

# Theory of Linear Models

Longhai Li

2026-01-11

# Preface

## Key Features

This text adopts a geometric approach to the statistical theory of linear models, aiming to provide a deeper understanding than standard algebraic treatments. Key features include:

- **Projection Perspective:** We prioritize the geometric interpretation of least squares, viewing estimation as a projection of the response vector onto a model subspace. This visual framework unifies diverse topics—from simple regression to complex ANOVA designs—under a single theoretical umbrella.
- **Interactive Visualizations:** Abstract concepts are brought to life through interactive 3D plots. Readers can rotate and inspect vector spaces, residual planes, and projection geometries to build a tangible intuition for high-dimensional operations.
- **Computational Integration:** Theory is seamlessly integrated with practice. The text provides implementation examples using R (and Python), demonstrating how theoretical matrix equations translate directly into computational code.
- **Rigorous Foundations:** While visually driven, the text maintains mathematical rigor, covering essential topics such as spectral theory, the generalized inverse and the multivariate normal distribution to ensure a solid theoretical grounding.

## Overview

This course is a rigorous examination of the general linear models using vector space theory, in particular the approach of regarding least square as projection. The topics includes: vector space; projection; matrix algebra; generalized inverses; quadratic forms; theory for point estimation; theory for hypothesis test; theory for non-full-rank models.

## Audience

This book is designed for graduate students and advanced undergraduate students in statistics, data science, and related quantitative fields. It serves as a bridge between applied regression analysis and the theoretical foundations of linear models. Researchers and practitioners seeking a deeper geometric and algebraic understanding of the statistical methods they use daily will also find this text valuable.

## Prerequisites

To get the most out of this book, readers should have a comfortable grasp of the following topics:

**Linear Algebra:** An elementary understanding of matrix operations is essential. You should be familiar with matrix multiplication, determinants, inversion, and the basic concepts of vector spaces (such as linear independence, basis vectors, and subspaces). While we review key spectral theory concepts (like eigenvalues and the singular value decomposition) in the early chapters, prior exposure to these ideas is helpful.

**Probability and Statistics:** A standard introductory course in probability and mathematical statistics is required. Readers should be familiar with random variables, expectation, variance, covariance, common probability distributions (especially the Normal distribution), and fundamental concepts of hypothesis testing and estimation.

# 1 Introduction

## 1.1 Multiple Linear Regression

Suppose we have observations on  $Y$  and  $X_j$ . The data can be represented in matrix form.

$$\underset{n \times 1}{y} = \underset{n \times p}{X} \underset{n \times 1}{\beta} + \underset{n \times 1}{\epsilon} \quad (1.1)$$

where the error terms are distributed as:

$$\epsilon \sim N_n(0, \sigma^2 I_n), \quad (1.2)$$

in which  $I_n$  is the identity matrix:

$$I_n = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} \quad (1.3)$$

The scalar equation for a single observation is:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \epsilon_i \quad (1.4)$$

## 1.2 Examples

### 1.2.1 Polynomial Regression

Polynomial regression fits a curved line to the data points but remains linear in the parameters ( $\beta$ ).

The model equation is:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_{p-1} x_i^{p-1} \quad (1.5)$$

### 1.2.2 Design Matrix Construction

The design matrix  $X$  is constructed by taking powers of the input variable.

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{p-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{p-1} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix} \quad (1.6)$$

### 1.2.3 One-Way ANOVA

ANOVA can be expressed as a linear model using categorical predictors (dummy variables).

Suppose we have 3 groups ( $G_1, G_2, G_3$ ) with observations:

$$Y_{ij} = \mu_i + \epsilon_{ij}, \quad \epsilon_{ij} \sim N(0, \sigma^2) \quad (1.7)$$

$$\begin{array}{c|c|c} G_1 & G_2 & G_3 \\ \hline Y_{11} & Y_{21} & Y_{31} \\ Y_{12} & Y_{22} & Y_{32} \end{array} \quad (1.8)$$

We construct the matrix  $X$  to select the group mean ( $\mu$ ) corresponding to the observation:

$$\underset{6 \times 1}{y} = \underset{6 \times 3}{X} \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{pmatrix} + \epsilon \quad (1.9)$$

$$\begin{bmatrix} Y_{11} \\ Y_{12} \\ Y_{21} \\ Y_{22} \\ Y_{31} \\ Y_{32} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} + \epsilon \quad (1.10)$$

### 1.2.4 Analysis of Covariance (ANCOVA)

ANCOVA combines continuous variables and categorical (dummy) variables in the same design matrix.

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} X_{1,\text{cont}} & 1 & 0 \\ X_{2,\text{cont}} & 1 & 0 \\ \vdots & 0 & 1 \\ X_{n,\text{cont}} & 0 & 1 \end{bmatrix} \beta + \epsilon \quad (1.11)$$

## 1.3 Least Squares Estimation

For the general linear model  $y = X\beta + \epsilon$ , the Least Squares estimator is:

$$\hat{\beta} = (X'X)^{-1}X'y \quad (1.12)$$

The predicted values ( $\hat{y}$ ) are obtained via the Projection Matrix (Hat Matrix)  $P_X$ :

$$\hat{y} = X\hat{\beta} = X(X'X)^{-1}X'y = P_X y \quad (1.13)$$

The residuals and Sum of Squared Errors are:

$$\hat{e} = y - \hat{y} \quad (1.14)$$

$$\text{SSE} = \|\hat{e}\|^2 \quad (1.15)$$

The coefficient of determination is:

$$R^2 = \frac{\text{SST} - \text{SSE}}{\text{SST}} \quad (1.16)$$

where  $\text{SST} = \sum (y_i - \bar{y})^2$ .

## 1.4 Geometric Perspective of Least Square Estimation

We align the coordinate system to the models for clarity:

1. **Reduced Model** ( $M_0$ ): Represented by the **X-axis** (labeled  $j_3$ ).
  - $\hat{y}_0$  is the projection of  $y$  onto this axis.
2. **Full Model** ( $M_1$ ): Represented by the **XY-plane** (the floor).
  - $\hat{y}_1$  is the projection of  $y$  onto this plane ( $z = 0$ ).
3. **Observed Data** ( $y$ ): A point in 3D space.

The “improvement” due to adding predictors is the distance between  $\hat{y}_0$  and  $\hat{y}_1$ .

The geometric perspective is not merely for intuition, but as the most robust framework for mastering linear models. This approach offers three distinct advantages:

- **Statistical Clarity:** Geometry provides the most natural path to understanding the properties of estimators. By viewing least square estimation as an orthogonal projection, the decomposition of sums of squares into independent components becomes visually obvious, demystifying how degrees of freedom relate to subspace dimensions rather than abstract algebraic constants. The sampling distribution of the sum squares become straightforward.
- **Computational Stability:** A geometric understanding is essential for implementing efficient and numerically stable algorithms. While the algebraic “Normal Equations”  $((X'X)^{-1}X'y)$  are theoretically valid, they are often computationally hazardous. The geometric approach leads directly to superior methods—such as QR and Singular Value Decompositions—that are the backbone of modern statistical software.
- **Generalizability:** The principles of projection and orthogonality extend far beyond the Gaussian linear model. These geometric insights provide the foundational intuition needed for tackling non-Gaussian optimization problems, including Generalized Linear Models (GLMs) and convex optimization, where solutions can often be viewed as projections onto convex sets.

Geometric Interpretation: Aligned View

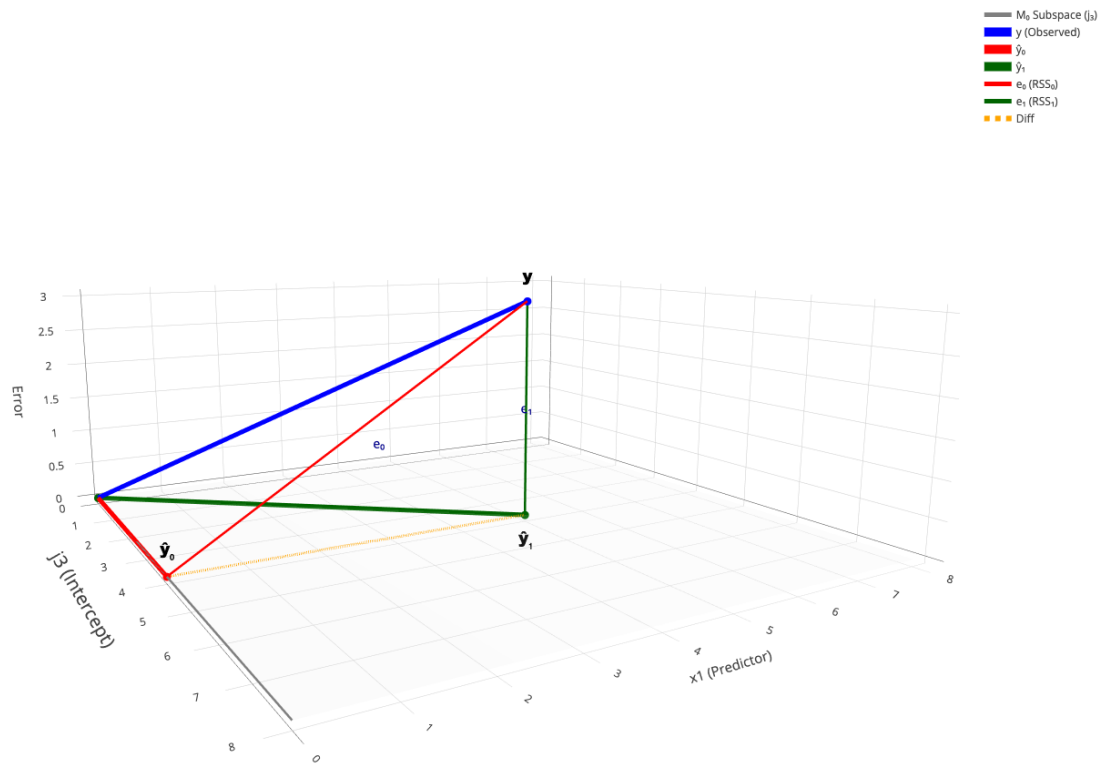


Figure 1.1: Geometric Interpretation: Projection onto Axis (M0) vs Plane (M1)

## 2 Projection in Vector Space

### 2.1 Vector and Projection onto a Line

#### 2.1.1 Vectors and Operations

The concept of a vector is fundamental to linear algebra and linear models. We begin by formally defining what a vector is in the context of Euclidean space.

**Definition 2.1** (Vector). A **vector**  $x$  is defined as a point in  $n$ -dimensional space ( $\mathbb{R}^n$ ). It is typically represented as a column vector containing  $n$  real-valued components:

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad (2.1)$$

Vectors are not just static points; they can be combined and manipulated. The two most basic geometric operations are addition and subtraction.

**Vector Arithmetic:** Vectors can be manipulated geometrically:

**Definition 2.2** (Vector Addition). The sum of two vectors  $x$  and  $y$  creates a new vector. The operation is performed component-wise, adding corresponding elements from each vector. Geometrically, this follows the “parallelogram rule” or the “head-to-tail” method, where you place the tail of  $y$  at the head of  $x$ .

$$x + y = \begin{pmatrix} x_1 + y_1 \\ \vdots \\ x_n + y_n \end{pmatrix} \quad (2.2)$$

**Definition 2.3** (Vector Subtraction). The difference  $d = y - x$  is the vector that “closes the triangle” formed by  $x$  and  $y$ . It represents the displacement vector that connects the tip of  $x$  to the tip of  $y$ , such that  $x + d = y$ .



### 2.1.2 Scalar Multiplication and Distance

In addition to combining vectors with each other, we can modify a single vector using a real number, known as a scalar.

**Definition 2.4** (Scalar Multiplication). Multiplying a vector by a scalar  $c$  scales its magnitude (length) without changing its line of direction. If  $c$  is positive, the direction remains the same; if  $c$  is negative, the direction is reversed.

$$cx = \begin{pmatrix} cx_1 \\ \vdots \\ cx_n \end{pmatrix} \quad (2.3)$$

We often need to quantify the “size” of a vector. This is done using the concept of length, or norm.

**Definition 2.5** (Euclidean Distance (Length)). The length (or norm) of a vector  $x = (x_1, \dots, x_n)^T$  corresponds to the straight-line distance from the origin to the point defined by  $x$ . It is defined as the square root of the sum of squared components:

$$\|x\|^2 = \sum_{i=1}^n x_i^2 \quad (2.4)$$

$$\|x\| = \sqrt{\sum_{i=1}^n x_i^2} \quad (2.5)$$

### 2.1.3 Angle and Inner Product

To understand the relationship between two vectors  $x$  and  $y$  beyond just their lengths, we must look at the angle between them. Consider the triangle formed by the vectors  $x$ ,  $y$ , and their difference  $y - x$ . By applying the classic **Law of Cosines** to this triangle, we can relate the geometric angle to the vector lengths.

**Theorem 2.1** (Law of Cosines). *For a triangle with sides  $a, b, c$  and angle  $\theta$  opposite to side  $c$ :*

$$c^2 = a^2 + b^2 - 2ab \cos \theta \quad (2.6)$$

Translating this geometric theorem into vector notation where the side lengths correspond to the norms of the vectors, we get:

$$\|y - x\|^2 = \|x\|^2 + \|y\|^2 - 2\|x\| \cdot \|y\| \cos \theta \quad (2.7)$$

This equation provides a critical link between the geometric angle  $\theta$  and the algebraic norms of the vectors.

#### Derivation of Inner Product

We can express the squared distance term  $\|y - x\|^2$  purely algebraically by expanding the components:

$$\|y - x\|^2 = \sum_{i=1}^n (x_i - y_i)^2 \quad (2.8)$$

$$= \sum_{i=1}^n (x_i^2 + y_i^2 - 2x_i y_i) \quad (2.9)$$

$$= \|x\|^2 + \|y\|^2 - 2 \sum_{i=1}^n x_i y_i \quad (2.10)$$

By comparing this expanded form with the result from the Law of Cosines derived previously, we can identify a corresponding interaction term. This term is so important that we give it a special name: the **Inner Product** (or dot product).

**Definition 2.6** (Inner Product). The inner product of two vectors  $x$  and  $y$  is defined as the sum of the products of their corresponding components:

$$x' y = \sum_{i=1}^n x_i y_i = \langle x, y \rangle \quad (2.11)$$

Thus, equating the geometric and algebraic forms yields the fundamental relationship:

$$x' y = \|x\| \cdot \|y\| \cos \theta \quad (2.12)$$

### 2.1.4 Coordinate (Scalar) Projection

The inner product allows us to calculate projections, which quantify how much of one vector “lies along” another. If we rearrange the cosine formula derived above, we can isolate the term that represents the length of the “shadow” cast by vector  $y$  onto vector  $x$ .

The length of this projection is given by:

$$\|y\| \cos \theta = \frac{x' y}{\|x\|} \quad (2.13)$$

This expression can be interpreted as the inner product of  $y$  with the normalized (unit) vector in the direction of  $x$ :

$$\text{Scalar Projection} = \left\langle \frac{x}{\|x\|}, y \right\rangle \quad (2.14)$$

### 2.1.5 Vector Projection Formula

The scalar projection only gives us a magnitude (a number). To define the projection as a vector in the same space, we need to multiply this scalar magnitude by the direction of the vector we are projecting onto.

**Definition 2.7** (Vector Projection). The projection of vector  $y$  onto vector  $x$ , denoted  $\hat{y}$ , is calculated as:

$$\text{Projection Vector} = (\text{Length}) \cdot (\text{Direction}) \quad (2.15)$$

$$\hat{y} = \left( \frac{x'y}{||x||} \right) \cdot \frac{x}{||x||} \quad (2.16)$$

This is often written compactly by combining the denominators:

$$\hat{y} = \frac{x'y}{||x||^2} x \quad (2.17)$$

### 2.1.6 Perpendicularity (Orthogonality)

A special case of the angle between vectors arises when  $\theta = 90^\circ$ . This geometric concept of perpendicularity is central to the theory of projections and least squares.

**Definition 2.8** (Perpendicularity). Two vectors are defined as **perpendicular** (or orthogonal) if the angle between them is  $90^\circ$  ( $\pi/2$ ).

Since  $\cos(90^\circ) = 0$ , the condition for orthogonality simplifies to the inner product being zero:

$$x'y = 0 \iff x \perp y \quad (2.18)$$

**Example 2.1** (Orthogonal Vectors). Consider two vectors in  $\mathbb{R}^2$ :  $x = (1, 1)'$  and  $y = (1, -1)'$ .

$$x'y = 1(1) + 1(-1) = 1 - 1 = 0 \quad (2.19)$$

Since their inner product is zero, these vectors are orthogonal to each other.

### 2.1.7 Projection onto a Line (Subspace)

We can generalize the concept of projecting onto a single vector to projecting onto the entire line (a 1-dimensional subspace) defined by that vector.

**Definition 2.9** (Line Spanned by a Vector). The line space  $L(x)$ , or the space spanned by a vector  $x$ , is defined as the set of all scalar multiples of  $x$ :

$$L(x) = \{cx \mid c \in \mathbb{R}\} \quad (2.20)$$

The projection of  $y$  onto  $L(x)$ , denoted  $\hat{y}$ , is defined by the geometric property that it is the closest point on the line to  $y$ . This implies that the error vector (or residual) must be perpendicular to the line itself.

**Definition 2.10** (Projection onto a Line). A vector  $\hat{y}$  is the projection of  $y$  onto the line  $L(x)$  if:

1.  $\hat{y}$  lies on the line  $L(x)$  (i.e.,  $\hat{y} = cx$  for some scalar  $c$ ).
2. The residual vector  $(y - \hat{y})$  is perpendicular to the direction vector  $x$ .

**Derivation:** To find the value of the scalar  $c$ , we apply the orthogonality condition:

$$(y - \hat{y}) \perp x \implies x'(y - cx) = 0 \quad (2.21)$$

Expanding this inner product gives:

$$x'y - c(x'x) = 0 \quad (2.22)$$

Solving for  $c$ , we obtain:

$$c = \frac{x'y}{||x||^2} \quad (2.23)$$

This confirms the formula derived previously using the inner product geometry. It shows that the least squares principle (shortest distance) leads to the same result as the geometric projection.

#### Alternative Forms of the Projection Formula

We can express the projection vector  $\hat{y}$  in several equivalent ways to highlight different geometric interpretations.

**Definition 2.11** (Forms of Projection). The projection of  $y$  onto the vector  $x$  is given by:

$$\hat{y} = \frac{x'y}{||x||^2}x = \left\langle y, \frac{x}{||x||} \right\rangle \frac{x}{||x||} \quad (2.24)$$

This second form separates the components into:

$$\text{Projection} = (\text{Scalar Projection}) \times (\text{Unit Direction}) \quad (2.25)$$

### 2.1.8 Projection Matrix ( $P_x$ )

In linear models, it is often more convenient to view projection as a linear transformation applied to the vector  $y$ . This allows us to define a **Projection Matrix**.

We can rewrite the formula for  $\hat{y}$  by factoring out  $y$ :

$$\hat{y} = \text{proj}(y|x) = x \frac{x'y}{||x||^2} = \frac{xx'}{||x||^2} y \quad (2.26)$$

This leads to the definition of the projection matrix  $P_x$ .

**Definition 2.12** (Projection Matrix onto a Single Vector). The matrix  $P_x$  that projects any vector  $y$  onto the line spanned by  $x$  is defined as:

$$P_x = \frac{xx'}{||x||^2} \quad (2.27)$$

Using this matrix, the projection is simply:

$$\hat{y} = P_x y \quad (2.28)$$

If  $x \in \mathbb{R}^p$ , then  $P_x$  is a  $p \times p$  symmetric matrix.

Let's apply these concepts to a concrete example.

**Example 2.2** (Numerical Projection). Let  $y = (1, 3)'$  and  $x = (1, 1)'$ . We want to find the projection of  $y$  onto  $x$ .

**Method 1: Using the Vector Formula** First, calculate the inner products:

$$x'y = 1(1) + 1(3) = 4 \quad (2.29)$$

$$||x||^2 = 1^2 + 1^2 = 2 \quad (2.30)$$

Now, apply the formula:

$$\hat{y} = \frac{4}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 2 \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \quad (2.31)$$

**Method 2: Using the Projection Matrix** Construct the matrix  $P_x$ :

$$P_x = \frac{1}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{pmatrix} \quad (2.32)$$

Multiply by  $y$ :

$$\hat{y} = P_x y = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{pmatrix} \begin{pmatrix} 1 \\ 3 \end{pmatrix} = \begin{pmatrix} 0.5(1) + 0.5(3) \\ 0.5(1) + 0.5(3) \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \quad (2.33)$$

**Example: Projection onto the Ones Vector ( $j_n$ )**

A very common operation in statistics is calculating the sample mean. This can be viewed geometrically as a projection onto a specific vector.

**Example 2.3** (Projection onto the Ones Vector). Let  $y = (y_1, \dots, y_n)'$  be a data vector. Let  $j_n = (1, 1, \dots, 1)'$  be a vector of all ones.

The projection of  $y$  onto  $j_n$  is:

$$\text{proj}(y|j_n) = \frac{j_n' y}{||j_n||^2} j_n \quad (2.34)$$

Calculating the components:

$$j_n' y = \sum_{i=1}^n y_i \quad (\text{Sum of observations}) \quad (2.35)$$

$$||j_n||^2 = \sum_{i=1}^n 1^2 = n \quad (2.36)$$

Substituting these back:

$$\hat{y} = \frac{\sum y_i}{n} j_n = \bar{y} j_n = \begin{pmatrix} \bar{y} \\ \vdots \\ \bar{y} \end{pmatrix} \quad (2.37)$$

Thus, replacing a data vector with its mean vector is geometrically equivalent to projecting the data onto the line spanned by the vector of ones.

**2.1.9 Pythagorean Theorem**

The Pythagorean theorem generalizes from simple geometry to vector spaces using the concept of orthogonality defined by the inner product.

**Theorem 2.2** (Pythagorean Theorem). *If two vectors  $x$  and  $y$  are orthogonal (i.e.,  $x \perp y$  or  $x'y = 0$ ), then the squared length of their sum is equal to the sum of their squared lengths:*

$$||x + y||^2 = ||x||^2 + ||y||^2 \quad (2.38)$$

*Proof.* We expand the squared norm using the inner product:

$$\begin{aligned} ||x + y||^2 &= (x + y)'(x + y) \\ &= x'x + x'y + y'x + y'y \\ &= ||x||^2 + 2x'y + ||y||^2 \end{aligned} \quad (2.39)$$

Since  $x \perp y$ , the inner product  $x'y = 0$ . Thus, the term  $2x'y$  vanishes, leaving:

$$||x + y||^2 = ||x||^2 + ||y||^2 \quad (2.40)$$

□

The proof after defining inner product to represent  $\cos(\theta)$  is trivial. Figure 2.1 shows a geometric proof of the fundamental Pythagorean Theorem (aka □□□□).

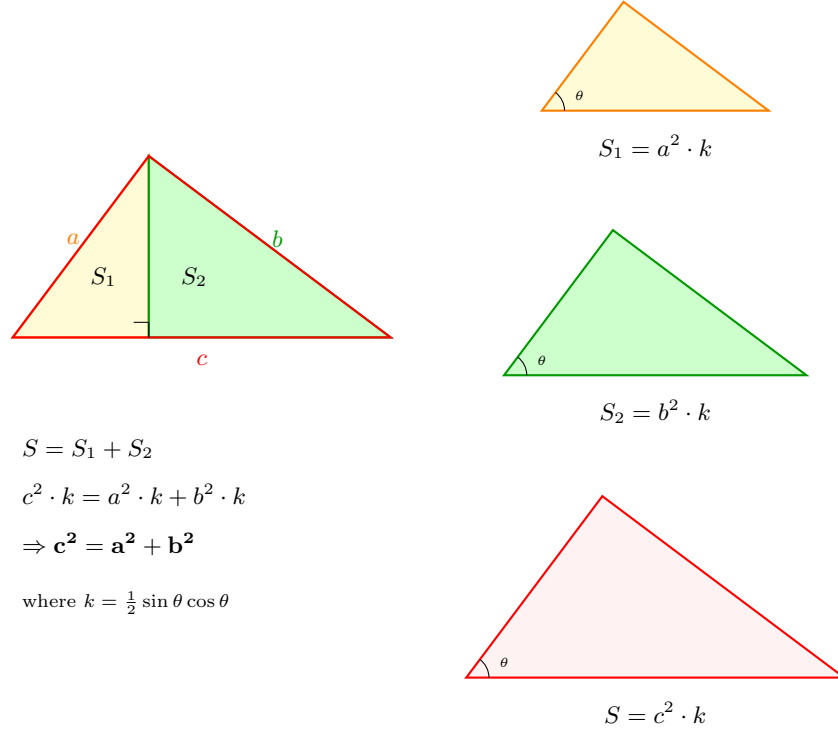


Figure 2.1: Proof of Pythagorean Theorem using Area Scaling

### 2.1.10 Least Square Property

One of the most important properties of the orthogonal projection is that it minimizes the distance between the vector  $y$  and the subspace (or line) onto which it is projected.

**Theorem 2.3** (Least Square Property). *Let  $\hat{y}$  be the projection of  $y$  onto the line  $L(x)$ . For any other vector  $y^*$  on the line  $L(x)$ , the distance from  $y$  to  $y^*$  is always greater than or equal to the distance from  $y$  to  $\hat{y}$ .*

$$||y - y^*|| \geq ||y - \hat{y}|| \quad (2.41)$$

*Proof.* Since both  $\hat{y}$  and  $y^*$  lie on the line  $L(x)$ , their difference  $(\hat{y} - y^*)$  also lies on  $L(x)$ . From the definition of projection, the residual  $(y - \hat{y})$  is orthogonal to the line  $L(x)$ . Therefore:

$$(y - \hat{y}) \perp (\hat{y} - y^*) \quad (2.42)$$

We can write the vector  $(y - y^*)$  as:

$$y - y^* = (y - \hat{y}) + (\hat{y} - y^*) \quad (2.43)$$

Applying the Pythagorean Theorem:

$$\|y - y^*\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - y^*\|^2 \quad (2.44)$$

Since  $\|\hat{y} - y^*\|^2 \geq 0$ , it follows that:

$$\|y - y^*\|^2 \geq \|y - \hat{y}\|^2 \quad (2.45)$$

□

## 2.2 Vector Space

We now generalize our discussion from lines to broader spaces.

**Definition 2.13** (Vector Space). A set  $V \subseteq \mathbb{R}^n$  is called a **Vector Space** if it is closed under vector addition and scalar multiplication:

1. **Closed under Addition:** If  $x_1 \in V$  and  $x_2 \in V$ , then  $x_1 + x_2 \in V$ .
2. **Closed under Scalar Multiplication:** If  $x \in V$ , then  $cx \in V$  for any scalar  $c \in \mathbb{R}$ .

It follows that the zero vector 0 must belong to any subspace (by choosing  $c = 0$ ).

### 2.2.1 Spanned Vector Space

The most common way to construct a vector space in linear models is by spanning it with a set of vectors.

**Definition 2.14** (Spanned Vector Space). Let  $x_1, \dots, x_p$  be a set of vectors in  $\mathbb{R}^n$ . The space spanned by these vectors, denoted  $L(x_1, \dots, x_p)$ , is the set of all possible linear combinations of them:

$$L(x_1, \dots, x_p) = \{r \mid r = c_1x_1 + \dots + c_px_p, \text{ for } c_i \in \mathbb{R}\} \quad (2.46)$$

### 2.2.2 Column Space and Row Space

When vectors are arranged into a matrix, we define specific spaces based on their columns and rows.

**Definition 2.15** (Column Space). For a matrix  $X = (x_1, \dots, x_p)$ , the **Column Space**, denoted  $\text{Col}(X)$ , is the vector space spanned by its columns:

$$\text{Col}(X) = L(x_1, \dots, x_p) \quad (2.47)$$

**Definition 2.16** (Row Space). The **Row Space**, denoted  $\text{Row}(X)$ , is the vector space spanned by the rows of the matrix  $X$ .



### 2.2.3 Linear Independence and Rank

Not all vectors in a spanning set contribute new dimensions to the space. This concept is captured by linear independence.

**Definition 2.17** (Linear Independence). A set of vectors  $x_1, \dots, x_p$  is said to be **Linearly Independent** if the only solution to the linear combination equation equal to zero is the trivial solution:

$$\sum_{i=1}^p c_i x_i = 0 \implies c_1 = c_2 = \dots = c_p = 0 \quad (2.48)$$

If there exist non-zero  $c_i$ 's such that sum is zero, the vectors are **Linearly Dependent**.

## 2.3 Rank of Matrices and Dim of Vector Space

**Definition 2.18** (Rank). The **Rank** of a matrix  $X$ , denoted  $\text{Rank}(X)$ , is the maximum number of linearly independent columns in  $X$ . This is equivalent to the dimension of the column space:

$$\text{Rank}(X) = \text{Dim}(\text{Col}(X)) \quad (2.49)$$

There are several fundamental properties regarding the rank of a matrix.

**Example 2.4** (Example of the Equality of Row and Col Rank). Consider the following  $3 \times 4$  matrix ( $n = 3, p = 4$ ):

$$X = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix} \quad (2.50)$$

Notice that the third row is the sum of the first two ( $r_3 = r_1 + r_2$ ).

**1. Row Rank and Basis  $U$**  The first two rows are linearly independent. We set the row rank  $r = 2$  and use these rows as our basis matrix  $U$  ( $2 \times 4$ ):

$$U = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix} \quad (2.51)$$

**2. Coefficient Matrix  $C$**  We express every row of  $X$  as a linear combination of the rows of  $U$ :

- Row 1:  $1 \cdot u_1 + 0 \cdot u_2$
- Row 2:  $0 \cdot u_1 + 1 \cdot u_2$
- Row 3:  $1 \cdot u_1 + 1 \cdot u_2$

These coefficients form the matrix  $C$  ( $3 \times 2$ ):

$$C = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix} \quad (2.52)$$

**3. The Decomposition ( $X = CU$ )** We verify that  $X$  is the product of  $C$  and  $U$ :

$$\underbrace{\begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix}}_{X \ (3 \times 4)} = \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}}_{C \ (3 \times 2)} \underbrace{\begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix}}_{U \ (2 \times 4)} \quad (2.53)$$

**4. Conclusion on Column Rank** The columns of  $X$  are linear combinations of the columns of  $C$ .

$$\text{Col}(X) \subseteq \text{Col}(C) \quad (2.54)$$

Since  $C$  has only 2 columns, the dimension of its column space (and thus  $X$ 's column space) cannot exceed 2.

$$\text{Dim}(\text{Col}(X)) \leq 2 \quad (2.55)$$

This confirms that Row Rank (2)  $\geq$  Column Rank. (By symmetry, they are equal).

**Theorem 2.4** (Row Rank equals Column Rank).

1. **Row Rank equals Column Rank:** *The dimension of the column space is equal to the dimension of the row space.*

$$\text{Dim}(\text{Col}(X)) = \text{Dim}(\text{Row}(X)) \implies \text{Rank}(X) = \text{Rank}(X') \quad (2.56)$$

2. **Bounds:** *For an  $n \times p$  matrix  $X$ :*

$$\text{Rank}(X) \leq \min(n, p) \quad (2.57)$$

### 2.3.1 Orthogonality to a Subspace

We can extend the concept of orthogonality from single vectors to entire subspaces.

**Definition 2.19** (Orthogonality to a Subspace). A vector  $y$  is orthogonal to a subspace  $V$  (denoted  $y \perp V$ ) if  $y$  is orthogonal to **every** vector  $x$  in  $V$ .

$$y \perp V \iff y'x = 0 \quad \forall x \in V \quad (2.58)$$

**Definition 2.20** (Orthogonal Complement). The set of all vectors that are orthogonal to a subspace  $V$  is called the **Orthogonal Complement** of  $V$ , denoted  $V^\perp$ .

$$V^\perp = \{y \in \mathbb{R}^n \mid y \perp V\} \quad (2.59)$$

### 2.3.2 Kernel (Null Space) and Image

For a matrix transformation defined by  $X$ , we define two key spaces: the Image (Column Space) and the Kernel (Null Space).

**Definition 2.21** (Image and Kernel).

1. **Image (Column Space):** The set of all possible outputs.

$$\text{Im}(X) = \text{Col}(X) = \{X\beta \mid \beta \in \mathbb{R}^p\} \quad (2.60)$$

2. **Kernel (Null Space):** The set of all inputs mapped to the zero vector.

$$\text{Ker}(X) = \{\beta \in \mathbb{R}^p \mid X\beta = 0\} \quad (2.61)$$

**Theorem 2.5** (Relationship between Kernel and Row Space). *The kernel of  $X$  is the orthogonal complement of the row space of  $X$ :*

$$\text{Ker}(X) = [\text{Row}(X)]^\perp \quad (2.62)$$

*Proof.* Let  $x \in \mathbb{R}^p$ .  $x \in \text{Ker}(X)$  if and only if  $Xx = 0$ . If we denote the rows of  $X$  as  $r'_1, \dots, r'_n$ , then the equation  $Xx = 0$  is equivalent to the system of equations:

$$\begin{pmatrix} r'_1 \\ \vdots \\ r'_n \end{pmatrix} x = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \iff r'_i x = 0 \text{ for all } i = 1, \dots, n \quad (2.63)$$

This means  $x$  is orthogonal to every row of  $X$ . Since the rows span the row space  $\text{Row}(X)$ , being orthogonal to every generator  $r_i$  implies  $x$  is orthogonal to the entire space  $\text{Row}(X)$ . Thus,  $\text{Ker}(X) = \{x \mid x \perp \text{Row}(X)\} = [\text{Row}(X)]^\perp$ .  $\square$

### 2.3.3 Nullity Theorem

There is a fundamental relationship between the dimensions of these spaces.

**Theorem 2.6** (Rank-Nullity Theorem). *For an  $n \times p$  matrix  $X$ :*

$$\text{Rank}(X) + \text{Nullity}(X) = p \quad (2.64)$$

where  $\text{Nullity}(X) = \text{Dim}(\text{Ker}(X))$ .

*Proof.* From the previous theorem, we established that the kernel is the orthogonal complement of the row space:

$$\text{Ker}(X) = [\text{Row}(X)]^\perp \quad (2.65)$$

Since the row space is a subspace of  $\mathbb{R}^p$ , the entire space can be decomposed into the direct sum of the row space and its orthogonal complement:

$$\mathbb{R}^p = \text{Row}(X) \oplus [\text{Row}(X)]^\perp = \text{Row}(X) \oplus \text{Ker}(X) \quad (2.66)$$

Taking the dimensions of these spaces:

$$\text{Dim}(\mathbb{R}^p) = \text{Dim}(\text{Row}(X)) + \text{Dim}(\text{Ker}(X)) \quad (2.67)$$

Substituting the definitions of Rank (dimension of row/column space) and Nullity:

$$p = \text{Rank}(X) + \text{Nullity}(X) \quad (2.68)$$

□

### Comparing Ranks via Kernel Containment

The Rank-Nullity Theorem provides a powerful and convenient tool for comparing the ranks of two matrices  $A$  and  $B$  (with the same number of columns) by inspecting their null spaces.

**Theorem 2.7** (Kernel Containment and Rank Inequality). *Let  $A$  and  $B$  be two matrices with  $p$  columns. If the kernel of  $A$  is contained within the kernel of  $B$ , then the rank of  $A$  is greater than or equal to the rank of  $B$ .*

$$\text{Ker}(A) \subseteq \text{Ker}(B) \implies \text{Rank}(A) \geq \text{Rank}(B) \quad (2.69)$$

*Proof.* From the subspace inclusion  $\text{Ker}(A) \subseteq \text{Ker}(B)$ , it follows that the dimension of the smaller space cannot exceed the dimension of the larger space:

$$\text{Nullity}(A) \leq \text{Nullity}(B) \quad (2.70)$$

Using the Rank-Nullity Theorem ( $\text{Rank} = p - \text{Nullity}$ ), we reverse the inequality:

$$p - \text{Nullity}(A) \geq p - \text{Nullity}(B) \quad (2.71)$$

$$\text{Rank}(A) \geq \text{Rank}(B) \quad (2.72)$$

□

### 2.3.4 Rank Inequalities

Understanding the bounds of the rank of matrix products is crucial for deriving properties of linear estimators.

**Theorem 2.8** (Rank of a Matrix Product). *Let  $X$  be an  $n \times p$  matrix and  $Z$  be a  $p \times k$  matrix. The rank of their product  $XZ$  is bounded by the rank of the individual matrices:*

$$\text{Rank}(XZ) \leq \min(\text{Rank}(X), \text{Rank}(Z)) \quad (2.73)$$

*Proof.* The columns of  $XZ$  are linear combinations of the columns of  $X$ . Thus, the column space of  $XZ$  is a subspace of the column space of  $X$ :

$$\text{Col}(XZ) \subseteq \text{Col}(X) \implies \text{Rank}(XZ) \leq \text{Rank}(X) \quad (2.74)$$

Similarly, the rows of  $XZ$  are linear combinations of the rows of  $Z$ . Thus, the row space of  $XZ$  is a subspace of the row space of  $Z$ :

$$\text{Row}(XZ) \subseteq \text{Row}(Z) \implies \text{Rank}(XZ) \leq \text{Rank}(Z) \quad (2.75)$$

□

### Rank and Invertible Matrices

Multiplying by an invertible (non-singular) matrix preserves the rank. This is a very useful property when manipulating linear equations.

**Theorem 2.9** (Rank with Non-Singular Multiplication). *Let  $A$  be an  $n \times n$  invertible matrix (i.e.,  $\text{Rank}(A) = n$ ) and  $X$  be an  $n \times p$  matrix. Then:*

$$\text{Rank}(AX) = \text{Rank}(X) \quad (2.76)$$

*Similarly, if  $B$  is a  $p \times p$  invertible matrix, then:*

$$\text{Rank}(XB) = \text{Rank}(X) \quad (2.77)$$

*Proof.* From the previous theorem, we know  $\text{Rank}(AX) \leq \text{Rank}(X)$ . Since  $A$  is invertible, we can write  $X = A^{-1}(AX)$ . Applying the theorem again:

$$\text{Rank}(X) = \text{Rank}(A^{-1}(AX)) \leq \text{Rank}(AX) \quad (2.78)$$

Thus,  $\text{Rank}(AX) = \text{Rank}(X)$ .

□

### 2.3.5 Rank of $X'X$ and $XX'$

The matrix  $X'X$  (the Gram matrix) appears in the normal equations for least squares ( $X'X\beta = X'y$ ). Its properties are closely tied to  $X$ .

**Theorem 2.10** (Rank of Gram Matrix). *For any real matrix  $X$ , the rank of  $X'X$  and  $XX'$  is the same as the rank of  $X$  itself:*

$$\text{Rank}(X'X) = \text{Rank}(X) \quad (2.79)$$

$$\text{Rank}(XX') = \text{Rank}(X) \quad (2.80)$$

*Proof.* We first show that the null space (kernel) of  $X$  is the same as the null space of  $X'X$ . If  $v \in \text{Ker}(X)$ , then  $Xv = 0 \implies X'Xv = 0 \implies v \in \text{Ker}(X'X)$ . Conversely, if  $v \in \text{Ker}(X'X)$ , then  $X'Xv = 0$ . Multiply by  $v'$ :

$$v'X'Xv = 0 \implies (Xv)'(Xv) = 0 \implies \|Xv\|^2 = 0 \implies Xv = 0 \quad (2.81)$$

So  $\text{Ker}(X) = \text{Ker}(X'X)$ . By the Rank-Nullity Theorem, since they have the same number of columns and same nullity, they must have the same rank.  $\square$

#### Column Space of $XX'$

Beyond just the rank, the column spaces themselves are related.

**Theorem 2.11** (Column Space Equivalence). *The column space of  $XX'$  is identical to the column space of  $X$ :*

$$\text{Col}(XX') = \text{Col}(X) \quad (2.82)$$

*Proof.*

1. **Forward ( $\subseteq$ ):** Let  $z \in \text{Col}(XX')$ . Then  $z = XX'w$  for some vector  $w$ . We can rewrite this as  $z = X(X'w)$ . Since  $z$  is a linear combination of columns of  $X$  (with coefficients  $X'w$ ),  $z \in \text{Col}(X)$ . Thus,  $\text{Col}(XX') \subseteq \text{Col}(X)$ .
2. **Equality via Rank:** From the previous theorem, we know that  $\text{Rank}(XX') = \text{Rank}(X)$ . Since  $\text{Col}(XX')$  is a subspace of  $\text{Col}(X)$  and they have the same finite dimension (Rank), the subspaces must be identical.

$\square$

**Implication:** This property ensures that for any  $y$ , the projection of  $y$  onto  $\text{Col}(X)$  lies in the same space as the projection onto  $\text{Col}(XX')$ . This is vital for the existence of solutions in generalized least squares.

## 2.4 Orthogonal Projection onto a Subspace

Let  $V$  be a subspace of  $\mathbb{R}^n$ . For any vector  $y \in \mathbb{R}^n$ , there exists a **unique** vector  $\hat{y} \in V$  such that the residual is orthogonal to the subspace:

$$(y - \hat{y}) \perp V \quad (2.83)$$

Equivalently:

$$\langle y - \hat{y}, v \rangle = 0 \quad \forall v \in V \quad (2.84)$$

### 2.4.1 Equivalence to Least Squares

The geometric definition of projection (orthogonality) is mathematically equivalent to the optimization problem of minimizing distance (least squares).

**Theorem 2.12** (Best Approximation Theorem (Least Squares Property)). *Let  $V$  be a subspace of  $\mathbb{R}^n$  and  $y \in \mathbb{R}^n$ . Let  $\hat{y}$  be the orthogonal projection of  $y$  onto  $V$ . Then  $\hat{y}$  is the closest point in  $V$  to  $y$ . That is, for any vector  $v \in V$  such that  $v \neq \hat{y}$ :*

$$\|y - \hat{y}\|^2 < \|y - v\|^2 \quad (2.85)$$

*Proof.* Let  $v$  be any vector in  $V$ . We can rewrite the difference vector  $y - v$  by adding and subtracting the projection  $\hat{y}$ :

$$y - v = (y - \hat{y}) + (\hat{y} - v) \quad (2.86)$$

Observe the properties of the two terms on the right-hand side:

1. **Residual:**  $(y - \hat{y})$  is orthogonal to  $V$  by definition.
2. **Difference in Subspace:** Since both  $\hat{y} \in V$  and  $v \in V$ , their difference  $(\hat{y} - v)$  is also in  $V$ .

Therefore, the two terms are orthogonal to each other:

$$(y - \hat{y}) \perp (\hat{y} - v) \quad (2.87)$$

Applying the Pythagorean Theorem:

$$\|y - v\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - v\|^2 \quad (2.88)$$

Since squared norms are non-negative, and  $\|\hat{y} - v\|^2 > 0$  (because  $v \neq \hat{y}$ ):

$$\|y - v\|^2 > \|y - \hat{y}\|^2 \quad (2.89)$$

The projection  $\hat{y}$  minimizes the squared error distance (and error distance itself).  $\square$

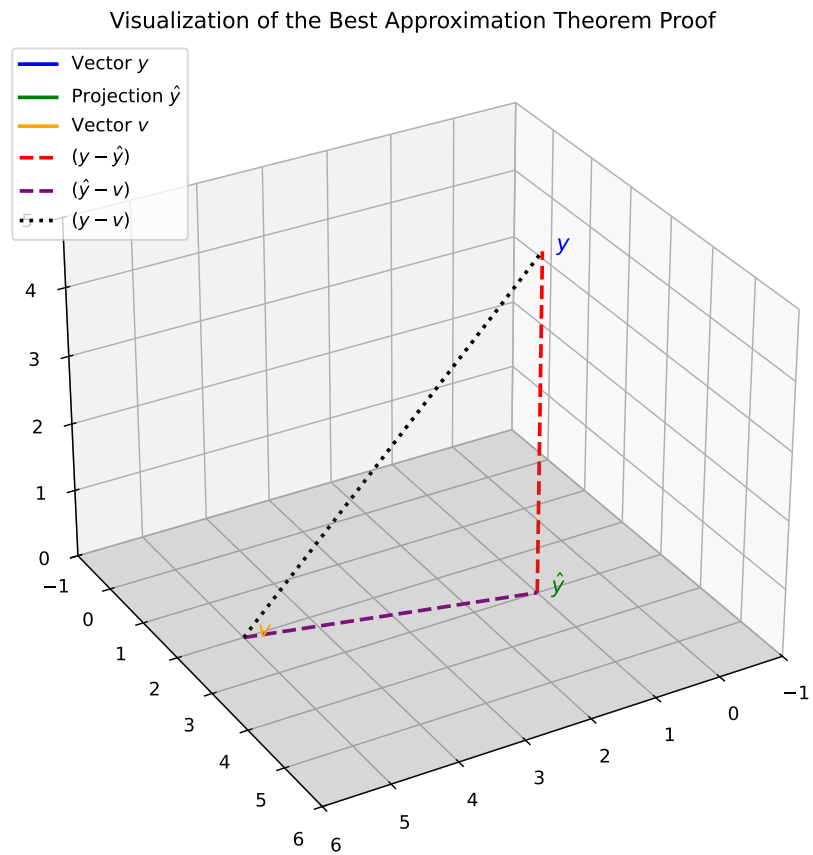


Figure 2.2: Visualization of the Best Approximation Theorem



## 2.4.2 Uniqueness of Projection

While the existence of a least-squares solution is guaranteed, we must also prove that there is only one such vector.

**Theorem 2.13** (Uniqueness of Orthogonal Projection). *For a given vector  $y$  and subspace  $V$ , the projection vector  $\hat{y}$  satisfying  $(y - \hat{y}) \perp V$  is unique.*

*Proof.* Assume there are two vectors  $\hat{y}_1 \in V$  and  $\hat{y}_2 \in V$  that both satisfy the orthogonality condition.

$$(y - \hat{y}_1) \perp V \quad \text{and} \quad (y - \hat{y}_2) \perp V \quad (2.90)$$

This means that for any  $v \in V$ , both inner products are zero:

$$\langle y - \hat{y}_1, v \rangle = 0 \quad (2.91)$$

$$\langle y - \hat{y}_2, v \rangle = 0 \quad (2.92)$$

Subtracting the second equation from the first:

$$\langle y - \hat{y}_1, v \rangle - \langle y - \hat{y}_2, v \rangle = 0 \quad (2.93)$$

Using the linearity of the inner product:

$$\langle (y - \hat{y}_1) - (y - \hat{y}_2), v \rangle = 0 \quad (2.94)$$

$$\langle \hat{y}_2 - \hat{y}_1, v \rangle = 0 \quad (2.95)$$

This equation holds for **all**  $v \in V$ . Since  $\hat{y}_1$  and  $\hat{y}_2$  are both in  $V$ , their difference  $d = \hat{y}_2 - \hat{y}_1$  must also be in  $V$ . We can therefore choose  $v = d = \hat{y}_2 - \hat{y}_1$ .

$$\langle \hat{y}_2 - \hat{y}_1, \hat{y}_2 - \hat{y}_1 \rangle = 0 \implies \|\hat{y}_2 - \hat{y}_1\|^2 = 0 \quad (2.96)$$

The only vector with a norm of zero is the zero vector itself.

$$\hat{y}_2 - \hat{y}_1 = 0 \implies \hat{y}_1 = \hat{y}_2 \quad (2.97)$$

Thus, the projection is unique. □

## 2.5 Projection via Orthonormal Basis ( $Q$ )

### 2.5.1 Orthonormal Basis

Before discussing projections onto general subspaces, we must formally define the coordinate system of a subspace, known as a basis.

**Definition 2.22** (Basis). A set of vectors  $\{x_1, \dots, x_k\}$  is a **Basis** for a vector space  $V$  if:

1. The vectors span the space:  $V = L(x_1, \dots, x_k)$ .

2. The vectors are linearly independent.

The number of vectors in a basis is unique and is defined as the **Dimension** of  $V$ .

Calculations become significantly simpler if we choose a basis with special geometric properties.

**Definition 2.23** (Orthonormal Basis). A basis  $\{q_1, \dots, q_k\}$  is called an **Orthonormal Basis** if:

1. **Orthogonal:** Each pair of vectors is perpendicular.

$$q'_i q_j = 0 \quad \text{for } i \neq j \quad (2.98)$$

2. **Normalized:** Each vector has unit length.

$$\|q_i\|^2 = q'_i q_i = 1 \quad (2.99)$$

Combining these, we write  $q'_i q_j = \delta_{ij}$  (Kronecker delta).

We now generalize the projection problem. Instead of projecting  $y$  onto a single line, we project it onto a subspace  $V$  of dimension  $k$ .

If we have an orthonormal basis  $\{q_1, \dots, q_k\}$  for  $V$ , the projection  $\hat{y}$  is simply the sum of the projections onto the individual basis vectors.

**Definition 2.24** (Projection Defined with Orthonormal Basis). The projection of  $y$  onto the subspace  $V = L(q_1, \dots, q_k)$  is:

$$\hat{y} = \sum_{i=1}^k \text{proj}(y|q_i) = \sum_{i=1}^k (q'_i y) q_i \quad (2.100)$$

Since the basis vectors are normalized, we do not need to divide by  $\|q_i\|^2$ .

**Theorem 2.14** (Projection via Orthonormal Basis). *Let  $\{q_1, \dots, q_k\}$  be an orthonormal basis for the subspace  $V \subseteq \mathbb{R}^n$ . The vector defined by the sum of individual projections:*

$$\hat{y} = \sum_{i=1}^k \langle y, q_i \rangle q_i \quad (2.101)$$

*is indeed the orthogonal projection of  $y$  onto  $V$ . That is, it satisfies  $(y - \hat{y}) \perp V$ .*

*Proof.* To prove this, we must check two conditions:

1.  $\hat{y} \in V$ : This is immediate because  $\hat{y}$  is a linear combination of the basis vectors  $\{q_1, \dots, q_k\}$ .

2.  $(y - \hat{y}) \perp V$ : It suffices to show that the error vector  $e = y - \hat{y}$  is orthogonal to every basis vector  $q_j$  (for  $j = 1, \dots, k$ ).

Let's calculate the inner product  $\langle y - \hat{y}, q_j \rangle$ :

$$\begin{aligned}\langle y - \hat{y}, q_j \rangle &= \langle y, q_j \rangle - \langle \hat{y}, q_j \rangle \\ &= \langle y, q_j \rangle - \left\langle \sum_{i=1}^k \langle y, q_i \rangle q_i, q_j \right\rangle \\ &= \langle y, q_j \rangle - \sum_{i=1}^k \langle y, q_i \rangle \underbrace{\langle q_i, q_j \rangle}_{\delta_{ij}}\end{aligned}\tag{2.102}$$

Since the basis is orthonormal,  $\langle q_i, q_j \rangle$  is 1 if  $i = j$  and 0 otherwise. Thus, the summation collapses to a single term where  $i = j$ :

$$\begin{aligned}\langle y - \hat{y}, q_j \rangle &= \langle y, q_j \rangle - \langle y, q_j \rangle \cdot 1 \\ &= 0\end{aligned}\tag{2.103}$$

Since  $(y - \hat{y})$  is orthogonal to every basis vector  $q_j$ , it is orthogonal to the entire subspace  $V$ . Thus,  $\hat{y}$  is the unique orthogonal projection.

□

## 2.5.2 Projection Matrix via Orthonormal Basis ( $Q$ )

### Matrix Form with Orthonormal Basis

We can express the summation formula for  $\hat{y}$  compactly using matrix notation.

Let  $Q$  be an  $n \times k$  matrix whose columns are the orthonormal basis vectors  $q_1, \dots, q_k$ .

$$Q = (q_1 \quad q_2 \quad \dots \quad q_k)\tag{2.104}$$

Properties of  $Q$ :

- $Q'Q = I_k$  (Identity matrix of size  $k \times k$ ).
- $QQ'$  is **not** necessarily  $I_n$  (unless  $k = n$ ).

**Definition 2.25** (Projection Matrix in Terms of  $Q$ ). The projection  $\hat{y}$  can be written as:

$$\hat{y} = (q_1 \quad \dots \quad q_k) \begin{pmatrix} q_1' y \\ \vdots \\ q_k' y \end{pmatrix} = Q(Q'y) = (QQ')y\tag{2.105}$$

Thus, the projection matrix  $P$  onto the subspace  $V$  is:

$$P = QQ'\tag{2.106}$$

## Properties of Projection Matrices

We have defined the projection matrix as  $P = X(X'X)^{-1}X'$  (or  $P = QQ'$  for orthonormal bases). All orthogonal projection matrices share two fundamental algebraic properties.

**Theorem 2.15** (Symmetry and Idempotence). *A square matrix  $P$  represents an orthogonal projection onto some subspace if and only if it satisfies:*

1. **Idempotence:**  $P^2 = P$  (Applying the projection twice is the same as applying it once).
2. **Symmetry:**  $P' = P$ .

*Proof:* If  $\hat{y} = Py$  is already in the subspace  $\text{Col}(X)$ , then projecting it again should not change it.

$$P(Py) = Py \implies P^2y = Py \quad \forall y \quad (2.107)$$

Thus,  $P^2 = P$ . □

## Example: ANOVA (Analysis of Variance)

One of the most common applications of projection is in Analysis of Variance (ANOVA). We can view the calculation of group means as a projection onto a subspace defined by group indicator variables.

**Example 2.5** (Finding Projection for One-way ANOVA). Consider a one-way ANOVA model with  $k$  groups:

$$y_{ij} = \mu_i + \epsilon_{ij} \quad (2.108)$$

where  $i \in \{1, \dots, k\}$  represents the group and  $j \in \{1, \dots, n_i\}$  represents the observation within the group. Let  $N = \sum_{i=1}^k n_i$  be the total number of observations.

**1. Matrix Definitions** We define the data vector  $y$  and the design matrix  $X$  as follows:

- **Data Vector ( $y$ ):** An  $N \times 1$  vector containing all observations stacked by group:

$$y = \begin{pmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ y_{21} \\ \vdots \\ y_{kn_k} \end{pmatrix} \quad (2.109)$$

- **Design Matrix ( $X$ ):** An  $N \times k$  matrix constructed from  $k$  column vectors,  $X = (x_1, x_2, \dots, x_k)$ . Each vector  $x_g$  is an **indicator variable** (dummy variable) for group  $g$ :

$$x_g = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \quad \leftarrow \text{Entries are 1 if observation belongs to group } g \quad (2.110)$$

**2. Orthogonality** These column vectors  $x_1, \dots, x_k$  are mutually orthogonal because no observation can belong to two groups at once. The dot product of any two distinct columns is zero:

$$\langle x_g, x_h \rangle = 0 \quad \text{for } g \neq h \quad (2.111)$$

This allows us to find the projection onto the column space of  $X$  by simply summing the projections onto each column individually.

**3. Calculating Individual Projections** For a specific group vector  $x_g$ , the projection is:

$$\text{proj}(y|x_g) = \frac{\langle y, x_g \rangle}{\langle x_g, x_g \rangle} x_g \quad (2.112)$$

We calculate the two scalar terms:

- **Denominator** ( $\langle x_g, x_g \rangle$ ): The sum of squared elements of  $x_g$ . Since  $x_g$  contains  $n_g$  ones and zeros elsewhere:

$$\langle x_g, x_g \rangle = \sum \mathbb{1}_{\{i=g\}}^2 = n_g \quad (2.113)$$

- **Numerator** ( $\langle y, x_g \rangle$ ): The dot product sums only the  $y$  values belonging to group  $g$ :

$$\langle y, x_g \rangle = \sum_{i,j} y_{ij} \cdot \mathbb{1}_{\{i=g\}} = \sum_{j=1}^{n_g} y_{gj} = y_g. \quad (\text{Group Total}) \quad (2.114)$$

**4. The Resulting Projection** Substituting these back into the formula gives the coefficient for the vector  $x_g$ :

$$\text{proj}(y|x_g) = \frac{y_g}{n_g} x_g = \bar{y}_g \cdot x_g \quad (2.115)$$

The total projection  $\hat{y}$  is the sum over all groups:

$$\hat{y} = \sum_{g=1}^k \bar{y}_g \cdot x_g \quad (2.116)$$

This confirms that the fitted value for any specific observation  $y_{ij}$  is simply its group mean  $\bar{y}_i$ .

### 2.5.3 Gram-Schmidt Process

To use the simplified formula  $P = QQ'$ , we need an orthonormal basis. The Gram-Schmidt process provides a method to construct such a basis from any set of linearly independent vectors.

**Gram-Schmidt Process** Given linearly independent vectors  $x_1, \dots, x_p$ :

1. **Step 1:** Normalize the first vector.

$$q_1 = \frac{x_1}{\|x_1\|} \quad (2.117)$$

2. **Step 2:** Project  $x_2$  onto  $q_1$  and subtract it to find the orthogonal component.

$$v_2 = x_2 - (x_2'q_1)q_1 \quad (2.118)$$

Then normalize:

$$q_2 = \frac{v_2}{||v_2||} \quad (2.119)$$

3. **Step k:** Subtract the projections onto all previous  $q$  vectors.

$$v_k = x_k - \sum_{j=1}^{k-1} (x_k'q_j)q_j \quad (2.120)$$

$$q_k = \frac{v_k}{||v_k||} \quad (2.121)$$

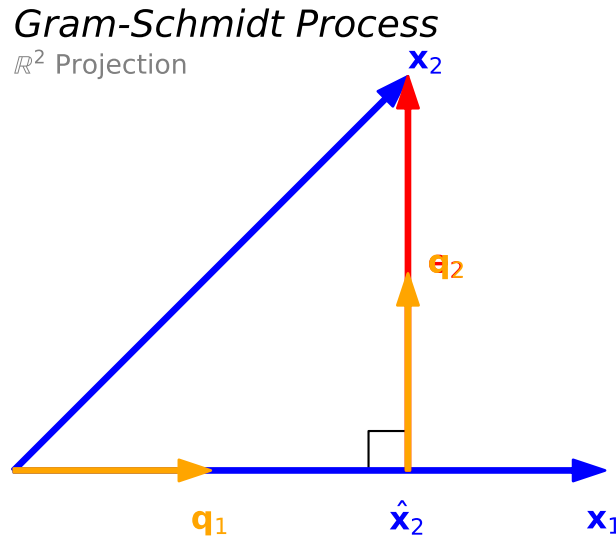


Figure 2.3: Gram-Schmidt Process: Projecting  $x_2$  onto  $x_1$

This process leads to the **QR Decomposition** of a matrix:  $X = QR$ , where  $Q$  is orthogonal and  $R$  is upper triangular.

## 2.6 Hat Matrix (Projection Matrix via $X$ )

### 2.6.1 Norm Equations

Let  $X = (x_1, \dots, x_p)$  be an  $n \times p$  matrix, where each column  $x_j$  is a predictor vector.

We want to project the target vector  $y$  onto the column space  $\text{Col}(X)$ . This is equivalent to finding a coefficient vector  $\beta \in \mathbb{R}^p$  such that the error vector (residual) is orthogonal to the entire subspace  $\text{Col}(X)$ .

$$y - X\beta \perp \text{Col}(X) \quad (2.122)$$

Since the columns of  $X$  span the subspace, the residual must be orthogonal to **every** column vector  $x_j$  individually:

$$y - X\beta \perp x_j \quad \text{for } j = 1, \dots, p \quad (2.123)$$

Writing this geometric condition as an algebraic dot product (where  $x_j'$  denotes the transpose):

$$x_j'(y - X\beta) = 0 \quad \text{for each } j \quad (2.124)$$

We can stack these  $p$  separate linear equations into a single matrix equation. Since the rows of  $X'$  are the columns of  $X$ , this becomes:

$$\begin{pmatrix} x_1' \\ \vdots \\ x_p' \end{pmatrix} (y - X\beta) = \mathbf{0} \implies X'(y - X\beta) = 0 \quad (2.125)$$

Finally, we distribute the matrix transpose and rearrange terms to solve for  $\beta$ :

$$\begin{aligned} X'y - X'X\beta &= 0 \\ X'X\beta &= X'y \end{aligned} \quad (2.126)$$

This system is known as the **Normal Equations**.

**Theorem 2.16** (Least Squares Estimator). *If  $X'X$  is invertible (i.e.,  $X$  has full column rank), the unique solution for  $\beta$  is:*

$$\hat{\beta} = (X'X)^{-1}X'y \quad (2.127)$$

### 2.6.2 Hat Matrix

Substituting the estimator  $\hat{\beta}$  back into the equation for  $\hat{y}$  gives us the projection matrix.

**Definition 2.26** (Hat Matrix). The projection of  $y$  onto  $\text{Col}(X)$  is given by:

$$\hat{y} = X\hat{\beta} = X(X'X)^{-1}X'y \quad (2.128)$$

Thus, the hat matrix  $H$  is defined as:

$$H = X(X'X)^{-1}X' \quad (2.129)$$

### 2.6.3 Equivalence of Hat Matrix and $QQ'$

If we use the QR decomposition such that  $X = QR$ , where the columns of  $Q$  form an orthonormal basis for  $\text{Col}(X)$ , the formula simplifies significantly.

Recall that for orthonormal columns,  $Q'Q = I$ . Substituting  $X = QR$  into the general formula:

$$\begin{aligned} H &= QR((QR)'(QR))^{-1}(QR)' \\ &= QR(R'Q'QR)^{-1}R'Q' \\ &= QR(\underbrace{R'Q'Q}_I R)^{-1}R'Q' \\ &= QR(R'R)^{-1}R'Q' \\ &= QR R^{-1}(R')^{-1}R'Q' \\ &= Q \underbrace{RR^{-1}}_I \underbrace{(R')^{-1}R'}_I Q' \\ &= QQ' \end{aligned} \quad (2.130)$$

This confirms that  $H = QQ'$  is consistent with the general formula  $H = X(X'X)^{-1}X'$ .

### 2.6.4 Properties of Hat Matrix

We revisit the properties of projection matrices in this general context.

**Theorem 2.17** (Properties of Hat Matrix). The matrix  $H = X(X'X)^{-1}X'$  satisfies:

1. **Symmetric:**  $H' = H$
2. **Idempotent:**  $H^2 = H$
3. **Trace:** The trace of a projection matrix equals the dimension of the subspace it projects onto.

$$\text{tr}(H) = \text{tr}(X(X'X)^{-1}X') = \text{tr}((X'X)^{-1}X'X) = \text{tr}(I_p) = p \quad (2.131)$$



## 2.7 Projection Defined with Orthogonal Projection Matrix

Projection don't have to be defined with a subspace or a matrix  $X$  as we discussed before. Projection matrix is a self-contained definition of the subspace it projects onto.

### 2.7.1 Orthogonal Projection Matrix

**Definition 2.27** (Orthogonal Projection Matrix). A square matrix  $P$  is called an **orthogonal projection matrix** if it satisfies two conditions:

1. **Symmetry:**  $P^\top = P$
2. **Idempotency:**  $P^2 = P$

**Theorem 2.18** (Projection onto Column Space). *If a matrix  $P$  is symmetric and idempotent, then  $P$  represents the orthogonal projection onto its column space,  $\text{Col}(P)$ .*

*Specifically, for any vector  $y$ , the vector  $\hat{y} = Py$  is the unique vector in  $\text{Col}(P)$  such that the residual  $e = y - \hat{y}$  is orthogonal to  $\text{Col}(P)$ .*

*Proof.* Let  $y \in \mathbb{R}^n$ . We decompose  $y$  as  $y = Py + (I - P)y$ . We must show that the residual term  $(I - P)y$  is orthogonal to any vector  $z \in \text{Col}(P)$ .

Since  $z \in \text{Col}(P)$ , there exists a vector  $x$  such that  $z = Px$ . The inner product between  $z$  and the residual is:

$$\langle z, (I - P)y \rangle = z^\top (I - P)y = (Px)^\top (I - P)y \quad (2.132)$$

Using the matrix transpose property  $(AB)^\top = B^\top A^\top$ , we rewrite Equation 2.132 as:

$$\langle z, (I - P)y \rangle = x^\top P^\top (I - P)y \quad (2.133)$$

Since  $P$  is symmetric ( $P^\top = P$ ), we can substitute  $P$  for  $P^\top$  in Equation 2.133:

$$\langle z, (I - P)y \rangle = x^\top P(I - P)y = x^\top (P - P^2)y \quad (2.134)$$

Finally, utilizing the idempotency of  $P$  (where  $P^2 = P$ ), the expression in Equation 2.134 simplifies to 0:

$$x^\top (P - P)y = x^\top (0)y = 0 \quad (2.135)$$

Since the inner product is 0, the residual is orthogonal to every vector in  $\text{Col}(P)$ . Thus,  $P$  is the orthogonal projector.  $\square$

## 2.7.2 Projection onto Complement Space

**Theorem 2.19** (Projection onto Orthogonal Complement). *Let  $P$  be an orthogonal projection matrix. The matrix  $M$  defined as:*

$$M = I - P \quad (2.136)$$

*is the orthogonal projection matrix onto the orthogonal complement of the column space of  $P$ , denoted  $\text{Col}(P)^\perp$ .*

*Proof. 1. Symmetry and Idempotency* Since  $P$  is a projection matrix,  $P^\top = P$  and  $P^2 = P$ . We verify these properties for  $M$ :

$$M^\top = (I - P)^\top = I - P^\top = I - P = M \quad (2.137)$$

$$M^2 = (I - P)(I - P) = I - 2P + P^2 = I - 2P + P = I - P = M \quad (2.138)$$

By Equation 2.137 and Equation 2.138,  $M$  is symmetric and idempotent, so it is an orthogonal projection matrix.

**2. Identifying the Subspace** By Theorem 2.18,  $M$  projects onto its own column space,  $\text{Col}(M)$ . A vector  $v$  is in  $\text{Col}(M)$  if and only if it is fixed by the projection ( $Mv = v$ ).

$$Mv = v \quad (2.139)$$

Substituting  $M = I - P$  into Equation 2.139 gives:

$$(I - P)v = v \quad (2.140)$$

Rearranging Equation 2.140, we find the condition for  $v$ :

$$v - Pv = v \implies Pv = 0 \quad (2.141)$$

The condition  $Pv = 0$  in Equation 2.141 implies that  $v$  belongs to the null space of  $P$ , denoted  $\text{Null}(P)$ . By the Fundamental Theorem of Linear Algebra for symmetric matrices, the null space is the orthogonal complement of the column space:

$$\text{Null}(P) = \text{Col}(P^\top)^\perp = \text{Col}(P)^\perp \quad (2.142)$$

Thus, the image of  $M$  is exactly  $\text{Col}(P)^\perp$ .  $\square$

**Exercise 2.1** (Column Space of the Hat Matrix). Let  $H = X(X^\top X)^{-1}X^\top$  be the hat matrix.

1. Prove that the column space of  $H$  is identical to the column space of  $X$ :

$$\text{Col}(H) = \text{Col}(X) \quad (2.143)$$

2. Using the result above, show that the column space of the residual maker matrix  $M = I - H$  is the orthogonal complement of  $\text{Col}(X)$ :

$$\text{Col}(M) = \text{Col}(X)^\perp \quad (2.144)$$

## Solutions

**1. Equivalence of Column Spaces** To prove  $\text{Col}(H) = \text{Col}(X)$ , we show inclusion in both directions.

- **Forward** ( $\text{Col}(H) \subseteq \text{Col}(X)$ ): By definition,  $H = X[(X^\top X)^{-1}X^\top]$ . Any column of  $H$  is a linear combination of the columns of  $X$  (weighted by the matrix in brackets). Therefore, any vector in the image of  $H$  must lie in  $\text{Col}(X)$ .
- **Reverse** ( $\text{Col}(X) \subseteq \text{Col}(H)$ ): Take any vector  $v \in \text{Col}(X)$ . By definition,  $v = Xb$  for some vector  $b$ . Apply  $H$  to  $v$ :

$$Hv = X(X^\top X)^{-1}X^\top(Xb) = X(X^\top X)^{-1}(X^\top X)b = X(I)b = Xb = v \quad (2.145)$$

Since  $Hv = v$ , the vector  $v$  lies in the column space of  $H$  (specifically, it is an eigenvector with eigenvalue 1).

Since both inclusions hold,  $\text{Col}(H) = \text{Col}(X)$ .

**2. Orthogonal Complements** From part 1, we know the subspaces are identical. Therefore, their orthogonal complements must also be identical:

$$\text{Col}(H)^\perp = \text{Col}(X)^\perp \quad (2.146)$$

We previously established in Theorem 2.19 that for any projection matrix  $P$ , the complement projection  $M = I - P$  projects onto  $\text{Col}(P)^\perp$ . Substituting  $H$  for  $P$ :

$$\text{Col}(M) = \text{Col}(H)^\perp \quad (2.147)$$

Combining these results gives the required equality:

$$\text{Col}(M) = \text{Col}(X)^\perp \quad (2.148)$$

## 2.8 Projection onto Nested Subspaces

### 2.8.1 Nested Models and Subspaces

In hypothesis testing (like comparing a null model to an alternative model), we often deal with nested subspaces.

**Definition 2.28** (Nested Models). Consider two models:

1. **Reduced Model** ( $M_0$ ):  $y \in \text{Col}(X_0)$
2. **Full Model** ( $M_1$ ):  $y \in \text{Col}(X_1)$

We say the models are nested if the column space of the reduced model is contained entirely within the column space of the full model:

$$\text{Col}(X_0) \subseteq \text{Col}(X_1) \quad (2.149)$$

Usually,  $X_1$  is constructed by adding columns to  $X_0$ :  $X_1 = [X_0, X_{\text{new}}]$ .

## 2.8.2 Projections onto Nested Subspaces

Let  $P_0$  be the projection matrix onto  $\text{Col}(X_0)$  and  $P_1$  be the projection matrix onto  $\text{Col}(X_1)$ . Since  $\text{Col}(X_0) \subseteq \text{Col}(X_1)$ , we have important relationships between these matrices.

**Theorem 2.20** (Composition of Projections). *If  $\text{Col}(P_0) \subseteq \text{Col}(P_1)$ , then:*

1.  $P_1 P_0 = P_0$  (Projecting onto the small space, then the large space, keeps you in the small space).
2.  $P_0 P_1 = P_0$  (Projecting onto the large space, then the small space, is the same as just projecting onto the small space).

*Proof. 1. Proof of  $P_1 P_0 = P_0$ :* For any vector  $y \in \mathbb{R}^n$ , the vector  $v = P_0 y$  lies in  $\text{Col}(X_0)$ . Since  $\text{Col}(X_0) \subseteq \text{Col}(X_1)$ , the vector  $v$  also lies in  $\text{Col}(X_1)$ . A projection matrix  $P_1$  acts as the identity operator for any vector already in its column space. Therefore,  $P_1 v = v$ . Substituting  $v = P_0 y$ , we get  $P_1 P_0 y = P_0 y$  for all  $y$ . Thus,  $P_1 P_0 = P_0$ .

**2. Proof of  $P_0 P_1 = P_0$ :** Take the transpose of the previous result ( $P_1 P_0 = P_0$ ).

$$(P_1 P_0)' = P_0' \quad (2.150)$$

Using the property that projection matrices are symmetric ( $P' = P$ ):

$$P_0' P_1' = P_0' \implies P_0 P_1 = P_0 \quad (2.151)$$

□

### Difference of Projections

The difference between the two projection matrices,  $P_1 - P_0$ , is itself a projection matrix.

**Theorem 2.21** (Difference Projection). *The matrix  $P_\Delta = P_1 - P_0$  is an orthogonal projection matrix onto the subspace  $\text{Col}(X_1) \cap \text{Col}(X_0)^\perp$ . This subspace represents the “extra” information in the full model that is orthogonal to the reduced model.*

**Properties:**

1. **Symmetric:**  $(P_1 - P_0)' = P_1 - P_0$ .
2. **Idempotent:**  $(P_1 - P_0)(P_1 - P_0) = P_1 - P_0 P_1 - P_1 P_0 + P_0 = P_1 - P_0 - P_0 + P_0 = P_1 - P_0$ .
3. **Orthogonality:**  $(P_1 - P_0)P_0 = P_1 P_0 - P_0 = P_0 - P_0 = 0$ .

*Proof. 1. Symmetry:* Since  $P_1$  and  $P_0$  are symmetric:  $(P_1 - P_0)' = P_1' - P_0' = P_1 - P_0$ .

**2. Idempotency:**

$$\begin{aligned} (P_1 - P_0)^2 &= (P_1 - P_0)(P_1 - P_0) \\ &= P_1^2 - P_1 P_0 - P_0 P_1 + P_0^2 \end{aligned} \quad (2.152)$$

Using the projection properties ( $P^2 = P$ ) and the nested property ( $P_1 P_0 = P_0$  and  $P_0 P_1 = P_0$ ):

$$= P_1 - P_0 - P_0 + P_0 = P_1 - P_0 \quad (2.153)$$

### 3. Orthogonality to $P_0$ :

$$(P_1 - P_0)P_0 = P_1P_0 - P_0^2 = P_0 - P_0 = 0 \quad (2.154)$$

Since  $(P_1 - P_0)$  is symmetric and idempotent, it is an orthogonal projection matrix. Since it is orthogonal to  $P_0$  (the space of  $M_0$ ) but is derived from  $P_1$ , it projects onto the subspace of  $M_1$  that is orthogonal to  $M_0$ .  $\square$

### 2.8.3 Decomposition of Projections and their Sum Squares

**Theorem 2.22** (Orthogonal Decomposition). *Let  $M_0 \subset M_1$  be two nested linear models with corresponding design matrices  $X_0$  and  $X_1$  such that  $\text{Col}(X_0) \subset \text{Col}(X_1)$ . Let  $P_0$  and  $P_1$  be the orthogonal projection matrices onto  $\text{Col}(X_0)$  and  $\text{Col}(X_1)$  respectively.*

*For any observation vector  $y$ , we have the decomposition:*

$$y = \underbrace{P_0 y}_{\hat{y}_0} + \underbrace{(P_1 - P_0)y}_{\hat{y}_1 - \hat{y}_0} + \underbrace{(I - P_1)y}_{y - \hat{y}_1} \quad (2.155)$$

#### **Geometric Interpretation:**

1.  $\hat{y}_0 \in \text{Col}(X_0)$ : The fit of the reduced model.
2.  $(\hat{y}_1 - \hat{y}_0) \in \text{Col}(X_0)^\perp \cap \text{Col}(X_1)$ : The additional fit provided by  $M_1$  over  $M_0$ .
3.  $(y - \hat{y}_1) \in \text{Col}(X_1)^\perp$ : The projection of  $y$  onto the **orthogonal complement** of  $\text{Col}(X_1)$ .

The three component vectors are mutually orthogonal. Consequently, their squared norms sum to the total squared norm:

$$\|y\|^2 = \|\hat{y}_0\|^2 + \|\hat{y}_1 - \hat{y}_0\|^2 + \|y - \hat{y}_1\|^2 \quad (2.156)$$

*Proof.* **1. Definitions** We define the three components as vectors  $v_1, v_2, v_3$ :

- $v_1 = \hat{y}_0 = P_0 y$ .
- $v_2 = \hat{y}_1 - \hat{y}_0 = (P_1 - P_0)y$ .
- $v_3 = y - \hat{y}_1 = (I - P_1)y$ .

– **Note:** Since  $P_1$  projects onto  $\text{Col}(X_1)$ , the matrix  $(I - P_1)$  projects onto the **orthogonal complement**  $\text{Col}(X_1)^\perp$ . Thus,  $v_3 \in \text{Col}(I - P_1)$ .

Note that since  $\text{Col}(X_0) \subset \text{Col}(X_1)$ , we have the property  $P_1 P_0 = P_0 P_1 = P_0$ . (Projecting onto the smaller subspace  $M_0$  is unchanged if we first project onto the enclosing subspace  $M_1$ ).

**2. Orthogonality of  $v_1$  and  $v_2$**  We check the inner product  $\langle v_1, v_2 \rangle = v_1' v_2$ :

$$\begin{aligned} v_1' v_2 &= (P_0 y)' (P_1 - P_0) y \\ &= y' P_0' (P_1 - P_0) y \\ &= y' (P_0 P_1 - P_0^2) y \quad (\text{Since } P_0 \text{ is symmetric}) \\ &= y' (P_0 - P_0) y \quad (\text{Since } P_0 P_1 = P_0 \text{ and } P_0^2 = P_0) \\ &= 0 \end{aligned} \quad (2.157)$$

**3. Orthogonality of  $(v_1 + v_2)$  and  $v_3$**  Note that  $v_1 + v_2 = P_1 y = \hat{y}_1$ . We check if the total fit  $\hat{y}_1$  is orthogonal to the residual  $v_3$ :

$$\begin{aligned}
\hat{y}_1' v_3 &= (P_1 y)' (I - P_1) y \\
&= y' P_1 (I - P_1) y \\
&= y' (P_1 - P_1^2) y \\
&= y' (P_1 - P_1) y \\
&= 0
\end{aligned} \tag{2.158}$$

Since  $\hat{y}_1$  is orthogonal to  $v_3$ , and  $\hat{y}_0$  is a component of  $\hat{y}_1$ , it follows that all three pieces are mutually orthogonal.

**4. Sum of Squares** By the Pythagorean theorem applied twice to these orthogonal vectors, the equality of squared norms follows immediately.  $\square$

**Example 2.6** (ANOVA Sum Squares). We apply the **Nested Model Theorem** ( $M_0 \subset M_1$ ) to the One-way ANOVA setting.

### 1. Notation and Definitions

Consider a dataset with  $k$  groups. Let  $i = 1, \dots, k$  index the groups, and  $j = 1, \dots, n_i$  index the observations within group  $i$ .

- $N$ : Total number of observations,  $N = \sum_{i=1}^k n_i$ .
- $y_{ij}$ : The  $j$ -th observation in the  $i$ -th group.
- $\bar{y}_i$ : The sample mean of group  $i$ .

$$\bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij} \tag{2.159}$$

- $\bar{y}_{..}$ : The grand mean of all observations.

$$\bar{y}_{..} = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij} \tag{2.160}$$

### 2. The Data and Projection Vectors

Table 2.1: ANOVA Vectors: Data, Null Model, and Full Model

Observation ( $y$ )	Null Projection ( $\hat{y}_0$ )	Full Projection ( $\hat{y}_1$ )
$\begin{pmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ \vdots \\ y_{k1} \\ \vdots \\ y_{kn_k} \end{pmatrix}$	$\begin{pmatrix} \bar{y}_{..} \\ \vdots \\ \bar{y}_{..} \\ \vdots \\ \bar{y}_{..} \\ \vdots \\ \bar{y}_{..} \end{pmatrix}$	$\begin{pmatrix} \bar{y}_{1.} \\ \vdots \\ \bar{y}_{1.} \\ \vdots \\ \bar{y}_{k.} \\ \vdots \\ \bar{y}_{k.} \end{pmatrix}$

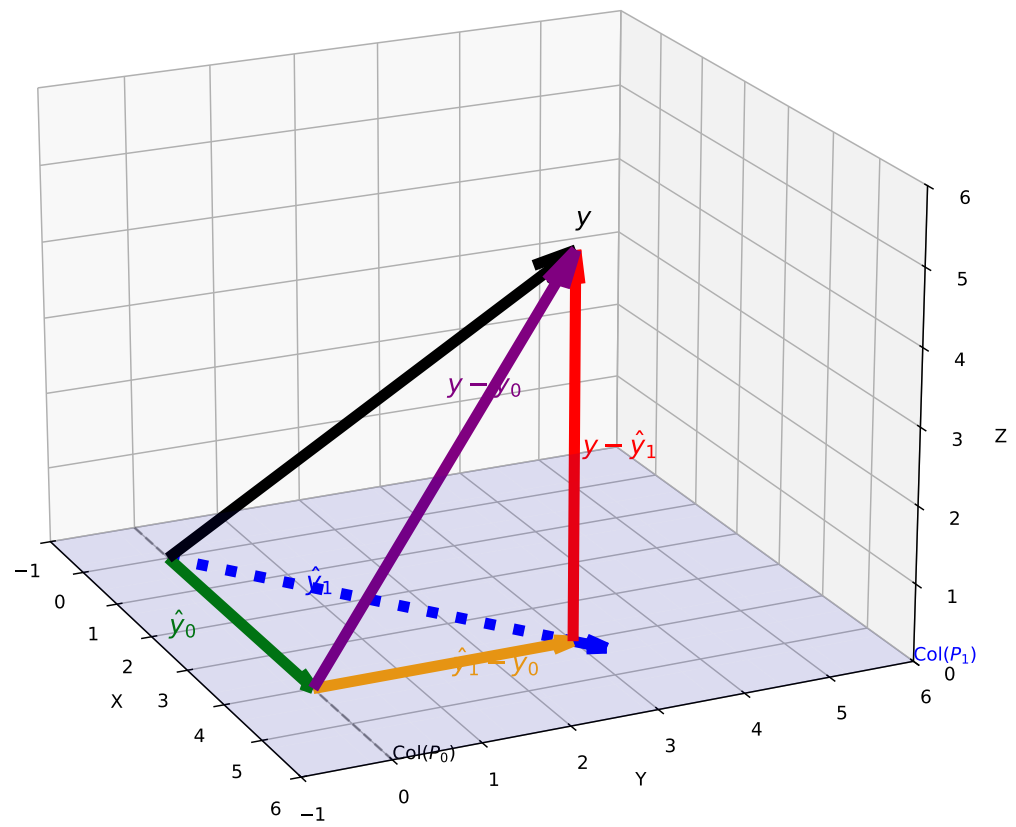


Figure 2.4: Illustration of Projections onto Nested Subspaces

### 3. Decomposition and Sum of Squares

Component	Notation	Definition	Vector Elements	Squared Norm (Sum of Squares)
<b>Null Proj.</b>	$\hat{y}_0$	$P_0 y$	Grand Mean ( $\bar{y}_{..}$ )	$\ \hat{y}_0\ ^2 = N\bar{y}_{..}^2$
<b>Full Proj.</b>	$\hat{y}_1$	$P_1 y$	Group Means ( $\bar{y}_{i.}$ )	$\ \hat{y}_1\ ^2 = \sum_{i=1}^k n_i \bar{y}_{i.}^2$

### 4. Geometric Justification of Shortcut Formulas

**A. Total Sum of Squares (SST)** Since  $\hat{y}_0 \perp (y - \hat{y}_0)$ , we have  $\|y\|^2 = \|\hat{y}_0\|^2 + \|y - \hat{y}_0\|^2$ :

$$\text{SST} = \|y - \hat{y}_0\|^2 = \|y\|^2 - \|\hat{y}_0\|^2 \quad (2.161)$$

$$\text{SST} = \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}^2 - N\bar{y}_{..}^2 \quad (2.162)$$

**B. Between Group Sum of Squares (SSB)** Since  $\hat{y}_0 \perp (\hat{y}_1 - \hat{y}_0)$ , we have  $\|\hat{y}_1\|^2 = \|\hat{y}_0\|^2 + \|\hat{y}_1 - \hat{y}_0\|^2$ :

$$\text{SSB} = \|\hat{y}_1 - \hat{y}_0\|^2 = \|\hat{y}_1\|^2 - \|\hat{y}_0\|^2 \quad (2.163)$$

$$\text{SSB} = \sum_{i=1}^k n_i \bar{y}_{i.}^2 - N\bar{y}_{..}^2 \quad (2.164)$$

**C. Within Group Sum of Squares (SSW)** Since  $\hat{y}_1 \perp (y - \hat{y}_1)$ , we have  $\|y\|^2 = \|\hat{y}_1\|^2 + \|y - \hat{y}_1\|^2$ :

$$\text{SSW} = \|y - \hat{y}_1\|^2 = \|y\|^2 - \|\hat{y}_1\|^2 \quad (2.165)$$

$$\text{SSW} = \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}^2 - \sum_{i=1}^k n_i \bar{y}_{i.}^2 \quad (2.166)$$

**Conclusion:**

$$\underbrace{\|y\|^2 - N\bar{y}_{..}^2}_{\text{SST}} = \underbrace{\left(\sum n_i \bar{y}_{i.}^2 - N\bar{y}_{..}^2\right)}_{\text{SSB}} + \underbrace{\left(\sum \sum y_{ij}^2 - \sum n_i \bar{y}_{i.}^2\right)}_{\text{SSW}} \quad (2.167)$$

### 5. Visualizing ANOVA Components in Data Space

## 2.9 Projections onto Orthogonal Subspaces

Finally, we consider the case where the entire space  $\mathbb{R}^n$  is decomposed into mutually orthogonal subspaces.



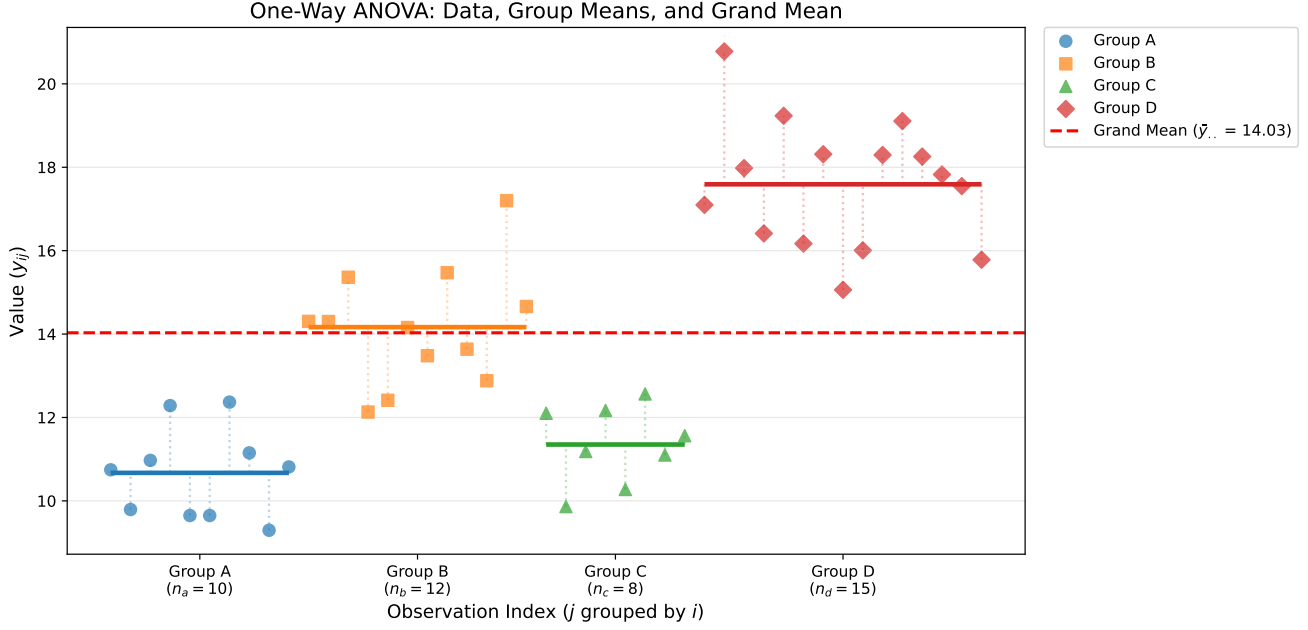


Figure 2.5: Visualization of Group Means vs. Grand Mean

**Theorem 2.23** (General Orthogonal Projections). *If  $\mathbb{R}^n$  is the direct sum of orthogonal subspaces  $V_1, V_2, \dots, V_k$ :*

$$\mathbb{R}^n = V_1 \oplus V_2 \oplus \dots \oplus V_k \quad (2.168)$$

where  $V_i \perp V_j$  for all  $i \neq j$ .

Then any vector  $y$  can be uniquely written as:

$$y = \hat{y}_1 + \hat{y}_2 + \dots + \hat{y}_k \quad (2.169)$$

where  $\hat{y}_i \in V_i$ .

Furthermore, each component  $\hat{y}_i$  is simply the projection of  $y$  onto the subspace  $V_i$ :

$$\hat{y}_i = P_i y \quad (2.170)$$

*Proof.* **1. Existence:** Since  $\mathbb{R}^n$  is the direct sum of  $V_1, \dots, V_k$ , by definition, any vector  $y \in \mathbb{R}^n$  can be written as a sum  $y = v_1 + \dots + v_k$  where  $v_i \in V_i$ .

**2. Uniqueness:** Suppose there are two such representations:  $y = \sum v_i = \sum w_i$ , with  $v_i, w_i \in V_i$ . Then  $\sum (v_i - w_i) = 0$ . Since subspaces in a direct sum are independent, the only way for the sum of elements to be zero is if each individual element is zero. Thus,  $v_i - w_i = 0 \implies v_i = w_i$ . The representation is unique. Let  $\hat{y}_i = v_i$ .

**3. Projection Property:** We claim that the  $i$ -th component  $\hat{y}_i$  is the orthogonal projection of  $y$  onto  $V_i$ . We must show that the residual  $(y - \hat{y}_i)$  is orthogonal to  $V_i$ .

$$y - \hat{y}_i = \sum_{j \neq i} \hat{y}_j \quad (2.171)$$

Let  $z$  be any vector in  $V_i$ . We calculate the inner product:

$$\langle y - \hat{y}_i, z \rangle = \left\langle \sum_{j \neq i} \hat{y}_j, z \right\rangle = \sum_{j \neq i} \langle \hat{y}_j, z \rangle \quad (2.172)$$

Since  $\hat{y}_j \in V_j$  and  $z \in V_i$ , and the subspaces are mutually orthogonal ( $V_j \perp V_i$  for  $j \neq i$ ), every term in the sum is zero. Therefore,  $(y - \hat{y}_i) \perp V_i$ . By the definition of orthogonal projection,  $\hat{y}_i = P_i y$ .  $\square$

This implies that the identity matrix can be decomposed into a sum of projection matrices:

$$I_n = P_1 + P_2 + \cdots + P_k \quad (2.173)$$

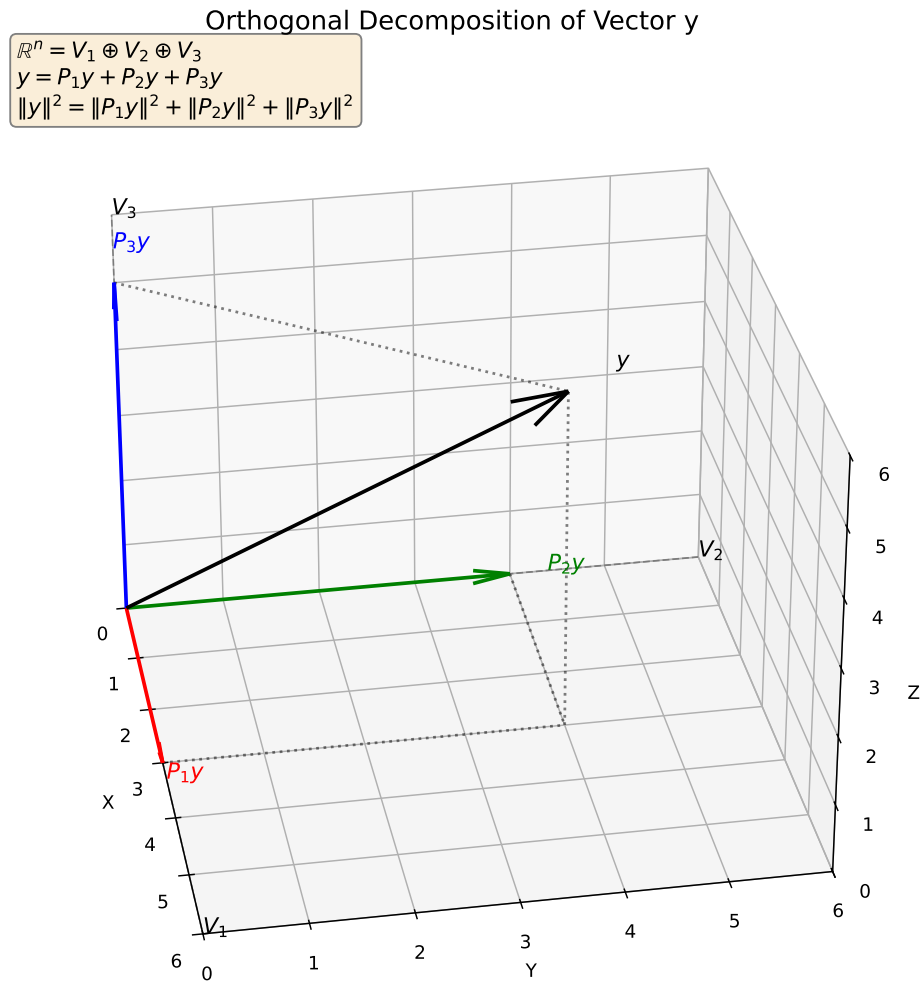


Figure 2.6: Orthogonal decomposition of vector  $y$  into subspaces

**Theorem 2.24** (Complete Orthogonal Decomposition of  $\mathbb{R}^n$ ). *Let  $P_0, P_1, \dots, P_k$  be a sequence of orthogonal projection matrices with nested column spaces:*

$$\text{Col}(P_0) \subseteq \text{Col}(P_1) \subseteq \dots \subseteq \text{Col}(P_k) \quad (2.174)$$

*Define the sequence of difference matrices  $\Delta P_i$  and their column spaces  $V_i$  as follows:*

$$\begin{aligned} \Delta P_0 &= P_0, & V_0 &= \text{Col}(\Delta P_0) \\ \Delta P_i &= P_i - P_{i-1} \quad (1 \leq i \leq k), & V_i &= \text{Col}(\Delta P_i) \\ \Delta P_{k+1} &= I - P_k, & V_{k+1} &= \text{Col}(\Delta P_{k+1}) \end{aligned}$$

**Conclusion:**

1. **Projection Property:** Each  $\Delta P_i$  is the orthogonal projection matrix onto  $V_i$  for  $i = 0, \dots, k+1$ .
2. **Mutual Orthogonality:** The collection  $\{\Delta P_i\}$  are mutually orthogonal operators:

$$\Delta P_i \Delta P_j = 0 \quad \text{for all } i \neq j \quad (2.175)$$

3. **Direct Sum Decomposition:** The vector space  $\mathbb{R}^n$  is the direct sum of these orthogonal subspaces:

$$\mathbb{R}^n = V_0 \oplus V_1 \oplus \dots \oplus V_{k+1} \quad (2.176)$$

*Proof.* **1. Proof that  $\Delta P_i$  is the Projection onto  $V_i$**  We must show each  $\Delta P_i$  is symmetric and idempotent.

- For  $\Delta P_0 = P_0$ : True by definition.
- For  $\Delta P_i$  ( $1 \leq i \leq k$ ):
  - **Symmetry:** Difference of symmetric matrices ( $P_i, P_{i-1}$ ) is symmetric.
  - **Idempotency:**  $(\Delta P_i)^2 = (P_i - P_{i-1})^2 = P_i^2 - P_i P_{i-1} - P_{i-1} P_i + P_{i-1}^2$ . Using nested properties ( $P_i P_{i-1} = P_{i-1}$ ), this simplifies to  $P_i - P_{i-1} = \Delta P_i$ .
- For  $\Delta P_{k+1} = I - P_k$ :
  - **Symmetry:**  $(I - P_k)' = I - P_k$ .
  - **Idempotency:**  $(I - P_k)^2 = I - 2P_k + P_k^2 = I - P_k$ .

**2. Proof of Mutual Orthogonality** We show  $\Delta P_j \Delta P_i = 0$  for  $i < j$ .

- **Case 1: Both indices  $\leq k$**  (i.e.,  $1 \leq i < j \leq k$ ):

$$(P_j - P_{j-1})(P_i - P_{i-1}) = P_j P_i - P_j P_{i-1} - P_{j-1} P_i + P_{j-1} P_{i-1} \quad (2.177)$$

Since  $\text{Col}(P_i) \subseteq \text{Col}(P_{j-1})$ , all terms reduce to  $P_i - P_{i-1} - P_i + P_{i-1} = 0$ .

- **Case 2: One index is the residual** ( $j = k+1$ ): We check  $\Delta P_{k+1} \Delta P_i = (I - P_k) \Delta P_i$  for any  $i \leq k$ . Since  $V_i \subseteq \text{Col}(P_k)$ , we have  $P_k \Delta P_i = \Delta P_i$ .

$$(I - P_k) \Delta P_i = \Delta P_i - P_k \Delta P_i = \Delta P_i - \Delta P_i = 0 \quad (2.178)$$

**3. Proof of Direct Sum** The sum of the difference matrices forms a telescoping series:

$$\sum_{j=0}^{k+1} \Delta P_j = P_0 + \sum_{i=1}^k (P_i - P_{i-1}) + (I - P_k) \quad (2.179)$$

$$= P_k + (I - P_k) = I \quad (2.180)$$

Since the identity operator  $I$  (which maps  $\mathbb{R}^n$  to itself) is the sum of mutually orthogonal projection operators, the space  $\mathbb{R}^n$  decomposes into the direct sum of their respective image subspaces  $V_i$ .  $\square$

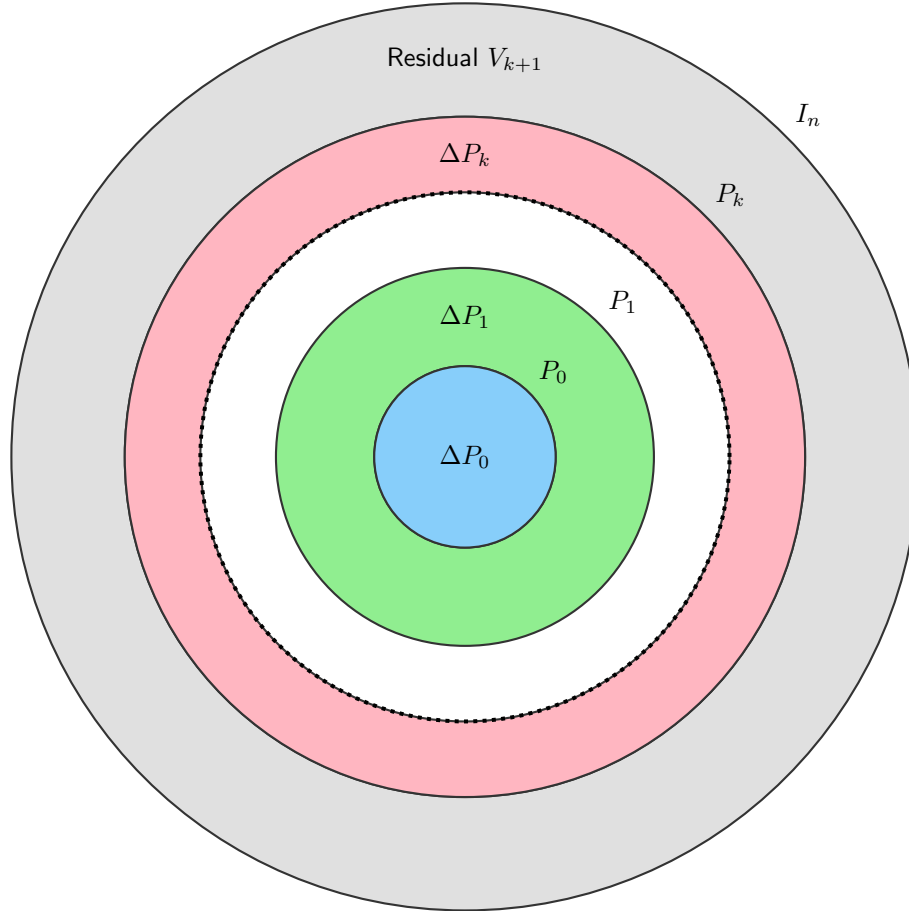


Figure 2.7: Venn Diagram of Nested Projections with Colored Increments

## 3 Matrix Algebra

This chapter covers a review of matrix algebra concepts essential for linear models, including eigenvalues, spectral decomposition, singular value decomposition.

### 3.1 Spectral Theory

#### 3.1.1 Eigenvalues and Eigenvectors

**Definition 3.1** (Eigenvalues and Eigenvectors). For a square matrix  $A$  ( $n \times n$ ), a scalar  $\lambda$  is an **eigenvalue** and a non-zero vector  $x$  is the corresponding **eigenvector** if:

$$Ax = \lambda x \iff (A - \lambda I_n)x = 0 \quad (3.1)$$

The eigenvalues are found by solving the characteristic equation:

$$|A - \lambda I_n| = 0 \quad (3.2)$$

#### 3.1.2 Spectral Decomposition

For symmetric matrices, we have a powerful decomposition theorem.

**Theorem 3.1** (Spectral Decomposition). *If  $A$  is a symmetric  $n \times n$  matrix, all its eigenvalues  $\lambda_1, \dots, \lambda_n$  are real. Furthermore, there exists an orthogonal matrix  $Q$  such that:*

$$A = Q\Lambda Q' = \sum_{i=1}^n \lambda_i q_i q_i' \quad (3.3)$$

where:

- $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  contains the eigenvalues.
- $Q = (q_1, \dots, q_n)$  contains the corresponding orthonormal eigenvectors ( $q_i' q_j = \delta_{ij}$ ).

**Explantion:** This allows us to view the transformation  $Ax$  as a rotation ( $Q'$ ), a scaling ( $\Lambda$ ), and a rotation back ( $Q$ ). For a symmetric matrix  $A$ , we can write the spectral decomposition as a product of the eigenvector matrix  $Q$  and eigenvalue matrix  $\Lambda$ :

$$\begin{aligned}
A &= Q\Lambda Q' \\
&= (q_1 \quad q_2 \quad \cdots \quad q_n) \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{pmatrix} \begin{pmatrix} q'_1 \\ q'_2 \\ \vdots \\ q'_n \end{pmatrix} \\
&= (\lambda_1 q_1 \quad \lambda_2 q_2 \quad \cdots \quad \lambda_n q_n) \begin{pmatrix} q'_1 \\ q'_2 \\ \vdots \\ q'_n \end{pmatrix} \\
&= \lambda_1 q_1 q'_1 + \lambda_2 q_2 q'_2 + \cdots + \lambda_n q_n q'_n \\
&= \sum_{i=1}^n \lambda_i q_i q'_i
\end{aligned} \tag{3.4}$$

where the eigenvectors  $q_i$  satisfy the orthogonality conditions:

$$q'_i q_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \tag{3.5}$$

And  $Q$  is an orthogonal matrix:  $Q'Q = QQ' = I_n$ .

### 3.1.3 Quadratic Form

**Definition 3.2.** A **quadratic form** in  $n$  variables  $x_1, x_2, \dots, x_n$  is a scalar function defined by a symmetric matrix  $A$ :

$$Q(x) = x'Ax = \sum_{i=1}^n \sum_{j=1}^n a_{ij}x_i x_j \tag{3.6}$$

### 3.1.4 Positive and Non-Negative Definite Matrices

**Definition 3.3** (Positive and Non-Negative Definite Matrices). A symmetric matrix  $A$  is **positive definite (p.d.)** if:

$$x'Ax > 0 \quad \forall x \neq 0 \tag{3.7}$$

It is **non-negative definite (n.n.d.)** if:

$$x'Ax \geq 0 \quad \forall x \tag{3.8}$$

**Theorem 3.2** (Properties of Definite Matrices). *Let  $A$  be a symmetric  $n \times n$  matrix with eigenvalues  $\lambda_1, \dots, \lambda_n$ .*

1. **Eigenvalue Characterization:**

- $A$  is p.d.  $\iff$  all  $\lambda_i > 0$ .
- $A$  is n.n.d.  $\iff$  all  $\lambda_i \geq 0$ .

## 2. Determinant and Inverse:

- If  $A$  is p.d., then  $|A| > 0$  and  $A^{-1}$  exists.
- If  $A$  is n.n.d. and singular, then  $|A| = 0$  (at least one  $\lambda_i = 0$ ).

## 3. Gram Matrices ( $B'B$ ): Let $B$ be an $n \times p$ matrix.

- If  $\text{rank}(B) = p$ , then  $B'B$  is p.d.
- If  $\text{rank}(B) < p$ , then  $B'B$  is n.n.d.

### 3.1.5 Properties of Symmetric Matrices

**Theorem 3.3** (Properties of Symmetric Matrices). Let  $A$  be a symmetric matrix with spectral decomposition  $A = Q\Lambda Q'$ . The following properties hold:

1. **Trace:**  $\text{tr}(A) = \sum \lambda_i$ .
2. **Determinant:**  $|A| = \prod \lambda_i$ .
3. **Singularity:**  $A$  is singular if and only if at least one  $\lambda_i = 0$ .
4. **Inverse:** If  $A$  is non-singular ( $\lambda_i \neq 0$ ), then  $A^{-1} = Q\Lambda^{-1}Q'$ .
5. **Powers:**  $A^k = Q\Lambda^kQ'$ .
  - Square Root:  $A^{1/2} = Q\Lambda^{1/2}Q'$  (if  $\lambda_i \geq 0$ ).
6. **Spectral Representation of Quadratic Forms:** The quadratic form  $x'Ax$  can be diagonalized using the eigenvectors of  $A$ :

$$x'Ax = x'Q\Lambda Q'x = y'\Lambda y = \sum_{i=1}^n \lambda_i y_i^2 \quad (3.9)$$

where  $y = Q'x$  represents a rotation of the coordinate system.

### 3.1.6 Spectral Representation of Projection Matrices

We revisit projection matrices in the context of eigenvalues.

**Theorem 3.4** (Eigenvalues of Projection Matrices). A symmetric matrix  $P$  is a projection matrix (idempotent,  $P^2 = P$ ) if and only if its eigenvalues are either 0 or 1.

$$P^2x = \lambda^2x \quad \text{and} \quad Px = \lambda x \implies \lambda^2 = \lambda \implies \lambda \in \{0, 1\} \quad (3.10)$$

For a projection matrix  $P$ :

- If  $x \in \text{Col}(P)$ ,  $Px = x$  (Eigenvalue 1).
- If  $x \perp \text{Col}(P)$ ,  $Px = 0$  (Eigenvalue 0).
- $\text{rank}(P) = \text{tr}(P) = \sum \lambda_i$  (Count of 1s).

**Example 3.1.** For  $P = \frac{1}{n} J_n J_n'$ , the rank is  $\text{tr}(P) = 1$ .

## 3.2 Singular Value Decomposition (SVD)

**Theorem 3.5** (Singular Value Decomposition (SVD)). *Let  $X$  be an  $n \times p$  matrix with rank  $r \leq \min(n, p)$ .  $X$  can be decomposed into the product of three matrices:*

$$X = U \mathbf{D} V' \quad (3.11)$$

### 1. Partitioned Matrix Form

$$X = \begin{pmatrix} U_1 & U_2 \end{pmatrix}_{n \times n} \begin{pmatrix} \Lambda_r & O_{r \times (p-r)} \\ O_{(n-r) \times r} & O_{(n-r) \times (p-r)} \end{pmatrix} \begin{pmatrix} V_1' \\ V_2' \end{pmatrix}_{p \times p} \quad (3.12)$$

### 2. Detailed Matrix Form

Expanding the diagonal matrix explicitly:

$$X = \begin{pmatrix} u_1 & \dots & u_n \end{pmatrix}_{n \times n} \left( \begin{array}{cccc|cc} \lambda_1 & 0 & \dots & 0 & & \\ 0 & \lambda_2 & \dots & 0 & & \\ \vdots & \vdots & \ddots & \vdots & & \\ 0 & 0 & \dots & \lambda_r & & \\ \hline & & & & O_{21} & \\ & & & & & O_{22} \end{array} \right) \begin{pmatrix} v_1' \\ \vdots \\ v_p' \end{pmatrix}_{p \times p} \quad (3.13)$$

### 3. Reduced Form

$$X = U_1 \Lambda_r V_1' = \sum_{i=1}^r \lambda_i u_i v_i' \quad (3.14)$$

**Properties:**

1. **Singular Values ( $\Lambda_r$ ):**  $\Lambda_r = \text{diag}(\lambda_1, \dots, \lambda_r)$  contains the singular values ( $\lambda_i > 0$ ), which are the square roots of the non-zero eigenvalues of  $X'X$ .
2. **Orthogonality:**
  - $U$  is  $n \times n$  orthogonal ( $U'U = I_n$ ).
  - $V$  is  $p \times p$  orthogonal ( $V'V = I_p$ ).



### 3.2.0.1 Connection to Gram Matrices

The matrices  $U$  and  $V$  provide the basis vectors (eigenvectors) for the Gram matrices of  $X$ .

1. **Right Singular Vectors ( $V$ ):** The columns of  $V$  are the eigenvectors of the Gram matrix  $X'X$ .

$$X'X = (U\Lambda V')'(U\Lambda V') = V\Lambda U'U\Lambda V' = V\Lambda^2 V' \quad (3.15)$$

- The eigenvalues of  $X'X$  are the squared singular values  $\lambda_i^2$ .

2. **Left Singular Vectors ( $U$ ):** The columns of  $U$  are the eigenvectors of the Gram matrix  $XX'$ .

$$XX' = (U\Lambda V')(U\Lambda V')' = U\Lambda V'V\Lambda U' = U\Lambda^2 U' \quad (3.16)$$

- The eigenvalues of  $XX'$  are also  $\lambda_i^2$  (for non-zero values).

### 3.2.0.2 Numerical Example

Consider the matrix  $X = \begin{pmatrix} 1 & 1 \\ 2 & 2 \end{pmatrix}$ .

1. **Compute  $X'X$  and find  $V$ :**

$$X'X = \begin{pmatrix} 1 & 2 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 2 & 2 \end{pmatrix} = \begin{pmatrix} 5 & 5 \\ 5 & 5 \end{pmatrix} \quad (3.17)$$

- Eigenvalues of  $X'X$ : Trace is 10, Determinant is 0. Thus,  $\mu_1 = 10, \mu_2 = 0$ .
- **Singular Values:**  $\lambda_1 = \sqrt{10}, \lambda_2 = 0$ .
- Eigenvector for  $\mu_1 = 10$ : Normalized  $v_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ .
- Eigenvector for  $\mu_2 = 0$ : Normalized  $v_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ .
- Therefore,  $V = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$ .

2. **Compute  $XX'$  and find  $U$ :**

$$XX' = \begin{pmatrix} 1 & 1 \\ 2 & 2 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 1 & 2 \end{pmatrix} = \begin{pmatrix} 2 & 4 \\ 4 & 8 \end{pmatrix} \quad (3.18)$$

- Eigenvalues are again 10 and 0.
- Eigenvector for  $\mu_1 = 10$ : Normalized  $u_1 = \frac{1}{\sqrt{5}} \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ .
- Eigenvector for  $\mu_2 = 0$ : Normalized  $u_2 = \frac{1}{\sqrt{5}} \begin{pmatrix} 2 \\ -1 \end{pmatrix}$ .
- Therefore,  $U = \frac{1}{\sqrt{5}} \begin{pmatrix} 1 & 2 \\ 2 & -1 \end{pmatrix}$ .

3. **Verification:**

$$X = \sqrt{10}u_1v_1' = \sqrt{10} \begin{pmatrix} \frac{1}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} \end{pmatrix} \left( \frac{1}{\sqrt{2}} \quad \frac{1}{\sqrt{2}} \right) = \begin{pmatrix} 1 & 1 \\ 2 & 2 \end{pmatrix} \quad (3.19)$$

### 3.3 Cholesky Decomposition

A symmetric matrix  $A$  has a Cholesky decomposition if and only if it is **non-negative definite** (i.e.,  $x'Ax \geq 0$  for all  $x$ ).

$$A = B'B \quad (3.20)$$

where  $B$  is an **upper triangular** matrix with non-negative diagonal entries.

#### 3.3.1 Matrix Representation of the Algorithm

To derive the algorithm, we equate the elements of  $A$  with the product of the lower triangular matrix  $B'$  and the upper triangular matrix  $B$ .

For a  $3 \times 3$  matrix, this looks like:

$$\underbrace{\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}}_A = \underbrace{\begin{pmatrix} b_{11} & 0 & 0 \\ b_{12} & b_{22} & 0 \\ b_{13} & b_{23} & b_{33} \end{pmatrix}}_{B'} \underbrace{\begin{pmatrix} b_{11} & b_{12} & b_{13} \\ 0 & b_{22} & b_{23} \\ 0 & 0 & b_{33} \end{pmatrix}}_B \quad (3.21)$$

Multiplying the matrices on the right yields the system of equations:

$$A = \begin{pmatrix} \mathbf{b_{11}^2} & b_{11}b_{12} & b_{11}b_{13} \\ b_{12}b_{11} & \mathbf{b_{12}^2 + b_{22}^2} & b_{12}b_{13} + b_{22}b_{23} \\ b_{13}b_{11} & b_{13}b_{12} + b_{23}b_{22} & \mathbf{b_{13}^2 + b_{23}^2 + b_{33}^2} \end{pmatrix} \quad (3.22)$$

By solving for the bolded diagonal terms and substituting known values from previous rows, we get the recursive algorithm.

#### 3.3.2 The Algorithm

1. **Row 1:** Solve for  $b_{11}$  using  $a_{11}$ , then solve the rest of the row ( $b_{1j}$ ) by division.

- $b_{11} = \sqrt{a_{11}}$
- $b_{1j} = a_{1j}/b_{11}$

2. **Row 2:** Solve for  $b_{22}$  using  $a_{22}$  and the known  $b_{12}$ , then solve  $b_{2j}$ .

- $b_{22} = \sqrt{a_{22} - b_{12}^2}$
- $b_{2j} = (a_{2j} - b_{12}b_{1j})/b_{22}$

3. **Row 3:** Solve for  $b_{33}$  using  $a_{33}$  and the known  $b_{13}$ ,  $b_{23}$ .

- $b_{33} = \sqrt{a_{33} - b_{13}^2 - b_{23}^2}$

### 3.3.3 Numerical Example

Consider the positive definite matrix  $A$ :

$$A = \begin{pmatrix} 4 & 2 & -2 \\ 2 & 10 & 2 \\ -2 & 2 & 6 \end{pmatrix} \quad (3.23)$$

We find  $B$  such that  $A = B'B$ :

**1. First Row of  $B$  ( $b_{11}, b_{12}, b_{13}$ ):**

- $b_{11} = \sqrt{4} = 2$
- $b_{12} = 2/2 = 1$
- $b_{13} = -2/2 = -1$

**2. Second Row of  $B$  ( $b_{22}, b_{23}$ ):**

- $b_{22} = \sqrt{10 - (1)^2} = \sqrt{9} = 3$
- $b_{23} = (2 - (1)(-1))/3 = 3/3 = 1$

**3. Third Row of  $B$  ( $b_{33}$ ):**

- $b_{33} = \sqrt{6 - (-1)^2 - (1)^2} = \sqrt{4} = 2$

**Result:**

$$B = \begin{pmatrix} 2 & 1 & -1 \\ 0 & 3 & 1 \\ 0 & 0 & 2 \end{pmatrix} \quad (3.24)$$

# 4 Multivariate Normal Distribution

## 4.1 Motivation

Consider the linear model:

$$y = X\beta + \epsilon, \quad \epsilon_i \sim N(0, \sigma^2) \quad (4.1)$$

We are often interested in the distributional properties of the response vector  $y$  and the residuals. Specifically, if  $y = (y_1, \dots, y_n)'$ , we need to understand its multivariate distribution.

$$\hat{y} = Py, \quad e = y - \hat{y} = (I_n - P)y \quad (4.2)$$

## 4.2 Random Vectors and Matrices

**Definition 4.1** (Random Vector and Matrix). A **Random Vector** is a vector whose elements are random variables. E.g.,

$$x_{k \times 1} = (x_1, x_2, \dots, x_k)^T \quad (4.3)$$

where  $x_1, \dots, x_k$  are each random variables.

A **Random Matrix** is a matrix whose elements are random variables. E.g.,  $X_{n \times k} = (x_{ij})$ , where  $x_{11}, \dots, x_{nk}$  are each random variables.

**Definition 4.2** (Expected Value). The expected value (population mean) of a random matrix (or vector) is the matrix (or vector) of expected values of its elements.

For  $X_{n \times k}$ :

$$E(X) = \begin{pmatrix} E(x_{11}) & \dots & E(x_{1k}) \\ \vdots & \ddots & \vdots \\ E(x_{n1}) & \dots & E(x_{nk}) \end{pmatrix} \quad (4.4)$$

$$E \left( \begin{pmatrix} x_1 \\ \vdots \\ x_k \end{pmatrix} \right) = \begin{pmatrix} E(x_1) \\ \vdots \\ E(x_k) \end{pmatrix} \quad (4.5)$$

**Definition 4.3** (Variance-Covariance Matrix). For a random vector  $x_{k \times 1} = (x_1, \dots, x_k)^T$ , the matrix is:

$$\text{Var}(x) = \Sigma_x = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1k} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{k1} & \sigma_{k2} & \dots & \sigma_{kk} \end{pmatrix} \quad (4.6)$$

Where:

- $\sigma_{ij} = \text{Cov}(x_i, x_j) = E[(x_i - \mu_i)(x_j - \mu_j)]$
- $\sigma_{ii} = \text{Var}(x_i) = E[(x_i - \mu_i)^2]$

In matrix notation:

$$\text{Var}(x) = E[(x - \mu_x)(x - \mu_x)^T] \quad (4.7)$$

Note:  $\text{Var}(x)$  is symmetric.

#### 4.2.1 Derivation of Covariance Matrix Structure

Expanding the vector multiplication for variance:

$$(x - \mu_x)(x - \mu_x)' \quad \text{where } \mu_x = (\mu_1, \dots, \mu_n)' \quad (4.8)$$

$$= \begin{pmatrix} x_1 - \mu_1 \\ \vdots \\ x_n - \mu_n \end{pmatrix} (x_1 - \mu_1, \dots, x_n - \mu_n) \quad (4.9)$$

This results in the matrix  $A = (a_{ij})$  where  $a_{ij} = (x_i - \mu_i)(x_j - \mu_j)$ . Taking expectations yields the covariance matrix elements  $\sigma_{ij}$ .

**Definition 4.4** (Covariance Matrix (Two Vectors)). For random vectors  $x_{k \times 1}$  and  $y_{n \times 1}$ , the covariance matrix is:

$$\text{Cov}(x, y) = E[(x - \mu_x)(y - \mu_y)^T] = \begin{pmatrix} \text{Cov}(x_1, y_1) & \dots & \text{Cov}(x_1, y_n) \\ \vdots & \ddots & \vdots \\ \text{Cov}(x_k, y_1) & \dots & \text{Cov}(x_k, y_n) \end{pmatrix} \quad (4.10)$$

Note that  $\text{Cov}(x, x) = \text{Var}(x)$ .

**Definition 4.5** (Correlation Matrix). The correlation matrix of a random vector  $x$  is:

$$\text{corr}(x) = \begin{pmatrix} 1 & \rho_{12} & \dots & \rho_{1k} \\ \vdots & \ddots & \vdots & \\ \rho_{k1} & \rho_{k2} & \dots & 1 \end{pmatrix} \quad (4.11)$$

where  $\rho_{ij} = \text{corr}(x_i, x_j)$ .

**Relationships:** Let  $V_x = \text{diag}(\text{Var}(x_1), \dots, \text{Var}(x_k))$ .

$$\Sigma_x = V_x^{1/2} \rho_x V_x^{1/2} \quad \text{and} \quad \rho_x = (V_x^{1/2})^{-1} \Sigma_x (V_x^{1/2})^{-1} \quad (4.12)$$

Similarly for two vectors:

$$\Sigma_{xy} = V_x^{1/2} \rho_{xy} V_y^{1/2} \quad (4.13)$$

### 4.3 Properties of Mean and Variance

We can derive several key algebraic properties for operations on random vectors.

1.  $E(X + Y) = E(X) + E(Y)$
2.  $E(AXB) = AE(X)B$  (In particular,  $E(AX) = A\mu_x$ )
3.  $\text{Cov}(x, y) = \text{Cov}(y, x)^T$
4.  $\text{Cov}(x + c, y + d) = \text{Cov}(x, y)$
5.  $\text{Cov}(Ax, By) = A\text{Cov}(x, y)B^T$ 
  - Special case for scalars:  $\text{Cov}(ax, by) = ab \cdot \text{Cov}(x, y)$
6.  $\text{Cov}(x_1 + x_2, y_1) = \text{Cov}(x_1, y_1) + \text{Cov}(x_2, y_1)$
7.  $\text{Var}(x + c) = \text{Var}(x)$
8.  $\text{Var}(Ax) = A\text{Var}(x)A^T$
9.  $\text{Var}(x_1 + x_2) = \text{Var}(x_1) + \text{Cov}(x_1, x_2) + \text{Cov}(x_2, x_1) + \text{Var}(x_2)$
10.  $\text{Var}(\sum x_i) = \sum \text{Var}(x_i)$  if independent.

*Proof.* **Property 5 (Covariance of Linear Transformation):**

$$\begin{aligned}\text{Cov}(Ax, By) &= E[(Ax - A\mu_x)(By - B\mu_y)^T] \\ &= AE[(x - \mu_x)(y - \mu_y)^T]B^T \\ &= A\text{Cov}(x, y)B^T\end{aligned}\tag{4.14}$$

**Property 2 (Expectation of Linear Transformation):**

To prove  $E(AXB) = AE(X)B$ : First consider  $E(Ax_j)$  where  $x_j$  is a column of  $X$ .

$$E(Ax_j) = E\begin{pmatrix} a'_1 x_j \\ \vdots \\ a'_n x_j \end{pmatrix} = \begin{pmatrix} E(a'_1 x_j) \\ \vdots \\ E(a'_n x_j) \end{pmatrix}\tag{4.15}$$

Since  $a_i$  are constants:

$$E(a'_i x_j) = E\left(\sum_{k=1}^p a_{ik} x_{kj}\right) = \sum_{k=1}^p a_{ik} E(x_{kj}) = a'_i E(x_j)\tag{4.16}$$

Thus  $E(Ax_j) = AE(x_j)$ . Applying this to all columns of  $X$ :

$$E(AX) = [E(Ax_1), \dots, E(Ax_m)] = [AE(x_1), \dots, AE(x_m)] = AE(X)\tag{4.17}$$

Similarly,  $E(XB) = E(X)B$ .

**Proof of Property 9 (Variance of Sum):**

$$\text{Var}(x_1 + x_2) = E[(x_1 + x_2 - \mu_1 - \mu_2)(x_1 + x_2 - \mu_1 - \mu_2)^T]\tag{4.18}$$

Let centered variables be denoted by differences.

$$= E[((x_1 - \mu_1) + (x_2 - \mu_2))((x_1 - \mu_1) + (x_2 - \mu_2))^T]\tag{4.19}$$

Expanding terms:

$$= E[(x_1 - \mu_1)(x_1 - \mu_1)^T + (x_1 - \mu_1)(x_2 - \mu_2)^T + (x_2 - \mu_2)(x_1 - \mu_1)^T + (x_2 - \mu_2)(x_2 - \mu_2)^T] \quad (4.20)$$

$$= \text{Var}(x_1) + \text{Cov}(x_1, x_2) + \text{Cov}(x_2, x_1) + \text{Var}(x_2) \quad (4.21)$$

□

## 4.4 The Multivariate Normal Distribution

### 4.4.1 Definition and Density

**Definition 4.6** (Independent Standard Normal). Let  $z = (z_1, \dots, z_n)'$  where  $z_i \sim N(0, 1)$  are independent. We say  $z \sim N_n(0, I_n)$ . The joint PDF is the product of marginals:

$$f(z) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{z_i^2}{2}} = \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} z^T z} \quad (4.22)$$

Properties:  $E(z) = 0$  and  $\text{Var}(z) = I_n$  (Covariance is 0 for  $i \neq j$ , Variance is 1).

**Definition 4.7** (Multivariate Normal Distribution). A random vector  $x$  ( $n \times 1$ ) has a **multivariate normal distribution** if it has the same distribution as:

$$x = A_{n \times p} z_{p \times 1} + \mu_{n \times 1} \quad (4.23)$$

where  $z \sim N_p(0, I_p)$ ,  $A$  is a matrix of constants, and  $\mu$  is a vector of constants. The moments are:

- $E(x) = \mu$
- $\text{Var}(x) = AA^T = \Sigma$

### 4.4.2 Geometric Interpretation

Using Spectral Decomposition,  $\Sigma = Q\Lambda Q'$ . We can view the transformation  $x = Az + \mu$  as:

1. Scaling by eigenvalues ( $\Lambda^{1/2}$ ).
2. Rotation by eigenvectors ( $Q$ ).
3. Shift by mean ( $\mu$ ).

### 4.4.3 Probability Density Function

If  $\Sigma$  is positive definite, the PDF exists. We use the change of variable formula for  $x = Az + \mu$ :

$$f_x(x) = f_z(g^{-1}(x)) \cdot |J| \quad (4.24)$$

where  $z = A^{-1}(x - \mu)$  and  $J = \det(A^{-1}) = |A|^{-1}$ .

$$f_x(x) = (2\pi)^{-p/2} |A|^{-1} \exp \left\{ -\frac{1}{2} (A^{-1}(x - \mu))^T (A^{-1}(x - \mu)) \right\} \quad (4.25)$$

Using  $|\Sigma| = |AA^T| = |A|^2$  and  $\Sigma^{-1} = (AA^T)^{-1}$ , we get:

$$f_x(x) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\} \quad (4.26)$$

### 4.4.4 Moment Generating Function

**Definition 4.8** (Moment Generating Function (MGF)). The MGF of a random vector  $x$  is  $M_x(t) = E(e^{t^T x})$ . For  $x = Az + \mu$ :

$$M_x(t) = E[e^{t^T(Az+\mu)}] = e^{t^T \mu} E[e^{(A^T t)^T z}] = e^{t^T \mu} M_z(A^T t) \quad (4.27)$$

Since  $M_z(u) = e^{u^T \mu/2}$ :

$$M_x(t) = e^{t^T \mu} \exp \left( \frac{1}{2} t^T (AA^T) t \right) = \exp \left( t^T \mu + \frac{1}{2} t^T \Sigma t \right) \quad (4.28)$$

Key Properties:

1. **Uniqueness:** Two random vectors with the same MGF have the same distribution.
2. **Independence:**  $y_1$  and  $y_2$  are independent iff  $M_y(t) = M_{y_1}(t_1)M_{y_2}(t_2)$ .

## 4.5 Construction and Linear Transformations

**Theorem 4.1** (Constructing MVN Random Vector). Let  $\mu \in \mathbb{R}^n$  and  $\Sigma$  be an  $n \times n$  symmetric non-negative definitive (n.n.d) matrix. Then there exists a multivariate normal distribution with mean  $\mu$  and covariance  $\Sigma$ .

*Proof.* Since  $\Sigma$  is n.n.d., there exists  $B$  such that  $\Sigma = BB^T$  (e.g., via Cholesky or Spetral Decomposition). Let  $z \sim N_n(0, I)$  and define  $x = Bz + \mu$ .  $\square$

**Theorem 4.2** (Linear Transformation Theorem). Let  $x \sim N_n(\mu, \Sigma)$ . Let  $y = Cx + d$  where  $C$  is  $r \times n$  and  $d$  is  $r \times 1$ . Then:

$$y \sim N_r(C\mu + d, C\Sigma C^T) \quad (4.29)$$

*Proof.*  $x = Az + \mu$  where  $AA^T = \Sigma$ .

$$y = C(Az + \mu) + d = (CA)z + (C\mu + d) \quad (4.30)$$

This fits the definition of MVN with mean  $C\mu + d$  and variance  $C\Sigma C^T$ .  $\square$



### 4.5.1 Important Corollaries of Theorem 4.2

**Corollary 4.1** (Marginals). *Any subvector of a multivariate normal vector is also multivariate normal.*

*Proof.* If we partition  $x = (x'_1, x'_2)'$ , we can use  $C = (I_r, 0)$  to show  $x_1 \sim N(\mu_1, \Sigma_{11})$ . □

**Corollary 4.2** (Univariate Combinations). *Any linear combination  $a^T x$  is univariate normal:*

$$a^T x \sim N(a^T \mu, a^T \Sigma a) \quad (4.31)$$

**Corollary 4.3** (Orthogonal Transformations). *If  $x \sim N(0, I_n)$  and  $Q$  is orthogonal ( $Q'Q = I$ ), then  $y = Q'x \sim N(0, I_n)$ .*

**Corollary 4.4** (Standardization). *If  $y \sim N_n(\mu, \Sigma)$  and  $\Sigma$  is positive definite:*

$$\Sigma^{-1/2}(y - \mu) \sim N_n(0, I_n) \quad (4.32)$$

*Proof.* Let  $z = \Sigma^{-1/2}(y - \mu)$ . Then  $\text{Var}(z) = \Sigma^{-1/2}\Sigma\Sigma^{-1/2} = I_n$ . □

## 4.6 Independence

**Theorem 4.3** (Independence in MVN). *Let  $y \sim N(\mu, \Sigma)$  be partitioned into  $y_1$  and  $y_2$ .*

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \quad (4.33)$$

*Then  $y_1$  and  $y_2$  are independent if and only if  $\Sigma_{12} = 0$  (zero covariance).*

*Proof.* **1. Independence  $\implies$  Covariance is 0:** This holds generally for any distribution.

$$\text{Cov}(y_1, y_2) = E[(y_1 - \mu_1)(y_2 - \mu_2)'] = 0 \quad (4.34)$$

**2. Covariance is 0  $\implies$  Independence:** This is specific to MVN. We use MGFs. If  $\Sigma_{12} = 0$ , the quadratic form in the MGF splits:

$$t^T \Sigma t = t_1^T \Sigma_{11} t_1 + t_2^T \Sigma_{22} t_2 \quad (4.35)$$

The MGF becomes:

$$M_y(t) = \exp(t_1^T \mu_1 + \frac{1}{2} t_1^T \Sigma_{11} t_1) \times \exp(t_2^T \mu_2 + \frac{1}{2} t_2^T \Sigma_{22} t_2) \quad (4.36)$$

$$M_y(t) = M_{y_1}(t_1) M_{y_2}(t_2) \quad (4.37)$$

Thus, they are independent. □

## 4.7 Signal-Noise Decomposition for Multivariate Normal Distribution

We can formalize the relationship between two random vectors  $y$  and  $x$  through a decomposition theorem that separates the systematic signal from the stochastic noise.

**Theorem 4.4** (Regression Decomposition Theorem). *Let the random vector  $V$  of dimension  $p \times 1$  be partitioned into two subvectors  $y$  ( $p_1 \times 1$ ) and  $x$  ( $p_2 \times 1$ ). Assume  $V$  follows a multivariate normal distribution:*

$$\begin{pmatrix} y \\ x \end{pmatrix} \sim N_p \left( \begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix}, \begin{pmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{pmatrix} \right) \quad (4.38)$$

The response vector  $y$  can be uniquely decomposed into a systematic component and a stochastic error:

$$y = m(x) + e \quad (4.39)$$

where we define the **Regression Coefficient Matrix**  $B$  and the components as:

$$B = \Sigma_{yx} \Sigma_{xx}^{-1} \quad (4.40)$$

$$m(x) = \mu_y + B(x - \mu_x) \quad (4.41)$$

$$e = y - m(x) \quad (4.42)$$

**Properties:**

1. **Independence:** The noise vector  $e$  is statistically independent of the predictor  $x$  (and consequently independent of  $m(x)$ ).
2. **Marginal Distributions:**
  - $m(x) \sim N_{p_1}(\mu_y, B\Sigma_{xx}B^T)$
  - $e \sim N_{p_1}(0, \Sigma_{yy} - B\Sigma_{xx}B^T)$
3. **Conditional Distribution:** Since  $y = m(x) + e$ , and  $e$  is independent of  $x$ , the conditional distribution is:

$$y|x \sim N_{p_1}(m(x), \Sigma_{y|x}) \quad (4.43)$$

where:

$$m(x) = \mu_y + B(x - \mu_x) = \mu_y + \Sigma_{yx} \Sigma_{xx}^{-1}(x - \mu_x) \quad (4.44)$$

$$\Sigma_{y|x} = \Sigma_{yy} - B\Sigma_{xx}B^T = \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \quad (4.45)$$

*Proof.* We define a transformation from the input vector  $V = \begin{pmatrix} y \\ x \end{pmatrix}$  to the target vector  $W = \begin{pmatrix} m(x) \\ e \end{pmatrix}$ .

Using the linear transformation  $W = CV + d$ :

$$\underbrace{\begin{pmatrix} m(x) \\ e \end{pmatrix}}_W = \underbrace{\begin{pmatrix} 0 & B \\ I & -B \end{pmatrix}}_C \underbrace{\begin{pmatrix} y \\ x \end{pmatrix}}_V + \underbrace{\begin{pmatrix} \mu_y - B\mu_x \\ -(\mu_y - B\mu_x) \end{pmatrix}}_d \quad (4.46)$$

### 1. Mean Vector

$$E[W] = CE[V] + d = \begin{pmatrix} 0 & B \\ I & -B \end{pmatrix} \begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix} + \begin{pmatrix} \mu_y - B\mu_x \\ -\mu_y + B\mu_x \end{pmatrix} = \begin{pmatrix} B\mu_x \\ \mu_y - B\mu_x \end{pmatrix} + \begin{pmatrix} \mu_y - B\mu_x \\ -\mu_y + B\mu_x \end{pmatrix} = \begin{pmatrix} \mu_y \\ 0 \end{pmatrix} \quad (4.47)$$

### 2. Covariance Matrix

We compute  $\text{Var}(W) = C\Sigma C^T$  directly:

$$\begin{aligned} C\Sigma C^T &= \begin{pmatrix} 0 & B \\ I & -B \end{pmatrix} \begin{pmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{pmatrix} \begin{pmatrix} 0 & I \\ B^T & -B^T \end{pmatrix} \\ &= \begin{pmatrix} B\Sigma_{xy} & B\Sigma_{xx} \\ \Sigma_{yy} - B\Sigma_{xy} & \Sigma_{yx} - B\Sigma_{xx} \end{pmatrix} \begin{pmatrix} 0 & I \\ B^T & -B^T \end{pmatrix} \\ &= \begin{pmatrix} B\Sigma_{xx}B^T & B\Sigma_{xy} - B\Sigma_{xx}B^T \\ \Sigma_{yx}B^T - B\Sigma_{xx}B^T & (\Sigma_{yy} - B\Sigma_{xy}) - (\Sigma_{yx} - B\Sigma_{xx})B^T \end{pmatrix} \\ &= \begin{pmatrix} B\Sigma_{xx}B^T & 0 \\ 0 & \Sigma_{yy} - B\Sigma_{xx}B^T \end{pmatrix} \end{aligned} \quad (4.48)$$

### 3. Conditional Distribution

We have established that  $y = m(x) + e$  where  $e$  is independent of  $x$ . To find the distribution of  $y$  conditional on  $x$ , we observe that  $m(x)$  becomes a constant vector when  $x$  is fixed, and the randomness comes solely from  $e$ :

$$E[y|x] = m(x) + E[e|x] = m(x) + 0 = m(x) \quad (4.49)$$

$$\text{Var}(y|x) = \text{Var}(m(x)|x) + \text{Var}(e|x) = 0 + \text{Var}(e) = \Sigma_{y|x} \quad (4.50)$$

Thus,  $y|x \sim N(m(x), \Sigma_{y|x})$ . □

## 4.7.1 Connections with Other Formulas

### 4.7.1.1 Rao-Blackwell Decomposition of Variance

The Law of Total Variance (Rao-Blackwell theorem) allows us to decompose the total variance of  $y$  into two orthogonal components based on the predictor  $x$ :

$$\text{Var}(y) = \underbrace{E[\text{Var}(y|x)]}_{\text{Unexplained (Noise)}} + \underbrace{\text{Var}[E(y|x)]}_{\text{Explained (Signal)}} \quad (4.51)$$

In the Multivariate Normal case, this decomposition perfectly aligns with our regression model  $y = m(x) + e$ .

#### Variance of Noise

This term represents the average variance remaining in  $y$  after accounting for  $x$ . It corresponds to the variance of the error term  $e$ :

$$E[\text{Var}(y|x)] = \text{Var}(e) = \Sigma_{yy} - B\Sigma_{xx}B^T \quad (4.52)$$

#### Variance of Signal

This term represents the variability of the conditional mean  $m(x)$  itself. Using the matrix  $B$ , this takes the quadratic form:

$$\text{Var}[E(y|x)] = \text{Var}[m(x)] = B\Sigma_{xx}B^T \quad (4.53)$$

#### Total Variance

Summing the Signal and Noise components recovers the total marginal variance of  $y$ :

$$\Sigma_{yy} = \underbrace{\Sigma_{yy} - B\Sigma_{xx}B^T}_{\text{Unexplained (Noise)}} + \underbrace{B\Sigma_{xx}B^T}_{\text{Explained (Signal)}} \quad (4.54)$$

#### 4.7.1.2 Connection to OLS Regression Estimators

In OLS regression, centering the data allows us to separate the intercept from the slopes. Let  $\mathbf{y}_c$  and  $\mathbf{X}_c$  be the centered response and design matrices (where  $\mathbf{X}_c$  **excludes the column of 1s**). Using this centered form, the total sum of squares decomposes exactly like the population variance:

$$\text{SST} = \text{SSR} + \text{SSE} \quad (4.55)$$

Comparing the sample quantities to their population counterparts:

1. **Regression Coefficients:**

$$\hat{\beta}^T = (\mathbf{X}_c^T \mathbf{X}_c)^{-1} \mathbf{X}_c^T \mathbf{y}_c \approx B \quad (4.56)$$

*Note:  $\hat{\beta}$  here represents only the slope coefficients, matching the dimensions of the covariance matrix  $\Sigma_{xx}$ .*

2. **Explained Variation (Signal):**

$$\text{SSR} = \hat{\beta}^T (\mathbf{X}_c^T \mathbf{X}_c) \hat{\beta} \approx (n-1) B \Sigma_{xx} B^T \quad (4.57)$$

3. **Unexplained Variation (Noise):**

$$\text{SSE} = \mathbf{y}_c^T \mathbf{y}_c - \hat{\beta}^T (\mathbf{X}_c^T \mathbf{X}_c) \hat{\beta} \approx (n-1) (\Sigma_{yy} - B \Sigma_{xx} B^T) \quad (4.58)$$

## 4.8 Partial and Multiple Correlation

**Definition 4.9** (Partial Correlation). The partial correlation between elements  $y_i$  and  $y_j$  given a set of variables  $x$  is derived from the conditional covariance matrix  $\Sigma_{y|x}$ :

$$\rho_{ij|x} = \frac{\sigma_{ij|x}}{\sqrt{\sigma_{ii|x} \sigma_{jj|x}}} \quad (4.59)$$

where  $\sigma_{ij|x}$  are elements of  $\Sigma_{y|x} = \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}$ .

**Definition 4.10** (Multiple Correlation ( $R^2$ )). For a scalar  $y$  and vector  $x$ , the squared multiple correlation is the proportion of variance of  $y$  explained by the conditional mean:

$$R_{y|x}^2 = \frac{\text{Var}(E(y|x))}{\text{Var}(y)} = \frac{\Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}}{\sigma_y^2} \quad (4.60)$$

Note: this definition is the population or theoretical  $R^2$ , which is estimated by adjusted  $R^2$  using sample in linear regression.

## 4.9 Examples

**Example 4.1** (Bivariate Normal). Let the random vector  $\begin{pmatrix} y \\ x \end{pmatrix}$  follow a bivariate normal distribution:

$$\begin{pmatrix} y \\ x \end{pmatrix} \sim N \left( \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \begin{pmatrix} 2 & 2 \\ 2 & 4 \end{pmatrix} \right) \quad (4.61)$$

Here,  $\mu_y = 1$ ,  $\mu_x = 2$ ,  $\Sigma_{yy} = 2$ ,  $\Sigma_{xx} = 4$ , and  $\Sigma_{yx} = 2$ .

**1. Finding the Regression Coefficient Matrix  $B$**  Using the population formula:

$$B = \Sigma_{yx} \Sigma_{xx}^{-1} = 2(4)^{-1} = 0.5 \quad (4.62)$$

**2. Finding the Conditional Mean  $m(x)$  (The Signal)** The systematic component represents the projection of  $y$  onto  $x$ :

$$\begin{aligned} m(x) &= \mu_y + B(x - \mu_x) \\ &= 1 + 0.5(x - 2) = 0.5x \end{aligned} \quad (4.63)$$

**3. Variance of the Signal  $\text{Var}(m(x))$**  Using the quadratic form established in the theorem:

$$\text{Var}(m(x)) = B \Sigma_{xx} B^T = 0.5(4)(0.5) = 1 \quad (4.64)$$

**4. Variance of the Noise  $\text{Var}(y|x)$  (The Residual)** By the Signal-Noise Decomposition:

$$\begin{aligned} \text{Var}(y|x) &= \Sigma_{yy} - \text{Var}(m(x)) \\ &= 2 - 1 = 1 \end{aligned} \quad (4.65)$$

Thus,  $y|x \sim N(m(x), 1)$ . The total variance (2) is split equally between signal (1) and noise (1).

**5. Multiple Correlation Coefficient ( $R^2$ )**

$$R^2 = \frac{\text{Var}(m(x))}{\Sigma_{yy}} = \frac{1}{2} = 0.5 \quad (4.66)$$

**Example 4.2** (Trivariate Normal with 2 Predictors). Let  $V = (y, x_1, x_2)' \sim N_3(\mu, \Sigma)$  with:

$$\mu = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} 10 & 3 & 4 \\ 3 & 2 & 1 \\ 4 & 1 & 4 \end{pmatrix} \quad (4.67)$$

We partition these into  $\Sigma_{yy} = 10$ ,  $\Sigma_{yx} = \begin{pmatrix} 3 & 4 \end{pmatrix}$ , and  $\Sigma_{xx} = \begin{pmatrix} 2 & 1 \\ 1 & 4 \end{pmatrix}$ .

**1. Finding the Regression Coefficient Matrix  $B$**

$$\Sigma_{xx}^{-1} = \frac{1}{7} \begin{pmatrix} 4 & -1 \\ -1 & 2 \end{pmatrix} \Rightarrow B = \Sigma_{yx} \Sigma_{xx}^{-1} = \begin{pmatrix} \frac{8}{7} & \frac{5}{7} \end{pmatrix} \quad (4.68)$$

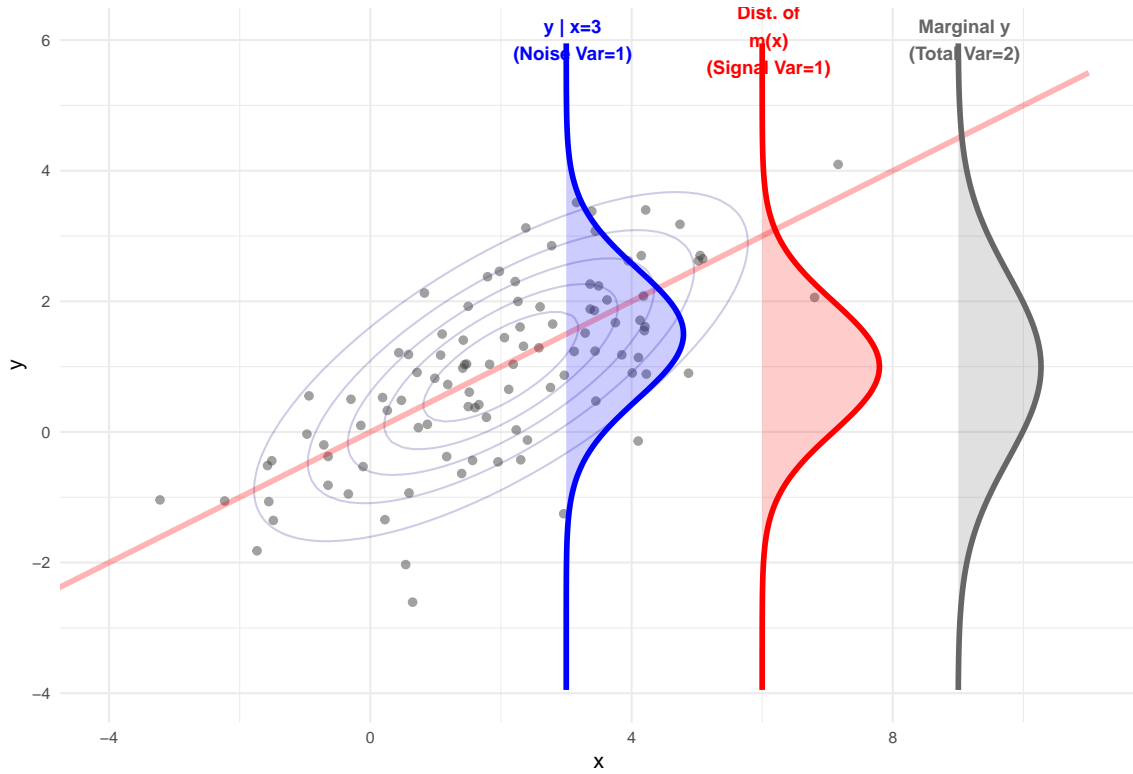


Figure 4.1: Illustration of Rao-Blackwell Variance Decomposition in Bivariate Normal

## 2. Finding the Conditional Mean $m(x)$ (The Signal)

$$m(x) = 1 + \frac{8}{7}(x_1 - 2) + \frac{5}{7}(x_2 - 3) \quad (4.69)$$

## 3. Variance of the Signal $\text{Var}(m(x))$

$$\text{Var}(m(x)) = B\Sigma_{xx}B^T = \begin{pmatrix} \frac{8}{7} & \frac{5}{7} \end{pmatrix} \begin{pmatrix} 3 \\ 4 \end{pmatrix} = \frac{44}{7} \approx 6.29 \quad (4.70)$$

## 4. Variance of the Noise $\text{Var}(y|x)$ (The Residual) Using the Signal-Noise Decomposition:

$$\Sigma_{y|x} = \Sigma_{yy} - \text{Var}(m(x)) = 10 - 6.29 = 3.71 \quad (4.71)$$

## 5. Multiple Correlation Coefficient ( $R^2$ )

$$R^2 = \frac{6.29}{10} = 0.629 \quad (4.72)$$

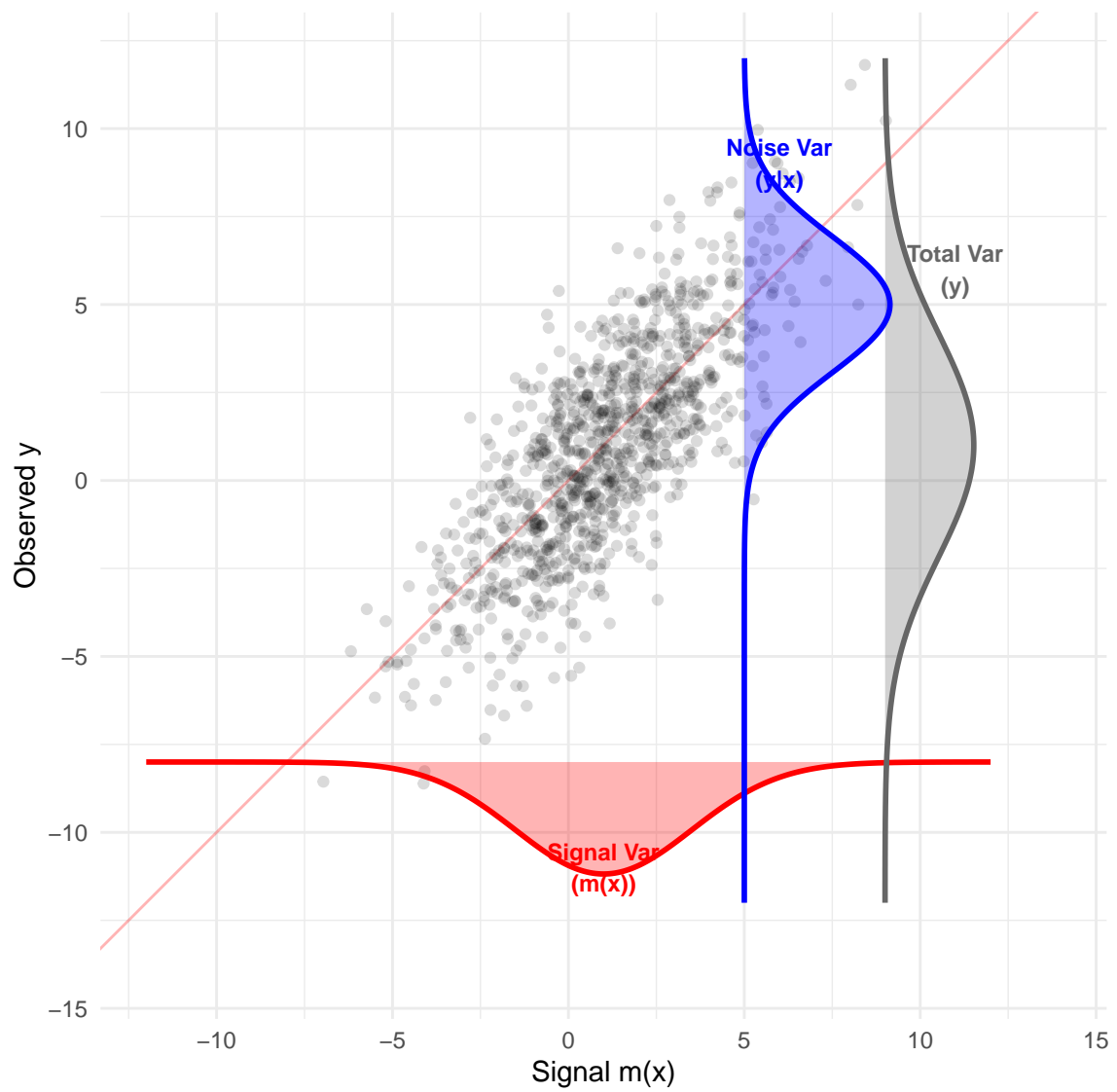


Figure 4.2: Signal-Noise Variance Decomposition in Multivariate Normal



## 5 Distribution of Quadratic Forms

This chapter covers the distribution of quadratic forms (sums of squares), which is crucial for hypothesis testing in linear models.

### 5.1 Quadratic Forms

A quadratic form is a polynomial with terms all of degree two.

**Definition 5.1** (Quadratic Form). Let  $y = (y_1, \dots, y_n)'$  be a random vector and  $A$  be a symmetric  $n \times n$  matrix. The scalar quantity  $y' Ay$  is called a **quadratic form** in  $y$ .

$$y' Ay = \sum_{i=1}^n \sum_{j=1}^n a_{ij} y_i y_j \quad (5.1)$$

**Examples:**

- **Squared Norm:** If  $A = I_n$ , then  $y' I_n y = y' y = \sum y_i^2 = ||y||^2$ .
- **Weighted Sum of Squares:** If  $A$  is diagonal with elements  $\lambda_i$ , then  $y' Ay = \sum \lambda_i y_i^2$ .
- **Projection Sum of Squares:** If  $P$  is a projection matrix,  $||Py||^2 = (Py)'(Py) = y' P' Py = y' Py$  (since  $P$  is symmetric and idempotent).

### 5.2 Mean of Quadratic Forms

We can find the expected value of a quadratic form without assuming normality.

**Lemma 5.1** (Mean of Simplified Quadratic Form). If  $y$  is a random vector with mean  $E(y) = \mu$  and covariance matrix  $\text{Var}(y) = I_n$ , then:

$$E(y' y) = \text{tr}(I_n) + \mu' \mu = n + \mu' \mu \quad (5.2)$$

*Proof.* Let us decompose  $y$  into its mean and a stochastic component:  $y = \mu + z$ , where  $E(z) = 0$  and  $\text{Var}(z) = E(z z') = I_n$ . Substituting this into the quadratic form:

$$\begin{aligned} y' y &= (\mu + z)'(\mu + z) \\ &= \mu' \mu + \mu' z + z' \mu + z' z \\ &= \mu' \mu + 2\mu' z + z' z \end{aligned} \quad (5.3)$$

Taking the expectation:

$$\begin{aligned}
 E(y'y) &= \mu'\mu + 2\mu'E(z) + E(z'z) \\
 &= \mu'\mu + 0 + E\left(\sum_{i=1}^n z_i^2\right)
 \end{aligned} \tag{5.4}$$

Since  $\text{Var}(z_i) = E(z_i^2) - (E(z_i))^2 = 1 - 0 = 1$ , we have  $E(\sum z_i^2) = \sum 1 = n$ . Thus,  $E(y'y) = n + \mu'\mu$ .  $\square$

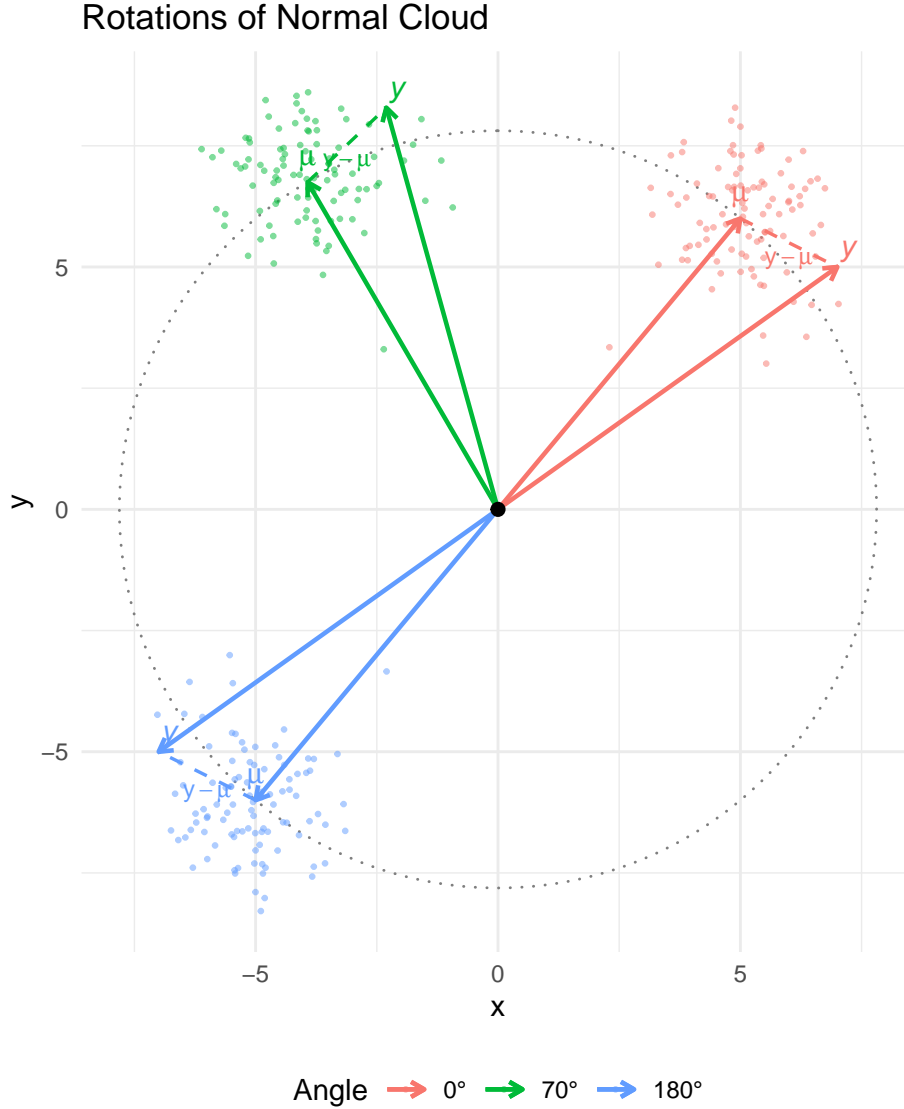


Figure 5.1: Illustration of the Mean and Distribution of Quadratic Forms

**Theorem 5.1** (Mean of Quadratic Form). *If  $y$  is a random vector with mean  $E(y) = \mu$  and covariance matrix  $\text{Var}(y) = \Sigma$ , and  $A$  is a symmetric matrix of constants, then:*

$$E(y'Ay) = \text{tr}(A\Sigma) + \mu'A\mu \tag{5.5}$$

*Proof.* We present three methods to derive the expectation of the quadratic form.

### Method 1: Using the Trace Trick

Using the fact that a scalar is equal to its own trace ( $\text{tr}(c) = c$ ) and the linearity of expectation:

$$\begin{aligned} E(y' Ay) &= E[\text{tr}(y' Ay)] \\ &= E[\text{tr}(Ayy')] \quad (\text{cyclic property of trace}) \\ &= \text{tr}(AE[yy']) \quad (\text{linearity of expectation}) \end{aligned} \tag{5.6}$$

Recall that the covariance matrix is defined as  $\Sigma = E[(y - \mu)(y - \mu)'] = E(yy') - \mu\mu'$ . Rearranging this gives the second moment:  $E(yy') = \Sigma + \mu\mu'$ . Substituting this back:

$$\begin{aligned} E(y' Ay) &= \text{tr}(A(\Sigma + \mu\mu')) \\ &= \text{tr}(A\Sigma) + \text{tr}(A\mu\mu') \\ &= \text{tr}(A\Sigma) + \text{tr}(\mu' A\mu) \quad (\text{cyclic property on second term}) \\ &= \text{tr}(A\Sigma) + \mu' A\mu \end{aligned} \tag{5.7}$$

### Method 2: Using Scalar Summation

We can express the quadratic form in scalar notation using the entries of  $A = (a_{ij})$ ,  $\Sigma = (\sigma_{ij})$ , and  $\mu = (\mu_i)$ :

$$\begin{aligned} E(y' Ay) &= E\left(\sum_{i=1}^n \sum_{j=1}^n a_{ij} y_i y_j\right) \\ &= \sum_{i=1}^n \sum_{j=1}^n a_{ij} E(y_i y_j) \\ &= \sum_{i=1}^n \sum_{j=1}^n a_{ij} (\sigma_{ij} + \mu_i \mu_j) \\ &= \sum_{i=1}^n \sum_{j=1}^n a_{ij} \sigma_{ji} + \sum_{i=1}^n \sum_{j=1}^n \mu_i a_{ij} \mu_j \quad (\text{since } \Sigma \text{ is symmetric, } \sigma_{ij} = \sigma_{ji}) \\ &= \text{tr}(A\Sigma) + \mu' A\mu \end{aligned} \tag{5.8}$$

### Method 3: Using Spectral Decomposition of A

Since  $A$  is symmetric, we use its spectral decomposition  $A = \sum_{i=1}^n \lambda_i q_i q_i'$ . Substituting this into the quadratic form:

$$y' Ay = y' \left( \sum_{i=1}^n \lambda_i q_i q_i' \right) y = \sum_{i=1}^n \lambda_i (q_i' y)^2 \tag{5.9}$$

Let  $w_i = q_i' y$ . This is a scalar random variable which is a linear transformation of  $y$ . Its properties are:

1. **Mean:**  $E(w_i) = q_i' E(y) = q_i' \mu$ .
2. **Variance:**  $\text{Var}(w_i) = \text{Var}(q_i' y) = q_i' \text{Var}(y) q_i = q_i' \Sigma q_i$ .

Using the relation  $E(w_i^2) = \text{Var}(w_i) + [E(w_i)]^2$ , we have:

$$E[(q'_i y)^2] = q'_i \Sigma q_i + (q'_i \mu)^2 \quad (5.10)$$

Summing over all  $i$  weighted by  $\lambda_i$ :

$$\begin{aligned} E(y' A y) &= \sum_{i=1}^n \lambda_i [q'_i \Sigma q_i + (q'_i \mu)^2] \\ &= \sum_{i=1}^n \text{tr}(\lambda_i q'_i \Sigma q_i) + \mu' \left( \sum_{i=1}^n \lambda_i q_i q'_i \right) \mu \\ &= \text{tr} \left( \Sigma \sum_{i=1}^n \lambda_i q_i q'_i \right) + \mu' A \mu \\ &= \text{tr}(\Sigma A) + \mu' A \mu \end{aligned} \quad (5.11)$$

□

*Remark (Geometric Interpretation via Sigma).* If we further decompose  $\Sigma = \sum_{j=1}^n \gamma_j v_j v'_j$  (where  $\gamma_j, v_j$  are eigenvalues/vectors of  $\Sigma$ ), the trace term becomes:

$$\text{tr}(A \Sigma) = \sum_{i=1}^n \sum_{j=1}^n \lambda_i \gamma_j (q'_i v_j)^2 \quad (5.12)$$

Here,  $(q'_i v_j)^2 = \cos^2(\theta_{ij})$  represents the alignment between the axes of the quadratic form ( $A$ ) and the axes of the data covariance ( $\Sigma$ ). The expectation is maximized when the eigenspaces of  $A$  and  $\Sigma$  align.

**Corollary 5.1** (Expectation with Projection Matrix). *Consider the special case where:*

1.  $P$  is a **projection matrix** (symmetric and idempotent,  $P^2 = P$ ).
2. The covariance is **spherical**:  $\Sigma = \sigma^2 I_n$ .

Then the expectation simplifies to:

$$E(y' P y) = \sigma^2 r + \|P \mu\|^2 \quad (5.13)$$

where  $r = \text{rank}(P) = \text{tr}(P)$ .

**Proof:** Using Theorem 5.1 with  $A = P$  and  $\Sigma = \sigma^2 I_n$ :

1. **Trace Term:**  $\text{tr}(P \Sigma) = \text{tr}(P(\sigma^2 I_n)) = \sigma^2 \text{tr}(P)$ . Since  $P$  is idempotent, its eigenvalues are either 0 or 1, so  $\text{tr}(P) = \text{rank}(P) = r$ .
2. **Mean Term:** Since  $P$  is symmetric and idempotent ( $P' P = P^2 = P$ ), we can rewrite the quadratic form:

$$\mu' P \mu = \mu' P' P \mu = (P \mu)' (P \mu) = \|P \mu\|^2 \quad (5.14)$$

**Example 5.1** (Expectation of Sum of Squares Decomposition (i.i.d. Case)). Consider a random vector  $y = (y_1, \dots, y_n)'$  with mean vector  $\mu_y = \mu j_n$  and covariance  $\Sigma = \sigma^2 I_n$ . We analyze the two components of the total sum of squares by projecting  $y$  onto the mean space ( $P_{j_n}$ ) and the residual space ( $I - P_{j_n}$ ).

## 1. The Projection Vectors

First, we write the explicit forms of the projected vectors using  $P_{j_n} = \frac{1}{n} j_n j'_n$ :

- **Mean Vector** ( $P_{j_n} y$ ): Projecting  $y$  onto the column space of  $j_n$  replaces every element with the sample mean  $\bar{y}$ .

$$P_{j_n} y = \bar{y} j_n = \begin{pmatrix} \bar{y} \\ \bar{y} \\ \vdots \\ \bar{y} \end{pmatrix} \quad (5.15)$$

- **Residual Vector** ( $(I - P_{j_n})y$ ): Subtracting the mean projection from  $y$  yields the deviations.

$$(I - P_{j_n})y = y - \bar{y} j_n = \begin{pmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix} \quad (5.16)$$

## 2. Expectations of Squared Norms

We now find the expectation of the squared length of these vectors using Corollary 5.1.

**Part A: Sum of Squares for Mean** The quadratic form is the squared norm of the projected mean vector:

$$y' P_{j_n} y = \|P_{j_n} y\|^2 = \sum_{i=1}^n \bar{y}^2 = n\bar{y}^2 \quad (5.17)$$

Applying the corollary with  $P = P_{j_n}$ :

- **Rank:**  $\text{tr}(P_{j_n}) = 1$ .
- **Mean:**  $P_{j_n} \mu_y = P_{j_n} (\mu j_n) = \mu j_n$ . The squared norm is  $n\mu^2$ .

$$E[\|P_{j_n} y\|^2] = \sigma^2(1) + n\mu^2 \quad (5.18)$$

**Part B: Sum of Squared Errors (SSE)** The quadratic form is the squared norm of the residual vector:

$$y' (I - P_{j_n}) y = \|(I - P_{j_n}) y\|^2 = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (5.19)$$

Applying the corollary with  $P = I - P_{j_n}$ :

- **Rank:**  $\text{tr}(I - P_{j_n}) = n - 1$ .
- **Mean:**  $(I - P_{j_n}) \mu_y = \mu_y - P_{j_n} \mu_y = \mu j_n - \mu j_n = 0$ . The squared norm is 0.

$$E[\|(I - P_{j_n}) y\|^2] = \sigma^2(n - 1) + 0 \quad (5.20)$$

**Conclusion** These results confirm the standard properties:  $E(\bar{y}^2) = \frac{\sigma^2}{n} + \mu^2$  and  $E(S^2) = \sigma^2$ .

**Example 5.2** (Expectation of Total Sum of Squares (Regression Case)). Consider now a regression setting where the mean of  $y$  depends on covariates (e.g.,  $\mu_i = \beta_0 + \beta_1 x_i$ ). The mean vector  $\mu_y$  is **not** proportional to  $j_n$ . We are interested in the expectation of the **Total Sum of Squares (SST)**.

**1. Identification** The SST measures the variation of  $y$  around the *global sample mean*  $\bar{y}$ , ignoring the covariates:

$$\text{SST} = \sum_{i=1}^n (y_i - \bar{y})^2 = y'(I - P_{j_n})y \quad (5.21)$$

This is the same quadratic form as Part B in the previous example, but the underlying mean  $\mu_y$  has changed.

**2. Calculation** We apply Corollary 5.1 with  $P = I - P_{j_n}$  and general  $\mu_y$ :

- **Rank Term:** Same as before,  $\text{tr}(I - P_{j_n}) = n - 1$ .
- **Mean Term:** The projection of the mean vector is no longer zero.

$$(I - P_{j_n})\mu_y = \mu_y - \bar{\mu}j_n = \begin{pmatrix} \mu_1 - \bar{\mu} \\ \vdots \\ \mu_n - \bar{\mu} \end{pmatrix} \quad (5.22)$$

where  $\bar{\mu} = \frac{1}{n} \sum \mu_i$  is the average of the true means. The squared norm is the sum of squared deviations of the true means:

$$\|(I - P_{j_n})\mu_y\|^2 = \sum_{i=1}^n (\mu_i - \bar{\mu})^2 \quad (5.23)$$

**Conclusion**

$$E(\text{SST}) = (n - 1)\sigma^2 + \sum_{i=1}^n (\mu_i - \bar{\mu})^2 \quad (5.24)$$

This shows that in regression, the SST estimates  $(n - 1)\sigma^2$  *plus* the variability introduced by the regression signal (the spread of the true means  $\mu_i$ ).

## 5.3 Non-central $\chi^2$ Distribution

To understand the distribution of quadratic forms under normality, we introduce the non-central chi-square distribution.

**Definition 5.2** (Non-central  $\chi^2$  Distribution). Let  $y \sim N_n(\mu, I_n)$ . The random variable  $V = y'y = \sum y_i^2$  follows a **non-central chi-square distribution** with  $n$  degrees of freedom and non-centrality parameter  $\lambda$ .

$$V \sim \chi^2(n, \lambda) \quad \text{where } \lambda = \frac{1}{2}\mu'\mu = \frac{1}{2}\|\mu\|^2 \quad (5.25)$$

**Note:** Some definitions of non-central  $\chi^2$  use  $\lambda = \mu'\mu$ . In this course, we use  $\lambda = \frac{1}{2}\mu'\mu$ .

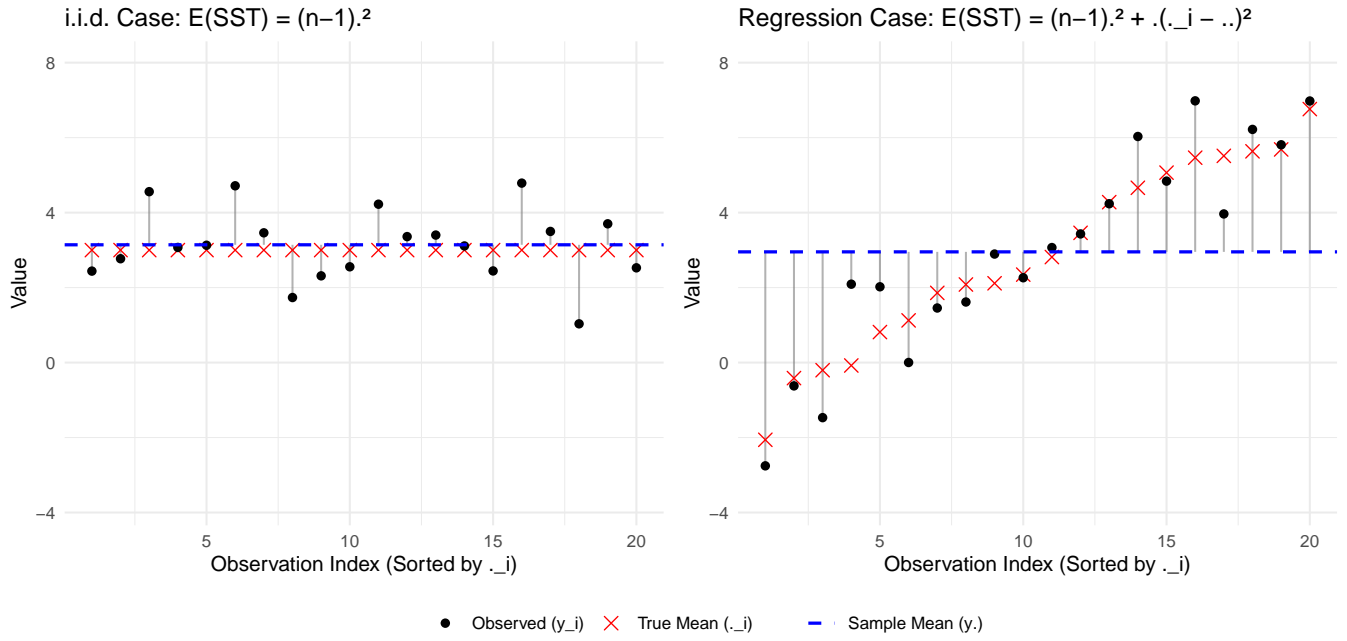


Figure 5.2: Comparison of SST components with increased variation in the true means. The vertical lines represent the deviations  $(y_i - \bar{y})$ . With  $\text{sd}(\mu_i) = 3$ , the regression case (right) shows significantly larger deviations, illustrating how the systematic spread of the means dominates the Total Sum of Squares.

### 5.3.1 Visualizing $\chi^2$ Distributions

Here is a plot visualizing the difference between central and non-central Chi-square distributions.

The density of the non-central chi-square distribution shifts to the right and becomes flatter as the non-centrality parameter  $\lambda$  increases.

### 5.3.2 Mean, Variance, and MGF

We summarize the key properties of the non-central chi-square distribution.

**Theorem 5.2** (Properties of Non-central Chi-square). *Let  $V \sim \chi^2(n, \lambda)$ . Then:*

1. **Mean:**  $E(V) = n + 2\lambda$
2. **Variance:**  $\text{Var}(V) = 2n + 8\lambda$
3. **Moment Generating Function (MGF):**

$$m_V(t) = \frac{\exp[-\lambda\{1 - 1/(1 - 2t)\}]}{(1 - 2t)^{n/2}} \quad \text{for } t < 1/2 \quad (5.26)$$

*Mean.* By definition,  $V \sim \chi^2(n, \lambda)$  is the distribution of  $y'y$  where  $y \sim N_n(\mu, I_n)$  and the non-centrality parameter is  $\lambda = \frac{1}{2}\mu'\mu$ . Applying Lemma 5.1 to the random vector  $y$ :

$$E(V) = E(y'y) = n + \mu'\mu = n + 2\lambda \quad (5.27)$$

## Chi-square Distributions

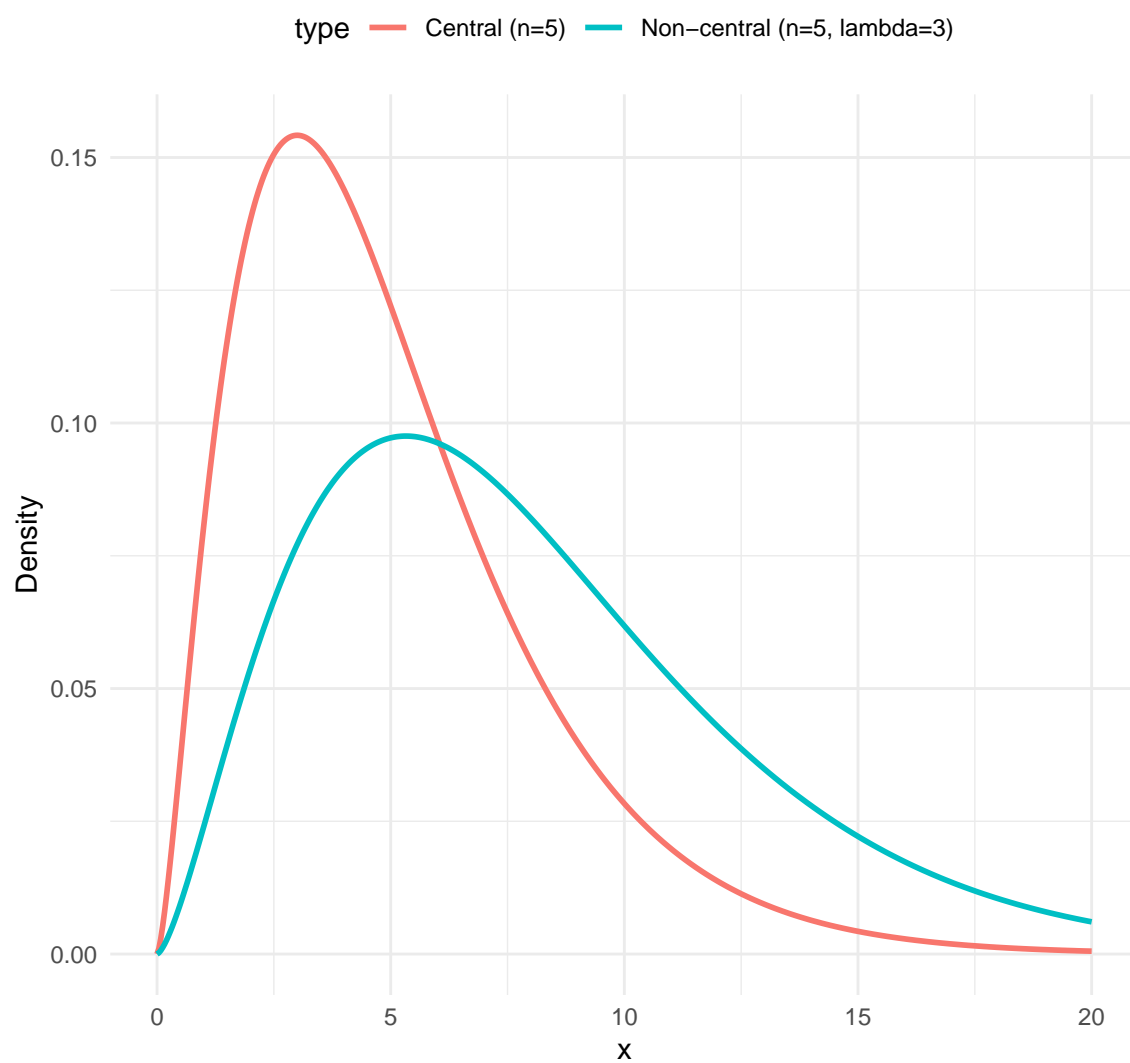


Figure 5.3: Central vs Non-central Chi-square Distribution



□

*MGF.* Since the components  $y_i$  of the vector  $y$  are independent  $N(\mu_i, 1)$ , and  $V = \sum_{i=1}^n y_i^2$ , the MGF of  $V$  is the product of the MGFs of each  $y_i^2$ :

$$m_V(t) = E[e^{t \sum y_i^2}] = \prod_{i=1}^n E[e^{t y_i^2}] \quad (5.28)$$

Consider a single component  $y_i \sim N(\mu_i, 1)$ . Its squared expectation is:

$$\begin{aligned} E[e^{t y_i^2}] &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{t y^2} e^{-\frac{1}{2}(y-\mu_i)^2} dy \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \exp \left\{ -\frac{1}{2} [(1-2t)y^2 - 2\mu_i y + \mu_i^2] \right\} dy \end{aligned} \quad (5.29)$$

Completing the square in the exponent for  $y$  (assuming  $t < 1/2$ ):

$$(1-2t)y^2 - 2\mu_i y + \mu_i^2 = (1-2t) \left( y - \frac{\mu_i}{1-2t} \right)^2 + \mu_i^2 - \frac{\mu_i^2}{1-2t} \quad (5.30)$$

The integral of the Gaussian kernel  $\exp\{-\frac{1}{2}(1-2t)(y - \dots)^2\}$  yields  $\sqrt{\frac{2\pi}{1-2t}}$ . The remaining constant term is:

$$\exp \left\{ -\frac{1}{2} \left( \mu_i^2 - \frac{\mu_i^2}{1-2t} \right) \right\} = \exp \left\{ \frac{\mu_i^2}{2} \left( \frac{1}{1-2t} - 1 \right) \right\} = \exp \left\{ \frac{\mu_i^2 t}{1-2t} \right\} \quad (5.31)$$

Thus, for a single component:

$$m_{y_i^2}(t) = (1-2t)^{-1/2} \exp \left( \frac{\mu_i^2 t}{1-2t} \right) \quad (5.32)$$

Multiplying the MGFs for all  $n$  components:

$$\begin{aligned} m_V(t) &= \prod_{i=1}^n (1-2t)^{-1/2} \exp \left( \frac{\mu_i^2 t}{1-2t} \right) \\ &= (1-2t)^{-n/2} \exp \left( \frac{t \sum \mu_i^2}{1-2t} \right) \end{aligned} \quad (5.33)$$

Substituting  $\lambda = \frac{1}{2} \sum \mu_i^2$  (so  $\sum \mu_i^2 = 2\lambda$ ):

$$m_V(t) = (1-2t)^{-n/2} \exp \left( \frac{2\lambda t}{1-2t} \right) \quad (5.34)$$

This form matches the theorem statement, noting that  $\frac{2\lambda t}{1-2t} = -\lambda(1 - \frac{1}{1-2t})$ . □

*Variance.* We use the **Cumulant Generating Function**,  $K_V(t) = \ln m_V(t)$ , as its derivatives yield the mean and variance directly:

$$K_V(t) = -\frac{n}{2} \ln(1-2t) + \frac{2\lambda t}{1-2t} \quad (5.35)$$

First derivative (Mean):

$$\begin{aligned} K'_V(t) &= -\frac{n}{2} \left( \frac{-2}{1-2t} \right) + 2\lambda \left[ \frac{1(1-2t) - t(-2)}{(1-2t)^2} \right] \\ &= \frac{n}{1-2t} + 2\lambda \frac{1}{(1-2t)^2} \end{aligned} \quad (5.36)$$

Second derivative (Variance):

$$\begin{aligned} K''_V(t) &= n(-1)(1-2t)^{-2}(-2) + 2\lambda(-2)(1-2t)^{-3}(-2) \\ &= \frac{2n}{(1-2t)^2} + \frac{8\lambda}{(1-2t)^3} \end{aligned} \quad (5.37)$$

Evaluating at  $t = 0$ :

$$\text{Var}(V) = K''_V(0) = 2n + 8\lambda \quad (5.38)$$

□

### 5.3.3 Additivity

**Theorem 5.3** (Additivity of Chi-square). *If  $v_1, \dots, v_k$  are independent random variables distributed as  $\chi^2(n_i, \lambda_i)$ , then their sum follows a chi-square distribution:*

$$\sum_{i=1}^k v_i \sim \chi^2 \left( \sum_{i=1}^k n_i, \sum_{i=1}^k \lambda_i \right) \quad (5.39)$$

*Proof.* **Method 1: Using MGFs**

The moment generating function of  $v_i \sim \chi^2(n_i, \lambda_i)$  is:

$$M_{v_i}(t) = \frac{\exp \left[ -\lambda_i \left( 1 - \frac{1}{1-2t} \right) \right]}{(1-2t)^{n_i/2}} \quad (5.40)$$

Since  $v_1, \dots, v_k$  are independent, the MGF of their sum  $V = \sum v_i$  is the product of their individual MGFs:

$$\begin{aligned} M_V(t) &= \prod_{i=1}^k M_{v_i}(t) \\ &= \prod_{i=1}^k \frac{\exp \left[ -\lambda_i \left( 1 - \frac{1}{1-2t} \right) \right]}{(1-2t)^{n_i/2}} \\ &= \frac{\exp \left[ -\sum \lambda_i \left( 1 - \frac{1}{1-2t} \right) \right]}{(1-2t)^{\sum n_i/2}} \end{aligned} \quad (5.41)$$

This is the MGF of a non-central chi-square distribution with degrees of freedom  $\sum n_i$  and non-centrality parameter  $\sum \lambda_i$ .

**Method 2: Geometric Interpretation**

Let  $v_i = ||y_i||^2$  where  $y_i \sim N_{n_i}(\mu_i, I_{n_i})$ . Since the vectors  $y_i$  are independent, we can stack them into a larger vector  $y = (y'_1, \dots, y'_k)'$ .

$$y \sim N_{\sum n_i}(\mu, I_{\sum n_i}) \quad \text{where } \mu = (\mu'_1, \dots, \mu'_k)' \quad (5.42)$$

The sum of squares is:

$$\sum v_i = \sum ||y_i||^2 = ||y||^2 \quad (5.43)$$

By definition,  $||y||^2$  follows a non-central chi-square distribution with degrees of freedom equal to the dimension of  $y$  ( $\sum n_i$ ) and non-centrality parameter  $\lambda = \frac{1}{2}||\mu||^2$ .

$$\lambda = \frac{1}{2} \sum_{i=1}^k ||\mu_i||^2 = \sum_{i=1}^k \lambda_i \quad (5.44)$$

□

### 5.3.4 Poisson Mixture Representation

**Theorem 5.4** (Poisson Mixture Representation). *Let  $v \sim \chi^2(n, \lambda)$  be a non-central chi-square random variable. Its probability density function can be represented as a Poisson-weighted sum of central chi-square density functions:*

$$f(v; n, \lambda) = \sum_{j=0}^{\infty} \left( \frac{e^{-\lambda} \lambda^j}{j!} \right) f(v; n + 2j, 0) \quad (5.45)$$

where  $f(v; \nu, 0)$  is the density of a central chi-square distribution with  $\nu$  degrees of freedom.

*Proof.* We use the Moment Generating Function (MGF) approach. The MGF of a non-central chi-square distribution  $v \sim \chi^2(n, \lambda)$  is:

$$M_v(t) = (1 - 2t)^{-n/2} \exp \left( \lambda \left[ \frac{1}{1 - 2t} - 1 \right] \right) \quad (5.46)$$

We can expand the exponential term using the power series  $e^x = \sum_{j=0}^{\infty} \frac{x^j}{j!}$ :

$$\begin{aligned} M_v(t) &= (1 - 2t)^{-n/2} e^{-\lambda} \exp \left( \frac{\lambda}{1 - 2t} \right) \\ &= e^{-\lambda} (1 - 2t)^{-n/2} \sum_{j=0}^{\infty} \frac{1}{j!} \left( \frac{\lambda}{1 - 2t} \right)^j \\ &= \sum_{j=0}^{\infty} \left( \frac{e^{-\lambda} \lambda^j}{j!} \right) (1 - 2t)^{-(n+2j)/2} \end{aligned} \quad (5.47)$$

Recognizing the terms:

1. The term in parentheses,  $P(J = j) = \frac{e^{-\lambda} \lambda^j}{j!}$ , is the probability mass function of a **Poisson** random variable  $J \sim \text{Poisson}(\lambda)$ .
2. The term  $(1 - 2t)^{-(n+2j)/2}$  is the MGF of a **central chi-square** distribution with  $n + 2j$  degrees of freedom.

Since the MGF of the mixture is the sum of the MGFs of the components weighted by the mixture probabilities, the density must follow the same mixture structure.  $\square$

*Remark.* This theorem implies a hierarchical model for generating a non-central chi-square variable:

1. Sample  $J \sim \text{Poisson}(\lambda)$ .
2. Given  $J = j$ , sample  $V \sim \chi^2(n + 2j, 0)$ .

This is particularly useful for numerical computation, as it allows the non-central CDF to be approximated by a finite sum of central chi-square CDFs.

## 5.4 Distribution of Quadratic Forms

### 5.4.1 MGF of Quadratic Forms

To determine the distribution of general quadratic forms  $y' Ay$ , we look at their MGF.

**Theorem 5.5** (MGF of Quadratic Form). *If  $y \sim N_p(\mu, \Sigma)$ , then the MGF of  $Q = y' Ay$  is:*

$$M_Q(t) = |I - 2tA\Sigma|^{-1/2} \exp\left(-\frac{1}{2}\mu'[I - (I - 2tA\Sigma)^{-1}]\Sigma^{-1}\mu\right) \quad (5.48)$$

### 5.4.2 Non-central $\chi^2$ of Quadratic Forms

We will prove a simplified version of Theorem 5.7 first.

**Theorem 5.6** (Distribution of Projected Spherical Normal). *If  $y \sim N_n(\mu, \sigma^2 I_n)$  and  $P_V$  is a projection matrix onto a subspace  $V$  of dimension  $r$ , then:*

$$\frac{1}{\sigma^2} y' P_V y = \frac{\|P_V y\|^2}{\sigma^2} \sim \chi^2\left(r, \frac{\|P_V \mu\|^2}{2\sigma^2}\right) \quad (5.49)$$

*This holds because  $\frac{1}{\sigma^2} P_V(\sigma^2 I) = P_V$ , which is idempotent.*

#### ! Crucial Theorem

This is one of the most important theorems in the course, establishing the fundamental conditions under which a quadratic form follows a chi-square distribution.

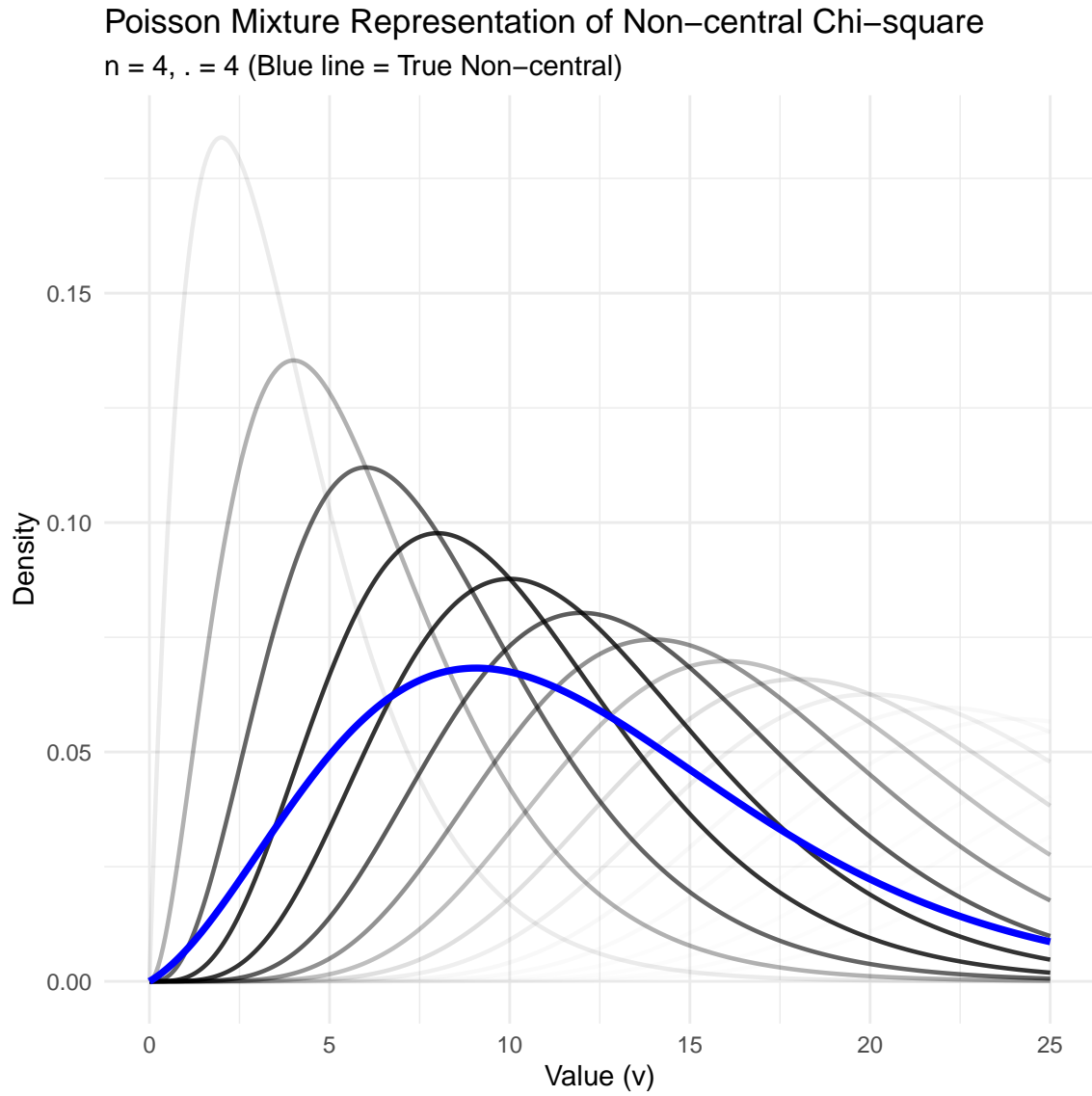


Figure 5.4: The non-central chi-square distribution as a Poisson mixture. The black curves represent central chi-square densities with  $df = n + 2j$ , with transparency (alpha) proportional to the Poisson weight  $P(J = j)$ . The solid blue line is the true non-central chi-square density.

*Proof.* **When**  $\sigma^2 = 1$

Let  $P_V$  be the projection matrix. We know  $P_V = QQ'$  where  $Q = (q_1, \dots, q_r)$  is an  $n \times r$  matrix with orthonormal columns ( $Q'Q = I_r$ ).

The projection of vector  $y$  onto the subspace  $V$  can be expressed using the orthonormal basis vectors:

$$P_V y = QQ' y = (q_1, \dots, q_r) \begin{pmatrix} q_1' y \\ \vdots \\ q_r' y \end{pmatrix} = \sum_{i=1}^r (q_i' y) q_i \quad (5.50)$$

The squared norm of the projection is:

$$y' P_V y = y' Q Q' y = (Q' y)' (Q' y) = \|Q' y\|^2 \quad (5.51)$$

Since  $y \sim N(\mu, I_n)$ , the linear transformation  $z = Q' y$  follows:

$$z \sim N(Q' \mu, Q' I_n Q) = N(Q' \mu, I_r) \quad (5.52)$$

Thus,  $z$  is a vector of  $r$  independent normal variables with variance 1. The sum of squares  $\|z\|^2$  is by definition non-central chi-square:

$$\|z\|^2 \sim \chi^2(r, \lambda) \quad (5.53)$$

where the non-centrality parameter is:

$$\lambda = \frac{1}{2} \|E(z)\|^2 = \frac{1}{2} \|Q' \mu\|^2 \quad (5.54)$$

Note that  $\|Q' \mu\|^2 = \mu' Q Q' \mu = \mu' P_V \mu = \|P_V \mu\|^2$ .

Thus,  $y' P_V y \sim \chi^2(r, \frac{1}{2} \|P_V \mu\|^2)$ .

**When**  $\sigma^2 \neq 1$

If  $y \sim N(\mu, \sigma^2 I_n)$ , we standardize by dividing by  $\sigma$ .

Let  $z = y/\sigma$ . Then  $z \sim N(\mu/\sigma, I_n)$ . Applying the previous result to  $z$ :

$$z' P_V z = \frac{y' P_V y}{\sigma^2} \sim \chi^2 \left( r, \frac{1}{2} \left\| P_V \frac{\mu}{\sigma} \right\|^2 \right) \quad (5.55)$$

which simplifies to:

$$\frac{\|P_V y\|^2}{\sigma^2} \sim \chi^2 \left( r, \frac{\|P_V \mu\|^2}{2\sigma^2} \right) \quad (5.56)$$

□

**! Important**

The term  $\|P_V y\|^2$  itself is **not** a standard chi-square variable; it is a scaled chi-square variable. Its mean is:

$$E(\|P_V y\|^2) = \sigma^2 \left( r + \frac{\|P_V \mu\|^2}{\sigma^2} \right) = r\sigma^2 + \|P_V \mu\|^2 \quad (5.57)$$

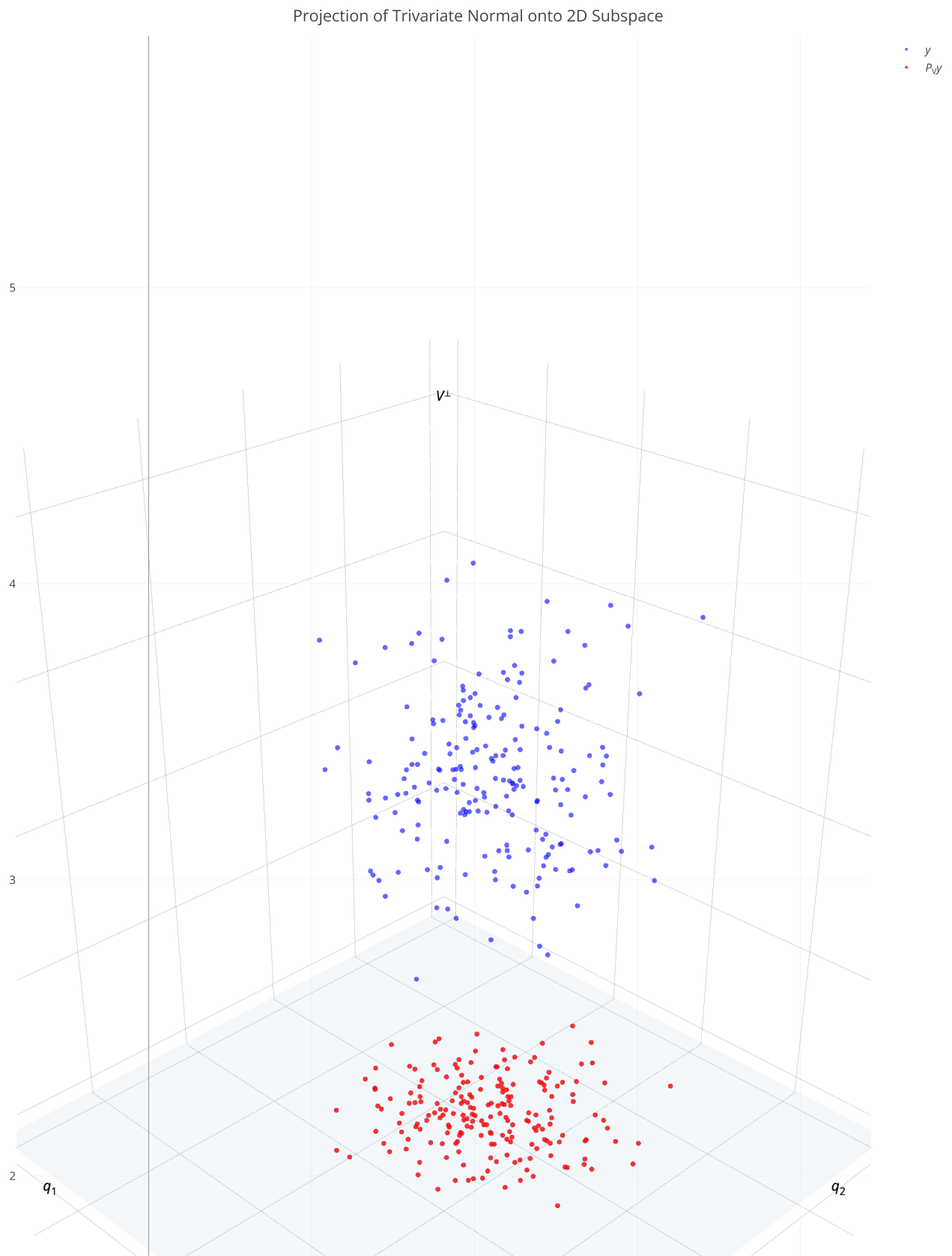


Figure 5.5: Visualization of Projected Trivariate Normal Cloud



### 5.4.3 Distribution of General Quadratic Forms

**Lemma 5.2** (Idempotent Matrix Property). *Let  $\Sigma$  be a positive definite matrix such that  $\Sigma = \Sigma^{1/2}\Sigma^{1/2}$ . The matrix  $A\Sigma$  is idempotent if and only if  $\Sigma^{1/2}A\Sigma^{1/2}$  is idempotent.*

*Proof.* ( $\Rightarrow$ ) Assume  $A\Sigma$  is idempotent, so  $A\Sigma A\Sigma = A\Sigma$ . Then:

$$\begin{aligned} (\Sigma^{1/2}A\Sigma^{1/2})^2 &= \Sigma^{1/2}A(\Sigma^{1/2}\Sigma^{1/2})A\Sigma^{1/2} \\ &= \Sigma^{1/2}(A\Sigma A)\Sigma^{1/2} \end{aligned} \quad (5.58)$$

From the assumption  $A\Sigma A\Sigma = A\Sigma$ , post-multiplying by  $\Sigma^{-1}$  gives  $A\Sigma A = A$ . Substituting this back:

$$\Sigma^{1/2}(A)\Sigma^{1/2} = \Sigma^{1/2}A\Sigma^{1/2} \quad (5.59)$$

( $\Leftarrow$ ) Assume  $\Sigma^{1/2}A\Sigma^{1/2}$  is idempotent. Then:

$$(\Sigma^{1/2}A\Sigma^{1/2})(\Sigma^{1/2}A\Sigma^{1/2}) = \Sigma^{1/2}A\Sigma^{1/2} \quad (5.60)$$

Expanding the left side:

$$\Sigma^{1/2}A(\Sigma^{1/2}\Sigma^{1/2})A\Sigma^{1/2} = \Sigma^{1/2}A\Sigma A\Sigma^{1/2} \quad (5.61)$$

Equating this to the right side:

$$\Sigma^{1/2}A\Sigma A\Sigma^{1/2} = \Sigma^{1/2}A\Sigma^{1/2} \quad (5.62)$$

Pre-multiply by  $\Sigma^{-1/2}$  and post-multiply by  $\Sigma^{1/2}$  (which exist since  $\Sigma$  is positive definite):

$$\begin{aligned} \Sigma^{-1/2}(\Sigma^{1/2}A\Sigma A\Sigma^{1/2})\Sigma^{1/2} &= \Sigma^{-1/2}(\Sigma^{1/2}A\Sigma^{1/2})\Sigma^{1/2} \\ I(A\Sigma A)\Sigma &= I(A)\Sigma \\ A\Sigma A\Sigma &= A\Sigma \end{aligned} \quad (5.63)$$

□

**Lemma 5.3** (Rank Invariance). *Under the conditions of Lemma 5.2, if  $A\Sigma$  is idempotent, then:*

$$\text{rank}(A\Sigma) = \text{rank}(\Sigma^{1/2}A\Sigma^{1/2}) = \text{tr}(A\Sigma) \quad (5.64)$$

*Proof.* Since  $A\Sigma$  and  $\Sigma^{1/2}A\Sigma^{1/2}$  are both idempotent (by Lemma 5.2), their ranks are equal to their traces.

Using the cyclic property of the trace operator ( $\text{tr}(XYZ) = \text{tr}(ZXY)$ ):

$$\begin{aligned} \text{rank}(A\Sigma) &= \text{tr}(A\Sigma) \\ &= \text{tr}(A\Sigma^{1/2}\Sigma^{1/2}) \\ &= \text{tr}(\Sigma^{1/2}A\Sigma^{1/2}) \\ &= \text{rank}(\Sigma^{1/2}A\Sigma^{1/2}) \end{aligned} \quad (5.65)$$

Alternatively, notice that  $A\Sigma$  is similar to  $\Sigma^{1/2}A\Sigma^{1/2}$ :

$$A\Sigma = \Sigma^{-1/2}(\Sigma^{1/2}A\Sigma^{1/2})\Sigma^{1/2} \quad (5.66)$$

Since similar matrices have the same rank, the equality holds. □

**Theorem 5.7** (Distribution of  $y' Ay$ ). Let  $y \sim N_p(\mu, \Sigma)$ . Let  $A$  be a symmetric matrix of rank  $r$ . Then  $y' Ay \sim \chi^2(r, \lambda)$  with  $\lambda = \frac{1}{2} \mu' A \mu$  **if and only if**  $A\Sigma$  is idempotent ( $A\Sigma A\Sigma = A\Sigma$ ).

**Special Case** ( $\Sigma = I$ ): If  $\Sigma = I$ , the condition simplifies to  $A$  being idempotent ( $A^2 = A$ ).

*Proof.* Let  $y^* = \Sigma^{-1/2} y$ , so  $y^* \sim N_n(\Sigma^{-1/2} \mu, I_n)$ . We rewrite the quadratic form:

$$y' Ay = y' \Sigma^{-1/2} (\Sigma^{1/2} A \Sigma^{1/2}) \Sigma^{-1/2} y = (y^*)' P_V y^* = \|P_V y^*\|^2 \quad (5.67)$$

Since  $A\Sigma$  is idempotent,  $P_V = \Sigma^{1/2} A \Sigma^{1/2}$  is a projection matrix with rank  $r$ . By the definition of the non-central chi-square,  $y' Ay \sim \chi^2(r, \frac{1}{2} \|P_V \Sigma^{-1/2} \mu\|^2)$ . The non-centrality parameter simplifies to  $\lambda = \frac{1}{2} \mu' A \mu$ .  $\square$

#### 5.4.4 Standardized Distance Distribution

**Corollary 5.2** (Standardized Distance Distribution). Suppose  $y \sim N_n(\mu, \Sigma)$ . Then the quadratic form representing the standardized distance from a constant vector  $\mu_0$  follows a non-central chi-square distribution:

$$(y - \mu_0)' \Sigma^{-1} (y - \mu_0) \sim \chi^2(n, \lambda = \frac{1}{2} (\mu - \mu_0)' \Sigma^{-1} (\mu - \mu_0)) \quad (5.68)$$

*Proof.* Let  $A = \Sigma^{-1}$ . Then  $A\Sigma = \Sigma^{-1} \Sigma = I_n$ , which is clearly idempotent. Alternatively, let  $z = \Sigma^{-1/2} (y - \mu_0)$ , then  $z \sim N_n(\Sigma^{-1/2} (\mu - \mu_0), I_n)$ . By the definition of chi-square,  $\|z\|^2 = (y - \mu_0)' \Sigma^{-1} (y - \mu_0)$  follows the stated distribution.  $\square$

#### ! Crucial Theorem

This is an important theorem we will use later.

### 5.5 Distributions of Projections of Spherical Normal

**Theorem 5.8** (Distribution of Projections). Let  $V$  be a  $k$ -dimensional subspace of  $\mathcal{R}^n$  with projection matrix  $P_V$ , and let  $y$  be a random vector in  $\mathcal{R}^n$  with mean  $E(y) = \mu$ . Then:

1.  $E(P_V y) = P_V \mu$ .
2. If  $\text{Var}(y) = \sigma^2 I_n$ , then  $\text{Var}(P_V y) = \sigma^2 P_V$  and  $E(\|P_V y\|^2) = \sigma^2 k + \|P_V \mu\|^2$ .
3. If  $y \sim N_n(\mu, \sigma^2 I_n)$ , then  $\frac{1}{\sigma^2} \|P_V y\|^2 = \frac{1}{\sigma^2} y' P_V y \sim \chi^2(k, \frac{1}{2\sigma^2} \|P_V \mu\|^2)$ .

*Proof.*

1. Since the projection operation is linear,  $E(P_V y) = P_V E(y) = P_V \mu$ .
2.  $\text{Var}(P_V y) = P_V \text{Var}(y) P_V^T = P_V \sigma^2 I_n P_V = \sigma^2 P_V$ . The expectation of the squared norm follows from the mean of a quadratic form:  $E(y' P_V y) = \text{tr}(P_V \sigma^2 I) + \mu' P_V \mu = \sigma^2 k + \|P_V \mu\|^2$ .
3. This is a special case of the general quadratic distribution theorem where  $A = \frac{1}{\sigma^2} P_V$  and  $A(\sigma^2 I) = P_V$ , which is idempotent.

□

**Theorem 5.9** (Orthogonal Projections). *Let  $V_1, \dots, V_k$  be mutually orthogonal subspaces with dimensions  $d_i$  and projection matrices  $P_i$ . If  $y \sim N_n(\mu, \sigma^2 I_n)$ , then:*

1. *The projections  $\hat{y}_i = P_i y$  are independent with  $\hat{y}_i \sim N(P_i \mu, \sigma^2 P_i)$ .*
2. *The squared norms  $\|\hat{y}_i\|^2$  are mutually independent.*
3.  *$\frac{1}{\sigma^2} \|\hat{y}_i\|^2 \sim \chi^2(d_i, \frac{1}{2\sigma^2} \|P_i \mu\|^2)$ .*

*Proof.*

1. For  $i \neq j$ ,  $\text{Cov}(P_i y, P_j y) = \sigma^2 P_i P_j = 0$  because orthogonal projection matrices satisfy  $P_i P_j = 0$ . Under normality, zero covariance implies independence.
2. Since  $\hat{y}_i$  are independent, any measurable functions of them, such as their squared norms, are also independent.
3. This follows directly from applying the projection distribution theorem to each independent subspace.

□

### 5.5.1 Independence of Forms

**Theorem 5.10** (Independence Conditions). *Suppose  $y \sim N_n(\mu, \Sigma)$ .*

- **Linear and Quadratic:**  *$By$  and  $y' Ay$  (where  $A$  is symmetric) are independent if and only if  $B\Sigma A = 0$ .*
- **Quadratic and Quadratic:**  *$y' Ay$  and  $y' By$  (where  $A, B$  are symmetric) are independent if and only if  $A\Sigma B = 0$ .*

*Proof.* If  $B\Sigma A = 0$ , the normal vectors  $By$  and  $Ay$  have zero covariance and are independent. Because  $By$  is independent of  $Ay$ , it is also independent of any measurable function of  $Ay$ , specifically  $y' Ay = \|Ay\|^2$  (if  $A$  is idempotent). □

### 5.5.2 Cochran's Theorem

**Theorem 5.11** (Cochran's Result). *Let  $y \sim N_n(\mu, \sigma^2 I)$  and  $y' y = \sum y' A_i y$ . The quadratic forms  $y' A_i y / \sigma^2$  are mutually independent  $\chi^2(r_i, \lambda_i)$  if and only if any one of the following holds:*

- *Each  $A_i$  is idempotent.*
- *$A_i A_j = 0$  for all  $i \neq j$ .*
- *$n = \sum r_i$ .*

## 5.6 Non-central Distributions Derived from Non-central $\chi^2$

We begin by defining two independent Chi-squared random variables that form the building blocks for statistical power analysis.

- **Non-central Component ( $X_1$ ):**  $X_1 \sim \chi^2(df_1, 2\lambda)$ . Here, we use the “natural” definition where  $\lambda$  is **half the sum of squared means**. (Note:  $R$  uses the non-centrality parameter as the full sum of squares, i.e.,  $2\lambda$ .)
- **Central Component ( $X_2$ ):**  $X_2 \sim \chi^2(df_2)$ .  $X_2$  often represents the **Noise Sum of Squares**,  $SSE_1$  of an adequate model, which is assume to follow a central  $\chi^2$ ,

We visualize these components as using the follow diagram.

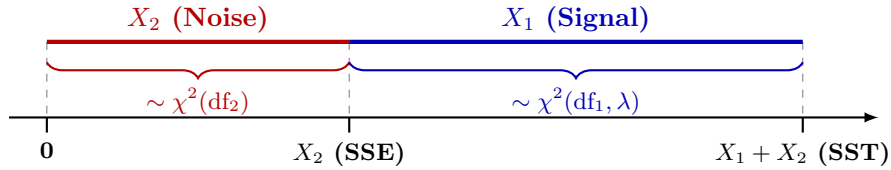


Figure 5.6: A diagram of two independent  $\chi^2$  random variables

### 5.6.1 The Non-central F-distribution $F(df_1, df_2, \lambda)$

**Definition 5.3** (Non-central F). Let  $X_1 \sim \chi^2(df_1, \lambda)$  and  $X_2 \sim \chi^2(df_2)$  be independent. The random variable  $F$  follows a **non-central F-distribution**:

$$F = \frac{X_1/df_1}{X_2/df_2} \sim F(df_1, df_2, \lambda) \quad (5.69)$$

- **Expectation:**
  - **Under  $H_0$  ( $\lambda = 0$ ):** Exact mean is  $\frac{df_2}{df_2 - 2}$  (for  $df_2 > 2$ ).
  - **Under  $H_1$  ( $\lambda \neq 0$ ):** Approximate mean is  $1 + \frac{2\lambda}{df_1}$ .

### 5.6.2 Type I Non-central Beta $Beta_1(df_1/2, df_2/2, \lambda)$

**Definition 5.4** (Type I Non-central Beta). The random variable  $B_I$  follows a **Type I non-central Beta distribution**, defined as the signal’s proportion of the total sum ( $R^2$ ):

$$B_I = \frac{X_1}{X_1 + X_2} \sim Beta_1\left(\frac{df_1}{2}, \frac{df_2}{2}, \lambda\right) \quad (5.70)$$

- **Relationship to F:**  $B_I = \frac{(df_1/df_2)F}{1 + (df_1/df_2)F}$
- **Expectation:**
  - **Under  $H_0$  ( $\lambda = 0$ ):** Exact mean is  $\frac{df_1}{df_1 + df_2}$ .
  - **Under  $H_1$  ( $\lambda \neq 0$ ):** Approximate mean is  $\frac{df_1 + 2\lambda}{df_1 + df_2 + 2\lambda}$ .

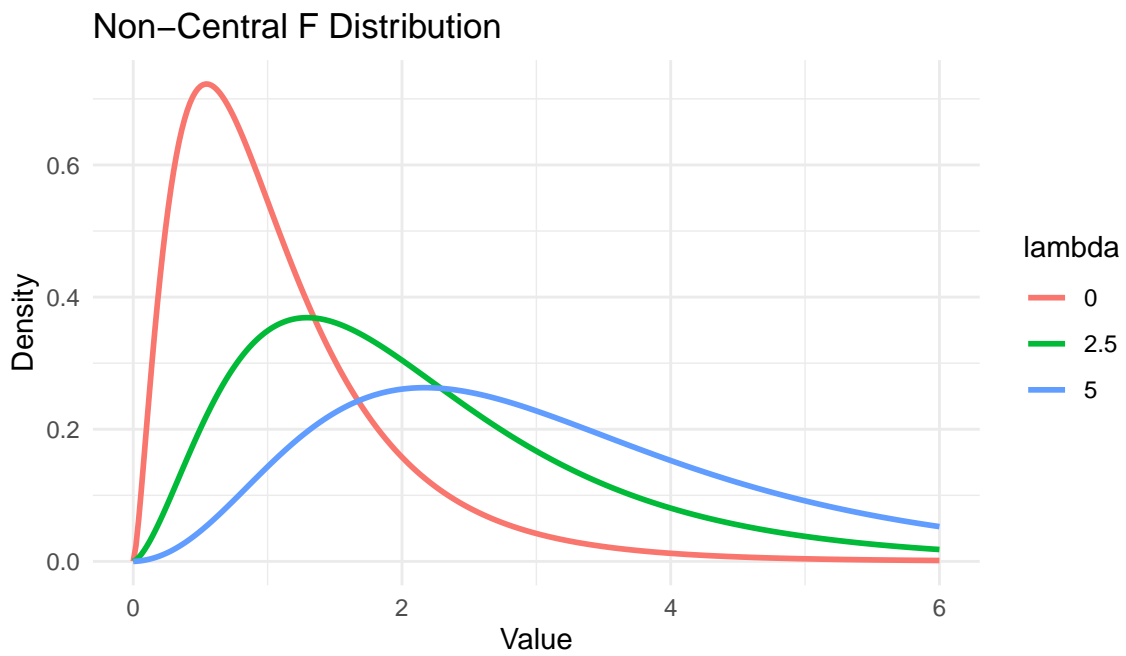


Figure 5.7: Densities of Non-Central F ( $\lambda$  defined as half-sum of squares).

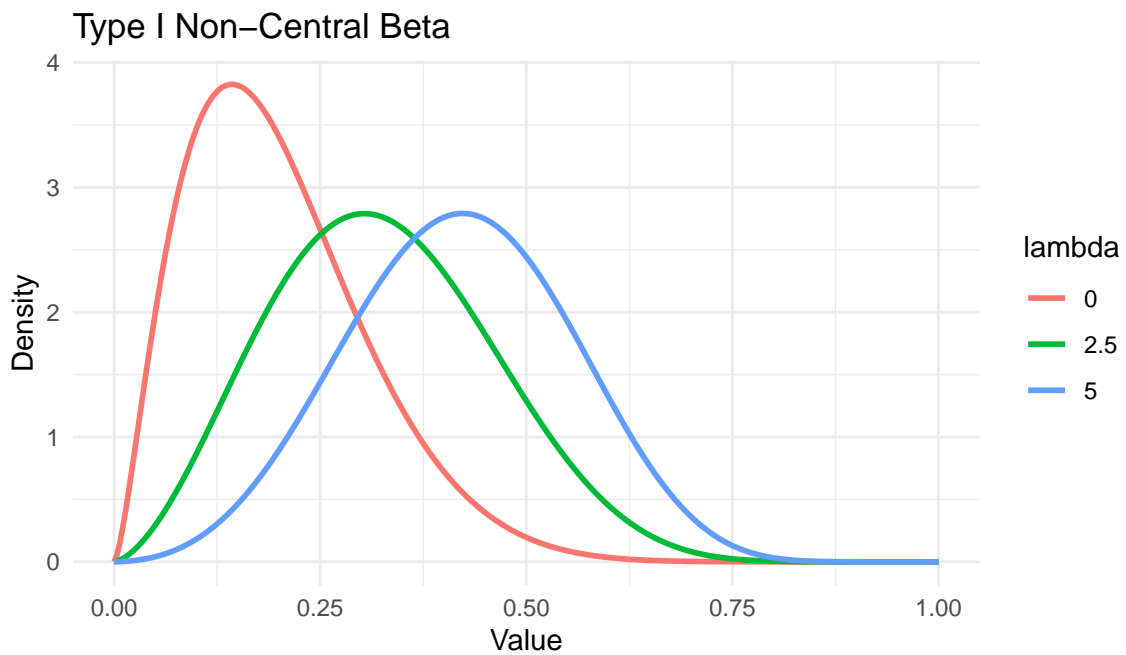


Figure 5.8: Densities of Type I Beta ( $R^2$ ).

### 5.6.3 Type II Non-central Beta $\text{Beta}_2(\text{df}_2/2, \text{df}_1/2, \lambda)$

**Definition 5.5** (Type II Non-central Beta).

$$B_{II} = \frac{X_2}{X_1 + X_2} = 1 - B_I \sim \text{Beta}_2\left(\frac{\text{df}_2}{2}, \frac{\text{df}_1}{2}, \lambda\right) \quad (5.71)$$

- **Relationship to F:**  $B_{II} = \frac{1}{1 + (\text{df}_1/\text{df}_2)F}$
- **Expectation:**
  - **Under  $H_0$  ( $\lambda = 0$ ):** Exact mean is  $\frac{\text{df}_2}{\text{df}_1 + \text{df}_2}$ .
  - **Under  $H_1$  ( $\lambda \neq 0$ ):** Approximate mean is  $\frac{\text{df}_2}{\text{df}_1 + \text{df}_2 + 2\lambda}$ .

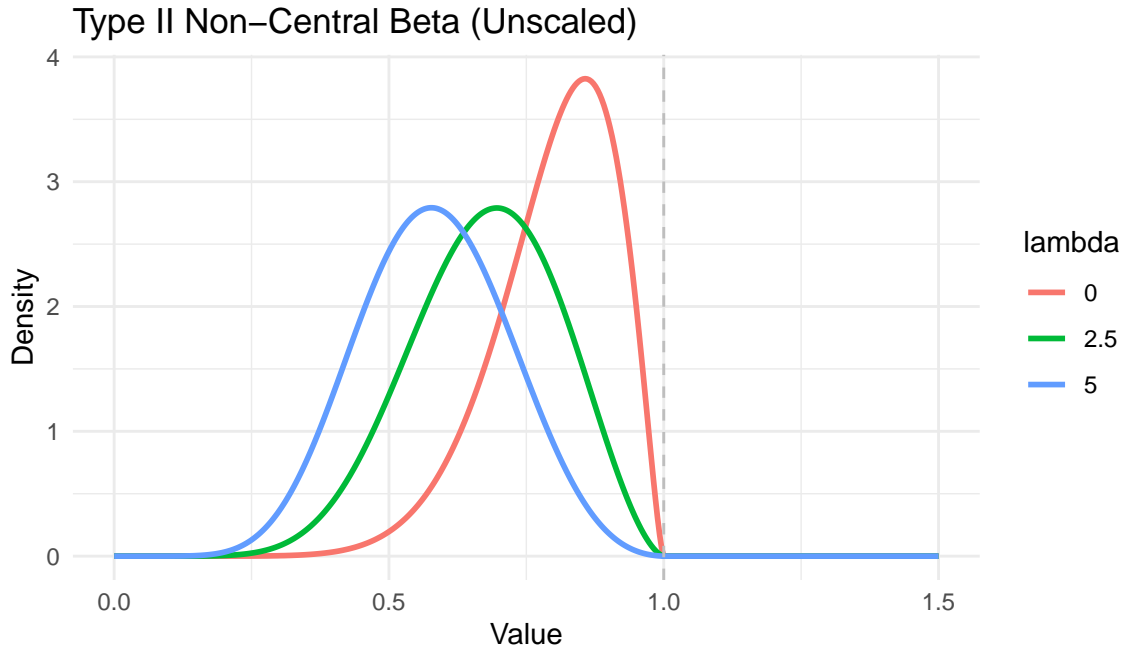


Figure 5.9: Densities of Type II Beta ( $SSE/SST$ ). Support is  $[0, 1]$ .

### 5.6.4 Scaled Type II Beta $\text{Scaled-Beta}_2(\text{df}_2/2, \text{df}_1/2, \lambda)$

**Definition 5.6** (Scaled Type II Beta).

$$S = \frac{X_2/\text{df}_2}{(X_1 + X_2)/(\text{df}_1 + \text{df}_2)} \sim \text{Scaled-Beta}_2 \quad (5.72)$$

- **Relationship to F:**  $S = \frac{\text{df}_1 + \text{df}_2}{\text{df}_2 + \text{df}_1 F}$
- **Expectation:**
  - **Under  $H_0$  ( $\lambda = 0$ ):** Exact mean is 1.
  - **Under  $H_1$  ( $\lambda \neq 0$ ):** Approximate mean is  $\frac{\text{df}_1 + \text{df}_2}{\text{df}_1 + \text{df}_2 + 2\lambda}$ .

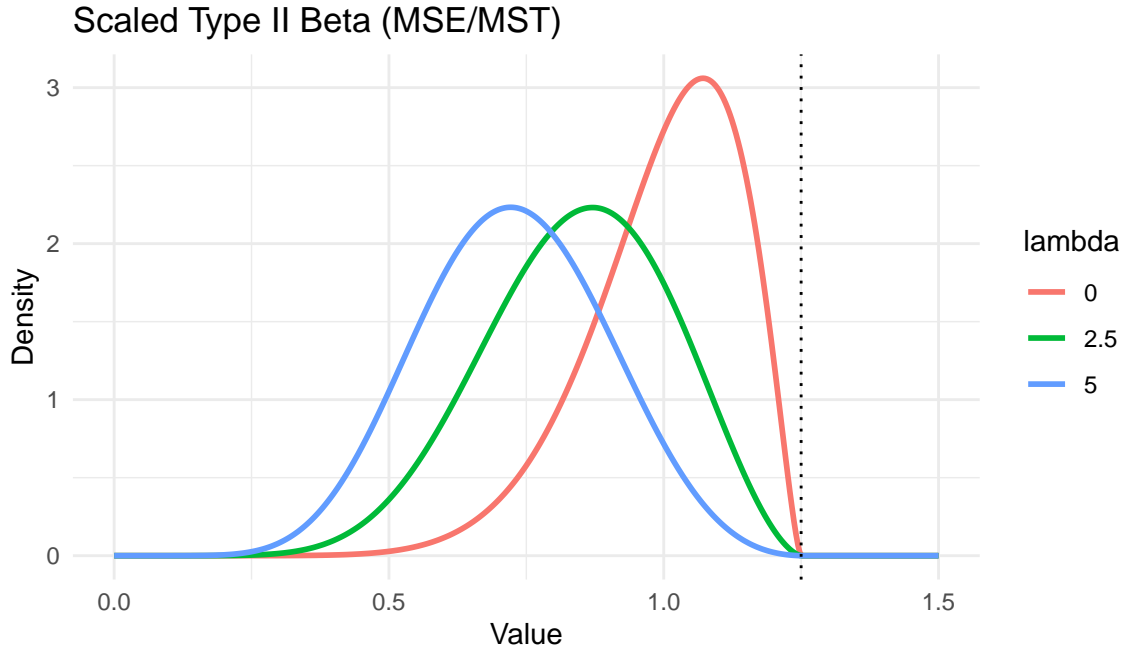


Figure 5.10: Densities of Scaled Type II Beta ( $MSE/MST$ ).

### 5.6.5 The Non-central t-distribution $t(df_2, \delta)$

**Definition 5.7** (Non-central t). Let  $Z \sim N(\delta, 1)$  and  $X_2 \sim \chi^2(df_2)$  be independent. The random variable  $T$  follows a **non-central t-distribution**:

$$T = \frac{Z}{\sqrt{X_2/df_2}} \sim t(df_2, \delta) \quad (5.73)$$

- **Relationship to F:**  $F = T^2$  (when  $df_1 = 1$ ). Note  $\delta^2 = 2\lambda$ .
- **Expectation:**
  - **Under  $H_0$  ( $\delta = 0$ ):** Exact mean is 0.
  - **Under  $H_1$  ( $\delta \neq 0$ ):** Approximate mean is  $\delta$ .

## 5.7 Example: Inference of the Mean of Normal Sample

Consider a random sample  $y \sim N_n(\mu j_n, \sigma^2 I_n)$ . We wish to test:

- $M_1$  (**Full Model**):  $\mu$  is unknown.
- $M_0$  (**Reduced Model**):  $\mu = \mu_0$ .

Let's define the transformed vector  $y^* = y - \mu_0 j_n$ . Note that  $y^* \sim N_n((\mu - \mu_0)j_n, \sigma^2 I_n)$ .

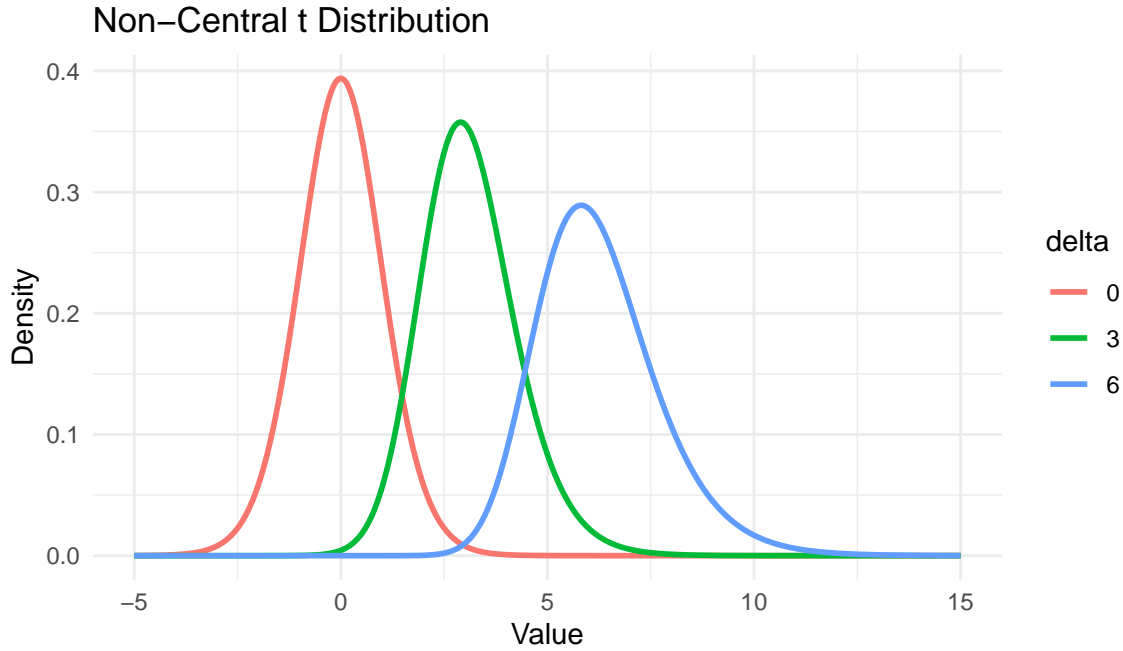


Figure 5.11: Densities of Non-Central t ( $df = 20$ ).

### 5.7.1 Sum of Squares and Their Distributions

We use the projection matrix  $P_{j_n} = \frac{1}{n}j_nj_n'$  and its complement  $(I_n - P_{j_n})$  to partition the transformed vector.

- **Total SSE ( $SSE_0$  for  $M_0$ ):**

$$SSE_0 = \|I_n y^*\|^2 = \sum_{i=1}^n (Y_i - \mu_0)^2 \quad (5.74)$$

This follows a non-central distribution with  $df_{\text{total}} = n$ :

$$\frac{SSE_0}{\sigma^2} \sim \chi^2(n, \lambda) \quad \text{where } \lambda = \frac{n(\mu - \mu_0)^2}{2\sigma^2} \quad (5.75)$$

- **Residual SSE ( $SSE_1$  for  $M_1$ ):**

$$SSE_1 = \|(I_n - P_{j_n})y^*\|^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad (5.76)$$

This captures the random noise (central component) with  $df_2 = n - 1$ :

$$\frac{SSE_1}{\sigma^2} \sim \chi^2(n - 1) \quad (5.77)$$

- **Difference SS ( $SS_{\text{diff}}$ ):**

$$SS_{\text{diff}} = \|P_{j_n} y^*\|^2 = n(\bar{Y} - \mu_0)^2 \quad (5.78)$$

This captures the signal (non-central component) with  $df_1 = 1$ :

$$\frac{SS_{\text{diff}}}{\sigma^2} \sim \chi^2(1, \lambda) \quad (5.79)$$



### 5.7.2 Distributions of Equivalent Statistics

We can construct five equivalent statistics to compare  $M_0$  and  $M_1$ .

- **The t-statistic ( $T$ ):**

$$T = \frac{\bar{Y} - \mu_0}{S/\sqrt{n}} \quad (5.80)$$

- **The F-statistic ( $F$ ):**

$$F = \frac{n(\bar{Y} - \mu_0)^2}{S^2} = T^2 \quad (5.81)$$

- **The Type I Beta statistic ( $B_I$ ):**

$$B_I = \frac{SS_{\text{diff}}}{SSE_0} = \frac{n(\bar{Y} - \mu_0)^2}{\sum (Y_i - \mu_0)^2} \quad (5.82)$$

- **The Type II Beta statistic ( $B_{II}$ ):**

$$B_{II} = \frac{SSE_1}{SSE_0} = \frac{\sum (Y_i - \bar{Y})^2}{\sum (Y_i - \mu_0)^2} = 1 - B_I \quad (5.83)$$

- **The Scaled Type II Beta statistic ( $S_{\text{scaled}}$ ):**

$$S_{\text{scaled}} = \frac{SSE_1/(n-1)}{SSE_0/n} = \left( \frac{n}{n-1} \right) B_{II} \quad (5.84)$$

### 5.7.3 Expectations Under $M_1$ and $M_0$

The table below contrasts the distributions and expected values of these statistics. We assume the sample size  $n$  is large enough for the mean of  $F$  to exist ( $n > 3$ ).

- **Degrees of Freedom:**  $df_1 = 1, df_2 = n - 1$ .
- **Non-centrality:**  $\delta = \frac{\sqrt{n}(\mu - \mu_0)}{\sigma}$  and  $\lambda = \frac{1}{2}\delta^2 = \frac{n(\mu - \mu_0)^2}{2\sigma^2}$ .

Table 5.1: Expected Values of Test Statistics Under Null and Alternative Hypotheses

Statistic	Distribution under $H_1$ ( $\mu \neq \mu_0$ )	Exact Mean under $H_0$ ( $\mu = \mu_0$ )	Approximate Mean under $H_1$
$T$	$t(n-1, \delta)$	0	$\frac{\sqrt{n}(\mu - \mu_0)}{\sigma}$
$F$	$F(1, n-1, \lambda)$	$\frac{n-1}{n-3} \approx 1$	$1 + \frac{\frac{\sigma}{n(\mu - \mu_0)^2}}{\sigma^2}$
$B_I$	$\text{Beta}_1\left(\frac{1}{2}, \frac{n-1}{2}, \lambda\right)$	$\frac{1}{n}$	$\frac{1/n + \frac{(\mu - \mu_0)^2}{\sigma^2}}{1 + \frac{(\mu - \mu_0)^2}{\sigma^2}}$
$B_{II}$	$\text{Beta}_2\left(\frac{n-1}{2}, \frac{1}{2}, \lambda\right)$	$\frac{n-1}{n}$	$\frac{\frac{(n-1)/n}{\sigma^2}}{1 + \frac{(\mu - \mu_0)^2}{\sigma^2}}$
$S_{\text{scaled}}$	$\text{Scaled-Beta}_2\left(\frac{n-1}{2}, \frac{1}{2}, \lambda\right)$	1	$\frac{1}{1 + \frac{(\mu - \mu_0)^2}{\sigma^2}}$

**Key Interpretation:** All statistics are functionally driven by the signal energy. Notably, for  $S_{\text{scaled}}$ , the sample size  $n$  cancels out in the approximate mean. This makes it a direct measure of the ratio between Noise Variance and Total Variance (Noise + Signal) in the population distributions, connected to the Rao-Blackwell decomposition of variances.

# 6 Estimation in Multiple Linear Regression

## 6.1 Linear Models and Least Square Estimator

### 6.1.1 Assumptions in Linear Models

Suppose that on a random sample of  $n$  units (patients, animals, trees, etc.) we observe a response variable  $Y$  and explanatory variables  $X_1, \dots, X_k$ . Our data are then  $(y_i, x_{i1}, \dots, x_{ik})$ ,  $i = 1, \dots, n$ , or in vector/matrix form  $y, x_1, \dots, x_k$  where  $y = (y_1, \dots, y_n)$  and  $x_j = (x_{1j}, \dots, x_{nj})^T$  or  $y, X$  where  $X = (x_1, \dots, x_k)$ .

Either by design or by conditioning on their observed values,  $x_1, \dots, x_k$  are regarded as vectors of known constants. The linear model in its classical form makes the following assumptions:

#### Assumptions on Linear Models

- **A1. (Additive Error)**  $y = \mu + e$  where  $e = (e_1, \dots, e_n)^T$  is an unobserved random vector with  $E(e) = 0$ . This implies that  $\mu = E(y)$  is the unknown mean of  $y$ .
- **A2. (Linearity)**  $\mu = \beta_1 x_1 + \dots + \beta_k x_k = X\beta$  where  $\beta_1, \dots, \beta_k$  are unknown parameters. This assumption says that  $E(y) = \mu \in \text{Col}(X)$  (lies in the column space of  $X$ ); i.e., it is a linear combination of explanatory vectors  $x_1, \dots, x_k$  with coefficients the unknown parameters in  $\beta = (\beta_1, \dots, \beta_k)^T$ . Note that it is linear in  $\beta_1, \dots, \beta_k$ , not necessarily in the  $x$ 's.
- **A3. (Independence)**  $e_1, \dots, e_n$  are independent random variables (and therefore so are  $y_1, \dots, y_n$ ).
- **A4. (Homoscedasticity)**  $e_1, \dots, e_n$  all have the same variance  $\sigma^2$ ; that is,  $\text{Var}(e_1) = \dots = \text{Var}(e_n) = \sigma^2$  which implies  $\text{Var}(y_1) = \dots = \text{Var}(y_n) = \sigma^2$ .
- **A5. (Normality)**  $e \sim N_n(0, \sigma^2 I_n)$ .

### 6.1.2 Matrix Formulation

The model can be written algebraically as:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}, \quad i = 1, \dots, n \quad (6.1)$$

Or in matrix notation:

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} \quad (6.2)$$

This is expressed compactly as:

$$y = X\beta + e \quad (6.3)$$

where  $X$  is the design matrix, and  $e \sim N_n(0, \sigma^2 I)$ . Alternatively:

$$y = \beta_0 j_n + \beta_1 x_1 + \cdots + \beta_k x_k + e \quad (6.4)$$

Taken together, all five assumptions can be stated more succinctly as:

$$y \sim N_n(X\beta, \sigma^2 I) \quad (6.5)$$

with the mean vector  $\mu_y = X\beta \in \text{Col}(X)$ .

#### A Note on Coefficients

The effect of a parameter depends upon what other explanatory variables are present in the model. For example,  $\beta_0$  and  $\beta_1$  in the model:

$$y = \beta_0 j_n + \beta_1 x_1 + \beta_2 x_2 + e \quad (6.6)$$

will typically be different than  $\beta_0^*$  and  $\beta_1^*$  in the model:

$$y = \beta_0^* j_n + \beta_1^* x_1 + e^* \quad (6.7)$$

In this context,  $\beta_0^*$  and  $\beta_1^*$  are the population-projected coefficients of the full model, that is,  $\beta_0^*$  and  $\beta_1^*$  are the parameters that can best approximate the full model.

#### ! Important

We will first consider the case that  $\text{rank}(X) = k + 1$ .

### 6.1.3 Least Squares Estimator of $\beta$ and Fitted Value $\hat{Y}$

**Definition 6.1** (Least Squares Estimator). The **Least Squares Estimator (LSE)** of  $\beta$ , denoted as  $\hat{\beta}$ , is the vector that minimizes the Sum of Squared Errors (SSE), which measures the discrepancy between the observed responses  $y$  and the fitted values  $X\hat{\beta}$ .

$$Q(\beta) = \sum_{i=1}^n (y_i - x_i^T \beta)^2 = (y - X\beta)'(y - X\beta) \quad (6.8)$$

We can derive the closed-form solution for  $\hat{\beta}$  using the geometry of projections discussed in previous chapters.

#### 1. Obtaining $\hat{Y}$

In the linear model  $y = X\beta + e$ , the systematic component (the mean  $E[y]$ ) is constrained to lie in the column space of  $X$ , denoted as  $\text{Col}(X)$ . We seek the vector in  $\text{Col}(X)$  that is “closest” to the observed data vector  $y$ . As

established in the theory of projections, this closest vector is the **orthogonal projection** of  $y$  onto  $\text{Col}(X)$ . Let  $\hat{y}$  denote this fitted value vector. Using the explicit formula for the projection matrix

$$H = X(X'X)^{-1}X', \quad (6.9)$$

we have:

$$\hat{y} = Hy = X(X'X)^{-1}X'y. \quad (6.10)$$

## 2. Obtaining $\hat{\beta}$ by Solving $X\beta = \hat{y}$

Since the fitted vector  $\hat{y}$  is a projection onto  $\text{Col}(X)$ , it must lie entirely within that column space. This guarantees that the linear system for the coefficients  $\hat{\beta}$  is consistent (has an exact solution):

$$X\hat{\beta} = \hat{y} \quad (6.11)$$

To isolate  $\hat{\beta}$ , we pre-multiply both sides by the left pseudo-inverse of  $X$ , which is  $(X'X)^{-1}X'$ :

$$\begin{aligned} (X'X)^{-1}X'(X\hat{\beta}) &= (X'X)^{-1}X'\hat{y} \\ \underbrace{(X'X)^{-1}(X'X)}_I \hat{\beta} &= (X'X)^{-1}X'\hat{y} \end{aligned} \quad (6.12)$$

This gives us the estimator expressed in terms of the fitted values:

$$\boxed{\hat{\beta} = (X'X)^{-1}X'\hat{y}} \quad (6.13)$$

However, we typically calculate the estimator from the observed data  $y$ . Recall that because  $\hat{y}$  is an orthogonal projection, the difference  $y - \hat{y}$  is orthogonal to  $X$ . This implies  $X'\hat{y} = X'y$ . Substituting this into the equation above yields the standard closed-form solution:

$$\boxed{\hat{\beta} = (X'X)^{-1}X'y} \quad (6.14)$$

### 6.1.4 Properties of the Estimator $\hat{\beta}$

**Theorem 6.1** (Unbiasedness of  $\hat{\beta}$ ). *If  $E(y) = X\beta$ , then  $\hat{\beta}$  is an unbiased estimator for  $\beta$ .*

*Proof.*

$$\begin{aligned} E(\hat{\beta}) &= E[(X'X)^{-1}X'y] \\ &= (X'X)^{-1}X'E(y) \quad [\text{using linearity of expectation}] \\ &= (X'X)^{-1}X'X\beta \\ &= \beta \end{aligned} \quad (6.15)$$

□

**Theorem 6.2** (Variance of  $\hat{\beta}$ ). If  $\text{Var}(y) = \sigma^2 I$ , the covariance matrix for  $\hat{\beta}$  is given by  $\sigma^2 (X'X)^{-1}$ .

*Proof.*

$$\begin{aligned}
 \text{Var}(\hat{\beta}) &= \text{Var}[(X'X)^{-1}X'y] \\
 &= (X'X)^{-1}X'\text{Var}(y)[(X'X)^{-1}X']' \quad [\text{using } \text{Var}(Ay) = A\text{Var}(y)A'] \\
 &= (X'X)^{-1}X'(\sigma^2 I)X(X'X)^{-1} \\
 &= \sigma^2 (X'X)^{-1}X'X(X'X)^{-1} \\
 &= \sigma^2 (X'X)^{-1}
 \end{aligned} \tag{6.16}$$

□

**Note:** These theorems require no assumption of normality.

## 6.2 Best Linear Unbiased Estimator (BLUE)

**Theorem 6.3** (Gauss-Markov Theorem). If  $E(y) = X\beta$  and  $\text{Var}(y) = \sigma^2 I$ , the least-squares estimators  $\hat{\beta}_j, j = 0, 1, \dots, k$  have minimum variance among all linear unbiased estimators.

*Proof.* We consider a linear estimator  $Ay$  of  $\beta$  and seek the matrix  $A$  for which  $Ay$  is a minimum variance unbiased estimator.

**1. Unbiasedness Condition:** In order for  $Ay$  to be an unbiased estimator of  $\beta$ , we must have  $E(Ay) = \beta$ . Using the assumption  $E(y) = X\beta$ , this is expressed as:

$$E(Ay) = AE(y) = AX\beta = \beta \tag{6.17}$$

which implies the condition  $AX = I_{k+1}$  since the relationship must hold for any  $\beta$ .

**2. Minimizing Variance:** The covariance matrix for the estimator  $Ay$  is:

$$\text{Var}(Ay) = A\text{Var}(y)A' = A(\sigma^2 I)A' = \sigma^2 AA' \tag{6.18}$$

We need to choose  $A$  (subject to  $AX = I$ ) so that the diagonal elements of  $AA'$  are minimized.

To relate  $Ay$  to  $\hat{\beta} = (X'X)^{-1}X'y$ , we define  $\hat{A} = (X'X)^{-1}X'$  and write  $A = (A - \hat{A}) + \hat{A}$ . Then:

$$AA' = [(A - \hat{A}) + \hat{A}][(A - \hat{A}) + \hat{A}]' \tag{6.19}$$

Expanding this, the cross terms vanish because  $(A - \hat{A})\hat{A}' = A\hat{A}' - \hat{A}\hat{A}'$ . Note that  $\hat{A}\hat{A}' = (X'X)^{-1}X'X(X'X)^{-1} = (X'X)^{-1}$ . Also,  $A\hat{A}' = AX(X'X)^{-1} = I(X'X)^{-1} = (X'X)^{-1}$  (since  $AX = I$ ). Thus,  $(A - \hat{A})\hat{A}' = 0$ .

The expansion simplifies to:

$$AA' = (A - \hat{A})(A - \hat{A})' + \hat{A}\hat{A}' \tag{6.20}$$

The matrix  $(A - \hat{A})(A - \hat{A})'$  is positive semidefinite, meaning its diagonal elements are non-negative. To minimize the diagonal of  $AA'$ , we must set  $A - \hat{A} = 0$ , which implies  $A = \hat{A}$ .

Thus, the minimum variance estimator is:

$$Ay = (X'X)^{-1}X'y = \hat{\beta} \tag{6.21}$$

□

### 6.2.1 Notes on Gauss-markov

1. **Distributional Generality:** The remarkable feature of the Gauss-Markov theorem is that it holds for *any* distribution of  $y$ ; normality is not required. The only assumptions used are linearity ( $E(y) = X\beta$ ) and homoscedasticity ( $\text{Var}(y) = \sigma^2 I$ ).
2. **Extension to All Linear Combinations:** The theorem extends beyond just the parameter vector  $\beta$  to any linear combination of the parameters.

**Corollary 6.1** (BLUE for All Linear Combinations). *If  $E(y) = X\beta$  and  $\text{Var}(y) = \sigma^2 I$ , the best linear unbiased estimator of the scalar  $a'\beta$  is  $a'\hat{\beta}$ , where  $\hat{\beta}$  is the least-squares estimator.*

*Proof.* Let  $\tilde{\beta} = Ay$  be any other linear unbiased estimator of  $\beta$ . The variance of the linear combination  $a'\tilde{\beta}$  is:

$$\frac{1}{\sigma^2} \text{Var}(a'\tilde{\beta}) = \frac{1}{\sigma^2} \text{Var}(a' Ay) = a' A A' a \quad (6.22)$$

From the proof of the Gauss-Markov theorem, we established that  $AA' = (A - \hat{A})(A - \hat{A})' + (X'X)^{-1}$  where  $\hat{A} = (X'X)^{-1}X'$ . Substituting this into the variance equation:

$$a' A A' a = a' (A - \hat{A})(A - \hat{A})' a + a' (X'X)^{-1} a \quad (6.23)$$

The term  $a' (A - \hat{A})(A - \hat{A})' a$  is a quadratic form with a positive semidefinite matrix, so it is always non-negative. Therefore:

$$a' A A' a \geq a' (X'X)^{-1} a = \frac{1}{\sigma^2} \text{Var}(a'\hat{\beta}) \quad (6.24)$$

The variance is minimized when  $A = \hat{A}$  (specifically when the first term is zero), proving that  $a'\hat{\beta}$  has the minimum variance among all linear unbiased estimators.  $\square$

3. **Scaling Invariance:** The predictions made by the model are invariant to the scaling of the explanatory variables.

**Theorem 6.4** (Scaling Explanatory Variables). *If  $x = (1, x_1, \dots, x_k)'$  and  $z = (1, c_1 x_1, \dots, c_k x_k)'$ , then the fitted values are identical:  $\hat{y} = \hat{\beta}' x = \hat{\beta}'_z z$ .*

*Proof.* Let  $D = \text{diag}(1, c_1, \dots, c_k)$  such that the design matrix is transformed to  $Z = XD$ . The LSE for the transformed data is:

$$\begin{aligned} \hat{\beta}_z &= (Z'Z)^{-1} Z'y = [(XD)'(XD)]^{-1} (XD)'y \\ &= D^{-1} (X'X)^{-1} (D')^{-1} D' X'y \\ &= D^{-1} (X'X)^{-1} X'y = D^{-1} \hat{\beta} \end{aligned} \quad (6.25)$$

. Then, the prediction is:

$$\hat{\beta}'_z z = (D^{-1} \hat{\beta})' (Dx) = \hat{\beta}' (D^{-1})' Dx = \hat{\beta}' x \quad (6.26)$$

$\square$

### 6.2.1.1 Limitations: Restriction to Unbiased Estimators

It is crucial to recognize that the Gauss-Markov theorem only guarantees optimality within the class of **linear** and **unbiased** estimators.

- **Assumption Sensitivity:** If the assumptions of linearity ( $E(y) = X\beta$ ) and homoscedasticity ( $\text{Var}(y) = \sigma^2 I$ ) do not hold,  $\hat{\beta}$  may be biased or may have a larger variance than other estimators.
- **Unbiasedness Constraint:** The theorem does not compare  $\hat{\beta}$  to biased estimators. It is possible for a biased estimator (e.g., shrinkage estimators) to have a smaller Mean Squared Error (MSE) than the BLUE by accepting some bias to significantly reduce variance. The LSE is only “best” (minimum variance) among those estimators that satisfy the unbiasedness constraint.

## 6.3 Estimator of Error Variance

We estimate  $\sigma^2$  by the residual mean square:

**Definition 6.2** (Residual Variance Estimator).

$$s^2 = \frac{1}{n - k - 1} \sum_{i=1}^n (y_i - x_i' \hat{\beta})^2 = \frac{\text{SSE}}{n - k - 1} \quad (6.27)$$

where  $\text{SSE} = (y - X\hat{\beta})'(y - X\hat{\beta})$ .

Alternatively, SSE can be written as:

$$\text{SSE} = y'y - \hat{\beta}' X'y \quad (6.28)$$

This is often useful for computation ( $y'y$  is the total sum of squares of the raw data).

### 6.3.1 Unbiasedness of $s^2$

**Theorem 6.5** (Unbiasedness of s-squared). *If  $s^2$  is defined as above, and if  $E(y) = X\beta$  and  $\text{Var}(y) = \sigma^2 I$ , then  $E(s^2) = \sigma^2$ .*

*Proof.* We use the Hat Matrix  $H = X(X'X)^{-1}X'$ , which projects  $y$  onto  $\text{Col}(X)$ . Thus,  $\hat{y} = Hy$ . The residuals are  $y - \hat{y} = (I - H)y$ . The Sum of Squared Errors is:

$$\text{SSE} = \|(I - H)y\|^2 = y'(I - H)'(I - H)y \quad (6.29)$$

Since  $H$  is symmetric and idempotent,  $(I - H)$  is also symmetric and idempotent. Thus:

$$\text{SSE} = y'(I - H)y \quad (6.30)$$



To find the expectation, we use the trace trick for quadratic forms:  $E[y' Ay] = \text{tr}(A \text{Var}(y)) + E[y]' A E[y]$ .

$$\begin{aligned} E(\text{SSE}) &= E[y'(I - H)y] \\ &= \text{tr}((I - H)\sigma^2 I) + (X\beta)'(I - H)(X\beta) \\ &= \sigma^2 \text{tr}(I - H) + \beta' X'(I - H)X\beta \end{aligned} \quad (6.31)$$

**Trace Term:**  $\text{tr}(I_n - H) = \text{tr}(I_n) - \text{tr}(H) = n - (k + 1)$ , since  $\text{tr}(H) = \text{tr}(X(X'X)^{-1}X') = \text{tr}((X'X)^{-1}X'X) = \text{tr}(I_{k+1}) = k + 1$ .

**Non-centrality Term:** Since  $HX = X$ , we have  $(I - H)X = 0$ . Therefore, the second term vanishes:  $\beta' X'(I - H)X\beta = 0$ .

Combining these:

$$E(\text{SSE}) = \sigma^2(n - k - 1) \quad (6.32)$$

Dividing by the degrees of freedom  $(n - k - 1)$ , we get  $E(s^2) = \sigma^2$ .  $\square$

## 6.4 Distributions Under Normality

If we add Assumption A5 ( $y \sim N_n(X\beta, \sigma^2 I)$ ), we can derive the exact sampling distributions.

**Corollary 6.2** (Estimated Covariance of Beta). *An unbiased estimator of  $\text{Cov}(\hat{\beta})$  is given by:*

$$\widehat{\text{Cov}}(\hat{\beta}) = s^2(X'X)^{-1} \quad (6.33)$$

**Theorem 6.6** (Sampling Distributions). *Under assumptions A1-A5:*

1.  $\hat{\beta} \sim N_{k+1}(\beta, \sigma^2(X'X)^{-1})$ .
2.  $(n - k - 1)s^2/\sigma^2 \sim \chi^2(n - k - 1)$ .
3.  $\hat{\beta}$  and  $s^2$  are independent.

*Proof. Part (i):* Since  $\hat{\beta} = (X'X)^{-1}X'y$  is a linear transformation of the normal vector  $y$ , it is also normally distributed. We already established its mean and variance in Theorem 6.1 and Theorem 6.2.

**Part (ii):** We showed  $\text{SSE} = y'(I - H)y$ . Since  $(I - H)$  is idempotent with rank  $n - k - 1$ , and  $(I - H)X\beta = 0$ , by the theory of quadratic forms in normal variables,  $\text{SSE}/\sigma^2 \sim \chi^2(n - k - 1)$ .

**Part (iii):**  $\hat{\beta}$  depends on  $Hy$  (or  $X'y$ ), while  $s^2$  depends on  $(I - H)y$ . Since  $H(I - H) = H - H^2 = 0$ , the linear forms defining the estimator and the residuals are orthogonal. For normal vectors, zero covariance implies independence.  $\square$

## 6.5 Maximum Likelihood Estimator (MLE)

**Theorem 6.7** (MLE for Linear Regression). *If  $y \sim N_n(X\beta, \sigma^2 I)$ , the Maximum Likelihood Estimators are:*

$$\hat{\beta}_{MLE} = (X'X)^{-1}X'y \quad (6.34)$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n}(y - X\hat{\beta})'(y - X\hat{\beta}) = \frac{SSE}{n} \quad (6.35)$$

*Proof.* The log-likelihood function is:

$$\ln L(\beta, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} (y - X\beta)'(y - X\beta) \quad (6.36)$$

Maximizing this with respect to  $\beta$  is equivalent to minimizing the quadratic term  $(y - X\beta)'(y - X\beta)$ , which yields the Least Squares Estimator. Differentiating with respect to  $\sigma^2$  and setting to zero yields  $\hat{\sigma}^2 = SSE/n$ .  $\square$

**Note:** The MLE for  $\sigma^2$  is biased (denominator  $n$ ), whereas  $s^2$  is unbiased (denominator  $n - k - 1$ ).

## 6.6 Linear Models in Centered Form

The regression model can be written in a centered form by subtracting the means of the explanatory variables:

$$y_i = \alpha + \beta_1(x_{i1} - \bar{x}_1) + \beta_2(x_{i2} - \bar{x}_2) + \cdots + \beta_k(x_{ik} - \bar{x}_k) + e_i \quad (6.37)$$

for  $i = 1, \dots, n$ , where the intercept term is adjusted:

$$\alpha = \beta_0 + \beta_1\bar{x}_1 + \beta_2\bar{x}_2 + \cdots + \beta_k\bar{x}_k \quad (6.38)$$

and  $\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$ .

### 6.6.1 Matrix Formulation

In matrix form, the equivalence between the original model and the centered model is:

$$y = X\beta + e = (j_n, X_c) \begin{pmatrix} \alpha \\ \beta_1 \end{pmatrix} + e \quad (6.39)$$

where  $\beta_1 = (\beta_1, \dots, \beta_k)^T$  represents the slope coefficients, and  $X_c$  is the centered design matrix:

$$X_c = (I - P_{j_n})X_1 \quad (6.40)$$

Here,  $X_1$  consists of the original columns of  $X$  excluding the intercept column.

To see the structure of  $X_c$ , we first calculate the projection of the data onto the intercept space,  $P_{j_n} X_1$ :

$$\begin{aligned}
P_{j_n} X_1 &= \frac{1}{n} j_n j_n' X_1 \\
&= \begin{pmatrix} 1/n & 1/n & \cdots & 1/n \\ 1/n & 1/n & \cdots & 1/n \\ \vdots & \vdots & \ddots & \vdots \\ 1/n & 1/n & \cdots & 1/n \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} \\
&= \begin{pmatrix} \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_k \\ \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_k \\ \vdots & \vdots & \ddots & \vdots \\ \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_k \end{pmatrix}
\end{aligned} \tag{6.41}$$

This results in a matrix where every row is the vector of column means. Subtracting this from  $X_1$  gives  $X_c$ :

$$\begin{aligned}
X_c &= X_1 - P_{j_n} X_1 \\
&= \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} - \begin{pmatrix} \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_k \\ \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_k \\ \vdots & \vdots & \ddots & \vdots \\ \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_k \end{pmatrix} \\
&= \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1k} - \bar{x}_k \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \cdots & x_{2k} - \bar{x}_k \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \cdots & x_{nk} - \bar{x}_k \end{pmatrix}
\end{aligned} \tag{6.42}$$

### 6.6.2 Estimation in Centered Form

Because the column space of the intercept  $j_n$  is orthogonal to the columns of  $X_c$  (since columns of  $X_c$  sum to zero), the cross-product matrix becomes block diagonal:

$$\begin{pmatrix} j_n' \\ X_c' \end{pmatrix} (j_n, X_c) = \begin{pmatrix} j_n' j_n & j_n' X_c \\ X_c' j_n & X_c' X_c \end{pmatrix} = \begin{pmatrix} n & 0 \\ 0 & X_c' X_c \end{pmatrix} \tag{6.43}$$

**Theorem 6.8** (Centered Estimators). *The least squares estimators for the centered parameters are:*

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} n & 0 \\ 0 & X_c' X_c \end{pmatrix}^{-1} \begin{pmatrix} j_n' y \\ X_c' y \end{pmatrix} = \begin{pmatrix} \bar{y} \\ (X_c' X_c)^{-1} X_c' y \end{pmatrix} \tag{6.44}$$

Thus:

1.  $\hat{\alpha} = \bar{y}$  (The sample mean of  $y$ ).
2.  $\hat{\beta}_1 = S_{xx}^{-1} S_{xy}$ , using the sample covariance notations.

Recovering the original intercept:

$$\hat{\beta}_0 = \hat{\alpha} - \hat{\beta}_1 \bar{x}_1 - \cdots - \hat{\beta}_k \bar{x}_k = \bar{y} - \hat{\beta}_1' \bar{x} \tag{6.45}$$

## 6.7 Sum of Squares Decomposition

We partition the total variation based on the orthogonal subspaces.

**Definition 6.3** (Sum of Squares Components). The total variation is decomposed as  $SST = SSR + SSE$ .

1. **Total Sum of Squares (SST):** The squared length of the centered response vector.

$$SST = \|y - \bar{y}j_n\|^2 = \|(I - P_{j_n})y\|^2 \quad (6.46)$$

2. **Regression Sum of Squares (SSR):** The variation explained by the regressors  $X_c$ .

$$SSR = \|\hat{y} - \bar{y}j_n\|^2 = \|P_{X_c}y\|^2 = \hat{\beta}_1' X_c' X_c \hat{\beta}_1 \quad (6.47)$$

3. **Sum of Squared Errors (SSE):** The residual variation.

$$SSE = \|y - \hat{y}\|^2 = \|(I - H)y\|^2 \quad (6.48)$$

### 6.7.1 3D Visualization of Decomposition of $y$

We partition the total variation in  $y$  based on the orthogonal subspaces.

1. **Space of the Mean:**  $L(j_n)$ , spanned by the intercept vector  $j_n$ .
2. **Space of the Regressors:**  $L(X_c)$ , spanned by the centered predictors  $X_c$ .
3. **Error Space:**  $\text{Col}(X)^\perp$ , orthogonal to the model space.

The vector  $y$  can be decomposed into three orthogonal components:

$$y = \bar{y}j_n + P_{X_c}y + (y - \hat{y}) \quad (6.49)$$

Visually, this corresponds to projecting the vector  $y$  onto three orthogonal axes.

#### Interactive Visualization:

We generate a cloud of 100 observations of  $y$  from  $N(\mu, \sigma = 1)$  where  $\mu = (5, 5, 0)$ . The projections onto the Model Plane ( $z = 0$ ) are highlighted in **red**, and the projections onto the error axis ( $z$ ) are in **yellow**.

### 6.7.1.1 Effect Exists (signal)

Scenario A: Effect Exists

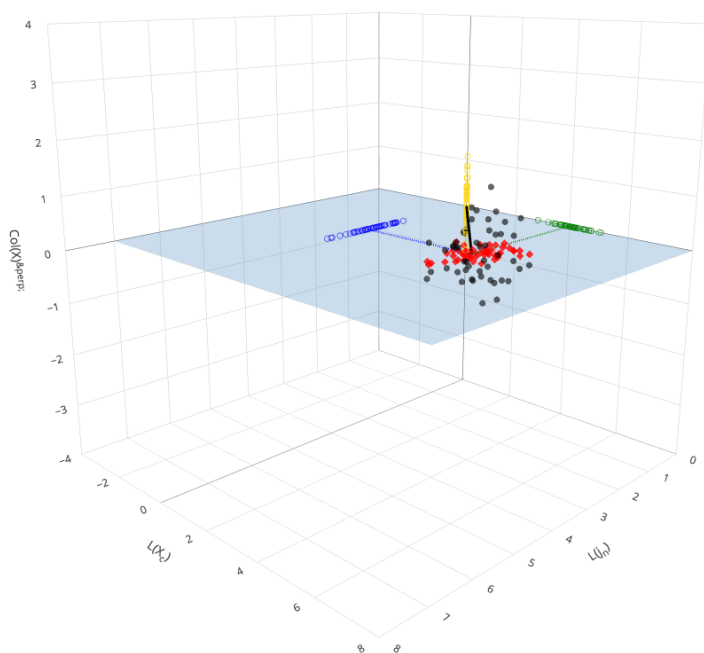


Figure 6.1: Scenario 1: Significant regression effect ( $\beta_1 \neq 0$ ). The mean vector projects significantly onto the predictor space.

### 6.7.1.2 No Effect (noise)

Scenario B: No Effect

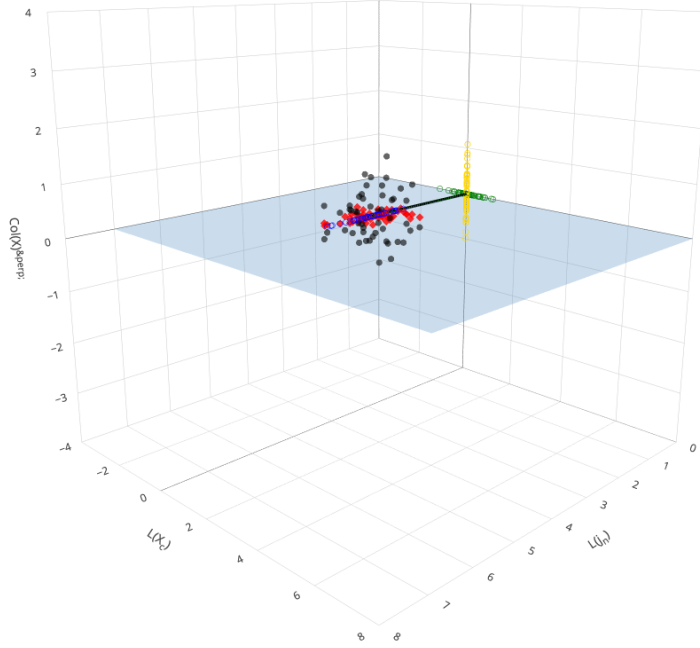


Figure 6.2: Scenario 2: No regression effect ( $\beta_1 = 0$ ). The mean vector lies purely on the intercept axis.

### 6.7.2 A Diagram to Show Decomposition of Sum of Squares

The decomposition of the total variation is visualized below. The total deviation (Orange) is the vector sum of the regression deviation (Green) and the residual error (Red).

### 6.7.3 Distribution of Sum of Squares

We apply the general theory of projections to the specific components defined in Definition 6.3.

**Theorem 6.9** (Distribution of Sum of Squares). *Let  $y \sim N(\mu, \sigma^2 I_n)$ , where  $\mu \in \text{Col}(X)$ . Consider the decomposition defined by the projection matrices  $P_{X_c}$  and  $M = I - H$ .*

- **Independence:** *The quadratic forms SSR and SSE are statistically independent because the subspaces  $L(X_c)$  and  $\text{Col}(X)^\perp$  are orthogonal.*

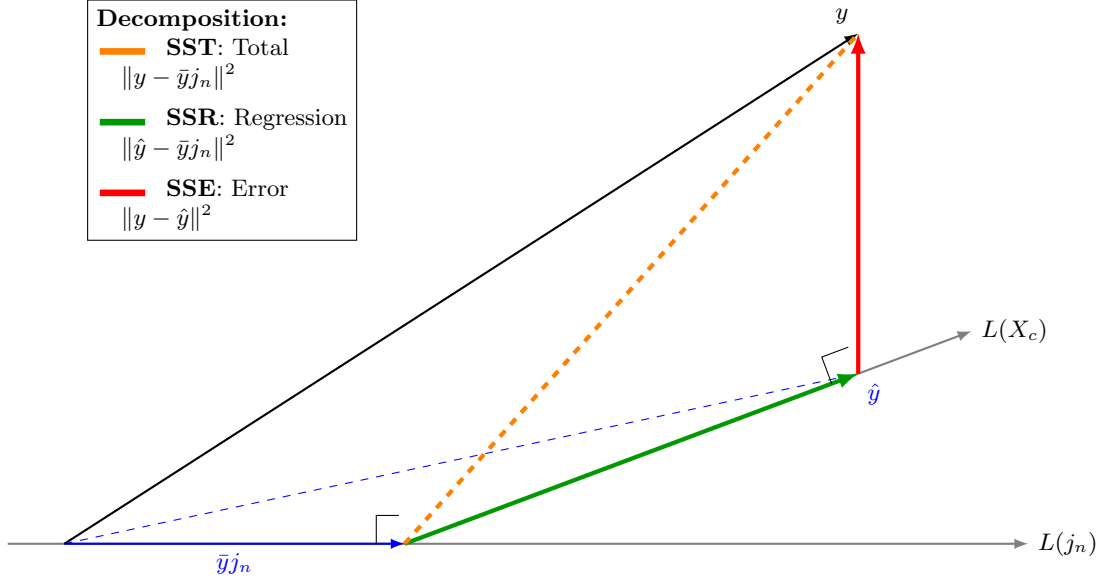


Figure 6.3: Geometric Decomposition:  $SST = SSR + SSE$

- **Distribution of SSE:** The scaled sum of squared errors follows a central Chi-squared distribution:

$$\frac{SSE}{\sigma^2} = \frac{\|(I - H)y\|^2}{\sigma^2} \sim \chi^2(n - k - 1) \quad (6.50)$$

**Mean:**

$$E[SSE] = \sigma^2(n - k - 1) \quad (6.51)$$

- **Distribution of SSR:** The scaled regression sum of squares follows a **non-central** Chi-squared distribution:

$$\frac{SSR}{\sigma^2} = \frac{\|P_{X_c}y\|^2}{\sigma^2} \sim \chi^2(k, \lambda) \quad (6.52)$$

**Mean:**

$$E[SSR] = \sigma^2 k + \|P_{X_c}\mu\|^2 \quad (6.53)$$

**Non-centrality Parameter ( $\lambda$ ):**

$$\lambda = \frac{1}{2\sigma^2} \|P_{X_c}\mu\|^2 \quad (6.54)$$

where

$$\|P_{X_c}\mu\|^2 = \|X_c\beta_1\|^2 = (X_c\beta_1)'(X_c\beta_1) = \beta_1'X_c'X_c\beta_1 \quad (6.55)$$

*Proof.* We apply Theorem 5.8 to the specific projection matrices identified in the definitions.

- **For SSE (Error Space):** SSE is defined by the projection matrix  $P_V = I - H$ .
  - **Dimension:** The rank of  $(I - H)$  is  $n - \text{rank}(X) = n - (k + 1) = n - k - 1$ .
  - **Non-centrality:** Since  $\mu \in \text{Col}(X)$ , the projection onto the orthogonal complement is zero:  $\|(I - H)\mu\|^2 = 0$ . Thus,  $\lambda = 0$ .

- **Expectation:** Using Part 2 of Theorem 5.8 ( $E(\|P_V y\|^2) = \sigma^2 \text{rank}(P_V) + \|P_V \mu\|^2$ ):

$$E[\text{SSE}] = \sigma^2(n - k - 1) + 0 = \sigma^2(n - k - 1) \quad (6.56)$$

- **For SSR (Regression Space):** SSR is defined by the projection matrix  $P_V = P_{X_c}$ .

- **Dimension:** The rank of  $P_{X_c}$  is  $(k + 1) - 1 = k$ .
- **Non-centrality:** The projection of  $\mu$  onto  $L(X_c)$  is  $P_{X_c} \mu$ .

$$\lambda = \frac{1}{2\sigma^2} \|P_{X_c} \mu\|^2 \quad (6.57)$$

- **Expectation:** Using Part 2 of Theorem 5.8:

$$E[\text{SSR}] = \sigma^2 k + \|P_{X_c} \mu\|^2 \quad (6.58)$$

This shows that while  $E[\text{SSE}]$  depends only on the noise variance and sample size,  $E[\text{SSR}]$  is inflated by the magnitude of the true regression signal  $\|P_{X_c} \mu\|^2$ .

□

## 6.8 F-test for Testing Overall Regression Effect

We wish to test whether the regression model provides any explanatory power beyond the simple intercept-only model.

### Hypotheses:

- **Null Hypothesis ( $H_0$ ):**  $\beta_1 = \beta_2 = \dots = \beta_k = 0$  (No regression effect). This implies  $\mu \in \text{span}(j_n)$  and the true signal variance  $\|X_c \beta_1\|^2 = 0$ .
- **Alternative Hypothesis ( $H_1$ ):** At least one  $\beta_j \neq 0$ .

### The F-statistic

We construct the test statistic using the ratio of the Mean Squares defined previously:

$$F = \frac{\text{MSR}}{\text{MSE}} = \frac{\text{SSR}/k}{\text{SSE}/(n - k - 1)} \quad (6.59)$$



## Understanding $F$ via Expectations

The logic of the F-test is transparent when we examine the expected values of the numerator and denominator:

$$\begin{aligned} E[\text{MSE}] &= \sigma^2 \\ E[\text{MSR}] &= \sigma^2 + \frac{\|X_c \beta_1\|^2}{k} \end{aligned} \quad (6.60)$$

- **If  $H_0$  is true:** The signal term is zero. Both Mean Squares estimate  $\sigma^2$  unbiasedly. We expect  $F \approx 1$ .
- **If  $H_1$  is true:** The numerator includes the positive term  $\frac{\|X_c \beta_1\|^2}{k}$ . We expect  $F > 1$ .

Therefore, we reject  $H_0$  for sufficiently large values of  $F$ . Specifically, we reject at level  $\alpha$  if  $F_{obs} > F_\alpha(k, n - k - 1)$ .

### 6.8.1 Distributional Theory

To derive the exact sampling distribution, we rely on the independence of the sums of squares (from Theorem 6.9) and the definition of the non-central F-distribution given in Definition 5.3.

**Theorem 6.10** (Distribution of Regression F-Statistic). *Under the assumption of normality, the regression F-statistic follows a **non-central F-distribution**:*

$$F \sim F(k, n - k - 1, \lambda) \quad (6.61)$$

The non-centrality parameter  $\lambda$  is determined by the ratio of the signal sum of squares to the error variance:

$$\lambda = \frac{\|X_c \beta_1\|^2}{2\sigma^2} \quad (6.62)$$

**Special Cases:**

1. **Under  $H_1$  (Signal exists):**  $\lambda > 0$ , so  $F$  follows the non-central distribution.
2. **Under  $H_0$  (No signal):**  $\beta_1 = 0 \implies \lambda = 0$ . The distribution collapses to the **central F-distribution**:

$$F \sim F(k, n - k - 1) \quad (6.63)$$

*Proof.* We identify the components from Definition 5.3:

1. **Numerator ( $X_1$ ):** Let  $X_1 = \text{SSR}/\sigma^2$ . From Theorem 6.9,  $X_1 \sim \chi^2(k, 2\lambda)$ .
2. **Denominator ( $X_2$ ):** Let  $X_2 = \text{SSE}/\sigma^2$ . From Theorem 6.9,  $X_2 \sim \chi^2(n - k - 1)$ .
3. **Independence:**  $X_1$  and  $X_2$  are independent.

Substituting these into the F-statistic:

$$F = \frac{\text{MSR}}{\text{MSE}} = \frac{(\text{SSR}/\sigma^2)/k}{(\text{SSE}/\sigma^2)/(n - k - 1)} = \frac{X_1/k}{X_2/(n - k - 1)} \quad (6.64)$$

By definition Definition 5.3, this ratio follows  $F(k, n - k - 1, \lambda)$ . □

### 6.8.2 Visualization of the Rejection Region

The following plot illustrates the central F-distribution (valid under  $H_0$ ) for  $k = 3$  predictors and  $n = 20$  observations ( $df_1 = 3, df_2 = 16$ ). An observed statistic of  $F = 2$  is marked, with the p-value represented by the shaded tail area.

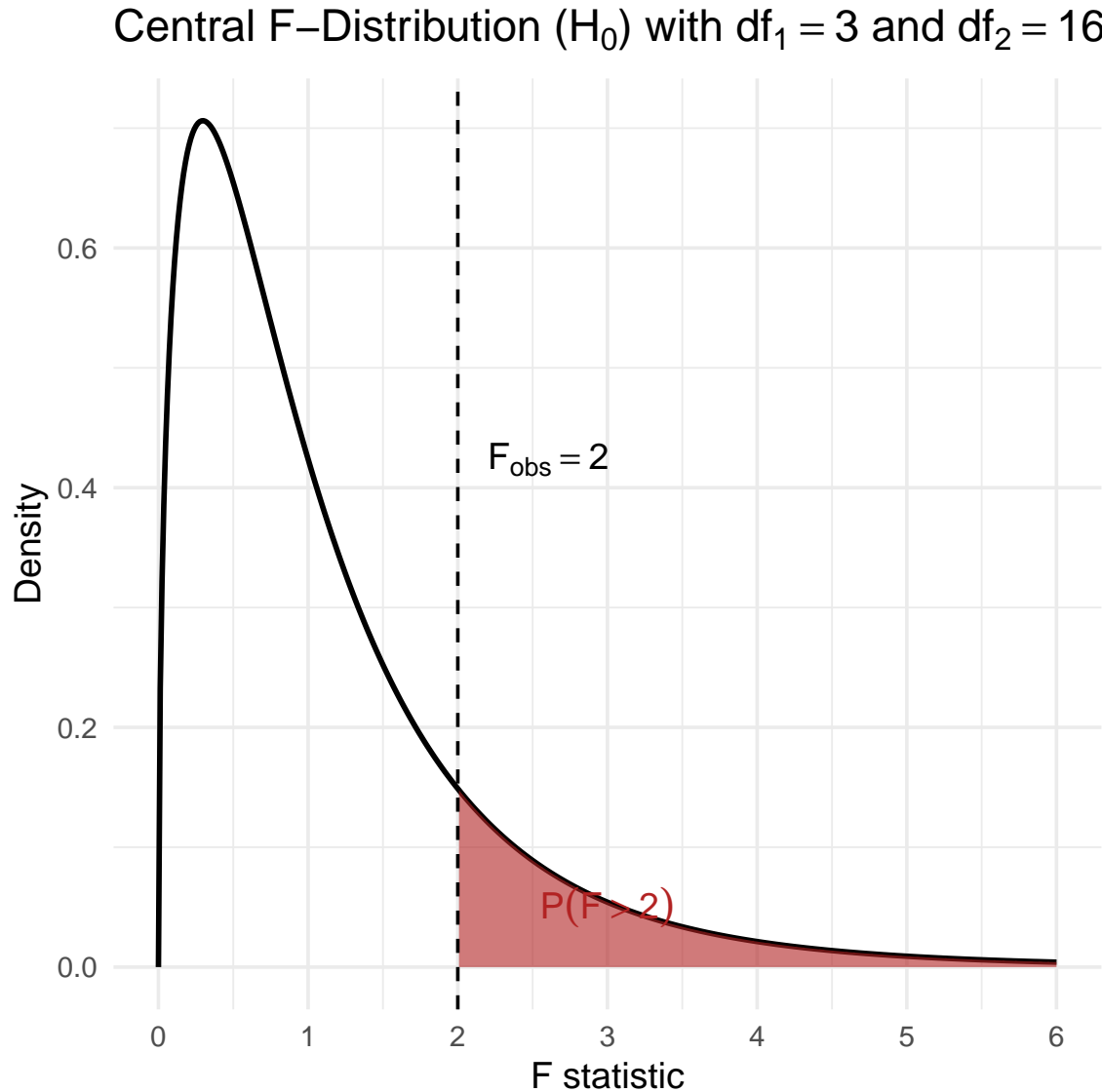


Figure 6.4: Probability Density Function of  $F(3, 16)$  under  $H_0$ . The shaded region represents the p-value.

## 6.9 Coefficient of Determination ( $R^2$ )

### 6.9.1 Definition

The  $R^2$  statistic measures the proportion of total variation explained by the regression model.

**Definition 6.4** (R-Squared).

$$R^2 = \frac{\text{SSR}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}} \quad (6.65)$$

Since  $0 \leq \text{SSE} \leq \text{SST}$ , it follows that  $0 \leq R^2 \leq 1$ .

### 6.9.2 Expectation and Bias

To understand the bias in  $R^2$ , it is more illuminating to analyze the expectation of the **unexplained variance**  $(1 - R^2)$ . This term represents the ratio of error sum of squares to the total sum of squares:

$$E[1 - R^2] = E\left[\frac{\text{SSE}}{\text{SST}}\right] \quad (6.66)$$

Using the first-order approximation  $E[X/Y] \approx E[X]/E[Y]$ , we examine the numerator and denominator separately:

$$\begin{aligned} E[\text{SSE}] &= \sigma^2(n - k - 1) \\ E[\text{SST}] &= \sigma^2(n - 1) + 2\sigma^2\lambda = \sigma^2\left((n - 1) + \frac{\|X_c\beta_1\|^2}{\sigma^2}\right) \end{aligned} \quad (6.67)$$

Substituting these back, we approximate the expected unexplained fraction:

$$E[1 - R^2] \approx \frac{\sigma^2(n - k - 1)}{\sigma^2\left((n - 1) + \frac{\|X_c\beta_1\|^2}{\sigma^2}\right)} = \frac{n - k - 1}{(n - 1) + \frac{\|X_c\beta_1\|^2}{\sigma^2}} \quad (6.68)$$

**Behavior under Null Hypothesis ( $H_0$ ):** When there is no true signal ( $\beta_1 = 0$ ), the term  $\frac{\|X_c\beta_1\|^2}{\sigma^2}$  vanishes. The expected proportion of unexplained variance becomes:

$$E[1 - R^2|H_0] \approx \frac{n - k - 1}{n - 1} \quad (6.69)$$

This result reveals the source of the bias:

1. Ideally, if predictors are noise, the model should explain nothing, and  $E[1 - R^2]$  should be 1.
2. Instead, the expected error ratio is **less than 1**, specifically scaled by  $\frac{n-k-1}{n-1}$ .
3. This scaling factor is exactly what the **Adjusted R-squared** ( $R_a^2$ ) attempts to correct by multiplying the observed ratio by the inverse  $\frac{n-1}{n-k-1}$ .

### 6.9.3 Exact Distribution

The  $R^2$  statistic follows the Type I Non-central Beta distribution derived from the ratio of independent Chi-squared variables.

**Theorem 6.11** (Distribution of R-Squared).

$$R^2 \sim \text{Beta}_1\left(\frac{k}{2}, \frac{n-k-1}{2}, \lambda\right) \quad (6.70)$$

where  $df_1 = k$  and  $df_2 = n - k - 1$ .

### 6.9.4 Adjusted R-squared ( $R_a^2$ )

To correct for the inflation of  $R^2$  due to model complexity ( $k$ ), we introduce the Adjusted  $R^2$ . This statistic penalizes the sum of squares by their degrees of freedom:

$$R_a^2 = 1 - \frac{\text{SSE}/(n-k-1)}{\text{SST}/(n-1)} = 1 - \frac{\text{MSE}}{\text{MST}} = 1 - (1 - R^2) \frac{n-1}{n-k-1} \quad (6.71)$$

**Expectation:**

Under  $H_0$ , since  $E[\text{MSE}] = E[\text{MST}] = \sigma^2$ , the estimator is asymptotically unbiased:

$$E[R_a^2|H_0] \approx 0 \quad (6.72)$$

**Variance and Stability:**

While  $R_a^2$  corrects the bias, it introduces instability. The variance of  $R_a^2$  under  $H_0$  can be derived from the variance of the Beta distribution:

$$\text{Var}(R_a^2|H_0) = \left(\frac{n-1}{n-k-1}\right)^2 \text{Var}(R^2|H_0) \quad (6.73)$$

Substituting  $\text{Var}(R^2|H_0) = \frac{2k(n-k-1)}{(n-1)^2(n+1)}$ , we obtain:

$$\text{Var}(R_a^2|H_0) = \frac{2k}{(n-k-1)(n+1)} \quad (6.74)$$

**Key Insight:**

As the model complexity  $k$  increases relative to  $n$ :

1. The denominator  $(n - k - 1)$  shrinks.
2. The variance  $\text{Var}(R_a^2)$  explodes.

This implies that for high-dimensional models (large  $k/n$ ),  $R_a^2$  becomes an extremely noisy estimator, often yielding large negative values even for null models.

### 6.9.5 Relationship with Rao-Blackwell Decomposition of Variances

The formula for the expected Adjusted  $R^2$  reveals a deep connection to the decomposition of variance in population quantities. Recall the Rao-Blackwell theorem (or Law of Total Variance), which decomposes the total variance of a single observation  $Y_i$  into the expected conditional variance (noise) and the variance of the conditional expectation (signal). Let  $\sigma_\mu^2$  denote the signal variance and  $\sigma^2$  denote the noise variance:

$$\text{Var}(Y_i) = E[\text{Var}(Y_i|x_{(i)})] + \text{Var}(E[Y_i|x_{(i)}]) \quad (6.75)$$

$$\sigma_Y^2 = \sigma^2 + \sigma_\mu^2 \quad (6.76)$$

In our derived expectation for  $R_a^2$ :

$$E[R_a^2] \approx \frac{\frac{\|X_c\beta_1\|^2}{n-1}}{\sigma^2 + \frac{\|X_c\beta_1\|^2}{n-1}} \quad (6.77)$$

The term in the numerator,  $\frac{\|X_c\beta_1\|^2}{n-1}$ , is precisely the **sample variance of the true means**  $\mu_i$ . Let  $\mu = X\beta$ . We can expand the centered signal vector  $X_c\beta_1$  to see this explicitly. Since  $\mu \in \text{Col}(X)$ , we know  $H\mu = \mu$ :

$$X_c\beta_1 = P_{X_c}\mu = (H - P_{j_n})\mu = H\mu - P_{j_n}\mu = \mu - \bar{\mu}j_n = \begin{pmatrix} \mu_1 - \bar{\mu} \\ \mu_2 - \bar{\mu} \\ \vdots \\ \mu_n - \bar{\mu} \end{pmatrix} \quad (6.78)$$

This vector represents the deviation of each observation's true mean from the grand mean. Consequently, the squared norm divided by degrees of freedom is:

$$\frac{\|X_c\beta_1\|^2}{n-1} = \frac{\sum_{i=1}^n (\mu_i - \bar{\mu})^2}{n-1} = \hat{\sigma}_\mu^2 \quad (6.79)$$

If we view the rows of the design matrix  $X$  as random draws  $x_{(1)}, \dots, x_{(n)}$  from a population of covariate vectors, this term estimates  $\text{Var}(x'_{(i)}\beta)$ , which is the variance of the signal component  $\sigma_\mu^2$ .

Thus,  $R_a^2$  can be interpreted as a method-of-moments estimator for the **proportion of variance explained by the signal** in the population:

$$E[R_a^2] \approx \frac{\sigma_\mu^2}{\sigma^2 + \sigma_\mu^2} = \frac{\text{Var}(E[Y_i|x_{(i)}])}{\text{Var}(Y_i)} \quad (6.80)$$

#### ! MSR Is Not a Variance Estimator

- Observing that  $E[\text{MST}] \approx \sigma^2 + \sigma_\mu^2$  and  $E[\text{MSE}] = \sigma^2$ , we can see that the difference  $\text{MST} - \text{MSE}$  provides a direct method-of-moments estimator for the variance of the signal itself ( $\sigma_\mu^2$ ).
- It is important to recognize that the commonly used **Mean Square Regression (MSR)**, defined as  $\text{SSR}/k$ , is **not** an estimator of the signal variance. Because  $E[\text{MSR}] = \sigma^2 + \frac{\|X_c\beta_1\|^2}{k}$ , it scales with the sample size  $n$  (via the squared norm) rather than converging to a population parameter. MSR is

designed for hypothesis testing (detecting *existence* of signal), not for estimating the *magnitude* of the signal variance.

## 6.10 Relationship between $R^2$ and $F$ Test

The proportion of *unexplained* variance,  $1 - R^2$ , follows the Type II Non-central Beta distribution.

**Theorem 6.12** (Distribution of Unexplained Variance).

$$1 - R^2 = \frac{SSE}{SST} \sim \text{Beta}_2 \left( \frac{n - k - 1}{2}, \frac{k}{2}, \lambda \right) \quad (6.81)$$

Both the standard and adjusted coefficients of determination are monotonic functions of the  $F$ -statistic ( $F = \text{MSR}/\text{MSE}$ ). We can derive these relationships directly from the definition of  $F$ :

$$F = \frac{\text{SSR}/k}{\text{SSE}/(n - k - 1)} \implies \frac{\text{SSR}}{\text{SSE}} = \frac{k}{n - k - 1} F \quad (6.82)$$

### 1. Standard $R^2$ :

Recall that  $R^2 = \frac{\text{SSR}}{\text{SST}} = \frac{\text{SSR}}{\text{SSR} + \text{SSE}}$ . Dividing the numerator and denominator by SSE:

$$R^2 = \frac{\text{SSR}/\text{SSE}}{(\text{SSR}/\text{SSE}) + 1} = \frac{\frac{k}{n - k - 1} F}{1 + \frac{k}{n - k - 1} F} \quad (6.83)$$

### 2. Adjusted $R_a^2$ :

Start with the definition involving the unexplained variance scaling:

$$R_a^2 = 1 - (1 - R^2) \frac{n - 1}{n - k - 1} \quad (6.84)$$

From the derivation above, note that the unexplained variance is:

$$1 - R^2 = \frac{1}{1 + \frac{k}{n - k - 1} F} = \frac{n - k - 1}{(n - k - 1) + kF} \quad (6.85)$$

Substituting this into the adjusted formula simplifies the expression:

$$R_a^2 = 1 - \left( \frac{n - k - 1}{(n - k - 1) + kF} \right) \frac{n - 1}{n - k - 1} = 1 - \frac{n - 1}{(n - k - 1) + kF} \quad (6.86)$$

### 6.10.1 Confidence Interval of Population $R^2$

Since we have identified that  $R_a^2$  estimates the population proportion of variance, let us define this target parameter formally as  $\rho^2$ :

$$\rho^2 = 1 - \frac{\sigma^2}{\sigma_Y^2} = \frac{\sigma_\mu^2}{\sigma_Y^2} \quad (6.87)$$

While  $R_a^2$  provides a point estimate, we can construct an exact confidence interval for  $\rho^2$  by exploiting the distribution of the  $F$ -statistic.

#### 1. The Link to Non-centrality:

Recall that the  $F$ -statistic follows a non-central distribution  $F(k, n - k - 1, \lambda)$ . The non-centrality parameter  $\lambda$  is directly related to the population  $\rho^2$ . Using the variance decomposition derived above:

$$\lambda = \frac{\|X_c \beta_1\|^2}{2\sigma^2} = \frac{n-1}{2} \left( \frac{\sigma_\mu^2}{\sigma^2} \right) \quad (6.88)$$

Substituting the signal-to-noise ratio  $\frac{\sigma_\mu^2}{\sigma^2} = \frac{\rho^2}{1-\rho^2}$ , we obtain a one-to-one mapping between  $\lambda$  and  $\rho^2$ :

$$\lambda(\rho^2) = \frac{n-1}{2} \left( \frac{\rho^2}{1-\rho^2} \right) \quad (6.89)$$

#### 2. Inverting the Test Statistic:

We find a confidence interval  $[\lambda_L, \lambda_U]$  for  $\lambda$  by “inverting” the observed  $F$ -statistic ( $F_{obs}$ ). We search for two specific non-central F-distributions: one where  $F_{obs}$  cuts off the upper  $\alpha/2$  tail, and one where it cuts off the lower  $\alpha/2$  tail.

- **Lower Bound ( $\lambda_L$ ):** The non-centrality parameter such that  $F_{obs}$  is the  $1 - \alpha/2$  quantile.
- **Upper Bound ( $\lambda_U$ ):** The non-centrality parameter such that  $F_{obs}$  is the  $\alpha/2$  quantile.

This concept is illustrated in the figure below.

#### 3. The Interval for $\rho^2$ :

Once  $[\lambda_L, \lambda_U]$  are found numerically, we map them back to the population  $R^2$  scale using the inverse relationship:

$$\rho^2 = \frac{2\lambda}{2\lambda + (n-1)} \quad (6.90)$$

This produces an exact confidence interval  $[\rho_L^2, \rho_U^2]$  for the proportion of variance explained by the model in the population.

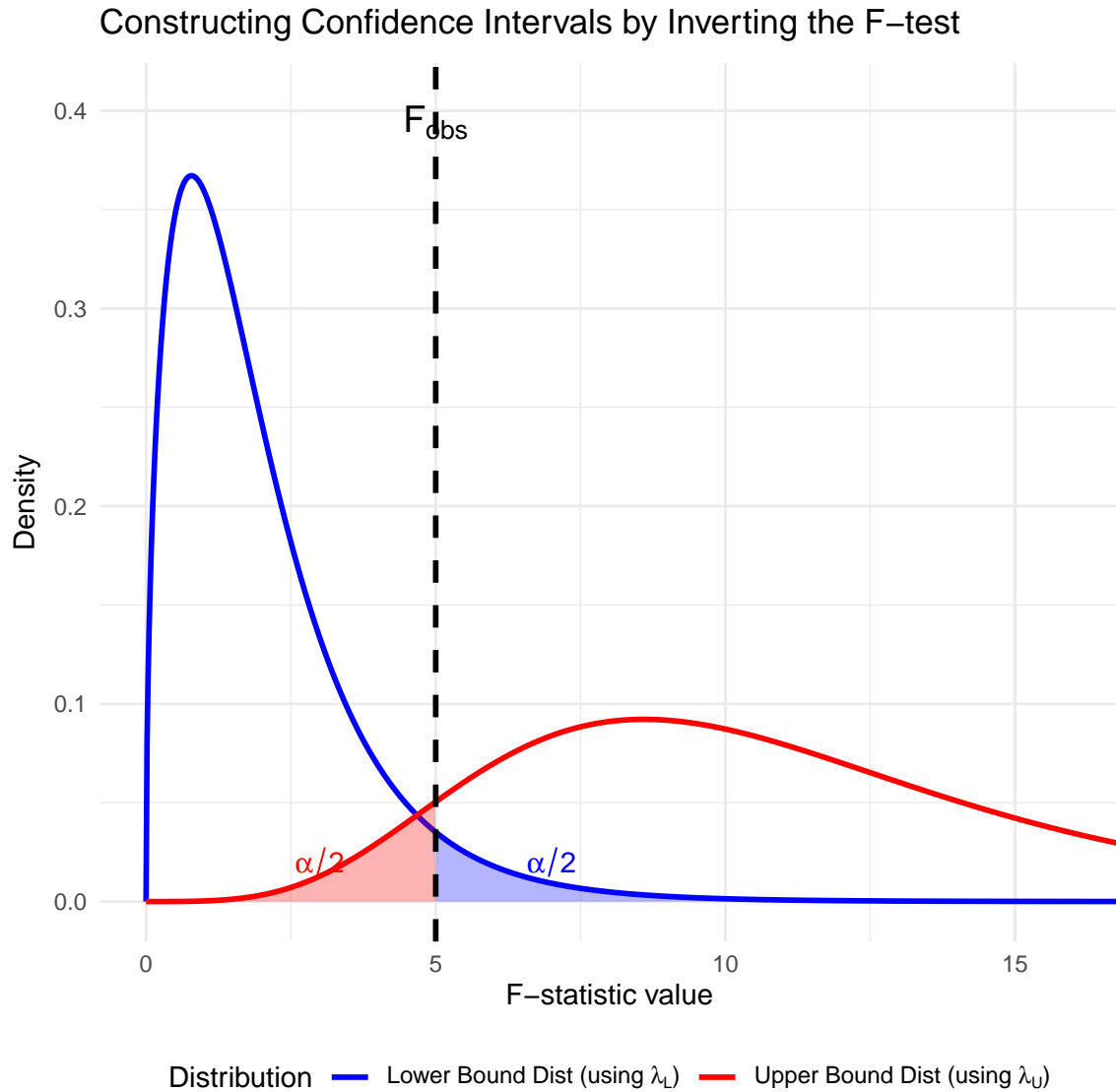


Figure 6.5: Illustration of constructing a confidence interval for the non-centrality parameter  $\lambda$  by inverting the F-test. The observed  $F_{obs}$  (dashed line) is the 97.5<sup>th</sup> percentile of the distribution defined by the lower bound  $\lambda_L$  (blue), and the 2.5<sup>th</sup> percentile of the distribution defined by the upper bound  $\lambda_U$  (red). The shaded areas each represent  $\alpha/2$ .



## 6.11 An Animation for Illustrating $R_a^2$ Under $H_0$ and $H_1$

We simulate a dataset with  $n = 30$  observations and consider a sequence of nested models adding groups of predictors.

### Predictor Groups:

1. **Step 1** ( $k = 1$ ): Add  $x_1$ . (Signal under  $H_1$ ).
2. **Step 2** ( $k = 6$ ): Add  $x_2, \dots, x_6$  (Noise).
3. **Step 3** ( $k = 11$ ): Add  $x_7, \dots, x_{11}$  (Noise).
4. **Step 4** ( $k = 20$ ): Add  $x_{12}, \dots, x_{20}$  (Noise).

### 6.11.0.1 Null Hypothesis ( $H_0$ )

Under  $H_0$ , the true coefficient for  $x_1$  is  $\beta_1 = 0$ . All predictors are noise.

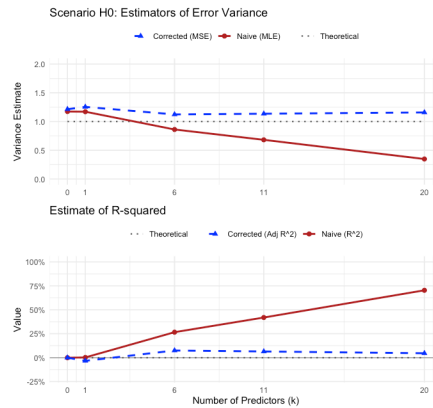


Figure 6.6: Simulation under  $H_0$ : As predictors are added (pure noise), standard R-squared increases while Adjusted R-squared and MSE remain stable.

### 6.11.0.2 Alternative Hypothesis ( $H_1$ )

Under  $H_1$ ,  $x_1$  is a true predictor ( $\beta_1 = 2$ ). The subsequent groups ( $x_2 \dots x_{20}$ ) remain noise.

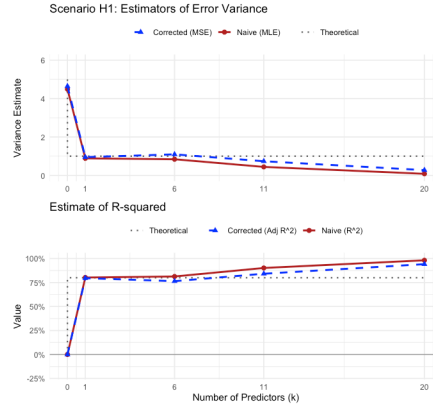


Figure 6.7: Simulation under H1: Adjusted R-squared correctly identifies the signal at  $k=1$ , then penalizes the subsequent noise predictors.

## 6.12 A Data Example with House Price Valuation

A real estate agency wants to refine their pricing model. They regress the selling price of houses ( $y$ ) on five predictors ( $X$ ): Size, Age, Bedrooms, Garage Capacity, and Lawn Size.

We assume the data has been collected and saved to `house_prices_5pred.csv`.

### 6.12.1 Visualize the Data

First, we load the dataset. We display the first 10 rows for PDF output, or a full paged table for HTML.

Table 6.1: First 10 rows of House Prices

Price	Size	Age	Beds	Garage	Lawn
425767	3092	18	5	1	325
336991	1802	37	2	1	687
528842	2701	49	2	0	261
399797	2745	0	5	2	554
427580	2143	1	5	3	296
478082	2754	26	4	0	833
295549	2039	17	2	3	194
335058	1758	11	3	1	111
461110	3191	58	2	2	286
204405	1298	41	2	0	813

## 6.12.2 Fit the Model

We will solve for the coefficients  $\hat{\beta}$  using three distinct methods.

### 6.12.2.1 Method 1: Naive Matrix Formula

This method solves the normal equations directly on the raw data:  $\hat{\beta} = (X'X)^{-1}X'y$ .

Matrix  $X'X$  (Cross-products of predictors):

	Intercept	Size	Age	Beds	Garage	Lawn
Intercept	25	59157	766	77	32	12049
Size	59157	153276601	1825278	183641	75356	27856192
Age	766	1825278	31762	2083	871	398142
Beds	77	183641	2083	273	95	38151
Garage	32	75356	871	95	74	13196
Lawn	12049	27856192	398142	38151	13196	7896317

Matrix  $X'y$  (Cross-products with response):

	[,1]
Intercept	9754828
Size	24994304201
Age	294127221
Beds	30527650
Garage	12400638
Lawn	4533697188

Solved Coefficients (Beta):

	Intercept	Size	Age	Beds	Garage	Lawn
[1,]	85085.08	142.1976	-624.3362	3446.122	-5014.307	-34.10539

### 6.12.2.2 Method 2: Centralized Formula

This method reduces multicollinearity issues. Formula:  $\hat{\beta}_{\text{slope}} = (X'_c X_c)^{-1} X'_c y_c$ .

Matrix  $X_c'X_c$  (Centered Sum of Squares):

	Size	Age	Beds	Garage	Lawn
Size	13294575	12708	1437	-365	-655116
Age	12708	8292	-276	-109	28961
Beds	1437	-276	36	-4	1040
Garage	-365	-109	-4	33	-2227
Lawn	-655116	28961	1040	-2227	2089181

Matrix  $X_c'y_c$  (Centered Cross-products):

	[,1]
Size	1911649801
Age	-4760709
Beds	482780
Garage	-85542
Lawn	-167739715

Solved Coefficients (Beta):

	Intercept	Size	Age	Beds	Garage	Lawn
[1,]	85085.08	142.1976	-624.3362	3446.122	-5014.307	-34.10539

### 6.12.2.3 Method 3: Using R's `lm` Function

This is the standard approach for practitioners.

(Intercept)	Size	Age	Beds	Garage	Lawn
85085.08052	142.19762	-624.33616	3446.12224	-5014.30654	-34.10539

## 6.12.3 Analysis of Variance (ANOVA)

We now evaluate the sources of variation to test the overall model significance.

### 6.12.3.1 1. Computing Sums of Squares

SST: 358921762209    SSR: 282617820807    SSE: 76303941402

### 6.12.3.2 2. Manual ANOVA Construction

We build the table manually using the sums of squares and degrees of freedom.

Table 6.2: Manual ANOVA Table

Source	DF	SS	MS	F_Statistic	P_Value
Regression (Model)	5	282617820807	56523564161	14.0746	0
Error (Residual)	19	76303941402	4015996916	NA	NA
Total	24	358921762209	14955073425	NA	NA

### 6.12.3.3 3. Standard R Output (summary and anova)

We display the standard `summary()` which provides the coefficients, t-tests, and the overall F-statistic found at the bottom. We also show `anova()` which gives the sequential sum of squares.

Call:

```
lm(formula = Price ~ ., data = df)
```

Residuals:

```
      Min       1Q   Median       3Q      Max
-126301  -37660    5319   40319   92283
```

Coefficients:

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 85085.08   75276.69   1.130   0.272
Size         142.20     17.74    8.015 1.63e-07 ***
Age          -624.34    888.03   -0.703   0.491
Beds         3446.12   13154.95    0.262   0.796
Garage      -5014.31   11826.13   -0.424   0.676
Lawn         -34.11     48.21   -0.707   0.488
```

---

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 63370 on 19 degrees of freedom

Multiple R-squared: 0.7874, Adjusted R-squared: 0.7315

F-statistic: 14.07 on 5 and 19 DF, p-value: 7.867e-06

Analysis of Variance Table

Response: Price

```
      Df    Sum Sq   Mean Sq F value    Pr(>F)
Size   1 2.7488e+11 2.7488e+11 68.4461 1.014e-07 ***
Age    1 5.2419e+09 5.2419e+09  1.3053  0.2674
Beds   1 1.1538e+08 1.1538e+08  0.0287  0.8672
Garage 1 3.7115e+08 3.7115e+08  0.0924  0.7644
Lawn   1 2.0100e+09 2.0100e+09  0.5005  0.4879
```

```
Residuals 19 7.6304e+10 4.0160e+09
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

#### 6.12.4 Coefficient of Determination and Variance Decomposition

We calculate  $R^2$  and Adjusted  $R^2$ , and then present them in a **Variance Decomposition Table**.

##### 6.12.4.1 1. Calculation

Standard  $R^2$ : 0.7874

Adjusted  $R^2$ : 0.7315

##### 6.12.4.2 2. Variance Decomposition Table

This table extends standard ANOVA. While ANOVA focuses on **Mean Squares (MS)** for hypothesis testing (is  $MSR > MSE$ ?), this table focuses on **Variance Components** ( $\hat{\sigma}^2$ ) for estimation (how much variance is Signal vs. Noise?).

- **Signal Variance** ( $\hat{\sigma}_\mu^2$ ): Estimated by  $MST - MSE$ . (Note:  $MSR$  is biased and overestimates signal).
- **Noise Variance** ( $\hat{\sigma}^2$ ): Estimated by  $MSE$ .
- **Total Variance** ( $\hat{\sigma}_Y^2$ ): Estimated by  $MST$ .

Table 6.3: Variance Decomposition Table: Estimating Signal vs. Noise

Component	DF	SS	MS	Value ( $\hat{\sigma}^2$ )	Proportion
Signal (Model)	5	282617820807	56523564161	10939076509	0.7315
Noise (Error)	19	76303941402	4015996916	4015996916	0.2685
Total (Y)	24	358921762209	14955073425	14955073425	1.0000

#### 6.12.5 Confidence Interval for Population $R^2$ ( $\rho^2$ )

We construct a 95% confidence interval for the population proportion of variance explained ( $\rho^2$ ).

##### 6.12.5.1 1. Manual Inversion Method

We solve for the non-centrality parameters  $\lambda_L$  and  $\lambda_U$  such that our observed  $F_{obs}$  corresponds to the appropriate quantiles.

- **Correction Note:** R's `ncp` parameter represents the *full* sum of squares. Therefore, the conversion to  $\rho^2$  for fixed predictors is  $\rho^2 = \frac{\lambda}{\lambda + n}$ .

Manual Calculation:

95% CI for Population  $\rho^2$ : [ 0.4674 , 0.8401 ]

### 6.12.5.2 2. Using R Package MBESS

The MBESS package automates this procedure. We use `Random.Predictors = FALSE` to match the fixed-predictor assumption used in our manual calculation.

```
$Lower.Conf.Limit.R2
```

```
[1] 0.467363
```

```
$Prob.Less.Lower
```

```
[1] 0.025
```

```
$Upper.Conf.Limit.R2
```

```
[1] 0.8400781
```

```
$Prob.Greater.Upper
```

```
[1] 0.025
```

## 6.13 Overfitting

Suppose the reduced model  $y = X_1\beta_1^* + e$  is true (i.e.,  $\beta_2 = 0$ ), but we fit the full model ( $\dagger$ ). Since the full model includes the true model as a special case, the estimator  $\hat{\beta}$  from the full model remains unbiased.

However, fitting the extraneous variables affects the variance.

**Theorem 6.13** (Variance Comparison). *Let  $\hat{\beta}_1$  be the estimator from the full model and  $\hat{\beta}_1^*$  be the estimator from the reduced model. Then:*

$$\text{Var}(\hat{\beta}_1) - \text{Var}(\hat{\beta}_1^*) = \sigma^2 AB^{-1}A^T \quad (6.91)$$

where  $A = (X_1^T X_1)^{-1} X_1^T X_2$  and  $B = X_2^T X_2 - X_2^T X_1 A$ . Since  $AB^{-1}A^T$  is positive semidefinite,  $\text{Var}(\hat{\beta}_1) \geq \text{Var}(\hat{\beta}_1^*)$ .

*Proof.* Using the inverse of a partitioned matrix, the top-left block of  $(X^T X)^{-1}$  corresponding to  $\beta_1$  is:

$$H^{11} = (X_1^T X_1)^{-1} + (X_1^T X_1)^{-1} X_1^T X_2 B^{-1} X_2^T X_1 (X_1^T X_1)^{-1} \quad (6.92)$$

Since  $\text{Var}(\hat{\beta}_1) = \sigma^2 H^{11}$  and  $\text{Var}(\hat{\beta}_1^*) = \sigma^2 (X_1^T X_1)^{-1}$ , the difference is the second term:

$$\text{Var}(\hat{\beta}_1) - \text{Var}(\hat{\beta}_1^*) = \sigma^2 AB^{-1}A^T \quad (6.93)$$

□

### 6.13.1 Summary

1. **Underfitting:** Reduces variance but introduces bias (unless variables are orthogonal).
2. **Overfitting:** Keeps estimators unbiased but increases variance.



# 7 Generalized Inverse

## 7.1 Generalized Inverses

### 7.1.1 Motivation

Consider the linear system  $X\beta = y$ . In  $\mathbb{R}^2$ , if  $X = [x_1, x_2]$  is invertible, the solution is unique:  $\beta = X^{-1}y$ . This satisfies  $X(X^{-1}y) = y$ . However, if  $X$  is not square or not invertible (e.g.,  $X$  is  $2 \times 3$ ),  $X\beta = y$  does not have a unique solution. We seek a matrix  $G$  such that  $\beta = Gy$  provides a solution whenever  $y \in C(X)$  (the column space of  $X$ ). Substituting  $\beta = Gy$  into the equation  $X\beta = y$ :

$$X(Gy) = y \quad \forall y \in C(X) \quad (7.1)$$

Since any  $y \in C(X)$  can be written as  $Xw$  for some vector  $w$ :

$$XGXw = Xw \quad \forall w \quad (7.2)$$

This implies the defining condition:

$$XGX = X \quad (7.3)$$

### 7.1.2 Definition of Generalized Inverse

**Definition 7.1** (Generalized Inverse). Let  $X$  be an  $n \times p$  matrix. A matrix  $X^-$  of size  $p \times n$  is called a **generalized inverse** of  $X$  if it satisfies:

$$XX^-X = X \quad (7.4)$$

**Example 7.1** (Examples of Generalized Inverse).

- **Example 1: Diagonal Matrix** If  $X = \text{diag}(\lambda_1, \lambda_2, 0, 0)$ , we can write it in matrix form as:

$$X = \begin{pmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (7.5)$$

A generalized inverse is obtained by inverting the non-zero elements:

$$X^- = \begin{pmatrix} \lambda_1^{-1} & 0 & 0 & 0 \\ 0 & \lambda_2^{-1} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (7.6)$$

- **Example 2: Row Vector** Let  $X = (1, 2, 3)$ . One possible generalized inverse is a column vector where the first element is the reciprocal of the first non-zero element of  $X$  (which is 1), and others are zero:

$$X^- = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad (7.7)$$

**Verification:**

$$XX^-X = (1, 2, 3) \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} (1, 2, 3) = (1) \cdot (1, 2, 3) = (1, 2, 3) = X \quad (7.8)$$

Other valid generalized inverses include  $\begin{pmatrix} 0 \\ 1/2 \\ 0 \end{pmatrix}$  or  $\begin{pmatrix} 0 \\ 0 \\ 1/3 \end{pmatrix}$ .

- **Example 3: Rank Deficient Matrix** Let  $A = \begin{pmatrix} 2 & 2 & 3 \\ 1 & 0 & 1 \\ 3 & 2 & 4 \end{pmatrix}$ . Note that Row 3 = Row 1 + Row 2, so

$\text{Rank}(A) = 2$ .

**Solution:** A generalized inverse can be found by locating a non-singular  $2 \times 2$  submatrix, inverting it, and padding the rest with zeros. Let's take the top-left minor  $M = \begin{pmatrix} 2 & 2 \\ 1 & 0 \end{pmatrix}$ . The inverse is  $M^{-1} = \frac{1}{-2} \begin{pmatrix} 0 & -2 \\ -1 & 2 \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ 0.5 & -1 \end{pmatrix}$ .

Placing this in the corresponding position in  $A^-$  and setting the rest to 0:

$$A^- = \begin{pmatrix} 0 & 1 & 0 \\ 0.5 & -1 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (7.9)$$

**Verification** ( $AA^-A = A$ ): First, compute  $AA^-$ :

$$AA^- = \begin{pmatrix} 2 & 2 & 3 \\ 1 & 0 & 1 \\ 3 & 2 & 4 \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 \\ 0.5 & -1 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{pmatrix} \quad (7.10)$$

Then multiply by  $A$ :

$$(AA^-)A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} 2 & 2 & 3 \\ 1 & 0 & 1 \\ 3 & 2 & 4 \end{pmatrix} = \begin{pmatrix} 2 & 2 & 3 \\ 1 & 0 & 1 \\ 3 & 2 & 4 \end{pmatrix} = A \quad (7.11)$$

### 7.1.3 A Procedure to Find a Generalized Inverse

If we can partition  $X$  (possibly after permuting rows/columns) such that  $R_{11}$  is a non-singular rank  $r$  submatrix:

$$X = \begin{pmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{pmatrix} \quad (7.12)$$

Then a generalized inverse is:

$$X^- = \begin{pmatrix} R_{11}^{-1} & 0 \\ 0 & 0 \end{pmatrix} \quad (7.13)$$

**Verification:**

$$\begin{aligned} XX^-X &= \begin{pmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{pmatrix} \begin{pmatrix} R_{11}^{-1} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{pmatrix} \\ &= \begin{pmatrix} I_r & 0 \\ R_{21}R_{11}^{-1} & 0 \end{pmatrix} \begin{pmatrix} R_{11} & R_{12} \\ R_{21} & R_{22} \end{pmatrix} \\ &= \begin{pmatrix} R_{11} & R_{12} \\ R_{21} & R_{21}R_{11}^{-1}R_{12} \end{pmatrix} \end{aligned} \quad (7.14)$$

Note that since  $\text{rank}(X) = \text{rank}(R_{11})$ , the rows of  $[R_{21}, R_{22}]$  are linear combinations of  $[R_{11}, R_{12}]$ , implying  $R_{22} = R_{21}R_{11}^{-1}R_{12}$ . Thus,  $XX^-X = X$ .

#### An Algorithm for Finding a Generalized Inverse

A systematic procedure to find a generalized inverse  $A^-$  for any matrix  $A$ :

1. Find any non-singular  $r \times r$  submatrix  $C$ , where  $r$  is the rank of  $A$ . It is not necessary for the elements of  $C$  to occupy adjacent rows and columns in  $A$ .
2. Find  $C^{-1}$  and  $(C^{-1})'$ .
3. Replace the elements of  $C$  in  $A$  with the elements of  $(C^{-1})'$ .
4. Replace all other elements in  $A$  with zeros.
5. Transpose the resulting matrix.

#### Matrix Visual Representation

$$\begin{array}{ccc} \begin{pmatrix} \times & \otimes & \times & \otimes \\ \times & \otimes & \times & \otimes \\ \times & \times & \times & \times \end{pmatrix} & \xrightarrow[\text{with } (C^{-1})']{\text{Replace } C} & \begin{pmatrix} \times & \triangle & \times & \triangle \\ \times & \triangle & \times & \triangle \\ \times & \times & \times & \times \end{pmatrix} & \xrightarrow[\text{Result}]{\text{Transpose}} & \begin{pmatrix} \times & \times & \times \\ \square & \square & \times \\ \times & \times & \times \\ \square & \square & \times \end{pmatrix} \\ \text{Original } A & & \text{Intermediate} & & \text{Final } A^- \end{array} \quad (7.15)$$

**Legend:**

- $\otimes$ : Elements of submatrix  $C$
- $\triangle$ : Elements of  $(C^{-1})'$
- $\square$ : Elements of  $C^{-1}$  (after transposition)
- $\times$ : Other elements (replaced by 0 in the final calculation)

### 7.1.4 Moore-Penrose Inverse

The Moore-Penrose inverse (denoted  $X^+$ ) is a unique generalized inverse defined via Singular Value Decomposition (SVD).

If  $X$  has SVD:

$$X = U \begin{pmatrix} \Lambda_r & 0 \\ 0 & 0 \end{pmatrix} V' \quad (7.16)$$

Then the Moore-Penrose inverse is:

$$X^+ = V \begin{pmatrix} \Lambda_r^{-1} & 0 \\ 0 & 0 \end{pmatrix} U' \quad (7.17)$$

where  $\Lambda_r = \text{diag}(\lambda_1, \dots, \lambda_r)$  contains the singular values. Unlike standard generalized inverses,  $X^+$  is unique.

#### Verification:

We verify that  $X^+$  satisfies the condition  $XX^+X = X$ .

##### 1. Substitute definitions:

$$XX^+X = \left[ U \begin{pmatrix} \Lambda_r & 0 \\ 0 & 0 \end{pmatrix} V' \right] \left[ V \begin{pmatrix} \Lambda_r^{-1} & 0 \\ 0 & 0 \end{pmatrix} U' \right] \left[ U \begin{pmatrix} \Lambda_r & 0 \\ 0 & 0 \end{pmatrix} V' \right] \quad (7.18)$$

##### 2. Apply orthogonality: Recall that $V'V = I$ and $U'U = I$ .

$$= U \begin{pmatrix} \Lambda_r & 0 \\ 0 & 0 \end{pmatrix} \underbrace{(V'V)}_I \begin{pmatrix} \Lambda_r^{-1} & 0 \\ 0 & 0 \end{pmatrix} \underbrace{(U'U)}_I \begin{pmatrix} \Lambda_r & 0 \\ 0 & 0 \end{pmatrix} V' \quad (7.19)$$

##### 3. Multiply diagonal matrices:

$$= U \left[ \begin{pmatrix} \Lambda_r & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \Lambda_r^{-1} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \Lambda_r & 0 \\ 0 & 0 \end{pmatrix} \right] V' \quad (7.20)$$

Since  $\Lambda_r \Lambda_r^{-1} \Lambda_r = I \cdot \Lambda_r = \Lambda_r$ :

$$= U \begin{pmatrix} \Lambda_r & 0 \\ 0 & 0 \end{pmatrix} V' = X \quad (7.21)$$

### 7.1.5 Solving Linear Systems with Generalized Inverse

We apply generalized inverses to solve systems of linear equations  $X\beta = c$  where  $X$  is  $n \times p$ .

**Definition 7.2** (Consistency and Solution). The system  $X\beta = c$  is consistent if and only if  $c \in \mathcal{C}(X)$  (the column space of  $X$ ). If consistent,  $\beta = X^-c$  is a solution.

**Proof:** If the system is consistent, there exists some  $b$  such that  $Xb = c$ . Using the definition  $XX^-X = X$ :

$$X(X^-c) = X(X^-Xb) = (XX^-X)b = Xb = c \quad (7.22)$$

Thus,  $X^-c$  is a solution. Note that the solution is not unique if  $X$  is not full rank.

**Example 7.2** (Examples of Solutions of Linear System with Generalized Inverse).

• **Example 1: Underdetermined System**

Let  $X = \begin{pmatrix} 1 & 2 & 3 \end{pmatrix}$  and we want to solve  $X\beta = 4$ .

**Solution 1:** Using the generalized inverse  $X^- = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$ :

$$\beta = X^- \cdot 4 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} 4 = \begin{pmatrix} 4 \\ 0 \\ 0 \end{pmatrix} \quad (7.23)$$

**Verification:**

$$X\beta = \begin{pmatrix} 1 & 2 & 3 \end{pmatrix} \begin{pmatrix} 4 \\ 0 \\ 0 \end{pmatrix} = 1(4) + 2(0) + 3(0) = 4 \quad \checkmark \quad (7.24)$$

**Solution 2:** Using another generalized inverse  $X^- = \begin{pmatrix} 0 \\ 0 \\ 1/3 \end{pmatrix}$ :

$$\beta = X^- \cdot 4 = \begin{pmatrix} 0 \\ 0 \\ 1/3 \end{pmatrix} 4 = \begin{pmatrix} 0 \\ 0 \\ 4/3 \end{pmatrix} \quad (7.25)$$

**Verification:**

$$X\beta = \begin{pmatrix} 1 & 2 & 3 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 4/3 \end{pmatrix} = 0 + 0 + 3(4/3) = 4 \quad \checkmark \quad (7.26)$$

• **Example 2: Overdetermined System**

Let  $X = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}$ . Solve  $X\beta = \begin{pmatrix} 2 \\ 4 \\ 6 \end{pmatrix} = c$ . Here  $c = 2X$ , so the system is consistent. Since  $X$  is a column vector,  $\beta$  is a scalar.

**Solution:** Using the generalized inverse  $X^- = \begin{pmatrix} 1 & 0 & 0 \end{pmatrix}$ :

$$\beta = X^-c = \begin{pmatrix} 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 2 \\ 4 \\ 6 \end{pmatrix} = 1(2) + 0(4) + 0(6) = 2 \quad (7.27)$$

**Verification:**

$$X\beta = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} (2) = \begin{pmatrix} 2 \\ 4 \\ 6 \end{pmatrix} = c \quad \checkmark \quad (7.28)$$

## 7.2 Least Squares for Non-full-rank $X$ with Generalized Inverse

### 7.2.1 Projection Matrix with Generalized Inverse of $X'X$

For the normal equations  $(X'X)\beta = X'y$ , a solution is given by:

$$\hat{\beta} = (X'X)^- X'y \quad (7.29)$$

The fitted values are

$$\hat{y} = X\hat{\beta} = X(X'X)^- X'y. \quad (7.30)$$

This  $\hat{y}$  represents the unique orthogonal projection of  $y$  onto  $\text{Col}(X)$ .

### 7.2.2 Invariance and Uniqueness of “the” Projection Matrix

**Theorem 7.1** (Transpose Property of Generalized Inverses).  $(X^-)'$  is a version of  $(X')^-$ . That is,  $(X^-)'$  is a generalized inverse of  $X'$ .

*Proof.* By definition, a generalized inverse  $X^-$  satisfies the property:

$$XX^-X = X \quad (7.31)$$

To verify that  $(X^-)'$  is a generalized inverse of  $X'$ , we need to show that it satisfies the condition  $AGA = A$  where  $A = X'$  and  $G = (X^-)'$ .

1. Start with the fundamental definition:

$$XX^-X = X \quad (7.32)$$

2. Take the transpose of both sides of the equation:

$$(XX^-X)' = X' \quad (7.33)$$

3. Apply the reverse order law for transposes,  $(ABC)' = C'B'A'$ :

$$X'(X^-)'X' = X' \quad (7.34)$$

Since substituting  $(X^-)'$  into the generalized inverse equation for  $X'$  yields  $X'$ ,  $(X^-)'$  is a valid generalized inverse of  $X'$ .  $\square$

**Lemma 7.1** (Invariance of Generalized Least Squares). For any version of the generalized inverse  $(X'X)^-$ , the matrix  $X'(X'X)^-X'$  is invariant and equals  $X'$ .

$$X'X(X'X)^-X' = X' \quad (7.35)$$

**Proof (using Projection):** Let  $P = X(X'X)^-X'$ . This is the projection matrix onto  $\mathcal{C}(X)$ . By definition of projection,  $Px = x$  for any  $x \in \text{Col}(X)$ . Since columns of  $X$  are in  $\text{Col}(X)$ ,  $PX = X$ . Taking the transpose:  $(PX)' = X' \implies X'P' = X'$ . Since projection matrices are symmetric ( $P = P'$ ),  $X'P = X'$ . Substituting  $P$ :  $X'X(X'X)^-X' = X'$ .

**Proof (Direct Matrix Manipulation):** Decompose  $y = X\beta + e$  where  $e \perp \text{Col}(X)$  (i.e.,  $X'e = 0$ ).

$$\begin{aligned} X'X(X'X)^-X'y &= X'X(X'X)^-X'(X\beta + e) \\ &= X'X(X'X)^-X'X\beta + X'X(X'X)^-X'e \end{aligned} \quad (7.36)$$

Using the property  $AA^-A = A$  (where  $A = X'X$ ), the first term becomes  $X'X\beta$ . The second term is 0 because  $X'e = 0$ . Thus, the expression simplifies to  $X'X\beta = X'(X\beta) = X'\hat{y}_{\text{proj}}$ . This implies the operator acts as  $X'$ .

**Theorem 7.2** (Properties of Projection Matrix  $P$ ). *Let  $P = X(X'X)^-X'$ . This matrix has the following properties:*

1. **Symmetry:**  $P = P'$ .
2. **Idempotence:**  $P^2 = P$ .

$$P^2 = X(X'X)^-X'X(X'X)^-X' = X(X'X)^-(X'X(X'X)^-X') \quad (7.37)$$

Using the identity from Lemma 7.1 ( $X'X(X'X)^-X' = X'$ ), this simplifies to:

$$X(X'X)^-X' = P \quad (7.38)$$

3. **Uniqueness:**  $P$  is unique and invariant to the choice of the generalized inverse  $(X'X)^-$ .

*Proof.* **Proof of Uniqueness:**

Let  $A$  and  $B$  be two different generalized inverses of  $X'X$ . Define  $P_A = XAX'$  and  $P_B = XBX'$ . From Lemma 7.1, we know that  $X'P_A = X'$  and  $X'P_B = X'$ .

Subtracting these two equations:

$$X'(P_A - P_B) = 0 \quad (7.39)$$

Taking the transpose, we get  $(P_A - P_B)X = 0$ . This implies that the columns of the difference matrix  $D = P_A - P_B$  are orthogonal to the columns of  $X$  (i.e.,  $D \perp \text{Col}(X)$ ).

However, by definition, the columns of  $P_A$  and  $P_B$  (and thus  $D$ ) are linear combinations of the columns of  $X$  (i.e.,  $D \in \text{Col}(X)$ ).

The only matrix that lies in the column space of  $X$  but is also *orthogonal* to the column space of  $X$  is the zero matrix. Therefore:

$$P_A - P_B = 0 \implies P_A = P_B \quad (7.40)$$

□

## 7.3 The Left Inverse View: Recovering $\hat{\beta}$ from $\hat{y}$

While the geometric properties of the linear model are most naturally established via the unique orthogonal projection  $\hat{y}$ , we require a functional mapping—a statistical “bridge”—to translate the distribution of these fitted values back into the parameter space of  $\hat{\beta}$ . This bridge is provided by the generalized left inverse.

### 7.3.1 The Generalized Left Inverse

To recover the parameter estimates directly from the fitted values, we define the generalized left inverse, denoted as  $X_{\text{left}}^-$ , such that:

$$\hat{\beta} = X_{\text{left}}^- \hat{y} \quad (7.41)$$

A standard choice for this operator, derived from the normal equations, is:

$$X_{\text{left}}^- = (X'X)^- X' \quad (7.42)$$

When  $X$  is full-rank, the  $X_{\text{left}}^-$  is unique, which is given by

$$X_{\text{left}}^- = (X'X)^{-1} X' \quad (7.43)$$

### 7.3.2 Verification of the Inverse Property

To verify that  $X_{\text{left}}^-$  acts as a valid generalized inverse of  $X$ , it must satisfy the condition  $XX_{\text{left}}^-X = X$ . Substituting our definition:

$$X \underbrace{[(X'X)^- X']}_{X_{\text{left}}^-} X = X(X'X)^-(X'X) \quad (7.44)$$

Using the property of generalized inverses for symmetric matrices where  $(X'X)(X'X)^-X' = X'$ , the transpose of this identity gives  $X(X'X)^-(X'X) = X$ . Thus, the condition holds:

$$XX_{\text{left}}^-X = X \quad (7.45)$$



### 7.3.3 Recovering the Estimator

We can now demonstrate that applying this left inverse to the fitted values  $\hat{y}$  yields the standard solution to the normal equations.

Substituting the projection formula  $\hat{y} = X(X'X)^{-}X'y$ :

$$\begin{aligned} X_{\text{left}}^{-}\hat{y} &= [(X'X)^{-}X'] [X(X'X)^{-}X'y] \\ &= (X'X)^{-} \underbrace{(X'X)(X'X)^{-}(X'X)}_{\text{Property } AA^{-}A=A} (X'X)^{-}X'y \end{aligned} \quad (7.46)$$

Simplifying using the generalized inverse property  $A^{-}AA^{-} = A^{-}$  (where  $A = X'X$ ):

$$\begin{aligned} X_{\text{left}}^{-}\hat{y} &= \underbrace{(X'X)^{-}(X'X)(X'X)^{-}}_{(X'X)^{-}} X'y \\ &= (X'X)^{-}X'y \end{aligned} \quad (7.47)$$

Thus, we recover the standard estimator used in the normal equations:

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y} \quad (7.48)$$

## 7.4 Non-full-rank Least Squares with QR Decomposition

When  $X$  has rank  $r < p$  (where  $X$  is  $n \times p$ ), we can derive the least squares estimator using partitioned matrices.

Assume the first  $r$  columns of  $X$  are linearly independent. We can partition  $X$  as:

$$X = Q(R_1, R_2) \quad (7.49)$$

where  $Q$  is an  $n \times r$  matrix with orthogonal columns ( $Q'Q = I_r$ ),  $R_1$  is an  $r \times r$  non-singular matrix, and  $R_2$  is  $r \times (p - r)$ .

The normal equations are:

$$X'X\beta = X'y \implies \begin{pmatrix} R_1' \\ R_2' \end{pmatrix} Q'Q(R_1, R_2)\beta = \begin{pmatrix} R_1' \\ R_2' \end{pmatrix} Q'y \quad (7.50)$$

Simplifying ( $Q'Q = I_r$ ):

$$\begin{pmatrix} R_1'R_1 & R_1'R_2 \\ R_2'R_1 & R_2'R_2 \end{pmatrix} \beta = \begin{pmatrix} R_1'Q'y \\ R_2'Q'y \end{pmatrix} \quad (7.51)$$

### 7.4.1 Constructing a Solution by Solving Normal Equations

One specific generalized inverse of  $X'X$  can be found by focusing on the non-singular block  $R_1'R_1$ :

$$(X'X)^- = \begin{pmatrix} (R_1'R_1)^{-1} & 0 \\ 0 & 0 \end{pmatrix} \quad (7.52)$$

Using this generalized inverse, the estimator  $\hat{\beta}$  becomes:

$$\hat{\beta} = (X'X)^- X'y = \begin{pmatrix} (R_1'R_1)^{-1} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} R_1'Q'y \\ R_2'Q'y \end{pmatrix} \quad (7.53)$$

$$\hat{\beta} = \begin{pmatrix} (R_1'R_1)^{-1} R_1'Q'y \\ 0 \end{pmatrix} = \begin{pmatrix} R_1^{-1}Q'y \\ 0 \end{pmatrix} \quad (7.54)$$

The fitted values are:

$$\hat{y} = X\hat{\beta} = Q(R_1, R_2) \begin{pmatrix} R_1^{-1}Q'y \\ 0 \end{pmatrix} = QR_1R_1^{-1}Q'y = QQ'y \quad (7.55)$$

This confirms that  $\hat{y}$  is the projection of  $y$  onto the column space of  $Q$  (which is the same as the column space of  $X$ ).

### 7.4.2 Constructing a Solution by Solving Reparametrized $\beta$

We can view the model as:

$$y = Q(R_1, R_2)\beta + \epsilon = Qb + \epsilon \quad (7.56)$$

where  $b = R_1\beta_1 + R_2\beta_2$ .

Since the columns of  $Q$  are orthogonal, the least squares estimate for  $b$  is simply:

$$\hat{b} = (Q'Q)^{-1}Q'y = Q'y \quad (7.57)$$

To find  $\beta$ , we solve the underdetermined system:

$$R_1\beta_1 + R_2\beta_2 = \hat{b} = Q'y \quad (7.58)$$

**Solution 1:** Set  $\beta_2 = 0$ . Then:

$$R_1\beta_1 = Q'y \implies \hat{\beta}_1 = R_1^{-1}Q'y \quad (7.59)$$

This yields the same result as the generalized inverse method above:  $\hat{\beta} = \begin{pmatrix} R_1^{-1}Q'y \\ 0 \end{pmatrix}$ .

**Solution 2:** Using the generalized inverse of  $R = (R_1, R_2)$ :

$$R^- = \begin{pmatrix} R_1^{-1} \\ 0 \end{pmatrix} \quad (7.60)$$

$$\hat{\beta} = R^-Q'y = \begin{pmatrix} R_1^{-1}Q'y \\ 0 \end{pmatrix} \quad (7.61)$$

This demonstrates that finding a solution to the normal equations using  $(X'X)^-$  is equivalent to solving the reparameterized system  $b = R\beta$ .