

L 2. Simple Linear Regression

- Fitted value and residuals
- Estimate of σ^2
- Distribution of $\hat{\beta}_1$
- Sum of Squares Phenomenon

Recall:

Simple Linear Regression Model

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \varepsilon_i \stackrel{i.i.d}{\sim} N(0, \sigma^2)$$

Least-square estimate $\hat{\beta}_0$ and $\hat{\beta}_1$, minimize $S(\beta_0, \beta_1) = \sum (y_i - \beta_0 - \beta_1 x_i)^2$

then the fitted value of the i^{th} observation is

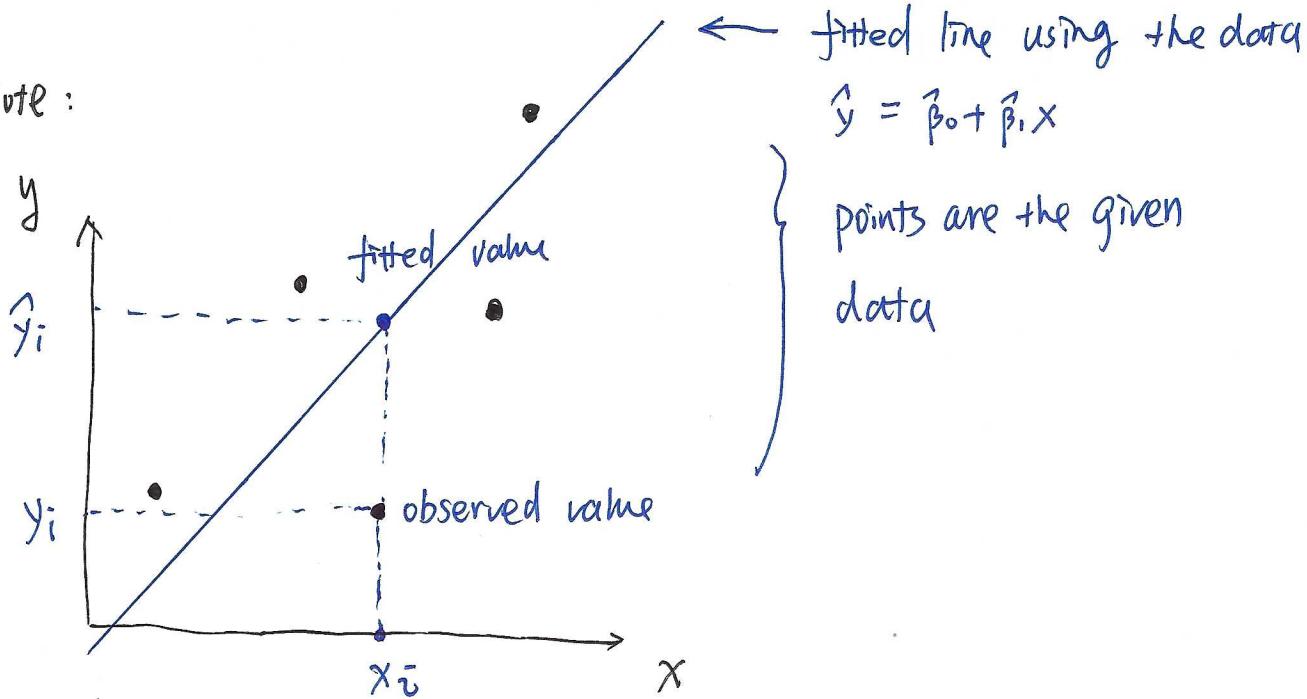
$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

For any given value $\underline{X=x}$, we would estimate $E(y) = \beta_0 + \beta_1 x$
doesn't have to be in the data

by

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad \leftarrow \begin{array}{l} \text{Best prediction of } y \text{ given} \\ X=x \text{ using the existing data} \end{array}$$

Note :



- $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ predicts $E(y|x=x) = b_0 + b_1 x$
- $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ predicts the value of y at $X=x$

What's the difference? The estimated value is the same, but the estimation error (variance) is different. will discuss

In L3.

I. Fitted Value and Residuals

Now we have $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

and fitted value $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

Then ε_i can be estimated by residuals

$$\hat{\varepsilon}_i = e_i = y_i - \hat{y}_i$$

Can you see
 e_i is unbiased estimate
of ε_i ?

Properties of \hat{y}_i and e_i

(1) The fitted residuals must sum to zero

$$\sum e_i = \sum (y_i - \hat{y}_i)$$

$$= \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = \underline{n(\bar{y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{x}) = 0}$$

— from LSE

Consequently

$$(2) \quad \sum e_i = \sum \hat{y}_i$$

(3) $\sum e_i^2$ is at minimum, there's no better fitted line to make it smaller — from LSE

$$(4) \quad \sum e_i x_i = 0 \quad — \text{from LSE}$$

$$(5) \quad \sum e_i \hat{y}_i = 0$$

II. Estimate of σ^2

Recall we have some unfinished business: unknown σ^2

If we use e_i to estimate ϵ_i , then we may use it to estimate σ^2 .

— recall in stats, estimate of σ^2 always starts with χ^2 distribution.

Sum of Squared Errors $SSE = \sum_{i=1}^n e_i^2 = \sum (y_i - \hat{y}_i)^2$

where $y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$, since it's somehow the sum of squares of normal distributions, we grab the result:

$$\frac{SSE}{\sigma^2} \sim \chi^2(n-2)$$

The proof is lengthy and not required, but I will provide it in the materials for the ones who're interested.

Therefore, $E\left(\frac{SSE}{\sigma^2}\right) = n-2 \Rightarrow E\left(\frac{SSE}{n-2}\right) = \sigma^2$

Mean Squared Errors

$$MSE = \frac{SSE}{n-2} = \frac{\sum e_i^2}{n-2}$$

Is the unbiased estimator of σ^2

III. Distribution of $\hat{\beta}_1$

Recall the facts we have known about $\hat{\beta}_1$:

$$E(\hat{\beta}_1) = \beta_1, \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}, \quad \text{where } \sigma^2 \text{ can be estimated by MSE.}$$

- When σ^2 is known (magic!)

$$Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

$$\text{so } \hat{\beta}_1 = \sum k_i y_i \sim N(\beta_1, \frac{\sigma^2}{S_{xx}})$$

- When σ^2 is unknown (common!), and we estimate it by MSE, then the studentized distribution becomes:

$$\frac{\hat{\beta}_1 - \beta_1}{S(\hat{\beta}_1)} \sim t(n-2)$$

Hint: $S(\hat{\beta}_1) = \sqrt{\frac{\text{MSE}}{\sum (x_i - \bar{x})^2}}$

Why do we like the distribution of $\hat{\beta}_1$?

It can give us the t-test for the significance of X

- Idea: $y = \beta_0 + \beta_1 x + \varepsilon$

If X is a "significant" factor to y ,
when x changes, y should change with it
significantly. β_1 is the slope for this change.

- Hypothesis: $H_0: \beta_1 = 0$ v.s. $H_1: \beta_1 \neq 0$

When H_1 is true, β_1 is significantly different
from 0, therefore a change in x cause significant
change in y — X is significant to Y .

- Test statistic: $t_0 = \frac{\hat{\beta}_1}{\sqrt{\frac{MSE}{S_{xx}}}}$

- Decision rule: When $|t_0| \geq t_{\alpha/2}$, $df=n-2$

or p-value $\leq \alpha$ \Rightarrow Reject \Rightarrow X is
 H_0 significant
to y .

IV. Sum of Squares Phenomenon - ANOVA and R²

- ANOVA: testing the hypothesis related to equality of more than one parameters
 - More useful in Multiple Linear Regression
 - Just to develop the basic math here
 - In SLR, it's also testing $H_0: \beta_1 = 0$

(a) Goal: The variation in y can be explained by the variation in regression line and the variation in the error. OR.

$$\text{Total Sum of Square in } Y = \frac{\text{Regression}}{\text{sum of squares}} + \frac{\text{Error}}{\text{sum of squares}}$$

$SST = SSR + SSE$
$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (\hat{y}_i - y_i)^2$

$$\begin{aligned}
 \text{Proof: } \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\
 &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})
 \end{aligned}$$

Note that

$$\begin{aligned}
 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) &= \sum_{i=1}^n (y_i - \hat{y}_i)\hat{y}_i - \sum_{i=1}^n (y_i - \hat{y}_i)\bar{y} \\
 &= \sum_{i=1}^n e_i \hat{y}_i - \bar{y} \sum_{i=1}^n e_i \\
 &= 0 - 0 \\
 &= 0
 \end{aligned}$$

so

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

or

$$SST = SSR + SSE$$

(b) F test for $H_0: \beta_1 = 0$

use distributions

	sum of squares	d.f.	mean squares
Regression	SSR	1	$MSR = SSR/1$
Error	SSE	$n-2$	$MSE = SSE/(n-2)$
Total	SST	$n-1$	

We know that $\frac{SSE}{\sigma^2} \sim \chi^2(n-2)$

and $\frac{SSR}{\sigma^2} \sim \chi^2(1)$

Ideas

- If the regression line can explain "more" of the variation and the error explains "less" variable variation, then it means including x in the line is significant
- This comparison can be done using the ratio $\frac{MSR}{MSE}$

Since $\frac{\chi^2(df_1)/df_1}{\chi^2(df_2)/df_2}$ gives a F distribution. so we can do test!

- To test $H_0: \beta_1 = 0$ v.s. $H_1: \beta_1 \neq 0$,

think how it affects MSR and MSE:

We know that $E(MSE) = \sigma^2$

We can prove that $E(MSR) = \sigma^2 + \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$

$$\text{Hint: Show } MSR = \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\text{then } E(MSR) = E(\hat{\beta}_1^2) \sum (x_i - \bar{x})^2$$

- When $\beta_1 = 0$: $E(MSR) = E(MSE) = \sigma^2$

- When $\beta_1 \neq 0$: $E(MSR) > E(MSE)$

So we are expecting $\frac{MSR}{MSE} >> 1$.

- The test statistic $F = \frac{MSR}{MSE} \sim F(1, n-2)$ under H_0

- Decision Rule: Reject H_0 when $F > F_{1-\alpha}$ ($df_1 = 1$, $df_2 = n-2$)
or when p-value $< \alpha$.

□

R Squared

A similar idea can be measured to see the goodness of fit of the model

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

- This is known as the " coefficient of determination"
- This measures how much variation in y (stated by $SST = \sum (y_i - \bar{y})^2$), is explained by the variation in regression (stated by $SSR = \sum (\hat{y}_i - \bar{y})^2$) and how much unexplainable part is contained in the variation of error (stated by $SSE = \sum (\hat{y}_i - y_i)^2$).
- It can be proved that

$$R^2 = r_{xy}^2 = \frac{S_{xy}}{\sqrt{S_{xx}} \sqrt{S_{yy}}} \quad \begin{aligned} &\text{the square of} \\ &\text{the sample correlation} \\ &\text{between } x \text{ and } y. \end{aligned}$$

- $0 \leq R^2 \leq 1$, closer to 1 indicates better fit.

Summary and Takeaways of this lecture.

Fitted value $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

Residual $e_i = y_i - \hat{y}_i$

Fitted line: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$

$$E(\hat{Y}_i) = \hat{\beta}_0 + \hat{\beta}_1 x_i = E(Y_i)$$

$\text{Var}(\hat{Y}_i)$ will be discussed later

$$E(e_i) = \epsilon_i$$

$$E(\bar{Y}) = \beta_0 + \beta_1 \bar{X}$$

sum of squared errors

$$\text{SSE} = \sum e_i^2 = \sum (y_i - \hat{y}_i)^2$$

$$\frac{\text{SSE}}{\sigma^2} \sim \chi^2(n-2)$$

mean squared errors

$$\text{MSE} = \frac{\text{SSE}}{n-2}$$

$$E(\text{MSE}) = \sigma^2$$

SSE measures the variation in errors;

MSE is the unbiased estimate of σ^2

when σ^2 is known:

$$\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{xx}})$$

when σ^2 is unknown

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\text{MSE}}{S_{xx}}}} \sim t(n-2)$$

Use the t distribution

of $\hat{\beta}_1$ to test:

$H_0: \beta_1 = 0$ v.s. $H_1: \beta_1 \neq 0$

for the significance of X to y .

ANOVA Test for $H_0: \beta_1 = 0$ v.s. $\beta_1 \neq 0$

	S.S.	d.f.	M.S	F stats
Regression	$SSR = \sum (\hat{y}_i - \bar{y})^2$	1	$MSR = \frac{SSR}{1}$	$\frac{MSR}{MSE}$
Error	$SSE = \sum (\hat{y}_i - y_i)^2$	$n-2$	$MSE = \frac{SSE}{n-2}$	
Total	$SST = \sum (y_i - \bar{y})^2$	$n-1$		

Reject H_0 if $F_{\text{stat}} > F_{1-\alpha}, df_1=1, df_2=n-2$

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}, \text{ closer to 1 better fit}$$



- ANOVA is essential comparing two models:

$$H_0: \beta_1 = 0 \Leftrightarrow y = \beta_0 + \varepsilon$$

$$H_1: \beta_1 \neq 0 \Leftrightarrow y = \beta_0 + \beta_1 x + \varepsilon$$

If reject H_0 , the model including x factor is significantly better than the null model.

$$\begin{aligned} \bullet \text{ In the simple regression case, } F_{\text{stat}} &= \frac{\hat{\beta}_1^2 \sum (x_i - \bar{x})^2}{MSE} = \left(\frac{\hat{\beta}_1}{\sqrt{\frac{MSE}{S_{xx}}}} \right)^2 \\ &= t_{\text{stat}}^2 \end{aligned}$$

the test result is often the same.