

Linear Regression

Professor Jie Peng, PhD

Department of Statistics

University of California, Davis

Recap: Fitted Values and Residuals in Matrix Form

- ▶ Fitted values vector: $n \times 1$ column vector:

$$\widehat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = \mathbf{H}\mathbf{Y},$$

where $\mathbf{H} := \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ is called the **hat matrix**.

- ▶ Residuals vector: $n \times 1$ column vector:

$$\mathbf{e} = \mathbf{Y} - \widehat{\mathbf{Y}} = (\mathbf{I}_n - \mathbf{H})\mathbf{Y}.$$

- ▶ Note that, fitted values vector $\widehat{\mathbf{Y}}$ and residuals vector \mathbf{e} are linear transformations of the observations vector \mathbf{Y} .

Recap: Hat Matrix

The hat matrix

$$\mathbf{H}_{n \times n} := \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

and the matrix

$$\mathbf{I}_n - \mathbf{H} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

are $n \times n$ **projection matrices**. Meaning that they satisfy:

- ▶ **Symmetric:** $\mathbf{H}' = \mathbf{H}$, $(\mathbf{I}_n - \mathbf{H})' = \mathbf{I}_n - \mathbf{H}$
- ▶ **Idempotent:** $\mathbf{H}^2 := \mathbf{H}\mathbf{H} = \mathbf{H}$, $(\mathbf{I}_n - \mathbf{H})^2 = \mathbf{I}_n - \mathbf{H}$.

Moreover, $\text{rank}(\mathbf{H}) = 2$, $\text{rank}(\mathbf{I}_n - \mathbf{H}) = n - 2$ (provided that X_i s are not all equal).

LS Estimation: Mean and Variance

Review: Linear Transformations of Random Vector

If \mathbf{Z} is an $r \times 1$ random vector, and \mathbf{A} is an $s \times r$ non-random matrix, then

$$\underset{s \times 1}{\mathbf{W}} = \underset{s \times r}{\mathbf{A}} \underset{r \times 1}{\mathbf{Z}}$$

is an $s \times 1$ random vector with

$$\mathbf{E}\{\mathbf{W}\} = \mathbf{E}\{\mathbf{AZ}\} = \mathbf{AE}\{\mathbf{Z}\}$$

$$\sigma^2\{\mathbf{W}\} = \sigma^2\{\mathbf{AZ}\} = \mathbf{A}\sigma^2\{\mathbf{Z}\}\mathbf{A}'$$

If further \mathbf{B} is a $t \times r$ non-random matrix, then

$$\text{Cov}(\mathbf{AZ}, \mathbf{BZ}) = \mathbf{A}\sigma^2\{\mathbf{Z}\}\mathbf{B}'$$

LS Estimation: Expectations

- ▶ LS estimator is unbiased:

$$\mathbf{E}\{\hat{\beta}\} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{E}\{\mathbf{Y}\} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta = \beta$$

- ▶ Expectation of the fitted values:

$$\mathbf{E}\{\widehat{\mathbf{Y}}\} = \mathbf{E}\{\mathbf{X}\hat{\beta}\} = \mathbf{X}\mathbf{E}\{\hat{\beta}\} = \mathbf{X}\beta = \mathbf{E}\{\mathbf{Y}\}$$

- ▶ Expectation of the residuals:

$$\mathbf{E}\{\mathbf{e}\} = \mathbf{E}\{\mathbf{Y} - \widehat{\mathbf{Y}}\} = \mathbf{E}\{\mathbf{Y}\} - \mathbf{E}\{\widehat{\mathbf{Y}}\} = \mathbf{0}_n$$

LS Estimation: Variance-Covariance Matrices

Variance-covariance of the LS estimator:

$$\begin{aligned}\sigma^2\{\hat{\beta}\} &= \sigma^2\{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}\} = ((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\sigma^2\{\mathbf{Y}\}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')' \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2 \begin{bmatrix} \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} & -\frac{\bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ -\frac{\bar{X}}{\sum_{i=1}^n (X_i - \bar{X})^2} & \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2} \end{bmatrix}\end{aligned}$$

What is the covariance between $\hat{\beta}_0$ and $\hat{\beta}_1$? What happens if

$\bar{X} = 0$?

- Variance-covariance of the fitted values:

$$\sigma^2\{\widehat{\mathbf{Y}}\} = \mathbf{H}\sigma^2\{\mathbf{Y}\}\mathbf{H}' = \sigma^2\mathbf{H}$$

- Variance-covariance of the residuals:

$$\sigma^2\{\mathbf{e}\} = (\mathbf{I}_n - \mathbf{H})\sigma^2\{\mathbf{Y}\}(\mathbf{I}_n - \mathbf{H})' = \sigma^2(\mathbf{I}_n - \mathbf{H})$$

Are residuals uncorrelated? Do they have the same variance?

How about the fitted values? What are the covariances

between the residuals and fitted values?

Simple Regression: Geometric Interpretation

Some Notations

- ▶ Let $\mathbf{1}_n$ denote the n-dimensional column vector of ones;
- ▶ Let $\mathbf{x} = (X_1, \dots, X_n)^T$ denote the n-dimensional column vector of X values;
- ▶ The design matrix \mathbf{X} for simple regression is formed by these two column vectors:

$$\mathbf{X} = (\mathbf{1}_n, \mathbf{x})$$

Column Space of the Design Matrix

- ▶ Let $\text{col}(\mathbf{X})$ denote the set of linear combinations of the two column vectors of \mathbf{X} :

$$\text{col}(\mathbf{X}) = \{\mathbf{v} \in \mathbb{R}^n : \text{there exists } c_0, c_1 \in \mathbb{R}, \text{ s.t., } \mathbf{v} = c_0 \mathbf{1}_n + c_1 \mathbf{x}\}.$$

- ▶ $\text{col}(\mathbf{X})$ is referred to as the *column space of the design matrix* \mathbf{X} .
- ▶ Note that $\text{col}(\mathbf{X})$ forms a **linear subspace** of \mathbb{R}^n .

Projection to $\text{col}(\mathbf{X})$ by Hat Matrix

The hat matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ “orthogonally projects” vectors in \mathbb{R}^n to $\text{col}(\mathbf{X})$ in the following sense:

For (any) $\mathbf{w} \in \mathbb{R}^n$:

- ▶ $\mathbf{H}\mathbf{w} \in \text{col}(\mathbf{X})$, i.e., we can find two constants $c_0, c_1 \in \mathbf{R}$ such that $\mathbf{H}\mathbf{w} = c_0\mathbf{1}_n + c_1\mathbf{x}$.
- ▶ $\mathbf{w} - \mathbf{H}\mathbf{w} \perp \text{col}(\mathbf{X})$, i.e., for any $\mathbf{v} \in \text{col}(\mathbf{X})$, $(\mathbf{w} - \mathbf{H}\mathbf{w})'\mathbf{v} = 0$.

Notes: “ \in ” is read as “belongs to”; “ \perp ” is read as “is orthogonal to”.

What is $\mathbf{H}\mathbf{X}$? What is $\mathbf{H}\mathbf{1}_n$, $\mathbf{H}\mathbf{x}$? What is $\mathbf{H}\mathbf{v}$ for $\mathbf{v} \in \text{col}(\mathbf{X})$?

Optional Reading: Proof of the Projection Properties of \mathbf{H}

- ▶ By definition of \mathbf{H} , $\mathbf{H}\mathbf{w} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{w}$. Now, let $\mathbf{c} := (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{w}$ (which is a 2×1 vector). We can see $\mathbf{H}\mathbf{w} = \mathbf{X}\mathbf{c} = c_0\mathbf{1}_n + c_1\mathbf{x}$, where c_0, c_1 are the 1st and 2nd elements of the vector \mathbf{c} .
- ▶ By definition of $\text{col}(\mathbf{X})$, for (any) $\mathbf{v} \in \text{col}(\mathbf{X})$, we can find a 2×1 vector \mathbf{c} such that $\mathbf{v} = \mathbf{X}\mathbf{c}$. So $(\mathbf{w} - \mathbf{H}\mathbf{w})'\mathbf{v} = (\mathbf{w} - \mathbf{H}\mathbf{w})'\mathbf{X}\mathbf{c} = \mathbf{w}'(\mathbf{I}_n - \mathbf{H})\mathbf{X}\mathbf{c}$. Note that, $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \mathbf{X} - \mathbf{H}\mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0}_{n \times 2}$. So $(\mathbf{w} - \mathbf{H}\mathbf{w})'\mathbf{v} = 0$.

Fitted Values and Residuals

- ▶ The fitted values vector $\widehat{\mathbf{Y}} = \mathbf{H}\mathbf{Y} = \hat{\beta}_0 \mathbf{1}_n + \hat{\beta}_1 \mathbf{x}$ belongs to $\text{col}(\mathbf{X})$.
- ▶ The residuals vector $\mathbf{e} = \mathbf{Y} - \mathbf{H}\mathbf{Y}$ is orthogonal to $\text{col}(\mathbf{X})$.
- ▶ Also note that, $\mathbf{1}_n, \mathbf{x} \in \text{col}(\mathbf{X})$.
- ▶ Therefore:

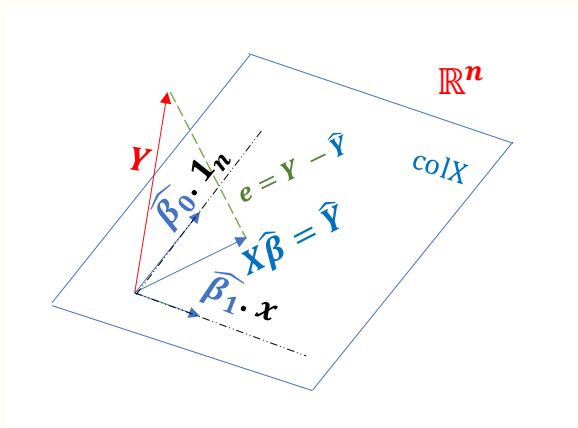
$$\langle \mathbf{e}, \mathbf{1}_n \rangle = \sum_{i=1}^n e_i = 0$$

$$\langle \mathbf{e}, \mathbf{x} \rangle = \sum_{i=1}^n x_i e_i = 0$$

$$\langle \mathbf{e}, \widehat{\mathbf{Y}} \rangle = \sum_{i=1}^n \hat{Y}_i e_i = 0$$

Geometric Interpretation of Regression

Figure: To regress Y onto X is to “orthogonally project” the response vector \mathbf{Y} onto the column space of the design matrix \mathbf{X}



Sums of Squares: Matrix Form

Error Sum of Squares

$$SSE = \sum_{i=1}^n e_i^2$$

can be expressed in matrix form as:

$$SSE = \mathbf{e}'\mathbf{e} = \mathbf{Y}'(\mathbf{I}_n - \mathbf{H})'(\mathbf{I}_n - \mathbf{H})\mathbf{Y} = \mathbf{Y}'(\mathbf{I}_n - \mathbf{H})\mathbf{Y}.$$

It can be shown that:

- ▶ $\mathbf{I}_n - \mathbf{H}$ is a projection matrix.
- ▶ $\text{rank}(\mathbf{I}_n - \mathbf{H}) = n - 2 = df(SSE)$.

Total Sum of Squares

$$SSTO = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n Y_i^2 - n(\bar{Y})^2$$

can be expressed in matrix form as:

$$SSTO = \mathbf{Y}'\mathbf{Y} - \frac{1}{n}\mathbf{Y}'\mathbf{J}_n\mathbf{Y} = \mathbf{Y}'\left(\mathbf{I}_n - \frac{1}{n}\mathbf{J}_n\right)\mathbf{Y},$$

where \mathbf{J}_n is the $n \times n$ matrix of ones. It can be shown that:

- ▶ $\mathbf{I}_n - \frac{1}{n}\mathbf{J}_n$ is a projection matrix.
- ▶ $\text{rank}(\mathbf{I}_n - \frac{1}{n}\mathbf{J}_n) = n - 1 = df(SSTO)$.

Regression Sum of Squares

$$SSR = \sum_{i=1}^n (\widehat{Y}_i - \bar{Y})^2$$

can be expressed in matrix form as:

$$\begin{aligned} SSR &= (\widehat{\mathbf{Y}} - \bar{\mathbf{Y}})' (\widehat{\mathbf{Y}} - \bar{\mathbf{Y}}), & \bar{\mathbf{Y}} &:= \frac{1}{n} \mathbf{J}_n \mathbf{Y} \\ &= \mathbf{Y}' \left(\mathbf{H} - \frac{1}{n} \mathbf{J}_n \right)' \left(\mathbf{H} - \frac{1}{n} \mathbf{J}_n \right) \mathbf{Y} \\ &= \mathbf{Y}' \left(\mathbf{H} - \frac{1}{n} \mathbf{J}_n \right) \mathbf{Y} \end{aligned}$$

It can be shown that:

- ▶ $\mathbf{H} - \frac{1}{n} \mathbf{J}_n$ is a projection matrix.
- ▶ $\text{rank}(\mathbf{H} - \frac{1}{n} \mathbf{J}_n) = 1 = df(SSR)$.

Review: Matrix Trace Operation and Properties

- ▶ $\mathbf{M} = (m_{ij})$ is an $s \times s$ square matrix, its trace is the summation of its diagonal elements: $\text{Tr}(\mathbf{M}) = \sum_{i=1}^s m_{ii}$. Specifically, a scalar c may be viewed as a 1×1 square matrix, and $\text{Tr}(c) = c$.
- ▶ Trace is a linear operator: For two square matrices \mathbf{A} and \mathbf{B} (of the same dimension) and a scalar c : $\text{Tr}(\mathbf{A} + \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B})$,
 $\text{Tr}(c \cdot \mathbf{A}) = c \cdot \text{Tr}(\mathbf{A})$.
- ▶ Consequently, for a random (square) matrix \mathbf{A} :
 $E(\text{Tr}(\mathbf{A})) = \text{Tr}(E(\mathbf{A}))$.
- ▶ If \mathbf{A} is an $s \times t$ matrix and \mathbf{B} is a $t \times s$ matrix, then $\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA})$.

Deriving Expectation of SSE

Note that: $SSE = \text{Tr}(SSE) = \text{Tr}(\mathbf{Y}'(\mathbf{I}_n - \mathbf{H})\mathbf{Y})$, so

$$\begin{aligned}E(SSE) &= E(\text{Tr}(\mathbf{Y}'(\mathbf{I}_n - \mathbf{H})\mathbf{Y})) = E(\text{Tr}((\mathbf{I}_n - \mathbf{H})\mathbf{Y}\mathbf{Y}')) \\&= \text{Tr}(E((\mathbf{I}_n - \mathbf{H})\mathbf{Y}\mathbf{Y}')) = \text{Tr}((\mathbf{I}_n - \mathbf{H})E(\mathbf{Y}\mathbf{Y}')) \\&= \text{Tr}((\mathbf{I}_n - \mathbf{H})(\sigma^2\mathbf{I}_n + \mathbf{X}\beta\beta'\mathbf{X}')) \\&= \sigma^2 \text{Tr}(\mathbf{I}_n - \mathbf{H}) + \text{Tr}((\mathbf{I}_n - \mathbf{H})\mathbf{X}\beta\beta'\mathbf{X}') \\&= (n - 2)\sigma^2\end{aligned}$$

The last equality is because $\text{Tr}(\mathbf{I}_n - \mathbf{H}) = \text{Tr}(\mathbf{I}_n) - \text{Tr}(\mathbf{H}) = n - 2$

and $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \mathbf{X} - \mathbf{H}\mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0}$. *Why?*

Can you derive $E(SSR)$ and $E(SSTO)$ in a similar fashion?

Optional Reading: Eigen-decomposition of Projection Matrices

- ▶ A projection matrix \mathbf{P} can be decomposed as $\mathbf{P} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$, where \mathbf{Q} is an orthogonal matrix formed by \mathbf{P} 's eigenvectors and $\mathbf{\Lambda}$ is a diagonal matrix consisting of its eigenvalues.
- ▶ Moreover, \mathbf{P} 's eigenvalues are either 1 or 0.
- ▶ Consequently, the number of its nonzero eigenvalues equals its trace equals its rank.
- ▶ In simple linear regression:

$$\text{rank}(\mathbf{H}) = \text{tr}(\mathbf{H}) = 2, \quad \text{rank}(\mathbf{I}_n - \mathbf{H}) = \text{tr}(\mathbf{I}_n - \mathbf{H}) = n - 2$$

Optional Reading: Sampling Distribution of SSE under Normal Error Model

- ▶ $\mathbf{I}_n - \mathbf{H}$ is a projection matrix with rank $n - 2 \implies$ its spectral decomposition looks like:

$$\mathbf{I}_n - \mathbf{H} = \mathbf{Q}^T \mathbf{\Lambda} \mathbf{Q},$$

where $\mathbf{\Lambda} = \text{diag}\{1, \dots, 1, 0, 0\}$ and \mathbf{Q} is an orthogonal matrix.

- ▶ $(\mathbf{I}_n - \mathbf{H})\mathbf{X} = \mathbf{0} \implies$

$$\mathbf{e} = (\mathbf{I}_n - \mathbf{H})\mathbf{Y} = (\mathbf{I}_n - \mathbf{H})(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) = (\mathbf{I}_n - \mathbf{H})\boldsymbol{\epsilon}$$

- ▶ $SSE = \mathbf{e}^T \mathbf{e} = \boldsymbol{\epsilon}^T (\mathbf{I}_n - \mathbf{H}) \boldsymbol{\epsilon} = (\mathbf{Q}\boldsymbol{\epsilon})^T \boldsymbol{\Lambda} (\mathbf{Q}\boldsymbol{\epsilon})$
- ▶ Let $\mathbf{z} = \mathbf{Q}\boldsymbol{\epsilon}$, then

$$SSE = \sum_{i=1}^n \lambda_i z_i^2 = \sum_{i=1}^{n-2} z_i^2.$$

where λ_i is the i th diagonal element of $\boldsymbol{\Lambda}$ and we have $\lambda_i = 1$ for $i = 1, \dots, n-2$ and $\lambda_i = 0$ for $i = n-1, n$.

- ▶ Moreover

$$\mathbf{E}(\mathbf{z}) = \mathbf{Q}\mathbf{E}\{\boldsymbol{\epsilon}\} = \mathbf{0}, \quad \sigma^2\{\mathbf{z}\} = \mathbf{Q}\sigma^2\{\boldsymbol{\epsilon}\}\mathbf{Q}^T = \sigma^2\mathbf{Q}\mathbf{Q}^T = \sigma^2\mathbf{I}_n$$

So under Normal error model, z_i s are i.i.d. $N(0, \sigma^2)$.

- ▶ Thus, by the definition of chi-squares distributions,

$$SSE \sim \sigma^2 \chi_{(n-2)}^2.$$

Confidence Interval for σ^2

Under the Normal error model, a $(1 - \alpha)100\%$ -confidence interval for σ^2 is :

$$\left[\frac{SSE}{\chi^2(1 - \alpha/2; n - 2)}, \frac{SSE}{\chi^2(\alpha/2; n - 2)} \right]$$

- $\chi^2(1 - \alpha/2; n - 2)$ and $\chi^2(\alpha/2; n - 2)$ are the $(1 - \alpha/2)100\%$ percentile and $(\alpha/2)100\%$ percentile of the chi-squares distribution with $n - 2$ degrees of freedom.

How to derive a $(1 - \alpha)100\%$ -confidence interval for the error standard deviation σ ? What is a 95% confidence interval for the error variance in the “Heights example”?

Optional Reading: Deriving C.I. for σ^2

This is through the pivotal quantity:

$$\frac{SSE}{\sigma^2}$$

and the fact that it follows $\chi^2_{(n-2)}$ distribution since $SSE \sim \sigma^2 \chi^2_{(n-2)}$.

Probability Density Curves of χ^2 Distributions

