

Types of Problems–Review

Notation: Let $\mathbf{A}^T = \mathbf{A}'$ be the transpose of \mathbf{A} .

0) Covariance and Expected Value = Mean, and the Multivariate Normal (MVN) Distribution:

Notation: Unless told otherwise, assume expectations exist and that conformable matrices and vectors are used.

The *population mean* of a random $n \times 1$ vector $\mathbf{x} = (x_1, \dots, x_n)^T$ is $E(\mathbf{x}) = \boldsymbol{\mu} = (E(x_1), \dots, E(x_n))^T$ and the $n \times n$ *population covariance matrix* $\text{Cov}(\mathbf{x}) = \boldsymbol{\Sigma}\mathbf{x} = E(\mathbf{x} - E(\mathbf{x}))(\mathbf{x} - E(\mathbf{x}))^T = (\sigma_{i,j})$ where $\text{Cov}(x_i, x_j) = \sigma_{i,j}$. The *population covariance matrix* of \mathbf{x} with \mathbf{y} is

$$\text{Cov}(\mathbf{x}, \mathbf{y}) = \boldsymbol{\Sigma}\mathbf{x}, \mathbf{y} = E[(\mathbf{x} - E(\mathbf{x}))(\mathbf{y} - E(\mathbf{y}))^T].$$

If \mathbf{X} and \mathbf{Y} are $n \times 1$ random vectors, \mathbf{a} a conformable constant vector, and \mathbf{A} and \mathbf{B} are conformable constant matrices, then

$$E(\mathbf{X} + \mathbf{Y}) = E(\mathbf{X}) + E(\mathbf{Y}), \quad E(\mathbf{a} + \mathbf{Y}) = \mathbf{a} + E(\mathbf{Y}), \quad \& \quad E(\mathbf{A}\mathbf{X}\mathbf{B}) = \mathbf{A}E(\mathbf{X})\mathbf{B}.$$

Also

$$\text{Cov}(\mathbf{a} + \mathbf{A}\mathbf{X}) = \text{Cov}(\mathbf{A}\mathbf{X}) = \mathbf{A}\text{Cov}(\mathbf{X})\mathbf{A}^T.$$

Note that $E(\mathbf{A}\mathbf{Y}) = \mathbf{A}E(\mathbf{Y})$ and $\text{Cov}(\mathbf{A}\mathbf{Y}) = \mathbf{A}\text{Cov}(\mathbf{Y})\mathbf{A}^T$.

If \mathbf{X} ($m \times 1$) and \mathbf{Y} ($n \times 1$) are random vectors, and \mathbf{A} and \mathbf{B} are conformable constant matrices, then

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

If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then $E(\mathbf{X}) = \boldsymbol{\mu}$, $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$, and $m_{\mathbf{X}}(\mathbf{t}) = \exp(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t})$.

If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and if \mathbf{A} is a $q \times p$ matrix, then $\mathbf{A}\mathbf{X} \sim N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$. If \mathbf{a} ($p \times 1$) and \mathbf{b} ($q \times 1$) are constant vectors, then $\mathbf{X} + \mathbf{a} \sim N_p(\boldsymbol{\mu} + \mathbf{a}, \boldsymbol{\Sigma})$ and $\mathbf{A}\mathbf{X} + \mathbf{b} \sim N_q(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$.

$$\text{Let } \mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \text{and } \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

The conditional distribution of a MVN is MVN. If $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then the conditional distribution of \mathbf{X}_1 given that $\mathbf{X}_2 = \mathbf{x}_2$ is multivariate normal with mean $\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$ and covariance matrix $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$. That is,

$$\mathbf{X}_1 | \mathbf{X}_2 = \mathbf{x}_2 \sim N_q(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}).$$

Notation:

$$\mathbf{X}_1 | \mathbf{X}_2 \sim N_q(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{X}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}).$$

1) Projection Matrices, Generalized Inverses, and the Column Space $C(\mathbf{X})$:

Let $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_m]$ be an $n \times m$ matrix. The space spanned by the columns of \mathbf{A} = column space of $\mathbf{A} = C(\mathbf{A})$.

Let $\mathbf{X} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_p]$ be an $n \times p$ matrix. Then $C(\mathbf{X}) = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{y} = \mathbf{X}\boldsymbol{\beta} \text{ for some } \boldsymbol{\beta} \in \mathbb{R}^p\}$.

One way to show $C(\mathbf{A}) = C(\mathbf{B})$ is to show that i) $\mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{y} \in C(\mathbf{B})$ and ii) $\mathbf{B}\mathbf{y} = \mathbf{A}\mathbf{x} \in C(\mathbf{A})$.

The *null space* of $\mathbf{A} = N(\mathbf{A}) = \{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{0}\} = \text{kernel of } \mathbf{A}$. The subspace $V^\perp = \{\mathbf{y} \in \mathbb{R}^k : \mathbf{y} \perp V\}$ is the *orthogonal complement* of V .

$N(\mathbf{A}^T) = [C(\mathbf{A})]^\perp$, so $N(\mathbf{A}) = [C(\mathbf{A}^T)]^\perp$.

A **generalized inverse** of an $m \times n$ matrix \mathbf{A} is any $n \times m$ matrix \mathbf{A}^- satisfying $\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A}$. Other names are conditional inverse, pseudo inverse, g-inverse, and p-inverse. Usually a generalized inverse is not unique, but if \mathbf{A}^{-1} exists, then $\mathbf{A}^- = \mathbf{A}^{-1}$ is unique. Notation: $\mathbf{G} := \mathbf{A}^-$ means \mathbf{G} is a generalized inverse of \mathbf{A} .

Let V be a subspace of \mathbb{R}^k . Then every $\mathbf{y} \in \mathbb{R}^k$ can be expressed uniquely as $\mathbf{y} = \mathbf{w} + \mathbf{z}$ where $\mathbf{w} \in V$ and $\mathbf{z} \in V^\perp$.

Let $\mathbf{X} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_p]$ be $n \times p$, and let $V = C(\mathbf{X}) = \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_p)$. Then the $n \times n$ matrix $\mathbf{P}_V = \mathbf{P}_\mathbf{X}$ is a **projection matrix** on $C(\mathbf{X})$ if $\mathbf{P}_\mathbf{X} \mathbf{y} = \mathbf{w} \ \forall \ \mathbf{y} \in \mathbb{R}^n$. (Here $\mathbf{y} = \mathbf{w} + \mathbf{z} = \mathbf{w}\mathbf{y} + \mathbf{z}\mathbf{y}$, so \mathbf{w} depends on \mathbf{y} .)

Projection Matrix Theorem: a) $\mathbf{P}_\mathbf{X}$ is unique.

b) $\mathbf{P}_\mathbf{X} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^-\mathbf{X}^T$ where $(\mathbf{X}^T\mathbf{X})^-$ is any generalized inverse of $\mathbf{X}^T\mathbf{X}$.

c) \mathbf{A} is a projection matrix on $C(\mathbf{A})$ iff \mathbf{A} is symmetric and idempotent. Hence $\mathbf{P}_\mathbf{X}$ is a projection matrix on $C(\mathbf{P}_\mathbf{X}) = C(\mathbf{X})$.

d) $\mathbf{I}_n - \mathbf{P}_\mathbf{X}$ is the projection matrix on $[C(\mathbf{X})]^\perp$.

e) $\mathbf{A} = \mathbf{P}_\mathbf{X}$ iff i) $\mathbf{y} \in C(\mathbf{X})$ implies $\mathbf{A}\mathbf{y} = \mathbf{y}$ and ii) $\mathbf{y} \perp C(\mathbf{X})$ implies $\mathbf{A}\mathbf{y} = \mathbf{0}$.

f) $\mathbf{P}_\mathbf{X}\mathbf{X} = \mathbf{X}$, and $\mathbf{P}_\mathbf{X}\mathbf{W} = \mathbf{W}$ if each column of $\mathbf{W} \in C(\mathbf{X})$.

g) $\mathbf{P}_\mathbf{X}\mathbf{v}_i = \mathbf{v}_i$.

h) If $C(\mathbf{X}_R)$ is a subspace of $C(\mathbf{X})$, then $\mathbf{P}_\mathbf{X}\mathbf{P}_{\mathbf{X}_R} = \mathbf{P}_{\mathbf{X}_R}\mathbf{P}_\mathbf{X} = \mathbf{P}_{\mathbf{X}_R}$.

i) $\text{rank}(\mathbf{P}_\mathbf{X}) = \text{tr}(\mathbf{P}_\mathbf{X}) = \text{rank}(\mathbf{X})$.

Note that \mathbf{P} is a projection matrix iff \mathbf{P} is symmetric and idempotent. Partition \mathbf{X} as $\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2]$, let \mathbf{P} be the projection matrix for $C(\mathbf{X})$ and let \mathbf{P}_1 be the projection matrix for $C(\mathbf{X}_1)$. Since $C(\mathbf{P}_1) = C(\mathbf{X}_1) \subseteq C(\mathbf{X})$, $\mathbf{P}\mathbf{P}_1 = \mathbf{P}_1$. Hence $\mathbf{P}_1\mathbf{P} = (\mathbf{P}\mathbf{P}_1)' = \mathbf{P}'_1 = \mathbf{P}_1$.

1a): Given small \mathbf{X} , be able to find the projection matrix \mathbf{P} for $C(\mathbf{X})$.

1b): Given small \mathbf{X} , be able to find $\text{rank}(\mathbf{X})$, a basis for $C(\mathbf{X})$, and $[C(\mathbf{X})]^\perp = \text{nullspace of } \mathbf{X}^T$.

1c): Be able to show that $\mathbf{G} := \mathbf{A}^-$.

2) Quadratic Forms $\mathbf{Y}'\mathbf{A}\mathbf{Y}$ and terms like $\mathbf{A}\mathbf{Y}$:

The matrix \mathbf{A} in a quadratic form $\mathbf{x}^T\mathbf{A}\mathbf{x}$ is **symmetric**. \mathbf{A} is **positive definite** ($\mathbf{A} > 0$) if $\mathbf{x}^T\mathbf{A}\mathbf{x} > 0 \ \forall \ \mathbf{x} \neq \mathbf{0}$. \mathbf{A} is **positive semidefinite** ($\mathbf{A} \geq 0$) if $\mathbf{x}^T\mathbf{A}\mathbf{x} \geq 0 \ \forall \ \mathbf{x}$.

Let \mathbf{A} be symmetric. If $\mathbf{A} \geq 0$ then the eigenvalues of \mathbf{A} are real and nonnegative. If $\mathbf{A} > 0$, let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$. If $\mathbf{A} > 0$, then $\lambda_n > 0$.

Theorem 2.5 (Seber and Lee Th. 1.5) expected value of a quadratic form: Let \mathbf{X} be a random vector with $E(\mathbf{X}) = \boldsymbol{\mu}$ and $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$. Then

$$E(\mathbf{X}^T\mathbf{A}\mathbf{X}) = \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) + [E(\mathbf{X})]^T\mathbf{A}E(\mathbf{X}) = \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) + \boldsymbol{\mu}^T\mathbf{A}\boldsymbol{\mu}.$$

Theorems 2.6 and 2.7: If $\mathbf{AY} \perp \mathbf{BY}$, then $f(\mathbf{AY}) \perp g(\mathbf{BY})$ where f and g are functions (such that $f(\mathbf{AY})$ only depends on \mathbf{A} and \mathbf{AY} and $g(\mathbf{BY})$ only depends on \mathbf{B} and \mathbf{BY}). Note that $\mathbf{Y}'\mathbf{AY} = \mathbf{Y}'\mathbf{A}'\mathbf{A}^-\mathbf{AY} = f(\mathbf{AY})$ (for a quadratic form \mathbf{A} is symmetric), $\mathbf{Y}'(\mathbf{I} - \mathbf{P})\mathbf{Y} = \|(\mathbf{I} - \mathbf{P})\mathbf{Y}\|^2$, and $\mathbf{Y}'\mathbf{PY} = \|\mathbf{PY}\|^2$ where the squared Euclidean norm $\|\mathbf{Z}\|^2 = \mathbf{Z}'\mathbf{Z}$.

Theorem 2.8. Let $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. a) Let $\mathbf{u} = \mathbf{AY}$ and $\mathbf{w} = \mathbf{BY}$. Then $\mathbf{AY} \perp \mathbf{BY}$ iff $\text{Cov}(\mathbf{u}, \mathbf{w}) = \mathbf{A}\boldsymbol{\Sigma}\mathbf{B}' = \mathbf{0}$ iff $\mathbf{B}\boldsymbol{\Sigma}\mathbf{A}' = \mathbf{0}$. Note that if $\boldsymbol{\Sigma} = \sigma^2\mathbf{I}_n$, then $\mathbf{AY} \perp \mathbf{BY}$ iff $\mathbf{AB}' = \mathbf{0}$ if $\mathbf{BA}' = \mathbf{0}$.

b) If \mathbf{A} is a symmetric $n \times n$ matrix, and \mathbf{B} is an $m \times n$ matrix, then $\mathbf{Y}'\mathbf{AY} \perp \mathbf{BY}$ iff $\mathbf{A}\boldsymbol{\Sigma}\mathbf{B}' = \mathbf{0}$ iff $\mathbf{B}\boldsymbol{\Sigma}\mathbf{A} = \mathbf{0}$.

Craig's Theorem: Let $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

a) If $\boldsymbol{\Sigma} > \mathbf{0}$, then $\mathbf{Y}'\mathbf{AY} \perp \mathbf{Y}'\mathbf{BY}$ iff $\mathbf{A}\boldsymbol{\Sigma}\mathbf{B} = \mathbf{0}$ iff $\mathbf{B}\boldsymbol{\Sigma}\mathbf{A} = \mathbf{0}$.

b) If $\boldsymbol{\Sigma} \geq \mathbf{0}$, then $\mathbf{Y}'\mathbf{AY} \perp \mathbf{Y}'\mathbf{BY}$ if $\mathbf{A}\boldsymbol{\Sigma}\mathbf{B} = \mathbf{0}$ (or if $\mathbf{B}\boldsymbol{\Sigma}\mathbf{A} = \mathbf{0}$).

c) If $\boldsymbol{\Sigma} \geq \mathbf{0}$, then $\mathbf{Y}'\mathbf{AY} \perp \mathbf{Y}'\mathbf{BY}$ iff

(*) $\boldsymbol{\Sigma}\mathbf{A}\boldsymbol{\Sigma}\mathbf{B}\boldsymbol{\Sigma} = \mathbf{0}$, $\boldsymbol{\Sigma}\mathbf{A}\boldsymbol{\Sigma}\mathbf{B}\boldsymbol{\mu} = \mathbf{0}$, $\boldsymbol{\Sigma}\mathbf{B}\boldsymbol{\Sigma}\mathbf{A}\boldsymbol{\mu} = \mathbf{0}$, and $\boldsymbol{\mu}'\mathbf{A}\boldsymbol{\Sigma}\mathbf{B}\boldsymbol{\mu} = 0$.

Note that if $\mathbf{A}\boldsymbol{\Sigma}\mathbf{B} = \mathbf{0}$, then (*) holds.

Theorem 2.13. If $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma} > \mathbf{0}$, then $\mathbf{Y}'\mathbf{AY} \sim \chi^2(\text{rank}(\mathbf{A}), \boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu}/2)$ iff $\mathbf{A}\boldsymbol{\Sigma}$ is idempotent.

Remark 1: If the theorem is for $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \mathbf{I})$ and $\mathbf{Z} \sim N_n(E(\mathbf{Z}), \sigma^2\mathbf{I})$, then use $\mathbf{Y} = \mathbf{Z}/\sigma \sim N_n(\boldsymbol{\mu} = E(\mathbf{Z})/\sigma, \mathbf{I})$.

Theorem 2.14. Let $\mathbf{A} = \mathbf{A}'$ be symmetric.

a) If $\mathbf{Y} \sim N_n(\mathbf{0}, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma}$ is a projection matrix, then $\mathbf{Y}'\mathbf{Y} \sim \chi^2(\text{rank}(\boldsymbol{\Sigma}))$ where $\text{rank}(\boldsymbol{\Sigma}) = \text{tr}(\boldsymbol{\Sigma})$.

b) If $\mathbf{Y} \sim N_n(\mathbf{0}, \mathbf{I})$, then $\mathbf{Y}'\mathbf{AY} \sim \chi_r^2$ iff \mathbf{A} is idempotent with $\text{rank}(\mathbf{A}) = \text{tr}(\mathbf{A}) = r$.

c) Let $\mathbf{Y} \sim N_n(\mathbf{0}, \sigma^2\mathbf{I})$. Then

$$\frac{\mathbf{Y}'\mathbf{AY}}{\sigma^2} \sim \chi_r^2 \quad \text{or} \quad \mathbf{Y}'\mathbf{AY} \sim \sigma^2 \chi_r^2$$

iff \mathbf{A} is idempotent of rank r .

d) If $\mathbf{Y} \sim N_n(\mathbf{0}, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma} > \mathbf{0}$, then $\mathbf{Y}'\mathbf{AY} \sim \chi_r^2$ iff $\mathbf{A}\boldsymbol{\Sigma}$ is idempotent with $\text{rank}(\mathbf{A}) = r = \text{rank}(\mathbf{A}\boldsymbol{\Sigma})$.

e) If $\mathbf{Y} \sim N_n(\mathbf{0}, \sigma^2\mathbf{I})$ then $\frac{\mathbf{Y}'\mathbf{Y}}{\sigma^2} \sim \chi^2\left(n, \frac{\boldsymbol{\mu}'\boldsymbol{\mu}}{2\sigma^2}\right)$.

f) If $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \mathbf{I})$ then $\mathbf{Y}'\mathbf{AY} \sim \chi^2(r, \boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu}/2)$ iff \mathbf{A} is idempotent with $\text{rank}(\mathbf{A}) = \text{tr}(\mathbf{A}) = r$.

g) If $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \sigma^2\mathbf{I})$ then $\frac{\mathbf{Y}'\mathbf{AY}}{\sigma^2} \sim \chi^2\left(r, \frac{\boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu}}{2\sigma^2}\right)$ iff \mathbf{A} is idempotent with $\text{rank}(\mathbf{A}) = \text{tr}(\mathbf{A}) = r$.

3) MLE: The following problem is typical. It is assumed that $\sigma > 0$ and $\boldsymbol{\beta} \in \mathbb{R}^p$.

Suppose $Y_i = \mathbf{x}_i'\boldsymbol{\beta} + \epsilon_i$ with $Q(\boldsymbol{\beta}) \geq 0$. Let c_n be a constant that does not depend on $\boldsymbol{\beta}$ or σ^2 . Suppose the likelihood function is

$$L(\boldsymbol{\beta}, \sigma^2) = c_n \frac{1}{\sigma^n} \exp\left(\frac{-1}{2\sigma^2}Q(\boldsymbol{\beta})\right).$$

a) Suppose that $\hat{\beta}_Q$ minimizes $Q(\beta)$. Show that $\hat{\beta}_Q$ is the MLE of β .

b) Then find the MLE $\hat{\sigma}^2$ of σ^2 .

Solution: a) For fixed $\sigma > 0$, $L(\beta, \sigma^2)$ is maximized by minimizing $Q(\beta) \geq 0$. So $\hat{\beta}_Q$ maximizes $L(\beta, \sigma^2)$ regardless of the value of $\sigma^2 > 0$. So $\hat{\beta}_Q$ is the MLE.

b) Let $Q = Q(\hat{\beta}_Q)$. Then the MLE $\hat{\sigma}^2$ is found by maximizing the profile likelihood, $L_p(\sigma^2) = L(\hat{\beta}_Q, \sigma^2) = c_n \frac{1}{\sigma^n} \exp\left(\frac{-1}{2\sigma^2}Q\right)$. Let $\tau = \sigma^2$. The $L_p(\tau) = c_n \frac{1}{\tau^{n/2}} \exp\left(\frac{-1}{2\tau}Q\right)$, and the log profile likelihood $\log L_p(\tau) = d - \frac{n}{2} \log(\tau) - \frac{Q}{2\tau}$. Thus

$$\frac{d \log L_p(\tau)}{d\tau} = \frac{-n}{2\tau} + \frac{Q}{2\tau^2} \stackrel{set}{=} 0$$

or $-n\tau + Q = 0$ or $\hat{\tau} = \hat{\sigma}^2 = Q/n$, unique. Then

$$\frac{d^2 \log L_p(\tau)}{d\tau^2} = \frac{n}{2\tau^2} - \frac{2Q}{2\tau^3} \Big|_{\hat{\tau}} = \frac{n}{2\tau^2} - \frac{2n\hat{\tau}}{2\hat{\tau}^3} = \frac{-n}{2\hat{\tau}^2} < 0$$

which proves that $\hat{\sigma}^2$ is the MLE of σ^2 .

Note: A negative second derivative shows that $\hat{\sigma}^2$ is a local max. The result that $\hat{\sigma}^2$ was the unique solution to setting the first derivative of the profile likelihood equal to zero makes $\hat{\sigma}^2$ the global max.

Common errors: Students use $Q(\beta)$ instead of $Q(\hat{\beta})$ in the profile likelihood. Students forget to write the word “unique.”

Variant: $Q(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|^2 = (\mathbf{Y} - \mathbf{X}\beta)^T(\mathbf{Y} - \mathbf{X}\beta)$ is the least squares criterion. Recognize that $Q(\beta)$ is minimized by $\hat{\beta} = \hat{\beta}_{OLS}$, and proceed as in the above problem.

Note: If the e_i are iid $N(0, \sigma^2)$ and least squares is used, then the MLE of β is the least squares estimator $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ and the MLE of σ^2 is

$$\hat{\sigma}_M^2 = \frac{n-p}{n} MSE = \frac{1}{n} \sum_{i=1}^n r_i^2.$$

4) LS Estimators for $p \leq 2$:

Given a least squares model with $p \leq 2$, derive or find the least squares estimator $\hat{\beta}$.

Tip: If the LS model is $Y_i = \mathbf{x}_i^T \beta + e_i$ for $i = 1, \dots, n$, then the LS criterion is $Q(\beta) = \sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2 = \sum_{i=1}^n r_i^2(\beta)$.

To derive the LS estimator, let $Q(\beta_1, \beta_2) = \sum_{i=1}^n (Y_i - \beta_1 - \beta_2 x_i)^2$ be the residual sum of squares where β_i vary on \mathbb{R} . Take the partial derivatives, set them to 0, and solve for the least squares estimators. If $p = 2$, we will assume 2nd derivatives do not need to be taken. If $p = 1$, show the solution is unique and show that the second derivative evaluated at $\hat{\beta}$ is positive. The β_i could be replaced by other symbols such as η_i .

Location model: $Y_i = \beta + e_i$ or $\mathbf{Y} = \mathbf{1}\beta + \mathbf{e}$. The parameter β could be replaced with μ or θ . The LS criterion $Q(\beta) = \sum_{i=1}^n (Y_i - \beta)^2$, and $\hat{\beta} = \bar{Y}$, the sample mean.

$$\text{Proof : } \frac{dQ(\beta)}{d\beta} = -2 \sum_{i=1}^n (Y_i - \beta).$$

Setting the derivative equal to 0 and calling the unique solution $\hat{\beta}$ gives $\sum_{i=1}^n Y_i = n\hat{\beta}$ or $\hat{\beta} = \bar{Y}$. The second derivative

$$\frac{d^2Q(\beta)}{d\beta^2} = 2n > 0,$$

hence $\hat{\beta}$ is the global minimizer.

Simple linear regression (SLR): $Y_i = \beta_1 + x_i\beta_2 + e_i$ or $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ where $\mathbf{X} = [\mathbf{1} \ \mathbf{x}]$ and $\boldsymbol{\beta} = (\beta_1 \ \beta_2)^T$. The LS criterion $Q(\beta_1, \beta_2) = \sum_{i=1}^n (Y_i - \beta_1 - x_i\beta_2)^2$.

The least squares (OLS) line is $\hat{Y} = \hat{\beta}_1 + \hat{\beta}_2 X$ where the slope

$$\hat{\beta}_2 \equiv \hat{\beta} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})Y_i}{\sum_{j=1}^n (X_j - \bar{X})^2} = \sum_{i=1}^n k_i Y_i$$

with

$$k_i = \frac{X_i - \bar{X}}{\sum_{j=1}^n (X_j - \bar{X})^2} = \frac{X_i - \bar{X}}{(n-1)S_X^2},$$

and the intercept $\hat{\beta}_1 \equiv \hat{\alpha} = \bar{Y} - \hat{\beta}_2 \bar{X}$.

By the **chain rule**,

$$\frac{\partial Q}{\partial \beta_1} = -2 \sum_{i=1}^n (Y_i - \beta_1 - \beta_2 X_i)$$

and

$$\frac{\partial^2 Q}{\partial \beta_1^2} = 2n.$$

Similarly,

$$\frac{\partial Q}{\partial \beta_2} = -2 \sum_{i=1}^n X_i (Y_i - \beta_1 - \beta_2 X_i)$$

and

$$\frac{\partial^2 Q}{\partial \beta_2^2} = 2 \sum_{i=1}^n X_i^2.$$

Setting the first partial derivatives to zero and calling the solutions $\hat{\beta}_1$ and $\hat{\beta}_2$ shows that the OLS estimators $\hat{\beta}_1$ and $\hat{\beta}_2$ satisfy the **normal equations**:

$$\sum_{i=1}^n Y_i = n\hat{\beta}_1 + \hat{\beta}_2 \sum_{i=1}^n X_i \quad \text{and}$$

$$\sum_{i=1}^n X_i Y_i = \hat{\beta}_1 \sum_{i=1}^n X_i + \hat{\beta}_2 \sum_{i=1}^n X_i^2.$$

The first equation gives $\hat{\beta}_1 = \bar{Y} - \hat{\beta}_2 \bar{X}$.

There are several equivalent formulas for the slope $\hat{\beta}_2$.

$$\hat{\beta}_2 \equiv \hat{\beta} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n X_i Y_i - \frac{1}{n}(\sum_{i=1}^n X_i)(\sum_{i=1}^n Y_i)}{\sum_{i=1}^n X_i^2 - \frac{1}{n}(\sum_{i=1}^n X_i)^2}$$

$$= \frac{\sum_{i=1}^n (X_i - \bar{X})Y_i}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n X_i Y_i - n\bar{X}\bar{Y}}{\sum_{i=1}^n X_i^2 - n(\bar{X})^2} = \hat{\rho}_{s_Y}/s_X.$$

Here the sample correlation $\hat{\rho} \equiv \hat{\rho}(X, Y) = \text{corr}(X, Y) =$

$$\frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(n-1)s_X s_Y} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

where the sample standard deviation

$$s_W = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (W_i - \bar{W})^2}$$

for $W = X, Y$. Notice that the term $n-1$ that occurs in the denominator of $\hat{\rho}$, s_Y^2 , and s_X^2 can be replaced by n as long as n is used in all 3 quantities.

SLR through the origin: $Y_i = x_i\beta + e_i$ or $Y = \mathbf{x}\beta + \mathbf{e}$. The LS criterion $Q(\beta) = \sum_{i=1}^n (Y_i - x_i\beta)^2$, and $\hat{\beta} = \sum_{i=1}^n x_i Y_i / \sum_{i=1}^n x_i^2$.

Known intercept: $Y_i = a + x_i\beta + e_i$ where the intercept a is known.
 $Q(\beta) = \sum_{i=1}^n (Y_i - a - x_i\beta)^2$.

Known slope: $Y_i = \beta + x_i b + e_i$ where the slope b is known.
 $Q(\beta) = \sum_{i=1}^n (Y_i - \beta - x_i b)^2$. Here, β may be replaced by α .

5) WLS:

For the WLS model $Y|\mathbf{x} = \mathbf{x}^T\boldsymbol{\beta} + e$ where the e_i are independent with $E(e_i) = 0$ and $V(e_i) = \sigma_i^2$. Hence $\mathbf{Y} = \mathbf{Y}|\mathbf{X} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ where $E(\mathbf{e}) = \mathbf{0}$ and $\text{Cov}(\mathbf{e}) = \text{diag}(\sigma_i^2)$.

An alternative model is $Y|\mathbf{x}^T\boldsymbol{\beta} = \mathbf{x}^T\boldsymbol{\beta} + \mathbf{u}$ where the u_i are independent with $E(u_i) = 0$ and $V(u_i) = \tau_i^2$. Hence $\mathbf{Y} = \mathbf{Y}|\mathbf{X}\boldsymbol{\beta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ where $E(\mathbf{u}) = \mathbf{0}$ and $\text{Cov}(\mathbf{u}) = \text{diag}(\tau_i^2)$.

6) Non-full rank linear models:

The **nonfull rank linear model** is $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ where \mathbf{X} has rank $r < p \leq n$, and \mathbf{X} is an $n \times p$ matrix.

Theorem 3.1. i) $\mathbf{P} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^- \mathbf{X}^T$ is the unique projection matrix on $C(\mathbf{X})$ and does not depend on the generalized inverse $(\mathbf{X}^T\mathbf{X})^-$.

ii) $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^- \mathbf{X}^T\mathbf{Y}$ does depend on $(\mathbf{X}^T\mathbf{X})^-$ and is not unique.

iii) $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{P}\mathbf{Y}$, $\mathbf{r} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I} - \mathbf{P})\mathbf{Y}$ and $RSS = \mathbf{r}^T\mathbf{r}$ are unique and so do not depend on $(\mathbf{X}^T\mathbf{X})^-$.

iv) $\hat{\boldsymbol{\beta}}$ is a solution to the *normal equations*: $\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y}$.

v) $\text{Rank}(\mathbf{P}) = r$ and $\text{rank}(\mathbf{I} - \mathbf{P}) = n - r$.

vi) If $\text{Cov}(\mathbf{Y}) = \text{Cov}(\mathbf{e}) = \sigma^2\mathbf{I}$, then $MSE = \frac{RSS}{n-r} = \frac{\mathbf{r}^T\mathbf{r}}{n-r}$ is an unbiased estimator of σ^2 .

vii) Let the columns of \mathbf{X}_1 form a basis for $C(\mathbf{X})$. For example, take r linearly independent columns of \mathbf{X} to form \mathbf{X}_1 . Then $\mathbf{P} = \mathbf{X}_1(\mathbf{X}_1^T\mathbf{X}_1)^{-1}\mathbf{X}_1^T$.

7) Estimability and the Gauss Markov Theorem:

Let \mathbf{a} and \mathbf{b} be constant vectors. Then $\mathbf{a}^T\boldsymbol{\beta}$ is **estimable** if there exists a linear unbiased estimator $\mathbf{b}^T\mathbf{Y}$ so $E(\mathbf{b}^T\mathbf{Y}) = \mathbf{a}^T\boldsymbol{\beta}$. Also, $\mathbf{a}^T\boldsymbol{\beta}$ is estimable iff $\mathbf{a}^T = \mathbf{b}^T\mathbf{X}$ iff $\mathbf{a} = \mathbf{X}^T\mathbf{b}$ iff $\mathbf{a} \in C(\mathbf{X}^T)$.

The linear estimator $\mathbf{a}^T\mathbf{Y}$ of $\mathbf{c}^T\boldsymbol{\theta}$ is the best linear unbiased estimator (BLUE) of $\mathbf{c}^T\boldsymbol{\theta}$ if $E(\mathbf{a}^T\mathbf{Y}) = \mathbf{c}^T\boldsymbol{\theta}$, and if for any other unbiased linear estimator $\mathbf{b}^T\mathbf{Y}$ of $\mathbf{c}^T\boldsymbol{\theta}$, $V(\mathbf{a}^T\mathbf{Y}) \leq V(\mathbf{b}^T\mathbf{Y})$. Note that $E(\mathbf{b}^T\mathbf{Y}) = \mathbf{c}^T\boldsymbol{\theta}$.

The next theorem shows that the least squares estimator of an estimable function $\mathbf{a}^T\boldsymbol{\beta}$ is $\mathbf{a}^T\hat{\boldsymbol{\beta}} = \mathbf{b}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{b}^T\mathbf{P}\mathbf{Y}$. Note that $\mathbf{b}^T\mathbf{Y}$ is also an unbiased estimator of $\mathbf{a}^T\boldsymbol{\beta}$ since $E(\mathbf{b}^T\mathbf{Y}) = \mathbf{b}^T\mathbf{X}\boldsymbol{\beta} = \mathbf{a}^T\boldsymbol{\beta}$.

Theorem 3.2 (see Seber and Lee Th 3.2) Let $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ where \mathbf{X} has rank $r \leq p \leq n$, $E(\mathbf{e}) = \mathbf{0}$, and $\text{Cov}(\mathbf{e}) = \sigma^2\mathbf{I}$.

a) The quantity $\mathbf{a}^T\boldsymbol{\beta}$ is estimable iff $\mathbf{a}^T = \mathbf{b}^T\mathbf{X}$ iff $\mathbf{a} = \mathbf{X}^T\mathbf{b}$ (for some constant vector \mathbf{b}) iff $\mathbf{a} \in C(\mathbf{X}^T)$.

b) Let $\hat{\boldsymbol{\theta}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ and $\boldsymbol{\theta} = \mathbf{X}\boldsymbol{\beta}$. Suppose there exists a constant vector \mathbf{c} such that $E(\mathbf{c}^T\hat{\boldsymbol{\theta}}) = \mathbf{c}^T\boldsymbol{\theta}$. Then among the class of linear unbiased estimators of $\mathbf{c}^T\boldsymbol{\theta}$, the least squares estimator $\mathbf{c}^T\hat{\boldsymbol{\theta}}$ is the unique BLUE.

c) **Gauss Markov Theorem:** If $\mathbf{a}^T\boldsymbol{\beta}$ is estimable and a least squares estimator $\hat{\boldsymbol{\beta}}$ is any solution to the normal equations $\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^T\mathbf{Y}$, then $\mathbf{a}^T\hat{\boldsymbol{\beta}}$ is the unique BLUE of $\mathbf{a}^T\boldsymbol{\beta}$.

Proof: a) If $\mathbf{a}^T\boldsymbol{\beta}$ is estimable, then $\mathbf{a}^T\boldsymbol{\beta} = E(\mathbf{b}^T\mathbf{Y}) = \mathbf{b}^T\mathbf{X}\boldsymbol{\beta}$ for all $\boldsymbol{\beta} \in \mathbb{R}^p$. Thus $\mathbf{a}^T = \mathbf{b}^T\mathbf{X}$ or $\mathbf{a} = \mathbf{X}^T\mathbf{b}$. Hence $\mathbf{a}^T\boldsymbol{\beta}$ is estimable iff $\mathbf{a}^T = \mathbf{b}^T\mathbf{X}$ iff $\mathbf{a} = \mathbf{X}^T\mathbf{b}$ iff $\mathbf{a} \in C(\mathbf{X}^T)$.

b) Since $\hat{\boldsymbol{\theta}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{P}\mathbf{Y}$, it follows that $E(\mathbf{c}^T\hat{\boldsymbol{\theta}}) = E(\mathbf{c}^T\mathbf{P}\mathbf{Y}) = \mathbf{c}^T\mathbf{P}\mathbf{X}\boldsymbol{\beta} = \mathbf{c}^T\mathbf{X}\boldsymbol{\beta} = \mathbf{c}^T\boldsymbol{\theta}$. Thus $\mathbf{c}^T\hat{\boldsymbol{\theta}} = \mathbf{c}^T\mathbf{P}\mathbf{Y} = (\mathbf{P}\mathbf{c})^T\mathbf{Y}$ is a linear unbiased estimator of $\mathbf{c}^T\boldsymbol{\theta}$. Let $\mathbf{d}^T\mathbf{Y}$ be any other linear unbiased estimator of $\mathbf{c}^T\boldsymbol{\theta}$. Hence $E(\mathbf{d}^T\mathbf{Y}) = \mathbf{d}^T\boldsymbol{\theta} = \mathbf{c}^T\boldsymbol{\theta}$ for all $\boldsymbol{\theta} \in C(\mathbf{X})$. So $(\mathbf{c} - \mathbf{d})^T\boldsymbol{\theta} = 0$ for all $\boldsymbol{\theta} \in C(\mathbf{X})$. Hence $(\mathbf{c} - \mathbf{d}) \in [C(\mathbf{X})]^\perp$ and $\mathbf{P}(\mathbf{c} - \mathbf{d}) = \mathbf{0}$, or $\mathbf{P}\mathbf{c} = \mathbf{P}\mathbf{d}$. Thus $V(\mathbf{c}^T\hat{\boldsymbol{\theta}}) = V(\mathbf{c}^T\mathbf{P}\mathbf{Y}) = V(\mathbf{d}^T\mathbf{P}\mathbf{Y}) = \sigma^2\mathbf{d}^T\mathbf{P}^T\mathbf{P}\mathbf{d} = \sigma^2\mathbf{d}^T\mathbf{P}\mathbf{d}$. Then $V(\mathbf{d}^T\mathbf{Y}) - V(\mathbf{c}^T\hat{\boldsymbol{\theta}}) = V(\mathbf{d}^T\mathbf{Y}) - V(\mathbf{d}^T\mathbf{P}\mathbf{Y}) = \sigma^2[\mathbf{d}^T\mathbf{d} - \mathbf{d}^T\mathbf{P}\mathbf{d}] = \sigma^2\mathbf{d}^T(\mathbf{I}_n - \mathbf{P})\mathbf{d} = \sigma^2\mathbf{d}^T(\mathbf{I}_n - \mathbf{P})^T(\mathbf{I}_n - \mathbf{P})\mathbf{d} = \mathbf{g}^T\mathbf{g} \geq 0$ with equality iff $\mathbf{g} = (\mathbf{I}_n - \mathbf{P})\mathbf{d} = \mathbf{0}$, or $\mathbf{d} = \mathbf{P}\mathbf{d} = \mathbf{P}\mathbf{c}$. Thus $\mathbf{c}^T\hat{\boldsymbol{\theta}}$ has minimum variance and is unique.

c) Since $\mathbf{a}^T\boldsymbol{\beta}$ is estimable, $\mathbf{a}^T\hat{\boldsymbol{\beta}} = \mathbf{b}^T\mathbf{X}\hat{\boldsymbol{\beta}}$. Then $\mathbf{a}^T\hat{\boldsymbol{\beta}} = \mathbf{b}^T\hat{\boldsymbol{\theta}}$ is the unique BLUE of $\mathbf{a}^T\boldsymbol{\beta} = \mathbf{b}^T\boldsymbol{\theta}$ by b).

Gauss Markov Theorem-Full Rank Case: Let $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ where \mathbf{X} is full rank, $E(\mathbf{e}) = \mathbf{0}$, and $\text{Cov}(\mathbf{e}) = \sigma^2\mathbf{I}$. Then $\mathbf{a}^T\hat{\boldsymbol{\beta}}$ is the unique BLUE of $\mathbf{a}^T\boldsymbol{\beta}$ for every constant $p \times 1$ vector \mathbf{a} .

Notation: $\boldsymbol{\beta}$ is “estimable” by $\hat{\boldsymbol{\beta}}$ for the full rank model, but not for the non-full rank model.

8) Hypothesis Testing:

Theorem 2.16. Let $\boldsymbol{\theta} = \mathbf{X}\boldsymbol{\eta} \in C(\mathbf{X})$ where $Y_i = \mathbf{x}_i^T\boldsymbol{\eta} + r_i(\boldsymbol{\eta})$ and the residual $r_i(\boldsymbol{\eta})$ depends on $\boldsymbol{\eta}$. The **least squares estimator** $\hat{\boldsymbol{\beta}}$ is the value of $\boldsymbol{\eta} \in \mathbb{R}^p$ that minimizes the **least squares criterion**

$$\sum_{i=1}^n r_i^2(\boldsymbol{\eta}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\eta}\|^2.$$

LS CLT (Least Squares Central Limit Theorem): Consider the MLR model $Y_i = \mathbf{x}_i^T\boldsymbol{\beta} + e_i$ and assume that the zero mean errors are iid with $E(e_i) = 0$ and $\text{VAR}(e_i) =$

σ^2 . Also assume that $\max_i(h_1, \dots, h_n) \rightarrow 0$ in probability as $n \rightarrow \infty$ and

$$\frac{\mathbf{X}^T \mathbf{X}}{n} \rightarrow \mathbf{W}^{-1}$$

as $n \rightarrow \infty$. Then the least squares (OLS) estimator $\hat{\boldsymbol{\beta}}$ satisfies

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \sigma^2 \mathbf{W}). \quad (1)$$

Partial F Test Theorem: Suppose $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{0}$ is true for the partial F test where \mathbf{L} is a full rank $r \times p$ matrix. Under the OLS full rank model, a)

$$F_R = \frac{1}{rMSE} (\mathbf{L}\hat{\boldsymbol{\beta}})^T [\mathbf{L}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} (\mathbf{L}\hat{\boldsymbol{\beta}}).$$

b) If $\mathbf{e} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$, then $F_R \sim F_{r, n-p}$.

c) For a large class of zero mean error distributions $rF_R \xrightarrow{D} \chi_r^2$.

d) The partial F test that rejects $H_0 : \mathbf{L}\boldsymbol{\beta} = \mathbf{0}$ if $F_R > F_{r, n-p}(1 - \delta)$ is a large sample right tail δ test for the OLS model for a large class of zero mean error distributions.

Assume H_0 is true. By the OLS CLT, $\sqrt{n}(\mathbf{L}\hat{\boldsymbol{\beta}} - \mathbf{L}\boldsymbol{\beta}) = \sqrt{n}\mathbf{L}\hat{\boldsymbol{\beta}} \xrightarrow{D} N_r(\mathbf{0}, \sigma^2 \mathbf{L}\mathbf{W}\mathbf{L}^T)$. Thus $\sqrt{n}(\mathbf{L}\hat{\boldsymbol{\beta}})^T (\sigma^2 \mathbf{L}\mathbf{W}\mathbf{L}^T)^{-1} \sqrt{n}\mathbf{L}\hat{\boldsymbol{\beta}} \xrightarrow{D} \chi_r^2$. Let $\hat{\sigma}^2 = MSE$ and $\hat{\mathbf{W}} = n(\mathbf{X}^T \mathbf{X})^{-1}$. Then

$$n(\mathbf{L}\hat{\boldsymbol{\beta}})^T [MSE \mathbf{L}n(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{L}^T]^{-1} \mathbf{L}\hat{\boldsymbol{\beta}} = rF_R \xrightarrow{D} \chi_r^2.$$

Partial F test: Let the full model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ with a constant β_1 in the model: $\mathbf{1}$ is the 1st column of \mathbf{X} . Let the reduced model $\mathbf{Y} = \mathbf{X}_R \boldsymbol{\beta}_R + \mathbf{e}$ also have a constant in the model where the columns of \mathbf{X}_R are a subset of k of the columns of \mathbf{X} . Let \mathbf{P}_R be the projection matrix on $C(\mathbf{X}_R)$ so $\mathbf{P}\mathbf{P}_R = \mathbf{P}_R$. Then $F_R = \frac{SSE(R) - SSE(F)}{rMSE(F)}$ where $r = df_R - df_F = p - k =$ number of predictors in the full model but not in the reduced model. $MSE = MSE(F) = SSE(F)/(n - p)$ where $SSE = SSE(F) = \mathbf{Y}(\mathbf{I} - \mathbf{P})\mathbf{Y}$. $SSE(R) - SSE(F) = \mathbf{Y}^T(\mathbf{P} - \mathbf{P}_R)\mathbf{Y}$ where $SSE(R) = \mathbf{Y}^T(\mathbf{I} - \mathbf{P}_R)\mathbf{Y}$.

Now assume $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$, and when H_0 is true, $\mathbf{Y} \sim N_n(\mathbf{X}_R \boldsymbol{\beta}_R, \sigma^2 \mathbf{I})$. Since $(\mathbf{I} - \mathbf{P})(\mathbf{P} - \mathbf{P}_R) = \mathbf{0}$, $[SSE(R) - SSE(F)] \perp MSE(F)$ by Craig's Theorem. When H_0 is true, $\boldsymbol{\mu} = \mathbf{X}_R \boldsymbol{\beta}_R$ and $\boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} = 0$ where $\mathbf{A} = (\mathbf{I} - \mathbf{P})$ or $\mathbf{A} = (\mathbf{P} - \mathbf{P}_R)$. Hence the noncentrality parameter is 0, and by Theorem 2.14 g), $SSE \sim \sigma^2 \chi_{n-p}^2$ and $SSE(R) - SSE(F) \sim \sigma^2 \chi_{p-k}^2$ since $rank(\mathbf{P} - \mathbf{P}_R) = tr(\mathbf{P} - \mathbf{P}_R) = p - k$. Hence under H_0 , $F_R \sim F_{p-k, n-p}$.

An ANOVA table for the partial F test is shown below, where $k = p_R$ is the number of predictors used by the reduced model, and $r = p - p_R = p - k$ is the number of predictors in the full model that are not in the reduced model.

Source	df	SS	MS	F
Reduced	$n - p_R$	$SSE(R) = \mathbf{Y}^T(\mathbf{I} - \mathbf{P}_R)\mathbf{Y}$	$MSE(R)$	$F_R = \frac{SSE(R) - SSE}{rMSE} =$
Full	$n - p$	$SSE = \mathbf{Y}^T(\mathbf{I} - \mathbf{P})\mathbf{Y}$	MSE	$\frac{\mathbf{Y}^T(\mathbf{P} - \mathbf{P}_R)\mathbf{Y}/r}{\mathbf{Y}^T(\mathbf{I} - \mathbf{P})\mathbf{Y}/(n - p)}$

The ANOVA F test is the special case where $k = 1$, $\mathbf{X}_R = \mathbf{1}$, $\mathbf{P}_R = \mathbf{P}_1$, and $SSE(R) - SSE(F) = SSTO - SSE = SSR$.

ANOVA table: $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ with a constant β_1 in the model: $\mathbf{1}$ is the 1st column of \mathbf{X} . $MS = SS/df$.

$SSTO = \mathbf{Y}^T(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^T)\mathbf{Y} = \sum_{i=1}^n (Y_i - \bar{Y})^2$, $SSE = \sum_{i=1}^n r_i^2$, $SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$,
 $SSTO = SSR + SSE$. SSTO is the SSE (residual sum of squares) for the location model $\mathbf{Y} = \mathbf{1}\beta_1 + \mathbf{e}$ that contains a constant but no nontrivial predictors. The location model has projection matrix $\mathbf{P}_1 = \mathbf{1}(\mathbf{1}^T\mathbf{1})^{-1}\mathbf{1}^T = \frac{1}{n}\mathbf{1}\mathbf{1}^T$. Hence $\mathbf{P}\mathbf{P}_1 = \mathbf{P}_1$ and $\mathbf{P}_1\mathbf{1} = \mathbf{1}$.

Source	df	SS	MS	F	p-value
Regression	p-1	$SSR = \mathbf{Y}^T(\mathbf{P} - \frac{1}{n}\mathbf{1}\mathbf{1}^T)\mathbf{Y}$	MSR	$F_0 = MSR/MSE$	for H_0 :
Residual	n-p	$SSE = \mathbf{Y}^T(\mathbf{I} - \mathbf{P})\mathbf{Y}$	MSE		$\beta_2 = \dots = \beta_p = 0$

The matrices in the quadratic forms for SSR and SSE are symmetric and idempotent and their product is $\mathbf{0}$. Hence if $\mathbf{e} \sim N_n(\mathbf{0}, \sigma^2\mathbf{I})$ so $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$, then $SSE \perp SSR$ by Craig's Theorem. If H_0 is true under normality, then $\mathbf{Y} \sim N_n(\mathbf{1}\beta_1, \sigma^2\mathbf{I})$, and by Theorem 2.14 g), $SSE \sim \sigma^2\chi_{n-p}^2$ and $SSR \sim \sigma^2\chi_{p-1}^2$ since $rank(\mathbf{I} - \mathbf{P}) = tr(\mathbf{I} - \mathbf{P}) = n-p$ and $rank(\mathbf{P} - \frac{1}{n}\mathbf{1}\mathbf{1}^T) = tr(\mathbf{P} - \frac{1}{n}\mathbf{1}\mathbf{1}^T) = p-1$. Hence under normality, $F_0 \sim F_{p-1, n-p}$.

9) Expected Value, Covariance Matrix and Large Sample Theory for least squares quantities:

For the full rank model, $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ where $E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$, $E(\mathbf{e}) = \mathbf{0}$ and $Cov(\mathbf{e}) = Cov(\mathbf{Y}) = \sigma^2\mathbf{I}$, $E(\mathbf{A}\mathbf{Y}) = \mathbf{A}\mathbf{X}\boldsymbol{\beta}$ and $Cov(\mathbf{A}\mathbf{Y}) = \sigma^2\mathbf{A}\mathbf{A}^T$.

$\mathbf{A} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ is used for $\hat{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y}$. $\mathbf{A} = \mathbf{I} - \mathbf{P} = \mathbf{I} - \mathbf{H}$ is used for the residual vector $\mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{A}\mathbf{Y}$. $\mathbf{A} = \mathbf{P} = \mathbf{H}$ is used for the vector of fitted values $\hat{\mathbf{Y}}$.

For the full rank Gaussian linear model, $\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$, and if \mathbf{A} is $k \times n$ with rank k , then $\mathbf{A}\mathbf{Y} \sim N_k(\mathbf{A}\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{A}\mathbf{A}^T)$.

If $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, \sigma^2\mathbf{W})$, and \mathbf{A} is $k \times p$ with rank k , then $\sqrt{n}(\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{A}\boldsymbol{\beta}) \xrightarrow{D} N_k(\mathbf{0}, \sigma^2\mathbf{A}\mathbf{W}\mathbf{A}^T)$.

The non-full rank model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ also has $E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$, $E(\mathbf{e}) = \mathbf{0}$, $Cov(\mathbf{e}) = Cov(\mathbf{Y}) = \sigma^2\mathbf{I}$, $E(\mathbf{A}\mathbf{Y}) = \mathbf{A}\mathbf{X}\boldsymbol{\beta}$ and $Cov(\mathbf{A}\mathbf{Y}) = \sigma^2\mathbf{A}\mathbf{A}^T$.

For the non-full rank model $\mathbf{A} = (\mathbf{X}^T\mathbf{X})^{-}\mathbf{X}^T$ is used for $\hat{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y}$ and $\mathbf{P} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-}\mathbf{X}^T$.

You should be able to handle the linear model written in different ways. The residual bootstrap model $\mathbf{Y}^* = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e}^*$ with $E(\mathbf{e}^*) = \mathbf{0}$ and $Cov(\mathbf{e}^*) = Cov(\mathbf{Y}^*) = \hat{\sigma}^2\mathbf{I}$. The parametric bootstrap model $\mathbf{Y}^* = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e}^*$ with $\mathbf{Y}^* \sim N_n(\mathbf{X}\hat{\boldsymbol{\beta}}, MSE\mathbf{I})$. In numerical linear algebra, the least squares solution to " $\mathbf{A}\mathbf{x} = \mathbf{b}$ " is of interest where the problem is actually the multiple linear regression model $\mathbf{b} = \mathbf{A}\mathbf{x} + \boldsymbol{\epsilon}$ where \mathbf{A} has full rank p , and we will assume that $E(\boldsymbol{\epsilon}) = \mathbf{0}$, and $Cov(\boldsymbol{\epsilon}) = \sigma^2\mathbf{I}_n$.

References:

Christensen, R. (2020), *Plane Answers to Complex Questions: the Theory of Linear Models*, 5th ed., Springer, New York, NY.

Graybill, F.A. (2000), *Theory and Application of the Linear Model*, Brooks/Cole, Pacific Grove, CA.

Olive, D.J. (2021), *Theory of Linear Models*, online course notes, see (<http://parker.ad.siu.edu/Olive/linmodbk.htm>).

Rencher, A.C., and Schaalje, G.B. (2008), *Linear Models in Statistics*, 2nd ed., Wiley, Hoboken, NJ.

Searle, S.R., and Gruber, M.H.J. (2017), *Linear Models*, 2nd ed., Wiley, Hoboken, NJ.

Seber, G.A.F., and Lee, A.J. (2003), *Linear Regression Analysis*, 2nd ed., Wiley, New York, NY.

Zimmerman, D.L. (2020), *Linear Model Theory: Exercises and Solutions*, Springer, New York, NY.