

Linear Regression: Reading Material

AERO 689: Introduction to Machine Learning for Aerospace Engineers

Dr. Raktim Bhattacharya

Table of contents

1	Introduction	5
1.1	Learning Objectives	6
2	Motivation: The Fuel Crisis Challenge	6
2.1	The \$100 Million Question	6
2.2	Real Impact	7
3	From Wind Tunnel to Flight: The Data Challenge	7
3.1	Traditional Approach: Empirical Models	7
3.2	Available Data Sources	8
3.2.1	1. Wind Tunnel Data	8
3.2.2	2. Flight Test Data	8
3.2.3	3. Operational Data	8
3.3	The Linear Regression Framework	9
4	Mathematical Foundation: The Linear Model	9
4.1	General Form	9
4.2	Vector Notation	10
5	Matrix Formulation	10
5.1	Design Matrix Structure	10
5.2	Aerospace Example: Drag Prediction	11
6	Key Insight: What “Linear” Means	12
6.1	Linear in Parameters, Not Features	12
6.2	Why Linearity in Parameters Matters	12
6.3	Aerospace Example	13
7	Real Data: Noisy Measurements	13
7.1	Understanding Measurement Noise	13

7.2	Characteristics of Real Wind Tunnel Data	14
7.3	Why the Error Term ϵ_i is Essential	14
8	The Optimization Problem	14
8.1	Objective: Minimize Squared Error	14
8.2	Why Square the Errors?	15
8.3	Matrix Form	15
8.4	The Optimization Goal	15
9	Derivation: Normal Equations	16
9.1	Step 1: Expand the Objective Function	16
9.2	Step 2: Take Derivative with Respect to β	16
9.3	Step 3: Set to Zero and Solve	16
9.4	Step 4: Closed-Form Solution	17
10	Geometric Interpretation: Understanding the Error	17
10.1	What is the Residual Vector?	17
10.2	Goal: Minimize the Length of the Error Vector	17
10.3	Column Space of X	18
10.4	The Fundamental Geometric Principle	18
11	Why Projection Minimizes Error	18
11.1	Visualizing the Projection	18
11.2	Why Perpendicular is Optimal	19
11.3	Mathematical Statement of Orthogonality	19
12	Deriving the Optimal Solution via Projection	19
12.1	Step 1: State the Orthogonality Condition	19
12.2	Step 2: Express in Terms of β	20
12.3	Step 3: Derive the Normal Equations	20
12.4	Step 4: Solve for β^*	20
13	The Projection Matrix	21
13.1	Definition	21
13.2	Computing Predictions	21
13.3	Key Properties	21
13.3.1	1. Idempotent (Projecting Twice = Projecting Once)	21
13.3.2	2. Symmetric	21
13.3.3	3. Projects onto $\text{col}(X)$	22
13.3.4	4. Residual Matrix	22
13.4	Aerospace Interpretation	22
14	When Direct Solution Fails	22
14.1	Problem 1: Singular Matrix (Non-Invertibility)	23

14.2 Problem 2: Computational Cost	23
14.3 Problem 3: Numerical Stability	24
14.4 Solutions	24
14.4.1 1. Regularization	24
14.4.2 2. Gradient Descent	24
14.4.3 3. QR Decomposition	24
14.4.4 4. Singular Value Decomposition (SVD)	25
15 Gradient Descent: Iterative Approach	25
15.1 The Algorithm	25
15.1.1 Steps	25
15.2 Computing the Gradient	26
15.3 Convergence	26
16 Gradient Descent Variants	26
16.1 Batch Gradient Descent	26
16.2 Stochastic Gradient Descent (SGD)	27
16.3 Mini-Batch Gradient Descent	27
16.4 Comparison Summary	28
17 Learning Rate Selection	28
17.1 The Role of Learning Rate	28
17.2 Too Small Learning Rate ($\eta \ll 1$)	29
17.3 Too Large Learning Rate ($\eta \gg 1$)	29
17.4 Finding the Right Learning Rate	29
17.5 Adaptive Learning Rates	29
17.5.1 1. Learning Rate Decay	30
17.5.2 2. Momentum	30
17.5.3 3. Adam (Adaptive Moment Estimation)	30
17.5.4 4. Line Search	31
18 Statistical Properties: Assumptions	31
18.1 Classical Linear Regression Assumptions	31
18.1.1 1. Linearity	31
18.1.2 2. Independence	31
18.1.3 3. Homoscedasticity (Constant Variance)	32
18.1.4 4. Normality	32
18.1.5 5. No Perfect Multicollinearity	33
19 Gauss-Markov Theorem: Implications for Practice	34
19.1 Statement of the Theorem	34
19.2 What Does BLUE Mean?	34
19.2.1 Best	34

19.2.2	Linear	34
19.2.3	Unbiased	35
19.3	Aerospace Context	35
19.4	When BLUE Doesn't Apply	35
20	Statistical Properties: Distribution	36
20.1	Distribution of the OLS Estimator	36
20.2	Proving Unbiasedness	36
20.3	Deriving the Covariance Matrix	37
20.4	Unpacking the Covariance Matrix	37
20.4.1	Diagonal Elements: Individual Coefficient Variance	38
20.4.2	Off-Diagonal Elements: Correlation Between Coefficients	38
20.4.3	Variance Inflation Factor (VIF)	38
20.5	Practical Implications	39
20.6	What If Normality Doesn't Hold?	39
20.7	The Role of the Central Limit Theorem	40
21	Estimating the Noise Variance	40
21.1	Residual Variance Estimator	40
21.2	Why Divide by $n - d - 1$?	41
21.2.1	Degrees of Freedom	41
21.2.2	Why Not Divide by n ?	41
21.2.3	Distribution of $\hat{\sigma}^2$	42
21.3	Aerospace Example	42
21.4	Relationship to Model Fit	43
21.5	Residual Standard Error vs. R-squared	44
22	Confidence Intervals for Coefficients	44
22.1	Understanding Coefficient Uncertainty	44
22.2	Standard Error	45
22.2.1	Factors Affecting Standard Error	45
22.2.2	Computing Standard Errors in Practice	46
22.3	Confidence Interval Formula	46
22.3.1	Why Use the t-Distribution?	46
22.3.2	Common Confidence Levels	47
22.4	Interpretation	47
22.4.1	Practical Interpretation for Engineers	48
22.5	Aerospace Example	48
22.5.1	Impact of Sample Size	49
23	Prediction Intervals vs. Confidence Intervals	49
23.1	The Two Types of Uncertainty	49
23.1.1	1. Parameter Uncertainty (Epistemic)	49

23.1.2 2. Irreducible Noise (Aleatoric)	50
23.2 Key Differences	50
23.3 Aerospace Example: Drag Prediction	51
23.4 Extrapolation Warning	52
24 Hypothesis Testing for Individual Coefficients	52
24.1 The Central Question	52
24.2 Setting Up the Test	52
24.2.1 Hypotheses	52
24.3 The Test Statistic	53
24.4 Distribution Under the Null	53
24.5 Making the Decision	53
24.5.1 Approach 1: P-value	53
24.5.2 Approach 2: Critical Value	54
24.6 Aerospace Example	54
24.7 Important Caveats	55
25 F-Test for Overall Model Significance	55
25.1 The Hypotheses	55
25.2 Decomposing Variance	55
25.3 The F-Statistic	56
25.4 Distribution Under the Null	56
25.5 Making the Decision	57
25.6 Relationship to R-squared	57
25.7 Aerospace Example	57
25.8 F-Test vs. Multiple t-Tests	58
25.9 F-Test for Model Comparison	59
26 Summary	59
26.1 Mathematical Foundations	59
26.2 Two Solution Approaches	60
26.3 Geometric Insight	60
26.4 Statistical Properties	60
26.5 Practical Considerations	60
26.6 Aerospace Applications	60

1 Introduction

Linear regression is one of the most fundamental and widely used techniques in machine learning and statistical modeling. In aerospace engineering, it plays a crucial role in predicting aircraft performance, analyzing wind tunnel data, and optimizing flight operations. This reading

material provides detailed explanations of the concepts presented in the lecture slides, with a focus on both mathematical rigor and practical aerospace applications.

1.1 Learning Objectives

By the end of this module, you should be able to:

1. **Apply linear regression to aerospace performance prediction:** Understand how to use linear regression to model relationships between flight parameters (e.g., angle of attack, Mach number) and performance metrics (e.g., drag coefficient, fuel consumption).
2. **Understand least squares method and gradient descent:** Master both the analytical (closed-form) and iterative (gradient descent) approaches to solving linear regression problems.
3. **Implement drag coefficient prediction from wind tunnel data:** Learn to process real experimental data and build predictive models that account for measurement noise and physical constraints.
4. **Validate models using aerospace-specific metrics:** Understand how to assess model quality using appropriate statistical measures and ensure predictions meet safety and certification requirements.

2 Motivation: The Fuel Crisis Challenge

2.1 The \$100 Million Question

Consider a real-world scenario that aerospace engineers face daily: an airline operates a fleet of 200 aircraft. The fuel cost alone exceeds \$50 million annually per aircraft type. With such enormous operational costs, even small improvements in fuel efficiency can translate to massive savings.

The Challenge: Airlines need to accurately predict fuel consumption for flight planning. Currently, most operations rely on simplified performance charts that were developed under idealized conditions. These charts provide conservative estimates but may not capture the full complexity of real-world flight operations.

The Machine Learning Opportunity: By using actual flight data collected from thousands of flights, we can build precise models that account for:

- Actual weather conditions encountered
- Real payload and weight variations
- Engine performance degradation over time
- Pilot technique variations

- Air traffic control routing constraints

2.2 Real Impact

The potential benefits of improved fuel consumption models are substantial:

- **1% fuel savings** = Over \$100 million annually across the industry
- **Better range predictions** = More efficient route planning and optimization
- **Accurate payload calculations** = Improved safety margins while maximizing revenue cargo capacity

Discussion Question: What factors do you think affect aircraft fuel consumption? Consider environmental factors (weather, altitude, temperature), operational factors (weight, speed, routing), and mechanical factors (engine condition, aerodynamic efficiency).

3 From Wind Tunnel to Flight: The Data Challenge

3.1 Traditional Approach: Empirical Models

Aerospace engineers have long used empirical models to characterize aircraft performance. One classic example is the **parabolic drag polar**:

$$C_D = C_{D_0} + KC_L^2$$

where:

- C_D is the drag coefficient
- C_{D_0} is the zero-lift drag coefficient (parasitic drag)
- K is the induced drag factor
- C_L is the lift coefficient

The Problem: This model assumes perfectly controlled conditions—smooth flow, steady state, clean configuration. In reality:

- Real flights encounter turbulence, wind shear, and varying atmospheric conditions
- Aircraft weight changes continuously as fuel is consumed
- Engines degrade over time, affecting performance
- Manufacturing tolerances mean each aircraft is slightly different

The Solution: Machine learning allows us to learn patterns directly from operational data, capturing complexities that simplified analytical models may miss.

3.2 Available Data Sources

Aerospace engineers can draw from three main sources of data:

3.2.1 1. Wind Tunnel Data

Characteristics:

- Highly controlled environment
- Precise measurements
- Limited range of conditions
- Expensive to collect
- May not capture full-scale Reynolds number effects

Best for: Understanding fundamental aerodynamic behavior, validating CFD simulations

3.2.2 2. Flight Test Data

Characteristics:

- Real flight conditions
- Instrumented aircraft
- Very expensive to collect (dedicated test aircraft, crew, facilities)
- Limited sample size
- High quality, well-documented

Best for: Aircraft certification, validating performance predictions, boundary exploration

3.2.3 3. Operational Data

Characteristics:

- Massive scale (thousands or millions of flights)
- Representative of actual operations
- Noisy (sensor errors, environmental variations)
- May lack detailed instrumentation
- Continuously collected

Best for: Statistical modeling, fleet-wide trends, operational optimization

3.3 The Linear Regression Framework

For this module, we'll focus on a specific problem: predicting the drag coefficient from flight parameters.

Goal: Predict C_D (drag coefficient) accurately

Input Features:

- α (angle of attack)
- M (Mach number)
- Re (Reynolds number)
- Potentially interaction terms and polynomial features

Output: Drag coefficient value for performance calculations

Why is this important? Accurate drag prediction is essential for:

- Fuel consumption estimation
- Range calculations
- Climb performance
- Flight envelope determination
- Control system design

4 Mathematical Foundation: The Linear Model

4.1 General Form

Linear regression models the relationship between input features and a continuous output variable. For a dataset with n samples and d features, the model is:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_d x_{id} + \epsilon_i$$

Notation Explained:

- y_i : The response variable (dependent variable) for sample i . In aerospace: this could be drag coefficient, fuel flow rate, or any measurable output.
- x_{ij} : The j -th feature (independent variable) of the i -th sample. Examples: Mach number, angle of attack, altitude, etc.
- β_j : Regression coefficients (parameters) that we need to learn from data. These quantify how each feature affects the response.
- β_0 : The intercept term. This represents the baseline value when all features are zero.

- ϵ_i : The error term (residual). This captures:

- Measurement noise
- Unmodeled physics
- Random variations
- Model approximation errors

4.2 Vector Notation

To work with all samples simultaneously, we use matrix-vector notation:

$$y = X\beta + \epsilon$$

where:

- $y \in \mathbb{R}^n$: Vector of all n response values
- $X \in \mathbb{R}^{n \times (d+1)}$: Design matrix (also called feature matrix)
- $\beta \in \mathbb{R}^{d+1}$: Vector of true parameters (including intercept)
- $\epsilon \in \mathbb{R}^n$: Vector of all error terms

The design matrix X is augmented with a column of ones to account for the intercept term.

5 Matrix Formulation

5.1 Design Matrix Structure

The design matrix organizes all our data into a structured format:

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1d} \\ 1 & x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix}$$

Understanding the structure:

- Each **row** represents one data sample (one observation, one flight, one wind tunnel measurement)
- Each **column** (except the first) represents one feature across all samples
- The **first column** of all ones corresponds to the intercept term β_0
- Dimensions: n rows \times $(d + 1)$ columns

The parameter vector and response vector are:

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

5.2 Aerospace Example: Drag Prediction

Let's make this concrete with a drag coefficient prediction problem.

Scenario: We have wind tunnel data measuring drag coefficient at various angles of attack and Mach numbers. We want to build a model that includes:

- Linear terms: α, M, Re
- Quadratic terms: α^2, M^2
- Interaction terms: $\alpha M, \alpha Re$, etc.

Design Matrix:

$$X = \begin{bmatrix} 1 & \alpha_1 & M_1 & Re_1 & \alpha_1^2 & M_1^2 & \alpha_1 M_1 & \cdots \\ 1 & \alpha_2 & M_2 & Re_2 & \alpha_2^2 & M_2^2 & \alpha_2 M_2 & \cdots \\ \vdots & \ddots \end{bmatrix}$$

Parameter Vector:

$$\beta = \begin{bmatrix} C_{D_0} \\ k_\alpha \\ k_M \\ k_{Re} \\ k_{\alpha^2} \\ k_{M^2} \\ k_{\alpha M} \\ \cdots \end{bmatrix}$$

Target Vector:

$$y = \begin{bmatrix} C_{D_1} \\ C_{D_2} \\ \cdots \\ C_{D_n} \end{bmatrix}$$

This formulation allows us to capture complex aerodynamic behavior while maintaining the linear regression framework.

6 Key Insight: What “Linear” Means

6.1 Linear in Parameters, Not Features

A common source of confusion: “linear regression” refers to linearity **in the parameters** β , **not** in the input features x .

The General Form:

$$y = \beta_0\phi_0(x) + \beta_1\phi_1(x) + \cdots + \beta_d\phi_d(x)$$

where $\phi_j(x)$ are **basis functions** that can be:

- Linear: $\phi_1(x) = x$
- Polynomial: $\phi_2(x) = x^2$, $\phi_3(x) = x^3$
- Trigonometric: $\phi_4(x) = \sin(x)$, $\phi_5(x) = \cos(x)$
- Exponential: $\phi_6(x) = e^x$
- Interactions: $\phi_7(x_1, x_2) = x_1x_2$
- Any other nonlinear transformation

Key Properties:

1. **Basis functions** $\phi_j(x)$ are **fixed** (you choose them before fitting)
2. **Coefficients** β_j are **what we solve for** (learned from data)
3. The model is **linear** in β_j because each β_j appears with power 1 and no products of different β_j terms
4. The model can be **nonlinear** in x because the basis functions can transform inputs arbitrarily

6.2 Why Linearity in Parameters Matters

Mathematical Benefit: Linearity in parameters means:

- The optimization problem is **convex** (has exactly one global minimum, no local minima)
- A **closed-form solution** exists (we can write down the exact answer)
- The solution is **unique** (if $X^T X$ is invertible)
- We can use **efficient linear algebra** algorithms

Practical Benefit: We can model complex, nonlinear physical phenomena while maintaining:

- Computational efficiency
- Guaranteed convergence
- Interpretable coefficients
- Statistical properties (confidence intervals, hypothesis tests)

6.3 Aerospace Example

Consider modeling the drag coefficient with this equation:

$$C_D = \beta_0 + \beta_1\alpha + \beta_2\alpha^2 + \beta_3M^2 + \beta_4(\alpha M)$$

Analysis:

- **Nonlinear in physical variables:** The relationship between C_D and (α, M) is nonlinear due to the α^2 , M^2 , and αM terms. This is a parabolic surface in the (α, M) space.
- **Linear in parameters:** The equation is a weighted sum of basis functions:

$$C_D = \beta_0 \cdot 1 + \beta_1 \cdot \alpha + \beta_2 \cdot \alpha^2 + \beta_3 \cdot M^2 + \beta_4 \cdot (\alpha M)$$

- **Effect of doubling β_2 :** If we double β_2 , the contribution of the α^2 term exactly doubles. This linear relationship in the parameters is what makes the problem tractable.

Matrix Representation:

$$X = \begin{bmatrix} 1 & \alpha_1 & \alpha_1^2 & M_1^2 & \alpha_1 M_1 \\ 1 & \alpha_2 & \alpha_2^2 & M_2^2 & \alpha_2 M_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

Each row processes one data point through all basis functions, creating a standard linear regression problem.

7 Real Data: Noisy Measurements

7.1 Understanding Measurement Noise

In practice, all experimental data contains measurement errors. For aerospace applications, these errors arise from:

1. **Sensor limitations:** Finite precision, drift, calibration errors
2. **Environmental factors:** Temperature variations, atmospheric turbulence
3. **Flow unsteadiness:** Turbulent fluctuations, vortex shedding
4. **Model simplification:** Real physics is more complex than our model
5. **Human factors:** Data recording errors, experimental setup variations

7.2 Characteristics of Real Wind Tunnel Data

When measuring drag coefficients in a wind tunnel, we typically observe:

- **Data scatter** around the true relationship
- **Heteroscedastic noise:** Measurement uncertainty often increases with angle of attack (larger forces → larger absolute errors)
- **Systematic biases:** Wall effects, support interference
- **Outliers:** Occasionally anomalous measurements due to flow separation or experimental issues

7.3 Why the Error Term ϵ_i is Essential

The error term in our model $y_i = x_i^T \beta^* + \epsilon_i$ is not just a mathematical convenience—it's a fundamental recognition that:

1. **No model is perfect:** Even the best physical model cannot capture every detail
2. **Measurements are imperfect:** Sensors have inherent limitations
3. **Random variations exist:** Physical processes have inherent stochasticity
4. **We seek expected behavior:** Our goal is to find the average relationship, not fit every noise fluctuation

Important: We want our model to capture the **signal** (true underlying relationship) without overfitting to the **noise** (random fluctuations). This is the essence of good statistical modeling.

8 The Optimization Problem

8.1 Objective: Minimize Squared Error

Given data $\{(x_i, y_i)\}_{i=1}^n$, we want to find parameters β that make our predictions $\hat{y}_i = x_i^T \beta$ as close as possible to the observed values y_i .

Residual for sample i :

$$r_i = y_i - \hat{y}_i = y_i - x_i^T \beta$$

Residual Sum of Squares (RSS):

$$\text{RSS}(\beta) = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - x_i^T \beta)^2$$

8.2 Why Square the Errors?

We could use different error metrics. Why do we square the errors?

1. **Penalizes large errors more:** An error of 2 contributes 4 to RSS, while two errors of 1 each contribute only 2. This makes the model sensitive to outliers.
2. **Mathematical tractability:** Squared errors lead to a smooth, differentiable objective function with a unique minimum.
3. **Statistical optimality:** Under Gaussian noise, least squares is the maximum likelihood estimator.
4. **Geometric interpretation:** Minimizing RSS is equivalent to finding the projection onto the column space of X .
5. **No sign cancellation:** Unlike simple sum of errors (where positive and negative errors cancel), squaring ensures all errors contribute positively.

8.3 Matrix Form

In matrix-vector notation:

$$\text{RSS}(\beta) = \|y - X\beta\|^2 = (y - X\beta)^T(y - X\beta)$$

This is the squared Euclidean norm of the residual vector $r = y - X\beta$.

8.4 The Optimization Goal

Our objective is to find the parameters that minimize RSS:

$$\beta^* = \arg \min_{\beta} \text{RSS}(\beta) = \arg \min_{\beta} \|y - X\beta\|^2$$

Physical Interpretation: In the context of aircraft performance modeling, we're finding the model parameters that minimize the total squared difference between predicted and actual drag coefficients across all wind tunnel measurements. This gives us the “best fit” in a least-squares sense.

9 Derivation: Normal Equations

We can derive the optimal solution using calculus. The approach is to expand the objective function, take the derivative with respect to β , set it to zero, and solve.

9.1 Step 1: Expand the Objective Function

Starting with:

$$\text{RSS}(\beta) = (y - X\beta)^T(y - X\beta)$$

Expand:

$$\text{RSS}(\beta) = y^T y - y^T X \beta - \beta^T X^T y + \beta^T X^T X \beta$$

Since $y^T X \beta$ is a scalar, it equals its transpose: $y^T X \beta = \beta^T X^T y$

Therefore:

$$\text{RSS}(\beta) = y^T y - 2\beta^T X^T y + \beta^T X^T X \beta$$

9.2 Step 2: Take Derivative with Respect to β

Using matrix calculus identities: $-\frac{\partial}{\partial \beta}(a^T \beta) = a - \frac{\partial}{\partial \beta}(\beta^T A \beta) = 2A\beta$ (when A is symmetric)

We get:

$$\frac{\partial \text{RSS}}{\partial \beta} = -2X^T y + 2X^T X \beta$$

9.3 Step 3: Set to Zero and Solve

At the minimum, the gradient must be zero:

$$-2X^T y + 2X^T X \beta^* = 0$$

Dividing by 2 and rearranging:

$$X^T X \beta^* = X^T y$$

These are called the **Normal Equations**.

9.4 Step 4: Closed-Form Solution

If $X^T X$ is invertible (which requires that the columns of X are linearly independent), we can solve for β^* :

$$\beta^* = (X^T X)^{-1} X^T y$$

This is the **ordinary least squares (OLS) solution**.

Verification that this is a minimum: The second derivative (Hessian matrix) is:

$$\frac{\partial^2 \text{RSS}}{\partial \beta^2} = 2X^T X$$

This is positive definite (assuming X has full column rank), confirming that we have found a minimum, not a maximum or saddle point.

10 Geometric Interpretation: Understanding the Error

10.1 What is the Residual Vector?

The residual (error) vector is defined as:

$$r = y - \hat{y} = y - X\beta$$

This is the vector of differences between:

- **Observed data y :** What we actually measured
- **Predictions $\hat{y} = X\beta$:** What our model predicts

10.2 Goal: Minimize the Length of the Error Vector

Our optimization problem can be stated as:

$$\min_{\beta} \|r\|^2 = \min_{\beta} \|y - X\beta\|^2$$

We want to make the residual vector as short as possible.

10.3 Column Space of X

Definition: The column space of X , denoted $\text{col}(X)$, is the set of all possible linear combinations of the columns of X :

$$\text{col}(X) = \{X\beta : \beta \in \mathbb{R}^{d+1}\}$$

Key Insight:

- Every prediction $\hat{y} = X\beta$ **must lie in** $\text{col}(X)$
- The observed data y typically does **not** lie in $\text{col}(X)$ (due to measurement noise)
- We seek the **best approximation** to y within $\text{col}(X)$

10.4 The Fundamental Geometric Principle

Question: What is the closest point in $\text{col}(X)$ to y ?

Answer: The **orthogonal projection** of y onto $\text{col}(X)$.

Why? The **Projection Theorem** from linear algebra states: The shortest distance from a point to a subspace is achieved by the perpendicular projection onto that subspace.

Mathematical Statement: The error is minimized when it is orthogonal to the column space:

$$r \perp \text{col}(X) \iff X^T r = 0$$

This orthogonality condition is the **geometric essence** of least squares.

11 Why Projection Minimizes Error

11.1 Visualizing the Projection

Imagine a 2D plane (representing $\text{col}(X)$) embedded in 3D space (representing \mathbb{R}^n). The data vector y is a point not on this plane.

Key Observations:

1. Any prediction \hat{y} must lie on the plane (in $\text{col}(X)$)
2. The error $r = y - \hat{y}$ is the vector from \hat{y} to y
3. We want to minimize $\|r\|$ (the length of this vector)

11.2 Why Perpendicular is Optimal

Consider any two points in $\text{col}(X)$: - \hat{y}_\perp : The perpendicular projection - \hat{y}_{other} : Any other point

By the **Pythagorean theorem**:

$$\|y - \hat{y}_{\text{other}}\|^2 = \|y - \hat{y}_\perp\|^2 + \|\hat{y}_\perp - \hat{y}_{\text{other}}\|^2$$

Since $\|\hat{y}_\perp - \hat{y}_{\text{other}}\|^2 \geq 0$, we have:

$$\|y - \hat{y}_{\text{other}}\|^2 \geq \|y - \hat{y}_\perp\|^2$$

Therefore, the perpendicular projection gives the smallest error.

11.3 Mathematical Statement of Orthogonality

The optimal prediction \hat{y}^* satisfies:

$$(y - \hat{y}^*) \perp \text{col}(X)$$

In matrix form:

$$X^T(y - \hat{y}^*) = 0$$

Understanding $X^T r = 0$:

- $X^T r$ is a vector of dot products: $[x_1^T r, x_2^T r, \dots, x_{d+1}^T r]^T$
- Each dot product equals zero: $x_j^T r = 0$
- This means r is perpendicular to **every column** of X
- Therefore, r is perpendicular to the **entire column space**

12 Deriving the Optimal Solution via Projection

Rather than using calculus, we can derive the OLS solution directly from the geometric principle of orthogonal projection.

12.1 Step 1: State the Orthogonality Condition

From geometry, the error must be perpendicular to $\text{col}(X)$:

$$r \perp \text{col}(X) \implies X^T r = 0$$

This is **the fundamental condition** for least squares optimality.

12.2 Step 2: Express in Terms of β

The optimal prediction is $\hat{y}^* = X\beta^*$ and the optimal residual is $r^* = y - \hat{y}^*$.

Substituting into the orthogonality condition:

$$X^T r^* = 0$$

$$X^T(y - \hat{y}^*) = 0$$

$$X^T(y - X\beta^*) = 0$$

12.3 Step 3: Derive the Normal Equations

Expanding:

$$X^T y - X^T X \beta^* = 0$$

Rearranging:

$$X^T X \beta^* = X^T y$$

These are the **Normal Equations**—exactly the same result we obtained using calculus!

12.4 Step 4: Solve for β^*

Assuming $X^T X$ is invertible:

$$\beta^* = (X^T X)^{-1} X^T y$$

Key Insight: We derived the OLS solution using **projection geometry** (orthogonality condition) instead of **calculus** (setting derivative to zero). Both paths lead to the same answer, providing two complementary perspectives:

- **Calculus view:** β^* minimizes the cost function
- **Geometric view:** β^* gives the perpendicular projection

13 The Projection Matrix

13.1 Definition

The **projection matrix** is defined as:

$$P = X(X^T X)^{-1} X^T$$

This matrix projects any vector in \mathbb{R}^n onto $\text{col}(X)$.

13.2 Computing Predictions

Using the projection matrix, we can write the predicted values as:

$$\hat{y} = X\beta^* = X(X^T X)^{-1} X^T y = Py$$

Interpretation: P directly maps observed data y to predictions \hat{y} , bypassing the need to explicitly compute β^* .

13.3 Key Properties

13.3.1 1. Idempotent (Projecting Twice = Projecting Once)

$$P^2 = P$$

Proof:

$$P^2 = X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = P$$

Interpretation: If a vector is already in $\text{col}(X)$, projecting it again doesn't change it.

13.3.2 2. Symmetric

$$P^T = P$$

Proof: Each factor is symmetric:

$$P^T = (X(X^T X)^{-1} X^T)^T = X((X^T X)^{-1})^T X^T = X(X^T X)^{-1} X^T = P$$

13.3.3 3. Projects onto $\text{col}(X)$

$$PX = X$$

Proof:

$$PX = X(X^T X)^{-1} X^T X = X$$

Interpretation: Any vector in $\text{col}(X)$ is unchanged by the projection.

13.3.4 4. Residual Matrix

$$I - P$$

This matrix projects onto the orthogonal complement of $\text{col}(X)$, giving the residuals:

$$r = (I - P)y$$

13.4 Aerospace Interpretation

In the context of drag coefficient prediction:

- Py : The component of observed drag that **can be explained** by our aerodynamic features (angle of attack, Mach number, etc.)
- $(I-P)y$: The component that **cannot be explained** (residual variance due to unmodeled effects, measurement noise, etc.)

The projection matrix decomposes the total variance into:

- **Explained variance** (signal captured by the model)
- **Unexplained variance** (noise or missing features)

14 When Direct Solution Fails

While the closed-form solution $\beta^* = (X^T X)^{-1} X^T y$ is elegant, it faces several practical challenges.

14.1 Problem 1: Singular Matrix (Non-Invertibility)

When it occurs:

1. More features than samples: $n < d + 1$
 - Not enough data to uniquely determine all parameters
 - Infinitely many solutions exist
 - Example: 50 wind tunnel measurements but 100 polynomial features
2. Multicollinearity: Highly correlated features
 - Example: Including both velocity and Mach number when temperature is constant (they're perfectly correlated)
 - $X^T X$ becomes nearly singular (very small eigenvalues)
 - Small numerical errors can cause huge changes in solution

Consequences:

- Cannot compute $(X^T X)^{-1}$
- Solution is non-unique or unstable

14.2 Problem 2: Computational Cost

Matrix inversion complexity: $O(d^3)$ operations

For high-dimensional problems:

- $d = 1000$ features: ~1 billion operations
- $d = 10000$ features: ~1 trillion operations

Memory requirements: Storing $X^T X$ requires $O(d^2)$ memory

When it matters:

- Large-scale machine learning (millions of features)
- Online learning (real-time updates)
- Embedded systems (limited computational resources)

14.3 Problem 3: Numerical Stability

Condition number: $\kappa(X^T X) = \frac{\sigma_{\max}}{\sigma_{\min}}$

- Large condition number \rightarrow small perturbations in data cause large changes in solution
- Floating-point arithmetic errors accumulate
- Results become unreliable

Example: If $\kappa = 10^{10}$ and we use 64-bit floating point (16 digits precision), we may lose all significant digits in the solution.

14.4 Solutions

14.4.1 1. Regularization

Add a penalty term to prevent overfitting and improve conditioning:

Ridge Regression (L2):

$$\beta^* = (X^T X + \lambda I)^{-1} X^T y$$

Lasso (L1):

$$\min_{\beta} \|y - X\beta\|^2 + \lambda \|\beta\|_1$$

14.4.2 2. Gradient Descent

Iteratively update parameters without matrix inversion:

- Complexity per iteration: $O(nd)$ instead of $O(d^3)$
- Can handle very large d
- Works for any differentiable loss function

14.4.3 3. QR Decomposition

Numerically stable direct method:

$$X = QR$$

- More stable than normal equations
- Still $O(nd^2)$ complexity

14.4.4 4. Singular Value Decomposition (SVD)

Most numerically stable method:

$$X = U\Sigma V^T$$

- Handles rank deficiency gracefully
- Reveals multicollinearity
- Gold standard for numerical stability

15 Gradient Descent: Iterative Approach

When direct solution is impractical, we can use iterative optimization methods. Gradient descent is the foundation of modern machine learning.

15.1 The Algorithm

Goal: Minimize $\text{RSS}(\beta) = \|y - X\beta\|^2$

Strategy: Start with an initial guess and repeatedly move in the direction of steepest descent.

15.1.1 Steps

1. **Initialize:** Choose starting parameters $\beta^{(0)}$ (typically random or zeros)
2. **Iterate:** For $t = 0, 1, 2, \dots$ until convergence:

$$\beta^{(t+1)} = \beta^{(t)} - \eta \nabla_{\beta} \text{RSS}(\beta^{(t)})$$

where $\eta > 0$ is the **learning rate** (step size)

3. **Stop:** When change is small enough or maximum iterations reached

15.2 Computing the Gradient

The gradient of RSS is:

$$\nabla_{\beta} \text{RSS} = \frac{\partial}{\partial \beta} \|y - X\beta\|^2 = -2X^T(y - X\beta)$$

Update Rule:

$$\beta^{(t+1)} = \beta^{(t)} + 2\eta X^T(y - X\beta^{(t)})$$

Intuition:

- The gradient ∇RSS points in the direction of steepest **increase**
- We move in the **opposite direction** (negative gradient) to decrease the error
- The learning rate η controls how big a step we take

15.3 Convergence

Under appropriate conditions (convexity, suitable learning rate):

$$\beta^{(t)} \rightarrow \beta^* \quad \text{as} \quad t \rightarrow \infty$$

The method converges to the same solution as the direct method, but arrives there iteratively.

16 Gradient Descent Variants

Different variants of gradient descent trade off between convergence speed and computational cost per iteration.

16.1 Batch Gradient Descent

Use all n samples in each iteration:

$$\beta^{(t+1)} = \beta^{(t)} - \eta \nabla_{\beta} \text{RSS}(\beta^{(t)})$$

where the gradient uses all data:

$$\nabla_{\beta} \text{RSS} = -2X^T(y - X\beta^{(t)}) = -2 \sum_{i=1}^n x_i(y_i - x_i^T \beta^{(t)})$$

Advantages:

- Stable, smooth convergence
- Gradient is exact (no sampling noise)
- Can use optimized linear algebra libraries

Disadvantages:

- Each iteration costs $O(nd)$
- Slow for very large datasets (n in millions)
- Entire dataset must fit in memory

When to use: Small to medium datasets, when you need stable convergence

16.2 Stochastic Gradient Descent (SGD)

Use one random sample i per iteration:

$$\beta^{(t+1)} = \beta^{(t)} + 2\eta x_i(y_i - x_i^T \beta^{(t)})$$

Advantages:

- Very fast updates: $O(d)$ per iteration
- Can process data online (streaming)
- May escape shallow local minima (for non-convex problems)
- Naturally handles huge datasets

Disadvantages:

- Noisy updates (high variance)
- Oscillates around minimum (never truly converges)
- Requires careful learning rate tuning
- May need learning rate decay

When to use: Very large datasets, online learning, when fast iteration is more important than stable convergence

16.3 Mini-Batch Gradient Descent

Use a random subset of b samples per iteration (typical: $b = 32, 64, 128, 256$):

$$\beta^{(t+1)} = \beta^{(t)} + \frac{2\eta}{b} \sum_{i \in \mathcal{B}_t} x_i(y_i - x_i^T \beta^{(t)})$$

where \mathcal{B}_t is a random mini-batch at iteration t .

Advantages:

- **Best balance** between speed and stability
- Vectorized operations (GPU-friendly)
- Reduced variance compared to SGD
- Faster than batch for large n

Disadvantages:

- One more hyperparameter (batch size)
- Still some oscillation

When to use: Default choice for modern machine learning—combines benefits of both batch and stochastic methods

16.4 Comparison Summary

Method	Samples per Iteration	Cost per Iteration	Convergence	Best For
Batch	All (n)	$O(nd)$	Smooth, stable	Small datasets
Stochastic	1	$O(d)$	Noisy, fast	Huge datasets, online
Mini-batch	b (32-256)	$O(bd)$	Balanced	Most applications

17 Learning Rate Selection

The learning rate η is one of the most critical hyperparameters in gradient descent.

17.1 The Role of Learning Rate

In the update rule:

$$\beta^{(t+1)} = \beta^{(t)} - \eta \nabla_{\beta} \text{RSS}(\beta^{(t)})$$

η controls **how far** we move in the direction of the negative gradient.

17.2 Too Small Learning Rate ($\eta \ll 1$)

Symptoms:

- Very slow progress toward minimum
- Requires many iterations to converge
- May time out before reaching optimum

Cost: Wasted computation time

Example: If optimal $\eta = 0.01$ but we use $\eta = 0.0001$, we need $100\times$ more iterations.

17.3 Too Large Learning Rate ($\eta \gg 1$)

Symptoms:

- Overshooting the minimum
- Oscillation around optimum
- Divergence (error increases instead of decreasing)
- Instability

Cost: Never converges, wasted effort

Example: Parameters “bounce” past the minimum on each side, never settling down.

17.4 Finding the Right Learning Rate

Rule of thumb: Start with $\eta \in \{0.001, 0.01, 0.1, 1.0\}$ and tune based on training curves.

Grid search: Try multiple values, plot RSS vs. iteration, choose the largest η that converges smoothly.

Diagnostic: Plot $\text{RSS}^{(t)}$ vs. t : - Decreasing smoothly \rightarrow good - Decreasing then oscillating \rightarrow too large - Flat or barely decreasing \rightarrow too small - Increasing \rightarrow way too large

17.5 Adaptive Learning Rates

Modern optimization uses sophisticated learning rate strategies:

17.5.1 1. Learning Rate Decay

Decrease η over time:

$$\eta_t = \frac{\eta_0}{1 + kt}$$

or exponential decay:

$$\eta_t = \eta_0 e^{-kt}$$

Rationale: Start with large steps (fast progress), then small steps (fine-tuning near minimum)

17.5.2 2. Momentum

Use exponentially weighted moving average of gradients:

$$v^{(t)} = \gamma v^{(t-1)} + \eta \nabla_{\beta} \text{RSS}(\beta^{(t)})$$

$$\beta^{(t+1)} = \beta^{(t)} - v^{(t)}$$

Benefits:

- Accelerates convergence in relevant directions
- Dampens oscillations
- Helps escape plateaus

17.5.3 3. Adam (Adaptive Moment Estimation)

Most popular modern optimizer. Combines: - Momentum (first moment) - Adaptive learning rates per parameter (second moment)

Update rule (simplified):

$$\begin{aligned} m^{(t)} &= \beta_1 m^{(t-1)} + (1 - \beta_1) g^{(t)} \\ v^{(t)} &= \beta_2 v^{(t-1)} + (1 - \beta_2) (g^{(t)})^2 \\ \beta^{(t+1)} &= \beta^{(t)} - \eta \frac{m^{(t)}}{\sqrt{v^{(t)}} + \epsilon} \end{aligned}$$

where $g^{(t)} = \nabla_{\beta} \text{RSS}(\beta^{(t)})$

Advantages:

- Works well with default parameters ($\beta_1 = 0.9$, $\beta_2 = 0.999$, $\eta = 0.001$)
- Adapts learning rate for each parameter
- Widely used in deep learning

17.5.4 4. Line Search

At each iteration, optimize the learning rate:

$$\eta_t = \arg \min_{\eta} \text{RSS}(\beta^{(t)} - \eta \nabla_{\beta} \text{RSS}(\beta^{(t)}))$$

Benefits: Guaranteed improvement at each step

Cost: Additional computation per iteration

18 Statistical Properties: Assumptions

The theoretical properties of linear regression rely on several key assumptions. Understanding these helps us know when the method will work well and when it might fail.

18.1 Classical Linear Regression Assumptions

18.1.1 1. Linearity

Assumption: The true relationship is linear in parameters:

$$y = X\beta_{\text{true}} + \epsilon$$

What it means: There exist true parameters β_{true} such that the data-generating process follows this form.

If violated:

- OLS estimates are biased
- Predictions may be systematically wrong
- **Solution:** Transform features, add polynomial or interaction terms, use nonlinear models

18.1.2 2. Independence

Assumption: Samples (x_i, y_i) are independent and identically distributed (i.i.d.).

What it means:

- Each observation is drawn independently from the same distribution
- Knowing one sample doesn't tell you about others

If violated:

- Time series data (autocorrelation)
- Spatial data (neighboring points are similar)
- Clustered data (measurements from same aircraft)

Consequences:

- Standard errors are wrong
- Confidence intervals unreliable
- **Solution:** Use time series models, spatial models, clustered standard errors

18.1.3 3. Homoscedasticity (Constant Variance)

Assumption: Error variance is constant across all observations:

$$\text{Var}(\epsilon_i) = \sigma^2 \quad \text{for all } i$$

What it means: The “noise level” is the same regardless of feature values.

Heteroscedasticity (violation): Variance depends on features

$$\text{Var}(\epsilon_i) = \sigma_i^2 \quad (\text{different for each } i)$$

Aerospace example:

- Sensor noise may increase with dynamic pressure
- Measurement error in drag coefficient may depend on angle of attack
- Flow unsteadiness increases near stall

If violated:

- OLS is still unbiased but not optimal (not minimum variance)
- Standard errors are wrong
- **Solution:** Weighted least squares, robust standard errors, transform the response

18.1.4 4. Normality

Assumption: Errors are normally distributed:

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \text{independently}$$

What it means: Random fluctuations follow a bell curve (Gaussian distribution).

Importance:

- Required for exact hypothesis tests and confidence intervals
- **Not required** for OLS estimation itself
- By Central Limit Theorem, normality becomes less critical for large n

If violated:

- OLS estimates still unbiased
- Hypothesis tests may be unreliable (especially for small n)
- **Solution:** Use robust inference, bootstrap, or transform data

18.1.5 5. No Perfect Multicollinearity

Assumption: Columns of X are linearly independent (full rank):

$$\text{rank}(X) = d + 1$$

What it means: No feature can be perfectly predicted by other features.

Perfect multicollinearity examples:

- Including both Celsius and Fahrenheit temperature
- Including velocity twice
- Sum of dummy variables equals 1 (dummy variable trap)

Near multicollinearity: Features are highly (but not perfectly) correlated

- Mach number and true airspeed (at constant altitude/temperature)
- Multiple polynomial terms of the same variable
- Multiple similar wind tunnel configurations

Consequences:

- $X^T X$ is singular or nearly singular
- Cannot compute $(X^T X)^{-1}$ or it's numerically unstable
- Coefficients have huge standard errors
- Small data changes cause large coefficient changes

Solutions:

- Remove redundant features
- Regularization (Ridge, Lasso)
- Principal Component Analysis (PCA)
- Domain knowledge to select meaningful features

19 Gauss-Markov Theorem: Implications for Practice

19.1 Statement of the Theorem

Gauss-Markov Theorem: Under assumptions 1-3 (linearity, independence, homoscedasticity), the OLS estimator $\beta^* = (X^T X)^{-1} X^T y$ is **BLUE** (Best Linear Unbiased Estimator).

19.2 What Does BLUE Mean?

19.2.1 Best

Minimum variance among all linear unbiased estimators.

What it means:

- OLS gives the most precise estimates possible (smallest uncertainty)
- No other linear unbiased method produces tighter confidence intervals
- In the space of linear unbiased estimators, OLS is optimal

Why it matters: For safety-critical aerospace applications, we want the most precise performance predictions possible.

19.2.2 Linear

Estimator is a linear function of the observations y :

$$\beta^* = Cy$$

for some matrix C (that may depend on X but not on y).

For OLS: $C = (X^T X)^{-1} X^T$

Why restrict to linear estimators?:

- Computationally tractable
- Well-understood statistical properties
- Note: Nonlinear estimators might have lower variance, but they're typically biased

19.2.3 Unbiased

Expected value equals true parameter:

$$E[\beta^*] = \beta_{\text{true}}$$

What it means:

- On average (over many datasets), estimates equal true values
- No systematic over- or under-estimation
- If we repeat the experiment many times, the average of our estimates converges to the truth

Why it matters: We want methods that are correct on average, not systematically biased.

19.3 Aerospace Context

Critical for certification:

Flight envelopes must be determined with:

- **Minimal uncertainty** (Best)
- **No systematic bias** (Unbiased)

The BLUE property ensures:

- Tightest possible bounds on performance predictions
- Maximum confidence in safety margins
- Statistically rigorous methods required for regulatory compliance (FAA and EASA require demonstrated performance with appropriate confidence levels and uncertainty quantification)

Example: When certifying maximum operating Mach number, we need the most precise drag predictions possible. Using OLS (under appropriate assumptions) guarantees we're using the optimal estimator.

19.4 When BLUE Doesn't Apply

BLUE only applies under assumptions 1-3. If:

- **Nonlinearity:** True relationship is nonlinear → OLS is biased
- **Heteroscedasticity:** Variance is not constant → OLS is unbiased but not minimum variance (use weighted least squares)
- **Dependence:** Samples are correlated → Standard errors are wrong (use robust methods)

Additionally, BLUE is restricted to **linear unbiased** estimators. In modern machine learning:

- We often accept **bias** for reduced **variance** (bias-variance tradeoff)
- Regularization (Ridge, Lasso) introduces bias but reduces variance
- For prediction, biased estimators can outperform OLS

20 Statistical Properties: Distribution

Understanding the distribution of our parameter estimates allows us to quantify uncertainty and perform statistical inference. This section provides a rigorous foundation for confidence intervals, hypothesis tests, and prediction intervals.

20.1 Distribution of the OLS Estimator

Under the **normality assumption** ($\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ independently), the OLS estimator has a multivariate normal distribution:

$$\beta^* \sim \mathcal{N}(\beta_{\text{true}}, \sigma^2(X^T X)^{-1})$$

20.2 Proving Unbiasedness

Let's rigorously derive that $E[\beta^*] = \beta_{\text{true}}$.

Starting from the OLS formula:

$$\beta^* = (X^T X)^{-1} X^T y$$

Substitute the true model $y = X\beta_{\text{true}} + \epsilon$:

$$\beta^* = (X^T X)^{-1} X^T (X\beta_{\text{true}} + \epsilon)$$

Distribute:

$$\beta^* = (X^T X)^{-1} X^T X\beta_{\text{true}} + (X^T X)^{-1} X^T \epsilon$$

Simplify (since $(X^T X)^{-1} X^T X = I$):

$$\beta^* = \beta_{\text{true}} + (X^T X)^{-1} X^T \epsilon$$

Take expectation (noting that X is fixed and $E[\epsilon] = 0$):

$$E[\beta^*] = \beta_{\text{true}} + (X^T X)^{-1} X^T E[\epsilon] = \beta_{\text{true}}$$

Conclusion: OLS is unbiased—on average, it recovers the true parameters.

Key requirement: This derivation only assumes $E[\epsilon] = 0$ and that X is fixed (or independent of ϵ). **Normality is not required for unbiasedness.**

20.3 Deriving the Covariance Matrix

Now let's derive the variance of the OLS estimator.

Starting from:

$$\beta^* = \beta_{\text{true}} + (X^T X)^{-1} X^T \epsilon$$

The estimation error is:

$$\beta^* - \beta_{\text{true}} = (X^T X)^{-1} X^T \epsilon$$

The covariance matrix is:

$$\text{Cov}(\beta^*) = E[(\beta^* - \beta_{\text{true}})(\beta^* - \beta_{\text{true}})^T]$$

Substitute the error expression:

$$\text{Cov}(\beta^*) = E[(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}]$$

Factor out constants (since X is non-random):

$$\text{Cov}(\beta^*) = (X^T X)^{-1} X^T E[\epsilon \epsilon^T] X (X^T X)^{-1}$$

Under the assumption $E[\epsilon \epsilon^T] = \sigma^2 I$ (independence and homoscedasticity):

$$\text{Cov}(\beta^*) = (X^T X)^{-1} X^T (\sigma^2 I) X (X^T X)^{-1}$$

Simplify:

$$\text{Cov}(\beta^*) = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

Key insight: The covariance depends on:

1. **Error variance** σ^2 : More noise \rightarrow higher variance in estimates
2. **Design matrix** X : The choice of features and sample points affects precision

20.4 Unpacking the Covariance Matrix

$$\text{Cov}(\beta^*) = \sigma^2 (X^T X)^{-1}$$

20.4.1 Diagonal Elements: Individual Coefficient Variance

$$\text{Var}(\beta_j^*) = \sigma^2[(X^T X)^{-1}]_{jj}$$

What affects this variance?

1. **Sample size:** As n increases, $X^T X$ grows, so $(X^T X)^{-1}$ shrinks \rightarrow variance decreases
2. **Feature spread:** If feature j has large variance across samples, the corresponding diagonal element of $(X^T X)^{-1}$ is smaller \rightarrow lower variance in β_j^*
3. **Multicollinearity:** If feature j is correlated with other features, $[(X^T X)^{-1}]_{jj}$ increases \rightarrow higher variance (harder to distinguish the individual effect of feature j)

20.4.2 Off-Diagonal Elements: Correlation Between Coefficients

$$\text{Cov}(\beta_j^*, \beta_k^*) = \sigma^2[(X^T X)^{-1}]_{jk}$$

Physical interpretation:

- If features j and k are correlated, their coefficient estimates are also correlated
- Positive covariance: Overestimating β_j tends to coincide with overestimating β_k
- Negative covariance: Overestimating β_j tends to coincide with underestimating β_k

Aerospace example: In a model with both Mach number M and dynamic pressure $q = \frac{1}{2}\rho V^2$:

- These features are highly correlated (both increase with velocity)
- Their coefficient estimates will be negatively correlated
- If we overestimate the Mach effect, we tend to underestimate the dynamic pressure effect (and vice versa)

20.4.3 Variance Inflation Factor (VIF)

The **Variance Inflation Factor** quantifies how much multicollinearity increases the variance of a coefficient estimate:

$$\text{VIF}_j = \frac{1}{1 - R_j^2}$$

where R_j^2 is the coefficient of determination from regressing feature j on all other features.

Interpretation:

- $VIF_j = 1$: No correlation with other features (ideal)
- $VIF_j = 5$: Variance is 5× larger than if features were uncorrelated
- $VIF_j > 10$: Severe multicollinearity (rule of thumb)

Relationship to covariance matrix:

$$\text{Var}(\beta_j^*) = \sigma^2 VIF_j \cdot \frac{1}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}$$

Engineering guidance: High VIF suggests:

- Remove redundant features
- Use regularization (Ridge regression)
- Collect more data with different feature combinations

20.5 Practical Implications

This distribution allows us to:

1. **Construct confidence intervals** for each coefficient
2. **Perform hypothesis tests** about parameters
3. **Quantify uncertainty** in predictions
4. **Design experiments** to minimize variance (optimal experimental design)

20.6 What If Normality Doesn't Hold?

Good news: Even without normality:

- OLS estimates remain **unbiased** (under linearity, independence)
- OLS is still **BLUE** (under homoscedasticity)
- **Asymptotically** (large n), the Central Limit Theorem implies approximate normality

When it matters:

- Small samples ($n < 30$): normality important for exact inference
- Large samples ($n > 100$): approximate normality holds regardless
- **Solution for small samples:** Use robust inference methods, bootstrap

20.7 The Role of the Central Limit Theorem

The **Central Limit Theorem (CLT)** provides powerful justification for using normal-theory inference even when errors aren't normally distributed.

CLT statement (simplified): For the OLS estimator,

$$\frac{\beta^* - \beta_{\text{true}}}{\text{SE}(\beta^*)} \xrightarrow{d} \mathcal{N}(0, I) \quad \text{as } n \rightarrow \infty$$

What this means:

1. Even if ϵ_i are not normal, the **distribution of β^* becomes approximately normal** for large n
2. The approximation improves with sample size
3. **Practical rule:** For $n > 30 - 50$, normal approximations are usually adequate

Why this happens:

β^* depends on the sum $X^T \epsilon = \sum_{i=1}^n x_i \epsilon_i$. The CLT says that sums of independent random variables (the ϵ_i) approach a normal distribution regardless of the individual distributions.

Engineering implications:

- With large wind tunnel datasets ($n > 100$), we can trust confidence intervals and hypothesis tests even if measurement errors aren't perfectly Gaussian
- For small samples (early flight tests with $n < 20$), normality of errors is more critical
- Bootstrap methods can verify the adequacy of normal approximations

21 Estimating the Noise Variance

In practice, we don't know the true noise variance σ^2 . We must estimate it from the data using the residuals.

21.1 Residual Variance Estimator

$$\hat{\sigma}^2 = \frac{1}{n-d-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{\text{RSS}}{n-d-1}$$

where:

- $\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ is the residual sum of squares
- n is the number of samples

- d is the number of features (not counting intercept)
- $d + 1$ is the total number of parameters (including intercept)

21.2 Why Divide by $n - d - 1$?

21.2.1 Degrees of Freedom

Definition: Degrees of freedom (df) = number of independent pieces of information available for estimating variance.

Calculation:

- Start with n observations
- Estimate $d + 1$ parameters ($\beta_0, \beta_1, \dots, \beta_d$)
- Each parameter estimate “uses up” one degree of freedom
- Remaining df: $n - (d + 1) = n - d - 1$

Physical interpretation:

The residuals $r_i = y_i - \hat{y}_i$ are not independent—they satisfy the normal equations, which impose $d + 1$ linear constraints. Therefore, only $n - d - 1$ of the residuals are truly “free” to vary.

21.2.2 Why Not Divide by n ?

If we divided by n , the estimator would be **biased**:

$$E \left[\frac{1}{n} \sum_{i=1}^n r_i^2 \right] < \sigma^2$$

This **underestimates** the true variance because the residuals are constrained to satisfy the normal equations (they’re not truly independent).

Mathematical proof of bias:

The expected value of RSS is:

$$E[\text{RSS}] = E[r^T r] = E[((I - P)y)^T ((I - P)y)]$$

where P is the projection matrix. After considerable algebra:

$$E[\text{RSS}] = (n - d - 1)\sigma^2$$

Therefore:

$$E\left[\frac{\text{RSS}}{n}\right] = \frac{n-d-1}{n}\sigma^2 < \sigma^2$$

But dividing by $n - d - 1$ gives an unbiased estimator:

$$E\left[\frac{\text{RSS}}{n-d-1}\right] = \sigma^2$$

21.2.3 Distribution of $\hat{\sigma}^2$

Under normality of errors, the scaled RSS has a chi-squared distribution:

$$\frac{\text{RSS}}{\sigma^2} \sim \chi_{n-d-1}^2$$

This implies:

$$\frac{(n-d-1)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-d-1}^2$$

Properties:

- **Mean:** $E[\hat{\sigma}^2] = \sigma^2$ (unbiased)
- **Variance:** $\text{Var}(\hat{\sigma}^2) = \frac{2\sigma^4}{n-d-1}$
- For large df, the distribution becomes approximately normal

Practical implication: The variance estimate itself has uncertainty. With small samples, $\hat{\sigma}^2$ can be quite variable.

21.3 Aerospace Example

Scenario: Fit drag coefficient model:

$$C_D = \beta_0 + \beta_1 \alpha + \beta_2 \alpha^2$$

using $n = 20$ wind tunnel data points.

Parameters:

- Number of features: $d = 2$ (we have α and α^2)
- Total parameters: $d + 1 = 3$ (including intercept β_0)
- Degrees of freedom: $n - d - 1 = 20 - 2 - 1 = 17$

Calculation:

Suppose after fitting, we find:

$$\text{RSS} = \sum_{i=1}^{20} (C_{D_i} - \hat{C}_{D_i})^2 = 0.0085$$

Variance estimate:

$$\hat{\sigma}^2 = \frac{0.0085}{17} = 0.0005$$

Standard deviation (root mean square error):

$$\hat{\sigma} = \sqrt{0.0005} = 0.0224$$

Interpretation: The typical prediction error for drag coefficient is about 0.022 (in C_D units, typically in counts where 1 count = 0.0001).

21.4 Relationship to Model Fit

Small $\hat{\sigma}^2$:

- Residuals are small
- Model fits data well
- Low prediction uncertainty
- May indicate good model or low measurement noise

Large $\hat{\sigma}^2$:

- Residuals are large
- Model doesn't fit well or data is very noisy
- High prediction uncertainty
- May indicate:
 - Missing important features
 - Wrong functional form
 - High measurement noise
 - Outliers

21.5 Residual Standard Error vs. R-squared

Two common measures of fit quality:

Residual Standard Error (RSE):

$$\text{RSE} = \sqrt{\hat{\sigma}^2} = \sqrt{\frac{\text{RSS}}{n - d - 1}}$$

- Has same units as response variable
- Absolute measure of fit quality
- Useful for prediction intervals

R-squared:

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} = 1 - \frac{\sum(y_i - \hat{y}_i)^2}{\sum(y_i - \bar{y})^2}$$

- Dimensionless (between 0 and 1)
- Relative measure (fraction of variance explained)
- Can be misleading with many features

Relationship: Both measure fit quality but RSE is often more interpretable for prediction purposes.

22 Confidence Intervals for Coefficients

Confidence intervals quantify the uncertainty in our parameter estimates, providing a range of plausible values for each coefficient.

22.1 Understanding Coefficient Uncertainty

Each estimated coefficient β_j^* is a **random variable**—it depends on the random sample we collected. If we repeated the experiment, we'd get different data and different estimates.

Key question: How uncertain are we about each coefficient?

22.2 Standard Error

The **standard error** of β_j^* measures its sampling variability:

$$\text{SE}(\beta_j^*) = \sqrt{\hat{\sigma}^2 [(X^T X)^{-1}]_{jj}}$$

Breaking down the formula:

1. $\hat{\sigma}^2$: Estimated error variance (data noise level)
2. $[(X^T X)^{-1}]_{jj}$: Depends on experimental design
3. **Square root**: Converts variance to standard deviation

Interpretation:

- Small SE \rightarrow precise estimate (high confidence)
- Large SE \rightarrow imprecise estimate (low confidence)

22.2.1 Factors Affecting Standard Error

1. Data noise ($\hat{\sigma}^2$):

- More noise \rightarrow larger SE
- Better measurement quality \rightarrow smaller SE
- Controllable through experimental design

2. Sample size (n):

More data \rightarrow smaller SE (typically $\text{SE} \propto 1/\sqrt{n}$)

Example: To halve the SE, need $4\times$ as much data

3. Feature spread:

Features with more spread \rightarrow smaller SE

Intuition: If angle of attack varies from 0° to 20° , we can estimate its effect more precisely than if it only varies from 0° to 2°

Mathematically: For simple linear regression,

$$\text{SE}(\beta_1) \propto \frac{1}{\sqrt{\sum(x_i - \bar{x})^2}}$$

4. Feature correlation (multicollinearity):

Correlated features \rightarrow larger SE

Why: Hard to distinguish individual effects when features move together

Quantified by VIF:

$$\text{SE}(\beta_j) \propto \sqrt{\text{VIF}_j}$$

22.2.2 Computing Standard Errors in Practice

Step 1: Compute $X^T X$ and invert it

Step 2: Extract diagonal elements of $(X^T X)^{-1}$

Step 3: Estimate error variance $\hat{\sigma}^2 = \text{RSS}/(n - d - 1)$

Step 4: Compute SE for each coefficient:

$$\text{SE}(\beta_j^*) = \sqrt{\hat{\sigma}^2 \cdot [(X^T X)^{-1}]_{jj}}$$

Most statistical software (R, Python statsmodels, MATLAB) reports standard errors automatically.

22.3 Confidence Interval Formula

A $100(1 - \gamma)\%$ **confidence interval** for β_j is:

$$\beta_j^* \pm t_{\gamma/2, n-d-1} \cdot \text{SE}(\beta_j^*)$$

where $t_{\gamma/2, n-d-1}$ is the critical value from the **t-distribution** with $n - d - 1$ degrees of freedom.

22.3.1 Why Use the t-Distribution?

The issue: We don't know σ^2 , only $\hat{\sigma}^2$

The solution: Account for uncertainty in $\hat{\sigma}^2$ by using the t-distribution instead of the normal distribution

The standardized coefficient:

$$\frac{\beta_j^* - \beta_j}{\text{SE}(\beta_j^*)} = \frac{\beta_j^* - \beta_j}{\sqrt{\hat{\sigma}^2 [(X^T X)^{-1}]_{jj}}}$$

follows a **t-distribution with $n - d - 1$ degrees of freedom**, not a standard normal.

Properties of t-distribution:

- Symmetric and bell-shaped like normal
- Heavier tails than normal (more probability in extremes)
- Approaches normal as df increases
- Accounts for additional uncertainty from estimating σ^2

22.3.2 Common Confidence Levels

95% confidence: Most common; $\gamma = 0.05$

- Critical value: $t_{0.025,n-d-1}$
- For large n : $t_{0.025,\infty} \approx 1.96 \approx 2$
- For $n = 20$, $d = 2$: $t_{0.025,17} \approx 2.11$
- For $n = 10$, $d = 2$: $t_{0.025,7} \approx 2.36$

90% confidence: Narrower interval; $\gamma = 0.10$

- Critical value: $t_{0.05,n-d-1}$
- For large n : $t_{0.05,\infty} \approx 1.645$

99% confidence: Wider interval; $\gamma = 0.01$

- Critical value: $t_{0.005,n-d-1}$
- For large n : $t_{0.005,\infty} \approx 2.576$

Tradeoff: Higher confidence \rightarrow wider interval \rightarrow less precision

22.4 Interpretation

What a 95% CI means:

If we repeated the experiment many times and computed a 95% CI each time, approximately 95% of these intervals would contain the true parameter value β_j .

Frequentist interpretation:

- The interval is random (varies with data)
- The true parameter is fixed (not random)
- In the long run, 95% of such intervals capture the truth

What it does NOT mean:

- “There’s a 95% probability that β_j is in this interval” (Bayesian interpretation)
- “95% of the data falls in this interval” (confusing with prediction interval)
- “We are 95% sure our estimate is correct” (certainty about a specific estimate)

22.4.1 Practical Interpretation for Engineers

What you CAN say:

“We are 95% confident the lift slope is between 0.097 and 0.113 per degree”

“The data is consistent with lift slopes in the range [0.097, 0.113]”

“If the true lift slope were outside this range, we would be surprised (only 5% chance of observing our data)”

22.5 Aerospace Example

Scenario: Fitting lift coefficient model:

$$C_L = \beta_0 + \beta_1 \alpha$$

Data: $n = 30$ wind tunnel measurements

Results:

- $\beta_1^* = 0.105$ (per degree)
- $SE(\beta_1^*) = 0.004$
- $n = 30, d = 1$, so $df = 28$
- $t_{0.025, 28} \approx 2.048$

95% Confidence Interval:

$$0.105 \pm 2.048 \times 0.004 = 0.105 \pm 0.008 = [0.097, 0.113]$$

Interpretation: We are 95% confident the true lift slope is between 0.097 and 0.113 per degree.

Engineering decision:

- For performance calculations, use $\beta_1^* = 0.105$ (best estimate)
- For conservative safety margins, use lower bound $\beta_1 = 0.097$
- The 95% CI ensures that performance predictions account for estimation uncertainty
- Width of interval (0.016) indicates good precision

22.5.1 Impact of Sample Size

What if we had only $n = 10$ measurements?

- Same estimate: $\beta_1^* = 0.105$
- Larger SE (less data): $\text{SE}(\beta_1^*) = 0.007$ (hypothetically)
- Smaller df: $n - d - 1 = 8$
- Larger critical value: $t_{0.025,8} = 2.306$

New 95% CI:

$$0.105 \pm 2.306 \times 0.007 = 0.105 \pm 0.016 = [0.089, 0.121]$$

Comparison:

- Interval is twice as wide (0.032 vs. 0.016)
- Less precision with smaller sample
- Need more data for tighter bounds

23 Prediction Intervals vs. Confidence Intervals

While confidence intervals quantify uncertainty in **parameter estimates**, prediction intervals quantify uncertainty in **future observations**. This distinction is critical for aerospace applications.

23.1 The Two Types of Uncertainty

When making predictions, we face two sources of uncertainty:

23.1.1 1. Parameter Uncertainty (Epistemic)

Source: We estimated β from finite data, so $\beta^* \neq \beta_{\text{true}}$

Captured by: Confidence intervals for $E[y | x]$

Behavior: Decreases as n increases (we learn the true parameters)

For the mean response at a point x_0 :

$$\hat{y}_0 = x_0^T \beta^*$$

The variance of this prediction is:

$$\text{Var}(\hat{y}_0) = \sigma^2 x_0^T (X^T X)^{-1} x_0$$

95% Confidence Interval for the mean:

$$x_0^T \beta^* \pm t_{0.025, n-d-1} \cdot \sqrt{\hat{\sigma}^2 x_0^T (X^T X)^{-1} x_0}$$

23.1.2 2. Irreducible Noise (Aleatoric)

Source: Even with perfect knowledge of β_{true} , individual observations have noise ϵ

Captured by: Prediction intervals for $y | x$

Behavior: Does NOT decrease with n (inherent randomness)

For a new observation at x_0 :

$$y_{\text{new}} = x_0^T \beta_{\text{true}} + \epsilon_{\text{new}}$$

The prediction error is:

$$y_{\text{new}} - \hat{y}_0 = (x_0^T \beta_{\text{true}} - x_0^T \beta^*) + \epsilon_{\text{new}}$$

The variance combines both sources:

$$\text{Var}(y_{\text{new}} - \hat{y}_0) = \sigma^2 (1 + x_0^T (X^T X)^{-1} x_0)$$

95% Prediction Interval:

$$x_0^T \beta^* \pm t_{0.025, n-d-1} \cdot \sqrt{\hat{\sigma}^2 (1 + x_0^T (X^T X)^{-1} x_0)}$$

23.2 Key Differences

Confidence Interval (for mean):

- **Question:** Where is the **true regression line**?
- **Uncertainty:** Only parameter estimation error
- **Width:** $\propto 1/\sqrt{n}$ (shrinks to zero as $n \rightarrow \infty$)
- **Formula:** $\sqrt{\hat{\sigma}^2 x_0^T (X^T X)^{-1} x_0}$

Prediction Interval (for individual observation):

- **Question:** Where will a **new observation** fall?
- **Uncertainty:** Parameter error PLUS inherent noise
- **Width:** Has a minimum (set by σ) even as $n \rightarrow \infty$
- **Formula:** $\sqrt{\hat{\sigma}^2 (1 + x_0^T (X^T X)^{-1} x_0)}$

Mathematical relationship:

$$\text{Prediction interval} > \text{Confidence interval}$$

The extra “1” inside the square root accounts for irreducible noise.

23.3 Aerospace Example: Drag Prediction

Scenario: Predict drag coefficient at $\alpha = 5^\circ$ using our fitted model

Fitted model: $C_D = 0.012 + 0.0003\alpha + 0.00002\alpha^2$

Given:

- $\hat{\sigma}^2 = 0.0001$ (from $n = 50$ points)
- $x_0 = [1, 5, 25]^T$ (for $\alpha = 5^\circ$)
- $x_0^T(X^T X)^{-1} x_0 = 0.05$ (computed from data)
- $t_{0.025, 47} \approx 2.01$

Point prediction:

$$\hat{C}_D = 0.012 + 0.0003(5) + 0.00002(25) = 0.0140$$

95% Confidence Interval (for mean drag at $\alpha = 5^\circ$):

$$0.0140 \pm 2.01\sqrt{0.0001 \times 0.05} = 0.0140 \pm 0.0014 = [0.0126, 0.0154]$$

95% Prediction Interval (for a new measurement at $\alpha = 5^\circ$):

$$0.0140 \pm 2.01\sqrt{0.0001 \times (1 + 0.05)} = 0.0140 \pm 0.0206 = [0.0120, 0.0346]$$

Interpretation:

- **CI:** We’re 95% confident the **true mean** drag at $\alpha = 5^\circ$ is between 0.0126 and 0.0154
- **PI:** A **new measurement** at $\alpha = 5^\circ$ will fall between 0.0120 and 0.0346 with 95% probability
- **Width ratio:** PI is $\sqrt{1.05/0.05} \approx 4.6$ times wider than CI

Engineering use cases:

- **CI:** For performance calculations requiring mean behavior (cruise drag, range estimates)
- **PI:** For safety margins accounting for measurement variability (test point acceptance criteria)

23.4 Extrapolation Warning

Both intervals **widen** as x_0 moves away from the center of the training data.

Why: The term $x_0^T(X^T X)^{-1}x_0$ increases with distance from data centroid

Implications:

- Predictions are most reliable **within the range** of training data
- **Extrapolation** (predicting outside data range) has much larger uncertainty
- For aerospace: Don't predict drag at $\alpha = 30^\circ$ if you only have data up to $\alpha = 15^\circ$

Practical guideline:

- **Interpolation:** Predicting within data range → reliable
- **Modest extrapolation:** Slightly beyond data → use with caution
- **Extreme extrapolation:** Far beyond data → unreliable and dangerous

24 Hypothesis Testing for Individual Coefficients

24.1 The Central Question

After estimating coefficients from data, we ask: **Is this coefficient significantly different from zero, or could it just be noise?**

This is crucial for:

- **Feature selection:** Which features actually matter?
- **Model simplification:** Can we remove unimportant features?
- **Physical interpretation:** Is this effect real or spurious?

24.2 Setting Up the Test

24.2.1 Hypotheses

Null hypothesis H_0 :

$$\beta_j = 0$$

Interpretation: Feature j has **no true effect** on the response. Any non-zero estimate we observed is just random sampling variation.

Alternative hypothesis H_1 :

$$\beta_j \neq 0$$

Interpretation: Feature j **does affect** the response. The estimated effect is real, not just noise.

24.3 The Test Statistic

We construct a **t-statistic**:

$$t_j = \frac{\beta_j^*}{\text{SE}(\beta_j^*)} = \frac{\text{Estimated coefficient}}{\text{Standard error of estimate}}$$

Interpretation: - Measures **how many standard errors** the estimate is away from zero - If $\beta_j = 0$ truly, we expect $\beta_j^* \approx 0$ (within sampling error) - Large $|t_j|$ means the estimate is far from zero \rightarrow unlikely if H_0 is true

24.4 Distribution Under the Null

If H_0 is true ($\beta_j = 0$), then:

$$t_j \sim t_{n-d-1}$$

The test statistic follows a **t-distribution** with $n - d - 1$ degrees of freedom.

This allows us to compute **p-values** and make decisions.

24.5 Making the Decision

24.5.1 Approach 1: P-value

P-value: The probability of observing a test statistic as extreme as (or more extreme than) what we got, if the null hypothesis were true.

$$\text{p-value} = P(|t| \geq |t_j| \mid H_0 \text{ is true})$$

Decision rule: - If $p < \alpha$ (significance level, typically 0.05): **Reject** H_0 (coefficient is significant) - If $p \geq \alpha$: **Fail to reject** H_0 (insufficient evidence that coefficient differs from zero)

Interpretation of p-values:

- $p < 0.001$: Very strong evidence against H_0 (highly significant)
- $p < 0.01$: Strong evidence

- $p < 0.05$: Moderate evidence (conventional cutoff)
- $p > 0.05$: Weak or no evidence

24.5.2 Approach 2: Critical Value

Critical value: The threshold $t_{\text{crit}} = t_{\alpha/2, n-d-1}$ such that:

$$P(|t| > t_{\text{crit}} \mid H_0) = \alpha$$

Decision rule: - If $|t_j| > t_{\text{crit}}$: **Reject H_0** - If $|t_j| \leq t_{\text{crit}}$: **Fail to reject H_0**

Equivalence: Both approaches give the same conclusion. The p-value approach is more informative because it quantifies the strength of evidence.

24.6 Aerospace Example

Scenario: Testing whether Reynolds number affects drag coefficient in our model:

$$C_D = \beta_0 + \beta_1 \alpha + \beta_2 \alpha^2 + \beta_3 \log(\text{Re})$$

Results for β_3 :

- $\beta_3^* = -0.0015$
- $\text{SE}(\beta_3^*) = 0.0008$
- $n = 50, d = 3, \text{df} = 46$

Test statistic:

$$t_3 = \frac{-0.0015}{0.0008} = -1.875$$

P-value (two-tailed, df = 46):

$$p \approx 0.067$$

Decision at $\alpha = 0.05$:

- $p = 0.067 > 0.05$: Fail to reject H_0
- **Conclusion:** Insufficient evidence that Reynolds number significantly affects drag (at 5% significance level)

Engineering interpretation:

- We might consider **removing** $\log(\text{Re})$ from the model to simplify it
- Or acknowledge that Reynolds number effects are weak in our data range
- Or collect more data to increase statistical power

24.7 Important Caveats

1. **Statistical significance vs practical significance:** A tiny effect can be statistically significant with enough data
2. **Multiple testing:** Testing many coefficients increases false positive rate (use Bonferroni correction or similar)
3. **Correlation:** In models with many correlated features, individual t-tests can be misleading
4. **Model assumptions:** These tests assume normality, homoscedasticity, etc.

25 F-Test for Overall Model Significance

While t-tests examine individual coefficients, the **F-test** assesses whether **any** of the features are useful for prediction. This is a test of the **overall model significance**.

25.1 The Hypotheses

Null hypothesis H_0 :

$$\beta_1 = \beta_2 = \dots = \beta_d = 0$$

Interpretation: **None** of the features have predictive power. The best model is just the intercept $y_i = \beta_0 + \epsilon_i$.

Alternative hypothesis H_1 :

At least one $\beta_j \neq 0$ for $j = 1, \dots, d$

Interpretation: At least one feature provides useful information.

25.2 Decomposing Variance

The F-test is based on decomposing the total variance in the response:

Total Sum of Squares (TSS):

$$\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2$$

Measures total variability in y around its mean.

Residual Sum of Squares (RSS):

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Measures unexplained variability (what the model missed).

Model Sum of Squares (MSS):

$$\text{MSS} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \text{TSS} - \text{RSS}$$

Measures explained variability (what the model captured).

Key identity:

$$\text{TSS} = \text{MSS} + \text{RSS}$$

Total variance = Explained + Unexplained

25.3 The F-Statistic

The test statistic compares explained variance per parameter to unexplained variance:

$$F = \frac{\text{MSS}/d}{\text{RSS}/(n - d - 1)} = \frac{\text{Mean Model Sum of Squares}}{\text{Mean Residual Sum of Squares}}$$

Intuition:

- **Numerator:** Average amount of variance explained per feature
- **Denominator:** Average unexplained variance (noise estimate)
- **Large F:** Features explain much more than noise → reject H_0
- **Small F:** Features explain similar amount to noise → fail to reject H_0

25.4 Distribution Under the Null

If H_0 is true (all $\beta_j = 0$ except intercept), then:

$$F \sim F_{d, n-d-1}$$

The test statistic follows an **F-distribution** with d and $n - d - 1$ degrees of freedom.

Properties of F-distribution:

- Always positive (ratio of squared quantities)
- Right-skewed
- Parameterized by two df: numerator df = d , denominator df = $n - d - 1$
- Mean ≈ 1 under H_0 (for large df)

25.5 Making the Decision

P-value:

$$p\text{-value} = P(F_{d,n-d-1} \geq F_{\text{observed}} \mid H_0)$$

Decision rule (at significance level α , typically 0.05):

- If $p < \alpha$: **Reject** $H_0 \rightarrow$ model is useful
- If $p \geq \alpha$: **Fail to reject** $H_0 \rightarrow$ model has no explanatory power

25.6 Relationship to R-squared

The F-statistic is closely related to R^2 :

$$R^2 = \frac{\text{MSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

We can rewrite the F-statistic as:

$$F = \frac{R^2/d}{(1 - R^2)/(n - d - 1)}$$

Interpretation: The F-test asks whether R^2 is significantly different from zero.

25.7 Aerospace Example

Scenario: Testing overall significance of drag model

$$C_D = \beta_0 + \beta_1 \alpha + \beta_2 \alpha^2 + \beta_3 M^2$$

Data: $n = 100$ wind tunnel measurements, $d = 3$ features

Results:

- TSS = 0.050 (total drag variability)
- RSS = 0.005 (residual variability)
- MSS = 0.045 (explained by model)

- $R^2 = 0.045/0.050 = 0.90$ (90% variance explained)

F-statistic:

$$F = \frac{0.045/3}{0.005/96} = \frac{0.015}{0.0000521} \approx 288$$

Critical value: $F_{0.05,3,96} \approx 2.70$

P-value: $p < 0.001$ (highly significant)

Conclusion:

- **Reject H_0 :** The model is highly significant
- At least one of $\{\alpha, \alpha^2, M^2\}$ provides predictive power
- The model explains 90% of drag variability, far more than noise

Engineering interpretation: The aerodynamic features (angle of attack and Mach number) are strong predictors of drag coefficient.

25.8 F-Test vs. Multiple t-Tests

Why not just look at individual t-tests?

1. **Joint significance:** Features may be individually weak but jointly strong
2. **Multiple testing problem:** With many features, some t-tests will be significant by chance
3. **Correlated features:** Individual t-tests can miss patterns when features are correlated

Complementary use:

- **F-test:** “Is the overall model useful?”
- **t-tests:** “Which specific features matter?”

Typical workflow:

1. Check F-test: If not significant, stop (model is useless)
2. If F-test is significant, examine individual t-tests
3. Remove non-significant features and refit

25.9 F-Test for Model Comparison

The F-test can also compare **nested models** (one model is a special case of another).

Example: Compare simple vs. complex drag models

Model 1 (simple): $C_D = \beta_0 + \beta_1\alpha + \beta_2\alpha^2$

Model 2 (complex): $C_D = \beta_0 + \beta_1\alpha + \beta_2\alpha^2 + \beta_3M + \beta_4M^2 + \beta_5\alpha M$

Question: Do the additional terms ($M, M^2, \alpha M$) significantly improve fit?

F-statistic for model comparison:

$$F = \frac{(\text{RSS}_1 - \text{RSS}_2)/q}{\text{RSS}_2/(n - d_2 - 1)}$$

where:

- RSS_1 = residual sum of squares for simpler model
- RSS_2 = residual sum of squares for complex model
- q = number of additional parameters in complex model
- d_2 = number of features in complex model

Distribution: $F \sim F_{q, n-d_2-1}$ under H_0 (simple model is adequate)

Decision: If $p < \alpha$, the complex model provides significant improvement.

26 Summary

Linear regression is a powerful and versatile tool for aerospace engineering applications. Key takeaways:

26.1 Mathematical Foundations

- Linear regression models relationships as $y = X\beta + \epsilon$
- “Linear” refers to parameters, not features—we can model nonlinear relationships
- The optimal solution minimizes squared error: $\min \|y - X\beta\|^2$

26.2 Two Solution Approaches

1. **Closed-form** (Normal Equations): $\beta^* = (X^T X)^{-1} X^T y$

- Exact, one-step solution
- Requires matrix inversion
- Can fail for large d or singular $X^T X$

2. **Gradient Descent**: Iterative optimization

- Works for any problem size
- Requires tuning (learning rate, iterations)
- Foundation of modern machine learning

26.3 Geometric Insight

- OLS finds the **orthogonal projection** of y onto $\text{col}(X)$
- Residuals are **perpendicular** to feature space: $X^T r = 0$
- Projection matrix $P = X(X^T X)^{-1} X^T$ separates signal from noise

26.4 Statistical Properties

- Under standard assumptions, OLS is **BLUE** (Best Linear Unbiased Estimator)
- Coefficient distribution: $\beta^* \sim \mathcal{N}(\beta_{\text{true}}, \sigma^2(X^T X)^{-1})$
- Enables **confidence intervals** and **hypothesis tests**

26.5 Practical Considerations

- Check assumptions (linearity, independence, homoscedasticity, normality)
- Watch for multicollinearity
- Use appropriate validation metrics
- Consider regularization for high-dimensional problems

26.6 Aerospace Applications

Linear regression is essential for:

- Performance prediction (drag, lift, fuel consumption)
- Wind tunnel data analysis
- Flight test analysis
- System identification
- Model validation and uncertainty quantification

The techniques you've learned form the foundation for more advanced machine learning methods used throughout aerospace engineering.