

Lecture Notes: Inference Under the Linear Model

Your Name/Course Name Here

April 22, 2025

1 From Least Squares Fitting to Statistical Inference

In our exploration of linear models so far, we've focused on the method of **Least Squares (LS)**. Given a design matrix \mathbf{X} (size $n \times (p + 1)$, where n is the number of observations and p is the number of predictor variables, plus one for an intercept) and a vector of observed responses \mathbf{Y} (size $n \times 1$), we found the coefficient vector $\hat{\boldsymbol{\beta}}$ that minimizes the sum of squared differences between the observed responses and the values predicted by the linear model. This LS estimator is given by:

$$\hat{\boldsymbol{\beta}} = \mathbf{A}\mathbf{Y}, \quad \text{where} \quad \mathbf{A} := (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

Assuming, of course, that $\mathbf{X}^\top \mathbf{X}$ is invertible, which typically holds if $n \geq p + 1$ and the columns of \mathbf{X} are linearly independent.

From this estimate, we defined two important vectors:

- The vector of **fitted values**: $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$. This represents the projection of the observed data \mathbf{Y} onto the column space of \mathbf{X} .
- The vector of **residuals**: $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}}$. This captures the part of the data *not* explained by the linear fit.

A fundamental geometric property we discovered is that the residual vector is orthogonal to the fitted values vector, $\mathbf{e} \perp \hat{\mathbf{Y}}$, and indeed, orthogonal to any vector in the column space of \mathbf{X} .

Remark 1.1 (Algebra, Not Statistics (Yet)). It's essential to remember that everything described above is derived from the algebraic goal of minimizing $\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$. We haven't needed to make any assumptions about probability distributions or random errors. These results hold for any given dataset (\mathbf{X}, \mathbf{Y}) .

Our goal now is to move beyond simply finding the "best fit" line or hyperplane. We want to perform **statistical inference**. This involves asking questions like:

- How uncertain is our estimate $\hat{\boldsymbol{\beta}}$? Can we construct confidence intervals for the true coefficients $\boldsymbol{\beta}$?
- Is a particular predictor variable significantly associated with the response? (i.e., can we test hypotheses like $H_0 : \beta_j = 0$?)
- How much variability in the response is inherent noise versus explained by the model?

To answer these questions, we need to introduce a probabilistic framework – the **linear model assumptions** – which describes how the data \mathbf{Y} are generated.

2 The Linear Model: Assumptions and Basic Properties

We now formally adopt the standard (or Gaussian) linear model.

Definition 2.1 (The Linear Model). The linear model assumes that the response vector \mathbf{Y} is generated according to:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad (1)$$

where:

- \mathbf{X} is the $n \times (p + 1)$ design matrix, treated as fixed and known.
- $\boldsymbol{\beta}$ is the $(p + 1) \times 1$ vector of *true, unknown* population coefficients.
- $\boldsymbol{\epsilon}$ is the $n \times 1$ vector of *unobserved random errors*, satisfying:

(LM1) **Zero Mean:** $\mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}$. (Equivalently, $\mathbb{E}[\epsilon_i] = 0$ for each $i = 1, \dots, n$).

(LM2) **Constant Variance (Homoscedasticity):** $\text{Var}(\epsilon_i) = \sigma^2$ for all i , where $\sigma^2 > 0$ is an unknown parameter.

(LM3) **Uncorrelated Errors:** $\text{Cov}(\epsilon_i, \epsilon_j) = 0$ for all $i \neq j$.

Assumptions (LM2) and (LM3) can be compactly written using the covariance matrix of the error vector: $\text{Cov}(\boldsymbol{\epsilon}) = \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top] - (\mathbb{E}[\boldsymbol{\epsilon}])(\mathbb{E}[\boldsymbol{\epsilon}])^\top = \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top] = \sigma^2 \mathbf{I}_n$, where \mathbf{I}_n is the $n \times n$ identity matrix.

Often, an additional assumption is made for exact inference:

(LM4) **Normality:** The errors ϵ_i are normally distributed. Combined with (LM1)-(LM3), this means $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$.

For the results in this section (calculating means and covariances), we only need assumptions (LM1)-(LM3). The normality assumption (LM4) will become crucial later when we discuss distributions of test statistics (like t-tests and F-tests).

Under these assumptions, \mathbf{Y} becomes a random vector. Let's find its mean and covariance matrix.

Mean of \mathbf{Y} : Using the linearity of expectation and (LM1):

$$\mathbb{E}[\mathbf{Y}] = \mathbb{E}[\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}] = \mathbb{E}[\mathbf{X}\boldsymbol{\beta}] + \mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{X}\boldsymbol{\beta} + \mathbf{0} = \mathbf{X}\boldsymbol{\beta}$$

(Note: $\mathbf{X}\boldsymbol{\beta}$ is considered a constant vector in this expectation, as \mathbf{X} is fixed and $\boldsymbol{\beta}$ is a fixed, albeit unknown, parameter vector).

Covariance of \mathbf{Y} : Using properties of covariance (adding a constant vector doesn't change covariance) and the definition $\text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n$:

$$\text{Cov}(\mathbf{Y}) = \text{Cov}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) = \text{Cov}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n$$

3 Statistical Properties of the LS Estimator $\hat{\beta}$

Now that \mathbf{Y} is a random vector, our LS estimator $\hat{\beta} = \mathbf{A}\mathbf{Y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ also becomes a random vector. We can investigate its statistical properties, specifically its mean and covariance matrix.

3.1 Mean of $\hat{\beta}$ (Unbiasedness)

Let's calculate the expected value of our estimator.

$$\begin{aligned}
 \mathbb{E}[\hat{\beta}] &= \mathbb{E}[\mathbf{A}\mathbf{Y}] \\
 &= \mathbf{A}\mathbb{E}[\mathbf{Y}] \quad (\text{since } \mathbf{A} \text{ is a constant matrix}) \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\beta) \quad (\text{substituting } \mathbb{E}[\mathbf{Y}] = \mathbf{X}\beta) \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{X})\beta \\
 &= \mathbf{I}_{p+1}\beta \\
 &= \beta
 \end{aligned}$$

This proves a fundamental result:

Proposition 3.1. *Under the linear model assumptions (LM1)-(LM3), the Least Squares estimator $\hat{\beta}$ is an **unbiased estimator** of the true coefficient vector β . That is, $\mathbb{E}[\hat{\beta}] = \beta$.*

This means that if we were to repeat our experiment many times and calculate $\hat{\beta}$ each time, the average of these estimates would converge to the true value β . This is a very desirable property for an estimator.

3.2 Covariance Matrix of $\hat{\beta}$

Next, let's find the covariance matrix of $\hat{\beta}$. Recall the property for a constant matrix \mathbf{B} and a random vector \mathbf{Z} : $\text{Cov}(\mathbf{B}\mathbf{Z}) = \mathbf{B} \text{Cov}(\mathbf{Z}) \mathbf{B}^\top$. We apply this with $\mathbf{Z} = \mathbf{Y}$ and $\mathbf{B} = \mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.

$$\begin{aligned}
 \text{Cov}(\hat{\beta}) &= \text{Cov}(\mathbf{A}\mathbf{Y}) \\
 &= \mathbf{A} \text{Cov}(\mathbf{Y}) \mathbf{A}^\top \\
 &= \mathbf{A}(\sigma^2 \mathbf{I}_n) \mathbf{A}^\top \quad (\text{substituting } \text{Cov}(\mathbf{Y}) = \sigma^2 \mathbf{I}_n) \\
 &= \sigma^2 \mathbf{A} \mathbf{A}^\top \\
 &= \sigma^2 \left[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right] \left[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \right]^\top \\
 &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \left((\mathbf{X}^\top)^\top ((\mathbf{X}^\top \mathbf{X})^{-1})^\top \right) \\
 &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} ((\mathbf{X}^\top \mathbf{X})^{-1}) \quad (\text{since } (\mathbf{X}^\top \mathbf{X})^{-1} \text{ is symmetric}) \\
 &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{X}) (\mathbf{X}^\top \mathbf{X})^{-1} \\
 &= \sigma^2 \mathbf{I}_{p+1} (\mathbf{X}^\top \mathbf{X})^{-1} \\
 &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}
 \end{aligned}$$

We have derived the covariance matrix of the LS estimator:

Proposition 3.2. *Under the linear model assumptions (LM1)-(LM3), the covariance matrix of the LS estimator $\hat{\beta}$ is given by:*

$$\boxed{\text{Cov}(\hat{\beta}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}}$$

Remark 3.3 (Interpreting $\text{Cov}(\hat{\beta})$). This $(p+1) \times (p+1)$ matrix is crucial for inference.

- The diagonal entries give the variances of the individual coefficient estimators: $\text{Var}(\hat{\beta}_j) = \sigma^2[(\mathbf{X}^\top \mathbf{X})^{-1}]_{jj}$. This measures the precision of each estimate.
- The off-diagonal entries give the covariances between different coefficient estimators: $\text{Cov}(\hat{\beta}_j, \hat{\beta}_k) = \sigma^2[(\mathbf{X}^\top \mathbf{X})^{-1}]_{jk}$. This tells us how the estimates tend to vary together.
- Notice the dependence on σ^2 : higher underlying noise variance leads to higher variance (less precision) in our estimates.
- The term $(\mathbf{X}^\top \mathbf{X})^{-1}$ reflects the influence of the experimental design or data structure. For example, multicollinearity (near linear dependence among columns of \mathbf{X}) tends to inflate the diagonal elements of $(\mathbf{X}^\top \mathbf{X})^{-1}$, increasing the variance of the corresponding coefficient estimates.

To use this covariance matrix for practical inference (like constructing confidence intervals or hypothesis tests), we need the value of σ^2 . Since σ^2 is typically unknown, we must estimate it from the data.

4 Estimating the Error Variance σ^2

4.1 Motivation and Definition

How can we estimate the underlying noise level σ^2 ? Intuitively, the residuals $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}}$ represent the discrepancy between our observations and the model's fit. The magnitude of these residuals should reflect the magnitude of the true errors $\boldsymbol{\epsilon}$.

A natural quantity to consider is the **Sum of Squared Residuals (SSR)**:

$$\text{SSR} = \|\mathbf{e}\|^2 = \sum_{i=1}^n e_i^2$$

Since $\sigma^2 = \text{Var}(\epsilon_i) = \mathbb{E}[\epsilon_i^2]$ (because $\mathbb{E}[\epsilon_i] = 0$), we might think SSR is related to $n\sigma^2$. However, the residuals e_i are not the same as the true errors ϵ_i . The residuals are calculated using the *estimated* coefficients $\hat{\beta}$, which depend on the data \mathbf{Y} . This process of estimating β "uses up" some information from the data.

Specifically, we estimated $p+1$ parameters (the components of β). It turns out that the appropriate divisor for SSR to get an unbiased estimate of σ^2 is not n , but $n - (p+1)$, the **residual degrees of freedom**.

Definition 4.1 (Unbiased Estimator of σ^2). The unbiased estimator of the error variance σ^2 , often denoted by $\hat{\sigma}^2$ or s^2 , is defined as:

$$\hat{\sigma}^2 := \frac{1}{n-p-1} \|\mathbf{e}\|^2 = \frac{1}{n-p-1} \sum_{i=1}^n e_i^2 = \frac{\text{SSR}}{n-p-1}$$

This quantity is also known as the **Residual Mean Square** (RMS or MSE).

We will now rigorously prove that this definition indeed yields an unbiased estimator.

Proposition 4.2. *Under the linear model assumptions (LM1)-(LM3), the estimator $\hat{\sigma}^2$ defined above is an unbiased estimator of σ^2 . That is, $\mathbb{E}[\hat{\sigma}^2] = \sigma^2$.*

We offer two distinct proofs of this important result. They provide different perspectives and utilize different mathematical tools.

4.2 Proof 1: Using Projection Matrices

This proof relies heavily on the geometric interpretation of least squares in terms of projections.

First Proof of Proposition 4.2. Let $M = \text{Im}(\mathbf{X})$ be the column space of the design matrix \mathbf{X} . This is the subspace of \mathbb{R}^n spanned by the columns of \mathbf{X} . Assuming \mathbf{X} has full column rank, the dimension of M is $p + 1$.

Recall that $\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ is the orthogonal projection of \mathbf{Y} onto M . Let $\mathbf{P} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ be the $n \times n$ projection matrix onto M . So, $\hat{\mathbf{Y}} = \mathbf{P}\mathbf{Y}$.

The residual vector is $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{P}\mathbf{Y} = (\mathbf{I}_n - \mathbf{P})\mathbf{Y}$. Let $\mathbf{Q} := \mathbf{I}_n - \mathbf{P}$. \mathbf{Q} is the projection matrix onto the orthogonal complement of M , denoted M^\perp . We know the following properties of \mathbf{Q} :

- It is symmetric: $\mathbf{Q}^\top = \mathbf{Q}$.
- It is idempotent: $\mathbf{Q}^2 = \mathbf{Q}\mathbf{Q} = \mathbf{Q}$.
- It annihilates vectors in M : $\mathbf{Q}\mathbf{X} = (\mathbf{I}_n - \mathbf{P})\mathbf{X} = \mathbf{X} - \mathbf{P}\mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0}$. (Since $\mathbf{P}\mathbf{X} = \mathbf{X}$ because columns of \mathbf{X} are in M).

Now, let's express the residual vector in terms of the true errors $\boldsymbol{\epsilon}$. Using the linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$:

$$\mathbf{e} = \mathbf{Q}\mathbf{Y} = \mathbf{Q}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) = \mathbf{Q}\mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\epsilon} = \mathbf{0} \cdot \boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\epsilon} = \mathbf{Q}\boldsymbol{\epsilon}$$

This crucial step shows that the residual vector \mathbf{e} is simply the projection of the true (unobserved) error vector $\boldsymbol{\epsilon}$ onto the subspace M^\perp , which is orthogonal to the space spanned by our predictors.

Now we compute the expected value of the Sum of Squared Residuals, $\mathbb{E}[\|\mathbf{e}\|^2]$.

$$\begin{aligned} \mathbb{E}[\|\mathbf{e}\|^2] &= \mathbb{E}[\|\mathbf{Q}\boldsymbol{\epsilon}\|^2] \\ &= \mathbb{E}[(\mathbf{Q}\boldsymbol{\epsilon})^\top (\mathbf{Q}\boldsymbol{\epsilon})] \\ &= \mathbb{E}[\boldsymbol{\epsilon}^\top \mathbf{Q}^\top \mathbf{Q} \boldsymbol{\epsilon}] \\ &= \mathbb{E}[\boldsymbol{\epsilon}^\top \mathbf{Q} \boldsymbol{\epsilon}] \quad (\text{using } \mathbf{Q}^\top = \mathbf{Q} \text{ and } \mathbf{Q}^2 = \mathbf{Q}) \end{aligned}$$

The term $\boldsymbol{\epsilon}^\top \mathbf{Q} \boldsymbol{\epsilon}$ is a quadratic form in the random vector $\boldsymbol{\epsilon}$. We can write it explicitly as $\sum_{i=1}^n \sum_{j=1}^n Q_{ij} \epsilon_i \epsilon_j$. By linearity of expectation:

$$\mathbb{E}[\boldsymbol{\epsilon}^\top \mathbf{Q} \boldsymbol{\epsilon}] = \mathbb{E} \left[\sum_{i=1}^n \sum_{j=1}^n Q_{ij} \epsilon_i \epsilon_j \right] = \sum_{i=1}^n \sum_{j=1}^n Q_{ij} \mathbb{E}[\epsilon_i \epsilon_j]$$

From the linear model assumptions (LM1)-(LM3), we know $\mathbb{E}[\epsilon_i] = 0$ and $\text{Cov}(\epsilon_i, \epsilon_j) = \mathbb{E}[\epsilon_i \epsilon_j]$. Thus:

$$\mathbb{E}[\epsilon_i \epsilon_j] = \begin{cases} \text{Var}(\epsilon_i) = \sigma^2, & \text{if } i = j \\ \text{Cov}(\epsilon_i, \epsilon_j) = 0, & \text{if } i \neq j \end{cases}$$

Substituting this into our sum:

$$\begin{aligned} \mathbb{E}[\|e\|^2] &= \sum_{i=1}^n \sum_{j=1}^n Q_{ij} \mathbb{E}[\epsilon_i \epsilon_j] \\ &= \sum_{i=1}^n Q_{ii} \mathbb{E}[\epsilon_i^2] + \sum_{i \neq j} Q_{ij} \mathbb{E}[\epsilon_i \epsilon_j] \\ &= \sum_{i=1}^n Q_{ii} (\sigma^2) + \sum_{i \neq j} Q_{ij} (0) \\ &= \sigma^2 \sum_{i=1}^n Q_{ii} \end{aligned}$$

The sum $\sum_{i=1}^n Q_{ii}$ is precisely the trace of the matrix \mathbf{Q} , denoted $\text{tr}(\mathbf{Q})$. So, we have $\mathbb{E}[\|e\|^2] = \sigma^2 \text{tr}(\mathbf{Q})$.

What is the trace of the projection matrix \mathbf{Q} ? The trace of any projection matrix is equal to the dimension of the subspace it projects onto. \mathbf{Q} projects onto M^\perp . We know $\dim(\mathbb{R}^n) = n$ and $\dim(M) = p+1$ (assuming \mathbf{X} has full rank). By the rank-nullity theorem or properties of orthogonal complements, $\dim(M^\perp) = \dim(\mathbb{R}^n) - \dim(M) = n - (p+1)$. Therefore, $\text{rank}(\mathbf{Q}) = \dim(M^\perp) = n - p - 1$. Since the trace of a projection matrix equals its rank, we have $\text{tr}(\mathbf{Q}) = n - p - 1$.

Substituting this back into our expectation calculation:

$$\mathbb{E}[\|e\|^2] = \sigma^2(n - p - 1)$$

Finally, we can find the expectation of our proposed estimator $\hat{\sigma}^2$:

$$\mathbb{E}[\hat{\sigma}^2] = \mathbb{E}\left[\frac{1}{n - p - 1} \|e\|^2\right] = \frac{1}{n - p - 1} \mathbb{E}[\|e\|^2] = \frac{1}{n - p - 1} \sigma^2(n - p - 1) = \sigma^2$$

This completes the proof that $\hat{\sigma}^2$ is an unbiased estimator of σ^2 . \square

Remark 4.3 (Degrees of Freedom). The quantity $n - p - 1$ is often called the **degrees of freedom for error** (or residual degrees of freedom). It represents the number of independent pieces of information in the data that are available for estimating the variance σ^2 , after having already used $p + 1$ degrees of freedom to estimate the coefficients in β .

4.3 Proof 2: Using a General Lemma about Expected Norms

An alternative, perhaps more abstract, proof uses a general result about the expected squared norm of a random vector.

Lemma 4.4. *Let \mathbf{Z} be a random vector in \mathbb{R}^k with mean $\mu_{\mathbf{Z}} = \mathbb{E}[\mathbf{Z}]$ and covariance matrix $\text{Cov}(\mathbf{Z})$. Then,*

$$\mathbb{E}[\|\mathbf{Z}\|^2] = \text{tr}(\mathbb{E}[\mathbf{Z}\mathbf{Z}^\top]) = \text{tr}(\text{Cov}(\mathbf{Z}) + \mu_{\mathbf{Z}}\mu_{\mathbf{Z}}^\top)$$

As a special case, if the mean is zero ($\mu_{\mathbf{Z}} = \mathbf{0}$), then $\mathbb{E}[\|\mathbf{Z}\|^2] = \text{tr}(\text{Cov}(\mathbf{Z}))$.

Proof of Lemma 4.4. We start with the definition of the squared Euclidean norm $\|\mathbf{Z}\|^2 = \mathbf{Z}^\top \mathbf{Z}$.

$$\begin{aligned}
\mathbb{E}[\|\mathbf{Z}\|^2] &= \mathbb{E}[\mathbf{Z}^\top \mathbf{Z}] \\
&\stackrel{(a)}{=} \mathbb{E}[\text{tr}(\mathbf{Z}^\top \mathbf{Z})] \quad (\text{Since } \mathbf{Z}^\top \mathbf{Z} \text{ is a } 1 \times 1 \text{ matrix, its trace is itself}) \\
&\stackrel{(b)}{=} \mathbb{E}[\text{tr}(\mathbf{Z} \mathbf{Z}^\top)] \quad (\text{Using the cyclic property of trace: } \text{tr}(AB) = \text{tr}(BA)) \\
&\stackrel{(c)}{=} \text{tr}(\mathbb{E}[\mathbf{Z} \mathbf{Z}^\top]) \quad (\text{Linearity of trace and expectation allows swapping them}) \\
&\stackrel{(d)}{=} \text{tr}(\text{Cov}(\mathbf{Z}) + \boldsymbol{\mu}_\mathbf{Z} \boldsymbol{\mu}_\mathbf{Z}^\top) \quad (\text{Using the definition } \text{Cov}(\mathbf{Z}) = \mathbb{E}[\mathbf{Z} \mathbf{Z}^\top] - \boldsymbol{\mu}_\mathbf{Z} \boldsymbol{\mu}_\mathbf{Z}^\top)
\end{aligned}$$

This establishes the lemma. The special case follows immediately by setting $\boldsymbol{\mu}_\mathbf{Z} = \mathbf{0}$. \square

Now we apply this lemma to prove Proposition 4.2.

Alternative Proof of Proposition 4.2. We want to compute $\mathbb{E}[\|\mathbf{e}\|^2]$. We apply Lemma 4.4 with the random vector $\mathbf{Z} = \mathbf{e}$. To do this, we first need the mean and covariance matrix of \mathbf{e} .

Recall from the first proof that $\mathbf{e} = \mathbf{Q}\mathbf{Y}$, where $\mathbf{Q} = \mathbf{I}_n - \mathbf{P}$. The mean of \mathbf{e} is:

$$\mathbb{E}[\mathbf{e}] = \mathbb{E}[\mathbf{Q}\mathbf{Y}] = \mathbf{Q}\mathbb{E}[\mathbf{Y}] = \mathbf{Q}(\mathbf{X}\boldsymbol{\beta}) = (\mathbf{Q}\mathbf{X})\boldsymbol{\beta} = \mathbf{0} \cdot \boldsymbol{\beta} = \mathbf{0}$$

So, the mean vector $\boldsymbol{\mu}_\mathbf{e}$ is the zero vector.

The covariance matrix of \mathbf{e} is:

$$\begin{aligned}
\text{Cov}(\mathbf{e}) &= \text{Cov}(\mathbf{Q}\mathbf{Y}) \\
&= \mathbf{Q} \text{Cov}(\mathbf{Y}) \mathbf{Q}^\top \\
&= \mathbf{Q}(\sigma^2 \mathbf{I}_n) \mathbf{Q}^\top \quad (\text{using } \text{Cov}(\mathbf{Y}) = \sigma^2 \mathbf{I}_n) \\
&= \sigma^2 \mathbf{Q} \mathbf{I}_n \mathbf{Q}^\top \\
&= \sigma^2 \mathbf{Q} \mathbf{Q}^\top \\
&= \sigma^2 \mathbf{Q} \mathbf{Q} \quad (\text{since } \mathbf{Q} \text{ is symmetric}) \\
&= \sigma^2 \mathbf{Q} \quad (\text{since } \mathbf{Q} \text{ is idempotent})
\end{aligned}$$

So, $\text{Cov}(\mathbf{e}) = \sigma^2 \mathbf{Q}$.

Now we apply the special case of Lemma 4.4 (since $\boldsymbol{\mu}_\mathbf{e} = \mathbf{0}$):

$$\begin{aligned}
\mathbb{E}[\|\mathbf{e}\|^2] &= \text{tr}(\text{Cov}(\mathbf{e})) \\
&= \text{tr}(\sigma^2 \mathbf{Q}) \\
&= \sigma^2 \text{tr}(\mathbf{Q}) \quad (\text{Linearity of trace})
\end{aligned}$$

As established in the first proof using properties of projection matrices, $\text{tr}(\mathbf{Q}) = n - p - 1$. Therefore,

$$\mathbb{E}[\|\mathbf{e}\|^2] = \sigma^2(n - p - 1)$$

Finally, the expectation of our estimator $\hat{\sigma}^2$ is:

$$\mathbb{E}[\hat{\sigma}^2] = \mathbb{E}\left[\frac{1}{n - p - 1} \|\mathbf{e}\|^2\right] = \frac{1}{n - p - 1} \mathbb{E}[\|\mathbf{e}\|^2] = \frac{1}{n - p - 1} \sigma^2(n - p - 1) = \sigma^2$$

This provides a second confirmation that $\hat{\sigma}^2$ is an unbiased estimator for σ^2 . \square

With these results – the unbiasedness of $\hat{\beta}$ and $\hat{\sigma}^2$, and the covariance matrix of $\hat{\beta}$ – we have laid the groundwork for constructing confidence intervals and hypothesis tests concerning the regression coefficients β , which are central tasks in statistical inference for linear models.