

# Scenario 1:

- ① Do the residuals indicate any problems with the fit?
- ② If densities ( $y$ ) is not reported exactly, how might this affect the fit?
- ③ What if the blocks were not measured in a random order?

Regression :

- ① Draw Scatterplot:  $\Rightarrow$  Identify which one is response variable, which one is explanatory variable.
- + Residual Plot  $\Rightarrow$  Analyze relationship : sign: positive / negative  
shape: linear / quadratic / etc.  
level of correlation: strong / moderately strong / weak

② Calculate Correlation : Def: strength of the linear association between two variables.  $[-1, 1]$

$$\text{Population mean: } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + \dots + x_n}{n}$$

$$\text{Population variance: } s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\text{Population correlation: } R \text{ (Empirical Correlation Coefficient)} \\ = \frac{\frac{1}{n} \sum_{i=1}^n (x_i \cdot y_i) - (\bar{x} \cdot \bar{y})}{s_x \cdot s_y}$$

$$\text{slope of regression: } \hat{\beta}_1 = \frac{s_y}{s_x} \cdot R$$

Interpretation: For each additional % point in explanatory variable ( $x$ ), we would expect the % point in response variable ( $y$ ) increase on average  $\hat{\beta}_1$  points.

$$\text{Intercept: } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Intuition: A regression line always passes through  $(\bar{x}, \bar{y})$

Interpretation: With no explanatory variable ( $x$ ), we expect on average to have  $\hat{\beta}_0$  of response variable ( $y$ ).

③ Fitting :

$\Rightarrow$  Residual: distance between observed and predicted

Goal: We want a line with small residuals

$$\text{Formula: } \hat{\epsilon}_i = y_i - \hat{\beta}_1 x_i - \hat{\beta}_0 \\ (\epsilon_i = y_i - \hat{y}_i)$$

$\Rightarrow$  least square : Def: = Minimizes the sum of the squares of data point to regression line.

$$\text{formula} = \min \sum_{i=1}^n (\text{Residual}_i)^2 \\ = \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Why Least Square :

- ① most commonly used
- ② Easier to compute
- ③ A residual twice as large as another is usually more than twice as bad.

$\Rightarrow$  least absolute deviation: Def: Minimizes the sum of the data point to regression line.

$$\text{formula} = \min \sum_{i=1}^n |\text{Residual}_i|$$

$$= \min_{\beta_0, \beta_1} \sum_{i=1}^n |y_i - \beta_0 - \beta_1 x_i|$$

$$\Rightarrow \sum_{i=1}^n h(\text{residuals}_i)$$

$$h(\text{residuals}_i) = \begin{cases} |x| & |x| > \tau \\ x^2 & |x| < \tau \end{cases}$$



Note: The statements are not causal, unless the study is a randomized controlled experiment.

## Scenario 2

① adding bands around least square line  $\Rightarrow$  make estimates  
(Prediction Interval + Confidence Interval)  
Note: several measurement of gain was taken for one density

### ④ Predicting:

- \* Fitted Value  $\Leftrightarrow$  Old Value
- New Value  $\Leftrightarrow$  New Value

### \* Definition:

$\Rightarrow$  Prediction : Give explanatory variable ( $x$ ) to predict response variable ( $y$ )  
(Plug in  $x$  to the model)

$\Rightarrow$  Extrapolation: Apply a model estimate to values outside of the realm of the original data ( $x=0$ )  $\leftarrow$  unreachable data

\* Least Square:  $\epsilon_i \sim N(0, \sigma^2)$   $\hat{\beta}_i$  is MLE

Least Absolute Value:  $\epsilon_i \sim \text{expl}$   $\hat{\beta}_i$  is MLE

### \* Requirement:

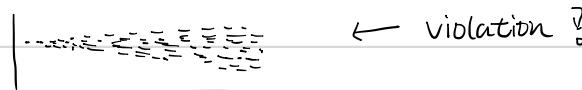
① Plot Residuals + Scatterplot, check linearity

② Plot QQ-plot + residuals histogram, Check if residuals plot  $\sim$  Normal

③ If use scatterplot to plot residuals, should only see scatter plots without clear pattern. ( $x \Rightarrow$  explanatory variable,  $y \Rightarrow$  residuals)

[homoscedasticity]  $\Rightarrow$  constant variability

[131] Residual Plot



$\leftarrow$  violation?

### \* $R^2$ :

$\Rightarrow$  Definition: The strength of the fit of a linear model

$\Rightarrow$  formula:  $= \left( \frac{\sum_{i=1}^n (x_i \cdot y_i) - (\bar{x} \cdot \bar{y})}{S_x \cdot S_y} \right)^2 \Rightarrow$  smaller than R

$\Rightarrow$  Interpretation: [what percent of variability in the response variable is explained by the model. The remainder of the variability is explained by variables not included in the model / due to randomness in data]

$R^2 \times 100\%$  of the variability in the response variable is explained by the model.

$\Rightarrow$  Adjusted  $R^2$ : (If we have more than 2 variables  $\Rightarrow$  gives more accurate interpretation)

### \* Categorical Data:

$\hookrightarrow$  reference level = (categorical data = 0) = (y = intercept)

$\hookrightarrow$  Use LS to test if X and Y are independent

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

Which region (northeast, midwest, west, or south) has the lowest poverty percentage?

$$H_0: \beta_1 = 0 \quad H_1: \beta_1 \neq 0$$

Hypothesis testing

probability  
of getting  
non-zero  
 $\beta_1$  value

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	9.50	0.87	10.94	0.00
region4midwest	slope { 0.03 based on reference level}	$\beta_0$	1.15	0.02
region4west	1.79	$\beta_1$	1.13	0.12
region4south	4.16	1.07	3.87	0.00

$$SE_{\hat{\beta}} = \frac{\sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$$\frac{|\hat{\beta} - \beta_1|}{SE_{\hat{\beta}}} \sim T_{n-2}$$

$$= \frac{| \text{point estimate} - \text{null} |}{SE}$$

#### ⑤ Outlier:

\* Distinguish if that outlier is influential / high leverage

Relationship between biological IQ & foster IQ

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	9.20760	9.29990	0.990	0.332
bioIQ	0.90144	0.09633	9.358	1.2e-09

Residual standard error: 7.729 on 25 degrees of freedom

n-2

Multiple R-squared: 0.7779, Adjusted R-squared: 0.769

F-statistic: 87.56 on 1 and 25 DF, p-value: 1.204e-09

Confidence Interval:

$$\text{Estimate} \pm SE \times t_{n-2}$$