# Theory of Linear Models

Longhai Li
University of Saskatchewan

2026-01-13

# 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 inverseand 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.

# Introduction

## Multiple Linear Regression

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

$$y_{n \times 1} = X_{n \times p}\beta + \epsilon_{n \times 1}$$

where the error terms are distributed as:

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

in which $I_n$ is the identity matrix:

$$I_n = \begin{pmatrix} 1 & 0 & ... & 0 \\ 0 & 1 & ... & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & ... & 1 \end{pmatrix}$$

The scalar equation for a single observation is:

$$Y_i = \beta_0 + \beta_1 X_{i1} + ... + \beta_p X_{ip} + \epsilon_i$$

## Examples

### Polynomial Regression

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

The model equation is:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + ... + \beta_{p-1} x_i^{p-1}$$

### Design Matrix Construction

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

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

**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:

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

$$\boxed{\begin{matrix} Y_{11} \\ Y_{12} \end{matrix}}^{G_1} \quad \boxed{\begin{matrix} Y_{21} \\ Y_{22} \end{matrix}}^{G_2} \quad \boxed{\begin{matrix} Y_{31} \\ Y_{32} \end{matrix}}^{G_3}$$

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

$$y_{6\times1} = X_{6\times3} \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{pmatrix} + \epsilon$$

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

**Analysis of Covariance (ANCOVA)**

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

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

## Least Squares Estimation

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

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

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

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

The residuals and Sum of Squared Errors are:

$$\hat{e} = y - \hat{y}$$

$$\text{SSE} = \| \hat{e} \|^2$$

The coefficient of determination is:

$$R^2 = \frac{\text{SST} - \text{SSE}}{\text{SST}}$$

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

## 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 $\hat{y}_1$.
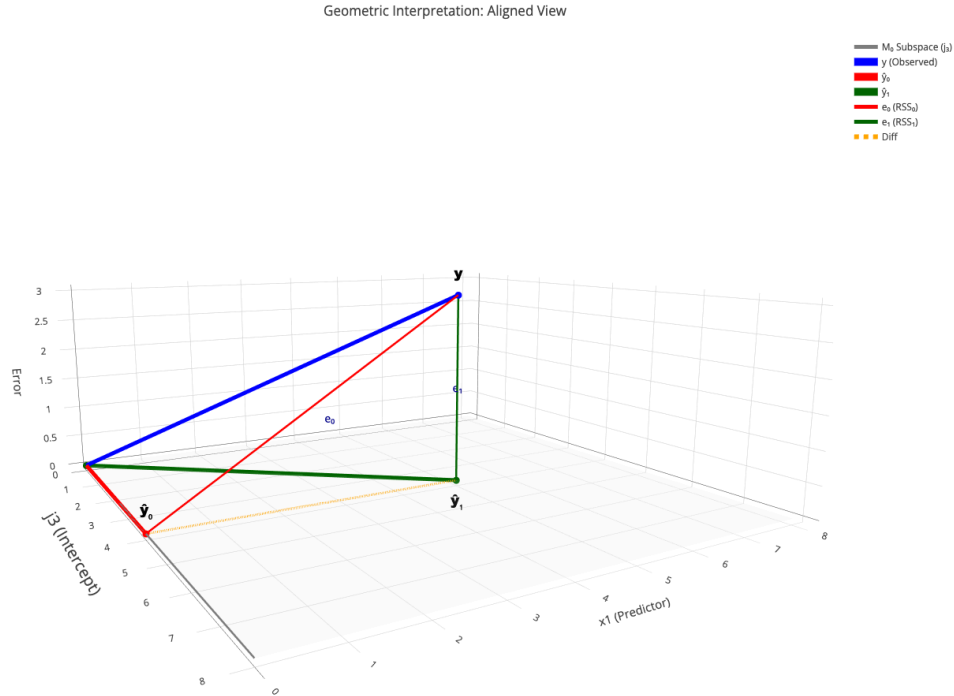


Figure 1: Geometric Interpretation: Projection onto Axis (M0) vs Plane (M1)

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.

## Projection in Vector Space

### Vector and Projection onto a Line

#### 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 0.1** (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:
>
> $$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

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 0.2** (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 $x$.

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

**Definition 0.3** (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 $x + d = y$.

**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 0.4** (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.

$$cx = \begin{pmatrix} cx_1 \\ \vdots \\ cx_n \end{pmatrix}$$

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

**Definition 0.5** (Euclidean Distance (Length)): The length (or norm) of a vector $x = (x_1, ..., x_n)^T$ 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:

$$\| x \|^2 = \sum_{i=1}^{n} x_i^2$$

$$\| x \| = \sqrt{\sum_{i=1}^{n} x_i^2}$$

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

$$c^2 = a^2 + b^2 - 2ab \cos \theta$$

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

$$\| y - x \|^2 = \| x \|^2 + \| y \|^2 - 2 \| x \| \cdot \| y \| \cos \theta$$

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:

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

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

$$= \parallel x \parallel^2 + \parallel y \parallel^2 - 2 \sum_{i=1}^{n} x_i y_i$$

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

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

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

$$x'y = \parallel x \parallel \cdot \parallel y \parallel \cos \theta$$

**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:

$$\parallel y \parallel \cos \theta = \frac{x'y}{\parallel x \parallel}$$

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

$$\text{Scalar Projection} = \left\langle \frac{x}{\parallel x \parallel}, y \right\rangle$$

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

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

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

This is often written compactly by combining the denominators:

$$\hat{y} = \frac{x'y}{\| x \|^2} x$$

**Perpendicularity (Orthogonality)**

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

**Definition 0.8** (Perpendicularity): Two vectors are defined as **perpendicular** (or orthogonal) if the angle between them is $90°$ ($\pi/2$).

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

$$x'y = 0 \Leftrightarrow x \perp y$$

**Example 0.1** (Orthogonal Vectors): Consider two vectors in $\mathbb{R}^2$: $x = (1,1)'$ and $y = (1,-1)'$.

$$x'y = 1(1) + 1(-1) = 1 - 1 = 0$$

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

**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 0.9** (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 $x$:

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

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 0.10** (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:

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

Expanding this inner product gives:

$$x'y - c(x'x) = 0$$

Solving for $c$, we obtain:

$$c = \frac{x'y}{\| x \|^2}$$

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

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

This second form separates the components into:

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

**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 $y$:

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

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

**Definition 0.12** (Projection Matrix onto a Single Vector): The matrix $P_x$ that projects any vector $y$ onto the line spanned by $x$ is defined as:

$$P_x = \frac{xx'}{\| x \|^2}$$

Using this matrix, the projection is simply:

$$\hat{y} = P_x y$$

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 0.2** (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:

$$x'y = 1(1) + 1(3) = 4$$

$$\| x \|^2 = 1^2 + 1^2 = 2$$

Now, apply the formula:

$$\hat{y} = \frac{4}{2}\begin{pmatrix} 1 \\ 1 \end{pmatrix} = 2\begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

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

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

Multiply by $y$:

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

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

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

$$\text{proj}\,(y \mid j_n) = \frac{j_n' y}{\parallel j_n \parallel^2} j_n$$

Calculating the components:

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

$$\parallel j_n \parallel^2 = \sum_{i=1}^{n} 1^2 = n$$

Substituting these back:

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

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.

**Pythagorean Theorem**

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

**Theorem 0.2** (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:

$$\parallel x + y \parallel^2 = \parallel x \parallel^2 + \parallel y \parallel^2$$

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

$$\parallel x + y \parallel^2 = (x + y)'(x + y)$$
$$= x'x + x'y + y'x + y'y$$
$$= \parallel x \parallel^2 + 2x'y + \parallel y \parallel^2$$

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

$$\parallel x + y \parallel^2 = \parallel x \parallel^2 + \parallel y \parallel^2$$

The proof after defining inner product to represent $\cos(\theta)$ is trivial. Figure 2 shows a geometric proof of the fundamental Pythagorean Theorem (aka 勾股定理).

$$S \;=\; S_1 \,+\, S_2$$

$$c^2 \cdot k \;=\; a^2 \cdot k \,+\, b^2 \cdot k$$

$$\Longrightarrow c^2 \;=\; a^2 \,+\, b^2$$

$$\text{where } c^2 = k \tfrac{b^2}{2} \sin\theta \cos\theta$$

Figure 2: Proof of Pythagorean Theorem using Area Scaling

**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 0.3** (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 $\hat{y}$.
>
> $$\| y - y^* \| \geq \| y - \hat{y} \|$$

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

$$(y - \hat{y}) \perp \left(\hat{y} - y^*\right)$$

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

$$y - y^* = (y - \hat{y}) + \left(\hat{y} - y^*\right)$$

Applying the Pythagorean Theorem:

$$\| y - y^* \|^2 = \| y - \hat{y} \|^2 + \| \hat{y} - y^* \|^2$$

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

$$\| y - y^* \|^2 \geq \| y - \hat{y} \|^2$$

## Vector Space

We now generalize our discussion from lines to broader spaces.

**Definition 0.13** (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 $c = 0$).

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

$$L(x_1, ..., x_p) = \{r \mid r = c_1 x_1 + ... + c_p x_p, \text{ for } c_i \in \mathbb{R}\}$$

**Column Space and Row Space**

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

**Definition 0.15** (Column Space): For a matrix $X = (x_1, ..., x_p)$, the **Column Space**, denoted Col $(X)$, is the vector space spanned by its columns:

$$\text{Col } (X) = L(x_1, ..., x_p)$$

**Definition 0.16** (Row Space): The **Row Space**, denoted Row $(X)$, is the vector space spanned by the rows of the matrix $X$.

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

$$\sum_{i=1}^{p} c_i x_i = 0 \implies c_1 = c_2 = ... = c_p = 0$$

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

## Rank of Matrices and Dim of Vector Space

**Definition 0.18** (Rank): The **Rank** of a matrix $X$, denoted Rank $(X)$, is the maximum number of linearly independent columns in $X$. This is equivalent to the dimension of the column space:

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

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

**Example 0.4** (Example of the Equality of Row and Col Rank): Consider the following $3 \times 4$ matrix ($n = 3, p = 4$):

$$X = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix}$$

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 $U$ ($2 \times 4$):

$$U = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix}$$

**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 $C$ ($3 \times 2$):

$$C = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}$$

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

$$\underbrace{\begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix}}_{X\,(3\times4)} = \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}}_{C\,(3\times2)} \underbrace{\begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix}}_{U\,(2\times4)}$$

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

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

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

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

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

**Theorem 0.4** (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.

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

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

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

## Orthogonality to a Subspace

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

**Definition 0.19** (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 $V$.

$$y \perp V \Leftrightarrow y'x = 0 \quad \forall x \in V$$

**Definition 0.20** (Orthogonal Complement): The set of all vectors that are orthogonal to a subspace $V$ is called the **Orthogonal Complement** of $V$, denoted $V^{\perp}$.

$$V^{\perp} = \{y \in \mathbb{R}^n \mid y \perp V\}$$

## 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 0.21** (Image and Kernel):

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

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

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

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

**Theorem 0.5** (Relationship between Kernel and Row Space): The kernel of $X$ is the orthogonal complement of the row space of $X$:

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

*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 $r_{1'}, ..., r_{n'}$, then the equation $Xx = 0$ is equivalent to the system of equations:

$$\begin{pmatrix} r_{1'} \\ \vdots \\ r_{n'} \end{pmatrix} x = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \Leftrightarrow r_{i'} x = 0 \ \text{ for all } \ i = 1, ..., n$$

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

**Nullity Theorem**

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

**Theorem 0.6** (Rank-Nullity Theorem): For an $n \times p$ matrix $X$:

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

where 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:

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

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:

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

Taking the dimensions of these spaces:

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

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

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

**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 0.7** (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 $B$.
>
> $$\text{Ker } (A) \subseteq \text{Ker } (B) \implies \text{Rank } (A) \geq \text{Rank } (B)$$

*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:

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

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

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

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

## Rank Inequalities

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

> **Theorem 0.8** (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:
>
> $$\text{Rank } (XZ) \leq \min(\text{Rank } (X), \text{Rank } (Z))$$

*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 $X$:

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

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 $Z$:

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

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

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

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

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

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

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

Thus, Rank $(AX) =$ Rank $(X)$.

**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 0.10** (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:

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

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

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

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

So Ker $(X) =$ 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.

**Column Space of $XX'$**

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

**Theorem 0.11** (Column Space Equivalence): The column space of $XX'$ is identical to the column space of $X$:

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

*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.

**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.

## 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:

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

Equivalently:

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

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

**Theorem 0.12** (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 $v \neq \hat{y}$:

$$\| y - \hat{y} \|^2 < \| y - v \|^2$$

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

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

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:

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

Applying the Pythagorean Theorem:

$$\| y - v \|^2 = \| y - \hat{y} \|^2 + \| \hat{y} - v \|^2$$

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

$$\| y - v \|^2 > \| y - \hat{y} \|^2$$

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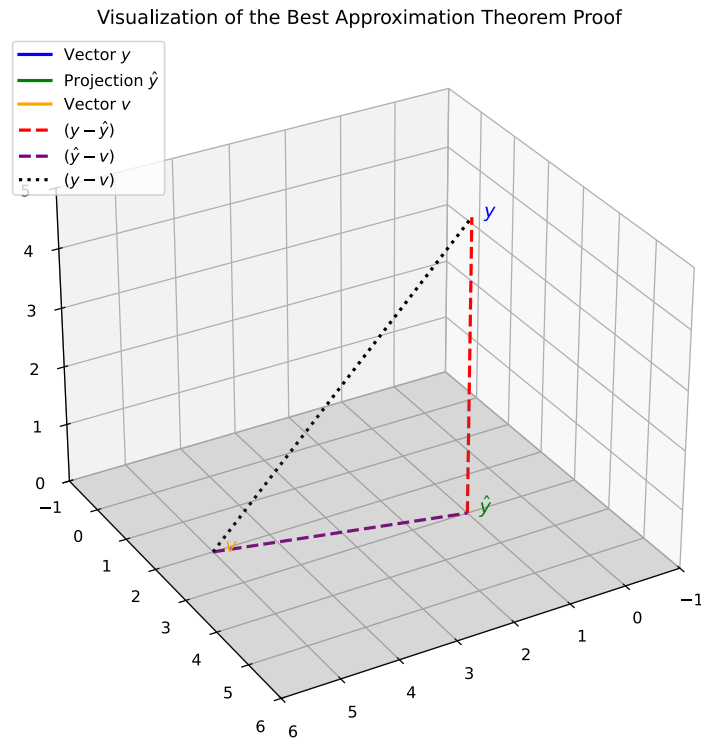


Figure 3: Visualization of the Best Approximation Theorem

**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 0.13** (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.

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

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

$$\langle y - \hat{y}_1, v \rangle = 0$$

$$\langle y - \hat{y}_2, v \rangle = 0$$

Subtracting the second equation from the first:

$$\langle y - \hat{y}_1, v \rangle - \langle y - \hat{y}_2, v \rangle = 0$$

Using the linearity of the inner product:

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

$$\langle \hat{y}_2 - \hat{y}_1, v \rangle = 0$$

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 $v = d = \hat{y}_2 - \hat{y}_1$.

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

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

$$\hat{y}_2 - \hat{y}_1 = 0 \implies \hat{y}_1 = \hat{y}_2$$

Thus, the projection is unique.

## Projection via Orthonormal Basis ($Q$)

**Orthonomal Basis**

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

**Definition 0.22** (Basis): A set of vectors $\{x_1, ..., x_k\}$ is a **Basis** for a vector space $V$ if:

1. The vectors span the space: $V = L(x_1, ..., x_k)$.
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 0.23** (Orthonormal Basis): A basis $\{q_1, ..., q_k\}$ is called an **Orthonormal Basis** if:

1. **Orthogonal:** Each pair of vectors is perpendicular.
$$q_{i'}q_j = 0 \quad \text{for} \quad i \neq j$$

2. **Normalized:** Each vector has unit length.
$$\| q_i \|^2 = q_{i'}q_i = 1$$

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 $\{q_1, ..., q_k\}$ for $V$, the projection $\hat{y}$ is simply the sum of the projections onto the individual basis vectors.

**Definition 0.24** (Projection Defined with Orthonormal Basis): The projection of $y$ onto the subspace $V = L(q_1, ..., q_k)$ is:
$$\hat{y} = \sum_{i=1}^{k} \text{proj}\ (y \mid q_i) = \sum_{i=1}^{k} (q_{i'}y)q_i$$

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

**Theorem 0.14** (Projection via Orthonormal Basis): Let $\{q_1, ..., q_k\}$ be an orthonormal basis for the subspace $V \subseteq \mathbb{R}^n$. The vector defined by the sum of individual projections:

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

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 $\{q_1, ..., q_k\}$.

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 $j = 1, ..., k$).

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

$$\langle y - \hat{y}, q_j \rangle = \langle y, q_j \rangle - \langle \hat{y}, q_j \rangle$$

$$= \langle y, q_j \rangle - \left\langle \sum_{i=1}^{k} \langle y, q_i \rangle q_i, q_j \right\rangle$$

$$= \langle y, q_j \rangle - \sum_{i=1}^{k} \langle y, q_i \rangle \underbrace{\langle q_i, q_j \rangle}_{\delta_{ij}}$$

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 $i = j$:

$$\langle y - \hat{y}, q_j \rangle = \langle y, q_j \rangle - \langle y, q_j \rangle \cdot 1$$
$$= 0$$

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.

**Projection Matrix via Orthonomal 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 $q_1, ..., q_k$.

$$Q = \begin{pmatrix} q_1 & q_2 & \cdots & q_k \end{pmatrix}$$

Properties of $Q$:

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

**Definition 0.25** (Projection Matrix in Terms of $Q$): The projection $\hat{y}$ can be written as:

$$\hat{y} = \begin{pmatrix} q_1 & \cdots & q_k \end{pmatrix} \begin{pmatrix} q_1{}'y \\ \vdots \\ q_{k'}y \end{pmatrix} = Q(Q'y) = (QQ')y$$

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

$$P = QQ'$$

**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 0.15** (Symmeticity 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 Col $(X)$, then projecting it again should not change it.

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

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

$$y_{ij} = \mu_i + \epsilon_{ij}$$

where $i \in \{1, ..., k\}$ represents the group and $j \in \{1, ..., n_i\}$ 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:

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

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

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

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

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

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:

$$\text{proj}\left(y \mid x_g\right) = \frac{\langle y, x_g \rangle}{\langle x_g, x_g \rangle} x_g$$

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:

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

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

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

**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 $x_1, ..., x_p$:

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

$$q_1 = \frac{x_1}{\| x_1 \|}$$

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

$$v_2 = x_2 - (x_2'q_1)q_1$$

   Then normalize:

$$q_2 = \frac{v_2}{\| v_2 \|}$$

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

$$v_k = x_k - \sum_{j=1}^{k-1}(x_k'q_j)q_j$$
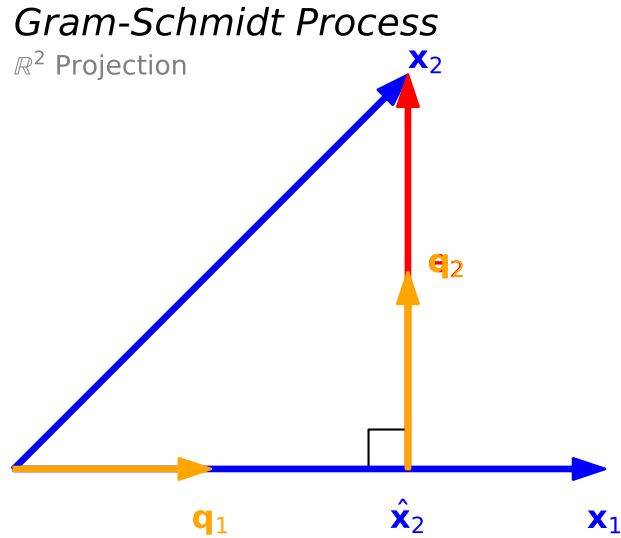
$$q_k = \frac{v_k}{\| v_k \|}$$

Figure 4: 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.

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

### Norm Equations

Let $X = (x_1, ..., x_p)$ 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 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 Col $(X)$.

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

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

$$y - X\beta \perp x_j \quad \text{for}\ \ j = 1, ..., p$$

Writing this geometric condition as an algebraic dot product (where $x_{j'}$ denotes the transpose):

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

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:

$$\begin{pmatrix} x_{1'} \\ \vdots \\ x_{p'} \end{pmatrix} (y - X\beta) = \mathbf{0} \implies X'(y - X\beta) = 0$$

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

$$X'y - X'X\beta = 0$$
$$X'X\beta = X'y$$

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

**Theorem 0.16** (Least Squares Estimator): If $X'X$ is invertible (i.e., $X$ has full column rank), the unique solution for $\beta$ is:

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

**Hat Matrix**

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

**Definition 0.26** (Hat Matrix): The projection of $y$ onto Col $(X)$ is given by:

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

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

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

**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 Col $(X)$, the formula simplifies significantly.

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

$$H = QR((QR)'(QR))^{-1}(QR)'$$
$$= QR(R'Q'QR)^{-1}R'Q'$$
$$= QR\left(R'\underbrace{Q'Q}_{I}R\right)^{-1}R'Q'$$
$$= QR(R'R)^{-1}R'Q'$$
$$= QRR^{-1}(R')^{-1}R'Q'$$
$$= Q\underbrace{RR^{-1}}_{I}\underbrace{(R')^{-1}R'}_{I}Q'$$
$$= QQ'$$

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

**Properties of Hat Matrix**

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

> **Theorem 0.17** (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.
>
> $$\text{tr}\,(H) = \text{tr}\,\left(X(X'X)^{-1}X'\right) = \text{tr}\,\left((X'X)^{-1}X'X\right) = \text{tr}\,\left(I_p\right) = p$$

## 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.

**Orthogonal Projection Matrix**

> **Definition 0.27** (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 0.18** (Projection onto Column Space): If a matrix $P$ is symmetric and idempotent, then $P$ represents the orthogonal projection onto its column space, Col $(P)$.

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

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

$$\langle z, (I - P)y \rangle = z^\top (I - P)y = (Px)^\top (I - P)y \tag{1}$$

Using the matrix transpose property $(AB)^\top = B^\top A^\top$, we rewrite Equation 1 as:

$$\langle z, (I - P)y \rangle = x^\top P^\top (I - P)y \tag{2}$$

Since $P$ is symmetric ($P^\top = P$), we can substitute $P$ for $P^\top$ in Equation 2:

$$\langle z, (I - P)y \rangle = x^\top P(I - P)y = x^\top (P - P^2)y \tag{3}$$

Finally, utilizing the idempotency of $P$ (where $P^2 = P$), the expression in Equation 3 simplifies to 0:

$$x^\top (P - P)y = x^\top (0)y = 0 \tag{4}$$

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

**Projection onto Complement Space**

**Theorem 0.19** (Projection onto Orthogonal Complement): Let $P$ be an orthogonal projection matrix. The matrix $M$ defined as:

$$M = I - P$$

is the orthogonal projection matrix onto the orthogonal complement of the column space of $P$, denoted 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 $M$:

$$M^\top = (I - P)^\top = I - P^\top = I - P = M \tag{5}$$

$$M^2 = (I - P)(I - P) = I - 2P + P^2 = I - 2P + P = I - P = M \tag{6}$$

By Equation 5 and Equation 6, $M$ is symmetric and idempotent, so it is an orthogonal projection matrix.

**2. Identifying the Subspace** By Theorem 0.18, $M$ projects onto its own column space, Col $(M)$. A vector $v$ is in Col $(M)$ if and only if it is fixed by the projection $(Mv = v)$.

$$Mv = v \tag{7}$$

Substituting $M = I - P$ into Equation 7 gives:

$$(I - P)v = v \tag{8}$$

Rearranging Equation 8, we find the condition for $v$:

$$v - Pv = v \implies Pv = 0 \tag{9}$$

The condition $Pv = 0$ in Equation 9 implies that $v$ belongs to the null space of $P$, denoted Null $(P)$. By the Fundamental Theorem of Linear Algebra for symmetric matrices, the null space is the orthogonal complement of the column space:

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

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

**Exercise 0.1** (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 $X$:

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

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

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

> **i Solutions**
>
> **1. Equivalence of Column Spaces** To prove Col $(H)$ = Col $(X)$, we show inclusion in both directions.
>
> - **Forward** (Col $(H) \subseteq$ Col $(X)$): By definition, $H = X\left[(X^\top X)^{-1} X^\top\right]$. 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 Col $(X)$.
>
> - **Reverse** (Col $(X) \subseteq$ Col $(H)$): Take any vector $v \in$ Col $(X)$. By definition, $v = Xb$ for some vector $b$. Apply $H$ to $v$:
>
> $$Hv = X(X^\top X)^{-1} X^\top (Xb) = X(X^\top X)^{-1}(X^\top X)b = X(I)b = Xb = v$$
>
> 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, Col $(H)$ = Col $(X)$.
>
> **2. Orthogonal Complements** From part 1, we know the subspaces are identical. Therefore, their orthogonal complements must also be identical:
>
> $$\text{Col } (H)^\perp = \text{Col } (X)^\perp$$
>
> We previously established in Theorem 0.19 that for any projection matrix $P$, the complement projection $M = I - P$ projects onto Col $(P)^\perp$. Substituting $H$ for $P$:
>
> $$\text{Col } (M) = \text{Col } (H)^\perp$$
>
> Combining these results gives the required equality:
>
> $$\text{Col } (M) = \text{Col } (X)^\perp$$

## Projection onto Nested Subspaces

**Nested Models and Subspaces**

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

**Definition 0.28** (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:

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

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

**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 0.20** (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 ($P_1 P_0 = P_0$).

$$(P_1 P_0)' = P_{0'}$$

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

$$P_{0'} P_{1'} = P_{0'} \implies P_0 P_1 = P_0$$

**Difference of Projections**

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

**Theorem 0.21** (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:**

$$(P_1 - P_0)^2 = (P_1 - P_0)(P_1 - P_0)$$
$$= P_1^2 - P_1 P_0 - P_0 P_1 + P_0^2$$

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

$$= P_1 - P_0 - P_0 + P_0 = P_1 - P_0$$

**3. Orthogonality to $P_0$:**

$$(P_1 - P_0)P_0 = P_1 P_0 - P_0^2 = P_0 - P_0 = 0$$

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 $M_0$.

**Decomposition of Projections and their Sum Squares**

**Theorem 0.22** (Orthogonal Decomposition): Let $M_0 \subset M_1$ be two nested linear models with corresponding design matrices $X_0$ and $X_1$ such that Col $(X_0) \subset$ Col $(X_1)$. Let $P_0$ and $P_1$ be the orthogonal projection matrices onto Col $(X_0)$ and Col $(X_1)$ respectively.

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

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

**Geometric Interpretation:**

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

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

$$\| y \|^2 = \| \hat{y}_0 \|^2 + \| \hat{y}_1 - \hat{y}_0 \|^2 + \| y - \hat{y}_1 \|^2$$

*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 Col $(X_1)$, the matrix $(I - P_1)$ projects onto the **orthogonal complement** Col $(X_1)^\perp$. Thus, $v_3 \in$ Col $(I - P_1)$.

Note that since Col $(X_0) \subset$ 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 $\langle v_1, v_2 \rangle = v_1' v_2$:

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

**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 $v_3$:

$$\hat{y}_{1}'v_3 = (P_1 y)'(I - P_1)y$$
$$= y'P_1(I - P_1)y$$
$$= y'(P_1 - P_1^2)y$$
$$= y'(P_1 - P_1)y$$
$$= 0$$

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.
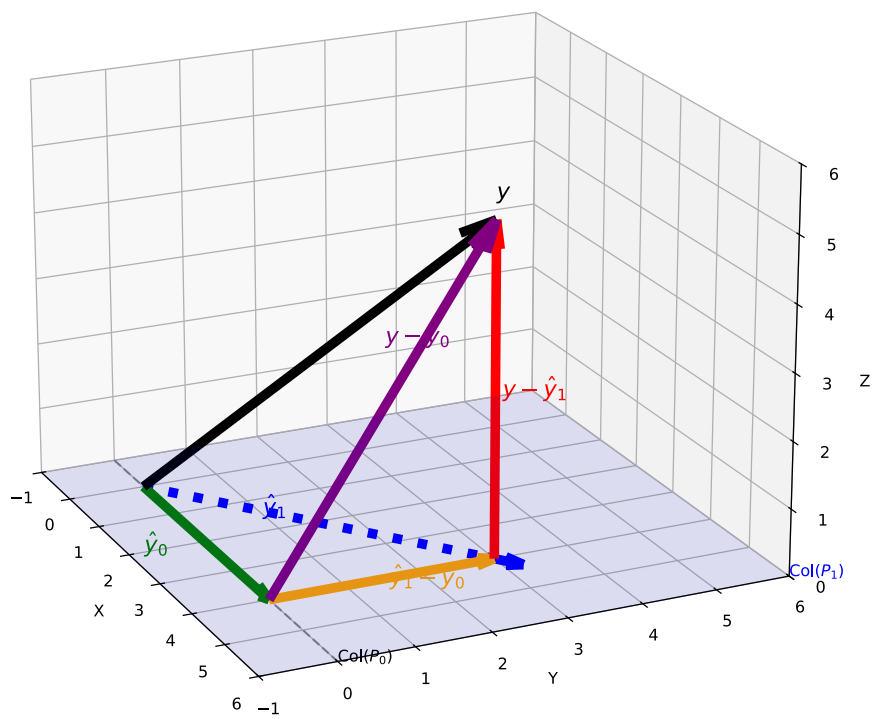
Figure 5: Illustration of Projections onto Nested Subspaces

**Example 0.6** (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 $i = 1, ..., k$ index the groups, and $j = 1, ..., n_i$ 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 $i$.

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

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

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

### 2. The Data and Projection Vectors

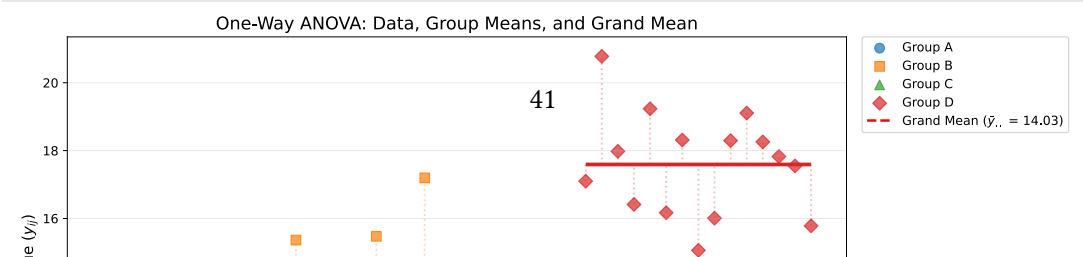Table 1: ANOVA Vectors: Data, Null Model, and Full Model

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

### 3. Decomposition and Sum of Squares

| Component | Notation | Definition | Vector Elements | Squared Norm (Sum of Squares) |
|---|---|---|---|---|
| **Null Proj.** | $\hat{y}_0$ | $P_0 y$ | Grand Mean $(\bar{y}_{..})$ | $\| \hat{y}_0 \|^2 = N\bar{y}_{..}^2$ |
| **Full Proj.** | $\hat{y}_1$ | $P_1 y$ | 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 $\| y \|^2 = \| \hat{y}_0 \|^2 + \| y - \hat{y}_0 \|^2$:



One-Way ANOVA: Data, Group Means, and Grand Mean

41

## Projections onto Orthogonal Subspaces

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

**Theorem 0.23** (General Orthogonal Projections): If $\mathbb{R}^n$ is the direct sum of orthogonal subspaces $V_1, V_2, ..., V_k$:

$$\mathbb{R}^n = V_1 \oplus V_2 \oplus ... \oplus V_k$$

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

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

$$y = \hat{y}_1 + \hat{y}_2 + ... + \hat{y}_k$$

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

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

$$\hat{y}_i = P_i y$$

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

$$y - \hat{y}_i = \sum_{j \neq i} \hat{y}_j$$

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

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

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, $\hat{y}_i = P_i y$.

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

$$I_n = P_1 + P_2 + ... + P_k$$

Orthogonal Decomposition of Vector y

$\mathbb{R}^n = V_1 \oplus V_2 \oplus V_3$
$y = P_1 y + P_2 y + P_3 y$
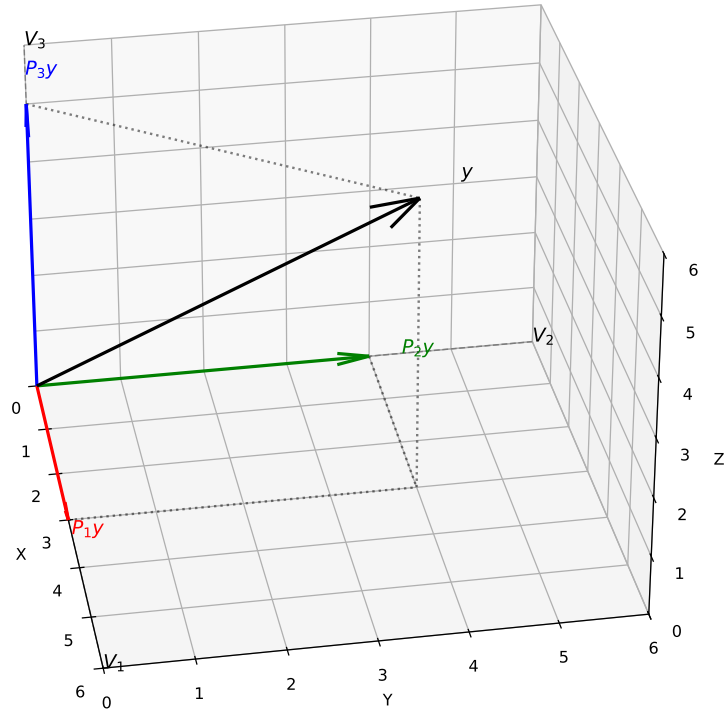$\|y\|^2 = \|P_1 y\|^2 + \|P_2 y\|^2 + \|P_3 y\|^2$

Figure 7: Orthogonal decomposition of vector y into subspaces

**Theorem 0.24** (Complete Orthogonal Decomposition of $\mathbb{R}^n$): Let $P_0, P_1, ..., P_k$ be a sequence of orthogonal projection matrices with nested column spaces:

$$\text{Col } (P_0) \subseteq \text{Col } (P_1) \subseteq ... \subseteq \text{Col } (P_k)$$

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

**Conclusion:**

1. **Projection Property:** Each $\Delta P_i$ is the orthogonal projection matrix onto $V_i$ for $i = 0, ..., k+1$.

2. **Mutual Orthogonality:** The collection $\{\Delta P_i\}$ are mutually orthogonal operators:

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

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

$$\mathbb{R}^n = V_0 \oplus V_1 \oplus ... \oplus V_{k+1}$$

*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., $1 \leq i < j \leq k$):

$$\left(P_j - P_{j-1}\right)(P_i - P_{i-1}) = P_j P_i - P_j P_{i-1} - P_{j-1} P_i + P_{j-1} P_{i-1}$$

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 $P_k \Delta P_i = \Delta P_i$.

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

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

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

$$= P_k + (I - P_k) = I$$

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 $V_i$.

**Residual** $V_{k+1}$

Figure 8: Venn Diagram of Nested Projections with Colored Increments