

1 Prove OLS Estimator is Unbiased

Substitute model into estimator:

$$\begin{aligned}\hat{\beta} &= (X^T X)^{-1} X^T (X\beta + \varepsilon) \\ &= (X^T X)^{-1} X^T X\beta + (X^T X)^{-1} X^T \varepsilon \\ &= \beta + (X^T X)^{-1} X^T \varepsilon\end{aligned}$$

Take expectation:

$$\mathbb{E}[\hat{\beta}] = \beta + (X^T X)^{-1} X^T \mathbb{E}[\varepsilon]$$

Since $\mathbb{E}[\varepsilon] = 0$,

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

✓ OLS is unbiased

2 Derive Covariance of OLS Estimator

$$\hat{\beta} - \beta = (X^T X)^{-1} X^T \varepsilon$$

Covariance:

$$\text{Var}(\hat{\beta}) = (X^T X)^{-1} X^T \text{Var}(\varepsilon) X (X^T X)^{-1}$$

Since $\text{Var}(\varepsilon) = \sigma^2 I$,

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

3 When is OLS BLUE? (Gauss–Markov Theorem)

OLS is **BLUE (Best Linear Unbiased Estimator)** if:

1. Linear model is correct
2. $E[\varepsilon] = 0$
3. $\text{Var}(\varepsilon) = \sigma^2 I$
4. No perfect multicollinearity
5. Errors uncorrelated

Then OLS has **minimum variance among all linear unbiased estimators**.

If heteroscedasticity exists \rightarrow OLS not efficient.

4 Why Multicollinearity Increases Variance

From covariance:

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

If predictors are highly correlated:

- $X^T X$ becomes nearly singular
- Small eigenvalues appear
- Inverse blows up

Since:

$$(X^T X)^{-1} = Q \Lambda^{-1} Q^T$$

Small eigenvalue $\lambda_i \Rightarrow 1/\lambda_i$ very large

Thus variance explodes.

5 Show Ridge Regression Shrinks Eigenvalues

Ridge estimator:

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

Eigen-decompose:

$$X^T X = Q \Lambda Q^T$$

Then:

$$X^T X + \lambda I = Q(\Lambda + \lambda I)Q^T$$

Inverse:

$$(X^T X + \lambda I)^{-1} = Q(\Lambda + \lambda I)^{-1} Q^T$$

Each eigenvalue transforms:

$$\lambda_i \rightarrow \frac{1}{\lambda_i + \lambda}$$

Since $\lambda > 0$,

Denominator increases \Rightarrow coefficients shrink.



6 Ridge vs Lasso (Geometric View)

OLS:

- Elliptical loss contours
- No constraint

Ridge:

- Constraint: $\|\beta\|_2 \leq t \|\beta\|_2^2 \leq t$
- Circular constraint region
- Shrinks continuously
- Rarely exact zero

Lasso:

- Constraint: $\|\beta\|_1 \leq t \|\beta\|_1 \leq t$
- Diamond-shaped region
- Corners \rightarrow sparsity
- Produces exact zeros

Thus:

- Ridge = shrinkage
- Lasso = feature selection

7 Linear Regression as MAP Estimate

Assume:

Likelihood:

$$y|X, \beta \sim \mathcal{N}(X\beta, \sigma^2 I)$$

Prior:

$$\beta \sim \mathcal{N}(0, \tau^2 I)$$

Posterior \propto Likelihood \times Prior

Taking negative log:

$$\frac{1}{2\sigma^2} \|y - X\beta\|^2 + \frac{1}{2\tau^2} \|\beta\|^2$$

Multiply by constant:

$$\|y - X\beta\|^2 + \lambda \|\beta\|^2$$

Where:

$$\lambda = \sigma^2 / \tau^2$$

This is Ridge regression.

Thus:

- OLS = MLE
- Ridge = MAP with Gaussian prior

8 Condition Number of $X^T X$

Condition number:

$$\kappa(X^T X) = \frac{\lambda_{\max}}{\lambda_{\min}}$$

If $\lambda_{\min} \rightarrow 0$:

- Ill-conditioned
- Numerical instability
- High variance

Multicollinearity \Rightarrow small eigenvalues \Rightarrow large condition number.

Ridge improves condition number:

$$\kappa_{ridge} = \frac{\lambda_{\max} + \lambda}{\lambda_{\min} + \lambda}$$

9 Prove Convexity of Squared Loss

Loss:

$$J(\beta) = \|y - X\beta\|^2$$

Expand:

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

Hessian:

$$\nabla^2 J(\beta) = 2X^T X$$

Since $X^T X$ is positive semidefinite:

$$z^T X^T X z = \|Xz\|^2 \geq 0$$

Therefore Hessian PSD \Rightarrow convex.

If full rank \Rightarrow strictly convex.

10 What Happens When p>>np >> np>>n?

Then:

- $X^T X$ is singular
- Infinite OLS solutions
- Overfitting guaranteed
- Interpolates data perfectly

Geometrically:

- Many hyperplanes fit data

Variance:

- Extremely high

Solution:

- Ridge
- Lasso
- Dimensionality reduction

Modern ML (deep learning) often works in $p \gg n_p \gg n_p \gg n$, but classical OLS theory fails.