

ST552: Linear Models and Variance Components

Mon/Wed 11:45am-1:00pm, Wed 1:30pm-2:45pm, SAS Hall 5270

Instructor: Dr Hua Zhou, hua_zhou@ncsu.edu

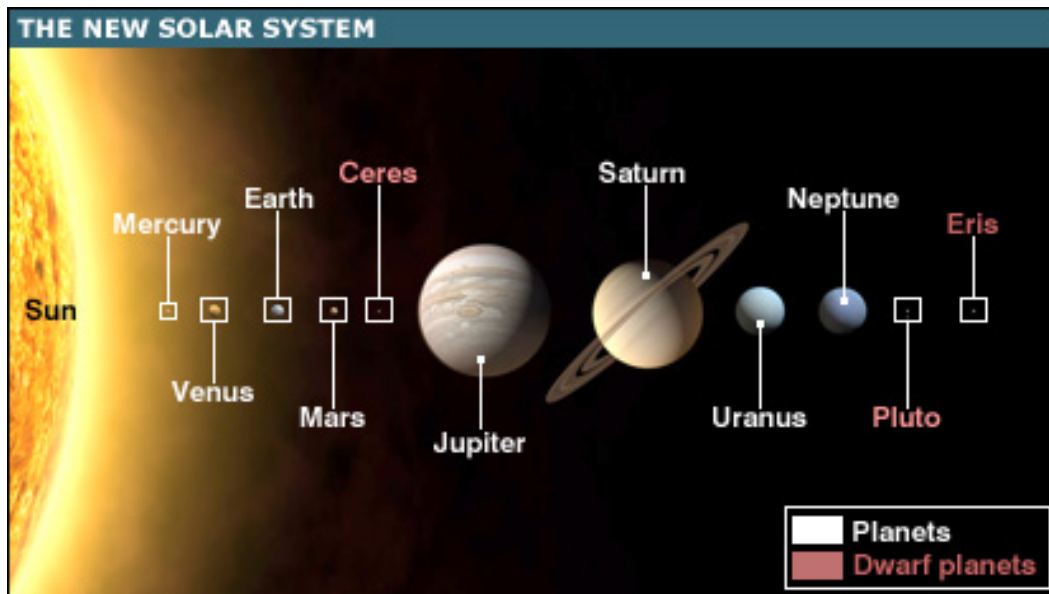
1 Lecture 1: Aug 21

Today

- Introduction
- Course logistics
- Read JM Appendix A and chapter 1
- Linear algebra review
- No afternoon session today

How Gauss became famous?





- 1801, *Dr Carl Friedrich Gauss*, 24; proved Fundamental Theorem of Algebra; wrote the book *Disquisitiones Arithmeticae*, which is still being studied today.
- 1801, Jan 1 - Feb 11 (41 days), astronomer Piazzi observed Ceres (a dwarf planet), which was then lost behind sun.
- 1801, Aug – Sep, futile search by top astronomers; Laplace claimed it unsolvable.
- 1801, Oct – Nov, young Gauss did calculations by *method of least squares*.
- 1801, Dec 31, astronomer von Zach relocated Ceres according to Gauss' calculation.
- 1802, *Summarische Übersicht der Bestimmung der Bahnen der beiden neuen Hauptplaneten angewandten Methoden*, considered the origin of linear algebra.
- 1807, Professor of Astronomy and Director of Göttingen Observatory in remainder of his life.
- 1809, *Theoria motus corporum coelestium in sectionibus conicis solem ambientum* (Theory of motion of the celestial bodies moving in conic sections around the Sun); birth of the Gaussian distribution, as an attempt to rationalize the method of least squares.

- 1810, Laplace consolidated importance of Gaussian distribution by proving the central limit theorem.
- 1829, Gauss-Markov Theorem. Under Gaussian error assumption (actually only uncorrelated and homoscedastic needed), least square solution is the best linear unbiased estimate (BLUE), i.e., it has the smallest variance and MSE among all linear unbiased estimators. Other estimators such as the James-Stein estimator may have smaller MSE, but they are *nonlinear*.

For more details

- <http://www.keplersdiscovery.com/Asteroid.html>
- Teets and Whitehead (1999)

ARTICLES

The Discovery of Ceres: How Gauss Became Famous

DONALD TEETS
KAREN WHITEHEAD
South Dakota School of Mines and Technology
Rapid City, SD 57701

“The Duke of Brunswick has discovered more in his country than a planet: a super-terrestrial spirit in a human body.”

These words, attributed to Laplace in 1801, refer to the accomplishment of Carl Friedrich Gauss in computing the orbit of the newly discovered planetoid *Ceres Ferdinandea* from extremely limited data. Indeed, although Gauss had already achieved some fame among mathematicians, it was his work on the Ceres orbit that “made Gauss a European celebrity—this a consequence of the popular appeal which astronomy has always enjoyed...” [2]. The story of Gauss’s work on this problem is a good one and is often told in biographical sketches of Gauss (e.g., [2], [3], [6]), but the mathematical details of how he solved the problem are invariably omitted from such historical works. We are left to wonder, how did he do it? *Just how did Gauss*

Gauss’ story

- Motivated by a real problem.
- Heuristic solution: method of least squares.
- Solution readily verifiable: Ceres was re-discovered!

- Algorithmic development: linear algebra, Gaussian elimination, FFT (fast Fourier transform).
- Theoretical justification: Gaussian distribution, Gauss-Markov theorem.

What is this course about?

This course focuses on the “theoretical” aspect of the method of least squares and statistical inference under the normal assumption. Read JM chapter 1 for a few specific examples of linear models.

A hierarchy of linear models

- The linear mean model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

where $\mathbf{E}(\mathbf{e}) = \mathbf{0}$. Only assumption is that errors have mean 0.

- Gauss-Markov model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

where $\mathbf{E}(\mathbf{e}) = \mathbf{0}$ and $\mathbf{Var}(\mathbf{e}) = \sigma^2\mathbf{I}$ (uncorrelated errors with constant variance).

- Aitken model or general linear model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

where $\mathbf{E}(\mathbf{e}) = \mathbf{0}$ and $\mathbf{Var}(\mathbf{e}) = \sigma^2\mathbf{V}$. \mathbf{V} is fixed and known.

- Linear models with joint normal errors: $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ or $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{V})$ with \mathbf{V} known.
- Variance components model: $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma_1^2\mathbf{V}_1 + \cdots + \sigma_r^2\mathbf{V}_r)$ with $\mathbf{V}_1, \dots, \mathbf{V}_r$ known.
- Multivariate linear model: $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$.
- Generalized linear models (GLMs). Logistic regression, probit regression, log-linear model (Poisson regression), ... Note the difference from the general linear model. GLMs are generalization of the *concept* of linear models. They are typically covered in a categorical data analysis course.

Syllabus

Check course website frequently for updates and announcements.

<http://hua-zhou.github.io/teaching/st552-2013fall/>

Lecture notes will be updated and posted after each lecture.

Vector and vector space

- A set of vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ are *linearly dependent* if there exist coefficients c_j such that $\sum_{j=1}^n c_j \mathbf{x}_j = \mathbf{0}$ and $\|\mathbf{c}\|_2 = \sum_{j=1}^n c_j^2 > 0$. They are *linearly independent* if $\sum_{j=1}^n c_j \mathbf{x}_j = \mathbf{0}$ implies $c_j = 0$ for all j .
- A *vector space* \mathcal{V} is a set of vectors that are closed under addition and scalar multiplication (closed under *axy* operation). Any vector space must contain the zero vector $\mathbf{0}$ (why?).
- A vector space \mathcal{V} is *generated* or *spanned* by a set of vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$, written as $\mathcal{V} = \text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, if any vector \mathbf{x} in the vector space is a linear combination of \mathbf{x}_i , $i = 1, \dots, n$.
- A set of linearly independent vectors that generate or span a space \mathcal{V} is called a *basis* of \mathcal{V} .
- Order and dimension. The *order* of a vectors space is simply the length of the vectors in that space. The *dimension* of a vector space is the maximum number of linearly independent vectors in that space.
- Two vector spaces \mathcal{V}_1 and \mathcal{V}_2 are *essentially disjoint* if the only element in $\mathcal{V}_1 \cap \mathcal{V}_2$ is the zero vector $\mathbf{0}$.
- If \mathcal{V}_1 and \mathcal{V}_2 are two vector spaces, then $\mathcal{V}_1 \cap \mathcal{V}_2$ are vector spaces.
- If \mathcal{V}_1 and \mathcal{V}_2 are two vector spaces of same order, then $\mathcal{V}_1 + \mathcal{V}_2 = \{\mathbf{v} : \mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2, \mathbf{v}_1 \in \mathcal{V}_1, \mathbf{v}_2 \in \mathcal{V}_2\}$ is a vector space. If furthermore \mathcal{V}_1 and \mathcal{V}_2 are essentially disjoint, the sum is called the *direct sum* and denoted by $\mathcal{V}_1 \oplus \mathcal{V}_2$.
- If \mathcal{V}_1 and \mathcal{V}_2 are two vector spaces, $\mathcal{V}_1 \cup \mathcal{V}_2$ is not necessarily a vector space. (Exercise: find a counter example).

- Affine spaces. Consider a system of m linear equations in variable $\mathbf{x} \in \mathbb{R}^n$

$$\begin{aligned} \mathbf{c}_1^T \mathbf{x} &= b_1 \\ &\vdots \\ \mathbf{c}_m^T \mathbf{x} &= b_m, \end{aligned}$$

where $\mathbf{c}_1, \dots, \mathbf{c}_m \in \mathbb{R}^n$ are linearly independent (and hence $m \leq n$). The set of solutions is called an *affine space*. The intersection of two affine spaces is an affine space (why?). If the zero vector $\mathbf{0}_n$ belongs to the affine space, i.e., $b_1 = \dots = b_m = 0$, then it is a vector space. Thus any affine space containing the origin $\mathbf{0}$ is a vector space, but other affine spaces are not vector spaces.

- If $m = 1$, the affine space is called a *hyperplane*. A hyperplane through the origin is an $(n - 1)$ -dimensional vector space.
- If $m = n - 1$, the affine space is a line. A line through the origin is a one-dimensional vector space.
- The mapping $\mathbf{x} \mapsto \mathbf{A}\mathbf{x} + \mathbf{b}$ is called an *affine function*. If $\mathbf{b} = \mathbf{0}$, it is called a *linear function*.

Orthogonality and orthogonalization

- Vector \mathbf{x}_1 is orthogonal to another vector \mathbf{x}_2 , denoted by $\mathbf{x}_1 \perp \mathbf{x}_2$, if $\langle \mathbf{x}_1, \mathbf{x}_2 \rangle = \mathbf{x}_1^T \mathbf{x}_2 = 0$. They are *orthonormal* if $\mathbf{x}_1 \perp \mathbf{x}_2$ and $\|\mathbf{x}_1\|_2 = \|\mathbf{x}_2\|_2 = 1$.
- The projection of a vector \mathbf{x}_2 onto the vector \mathbf{x}_1 is

$$\hat{\mathbf{x}}_2 = \frac{\langle \mathbf{x}_2, \mathbf{x}_1 \rangle}{\|\mathbf{x}_1\|_2^2} \mathbf{x}_1 = \langle \tilde{\mathbf{x}}_1, \mathbf{x}_2 \rangle \tilde{\mathbf{x}}_1,$$

where $\tilde{\mathbf{x}}_1 = \mathbf{x}_1 / \|\mathbf{x}_1\|_2$ is the normalized \mathbf{x}_1 vector.

- Gram-Schmidt transformation: orthonormalize two vectors \mathbf{x}_1 and \mathbf{x}_2 .

$$\begin{aligned} \tilde{\mathbf{x}}_1 &= \frac{1}{\|\mathbf{x}_1\|_2} \mathbf{x}_1 \\ \tilde{\mathbf{x}}_2 &= \frac{1}{\|\mathbf{x}_2 - \langle \tilde{\mathbf{x}}_1, \mathbf{x}_2 \rangle \tilde{\mathbf{x}}_1\|_2} (\mathbf{x}_2 - \langle \tilde{\mathbf{x}}_1, \mathbf{x}_2 \rangle \tilde{\mathbf{x}}_1) \end{aligned}$$

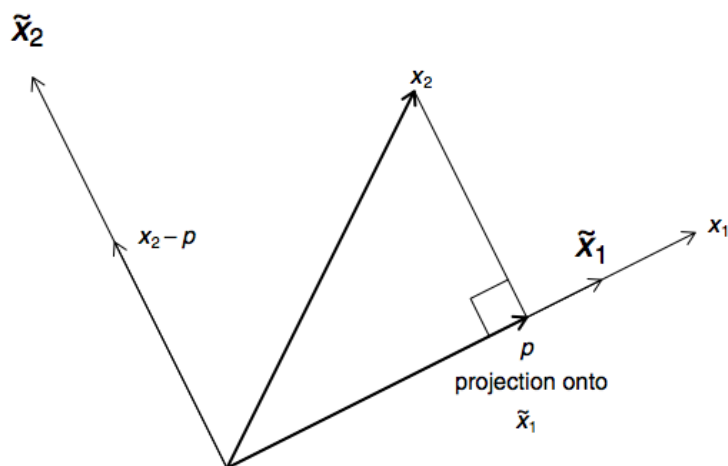


Fig. 2.2. Orthogonalization of x_1 and x_2

- Gram-Schmidt algorithm orthonormalizes a set of vectors.
- A set of nonzero, mutually orthogonal vectors are linearly independent. Proof by contradiction.
- Two vector spaces \mathcal{V}_1 and \mathcal{V}_2 are *orthogonal*, written $\mathcal{V}_1 \perp \mathcal{V}_2$, if each vector in \mathcal{V}_1 is orthogonal to every vector in \mathcal{V}_2 .
- The intersection of two orthogonal vector spaces consists only of the zero vector $\mathbf{0}$.
- If $\mathcal{V}_1 \perp \mathcal{V}_2$ and $\mathcal{V}_1 \oplus \mathcal{V}_2 = \mathbb{R}^n$, then \mathcal{V}_2 is the *orthogonal complement* of \mathcal{V}_1 . This is written as $\mathcal{V}_2 = \mathcal{V}_1^\perp$.

Rank

Assume $\mathbf{A} \in \mathbb{R}^{m \times n}$.

- $\text{rank}(\mathbf{A})$ is the maximum number of linearly independent rows of a matrix.
- Alternatively, $\text{rank}(\mathbf{A})$ is the maximum number of linearly independent columns of a matrix. (Exercise: show that these two definitions are equivalent. Hint: form the full rank partition of \mathbf{A} .)
- $\text{rank}(\mathbf{A}) \leq \min\{m, n\}$.

- A matrix is *full rank* if $\text{rank}(\mathbf{A}) = \min\{m, n\}$. It is *full row rank* if $\text{rank}(\mathbf{A}) = m$. It is *full column rank* if $\text{rank}(\mathbf{A}) = n$.
- A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *singular* if $\text{rank}(\mathbf{A}) < n$ and *non-singular* if $\text{rank}(\mathbf{A}) = n$.
- $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^T) = \text{rank}(\mathbf{A}^T \mathbf{A}) = \text{rank}(\mathbf{A} \mathbf{A}^T)$. (Exercise: show it.)
- $\text{rank}(\mathbf{A} \mathbf{B}) \leq \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})\}$. (Hint: Columns of $\mathbf{A} \mathbf{B}$ are spanned by columns of \mathbf{A} and rows of $\mathbf{A} \mathbf{B}$ are spanned by rows of \mathbf{B} .)
- $\text{rank}(\mathbf{A} \mathbf{B}) = \text{rank}(\mathbf{A})$ if \mathbf{B} is square and of full rank. More general, pre-multiplying by a matrix with full column rank or post-multiplying by a matrix with full row rank does not change rank. (Exercise: show it.)
- $\text{rank}(\mathbf{A} + \mathbf{B}) \leq \text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B})$.

$$\mathbf{A} + \mathbf{B} = (\mathbf{A} \ \mathbf{B}) \begin{pmatrix} \mathbf{I}_n \\ \mathbf{I}_n \end{pmatrix}.$$

- If $\mathbf{A} \mathbf{x} = \mathbf{0}_m$ for some $\mathbf{x} \neq \mathbf{0}_n$, then $\text{rank}(\mathbf{A}) \leq n - 1$.

2 Lecture 2: Aug 26

Last time

- How Gauss became famous? Method of least squares, Gauss-Markov theorem, Gaussian distribution, ...
- What do we mean by “linear models”? Different layers of assumptions: linear mean model, Gauss-Markov model, Aitken or general linear model, Gaussian assumption, variance components (linear mixed) model, multivariate linear model, GLM, ...
- Course logistics
- Linear algebra: vector and vector space, rank of a matrix

Today

- Continue review of linear algebra
- TA office hours are posted
- Homework 1 is posted, due Sep 4
- Reach JM chapter 2

Column space and null space

Assume $\mathbf{A} \in \mathbb{R}^{m \times n}$.

- The *column space* of a matrix \mathbf{A} , denoted by $\mathcal{C}(\mathbf{A})$, is the vector space (of order m) spanned by the columns of the matrix. Other names: *range* of \mathbf{A} , or the *manifold* of \mathbf{A} .
- The *null space* of a matrix \mathbf{A} , denoted by $\mathcal{N}(\mathbf{A})$, is the vector space (of order n) $\{\mathbf{y} : \mathbf{A}\mathbf{y} = \mathbf{0}\}$.
- $\dim(\mathcal{C}(\mathbf{A})) = r$ and $\dim(\mathcal{N}(\mathbf{A})) = n - r$, where $r = \text{rank}(\mathbf{A})$.
See JM Theorem A.1 for the proof.
Interpretation: “dimension of column space + dimension of null space = #

columns”

Misinterpretation: Columns space and null space are orthogonal complement to each other. They are of different orders in general! Next result gives the correct statement.

- $\mathcal{C}(\mathbf{A})$ and $\mathcal{N}(\mathbf{A}^T)$ are orthogonal complement to each other in \mathbb{R}^m : $\mathcal{C}(\mathbf{A}) = \mathcal{N}(\mathbf{A}^T)^\perp$.

Proof: (1) Orthogonality is trivial. (2) (Essentially disjoint) For any $\mathbf{v} \in \mathcal{C}(\mathbf{A}) \cap \mathcal{N}(\mathbf{A}^T)$, there exists a \mathbf{c} such that $\mathbf{A}\mathbf{c} = \mathbf{v}$ and $\mathbf{A}^T\mathbf{v} = \mathbf{0}_n$. Therefore $\mathbf{A}^T\mathbf{A}\mathbf{c} = \mathbf{A}^T\mathbf{v} = \mathbf{0}_n$ and thus $\mathbf{c}^T\mathbf{A}^T\mathbf{A}\mathbf{c} = 0$, implying that $\mathbf{v} = \mathbf{A}\mathbf{c} = \mathbf{0}$. Therefore $\mathcal{C}(\mathbf{A}) \cap \mathcal{N}(\mathbf{A}^T) = \{\mathbf{0}\}$. (3) Since $\dim(\mathcal{C}(\mathbf{A})) = \text{rank}(\mathbf{A})$ and $\dim(\mathcal{N}(\mathbf{A}^T)) = m - \text{rank}(\mathbf{A})$, we have $\mathcal{C}(\mathbf{A}) \cup \mathcal{N}(\mathbf{A}^T) = \mathbb{R}^m$.

- $\mathcal{C}(\mathbf{A}\mathbf{B}) \subset \mathcal{C}(\mathbf{A})$.
- If $\mathcal{C}(\mathbf{B}) \subset \mathcal{C}(\mathbf{A})$, then there exists a matrix \mathbf{C} such that $\mathbf{B} = \mathbf{A}\mathbf{C}$.
- $\mathcal{C}(\mathbf{A}\mathbf{A}^T) = \mathcal{C}(\mathbf{A})$ and thus $\text{rank}(\mathbf{A}\mathbf{A}^T) = \text{rank}(\mathbf{A})$. (Exercise: show it.)
“Multiplication by its transpose does not change the rank.”

Trace of a square matrix

Assume $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a square matrix.

- $\text{tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii}$ is the sum of diagonal entries.
- $\text{tr}(\mathbf{A}^T) = \text{tr}(\mathbf{A})$.
- $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$.
- $\text{tr}(c\mathbf{A}) = c\text{tr}(\mathbf{A})$, where c is a scalar.
- Invariance of trace function under cyclic permutation: $\text{tr}(\mathbf{A}\mathbf{B}) = \text{tr}(\mathbf{B}\mathbf{A})$ for $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B}^{n \times m}$. In general, $\text{tr}(\mathbf{A}_1 \cdots \mathbf{A}_k) = \text{tr}(\mathbf{A}_{j+1} \cdots \mathbf{A}_k \mathbf{A}_1 \cdots \mathbf{A}_j)$, $j = 1, \dots, k-1$, for matrices of compatible sizes.

Inner (dot) product between two matrices

- For $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$, $\langle \mathbf{A}, \mathbf{B} \rangle := \text{tr}(\mathbf{A}^T \mathbf{B}) = \text{tr}(\mathbf{B} \mathbf{A}^T)$. Analog of the inner product between two vectors $\langle \mathbf{x}, \mathbf{y} \rangle$.
- $\langle \mathbf{A}, \mathbf{B} \rangle = \langle \mathbf{A}^T, \mathbf{B}^T \rangle$.
- We say \mathbf{A} is orthogonal to \mathbf{B} if $\langle \mathbf{A}, \mathbf{B} \rangle = 0$.
- *Trace norm* (or *Frobenius norm*, or *Euclidean norm*) is the norm on the space of $m \times n$ matrices induced by the inner product: $\|\mathbf{A}\|_F = (\langle \mathbf{A}, \mathbf{A} \rangle)^{1/2} = \text{tr}(\mathbf{A}^T \mathbf{A})^{1/2} = (\sum_{i,j} a_{ij}^2)^{1/2}$.
- Cauchy-Schwartz: $\langle \mathbf{A}, \mathbf{B} \rangle \leq \langle \mathbf{A}, \mathbf{A} \rangle^{1/2} \langle \mathbf{B}, \mathbf{B} \rangle^{1/2}$.
Proof: Expand $\text{tr}((\mathbf{A} - x\mathbf{B})^T (\mathbf{A} - x\mathbf{B})) \geq 0$.

Matrix inverses

Assume $\mathbf{A} \in \mathbb{R}^{m \times n}$.

- The *Moore-Penrose inverse* of \mathbf{A} is a matrix $\mathbf{A}^+ \in \mathbb{R}^{n \times m}$ with following properties
 - (a) $\mathbf{A} \mathbf{A}^+ \mathbf{A} = \mathbf{A}$. (*Generalized inverse*, g_1 *inverse*, or *inner pseudo-inverse*)
 - (b) $\mathbf{A}^+ \mathbf{A} \mathbf{A}^+ = \mathbf{A}^+$. (*Outer pseudo-inverse*. Any g_1 inverse that satisfies this condition is called a g_2 *inverse*, or *reflexive generalized inverse* and is denoted by \mathbf{A}^* .)
 - (c) $\mathbf{A}^+ \mathbf{A}$ is symmetric.
 - (d) $\mathbf{A} \mathbf{A}^+$ is symmetric.
- \mathbf{A}^+ exists and is unique for any matrix \mathbf{A} .
- *Generalized inverse* (or g_1 *inverse*, denoted by \mathbf{A}^- or \mathbf{A}^g): property (a).
- g_2 *inverse* (denoted by \mathbf{A}^*): properties (a)+(b).
- *Moore-Penrose inverse* (denoted by \mathbf{A}^+): properties (a)+(b)+(c)+(d).
- If \mathbf{A} is square and full rank, then the generalized inverse is unique and denoted by \mathbf{A}^{-1} (inverse).

- How to find a generalized inverse (may be non-unique) for a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ of rank r ? Permute rows and columns to form the *full rank partitioning*

$$\mathbf{PAQ} = \begin{pmatrix} \mathbf{C}_{r \times r} & \mathbf{D}_{r \times (n-r)} \\ \mathbf{E}_{(m-r) \times r} & \mathbf{F}_{(m-r) \times (n-r)} \end{pmatrix},$$

where $\mathbf{C} \in \mathbb{R}^{r \times r}$ is of full rank. Then the matrix

$$\mathbf{Q} \begin{pmatrix} \mathbf{C}_{r \times r}^{-1} & \mathbf{0}_{r \times (m-r)} \\ \mathbf{0}_{(n-r) \times r} & \mathbf{0}_{(n-r) \times (m-r)} \end{pmatrix} \mathbf{P}$$

is a generalized inverse of \mathbf{A} . See JM Result A.11 (p247) for the argument.

- In practice, the Moore-Penrose inverse \mathbf{A}^+ is easily computed from the singular value decomposition (SVD) of \mathbf{A} .
- For any nonzero k , $(1/k)\mathbf{A}^-$ is a generalized inverse of $k\mathbf{A}$.
- $(\mathbf{A}^-)^T$ is a generalized inverse of \mathbf{A}^T .
- $\mathcal{C}(\mathbf{A}) = \mathcal{C}(\mathbf{AA}^-)$ and $\mathcal{C}(\mathbf{A}^T) = \mathcal{C}((\mathbf{A}^- \mathbf{A})^T)$.
 $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{AA}^-) = \text{rank}(\mathbf{A}^- \mathbf{A})$.

“Multiplication by generalized inverse does not change rank.”

Proof. We already know $\mathcal{C}(\mathbf{A}) \supset \mathcal{C}(\mathbf{AA}^-)$. Now since $\mathbf{A} = \mathbf{AA}^- \mathbf{A}$, we also have $\mathcal{C}(\mathbf{A}) \subset \mathcal{C}(\mathbf{A}^- \mathbf{A})$. □

- $\text{rank}(\mathbf{A}^-) \geq \text{rank}(\mathbf{A})$. “Generalized inverse has equal or a larger rank than original matrix.”

Proof. $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{AA}^- \mathbf{A}) \leq \text{rank}(\mathbf{AA}^-) = \text{rank}(\mathbf{A}^-)$. □

System of linear equations

$\mathbf{Ax} = \mathbf{b}$ where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$.

- When is there a solution? The following statements are equivalent.
 1. The linear system $\mathbf{Ax} = \mathbf{b}$ has a solution (*consistent*)
 2. $\mathbf{b} \in \mathcal{C}(\mathbf{A})$.

3. $\text{rank}(\mathbf{A}, \mathbf{b}) = \text{rank}(\mathbf{A})$.
4. $\mathbf{A}\mathbf{A}^{-}\mathbf{b} = \mathbf{b}$.

Proof. Equivalence between 1, 2 and 3 is trivial. 4 implies 1: apparently $\tilde{\mathbf{x}} = \mathbf{A}^{-}\mathbf{b}$ is a solution to $\mathbf{A}\mathbf{x} = \mathbf{b}$. 1 implies 4: if $\tilde{\mathbf{x}}$ is a solution, then $\mathbf{b} = \mathbf{A}\tilde{\mathbf{x}} = \mathbf{A}\mathbf{A}^{-}\mathbf{A}\tilde{\mathbf{x}} = \mathbf{A}\mathbf{A}^{-}\mathbf{b}$. \square

The last equivalence gives some intuition why \mathbf{A}^{-} is called an inverse.

- How to characterize the solution set? Let's first study the homogenous case $\mathbf{A}\mathbf{x} = \mathbf{0}$, which is always consistent (why?).
 $\tilde{\mathbf{x}}$ is a solution to $\mathbf{A}\mathbf{x} = \mathbf{0}$ if and only if

$$\tilde{\mathbf{x}} = (\mathbf{I}_n - \mathbf{A}^{-}\mathbf{A})\mathbf{q}$$

for some $\mathbf{q} \in \mathbb{R}^n$.

Proof. "If": Apparently $(\mathbf{I}_n - \mathbf{A}^{-}\mathbf{A})\mathbf{q}$ is a solution regardless value of \mathbf{q} since $\mathbf{A}(\mathbf{I}_n - \mathbf{A}^{-}\mathbf{A}) = \mathbf{A} - \mathbf{A} = \mathbf{0}_{m \times n}$. "Only if": If $\tilde{\mathbf{x}}$ is a solution, then $\tilde{\mathbf{x}} = (\mathbf{I}_n - \mathbf{A}^{-}\mathbf{A})\mathbf{q}$ by taking $\mathbf{q} = \tilde{\mathbf{x}}$. \square

- Rephrasing above result we have $\mathcal{N}(\mathbf{A}) = \mathcal{C}(\mathbf{I}_n - \mathbf{A}^{-}\mathbf{A})$.

- Now we study the inhomogeneous case.

If $\mathbf{A}\mathbf{x} = \mathbf{b}$ is consistent, then $\tilde{\mathbf{x}}$ is a solution to $\mathbf{A}\mathbf{x} = \mathbf{b}$ if and only if

$$\tilde{\mathbf{x}} = \mathbf{A}^{-}\mathbf{b} + (\mathbf{I}_n - \mathbf{A}^{-}\mathbf{A})\mathbf{q}$$

for some $\mathbf{q} \in \mathbb{R}^n$.

Interpretation: "a specific solution" + "a vector in the null space of \mathbf{A} ".

Proof.

$$\begin{aligned} \mathbf{A}\mathbf{x} = \mathbf{b} &\Leftrightarrow \mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{A}^{-}\mathbf{b} \Leftrightarrow \mathbf{A}(\mathbf{x} - \mathbf{A}^{-}\mathbf{b}) = \mathbf{0} \\ &\Leftrightarrow \mathbf{x} - \mathbf{A}^{-}\mathbf{b} = (\mathbf{I}_n - \mathbf{A}^{-}\mathbf{A})\mathbf{q} \Leftrightarrow \mathbf{x} = \mathbf{A}^{-}\mathbf{b} + (\mathbf{I}_n - \mathbf{A}^{-}\mathbf{A})\mathbf{q}. \end{aligned}$$

\square

- $\mathbf{A}\mathbf{x} = \mathbf{b}$ is consistent for *all* \mathbf{b} if and only if \mathbf{A} has full row rank.

Proof. “If”: $\dim(\mathcal{C}(\mathbf{A})) = \text{rank}(\mathbf{A}) = m$. Thus $\mathcal{C}(\mathbf{A}) = \mathbb{R}^m$ and contains any $\mathbf{b} \in \mathbb{R}^m$. “Only if”: $\mathbf{A}\mathbf{A}^-\mathbf{b} = \mathbf{b}$ for any \mathbf{b} . Take $\mathbf{b} = \mathbf{e}_i$ gives $\mathbf{A}\mathbf{A}^- = \mathbf{I}_m$. Thus $m \geq \text{rank}(\mathbf{A}) \geq \text{rank}(\mathbf{A}\mathbf{A}^-) = m$. \square

- If a system is consistent, its solution is unique if and only if \mathbf{A} has full column rank.

Proof. In previous proof, we see there is a one-to-one correspondence between the solution set to the inhomogenous system $\mathbf{A}\mathbf{x} = \mathbf{b}$ and the solution set to the homogeneous system $\mathbf{A}\mathbf{x} = \mathbf{0}$. Now \mathbf{A} has full column rank if and only if $\dim(\mathcal{N}(\mathbf{A})) = n - \dim(\mathcal{C}(\mathbf{A}^T)) = n - \text{rank}(\mathbf{A}) = n - n = 0$ if and only if there is a unique solution to $\mathbf{A}\mathbf{x} = \mathbf{0}$ if and only if there is a unique solution to $\mathbf{A}\mathbf{x} = \mathbf{b}$. \square

- If \mathbf{A} has full row and column rank, then \mathbf{A} is non-singular and the unique solution is $\mathbf{A}^{-1}\mathbf{b}$.

3 Lecture 3: Aug 28

Last time

- Linear algebra: column and null spaces $\mathcal{C}(\mathbf{A}) = \mathcal{N}(\mathbf{A}^T)^\perp$, trace (invariance under cyclic permutation), generalized matrix inverse $\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A}$, $\mathcal{C}(\mathbf{A}) = \mathcal{C}(\mathbf{A}\mathbf{A}^-)$, consistency of a linear system $\mathbf{A}\mathbf{A}^-\mathbf{b} = \mathbf{b}$, solution set of a consistent linear system $\mathbf{A}^-\mathbf{b} + (\mathbf{I}_n - \mathbf{A}^-\mathbf{A})\mathbf{q}$, $\mathcal{N}(\mathbf{A}) = \mathcal{C}(\mathbf{I}_n - \mathbf{A}^-\mathbf{A})$, $\mathbf{A}\mathbf{x} = \mathbf{b}$ is consistent for all \mathbf{b} if and only if \mathbf{A} has full row rank, A consistent system $\mathbf{A}\mathbf{x} = \mathbf{b}$ has a unique solution if and only if \mathbf{A} has full column rank, ...

Today

- Idempotent matrix, projection and orthogonal projection
- Method of least squares
- Announcement
 - Change of TA office hours
 - Pre-lecture notes
 - Questions on homework 1?
 - Make use of comments on course webpages

Idempotent matrix

Assume $\mathbf{A} \in \mathbb{R}^{n \times n}$.

- A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *idempotent* if and only if $\mathbf{A}^2 = \mathbf{A}$.
- Any idempotent matrix \mathbf{A} is a generalized inverse of itself.
- The only idempotent matrix of full rank is \mathbf{I} .

Proof. Since \mathbf{A} has full rank, the inverse \mathbf{A}^{-1} exists. Then $\mathbf{A} = \mathbf{A}^{-1}\mathbf{A}\mathbf{A} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$. \square

Interpretation: all idempotent matrices are singular except the identity matrix.

- \mathbf{A} is idempotent if and only if \mathbf{A}^T is idempotent if and only if $\mathbf{I}_n - \mathbf{A}$ is idempotent.
- If $\mathbf{A}^2 = k\mathbf{A}$ for some nonzero scalar k , then $\text{tr}(\mathbf{A}) = k \text{rank}(\mathbf{A})$.

Proof. Suppose $\text{rank}(\mathbf{A}) = r$. Choose a set of r linearly independent columns and form matrix $\mathbf{B} \in \mathbb{R}^{n \times r}$, which has full column rank. There exists matrix $\mathbf{C} \in \mathbb{R}^{r \times n}$ such that $\mathbf{A} = \mathbf{BC}$. \mathbf{C} must have full row rank, otherwise \mathbf{A} cannot have rank r . Then

$$\mathbf{B}(\mathbf{CB})\mathbf{C} = \mathbf{A}^2 = k\mathbf{A} = k\mathbf{BC} = \mathbf{B}(k\mathbf{I}_r)\mathbf{C}.$$

Since \mathbf{B} has full column rank, we must have $(\mathbf{CB})\mathbf{C} = (k\mathbf{I}_r)\mathbf{C}$ (why?). Since \mathbf{C} has full row rank, we must have $\mathbf{CB} = k\mathbf{I}_r$. Then $\text{tr}(\mathbf{A}) = \text{tr}(\mathbf{BC}) = \text{tr}(\mathbf{CB}) = \text{tr}(k\mathbf{I}_r) = kr = k \text{rank}(\mathbf{A})$. \square

- For any idempotent matrix \mathbf{A} , $\text{tr}(\mathbf{A}) = \text{rank}(\mathbf{A})$.
- For any idempotent matrix \mathbf{A} ,

$$\text{rank}(\mathbf{I}_n - \mathbf{A}) = \text{tr}(\mathbf{I}_n - \mathbf{A}) = n - \text{rank}(\mathbf{A}).$$

So $\text{rank}(\mathbf{A}) + \text{rank}(\mathbf{I} - \mathbf{A}) = n$ for any idempotent matrix \mathbf{A} .

- For a *general* matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, the matrices $\mathbf{A}^- \mathbf{A}$ and $\mathbf{A} \mathbf{A}^-$ are idempotent and

$$\begin{aligned} \text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^- \mathbf{A}) &= \text{rank}(\mathbf{A} \mathbf{A}^-) = \text{tr}(\mathbf{A}^- \mathbf{A}) = \text{tr}(\mathbf{A} \mathbf{A}^-) \\ \text{rank}(\mathbf{I}_n - \mathbf{A}^- \mathbf{A}) &= \text{tr}(\mathbf{I}_n - \mathbf{A}^- \mathbf{A}) = n - \text{rank}(\mathbf{A}) \\ \text{rank}(\mathbf{I}_m - \mathbf{A} \mathbf{A}^-) &= \text{tr}(\mathbf{I}_m - \mathbf{A} \mathbf{A}^-) = m - \text{rank}(\mathbf{A}). \end{aligned}$$

- $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times m}$. Then the following statements are equivalent.
 1. \mathbf{B} is a generalized inverse of \mathbf{A} .
 2. \mathbf{AB} is idempotent and $\text{rank}(\mathbf{AB}) = \text{rank}(\mathbf{A})$.
 3. \mathbf{BA} is idempotent and $\text{rank}(\mathbf{BA}) = \text{rank}(\mathbf{A})$.

Proof. Previous result shows 1 implies 2. We need to show 2 implies 1. If $\text{rank}(\mathbf{AB}) = \text{rank}(\mathbf{A})$, then $\mathcal{C}(\mathbf{A}) = \mathcal{C}(\mathbf{AB})$ and $\mathbf{A} = \mathbf{ABT}$ for some matrix \mathbf{T} . Thus by idempotency of \mathbf{AB} , $\mathbf{ABA} = \mathbf{ABABT} = \mathbf{ABT} = \mathbf{A}$. Equivalence between 1 and 3 is shown in a similar way. \square

Projection

- A matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ is a *projection* onto a vector space \mathcal{V} if and only if (a) \mathbf{P} is idempotent, (b) $\mathbf{P}\mathbf{x} \in \mathcal{V}$ for any $\mathbf{x} \in \mathbb{R}^n$, and (c) $\mathbf{P}\mathbf{z} = \mathbf{z}$ for any $\mathbf{z} \in \mathcal{V}$.
- Any idempotent matrix \mathbf{P} is a projection onto its own column space $\mathcal{C}(\mathbf{P})$.

Proof. Property (a) is free. Property (b) is trivial since $\mathbf{P}\mathbf{x} \in \mathcal{C}(\mathbf{P})$ for any \mathbf{x} . For property (c), note $\mathbf{P}\mathbf{P} = \mathbf{P}$ says $\mathbf{P}\mathbf{p}_i = \mathbf{p}_i$ for each column \mathbf{p}_i of \mathbf{P} . Thus $\mathbf{P}\mathbf{z} = \mathbf{z}$ for any $\mathbf{z} \in \mathcal{C}(\mathbf{P})$. \square

- $\mathbf{A}\mathbf{A}^-$ is a projection onto the column space $\mathcal{C}(\mathbf{A})$.

This gives a recipe for finding projections onto the column space of a matrix. Recall in last class we showed $\mathcal{C}(\mathbf{A}\mathbf{A}^-) = \mathcal{C}(\mathbf{A})$. Therefore the proof is trivial.

Direct proof. (a) Idempotent: $\mathbf{A}\mathbf{A}^-\mathbf{A}\mathbf{A}^- = \mathbf{A}\mathbf{A}^-$ by definition of generalized inverse. (b) $\mathbf{A}\mathbf{A}^-\mathbf{v} = \mathbf{A}(\mathbf{A}^-\mathbf{v}) \in \mathcal{C}(\mathbf{A})$ for any \mathbf{v} . (c) Let $\mathbf{z} \in \mathcal{C}(\mathbf{A})$, then $\mathbf{z} = \mathbf{A}\mathbf{c}$ for some \mathbf{c} . Thus $\mathbf{A}\mathbf{A}^-\mathbf{z} = \mathbf{A}\mathbf{A}^-\mathbf{A}\mathbf{c} = \mathbf{A}\mathbf{c} = \mathbf{z}$. \square

- $\mathbf{I}_n - \mathbf{A}^-\mathbf{A}$ is a projection onto the null space $\mathcal{N}(\mathbf{A})$.

This gives a recipe for finding projection onto the null space of a matrix. Recall in last class we showed $\mathcal{C}(\mathbf{I}_n - \mathbf{A}^-\mathbf{A}) = \mathcal{N}(\mathbf{A})$. Therefore the proof is trivial.

Direct proof. (a) $\mathbf{I}_n - \mathbf{A}^-\mathbf{A}$ is idempotent because $\mathbf{A}^-\mathbf{A}$ is idempotent. (b) For any \mathbf{x} , $\mathbf{A}(\mathbf{I}_n - \mathbf{A}^-\mathbf{A})\mathbf{x} = (\mathbf{A} - \mathbf{A})\mathbf{x} = \mathbf{0}$. That is $\mathbf{x} \in \mathcal{N}(\mathbf{A})$. (c) For any $\mathbf{z} \in \mathcal{N}(\mathbf{A})$, $(\mathbf{I}_n - \mathbf{A}^-\mathbf{A})\mathbf{z} = \mathbf{z} - \mathbf{A}^-\mathbf{A}\mathbf{z} = \mathbf{z} - \mathbf{0} = \mathbf{z}$. \square

- In general, the projections onto a vector space are not unique.

- JM Example A.7. Let $\mathbf{A} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \in \mathbb{R}^{2 \times 1}$. We first find generalized inverse $\mathbf{G} = (u, v) \in \mathbb{R}^{1 \times 2}$ of \mathbf{A} . Definition of generalized inverse requires $\mathbf{A}\mathbf{G}\mathbf{A} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} (u, v) \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \mathbf{A} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$, i.e., $v = (1 - u)/2$. Thus $\mathbf{G}_u = (u, (1 - u)/2)$ is a generalized inverse of \mathbf{A} for any value of u . Thus

$$\mathbf{A}\mathbf{G}_u = \begin{pmatrix} 1 \\ 2 \end{pmatrix} (u, (1 - u)/2) = \begin{pmatrix} u & (1 - u)/2 \\ 2u & 1 - u \end{pmatrix}$$

is a projection onto $\mathcal{C}(\mathbf{A})$ for any u . Taking $u = 0$ gives a projection

$$\mathbf{A}\mathbf{G}_0 = \begin{pmatrix} 0 & 1/2 \\ 0 & 1 \end{pmatrix}.$$

$\mathbf{A}\mathbf{G}_0\mathbf{x} = \begin{pmatrix} x_2/2 \\ x_2 \end{pmatrix}$ for any $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$. That is $\mathbf{A}\mathbf{G}_0$ projects in \mathbb{R}^2 points vertically to the line $\mathcal{C}(\mathbf{A})$. Taking $u = 1$ gives

$$\mathbf{A}\mathbf{G}_1 = \begin{pmatrix} 1 & 0 \\ 2 & 0 \end{pmatrix},$$

which projects points in \mathbb{R}^2 horizontally to the line. What if we require the projection $\mathbf{A}\mathbf{G}_u$ to be symmetric? Then $2u = (1 - u)/2$ suggests $u = 1/5$ and

$$\mathbf{A}\mathbf{G}_{1/5} = \begin{pmatrix} 1/5 & 2/5 \\ 2/5 & 4/5 \end{pmatrix},$$

which projects points in \mathbb{R}^2 to the closest point on the line. This is an instance of *orthogonal projection*.

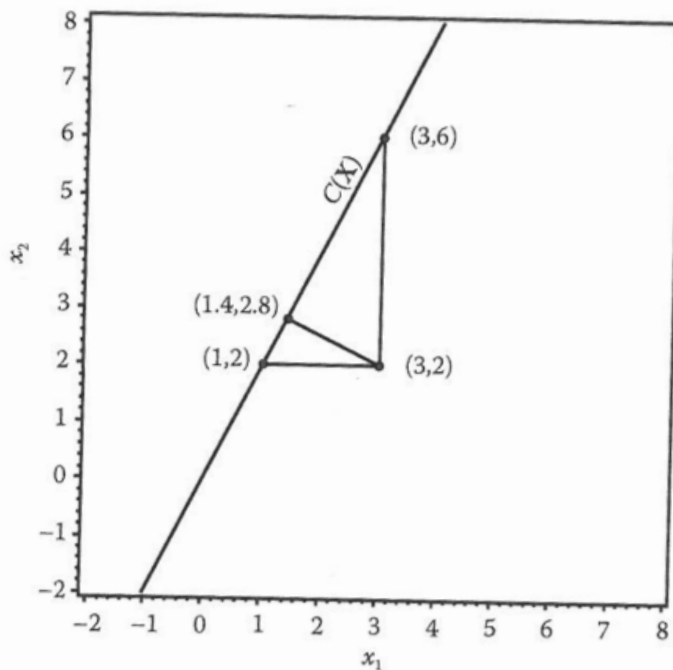


Figure A.1: Three projections onto a column space.

Orthogonal projection

- A symmetric, idempotent matrix \mathbf{P} that projects onto a vector space \mathcal{V} is unique.

Proof. Let \mathbf{P} and \mathbf{Q} be any two symmetric, idempotent matrices that project to the same vector space. \mathbf{Q} projects onto $\mathcal{C}(\mathbf{Q})$, thus \mathbf{P} projects onto $\mathcal{C}(\mathbf{Q})$ too. Therefore $\mathbf{P}\mathbf{q}_i = \mathbf{q}_i$ for each column \mathbf{q}_i of \mathbf{Q} . That is $\mathbf{P}\mathbf{Q} = \mathbf{Q}$. Similarly $\mathbf{Q}\mathbf{P} = \mathbf{P}$. Thus

$$\begin{aligned} \langle \mathbf{P} - \mathbf{Q}, \mathbf{P} - \mathbf{Q} \rangle &= \text{tr}((\mathbf{P} - \mathbf{Q})^T(\mathbf{P} - \mathbf{Q})) \\ &= \text{tr}(\mathbf{P}^T\mathbf{P} - \mathbf{P}^T\mathbf{Q} - \mathbf{Q}^T\mathbf{P} + \mathbf{Q}^T\mathbf{Q}) \\ &= \text{tr}(\mathbf{P}) - \text{tr}(\mathbf{P}\mathbf{Q}) - \text{tr}(\mathbf{Q}\mathbf{P}) + \text{tr}(\mathbf{Q}) \\ &= \text{tr}(\mathbf{P}) - \text{tr}(\mathbf{Q}) - \text{tr}(\mathbf{P}) + \text{tr}(\mathbf{Q}) \\ &= 0. \end{aligned}$$

\mathbf{P} and \mathbf{Q} must be equal. □

- Any symmetric, idempotent matrix is called an *orthogonal projection*. The term orthogonal comes from the fact that for any vector \mathbf{y} , the residual $\mathbf{y} - \mathbf{P}\mathbf{y}$ is orthogonal to the vector space \mathbf{P} is projecting onto. That is, for any $\mathbf{v} \in \mathcal{V}$, $\langle \mathbf{y} - \mathbf{P}\mathbf{y}, \mathbf{v} \rangle = \mathbf{y}^T(\mathbf{I}_n - \mathbf{P}^T)\mathbf{v} = \mathbf{y}^T(\mathbf{I}_n - \mathbf{P})\mathbf{v} = \mathbf{y}^T(\mathbf{v} - \mathbf{v}) = 0$.
- Many books use the term “projection” in the sense of orthogonal projection.
- If a symmetric, idempotent matrix \mathbf{P} projects onto \mathcal{V} , then $\mathbf{I} - \mathbf{P}$ projects onto the orthogonal complement \mathcal{V}^\perp .

Proof. Since \mathbf{P} projects onto $\mathcal{C}(\mathbf{P})$, we need to show $\mathbf{I} - \mathbf{P}$ projects onto $\mathcal{C}(\mathbf{P})^\perp = \mathcal{N}(\mathbf{P}^T) = \mathcal{N}(\mathbf{P})$. But \mathbf{P} is a generalized inverse of itself and thus $\mathbf{I} - \mathbf{P} = \mathbf{I} - \mathbf{P}\mathbf{P}$ projects to the $\mathcal{C}(\mathbf{I} - \mathbf{P}\mathbf{P}) = \mathcal{N}(\mathbf{P})$. Symmetry is trivial. □

Method of least squares

Given $\mathbf{y} \in \mathbb{R}^n$ and $\mathbf{X} \in \mathbb{R}^{n \times p}$, we want to find a vector $\mathbf{b} \in \mathbb{R}^p$ such that $\mathbf{X}\mathbf{b}$ approximates \mathbf{y} well. In other words, we want to approximate $\mathbf{y} \in \mathbb{R}^n$ by a linear combination of vectors $\mathbf{x}_j \in \mathbb{R}^n$, $j = 1, \dots, p$. Note the method of least squares only concerns approximation and has nothing to do with randomness and estimation.

- Suppose the linear system $\mathbf{X}\mathbf{b} = \mathbf{y}$ is consistent, we just need to find a solution to the linear system, which takes the general form $\mathbf{X}^{-}\mathbf{y} + (\mathbf{I}_p - \mathbf{X}^{-}\mathbf{X})\mathbf{q}$, $\mathbf{q} \in \mathbb{R}^p$. Recall that $\mathbf{I}_p - \mathbf{X}^{-}\mathbf{X}$ is a projection onto the null space $\mathcal{N}(\mathbf{X})$.

- Suppose the linear system is inconsistent. The *method of least squares* (due to Gauss) seeks \mathbf{b} that minimizes the Euclidean norm of the residual vector

$$Q(\mathbf{b}) = \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 = (\mathbf{y} - \mathbf{X}\mathbf{b})^T(\mathbf{y} - \mathbf{X}\mathbf{b}) = \mathbf{y}^T\mathbf{y} - 2\mathbf{b}^T\mathbf{X}^T\mathbf{y} + \mathbf{b}^T\mathbf{X}^T\mathbf{X}\mathbf{b}. \quad (1)$$

- Normal equation. To find the minimum, we take derivative and set the gradient to $\mathbf{0}$

$$\nabla Q(\mathbf{b}) = -2\mathbf{X}^T\mathbf{y} + 2\mathbf{X}^T\mathbf{X}\mathbf{b} = \mathbf{0}_p.$$

This leads to the *normal equation*

$$\mathbf{X}^T\mathbf{X}\mathbf{b} = \mathbf{X}^T\mathbf{y}. \quad (2)$$

- Is there a solution to the normal equation?

Normal equation is always consistent and thus admits at least one solution $\hat{\mathbf{b}}$. (HW1: show that $\mathcal{C}(\mathbf{X}^T) = \mathcal{C}(\mathbf{X}^T\mathbf{X})$.)

- Is the solution to the normal equation the minimizer to the least squares criterion?

Any solution $\hat{\mathbf{b}}$ to the normal equation (2) minimizes the least squares criterion (1).

Optimization argument: Any stationarity point (points with zero gradient vector) of a convex function is a global minimum. Now the least squares criterion is convex because the Hessian $\nabla^2 Q(\mathbf{b}) = \mathbf{X}^T\mathbf{X}$ is positive semidefinite. Therefore any solution to the normal equation is a stationarity point and thus a global minimum.

Direct argument: Let $\hat{\mathbf{b}}$ be a solution to the normal equation. For arbitrary

$\mathbf{b} \in \mathbb{R}^p$,

$$\begin{aligned}
Q(\mathbf{b}) - Q(\hat{\mathbf{b}}) &= -2(\mathbf{b} - \hat{\mathbf{b}})^T \mathbf{X}^T \mathbf{y} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b} - \hat{\mathbf{b}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{b}} \\
&= -2(\mathbf{b} - \hat{\mathbf{b}})^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{b}} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b} - \hat{\mathbf{b}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{b}} \\
&= \hat{\mathbf{b}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{b}} + \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{b} - 2\mathbf{b}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{b}} \\
&= (\mathbf{b} - \hat{\mathbf{b}})^T \mathbf{X}^T \mathbf{X} (\mathbf{b} - \hat{\mathbf{b}}) \\
&= \|\mathbf{X}(\mathbf{b} - \hat{\mathbf{b}})\|_2^2 \\
&\geq 0.
\end{aligned}$$

- The direct argument also reveals that the *fitted values* $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{b}}$ is invariant to the choice of the solution to the normal equation.
- Now we know the normal equation is always consistent and we want to find solution(s). In general the solution can be represented as

$$\hat{\mathbf{b}}(\mathbf{q}) = (\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{y} + [\mathbf{I}_p - (\mathbf{X}^T \mathbf{X})^{-} (\mathbf{X}^T \mathbf{X})] \mathbf{q}, \quad (3)$$

where $\mathbf{q} \in \mathbb{R}^p$ is arbitrary. One specific solution is

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{y}$$

with corresponding fitted values

$$\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{y}.$$

- When is the least squares solution unique?

The least squares solution is unique if and only if \mathbf{X} has full column rank. The solution is given by $\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.

Proof. The solution to normal equation is unique if and only if $\mathbf{X}^T \mathbf{X}$ has full (column) rank. Therefore $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{X}^T \mathbf{X}) = p$. \square

- $(\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T$ is a generalized inverse of \mathbf{X} .

Proof. By definition of generalized inverse, $\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{X} = \mathbf{X}^T \mathbf{X}$. Rearranging gives

$$\mathbf{X}^T \mathbf{X} [(\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{X} - \mathbf{I}_p] = \mathbf{0}_{p \times p}.$$

Each column of the matrix $\mathbf{X}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} - \mathbf{I}_p]$ is in $\mathcal{N}(\mathbf{X}^T)$ and also in $\mathcal{C}(\mathbf{X})$. Since $\mathcal{C}(\mathbf{X}) \cap \mathcal{N}(\mathbf{X}^T) = \{\mathbf{0}\}$, that column must be $\mathbf{0}$. Therefore $\mathbf{X}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} - \mathbf{I}_p] = \mathbf{0}_{n \times p}$. That is $\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} = \mathbf{X}$. \square

- The least squares solution has same format (3) regardless the consistency of the system $\mathbf{X}\mathbf{b} = \mathbf{y}$, since we can replace $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is a version of \mathbf{X}^- from the preceding result.

4 Lecture 4: Sep 4

Announcement

- HW1 is due today @ 11:59pm; TA is available answering questions in the afternoon session
- Read JM chapter 3
- Instructor out of town next week; no classes.

Last time

- Idempotent matrix
- Projection
 - any idempotent matrix \mathbf{P} projects onto $\mathcal{C}(\mathbf{P})$
 - $\mathbf{A}\mathbf{A}^-$ projects onto $\mathcal{C}(\mathbf{A})$
 - $\mathbf{I} - \mathbf{A}^- \mathbf{A}$ projects onto $\mathcal{N}(\mathbf{A})$
- Orthogonal projection (symmetric idempotent matrix): uniqueness and orthogonality; we will see $\mathbf{A}(\mathbf{A}^T \mathbf{A})^- \mathbf{A}^T$ is the orthogonal projection onto $\mathcal{C}(\mathbf{A})$
- Method of least squares:
 - least squares criterion: $Q(\mathbf{b}) = \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2$
 - normal equation: $\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{y}$
 - normal equation is always consistent
 - any solution to the normal equation is a least squares solution (that minimizes the least squares criterion)
 - least squares solution takes the general form

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{y} + (\mathbf{I}_p - (\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{X}) \mathbf{q},$$

where $\mathbf{q} \in \mathbb{R}^p$ is arbitrary

- the least squares solution is unique if and only if \mathbf{X} has full column rank
- $(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T$ is a generalized inverse of \mathbf{X} ; the least squares solution is the same whether the system $\mathbf{X}\mathbf{b} = \mathbf{y}$ is consistent or not

Today

- Geometry of least squares solution
- Gram-Schmidt orthogonalization
- Reparameterizations
- linear mean model
- random vectors
- estimable functions

Geometry of the least squares solution

- $\mathbf{X}(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T$ is *the* orthogonal projection onto $\mathcal{C}(\mathbf{X})$.

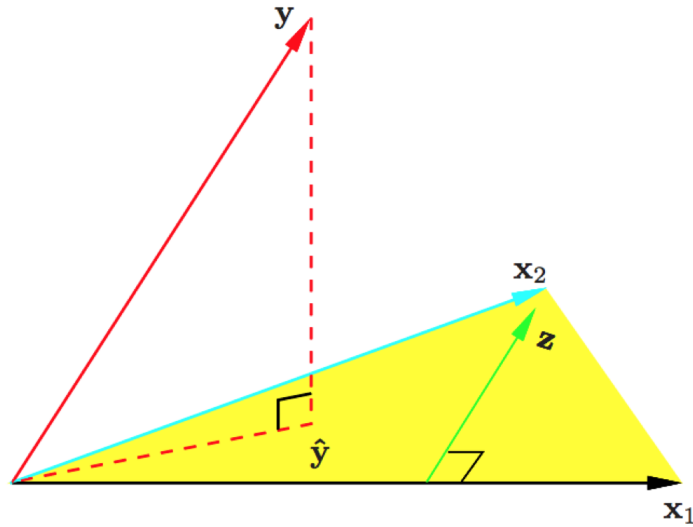
Proof. We showed that $(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T$ is a generalized inverse of \mathbf{X} in HW1. Therefore $\mathbf{X}(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T$ is a projection onto $\mathcal{C}(\mathbf{X})$. We only need to show the symmetry, which follows from the fact that transpose of $(\mathbf{X}^T \mathbf{X})^-$ is a generalized inverse of $(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$. \square

- Since orthogonal projection is unique, $\mathbf{X}(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T$ is invariant to the choice of the generalized inverse $(\mathbf{X}^T \mathbf{X})^-$ and thus can be denoted by $\mathbf{P}_\mathbf{X}$.
- Whichever least squares solution $\hat{\mathbf{b}}$ we use, the fitted values $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{b}}$ is the same since

$$\begin{aligned}\hat{\mathbf{y}} &= \mathbf{X}\hat{\mathbf{b}} \\ &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{y} + \mathbf{X}(\mathbf{I}_p - (\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{X})\mathbf{q} \\ &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{y} \\ &= \mathbf{P}_\mathbf{X} \mathbf{y}\end{aligned}$$

and the orthogonal projection is unique.

- Geometry: The fitted value from the least squares solution $\hat{\mathbf{y}} = \mathbf{P}_\mathbf{X} \mathbf{y}$ is the orthogonal projection of the response vector \mathbf{y} onto the column space $\mathcal{C}(\mathbf{X})$.



- $I_n - P_X$ is the orthogonal projection onto $\mathcal{N}(X^T)$.
- Decomposition of \mathbf{y} :

$$\mathbf{y} = P_X \mathbf{y} + (I_n - P_X) \mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{e}},$$

where $\hat{\mathbf{y}} \perp \hat{\mathbf{e}}$ and

$$\|\mathbf{y}\|_2^2 = \|\hat{\mathbf{y}}\|_2^2 + \|\hat{\mathbf{e}}\|_2^2.$$

The Pythagorean theorem follows from the orthogonality between $\mathcal{C}(X)$ and $\mathcal{N}(X^T)$ since $\|\mathbf{y}\|_2^2 = \mathbf{y}^T \mathbf{y} = (\hat{\mathbf{y}} + \hat{\mathbf{e}})^T (\hat{\mathbf{y}} + \hat{\mathbf{e}}) = \hat{\mathbf{y}}^T \hat{\mathbf{y}} + \hat{\mathbf{y}}^T \hat{\mathbf{e}} + \hat{\mathbf{e}}^T \hat{\mathbf{y}} + \hat{\mathbf{e}}^T \hat{\mathbf{e}} = \|\hat{\mathbf{y}}\|_2^2 + \|\hat{\mathbf{e}}\|_2^2$.

- Example: simple linear regression (predict y_i from intercept and one predictor: $y_i \approx b_0 + x_i b_1$).

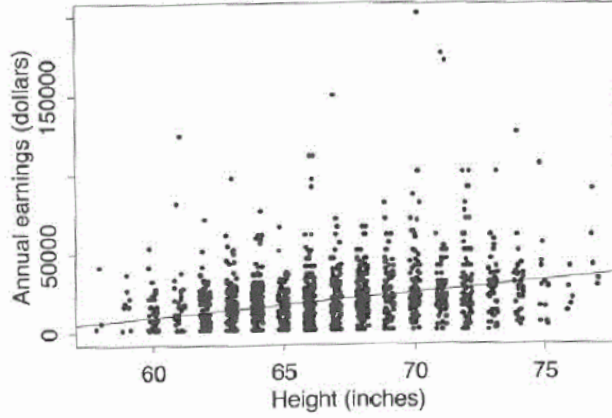


Fig. 4.3 Earnings vs. height for a random sample of adult Americans in 1990. The heights have been jittered slightly so that the points do not overlap.

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}.$$

Assume (x_1, \dots, x_n) is not a constant vector. The Gramian matrix is

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{pmatrix}$$

and its inverse is

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{n \sum_i (x_i - \bar{x})^2} \begin{pmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & n \end{pmatrix}.$$

The (unique) least squares solution is

$$\begin{aligned} \hat{\mathbf{b}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \frac{1}{n \sum_i (x_i - \bar{x})^2} \begin{pmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & n \end{pmatrix} \begin{pmatrix} \mathbf{1}_n^T \\ \mathbf{x}^T \end{pmatrix} \mathbf{y} \\ &= \frac{1}{n \sum_i (x_i - \bar{x})^2} \begin{pmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & n \end{pmatrix} \begin{pmatrix} \sum_i y_i \\ \sum_i x_i y_i \end{pmatrix} \\ &= \frac{1}{n \sum_i (x_i - \bar{x})^2} \begin{pmatrix} (\sum_i x_i^2)(\sum_i y_i) - (\sum_i x_i)(\sum_i x_i y_i) \\ n \sum_i x_i y_i - (\sum_i x_i)(\sum_i y_i) \end{pmatrix} \\ &= \begin{pmatrix} \bar{y} - \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} \bar{x} \\ \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} \end{pmatrix}. \end{aligned}$$

The fitted values are

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{b}} = (\bar{y} - \hat{b}_1\bar{x})\mathbf{1}_n + \hat{b}_1\mathbf{x} = \bar{y}\mathbf{1}_n - \hat{b}_1(\mathbf{x} - \bar{x}\mathbf{1}_n).$$

That is $\hat{y}_i = \bar{y} + \hat{b}_1(x_i - \bar{x})$. The residuals are $\hat{e}_i = (y_i - \bar{y}) - \hat{b}_1(x_i - \bar{x})$. And

$$\begin{aligned} \langle \hat{\mathbf{y}}, \hat{\mathbf{e}} \rangle &= \sum_i \hat{y}_i \hat{e}_i \\ &= \sum_i \left[\bar{y} + \hat{b}_1(x_i - \bar{x}) \right] \left[(y_i - \bar{y}) - \hat{b}_1(x_i - \bar{x}) \right] \\ &= \hat{b}_1 \left[\sum_i (x_i - \bar{x})(y_i - \bar{y}) - \hat{b}_1 \sum_i (x_i - \bar{x})^2 \right] \\ &= 0. \end{aligned}$$

Gram-Schmidt and QR decomposition

- Let $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$ be two non-zero, linearly independent vectors. Then $\mathbf{P}_{\mathcal{C}\{\mathbf{x}_1\}} = \mathbf{x}_1(\mathbf{x}_1^T \mathbf{x}_1)^{-1} \mathbf{x}_1^T = \mathbf{x}_1 \mathbf{x}_1^T / \|\mathbf{x}_1\|_2^2$ and

$$\mathbf{x}_2 - \mathbf{P}_{\mathbf{x}_1} \mathbf{x}_2 = \mathbf{x}_2 - \frac{\langle \mathbf{x}_2, \mathbf{x}_1 \rangle}{\|\mathbf{x}_1\|_2^2} \mathbf{x}_1$$

is orthogonal to \mathbf{x}_1 . This is the Gram-Schmidt orthogonalization of two vectors.

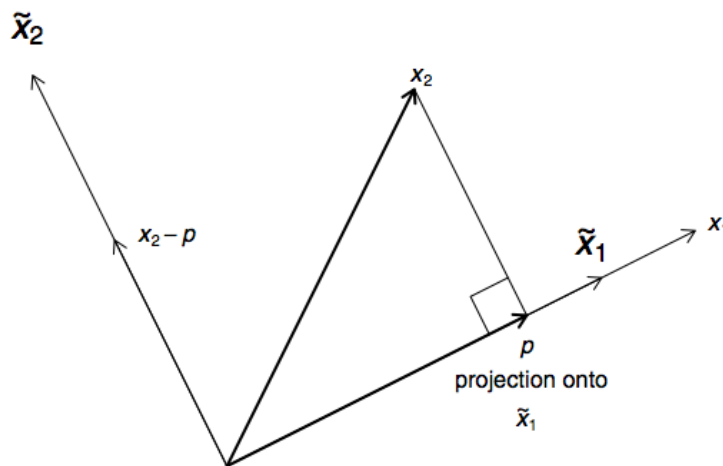


Fig. 2.2. Orthogonalization of \mathbf{x}_1 and \mathbf{x}_2

- Gram-Schmidt algorithm orthonormalizes a set of non-zero, *linearly independent* vectors $\mathbf{x}_1, \dots, \mathbf{x}_p$. Initialize $\mathbf{q}_1 = \mathbf{x}_1 / \|\mathbf{x}_1\|_2$; then for $j = 2, \dots, p$,

$$\begin{aligned}\mathbf{v}_j &= \mathbf{x}_j - \mathbf{P}_{\mathcal{C}\{\mathbf{q}_1, \dots, \mathbf{q}_{j-1}\}} \mathbf{x}_j = \mathbf{x}_j - \sum_{k=1}^{j-1} \langle \mathbf{x}_j, \mathbf{q}_k \rangle \mathbf{q}_k \\ \mathbf{q}_j &= \mathbf{v}_j / \|\mathbf{v}_j\|_2\end{aligned}$$

- For $j = 1, \dots, q$, $\mathcal{C}\{\mathbf{x}_1, \dots, \mathbf{x}_j\} = \mathcal{C}\{\mathbf{q}_1, \dots, \mathbf{q}_j\}$ and $\mathbf{q}_j \perp \mathcal{C}\{\mathbf{x}_1, \dots, \mathbf{x}_{j-1}\}$.
- Collectively, we have $\mathbf{X} = \mathbf{Q}\mathbf{R}$, where
 - $\mathbf{Q} \in \mathbb{R}^{n \times p}$ has orthonormal columns \mathbf{q}_j and thus $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}_p$.
 - $\mathbf{R} = \mathbf{Q}^T \mathbf{X} \in \mathbb{R}^{p \times p}$ has entries $r_{kj} = \langle \mathbf{q}_k, \mathbf{x}_j \rangle$, which are available from the algorithm. Note $r_{kj} = 0$ for $k > j$. Thus \mathbf{R} is upper triangular.

This is called the *QR decomposition* of \mathbf{X} .

- The original normal equation $\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{y}$ becomes $\mathbf{R}^T \mathbf{R} \mathbf{b} = \mathbf{R}^T \mathbf{Q} \mathbf{y}$, which is easy to solve (why?). Therefore QR decomposition by Gram-Schmidt offers a practical algorithm for solving the least squares problem.
- $\mathbf{X}^T \mathbf{X} = \mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R} = \mathbf{R}^T \mathbf{R}$ is the *Cholesky decomposition* of $\mathbf{X}^T \mathbf{X}$.

Reparameterization

- Example: We want to predict weight y_i by height x_i (and intercept). Intuitively it should not matter whether we record weight in pounds or kilograms.

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \quad \mathbf{W} = \begin{pmatrix} 1 & w_1 \\ 1 & w_2 \\ \vdots & \vdots \\ 1 & w_n \end{pmatrix},$$

where $w_i = 0.4536x_i$ (1 pound = 0.4536 kilogram). Note

$$\mathbf{W} = \mathbf{X} \begin{pmatrix} 1 & 0 \\ 0 & 0.4536 \end{pmatrix}, \quad \mathbf{X} = \mathbf{W} \begin{pmatrix} 1 & 0 \\ 0 & 2.2046 \end{pmatrix}.$$

- Example: We want to predict salary y_i by gender x_i (and intercept). Intuitively it should not matter whether we record gender as $x_i = 1_{\{\text{male}\}}$ or $w_i = 1_{\{\text{female}\}} = 1 - x_i$ or even keeping both x_i and w_i as predictors.

$$\mathbf{X}_1 = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \quad \mathbf{X}_2 = \begin{pmatrix} 1 & 1 - x_1 \\ 1 & 1 - x_2 \\ \vdots & \vdots \\ 1 & 1 - x_n \end{pmatrix}, \quad \mathbf{X}_3 = \begin{pmatrix} 1 & x_1 & 1 - x_1 \\ 1 & x_2 & 1 - x_2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & 1 - x_n \end{pmatrix}.$$

Note

$$\mathbf{X}_2 = \mathbf{X}_1 \begin{pmatrix} 1 & 1 \\ 0 & -1 \end{pmatrix}, \quad \mathbf{X}_3 = \mathbf{X}_1 \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \end{pmatrix}, \dots$$

- Two linear models $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$, where $\mathbf{X} \in \mathbb{R}^{n \times p}$, and $\mathbf{y} = \mathbf{W}\mathbf{c} + \mathbf{e}$, where $\mathbf{W} \in \mathbb{R}^{n \times q}$, are equivalent or *reparameterizations* of each other if $\mathcal{C}(\mathbf{X}) = \mathcal{C}(\mathbf{W})$.
- If $\mathcal{C}(\mathbf{X}) = \mathcal{C}(\mathbf{W})$, then $\mathbf{P}_\mathbf{X} = \mathbf{P}_\mathbf{W}$.

Proof. $\mathbf{P}_\mathbf{X}$ is an orthogonal projection onto $\mathcal{C}(\mathbf{X})$. $\mathbf{P}_\mathbf{W}$ is an orthogonal projection onto $\mathcal{C}(\mathbf{W})$. But we know orthogonal projection onto a vector space is unique. \square

- The fitted values $\hat{\mathbf{y}}$ and residuals $\hat{\mathbf{e}}$ of two equivalent linear models (parameterizations) are same. Hence $\mathbf{P}_\mathbf{X} = \mathbf{P}_\mathbf{W}$.

Proof. $\hat{\mathbf{y}} = \mathbf{P}_\mathbf{X}\mathbf{y}$ and $\hat{\mathbf{e}} = (\mathbf{I}_n - \mathbf{P}_\mathbf{X})\mathbf{y}$. Since $\mathbf{P}_\mathbf{X} = \mathbf{P}_\mathbf{W}$, we see the fitted values and residuals must be same. \square

- Translation between solutions of equivalent models.

If $\mathcal{C}(\mathbf{X}) = \mathcal{C}(\mathbf{W})$, then there exist $\mathbf{T} \in \mathbb{R}^{p \times q}$ and $\mathbf{S} \in \mathbb{R}^{q \times p}$ such that $\mathbf{W} = \mathbf{X}\mathbf{T}$ and $\mathbf{X} = \mathbf{W}\mathbf{S}$.

- If $\hat{\mathbf{c}}$ is a solution to the normal equation $\mathbf{W}^T\mathbf{W}\mathbf{c} = \mathbf{W}^T\mathbf{y}$, then $\hat{\mathbf{b}} = \mathbf{T}\hat{\mathbf{c}}$ is a solution to the normal equation $\mathbf{X}^T\mathbf{X}\mathbf{b} = \mathbf{X}^T\mathbf{y}$.

$$\textit{Proof. } \mathbf{X}^T\mathbf{X}\hat{\mathbf{b}} = \mathbf{X}^T\mathbf{X}\mathbf{T}\hat{\mathbf{c}} = \mathbf{X}^T\mathbf{W}\hat{\mathbf{c}} = \mathbf{X}^T\mathbf{P}_\mathbf{W}\mathbf{y} = \mathbf{X}^T\mathbf{P}_\mathbf{X}\mathbf{y} = (\mathbf{P}_\mathbf{X}\mathbf{X})^T\mathbf{y} = \mathbf{X}^T\mathbf{y} \quad \square$$

- Similarly, if $\hat{\mathbf{b}}$ is a solution to the normal equation $\mathbf{X}^T\mathbf{X}\mathbf{b} = \mathbf{X}^T\mathbf{y}$, then $\hat{\mathbf{c}} = \mathbf{S}\hat{\mathbf{b}}$ is a solution to the normal equation $\mathbf{W}^T\mathbf{W}\mathbf{c} = \mathbf{W}^T\mathbf{y}$.

Linear mean model

So far we have considered the method of least squares, which concerns the approximation of a vector $\mathbf{y} \in \mathbb{R}^n$ by p vectors $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^n$. Next we consider the linear mean model

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e},$$

where the random errors \mathbf{e} are assumed to have mean $\mathbb{E}(\mathbf{e}) = \mathbf{0}$.

Random vectors

A brief review of random vectors.

- Let $\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$ be a vector of n random variables with mean $\mathbb{E}(y_i) = \mu_i$, variance $\text{Var}(y_i) = \mathbb{E}(y_i - \mathbb{E}y_i)^2 = \sigma_i^2$, and covariance $\text{Cov}(y_i, y_j) = \mathbb{E}[(y_i - \mathbb{E}y_i)(y_j - \mathbb{E}y_j)] = \sigma_{ij} = \sigma_{ji}$. Collectively, we write

$$\mathbb{E}(\mathbf{y}) = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} = \boldsymbol{\mu}, \quad \text{Cov}(\mathbf{y}) = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{12} & \sigma_{22} & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_{nn} \end{pmatrix}.$$

- Let $\mathbf{y} \in \mathbb{R}^n$ be a random vector, then

$$\text{Cov}(\mathbf{y}) = \mathbb{E}(\mathbf{y}\mathbf{y}^T) - (\mathbb{E}\mathbf{y})(\mathbb{E}\mathbf{y})^T.$$

Special case: $\text{Var}(y) = \mathbb{E}y^2 - (\mathbb{E}y)^2$.

- If $\mathbf{y} \in \mathbb{R}^n$ is a random vector, $\mathbf{A} \in \mathbb{R}^{r \times n}$ is a constant matrix, and $\mathbf{b} \in \mathbb{R}^r$ is a constant vector, then

$$\begin{aligned} \mathbb{E}(\mathbf{A}\mathbf{y} + \mathbf{b}) &= \mathbf{A}\mathbb{E}(\mathbf{y}) + \mathbf{b} \\ \text{Cov}(\mathbf{A}\mathbf{y} + \mathbf{b}) &= \mathbf{A}\text{Cov}(\mathbf{y})\mathbf{A}^T. \end{aligned}$$

Special case: $\mathbb{E}(ay + b) = a\mathbb{E}(y) + b$ and $\text{Var}(ay + b) = a^2\text{Var}(y)$ for random variable y .

- For example, under the linear mean model, $\mathbf{E}\mathbf{y} = \mathbf{E}(\mathbf{X}\mathbf{b} + \mathbf{e}) = \mathbf{X}\mathbf{b} + \mathbf{E}(\mathbf{e}) = \mathbf{X}\mathbf{b}$.

- The covariance between two random vectors $\mathbf{x} \in \mathbb{R}^m$ and $\mathbf{y} \in \mathbb{R}^n$ is defined as

$$\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbf{E}[(\mathbf{x} - \mathbf{E}\mathbf{x})(\mathbf{y} - \mathbf{E}\mathbf{y})^T] \in \mathbb{R}^{m \times n}.$$

Note $\text{Cov}(\mathbf{y}, \mathbf{y}) = \text{Cov}(\mathbf{y})$.

- Let $\mathbf{x} \in \mathbb{R}^m$ and $\mathbf{y} \in \mathbb{R}^n$ be two random vectors and $\mathbf{A} \in \mathbb{R}^{p \times m}$ and $\mathbf{B} \in \mathbb{R}^{q \times n}$ be two constant matrices. Then

$$\text{Cov}(\mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{y}) = \mathbf{A}\text{Cov}(\mathbf{x}, \mathbf{y})\mathbf{B}^T.$$

Special case: $\text{Cov}(ax, by) = ab\text{Cov}(x, y)$ for random variables x, y .

- Let $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$ be two random vectors and $\mathbf{A} \in \mathbb{R}^{p \times m}$, $\mathbf{B} \in \mathbb{R}^{p \times n}$ are two constant matrices. Then

$$\begin{aligned} & \text{Cov}(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y}) \\ &= \mathbf{A}\text{Cov}(\mathbf{x})\mathbf{A}^T + \mathbf{A}\text{Cov}(\mathbf{x}, \mathbf{y})\mathbf{B}^T + \mathbf{B}\text{Cov}(\mathbf{y}, \mathbf{x})\mathbf{A}^T + \mathbf{B}\text{Cov}(\mathbf{y})\mathbf{B}^T. \end{aligned}$$

Special case: $\text{Var}(ax + by) = a^2\text{Var}(x) + 2ab\text{Cov}(x, y) + b^2\text{Var}(y)$ for random variables x, y .

- Expectation of quadratic form. Let $\mathbf{x} \in \mathbb{R}^n$ be a random vector with mean $\mathbf{E}(\mathbf{x}) = \boldsymbol{\mu}$ and covariance $\text{Cov}(\mathbf{x}) = \boldsymbol{\Omega}$. $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a constant matrix. Then

$$\mathbf{E}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = \text{tr}(\mathbf{A}\boldsymbol{\Omega}) + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}.$$

Proof. $\mathbf{E}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = \mathbf{E}\text{tr}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = \mathbf{E}\text{tr}(\mathbf{A} \mathbf{x} \mathbf{x}^T) = \text{tr}\mathbf{E}(\mathbf{A} \mathbf{x} \mathbf{x}^T) = \text{tr}\mathbf{A}\mathbf{E}(\mathbf{x} \mathbf{x}^T) = \text{tr}(\mathbf{A}(\boldsymbol{\Omega} + \boldsymbol{\mu}\boldsymbol{\mu}^T)) = \text{tr}(\mathbf{A}\boldsymbol{\Omega}) + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}. \quad \square$

Special case: $\mathbf{E}(ax^2) = a\text{Var}(x) + a[\mathbf{E}(x)]^2$ for a random variable x .

5 Lecture 5: Sep 16

Announcement

- HW1 returned.
- HW2 is posted and it is due next Wed, Sep 25.
<http://hua-zhou.github.io/teaching/st552-2013fall/ST552-2013-HW2.pdf>
Recitations in afternoon sessions of Sep 18 and Sep 25.
- Read JM chapter 3.

Last time

- Geometry of least squares solution: $\mathbf{P}_X = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$, fitted values $\hat{\mathbf{y}} = \mathbf{P}_X \mathbf{y}$, residuals $\hat{\mathbf{e}} = (\mathbf{I}_n - \mathbf{P}_X) \mathbf{y}$, $\hat{\mathbf{y}} \perp \hat{\mathbf{e}}$, $\|\mathbf{y}\|_2^2 = \|\hat{\mathbf{y}}\|_2^2 + \|\hat{\mathbf{e}}\|_2^2$
- Gram-Schmidt orthogonalization and QR decomposition: GS as a series of linear regressions, $\mathbf{X} = \mathbf{QR}$ as a practical way to solve normal equation, see HW2 for the case of non-full rank \mathbf{X}
- Reparameterizations (equivalent models): $\mathcal{C}(\mathbf{X}) = \mathcal{C}(\mathbf{W})$, $\mathbf{P}_X = \mathbf{P}_W$, translation between solutions of equivalent models
- Linear mean model: $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$, where $\mathbf{E}\mathbf{e} = \mathbf{0}$
- Random vectors

Today

- Positive (semi)definite matrix
- Estimable functions

Positive (semi)definite matrix

Assume $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric.

- A real symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *positive semi-definite* (or *nonnegative definite*, or p.s.d.) if $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ for all \mathbf{x} . Notation $\mathbf{A} \succeq \mathbf{0}_{n \times n}$.

- E.g., the *Gramian matrix* $\mathbf{X}^\top \mathbf{X}$ or $\mathbf{X} \mathbf{X}^\top$.
- The notation $\mathbf{A} \succeq \mathbf{0}_{n \times n}$ means $a_{ij} \geq 0$ for all i, j .
- If inequality is strict for all $\mathbf{x} \neq \mathbf{0}$, then \mathbf{A} is *positive definite*. $\mathbf{A} \succ \mathbf{0}_{n \times n}$.
- The notation $\mathbf{A} > \mathbf{0}_{n \times n}$ means $a_{ij} > 0$ for all i, j .
- We write $\mathbf{A} \succeq \mathbf{B}$ means $\mathbf{A} - \mathbf{B} \succeq \mathbf{0}_{n \times n}$.
- If $\mathbf{A} \succeq \mathbf{B}$, then $\det(\mathbf{A}) \geq \det(\mathbf{B})$ with equality if and only if $\mathbf{A} = \mathbf{B}$.
- Cholesky decomposition. Each positive semidefinite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ can be factorized as $\mathbf{A} = \mathbf{L} \mathbf{L}^\top$ for some lower triangular matrix $\mathbf{L} \in \mathbb{R}^{n \times n}$ with nonnegative diagonal entries.
- $\mathbf{A} \in \mathbb{R}^{n \times n}$ is positive semidefinite if and only if \mathbf{A} is a covariance matrix of a random vector.

Proof. “If part”: let $\mathbf{A} = \text{Cov}(\mathbf{x})$ for some random vector \mathbf{x} . Then for any constant \mathbf{c} of same length as \mathbf{x} , $\mathbf{c}^\top \mathbf{A} \mathbf{c} = \mathbf{c}^\top \text{Cov}(\mathbf{x}) \mathbf{c} = \text{Var}(\mathbf{c}^\top \mathbf{x}) \geq 0$. “Only if part”: let $\mathbf{A} = \mathbf{L} \mathbf{L}^\top$ be the Cholesky decomposition and \mathbf{x} a vector of iid standard normals. Then $\mathbf{L} \mathbf{x}$ has covariance matrix $\mathbf{L} \text{Cov}(\mathbf{x}) \mathbf{L}^\top = \mathbf{L} \mathbf{I}_n \mathbf{L}^\top = \mathbf{A}$. □

Estimable function

Assume the linear mean model: $\mathbf{y} = \mathbf{X} \mathbf{b} + \mathbf{e}$, $\text{E}(\mathbf{e}) = \mathbf{0}$. One main interest is estimation of the underlying parameter \mathbf{b} . Can \mathbf{b} be estimated or what functions of \mathbf{b} can be estimated?

- A parametric function $\mathbf{\Lambda} \mathbf{b}$, $\mathbf{\Lambda} \in \mathbb{R}^{m \times p}$, is said to be (linearly) *estimable* if there exists an *affinely unbiased estimator* of $\mathbf{\Lambda} \mathbf{b}$ for all $\mathbf{b} \in \mathbb{R}^p$. That is there exist constants $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{c} \in \mathbb{R}^m$ such that $\text{E}(\mathbf{A} \mathbf{y} + \mathbf{c}) = \mathbf{\Lambda} \mathbf{b}$ for all \mathbf{b} .
- Theorem: Assuming the linear mean model, the parametric function $\mathbf{\Lambda} \mathbf{b}$ is (linearly) estimable if and only if $\mathcal{C}(\mathbf{\Lambda}^\top) \subset \mathcal{C}(\mathbf{X}^\top)$, or equivalently $\mathcal{N}(\mathbf{X}) \subset \mathcal{N}(\mathbf{\Lambda})$.
 “ $\mathbf{\Lambda} \mathbf{b}$ is estimable \Leftrightarrow the row space of $\mathbf{\Lambda}$ is contained in the row space of $\mathbf{X} \Leftrightarrow$ the null space of \mathbf{X} is contained in the null space of $\mathbf{\Lambda}$.”

Proof. Let $\mathbf{A}\mathbf{y} + \mathbf{c}$ be an affine estimator of $\mathbf{\Lambda}\mathbf{b}$. Unbiasedness requires

$$\mathbb{E}(\mathbf{A}\mathbf{y} + \mathbf{c}) = \mathbf{A}\mathbb{E}(\mathbf{y}) + \mathbf{c} = \mathbf{A}\mathbf{X}\mathbf{b} + \mathbf{c} = \mathbf{\Lambda}\mathbf{b}$$

for all $\mathbf{b} \in \mathbb{R}^p$. Taking the special value $\mathbf{b} = \mathbf{0}$ shows that $\mathbf{c} = \mathbf{0}$. Thus $(\mathbf{A}\mathbf{X} - \mathbf{\Lambda})\mathbf{b} = \mathbf{0}$ for all \mathbf{b} . Now taking special values $\mathbf{b} = \mathbf{e}_i$ shows that the columns of the matrix $\mathbf{A}\mathbf{X} - \mathbf{\Lambda}$ are all zeros. This implies $\mathbf{A}\mathbf{X} = \mathbf{\Lambda}$. Therefore the matrix \mathbf{A} exists if and only if rows of $\mathbf{\Lambda}$ are linear combinations of the rows of \mathbf{X} , that is, if and only if $\mathcal{C}(\mathbf{\Lambda}^T) \subset \mathcal{C}(\mathbf{X}^T)$. \square

- Corollary: $\mathbf{X}\mathbf{b}$ is estimable.
“Expected value of any observation $\mathbb{E}(y_i)$ and their linear combinations are estimable.”
- Corollary: If \mathbf{X} has full column rank, then any linear combinations of \mathbf{b} are estimable.
- If $\mathbf{\Lambda}\mathbf{b}$ is (linearly) estimable, then its *least squares estimator* $\mathbf{\Lambda}\hat{\mathbf{b}}$ is invariant to the choice of the least squares solution $\hat{\mathbf{b}}$.

Proof. Let $\hat{\mathbf{b}}_1, \hat{\mathbf{b}}_2$ be two least squares solutions. Then $\hat{\mathbf{b}}_1 - \hat{\mathbf{b}}_2 \in \mathcal{N}(\mathbf{X}^T\mathbf{X}) = \mathcal{N}(\mathbf{X}) \subset \mathcal{N}(\mathbf{\Lambda})$. Hence $\mathbf{\Lambda}(\hat{\mathbf{b}}_1 - \hat{\mathbf{b}}_2) = \mathbf{0}$, that is $\mathbf{\Lambda}\hat{\mathbf{b}}_1 = \mathbf{\Lambda}\hat{\mathbf{b}}_2$. \square

- The least squares estimator $\mathbf{\Lambda}\hat{\mathbf{b}}$ is a linearly unbiased estimator of $\mathbf{\Lambda}\mathbf{b}$.

Proof. The least squares solution takes the general form

$$\hat{\mathbf{b}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} + [\mathbf{I}_p - (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}]\mathbf{q}$$

where $\mathbf{q} \in \mathbb{R}^p$ is arbitrary. Thus the least squares estimator

$$\begin{aligned} \mathbf{\Lambda}\hat{\mathbf{b}} &= \mathbf{\Lambda}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} + \mathbf{\Lambda}[\mathbf{I}_p - (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}]\mathbf{q} \\ &= \mathbf{\Lambda}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \end{aligned}$$

is a linear function of \mathbf{y} . Now

$$\begin{aligned} \mathbb{E}(\mathbf{\Lambda}\hat{\mathbf{b}}) &= \mathbf{\Lambda}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbb{E}(\mathbf{y}) \\ &= \mathbf{\Lambda}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}\mathbf{b} \\ &= \mathbf{\Lambda}\mathbf{b}, \end{aligned}$$

since $\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^-$ is a projection onto $\mathcal{C}(\mathbf{X}^T \mathbf{X}) = \mathcal{C}(\mathbf{X}^T)$ and $\mathcal{C}(\mathbf{\Lambda}^T) \subset \mathcal{C}(\mathbf{X}^T)$. Therefore the least squares estimator is unbiased. \square

- General recipes for showing estimability of $\mathbf{\Lambda} \mathbf{b}$:
 1. Identify each row of $\mathbf{\Lambda}$ as a linear combination of rows of \mathbf{X} .
 2. Find a projection \mathbf{P} onto $\mathcal{C}(\mathbf{X}^T)$, such as $\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^-$ or $\mathbf{P}_{\mathbf{X}^T} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^- \mathbf{X}$. $\mathbf{\Lambda} \mathbf{b}$ is estimable if and only if $\mathbf{P} \mathbf{\Lambda}^T = \mathbf{\Lambda}^T$.
 3. Find a basis $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_{p-r}) \in \mathbb{R}^{p \times (p-r)}$ of $\mathcal{N}(\mathbf{X})$. $\mathbf{\Lambda} \mathbf{b}$ is estimable if and only if $\mathbf{\Lambda} \mathbf{B} = \mathbf{0}_{m \times (p-r)}$.

One-way ANOVA analysis

- Example: We want to study whether the tips (percentage of meal cost) at a restaurant depends on the gender of the service person (waiter or waitress)? Data may be presented as

Waiter	Waitress
0.19, 0.15, 0.10	0.15, 0.14, 0.20

Responses: tip. Factors: gender of service person (waiter, waitress).

- Model: $y_{ij} = \mu + \alpha_i + e_{ij}$, $i = 1, \dots, a$ (a groups), $j = 1, \dots, n_i$ (i -th group has n_i observations).

$$\mathbb{E}(\mathbf{y}) = \mathbf{X} \mathbf{b} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & & & \\ & \mathbf{1}_{n_2} & & & \\ & & \mathbf{1}_{n_2} & & \\ & & & \ddots & \\ & & & & \mathbf{1}_{n_a} \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_a \end{pmatrix}.$$

In total we have $n = \sum_{i=1}^a n_i$ observations and $p = a + 1$ parameters.

- \mathbf{X} has rank $r = a$ and thus $\dim(\mathcal{N}(\mathbf{X})) = 1$. An obvious basis for $\mathcal{N}(\mathbf{X})$ is

$$\begin{pmatrix} 1 \\ -\mathbf{1}_a \end{pmatrix}.$$

Therefore $\boldsymbol{\lambda}^T \mathbf{b}$ is estimable if and only if $\boldsymbol{\lambda}^T \begin{pmatrix} 1 \\ -\mathbf{1}_a \end{pmatrix} = \lambda_0 - \sum_{i=1}^a \lambda_i = 0$.

- Examples of estimable functions:

- grand mean: $\mu + \bar{\alpha}$, where $\bar{\alpha} = a^{-1} \sum_{i=1}^a \alpha_i$
- cell means: $\mu + \alpha_i$
- *contrast*: $\sum_{i=1}^a d_i \alpha_i$ where $\sum_{i=1}^a d_i = 0$

Examples of non-estimable function:

- individual parameters: μ, α_i
- $\alpha_i + \alpha_j$

- To find the least squares solution, we form the normal equation

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \begin{pmatrix} n & n_1 & n_2 & \cdots & n_a \\ n_1 & n_1 & & & \\ n_2 & & n_2 & & \\ \vdots & & & \ddots & \\ n_a & & & & n_a \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_a \end{pmatrix} = \begin{pmatrix} n\bar{y}_{..} \\ n_1\bar{y}_{1.} \\ n_2\bar{y}_{2.} \\ \vdots \\ n_a\bar{y}_{a.} \end{pmatrix} = \mathbf{X}^T \mathbf{y}.$$

A generalized inverse of $\mathbf{X}^T \mathbf{X}$ is (see previous notes for finding generalized inverses)

$$(\mathbf{X}^T \mathbf{X})^- = \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 \\ 0 & n_1^{-1} & & & \\ 0 & & n_2^{-1} & & \\ \vdots & & & \ddots & \\ 0 & & & & n_a^{-1} \end{pmatrix}.$$

Therefore the least squares solution takes the general form

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{y} + [\mathbf{I}_p - (\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{X}] \mathbf{q} = \begin{pmatrix} 0 \\ \bar{y}_{1.} \\ \bar{y}_{2.} \\ \vdots \\ \bar{y}_{a.} \end{pmatrix} + z \begin{pmatrix} 1 \\ -1 \\ -1 \\ \vdots \\ -1 \end{pmatrix}$$

and the least squares estimator of an estimable function $\boldsymbol{\lambda}^T \mathbf{b}$ is

$$\boldsymbol{\lambda}^T \begin{pmatrix} 0 \\ \bar{y}_{1\cdot} \\ \bar{y}_{2\cdot} \\ \vdots \\ \bar{y}_{a\cdot} \end{pmatrix} = \sum_{i=1}^a \lambda_i \bar{y}_{i\cdot}.$$

For example, least squares estimators of $\mu + \alpha_i$ and $\sum_{i=1}^a d_i \alpha_i$ are \bar{y}_i and $\sum_{i=1}^a d_i \bar{y}_i$ respectively.

6 Lecture 6: Sep 18

Announcement

- HW2 recitation this afternoon
- Read JM chapter 4

Last time

- Positive (semi)definite matrix
- Estimability: $\mathbf{\Lambda}\mathbf{b}$ is estimable if and only if $\mathcal{C}(\mathbf{\Lambda}^T) \subset \mathcal{C}(\mathbf{X}^T)$, or equivalently $\mathcal{N}(\mathbf{X}) \subset \mathcal{N}(\mathbf{\Lambda})$, or equivalently $\mathbf{\Lambda}^T \perp \mathcal{N}(\mathbf{X})$
- The least squares estimator $\mathbf{\Lambda}\hat{\mathbf{b}}$ is a linearly unbiased estimate for an estimable function $\mathbf{\Lambda}\mathbf{b}$ and invariant to the choice of least squares solution $\hat{\mathbf{b}}$
- One-way ANOVA

Today

- Two-way ANOVA without and with interaction
- Estimability under reparameterizations

Two-way ANOVA without interaction

- Example: We want to study whether the tips (percentage of meal cost) at a restaurant depends on the gender of the service person (waitress or waiter) and the gender of paying customer? Data may be presented as

	Waiter	Waitress
Male customer	0.19, 0.15, 0.10	0.15, 0.14, 0.20
Female customer	0.16, 0.18, 0.14	0.13, 0.15, 0.19

Responses: tip. Factors: gender of service person (waiter, waitress), gender of paying customer (male, female).

- Model $y_{ijk} = \mu + \alpha_i + \beta_j + e_{ijk}$, $i = 1, \dots, a$ (a levels in factor 1), $j = 1, \dots, b$ (b levels in factor 2), and $k = 1, \dots, n_{ij}$ (n_{ij} observations in the (i, j) -th cell). In total we have $n = \sum_{i,j} n_{ij}$ observations and $p = a + b + 1$ parameters. For simplicity, we consider the case without replicates, i.e., $n_{ij} = 1$. Note adding more replicates to each cell does *not* change the rank of \mathbf{X} .

$$\mathbf{E}(\mathbf{y}) = \mathbf{X}\mathbf{b} = \begin{pmatrix} \mathbf{1}_b & \mathbf{1}_b & & & & & \mathbf{I}_b \\ \mathbf{1}_b & & \mathbf{1}_b & & & & \mathbf{I}_b \\ \vdots & & & \ddots & & & \\ \mathbf{1}_b & & & & \mathbf{1}_b & & \mathbf{I}_b \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_a \\ \beta_1 \\ \vdots \\ \beta_b \end{pmatrix}.$$

- $\mathbf{X} \in \mathbb{R}^{ab \times (1+a+b)}$ has rank $r = a + b - 1$ and $\dim(\mathcal{N}(\mathbf{X})) = 2$. An obvious basis for $\mathcal{N}(\mathbf{X})$ is

$$\left\{ \begin{pmatrix} 1 \\ -\mathbf{1}_a \\ \mathbf{0}_b \end{pmatrix}, \begin{pmatrix} 1 \\ \mathbf{0}_a \\ -\mathbf{1}_b \end{pmatrix} \right\}.$$

Therefore $\boldsymbol{\lambda}^T \mathbf{b}$ is estimable if and only if $\boldsymbol{\lambda}^T \begin{pmatrix} 1 \\ -\mathbf{1}_a \\ \mathbf{0}_b \end{pmatrix} = \lambda_0 - \sum_{i=1}^a \lambda_i = 0$ and

$$\boldsymbol{\lambda}^T \begin{pmatrix} 1 \\ \mathbf{0}_a \\ -\mathbf{1}_b \end{pmatrix} = \lambda_0 - \sum_{j=1}^b \lambda_{a+j} = 0.$$

- Examples of estimable functions:
 - grand mean: $\mu + \bar{\alpha} + \bar{\beta}$, where $\bar{\alpha} = a^{-1} \sum_{i=1}^a \alpha_i$ and $\bar{\beta} = b^{-1} \sum_{j=1}^b \beta_j$
 - cell means: $\mu + \alpha_i + \beta_j$
 - contrast $\sum_{i=1}^a d_i \alpha_i$ where $\sum_{i=1}^a d_i = 0$
 - contrast $\sum_{j=1}^b f_j \beta_j$ where $\sum_{j=1}^b f_j = 0$
- Examples of non-estimable functions:

- individual parameters: μ, α_i, β_j
- marginal means: $\mu + \alpha_i, \mu + \beta_j$
- $\alpha_i - \beta_j$

- To find the least squares estimators for these estimable functions, see HW2.

Two-way ANOVA with interaction

- Example: Responses: tip. Factors: gender of service person (waiter, waitress), gender of paying customer (male, female), and their interaction.
- Model $y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + e_{ijk}$, $i = 1, \dots, a$ (a levels in factor 1), $j = 1, \dots, b$ (b levels in factor 2), and $k = 1, \dots, n_{ij}$ (n_{ij} observations in the (i, j) -th cell). In total we have $n = \sum_{i,j} n_{ij}$ observations and $p = 1 + a + b + ab$ parameters.
- We consider the simple case without replicates, i.e., $n_{ij} = 1$. Note adding more replicates to each cell does *not* change the rank of \mathbf{X} .

$$E(\mathbf{y}) = \mathbf{X}\mathbf{b} = \begin{pmatrix} \mathbf{1}_b & \mathbf{1}_b & & & \mathbf{I}_b & \mathbf{I}_b & & & \\ \mathbf{1}_b & & \mathbf{1}_b & & \mathbf{I}_b & & \mathbf{I}_b & & \\ \vdots & & & \ddots & \vdots & & & \ddots & \\ \mathbf{1}_b & & & & \mathbf{1}_b & \mathbf{I}_b & & & \mathbf{I}_b \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_a \\ \beta_1 \\ \vdots \\ \beta_b \\ \gamma_{11} \\ \gamma_{12} \\ \vdots \\ \gamma_{a,(b-1)} \\ \gamma_{ab} \end{pmatrix}.$$

- $\mathbf{X} \in \mathbb{R}^{ab \times (1+a+b+ab)}$ has rank $r = ab$ and $\dim(\mathcal{N}(\mathbf{X})) = 1 + a + b$. A basis for $\mathcal{N}(\mathbf{X})$ is

$$\left\{ \begin{pmatrix} -1 \\ \mathbf{1}_a \\ \mathbf{0}_b \\ \mathbf{0}_a \otimes \mathbf{0}_b \end{pmatrix}, \begin{pmatrix} 0 \\ -\mathbf{e}_i \\ \mathbf{0}_b \\ \mathbf{e}_i \otimes \mathbf{1}_b \end{pmatrix}, i = 1, \dots, a, \begin{pmatrix} 0 \\ \mathbf{0}_a \\ -\mathbf{e}_j \\ \mathbf{1}_a \otimes \mathbf{e}_j \end{pmatrix}, j = 1, \dots, b \right\}.$$

Therefore $\boldsymbol{\lambda}^T \mathbf{b}$ is estimable if and only if $\boldsymbol{\lambda}$ is orthogonal to all these basis vectors.

- Examples of estimable functions:
 - grand mean: $\mu + \bar{\alpha} + \bar{\beta} + \bar{\gamma}$.
 - cell means: $\mu + \alpha_i + \beta_j + \gamma_{ij}$
 - main effect differences: $(\alpha_i + \bar{\gamma}_{i\cdot}) - (\alpha_{i'} + \bar{\gamma}_{i'\cdot}), (\beta_j + \bar{\gamma}_{\cdot j}) - (\beta_{j'} + \bar{\gamma}_{\cdot j'})$
 - Interaction effect: $\gamma_{ij} - \gamma_{ij'} - \gamma_{i'j} + \gamma_{i'j'}$
- Examples of non-estimable functions:
 - individual parameters: $\mu, \alpha_i, \beta_j, \gamma_{ij}$
 - marginal means: $\mu + \alpha_i + \bar{\gamma}_{i\cdot}$, where $\bar{\gamma}_{i\cdot} = b^{-1} \sum_{j=1}^b \gamma_{ij}$
 - marginal means: $\mu + \beta_j + \bar{\gamma}_{\cdot j}$, where $\bar{\gamma}_{\cdot j} = a^{-1} \sum_{i=1}^a \gamma_{ij}$

Estimability under reparameterizations

- Recall that two linear models $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$ and $\mathbf{y} = \mathbf{W}\mathbf{c} + \mathbf{e}$ are equivalent (*reparameterizations*) if $\mathcal{C}(\mathbf{X}) = \mathcal{C}(\mathbf{W})$. There exist transformations \mathbf{S} and \mathbf{T} such that

$$\mathbf{X} = \mathbf{W}\mathbf{S} \quad \text{and} \quad \mathbf{W} = \mathbf{X}\mathbf{T}.$$

Thus

$$\begin{aligned} \mathbb{E}(\mathbf{y}) &= \mathbf{X}\mathbf{b} = \mathbf{W}\mathbf{S}\mathbf{b} \\ &= \mathbf{W}\mathbf{c} = \mathbf{X}\mathbf{T}\mathbf{c}. \end{aligned}$$

- We already know how to translate least square solutions between equivalent models:
 - If $\hat{\mathbf{c}}$ solves $\mathbf{W}^T \mathbf{W}\mathbf{c} = \mathbf{W}^T \mathbf{y}$, then $\hat{\mathbf{b}} = \mathbf{T}\hat{\mathbf{c}}$ solves $\mathbf{X}^T \mathbf{X}\mathbf{b} = \mathbf{X}^T \mathbf{y}$.
 - If $\hat{\mathbf{b}}$ solves $\mathbf{X}^T \mathbf{X}\mathbf{b} = \mathbf{X}^T \mathbf{y}$, then $\hat{\mathbf{c}} = \mathbf{S}\hat{\mathbf{b}}$ solves $\mathbf{W}^T \mathbf{W}\mathbf{c} = \mathbf{W}^T \mathbf{y}$.
- Estimability under reparameterization

- If $\Lambda\mathbf{c}$ is estimable under model $\mathbf{y} = \mathbf{W}\mathbf{c} + \mathbf{e}$, then $\Lambda\mathbf{S}\mathbf{b}$ is estimable under model $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$.

Proof. Since $\Lambda\mathbf{c}$ is estimable, $\mathcal{C}(\Lambda^T) \subset \mathcal{C}(\mathbf{W}^T)$ and there exists a matrix \mathbf{A} such that $\Lambda^T = \mathbf{W}^T\mathbf{A}$. Thus $\mathbf{S}^T\Lambda^T = \mathbf{S}^T\mathbf{W}^T\mathbf{A} = \mathbf{X}^T\mathbf{A}$. In other words, $\mathcal{C}(\mathbf{S}^T\Lambda^T) \subset \mathcal{C}(\mathbf{X}^T)$. Therefore $\Lambda\mathbf{S}\mathbf{b}$ is estimable. \square

- Similarly, if $\Lambda\mathbf{b}$ is estimable under model $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$, then $\Lambda\mathbf{T}\mathbf{c}$ is estimable under model $\mathbf{y} = \mathbf{W}\mathbf{c} + \mathbf{e}$.

- Different parameterizations of one-way ANOVA analysis

- An over-parameterized model.

$$\mathbf{E}(\mathbf{y}) = \mathbf{X}\mathbf{b} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & & & \\ & \mathbf{1}_{n_2} & & & \\ & \vdots & & \ddots & \\ & & & & \mathbf{1}_{n_a} \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_a \end{pmatrix}.$$

Individual parameters in \mathbf{b} are *not* estimable since \mathbf{X} does not have full column rank. Certain functions are estimable, such as the grand mean $\mu + \bar{\alpha}$, cell means $\mu + \alpha_i$, and contrasts $\sum_{i=1}^a d_i\alpha_i$ ($\sum_{i=1}^a d_i = 0$).

- A full rank parameterization (deleting the first column of \mathbf{X}): *cell means model*

$$\mathbf{E}(\mathbf{y}) = \mathbf{Z}\boldsymbol{\mu} = \begin{pmatrix} \mathbf{1}_{n_1} & & & \\ & \mathbf{1}_{n_2} & & \\ & & \ddots & \\ & & & \mathbf{1}_{n_a} \end{pmatrix} \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_a \end{pmatrix}.$$

Note

$$\mathbf{X} = \mathbf{Z} \begin{pmatrix} 1 & 1 & & \\ \vdots & & \ddots & \\ 1 & & & 1 \end{pmatrix} \text{ and } \boldsymbol{\mu} = \begin{pmatrix} \mu + \alpha_1 \\ \vdots \\ \mu + \alpha_a \end{pmatrix}.$$

Each element of $\boldsymbol{\mu}$ (cell means) is estimable since \mathbf{Z} has full column rank.

- A full rank parameterization (deleting the last column of \mathbf{X}): *reference cell model*

$$\mathbf{E}(\mathbf{y}) = \mathbf{W}\mathbf{c} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & & & \\ & & \mathbf{1}_{n_2} & & \\ & & & \ddots & \\ & & & & \mathbf{1}_{n_{a-1}} \\ \mathbf{1}_{n_a} & & & & \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_a \end{pmatrix}.$$

Note

$$\mathbf{X} = \mathbf{W} \begin{pmatrix} 1 & & & & 1 \\ & 1 & & & -1 \\ & & \ddots & & \vdots \\ & & & 1 & -1 \end{pmatrix} \text{ and } \mathbf{c} = \begin{pmatrix} \mu + \alpha_a \\ \alpha_1 - \alpha_a \\ \vdots \\ \alpha_{a-1} - \alpha_a \end{pmatrix}.$$

Each element of \mathbf{c} is estimable since \mathbf{W} has full column rank.

- A full rank parameterization: *difference from the mean model*

$$\mathbf{E}(\mathbf{y}) = \mathbf{U}\boldsymbol{\delta} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & & & \\ & & \mathbf{1}_{n_2} & & \\ & & & \ddots & \\ & & & & \mathbf{1}_{n_{a-1}} \\ \mathbf{1}_{n_a} & -\mathbf{1}_{n_a} & -\mathbf{1}_{n_a} & \cdots & -\mathbf{1}_{n_a} \end{pmatrix} \begin{pmatrix} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_a \end{pmatrix}.$$

Note

$$\mathbf{X} = \mathbf{U} \begin{pmatrix} 1 & a^{-1} & a^{-1} & & a^{-1} \\ & 1 - a^{-1} & -a^{-1} & & -a^{-1} \\ & -a^{-1} & 1 - a^{-1} & & -a^{-1} \\ & \vdots & & \ddots & \\ -a^{-1} & -a^{-1} & & & 1 - a^{-1} \end{pmatrix} \text{ and } \boldsymbol{\delta} = \begin{pmatrix} \mu + \bar{\alpha} \\ \alpha_1 - \bar{\alpha} \\ \alpha_2 - \bar{\alpha} \\ \vdots \\ \alpha_a - \bar{\alpha} \end{pmatrix}.$$

Each element of $\boldsymbol{\delta}$ is estimable since \mathbf{U} has full column rank.

Review: correlation and independence

- Two random variables $x, y \in \mathbb{R}$ are *uncorrelated* if $\text{Cov}(x, y) = 0$.

- Two random vectors $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$ are *uncorrelated* if $\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbf{0}_{m \times n}$.
- If \mathbf{x} and \mathbf{y} are uncorrelated random vectors, then

$$\text{Cov}(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y}) = \mathbf{A}\text{Cov}(\mathbf{x})\mathbf{A}^T + \mathbf{B}\text{Cov}(\mathbf{y})\mathbf{B}^T.$$

Special case: If x, y are uncorrelated random variables, then $\text{Cov}(ax + by) = a^2\text{Var}(x) + b^2\text{Var}(y)$.

- We say random vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$ are (mutually) independent if their joint density factorizes as product of marginal densities

$$f_{\mathbf{X}_1, \dots, \mathbf{X}_n}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n f_{\mathbf{X}_i}(\mathbf{x}_i)$$

for all values of $\mathbf{x}_1, \dots, \mathbf{x}_n$.

- Mutual independence implies *pairwise independence*. The converse is not true in general. (HW3)
- If $\mathbf{X}_1, \dots, \mathbf{X}_n$ are (mutually) independent, then for any functions g_1, \dots, g_n

$$\mathbb{E} \prod_{i=1}^n g_i(\mathbf{X}_i) = \prod_{i=1}^n \mathbb{E} g_i(\mathbf{X}_i).$$

- If two random vectors \mathbf{x} and \mathbf{y} are independent, then they are uncorrelated. The converse is not true in general. (HW3)

Proof. Corollary of the preceding result. □

7 Lecture 7: Sep 23

Announcement

- Reminder: HW2 due this Wed
- HW3 posted and due Oct 2
<http://hua-zhou.github.io/teaching/st552-2013fall/ST552-2013-HW3.pdf>
- HW 2/3 recitation this Wed afternoon

Last time

- Estimability of two-way ANOVA without and with interaction
- Estimability under reparameterizations
- 4 commonly used parameterizations of one-way ANOVA
- Review: independence and correlation
- Gauss-Markov model: introduction

Today

- Gauss-Markov theorem
- Least squares estimator of σ^2
- Underfitting and overfitting
- Aitken model

Gauss-Markov model and Gauss-Markov theorem (JM 4.2)

- Assumptions of Gauss-Markov model:

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e},$$

where $E(\mathbf{e}) = \mathbf{0}$ and $\text{Cov}(\mathbf{e}) = \sigma^2 \mathbf{I}_n$. In words, errors have zero mean, constant variance (homoskedasticity), and are uncorrelated. Equivalently, $E(\mathbf{y}) = \mathbf{X}\mathbf{b}$ and $\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{I}_n$.

- Under linear mean model, we showed that the *least squares estimator* $\Lambda\hat{\mathbf{b}}$, where $\hat{\mathbf{b}}$ is any least squares solution to the normal equation, is a linearly unbiased estimator of an estimable function $\Lambda\mathbf{b}$ ($\mathcal{C}(\Lambda^T) \subset \mathcal{C}(\mathbf{X}^T)$).
- Under the Gauss-Markov assumptions, we can compute the covariance matrix of the least squares estimator

$$\begin{aligned}
\text{Cov}(\Lambda\hat{\mathbf{b}}) &= \text{Cov}(\Lambda(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}) \\
&= \Lambda(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\text{Cov}(\mathbf{y})\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\Lambda^T \\
&= \sigma^2\Lambda(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\Lambda^T \\
&= \sigma^2\Lambda(\mathbf{X}^T\mathbf{X})^{-1}\Lambda^T.
\end{aligned}$$

The last equality is because $\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}$ is a projection onto $\mathcal{C}(\mathbf{X}^T\mathbf{X}) = \mathcal{C}(\mathbf{X}^T) \supset \mathcal{C}(\Lambda^T)$.

- The celebrated Gauss-Markov theorem states that the least squares estimator $\Lambda\hat{\mathbf{b}}$ has the smallest variance among all linear unbiased estimators of $\Lambda\mathbf{b}$. Why smaller variance is better?
- Gauss-Markov Theorem (vector version). Under Gauss-Markov assumptions, if $\Lambda\mathbf{b}$ is estimable, then the least squares estimator $\Lambda\hat{\mathbf{b}}$ is the *best (minimum variance) affine unbiased estimator* (MVAUE). That is

$$\text{Cov}(\hat{\boldsymbol{\theta}}) \succeq \text{Cov}(\Lambda\hat{\mathbf{b}})$$

for any affine unbiased estimator $\hat{\boldsymbol{\theta}}$ of $\Lambda\mathbf{b}$.

Proof. Let $\hat{\boldsymbol{\theta}} = \mathbf{C}\mathbf{y} + \mathbf{d}$. Unbiasedness requires $\text{E}(\hat{\boldsymbol{\theta}}) = \mathbf{C}\mathbf{X}\mathbf{b} + \mathbf{d} = \Lambda\mathbf{b}$ for all \mathbf{b} . Taking $\mathbf{b} = \mathbf{0}$ shows that $\mathbf{d} = \mathbf{0}$. Thus $\hat{\boldsymbol{\theta}}$ must be a linear estimator.

$$\begin{aligned}
&\text{Cov}(\hat{\boldsymbol{\theta}}) \\
&= \text{Cov}(\mathbf{C}\mathbf{y}) \\
&= \text{Cov}(\Lambda\hat{\mathbf{b}} + \mathbf{C}\mathbf{y} - \Lambda\hat{\mathbf{b}}) \\
&= \text{Cov}(\Lambda\hat{\mathbf{b}}) + \text{Cov}(\mathbf{C}\mathbf{y} - \Lambda\hat{\mathbf{b}}) + \text{Cov}(\Lambda\hat{\mathbf{b}}, \mathbf{C}\mathbf{y} - \Lambda\hat{\mathbf{b}}) + \text{Cov}(\mathbf{C}\mathbf{y} - \Lambda\hat{\mathbf{b}}, \Lambda\hat{\mathbf{b}}).
\end{aligned}$$

The covariance terms vanish because

$$\begin{aligned}
& \text{Cov}(\Lambda \hat{\mathbf{b}}, \mathbf{C}\mathbf{y} - \Lambda \hat{\mathbf{b}}) \\
&= \text{E}(\Lambda \hat{\mathbf{b}} - \Lambda \mathbf{b})(\mathbf{C}\mathbf{y} - \Lambda \hat{\mathbf{b}})^T \\
&= \text{E}(\Lambda \hat{\mathbf{b}}\mathbf{y}^T \mathbf{C}^T) - \text{E}(\Lambda \hat{\mathbf{b}}\hat{\mathbf{b}}^T \Lambda^T) - \text{E}(\Lambda \mathbf{b}\mathbf{y}^T \mathbf{C}^T) + \text{E}(\Lambda \mathbf{b}\hat{\mathbf{b}}^T \Lambda^T) \\
&= \text{E}(\Lambda \hat{\mathbf{b}}\mathbf{y}^T \mathbf{C}^T) - \text{E}(\Lambda \hat{\mathbf{b}}\hat{\mathbf{b}}^T \Lambda^T) - \Lambda \mathbf{b}\mathbf{b}^T \Lambda^T + \Lambda \mathbf{b}\mathbf{b}^T \Lambda^T \\
&= \text{E}(\Lambda \hat{\mathbf{b}}\mathbf{y}^T \mathbf{C}^T) - \text{E}(\Lambda \hat{\mathbf{b}}\hat{\mathbf{b}}^T \Lambda^T) \\
&= \text{E}[\Lambda(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\mathbf{y}^T \mathbf{C}^T] - \text{E}[\Lambda(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\mathbf{y}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \Lambda^T] \\
&= \Lambda(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{E}(\mathbf{y}\mathbf{y}^T) \mathbf{C}^T - \Lambda(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{E}(\mathbf{y}\mathbf{y}^T) \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \Lambda^T \\
&= \Lambda(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{E}(\mathbf{y}\mathbf{y}^T) [\mathbf{C}^T - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \Lambda^T] \\
&= \Lambda(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\sigma^2 \mathbf{I} + \mathbf{X}\mathbf{b}\mathbf{b}^T \mathbf{X}^T) [\mathbf{C}^T - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \Lambda^T] \\
&= [\sigma^2 \Lambda(\mathbf{X}^T \mathbf{X})^{-1} + \Lambda(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\mathbf{b}\mathbf{b}^T] [\mathbf{X}^T \mathbf{C}^T - \mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \Lambda^T] \\
&= [\sigma^2 \Lambda(\mathbf{X}^T \mathbf{X})^{-1} + \Lambda(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\mathbf{b}\mathbf{b}^T] [\mathbf{X}^T \mathbf{C}^T - \Lambda^T] \text{(why?)} \\
&= [\sigma^2 \Lambda(\mathbf{X}^T \mathbf{X})^{-1} + \Lambda(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\mathbf{b}\mathbf{b}^T] \mathbf{0}_{p \times m} \text{(why?)} \\
&= \mathbf{0}_{m \times m}.
\end{aligned}$$

Therefore $\text{Cov}(\hat{\boldsymbol{\theta}}) - \text{Cov}(\Lambda \hat{\mathbf{b}}) = \text{Cov}(\mathbf{C}\mathbf{y} - \Lambda \hat{\mathbf{b}}) \succeq \mathbf{0}_{m \times m}$. □

- En route, we showed that any affine unbiased estimator of $\Lambda \mathbf{b}$ must be a linear estimator. Therefore we also say that the least squares estimator is BLUE (best linear unbiased estimator).
- En route, we also showed that $\text{Cov}(\Lambda \hat{\mathbf{b}}, \mathbf{C}\mathbf{y} - \Lambda \hat{\mathbf{b}}) = \mathbf{0}_{m \times m}$ for any unbiased estimator $\mathbf{C}\mathbf{y}$ of $\Lambda \hat{\mathbf{b}}$. Therefore the least squares estimator $\Lambda \hat{\mathbf{b}}$ is uncorrelated with any (linearly) unbiased estimators of $\mathbf{0}_m$.
- When \mathbf{X} has full column rank, then the least squares solution

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

is BLUE for estimating \mathbf{b} .

- A drawback of the criterion of best (minimum variance) affine unbiased estimator (MVAUE) is that it is not operational. We cannot minimize a matrix. However, we can minimize a scalar function of matrix: trace, determinant, largest eigenvalue and so on. The trace criterion is the most convenient.

- The *affine minimum-trace unbiased estimator* (MTAUE) of an estimable function $\Lambda \mathbf{b}$ is an affine unbiased estimator of $\Lambda \mathbf{b}$, say $\widehat{\Lambda \mathbf{b}}$, such that,

$$\text{tr}(\text{Cov}(\hat{\boldsymbol{\theta}})) \geq \text{tr}(\text{Cov}(\widehat{\Lambda \mathbf{b}}))$$

for all affine unbiased estimator $\hat{\boldsymbol{\theta}}$ of $\Lambda \mathbf{b}$.

- The best (minimum variance) affine unbiased estimator (MVAUE) is also an affine minimum-trace unbiased estimator (MTAUE). The converse is not true in general (HW3). If we can find an MTAUE and show that it is unique, then it must be the MVAUE unless it does not exist.
- A mechanic way to *derive* MVAUE is

1. Let $\hat{\boldsymbol{\theta}} = \mathbf{A}\mathbf{y} + \mathbf{c}$ be an affine estimator of $\Lambda \mathbf{b}$
2. Unbiasedness imposes $\mathbf{A}\mathbf{X} = \Lambda$ and $\mathbf{c} = \mathbf{0}$
3. Minimize $\text{tr}(\text{Cov}(\hat{\boldsymbol{\theta}})) = \sigma^2 \text{tr}(\mathbf{A}\mathbf{A}^T)$ subject to $\mathbf{A}\mathbf{X} = \Lambda$ to determine \mathbf{A} . This yields MTAUE.
4. Show that MTAUE is MVAUE by showing that MTAUE is unique and MVAUE exists or by direct argument

Least squares estimator of σ^2 (JM 4.3)

- Since σ^2 is a quadratic concept, we consider estimation of σ^2 by a quadratic function of \mathbf{y} . That is, a function of form

$$\mathbf{y}^T \mathbf{A} \mathbf{y}.$$

Any estimator of this form is called a *quadratic estimator*. If in addition \mathbf{A} is positive (semi)definite, then it is called a *quadratic and positive* estimator.

- An estimator $\hat{\sigma}^2$ of σ^2 is *unbiased* if

$$\text{E}(\hat{\sigma}^2) = \sigma^2$$

for all $\mathbf{b} \in \mathbb{R}^p$ and $\sigma^2 > 0$.

- Under Gauss-Markov assumptions ($E(\mathbf{y}) = \mathbf{X}\mathbf{b}$ and $\text{Cov}(\mathbf{e}) = \sigma^2\mathbf{I}$), the quadratic estimator

$$\hat{\sigma}^2 = \frac{1}{n-r} \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y} = \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2}{n-r} = \frac{\|\hat{\mathbf{e}}\|_2^2}{n-r} = \frac{\text{SSE}}{n-r},$$

where $r = \text{rank}(\mathbf{X})$, is unbiased.

Proof.

$$\begin{aligned} E(\hat{\sigma}^2) &= \frac{1}{n-r} E[\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}] \\ &= \frac{1}{n-r} \text{tr}[(\mathbf{I} - \mathbf{P}_X) \text{Cov}(\mathbf{y})] + \frac{1}{n-r} E(\mathbf{y})^T (\mathbf{I} - \mathbf{P}_X) E(\mathbf{y}) \\ &= \frac{\sigma^2}{n-r} \text{tr}(\mathbf{I} - \mathbf{P}_X) + \frac{1}{n-r} (\mathbf{X}\mathbf{b})^T (\mathbf{I} - \mathbf{P}_X) (\mathbf{X}\mathbf{b}) \\ &= \frac{\sigma^2}{n-r} \text{tr}(\mathbf{I} - \mathbf{P}_X) \\ &= \frac{\sigma^2}{n-r} [\text{rank}(\mathbf{I}) - \text{rank}(\mathbf{P}_X)] \\ &= \sigma^2. \end{aligned}$$

□

Remark: $\hat{\sigma}^2 = \frac{1}{n-r} \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}$ is called the *least squares estimator of σ^2* .

- To show that the least squares estimator of σ^2 is best (minimum variance), we need further assumptions on the third and fourth moment of \mathbf{e} . We will do this later under the *normal* linear regression model $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{I})$.

Underfitting and misspecification (JM 4.4.1)

- Underfitting (also called misspecification) means omitting predictors in the true model.
- Let $\mathbf{y} = \mathbf{X}\mathbf{b} + \boldsymbol{\eta} + \mathbf{e}$ represent the true model, where $\boldsymbol{\eta}$ includes the omitted variables and their coefficients.

- The bias of a least squares estimator $\Lambda \mathbf{b}$ for an estimable function $\Lambda \mathbf{b}$ is

$$\begin{aligned}
\mathbb{E}(\Lambda \hat{\mathbf{b}}) - \Lambda \mathbf{b} &= \Lambda (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{y}) - \Lambda \mathbf{b} \\
&= \Lambda (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \mathbf{b} + \boldsymbol{\eta}) - \Lambda \mathbf{b} \\
&= \Lambda \mathbf{b} + \Lambda (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\eta} - \Lambda \mathbf{b} \\
&= \Lambda (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\eta} \\
&= \mathbf{A}^T \mathbf{P}_X \boldsymbol{\eta}
\end{aligned}$$

where $\Lambda^T = \mathbf{X}^T \mathbf{A}$. \mathbf{A} exists since $\mathcal{C}(\Lambda^T) \subset \mathcal{C}(\mathbf{X}^T)$ due to estimability.

If $\boldsymbol{\eta} \perp \mathcal{C}(\mathbf{X})$, then the bias is zero.

- The bias of the least squares estimator of σ^2 is

$$\begin{aligned}
&(n-r)^{-1} \mathbb{E}[\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}] - \sigma^2 \\
&= (n-r)^{-1} \text{tr}[(\mathbf{I} - \mathbf{P}_X) \sigma^2 \mathbf{I}] + (n-r)^{-1} (\mathbf{X} \mathbf{b} + \boldsymbol{\eta})^T (\mathbf{I} - \mathbf{P}_X) (\mathbf{X} \mathbf{b} + \boldsymbol{\eta}) - \sigma^2 \\
&= \sigma^2 + (n-r)^{-1} (\mathbf{X} \mathbf{b} + \boldsymbol{\eta})^T (\mathbf{I} - \mathbf{P}_X) (\mathbf{X} \mathbf{b} + \boldsymbol{\eta}) - \sigma^2 \\
&= (n-r)^{-1} (\mathbf{X} \mathbf{b} + \boldsymbol{\eta})^T (\mathbf{I} - \mathbf{P}_X) (\mathbf{X} \mathbf{b} + \boldsymbol{\eta}) \\
&= (n-r)^{-1} \boldsymbol{\eta}^T (\mathbf{I} - \mathbf{P}_X) \boldsymbol{\eta}.
\end{aligned}$$

If $\boldsymbol{\eta} \in \mathcal{C}(\mathbf{X})$, then the bias is 0.

- A decomposition of the misspecification $\boldsymbol{\eta}$:

$$\mathbf{y} = \mathbf{X} \mathbf{b} + \boldsymbol{\eta} + \mathbf{e} = \mathbf{X} \mathbf{b} + \mathbf{P}_X \boldsymbol{\eta} + (\mathbf{I} - \mathbf{P}_X) \boldsymbol{\eta} + \mathbf{e}.$$

The piece $\mathbf{P}_X \boldsymbol{\eta}$ affects estimation of $\Lambda \mathbf{b}$; the other piece $(\mathbf{I} - \mathbf{P}_X) \boldsymbol{\eta}$ affects estimation of σ^2 .

- In presence of underfitting (misspecification), unbiasedness in estimating both $\Lambda \mathbf{b}$ and σ^2 can be achieved only when $\boldsymbol{\eta} = \mathbf{0}_n$.

Overfitting and multicollinearity (JM 4.4.2)

- Overfitting means including unnecessary predictors in the model.

$$\mathbf{y} = \mathbf{X}_1 \mathbf{b}_1 + \mathbf{X}_2 \mathbf{b}_2 + \mathbf{e},$$

where the second group of predictors are redundant. That is $\mathbf{b}_2 = \mathbf{0}$.

- For simplicity, assume both \mathbf{X}_1 and $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$ have full column rank.
- Using only \mathbf{X}_1 , the least squares estimator of \mathbf{b}_1 is

$$\tilde{\mathbf{b}}_1 = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y}$$

with

$$\begin{aligned} \mathbb{E}(\tilde{\mathbf{b}}_1) &= \mathbf{b}_1 \\ \text{Cov}(\tilde{\mathbf{b}}_1) &= \sigma^2 (\mathbf{X}_1^T \mathbf{X}_1)^{-1}. \end{aligned}$$

- Using both \mathbf{X}_1 and \mathbf{X}_2 (overfitting), the least squares estimator of $\begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix}$ is

$$\begin{aligned} \begin{pmatrix} \hat{\mathbf{b}}_1 \\ \hat{\mathbf{b}}_2 \end{pmatrix} &= \left[\begin{pmatrix} \mathbf{X}_1^T \\ \mathbf{X}_2^T \end{pmatrix} (\mathbf{X}_1 \ \mathbf{X}_2) \right]^{-1} \begin{pmatrix} \mathbf{X}_1^T \\ \mathbf{X}_2^T \end{pmatrix} \mathbf{y} \\ &= \begin{pmatrix} \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{X}_2 \\ \mathbf{X}_2^T \mathbf{X}_1 & \mathbf{X}_2^T \mathbf{X}_2 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}_1^T \mathbf{y} \\ \mathbf{X}_2^T \mathbf{y} \end{pmatrix} \end{aligned}$$

with

$$\begin{aligned} \mathbb{E} \begin{pmatrix} \hat{\mathbf{b}}_1 \\ \hat{\mathbf{b}}_2 \end{pmatrix} &= \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{0} \end{pmatrix} \\ \text{Cov} \begin{pmatrix} \hat{\mathbf{b}}_1 \\ \hat{\mathbf{b}}_2 \end{pmatrix} &= \sigma^2 \begin{pmatrix} \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{X}_2 \\ \mathbf{X}_2^T \mathbf{X}_1 & \mathbf{X}_2^T \mathbf{X}_2 \end{pmatrix}^{-1}. \end{aligned}$$

The corresponding block for $\text{Cov}(\hat{\mathbf{b}}_1)$ is (see HW3)

$$\begin{aligned} \text{Cov}(\hat{\mathbf{b}}_1) &= \sigma^2 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \\ &\quad + \sigma^2 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2 [\mathbf{X}_2^T (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1}) \mathbf{X}_2]^{-1} \mathbf{X}_2^T \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1}. \end{aligned}$$

- Therefore overfitting does not change the unbiasedness of least squares estimators. It inflates their variance.
- Also overfitting does not change the unbiasedness of the least squares estimator of σ^2 since

$$\begin{aligned} \mathbb{E}(\hat{\sigma}_1^2) &= \frac{1}{n - r_1} \mathbb{E}[\mathbf{y}^T (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1}) \mathbf{y}] = \sigma^2 \\ \mathbb{E}(\hat{\sigma}^2) &= \frac{1}{n - r} \mathbb{E}[\mathbf{y}^T (\mathbf{I} - \mathbf{P}_{\mathbf{X}}) \mathbf{y}] = \sigma^2, \end{aligned}$$

where $r_1 = \text{rank}(\mathbf{X}_1)$ and $r = \text{rank}(\mathbf{X})$.

- Under orthogonality between two sets of predictors $\mathbf{X}_1^T \mathbf{X}_2 = \mathbf{0}_{p_1 \times p_2}$, we have $\tilde{\mathbf{b}}_1 = \hat{\mathbf{b}}_1$ and $\text{Cov}(\tilde{\mathbf{b}}_1) = \text{Cov}(\hat{\mathbf{b}}_1)$.
- When \mathbf{X}_2 gets close to \mathbf{X}_1 (*multicollinearity*), the second term of $\text{Cov}(\hat{\mathbf{b}}_1)$ (extra variance) explodes.

8 Lecture 8: Sep 25

Announcement

- HW2 due today
- HW3 due next Wed Oct 2
- HW 2/3 recitation this afternoon
- HW4 will be posted this week (last HW before Midterm 1)

Last time

- Gauss-Markov model ($E(\mathbf{y}) = \mathbf{X}\mathbf{b}$ and $\text{Cov}(\mathbf{y}) = \sigma^2\mathbf{I}$) and Gauss-Markov theorem (least squares estimators $\Lambda\hat{\mathbf{b}}$ are MVAUE)
- Least squares estimator of σ^2 , $\hat{\sigma}^2 = (n - r)^{-1}\mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{y}$, is unbiased
- Underfitting (causing bias) and overfitting (inflating variance)
- Aitken model: introduction

Today

- Aitken model

Aitken model (JM 4.5)

- In the *Aitken model*, we relax the variance assumption $\text{Cov}(\mathbf{e}) = \sigma^2\mathbf{I}$ to $\text{Cov}(\mathbf{e}) = \sigma^2\mathbf{V}$, where \mathbf{V} is a fixed, positive semidefinite matrix.
- Example: Heteroskedasticity model $\mathbf{V} = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$. For instance, $\sigma_i^2 = \sigma^2 x_i^2$, $x_i \neq 0$.

- Example: Equicorrelation model

$$\text{Cov}(\mathbf{y}) = \text{Cov}(\mathbf{e}) = \sigma^2 \mathbf{I} + \tau^2 \mathbf{1}_n \mathbf{1}_n^T = \begin{pmatrix} \sigma^2 + \tau^2 & \tau^2 & \tau^2 & \dots & \tau^2 \\ \tau^2 & \sigma^2 + \tau^2 & & \dots & \tau^2 \\ \tau^2 & & \tau^2 & \sigma^2 + \tau^2 & \\ \vdots & & & \ddots & \\ \tau^2 & \tau^2 & & & \sigma^2 + \tau^2 \end{pmatrix}$$

- Example: AR(1) (autoregressive model) $e_i = \rho e_{i-1} + a_i$ where $\text{Cov}(\mathbf{a}) = \sigma_a^2 \mathbf{I}$.

$$\text{Cov}(\mathbf{y}) = \text{Cov}(\mathbf{e}) = \frac{\sigma_a^2}{1 - \rho^2} \begin{pmatrix} 1 & \rho & \rho^2 & \dots & \rho^{n-1} \\ \rho & 1 & \rho & \dots & \rho^{n-2} \\ \rho^2 & \rho & 1 & & \\ \vdots & & & \ddots & \\ \rho^{n-1} & \rho^{n-2} & & & 1 \end{pmatrix} = \frac{\sigma_a^2}{1 - \rho^2} (\rho^{|i-j|})_{i,j}.$$

Estimability under Aitken model

- Under the Aitken model ($\text{E}(\mathbf{y}) = \mathbf{X}\mathbf{b}$, $\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{V}$), a linear function $\mathbf{\Lambda}\mathbf{b}$ is estimable if and only if $\mathcal{C}(\mathbf{\Lambda}^T) \subset \mathcal{C}(\mathbf{X}^T)$.

Remark: no assumption about singularity of \mathbf{V} is needed.

Proof. The previous proof for the linear mean model case ($\text{E}(\mathbf{y}) = \mathbf{X}\mathbf{b}$) applies here since no assumption about the second moment were used. (Aitken model is just a special linear mean model.) \square

Review: method of Lagrangian multipliers for optimization with equality constraints (JM A.2)

$\min_U f(\mathbf{x})$, $U \subset \mathbb{R}^n$, subject to constraints $g_i(\mathbf{x}) = 0$ for $i = 1, \dots, m$. $g : \mathbb{R}^n \mapsto \mathbb{R}^m$.

- Lagrange multiplier theory. *Lagrangian* function $L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T g(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x})$. Strategy for finding equality constrained minimum. Find the stationary point $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ of the Lagrangian, $\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{0}_n$ for all $\boldsymbol{\lambda} \in \mathbb{R}^m$ together with $g(\mathbf{x}) = \mathbf{0}_m$.

- (Necessary conditions) Assume conditions (i) $g(\mathbf{y}) = \mathbf{0}_m$, (2) f and g are differentiable in some n -ball $B(\mathbf{y})$, (iii) $Dg(\mathbf{y}) \in \mathbb{R}^{m \times n}$ is continuous at \mathbf{y} , (iv) $Dg(\mathbf{y})$ has full row rank, (v) $f(\mathbf{x}) \geq f(\mathbf{y})$ for any $\mathbf{x} \in B(\mathbf{y})$ satisfying $g(\mathbf{x}) = \mathbf{0}_m$ (\mathbf{y} a local minimum subject to constraints). Then there exists $\boldsymbol{\lambda} \in \mathbb{R}^m$ satisfying $\nabla f(\mathbf{y}) + Dg(\mathbf{y})^\top \boldsymbol{\lambda} = \mathbf{0}_n$, i.e., $(\mathbf{y}, \boldsymbol{\lambda})$ is a stationarity point of the Lagrangian $L(\mathbf{x}, \boldsymbol{\lambda})$. In other words, there exists $\boldsymbol{\lambda} \in \mathbb{R}^m$, such that $\nabla L(\mathbf{y}, \boldsymbol{\lambda}) = \mathbf{0}_{m+n}$.
- (Sufficient conditions) (i) f twice differentiable at \mathbf{y} , (ii) g twice differentiable at \mathbf{y} , (iii) the Jacobian matrix $Dg(\mathbf{y}) \in \mathbb{R}^{m \times n}$ has full row rank m , (iv) it is a stationarity point of the Lagrangian at a given $\boldsymbol{\lambda} \in \mathbb{R}^m$, (v) $\mathbf{u}^\top d^2 f(\mathbf{y}) \mathbf{u} > 0$ for all $\mathbf{u} \neq \mathbf{0}_n$ satisfying $Dg(\mathbf{y}) \mathbf{u} = \mathbf{0}_m$. Then \mathbf{y} is a strict local minimum under constraint $g(\mathbf{y}) = \mathbf{0}_m$.
- Check condition (v). Condition (v) is equivalent to the “bordered determinantal criterion”

$$(-1)^m \det \begin{pmatrix} \mathbf{0}_{m \times m} & \mathbf{B}_r \\ \mathbf{B}_r^\top & \mathbf{A}_{rr} \end{pmatrix} > 0$$

for $r = m + 1, \dots, n$, where \mathbf{A}_{rr} is the top left r -by- r block of $d^2 f(\mathbf{y}) + \sum_{i=1}^m \lambda_i d^2 g_i(\mathbf{y})$, $\mathbf{B}_r \in \mathbb{R}^{m \times r}$ is the first r columns of $Dg(\mathbf{y})$.

- (Sufficient condition for a global minimum) Lagrangian first order condition + convexity of the Lagrangian on U for some $(\mathbf{y}, \boldsymbol{\lambda})$.
- Example: Linearly constrained least squares solution. $\min \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$ subject to linear constrained $\mathbf{V}\boldsymbol{\beta} = \mathbf{d}$. Form the Lagrangian

$$L(\boldsymbol{\beta}, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \boldsymbol{\lambda}^\top (\mathbf{V}\boldsymbol{\beta} - \mathbf{d}).$$

Stationary condition says

$$\begin{aligned} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} - \mathbf{X}^\top \mathbf{y} + \mathbf{V}^\top \boldsymbol{\lambda} &= \mathbf{0}_p \\ \mathbf{V}\boldsymbol{\beta} &= \mathbf{d} \end{aligned}$$

or equivalently

$$\begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{V}^\top \\ \mathbf{V} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^\top \mathbf{y} \\ \mathbf{d} \end{pmatrix}.$$

Review: derivatives of trace functions

Let $\mathbf{X} \in \mathbb{R}^{m \times n}$. \mathbf{A} and \mathbf{B} are constant matrices of compatible dimensions. Then

$$\begin{aligned}\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{A}\mathbf{X}) &= \mathbf{A}^T \\ \frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{X}\mathbf{A}\mathbf{X}^T\mathbf{B}) &= \mathbf{B}^T\mathbf{X}\mathbf{A}^T + \mathbf{B}\mathbf{X}\mathbf{A} \\ \frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{X}\mathbf{A}\mathbf{X}\mathbf{B}) &= \mathbf{B}^T\mathbf{X}^T\mathbf{A}^T + \mathbf{A}^T\mathbf{X}^T\mathbf{B}^T.\end{aligned}$$

Proof: check elementwise.

Aitken theorem and generalized least squares (no constraints, non-singular \mathbf{V})

- (Aitken theorem) Assume Aitken model ($\mathbf{E}(\mathbf{y}) = \mathbf{X}\mathbf{b}$, $\text{Cov}(\mathbf{y}) = \sigma^2\mathbf{V}$ and $\mathbf{V} \succ \mathbf{0}_{n \times n}$). The best (minimum variance) affine unbiased estimator (MVAUE) of an estimable function $\mathbf{\Lambda}\mathbf{b}$ is

$$\widehat{\mathbf{\Lambda}\mathbf{b}} = \mathbf{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}$$

with variance matrix

$$\text{Cov}(\mathbf{\Lambda}\mathbf{b}) = \sigma^2\mathbf{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{\Lambda}^T.$$

Proof. Since we are told the answer. We can go ahead checking unbiasedness and minimum variance property directly. But let's proceed in a more constructive way. First derive the estimator as a minimum trace affine unbiased estimator (MTAUE) and then show that it is MVAUE.

Let $\widehat{\mathbf{\Lambda}\mathbf{b}} = \mathbf{A}\mathbf{y} + \mathbf{c}$ be an affine estimator for $\mathbf{\Lambda}\mathbf{b}$. Unbiasedness requires

$$\mathbf{E}(\widehat{\mathbf{\Lambda}\mathbf{b}}) = \mathbf{A}\mathbf{X}\mathbf{b} + \mathbf{c} = \mathbf{\Lambda}\mathbf{b} \text{ for all } \mathbf{b}.$$

This implies $\mathbf{c} = \mathbf{0}$ and $\mathbf{A}\mathbf{X} = \mathbf{\Lambda}$. So $\widehat{\mathbf{\Lambda}\mathbf{b}}$ is linear and its variance is

$$\text{Cov}(\widehat{\mathbf{\Lambda}\mathbf{b}}) = \text{Cov}(\mathbf{A}\mathbf{y}) = \sigma^2\mathbf{A}\mathbf{V}\mathbf{A}^T.$$

To find a MTAUE, we consider the minimization problem

$$\begin{aligned}\text{minimize} & \quad \frac{1}{2}\text{tr}(\mathbf{A}\mathbf{V}\mathbf{A}^T) \\ \text{subject to} & \quad \mathbf{A}\mathbf{X} = \mathbf{\Lambda}.\end{aligned}$$

The relevant Lagrangian function is

$$\psi(\mathbf{A}, \mathbf{L}) = \frac{1}{2} \text{tr}(\mathbf{A}\mathbf{V}\mathbf{A}^T) - \text{tr}(\mathbf{L}^T(\mathbf{A}\mathbf{X} - \mathbf{\Lambda}))$$

where $\mathbf{L} \in \mathbb{R}^{m \times p}$ is a matrix of Lagrange multipliers. Differentiating ψ with respect to \mathbf{A} yields

$$\begin{aligned}\mathbf{A}\mathbf{V} &= \mathbf{L}\mathbf{X}^T \\ \mathbf{A}\mathbf{X} &= \mathbf{\Lambda}.\end{aligned}$$

From the first equation, $\mathbf{A} = \mathbf{L}\mathbf{X}^T\mathbf{V}^{-1}$. Substitution into the second equation we have

$$\mathbf{L}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X} = \mathbf{\Lambda}.$$

This system is always consistent (why? show $\mathcal{C}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X}) = \mathcal{C}(\mathbf{X}^T) \supset \mathcal{C}(\mathbf{\Lambda}^T)$ in HW4) so $\mathbf{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^-$ is a solution to \mathbf{L} (may not be unique). Therefore

$$\mathbf{A} = \mathbf{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^- \mathbf{X}^T\mathbf{V}^{-1}$$

is a solution to the constrained minimization problem. Actually it is invariant to the choice of generalized inverse $(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^-$ (HW4). Hence it is unique. Therefore the estimator

$$\widehat{\mathbf{\Lambda}}\mathbf{b} = \mathbf{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^- \mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}$$

with variance

$$\begin{aligned}\text{Cov}(\widehat{\mathbf{\Lambda}}\mathbf{b}) &= \sigma^2 \mathbf{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^- \mathbf{X}^T\mathbf{V}^{-1}\mathbf{V}\mathbf{V}^{-1}\mathbf{X}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^- \mathbf{\Lambda}^T \\ &= \sigma^2 \mathbf{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^- \mathbf{\Lambda}^T\end{aligned}$$

is the unique MTAUE. If MVAUE exists, this is it.

Now we have a good candidate, namely the MTAUE. It becomes a routine to check that it is actually the MVAUE. Consider an arbitrary affine unbiased estimator

$$[\mathbf{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^- \mathbf{X}^T\mathbf{V}^{-1} + \mathbf{C}]\mathbf{y} + \mathbf{d}.$$

Unbiasedness imposes $\mathbf{C}\mathbf{X} = \mathbf{0}_{m \times p}$ and $\mathbf{d} = \mathbf{0}_m$. Its covariance is

$$\begin{aligned} & \sigma^2[\boldsymbol{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1} + \mathbf{C}]\mathbf{V}[\boldsymbol{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1} + \mathbf{C}]^T \\ &= \sigma^2\boldsymbol{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\boldsymbol{\Lambda}^T + \sigma^2\mathbf{C}\mathbf{V}\mathbf{C}^T \\ &\succeq \text{Cov}(\widehat{\boldsymbol{\Lambda}\mathbf{b}}). \end{aligned}$$

Therefore the MTAUE is indeed the MVAUE. \square

- $\widehat{\boldsymbol{\Lambda}\mathbf{b}} = \boldsymbol{\Lambda}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}$ is called the *generalized least squares estimator* of an estimable function $\boldsymbol{\Lambda}\mathbf{b}$.
- The same good candidate can also be obtained by transforming the Aitken model back to the Gauss-Markov model (as done in the textbook).

Aitken model with linear constraints ($\mathbf{V} \succ \mathbf{0}_{n \times n}$)

- In many applications, researchers want to fit linear model with some constraints on the parameter \mathbf{b} . We investigate the estimability issue and try to find the MVAUE assuming that \mathbf{V} is non-singular.
- (Estimability under linear constraints) Consider the Aitken model ($\mathbb{E}(\mathbf{y}) = \mathbf{X}\mathbf{b}$ and $\text{Cov}(\mathbf{y}) = \sigma^2\mathbf{V}$) with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$. A linear function $\boldsymbol{\Lambda}\mathbf{b}$ is estimable if and only if $\mathcal{C}(\boldsymbol{\Lambda}^T) \subset \mathcal{C}((\mathbf{X}^T, \mathbf{R}^T))$.

Remark 1: Apparently there are *more* estimable functions with constraints than without constraints.

Remark 2: No assumption on the singularity of \mathbf{V} is needed here.

Proof. Aitken model with linear constraints is equivalent to the following expanded unconstrained Aitken model

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{r} \end{pmatrix} = \begin{pmatrix} \mathbf{X} \\ \mathbf{R} \end{pmatrix} \mathbf{b} + \begin{pmatrix} \mathbf{e} \\ \mathbf{0} \end{pmatrix},$$

where the error term has first two moments

$$\mathbb{E} \begin{pmatrix} \mathbf{e} \\ \mathbf{0} \end{pmatrix} = \mathbf{0}, \quad \text{Cov} \begin{pmatrix} \mathbf{e} \\ \mathbf{0} \end{pmatrix} = \sigma^2 \begin{pmatrix} \mathbf{V} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}.$$

Thus the estimability result follows from that for this expanded Aitken model without constraints. \square

- (MVAUE under linear constraints) Consider the Aitken model ($E(\mathbf{y}) = \mathbf{X}\mathbf{b}$ and $\text{Cov}(\mathbf{y}) = \sigma^2\mathbf{V}$) with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$ and assume $\mathbf{V} \succ \mathbf{0}_{n \times n}$. The best (minimum variance) affine unbiased estimator (MVAUE) of an estimable function $\mathbf{\Lambda}\mathbf{b}$ is

$$\widehat{\mathbf{\Lambda}\mathbf{b}} = \mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y} + \mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{R}^T(\mathbf{R}\mathbf{G}^{-1}\mathbf{R}^T)^{-1}(\mathbf{r} - \mathbf{R}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y})$$

where

$$\mathbf{G} = \mathbf{X}^T\mathbf{V}^{-1}\mathbf{X} + \mathbf{R}^T\mathbf{R}.$$

It has variance matrix

$$\text{Cov}(\widehat{\mathbf{\Lambda}\mathbf{b}}) = \sigma^2\mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{\Lambda}^T - \sigma^2\mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{R}^T(\mathbf{R}\mathbf{G}^{-1}\mathbf{R}^T)^{-1}\mathbf{R}\mathbf{G}^{-1}\mathbf{\Lambda}^T.$$

Remark 1: Notice that the MVAUE is indeed affine and not linear anymore under linear constraints.

Remark 2: Notice the variance is reduced in presence of constraints.

Proof. We again proceed in the constructive way: first find MTAUE and then show it is MVAUE.

Let $\widehat{\mathbf{\Lambda}\mathbf{b}} = \mathbf{A}\mathbf{y} + \mathbf{c}$ be an affine estimator of $\mathbf{\Lambda}\mathbf{b}$. Unbiasedness imposes

$$E(\widehat{\mathbf{\Lambda}\mathbf{b}}) = \mathbf{A}\mathbf{X}\mathbf{b} + \mathbf{c} = \mathbf{\Lambda}\mathbf{b} \quad \text{for all } \mathbf{b} \text{ satisfying } \mathbf{R}\mathbf{b} = \mathbf{r}.$$

Substituting the general solution $\mathbf{b}^* = \mathbf{R}^{-1}\mathbf{r} + (\mathbf{I} - \mathbf{R}^{-1}\mathbf{R})\mathbf{q}$ shows

$$(\mathbf{\Lambda} - \mathbf{A}\mathbf{X})[\mathbf{R}^{-1}\mathbf{r} + (\mathbf{I} - \mathbf{R}^{-1}\mathbf{R})\mathbf{q}] = \mathbf{c} \quad \text{for all } \mathbf{q},$$

which implies

$$\begin{aligned} (\mathbf{\Lambda} - \mathbf{A}\mathbf{X})\mathbf{R}^{-1}\mathbf{r} &= \mathbf{c} \\ (\mathbf{\Lambda} - \mathbf{A}\mathbf{X})(\mathbf{I} - \mathbf{R}^{-1}\mathbf{R}) &= \mathbf{0}_{m \times p}. \end{aligned}$$

The second equation says $\mathcal{C}((\mathbf{\Lambda} - \mathbf{A}\mathbf{X})^T) \subset \mathcal{C}(\mathbf{R}^T)$. There exists a matrix \mathbf{B} such that $\mathbf{\Lambda} - \mathbf{A}\mathbf{X} = \mathbf{B}\mathbf{R}$ and $\mathbf{c} = \mathbf{B}\mathbf{R}\mathbf{R}^{-1}\mathbf{r} = \mathbf{B}\mathbf{r}$. Therefore the affine estimator takes the form

$$\widehat{\mathbf{\Lambda}\mathbf{b}} = \mathbf{A}\mathbf{y} + \mathbf{B}\mathbf{r}$$

and the unbiasedness condition boils down to $\mathbf{A}\mathbf{X} + \mathbf{B}\mathbf{R} = \mathbf{\Lambda}$.

To find a MTAUE, we consider the minimization problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2}\text{tr}(\mathbf{A}\mathbf{V}\mathbf{A}^T) \\ & \text{subject to} && \mathbf{A}\mathbf{X} + \mathbf{B}\mathbf{R} = \mathbf{\Lambda}. \end{aligned}$$

The relevant Lagrangian function is

$$\psi(\mathbf{A}, \mathbf{B}, \mathbf{L}) = \frac{1}{2}\text{tr}\mathbf{A}\mathbf{V}\mathbf{A}^T - \text{tr}\mathbf{L}^T(\mathbf{A}\mathbf{X} + \mathbf{B}\mathbf{R} - \mathbf{\Lambda}),$$

where \mathbf{L} is a matrix of Lagrangian multipliers. Differentiating with respect to \mathbf{A} and \mathbf{B} gives the first order conditions

$$\begin{aligned} \mathbf{A}\mathbf{V} &= \mathbf{L}\mathbf{X}^T \\ \mathbf{R}\mathbf{L}^T &= \mathbf{0} \\ \mathbf{A}\mathbf{X} + \mathbf{B}\mathbf{R} &= \mathbf{\Lambda}. \end{aligned}$$

Step 1: Express \mathbf{L} in terms of \mathbf{A} . From the first equation, $\mathbf{A} = \mathbf{L}\mathbf{X}^T\mathbf{V}^{-1}$. Thus

$$\mathbf{A}\mathbf{X} = \mathbf{L}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X} = \mathbf{L}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{X} + \mathbf{R}^T\mathbf{R}) = \mathbf{L}\mathbf{G}$$

by the second equation, where $\mathbf{G} = \mathbf{X}^T\mathbf{V}^{-1}\mathbf{X} + \mathbf{R}^T\mathbf{R}$ is p.s.d. It can be shown that $\mathcal{C}(\mathbf{G}) \supset \mathcal{C}(\mathbf{X}^T) \cup \mathcal{C}(\mathbf{R}^T)$ (HW4). Therefore $\mathbf{A}\mathbf{X} = \mathbf{L}\mathbf{G}$ as an equation in \mathbf{L} is always consistent (since $\mathbf{G}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{A}^T = \mathbf{X}^T\mathbf{A}^T$) and $\mathbf{A}\mathbf{X}\mathbf{G}^{-1}$ is one solution to \mathbf{L} .

Step 2: Solve for \mathbf{B} . Post-multiply both sides of the third equation by $\mathbf{G}^{-1}\mathbf{R}^T$ we get $\mathbf{B}\mathbf{R}\mathbf{G}^{-1}\mathbf{R}^T = \mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{R}^T$. As a linear equation in \mathbf{B} , it is always consistent if we can check

$$(\mathbf{R}\mathbf{G}^{-1}\mathbf{R}^T)(\mathbf{R}\mathbf{G}^{-1}\mathbf{R}^T)^{-}\mathbf{R}\mathbf{G}^{-1}\mathbf{\Lambda}^T = \mathbf{R}\mathbf{G}^{-1}\mathbf{\Lambda}^T.$$

We show this by writing $\mathbf{\Lambda}^T = (\mathbf{X}^T, \mathbf{R}^T)\mathbf{T}$ for some transformation matrix \mathbf{T} . Such \mathbf{T} exists because of the estimability of $\mathbf{\Lambda}\mathbf{b}$. Then it only remains to show that $\mathcal{C}(\mathbf{R}\mathbf{G}^{-1}\mathbf{R}^T) \supset \mathcal{C}(\mathbf{R}\mathbf{G}^{-1}\mathbf{X}^T)$ (HW4).

Therefore the linear system is always consistent and one specific solution to \mathbf{B} is $\mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{R}^T(\mathbf{R}\mathbf{G}^{-1}\mathbf{R}^T)^{-}$.

Step 3: Assemble pieces to get the minimizing \mathbf{A}

$$\begin{aligned}
 \mathbf{A} &= \mathbf{L}\mathbf{X}^T\mathbf{V}^{-1} = \mathbf{A}\mathbf{X}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1} \\
 &= (\mathbf{\Lambda} - \mathbf{B}\mathbf{R})\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1} \\
 &= \mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1} - \mathbf{B}\mathbf{R}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1} \\
 &= \mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1} - \mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{R}^T(\mathbf{R}\mathbf{G}^{-1}\mathbf{R}^T)^{-1}\mathbf{R}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1}.
 \end{aligned}$$

Therefore the MTAUE is

$$\begin{aligned}
 \widehat{\mathbf{\Lambda}}\mathbf{b} &= \mathbf{A}\mathbf{y} + \mathbf{B}\mathbf{r} \\
 &= \mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y} + \mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{R}^T(\mathbf{R}\mathbf{G}^{-1}\mathbf{R}^T)^{-1}(\mathbf{r} - \mathbf{R}\mathbf{G}^{-1}\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}),
 \end{aligned}$$

which serves a good candidate for MVAUE. The last two (routine) steps are to derive the variance formula for this MTAUE and show that it is indeed the MVAUE. They are left as exercises in HW4. \square

9 Lecture 9: Sep 30

Announcement

- HW2 returned (77 ± 15). For comparison: HW1 (89 ± 4).
- HW3 due this Wed Oct 2
- HW 3/4 recitation this Wednesday afternoon
- HW4 is posted and due Oct 9 (HW4 questions are covered in Midterm 1)
- Plan for this week: finish Aitken model and Midterm 1 review
- Please send your specific questions (HWs, lecture notes, textbook, ...) to me (hua_zhou@ncsu.edu) by this Saturday; so I can go over during classes (Oct 2 or Oct 7)

Last time

- Aitken model (no constraints, non-singular \mathbf{V}):
MVAUE is $\widehat{\mathbf{\Lambda b}} = \mathbf{\Lambda}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}$ with variance $\sigma^2 \mathbf{\Lambda}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{\Lambda}^T$
- Aitken model (linear constraints $\mathbf{Rb} = \mathbf{r}$, non-singular \mathbf{V}):
Estimability: $\mathcal{C}(\mathbf{\Lambda}^T) \subset \mathcal{C}((\mathbf{X}^T, \mathbf{R}^T))$, MVAUE (via MTAUE).

Today

- Aitken model: imposing constraints for a unique solution
- Aitken model with singular \mathbf{V}

Aitken model with linear constraints: special case $\mathcal{C}(\mathbf{R}^T) \cap \mathcal{C}(\mathbf{X}^T) = \{\mathbf{0}\}$

- As we saw, incorporating constraints increases the class the linear functions $\mathbf{\Lambda b}$ that are estimatable. When $\begin{pmatrix} \mathbf{X} \\ \mathbf{R} \end{pmatrix}$ has full column rank p , we can estimate every single parameter in \mathbf{b} !

- In many applications, we purposely add constraints which are linearly independent of rows of \mathbf{X} . That is $\mathcal{C}(\mathbf{R}^T) \cap \mathcal{C}(\mathbf{X}^T) = \{\mathbf{0}\}$. In this case, the MVAUE takes a simpler form.
- Estimability: $\mathcal{C}(\mathbf{\Lambda}^T) \subset \mathcal{C}((\mathbf{X}^T, \mathbf{R}^T))$. Note $\text{rank}((\mathbf{X}^T, \mathbf{R}^T)) = \text{rank}(\mathbf{X}^T) + \text{rank}(\mathbf{R}^T)$ when $\mathcal{C}(\mathbf{R}^T) \cap \mathcal{C}(\mathbf{X}^T) = \{\mathbf{0}\}$.
- For example, in one-way ANOVA

$$\mathbf{X} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & & & \\ & \mathbf{1}_{n_2} & \mathbf{1}_{n_2} & & \\ & \vdots & & \ddots & \\ & \mathbf{1}_{n_a} & & & \mathbf{1}_{n_a} \end{pmatrix},$$

where $\text{rank}(\mathbf{X}) = a$. Three commonly used constraints are

1. $\alpha_a = 0$
2. $\sum_{i=1}^a \alpha_i = 0$
3. $\sum_{i=1}^a n_i \alpha_i = 0$.

Each of them will make every single parameter in $(\mu, \alpha_1, \dots, \alpha_a)$ estimable.

- (MVAUE) Consider the Aitken model ($E(\mathbf{y}) = \mathbf{X}\mathbf{b}$ and $\text{Cov}(\mathbf{y}) = \sigma^2\mathbf{V}$) with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$ where $\mathcal{C}(\mathbf{R}^T) \cap \mathcal{C}(\mathbf{X}^T) = \{\mathbf{0}\}$. Assume $\mathbf{V} \succ \mathbf{0}_{n \times n}$. The best (minimum variance) affine unbiased estimator (MVAUE) of an estimable function $\mathbf{\Lambda}\mathbf{b}$ is

$$\widehat{\mathbf{\Lambda}\mathbf{b}} = \mathbf{\Lambda}\mathbf{G}^{-1}(\mathbf{X}^T\mathbf{V}^{-1}\mathbf{y} + \mathbf{R}^T\mathbf{r}),$$

where

$$\mathbf{G} = \mathbf{X}^T\mathbf{V}^{-1}\mathbf{X} + \mathbf{R}^T\mathbf{R}.$$

It has variance matrix

$$\text{Cov}(\widehat{\mathbf{\Lambda}\mathbf{b}}) = \sigma^2\mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{\Lambda}^T - \sigma^2\mathbf{\Lambda}\mathbf{G}^{-1}\mathbf{R}^T\mathbf{R}\mathbf{G}^{-1}\mathbf{\Lambda}^T.$$

Remark: note $\mathcal{C}(\mathbf{R}^T) \cap \mathcal{C}(\mathbf{X}^T) = \{\mathbf{0}\}$ is equivalent to $\text{rank}(\mathbf{X}^T, \mathbf{R}^T) = \text{rank}(\mathbf{X}^T) + \text{rank}(\mathbf{R}^T)$.

Proof. Do HW4 Q3 and you will find some cancellations occur in the previous more general result. □

Aitken model with linear constraints and a singular V

Now we are in a position to consider the most general case: Aitken model with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$ and a possibly singular \mathbf{V} .

- First we notice that, with a singular \mathbf{V} , not all values of \mathbf{y} are possible. For example, $y_i = \mu + e_i$, $i = 1, 2, 3, 4$, where $\text{Var}(e_1) = \text{Var}(e_2) = 0$ and $\text{Var}(e_i) > 0$ for $i = 3, 4$, is a linear model with $\text{rank}(\mathbf{V}) = 2$. Any \mathbf{y} with $y_1 \neq y_2$ is inconsistent with this model.
- We can write the constrained linear model as an unconstrained one

$$\mathbf{y}_e = \mathbf{X}_e \mathbf{b} + \mathbf{u}, \quad \mathbf{E}(\mathbf{u}) = \mathbf{0}, \quad \text{Cov}(\mathbf{u}) = \sigma^2 \mathbf{V}_e,$$

where

$$\mathbf{y}_e = \begin{pmatrix} \mathbf{y} \\ \mathbf{r} \end{pmatrix}, \quad \mathbf{X}_e = \begin{pmatrix} \mathbf{X} \\ \mathbf{R} \end{pmatrix}, \quad \mathbf{V}_e = \begin{pmatrix} \mathbf{V} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}.$$

- (Consistency) Aitken model ($\mathbf{E}(\mathbf{y}) = \mathbf{X}\mathbf{b}$, $\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{V}$) with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$ is consistent if and only if

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{r} \end{pmatrix} \in \mathcal{C} \begin{pmatrix} \mathbf{X} & \mathbf{V} \\ \mathbf{R} & \mathbf{0} \end{pmatrix}.$$

Proof. Note the model $\mathbf{y}_e = \mathbf{X}_e \mathbf{b} + \mathbf{u}$ is equivalent to the model $\mathbf{y}_e = \mathbf{X}_e \mathbf{b} + \mathbf{V}_e \mathbf{V}_e^- \mathbf{u}$ since

$$\mathbf{E}(\mathbf{u}) = \mathbf{V}_e \mathbf{V}_e^- \mathbf{E}(\mathbf{u}) = \mathbf{0}_n, \quad \text{Cov}(\mathbf{u}) = \sigma^2 \mathbf{V}_e \mathbf{V}_e^- \mathbf{V}_e \mathbf{V}_e^- \mathbf{V}_e = \sigma^2 \mathbf{V}_e.$$

Therefore the linear model is consistent if and only if $\mathbf{y}_e \in \mathcal{C}((\mathbf{X}_e, \mathbf{V}_e \mathbf{V}_e^-)) = \mathcal{C}((\mathbf{X}_e, \mathbf{V}_e)) = \mathcal{C} \left(\begin{pmatrix} \mathbf{X} & \mathbf{V} \\ \mathbf{R} & \mathbf{0} \end{pmatrix} \right)$. \square

- (Estimability) $\mathbf{\Lambda}\mathbf{b}$ is estimable if and only if

$$\mathcal{C}(\mathbf{\Lambda}^T) \subset \mathcal{C}((\mathbf{X}^T, \mathbf{R}^T)).$$

- (MVAUE) Consider the Aitken model ($E(\mathbf{y}) = \mathbf{X}\mathbf{b}$, $\text{Cov}(\mathbf{y}) = \sigma^2\mathbf{V}$) with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$. The best (minimum) variance affine unbiased estimator (MVAUE) of an estimable function $\mathbf{\Lambda}\mathbf{b}$ is

$$\widehat{\mathbf{\Lambda}\mathbf{b}} = \mathbf{\Lambda}(\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- \mathbf{X}_e^T \mathbf{V}_0^- \mathbf{y}_e,$$

where $\mathbf{V}_0 = \mathbf{V}_e + \mathbf{X}_e \mathbf{X}_e^T$, with variance

$$\text{Cov}(\widehat{\mathbf{\Lambda}\mathbf{b}}) = \sigma^2 \mathbf{\Lambda}[(\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- - \mathbf{I}_p] \mathbf{\Lambda}^T.$$

Remark: Since \mathbf{y}_e involves both \mathbf{y} and \mathbf{r} , the MVAUE is affine in general.

Proof. Let's try to derive the MTAUE first, and then show that the MTAUE is MVAUE. Let $\mathbf{A}\mathbf{y}_e + \mathbf{c}$ be an affine unbiased estimator of an estimable function $\mathbf{\Lambda}\mathbf{b}$. Unbiasedness requires

$$\mathbf{A}\mathbf{X}_e \mathbf{b} + \mathbf{c} = \mathbf{\Lambda}\mathbf{b} \quad \text{for all } \mathbf{b} \in \mathbb{R}^p,$$

which implies

$$\mathbf{A}\mathbf{X}_e = \mathbf{\Lambda} \text{ and } \mathbf{c} = \mathbf{0}.$$

We seek the MTAUE by solving

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \text{tr} \mathbf{A} \mathbf{V}_e \mathbf{A}^T \\ & \text{subject to} && \mathbf{A}\mathbf{X}_e = \mathbf{\Lambda}. \end{aligned}$$

The relevant Lagrangian is

$$\psi(\mathbf{A}, \mathbf{L}) = \frac{1}{2} \text{tr} \mathbf{A} \mathbf{V}_e \mathbf{A}^T - \text{tr}(\mathbf{L}^T (\mathbf{A}\mathbf{X}_e - \mathbf{\Lambda})),$$

where \mathbf{L} is a matrix of Lagrangian multipliers. Differentiating with respect to \mathbf{A} yields the first order conditions

$$\begin{aligned} \mathbf{A} \mathbf{V}_e &= \mathbf{L} \mathbf{X}_e^T \\ \mathbf{A} \mathbf{X}_e &= \mathbf{\Lambda}, \end{aligned}$$

or, as a matrix equation,

$$\begin{pmatrix} \mathbf{V}_e & \mathbf{X}_e \\ \mathbf{X}_e^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{A}^T \\ -\mathbf{L}^T \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{\Lambda}^T \end{pmatrix}.$$

By HW4 Q5(d), a generalized inverse of the matrix on left side is

$$\begin{pmatrix} \mathbf{V}_e & \mathbf{X}_e \\ \mathbf{X}_e^T & \mathbf{0} \end{pmatrix}^- = \begin{pmatrix} \mathbf{V}_0^- - \mathbf{V}_0^- \mathbf{X}_e \mathbf{C}^- \mathbf{X}_e^T \mathbf{V}_0^- & \mathbf{V}_0^- \mathbf{X}_e \mathbf{C}^- \\ \mathbf{C}^- \mathbf{X}_e^T \mathbf{V}_0^- & -\mathbf{C}^- + \mathbf{C} \mathbf{C}^- \end{pmatrix},$$

where $\mathbf{V}_0 = \mathbf{V}_e + \mathbf{X}_e \mathbf{X}_e^T$ and $\mathbf{C} = \mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e$, and

$$\begin{pmatrix} \mathbf{V}_e & \mathbf{X}_e \\ \mathbf{X}_e^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{V}_e & \mathbf{X}_e \\ \mathbf{X}_e^T & \mathbf{0} \end{pmatrix}^- = \begin{pmatrix} \mathbf{V}_0 \mathbf{V}_0^- & \mathbf{0} \\ \mathbf{0} & \mathbf{C} \mathbf{C}^- \end{pmatrix}.$$

By estimability, $\mathcal{C}(\boldsymbol{\Lambda}^T) \subset \mathcal{C}(\mathbf{X}_e^T) = \mathcal{C}(\mathbf{C})$ (In HW4 Q5(c), we also show $\mathcal{C}(\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e) = \mathcal{C}(\mathbf{X}_e^T)$) and thus $\mathbf{C} \mathbf{C}^- \boldsymbol{\Lambda}^T = \boldsymbol{\Lambda}^T$. Therefore the linear system is consistent and the solution to \mathbf{A} takes the general form

$$\mathbf{A} = \boldsymbol{\Lambda} (\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- \mathbf{X}_e^T \mathbf{V}_0^- + \mathbf{Q} (\mathbf{I} - \mathbf{V}_0 \mathbf{V}_0^-)$$

where \mathbf{Q} is arbitrary. Therefore the MTAUE is

$$\begin{aligned} \widehat{\boldsymbol{\Lambda} \mathbf{b}} &= \mathbf{A} \mathbf{y}_e \\ &= \boldsymbol{\Lambda} (\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- \mathbf{X}_e^T \mathbf{V}_0^- \mathbf{y}_e + \mathbf{Q} (\mathbf{I} - \mathbf{V}_0 \mathbf{V}_0^-) \mathbf{y}_e \\ &= \boldsymbol{\Lambda} (\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- \mathbf{X}_e^T \mathbf{V}_0^- \mathbf{y}_e. \end{aligned}$$

The last equality follows from $\mathcal{C}(\mathbf{V}_0) = \mathcal{C}(\mathbf{V}_e + \mathbf{X}_e \mathbf{X}_e^T) = \mathcal{C}((\mathbf{X}_e, \mathbf{V}_e))$ (HW4 Q5(a)) and the consistency $\mathbf{y}_e \in \mathcal{C}((\mathbf{X}_e, \mathbf{V}_e))$.

The variance of the MTAUE is

$$\begin{aligned} \text{Cov}(\widehat{\boldsymbol{\Lambda} \mathbf{b}}) &= \sigma^2 \boldsymbol{\Lambda} (\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- \mathbf{X}_e^T \mathbf{V}_0^- \mathbf{V}_e \mathbf{V}_0^- \mathbf{X}_e (\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- \boldsymbol{\Lambda}^T \\ &= \sigma^2 \boldsymbol{\Lambda} (\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- \mathbf{X}_e^T \mathbf{V}_0^- (\mathbf{V}_0 - \mathbf{X}_e \mathbf{X}_e^T) \mathbf{V}_0^- \mathbf{X}_e (\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- \boldsymbol{\Lambda}^T \\ &= \sigma^2 \boldsymbol{\Lambda} (\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- \boldsymbol{\Lambda}^T - \sigma^2 \boldsymbol{\Lambda} \boldsymbol{\Lambda}^T \\ &= \sigma^2 \boldsymbol{\Lambda} [(\mathbf{X}_e^T \mathbf{V}_0^- \mathbf{X}_e)^- - \mathbf{I}] \boldsymbol{\Lambda}^T. \end{aligned}$$

Last order of business is to show that the MTAUE is indeed MVAUE. Here is another useful device for proving this. If $\boldsymbol{\Lambda} \mathbf{b}$ is estimable, then $\mathbf{c}^T \boldsymbol{\Lambda} \mathbf{b}$ is estimable (why?) and $\mathbf{c}^T \widehat{\boldsymbol{\Lambda} \mathbf{b}}$ is the MTAUE for estimating $\mathbf{c}^T \boldsymbol{\Lambda} \mathbf{b}$. Let $\hat{\boldsymbol{\theta}}$ be another affine unbiased estimator of $\boldsymbol{\Lambda} \mathbf{b}$. Then $\mathbf{c}^T \hat{\boldsymbol{\theta}}$ is an affine unbiased estimator of $\mathbf{c}^T \boldsymbol{\Lambda} \mathbf{b}$. Thus

$$\text{Var}(\mathbf{c}^T \widehat{\boldsymbol{\Lambda} \mathbf{b}}) \leq \text{Var}(\mathbf{c}^T \hat{\boldsymbol{\theta}}),$$

which implies that

$$\mathbf{c}^T \text{Cov}(\widehat{\boldsymbol{\Lambda}\mathbf{b}}) \mathbf{c}^T \leq \mathbf{c}^T \text{Cov}(\widehat{\boldsymbol{\theta}}) \mathbf{c}.$$

Since \mathbf{c} is arbitrary, we have shown that $\text{Cov}(\widehat{\boldsymbol{\Lambda}\mathbf{b}}) \preceq \text{Cov}(\widehat{\boldsymbol{\theta}})$. □

- (Conflict with previous result?) The special case $\mathbf{V} \succ \mathbf{0}_{n \times n}$ does not recover our previous result (Aitken with linear constraints and non-singular \mathbf{V}) obviously. We will take an alternative approach below – transforming a singular Aitken model to a non-singular one.

Before that we first review some basic results on determinant, eigenvalues and eigenvectors.

10 Lecture 10: Oct 2

Announcement

- HW3 due today
- HW 3/4 recitation this afternoon
- HW4 is due next Wed Oct 9 (covered in Midterm 1)
- Please send your specific questions (HWs, lecture notes, textbook, ...) to me (hua_zhou@ncsu.edu) by this Saturday

Last time

- Aitken model: imposing constraints for a unique solution
- Aitken model with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$ and a possibly singular \mathbf{V}

Today

- Generalized least squares

A result on the invariance to the choice of generalized inverse

We often need to show that a product $\mathbf{B}\mathbf{A}^{-}\mathbf{C}$ is invariant to the choice of the generalized inverse. Let's summarize a useful result here.

- If $\mathcal{C}(\mathbf{B}^T) \subset \mathcal{C}(\mathbf{A}^T)$ and $\mathcal{C}(\mathbf{C}) \subset \mathcal{C}(\mathbf{A})$, then the product $\mathbf{B}\mathbf{A}^{-}\mathbf{C}$ is invariant to the choice of the generalized inverse. In other words, it is unique.

Proof. Since $\mathcal{C}(\mathbf{B}^T) \subset \mathcal{C}(\mathbf{A}^T)$, there exists a matrix \mathbf{L} such that $\mathbf{B} = \mathbf{L}\mathbf{A}$. Since $\mathcal{C}(\mathbf{C}) \subset \mathcal{C}(\mathbf{A})$, there exists a matrix \mathbf{R} such that $\mathbf{C} = \mathbf{A}\mathbf{R}$. Then

$$\mathbf{B}\mathbf{A}^{-}\mathbf{C} = \mathbf{L}\mathbf{A}\mathbf{A}^{-}\mathbf{A}\mathbf{R} = \mathbf{L}\mathbf{A}\mathbf{R}.$$

It is also invariant to the choice of the transformation matrix since $\mathbf{L}_1\mathbf{A}\mathbf{R} - \mathbf{L}_2\mathbf{A}\mathbf{R} = \mathbf{B}\mathbf{R} - \mathbf{B}\mathbf{R} = \mathbf{0}$ for any two candidate transformation matrices \mathbf{L}_1 and \mathbf{L}_2 . Similarly it is invariant to the choice of transformation matrix \mathbf{R} too. \square

- Applications: uniqueness of $\mathbf{P}_X = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$, uniqueness of all MVAUEs we derived so far, uniqueness of the variance of all MVAUEs we derived so far, ...

Generalized least squares

- Gauss invented method of least squares as an approximation method. The least squares solution (magically) turns out to be MVAUE.
- For Aitken model, we approached MVAUE through MTAUE. Are the MVAUEs solutions to some generalized least squares criterion? The answer is affirmative. This coincidence is surprising but not trivial.
- (Generalized least squares for Aitken with linear constraints) Consider the Aitken model, $E(\mathbf{y}) = \mathbf{X}\mathbf{b}$ and $\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{V}$, with a possibly singular \mathbf{V} and linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$. Assume

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{r} \end{pmatrix} \in \mathcal{C} \left(\begin{pmatrix} \mathbf{X} & \mathbf{V} \\ \mathbf{R} & \mathbf{0} \end{pmatrix} \right) \quad (\text{consistency})$$

and $\mathcal{C}(\mathbf{\Lambda}^T) \subset \mathcal{C}((\mathbf{X}^T, \mathbf{R}^T))$ (estimability). Then the best (minimum) variance affine unbiased estimator (MVAUE) of $\mathbf{\Lambda}\mathbf{b}$ is $\mathbf{\Lambda}\hat{\mathbf{b}}$, where $\hat{\mathbf{b}}$ minimizes the generalized least squares criterion

$$\begin{pmatrix} \mathbf{y} - \mathbf{X}\mathbf{b} \\ \mathbf{r} - \mathbf{R}\mathbf{b} \end{pmatrix}^T \begin{pmatrix} \mathbf{V} + \mathbf{X}\mathbf{X}^T & \mathbf{X}\mathbf{R}^T \\ \mathbf{R}\mathbf{X}^T & \mathbf{R}\mathbf{R}^T \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{y} - \mathbf{X}\mathbf{b} \\ \mathbf{r} - \mathbf{R}\mathbf{b} \end{pmatrix}.$$

Remark 1: The generalized least squares solution $\hat{\mathbf{b}}$ may not be unique.

Remark 2: One implication of this connection is that we can obtain the MVAUE by solving the generalized least squares problem. There are good algorithms for minimizing a generalized least squares criterion.

Proof. Consider the augmented linear model $E(\mathbf{y}_e) = \mathbf{X}_e \mathbf{b}$, $\text{Cov}(\mathbf{y}_e) = \sigma^2 \mathbf{V}_e$ where

$$\mathbf{y}_e = \begin{pmatrix} \mathbf{y} \\ \mathbf{r} \end{pmatrix}, \quad \mathbf{X}_e = \begin{pmatrix} \mathbf{X} \\ \mathbf{R} \end{pmatrix}, \quad \mathbf{V}_e = \begin{pmatrix} \mathbf{V} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}.$$

Then generalized least squares criterion is

$$(\mathbf{y}_e - \mathbf{X}_e \mathbf{b})^T \mathbf{V}_0^- (\mathbf{y}_e - \mathbf{X}_e \mathbf{b}),$$

where $\mathbf{V}_0 = \mathbf{V}_e + \mathbf{X}_e \mathbf{X}_e^T$. We first show that this expression is invariant to the choice of the generalized inverse \mathbf{V}_0^- . This is true because $\mathbf{y}_e - \mathbf{X}_e \mathbf{b} \in \mathcal{C}((\mathbf{X}_e, \mathbf{V}_e)) = \mathcal{C}(\mathbf{V}_0)$ by consistency. Therefore without loss of generality we can use the Moore-Penrose inverse \mathbf{V}_0^+ , which is symmetric and positive semidefinite, in the criterion.

Let $\mathbf{V}_0^+ = \mathbf{L}\mathbf{L}^T$, where \mathbf{L} can be the Cholesky factor or the symmetric square root. Then the generalized least squares criterion becomes

$$\|\mathbf{L}^T (\mathbf{y}_e - \mathbf{X}_e \mathbf{b})\|_2^2.$$

The least squares solution (may be non-unique) to this regular least squares problem is

$$\hat{\mathbf{b}} = (\mathbf{X}_e^T \mathbf{L}\mathbf{L}^T \mathbf{X}_e)^- \mathbf{X}_e^T \mathbf{L}\mathbf{L}^T \mathbf{y}_e = (\mathbf{X}_e^T \mathbf{V}_0^+ \mathbf{X}_e)^- \mathbf{X}_e^T \mathbf{V}_0^+ \mathbf{y}_e.$$

The final step is to show that

$$\Lambda \hat{\mathbf{b}} = \Lambda (\mathbf{X}_e^T \mathbf{V}_0^+ \mathbf{X}_e)^- \mathbf{X}_e^T \mathbf{V}_0^+ \mathbf{y}_e$$

is invariant to the choice of the generalized inverse $(\mathbf{X}_e^T \mathbf{V}_0^+ \mathbf{X}_e)^-$ and we can replace \mathbf{V}_0^+ by any generalized inverse \mathbf{V}_0^- . We use the result summarized in the previous section.

- To show that we can replace the Moore-Penrose inverse \mathbf{V}_0^+ in $\mathbf{X}_e^T \mathbf{V}_0^+ \mathbf{X}_e$ by any generalized inverse \mathbf{V}_0^- , we note $\mathcal{C}(\mathbf{V}_0) = \mathcal{C}((\mathbf{X}_e, \mathbf{V}_e)) \supset \mathcal{C}(\mathbf{X}_e)$.
- To show that we can replace the Moore-Penrose inverse \mathbf{V}_0^+ in $\mathbf{X}_e^T \mathbf{V}_0^+ \mathbf{y}_e$ by any generalized inverse \mathbf{V}_0^- , we note $\mathcal{C}(\mathbf{V}_0) \supset \mathcal{C}(\mathbf{X}_e)$ and $\mathbf{y}_e \in \mathcal{C}((\mathbf{X}_e, \mathbf{V}_e)) = \mathcal{C}(\mathbf{V}_0)$ by consistency.
- To show that $\Lambda \hat{\mathbf{b}}$ is invariant to the choice of $(\mathbf{X}_e^T \mathbf{V}_0^+ \mathbf{X}_e)^-$, we note

$$\mathcal{C}(\mathbf{X}_e^T \mathbf{V}_0^+ \mathbf{X}_e) = \mathcal{C}(\mathbf{X}_e^T) \supset \mathcal{C}(\Lambda^T).$$

Thus we have shown that $\Lambda \hat{\mathbf{b}}$ coincides with the MVAUE we derived earlier. \square

- In absence of constraints, the generalized least squares criterion reduces to

$$(\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{V} + \mathbf{X}\mathbf{X}^T)^- (\mathbf{y} - \mathbf{X}\mathbf{b}).$$

11 Lecture 11: Oct 7

Announcement

- HW3 returned (92 ± 11)
- HW4 deadline extended to next Wed Oct 16
- Midterm 1, Oct 9 @ 11:45AM-1PM
- No afternoon session this Wed Oct 9
- No office hours today (instructor out for conference).
Makeup office hours: Tue Oct 8 @ 1-3PM
- No office hours Wed Oct 9
- No class and office hours next Monday Oct 14 (instructor out of town)

Last time

- Generalized least squares

Today

- Q&A
- HW4

Q&A

- After all the lectures, I am getting confused and mixed up with some definitions. I wonder if you can give us a summary of what are the ways to get (BLUE, invertible, estimable, consistency, MTAUE, MVAUE, etc.).
- HW5 question1 (b), and question 3.
- I made a mistake in HW2 about writing nested design model (last question from book). I wonder if there are more difficult models we should know how to write out. I just cannot come up with other examples.

- When we prove the unbiasedness of the $\Lambda \hat{\mathbf{b}}$, we plug in the least square solution of \mathbf{b} , then we have $\Lambda(\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{X} \mathbf{y} + \Lambda[\mathbf{I} - (\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{X}] \mathbf{q}$, then the second term is vanished. My thought is: $[\mathbf{I} - (\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{X}]$ is projection onto null space of $\mathbf{X}^T \mathbf{X}$. Because $\mathcal{N}(\mathbf{X}^T \mathbf{X}) = \mathcal{N}(\mathbf{X}^T)$, and because $\Lambda \mathbf{b}$ is estimable, so $\mathcal{C}(\Lambda^T) \subset \mathcal{C}(\mathbf{X}^T) \Leftrightarrow \mathcal{N}(\mathbf{X}^T) \subset \mathcal{N}(\Lambda^T)$. My problem is: we can take transpose, which is $[[\mathbf{I} - (\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{X}] \Lambda^T]^T$, but we only have $\mathcal{N}(\mathbf{X}^T) \subset \mathcal{N}(\Lambda^T)$, I cannot see projection onto $\mathcal{N}(\mathbf{X}^T)$ is also projection onto $\mathcal{N}(\Lambda^T)$? (since $\mathcal{N}(\Lambda^T)$ contains $\mathcal{N}(\mathbf{X}^T)$)
- Can you explain the jointly non-estimable function. And why the number of them is fixed, $p - r$?
- Reparameterization for one-way ANOVA, since we can reparameterize X to be full rank, I am thinking if that matrix is equivalent to one-way ANOVA without interception? (maybe this is trivial)
- For the overfitting part, we have $\text{Cov}(\hat{\mathbf{b}}) = \text{Cov}((\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{y})$, which is

$$\sigma^2(\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{X} ((\mathbf{X}^T \mathbf{X})^{-})^T.$$

my question is how do we come up with $\sigma^2(\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{X} ((\mathbf{X}^T \mathbf{X})^{-})^T = \sigma^2(\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-} = \sigma^2(\mathbf{X}^T \mathbf{X})^{-}$? Is that because the projection. It is like the property of \mathbf{G}^+ inverse?

- For the Lagrangian multiplier in linearly constrained least squares solution, what is that constrained do or mean? $\mathbf{V} \boldsymbol{\beta} = \mathbf{d}$ what is \mathbf{d} ?
- How does that come out for the consistency of Aitken model with linear constraints?
- Conflict, yes, I got confused about the Aitken model with linear constraints and non-singular V / Aitken model with linear constraints and singular V. Why the first MVAUE is non linear , the second is linear?
- Can you go over how to prove column spaces are equal using the null space property?
- Homework 1: 2 (d).
- Homework 2: 1 (JM 2.11) (f), (g) and (h);

12 Lecture 12: Oct 16

Announcement

- HW4 due today.
- HW5 will be posted today and due next Fri Oct 25 (?).
- Midterm 1 returned (80 ± 14).

Today

- Finish Chapter 4: Aitken model with linear constraints $\mathbf{Rb} = \mathbf{r}$ and a possibly singular \mathbf{V} – alternative approach
- Chapter 5: normal and related distributions

Review: determinant (JM A.6)

Assume $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a square matrix.

- The *determinant* of a matrix $\det(\mathbf{A}) = |\mathbf{A}| = \sum (-1)^{\phi(j_1, \dots, j_n)} \prod_{i=1}^n a_{ij_i}$, where $\phi(j_1, \dots, j_n)$ is the number of transpositions to change $(1, \dots, n)$ to (j_1, \dots, j_n) .
- Example,

$$\det \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} = a_{11}a_{22} - a_{12}a_{21}.$$

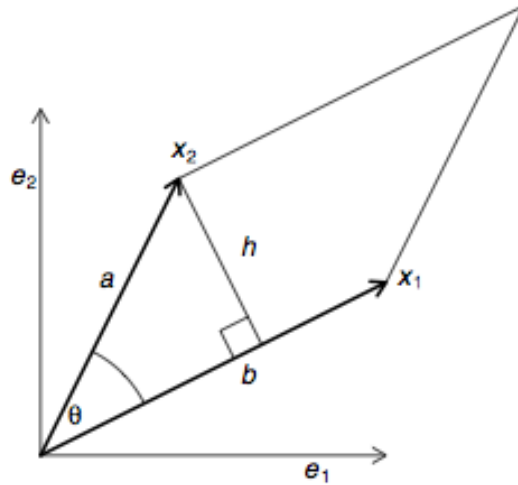


Fig. 3.1. Volume (Area) of Region Determined by \mathbf{x}_1 and \mathbf{x}_2

- Interpretation of the (absolute value of) determinant as the volume of the parallelogram defined by the columns of the matrix:

$$\begin{aligned}
 \text{area} &= bh = \|\mathbf{x}_1\| \|\mathbf{x}_2\| \sin(\theta) \\
 &= \|\mathbf{x}_1\| \|\mathbf{x}_2\| \sqrt{1 - \left(\frac{\langle \mathbf{x}_1, \mathbf{x}_2 \rangle}{\|\mathbf{x}_1\| \|\mathbf{x}_2\|} \right)^2} \\
 &= \sqrt{\|\mathbf{x}_1\|^2 \|\mathbf{x}_2\|^2 - (\langle \mathbf{x}_1, \mathbf{x}_2 \rangle)^2} \\
 &= \sqrt{(x_{11}^2 + x_{12}^2)(x_{21}^2 + x_{22}^2) - (x_{11}x_{21} - x_{12}x_{22})^2} \\
 &= |x_{11}x_{22} - x_{12}x_{21}| \\
 &= |\det(\mathbf{X})|.
 \end{aligned}$$

- Another interpretation of the determinant is the volume changing factor when operating on a set in \mathbb{R}^n . $\text{volume}(f(S)) = |\det(\mathbf{A})| \text{volume}(f(S))$ where $f : \mathbb{R}^n \mapsto \mathbb{R}^n$ is the linear mapping defined by \mathbf{A} .
- The determinant of a lower or upper triangular matrix \mathbf{A} is the product of the diagonal elements $\prod_{i=1}^n a_{ii}$.
- Determinant of a singular matrix is 0.
- Determinant of an orthogonal matrix is 1 (rotation) or -1 (reflection).

- $\det(\mathbf{A}^T) = \det(\mathbf{A})$.
- $\det(\mathbf{A}^{-1}) = 1/\det(\mathbf{A})$.
- $\det(c\mathbf{A}) = c^n \det(\mathbf{A})$.
- $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$.
- For \mathbf{A} and \mathbf{D} square and nonsingular,

$$\det \left(\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} \right) = \det(\mathbf{A}) \det(\mathbf{D} - \mathbf{CA}^{-1}\mathbf{B}) = \det(\mathbf{D}) \det(\mathbf{A} - \mathbf{BD}^{-1}\mathbf{C}).$$

Proof. Take determinant on the both sides of the matrix identity

$$\begin{pmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{CA}^{-1}\mathbf{B} \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{CA}^{-1} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} \begin{pmatrix} \mathbf{I} & -\mathbf{A}^{-1}\mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}.$$

□

Review: eigenvalues and eigenvectors (JM A.6)

Assume $\mathbf{A} \in \mathbb{R}^{n \times n}$ a square matrix.

- *Eigenvalues* are defined as roots of the characteristic equation $\det(\lambda \mathbf{I}_n - \mathbf{A}) = 0$.
- \mathbf{A} is singular if and only if it has at least one 0 eigenvalue.
- If λ is an eigenvalue of \mathbf{A} , then there exist non-zero vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ such that $\mathbf{Ax} = \lambda \mathbf{x}$ and $\mathbf{y}^T \mathbf{A} = \lambda \mathbf{y}^T$. \mathbf{x} and \mathbf{y} are called the (column) *eigenvector* and *row eigenvector* of \mathbf{A} associated with the eigenvalue λ .
- Eigenvectors associated with distinct eigenvalues are linearly independent.

Proof. Let $\mathbf{Ax}_1 = \lambda_1 \mathbf{x}_1$, $\mathbf{Ax}_2 = \lambda_2 \mathbf{x}_2$, and $\lambda_1 \neq \lambda_2$. Assume that \mathbf{x}_1 and \mathbf{x}_2 are linearly dependent such that $\mathbf{x}_2 = \alpha \mathbf{x}_1$. Then

$$\alpha \lambda_1 \mathbf{x}_1 = \alpha \mathbf{Ax}_1 = \mathbf{Ax}_2 = \lambda_2 \mathbf{x}_2 = \alpha \lambda_2 \mathbf{x}_1.$$

That is $\alpha(\lambda_1 - \lambda_2)\mathbf{x}_1 = \mathbf{0}$. Since $\alpha \neq 0$ and $\lambda_1 \neq \lambda_2$, we have $\mathbf{x}_1 = \mathbf{0}$, a contradiction. □

- Eigenvalues of an upper or lower triangular matrix are its diagonal entries:
 $\lambda_i = a_{ii}$.

- Eigenvalues of an idempotent matrix are either 0 or 1.

Proof. If $\mathbf{A}\mathbf{A} = \mathbf{A}$ and $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$, then

$$\lambda\mathbf{x} = \mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{A}\mathbf{x} = \lambda^2\mathbf{x}.$$

Hence $\lambda = \lambda^2$, i.e., $\lambda = 0$ or 1 . □

- Eigenvalues of an orthogonal matrix have complex modulus 1.
- In most statistical applications, we deal with eigenvalues/eigenvectors of symmetric matrices.

The eigenvalues and eigenvectors of a real *symmetric* matrix are real.

- Eigenvectors associated with distinct eigenvalues of a symmetry matrix are orthogonal to each other.
- Eigen-decomposition of a symmetric matrix: $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, where

- $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$

- columns of \mathbf{U} are the eigenvectors which are (or can be chosen to be) mutually orthonormal

- A real symmetric matrix is positive semidefinite (positive definite) if and only if all eigenvalues are nonnegative (positive).

Proof. If \mathbf{A} is positive definite and $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$, then $\mathbf{x}^T\mathbf{A}\mathbf{x} = \lambda\mathbf{x}^T\mathbf{x}$. Now $\mathbf{x}^T\mathbf{A}\mathbf{x} > 0$ and $\mathbf{x}^T\mathbf{x} > 0$ imply $\lambda > 0$. □

- If \mathbf{A} has r non-zero eigenvalues, then $\text{rank}(\mathbf{A}) \geq r$.

If \mathbf{A} is symmetric and has r non-zero eigenvalues, then $\text{rank}(\mathbf{A}) = r$.

If \mathbf{A} is idempotent and has r eigenvalues equal to 1, then $\text{rank}(\mathbf{A}) = \text{tr}(\mathbf{A}) = r$.

- If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is positive definite and has eigen-decomposition

$$\mathbf{A} = \mathbf{U} \text{diag}(\lambda_1, \dots, \lambda_n) \mathbf{U}^T,$$

where the eigenvalues λ_i are positive. Then the inverse is

$$\mathbf{A}^{-1} = \mathbf{U} \text{diag}(\lambda_1^{-1}, \dots, \lambda_n^{-1}) \mathbf{U}^T.$$

- If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is psd with rank $r < n$ and has eigen-decomposition

$$\mathbf{A} = \mathbf{U} \text{diag}(\lambda_1, \dots, \lambda_r, 0, \dots, 0) \mathbf{U}^T,$$

where $\lambda_i > 0$, $i = 1, \dots, r$, are positive eigenvalues. Then the Moore-Penrose inverse is

$$\mathbf{A}^+ = \mathbf{U} \text{diag}(\lambda_1^{-1}, \dots, \lambda_r^{-1}, 0, \dots, 0) \mathbf{U}^T = \mathbf{U}_r \text{diag}(\lambda_1^{-1}, \dots, \lambda_r^{-1}) \mathbf{U}_r^T,$$

where \mathbf{U}_r contains the first r columns of \mathbf{U} . Note \mathbf{A}^+ is psd too.

- If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is psd, then

$$\mathbf{A}^{1/2} = \mathbf{U} \text{diag}(\lambda_1^{1/2}, \dots, \lambda_r^{1/2}, 0, \dots, 0) \mathbf{U}^T = \mathbf{U}_r \text{diag}(\lambda_1^{1/2}, \dots, \lambda_r^{1/2}) \mathbf{U}_r^T$$

is a symmetric square root of \mathbf{A} and $\mathbf{A}^{1/2}$ is psd too.

- $\mathbf{A} \in \mathbb{R}^{n \times n}$ a square matrix (not required to be symmetric), then $\text{tr}(\mathbf{A}) = \sum_i \lambda_i$ and $\det(\mathbf{A}) = \prod_i \lambda_i$.

Aitken model with linear constraints and a singular \mathbf{V} (alternative approach)

Here we derive an alternative expression for the MVAUE of Aitken model, $\mathbf{E}(\mathbf{y}) = \mathbf{X}\mathbf{b}$ and $\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{V}$, with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$ and a possibly singular \mathbf{V} .

- First we show that a singular \mathbf{V} essentially imposes linear constraints on \mathbf{b} .

Suppose \mathbf{V} is positive semi-definite with rank s . From eigen-decomposition,

$$\begin{aligned}
\mathbf{V} &= \mathbf{U}\mathbf{D}\mathbf{U}^T \\
&= (\mathbf{S}_{n \times s}, \mathbf{T}_{n \times (n-s)}) \begin{pmatrix} d_1 & & & & \\ & \ddots & & & \\ & & d_s & & \\ & & & 0 & \\ & & & & \ddots \\ & & & & & 0 \end{pmatrix} \begin{pmatrix} \mathbf{S}_{n \times s}^T \\ \mathbf{T}_{n \times (n-s)}^T \end{pmatrix} \\
&= \mathbf{S} \text{diag}(d_1, \dots, d_s) \mathbf{S}^T \\
&= \mathbf{S} \mathbf{D}_s \mathbf{S}^T,
\end{aligned}$$

where d_1, \dots, d_s are positive eigenvalues, $\mathbf{S}^T \mathbf{S} = \mathbf{I}_s$, $\mathbf{T}^T \mathbf{T} = \mathbf{I}_{n-s}$, and $\mathbf{S}^T \mathbf{T} = \mathbf{0}_{s \times (n-s)}$. $\mathcal{C}(\mathbf{S}) = \mathcal{C}(\mathbf{V})$ and $\mathcal{C}(\mathbf{T}) = \mathcal{N}(\mathbf{V}^T) = \mathcal{N}(\mathbf{V})$. Note that

$$\mathbf{V}^+ = \mathbf{S} \mathbf{D}_s^{-1} \mathbf{S}^T,$$

is the Moore-Penrose inverse of \mathbf{V} .

Let's transform the original model $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$ by $\begin{pmatrix} \mathbf{S}^T \\ \mathbf{T}^T \end{pmatrix}$:

$$\begin{aligned}
\mathbf{S}^T \mathbf{y} &= \mathbf{S}^T \mathbf{X}\mathbf{b} + \mathbf{S}^T \mathbf{e} = \mathbf{S}^T \mathbf{X}\mathbf{b} + \mathbf{u} \\
\mathbf{T}^T \mathbf{y} &= \mathbf{T}^T \mathbf{X}\mathbf{b} + \mathbf{T}^T \mathbf{e} = \mathbf{T}^T \mathbf{X}\mathbf{b} + \mathbf{v},
\end{aligned}$$

where

$$\begin{aligned}
\mathbf{E}(\mathbf{u}) &= \mathbf{0}_s, & \text{Cov}(\mathbf{u}) &= \sigma^2 \mathbf{D}_s \\
\mathbf{E}(\mathbf{v}) &= \mathbf{0}_{n-s}, & \text{Cov}(\mathbf{v}) &= \mathbf{0}_{(n-s) \times (n-s)}.
\end{aligned}$$

This shows that the singularity of \mathbf{V} essentially introduces linear constraints to the parameters

$$\mathbf{T}^T \mathbf{y} = \mathbf{T}^T \mathbf{X}\mathbf{b}.$$

- Next we add explicit constraints. $\mathbf{R}\mathbf{b} = \mathbf{r}$ is the explicit constraints and singularity of \mathbf{V} imposes implicit constraints $\mathbf{T}^T \mathbf{X}\mathbf{b} = \mathbf{T}^T \mathbf{y}$. Collectively we write this as

$$\mathbf{R}_0 \mathbf{b} = \mathbf{r}_0,$$

where

$$\mathbf{R}_0 = \begin{pmatrix} \mathbf{T}^T \mathbf{X} \\ \mathbf{R} \end{pmatrix} \quad \text{and} \quad \mathbf{r}_0 = \begin{pmatrix} \mathbf{T}^T \mathbf{y} \\ \mathbf{r} \end{pmatrix}.$$

- In summary, Aitken model with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$ and a singular \mathbf{V} is equivalent to the constrained model

$$\mathbf{S}^T \mathbf{y} = \mathbf{S}^T \mathbf{X} \mathbf{b} + \mathbf{u},$$

with $E(\mathbf{u}) = \mathbf{0}_n$, $\text{Cov}(\mathbf{u}) = \sigma^2 \mathbf{D}_s$, and linear constraints

$$\mathbf{R}_0 \mathbf{b} = \mathbf{r}_0.$$

- (Estimability and MVAUE) Consider the Aitken model ($E(\mathbf{y}) = \mathbf{X}\mathbf{b}$, $\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{V}$) with linear constraints $\mathbf{R}\mathbf{b} = \mathbf{r}$. A linear function $\mathbf{\Lambda}\mathbf{b}$ is estimable if and only if $\mathcal{C}(\mathbf{\Lambda}^T) \subset \mathcal{C}((\mathbf{X}^T, \mathbf{R}^T))$. And the best (minimum) variance affine unbiased estimator (MVAUE) is

$$\widehat{\mathbf{\Lambda}\mathbf{b}} = \mathbf{\Lambda} \mathbf{G}^{-1} \mathbf{X}^T \mathbf{V}^+ \mathbf{y} + \mathbf{\Lambda} \mathbf{G}^{-1} \mathbf{R}_0^T (\mathbf{R}_0 \mathbf{G}^{-1} \mathbf{R}_0^T)^{-1} (\mathbf{r}_0 - \mathbf{R}_0 \mathbf{G}^{-1} \mathbf{X}^T \mathbf{V}^+ \mathbf{y}),$$

where

$$\mathbf{R}_0 = \begin{pmatrix} \mathbf{T}^T \mathbf{X} \\ \mathbf{R} \end{pmatrix}, \quad \mathbf{r}_0 = \begin{pmatrix} \mathbf{T}^T \mathbf{y} \\ \mathbf{r} \end{pmatrix}, \quad \mathbf{G} = \mathbf{X}^T \mathbf{V}^+ \mathbf{X} + \mathbf{R}_0^T \mathbf{R}_0,$$

and \mathbf{T} is a matrix of maximum rank such that $\mathbf{V}\mathbf{T} = \mathbf{0}$. The variance matrix of $\widehat{\mathbf{\Lambda}\mathbf{b}}$ is

$$\text{Cov}(\widehat{\mathbf{\Lambda}\mathbf{b}}) = \sigma^2 \mathbf{\Lambda} \mathbf{G}^{-1} \mathbf{\Lambda}^T - \sigma^2 \mathbf{\Lambda} \mathbf{G}^{-1} \mathbf{R}_0^T (\mathbf{R}_0 \mathbf{G}^{-1} \mathbf{R}_0^T)^{-1} \mathbf{R}_0 \mathbf{G}^{-1} \mathbf{\Lambda}^T.$$

Proof. Use previous result for Aitken model with linear constraints and a non-singular \mathbf{V} . □

Normal distribution (JM 5.2)

- So far we studied the following linear models
 - Method of least squares: find \mathbf{b} that minimizes $\|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2$ (approximation)

- Linear mean model: $E(\mathbf{y}) = \mathbf{X}\mathbf{b}$ (unbiasedness, estimability)
- Gauss-Markov model: $E(\mathbf{y}) = \mathbf{X}\mathbf{b}$, $\text{Cov}(\mathbf{y}) = \sigma^2\mathbf{I}$ (MVAUE, BLUE)
- Aitken model: $E(\mathbf{y}) = \mathbf{X}\mathbf{b}$, $\text{Cov}(\mathbf{y}) = \sigma^2\mathbf{V}$ (MVAUE, BLUE)

Next we add distribution assumption $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{V})$ so we can do more: confidence interval, testing, best estimator of σ^2, \dots

- Moment generating function (mgf) of random variables.

- Assume $X \sim F_X(t) = P(X \leq t)$. If $E(e^{tX}) < \infty$ for all t in a neighborhood of 0, i.e., $|t| < \delta$ for some $\delta > 0$, then the function

$$m_X(t) = E(e^{tX})$$

is defined for all t such that $|t| < \delta$. $m_X(t)$ is the *moment generating function* (mgf) of X .

- If $m_X(t)$ exists, then $E(|X|^j) < \infty$ for $j = 0, 1, \dots$
- Differentiating mgf and setting $t = 0$ yields moments: $m_X(0) = 1$, $m'_X(0) = E(X)$, $m''_X(0) = E(X^2)$, ...
In general, $m_X^{(j)}(0) = E(X^j)$, $j = 0, 1, 2, \dots$
- (mgf determines distribution) If X_1 and X_2 are random variables with mgf $m_{X_1}(t)$ and $m_{X_2}(t)$, then $m_{X_1}(t) = m_{X_2}(t)$ for all t in a neighborhood of 0 if and only if $F_{X_1}(t) = F_{X_2}(t)$ for all t .
- If X_1, X_2, \dots, X_n are independent with mgfs $m_{X_1}(t), \dots, m_{X_n}(t)$ and

$$Y = a_0 + \sum_{i=1}^n a_i X_i,$$

then

$$m_Y(t) = e^{ta_0} \prod_{i=1}^n m_{X_i}(a_i t).$$

- Moment generating function (mgf) of random vectors.

- Let $\mathbf{X} = (X_1, \dots, X_p)^T \in \mathbb{R}^p$ be a random vector. Then

$$m_{\mathbf{X}}(\mathbf{t}) = E(e^{\mathbf{t}^T \mathbf{X}}) = E(e^{t_1 X_1 + \dots + t_p X_p}), \quad \mathbf{t} \in \mathbb{R}^p$$

is the mgf of \mathbf{X} provided that $E(e^{\mathbf{t}^T \mathbf{X}}) < \infty$ for all \mathbf{t} such that $\|\mathbf{t}\| < \delta$ for some $\delta > 0$.

- If $m_{\mathbf{X}}(\mathbf{t})$ exists, then $E(|X_1|^{a_1} \cdots |X_p|^{a_p}) < \infty$ for $a_j = 0, 1, 2, \dots$. Especially $E(|X_j|^{a_j}) < \infty$, $a_j = 0, 1, 2, \dots$. Also

$$m_{X_j}(t) = m_{\mathbf{X}}(t\mathbf{e}_j) = m_{\mathbf{X}}((0, \dots, t, \dots, 0)^T).$$

- Differentiating mgf and setting $\mathbf{t} = \mathbf{0}$ yields moments. Let $\mathbf{n} = (n_1, \dots, n_p)$ and $\mathbf{t}^{\mathbf{n}} = \prod_{i=1}^p t_i^{n_i}$. Then

$$\left. \frac{\partial^{\sum_{i=1}^p n_i}}{\partial \mathbf{t}^{\mathbf{n}}} m_{\mathbf{X}}(\mathbf{t}) \right|_{\mathbf{t}=\mathbf{0}} = E\left(\prod_{i=1}^p X_i^{n_i}\right).$$

The first two moments are

$$\begin{aligned} \nabla m_{\mathbf{X}}(\mathbf{t})|_{\mathbf{t}=\mathbf{0}} &= E(\mathbf{X}) \\ d^2 m_{\mathbf{X}}(\mathbf{t})|_{\mathbf{t}=\mathbf{0}} &= E(\mathbf{X}\mathbf{X}^T). \end{aligned}$$

- (mgf determines distribution) $m_{\mathbf{X}}(\mathbf{t}) = m_{\mathbf{Y}}(\mathbf{t})$ for all $\|\mathbf{t}\| < \delta$ ($\delta > 0$) if and only if $F_{\mathbf{X}}(\mathbf{t}) = F_{\mathbf{Y}}(\mathbf{t})$ for all \mathbf{t} .

- Suppose that \mathbf{X} is partitioned as $\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_m \end{pmatrix}$ and has mgf $m_{\mathbf{X}}(\mathbf{t}) = E(e^{\mathbf{t}^T \mathbf{X}})$. Then $\mathbf{X}_1, \dots, \mathbf{X}_m$ are independent if and only if

$$m_{\mathbf{X}}(\mathbf{t}) = \prod_{i=1}^m m_{\mathbf{X}_i}(\mathbf{t}_i)$$

for all $\mathbf{t} = \begin{pmatrix} \mathbf{t}_1 \\ \vdots \\ \mathbf{t}_m \end{pmatrix}$ in a neighborhood of $\mathbf{0}$.

- If $\mathbf{X}_1, \dots, \mathbf{X}_n$ are independent random vectors and

$$\mathbf{Y} = \mathbf{A}_0 + \sum_{i=1}^n \mathbf{A}_i \mathbf{X}_i,$$

then

$$m_{\mathbf{Y}}(\mathbf{t}) = e^{\mathbf{t}^T \mathbf{A}_0} \prod_{i=1}^n m_{\mathbf{X}_i}(\mathbf{A}_i^T \mathbf{t}).$$

- A random variable Z has a *standard normal distribution*, denoted $Z \sim N(0, 1)$, if

$$F_Z(t) = P(Z \leq t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz,$$

or equivalently Z has density

$$f_Z(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}, \quad -\infty < z < \infty.$$

or equivalently

$$m_Z(t) = E(e^{tZ}) = e^{t^2/2}, \quad -\infty < t < \infty.$$

- Non-standard normal random variable.
 - Definition 1: A random variable X has *normal distribution* with mean μ and variance σ^2 , denoted $X \sim N(\mu, \sigma^2)$, if

$$X \stackrel{D}{=} \mu + \sigma Z,$$

where $Z \sim N(0, 1)$.

- Definition 2: $X \sim N(\mu, \sigma^2)$ if

$$m_X(t) = E(e^{tX}) = e^{t\mu + \sigma^2 t^2/2}, \quad -\infty < t < \infty.$$

- In both definitions, $\sigma^2 = 0$ is allowed. If $\sigma^2 > 0$, it has a density

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}, \quad -\infty < x < \infty.$$

- The *standard multivariate normal* is a vector of independent standard normals, denoted $\mathbf{Z} \sim N_p(\mathbf{0}_p, \mathbf{I}_p)$. The joint density is

$$f_{\mathbf{Z}}(\mathbf{z}) = \frac{1}{(2\pi)^{p/2}} e^{-\sum_{i=1}^p z_i^2/2}.$$

The mgf is

$$m_{\mathbf{Z}}(\mathbf{t}) = \prod_{i=1}^p m_{Z_i}(t_i) = \prod_{i=1}^p e^{t_i^2/2} = e^{\sum_{i=1}^p t_i^2/2} = e^{\mathbf{t}^T \mathbf{t}/2}.$$

- Consider the affine transformation $\mathbf{X} = \boldsymbol{\mu} + \mathbf{AZ}$, where $\mathbf{Z} \sim N_p(\mathbf{0}_p, \mathbf{I}_p)$. \mathbf{X} has mean and variance

$$\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}, \quad \text{Var}(\mathbf{X}) = \mathbf{AA}^T$$

and the moment generating function is

$$m_{\mathbf{X}}(\mathbf{t}) = \mathbb{E}(e^{\mathbf{t}^T(\boldsymbol{\mu} + \mathbf{AZ})}) = e^{\mathbf{t}^T \boldsymbol{\mu}} \mathbb{E}e^{\mathbf{t}^T \mathbf{AZ}} = e^{\mathbf{t}^T \boldsymbol{\mu} + \mathbf{t}^T \mathbf{AA}^T \mathbf{t} / 2}.$$

- $\mathbf{X} \in \mathbb{R}^p$ has a *multivariate normal distribution* with mean $\boldsymbol{\mu} \in \mathbb{R}^p$ and covariance $\mathbf{V} \in \mathbb{R}^{p \times p}$ (\mathbf{V} psd), denoted $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$, if its mgf takes the form

$$m_{\mathbf{X}}(\mathbf{t}) = e^{\mathbf{t}^T \boldsymbol{\mu} + \mathbf{t}^T \mathbf{V} \mathbf{t} / 2}, \quad \mathbf{t} \in \mathbb{R}^p.$$

Remarks

- We already see any affine transform of a multivariate standard normal is normal. Conversely, any multivariate normal $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$ is also an affine transform of a multivariate standard normal. Just take \mathbf{A} to be the Cholesky factor or symmetric square root of \mathbf{V} .
- \mathbf{V} can be singular, in which case the density does not exist. Suppose $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$ and \mathbf{V} is psd with rank $s < p$. From eigen-decomposition

$$\begin{aligned} \mathbf{V} &= \mathbf{UDU}^T \\ &= (\mathbf{S}_{p \times s}, \mathbf{T}_{p \times (p-s)}) \begin{pmatrix} d_1 & & & & & \\ & \ddots & & & & \\ & & d_s & & & \\ & & & 0 & & \\ & & & & \ddots & \\ & & & & & 0 \end{pmatrix} \begin{pmatrix} \mathbf{S}_{p \times s}^T \\ \mathbf{T}_{p \times (p-s)}^T \end{pmatrix} \\ &= \mathbf{S} \text{diag}(d_1, \dots, d_s) \mathbf{S}^T \\ &= \mathbf{SD}_s \mathbf{S}^T, \end{aligned}$$

where d_1, \dots, d_s are positive eigenvalues, $\mathbf{S}^T \mathbf{S} = \mathbf{I}_s$, $\mathbf{T}^T \mathbf{T} = \mathbf{I}_{p-s}$, and $\mathbf{S}^T \mathbf{T} = \mathbf{0}_{s \times (p-s)}$. Then

$$\begin{aligned} \mathbf{S}^T \mathbf{X} &\sim N_s(\mathbf{S}^T \boldsymbol{\mu}, \mathbf{D}_s) \\ \mathbf{T}^T \mathbf{X} &= \mathbf{T}^T \boldsymbol{\mu}. \end{aligned}$$

The second set of equations indicates that a singular \mathbf{V} imposes constraints on the values \mathbf{X} can take. The probability mass lies in a subspace of dimension s .

- If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$ and \mathbf{V} is non-singular, then

- $\mathbf{V} = \mathbf{A}\mathbf{A}^T$ for some non-singular \mathbf{A}

- $\mathbf{A}^{-1}(\mathbf{X} - \boldsymbol{\mu}) \sim N_p(\mathbf{0}_p, \mathbf{I}_p)$

- The density of \mathbf{X} is

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\mathbf{V}|^{1/2}} e^{-(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{V}^{-1} (\mathbf{x} - \boldsymbol{\mu}) / 2}.$$

Proof. The first fact follows by taking \mathbf{A} to be the Cholesky factor or square root of \mathbf{V} . The second fact is trivial. For the third fact, we use the change of variable formula. Let $\mathbf{Z} \sim N_p(\mathbf{0}_p, \mathbf{I}_p)$. By the change of variable formula $\mathbf{X} = \boldsymbol{\mu} + \mathbf{A}\mathbf{Z}$ has density

$$\begin{aligned} f_{\mathbf{X}}(\mathbf{x}) &= \frac{1}{(2\pi)^{p/2}} e^{-(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{A}^{-T} \mathbf{A}^{-1} (\mathbf{x} - \boldsymbol{\mu}) / 2} |\det(\mathbf{A}^{-1})| \\ &= \frac{1}{(2\pi)^{p/2} |\mathbf{V}|^{1/2}} e^{-(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{V}^{-1} (\mathbf{x} - \boldsymbol{\mu}) / 2}. \end{aligned}$$

□

- (Any affine transform of normal is normal) If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$ and $\mathbf{Y} = \mathbf{a} + \mathbf{B}\mathbf{X}$, where $\mathbf{a} \in \mathbb{R}^q$ and $\mathbf{B} \in \mathbb{R}^{q \times p}$, then $\mathbf{Y} \sim N_q(\mathbf{a} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\mathbf{V}\mathbf{B}^T)$.

Proof. Check the mgf. □

- (Marginal of normal is normal) If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$, then any subvector of \mathbf{X} is normal too.

Proof. Immediate corollary of the preceding result. □

- A convenient fact about normal random variables/vectors is that zero correlation/covariance implies independence.

If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$ and is partitioned as

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_m \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \vdots \\ \boldsymbol{\mu}_m \end{pmatrix}, \quad \mathbf{V} = \begin{pmatrix} \mathbf{V}_{11} & \cdots & \mathbf{V}_{1m} \\ \vdots & & \vdots \\ \mathbf{V}_{m1} & \cdots & \mathbf{V}_{mm} \end{pmatrix},$$

then $\mathbf{X}_1, \dots, \mathbf{X}_m$ are jointly independent if and only if $\mathbf{V}_{ij} = \mathbf{0}$ for all $i \neq j$.

Proof. If $\mathbf{X}_1, \dots, \mathbf{X}_m$ are jointly independent, then $\mathbf{V}_{ij} = \text{Cov}(\mathbf{X}_i, \mathbf{X}_j) = \text{E}(\mathbf{X}_i - \boldsymbol{\mu}_i)(\mathbf{X}_j - \boldsymbol{\mu}_j)^T = \text{E}(\mathbf{X}_i - \boldsymbol{\mu}_i)\text{E}(\mathbf{X}_j - \boldsymbol{\mu}_j)^T = \mathbf{0}_{p_i}\mathbf{0}_{p_j}^T = \mathbf{0}_{p_i \times p_j}$. Conversely, if $\mathbf{V}_{ij} = \mathbf{0}$ for all $i \neq j$, then the mgf of $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_m)^T$ is

$$\begin{aligned} m_{\mathbf{X}}(\mathbf{t}) &= e^{\mathbf{t}^T \boldsymbol{\mu} + \mathbf{t}^T \mathbf{V} \mathbf{t} / 2} \\ &= e^{\sum_{i=1}^m \mathbf{t}_i^T \boldsymbol{\mu}_i + \sum_{i=1}^m \mathbf{t}_i^T \mathbf{V}_i \mathbf{t}_i / 2} \\ &= m_{\mathbf{X}_1}(\mathbf{t}_1) \cdots m_{\mathbf{X}_m}(\mathbf{t}_m). \end{aligned}$$

Thus $\mathbf{X}_1, \dots, \mathbf{X}_m$ are jointly independent. \square

- Let $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$, $\mathbf{Y}_1 = \mathbf{a}_1 + \mathbf{B}_1 \mathbf{X}$, and $\mathbf{Y}_2 = \mathbf{a}_2 + \mathbf{B}_2 \mathbf{X}$. Then \mathbf{Y}_1 and \mathbf{Y}_2 are independent if and only if $\mathbf{B}_1 \mathbf{V} \mathbf{B}_2^T = \mathbf{0}$.

Proof. Note $\text{Cov}(\mathbf{Y}_1, \mathbf{Y}_2) = \mathbf{B}_1 \text{Cov}(\mathbf{X}) \mathbf{B}_2^T = \mathbf{B}_1 \mathbf{V} \mathbf{B}_2^T$. \square

Chi-square and related distributions (JM 5.3)

- Let $\mathbf{Z} \sim N_p(\mathbf{0}_p, \mathbf{I}_p)$. Then $U = \|\mathbf{Z}\|_2^2 = \sum_{i=1}^p Z_i^2$ has the *chi-square distribution* with p degrees of freedom, denoted by $U \sim \chi_p^2$.

– The mgf of U is

$$\begin{aligned} m_U(t) &= \text{E}(e^{tU}) = \text{E}\left(e^{t \sum_{i=1}^p Z_i^2}\right) \\ &= \prod_{i=1}^p \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{tz_i^2 - z_i^2/2} dz_i \\ &= \prod_{i=1}^p \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z_i^2}{2(1-2t)}} dz_i \\ &= \prod_{i=1}^p (1-2t)^{-1/2} \\ &= (1-2t)^{-p/2}. \end{aligned}$$

– The first two moments of U are

$$\mathbb{E}(U) = p, \quad \text{Var}(U) = 2p.$$

– The density of U is

$$f_U(u) = \frac{u^{(p-2)/2} e^{-u/2}}{\Gamma(p/2) 2^{p/2}}, \quad u > 0.$$

- Let $J \sim \text{Poisson}(\phi)$ and $U|J = j \sim \chi_{p+2j}^2$. Then the unconditional distribution of U is the *noncentral chi-square distribution* with noncentrality parameter ϕ , denoted by $U \sim \chi_p^2(\phi)$.

– The density of U is

$$f_U(u) = \sum_{j=0}^{\infty} e^{-\phi} \frac{\phi^j}{j!} \times \frac{u^{(p+2j-2)/2} e^{-u/2}}{\Gamma((p+2j)/2) 2^{(p+2j)/2}}, \quad u > 0.$$

– The mgf of U is

$$\begin{aligned} m_U(t) &= \mathbb{E}(e^{tU}) = \mathbb{E}[\mathbb{E}(e^{tU}|J)] \\ &= \mathbb{E}(1 - 2t)^{-(p+2J)/2} \\ &= (1 - 2t)^{-p/2} \mathbb{E}[(1 - 2t)^{-J}] \\ &= (1 - 2t)^{-p/2} e^{2\phi t/(1-2t)}. \end{aligned}$$

– The first two moments of U are

$$\mathbb{E}(U) = p + 2\phi, \quad \text{Var}(U) = 2p + 8\phi.$$

- If $U_i \sim \chi_{p_i}^2(\phi_i)$, $i = 1, \dots, m$, are jointly independent, then $U = \sum_{i=1}^m U_i \sim \chi_p(\phi)$, where $p = \sum_{i=1}^m p_i$ and $\phi = \sum_{i=1}^m \phi_i$.

Proof. The mgf of U is

$$m_U(t) = \prod_{i=1}^m m_{U_i}(t) = \prod_{i=1}^m (1 - 2t)^{-p_i/2} e^{2t\phi_i/(1-2t)} = (1 - 2t)^{-p/2} e^{2t\phi/(1-2t)}.$$

□

- If $X \sim N(\mu, 1)$, then $U = X^2 \sim \chi_1^2(\mu^2/2)$.

Proof. The mgf of U is

$$\begin{aligned}
m_U(t) &= \mathbb{E}(e^{tU}) = \mathbb{E}(e^{tX^2}) \\
&= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{tx^2 - (x-\mu)^2/2} \\
&= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1-2t}{2}(x - \frac{\mu}{1-2t})^2 - \mu^2 + \frac{\mu^2}{1-2t}} \\
&= (1-2t)^{-1/2} \times e^{(\mu^2/2)2t/(1-2t)},
\end{aligned}$$

which matches that of $\chi_1^2(\mu^2/2)$. □

- If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{I}_p)$, then $W = \|\mathbf{X}\|_2^2 = \sum_{i=1}^p X_i^2 \sim \chi_p^2(\frac{1}{2}\boldsymbol{\mu}^T \boldsymbol{\mu})$.

Proof. Immediate corollary of the preceding two results. □

- If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$ and \mathbf{V} is non-singular. Then $W = \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \sim \chi_p^2(\frac{1}{2}\boldsymbol{\mu}^T \mathbf{V}^{-1} \boldsymbol{\mu})$.

Proof. Let $\mathbf{V} = \mathbf{A}\mathbf{A}^T$. Then $\mathbf{Z} = \mathbf{A}^{-1}\mathbf{X} \sim N_p(\mathbf{A}^{-1}\boldsymbol{\mu}, \mathbf{I}_p)$. And

$$W = \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} = \mathbf{Z}^T \mathbf{Z} \sim \chi_p^2(\frac{1}{2}\boldsymbol{\mu}^T \mathbf{V}^{-1} \boldsymbol{\mu}).$$

□

- (Stochastic monotonicity of non-central chi-square) Let $U \sim \chi_p^2(\phi)$. Then $P(U > c)$ is increasing in ϕ for fixed p and $c > 0$.

Proof. For an algebraic proof, see JM Result 5.11 (p106). For a probabilistic proof, we fix p and assume $\phi_2 > \phi_1 > 0$. Let X_1, X_2, \dots be a sequence of independent standard normals, $J_1 \sim \text{Poisson}(\phi_1)$, $\Delta J \sim \text{Poisson}(\phi_2 - \phi_1)$, and $X_1, X_2, \dots, J_1, \Delta J$ are jointly independent. Let $J_2 = J_1 + \Delta J \sim \text{Poisson}(\phi_2)$. Define

$$\begin{aligned}
U_1 | J_1 &= \sum_{i=1}^{p+2J_1} X_i^2 \\
U_2 | J_2 &= \sum_{i=1}^{p+2J_2} X_i^2.
\end{aligned}$$

Then $U_1 \sim \chi_p^2(\phi_1)$, $U_2 \sim \chi_p^2(\phi_2)$, and $U_2 \geq U_1$ almost surely. Hence $P(U_2 > c) \geq P(U_1 > c)$. □

- Let $U_1 \sim \chi_{p_1}^2$ and $U_2 \sim \chi_{p_2}^2$ be independent. Then $F = \frac{U_1/p_1}{U_2/p_2}$ has the *F-distribution* with p_1 and p_2 degrees of freedom, denoted by $F \sim F_{p_1, p_2}$.

– Then density of F is

$$f_F(f) = \frac{\Gamma\left(\frac{p_1+p_2}{2}\right) \left(\frac{p_2}{p_1}\right)^{p_1/2}}{\Gamma\left(\frac{p_1}{2}\right)\Gamma\left(\frac{p_2}{2}\right)} f^{p_1/2-1} \left(1 + \frac{p_1}{p_2}f\right)^{-(p_1+p_2)/2}, \quad f > 0.$$

– Not all moments of F exist.

$$\begin{aligned} E(F) &= \frac{p_2}{p_2 - 2} \text{ for } p_2 > 2 \\ \text{Var}(F) &= \frac{2p_2^2(p_1 + p_2 - 2)}{p_1(p_2 - 2)^2(p_2 - 4)} \text{ for } p_2 > 4. \end{aligned}$$

- Let $U_1 \sim \chi_{p_1}^2(\phi)$ and $U_2 \sim \chi_{p_2}^2$ be independent. Then $F = \frac{U_1/p_1}{U_2/p_2}$ has the *noncentral F-distribution* with p_1 and p_2 degrees of freedom and *noncentrality* parameter ϕ , denoted by $F \sim F_{p_1, p_2}(\phi)$.
- (Stochastic monotonicity of noncentral F) Let $W \sim F_{p_1, p_2}(\phi)$. Then $P(W > c)$ is strictly increasing in ϕ for fixed p_1 and p_2 .

Proof. It follows from stochastic monotonicity of non-central chi-square. \square

- Let $U \sim N(\mu, 1)$ and $V \sim \chi_k^2$ be independent. Then $T = U/\sqrt{V/k}$ has the *noncentral Student's t-distribution* with k degrees of freedom and noncentrality parameter μ , denoted $T \sim t_k(\mu)$.

If $T \sim t_k(\mu)$, then $T^2 = F_{1, k}(\mu^2/2)$.

If $\mu = 0$, the distribution is the *Student's t distribution*, denoted by $T \sim t_k$, and has density

$$f_T(t) = \frac{\Gamma((k+1)/2)}{\Gamma(k/2)\sqrt{\pi k}} \times (1 + t^2/k)^{-(k+1)/2}.$$

13 Lecture 13: Oct 21

Announcement

- HW5 posted and due this Fri Oct 25.

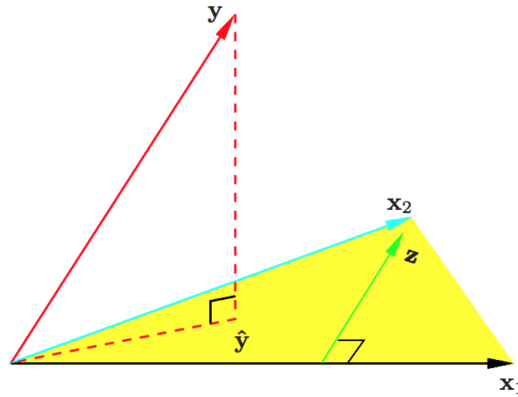
Last time

- Multivariate normal $N_p(\boldsymbol{\mu}, \mathbf{V})$, mgf $m(\mathbf{t}) = e^{\mathbf{t}^T \boldsymbol{\mu} + \mathbf{t}^T \mathbf{V} \mathbf{t} / 2}$, correlation and independence, existence of density when \mathbf{V} is non-singular.
- (Central) Chi-square distribution: If $\mathbf{Z} \sim N_p(\mathbf{0}_p, \mathbf{I}_p)$, then $U = \|\mathbf{Z}\|_2^2 \sim \chi_p^2$.
- Noncentral chi-square distribution:
 - If $\mathbf{Z} \sim N_p(\boldsymbol{\mu}, \mathbf{I}_p)$, then $U = \|\mathbf{Z}\|_2^2 \sim \chi_p^2(\boldsymbol{\mu}^T \boldsymbol{\mu} / 2)$.
 - If $\mathbf{Z} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$, then $U = \mathbf{Z}^T \mathbf{V}^{-1} \mathbf{Z} \sim \chi_p^2(\boldsymbol{\mu}^T \mathbf{V}^{-1} \boldsymbol{\mu} / 2)$.
 - $U \sim \chi_p^2(\phi)$ is stochastic monotone in p and ϕ .
 - If $U_i \sim \chi_{p_i}^2(\phi_i)$ are independent, then $\sum_i U_i \sim \chi_{\sum_i p_i}^2(\sum_i \phi_i)$.
- (Central) F distribution: If $U_1 \sim \chi_{p_1}^2$ and $U_2 \sim \chi_{p_2}^2$, then $(U_1/p_1)/(U_2/p_2) \sim F_{p_1, p_2}$.
- Noncentral F distribution: If $U_1 \sim \chi_{p_1}^2(\phi)$ and $U_2 \sim \chi_{p_2}^2$, then $F = (U_1/p_1)/(U_2/p_2) \sim F_{p_1, p_2}(\phi)$. $F_{p_1, p_2}(\phi)$ is stochastic monotone in ϕ .
- Noncentral t distribution: If $U \sim N(\mu, 1)$ and $V \sim \chi_p^2$, then $T = U/\sqrt{V/p} \sim t_p(\mu)$.

Today

- Distribution of quadratic forms (JM 5.4)
- Cochran's theorem (JM 5.5)

Distribution of quadratic forms (JM 5.4)



- Motivation: From the geometry of least squares problem, we know

$$\mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{e}} = \mathbf{P}_X \mathbf{y} + (\mathbf{I} - \mathbf{P}_X) \mathbf{y}$$

and

$$\|\mathbf{y}\|_2^2 = \|\hat{\mathbf{y}}\|_2^2 + \|\hat{\mathbf{e}}\|_2^2 = \mathbf{y}^T \mathbf{P}_X \mathbf{y} + \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}.$$

Now assuming $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I}_n)$, we would like to know the distribution of the sums of squares $\mathbf{y}^T \mathbf{P}_X \mathbf{y}$ and $\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}$.

- A symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is idempotent with rank s if and only if there exists a matrix $\mathbf{G} \in \mathbb{R}^{n \times s}$ with orthonormal columns, that is $\mathbf{G}^T \mathbf{G} = \mathbf{I}_s$, such that $\mathbf{A} = \mathbf{G}\mathbf{G}^T$.

Proof. The “if part” is easy. For the “only if part”, recall that the eigenvalues of an idempotent matrix are either 1 or 0 and $\text{rank}(\mathbf{A}) = s$ equals the number of nonzero eigenvalues. Thus by the eigen-decomposition,

$$\mathbf{A} = (\mathbf{Q}_1, \mathbf{Q}_2) \begin{pmatrix} \mathbf{I}_s & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{Q}_1^T \\ \mathbf{Q}_2^T \end{pmatrix} = \mathbf{Q}_1 \mathbf{Q}_1^T,$$

where $\mathbf{Q}_1 \in \mathbb{R}^{n \times s}$ and $\mathbf{Q}_2 \in \mathbb{R}^{n \times (n-s)}$. □

- Let $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{I}_p)$ and \mathbf{A} be symmetric and idempotent with rank s . Then $\mathbf{X}^T \mathbf{A} \mathbf{X} \sim \chi_s^2(\frac{1}{2} \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu})$.

Proof. By the preceding result, $\mathbf{A} = \mathbf{G}\mathbf{G}^T$, where $\mathbf{G} \in \mathbb{R}^{p \times s}$ and $\mathbf{G}\mathbf{G}^T = \mathbf{I}_s$. Then $\mathbf{G}^T \mathbf{X} \sim N_s(\mathbf{G}^T \boldsymbol{\mu}, \mathbf{I}_s)$ and we have

$$\mathbf{X}^T \mathbf{A} \mathbf{X} = (\mathbf{G}^T \mathbf{X})^T \mathbf{G}^T \mathbf{X} = \|\mathbf{G}^T \mathbf{X}\|_2^2 \sim \chi_s^2 \left(\frac{1}{2} \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} \right).$$

□

- (General case) Let $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$ with \mathbf{V} non-singular, and $\mathbf{A} \in \mathbb{R}^{p \times p}$ be symmetric. If $\mathbf{A}\mathbf{V}$ is idempotent with rank s , then $\mathbf{X}^T \mathbf{A} \mathbf{X} \sim \chi_s^2(\frac{1}{2} \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu})$.

Proof. Let $\mathbf{V} = \boldsymbol{\Gamma}\boldsymbol{\Gamma}^T$, where $\boldsymbol{\Gamma}$ can be the Cholesky factor or symmetric square root of \mathbf{V} , and $\mathbf{Y} = \boldsymbol{\Gamma}^{-1} \mathbf{X} \sim N_p(\boldsymbol{\Gamma}^{-1} \boldsymbol{\mu}, \mathbf{I}_p)$. Note

$$\mathbf{X}^T \mathbf{A} \mathbf{X} = \mathbf{Y}^T \boldsymbol{\Gamma}^T \mathbf{A} \boldsymbol{\Gamma} \mathbf{Y}.$$

The preceding results applies if $\boldsymbol{\Gamma}^T \mathbf{A} \boldsymbol{\Gamma}$ is symmetric (trivial) and idempotent. Idempotency holds since

$$(\boldsymbol{\Gamma}^T \mathbf{A} \boldsymbol{\Gamma})(\boldsymbol{\Gamma}^T \mathbf{A} \boldsymbol{\Gamma}) = \boldsymbol{\Gamma}^T \mathbf{A} \mathbf{V} \mathbf{A} \boldsymbol{\Gamma} = \boldsymbol{\Gamma}^T \mathbf{A} \boldsymbol{\Gamma}.$$

The last equation is true since $\mathbf{A}\mathbf{V}$ is a projection onto $\mathcal{C}(\mathbf{A}\mathbf{V}) = \mathcal{C}(\mathbf{A})$. Also note that $\text{rank}(\boldsymbol{\Gamma}^T \mathbf{A} \boldsymbol{\Gamma}) = \text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}\mathbf{V}) = s$. Therefore by the preceding result, $\mathbf{X}^T \mathbf{A} \mathbf{X} \sim \chi_s^2(\phi)$ with $\phi = \frac{1}{2} \boldsymbol{\mu}^T \boldsymbol{\Gamma}^{-T} \boldsymbol{\Gamma}^T \mathbf{A} \boldsymbol{\Gamma} \boldsymbol{\Gamma}^{-1} \boldsymbol{\mu} = \frac{1}{2} \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}$. □

- Consider the normal linear model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I}_n)$.

– Using $\mathbf{A} = (1/\sigma^2)(\mathbf{I} - \mathbf{P}_X)$, we have

$$\text{SSE}/\sigma^2 = \|\hat{\boldsymbol{\epsilon}}\|_2^2/\sigma^2 = \mathbf{y}^T \mathbf{A} \mathbf{y} \sim \chi_{n-r}^2,$$

where $r = \text{rank}(\mathbf{X})$. Note the noncentrality parameter is

$$\phi = \frac{1}{2} (\mathbf{X}\mathbf{b})^T (1/\sigma^2) (\mathbf{I} - \mathbf{P}_X) (\mathbf{X}\mathbf{b}) = 0 \text{ for all } \mathbf{b}.$$

– Using $\mathbf{A} = (1/\sigma^2) \mathbf{P}_X$, we have

$$\text{SSR}/\sigma^2 = \|\hat{\mathbf{y}}\|_2^2/\sigma^2 = \mathbf{y}^T \mathbf{A} \mathbf{y} \sim \chi_r^2(\phi)$$

with the noncentrality parameter

$$\phi = \frac{1}{2} (\mathbf{X}\mathbf{b})^T (1/\sigma^2) \mathbf{P}_X (\mathbf{X}\mathbf{b}) = \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{b}\|_2^2.$$

– The joint distribution of $\hat{\mathbf{y}}$ and $\hat{\mathbf{e}}$ is

$$\begin{pmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{e}} \end{pmatrix} = \begin{pmatrix} \mathbf{P}_X \\ \mathbf{I}_n - \mathbf{P}_X \end{pmatrix} \mathbf{y} \sim N_{2n} \left(\begin{pmatrix} \mathbf{X}\mathbf{b} \\ \mathbf{0}_n \end{pmatrix}, \begin{pmatrix} \sigma^2 \mathbf{P}_X & \mathbf{0} \\ \mathbf{0} & \sigma^2 (\mathbf{I} - \mathbf{P}_X) \end{pmatrix} \right).$$

So $\hat{\mathbf{y}}$ is independent of $\hat{\mathbf{e}}$. Thus $\|\hat{\mathbf{y}}\|_2^2$ is independent of $\|\hat{\mathbf{e}}\|_2^2$ and

$$F = \frac{\|\hat{\mathbf{y}}\|_2^2/r}{\|\hat{\mathbf{e}}\|_2^2/(n-r)} \sim F_{r, n-r} \left(\frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{b}\|_2^2 \right).$$

- (Independence between two linear forms of a multivariate normal) In last lecture we showed the following result.

Let $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$. Then $\mathbf{a}_1 + \mathbf{B}_1\mathbf{X}$ and $\mathbf{a}_2 + \mathbf{B}_2\mathbf{X}$ are independent if and only if $\mathbf{B}_1\mathbf{V}\mathbf{B}_2^T = \mathbf{0}$.

- (Independence between linear and quadratic forms of a multivariate normal) Let $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$ and \mathbf{A} is symmetric with rank s . If $\mathbf{B}\mathbf{V}\mathbf{A} = \mathbf{0}$, then $\mathbf{B}\mathbf{X}$ and $\mathbf{X}^T\mathbf{A}\mathbf{X}$ are independent.

Proof. By eigen-decomposition, $\mathbf{A} = \mathbf{Q}_1\boldsymbol{\Lambda}_1\mathbf{Q}_1^T$, where $\mathbf{Q}_1^T\mathbf{Q}_1 = \mathbf{I}_s$ and $\boldsymbol{\Lambda}_1 \in \mathbb{R}^{s \times s}$ is non-singular. Consider the joint distribution

$$\begin{pmatrix} \mathbf{B}\mathbf{X} \\ \mathbf{Q}_1^T\mathbf{X} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{B}\boldsymbol{\mu} \\ \mathbf{Q}_1^T\boldsymbol{\mu} \end{pmatrix}, \begin{pmatrix} \mathbf{B}\mathbf{V}\mathbf{B}^T & \mathbf{B}\mathbf{V}\mathbf{Q}_1 \\ \mathbf{Q}_1^T\mathbf{V}\mathbf{B}^T & \mathbf{Q}_1^T\mathbf{V}\mathbf{Q}_1 \end{pmatrix} \right).$$

By hypothesis

$$\mathbf{B}\mathbf{V}\mathbf{A} = \mathbf{B}\mathbf{V}\mathbf{Q}_1\boldsymbol{\Lambda}_1\mathbf{Q}_1^T = \mathbf{0}.$$

Post-multiplying both sides by $\mathbf{Q}_1\boldsymbol{\Lambda}_1^{-1}$ gives $\mathbf{B}\mathbf{V}\mathbf{Q}_1 = \mathbf{0}$, which implies that $\mathbf{B}\mathbf{X}$ is independent of both $\mathbf{Q}_1^T\mathbf{X}$ and $\mathbf{X}^T\mathbf{Q}_1\boldsymbol{\Lambda}_1\mathbf{Q}_1^T\mathbf{X} = \mathbf{X}^T\mathbf{A}\mathbf{X}$. \square

- (Independence between two quadratic forms of a multivariate normal) Let $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$, \mathbf{A} be symmetric with rank r , and \mathbf{B} be symmetric with rank s . If $\mathbf{B}\mathbf{V}\mathbf{A} = \mathbf{0}$, then $\mathbf{X}^T\mathbf{A}\mathbf{X}$ and $\mathbf{X}^T\mathbf{B}\mathbf{X}$ are independent.

Proof. Again by eigen-decomposition,

$$\begin{aligned} \mathbf{A} &= \mathbf{Q}_1\boldsymbol{\Lambda}_1\mathbf{Q}_1^T, & \text{where } \mathbf{Q}_1 &\in \mathbb{R}^{p \times r}, \boldsymbol{\Lambda}_1 \in \mathbb{R}^{r \times r} \text{ nonsingular} \\ \mathbf{B} &= \mathbf{Q}_2\boldsymbol{\Lambda}_2\mathbf{Q}_2^T, & \text{where } \mathbf{Q}_2 &\in \mathbb{R}^{p \times s}, \boldsymbol{\Lambda}_2 \in \mathbb{R}^{s \times s} \text{ nonsingular.} \end{aligned}$$

Now consider the joint distribution

$$\begin{pmatrix} \mathbf{Q}_1^T \mathbf{X} \\ \mathbf{Q}_2^T \mathbf{X} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{Q}_1^T \boldsymbol{\mu} \\ \mathbf{Q}_2^T \boldsymbol{\mu} \end{pmatrix}, \begin{pmatrix} \mathbf{Q}_1^T \mathbf{V} \mathbf{Q}_2 & \mathbf{Q}_1^T \mathbf{V} \mathbf{Q}_2 \\ \mathbf{Q}_2^T \mathbf{V} \mathbf{Q}_1 & \mathbf{Q}_2^T \mathbf{V} \mathbf{Q}_2 \end{pmatrix} \right).$$

By hypothesis

$$\mathbf{B} \mathbf{V} \mathbf{A} = \mathbf{Q}_2 \boldsymbol{\Lambda}_2 \mathbf{Q}_2^T \mathbf{V} \mathbf{Q}_1 \boldsymbol{\Lambda}_1 \mathbf{Q}_1^T = \mathbf{0}.$$

Pre-multiplying both sides by $\boldsymbol{\Lambda}_2^{-1} \mathbf{Q}_2^T$ and then post-multiplying both sides by $\mathbf{Q}_1 \boldsymbol{\Lambda}_1^{-1}$ gives

$$\mathbf{Q}_2^T \mathbf{V} \mathbf{Q}_1 = \mathbf{0}.$$

Therefore $\mathbf{Q}_1^T \mathbf{X}$ is independent of $\mathbf{Q}_2^T \mathbf{X}$, which implies $\mathbf{X}^T \mathbf{A} \mathbf{X} = \mathbf{X}^T \mathbf{Q}_1 \boldsymbol{\Lambda}_1 \mathbf{Q}_1^T$ is independent of $\mathbf{X}^T \mathbf{B} \mathbf{X} = \mathbf{X}^T \mathbf{Q}_2 \boldsymbol{\Lambda}_2 \mathbf{Q}_2^T$. \square

14 Lecture 14: Oct 23

Announcement

- HW5 due this Fri Oct 25.
- Homework session this afternoon.

Last time

- Distribution of quadratic forms (JM 5.4)
 - $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \mathbf{V})$ with \mathbf{V} non-singular, $\mathbf{A} \in \mathbb{R}^{p \times p}$ be symmetric, and $\mathbf{A}\mathbf{V}$ is idempotent with rank s , then $\mathbf{X}^T \mathbf{A} \mathbf{X} \sim \chi_s^2(\frac{1}{2} \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu})$.
 - See HW5 for a more general version.
- Independence between linear and quadratic forms of a normal $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{V})$
 - $(\mathbf{a} + \mathbf{A}\mathbf{X}) \perp (\mathbf{b} + \mathbf{B}\mathbf{X})$ if $\mathbf{A}\mathbf{V}\mathbf{B}^T = \mathbf{0}$.
 - $\mathbf{B}\mathbf{X} \perp \mathbf{X}^T \mathbf{A} \mathbf{X}$ if \mathbf{A} is symmetric and $\mathbf{B}\mathbf{V}\mathbf{A} = \mathbf{0}$.
 - $\mathbf{X}^T \mathbf{A} \mathbf{X} \perp \mathbf{X}^T \mathbf{B} \mathbf{X}$ if \mathbf{A} and \mathbf{B} are symmetric and $\mathbf{B}\mathbf{V}\mathbf{A} = \mathbf{0}$.

Today

- Cochran's theorem (JM 5.5)
- Statistical inference for normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I})$

Cochran theorem (JM 5.5)

- We already saw that, under the normal linear model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I})$, the sum of squares $\|\hat{\mathbf{y}}\|_2^2 = \|\mathbf{y}^T \mathbf{P}_X \mathbf{y}\|_2^2$ and $\|\hat{\mathbf{e}}\|_2^2 = \|\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}\|_2^2$ are independent chi-square distributions.

The Cochran theorem deals with the distributions of more general sum of squares from normal linear model.

- (Cochran theorem) Let $\mathbf{y} \sim N_n(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_n)$ and \mathbf{A}_i , $i = 1, \dots, k$, be symmetric idempotent matrix with rank s_i . If $\sum_{i=1}^k \mathbf{A}_i = \mathbf{I}_n$, then $(1/\sigma^2) \mathbf{y}^T \mathbf{A}_i \mathbf{y}$ are independent $\chi_{s_i}^2(\phi_i)$, with $\phi_i = \frac{1}{2\sigma^2} \boldsymbol{\mu}^T \mathbf{A}_i \boldsymbol{\mu}$ and $\sum_{i=1}^k s_i = n$.

Proof. Since \mathbf{A}_i is symmetric and idempotent with rank s_i , $\mathbf{A}_i = \mathbf{Q}_i \mathbf{Q}_i^T$ with $\mathbf{Q}_i \in \mathbb{R}^{n \times s_i}$ and $\mathbf{Q}_i^T \mathbf{Q}_i = \mathbf{I}_{s_i}$. Define $\mathbf{Q} = (\mathbf{Q}_1, \dots, \mathbf{Q}_k) \in \mathbb{R}^{n \times \sum_{i=1}^k s_i}$. Note

$$\begin{aligned} \mathbf{Q}^T \mathbf{Q} &= \mathbf{I}_{\sum_{i=1}^k s_i} \\ \mathbf{Q} \mathbf{Q}^T &= \sum_{i=1}^k \mathbf{Q}_i \mathbf{Q}_i^T = \sum_{i=1}^k \mathbf{A}_i = \mathbf{I}_n. \end{aligned}$$

From the second equation, we have

$$\sum_{i=1}^k s_i = \sum_{i=1}^k \text{rank}(\mathbf{A}_i) = \sum_{i=1}^k \text{tr}(\mathbf{A}_i) = \text{tr}(\mathbf{I}_n) = n.$$

Now

$$\mathbf{Q}^T \mathbf{y} = \begin{pmatrix} \mathbf{Q}_1^T \mathbf{y} \\ \vdots \\ \mathbf{Q}_k^T \mathbf{y} \end{pmatrix} \sim N_n \left(\begin{pmatrix} \mathbf{Q}_1^T \boldsymbol{\mu} \\ \vdots \\ \mathbf{Q}_k^T \boldsymbol{\mu} \end{pmatrix}, \sigma^2 \mathbf{I}_n \right),$$

implying that $\mathbf{Q}_i^T \mathbf{y} \sim N_{s_i}(\mathbf{Q}_i^T \boldsymbol{\mu}, \sigma^2 \mathbf{I}_{s_i})$ are jointly independent. Therefore $(1/\sigma^2) \mathbf{y}^T \mathbf{A}_i \mathbf{y} = (1/\sigma^2) \|\mathbf{Q}_i^T \mathbf{y}\|_2^2 \sim \chi_{s_i}^2(\frac{1}{2\sigma^2} \boldsymbol{\mu}^T \mathbf{A}_i \boldsymbol{\mu})$ are jointly independent. \square

- Application to the one-way ANOVA: $y_{ij} = \mu + \alpha_i + e_{ij}$, or $\mathbf{y} = N_n(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I}_n)$ where

$$\mathbf{X}\mathbf{b} = \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & & & \\ & \mathbf{1}_{n_2} & & & \\ & \vdots & & \ddots & \\ & & & & \mathbf{1}_{n_a} \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_a \end{pmatrix}.$$

Note

$$\begin{aligned}
\mathbf{P}_1 &= n^{-1} \mathbf{1} \mathbf{1}^T \\
\mathbf{P}_X &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\
&= \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & & & \\ \mathbf{1}_{n_2} & & \mathbf{1}_{n_2} & & \\ \vdots & & & \ddots & \\ \mathbf{1}_{n_a} & & & & \mathbf{1}_{n_a} \end{pmatrix} \begin{pmatrix} 0 & & & & \\ & n_1^{-1} & & & \\ & & n_2^{-1} & & \\ & & & \ddots & \\ & & & & n_a^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & & & \\ \mathbf{1}_{n_2} & & \mathbf{1}_{n_2} & & \\ \vdots & & & \ddots & \\ \mathbf{1}_{n_a} & & & & \mathbf{1}_{n_a} \end{pmatrix}^T \\
&= \begin{pmatrix} n_1^{-1} \mathbf{1}_{n_1} \mathbf{1}_{n_1}^T & & & & \\ & n_2^{-1} \mathbf{1}_{n_2} \mathbf{1}_{n_2}^T & & & \\ & & \ddots & & \\ & & & & n_a^{-1} \mathbf{1}_{n_a} \mathbf{1}_{n_a}^T \end{pmatrix}.
\end{aligned}$$

Define

$$\begin{aligned}
\mathbf{A}_1 &= \mathbf{P}_1 \\
\mathbf{A}_2 &= \mathbf{P}_X - \mathbf{P}_1 \\
\mathbf{A}_3 &= \mathbf{I}_n - \mathbf{P}_X.
\end{aligned}$$

with corresponding quadratic forms

$$\begin{aligned}
\text{SSM} &= \mathbf{y}^T \mathbf{A}_1 \mathbf{y} = n \bar{y}^2, \quad (1/\sigma^2) \mathbf{y}^T \mathbf{A}_1 \mathbf{y} \sim \chi_1^2(\phi_1), \\
\phi_1 &= \frac{1}{2\sigma^2} (\mathbf{X} \mathbf{b})^T \mathbf{A}_1 (\mathbf{X} \mathbf{b}) = \frac{n(\mu + \bar{\alpha})^2}{2\sigma^2} \\
\text{SSA}_{\text{cfm}} &= \mathbf{y}^T \mathbf{A}_2 \mathbf{y} = \sum_{i=1}^a n_i \bar{y}_i^2 - n \bar{y}^2, \quad (1/\sigma^2) \mathbf{y}^T \mathbf{A}_2 \mathbf{y} \sim \chi_{a-1}^2(\phi_2), \\
\phi_2 &= \frac{1}{2\sigma^2} (\mathbf{X} \mathbf{b})^T \mathbf{A}_2 (\mathbf{X} \mathbf{b}) = \frac{\sum_{i=1}^a n_i (\alpha_i - \bar{\alpha})^2}{2\sigma^2} \\
\text{SSE} &= \mathbf{y}^T \mathbf{A}_3 \mathbf{y} = \mathbf{y}^T \mathbf{y} - n \bar{y}^2 = \sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2, \quad (1/\sigma^2) \mathbf{y}^T \mathbf{A}_3 \mathbf{y} \sim \chi_{n-a}^2(\phi_3), \\
\phi_3 &= \frac{1}{2\sigma^2} (\mathbf{X} \mathbf{b})^T \mathbf{A}_3 (\mathbf{X} \mathbf{b}) = 0.
\end{aligned}$$

We have the classical ANOVA table

Source	df	Projection	SS	Noncentrality
Mean	1	\mathbf{P}_1	SSM= $n\bar{y}^2$	$\frac{1}{2\sigma^2}n(\mu + \bar{\alpha})^2$
Group	$a - 1$	$\mathbf{P}_X - \mathbf{P}_1$	SSA _{cfm} = $\sum_{i=1}^a n_i \bar{y}_i^2 - n\bar{y}^2$	$\frac{1}{2\sigma^2} \sum_{i=1}^a n_i (\alpha_i - \bar{\alpha})^2$
Error	$n - a$	$\mathbf{I} - \mathbf{P}_X$	SSE= $\sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$	0
Total	n	\mathbf{I}	SST= $\sum_i \sum_j y_{ij}^2$	$\frac{1}{2\sigma^2} \sum_{i=1}^a n_i (\mu + \alpha_i)^2$

Estimation under the normal Gauss-Markov model (JM 6.2)

Assume normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{I})$.

- The density of $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{I})$ is

$$\begin{aligned} f(\mathbf{y} \mid \mathbf{b}, \sigma^2) &= \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2}(\mathbf{y}-\mathbf{X}\mathbf{b})^T(\mathbf{y}-\mathbf{X}\mathbf{b})} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2}\mathbf{y}^T\mathbf{y} + \frac{1}{\sigma^2}\mathbf{y}^T\mathbf{X}\mathbf{b} - \frac{1}{2\sigma^2}\mathbf{b}^T\mathbf{X}^T\mathbf{X}\mathbf{b}}. \end{aligned}$$

- Recall Theorem 6.2.25 of Casella and Berger (2001).

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be iid observations from an exponential family with pdf or pmf of the form

$$f(\mathbf{x} \mid \boldsymbol{\theta}) = h(\mathbf{x})c(\boldsymbol{\theta}) \exp\left(\sum_{j=1}^k w(\theta_j)t_j(\mathbf{x})\right),$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$. Then the statistics

$$T(\mathbf{X}) = \left(\sum_{i=1}^n t_1(\mathbf{X}_i), \dots, \sum_{i=1}^n t_k(\mathbf{X}_i)\right)$$

is complete as long as the parameter space Θ contains an open set in \mathbb{R}^k .

- Also by Theorem 6.2.28 of Casella and Berger (2001):

If a minimal sufficient statistic exists, then any complete statistic is also a minimal sufficient statistic.

- From the density of the normal Gauss-Markov model, we see

$$T(\mathbf{y}) = (\mathbf{y}^T\mathbf{y}, \mathbf{X}^T\mathbf{y})$$

is a complete and minimal sufficient statistic for (σ^2, \mathbf{b}) .

- Recall the Rao-Blackwell theorem (Casella and Berger, 2001, Theorem 7.3.17). Any function of a sufficient statistic has the smallest variance among all unbiased estimators of its expectation.
- Under the normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{I})$, the least squares estimator $\mathbf{\Lambda}\hat{\mathbf{b}}$ of an estimable function $\mathbf{\Lambda}\mathbf{b}$ has the smallest variance among *all* unbiased estimators.

Remark: The least squares estimator $\mathbf{\Lambda}\hat{\mathbf{b}}$ is also called the MVUE (minimum variance unbiased estimator) under normal assumption.

Proof. The least squares estimator $\mathbf{\Lambda}\hat{\mathbf{b}} = \mathbf{\Lambda}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$ is a function of the complete sufficient statistic $T(\mathbf{y}) = (\mathbf{y}^T\mathbf{y}, \mathbf{X}^T\mathbf{y})$ and has expectation $E(\mathbf{\Lambda}\hat{\mathbf{b}}) = \mathbf{\Lambda}\mathbf{b}$. \square

- (MLE) Under the normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{I})$, the maximum likelihood estimator (MLE) of (\mathbf{b}, σ^2) is

$$\left(\hat{\mathbf{b}}, \frac{\text{SSE}}{n} \right),$$

where $\hat{\mathbf{b}}$ is any least squares solution (solution to the normal equation) and $\text{SSE} = \mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{y}$.

Remark: The MLE for σ^2 is biased.

Proof. To maximize the density

$$f(\mathbf{y} | \mathbf{b}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2}(\mathbf{y}-\mathbf{X}\mathbf{b})^T(\mathbf{y}-\mathbf{X}\mathbf{b})}$$

is equivalent to maximizing the log-likelihood

$$L(\mathbf{b}, \sigma^2 | \mathbf{y}) = \ln f(\mathbf{y} | \mathbf{b}, \sigma^2) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\mathbf{b})^T(\mathbf{y} - \mathbf{X}\mathbf{b}).$$

For any σ^2 , the quadratic form is minimized by any least squares solution (as we showed before) and optimal value SSE. To find the maximizing σ^2 , we set the derivative to 0

$$-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4}\text{SSE} = 0,$$

which yields $\hat{\sigma}^2 = \text{SSE}/n$. \square

- Under the normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{I})$, the MLE of an estimable function $\mathbf{A}\mathbf{b}$ is $\mathbf{A}\hat{\mathbf{b}}$, where $\hat{\mathbf{b}}$ is any least squares solution.

Proof. It follows from the invariance property of the MLE (Casella and Berger, 2001, Theorem 7.2.10). □

15 Lecture 15: Oct 28

Announcement

- HW5 due today.
- HW6 posted today and due next Monday.
- HW7 will be posted this Wed and due next Wed.
- HW5-7 are covered in Midterm 2.

Last time

- Cochran's theorem: Let $\mathbf{y} \sim N_n(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_n)$ and \mathbf{A}_i , $i = 1, \dots, k$, be symmetric idempotent matrix with rank s_i . If $\sum_{i=1}^k \mathbf{A}_i = \mathbf{I}_n$, then $(1/\sigma^2)\mathbf{y}^T \mathbf{A}_i \mathbf{y}$ are independent $\chi_{s_i}^2(\phi_i)$, with $\phi_i = \frac{1}{2\sigma^2} \boldsymbol{\mu}^T \mathbf{A}_i \boldsymbol{\mu}$ and $\sum_{i=1}^k s_i = n$.
- Application of Cochran's theorem to the one-way ANOVA model.
- Statistical inference for normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I})$.
 - $(\mathbf{y}^T \mathbf{y}, \mathbf{X}^T \mathbf{y})$ is a complete and minimal sufficient statistic for (σ^2, \mathbf{b}) .
 - The least squares estimator $\boldsymbol{\Lambda} \hat{\mathbf{b}}$ of an estimable function $\boldsymbol{\Lambda} \mathbf{b}$ is not only MVAUE, but also MVUE under normal Gauss-Markov assumption.
 - MLE of (\mathbf{b}, σ^2) is $(\hat{\mathbf{b}}, \text{SSE}/n)$, where $\hat{\mathbf{b}}$ is any least squares solution and $\text{SSE} = \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}$.
 - MLE of any estimable function $\boldsymbol{\Lambda} \mathbf{b}$ is $\boldsymbol{\Lambda} \hat{\mathbf{b}}$.

Today

- Testing general linear hypothesis: first principles test
- Testing general linear hypothesis: LRT

General linear hypothesis (JM 6.3)

We assume normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{I})$.

- In many applications, we are interested in testing a general linear hypothesis

$$H_0 : \mathbf{R}\mathbf{b} = \mathbf{r} \quad \text{vs} \quad H_A : \mathbf{R}\mathbf{b} \neq \mathbf{r},$$

where $\mathbf{R} \in \mathbb{R}^{s \times p}$ and $\mathbf{r} \in \mathbb{R}^s$.

- Examples

- $H_A : b_j = 0$
- $H_A : b_1 = b_2 = b_3 = 0$
- $H_A : b_2 = b_3 = b_4$
- $H_A : b_1 + b_3 = 1, b_2 = 3$
- $H_A : \mathbf{b} \in \mathcal{C}(\mathbf{B})$

- We say a general linear hypothesis $H_0 : \mathbf{R}\mathbf{b} = \mathbf{r}$ is *testable* if \mathbf{R} has full row rank s and $\mathbf{R}\mathbf{b}$ is estimable.
- Example (testing the interaction in two-way ANOVA): Consider the two-way ANOVA model with interaction

$$y_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij} + e_{ijk}.$$

In Lecture 6 (Sep 18), we showed that the interaction effects

$$\gamma_{ij} - \gamma_{ij'} - \gamma_{i'j} + \gamma_{i'j'}$$

are estimable. There are $(a-1)(b-1)$ linearly independent interaction effects. Thus $s = (a-1)(b-1)$ in the linear hypothesis $H_0 : \mathbf{R}\mathbf{b} = \mathbf{0}_s$ for testing interaction effects.

Recall that the design matrix for a two-way ANOVA with interaction has rank ab and that for a two-way ANOVA without interaction has rank $a+b-1$. The difference is $ab - (a+b-1) = (a-1)(b-1)$.

First principles test for a general linear hypothesis (JM 6.3)

- Consider testing a testable general linear hypothesis $\mathbf{Rb} = \mathbf{r}$ under the normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{Xb}, \sigma^2 \mathbf{I})$.
- Since \mathbf{Rb} is estimable, the MVAUE (also the MVUE and MLE)

$$\mathbf{R}\hat{\mathbf{b}} = \mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

is a multivariate normal with mean

$$\mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Xb} = \mathbf{Rb}$$

and covariance

$$\sigma^2 \mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T = \sigma^2 \mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T =: \sigma^2 \mathbf{H}.$$

Hence

$$\mathbf{R}\hat{\mathbf{b}} - \mathbf{r} \sim N_s(\mathbf{Rb} - \mathbf{r}, \sigma^2 \mathbf{H}).$$

- Note

$$\begin{aligned} \text{rank}(\mathbf{H}) &= \text{rank}(\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T) \\ &\geq \text{rank}(\mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T) = \text{rank}(\mathbf{R}^T) = s, \end{aligned}$$

thus \mathbf{H} is non-singular.

- By previous result (Lecture 12),

$$(\mathbf{R}\hat{\mathbf{b}} - \mathbf{r})^T (\sigma^2 \mathbf{H})^{-1} (\mathbf{R}\hat{\mathbf{b}} - \mathbf{r}) \sim \chi_s^2(\phi)$$

with noncentrality parameter

$$\phi = \frac{1}{2} (\mathbf{Rb} - \mathbf{r})^T (\sigma^2 \mathbf{H})^{-1} (\mathbf{Rb} - \mathbf{r}).$$

- Since

$$[\mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] (\sigma^2 \mathbf{I}) (\mathbf{I} - \mathbf{P}_X) = \mathbf{0},$$

the quadratic form $(\mathbf{R}\hat{\mathbf{b}} - \mathbf{r})^T (\sigma^2 \mathbf{H})^{-1} (\mathbf{R}\hat{\mathbf{b}} - \mathbf{r})$ is independent of the quadratic form $\text{SSE} = \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}$.

- Therefore the ratio is a F random variable

$$F = \frac{(\mathbf{R}\hat{\mathbf{b}} - \mathbf{r})^T \mathbf{H}^{-1}(\mathbf{R}\hat{\mathbf{b}} - \mathbf{r})/s}{\text{SSE}/(n-r)} \sim F_{s, n-r}(\phi).$$

Under $H_0 : \mathbf{R}\mathbf{b} = \mathbf{r}$, $\phi = 0$ and the ratio is a central F distribution.

Under $H_A : \mathbf{R}\mathbf{b} \neq \mathbf{r}$, $\phi > 0$ and the ratio is a noncentral F distribution.

- Consider one-way ANOVA model

$$y_{ij} = \mu + \alpha_i + e_{ij}.$$

We want test the hypothesis that all effects α_i are equal. The two hypotheses

$$H_0 : \alpha_1 - \alpha_2 = 0, \alpha_1 - \alpha_3 = 0, \dots, \alpha_a - \alpha_a = 0$$

and

$$H_0 : \alpha_1 - \alpha_2 = 0, \alpha_2 - \alpha_3 = 0, \dots, \alpha_{a-1} - \alpha_a = 0$$

are two logically equivalent ways of expressing the test. Will they lead to the same test procedure?

- Two linear hypotheses $\mathbf{R}_1\mathbf{b} = \mathbf{r}_1$ and $\mathbf{R}_2\mathbf{b} = \mathbf{r}_2$ are *equivalent* if

$$\{\mathbf{b} : \mathbf{R}_1\mathbf{b} = \mathbf{r}_1\} = \{\mathbf{b} : \mathbf{R}_2\mathbf{b} = \mathbf{r}_2\}.$$

- (Invariance of the first principles test under equivalent linear hypotheses) Suppose two linear hypotheses $\mathbf{R}_1\mathbf{b} = \mathbf{r}_1$ and $\mathbf{R}_2\mathbf{b} = \mathbf{r}_2$ are *equivalent* and both \mathbf{R}_1 and \mathbf{R}_2 have full row rank s . The points in the set $\{\mathbf{b} : \mathbf{R}_1\mathbf{b} = \mathbf{r}_1\}$ are characterized by

$$\mathbf{R}_1^- \mathbf{r}_1 + (\mathbf{I} - \mathbf{R}_1^- \mathbf{R}_1)\mathbf{q}.$$

Now

$$\mathbf{R}_2[\mathbf{R}_1^- \mathbf{r}_1 + (\mathbf{I} - \mathbf{R}_1^- \mathbf{R}_1)\mathbf{q}] = \mathbf{r}_2$$

for all \mathbf{q} implies that $\mathbf{R}_2\mathbf{R}_1^- \mathbf{r}_1 = \mathbf{r}_2$ (taking $\mathbf{q} = \mathbf{0}$) and

$$\mathbf{R}_2(\mathbf{I} - \mathbf{R}_1^- \mathbf{R}_1) = \mathbf{0}_{s \times p}.$$

Thus $\mathcal{C}(\mathbf{R}_2^T) \subset \mathcal{C}(\mathbf{R}_1^T)$. Reversing the roles of $(\mathbf{R}_1, \mathbf{b}_1)$ and $(\mathbf{R}_2, \mathbf{b}_2)$ shows that $\mathcal{C}(\mathbf{R}_1^T) \subset \mathcal{C}(\mathbf{R}_2^T)$. Hence $\mathcal{C}(\mathbf{R}_1^T) = \mathcal{C}(\mathbf{R}_2^T)$.

$\mathbf{R}_2 = \mathbf{T}\mathbf{R}_1$ for some transformation matrix $\mathbf{T} \in \mathbb{R}^{s \times s}$. \mathbf{T} must be non-singular, otherwise $\text{rank}(\mathbf{R}_2) \leq \text{rank}(\mathbf{T}) < s$, a contradiction with \mathbf{R}_2 has full row rank. Also we have

$$\mathbf{R}_2 \mathbf{R}_1^{-1} \mathbf{r}_1 = \mathbf{T} \mathbf{R}_1 \mathbf{R}_1^{-1} \mathbf{r}_1 = \mathbf{T} \mathbf{r}_1 = \mathbf{r}_2.$$

Thus the numerator of the F statistic for testing $H_0 : \mathbf{R}_2 \mathbf{b} = \mathbf{r}_2$

$$\begin{aligned} & (\mathbf{R}_2 \mathbf{b} - \mathbf{r}_2)^T [\mathbf{R}_2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}_2^T]^{-1} (\mathbf{R}_2 \mathbf{b} - \mathbf{r}_2) \\ &= (\mathbf{R}_1 \mathbf{b} - \mathbf{r}_1)^T \mathbf{T}^T [\mathbf{T} \mathbf{R}_1 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}_1^T \mathbf{T}^T]^{-1} \mathbf{T} (\mathbf{R}_1 \mathbf{b} - \mathbf{r}_1) \\ &= (\mathbf{R}_1 \mathbf{b} - \mathbf{r}_1)^T [\mathbf{R}_1 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}_1^T]^{-1} (\mathbf{R}_1 \mathbf{b} - \mathbf{r}_1). \end{aligned}$$

is same as that for testing $H_0 : \mathbf{R}_1 \mathbf{b} = \mathbf{r}_1$.

Likelihood ratio test (LRT) for a general linear hypothesis (JM 6.4-6.5)

- We derive the LRT for testing the general linear hypothesis $H_0 : \mathbf{R}\mathbf{b} = \mathbf{r}$.
- The log-likelihood function is

$$\begin{aligned} L(\mathbf{b}, \sigma^2 | \mathbf{y}) &= -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b}) \\ &= -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} Q(\mathbf{b}), \end{aligned}$$

where $Q(\mathbf{b}) = (\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b})$. As derived in the last lecture, for any \mathbf{b} , the maximizing σ^2 is given by

$$\hat{\sigma}^2 = \frac{Q(\mathbf{b})}{n}$$

and the resultant log-likelihood is

$$-\frac{n}{2} \ln[2\pi Q(\mathbf{b})/n] - \frac{n}{2}.$$

We only need to derive the constrained and unconstrained MLE for \mathbf{b} . That is the constrained and unconstrained maximizers of $Q(\mathbf{b})$.

- The unconstrained MLE is the least squares solution $\hat{\mathbf{b}}$, i.e., any solution to the normal equation.
- To derive the constrained MLE, we need to solve the minimization problem

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{b})^T(\mathbf{y} - \mathbf{X}\mathbf{b}) \\ & \text{subject to} \quad \mathbf{R}\mathbf{b} = \mathbf{r}. \end{aligned}$$

Setting the gradient of the Lagrangian function

$$\frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{b})^T(\mathbf{y} - \mathbf{X}\mathbf{b}) + \boldsymbol{\lambda}^T(\mathbf{R}\mathbf{b} - \mathbf{r})$$

to zero leads to the linear equation

$$\begin{pmatrix} \mathbf{X}^T\mathbf{X} & \mathbf{R}^T \\ \mathbf{R} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{b} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T\mathbf{y} \\ \mathbf{r} \end{pmatrix}.$$

This equation is always consistent since $\mathbf{X}^T\mathbf{y} \in \mathcal{C}(\mathbf{X}^T) = \mathcal{C}(\mathbf{X}^T\mathbf{X})$ and $\mathbf{r} = \mathbf{R}\mathbf{b} \in \mathcal{C}(\mathbf{R})$. By HW4 Q5(e), a generalized inverse of bordered Gramian matrix is

$$\begin{aligned} & \begin{pmatrix} \mathbf{X}^T\mathbf{X} & \mathbf{R}^T \\ \mathbf{R} & \mathbf{0} \end{pmatrix}^- \\ &= \begin{pmatrix} (\mathbf{X}^T\mathbf{X})^- - (\mathbf{X}^T\mathbf{X})^- \mathbf{R}^T \mathbf{H}^{-1} \mathbf{R} (\mathbf{X}^T\mathbf{X})^- & (\mathbf{X}^T\mathbf{X})^- \mathbf{R}^T \mathbf{H}^{-1} \\ \mathbf{H}^{-1} \mathbf{R} (\mathbf{X}^T\mathbf{X})^- & -\mathbf{H}^{-1} \end{pmatrix}, \end{aligned}$$

where $\mathbf{H} = \mathbf{R}(\mathbf{X}^T\mathbf{X})^- \mathbf{R}^T$ is non-singular under the assumption \mathbf{R} has full row rank. Therefore a solution is

$$\begin{aligned} \hat{\mathbf{b}}_0 &= [(\mathbf{X}^T\mathbf{X})^- - (\mathbf{X}^T\mathbf{X})^- \mathbf{R}^T \mathbf{H}^{-1} \mathbf{R} (\mathbf{X}^T\mathbf{X})^-] \mathbf{X}^T \mathbf{y} + (\mathbf{X}^T\mathbf{X})^- \mathbf{R}^T \mathbf{H}^{-1} \mathbf{r} \\ &= \hat{\mathbf{b}} - (\mathbf{X}^T\mathbf{X})^- \mathbf{R}^T \mathbf{H}^{-1} \mathbf{R} \hat{\mathbf{b}} + (\mathbf{X}^T\mathbf{X})^- \mathbf{R}^T \mathbf{H}^{-1} \mathbf{r} \\ &= \hat{\mathbf{b}} - (\mathbf{X}^T\mathbf{X})^- \mathbf{R}^T \mathbf{H}^{-1} (\mathbf{R}\hat{\mathbf{b}} - \mathbf{r}). \end{aligned}$$

- The change in the SSEs of constrained and unconstrained models is

$$\begin{aligned} & Q(\hat{\mathbf{b}}_0) - Q(\hat{\mathbf{b}}) \\ &= (\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}_0)^T(\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}_0) - (\mathbf{y} - \mathbf{X}\hat{\mathbf{b}})^T(\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}) \\ &= 2\mathbf{y}^T\mathbf{X}(\hat{\mathbf{b}} - \hat{\mathbf{b}}_0) - 2\hat{\mathbf{b}}\mathbf{X}^T\mathbf{X}\hat{\mathbf{b}} + 2\hat{\mathbf{b}}^T\mathbf{X}^T\mathbf{X}\hat{\mathbf{b}}_0 + (\hat{\mathbf{b}}_0 - \hat{\mathbf{b}})^T\mathbf{X}^T\mathbf{X}(\hat{\mathbf{b}}_0 - \hat{\mathbf{b}}) \\ &= 2(\mathbf{y} - \mathbf{X}\hat{\mathbf{b}})^T\mathbf{X}(\hat{\mathbf{b}} - \hat{\mathbf{b}}_0) + (\hat{\mathbf{b}}_0 - \hat{\mathbf{b}})^T\mathbf{X}^T\mathbf{X}(\hat{\mathbf{b}}_0 - \hat{\mathbf{b}}) \\ &= 0 + (\mathbf{R}\hat{\mathbf{b}} - \mathbf{r})^T \mathbf{H}^{-1} \mathbf{R} (\mathbf{X}^T\mathbf{X})^- \mathbf{X}^T \mathbf{X} (\mathbf{X}^T\mathbf{X})^- \mathbf{R}^T \mathbf{H}^{-1} (\mathbf{R}\hat{\mathbf{b}} - \mathbf{r}) \\ &= (\mathbf{R}\hat{\mathbf{b}} - \mathbf{r})^T \mathbf{H}^{-1} (\mathbf{R}\hat{\mathbf{b}} - \mathbf{r}). \end{aligned}$$

- The LRT rejects when

$$\frac{\max_{\Omega_0} L(\mathbf{b}, \sigma^2)}{\max_{\Omega} L(\mathbf{b}, \sigma^2)} = \left(\frac{Q(\hat{\mathbf{b}}_0)}{Q(\hat{\mathbf{b}})} \right)^{-n/2}$$

is small, or equivalently when

$$\frac{Q(\hat{\mathbf{b}}_0) - Q(\hat{\mathbf{b}})}{Q(\hat{\mathbf{b}})}$$

is large, or equivalently when

$$\frac{(\mathbf{R}\hat{\mathbf{b}} - \mathbf{r})^T \mathbf{H}^{-1} (\mathbf{R}\hat{\mathbf{b}} - \mathbf{r}) / s}{\text{SSE} / (n - r)}$$

is large.

Therefore the LRT is same as the first principles test!

16 Lecture 16: Oct 30

Announcement

- HW4 returned.
- HW4 Q5(e) mystery resolved: <http://hua-zhou.github.io/teaching/st552-2013fall/2013/10/29/hw4-Q5e-mystery-resolved.html>
- HW6 due next Monday (?)
- No HW session this afternoon.
- HW7 posted today and due next Wed.
- Midterm 2 covers Chapter 5-7 and HW 5-7.

Last time

- Testing general linear hypothesis $H_0 : \mathbf{R}\mathbf{b} = \mathbf{r}$: first principles test

$$F = \frac{(\mathbf{R}\hat{\mathbf{b}} - \mathbf{r})^T \mathbf{H}^{-1} (\mathbf{R}\hat{\mathbf{b}} - \mathbf{r}) / s}{\text{SSE} / (n - r)} \sim F_{s, n-r}(\phi).$$

- Testing general linear hypothesis: LRT = first principles test
- The first principles test (and LRT) is same under equivalent linear hypotheses.

Today

- t test and confidence interval
- simultaneous confidence intervals and multiple comparison

t test and confidence interval

Assume the normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I})$.

- If $\mathbf{r}^T \mathbf{b}$ is estimable, then the least squares estimator (MVAUE, MVUE, MLE) is

$$\mathbf{r}^T \hat{\mathbf{b}} \sim N(\mathbf{r}^T \mathbf{b}, \sigma^2 \mathbf{r}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}),$$

or equivalently

$$\frac{\mathbf{r}^T \hat{\mathbf{b}} - \mathbf{r}^T \mathbf{b}}{\sigma \sqrt{\mathbf{r}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}}} \sim N(0, 1).$$

- We estimate σ^2 by the usual $\hat{\sigma}^2 = \text{SSE}/(n - r)$. Then

$$\frac{\mathbf{r}^T \hat{\mathbf{b}} - \mathbf{r}^T \mathbf{b}}{\hat{\sigma} \sqrt{\mathbf{r}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}}} = \frac{(\mathbf{r}^T \hat{\mathbf{b}} - \mathbf{r}^T \mathbf{b}) / \sqrt{\sigma^2 \mathbf{r}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}}}{\text{SSE} / [\sigma^2 (n - r)]} \sim t_{n-r}.$$

- A $(1 - \alpha)$ confidence interval for $\mathbf{r}^T \mathbf{b}$ is given by

$$\mathbf{r}^T \hat{\mathbf{b}} \pm t_{n-r, \alpha/2} \hat{\sigma} \sqrt{\mathbf{r}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}},$$

where $t_{n-r, \alpha/2}$ is the critical value of a central t distribution with $n - r$ degrees of freedom. We are assured that

$$\mathbf{P}(\mathbf{r}^T \hat{\mathbf{b}} - t_{n-r, \alpha/2} \hat{\sigma} \sqrt{\mathbf{r}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}} \leq \mathbf{r}^T \mathbf{b} \leq \mathbf{r}^T \hat{\mathbf{b}} + t_{n-r, \alpha/2} \hat{\sigma} \sqrt{\mathbf{r}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}}) = 1 - \alpha.$$

- For two-sided testing

$$H_0 : \mathbf{r}^T \mathbf{b} = r \quad \text{vs} \quad H_A : \mathbf{r}^T \mathbf{b} \neq r.$$

The t test rejects H_0 when

$$\left| \frac{\mathbf{r}^T \hat{\mathbf{b}} - r}{\hat{\sigma} \sqrt{\mathbf{r}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}}} \right| > t_{n-r, \alpha/2}.$$

One-sided tests are carried out similarly using critical value $t_{n-r, \alpha}$.

Simultaneous confidence intervals/multiple comparison (JM 6.6)

Assume the normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I})$.

- Now we consider estimating an estimable function \mathbf{Rb} where \mathbf{R} has full row rank s . The least squares estimator (also MVAUE, MVUE and MLE) is

$$\mathbf{R}\hat{\mathbf{b}} \sim N_s(\mathbf{Rb}, \sigma^2 \mathbf{H}),$$

where

$$\mathbf{H} = \mathbf{R}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T.$$

- The $1 - \alpha$ confidence interval for the i -th linear function $\mathbf{r}_i^T \mathbf{b}$ is $[l_i, u_i]$, where

$$\begin{aligned} l_i &= \mathbf{r}_i^T \hat{\mathbf{b}} - t_{n-r, \alpha/2} \hat{\sigma} \sqrt{\mathbf{r}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}_i} = \mathbf{r}_i^T \hat{\mathbf{b}} - t_{n-r, \alpha/2} \hat{\sigma} \sqrt{h_{ii}} \\ u_i &= \mathbf{r}_i^T \hat{\mathbf{b}} + t_{n-r, \alpha/2} \hat{\sigma} \sqrt{\mathbf{r}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}_i} = \mathbf{r}_i^T \hat{\mathbf{b}} + t_{n-r, \alpha/2} \hat{\sigma} \sqrt{h_{ii}}. \end{aligned}$$

That means

$$\mathbf{P}(l_i \leq \mathbf{r}_i^T \mathbf{b} \leq u_i) = 1 - \alpha \quad \text{for } i = 1, \dots, s,$$

and thus

$$\mathbf{P}(l_i \leq \mathbf{r}_i^T \mathbf{b} \leq u_i \text{ for all } i = 1, \dots, s) \leq 1 - \alpha,$$

where the inequality can be strict. In other words, our joint confidence intervals don't have the right coverage ☹

Bonferroni correction

- (Bonferroni inequalities) Let E_i be a collection of events. Then

$$\begin{aligned} \mathbf{P}(\cup_j E_j) &= \mathbf{P}(\text{at least one } E_j) \leq \sum_j \mathbf{P}(E_j) \\ \mathbf{P}(\cap_j E_j) &= \mathbf{P}(\text{all } E_j) \geq 1 - \sum_j \mathbf{P}(E_j^c). \end{aligned}$$

- Bonferroni method for simultaneous confidence intervals. Set

$$\begin{aligned} l_i &= \mathbf{r}_i^T \hat{\mathbf{b}} - t_{n-r, \alpha/(2s)} \hat{\sigma} \sqrt{\mathbf{r}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}_i} = \mathbf{r}_i^T \hat{\mathbf{b}} - t_{n-r, \alpha/2} \hat{\sigma} \sqrt{h_{ii}} \\ u_i &= \mathbf{r}_i^T \hat{\mathbf{b}} + t_{n-r, \alpha/(2s)} \hat{\sigma} \sqrt{\mathbf{r}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}_i} = \mathbf{r}_i^T \hat{\mathbf{b}} + t_{n-r, \alpha/2} \hat{\sigma} \sqrt{h_{ii}}. \end{aligned}$$

Then

$$\begin{aligned}
& \mathbf{P}(l_i \leq \mathbf{r}_i^T \mathbf{b} \leq u_i \text{ for all } i = 1, \dots, s) \\
& \geq 1 - \sum_{i=1}^s \mathbf{P}(\mathbf{r}_i^T \mathbf{b} \notin [l_i, u_i]) \\
& = 1 - s(\alpha/s) \\
& = 1 - \alpha.
\end{aligned}$$

Now our joint confidence intervals have the right coverage \odot

- Obviously we can utilize these simultaneous intervals for testing $H_0 : \mathbf{Rb} = \mathbf{r}$. Simply reject H_0 when $|\mathbf{r}_i^T \hat{\mathbf{b}} - r_i| > t_{n-r, \alpha/2s} \hat{\sigma} \sqrt{h_{ii}}$ for any i .

Scheffé's method

- Scheffé's method constructs simultaneous confidence intervals

$$[\mathbf{u}^T \mathbf{R}\hat{\mathbf{b}} - c\hat{\sigma}\sqrt{\mathbf{u}^T \mathbf{H}\mathbf{u}}, \mathbf{u}^T \mathbf{R}\hat{\mathbf{b}} + c\hat{\sigma}\sqrt{\mathbf{u}^T \mathbf{H}\mathbf{u}}]$$

for *all* possible linear combinations $\mathbf{u}^T \mathbf{Rb}$ of the estimable function \mathbf{Rb} . The universal constant c needs to be chosen so that we have the right coverage.

- Observe that

$$\begin{aligned}
& \mathbf{P}(\mathbf{u}^T \mathbf{Rb} \in [\mathbf{u}^T \mathbf{R}\hat{\mathbf{b}} - c\hat{\sigma}\sqrt{\mathbf{u}^T \mathbf{H}\mathbf{u}}, \mathbf{u}^T \mathbf{R}\hat{\mathbf{b}} + c\hat{\sigma}\sqrt{\mathbf{u}^T \mathbf{H}\mathbf{u}}] \text{ for all } \mathbf{u}) \\
& = \mathbf{P}\left(\frac{|\mathbf{u}^T \mathbf{R}\hat{\mathbf{b}} - \mathbf{u}^T \mathbf{Rb}|}{\hat{\sigma}\sqrt{\mathbf{u}^T \mathbf{H}\mathbf{u}}} \leq c \text{ for all } \mathbf{u}\right) \\
& = \mathbf{P}\left(\max_{\mathbf{u}} \frac{|\mathbf{u}^T \mathbf{R}\hat{\mathbf{b}} - \mathbf{u}^T \mathbf{Rb}|}{\hat{\sigma}\sqrt{\mathbf{u}^T \mathbf{H}\mathbf{u}}} \leq c\right) \\
& = \mathbf{P}\left(\max_{\mathbf{u}} \frac{[\mathbf{u}^T (\mathbf{R}\hat{\mathbf{b}} - \mathbf{Rb})]^2}{\hat{\sigma}^2 \mathbf{u}^T \mathbf{H}\mathbf{u}} \leq c^2\right) \\
& = \mathbf{P}\left(\frac{(\mathbf{R}\hat{\mathbf{b}} - \mathbf{Rb})^T \mathbf{H}^{-1} (\mathbf{R}\hat{\mathbf{b}} - \mathbf{Rb})}{\hat{\sigma}^2} \leq c^2\right) \\
& = \mathbf{P}\left(\frac{(\mathbf{R}\hat{\mathbf{b}} - \mathbf{Rb})^T \mathbf{H}^{-1} (\mathbf{R}\hat{\mathbf{b}} - \mathbf{Rb})}{s\hat{\sigma}^2} \leq \frac{c^2}{s}\right).
\end{aligned}$$

The key step (the fourth equality) is due to a generalized Cauchy-Schwartz inequality (HW7). Since

$$\frac{(\mathbf{R}\hat{\mathbf{b}} - \mathbf{R}\mathbf{b})^T \mathbf{H}^{-1}(\mathbf{R}\hat{\mathbf{b}} - \mathbf{R}\mathbf{b})}{s\hat{\sigma}^2} \sim F_{s,n-r},$$

we choose c such that $c^2/s = F_{s,n-r,\alpha}$, i.e.,

$$c = \sqrt{sF_{s,n-r,\alpha}}.$$

- Scheffé's simultaneous confidence intervals are just inversion of the F test (first principles test, LRT) and are invariant under equivalent linear hypotheses.

Tukey's method

- Consider the *balanced* one-way ANOVA model: $y_{ij} = \mu + \alpha_i + e_{ij}$, where $n_1 = \dots = n_a = n$.
- We are interested constructing simultaneous confidence intervals for *all* pairwise differences of treatment effects $\alpha_i - \alpha_j$.
- The key observation is

$$\max_{i,j} (x_i - x_j) = \max_i x_i - \min_i x_i.$$

for any real numbers x_i .

- Under balanced one-way ANOVA

$$\bar{y}_i. \sim N(\mu + \alpha_i, \sigma^2/n)$$

for $i = 1, \dots, a$ and are independent. So

$$Z_i = \frac{\bar{y}_i. - (\mu + \alpha_i)}{\sigma/\sqrt{n}}$$

are independent standard normals. Then

$$\begin{aligned} & \max_{i,j} \frac{(\bar{y}_i. - \bar{y}_j.) - (\alpha_i - \alpha_j)}{\sigma/\sqrt{n}} \\ &= \max_i \frac{\bar{y}_i. - (\mu + \alpha_i)}{\sigma/\sqrt{n}} - \min_i \frac{\bar{y}_i. - (\mu + \alpha_i)}{\sigma/\sqrt{n}} \\ &= \max_i Z_i - \min_i Z_i. \end{aligned}$$

Also recall that

$$\frac{\text{SSE}}{\sigma^2} = \frac{a(n-1)\hat{\sigma}^2}{\sigma^2} \sim \chi_{a(n-1)}^2$$

and is independent of Z_i .

- We tabulate the distribution of the random variable

$$W = \frac{\max_j Z_j - \min_j Z_j}{\sqrt{U/\nu}},$$

where Z_j are k iid standard normals and independent of $U \sim \chi_\nu^2$. Let $w_{k,\nu,\alpha}$ be the critical values of this distribution.

Then

$$\begin{aligned} & 1 - \alpha \\ = & \mathbf{P} \left(\frac{\max_i (\bar{y}_i - (\mu + \alpha_i)) / (\sigma / \sqrt{n}) - \min_i (\bar{y}_i - (\mu + \alpha_i)) / (\sigma / \sqrt{n})}{\sqrt{\hat{\sigma}^2 / \sigma^2}} \leq w_{a,a(n-1),\alpha} \right) \\ = & \mathbf{P} \left(\max_i (\bar{y}_i - (\mu + \alpha_i)) - \min_i (\bar{y}_i - (\mu + \alpha_i)) \leq \frac{\hat{\sigma}}{\sqrt{n}} w_{a,a(n-1),\alpha} \right) \\ = & \mathbf{P} \left(\max_{i,j} [(\bar{y}_i - \bar{y}_j) - (\alpha_i - \alpha_j)] \leq \frac{\hat{\sigma}}{\sqrt{n}} w_{a,a(n-1),\alpha} \right) \\ = & \mathbf{P} \left(|(\bar{y}_i - \bar{y}_j) - (\alpha_i - \alpha_j)| \leq \frac{\hat{\sigma}}{\sqrt{n}} w_{a,a(n-1),\alpha} \text{ for all } i, j \right) \\ = & \mathbf{P} \left((\bar{y}_i - \bar{y}_j) - \frac{\hat{\sigma}}{\sqrt{n}} w_{a,a(n-1),\alpha} \leq \alpha_i - \alpha_j \leq (\bar{y}_i - \bar{y}_j) + \frac{\hat{\sigma}}{\sqrt{n}} w_{a,a(n-1),\alpha} \text{ for all } i, j \right). \end{aligned}$$

Therefore

$$\left[(\bar{y}_i - \bar{y}_j) - \frac{\hat{\sigma}}{\sqrt{n}} w_{a,a(n-1),\alpha}, (\bar{y}_i - \bar{y}_j) + \frac{\hat{\sigma}}{\sqrt{n}} w_{a,a(n-1),\alpha} \right]$$

are level $1 - \alpha$ simultaneous confidence intervals for all pairwise differences of treatment means $\alpha_i - \alpha_j$.

- (Extension to contrasts) Any linear combination $\mathbf{u}^T \boldsymbol{\alpha} = \sum_{i=1}^a u_i \alpha_i$, where $\sum_{i=1}^a u_i = 0$, is called a *contrast*. Tukey's method can be extended to simultaneous confidence intervals for all contrasts in balanced one-way ANOVA (HW7)

$$\left[\sum_i u_i \bar{y}_i - \frac{\hat{\sigma}}{\sqrt{n}} w_{a,a(n-1),\alpha} \times \frac{1}{2} \sum_i |u_i|, \sum_i u_i \bar{y}_i + \frac{\hat{\sigma}}{\sqrt{n}} w_{a,a(n-1),\alpha} \times \frac{1}{2} \sum_i |u_i| \right].$$

Sequential sum of squares (JM 7.3)

Assume the normal Gauss-Markov model $\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{I})$.

- Partition the design matrix as

$$\mathbf{X} = (\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_k),$$

where $\mathbf{X}_j \in \mathbb{R}^{n \times p_j}$ and $\sum_{j=0}^k p_j = p$.

- Define

$$\begin{aligned}\mathbf{A}_0 &= \mathbf{P}_{\mathbf{X}_0} \\ \mathbf{A}_j &= \mathbf{P}_{(\mathbf{X}_0, \dots, \mathbf{X}_j)} - \mathbf{P}_{(\mathbf{X}_0, \dots, \mathbf{X}_{j-1})}, \quad j = 1, \dots, k, \\ \mathbf{A}_{k+1} &= \mathbf{I}_n - \mathbf{P}_{\mathbf{X}}.\end{aligned}$$

Let

$$\begin{aligned}r_j &= \text{rank}(\mathbf{A}_j) = \text{rank}(\mathbf{P}_{(\mathbf{X}_0, \dots, \mathbf{X}_j)}) - \text{rank}(\mathbf{P}_{(\mathbf{X}_0, \dots, \mathbf{X}_{j-1})}) \\ &= \text{rank}((\mathbf{X}_0, \dots, \mathbf{X}_j)) - \text{rank}((\mathbf{X}_0, \dots, \mathbf{X}_{j-1}))\end{aligned}$$

be the extra rank contribution from predictors in \mathbf{X}_j .

- \mathbf{A}_j are orthogonal projections and

$$\sum_{j=0}^{k+1} \mathbf{A}_j = \mathbf{I}_n.$$

Therefore we can apply the Cochran's theorem,

$$\frac{\text{SS}_j}{\sigma^2} = \frac{1}{\sigma^2} \mathbf{y}^T \mathbf{A}_j \mathbf{y} \sim \chi_{r_j}^2(\phi_j), \quad j = 0, \dots, k+1,$$

with noncentrality parameter

$$\phi_j = \frac{1}{2\sigma^2} (\mathbf{X}\mathbf{b})^T \mathbf{A}_j (\mathbf{X}\mathbf{b})$$

and are independent.

- ANOVA table for sequential SS (type I SS in SAS). See Table 1.
- Example SAS output. Two-way ANOVA with interaction.

Source	DF	Projection	SS	Noncentrality
\mathbf{b}_0	$r(\mathbf{X}_0)$	$A_0 \mathbf{P}_{\mathbf{X}_0}$	$\mathbf{y}^T A_0 \mathbf{y}$	$(2\sigma^2)^{-1} (\mathbf{X}\mathbf{b})^T A_0 (\mathbf{X}\mathbf{b})$
\mathbf{b}_1 after \mathbf{b}_0	$r(\mathbf{X}_0, \mathbf{X}_1) - r(\mathbf{X}_0)$	$A_1 = \mathbf{P}_{(\mathbf{X}_0, \mathbf{X}_1)} - \mathbf{P}_{\mathbf{X}_0}$	$\mathbf{y}^T A_1 \mathbf{y}$	$(2\sigma^2)^{-1} (\mathbf{X}\mathbf{b})^T A_1 (\mathbf{X}\mathbf{b})$
...				
\mathbf{b}_j after $\mathbf{b}_0, \dots, \mathbf{b}_{j-1}$	$r(\mathbf{X}_0, \dots, \mathbf{X}_j) - r(\mathbf{X}_0, \dots, \mathbf{X}_{j-1})$	$A_j = \mathbf{P}_{(\mathbf{X}_0, \dots, \mathbf{X}_j)} - \mathbf{P}_{(\mathbf{X}_0, \dots, \mathbf{X}_{j-1})}$	$\mathbf{y}^T A_j \mathbf{y}$	$(2\sigma^2)^{-1} (\mathbf{X}\mathbf{b})^T A_j (\mathbf{X}\mathbf{b})$
...				
\mathbf{b}_k after $\mathbf{b}_0, \dots, \mathbf{b}_{k-1}$	$r(\mathbf{X}_0, \dots, \mathbf{X}_k) - r(\mathbf{X}_0, \dots, \mathbf{X}_{k-1})$	$A_k = \mathbf{P}_{(\mathbf{X}_0, \dots, \mathbf{X}_k)} - \mathbf{P}_{(\mathbf{X}_0, \dots, \mathbf{X}_{k-1})}$	$\mathbf{y}^T A_k \mathbf{y}$	$(2\sigma^2)^{-1} (\mathbf{X}\mathbf{b})^T A_k (\mathbf{X}\mathbf{b})$
Error	$n - r(\mathbf{X})$	$A_{k+1} = \mathbf{I}_n - \mathbf{P}_{\mathbf{X}}$	$\mathbf{y}^T (\mathbf{I}_n - \mathbf{P}_{\mathbf{X}}) \mathbf{y}$	0
Total	n	\mathbf{I}	$\mathbf{y}^T \mathbf{y}$	$(2\sigma^2)^{-1} (\mathbf{X}\mathbf{b})^T (\mathbf{X}\mathbf{b})$

Table 1: ANOVA table for sequential SS (type I SS in SAS).

```

proc glm data=bread;
  class height width;
  model sales=height width height*width;

```

The GLM Procedure
 Dependent Variable: sales

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	5	1580.000000	316.000000	30.58	0.0003
Error	6	62.000000	10.333333		
Corrected Total	11	1642.000000			

Source	DF	Type I SS	Mean Square	F Value	Pr > F
height	2	1544.000000	772.000000	74.71	<.0001
width	1	12.000000	12.000000	1.16	0.3226
height*width	2	24.000000	12.000000	1.16	0.3747

Source	DF	Type III SS	Mean Square	F Value	Pr > F
height	2	1544.000000	772.000000	74.71	<.0001
width	1	12.000000	12.000000	1.16	0.3226
height*width	2	24.000000	12.000000	1.16	0.3747

17 Lecture 17: Nov 4

Announcement

- HW6 due this Wed
- HW7 due this Wed.
- Midterm 2 change of time?

Last time

- t test and confidence interval

$$\mathbf{r}^T \hat{\mathbf{b}} \pm t_{n-r, \alpha/2} \hat{\sigma} \sqrt{\mathbf{r}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}}.$$

- Simultaneous confidence intervals and multiple comparison

- Bonferroni for s linear hypotheses $\mathbf{r}_i^T \mathbf{b}$

$$\mathbf{r}_i^T \hat{\mathbf{b}} \pm t_{n-r, \alpha/(2s)} \hat{\sigma} \sqrt{\mathbf{r}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{r}_i}.$$

- Scheffé for all linear combinations $\mathbf{u}^T \mathbf{R} \mathbf{b}$

$$\mathbf{u}^T \mathbf{R} \hat{\mathbf{b}} \pm c \hat{\sigma} \sqrt{\mathbf{u}^T \mathbf{R} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{R}^T \mathbf{u}}.$$

where $c = \sqrt{s F_{s, n-r, \alpha}}$.

- Tukey for all pairwise differences between treatment effects $\alpha_i - \alpha_j$ in balanced one-way ANOVA

$$(\bar{y}_i - \bar{y}_j) \pm \frac{\hat{\sigma}}{\sqrt{n}} w_{a, a(n-1), \alpha},$$

where $w_{a, a(n-1), \alpha}$ is critical value of $W = (\max_j Z_j - \min_j Z_j) / \sqrt{U / (an - a)}$,
 $U \sim \chi_{a(n-1)}^2$.

- Sequential SS (type I SS).

Today

- Testing under the Aitken model $\mathbf{y} \sim N(\mathbf{X} \mathbf{b}, \sigma^2 \mathbf{V})$.

Testing under the Aitken model

Assume the Aitken model $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{V})$, where \mathbf{V} can be singular.

- Consider testing a reduced model versus a full model

$$H_0 : \mathbf{y} \sim (\mathbf{X}_0\mathbf{c}, \sigma^2\mathbf{V}) \quad \text{vs} \quad H_A : \mathbf{y} \sim (\mathbf{X}\mathbf{b}, \sigma^2\mathbf{V}),$$

where $\mathcal{C}(\mathbf{X}_0) \subset \mathcal{C}(\mathbf{X})$. Equivalently, we are testing the column space hypothesis

$$H_0 : \mathbf{b} \in \mathcal{C}(\mathbf{B}) \quad \text{vs} \quad H_A : \mathbf{b} \in \mathbb{R}^p$$

for some \mathbf{B} . It's a special case of the more general linear hypothesis $\mathbf{R}\mathbf{b} = \mathbf{r}$.

- The MVAUE of the estimable function $\mathbf{X}\mathbf{b}$ is

$$\widehat{\Lambda}\mathbf{b} = \mathbf{X}(\mathbf{X}^T\mathbf{V}_1^-\mathbf{X})^-\mathbf{X}^T\mathbf{V}_1^-\mathbf{y},$$

where $\mathbf{V}_1 = \mathbf{V} + \mathbf{X}\mathbf{X}^T$. We also showed that

$$\widehat{\mathbf{X}}\mathbf{b} = \mathbf{X}\hat{\mathbf{b}},$$

where $\hat{\mathbf{b}} = (\mathbf{X}^T\mathbf{V}_1^-\mathbf{X})^-\mathbf{X}^T\mathbf{V}_1^-\mathbf{y}$ minimizes the generalized least squares criterion $(\mathbf{y} - \mathbf{X}\mathbf{b})^T\mathbf{V}_1^-(\mathbf{y} - \mathbf{X}\mathbf{b})$. Therefore the SSE under the full model is

$$\begin{aligned} \text{SSE} &= (\mathbf{y} - \mathbf{X}\hat{\mathbf{b}})^T\mathbf{V}_1^-(\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}) \\ &= \mathbf{y}^T[\mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{V}_1^-\mathbf{X})^-\mathbf{X}^T\mathbf{V}_1^-]^T\mathbf{V}_1^-[\mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{V}_1^-\mathbf{X})^-\mathbf{X}^T\mathbf{V}_1^-]\mathbf{y} \\ &= \mathbf{y}^T[\mathbf{V}_1^- - \mathbf{V}_1^-\mathbf{X}(\mathbf{X}^T\mathbf{V}_1^-\mathbf{X})^-\mathbf{X}^T\mathbf{V}_1^-]\mathbf{y} \\ &=: \mathbf{y}^T\mathbf{A}_1\mathbf{y}. \end{aligned}$$

The third equality uses HW4 Q5 (c) $\mathcal{C}(\mathbf{X}^T\mathbf{V}_1^-\mathbf{X}) = \mathcal{C}(\mathbf{X}^T)$.

We observe that SSE is invariant to the choice of the generalized inverses \mathbf{V}_1^- and $(\mathbf{X}^T\mathbf{V}_1^-\mathbf{X})^-$ (why?). Without loss of generality, we may assume \mathbf{A}_1 is symmetric by using the Moore-Penrose inverses throughout.

- Similarly, the SSE under the reduced model is

$$\begin{aligned} \text{SSE}_0 &= \mathbf{y}^T[\mathbf{I} - \mathbf{X}_0(\mathbf{X}_0^T\mathbf{V}_0^-\mathbf{X}_0)^-\mathbf{X}_0^T\mathbf{V}_0^-]^T\mathbf{V}_0^-[\mathbf{I} - \mathbf{X}_0(\mathbf{X}_0^T\mathbf{V}_0^-\mathbf{X}_0)^-\mathbf{X}_0^T\mathbf{V}_0^-]\mathbf{y} \\ &= \mathbf{y}^T[\mathbf{V}_0^- - \mathbf{V}_0^-\mathbf{X}_0(\mathbf{X}_0^T\mathbf{V}_0^-\mathbf{X}_0)^-\mathbf{X}_0^T\mathbf{V}_0^-]\mathbf{y} \\ &=: \mathbf{y}^T\mathbf{A}_0\mathbf{y}, \end{aligned}$$

where $\mathbf{V}_0 = \mathbf{V} + \mathbf{X}_0\mathbf{X}_0^T$.

- A reasonable test for a general linear hypothesis may be based on the change in SSEs under the constrained and full models

$$\frac{\text{SSE}_0 - \text{SSE}}{\text{SSE}}.$$

We need some preparatory results first.

- Assume $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{V})$, where $\boldsymbol{\mu} \in \mathcal{C}(\mathbf{V})$ and \mathbf{V} is an orthogonal projection matrix. Then

$$\mathbf{X}^T \mathbf{X} \sim \chi_{\text{rank}(\mathbf{V})}^2(\boldsymbol{\mu}^T \boldsymbol{\mu}/2).$$

Proof. $\mathbf{V} = \mathbf{O}\mathbf{O}^T$, where $\mathbf{O} \in \mathbb{R}^{n \times \text{rank}(\mathbf{V})}$ and $\mathbf{O}^T \mathbf{O} = \mathbf{I}_{\text{rank}(\mathbf{V})}$. Note $\mathcal{C}(\mathbf{V}) = \mathcal{C}(\mathbf{O})$. Let $\boldsymbol{\mu} = \mathbf{O}\mathbf{b}$ for some \mathbf{b} . Then $\mathbf{X} = \mathbf{O}\mathbf{Z}$, where $\mathbf{Z} \sim N(\mathbf{b}, \mathbf{I}_{\text{rank}(\mathbf{V})})$. Thus

$$\mathbf{X}^T \mathbf{X} = \mathbf{Z}^T \mathbf{O}^T \mathbf{O} \mathbf{Z} = \mathbf{Z}^T \mathbf{Z} \sim \chi_{\text{rank}(\mathbf{V})}^2(\boldsymbol{\mu}^T \boldsymbol{\mu}/2).$$

□

- Assume $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{V})$, \mathbf{V} possibly singular and \mathbf{A} is symmetric. Then

$$\mathbf{X}^T \mathbf{A} \mathbf{X} \sim \chi_{\text{tr}(\mathbf{A}\mathbf{V})}^2(\boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}/2)$$

if (1) $\mathbf{V}\mathbf{A}\mathbf{V}\mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{A}\mathbf{V}$, (2) $\boldsymbol{\mu}^T \mathbf{A}\mathbf{V}\mathbf{A}\boldsymbol{\mu} = \boldsymbol{\mu}^T \mathbf{A}\boldsymbol{\mu}$, and (3) $\mathbf{V}\mathbf{A}\mathbf{V}\mathbf{A}\boldsymbol{\mu} = \mathbf{V}\mathbf{A}\boldsymbol{\mu}$.

Remark: Previously (in Lecture 13) we showed the special case where \mathbf{V} is nonsingular, \mathbf{A} is symmetric, and $\mathbf{A}\mathbf{V}$ is idempotent.

Proof. $\mathbf{X} = \boldsymbol{\mu} + \mathbf{e}$, where $\mathbf{e} \sim N(\mathbf{0}, \mathbf{V})$ and thus $\mathbf{e} \in \mathcal{C}(\mathbf{V})$. Then

$$\begin{aligned} & \mathbf{X}^T \mathbf{A} \mathbf{V} \mathbf{A} \mathbf{X} \\ &= (\boldsymbol{\mu} + \mathbf{e})^T (\mathbf{A}\mathbf{V}\mathbf{A}) (\boldsymbol{\mu} + \mathbf{e}) \\ &= \boldsymbol{\mu}^T \mathbf{A}\mathbf{V}\mathbf{A}\boldsymbol{\mu} + 2\mathbf{e}^T \mathbf{A}\mathbf{V}\mathbf{A}\boldsymbol{\mu} + \mathbf{e}^T \mathbf{A}\mathbf{V}\mathbf{A}\mathbf{e} \\ &= \boldsymbol{\mu}^T \mathbf{A}\boldsymbol{\mu} + 2\mathbf{e}^T \mathbf{A}\boldsymbol{\mu} + \mathbf{e}^T \mathbf{A}\mathbf{e} \\ &= \mathbf{X}^T \mathbf{A} \mathbf{X}. \end{aligned}$$

Let $\mathbf{V} = \mathbf{Q}\mathbf{Q}^T$ where \mathbf{Q} can be the Cholesky factor or symmetric square root. Note $\mathcal{C}(\mathbf{V}) = \mathcal{C}(\mathbf{Q})$. Then

$$\mathbf{X}^T \mathbf{A} \mathbf{X} = \mathbf{X}^T \mathbf{A} \mathbf{V} \mathbf{A} \mathbf{X} = (\mathbf{X}^T \mathbf{A} \mathbf{Q})(\mathbf{Q}^T \mathbf{A} \mathbf{X})$$

and

$$\mathbf{Q}^T \mathbf{A} \mathbf{X} \sim N(\mathbf{Q}^T \mathbf{A} \boldsymbol{\mu}, \mathbf{Q}^T \mathbf{A} \mathbf{V} \mathbf{A} \mathbf{Q}).$$

By assumption (1), the covariance matrix is an orthogonal projection with rank

$$\begin{aligned} \text{rank}(\mathbf{Q}^T \mathbf{A} \mathbf{V} \mathbf{A} \mathbf{Q}) &= \text{rank}(\mathbf{Q}^T \mathbf{A} \mathbf{Q}) \\ &= \text{tr}(\mathbf{Q}^T \mathbf{A} \mathbf{Q}) = \text{tr}(\mathbf{A} \mathbf{Q} \mathbf{Q}^T) = \text{tr}(\mathbf{A} \mathbf{V}). \end{aligned}$$

We also need to check $\mathbf{Q}^T \mathbf{A} \boldsymbol{\mu} \in \mathcal{C}(\mathbf{Q}^T \mathbf{A} \mathbf{V} \mathbf{A} \mathbf{Q})$. This is true because

$$\begin{aligned} &\mathbf{Q}^T \mathbf{A} \mathbf{V} \mathbf{A} \mathbf{Q} (\mathbf{Q}^T \mathbf{A} \boldsymbol{\mu}) \\ &= \mathbf{Q}^T \mathbf{A} \mathbf{V} \mathbf{A} \mathbf{V} \mathbf{A} \boldsymbol{\mu} \\ &= \mathbf{Q}^T \mathbf{A} \mathbf{V} \mathbf{A} \boldsymbol{\mu} \\ &= \mathbf{Q}^T \mathbf{A} \boldsymbol{\mu}. \end{aligned}$$

Thus the result follows from the preceding result. \square

- If $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{V})$, \mathbf{V} possibly singular, and (1) $\mathbf{V} \mathbf{A} \mathbf{V} \mathbf{B} \mathbf{V} = \mathbf{0}$, (2) $\mathbf{V} \mathbf{A} \mathbf{V} \mathbf{B} \boldsymbol{\mu} = \mathbf{0}$, (3) $\mathbf{V} \mathbf{B} \mathbf{V} \mathbf{A} \boldsymbol{\mu} = \mathbf{0}$, (4) $\boldsymbol{\mu}^T \mathbf{A} \mathbf{V} \mathbf{B} \boldsymbol{\mu} = \mathbf{0}$, then $\mathbf{X}^T \mathbf{A} \mathbf{X}$ and $\mathbf{X}^T \mathbf{B} \mathbf{X}$ are independent.

Remark: Previously (in Lecture 13) we showed the special case where \mathbf{V} is nonsingular and \mathbf{A} and \mathbf{B} are symmetric.

Proof. Let $\mathbf{V} = \mathbf{Q} \boldsymbol{\Lambda}_s \mathbf{Q}^T$, where $\boldsymbol{\Lambda}_s$ is diagonal containing all nonzero eigenvalues of \mathbf{V} . Then $\mathbf{X} = \boldsymbol{\mu} + \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z}$, where $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$, and

$$\begin{aligned} \mathbf{X}^T \mathbf{A} \mathbf{X} &= (\boldsymbol{\mu} + \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z})^T \mathbf{A} (\boldsymbol{\mu} + \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z}) \\ &= \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} + 2\boldsymbol{\mu}^T \mathbf{A} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z} + \mathbf{Z}^T \boldsymbol{\Lambda}_s^{1/2} \mathbf{Q}^T \mathbf{A} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z} \\ \mathbf{X}^T \mathbf{B} \mathbf{X} &= (\boldsymbol{\mu} + \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z})^T \mathbf{B} (\boldsymbol{\mu} + \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z}) \\ &= \boldsymbol{\mu}^T \mathbf{B} \boldsymbol{\mu} + 2\boldsymbol{\mu}^T \mathbf{B} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z} + \mathbf{Z}^T \boldsymbol{\Lambda}_s^{1/2} \mathbf{Q}^T \mathbf{B} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z}. \end{aligned}$$

Now we check independence between the summands of the two expressions.

– From assumption (4),

$$\begin{aligned} & \boldsymbol{\mu}^T \mathbf{A} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Q}^T \mathbf{B} \boldsymbol{\mu} = \boldsymbol{\mu}^T \mathbf{A} \mathbf{V} \mathbf{B} \boldsymbol{\mu} = 0 \\ \Rightarrow & \boldsymbol{\mu}^T \mathbf{A} \mathbf{Q} \mathbf{Z} \perp \boldsymbol{\mu}^T \mathbf{B} \mathbf{Q} \mathbf{Z}. \end{aligned}$$

– From assumption (3),

$$\begin{aligned} & \mathbf{V} \mathbf{B} \mathbf{V} \mathbf{A} \boldsymbol{\mu} = \mathbf{0} \\ \Rightarrow & \mathbf{Q} \boldsymbol{\Lambda}_s \mathbf{Q}^T \mathbf{B} \mathbf{V} \mathbf{A} \boldsymbol{\mu} = \mathbf{0} \\ \Rightarrow & \mathbf{Q}^T \mathbf{B} \mathbf{V} \mathbf{A} \boldsymbol{\mu} = \mathbf{0} \\ \Rightarrow & \boldsymbol{\mu}^T \mathbf{A} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Q}^T \mathbf{B} \mathbf{Q} = \mathbf{0}^T \\ \Rightarrow & \boldsymbol{\mu}^T \mathbf{A} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z} \perp \mathbf{Z}^T \boldsymbol{\Lambda}_s^{1/2} \mathbf{Q}^T \mathbf{B} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z}. \end{aligned}$$

– Similarly, by assumption (2), $\boldsymbol{\mu}^T \mathbf{B} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z} \perp \mathbf{Z}^T \boldsymbol{\Lambda}_s^{1/2} \mathbf{Q}^T \mathbf{A} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z}$.

– From assumption (1),

$$\begin{aligned} & \mathbf{V} \mathbf{A} \mathbf{V} \mathbf{B} = \mathbf{0} \\ = & \mathbf{Q} \boldsymbol{\Lambda}_s \mathbf{Q}^T \mathbf{A} \mathbf{V} \mathbf{B} = \mathbf{0} \\ \Rightarrow & \mathbf{Q}^T \mathbf{A} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Q}^T \mathbf{B} \mathbf{Q} = \mathbf{0} \\ \Rightarrow & \mathbf{Z}^T \boldsymbol{\Lambda}_s^{1/2} \mathbf{Q}^T \mathbf{A} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z} \perp \mathbf{Z}^T \boldsymbol{\Lambda}_s^{1/2} \mathbf{Q}^T \mathbf{B} \mathbf{Q} \boldsymbol{\Lambda}_s^{1/2} \mathbf{Z}. \end{aligned}$$

Therefore, we have $\mathbf{X}^T \mathbf{A} \mathbf{X} \perp \mathbf{X}^T \mathbf{B} \mathbf{X}$. □

18 Lecture 18: Nov 6

Announcement

- HW6 due this Wed.
- HW7 due next Mon.
- Midterm 2 changed to next Wed Nov 13 @ 11:45AM-1PM.
- Q&A for basic exam questions? (if interested, submit your questions by 11/27).

Last time

Testing under Aitken model $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{V})$ – preparatory results. Note correction to the proof about independence of quadratic forms result (thanks to Xue Feng).

Today

Testing under Aitken model $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{V})$.

Testing under the Aitken model (cont'd)

Test $H_0 : \mathbf{y} \sim (\mathbf{X}_0\mathbf{c}, \sigma^2\mathbf{V})$ vs $H_A : \mathbf{y} \sim (\mathbf{X}\mathbf{b}, \sigma^2\mathbf{V})$, \mathbf{V} possibly singular.

- Now we can state and prove the results related to testing under Aitken model.
Recall the notation

$$\begin{aligned}\text{SSE} &= \mathbf{y}^T \mathbf{A}_1 \mathbf{y} \\ \text{SSE}_0 &= \mathbf{y}^T \mathbf{A}_0 \mathbf{y},\end{aligned}$$

where

$$\begin{aligned}\mathbf{A}_1 &= \mathbf{V}_1^- - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^- \mathbf{X}^T \mathbf{V}_1^-, & \mathbf{V}_1 &= \mathbf{V} + \mathbf{X} \mathbf{X}^T \\ \mathbf{A}_0 &= \mathbf{V}_0^- - \mathbf{V}_0^- \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{V}_0^- \mathbf{X}_0)^- \mathbf{X}_0^T \mathbf{V}_0^-, & \mathbf{V}_0 &= \mathbf{V} + \mathbf{X}_0 \mathbf{X}_0^T.\end{aligned}$$

- **Theorem**

1. $\frac{1}{\sigma^2} \mathbf{y}^T \mathbf{A}_1 \mathbf{y} \sim \chi_{\text{tr}(\mathbf{A}_1 \mathbf{V})}^2$.

2. If $\mathbf{X}\mathbf{b} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})$, then

$$\frac{1}{\sigma^2} \mathbf{y}^T (\mathbf{A}_0 - \mathbf{A}_1) \mathbf{y} \sim \chi_{\text{tr}((\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V})}^2 (\mathbf{b}^T \mathbf{X}^T \mathbf{A}_0 \mathbf{X} \mathbf{b} / 2).$$

3. If $\mathbf{X}\mathbf{b} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})$, then $\mathbf{y}^T \mathbf{A}_1 \mathbf{y}$ is independent of $\mathbf{y}^T (\mathbf{A}_0 - \mathbf{A}_1) \mathbf{y}$.

• The proof involves checking various conditions in the preparatory results.

$$\begin{aligned} - \mathbf{V} \mathbf{A}_1 &= \mathbf{V}_1 \mathbf{A}_1 \\ \mathbf{V} \mathbf{A}_0 &= \mathbf{V}_0 \mathbf{A}_0. \end{aligned}$$

Proof.

$$\begin{aligned} \mathbf{V} \mathbf{A}_1 &= (\mathbf{V}_1 - \mathbf{X} \mathbf{X}^T) \mathbf{A}_1 \\ &= \mathbf{V}_1 \mathbf{A}_1 - \mathbf{X} \mathbf{X}^T [\mathbf{V}_1^- - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_1^-] \\ &= \mathbf{V}_1 \mathbf{A}_1 - \mathbf{X} \mathbf{X}^T \mathbf{V}_1^- + \mathbf{X} \mathbf{X}^T \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_1^- \\ &= \mathbf{V}_1 \mathbf{A}_1 - \mathbf{X} \mathbf{X}^T \mathbf{V}_1^- + \mathbf{X} \mathbf{X}^T \mathbf{V}_1^- \\ &= \mathbf{V}_1 \mathbf{A}_1. \end{aligned}$$

Similarly $\mathbf{V} \mathbf{A}_0 = \mathbf{V}_0 \mathbf{A}_0$. □

$$\begin{aligned} - (1) \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_1 \mathbf{V} &= \mathbf{V} \mathbf{A}_1 \mathbf{V} \\ (2) \mathbf{b}^T \mathbf{X}^T \mathbf{A}_1 \mathbf{V} \mathbf{A}_1 \mathbf{X} \mathbf{b} &= \mathbf{b}^T \mathbf{X}^T \mathbf{A}_1 \mathbf{X} \mathbf{b} \\ (3) \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_1 \mathbf{X} \mathbf{b} &= \mathbf{V} \mathbf{A}_1 \mathbf{X} \mathbf{b}. \end{aligned}$$

Remark: The first part of the theorem is checked by this result. Note that the centrality parameter is $\mathbf{b}^T \mathbf{X}^T \mathbf{A}_1 \mathbf{X} \mathbf{b} = 0$ since

$$\mathbf{A}_1 \mathbf{X} = [\mathbf{V}_1^- - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_1^-] \mathbf{X} = \mathbf{0}.$$

Proof. For (1), by previous fact,

$$\begin{aligned} \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_1 \mathbf{V} &= \mathbf{V} \mathbf{A}_1 \mathbf{V}_1 \mathbf{A}_1 \mathbf{V} \\ &= \mathbf{V} [\mathbf{I} - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{V}_1^- \mathbf{V}_1 \mathbf{V}_1^- [\mathbf{I} - \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_1^-] \mathbf{V} \\ &= \mathbf{V} [\mathbf{I} - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{V}_1^- [\mathbf{V}_1 \mathbf{V}_1^- - \mathbf{V}_1 \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_1^-] \mathbf{V} \\ &= \mathbf{V} [\mathbf{I} - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{V}_1^- [\mathbf{V}_1 \mathbf{V}_1^- - \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_1^-] \mathbf{V} \\ &= \mathbf{V} [\mathbf{I} - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{V}_1^- [\mathbf{V}_1 \mathbf{V}_1^- \mathbf{V} - \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_1^- \mathbf{V}] \\ &= \mathbf{V} [\mathbf{I} - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{V}_1^- [\mathbf{V} - \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_1^- \mathbf{V}] \\ &= \mathbf{V} [\mathbf{I} - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{V}_1^- [\mathbf{I} - \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_1^-] \mathbf{V} \\ &= \mathbf{V} \mathbf{A}_1 \mathbf{V}. \end{aligned}$$

For (2),

$$\begin{aligned}
& \mathbf{b}^T \mathbf{X}^T \mathbf{A}_1 \mathbf{V} \mathbf{A}_1 \mathbf{X} \mathbf{b} \\
&= \mathbf{b}^T \mathbf{X}^T \mathbf{A}_1 \mathbf{V}_1 \mathbf{A}_1 \mathbf{X} \mathbf{b} \\
&= \mathbf{b}^T \mathbf{X}^T \mathbf{A}_1 \mathbf{X} \mathbf{b} \quad (\text{since } \mathcal{C}(\mathbf{X}) \subset \mathcal{C}(\mathbf{V}_1)).
\end{aligned}$$

For (3),

$$\begin{aligned}
& \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_1 \mathbf{X} \mathbf{b} \\
&= \mathbf{V}_1 \mathbf{A}_1 \mathbf{V}_1 \mathbf{A}_1 \mathbf{X} \mathbf{b} \\
&= \mathbf{V}_1 \mathbf{A}_1 \mathbf{X} \mathbf{b} \quad (\text{since } \mathcal{C}(\mathbf{X}) \subset \mathcal{C}(\mathbf{V}_1)) \\
&= \mathbf{V} \mathbf{A}_1 \mathbf{X} \mathbf{b}.
\end{aligned}$$

□

$$- \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_0 \mathbf{V} = \mathbf{V} \mathbf{A}_1 \mathbf{V}.$$

Proof.

$$\begin{aligned}
& \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_0 \mathbf{V} \\
&= \mathbf{V} \mathbf{A}_1 \mathbf{V}_0 [\mathbf{V}_0^- - \mathbf{V}_0^- \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{V}_0^- \mathbf{X}_0)^- \mathbf{X}_0^T \mathbf{V}_0^-] \mathbf{V} \\
&= \mathbf{V} \mathbf{A}_1 [\mathbf{V} - \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{V}_0^- \mathbf{X}_0)^- \mathbf{X}_0^T \mathbf{V}_0^- \mathbf{V}] \\
&= \mathbf{V} \mathbf{A}_1 \mathbf{V} \\
&= \mathbf{V} \mathbf{A} \mathbf{V}.
\end{aligned}$$

The third equality is because

$$\begin{aligned}
& \mathbf{A}_1 \mathbf{X}_0 \\
&= [\mathbf{V}_1^- - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^- \mathbf{X}^T \mathbf{V}_1^-] \mathbf{X}_0 \\
&= \mathbf{V}_1^- \mathbf{X}_0 - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^- \mathbf{X}^T \mathbf{V}_1^- \mathbf{X}_0 \\
&= \mathbf{V}_1^- \mathbf{X}_0 - \mathbf{V}_1^- \mathbf{X} (\mathbf{X}^T \mathbf{V}_1^- \mathbf{X})^- \mathbf{X}^T \mathbf{V}_1^- \mathbf{X} \mathbf{T} \quad (\mathcal{C}(\mathbf{X}_0) \subset \mathcal{C}(\mathbf{X})) \\
&= \mathbf{V}_1^- \mathbf{X}_0 - \mathbf{V}_1^- \mathbf{X}_0 \\
&= \mathbf{0}.
\end{aligned}$$

□

- (1) $\mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V} = \mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V}$.
- (2) $\mathbf{b}^T \mathbf{X}^T (\mathbf{A}_0 - \mathbf{A}_1) \mathbf{V} (\mathbf{A}_0 - \mathbf{A}_1) \mathbf{X} \mathbf{b} = \mathbf{b}^T \mathbf{X}^T (\mathbf{A}_0 - \mathbf{A}_1) \mathbf{X} \mathbf{b} = \mathbf{b}^T \mathbf{X}^T \mathbf{A}_0 \mathbf{X} \mathbf{b}$.
- (3) If $\mathbf{X} \mathbf{b} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})$, then $\mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1) \mathbf{V} (\mathbf{A}_0 - \mathbf{A}_1) \mathbf{X} \mathbf{b} = \mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1) \mathbf{X} \mathbf{b}$.

Remark: The second part of the theorem is checked by this result. The non-centrality parameter is

$$\frac{1}{2} \mathbf{b}^T \mathbf{X}^T (\mathbf{A}_0 - \mathbf{A}_1) \mathbf{X} \mathbf{b} = \frac{1}{2} \mathbf{b}^T \mathbf{X}^T \mathbf{A}_0 \mathbf{X} \mathbf{b}.$$

and the degrees of freedom is $\text{tr}((\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V})$.

Proof. For (1),

$$\begin{aligned} & \mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V} \\ &= \mathbf{V} \mathbf{A}_0 \mathbf{V} \mathbf{A}_0 \mathbf{V} - 2\mathbf{V} \mathbf{A}_0 \mathbf{V} \mathbf{A}_1 \mathbf{V} + \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_1 \mathbf{V} \\ &= \mathbf{V} \mathbf{A}_0 \mathbf{V} - 2\mathbf{V} \mathbf{A}_1 \mathbf{V} + \mathbf{V} \mathbf{A}_1 \mathbf{V} \\ &= \mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V}. \end{aligned}$$

For (2), recall that $\mathbf{A}_1 \mathbf{X} = \mathbf{0}$. Then (2) follows from (same proof as the preceding result part (1))

$$\mathbf{V}_0 \mathbf{A}_0 \mathbf{V}_0 \mathbf{A}_0 \mathbf{V}_0 = \mathbf{V}_0 \mathbf{A}_0 \mathbf{V}_0.$$

For (3), since $\mathbf{A}_1 \mathbf{X} = \mathbf{0}$, enough to show

$$\mathbf{V} \mathbf{A}_0 \mathbf{V} \mathbf{A}_0 \mathbf{X} \mathbf{b} - \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_0 \mathbf{X} \mathbf{b} = \mathbf{V} \mathbf{A}_0 \mathbf{X} \mathbf{b}.$$

With $\mathbf{V} \mathbf{A}_0 \mathbf{V} \mathbf{A}_0 \mathbf{X} \mathbf{b} = \mathbf{V} \mathbf{A}_0 \mathbf{X} \mathbf{b}$, it's enough to show

$$\mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_0 \mathbf{X} \mathbf{b} = \mathbf{0}.$$

This is verified by

$$\begin{aligned} & \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_0 \mathbf{X} \mathbf{b} \\ &= \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{A}_0 \mathbf{V} \mathbf{c} \quad (\text{because } \mathbf{X} \mathbf{b} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V}) = \mathcal{C}(\mathbf{V}_0) \text{ and } \mathbf{A}_0 \mathbf{X}_0 = \mathbf{0}) \\ &= \mathbf{V} \mathbf{A}_1 \mathbf{V} \mathbf{c} \quad (\text{preceding result}) \\ &= \mathbf{V} \mathbf{A}_1 \mathbf{X} \mathbf{b} \quad (\text{because } \mathbf{A}_1 \mathbf{X}_0 = \mathbf{0}) \\ &= \mathbf{0}. \quad (\text{because } \mathbf{A}_1 \mathbf{X} = \mathbf{0}) \end{aligned}$$

□

- (1) $\mathbf{V}\mathbf{A}_1\mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V} = \mathbf{0}$
- (2) If $\mathbf{X}\mathbf{b} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})$, then $\mathbf{V}\mathbf{A}_1\mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{X}\mathbf{b} = \mathbf{0}$.
- (3) $\mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V}\mathbf{A}_1\mathbf{X}\mathbf{b} = \mathbf{0}$
- (4) $\mathbf{b}^T\mathbf{X}^T(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V}\mathbf{A}_1\mathbf{X}\mathbf{b} = 0$.

Remark: The third part of the theorem (independence) is checked by this result.

Proof. For (1),

$$\begin{aligned}
& \mathbf{V}\mathbf{A}_1\mathbf{V}(\mathbf{A}_0 - \mathbf{A}_1)\mathbf{V} \\
&= \mathbf{V}\mathbf{A}_1\mathbf{V}\mathbf{A}_0\mathbf{V} - \mathbf{V}\mathbf{A}_1\mathbf{V}\mathbf{A}_1\mathbf{V} \\
&= \mathbf{V}\mathbf{A}_1\mathbf{V} - \mathbf{V}\mathbf{A}_1\mathbf{V} \\
&= \mathbf{0}.
\end{aligned}$$

(2) follows from the proof in the preceding result.

(3) and (4) are trivial since $\mathbf{A}_1\mathbf{X} = \mathbf{0}$. □

- Before we state the F-test, we show that the constrained model is consistent if and only if $\mathbf{X}\mathbf{b} \in \mathcal{C}(\mathbf{X}_0)$.

1. (Null model) If $\mathbf{X}\mathbf{b} \in \mathcal{C}(\mathbf{X}_0)$, then $\mathbf{P}(\mathbf{Y} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})) = 1$ and

$$\mathbf{b}^T\mathbf{X}^T\mathbf{A}_0\mathbf{X}\mathbf{b} = 0.$$

2. (Alternative model) If $\mathbf{X}\mathbf{b} \notin \mathcal{C}(\mathbf{X}_0)$, then either

- $\mathbf{X}\mathbf{b} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})$, $\mathbf{P}(\mathbf{Y} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})) = 1$, and $\mathbf{b}^T\mathbf{X}^T\mathbf{A}_0\mathbf{X}\mathbf{b} > 0$, or
- $\mathbf{X}\mathbf{b} \notin \mathcal{C}(\mathbf{X}_0, \mathbf{V})$ and $\mathbf{P}(\mathbf{Y} \notin \mathcal{C}(\mathbf{X}_0, \mathbf{V})) = 1$.

Proof. For Part 1, note $\mathbf{Y} = \mathbf{X}\mathbf{b} + \mathbf{e}$, where $\mathbf{e} \in \mathcal{C}(\mathbf{V})$ and $\mathbf{A}_0\mathbf{X}_0 = \mathbf{0}$. For part 2

- If $\mathbf{X}\mathbf{b} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})$, we need to show $\mathbf{b}^T\mathbf{X}^T\mathbf{A}_0\mathbf{X}\mathbf{b} > 0$. Here is a proof thanks to Po Ning. Suppose

$$\mathbf{b}^T\mathbf{X}^T\mathbf{A}_0\mathbf{X}\mathbf{b} = \mathbf{b}^T\mathbf{X}^T[\mathbf{V}_0^- - \mathbf{V}_0^-\mathbf{X}_0(\mathbf{X}_0^T\mathbf{V}_0^-\mathbf{X}_0)^-\mathbf{X}_0^T\mathbf{V}_0^-]\mathbf{X}\mathbf{b} = 0.$$

Note this quantity is invariant to the choice of generalized inverse \mathbf{V}_0^- due to the fact $\mathbf{X}\mathbf{b} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V}) = \mathcal{C}(\mathbf{V}_0)$ and $\mathcal{C}(\mathbf{X}_0) \subset \mathcal{C}(\mathbf{V}_0)$. Without loss of

generality, we can use \mathbf{V}_0^+ , which is psd and thus can be decomposed as $\mathbf{V}^+ = \mathbf{L}\mathbf{L}^T$, $\mathbf{L} \in \mathbb{R}^{n \times r}$ full column rank. Therefore

$$\begin{aligned}
& \mathbf{b}^T \mathbf{X}^T [\mathbf{V}_0^- - \mathbf{V}_0^- \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{V}_0^- \mathbf{X}_0)^- \mathbf{X}_0^T \mathbf{V}_0^-] \mathbf{X} \mathbf{b} \\
&= \mathbf{b}^T \mathbf{X}^T [\mathbf{V}_0^+ - \mathbf{V}_0^+ \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{V}_0^+ \mathbf{X}_0)^- \mathbf{X}_0^T \mathbf{V}_0^+] \mathbf{X} \mathbf{b} \\
&= \mathbf{b}^T \mathbf{X}^T [\mathbf{L}\mathbf{L}^T - \mathbf{L}\mathbf{L}^T \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{V}_0^+ \mathbf{X}_0)^- \mathbf{X}_0^T \mathbf{L}\mathbf{L}^T] \mathbf{X} \mathbf{b} \\
&= \mathbf{b}^T \mathbf{X}^T \mathbf{L} [\mathbf{I} - \mathbf{L}^T \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{L}\mathbf{L}^T \mathbf{X}_0)^- \mathbf{X}_0^T \mathbf{L}] \mathbf{L}^T \mathbf{X} \mathbf{b} \\
&= \mathbf{b}^T \mathbf{X}^T \mathbf{L} (\mathbf{I} - \mathbf{P}_Z) \mathbf{L}^T \mathbf{X} \mathbf{b} \quad (\text{let } \mathbf{Z} = \mathbf{L}^T \mathbf{X}_0) \\
&= \mathbf{b}^T \mathbf{X}^T \mathbf{L} (\mathbf{I} - \mathbf{P}_Z) (\mathbf{I} - \mathbf{P}_Z) \mathbf{L}^T \mathbf{X} \mathbf{b} \\
&= \|(\mathbf{I} - \mathbf{P}_Z) \mathbf{L}^T \mathbf{X} \mathbf{b}\|_2^2 \\
&= 0,
\end{aligned}$$

implying that $(\mathbf{I} - \mathbf{P}_Z) \mathbf{L}^T \mathbf{X} \mathbf{b} = \mathbf{0}$ and thus

$$\mathbf{L}^T \mathbf{X} \mathbf{b} \in \mathcal{C}(\mathbf{Z}) = \mathcal{C}(\mathbf{L}^T \mathbf{X}_0).$$

Then what?

- If $\mathbf{X} \mathbf{b} \notin \mathcal{C}(\mathbf{X}_0, \mathbf{V})$, then $\mathbf{Y} = \mathbf{X} \mathbf{b} + \mathbf{e}$ cannot be in $\mathcal{C}(\mathbf{X}_0, \mathbf{V})$ since $\mathbf{e} \in \mathcal{C}(\mathbf{V})$. If it does, then $\mathbf{X} \mathbf{b} = \mathbf{Y} - \mathbf{e} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})$, a contradiction. Therefore $\mathbf{P}(\mathbf{Y} \notin \mathcal{C}(\mathbf{X}_0, \mathbf{V})) = 1$.

□

- (F test under Aitken model) The test rejects $H_0 : \mathbf{y} \sim N(\mathbf{X}_0 \mathbf{c}, \sigma^2 \mathbf{V})$
 - if $\mathbf{y} \notin \mathcal{C}(\mathbf{X}_0, \mathbf{V})$, or
 - if $\mathbf{y} \in \mathcal{C}(\mathbf{X}_0, \mathbf{V})$ and

$$\frac{(\text{SSE}_0 - \text{SSE}) / \text{tr}((\mathbf{A}_0 - \mathbf{A}_1) \mathbf{V})}{\text{SSE} / \text{tr}(\mathbf{A}_1 \mathbf{V})} > F_{\text{tr}((\mathbf{A}_0 - \mathbf{A}_1) \mathbf{V}), \text{tr}(\mathbf{A}_1 \mathbf{V}), \alpha}.$$

It is a level- α test because $\mathbf{P}(\mathbf{Y} \notin \mathcal{C}(\mathbf{X}_0, \mathbf{V})) = 0$ under H_0 . The power of the test is always greater than α since if $\mathbf{b}^T \mathbf{X}^T \mathbf{A}_0 \mathbf{X} \mathbf{b} = 0$, then the test will reject H_0 with probability 1.

19 Lecture 19: Nov 11

Announcement

- HW6 returned.
- HW7 due today.
- Midterm 2 this Wed Nov 13 @ 11:45AM-1PM.

Last time

Testing under Aitken model $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{V})$ – main results and proofs.

Today

- Q&A
- Estimation of σ^2 .

Q&A for midterm 2

- Notes p101 about independence of numerator and denominator of the F statistic.
- Notes p103 about cancellation of \mathbf{T} .

MVQUE of σ^2

- So far we have focused on estimation and testing of \mathbf{b} . What's a good estimator of σ^2 ? We only know that the *least squares estimator*

$$\hat{\sigma}^2 = \text{SSE}/(n - \text{rank}(\mathbf{X}))$$

is unbiased under the Gauss-Markov assumption.

With additional normality assumption, does this have minimum variance within certain class of estimators? Does this have minimum mean square error (MSE) within certain class of estimators?

- Since σ^2 is a quadratic concept, we consider estimation of σ^2 by a quadratic function of \mathbf{y} .

We call any estimator $\mathbf{y}^T \mathbf{A} \mathbf{y}$, where \mathbf{A} is symmetric and non-stochastic, a *quadratic estimator*.

- Any estimator $\hat{\sigma}^2$ such that

$$E(\hat{\sigma}^2) = \sigma^2 \quad \text{for all } \mathbf{b} \in \mathbb{R}^p \text{ and } \sigma^2 > 0$$

is called an *unbiased estimator* of σ^2 .

- The *minimum variance unbiased estimator* (MVQUE, also called BQUE) of σ^2 is a quadratic unbiased estimator of σ^2 , say $\hat{\sigma}^2$, such that

$$\text{Var}(\hat{\sigma}^2) \leq \text{Var}(\hat{\tau}^2)$$

for all quadratic unbiased estimators $\hat{\tau}^2$ of $\hat{\sigma}^2$.

- Since we are dealing with quadratic function of a multivariate normal, we recall a useful result (shown in HW5).

Let $\mathbf{Y} \sim N(\boldsymbol{\mu}, \mathbf{V})$, \mathbf{V} nonsingular, and let $U = \mathbf{Y}^T \mathbf{A} \mathbf{Y}$ for \mathbf{A} symmetric.

1. The mgf of U is

$$\begin{aligned} m_U(t) &= |\mathbf{I} - 2t\mathbf{A}\mathbf{V}|^{-1/2} e^{-\boldsymbol{\mu}^T \mathbf{V}^{-1} \boldsymbol{\mu}/2 + \boldsymbol{\mu}^T (\mathbf{V} - 2t\mathbf{V}\mathbf{A}\mathbf{V})^{-1} \boldsymbol{\mu}/2} \\ &= |\mathbf{I} - 2t\mathbf{A}\mathbf{V}|^{-1/2} e^{-t\boldsymbol{\mu}^T (\mathbf{I} - 2t\mathbf{A}\mathbf{V})^{-1} \mathbf{A}\boldsymbol{\mu}}. \end{aligned}$$

2. The mean and variance of U is

$$\begin{aligned} EU &= \text{tr}(\mathbf{A}\mathbf{V}) + \boldsymbol{\mu}^T \mathbf{A}\boldsymbol{\mu} \\ \text{Var } U &= 2\text{tr}(\mathbf{A}\mathbf{V})^2 + 4\boldsymbol{\mu}^T \mathbf{A}\mathbf{V}\mathbf{A}\boldsymbol{\mu}. \end{aligned}$$

- Under Gauss-Markov normal model $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I})$, the least squares estimator

$$\hat{\sigma}^2 = \frac{\text{SSE}}{n - \text{rank}(\mathbf{X})} = \frac{\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}}{n - \text{rank}(\mathbf{X})}$$

is MVQUE.

Proof. Let $\hat{\sigma}^2 = \mathbf{y}^T \mathbf{A} \mathbf{y}$ be a quadratic estimator of σ^2 . Then

$$\begin{aligned}\hat{\sigma}^2 &= (\mathbf{X}\mathbf{b} + \mathbf{e})^T \mathbf{A} (\mathbf{X}\mathbf{b} + \mathbf{e}) \\ &= \mathbf{b}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{b} + 2\mathbf{b}^T \mathbf{X}^T \mathbf{A} \mathbf{e} + \mathbf{e}^T \mathbf{A} \mathbf{e}\end{aligned}$$

and

$$\mathbb{E}(\hat{\sigma}^2) = \mathbf{b}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{b} + \sigma^2 \text{tr} \mathbf{A} = \sigma^2$$

for all \mathbf{b} and $\sigma^2 > 0$ if and only if

$$\mathbf{X}^T \mathbf{A} \mathbf{X} = \mathbf{0}_{p \times p} \quad \text{and} \quad \text{tr} \mathbf{A} = 1.$$

The variance of $\hat{\sigma}^2$ is (why?)

$$\begin{aligned}\text{Var}(\hat{\sigma}^2) &= \text{Var}(\mathbf{y}^T \mathbf{A} \mathbf{y}) \\ &= 2\sigma^4 (\text{tr} \mathbf{A}^2 + 2\boldsymbol{\gamma}^T \mathbf{X}^T \mathbf{A}^2 \mathbf{X} \boldsymbol{\gamma}),\end{aligned}$$

where $\boldsymbol{\gamma} = \mathbf{b}/\sigma$. To seek the MVQUE, we solve the optimization problem

$$\begin{aligned}\text{minimize} \quad & \text{tr} \mathbf{A}^2 + 2\boldsymbol{\gamma}^T \mathbf{X}^T \mathbf{A}^2 \mathbf{X} \boldsymbol{\gamma} \\ \text{subject to} \quad & \mathbf{X}^T \mathbf{A} \mathbf{X} = \mathbf{0} \quad \text{and} \quad \text{tr} \mathbf{A} = 1.\end{aligned}$$

(It turns out the optimal \mathbf{A} is independent of the unknown parameter $\boldsymbol{\gamma}$.) We form the Lagrangian

$$L(\mathbf{A}, \lambda, \mathbf{L}) = \frac{1}{2} \text{tr} \mathbf{A}^2 + \boldsymbol{\gamma}^T \mathbf{X}^T \mathbf{A}^2 \mathbf{X} \boldsymbol{\gamma} - \lambda (\text{tr} \mathbf{A} - 1) - \text{tr}(\mathbf{L}^T \mathbf{X}^T \mathbf{A} \mathbf{X}),$$

where λ and $\mathbf{L} \in \mathbb{R}^{p \times p}$ are Lagrange multipliers. Since $\mathbf{X}^T \mathbf{A} \mathbf{X}$ is symmetric, we assume \mathbf{L} is symmetric too. Setting derivative of the Lagrangian to 0 gives

$$\begin{aligned}\mathbf{A} + \mathbf{X} \boldsymbol{\gamma} \boldsymbol{\gamma}^T \mathbf{X}^T \mathbf{A} + \mathbf{A} \mathbf{X} \boldsymbol{\gamma} \boldsymbol{\gamma}^T \mathbf{X}^T - \lambda \mathbf{I} &= \mathbf{X} \mathbf{L} \mathbf{X}^T \\ \mathbf{X}^T \mathbf{A} \mathbf{X} &= \mathbf{0}_{p \times p} \\ \text{tr} \mathbf{A} &= 1\end{aligned}$$

(Check Lecture 8 p56 for derivatives of trace functions.) Pre- and post-multiplying the first equation by $\mathbf{P}_{\mathbf{X}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ yields

$$-\lambda \mathbf{P}_{\mathbf{X}} = \mathbf{X}^T \mathbf{L} \mathbf{X}.$$

Substitution back to the first equation shows

$$\mathbf{A} = \lambda(\mathbf{I} - \mathbf{P}_X) - (\mathbf{X}\boldsymbol{\gamma}\boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{A} + \mathbf{A}\mathbf{X}\boldsymbol{\gamma}\boldsymbol{\gamma}^T\mathbf{X}^T).$$

Taking trace on both sides gives $\text{tr}(\mathbf{A}) = \lambda(n - r) = 1$, i.e.,

$$\lambda = \frac{1}{n - r}.$$

Taking square on both sides gives

$$\begin{aligned} \mathbf{A}^2 &= \lambda^2(\mathbf{I} - \mathbf{P}_X) - \lambda\mathbf{A}\mathbf{X}\boldsymbol{\gamma}\boldsymbol{\gamma}^T\mathbf{X}^T \\ &\quad - \lambda\mathbf{X}\boldsymbol{\gamma}\boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{A} + \mathbf{X}\boldsymbol{\gamma}\boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{A}^2\mathbf{X}\boldsymbol{\gamma}\boldsymbol{\gamma}^T\mathbf{X}^T \\ &\quad + \mathbf{A}\mathbf{X}\boldsymbol{\gamma}\boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{X}\boldsymbol{\gamma}\boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{A}. \end{aligned}$$

Taking trace on both sides

$$\text{tr}(\mathbf{A}^2) = \frac{1}{n - r} + 2(\boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{A}^2\mathbf{X}\boldsymbol{\gamma})(\boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{X}\boldsymbol{\gamma}).$$

Thus the objective function can be expressed as

$$\begin{aligned} &\text{tr}(\mathbf{A}^2) + 2\boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{A}^2\mathbf{X}\boldsymbol{\gamma} \\ &= \frac{1}{n - r} + 2(\boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{A}^2\mathbf{X}\boldsymbol{\gamma})(1 + \boldsymbol{\gamma}^T\mathbf{X}^T\mathbf{X}\boldsymbol{\gamma}), \end{aligned}$$

which is minimized by taking $\mathbf{A}\mathbf{X}\boldsymbol{\gamma} = \mathbf{0}$. Therefore

$$\mathbf{A} = \lambda(\mathbf{I} - \mathbf{P}_X) = \frac{\mathbf{I} - \mathbf{P}_X}{n - r}.$$

□

20 Lecture 20: Nov 18

Last time

Under Gauss-Markov normal model, the least squares estimator $\hat{\sigma}^2 = \text{SSE}/(n - r)$ is MVQUE.

Today

- Best quadratic invariant estimator of σ^2 .
- Variance components and mixed model (Chapter 8): introduction.

Best quadratic invariant estimation of σ^2

- In this section, we seek the quadratic estimator of σ^2 with smallest mean squared error (MSE), within the class of invariant quadratic estimators. Recall that the MSE of an estimator of σ^2 is

$$\begin{aligned} E(\hat{\sigma}^2 - \sigma^2)^2 &= E[\hat{\sigma}^2 - E(\hat{\sigma}^2) + E(\hat{\sigma}^2) - \sigma^2]^2 \\ &= \text{Var}(\hat{\sigma}^2) + [E(\hat{\sigma}^2) - \sigma^2]^2 \\ &= \text{Variance} + \text{Bias}^2. \end{aligned}$$

- Assume $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \sigma^2\mathbf{V})$, \mathbf{V} nonsingular. Let $\mathbf{y}^T \mathbf{A}\mathbf{y}$ be a quadratic estimator. In HW5, we showed that its first two moments are

$$\begin{aligned} E(\mathbf{y}^T \mathbf{A}\mathbf{y}) &= \mathbf{b}^T \mathbf{X}^T \mathbf{A}\mathbf{X}\mathbf{b} + \sigma^2 \text{tr}(\mathbf{A}\mathbf{V}) \\ \text{Var}(\mathbf{y}^T \mathbf{A}\mathbf{y}) &= 4\sigma^2 \mathbf{b}^T \mathbf{X}^T \mathbf{A}\mathbf{V}\mathbf{A}\mathbf{X}\mathbf{b} + 2\sigma^4 \text{tr}(\mathbf{A}\mathbf{V})^2 \end{aligned}$$

and its mgf is

$$\begin{aligned} m(t) &= |\mathbf{I} - 2t\mathbf{A}\mathbf{V}|^{-1/2} e^{-\mathbf{b}^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}\mathbf{b}/2 + \mathbf{b}^T \mathbf{X}^T (\mathbf{V} - 2t\mathbf{V}\mathbf{A}\mathbf{V})^{-1} \mathbf{X}\mathbf{b}/2} \\ &= |\mathbf{I} - 2t\mathbf{A}\mathbf{V}|^{-1/2} e^{-t\mathbf{b}^T \mathbf{X}^T (\mathbf{I} - 2t\mathbf{A}\mathbf{V})^{-1} \mathbf{A}\mathbf{X}\mathbf{b}}. \end{aligned}$$

Thus distribution is independent of \mathbf{b} if and only if $\mathbf{A}\mathbf{X} = \mathbf{0}_{n \times p}$.

Proof. The “if” part is trivial from the mgf. For the “only if” part, $\mathbf{b}^T \mathbf{X}^T \mathbf{A}\mathbf{V}\mathbf{A}\mathbf{X}\mathbf{b} = 0$ for all \mathbf{b} implies $\mathbf{X}^T \mathbf{A}\mathbf{V}\mathbf{A}\mathbf{X} = \mathbf{0}_{p \times p}$. Since \mathbf{V} is nonsingular, $\mathbf{A}\mathbf{X} = \mathbf{0}$. \square

- We say a quadratic estimator $\mathbf{y}^T \mathbf{A} \mathbf{y}$ is *invariant* under translation of \mathbf{b} if $\mathbf{A} \mathbf{X} = \mathbf{0}$.

A quadratic invariant estimator $\hat{\sigma}^2$ is the *best quadratic invariant estimator* of σ^2 if

$$\mathbb{E}(\hat{\sigma}^2 - \sigma^2)^2 \leq \mathbb{E}(\hat{\tau}^2 - \sigma^2)^2$$

for all quadratic invariant estimator $\hat{\tau}^2$ of σ^2 .

- Under Gauss-Markov normal model $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \sigma^2 \mathbf{I})$, the best quadratic invariant estimator of σ^2 is

$$\hat{\sigma}^2 = \frac{\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}}{n - r + 2} = \frac{\text{SSE}}{n - r + 2}.$$

Remark: By introducing bias (shrinkage) to the unbiased estimator $\text{SSE}/(n-r)$, we achieve optimal MSE.

Proof. Let $\hat{\sigma}^2 = \mathbf{y}^T \mathbf{A} \mathbf{y}$ be a quadratic estimator. Invariance imposes $\mathbf{A} \mathbf{X} = \mathbf{0}$. Thus

$$\hat{\sigma}^2 = \mathbf{y}^T \mathbf{A} \mathbf{y} = (\mathbf{X}\mathbf{b} + \mathbf{e})^T \mathbf{A} (\mathbf{X}\mathbf{b} + \mathbf{e}) = \mathbf{e}^T \mathbf{A} \mathbf{e}$$

and the MSE is (why?)

$$\mathbb{E}(\hat{\sigma}^2 - \sigma^2)^2 = 2\sigma^4 \text{tr} \mathbf{A}^2 + \sigma^4 (1 - \text{tr} \mathbf{A})^2.$$

We seek the best quadratic invariant estimator by solving

$$\begin{aligned} & \text{minimize} && 2\text{tr} \mathbf{A}^2 + (1 - \text{tr} \mathbf{A})^2 \\ & \text{subject to} && \mathbf{A} \mathbf{X} = \mathbf{0}. \end{aligned}$$

We form the Lagrangian

$$L(\mathbf{A}, \mathbf{L}) = \text{tr} \mathbf{A}^2 + \frac{1}{2} (1 - \text{tr} \mathbf{A})^2 - \text{tr}(\mathbf{L}^T \mathbf{A} \mathbf{X})$$

and set its gradient to zero

$$\begin{aligned} 2\mathbf{A} - (1 - \text{tr} \mathbf{A}) \mathbf{I}_n &= \mathbf{L} \mathbf{X}^T \\ \mathbf{A} \mathbf{X} &= \mathbf{0}. \end{aligned}$$

Pre-multiplying the first equation by \mathbf{A} yields

$$2\mathbf{A}^2 - (1 - \text{tr}\mathbf{A})\mathbf{A} = \mathbf{A}\mathbf{L}\mathbf{X}^T.$$

Post-multiplying by $\mathbf{P}_\mathbf{X} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ shows

$$\mathbf{A}\mathbf{L}\mathbf{X}^T = \mathbf{0}.$$

Substitution back gives

$$\mathbf{A}^2 = \frac{1 - \text{tr}\mathbf{A}}{2}\mathbf{A}.$$

Let $\text{rank}(\mathbf{A}) = \rho$. By a result shown in Lecture 3 (p16),

$$\text{tr}\mathbf{A} = \frac{1 - \text{tr}\mathbf{A}}{2}\text{rank}(\mathbf{A}) = \frac{1 - \text{tr}\mathbf{A}}{2}\rho.$$

Thus

$$\text{tr}(\mathbf{A}) = \frac{\rho}{\rho + 2} \quad \text{and} \quad \text{tr}(\mathbf{A}^2) = \frac{\rho}{(\rho + 2)^2}.$$

The objective function becomes

$$\frac{2}{\rho + 2}.$$

Since $\mathbf{A}\mathbf{X} = \mathbf{0}$, the maximizing ρ is $n - \text{rank}(\mathbf{X}) = n - r$. Note $(n - r + 2)\mathbf{A}$ is symmetric, idempotent, and orthogonal to $\mathcal{C}(\mathbf{X})$, thus equals $\mathbf{I} - \mathbf{P}_\mathbf{X}$. Therefore

$$\mathbf{A} = (\mathbf{I} - \mathbf{P}_\mathbf{X})/(n - r + 2).$$

□

- In summary, we have the following estimators for σ^2
 - MVQUE (aka least square estimator): $\text{SSE}/(n - r)$
 - MLE: SSE/n
 - Best quadratic invariant estimate: $\text{SSE}/(n - r + 2)$

Variance components and mixed models

- Traditionally, linear models have been divided into three categories:
 - *Fixed effects model*: $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$, where \mathbf{b} is fixed.
 - *Random effects model*: $\mathbf{y} = \mathbf{Z}\mathbf{u} + \mathbf{e}$, where \mathbf{u} is random.
 - *Mixed effects model*: $\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e}$, where \mathbf{b} is fixed and \mathbf{u} is random.
- In a mixed effects model

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

- $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a design matrix for fixed effects $\mathbf{b} \in \mathbb{R}^p$.
- $\mathbf{Z} \in \mathbb{R}^{n \times q}$ is a design matrix for random effects $\mathbf{u} \in \mathbb{R}^q$.
- The most general assumption is $\mathbf{e} \in N(\mathbf{0}_n, \mathbf{R})$, $\mathbf{u} \in N(\mathbf{0}_q, \mathbf{G})$, and \mathbf{e} is independent of \mathbf{u} .

In many applications, $\mathbf{e} \sim N(\mathbf{0}_n, \sigma^2 \mathbf{I}_n)$ and

$$\mathbf{Z}\mathbf{u} = \begin{pmatrix} \mathbf{Z}_1 & \cdots & \mathbf{Z}_m \end{pmatrix} \begin{pmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_m \end{pmatrix} = \mathbf{Z}_1\mathbf{u}_1 + \cdots + \mathbf{Z}_m\mathbf{u}_m,$$

where $\mathbf{u}_i \sim N(\mathbf{0}_{q_i}, \sigma_i^2 \mathbf{I}_{q_i})$, $\sum_{i=1}^m q_i = q$. \mathbf{e} and \mathbf{u}_i , $i = 1, \dots, m$, are jointly independent. Then the covariance of responses \mathbf{y}

$$\mathbf{V}(\sigma^2, \sigma_1^2, \dots, \sigma_m^2) = \sigma^2 \mathbf{I} + \sum_{i=1}^m \sigma_i^2 \mathbf{Z}_i \mathbf{Z}_i^T$$

is a function of the *variance components* $(\sigma^2, \sigma_1^2, \dots, \sigma_m^2)$.

- Primary goal of the mixed model (aka variance components model) is to
 - estimation and testing of the fixed effects \mathbf{b}
 - estimation and testing of the variance components $(\sigma^2, \sigma_1^2, \dots, \sigma_m^2)$
 - prediction

- Example: One-way ANOVA with random treatment effects (JM 8.2)

$$y_{ij} = \mu + \alpha_i + e_{ij}, \quad i = 1, \dots, a, j = 1, \dots, n_i, \sum_{i=1}^a n_i = n.$$

So far we considered the fixed effects model where the treatment effects α_i are assumed to be constant. In random effects model, the common assumption is

$$\begin{aligned} \alpha_i &\sim N(0, \sigma_a^2) \\ e_{ij} &\sim N(0, \sigma^2) \end{aligned}$$

are jointly independent. Equivalently the model is

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e} \\ &= \mathbf{1}_n\mu + \begin{pmatrix} \mathbf{1}_{n_1} & & & \\ & \mathbf{1}_{n_2} & & \\ & & \ddots & \\ & & & \mathbf{1}_{n_a} \end{pmatrix} \mathbf{u} + \mathbf{e}, \end{aligned}$$

where $\mathbf{u} \sim N(\mathbf{0}_a, \sigma_a^2 \mathbf{I}_a)$ and $\mathbf{e} \sim N(\mathbf{0}_n, \sigma^2 \mathbf{I}_n)$ are independent. The covariance matrix is

$$\mathbf{V} = \text{BlkDiag}(\mathbf{V}_1, \dots, \mathbf{V}_a) = \begin{pmatrix} \mathbf{V}_1 & & & \\ & \mathbf{V}_2 & & \\ & & \ddots & \\ & & & \mathbf{V}_a \end{pmatrix},$$

where $\mathbf{V}_i = \sigma^2 \mathbf{I}_{n_i} + \sigma_a^2 \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T$.

The model parameters are $(\mu, \sigma^2, \sigma_a^2)$.

- Example: Two-way ANOVA with random effects

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + e_{ijk}, \quad i = 1, \dots, a, j = 1, \dots, b, k = 1, \dots, n_{ij},$$

where

$$\begin{aligned} \alpha_i &\sim N(0, \sigma_a^2), & \beta_j &\sim N(0, \sigma_b^2) \\ \gamma_{ij} &\sim N(0, \sigma_c^2), & e_{ijk} &\sim N(0, \sigma^2) \end{aligned}$$

are jointly independent.

Equivalently the model is

$$\mathbf{y} = \mathbf{1}_n \mu + \begin{pmatrix} \mathbf{1}_{n_{11}} \\ \mathbf{1}_{n_{12}} \\ \vdots \\ \mathbf{1}_{n_{1b}} \\ \vdots \\ \mathbf{1}_{n_{a1}} \\ \mathbf{1}_{n_{a2}} \\ \vdots \\ \mathbf{1}_{n_{ab}} \end{pmatrix} \mathbf{u}_a + \begin{pmatrix} \mathbf{1}_{n_{11}} & & & & \\ & \mathbf{1}_{n_{12}} & & & \\ & & \ddots & & \\ & & & \mathbf{1}_{1b} & \\ \vdots & & & \vdots & \\ \mathbf{1}_{n_{a1}} & & & & \\ & \mathbf{1}_{n_{a2}} & & & \\ & & \ddots & & \\ & & & & \mathbf{1}_{ab} \end{pmatrix} \mathbf{u}_b + \begin{pmatrix} \mathbf{1}_{n_{11}} & & & & & & & \\ & \mathbf{1}_{n_{12}} & & & & & & \\ & & \ddots & & & & & \\ & & & \mathbf{1}_{1b} & & & & \\ & & & & \ddots & & & \\ & & & & & \mathbf{1}_{n_{a1}} & & \\ & & & & & & \mathbf{1}_{n_{a2}} & \\ & & & & & & & \ddots \\ & & & & & & & & \mathbf{1}_{ab} \end{pmatrix} \mathbf{u}_c + \mathbf{e},$$

where

$$\begin{aligned}
 \mathbf{u}_a &\sim N(\mathbf{0}_a, \sigma_a^2 \mathbf{I}_a), & \mathbf{u}_b &\sim N(\mathbf{0}_b, \sigma_b^2 \mathbf{I}_b), \\
 \mathbf{u}_c &\sim N(\mathbf{0}_{ab}, \sigma_c^2 \mathbf{I}_{ab}), & \mathbf{e} &\sim N(\mathbf{0}_n, \sigma^2 \mathbf{I}_n)
 \end{aligned}$$

are jointly independent.

The model parameters are $(\mu, \sigma^2, \sigma_a^2, \sigma_b^2, \sigma_c^2)$.

- We may also have two-way ANOVA mixed model, where one factor has fixed effects and the other factor has random effects (JM Example 8.1 p190).

21 Lecture 21: Nov 20

Announcement

- HW7 returned.
- Midterm 2 returned (85.63 ± 6.32).
- HW8 posted and due Mon Dec 2.
- No afternoon session today.
- Email your comments (compliments?) on TA to me.

Last time

- Best (minimal MSE) quadratic invariant estimator of σ^2 .
- Variance components and mixed model (Chapter 8): introduction.

Today

- Variance component estimation: MLE.
- Variance component estimation: REML.

Variance component estimation: MLE (JM 8.4.1)

In this section, we pursue MLE for the variance components model

$$\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \mathbf{V}),$$

where

$$\mathbf{V} = \sum_{i=0}^m \sigma_i^2 \mathbf{V}_i,$$

with \mathbf{V}_i psd and \mathbf{V} nonsingular. Parameters are \mathbf{b} and $(\sigma_0^2, \sigma_1^2, \dots, \sigma_m^2)$.

- Assume \mathbf{V} is nonsingular, then the log-likelihood function is

$$\begin{aligned} L(\mathbf{b}, \sigma_0^2, \sigma_1^2, \dots, \sigma_m^2) \\ = -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \ln \det(\mathbf{V}) - \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{b})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b}). \end{aligned}$$

- In general we need some iterative optimization algorithm (Fisher scoring, EM algorithm) to find the MLE.
- The following facts are useful when deriving the derivatives of log-likelihood.

1. The partial derivative of a matrix $\mathbf{B} = (b_{ij})$ with respect to a scalar parameter θ is the matrix

$$\frac{\partial}{\partial \theta} \mathbf{B} = \left(\frac{\partial}{\partial \theta} b_{ij} \right).$$

2. Because the trace function is linear,

$$\frac{\partial}{\partial \theta} \text{tr}(\mathbf{B}) = \text{tr} \left(\frac{\partial}{\partial \theta} \mathbf{B} \right).$$

3. The product rule of differentiation implies

$$\frac{\partial}{\partial \theta} (\mathbf{B}\mathbf{C}) = \left(\frac{\partial}{\partial \theta} \mathbf{B} \right) \mathbf{C} + \mathbf{B} \left(\frac{\partial}{\partial \theta} \mathbf{C} \right).$$

4. The derivative of a matrix inverse is

$$\frac{\partial}{\partial \theta} \mathbf{B}^{-1} = -\mathbf{B}^{-1} \left(\frac{\partial}{\partial \theta} \mathbf{B} \right) \mathbf{B}^{-1}$$

Proof. Solving for $\frac{\partial}{\partial \theta} \mathbf{B}^{-1}$ in

$$\begin{aligned} \mathbf{0} &= \frac{\partial}{\partial \theta} \mathbf{I} \\ &= \frac{\partial}{\partial \theta} (\mathbf{B}^{-1} \mathbf{B}) \\ &= \left(\frac{\partial}{\partial \theta} \mathbf{B}^{-1} \right) \mathbf{B} + \mathbf{B}^{-1} \left(\frac{\partial}{\partial \theta} \mathbf{B} \right). \end{aligned}$$

□

5. If \mathbf{B} is a square nonsingular matrix, then

$$\begin{aligned}\frac{\partial}{\partial \theta} \det \mathbf{B} &= \det \mathbf{B} \cdot \operatorname{tr} \left(\mathbf{B}^{-1} \frac{\partial}{\partial \theta} \mathbf{B} \right), \\ \frac{\partial}{\partial \theta} \ln \det \mathbf{B} &= \operatorname{tr} \left(\mathbf{B}^{-1} \frac{\partial}{\partial \theta} \mathbf{B} \right).\end{aligned}$$

- The score, observed information matrix, and expected (Fisher) information matrix of the variance components model are listed below (HW8).

1. Score (gradient) vector is

$$\begin{aligned}\frac{\partial}{\partial b_i} L &= \mathbf{e}_i^T \mathbf{X}^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \mathbf{b}) \\ \frac{\partial}{\partial \sigma_i^2} L &= -\frac{1}{2} \operatorname{tr}(\mathbf{V}^{-1} \mathbf{V}_i) + \frac{1}{2} (\mathbf{y} - \mathbf{X} \mathbf{b})^T \mathbf{V}^{-1} \mathbf{V}_i \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \mathbf{b}).\end{aligned}$$

2. The observed information matrix has entries

$$\begin{aligned}-\frac{\partial^2}{\partial b_i \partial b_j} L &= \mathbf{e}_i^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \mathbf{e}_j \\ -\frac{\partial^2}{\partial \sigma_i^2 \partial b_j} L &= \mathbf{e}_j^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{V}_i \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \mathbf{b}) \\ -\frac{\partial^2}{\partial \sigma_i^2 \partial \sigma_j^2} L &= -\frac{1}{2} \operatorname{tr}(\mathbf{V}^{-1} \mathbf{V}_i \mathbf{V}^{-1} \mathbf{V}_j) + (\mathbf{y} - \mathbf{X} \mathbf{b})^T \mathbf{V}^{-1} \mathbf{V}_i \mathbf{V}^{-1} \mathbf{V}_j \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \mathbf{b}).\end{aligned}$$

3. The expected (Fisher) information matrix has entries

$$\begin{aligned}\mathbb{E} \left(-\frac{\partial^2}{\partial b_i \partial b_j} L \right) &= \mathbf{e}_i^T \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \mathbf{e}_j \\ \mathbb{E} \left(-\frac{\partial^2}{\partial \sigma_i^2 \partial b_j} L \right) &= 0 \\ \mathbb{E} \left(-\frac{\partial^2}{\partial \sigma_i^2 \partial \sigma_j^2} L \right) &= \frac{1}{2} \operatorname{tr}(\mathbf{V}^{-1} \mathbf{V}_i \mathbf{V}^{-1} \mathbf{V}_j)\end{aligned}$$

and thus has a block diagonal form

$$\begin{pmatrix} \mathbb{E}(-d_{\mathbf{b}}^2 L) & \mathbf{0}_p \\ \mathbf{0}_p^T & \mathbb{E}(-d_{\boldsymbol{\sigma}^2}^2 L) \end{pmatrix}.$$

- Setting the score (gradient) vector to zero gives the *likelihood equation* for \mathbf{b} and σ_i^2

$$\begin{aligned}\mathbf{X}^T \mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\mathbf{b}) &= \mathbf{0} \\ (\mathbf{y} - \mathbf{X}\mathbf{b})^T \mathbf{V}^{-1} \mathbf{V}_i \mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\mathbf{b}) &= \text{tr}(\mathbf{V}^{-1} \mathbf{V}_i), \quad i = 0, \dots, m.\end{aligned}$$

Root of the likelihood equation is the MLE.

- Example: One-way ANOVA with random effects

$$y_{ij} = \mu + \alpha_i + e_{ij}, \quad i = 1, \dots, a, j = 1, \dots, n_i,$$

where $\alpha_i \sim N(0, \sigma_1^2)$ and $e_{ij} \sim N(0, \sigma_0^2)$ are independent. Here

$$\mathbf{X} = \mathbf{1}_n, \quad \mathbf{V}_0 = \mathbf{I}_n, \quad \mathbf{V}_1 = \text{BlkDiag}(\sigma_1^2 \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T),$$

and

$$\begin{aligned}\mathbf{V} &= \text{BlkDiag}(\sigma_0^2 \mathbf{I}_{n_i} + \sigma_1^2 \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T), \\ \mathbf{V}^{-1} &= \text{BlkDiag}\left(\sigma_0^{-2} \mathbf{I}_{n_i} - \sigma_0^{-2} \frac{\sigma_1^2}{\sigma_0^2 + n_i \sigma_1^2} \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T\right) \\ \mathbf{V}^{-2} &= \text{BlkDiag}\left(\sigma_0^{-4} \mathbf{I}_{n_i} + \sigma_0^{-4} \frac{\sigma_1^2 (n_i - 2)}{\sigma_0^2 + n_i \sigma_1^2} \mathbf{1}_{n_i} \mathbf{1}_{n_i}^T\right).\end{aligned}$$

Setting the score vector to 0 gives the likelihood equation (check it! it might be wrong)

$$\begin{aligned}\sigma_0^{-2} \sum_i \left(n_i - \frac{\sigma_1^2 n_i^2}{\sigma_0^2 + n_i \sigma_1^2}\right) \mu &= \sigma_0^{-2} \sum_i \left(1 - \frac{\sigma_1^2 n_i}{\sigma_0^2 + n_i \sigma_1^2}\right) \sum_j y_{ij} \\ \sigma_0^{-2} \sum_i \left(n_i - \frac{\sigma_1^2 n_i}{\sigma_0^2 + n_i \sigma_1^2}\right) &= \sigma_0^{-4} \sum_i \left(\sum_j (y_{ij} - \mu)^2 + \frac{\sigma_1^2 (n_i - 2)}{\sigma_0^2 + n_i \sigma_1^2} \left(\sum_j y_{ij} - n_i \mu\right)^2\right) \\ \sigma_0^{-2} \sigma_1^2 \sum_i \left(n_i - \frac{\sigma_1^2 n_i^2}{\sigma_0^2 + n_i \sigma_1^2}\right) &= \sigma_0^{-4} \sigma_1^2 \sum_i \frac{\sigma_0^2 - n_i \sigma_1^2}{\sigma_0^2 + n_i \sigma_1^2} \left(\sum_j y_{ij} - n_i \mu\right)^2.\end{aligned}$$

In the un-balanced case, there is no analytical solution to the likelihood equation. In the balanced case, analytical solution is available (HW8).

MLE under singular \mathbf{V}

- Suppose $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \mathbf{V})$, where \mathbf{V} is singular. MLE depends on a density, which does not exist for singular \mathbf{V} \odot
- Assume $r = \text{rank}(\mathbf{V}) < n$, let $\mathbf{B} \in \mathbb{R}^{n \times r}$ such that $\mathcal{C}(\mathbf{B}) = \mathcal{C}(\mathbf{V})$. Then

$$\mathbf{B}^T \mathbf{Y} \sim N_r(\mathbf{B}^T \boldsymbol{\mu}, \mathbf{B}^T \mathbf{V} \mathbf{B}),$$

where the covariance matrix is nonsingular since

$$\begin{aligned} \text{rank}(\mathbf{B}^T \mathbf{V} \mathbf{B}) &= \text{rank}(\mathbf{B}^T \mathbf{U}_r \mathbf{D}_r \mathbf{U}_r^T \mathbf{B}) \\ &= \text{rank}(\mathbf{U}_r^T \mathbf{B}) = \text{rank}(\mathbf{U}_r^T \mathbf{U}_r \mathbf{T}) = \text{rank}(\mathbf{U}_r^T \mathbf{U}_r) = r. \end{aligned}$$

- The log-likelihood of $\mathbf{B}^T \mathbf{Y}$ is

$$\begin{aligned} & -\frac{r}{2} \ln(2\pi) - \frac{1}{2} \ln \det(\mathbf{B}^T \mathbf{V} \mathbf{B}) - \frac{1}{2} (\mathbf{B}^T \mathbf{y} - \mathbf{B}^T \boldsymbol{\mu})^T (\mathbf{B}^T \mathbf{V} \mathbf{B})^{-1} (\mathbf{B}^T \mathbf{y} - \mathbf{B}^T \boldsymbol{\mu}) \\ &= -\frac{r}{2} \ln(2\pi) - \frac{1}{2} \ln \det(\mathbf{B}^T \mathbf{V} \mathbf{B}) - \frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T \mathbf{B} (\mathbf{B}^T \mathbf{V} \mathbf{B})^{-1} \mathbf{B}^T (\mathbf{y} - \boldsymbol{\mu}). \end{aligned}$$

- We now show that maximization of the log-density does *not* depend on choice of \mathbf{B} .

Let $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{n \times r}$ be such that $\mathcal{C}(\mathbf{B}_1) = \mathcal{C}(\mathbf{B}_2) = \mathcal{C}(\mathbf{V})$ with $r = \text{rank}(\mathbf{V})$. Then

1. $\det(\mathbf{B}_1^T \mathbf{V} \mathbf{B}_1) = c \det(\mathbf{B}_2^T \mathbf{V} \mathbf{B}_2)$ for some constant c .
2. $\mathbf{B}_1 (\mathbf{B}_1^T \mathbf{V} \mathbf{B}_1)^{-1} \mathbf{B}_1^T = \mathbf{B}_2 (\mathbf{B}_2^T \mathbf{V} \mathbf{B}_2)^{-1} \mathbf{B}_2^T$.

Proof. Since $\mathcal{C}(\mathbf{B}_1) = \mathcal{C}(\mathbf{B}_2)$, $\mathbf{B}_1 = \mathbf{B}_2 \mathbf{T}$ for some nonsingular transformation matrix \mathbf{T} . Then

$$\det(\mathbf{B}_1^T \mathbf{V} \mathbf{B}_1) = \det(\mathbf{T}^T \mathbf{B}_2^T \mathbf{V} \mathbf{B}_2 \mathbf{T}) = [\det(\mathbf{T})]^2 \det(\mathbf{B}_2^T \mathbf{V} \mathbf{B}_2)$$

and

$$\mathbf{B}_1 (\mathbf{B}_1^T \mathbf{V} \mathbf{B}_1)^{-1} \mathbf{B}_1^T = \mathbf{B}_2 \mathbf{T} (\mathbf{T}^T \mathbf{B}_2^T \mathbf{V} \mathbf{B}_2 \mathbf{T})^{-1} \mathbf{T}^T \mathbf{B}_2^T = \mathbf{B}_2 (\mathbf{B}_2^T \mathbf{V} \mathbf{B}_2)^{-1} \mathbf{B}_2^T.$$

□

22 Lecture 22: Nov 25

Announcement

- HW8 due Mon Dec 2.
- Email your comments (compliments?) on TA to me.

Last time

- Variance component estimation: MLE.

Today

- Variance component estimation: REML.
- Variance component estimation: method of moment.
- Variance component testing: exact F tests.

Variance component estimation: REML (JM 8.4.2)

- Again we consider estimating the variance components $(\sigma_0^2, \dots, \sigma_m^2)$ of the model

$$\mathbf{y} \sim N_n(\mathbf{X}\mathbf{b}, \mathbf{V}),$$

where

$$\mathbf{V} = \sum_{i=0}^m \sigma_i^2 \mathbf{V}_i,$$

with \mathbf{V}_i psd and \mathbf{V} nonsingular. Parameters are \mathbf{b} and variance components $(\sigma_0^2, \dots, \sigma_m^2)$.

- *Restricted (or residual) maximum likelihood* (REML) estimation involves finding the MLE of variance components from the distribution of residuals. This allows for estimation of the variance components without complication of the fixed effects.

- It is not clear how to define residuals since \mathbf{V} is unknown.

Let \mathbf{M} be any projection (not necessarily orthogonal) onto $\mathcal{C}(\mathbf{X})$. Then residuals can be defined as

$$(\mathbf{I} - \mathbf{M})\mathbf{y}.$$

The distribution of the residuals is a singular normal

$$(\mathbf{I} - \mathbf{M})\mathbf{Y} \sim N(\mathbf{0}, (\mathbf{I} - \mathbf{M})\mathbf{V}(\mathbf{I} - \mathbf{M})^T)$$

and, due to nonsingularity of \mathbf{V} ,

$$\mathcal{C}((\mathbf{I} - \mathbf{M})\mathbf{V}(\mathbf{I} - \mathbf{M})^T) = \mathcal{C}(\mathbf{I} - \mathbf{M})$$

Suppose $s = \text{rank}(\mathbf{X}) = \text{rank}(\mathbf{M})$, then

$$\text{rank}(\mathbf{I} - \mathbf{M}) = n - s.$$

- We use the result in last section to do MLE on this singular normal distribution. Let $\mathbf{B} \in \mathbb{R}^{n \times (n-s)}$ such that $\mathcal{C}(\mathbf{B}) = \mathcal{C}(\mathbf{I} - \mathbf{M})$. Then

$$\mathbf{B}^T(\mathbf{I} - \mathbf{M})\mathbf{Y} \sim N_{n-s}(\mathbf{0}, \mathbf{B}^T(\mathbf{I} - \mathbf{M})\mathbf{V}(\mathbf{I} - \mathbf{M})^T\mathbf{B})$$

and the MLE of variance components maximizes the log-likelihood

$$\begin{aligned} & -\frac{n-s}{2} \ln(2\pi) - \frac{1}{2} \ln \det(\mathbf{B}^T(\mathbf{I} - \mathbf{M})\mathbf{V}(\mathbf{I} - \mathbf{M})^T\mathbf{B}) \\ & - \frac{1}{2} \mathbf{y}^T(\mathbf{I} - \mathbf{M})^T\mathbf{B}[\mathbf{B}^T(\mathbf{I} - \mathbf{M})\mathbf{V}(\mathbf{I} - \mathbf{M})^T\mathbf{B}]^{-1}\mathbf{B}^T(\mathbf{I} - \mathbf{M})\mathbf{y}. \end{aligned}$$

In last section we saw that maximization of this log-likelihood does not depend on choice of \mathbf{B} . Does this depend on choice of projection matrix \mathbf{M} ?

- $\mathcal{C}((\mathbf{I} - \mathbf{M})^T\mathbf{B}) = \mathcal{N}(\mathbf{X}^T)$.

Proof. Clearly $\mathbf{B}^T(\mathbf{I} - \mathbf{M})\mathbf{X} = \mathbf{B}^T(\mathbf{X} - \mathbf{X}) = \mathbf{0}_{(n-s) \times p}$. So $\mathcal{C}((\mathbf{I} - \mathbf{M})^T\mathbf{B}) \subset \mathcal{N}(\mathbf{X}^T)$. It is enough to show $\text{rank}((\mathbf{I} - \mathbf{M})^T\mathbf{B}) = n - s = \dim(\mathcal{N}(\mathbf{X}^T))$. Since $\mathcal{C}(\mathbf{B}) = \mathcal{C}(\mathbf{I} - \mathbf{M})$ and $\mathbf{I} - \mathbf{M}$ is idempotent, $(\mathbf{I} - \mathbf{M})\mathbf{B} = \mathbf{B}$. Thus

$$\text{rank}((\mathbf{I} - \mathbf{M})^T\mathbf{B}) = \text{rank}(\mathbf{B}^T(\mathbf{I} - \mathbf{M})) = \text{rank}(\mathbf{B}^T) = n - s.$$

□

Remark: This result shows that for any projection \mathbf{M} and \mathbf{B} ,

$$(\mathbf{I} - \mathbf{M})^T \mathbf{B} = (\mathbf{I} - \mathbf{P}_X) \mathbf{T}$$

for some transformation matrix $\mathbf{T} \in \mathbb{R}^{n \times (n-s)}$ with full column rank and $\mathcal{C}(\mathbf{T}) = \mathcal{N}(\mathbf{X}^T)$ (otherwise contradicting with $\text{rank}((\mathbf{I} - \mathbf{M})^T \mathbf{B}) = n - s$). Therefore the MLE estimation of variance components does not depend on choice of either \mathbf{M} or \mathbf{B} .

- It turns out we even do *not* have to make this choice! The entire REML procedure can be carried out using the original data.
- An alternative definition of REML is to do MLE from

$$\mathbf{B}^T \mathbf{Y} \sim N_{n-s}(\mathbf{0}, \mathbf{B}^T \mathbf{V} \mathbf{B}),$$

where $\mathbf{B} \in \mathbb{R}^{n \times (n-s)}$ is any basis of $\mathcal{N}(\mathbf{X}^T)$. That is $\mathbf{B}^T \mathbf{X} = \mathbf{0}_{(n-s) \times p}$ and $\text{rank}(\mathbf{B}) = n - s$.

Essentially we find a basis \mathbf{B} of $\mathcal{N}(\mathbf{X}^T)$ and estimate the variance components from $\mathbf{B}^T \mathbf{Y}$. This is equivalent to the above procedure of choosing a projection \mathbf{M} and then doing MLE from residuals (why?). Hence the name *restricted maximum likelihood* (REML).

- Thus we are doing MLE of variance components from the log-likelihood

$$-\frac{n-s}{2} \ln(2\pi) - \frac{1}{2} \ln \det(\mathbf{B}^T \mathbf{V} \mathbf{B}) - \frac{1}{2} \mathbf{y}^T \mathbf{B} [\mathbf{B}^T \mathbf{V} \mathbf{B}]^{-1} \mathbf{B}^T \mathbf{y}.$$

Setting gradient to zero shows REML of the variance components has to satisfy the likelihood equation

$$\begin{aligned} & \mathbf{y}^T \mathbf{B} (\mathbf{B}^T \mathbf{V} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{V}_i \mathbf{B} (\mathbf{B}^T \mathbf{V} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{y} \\ &= \text{tr}((\mathbf{B}^T \mathbf{V} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{V}_i \mathbf{B}) \\ &= \text{tr}(\mathbf{B} (\mathbf{B}^T \mathbf{V} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{V}_i), \quad i = 0, \dots, m. \end{aligned}$$

It turns out the likelihood equation is independent of the choice of basis \mathbf{B} .

- $\mathbf{B} (\mathbf{B}^T \mathbf{V} \mathbf{B})^{-1} \mathbf{B}^T = \mathbf{V}^{-1} - \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}$.

Proof. Since \mathbf{V} is pd, the symmetric square root $\mathbf{V}^{1/2}$ is pd too. It suffices to show

$$\mathbf{V}^{1/2} \mathbf{B} (\mathbf{B}^T \mathbf{V} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{V}^{1/2} = \mathbf{I} - \mathbf{V}^{-1/2} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1/2}.$$

The left side is the orthogonal projection onto $\mathcal{C}(\mathbf{V}^{1/2} \mathbf{B}) = \mathcal{C}(\mathbf{B}) = \mathcal{N}(\mathbf{X}^T)$. The right side is the orthogonal projection onto $\mathcal{N}(\mathbf{X}^T \mathbf{V}^{-1/2}) = \mathcal{N}(\mathbf{X}^T)$. Since orthogonal projection to a vector space is unique, proof is done. \square

- In summary, let

$$\mathbf{A} := \mathbf{V}^{-1} - \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1}.$$

Then any solution $(\hat{\sigma}_0^2, \dots, \hat{\sigma}_m^2)$ to the (restricted) likelihood equation

$$\mathbf{y}^T \mathbf{A} \mathbf{V}_i \mathbf{A} \mathbf{y} = \text{tr}(\mathbf{A} \mathbf{V}_i), \quad i = 0, 1, \dots, m.$$

is the REML.

- Example: REML for balanced one-way ANOVA with random effects (HW8).

Variance component estimation: method of moment (JM 8.4.3)

Also called the ANOVA method (JM 8.4.3).

- Again we consider estimating the variance components $(\sigma_0^2, \dots, \sigma_m^2)$ of the model

$$\mathbf{y} \sim N_n(\mathbf{X} \mathbf{b}, \mathbf{V}),$$

where

$$\mathbf{V} = \sum_{i=0}^m \sigma_i^2 \mathbf{V}_i,$$

with \mathbf{V}_i psd and \mathbf{V} nonsingular.

- Recall

$$\begin{aligned} \mathbb{E}(\mathbf{y}^T \mathbf{A} \mathbf{y}) &= \text{tr}(\mathbf{A} \mathbf{V}) + \mathbf{b}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{b} \\ &= \sum_{i=0}^m \sigma_i^2 \text{tr}(\mathbf{A} \mathbf{V}_i) + \mathbf{b}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{b}. \end{aligned}$$

- Choose \mathbf{A}_i , $i = 0, \dots, m$, such that $\mathbf{X}^T \mathbf{A}_i \mathbf{X} = \mathbf{0}$. This leads to a system of equations by setting $E(\mathbf{y}^T \mathbf{A}_i \mathbf{y})$ to the observed $\mathbf{y}^T \mathbf{A}_i \mathbf{y}$

$$\begin{pmatrix} \text{tr}(\mathbf{A}_0 \mathbf{V}_0) & \text{tr}(\mathbf{A}_0 \mathbf{V}_1) & \cdots & \text{tr}(\mathbf{A}_0 \mathbf{V}_m) \\ \vdots & & & \vdots \\ \text{tr}(\mathbf{A}_m \mathbf{V}_0) & \text{tr}(\mathbf{A}_m \mathbf{V}_1) & \cdots & \text{tr}(\mathbf{A}_m \mathbf{V}_m) \end{pmatrix} \begin{pmatrix} \sigma_0^2 \\ \vdots \\ \sigma_m^2 \end{pmatrix} = \begin{pmatrix} \mathbf{y}^T \mathbf{A}_0 \mathbf{y} \\ \vdots \\ \mathbf{y}^T \mathbf{A}_m \mathbf{y} \end{pmatrix},$$

which can be solved for the variance components.

- Example: MoM for *unbalanced* one-way ANOVA with random effects (HW8). Choose $\mathbf{A}_0 = \mathbf{P}_Z - \mathbf{P}_1$ and $\mathbf{A}_1 = \mathbf{I} - \mathbf{P}_Z$.

Variance components testing: Wald's exact F test

Proposed by Seely and El Bassiouni (1983).

- Consider the mixed linear model:

$$\mathbf{Y} = \mathbf{X}\mathbf{b} + \mathbf{Z}_1\mathbf{u}_1 + \mathbf{Z}_2\mathbf{u}_2 + \mathbf{e},$$

where $\mathbf{Z}_1 \in \mathbb{R}^{n \times q_1}$ and $\mathbf{Z}_2 \in \mathbb{R}^{n \times q_2}$. \mathbf{u}_1 , \mathbf{u}_2 , and \mathbf{e} are independent random vectors with

$$\mathbf{u}_1 \sim N(\mathbf{0}_{q_1}, \mathbf{R}), \quad \mathbf{u}_2 \sim N(\mathbf{0}_{q_2}, \sigma_2^2 \mathbf{I}_{q_2}), \quad \mathbf{e} \sim N(\mathbf{0}_n, \sigma_0^2 \mathbf{I}_n).$$

Equivalently, we have the variance component model

$$\mathbf{Y} \sim N_n(\mathbf{X}\mathbf{b}, \mathbf{V}),$$

where

$$\mathbf{V} = \sigma_0^2 \mathbf{I}_n + \mathbf{Z}_1 \mathbf{R} \mathbf{Z}_1^T + \sigma_2^2 \mathbf{Z}_2 \mathbf{Z}_2^T.$$

- We are interested testing $H_0 : \sigma_2^2 = 0$ vs $H_A : \sigma_2^2 > 0$.
- Idea of the Wald test is to treat \mathbf{u}_1 and \mathbf{u}_2 as fixed effects and use the change in SSE as the test statistic.

- Let $\mathbf{P}_1 = \mathbf{P}_{\mathbf{X}, \mathbf{Z}_1}$ be the orthogonal projection onto $\mathcal{C}(\mathbf{X}, \mathbf{Z}_1)$.
Let $\mathbf{P}_2 = \mathbf{P}_{\mathbf{X}, \mathbf{Z}_1, \mathbf{Z}_2}$ be the orthogonal projection onto $\mathcal{C}(\mathbf{X}, \mathbf{Z}_1, \mathbf{Z}_2)$.

Then

$$\text{SSE}_1 = \mathbf{y}^T (\mathbf{I} - \mathbf{P}_1) \mathbf{y}$$

$$\text{SSE}_2 = \mathbf{y}^T (\mathbf{I} - \mathbf{P}_2) \mathbf{y}$$

and

$$\text{SSE}_1 - \text{SSE}_2 = \mathbf{y}^T (\mathbf{P}_2 - \mathbf{P}_1) \mathbf{y}.$$

Ideally we would like to conduct a F test based on the ratio

$$\frac{\text{SSE}_1 - \text{SSE}_2}{\text{SSE}_2}.$$

By previous result (p90 of notes), it is enough to show the following facts.

1. $\sigma_0^{-2} (\mathbf{I} - \mathbf{P}_2) \mathbf{V}$ is idempotent with $\text{rank}((\mathbf{I} - \mathbf{P}_2) \mathbf{V}) = \text{rank}(\mathbf{I} - \mathbf{P}_2)$.

Proof. Just note that

$$\begin{aligned} & (\mathbf{I} - \mathbf{P}_2) \mathbf{V} (\mathbf{I} - \mathbf{P}_2) \\ &= \mathbf{V} - \mathbf{V} \mathbf{P}_2 - \mathbf{P}_2 \mathbf{V} + \mathbf{P}_2 \mathbf{V} \mathbf{P}_2 \\ &= \mathbf{V} - (\sigma_0^2 \mathbf{P}_2 + \mathbf{Z}_1 \mathbf{R} \mathbf{Z}_1^T + \sigma_2^2 \mathbf{Z}_2 \mathbf{Z}_2^T) \\ &= \sigma_0^2 (\mathbf{I} - \mathbf{P}_2). \end{aligned}$$

□

2. $\sigma_0^{-2} (\mathbf{P}_2 - \mathbf{P}_1) \mathbf{V}$ is idempotent under the null hypothesis $\sigma_2^2 = 0$ with $\text{rank}(\mathbf{P}_2 - \mathbf{P}_1) \mathbf{V} = \text{rank}(\mathbf{P}_2 - \mathbf{P}_1)$.

Proof. Just note that, when $\sigma_2^2 = 0$,

$$\begin{aligned} & (\mathbf{P}_2 - \mathbf{P}_1) \mathbf{V} (\mathbf{P}_2 - \mathbf{P}_1) \\ &= \mathbf{P}_2 \mathbf{V} \mathbf{P}_2 - \mathbf{P}_2 \mathbf{V} \mathbf{P}_1 - \mathbf{P}_1 \mathbf{V} \mathbf{P}_2 + \mathbf{P}_1 \mathbf{V} \mathbf{P}_1 \\ &= (\sigma_0^2 \mathbf{P}_2 + \mathbf{Z}_1 \mathbf{R} \mathbf{Z}_1^T) - (\sigma_0^2 \mathbf{P}_1 + \mathbf{Z}_1 \mathbf{R} \mathbf{Z}_1^T) \\ &\quad - (\sigma_0^2 \mathbf{P}_1 + \mathbf{Z}_1 \mathbf{R} \mathbf{Z}_1^T) + (\sigma_0^2 \mathbf{P}_1 + \mathbf{Z}_1 \mathbf{R} \mathbf{Z}_1^T) \\ &= \sigma_0^2 (\mathbf{P}_2 - \mathbf{P}_1). \end{aligned}$$

□

3. $(\mathbf{P}_2 - \mathbf{P}_1)\mathbf{V}(\mathbf{I} - \mathbf{P}_1) = \mathbf{0}$.

Proof. Just note that

$$\begin{aligned}
& (\mathbf{P}_2 - \mathbf{P}_1)\mathbf{V}(\mathbf{I} - \mathbf{P}_1) \\
&= \mathbf{P}_2\mathbf{V} - \mathbf{P}_2\mathbf{V}\mathbf{P}_2 - \mathbf{P}_1\mathbf{V} + \mathbf{P}_1\mathbf{V}\mathbf{P}_2 \\
&= (\sigma_0^2\mathbf{P}_2 + \mathbf{Z}_1\mathbf{R}\mathbf{Z}_1^T + \sigma_2^2\mathbf{Z}_2\mathbf{Z}_2^T) - (\sigma_0^2\mathbf{P}_2 + \mathbf{Z}_1\mathbf{R}\mathbf{Z}_1^T + \sigma_2^2\mathbf{Z}_2\mathbf{Z}_2^T) \\
&\quad - (\sigma_0^2\mathbf{P}_1 + \mathbf{Z}_1\mathbf{R}\mathbf{Z}_1^T + \sigma_2^2\mathbf{P}_1\mathbf{Z}_2\mathbf{Z}_2^T) + (\sigma_0^2\mathbf{P}_1 + \mathbf{Z}_1\mathbf{R}\mathbf{Z}_1^T + \sigma_2^2\mathbf{P}_1\mathbf{Z}_2\mathbf{Z}_2^T) \\
&= \mathbf{0}.
\end{aligned}$$

□

- (Wald's F test) The test reject $H_0 : \sigma_2^2 = 0$ when

$$\frac{(\text{SSE}_1 - \text{SSE}_2)/\text{rank}(\mathbf{P}_2 - \mathbf{P}_1)}{\text{SSE}_1/\text{rank}(\mathbf{I} - \mathbf{P}_2)} > F_{\text{rank}(\mathbf{P}_2 - \mathbf{P}_1), \text{rank}(\mathbf{I} - \mathbf{P}_2), \alpha}.$$

Remark: When $\mathcal{C}(\mathbf{X}, \mathbf{Z}_1, \mathbf{Z}_2) = \mathcal{C}(\mathbf{X}, \mathbf{Z}_1)$, the test cannot be performed since the numerator is 0/0.

Remark: Under the alternative hypothesis $\sigma_2^2 > 0$, the numerator is not a χ^2 random variable anymore but can be numerically evaluated by a method of Davies (1980).

- Extension of Wald's F test to the case

$$\mathbf{V} = \sigma_0^2\mathbf{\Sigma} + \mathbf{Z}_1\mathbf{R}\mathbf{Z}_1^T + \sigma_2^2\mathbf{Z}_2\mathbf{Z}_2^T, \tag{4}$$

where $\mathbf{\Sigma} \in \mathbb{R}^{n \times n}$ is a known psd matrix. Let $r = \text{rank}(\mathbf{\Sigma})$.

Given eigen-decomposition $\mathbf{\Sigma} = \mathbf{U}_r\mathbf{D}_r\mathbf{U}_r^T$, define $\mathbf{T} = \mathbf{D}_r^{-1/2}\mathbf{U}_r^T$. Thus

$$\mathbf{T}\mathbf{Y} \sim N(\mathbf{T}\mathbf{X}\mathbf{b}, \sigma_0^2\mathbf{I}_r + (\mathbf{T}\mathbf{Z}_1)\mathbf{R}(\mathbf{T}\mathbf{Z}_1)^T + \sigma_2^2(\mathbf{T}\mathbf{Z}_2)(\mathbf{T}\mathbf{Z}_2)^T).$$

Then Wald's F test can be applied to $\mathbf{T}\mathbf{Y}$.

Variance components testing: Öfversten's exact F test

Proposed by Öfversten (1993) and Christensen (1996).

- Consider the mixed linear model:

$$\mathbf{Y} = \mathbf{X}\mathbf{b} + \mathbf{Z}_1\mathbf{u}_1 + \mathbf{Z}_2\mathbf{u}_2 + \mathbf{e},$$

where $\mathbf{Z}_1 \in \mathbb{R}^{n \times q_1}$ and $\mathbf{Z}_2 \in \mathbb{R}^{n \times q_2}$. \mathbf{u}_1 , \mathbf{u}_2 , and \mathbf{e} are independent random vectors with

$$\mathbf{u}_1 \sim N(\mathbf{0}_{q_1}, \sigma_1^2 \mathbf{I}_{q_1}), \quad \mathbf{u}_2 \sim N(\mathbf{0}_{q_2}, \sigma_2^2 \mathbf{I}_{q_2}), \quad \mathbf{e} \sim N(\mathbf{0}_n, \sigma_0^2 \mathbf{I}_n).$$

Equivalently, we have the variance component model

$$\mathbf{Y} \sim N_n(\mathbf{X}\mathbf{b}, \mathbf{V}),$$

where

$$\mathbf{V} = \sigma_0^2 \mathbf{I}_n + \sigma_1^2 \mathbf{Z}_1 \mathbf{Z}_1^T + \sigma_2^2 \mathbf{Z}_2 \mathbf{Z}_2^T.$$

- We are interested in testing $\mathbf{H}_0 : \sigma_1^2 = 0$ vs $H_A : \sigma_1^2 > 0$.
- The idea is to massage the problem to the case (4) considered above. First perform QR (Gram-Schmidt) on the matrix $\begin{pmatrix} \mathbf{X} & \mathbf{Z}_2 & \mathbf{Z}_1 & \mathbf{I}_n \end{pmatrix}$ to obtain an orthonormal basis $\begin{pmatrix} \mathbf{C}_1 & \mathbf{C}_2 & \mathbf{C}_3 & \mathbf{C}_4 \end{pmatrix}$ of \mathbb{R}^n , where
 - \mathbf{C}_1 is an orthonormal basis of $\mathcal{C}(\mathbf{X})$
 - \mathbf{C}_2 is an orthonormal basis of $\mathcal{C}(\mathbf{X}, \mathbf{Z}_2) - \mathcal{C}(\mathbf{X})$
 - \mathbf{C}_3 is an orthonormal basis of $\mathcal{C}(\mathbf{X}, \mathbf{Z}_2, \mathbf{Z}_1) - \mathcal{C}(\mathbf{X}, \mathbf{Z}_2)$
 - \mathbf{C}_4 is an orthonormal basis of $\mathcal{C}(\mathbf{X}, \mathbf{Z}_2, \mathbf{Z}_1)^\perp$.

Then we choose λ and matrix \mathbf{K} such that

$$\mathbf{C}_2^T \mathbf{Y} + \mathbf{K} \mathbf{C}_4^T \mathbf{Y} \sim N(\mathbf{0}, (\sigma_2^2 + \sigma_0^2/\lambda) \mathbf{C}_2^T \mathbf{Z}_2 \mathbf{Z}_2^T \mathbf{C}_2 + \sigma_1^2 \mathbf{C}_2^T \mathbf{Z}_1 \mathbf{Z}_1^T \mathbf{C}_2),$$

which is of form (4).

- If $\mathbf{C}_2^T \mathbf{Z}_1 = \mathbf{0}$, e.g., when $\mathcal{C}(\mathbf{Z}_1) \subset \mathcal{C}(\mathbf{X})$, then this test cannot be performed.
- If $\mathbf{C}_2^T \mathbf{Z}_2 \mathbf{Z}_2^T \mathbf{C}_2 = \lambda \mathbf{I}$, note

$$\begin{aligned} \mathbf{C}_2^T \mathbf{Y} &\sim N(\mathbf{0}, \sigma_0^2 \mathbf{I} + \sigma_2^2 \mathbf{C}_2^T \mathbf{Z}_2 \mathbf{Z}_2^T \mathbf{C}_2 + \sigma_1^2 \mathbf{C}_2^T \mathbf{Z}_1 \mathbf{Z}_1^T \mathbf{C}_2) \\ &= N(\mathbf{0}, (\sigma_0^2 + \lambda \sigma_2^2) \mathbf{I} + \sigma_1^2 \mathbf{C}_2^T \mathbf{Z}_1 \mathbf{Z}_1^T \mathbf{C}_2). \end{aligned}$$

Then the ordinary Wald's F test can be applied without using the $\mathbf{K} \mathbf{C}_4^T \mathbf{Y}$ piece as long as $\mathbf{C}_2^T \mathbf{Z}_1 \neq \mathbf{0}$.

- In general, $\mathbf{C}_2^T \mathbf{Z}_2 \mathbf{Z}_2^T \mathbf{C}_2 \neq \lambda \mathbf{I}$, then the test requires $\mathbf{K} \mathbf{C}_4^T \mathbf{Y}$. Note that

$$\mathbf{C}_2^T \mathbf{V} \mathbf{C}_4 = \mathbf{C}_2^T (\sigma_0^2 \mathbf{I}_n + \sigma_1^2 \mathbf{Z}_1 \mathbf{Z}_1^T + \sigma_2^2 \mathbf{Z}_2 \mathbf{Z}_2^T) \mathbf{C}_4 = \mathbf{0}.$$

Therefore $\mathbf{C}_2^T \mathbf{Y}$ is independent of $\mathbf{K} \mathbf{C}_4^T \mathbf{Y}$. We simply pick \mathbf{K} such that

$$\mathbf{K} \mathbf{C}_4^T \mathbf{Y} \sim N(\mathbf{0}, \sigma_0^2 (\lambda^{-1} \mathbf{C}_2^T \mathbf{Z}_2 \mathbf{Z}_2^T \mathbf{C}_2 - \mathbf{I})).$$

That is

$$\mathbf{K} \mathbf{K}^T = \lambda^{-1} \mathbf{C}_2^T \mathbf{Z}_2 \mathbf{Z}_2^T \mathbf{C}_2 - \mathbf{I}.$$

Apparently we need to choose λ such that $\lambda^{-1} \mathbf{C}_2^T \mathbf{Z}_2 \mathbf{Z}_2^T \mathbf{C}_2 - \mathbf{I}$ is psd. Let $\mathbf{C}_2^T \mathbf{Z}_2 \mathbf{Z}_2^T \mathbf{C}_2 = \mathbf{W} \mathbf{\Lambda} \mathbf{W}^T = \mathbf{W} \text{diag}(\lambda_i) \mathbf{W}^T$ be the eigendecomposition. Setting λ to be the smallest eigenvalue yields

$$\mathbf{K} = \mathbf{W} \text{diag}(\sqrt{\lambda_i / \lambda - 1}).$$

23 Lecture 23: Dec 2

Announcement

- HW8 due this Wed (?)
- Course evaluation: <https://classeval.ncsu.edu/>

Last time

- Variance component estimation: REML.
- Variance component estimation: method of moment/ANOVA method.
- Variance component testing: exact F tests.

Today

- Variance component testing: eLRT, eRLRT, score.

Exact LRT and RLRT with one variance component

This is proposed by Crainiceanu and Ruppert (2004). Slight notation change below.

- Consider the variance component model

$$\mathbf{Y} \sim N_n(\mathbf{X}\mathbf{b}, \mathbf{V}),$$

where

$$\mathbf{V} = \sigma_0^2 \mathbf{I}_n + \sigma_1^2 \mathbf{V}_1.$$

Let $\lambda = \sigma_1^2/\sigma_0^2$ be the signal-to-noise ratio, and write the covariance as

$$\mathbf{V} = \sigma_0^2(\mathbf{I}_n + \lambda \mathbf{V}_1) = \sigma_0^2 \mathbf{V}_\lambda.$$

The model parameters are $(\mathbf{b}, \sigma_0^2, \lambda)$. Denote $s = \text{rank}(\mathbf{X})$.

- Testing $H_0 : \sigma_1^2 = 0$ vs $H_A : \sigma_1^2 > 0$ is equivalent to testing $H_0 : \lambda = 0$ vs $H_A : \lambda > 0$.

- The log-likelihood function is

$$L(\mathbf{b}, \sigma_0^2, \lambda) = -\frac{n}{2} \ln \sigma_0^2 - \frac{1}{2} \ln \det(\mathbf{V}_\lambda) - \frac{1}{2\sigma_0^2} (\mathbf{y} - \mathbf{X}\mathbf{b})^T \mathbf{V}_\lambda^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b}).$$

- The likelihood ratio test (LRT) statistic is

$$\text{LRT} = 2 \sup_{H_A} L(\mathbf{b}, \sigma_0^2, \lambda) - 2 \sup_{H_0} L(\mathbf{b}, \sigma_0^2, \lambda).$$

- Under the null $\lambda = 0$, it is a regular linear model with MLE

$$\begin{aligned} \hat{\mathbf{b}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ \hat{\sigma}_0^2 &= \frac{\text{SSE}}{n} = \frac{\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}}{n} =: \frac{\mathbf{y}^T \mathbf{A}_0 \mathbf{y}}{n}. \end{aligned}$$

Thus

$$2 \sup_{H_0} L(\mathbf{b}, \sigma_0^2, \lambda) = -n \ln \mathbf{y}^T \mathbf{A}_0 \mathbf{y} + n \ln n - n.$$

- Under the alternative, for fixed $\lambda > 0$, the profile likelihood maximizers are

$$\begin{aligned} \hat{\mathbf{b}}(\lambda) &= (\mathbf{X}^T \mathbf{V}_\lambda^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_\lambda^{-1} \mathbf{y} \\ \hat{\sigma}_0^2(\lambda) &= \frac{\mathbf{y}^T [\mathbf{V}_\lambda^{-1} - \mathbf{V}_\lambda^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}_\lambda^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_\lambda^{-1}] \mathbf{y}}{n} =: \frac{\mathbf{y}^T \mathbf{A}_\lambda \mathbf{y}}{n}. \end{aligned}$$

Thus

$$2 \sup_{H_A} L(\mathbf{b}, \sigma_0^2, \lambda) = \sup_{\lambda \geq 0} -n \ln \mathbf{y}^T \mathbf{A}_\lambda \mathbf{y} - \ln \det(\mathbf{V}_\lambda) + n \ln n - n.$$

- Therefore the LRT statistic is

$$\text{LRT} = \sup_{\lambda \geq 0} n \ln \mathbf{y}^T \mathbf{A}_0 \mathbf{y} - n \ln \mathbf{y}^T \mathbf{A}_\lambda \mathbf{y} - \ln \det(\mathbf{V}_\lambda),$$

where

$$\begin{aligned} \mathbf{A}_0 &= \mathbf{I} - \mathbf{P}_X \\ \mathbf{A}_\lambda &= \mathbf{V}_\lambda^{-1} - \mathbf{V}_\lambda^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}_\lambda^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_\lambda^{-1}. \end{aligned}$$

- First obtain an orthonormal basis $(\mathbf{Q}_0 \ \mathbf{Q}_1 \ \mathbf{Q}_2)$ of \mathbb{R}^n such that

- $\mathbf{Q}_0 \in \mathbb{R}^{n \times s}$ is an orthonormal basis of $\mathcal{C}(\mathbf{X})$,
- $\mathbf{Q}_1 \in \mathbb{R}^{n \times k}$ is an orthonormal basis from the eigendecomposition

$$\mathbf{A}_0 \mathbf{V}_1 \mathbf{A}_0 = \mathbf{Q}_1 \text{diag}(\mu_1, \dots, \mu_k) \mathbf{Q}_1^T,$$

where $k = \text{rank}(\mathbf{A}_0 \mathbf{V}_1 \mathbf{A}_0)$, and

- $\mathbf{Q}_2 \in \mathbb{R}^{n \times (n-s-k)}$ is an orthonormal basis of $\mathcal{C}(\mathbf{Q}_0, \mathbf{Q}_1)^\perp = \mathcal{C}(\mathbf{X}, \mathbf{Q}_1)^\perp$.

Define $\mathbf{Q} = (\mathbf{Q}_1, \mathbf{Q}_2) \in \mathbb{R}^{n \times (n-s)}$, an orthonormal basis of $\mathcal{C}(\mathbf{X})^\perp = \mathcal{N}(\mathbf{X}^T)$.

Then the following results hold.

1. $\mathbf{A}_0 = \mathbf{Q} \mathbf{Q}^T$.

Proof. Both sides are the orthogonal projection onto $\mathcal{N}(\mathbf{X}^T)$. □

2. $\mathbf{Q}^T \mathbf{V}_\lambda \mathbf{Q} = \text{diag}(1 + \lambda\mu_1, \dots, 1 + \lambda\mu_k, 1, \dots, 1)$.

Proof. Note $\mathbf{A}_0 \mathbf{Q} = (\mathbf{I} - \mathbf{P}_X) \mathbf{Q} = \mathbf{Q}$. Then

$$\begin{aligned} & \mathbf{Q}^T \mathbf{V}_\lambda \mathbf{Q} \\ &= \mathbf{Q}^T \mathbf{A}_0 \mathbf{V}_\lambda \mathbf{A}_0 \mathbf{Q} \\ &= \mathbf{Q}^T \mathbf{A}_0 (\mathbf{I}_n + \lambda \mathbf{V}_1) \mathbf{A}_0 \mathbf{Q} \\ &= \mathbf{Q}^T \mathbf{A}_0 \mathbf{Q} + \lambda \mathbf{Q}^T \mathbf{Q}_1 \text{diag}(\mu_1, \dots, \mu_k) \mathbf{Q}_1^T \mathbf{Q} \\ &= \mathbf{Q}^T \mathbf{Q} + \lambda \begin{pmatrix} \mathbf{I}_k \\ \mathbf{0} \end{pmatrix} \text{diag}(\mu_1, \dots, \mu_k) \begin{pmatrix} \mathbf{I}_k & \mathbf{0} \end{pmatrix} \\ &= \mathbf{I}_{n-s} + \lambda \text{diag}(\mu_1, \dots, \mu_k, 0, \dots, 0) \\ &= \text{diag}(1 + \mu_1, \dots, 1 + \mu_k, 1, \dots, 1). \end{aligned}$$

□

3. $\mathbf{A}_\lambda = \mathbf{Q} \text{diag}((1 + \mu_1)^{-1}, \dots, (1 + \mu_k)^{-1}, 1, \dots, 1) \mathbf{Q}^T$.

Proof. Since \mathbf{Q} form a basis of $\mathcal{N}(\mathbf{X}^T)$, by a previous result (p143),

$$\mathbf{A}_\lambda = \mathbf{V}_\lambda^{-1} - \mathbf{V}_\lambda^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{V}_\lambda^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_\lambda^{-1} = \mathbf{Q} (\mathbf{Q}^T \mathbf{V}_\lambda \mathbf{Q})^{-1} \mathbf{Q}^T.$$

Then substitute the result 2 for $\mathbf{Q}^T \mathbf{V}_\lambda \mathbf{Q}$. □

4. Under the alternative $\lambda > 0$,

$$\sigma_0^{-1} \mathbf{Q}^T \mathbf{y} \sim N(\mathbf{0}_{n-s}, \text{diag}(1 + \lambda\mu_1, \dots, 1 + \lambda\mu_k, 1, \dots, 1)).$$

Under the null $\lambda = 0$,

$$\sigma_0^{-1} \mathbf{Q}^T \mathbf{y} \sim N(\mathbf{0}_{n-s}, \mathbf{I}_{n-s}).$$

Proof. Since $\mathcal{C}(\mathbf{Q}) = \mathcal{N}(\mathbf{X}^T)$, $\mathbf{Q}^T \mathbf{X} \mathbf{b} = \mathbf{0}_{n-s}$. Covariance follows from result 2. \square

5. Let the eigenvalues of \mathbf{V}_1 be (ξ_1, \dots, ξ_ℓ) , where $\ell = \text{rank}(\mathbf{V}_1)$. Then

$$\ln \det(\mathbf{V}_\lambda) = \sum_{i=1}^{\ell} \ln(1 + \lambda\xi_i).$$

Proof. This is trivial. \square

- Putting things together, under the null,

$$\begin{aligned} \text{LRT} &= \sup_{\lambda \geq 0} n \ln \mathbf{y}^T \mathbf{A}_0 \mathbf{y} - n \ln \mathbf{y}^T \mathbf{A}_\lambda \mathbf{y} - \ln \det(\mathbf{V}_\lambda) \\ &= \sup_{\lambda \geq 0} n \ln \mathbf{y}^T \mathbf{Q} \mathbf{Q}^T \mathbf{y} \\ &\quad - n \ln \mathbf{y}^T \mathbf{Q} \text{diag}((1 + \lambda\mu_1)^{-1}, \dots, (1 + \lambda\mu_k)^{-1}, 1, \dots, 1) \mathbf{Q}^T \mathbf{y} \\ &\quad - \ln \det(\mathbf{V}_\lambda) \\ &\stackrel{\mathcal{D}}{=} \sup_{\lambda \geq 0} n \ln \frac{\sum_{i=1}^{n-s} w_i^2}{\sum_{i=1}^k \frac{w_i^2}{1 + \lambda\mu_i} + \sum_{i=k+1}^{n-s} w_i^2} - \sum_{i=1}^{\ell} \ln(1 + \lambda\xi_i), \end{aligned}$$

where w_i are $(n - s)$ independent standard normals.

- The null distribution can be obtained from computer simulation:

Given \mathbf{y} , \mathbf{X} , \mathbf{V}_1 , and simulation replicates B

Eigen-decomposition: $\mathbf{V}_1 = \mathbf{U}_\ell \text{diag}(\xi_1, \dots, \xi_\ell) \mathbf{U}_\ell^T$

Regress $\mathbf{U}_\ell \text{diag}(\sqrt{\xi_1}, \dots, \sqrt{\xi_\ell})$ on \mathbf{X} and obtain residuals $\mathbf{A}_0 \mathbf{V}_1^{1/2}$

Obtain eigenvalues (μ_1, \dots, μ_k) of $\mathbf{A}_0 \mathbf{V}_1 \mathbf{A}_0$

for $b = 1$ to B do

Simulate (w_1, \dots, w_{n-s}) independent standard normals

Find and record the maximal value of

$$f(\lambda) = n \ln \frac{\sum_{i=1}^{n-s} w_i^2}{\sum_{i=1}^k \frac{w_i^2}{1+\lambda\mu_i} + \sum_{i=k+1}^{n-s} w_i^2} - \sum_{i=1}^l \ln(1 + \lambda\xi_i)$$

over $\lambda \in [0, \infty)$

end for

Remark: The maximization task is subject to an MM algorithm.

- 0 is a local maximum of $f(\lambda)$ if $f'(\lambda) \leq 0$. Therefore the probability of having a local maximum at 0 is

$$\text{Prob} \left(\frac{\sum_{i=1}^k \mu_i w_i^2}{\sum_{i=1}^{n-s} w_i^2} \leq \frac{1}{n} \sum_{i=1}^l \xi_i \right),$$

which provides a good approximation of the point mass at 0 of the null distribution of LRT.

- Same derivation can be carried out for the restricted (residual) LRT, in which case

$$\text{RLRT} \stackrel{\mathcal{D}}{=} \sup_{\lambda \geq 0} (n-s) \ln \frac{\sum_{i=1}^{n-s} w_i^2}{\sum_{i=1}^k \frac{w_i^2}{1+\lambda\mu_i} + \sum_{i=k+1}^{n-s} w_i^2} - \sum_{i=1}^l \ln(1 + \lambda\xi_i)$$

under the null $\lambda = 0$.

(Rao) score test with one variance component

- Again we consider the variance component model

$$\mathbf{Y} \sim N_n(\mathbf{X}\mathbf{b}, \mathbf{V}),$$

where

$$\mathbf{V} = \sigma_0^2 \mathbf{I}_n + \sigma_1^2 \mathbf{V}_1.$$

- We develop an exact score test for $H_0 : \sigma_1^2 = 0$ vs $H_A : \sigma_1^2 > 0$.

Score test avoids the maximization step in simulating the null distribution of test statistic.

- The log-likelihood function is

$$L(\mathbf{b}, \sigma_0^2, \lambda) = -\frac{1}{2} \ln \det(\mathbf{V}) - \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{b})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b}).$$

and the first derivative with respect to σ_1^2 is

$$\frac{\partial}{\partial \sigma_1^2} L = -\frac{1}{2} \text{tr}(\mathbf{V}^{-1} \mathbf{V}_1) + \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{b})^T \mathbf{V}^{-1} \mathbf{V}_1 \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b}).$$

The information matrix relevant to variance components are

$$\begin{aligned} \text{E} \left(-\frac{\partial^2}{\partial \sigma_0^2 \partial \sigma_0^2} L \right) &= \frac{1}{2} \text{tr}(\mathbf{V}^{-2}) \\ \text{E} \left(-\frac{\partial^2}{\partial \sigma_0^2 \partial \sigma_1^2} L \right) &= \text{E} \left(-\frac{\partial^2}{\partial \sigma_1^2 \partial \sigma_0^2} L \right) = \frac{1}{2} \text{tr}(\mathbf{V}^{-2} \mathbf{V}_1) \\ \text{E} \left(-\frac{\partial^2}{\partial \sigma_1^2 \partial \sigma_1^2} L \right) &= \frac{1}{2} \text{tr}(\mathbf{V}^{-1} \mathbf{V}_1 \mathbf{V}^{-1} \mathbf{V}_1). \end{aligned}$$

- The (Rao) score statistic is based on

$$\mathbf{I}_{\sigma_1^2, \sigma_1^2}^{-1} \left(\frac{\partial}{\partial \sigma_1^2} L \right)^2$$

evaluated at the MLE under the null.

- We evaluate the partial derivatives at the MLE under the null

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \hat{\sigma}_0^2 = \frac{\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}}{n}.$$

That is

$$\begin{aligned} D_1 &:= \frac{\partial}{\partial \sigma_1^2} L(\hat{\mathbf{b}}, \hat{\sigma}_0^2) \\ &= -\frac{n \text{tr}(\mathbf{V}_1)}{2 \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}} + \frac{n^2 \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{V}_1 (\mathbf{I} - \mathbf{P}_X) \mathbf{y}}{2 [\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}]^2} \\ &= \frac{-n \text{tr}(\mathbf{V}_1) [\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}] + n^2 \mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{V}_1 (\mathbf{I} - \mathbf{P}_X) \mathbf{y}}{2 [\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}]^2} \end{aligned}$$

$$\begin{aligned} J_{00} &:= \text{E} \left(-\frac{\partial^2}{\partial \sigma_0^2 \partial \sigma_0^2} L(\hat{\mathbf{b}}, \hat{\sigma}_0^2) \right) = \frac{n^3}{2 [\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}]^2} \\ J_{01} &= J_{10} := \text{E} \left(-\frac{\partial^2}{\partial \sigma_0^2 \partial \sigma_1^2} L(\hat{\mathbf{b}}, \hat{\sigma}_0^2) \right) = \frac{n^2 \text{tr}(\mathbf{V}_1)}{2 [\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}]^2} \\ J_{11} &:= \text{E} \left(-\frac{\partial^2}{\partial \sigma_1^2 \partial \sigma_1^2} L(\hat{\mathbf{b}}, \hat{\sigma}_0^2) \right) = \frac{n^2 \text{tr}(\mathbf{V}_1^2)}{2 [\mathbf{y}^T (\mathbf{I} - \mathbf{P}_X) \mathbf{y}]^2}, \end{aligned}$$

from which we form the score statistic

$$\begin{aligned}
T &= \begin{cases} (J_{11} - J_{10}J_{00}^{-1}J_{01})^{-1}D_1^2 & D_1 \geq 0 \\ 0 & D_1 < 0 \end{cases} \\
&= \begin{cases} \frac{\left[\frac{-n\text{tr}(\mathbf{V}_1) + n^2 \frac{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{V}_1(\mathbf{I} - \mathbf{P}_X)\mathbf{y}}{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{y}} \right]^2}{2[n^2\text{tr}(\mathbf{V}_1^2) - n\text{tr}(\mathbf{V}_1)^2]} & \frac{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{V}_1(\mathbf{I} - \mathbf{P}_X)\mathbf{y}}{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{y}} \geq \frac{\text{tr}(\mathbf{V}_1)}{n} \\ 0 & \frac{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{V}_1(\mathbf{I} - \mathbf{P}_X)\mathbf{y}}{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{y}} < \frac{\text{tr}(\mathbf{V}_1)}{n} \end{cases}.
\end{aligned}$$

Essentially the score test rejects when

$$T' = \max \left\{ \frac{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{V}_1(\mathbf{I} - \mathbf{P}_X)\mathbf{y}}{\mathbf{y}^T(\mathbf{I} - \mathbf{P}_X)\mathbf{y}}, \frac{\text{tr}(\mathbf{V}_1)}{n} \right\}$$

is large.

- Let the eigen-decomposition of $(\mathbf{I} - \mathbf{P}_X)\mathbf{V}_1(\mathbf{I} - \mathbf{P}_X)$ be

$$(\mathbf{I} - \mathbf{P}_X)\mathbf{V}_1(\mathbf{I} - \mathbf{P}_X) = \mathbf{Q}_1 \text{diag}(\mu_1, \dots, \mu_k) \mathbf{Q}_1^T,$$

\mathbf{Q}_2 be an orthonormal basis of $\mathcal{C}(\mathbf{X}, \mathbf{Q}_1)^\perp$, and $\mathbf{Q} = (\mathbf{Q}_1, \mathbf{Q}_2) \in \mathbb{R}^{n \times (n-s)}$. Then

$$\begin{aligned}
T' &= \max \left\{ \frac{\mathbf{y}^T \mathbf{Q} \text{diag}(\mu_1, \dots, \mu_k, 0, \dots, 0) \mathbf{Q}^T \mathbf{y}}{\mathbf{y}^T \mathbf{Q} \mathbf{Q}^T \mathbf{y}}, \frac{\text{tr}(\mathbf{V}_1)}{n} \right\} \\
&=^{\mathcal{D}} \max \left\{ \frac{\sum_{i=1}^k \mu_k w_i^2}{\sum_{i=1}^{n-s} w_i^2}, \frac{\text{tr}(\mathbf{V}_1)}{n} \right\},
\end{aligned}$$

where w_i are $n - s$ independent standard normals.

- The null distribution can be obtained from computer simulation:

Given $\mathbf{y}, \mathbf{X}, \mathbf{V}_1$, and simulation replicates B

Eigen-decomposition: $\mathbf{V}_1 = \mathbf{U}_\ell \text{diag}(\xi_1, \dots, \xi_\ell) \mathbf{U}_\ell^T$

Regress $\mathbf{U}_\ell \text{diag}(\sqrt{\xi_1}, \dots, \sqrt{\xi_\ell})$ on \mathbf{X} to obtain residuals $(\mathbf{I} - \mathbf{P}_X)\mathbf{V}_1^{1/2}$

Obtain eigenvalues (μ_1, \dots, μ_k) of $(\mathbf{I} - \mathbf{P}_X)\mathbf{V}_1(\mathbf{I} - \mathbf{P}_X)$

for $b = 1$ to B do

 Simulate (w_1, \dots, w_{n-s}) independent standard normals

 Record the value

$$\max \left\{ \frac{\sum_{i=1}^k \mu_k w_i^2}{\sum_{i=1}^{n-s} w_i^2}, \frac{\text{tr}(\mathbf{V}_1)}{n} \right\}$$

end for

- The probability mass of T' at lower boundary $\text{tr}(\mathbf{V}_1)/n$ is

$$\text{Prob} \left(\frac{\sum_{i=1}^k \mu_i w_i^2}{\sum_{i=1}^{n-s} w_i^2} \leq \frac{\text{tr}(\mathbf{V}_1)}{n} \right),$$

same as the point mass at 0 of the null distribution of eLRT.

Wald test with one variance component

TODO.

24 Lecture 24: Dec 4

Announcement

- HW8 due today.
- Course evaluation: <https://classeval.ncsu.edu/>
- No OH from next week. Reachable by email.
- Final exam: Friday, Dec 13 @ 8A-11A, SAS 5270.

Last time

- Variance component testing: eLRT, eRLRT, score test.

Today

- Variance component model: testing with two or more variance components.
- (Last topic!) Prediction: BLP and BLUP.

Testing with two or more variance components

- First we extend the test with one variance component to a slightly more general case $\mathbf{y} \sim N(\mathbf{X}\mathbf{b}, \mathbf{V})$ with

$$\mathbf{V} = \sigma_0^2 \mathbf{V}_0 + \sigma_1^2 \mathbf{V}_1, \quad (5)$$

where $\mathbf{V}_0 \in \mathbb{R}^{n \times n}$ is a known psd matrix. Let $r = \text{rank}(\mathbf{V}_0)$.

Given eigen-decomposition $\mathbf{V}_0 = \mathbf{U}_r \mathbf{D}_r \mathbf{U}_r^T$, define $\mathbf{T} = \mathbf{D}_r^{-1/2} \mathbf{U}_r^T \in \mathbb{R}^{r \times n}$. Then

$$\mathbf{T}\mathbf{Y} \sim N(\mathbf{T}\mathbf{X}\mathbf{b}, \sigma_0^2 \mathbf{I}_r + \sigma_1^2 \mathbf{T}\mathbf{V}_1 \mathbf{T}^T)$$

and the eLRT, eRLRT or score test can be applied to $\mathbf{T}\mathbf{Y}$.

- Now we consider the linear model with two variance components

$$\mathbf{Y} \sim N_n(\mathbf{X}\mathbf{b}, \mathbf{V}),$$

where

$$\mathbf{V} = \sigma_0^2 \mathbf{I}_n + \sigma_1^2 \mathbf{V}_1 + \sigma_2^2 \mathbf{V}_2.$$

- We are interested in testing $H_0 : \sigma_2^2 = 0$ vs $H_A : \sigma_2^2 > 0$.
- The idea is to massage the problem to the case (5) above.
- First perform QR (Gram-Schmidt) on the matrix $(\mathbf{X}, \mathbf{V}_1, \mathbf{V}_2, \mathbf{I}_n)$ to obtain an orthonormal basis $(\mathbf{Q}_0, \mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3)$ of \mathbb{R}^n , where
 - \mathbf{Q}_0 is an orthonormal basis of $\mathcal{C}(\mathbf{X})$
 - \mathbf{Q}_1 is an orthonormal basis of $\mathcal{C}(\mathbf{X}, \mathbf{V}_1) - \mathcal{C}(\mathbf{X})$
 - \mathbf{Q}_2 is an orthonormal basis of $\mathcal{C}(\mathbf{X}, \mathbf{V}_1, \mathbf{V}_2) - \mathcal{C}(\mathbf{X}, \mathbf{V}_1)$
 - \mathbf{Q}_3 is an orthonormal basis of $\mathcal{C}(\mathbf{X}, \mathbf{V}_1, \mathbf{V}_2)^\perp$.
- If $\text{rank}(\mathbf{Q}_2) > 0$, that is $\mathcal{C}(\mathbf{X}, \mathbf{V}_1) \subsetneq \mathcal{C}(\mathbf{X}, \mathbf{V}_1, \mathbf{V}_2)$, then

$$\mathbf{Q}_2^T \mathbf{Y} \sim N(\mathbf{0}, \sigma_0^2 \mathbf{I} + \sigma_2^2 \mathbf{Q}_2^T \mathbf{V}_2 \mathbf{Q}_2)$$

and the exact F test, eLRT, eRLRT and score test can be applied to $\mathbf{Q}_2^T \mathbf{y}$.

Is this beneficial to add the $\mathbf{Q}_3^T \mathbf{Y}$ component???

- If $\text{rank}(\mathbf{Q}_2) = 0$, that is $\mathcal{C}(\mathbf{X}, \mathbf{V}_1) = \mathcal{C}(\mathbf{X}, \mathbf{V}_1, \mathbf{V}_2)$, then we would like to choose λ and matrix \mathbf{K} such that

$$\mathbf{Q}_1^T \mathbf{Y} + \mathbf{K} \mathbf{Q}_3^T \mathbf{Y} \sim N(\mathbf{0}, (\sigma_1^2 + \sigma_0^2/\lambda) \mathbf{Q}_1^T \mathbf{V}_1 \mathbf{Q}_1 + \sigma_2^2 \mathbf{Q}_1^T \mathbf{V}_2 \mathbf{Q}_1),$$

which is of form (5). We consider following situations.

- If $\mathbf{Q}_1^T \mathbf{V}_2 = \mathbf{0}$, e.g., when $\mathcal{C}(\mathbf{V}_2) \subset \mathcal{C}(\mathbf{X})$, then this test cannot be performed.
- If $\mathbf{Q}_1^T \mathbf{V}_1 \mathbf{Q}_1 = \lambda \mathbf{I}$, note

$$\begin{aligned} \mathbf{Q}_1^T \mathbf{Y} &\sim N(\mathbf{0}, \sigma_0^2 \mathbf{I} + \sigma_1^2 \mathbf{Q}_1^T \mathbf{V}_1 \mathbf{Q}_1 + \sigma_2^2 \mathbf{Q}_1^T \mathbf{V}_2 \mathbf{Q}_1) \\ &= N(\mathbf{0}, (\sigma_0^2 + \lambda \sigma_1^2) \mathbf{I} + \sigma_2^2 \mathbf{Q}_1^T \mathbf{V}_2 \mathbf{Q}_1). \end{aligned}$$

Then the ordinary tests (F test, eLRT, eRLRT, score) for one variance component can be applied without using the $\mathbf{K} \mathbf{Q}_3^T \mathbf{y}$ piece as long as $\mathbf{Q}_1^T \mathbf{V}_2 \neq \mathbf{0}$.

- In general, $\mathbf{Q}_1^T \mathbf{V}_1 \mathbf{Q}_1 \neq \lambda \mathbf{I}$, then the test requires the $\mathbf{K} \mathbf{Q}_3^T \mathbf{y}$ term, which has distribution

$$\mathbf{K} \mathbf{Q}_3^T \mathbf{Y} \sim N(\mathbf{0}, \sigma_0^2 \mathbf{K} \mathbf{K}^T).$$

Note that

$$\mathbf{Q}_1^T \mathbf{V} \mathbf{Q}_3 = \mathbf{Q}_1^T (\sigma_0^2 \mathbf{I}_n + \sigma_1^2 \mathbf{V}_1 + \sigma_2^2 \mathbf{V}_2) \mathbf{Q}_3 = \mathbf{0}.$$

Therefore $\mathbf{Q}_1^T \mathbf{Y} \perp \mathbf{K} \mathbf{Q}_3^T \mathbf{Y}$. We simply pick \mathbf{K} such that

$$\mathbf{K} \mathbf{Q}_3^T \mathbf{Y} \sim N(\mathbf{0}, \sigma_0^2 (\lambda^{-1} \mathbf{Q}_1^T \mathbf{V}_2 \mathbf{Q}_1 - \mathbf{I})).$$

That is

$$\mathbf{K} \mathbf{K}^T = \lambda^{-1} \mathbf{Q}_1^T \mathbf{V}_2 \mathbf{Q}_1 - \mathbf{I}.$$

Apparently we need to choose λ such that $\lambda^{-1} \mathbf{Q}_1^T \mathbf{V}_2 \mathbf{Q}_1 - \mathbf{I}$ is psd. Let

$$\mathbf{Q}_1^T \mathbf{V}_2 \mathbf{Q}_1 = \mathbf{W} \mathbf{\Lambda} \mathbf{W}^T = \mathbf{W} \text{diag}(\lambda_i) \mathbf{W}^T$$

be the eigendecomposition. Setting λ to be the smallest eigenvalue yields

$$\mathbf{K} = \mathbf{W} \text{diag}(\sqrt{\lambda_i / \lambda - 1}).$$

Best linear prediction (BLP)

- Given data (y, x_1, \dots, x_p) , regression can be thought of one way to predict y from x_1, \dots, x_p .
- A reasonable criterion is to choose predictor $f(\mathbf{x})$ such that the mean squared error

$$\text{MSE} = \text{E}[y - f(\mathbf{x})]^2$$

is minimized. Here the expectation is wrt the joint distribution of (y, \mathbf{x}) .

- Let $m(\mathbf{x}) = \text{E}(y|\mathbf{x})$. Then for any other predictor $f(\mathbf{x})$,

$$\text{E}[y - m(\mathbf{x})]^2 \leq \text{E}[y - f(\mathbf{x})]^2.$$

That is $m(\mathbf{x})$ is the best predictor of y .

Proof.

$$\begin{aligned} & \text{E}[y - f(\mathbf{x})]^2 \\ &= \text{E}[y - m(\mathbf{x}) + m(\mathbf{x}) - f(\mathbf{x})]^2 \\ &= \text{E}[y - m(\mathbf{x})]^2 + \text{E}[m(\mathbf{x}) - f(\mathbf{x})]^2 + 2\text{E}[y - m(\mathbf{x})][m(\mathbf{x}) - f(\mathbf{x})]. \end{aligned}$$

But the cross term vanishes

$$\begin{aligned}
& \mathbb{E}[y - m(\mathbf{x})][m(\mathbf{x}) - f(\mathbf{x})] \\
&= \mathbb{E}\{\mathbb{E}[y - m(\mathbf{x})][m(\mathbf{x}) - f(\mathbf{x}) \mid \mathbf{x}]\} \\
&= \mathbb{E}\{[m(\mathbf{x}) - f(\mathbf{x})]\mathbb{E}[y - m(\mathbf{x}) \mid \mathbf{x}]\} \\
&= \mathbb{E}\{[m(\mathbf{x}) - f(\mathbf{x})]0\} \\
&= 0.
\end{aligned}$$

Therefore $\mathbb{E}[y - f(\mathbf{x})]^2 \geq \mathbb{E}[y - m(\mathbf{x})]^2$ for any $f(\mathbf{x})$. \square

- In order to use this result, we need to know the joint distribution of (y, \mathbf{x}) . This is often unrealistic \ominus If only the first two moments (means, variances, and covariances) are known, then we can find the *best linear predictor* (BLP) of y .
- Assume

$$\mathbb{E} \begin{pmatrix} y \\ \mathbf{x} \end{pmatrix} = \begin{pmatrix} \mu_y \\ \boldsymbol{\mu}_x \end{pmatrix}, \quad \text{Cov} \begin{pmatrix} \sigma_{yy} & \boldsymbol{\Sigma}_{xy}^T \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{pmatrix}.$$

Let $\boldsymbol{\beta}^*$ be a solution of $\boldsymbol{\Sigma}_{xx}\boldsymbol{\beta} = \boldsymbol{\Sigma}_{xy}$. Then

$$\widehat{\mathbb{E}}(y|\mathbf{x}) := \mu_y + (\mathbf{x} - \boldsymbol{\mu}_x)^T \boldsymbol{\beta}^*$$

is *the* best linear predictor (BLP) of y .

Proof. Let $f(\mathbf{x}) = \alpha + \mathbf{x}^T \boldsymbol{\beta}$ be an arbitrary linear predictor. Then we find $\alpha, \boldsymbol{\beta}$ by minimizing the MSE

$$\begin{aligned}
& \mathbb{E}[y - f(\mathbf{x})]^2 \\
&= \mathbb{E}(y - \alpha - \mathbf{x}^T \boldsymbol{\beta})^2 \\
&= \mathbb{E}(y - \alpha)^2 + \boldsymbol{\beta}^T \mathbb{E}(\mathbf{x}\mathbf{x}^T) \boldsymbol{\beta} - 2\mathbb{E}[(y - \alpha)(\mathbf{x}^T \boldsymbol{\beta})] \\
&= \mathbb{E}(y - \alpha)^2 + \boldsymbol{\beta}^T \mathbb{E}(\mathbf{x}\mathbf{x}^T) \boldsymbol{\beta} - 2\mathbb{E}(y\mathbf{x}^T) \boldsymbol{\beta} + 2\alpha \mathbb{E}(\mathbf{x}^T) \boldsymbol{\beta} \\
&= \mathbb{E}(y - \alpha)^2 + \boldsymbol{\beta}^T \mathbb{E}(\mathbf{x}\mathbf{x}^T) \boldsymbol{\beta} - 2\mathbb{E}(y\mathbf{x})^T \boldsymbol{\beta} + 2\alpha \boldsymbol{\mu}_x^T \boldsymbol{\beta}.
\end{aligned}$$

Setting derivatives to 0 gives

$$\begin{aligned}
\frac{\partial}{\partial \alpha} \mathbb{E}[y - f(\mathbf{x})]^2 &= 2(\alpha - \mu_y) + 2\boldsymbol{\mu}_x^T \boldsymbol{\beta} = 0 \\
\nabla \boldsymbol{\beta} \mathbb{E}[y - f(\mathbf{x})]^2 &= 2\mathbb{E}(\mathbf{x}\mathbf{x}^T) \boldsymbol{\beta} - 2\mathbb{E}(y\mathbf{x}) + 2\alpha \boldsymbol{\mu}_x = \mathbf{0}_p.
\end{aligned}$$

From the first equation, $\alpha = \mu_y - \boldsymbol{\mu}_x^T \boldsymbol{\beta}$. Substitution into the second equation yields

$$(\mathbf{E}\mathbf{x}\mathbf{x}^T - \boldsymbol{\mu}_x\boldsymbol{\mu}_x^T)\boldsymbol{\beta} = \mathbf{E}(\mathbf{x}y) - \boldsymbol{\mu}_x\mu_y.$$

That is

$$\boldsymbol{\Sigma}_{\mathbf{x}\mathbf{x}}\boldsymbol{\beta} = \boldsymbol{\Sigma}_{\mathbf{x}y}.$$

Therefore the optimal α is $\alpha^* = \mu_y - \boldsymbol{\mu}_x^T \boldsymbol{\beta}^*$, where $\boldsymbol{\beta}^*$ is any solution to above equation. And the BLP is

$$\alpha^* + \mathbf{x}^T \boldsymbol{\beta}^* = \mu_y + (\mathbf{x} - \boldsymbol{\mu}_x)^T \boldsymbol{\beta}^*.$$

Because the criterion function is a convex function, the stationary condition is both necessary and sufficient for the global minima. Therefore any BLP must be of this form. \square

Best linear unbiased prediction (BLUP)

- Consider random variables y_0, y_1, \dots, y_n , and we are interested predicting y_0 given data y_1, \dots, y_n . If we know the mean, variances, and covariances, then we can use above theory to find the BLP of y_0 .
- In practice, we don't know the means $\mu_0 = \mathbf{E}y_0$, $\boldsymbol{\mu}_y = \mathbf{E}\mathbf{y}$ most of time \odot
Let's impose a linear (Aitken) model for the means μ_i

$$\mathbf{E} \begin{pmatrix} \mathbf{y} \\ y_0 \end{pmatrix} = \begin{pmatrix} \mathbf{X}\mathbf{b} \\ \mathbf{x}_0^T \mathbf{b} \end{pmatrix}, \quad \text{Cov} \begin{pmatrix} \mathbf{y} \\ y_0 \end{pmatrix} = \begin{pmatrix} \mathbf{V} & \mathbf{V}_{\mathbf{y}y_0} \\ \mathbf{V}_{\mathbf{y}y_0}^T & V_{y_0y_0} \end{pmatrix}. \quad (6)$$

- If \mathbf{b} is known, then the BLP of y_0 is

$$\mathbf{x}_0^T \mathbf{b} + (\mathbf{y} - \mathbf{X}\mathbf{b})^T \boldsymbol{\beta}^*,$$

where $\boldsymbol{\beta}^*$ is a solution to $\mathbf{V}\boldsymbol{\beta} = \mathbf{V}_{\mathbf{y}y_0}$.

- If \mathbf{b} is unknown, then the hope is to find the *best linear unbiased predictor* (BLUP).

We call a linear predictor $f(\mathbf{y}) = a_0 + \mathbf{a}^T \mathbf{y}$ of y_0 BLUP if

1. it is unbiased, i.e.,

$$E f(\mathbf{y}) = a_0 + \mathbf{a}^T \mathbf{X} \mathbf{b} = \mathbf{x}_0^T \mathbf{b} = E y_0$$

for all \mathbf{b} , and

2. for any other linear unbiased predictor $b_0 + \mathbf{b}^T \mathbf{y}$,

$$E(y_0 - a_0 - \mathbf{a}^T \mathbf{y})^2 \leq E(y_0 - b_0 - \mathbf{b}^T \mathbf{y})^2.$$

- Theorem: Under the Aitken model (6) and assume $\mathbf{V}_{y_0} \in \mathcal{C}(\mathbf{V}, \mathbf{X})$ and $\mathbf{x}_0 \in \mathcal{C}(\mathbf{X}^T)$, the BLUP of y_0 is

$$\mathbf{x}_0^T \hat{\mathbf{b}}_{\text{GLS}} + (\mathbf{y} - \mathbf{X} \hat{\mathbf{b}}_{\text{GLS}})^T \boldsymbol{\beta}_*,$$

where $\boldsymbol{\beta}_*$ is a solution of $(\mathbf{V} + \mathbf{X} \mathbf{X}^T) \boldsymbol{\beta} = \mathbf{V}_{y_0}$ and $\mathbf{X} \hat{\mathbf{b}}_{\text{GLS}}$ is the BLUE of $\mathbf{X} \mathbf{b}$.

Remark 1: We don't assume \mathbf{V} is nonsingular. For nonsingular \mathbf{V} , we can take $\boldsymbol{\beta}_* = \mathbf{V}^{-1} \mathbf{V}_{y_0}$.

Remark 2: Both $\hat{\mathbf{b}}_{\text{GLS}}$ and $\boldsymbol{\beta}_*$ depend crucially on \mathbf{V} .

Proof. Let $a_0 + \mathbf{a}^T \mathbf{y}$ be arbitrary linear predictor of y_0 . Unbiasedness requires

$$a_0 + \mathbf{a}^T \mathbf{X} \mathbf{b} = \mathbf{x}_0^T \mathbf{b}$$

for all \mathbf{b} . Thus $a_0 = 0$ and $\mathbf{a}^T \mathbf{X} = \mathbf{x}_0^T$. We need to solve the constrained optimization problem

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} E(y_0 - \mathbf{a}^T \mathbf{y})^2 = \frac{1}{2} \mathbf{a}^T E(\mathbf{y} \mathbf{y}^T) \mathbf{a} - E(y_0 \mathbf{a}^T \mathbf{y}) + \frac{1}{2} E y_0^2 \\ & \text{subject to} \quad \mathbf{X}^T \mathbf{a} = \mathbf{x}_0. \end{aligned}$$

Setting the gradient of the Lagrangian

$$L(\mathbf{a}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{a}^T E(\mathbf{y} \mathbf{y}^T) \mathbf{a} - E(y_0 \mathbf{a}^T \mathbf{y}) + \frac{1}{2} E y_0^2 + \boldsymbol{\lambda}^T (\mathbf{X}^T \mathbf{a} - \mathbf{x}_0)$$

to zero yields equations

$$\begin{aligned} \mathbf{E}(\mathbf{y} \mathbf{y}^T) \mathbf{a} - \mathbf{E}(y_0 \mathbf{y}) + \mathbf{X} \boldsymbol{\lambda} &= \mathbf{0}_n \\ \mathbf{X}^T \mathbf{a} &= \mathbf{x}_0. \end{aligned}$$

Adding and subtracting $(\mathbf{E}\mathbf{y}\mathbf{E}\mathbf{y}^T)\mathbf{a} = \mathbf{X}\mathbf{b}\mathbf{b}^T\mathbf{X}^T\mathbf{a} = \mathbf{X}\mathbf{b}\mathbf{b}^T\mathbf{x}_0 = \mathbf{E}\mathbf{y}\mathbf{E}\mathbf{y}_0$ to the first equation shows $\mathbf{V}\mathbf{a} - \mathbf{V}_{\mathbf{y}\mathbf{y}_0} + \mathbf{X}\boldsymbol{\lambda} = \mathbf{0}_n$. In matrix notation,

$$\begin{pmatrix} \mathbf{V} & \mathbf{X} \\ \mathbf{X}^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} \mathbf{V}_{\mathbf{y}\mathbf{y}_0} \\ \mathbf{x}_0 \end{pmatrix}.$$

By HW4 5(d), solution for the optimal \mathbf{a} is

$$\mathbf{a}_* = [\mathbf{V}_0^- - \mathbf{V}_0^- \mathbf{X}(\mathbf{X}^T \mathbf{V}_0^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_0^-] \mathbf{V}_{\mathbf{y}\mathbf{y}_0} + \mathbf{V}_0^- \mathbf{X}(\mathbf{X}^T \mathbf{V}_0^- \mathbf{X})^{-1} \mathbf{x}_0,$$

where $\mathbf{V}_0 = \mathbf{V} + \mathbf{X}\mathbf{X}^T$. Thus the BLUP is

$$\begin{aligned} \mathbf{a}_*^T \mathbf{y} &= \mathbf{V}_{\mathbf{y}\mathbf{y}_0}^T [\mathbf{V}_0^- - \mathbf{V}_0^- \mathbf{X}(\mathbf{X}^T \mathbf{V}_0^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_0^-] \mathbf{y} + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{V}_0^- \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}_0^- \mathbf{y} \\ &= \mathbf{V}_{\mathbf{y}\mathbf{y}_0}^T \mathbf{V}_0^- (\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}_{\text{GLS}}) + \mathbf{x}_0^T \hat{\mathbf{b}}_{\text{GLS}} \\ &= \boldsymbol{\beta}_*^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}_{\text{GLS}}) + \mathbf{x}_0^T \hat{\mathbf{b}}_{\text{GLS}}. \end{aligned}$$

□

- The prediction variance of BLUP is (TODO)

$$\mathbb{E}(y_0 - \mathbf{a}_*^T \mathbf{y})^2 = V_{y_0 y_0} - 2\mathbf{a}_*^T \mathbf{V}_{\mathbf{y}\mathbf{y}_0} + \mathbf{a}_*^T \mathbf{V} \mathbf{a}_*. \quad ???$$

- Example (BLUP in Gauss-Markov linear model): $\mathbf{V} = \sigma^2 \mathbf{I}$ and $\mathbf{V}_{\mathbf{y}\mathbf{y}_0} = \mathbf{0}_p$. Thus $\boldsymbol{\beta}_* = \mathbf{0}_p$ and the BLUP for y_0 is $\mathbf{x}_0^T \hat{\mathbf{b}}$, which is also the BLUE of $\mathbf{x}_0^T \mathbf{b}$.

Mixed model equation (MME)

- Consider the mixed effects model

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

- $\mathbf{X} \in \mathbb{R}^{n \times p}$ is a design matrix for fixed effects $\mathbf{b} \in \mathbb{R}^p$.
- $\mathbf{Z} \in \mathbb{R}^{n \times q}$ is a design matrix for random effects $\mathbf{u} \in \mathbb{R}^q$.
- The most general assumption is $\mathbf{e} \in N(\mathbf{0}_n, \mathbf{R})$, $\mathbf{u} \in N(\mathbf{0}_q, \mathbf{G})$, and \mathbf{e} is independent of \mathbf{u} . That is

$$\begin{pmatrix} \mathbf{u} \\ \mathbf{e} \end{pmatrix} \sim N \left(\mathbf{0}_{q+n}, \begin{pmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{pmatrix} \right).$$

Assume \mathbf{G} and \mathbf{R} are nonsingular.

- We already know that the BLUE of \mathbf{Xb} is

$$\mathbf{X}\hat{\mathbf{b}}_{\text{GLS}} = \mathbf{X}[\mathbf{X}^T(\mathbf{R} + \mathbf{ZGZ}^T)^{-1}\mathbf{X}]^{-1}\mathbf{X}^T(\mathbf{R} + \mathbf{ZGZ}^T)^{-1}\mathbf{y}.$$

- We can apply the previous theorem to derive the BLUP of \mathbf{u} . Note

$$\text{E} \begin{pmatrix} \mathbf{y} \\ \mathbf{u} \end{pmatrix} = \begin{pmatrix} \mathbf{Xb} \\ \mathbf{0}_q \end{pmatrix}, \quad \text{Cov} \begin{pmatrix} \mathbf{y} \\ \mathbf{u} \end{pmatrix} = \begin{pmatrix} \mathbf{R} + \mathbf{ZGZ}^T & \mathbf{ZG} \\ \mathbf{GZ}^T & \mathbf{R} \end{pmatrix}.$$

Therefore the BLUP for \mathbf{u} is

$$\mathbf{GZ}^T(\mathbf{R} + \mathbf{ZGZ}^T)^{-1}(\mathbf{y} - \mathbf{X}\hat{\mathbf{b}}_{\text{GLS}}).$$

- It turns out both BLUE of \mathbf{Xb} and BLUP of \mathbf{u} can be obtained simultaneously by solving a so-called mixed model equation (MME).

Mixed model equation (MME) defined as

$$\begin{pmatrix} \mathbf{X}^T\mathbf{R}^{-1}\mathbf{X} & \mathbf{X}^T\mathbf{R}^{-1}\mathbf{Z} \\ \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{X} & \mathbf{G}^{-1} + \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{Z} \end{pmatrix} \begin{pmatrix} \mathbf{b} \\ \mathbf{u} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T\mathbf{R}^{-1}\mathbf{y} \\ \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{y} \end{pmatrix}$$

is a generalization of the normal equation for fixed effects model.

- Theorem: Let $(\hat{\mathbf{b}}, \hat{\mathbf{u}})$ be a solution to MME. Then $\mathbf{X}\hat{\mathbf{b}}$ is the BLUE of \mathbf{Xb} and $\hat{\mathbf{u}}$ is the BLUP of \mathbf{u} .

Proof. Let

$$\mathbf{V} = \text{Cov}(\mathbf{y}) = \mathbf{R} + \mathbf{ZGZ}^T.$$

Then $\mathbf{X}\hat{\mathbf{b}}$ is a BLUE of \mathbf{Xb} if $\hat{\mathbf{b}}$ is a solution to

$$\mathbf{X}^T\mathbf{V}^{-1}\mathbf{Xb} = \mathbf{X}^T\mathbf{V}^{-1}\mathbf{y}.$$

By the binomial inversion formula (HW1), we have

$$\mathbf{V}^{-1} = \mathbf{R}^{-1} - \mathbf{R}^{-1}\mathbf{Z}(\mathbf{G}^{-1} + \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{Z})^{-1}\mathbf{Z}^T\mathbf{R}^{-1}.$$

The MME says

$$\begin{aligned} \mathbf{X}^T\mathbf{R}^{-1}\mathbf{Xb} + \mathbf{X}^T\mathbf{R}^{-1}\mathbf{Zu} &= \mathbf{X}^T\mathbf{R}^{-1}\mathbf{y} \\ \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{Xb} + (\mathbf{G}^{-1} + \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{Z})\mathbf{u} &= \mathbf{Z}^T\mathbf{R}^{-1}\mathbf{y}. \end{aligned}$$

From the second equation we solve for

$$\mathbf{u} = (\mathbf{G}^{-1} + \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{X} \mathbf{b})$$

and substituting into the first equation shows

$$\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \hat{\mathbf{b}} = \mathbf{X}^T \mathbf{V}^{-1} \mathbf{y}.$$

Thus $\hat{\mathbf{b}}$ is a generalized least squares solution and $\mathbf{X} \hat{\mathbf{b}}$ is the BLUE of $\mathbf{X} \mathbf{b}$.

To show that $\hat{\mathbf{u}}$ is BLUP, only need to show that

$$(\mathbf{G}^{-1} + \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{R}^{-1} = \mathbf{G} \mathbf{Z}^T (\mathbf{R} + \mathbf{Z} \mathbf{G} \mathbf{Z}^T)^{-1},$$

or equivalently

$$\mathbf{Z}^T \mathbf{R}^{-1} (\mathbf{R} + \mathbf{Z} \mathbf{G} \mathbf{Z}^T) = (\mathbf{G}^{-1} + \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{Z}) \mathbf{G} \mathbf{Z}^T,$$

which is obvious. □

References

- Casella, G. and Berger, R. (2001). *Statistical Inference*. Duxbury Resource Center.
- Christensen, R. (1996). Exact tests for variance components. *Biometrics*, 52(1):309–314.
- Crainiceanu, C. M. and Ruppert, D. (2004). Likelihood ratio tests in linear mixed models with one variance component. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 66(1):165–185.
- Davies, R. B. (1980). The distribution of linear combinations of χ^2 random variables. *Applied Statistics*, 29:323–333.
- Öfversten, J. (1993). Exact tests for variance components in unbalanced mixed linear models. *Biometrics*, 49(1):45–57.
- Seely, J. F. and El Bassiouni, Y. (1983). Applying Wald’s variance component test. *Ann. Statist.*, 11(1):197–201.
- Teets, D. and Whitehead, K. (1999). The discovery of Ceres: how Gauss became famous. *Math. Mag.*, 72(2):83–93.