

# Linear Regression Q&A

## Advanced Theoretical Analysis of Linear Regression

We consider the classical linear model:  $\mathbf{y} = \mathbf{X}\beta + \varepsilon$

where:

- $\mathbf{X} \in \mathbb{R}^{n \times p}$
- $\beta \in \mathbb{R}^p$
- $\varepsilon \sim (\mathbf{0}, \sigma^2 \mathbf{I})$

Assume  $\mathbf{X}$  has full column rank unless stated otherwise

### 1. Prove that OLS estimator is unbiased.

OLS estimator:  $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

Substitute  $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ :

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

Take expectation:  $E[\hat{\beta}] = \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E[\varepsilon]$

Since  $E[\varepsilon] = \mathbf{0}$ :

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

✓ OLS is unbiased.

### 2. Derive covariance of OLS estimator.

From:  $\hat{\beta} = \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon$

$$\begin{aligned}\text{Variance: } \text{Var}(\hat{\beta}) &= \text{Var}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Var}(\varepsilon) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}\end{aligned}$$

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

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

- ✓ Covariance depends on  $(\mathbf{X}^T \mathbf{X})^{-1}$

### 3. When is OLS BLUE? (Gauss-Markov Theorem)

OLS is BLUE (Best Linear Unbiased Estimator) if:

- Model is linear in parameters
- $E[\varepsilon] = \mathbf{0}$
- $\text{Var}(\varepsilon) = \sigma^2 \mathbf{I}$
- No perfect multicollinearity

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

- ✓ No normality assumption required.

### 4. Why does multicollinearity increase variance?

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

If features are highly correlated:

- $\mathbf{X}^T \mathbf{X}$  becomes nearly singular
- Its inverse contains very large values
- Variance becomes very large

- ✓ Multicollinearity → unstable coefficient estimates

### 5. Show ridge regression shrinks eigenvalues.

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

Eigen-decomposition:  $\mathbf{X}^T \mathbf{X} = \mathbf{Q} \Lambda \mathbf{Q}^T$

Then:  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} = \mathbf{Q} (\Lambda + \lambda \mathbf{I}) \mathbf{Q}^T$

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

Each eigenvalue  $\lambda_i$  becomes:  $1 / (\lambda_i + \lambda)$

- ✓ Ridge increases eigenvalues
- ✓ Reduces variance
- ✓ Stabilizes matrix inversion

### 6. Compare ridge vs lasso geometrically.

**Ridge:**  $\min \|y - X\beta\|^2 + \lambda\|\beta\|^2$

**Lasso:**  $\min \|y - X\beta\|^2 + \lambda\|\beta\|_1$

Geometrically:

- Ridge constraint region = circle (L2 ball)
- Lasso constraint region = diamond (L1 ball)

Because the diamond has sharp corners, the solution often occurs on axes  $\rightarrow$  exact zeros.

- ✓ Lasso performs feature selection
- ✓ Ridge shrinks coefficients but rarely sets them to zero

## 7. Derive linear regression as MAP estimate.

Assume:

Likelihood:  $y | X, \beta \sim N(X\beta, \sigma^2 I)$

Prior (Gaussian):  $\beta \sim N(\mathbf{0}, \tau^2 I)$

Posterior maximization:  $\log p(\beta | y) \propto -\|y - X\beta\|^2 - \lambda\|\beta\|^2$

This equals the Ridge objective.

- ✓ OLS = MLE
- ✓ Ridge = MAP with Gaussian prior
- ✓ Lasso = MAP with Laplace prior

## 8. Analyze condition number of $X^T X$ .

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

If smallest eigenvalue  $\approx 0$ :

$$\kappa \rightarrow \infty$$

Large condition number  $\Rightarrow$  numerical instability.

Multicollinearity  $\rightarrow$  small eigenvalues  $\rightarrow$  large condition number.

- ✓ Ridge improves conditioning by adding  $\lambda$ .

## 9. Prove convexity of squared loss.

Loss:  $L(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2$

Gradient:  $\nabla L = -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta)$

Hessian:  $\nabla^2 L = 2\mathbf{X}^T\mathbf{X}$

Since  $\mathbf{X}^T\mathbf{X}$  is positive semidefinite:  $\nabla^2 L \geq \mathbf{0}$

- ✓ Squared loss is convex
- ✓ If full rank  $\rightarrow$  strictly convex  $\rightarrow$  unique global minimum

## 10. What happens when $p > n$ ?

When number of features exceeds number of samples:

- $\mathbf{X}^T\mathbf{X}$  is singular
- Infinite OLS solutions
- Model interpolates training data
- Severe overfitting
- High variance

### Solutions:

- Ridge regression
- Lasso
- Dimensionality reduction

### Modern Interpretation:

This setting is common in high-dimensional statistics and deep learning.