

Linear Models: Estimation and Inference under Normality

Lecture Notes Supplement

Your Name/Course Name Here

April 27, 2025

1 Estimating Linear Combinations of Coefficients

In our study of the linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, we've focused on estimating the entire vector of coefficients $\boldsymbol{\beta} \in \mathbb{R}^{p+1}$. Often, however, we are interested in estimating a specific linear combination of these coefficients.

Definition 1.1 (Linear Combination). Let $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^\top$ be the vector of regression coefficients, and let $\mathbf{a} = (a_0, a_1, \dots, a_p)^\top \in \mathbb{R}^{p+1}$ be a fixed, known vector of constants. The scalar quantity

$$\theta := \mathbf{a}^\top \boldsymbol{\beta} = \sum_{j=0}^p a_j \beta_j \quad (1)$$

is called a **linear combination** of the elements of $\boldsymbol{\beta}$.

Remark 1.2 (Motivation). Why study linear combinations?

- **Individual Coefficients:** Estimating β_j itself corresponds to choosing \mathbf{a} with $a_j = 1$ and all other $a_k = 0$.
- **Differences/Contrasts:** Comparing two coefficients, say $\beta_1 - \beta_2$, corresponds to $a_1 = 1, a_2 = -1$, and others zero.
- **Predictions:** For a new set of predictor values $\mathbf{x}_* = (1, x_{*1}, \dots, x_{*p})^\top$, the expected response is $\mathbb{E}[Y_*] = \mathbf{x}_*^\top \boldsymbol{\beta}$, which is a linear combination with $\mathbf{a} = \mathbf{x}_*$.

Given the least squares (LS) estimator $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ for $\boldsymbol{\beta}$, a natural estimator for $\theta = \mathbf{a}^\top \boldsymbol{\beta}$ is obtained by simply plugging in $\hat{\boldsymbol{\beta}}$:

Definition 1.3 (LS Estimator of a Linear Combination). The LS estimator for $\theta = \mathbf{a}^\top \boldsymbol{\beta}$ is

$$\hat{\theta} := \mathbf{a}^\top \hat{\boldsymbol{\beta}} = \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}. \quad (2)$$

Notice that $\hat{\theta}$ is a linear function of the observation vector \mathbf{Y} . We can write it as $\hat{\theta} = \mathbf{c}^\top \mathbf{Y}$, where the vector $\mathbf{c} \in \mathbb{R}^n$ is given by

$$\mathbf{c} := \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}. \quad (3)$$

It's worth noting the dimensions: \mathbf{a} is in \mathbb{R}^{p+1} , while \mathbf{c} is in \mathbb{R}^n .

Let's examine the properties of this estimator under the standard **linear model assumptions**:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \text{where } \mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0} \text{ and } \text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n. \quad (4)$$

These assumptions imply $\mathbb{E}[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta}$ and $\text{Cov}(\mathbf{Y}) = \sigma^2 \mathbf{I}_n$.

Mean of $\hat{\theta}$: Using the linearity of expectation and the unbiasedness of $\hat{\beta}$ ($\mathbb{E}[\hat{\beta}] = \beta$), we have:

$$\mathbb{E}[\hat{\theta}] = \mathbb{E}[\mathbf{a}^\top \hat{\beta}] = \mathbf{a}^\top \mathbb{E}[\hat{\beta}] = \mathbf{a}^\top \beta = \theta.$$

Thus, $\hat{\theta}$ is an **unbiased** estimator for θ .

Variance of $\hat{\theta}$: Using the linearity property of variance for vector transformations ($V(\mathbf{A}\mathbf{Y}) = \mathbf{A}\text{Cov}(\mathbf{Y})\mathbf{A}^\top$) or directly for the linear combination $\hat{\theta} = \mathbf{c}^\top \mathbf{Y}$:

$$\begin{aligned}\text{Var}(\hat{\theta}) &= \text{Var}(\mathbf{c}^\top \mathbf{Y}) \\ &= \mathbf{c}^\top \text{Cov}(\mathbf{Y}) \mathbf{c} \\ &= \mathbf{c}^\top (\sigma^2 \mathbf{I}_n) \mathbf{c} \\ &= \sigma^2 \mathbf{c}^\top \mathbf{c}.\end{aligned}$$

We can express this variance in terms of \mathbf{a} as well:

$$\begin{aligned}\text{Var}(\hat{\theta}) &= \sigma^2 \mathbf{c}^\top \mathbf{c} = \sigma^2 \left(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} \right)^\top \left(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} \right) \\ &= \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} \\ &= \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}.\end{aligned}$$

So, $\hat{\theta}$ is a linear, unbiased estimator with variance $\sigma^2 \mathbf{c}^\top \mathbf{c} = \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}$.

Is $\hat{\theta}$ the "Best" Linear Unbiased Estimator? We found *an* unbiased linear estimator $\hat{\theta}$. But could there be another estimator, say $\tilde{\theta}$, which is also linear in \mathbf{Y} and unbiased for θ , but somehow "better" than $\hat{\theta}$? First, we need to define what "better" means. A standard criterion is the Mean Squared Error (MSE).

Definition 1.4 (Mean Squared Error (MSE)). The MSE of an estimator $\hat{\phi}$ for a parameter ϕ is

$$\text{MSE}(\hat{\phi}) := \mathbb{E}_\phi \left[(\hat{\phi} - \phi)^2 \right],$$

where the expectation is taken with respect to the distribution governed by the true parameter ϕ . Note that the MSE generally depends on the true value of ϕ .

We say an estimator $\hat{\phi}$ is **better** than another estimator $\tilde{\phi}$ if $\text{MSE}(\hat{\phi}) \leq \text{MSE}(\tilde{\phi})$ for all possible values of the parameter ϕ , with strict inequality for at least one value.

Bias-Variance Decomposition of MSE: A fundamental property of MSE is its decomposition into variance and squared bias. For any estimator $\hat{\phi}$ of ϕ :

$$\begin{aligned}\text{MSE}(\hat{\phi}) &= \mathbb{E} \left[(\hat{\phi} - \phi)^2 \right] \\ &= \mathbb{E} \left[((\hat{\phi} - \mathbb{E}[\hat{\phi}]) + (\mathbb{E}[\hat{\phi}] - \phi))^2 \right] \\ &= \mathbb{E} \left[(\hat{\phi} - \mathbb{E}[\hat{\phi}])^2 \right] + \mathbb{E} \left[(\mathbb{E}[\hat{\phi}] - \phi)^2 \right] + 2\mathbb{E} \left[(\hat{\phi} - \mathbb{E}[\hat{\phi}])(\mathbb{E}[\hat{\phi}] - \phi) \right] \\ &= \mathbb{E} \left[(\hat{\phi} - \mathbb{E}[\hat{\phi}])^2 \right] + (\mathbb{E}[\hat{\phi}] - \phi)^2 + 2(\mathbb{E}[\hat{\phi}] - \phi) \underbrace{\mathbb{E}[\hat{\phi} - \mathbb{E}[\hat{\phi}]]}_{=0} \\ &= \underbrace{\mathbb{E} \left[(\hat{\phi} - \mathbb{E}[\hat{\phi}])^2 \right]}_{\text{Var}(\hat{\phi})} + \underbrace{(\mathbb{E}[\hat{\phi}] - \phi)^2}_{(\text{bias}(\hat{\phi}))^2}\end{aligned}\tag{5}$$

where $\text{bias}(\hat{\phi}) = \mathbb{E}[\hat{\phi}] - \phi$.

Crucially, for an **unbiased** estimator, the bias term is zero, so $\text{bias}(\hat{\phi}) = 0$. In this case, the decomposition simplifies beautifully:

$$\text{MSE}(\hat{\phi}) = \text{Var}(\hat{\phi}) \quad (\text{if } \hat{\phi} \text{ is unbiased}).$$

Therefore, if we restrict our search to *unbiased* estimators, finding the "best" estimator in terms of MSE is equivalent to finding the estimator with the *minimum variance*.

This leads us to a central question: Among all linear unbiased estimators for $\theta = \mathbf{a}^\top \boldsymbol{\beta}$, does the LS estimator $\hat{\theta} = \mathbf{a}^\top \hat{\boldsymbol{\beta}}$ have the smallest variance? The answer is yes, as formalized by the famous Gauss-Markov Theorem.

2 The Gauss-Markov Theorem

This theorem is a cornerstone of linear regression theory. It provides a powerful justification for using the least squares estimator, relying only on the assumptions about the first and second moments of the errors, not on any specific distribution like the normal distribution.

Theorem 2.1 (Gauss-Markov). *Assume the linear model holds: $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ with $\mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}$ and $\text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n$. Let $\theta = \mathbf{a}^\top \boldsymbol{\beta}$ be any linear combination of the coefficients. Let $\hat{\theta} = \mathbf{a}^\top \hat{\boldsymbol{\beta}}$ be the least squares estimator of θ , where $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$.*

*Consider any other estimator $\tilde{\theta}$ that is also **linear** in \mathbf{Y} (i.e., $\tilde{\theta} = \mathbf{d}^\top \mathbf{Y}$ for some vector $\mathbf{d} \in \mathbb{R}^n$) and **unbiased** for θ (i.e., $\mathbb{E}[\tilde{\theta}] = \theta$ for all possible $\boldsymbol{\beta}$).*

Then, the variance of the LS estimator is less than or equal to the variance of any other such linear unbiased estimator:

$$\text{Var}(\hat{\theta}) \leq \text{Var}(\tilde{\theta}) \quad \text{for all } \boldsymbol{\beta}.$$

For this reason, $\hat{\theta}$ is called the **Best Linear Unbiased Estimator (BLUE)** of θ .

Proof. Let $\hat{\theta} = \mathbf{c}^\top \mathbf{Y}$ where $\mathbf{c} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}$, as defined earlier. Let $\tilde{\theta} = \mathbf{d}^\top \mathbf{Y}$ be any other linear unbiased estimator for θ .

We can write the vector \mathbf{d} in terms of \mathbf{c} and a difference vector $\boldsymbol{\Delta}$:

$$\mathbf{d} = \mathbf{c} + \boldsymbol{\Delta}, \quad \text{where } \boldsymbol{\Delta} = \mathbf{d} - \mathbf{c} \in \mathbb{R}^n.$$

Now, let's use the unbiasedness condition for $\tilde{\theta}$: $\mathbb{E}[\tilde{\theta}] = \theta$ must hold for all $\boldsymbol{\beta}$.

$$\begin{aligned} \theta &= \mathbb{E}[\tilde{\theta}] = \mathbb{E}[\mathbf{d}^\top \mathbf{Y}] \\ &= \mathbb{E}[(\mathbf{c} + \boldsymbol{\Delta})^\top \mathbf{Y}] \\ &= \mathbb{E}[\mathbf{c}^\top \mathbf{Y}] + \mathbb{E}[\boldsymbol{\Delta}^\top \mathbf{Y}] \\ &= \mathbb{E}[\hat{\theta}] + \boldsymbol{\Delta}^\top \mathbb{E}[\mathbf{Y}] \quad (\text{since } \boldsymbol{\Delta} \text{ is constant}) \\ &= \theta + \boldsymbol{\Delta}^\top (\mathbf{X}\boldsymbol{\beta}) \quad (\text{since } \hat{\theta} \text{ is unbiased and } \mathbb{E}[\mathbf{Y}] = \mathbf{X}\boldsymbol{\beta}) \end{aligned}$$

Comparing the start and end of this chain of equalities, we must have:

$$\boldsymbol{\Delta}^\top \mathbf{X}\boldsymbol{\beta} = 0 \quad \text{for all } \boldsymbol{\beta} \in \mathbb{R}^{p+1}.$$

The only way this can hold for all possible vectors $\boldsymbol{\beta}$ is if the vector multiplying $\boldsymbol{\beta}$ is the zero vector. That is,

$$\boldsymbol{\Delta}^\top \mathbf{X} = \mathbf{0}^\top. \tag{6}$$

This condition captures the implication of unbiasedness for the difference vector Δ .

Now, let's use this condition to examine the relationship between Δ and \mathbf{c} . Recall $\mathbf{c} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}$. Consider the inner product $\Delta^\top \mathbf{c}$:

$$\Delta^\top \mathbf{c} = \Delta^\top \left(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} \right) = \underbrace{(\Delta^\top \mathbf{X})}_{=\mathbf{0}^\top} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} = \mathbf{0}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a} = 0.$$

So, the vectors Δ and \mathbf{c} are orthogonal: $\Delta^\top \mathbf{c} = 0$.

Finally, let's calculate the variance of $\tilde{\theta}$:

$$\begin{aligned} \text{Var}(\tilde{\theta}) &= \text{Var}(\mathbf{d}^\top \mathbf{Y}) \\ &= \mathbf{d}^\top \text{Cov}(\mathbf{Y}) \mathbf{d} \\ &= \mathbf{d}^\top (\sigma^2 \mathbf{I}_n) \mathbf{d} \\ &= \sigma^2 \mathbf{d}^\top \mathbf{d} \\ &= \sigma^2 (\mathbf{c} + \Delta)^\top (\mathbf{c} + \Delta) \\ &= \sigma^2 (\mathbf{c}^\top \mathbf{c} + \mathbf{c}^\top \Delta + \Delta^\top \mathbf{c} + \Delta^\top \Delta) \\ &= \sigma^2 (\mathbf{c}^\top \mathbf{c} + 0 + 0 + \Delta^\top \Delta) \quad (\text{since } \Delta^\top \mathbf{c} = \mathbf{c}^\top \Delta = 0) \\ &= \sigma^2 (\mathbf{c}^\top \mathbf{c} + \Delta^\top \Delta) \end{aligned}$$

We know that $\text{Var}(\hat{\theta}) = \sigma^2 \mathbf{c}^\top \mathbf{c}$. Therefore,

$$\text{Var}(\tilde{\theta}) = \text{Var}(\hat{\theta}) + \sigma^2 \Delta^\top \Delta.$$

Since $\Delta^\top \Delta = \|\Delta\|^2$ is the squared Euclidean norm of Δ , it must be non-negative ($\Delta^\top \Delta \geq 0$). Thus,

$$\text{Var}(\tilde{\theta}) \geq \text{Var}(\hat{\theta}).$$

Equality holds if and only if $\Delta^\top \Delta = 0$, which implies $\Delta = \mathbf{0}$, meaning $\mathbf{d} = \mathbf{c}$. This confirms that the LS estimator $\hat{\theta}$ has the minimum variance among all linear unbiased estimators. \square

Remark 2.2. The Gauss-Markov theorem is truly remarkable. It tells us that without making any assumptions about the *shape* of the error distribution (like normality), as long as the errors are uncorrelated and have constant variance, the least squares procedure naturally yields estimators that are optimal in the class of linear unbiased estimators, judged by the criterion of minimum variance (or equivalently, minimum MSE within this class).

3 Toward Inference: The Need for Distributional Assumptions

We have established the optimality of LS estimators in terms of minimum variance among linear unbiased estimators (point estimation). However, many statistical tasks go beyond point estimation. For example, we often want to:

- Construct **confidence intervals** for coefficients β_j or linear combinations $\theta = \mathbf{a}^\top \boldsymbol{\beta}$.
- Perform **hypothesis tests**, such as testing if a particular coefficient is zero ($H_0 : \beta_j = 0$).

To perform these tasks, we need more than just the mean and variance of our estimators; we need to know their *sampling distribution*. This typically requires making stronger assumptions about the distribution of the error term ϵ . The most common and mathematically tractable assumption is that the errors follow a multivariate normal distribution.

Before diving into the normal linear model, let's briefly review some key concepts about multivariate distributions, focusing on the multivariate normal.

3.1 Review of Multivariate Distributions

Consider a random vector $\mathbf{Z} = (Z_1, \dots, Z_k)^\top$.

Joint CDF: The joint cumulative distribution function (CDF) is defined as

$$F_{\mathbf{Z}}(z_1, \dots, z_k) := P(Z_1 \leq z_1, \dots, Z_k \leq z_k).$$

The CDF always exists and uniquely determines the distribution of \mathbf{Z} .

Independence: The components Z_1, \dots, Z_k are mutually (statistically) independent if and only if the joint CDF factors into the product of the marginal CDFs:

$$F_{\mathbf{Z}}(z_1, \dots, z_k) = P(Z_1 \leq z_1) \cdots P(Z_k \leq z_k) = F_{Z_1}(z_1) \cdots F_{Z_k}(z_k)$$

for all $z_1, \dots, z_k \in \mathbb{R}$.

Joint PDF: If the joint CDF is sufficiently smooth, we can define the joint probability density function (PDF) via differentiation:

$$f_{\mathbf{Z}}(z_1, \dots, z_k) = \frac{\partial^k}{\partial z_1 \cdots \partial z_k} F_{\mathbf{Z}}(z_1, \dots, z_k).$$

When the PDF exists, it provides an equivalent characterization of the distribution, satisfying

$$F_{\mathbf{Z}}(z_1, \dots, z_k) = \int_{-\infty}^{z_1} \cdots \int_{-\infty}^{z_k} f_{\mathbf{Z}}(u_1, \dots, u_k) du_k \cdots du_1,$$

and $f_{\mathbf{Z}}(\mathbf{z}) \geq 0$ with $\int_{\mathbb{R}^k} f_{\mathbf{Z}}(\mathbf{z}) d\mathbf{z} = 1$. For independent continuous random variables, the joint PDF is the product of the marginal PDFs: $f_{\mathbf{Z}}(\mathbf{z}) = f_{Z_1}(z_1) \cdots f_{Z_k}(z_k)$.

3.2 The Multivariate Normal Distribution

The multivariate normal (MVN) distribution is arguably the most important multivariate distribution in statistics, partly due to the Central Limit Theorem and its mathematical tractability.

Definition 3.1 (Multivariate Normal Distribution). A random vector $\mathbf{W} = (W_1, \dots, W_k)^\top$ is said to have a **multivariate normal distribution** if it can be represented as

$$\mathbf{W} \stackrel{d}{=} \boldsymbol{\mu} + \mathbf{A}\mathbf{Z} \tag{7}$$

where:

- $\boldsymbol{\mu} \in \mathbb{R}^k$ is a constant vector (the mean).

- $\mathbf{A} \in \mathbb{R}^{k \times l}$ is a constant matrix.
- $\mathbf{Z} = (Z_1, \dots, Z_l)^\top$ is a random vector whose components Z_i are independent and identically distributed (i.i.d.) standard normal random variables, $Z_i \sim \mathcal{N}(0, 1)$.

The symbol $\stackrel{d}{=}$ means "equal in distribution".

Properties of the Multivariate Normal Distribution:

1. **Mean and Covariance:** If $\mathbf{W} \stackrel{d}{=} \boldsymbol{\mu} + \mathbf{A}\mathbf{Z}$ as in Definition 3.1, we can calculate its mean vector and covariance matrix:

$$\begin{aligned}\mathbb{E}[\mathbf{W}] &= \mathbb{E}[\boldsymbol{\mu} + \mathbf{A}\mathbf{Z}] = \boldsymbol{\mu} + \mathbf{A}\mathbb{E}[\mathbf{Z}] = \boldsymbol{\mu} + \mathbf{A}\mathbf{0} = \boldsymbol{\mu}. \\ \text{Cov}(\mathbf{W}) &= \text{Cov}(\boldsymbol{\mu} + \mathbf{A}\mathbf{Z}) = \text{Cov}(\mathbf{A}\mathbf{Z}) \\ &= \mathbf{A}\text{Cov}(\mathbf{Z})\mathbf{A}^\top \\ &= \mathbf{A}(\mathbf{I}_l)\mathbf{A}^\top \quad (\text{since } Z_i \text{ are i.i.d. } \mathcal{N}(0, 1)) \\ &= \mathbf{A}\mathbf{A}^\top.\end{aligned}$$

Let $\mathbf{V} = \mathbf{A}\mathbf{A}^\top$. Then \mathbf{V} is the covariance matrix of \mathbf{W} . Note that \mathbf{V} is always symmetric and positive semi-definite. The representation $\mathbf{W} \stackrel{d}{=} \boldsymbol{\mu} + \mathbf{A}\mathbf{Z}$ implies $\mathbb{E}[\mathbf{W}] = \boldsymbol{\mu}$ and $\text{Cov}(\mathbf{W}) = \mathbf{V} = \mathbf{A}\mathbf{A}^\top$. The mean vector $\boldsymbol{\mu}$ is unique, but the matrix \mathbf{A} is not unique for a given \mathbf{V} (e.g., if \mathbf{U} is an $l \times l$ orthogonal matrix, $\mathbf{A}' = \mathbf{A}\mathbf{U}$ gives $\mathbf{A}'(\mathbf{A}')^\top = \mathbf{A}\mathbf{U}\mathbf{U}^\top\mathbf{A}^\top = \mathbf{A}\mathbf{A}^\top = \mathbf{V}$).

2. **Density Function (when \mathbf{V} is invertible):** If $l = k$ and the matrix $\mathbf{A}_{k \times k}$ is invertible (equivalently, has linearly independent columns), then the covariance matrix $\mathbf{V} = \mathbf{A}\mathbf{A}^\top$ is positive definite (and thus invertible). In this case, the random vector \mathbf{W} has a probability density function (PDF) given by:

$$f_{\mathbf{W}}(\mathbf{w}) = (2\pi)^{-k/2} |\det(\mathbf{V})|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu})^\top \mathbf{V}^{-1}(\mathbf{w} - \boldsymbol{\mu})\right), \quad \mathbf{w} \in \mathbb{R}^k.$$

If \mathbf{V} is singular (not invertible), the distribution does not have a density with respect to the Lebesgue measure on \mathbb{R}^k ; the distribution is concentrated on a lower-dimensional affine subspace.

3. **Notation:** Regardless of whether \mathbf{V} is invertible, the distribution of \mathbf{W} is completely determined by its mean vector $\boldsymbol{\mu}$ and covariance matrix \mathbf{V} . We use the notation

$$\mathbf{W} \sim \mathcal{N}_k(\boldsymbol{\mu}, \mathbf{V})$$

to denote that \mathbf{W} follows a k -dimensional multivariate normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix \mathbf{V} .

4. **Linear Combinations are Normal:** If $\mathbf{W} \sim \mathcal{N}_k(\boldsymbol{\mu}, \mathbf{V})$ and $\mathbf{c} \in \mathbb{R}^k$ is a constant vector, then the scalar random variable $Y = \mathbf{c}^\top \mathbf{W}$ has a univariate normal distribution:

$$\mathbf{c}^\top \mathbf{W} \sim \mathcal{N}(\mathbf{c}^\top \boldsymbol{\mu}, \mathbf{c}^\top \mathbf{V} \mathbf{c}).$$

This is a crucial property. It follows directly from the definition: if $\mathbf{W} = \boldsymbol{\mu} + \mathbf{A}\mathbf{Z}$, then $\mathbf{c}^\top \mathbf{W} = \mathbf{c}^\top \boldsymbol{\mu} + (\mathbf{c}^\top \mathbf{A})\mathbf{Z}$. Since $(\mathbf{c}^\top \mathbf{A})\mathbf{Z}$ is a linear combination of independent standard

normals, it is itself normally distributed (possibly with variance 0 if $\mathbf{c}^\top \mathbf{A} = \mathbf{0}^\top$). The mean is $\mathbf{c}^\top \boldsymbol{\mu}$ and the variance is $(\mathbf{c}^\top \mathbf{A})(\mathbf{c}^\top \mathbf{A})^\top = \mathbf{c}^\top (\mathbf{A}\mathbf{A}^\top) \mathbf{c} = \mathbf{c}^\top \mathbf{V} \mathbf{c}$. *Example:* Taking $\mathbf{c} = (0, \dots, 0, 1, 0, \dots, 0)^\top$ (with 1 in the j -th position) shows that each component W_j is univariate normal:

$$W_j \sim \mathcal{N}(\mu_j, V_{jj}),$$

where V_{jj} is the j -th diagonal element of \mathbf{V} .

5. **Characterization via Linear Combinations:** The previous property has a converse. A random vector \mathbf{W} is multivariate normal if and only if every linear combination $\mathbf{c}^\top \mathbf{W}$ is univariate normal. Formally:

$$\mathbf{W} \sim \mathcal{N}_k(\boldsymbol{\mu}, \mathbf{V}) \iff \mathbf{c}^\top \mathbf{W} \sim \mathcal{N}(\mathbf{c}^\top \boldsymbol{\mu}, \mathbf{c}^\top \mathbf{V} \mathbf{c}) \quad \forall \mathbf{c} \in \mathbb{R}^k.$$

(Here $\boldsymbol{\mu} = \mathbb{E}[\mathbf{W}]$ and $\mathbf{V} = \text{Cov}(\mathbf{W})$). This provides an alternative way to define or check for multivariate normality.

6. **Linear Transformations Preserve Normality:** If $\mathbf{W} \sim \mathcal{N}_k(\boldsymbol{\mu}, \mathbf{V})$ and $\mathbf{C} \in \mathbb{R}^{m \times k}$ is a constant matrix, then the transformed vector $\mathbf{Y} = \mathbf{C}\mathbf{W}$ is also multivariate normal:

$$\mathbf{C}\mathbf{W} \sim \mathcal{N}_m(\mathbf{C}\boldsymbol{\mu}, \mathbf{C}\mathbf{V}\mathbf{C}^\top).$$

This follows because if $\mathbf{W} = \boldsymbol{\mu} + \mathbf{A}\mathbf{Z}$, then $\mathbf{C}\mathbf{W} = \mathbf{C}\boldsymbol{\mu} + (\mathbf{C}\mathbf{A})\mathbf{Z}$. This fits the definition of MVN with mean $\mathbf{C}\boldsymbol{\mu}$ and matrix $\mathbf{A}' = \mathbf{C}\mathbf{A}$. The covariance is $\mathbf{A}'(\mathbf{A}')^\top = (\mathbf{C}\mathbf{A})(\mathbf{C}\mathbf{A})^\top = \mathbf{C}(\mathbf{A}\mathbf{A}^\top)\mathbf{C}^\top = \mathbf{C}\mathbf{V}\mathbf{C}^\top$.

7. **Sums of Independent Normals are Normal:** If $\mathbf{W}^{(j)} \sim \mathcal{N}_k(\boldsymbol{\mu}^{(j)}, \mathbf{V}^{(j)})$ for $j = 1, \dots, p$, are independent random vectors, and d_1, \dots, d_p are scalar constants, then their linear combination is also multivariate normal:

$$\sum_{j=1}^p d_j \mathbf{W}^{(j)} \sim \mathcal{N}_k \left(\sum_{j=1}^p d_j \boldsymbol{\mu}^{(j)}, \sum_{j=1}^p d_j^2 \mathbf{V}^{(j)} \right).$$

8. **Independence and Zero Covariance:** Let $\mathbf{W} \sim \mathcal{N}_k(\boldsymbol{\mu}, \mathbf{V})$. Partition \mathbf{W} into two sub-vectors, $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$, corresponding to disjoint subsets of indices $\mathcal{I}_1, \mathcal{I}_2 \subseteq \{1, \dots, k\}$. That is, $\mathbf{W}^{(1)} = (W_i : i \in \mathcal{I}_1)$ and $\mathbf{W}^{(2)} = (W_j : j \in \mathcal{I}_2)$. A remarkable property of the MVN distribution is that $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$ are **statistically independent** if and only if their cross-covariance is zero:

$$\text{Cov}(W_i, W_j) = 0 \quad \forall i \in \mathcal{I}_1, j \in \mathcal{I}_2 \iff \mathbf{W}^{(1)} \text{ and } \mathbf{W}^{(2)} \text{ are independent.}$$

In general, zero covariance does not imply independence, but it does for jointly normal random variables. This property is extremely useful in linear models.

4 Distributions Related to the Normal

Several important probability distributions arise from transformations of normal random variables. These are fundamental for constructing confidence intervals and test statistics in the context of the normal linear model.

Definition 4.1 (Chi-square (χ^2) Distribution). If Z_1, Z_2, \dots, Z_k are independent and identically distributed (i.i.d.) standard normal random variables, $Z_i \sim \mathcal{N}(0, 1)$, then the distribution of the sum of their squares,

$$Q = \sum_{j=1}^k Z_j^2,$$

is called the **Chi-square distribution** with k degrees of freedom. We denote this by $Q \sim \chi_k^2$. (In R: ‘pchisq()’, ‘qchisq()’, ‘rchisq()’, ‘dchisq()’).

Proposition 4.2 (Expected Value of χ_k^2). If $Q \sim \chi_k^2$, then $\mathbb{E}[Q] = k$.

Proof. Let $Q = \sum_{i=1}^k Z_i^2$ where $Z_i \stackrel{iid}{\sim} \mathcal{N}(0, 1)$. By linearity of expectation:

$$\mathbb{E}[Q] = \mathbb{E}\left[\sum_{i=1}^k Z_i^2\right] = \sum_{i=1}^k \mathbb{E}[Z_i^2].$$

For a standard normal random variable Z_i , we know $\mathbb{E}[Z_i] = 0$ and $\text{Var}(Z_i) = 1$. Since $\text{Var}(Z_i) = \mathbb{E}[Z_i^2] - (\mathbb{E}[Z_i])^2$, we have

$$\mathbb{E}[Z_i^2] = \text{Var}(Z_i) + (\mathbb{E}[Z_i])^2 = 1 + 0^2 = 1.$$

Therefore, $\mathbb{E}[Q] = \sum_{i=1}^k 1 = k$. □

Proposition 4.3 (Quadratic Forms in Normal Variables). Let $\mathbf{Z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n)$ (a vector of n i.i.d. standard normal variables). Let \mathbf{P} be an $n \times n$ symmetric ($\mathbf{P}^\top = \mathbf{P}$) and idempotent ($\mathbf{P}^2 = \mathbf{P}$) matrix with $\text{rank}(\mathbf{P}) = r$. Then the quadratic form $\mathbf{Z}^\top \mathbf{P} \mathbf{Z}$ follows a Chi-square distribution with r degrees of freedom:

$$\mathbf{Z}^\top \mathbf{P} \mathbf{Z} = \|\mathbf{P} \mathbf{Z}\|^2 \sim \chi_r^2.$$

Proof Sketch. Since \mathbf{P} is symmetric and idempotent, it represents an orthogonal projection onto a subspace of dimension $r = \text{rank}(\mathbf{P}) = \text{tr}(\mathbf{P})$. Such a matrix has eigenvalues that are either 0 or 1, with exactly r eigenvalues equal to 1. By the spectral theorem, $\mathbf{P} = \mathbf{U} \mathbf{D} \mathbf{U}^\top$ where \mathbf{U} is orthogonal and \mathbf{D} is diagonal with r ones and $n - r$ zeros. Let $\mathbf{W} = \mathbf{U}^\top \mathbf{Z}$. Since \mathbf{U}^\top is orthogonal, $\mathbf{W} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{U}^\top \mathbf{I}_n \mathbf{U}) = \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n)$. Then

$$\mathbf{Z}^\top \mathbf{P} \mathbf{Z} = \mathbf{Z}^\top \mathbf{U} \mathbf{D} \mathbf{U}^\top \mathbf{Z} = (\mathbf{U}^\top \mathbf{Z})^\top \mathbf{D} (\mathbf{U}^\top \mathbf{Z}) = \mathbf{W}^\top \mathbf{D} \mathbf{W} = \sum_{i=1}^n D_{ii} W_i^2 = \sum_{i=1}^r W_i^2,$$

where the sum is taken over the r indices i for which $D_{ii} = 1$. Since W_i are i.i.d. $\mathcal{N}(0, 1)$, this sum is, by definition, a χ_r^2 random variable. (The original note uses $\mathbb{E}\|\mathbf{P} \mathbf{Z}\|^2 = \text{tr}(\text{Cov}[\mathbf{P} \mathbf{Z}]) = \text{tr}(\mathbf{P} \mathbf{I} \mathbf{P}^\top) = \text{tr}(\mathbf{P}) = r$. While this correctly calculates the expected value, it doesn't fully prove the distribution is χ_r^2 . The spectral decomposition argument is more complete). □

Definition 4.4 (Student's t -distribution). If $Z \sim \mathcal{N}(0, 1)$ and $V \sim \chi_k^2$ are independent random variables, then the distribution of the ratio

$$T = \frac{Z}{\sqrt{V/k}}$$

is called the **t -distribution** with k degrees of freedom. We denote this by $T \sim t_k$. (In R: ‘pt()’, ‘qt()’, ‘rt()’, ‘dt()’).

Definition 4.5 (*F-distribution*). If $V_1 \sim \chi_{k_1}^2$ and $V_2 \sim \chi_{k_2}^2$ are independent random variables, then the distribution of the ratio of scaled Chi-square variables

$$F = \frac{V_1/k_1}{V_2/k_2}$$

is called the **F-distribution** with k_1 (numerator) and k_2 (denominator) degrees of freedom. We denote this by $F \sim F_{k_1, k_2}$. (In R: ‘pf()’, ‘qf()’, ‘rf()’, ‘df()’).

5 Inference Under the Normal Linear Model

We now combine the linear model structure with the assumption of normally distributed errors. This allows us to derive exact sampling distributions for our estimators, paving the way for confidence intervals and hypothesis tests.

Recall the standard linear model assumptions (4):

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}, \quad \text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n.$$

We now add the assumption that the errors are multivariate normal.

Definition 5.1 (Normal Linear Model). The **normal linear model** assumes:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n). \quad (8)$$

This implies that the observations \mathbf{Y} are also multivariate normal:

$$\mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n).$$

Under this model, we can determine the exact distributions of the LS estimator $\hat{\boldsymbol{\beta}}$ and the variance estimator $\hat{\sigma}^2$.

5.1 Distribution of the LS Estimator $\hat{\boldsymbol{\beta}}$

Recall the formula for the LS estimator:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}.$$

Let $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$, which is a $(p+1) \times n$ matrix. We can write $\hat{\boldsymbol{\beta}}$ in terms of the error vector $\boldsymbol{\epsilon}$:

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \mathbf{A}\mathbf{Y} \\ &= \mathbf{A}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) \\ &= \mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbf{A}\boldsymbol{\epsilon} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} + \mathbf{A}\boldsymbol{\epsilon} \\ &= \mathbf{I}_{p+1}\boldsymbol{\beta} + \mathbf{A}\boldsymbol{\epsilon} \\ &= \boldsymbol{\beta} + \mathbf{A}\boldsymbol{\epsilon}. \end{aligned}$$

Since $\boldsymbol{\epsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, and $\hat{\boldsymbol{\beta}}$ is a linear transformation of $\boldsymbol{\epsilon}$ shifted by a constant vector $\boldsymbol{\beta}$, it follows from Property 6 of MVN distributions that $\hat{\boldsymbol{\beta}}$ is also multivariate normal.

We already know its mean and covariance matrix from our earlier work (which did not require normality):

$$\mathbb{E}[\hat{\beta}] = \beta$$

$$\text{Cov}(\hat{\beta}) = \text{Cov}(\beta + \mathbf{A}\epsilon) = \text{Cov}(\mathbf{A}\epsilon) = \mathbf{A}\text{Cov}(\epsilon)\mathbf{A}^\top = \mathbf{A}(\sigma^2\mathbf{I}_n)\mathbf{A}^\top = \sigma^2\mathbf{A}\mathbf{A}^\top.$$

Calculating $\mathbf{A}\mathbf{A}^\top$:

$$\mathbf{A}\mathbf{A}^\top = \left((\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\right)\left((\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\right)^\top = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1} = (\mathbf{X}^\top\mathbf{X})^{-1}.$$

So, $\text{Cov}(\hat{\beta}) = \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}$.

Combining these results, we have the sampling distribution of the LS estimator under the normal linear model:

$$\hat{\beta} \sim \mathcal{N}_{p+1}\left(\beta, \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}\right). \quad (9)$$

5.2 Distribution of the Variance Estimator $\hat{\sigma}^2$

Recall the definition of the unbiased variance estimator:

$$\hat{\sigma}^2 = \frac{\|\mathbf{e}\|^2}{n-p-1} = \frac{SSE}{n-p-1},$$

where $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{P}_\mathbf{X}\mathbf{Y} = (\mathbf{I}_n - \mathbf{P}_\mathbf{X})\mathbf{Y}$ are the residuals, and $\mathbf{P}_\mathbf{X} = \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top$ is the projection matrix onto the column space of \mathbf{X} . Let $\mathbf{Q} = \mathbf{I}_n - \mathbf{P}_\mathbf{X}$. Then $\mathbf{e} = \mathbf{Q}\mathbf{Y}$.

We can express the residuals in terms of the errors ϵ :

$$\mathbf{e} = \mathbf{Q}\mathbf{Y} = \mathbf{Q}(\mathbf{X}\beta + \epsilon) = \mathbf{Q}\mathbf{X}\beta + \mathbf{Q}\epsilon.$$

Since \mathbf{Q} projects onto the space orthogonal to the column space of \mathbf{X} , we have $\mathbf{Q}\mathbf{X} = \mathbf{0}$. Thus,

$$\mathbf{e} = \mathbf{Q}\epsilon.$$

The sum of squared errors (SSE) is $\|\mathbf{e}\|^2 = \|\mathbf{Q}\epsilon\|^2 = \epsilon^\top \mathbf{Q}^\top \mathbf{Q} \epsilon$. The matrix \mathbf{Q} is symmetric ($\mathbf{Q}^\top = \mathbf{Q}$) and idempotent ($\mathbf{Q}^2 = \mathbf{Q}$) because it's an orthogonal projection matrix. Its rank is $\text{rank}(\mathbf{Q}) = \text{tr}(\mathbf{Q}) = \text{tr}(\mathbf{I}_n - \mathbf{P}_\mathbf{X}) = \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{P}_\mathbf{X}) = n - \text{tr}(\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top) = n - \text{tr}((\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X}) = n - \text{tr}(\mathbf{I}_{p+1}) = n - (p+1)$.

Now, consider the random vector $\mathbf{Z} = \epsilon/\sigma$. Since $\epsilon \sim \mathcal{N}_n(\mathbf{0}, \sigma^2\mathbf{I}_n)$, we have $\mathbf{Z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n)$. The SSE can be written as:

$$SSE = \|\mathbf{e}\|^2 = \|\mathbf{Q}\epsilon\|^2 = \|\sigma\mathbf{Q}(\epsilon/\sigma)\|^2 = \sigma^2\|\mathbf{Q}\mathbf{Z}\|^2.$$

We apply Proposition 4.3 with the matrix $\mathbf{P} = \mathbf{Q}$ and vector $\mathbf{Z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n)$. Since \mathbf{Q} is symmetric, idempotent, and has rank $r = n - p - 1$, we conclude that

$$\frac{SSE}{\sigma^2} = \frac{\|\mathbf{e}\|^2}{\sigma^2} = \|\mathbf{Q}\mathbf{Z}\|^2 \sim \chi_{n-p-1}^2.$$

Relating this back to $\hat{\sigma}^2$:

$$\frac{(n-p-1)\hat{\sigma}^2}{\sigma^2} = \frac{SSE}{\sigma^2} \sim \chi_{n-p-1}^2.$$

This gives us the sampling distribution of the variance estimator $\hat{\sigma}^2$. It follows a scaled Chi-square distribution.

5.3 Joint Distribution of $\hat{\beta}$ and $\hat{\sigma}^2$

A crucial result for constructing t -statistics and F -statistics is that, under the normal linear model, the LS estimator $\hat{\beta}$ is statistically independent of the variance estimator $\hat{\sigma}^2$.

To show this, we will show that $\hat{\beta}$ is independent of the residual vector \mathbf{e} , since $\hat{\sigma}^2$ is solely a function of \mathbf{e} (specifically, $\|\mathbf{e}\|^2$).

We know $\hat{\beta} = \mathbf{A}\mathbf{Y}$ and $\mathbf{e} = \mathbf{Q}\mathbf{Y}$, where $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ and $\mathbf{Q} = \mathbf{I}_n - \mathbf{P}_\mathbf{X}$. Consider the joint vector:

$$\begin{pmatrix} \hat{\beta} \\ \mathbf{e} \end{pmatrix} = \begin{pmatrix} \mathbf{A} \\ \mathbf{Q} \end{pmatrix} \mathbf{Y}.$$

Since $\mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$, this stacked vector is also multivariate normal because it's a linear transformation of \mathbf{Y} .

Now, we examine the covariance between the blocks $\hat{\beta}$ and \mathbf{e} :

$$\begin{aligned} \text{Cov}(\hat{\beta}, \mathbf{e}) &= \text{Cov}(\mathbf{A}\mathbf{Y}, \mathbf{Q}\mathbf{Y}) \\ &= \mathbf{A}\text{Cov}(\mathbf{Y})\mathbf{Q}^\top \\ &= \mathbf{A}(\sigma^2 \mathbf{I}_n)\mathbf{Q} \quad (\text{since } \mathbf{Q} \text{ is symmetric}) \\ &= \sigma^2 \mathbf{A}\mathbf{Q} \end{aligned}$$

Let's compute $\mathbf{A}\mathbf{Q}$:

$$\begin{aligned} \mathbf{A}\mathbf{Q} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{I}_n - \mathbf{P}_\mathbf{X}) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{I}_{p+1} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \\ &= \mathbf{0}_{(p+1) \times n} \quad (\text{The zero matrix}) \end{aligned}$$

Therefore, $\text{Cov}(\hat{\beta}, \mathbf{e}) = \mathbf{0}$.

Since $(\hat{\beta}, \mathbf{e})$ are jointly multivariate normal and their covariance is zero, we can invoke Property 8 of MVN distributions to conclude that $\hat{\beta}$ and \mathbf{e} are **statistically independent**.

Because $\hat{\sigma}^2 = \|\mathbf{e}\|^2/(n - p - 1)$ is a function of \mathbf{e} only, it follows that $\hat{\beta}$ and $\hat{\sigma}^2$ are also **statistically independent**.

Theorem 5.2 (Independence of $\hat{\beta}$ and $\hat{\sigma}^2$). *Under the normal linear model (8), the least squares estimator $\hat{\beta}$ is statistically independent of the residual variance estimator $\hat{\sigma}^2$.*

This independence is fundamental. It allows us, for example, to form t -statistics like

$$\frac{(\mathbf{a}^\top \hat{\beta} - \mathbf{a}^\top \beta)}{\sqrt{\hat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}}$$

where the numerator (related to $\hat{\beta}$) is independent of the denominator (related to $\hat{\sigma}^2$), leading to a t -distribution (after appropriate scaling), which forms the basis for confidence intervals and hypothesis tests for $\mathbf{a}^\top \beta$.