Week 8     { Linear Regression }

Linear Model     $Y \sim X$

$\Rightarrow$ parametric modeling

Linear Regression

ANOVA

$Y$ : conti num
$X$ : cate / conti
     ∨    ∨

$\begin{cases} Y : \text{conti numerical} \\ X : \text{categorical} \end{cases} \Rightarrow H_0 : \mu_A^A = \mu_B^B = \mu_C^W$
                          $\Rightarrow H_a : \quad ''$

Goal : Find a relationship between $X \cdot Y$
            & quantify its relationship

           MP

          * prediction

1-way ANOVA

1) $Y$ : conti.    $X$ : categorical

2) Normality for each level

3) Equal var

4) indep samples

Linear Regression

1) $X, Y$  linear relationship

2) Normality ← conditionmy on $X$

3) Equal var

4) indep samples

$$E(Y) = \beta_0 + \beta_1 X$$

$Y$: Final score

$X$: # of hours study

(Assumption 1)

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

mean of $Y$ for if $X = x$
$\rightarrow$ explained by $x$ (model)

Error Individual

can't be even

1.5

Estimated reg. line

① Using Samples,

$\hat{\beta_0}$ $\hat{\beta_1}$

goal = Estimate $\beta_0$ $\beta_1$

$$\hat{Y} = \hat{\beta_0} + \hat{\beta_1} X$$

prediction

$\hat{y} = 70$.

Estima : $\hat{\beta_0}$ $\hat{\beta_1}$

$$[ H_0: \beta_1 = 0 \quad vs. \quad H_a: \beta_1 \neq 0 ]$$

$\begin{cases} X, Y \\ no\ relationship \end{cases}$

sig. relat

: inferential problem

## Heteroscedasticity

income

$$E(Y) = \beta_0 + \beta_1 X$$

a

age

## Non-normality

in

$$E(Y) = \beta_0 + \beta_1 X$$

age

## Estimation



$\boxed{y}$ $\bar{y}$

$x \to x+1$
$\hat{\beta_1}$

$Y = \hat{\beta_0} + \hat{\beta_1} X$

Least Square Estimation (LSE)
or Ordinary Least Square (OLS)

$\Downarrow$

$\left( \text{minimize} \dfrac{\text{squared}}{\text{vertical distance}} \right)$

---

$Y$: cirrhosis
$X$: alcohol

$Y \sim \alpha X_1, \alpha X_2 \cdots \alpha X_p$
$\qquad \beta_1 \quad \beta_2 \qquad \beta_5$

— write the estimated regression line

$$\hat{y} = -5.99 + 1.9779 \times \text{alcohol}$$
$$\hat{\beta_0} \qquad \hat{\beta_1}$$

$\tilde{I}$ ① Test model significance (F-test)

$H_0: \beta_1 = 0$   $H_a: \beta_1 \ne 0$   $H_0$: linear model is not useful
$(\beta_1 = \beta_2 = \beta_3 = \beta_4 = \beta_5 = 0)$   $H_a$: at least one $\beta \ne 0$

② Test individual $X$ significance (t-test)

$H_0: \beta_1 = 0$ vs. $H_a: \beta_1 \ne 0$

③ Interpretation of $\hat{\beta_1}$



④ $\boxed{R^2} = \dfrac{SS_{model}}{SS_{Total}}$

$SS_{Total} = SS_{model} + SS_{Error}$

⑤ Model diagnostic
(0)

model assumption

Influential point diagnosis

Exploratory Analysis

Scatter plot

< correlation ⓟ : pearson correlation

$\underline{X} \cdot \underline{Y}$

(1) range.    sign + −        $|\rho|$

$-1 \leq \rho \leq 1$

(2) Only measure (linearity)

$\boxed{\rho}$

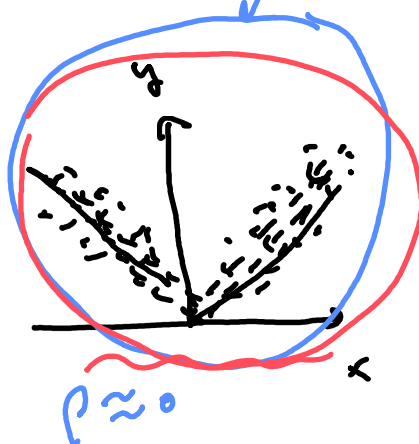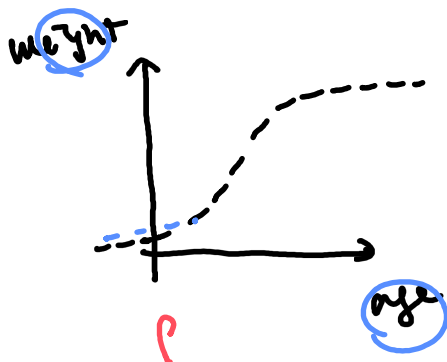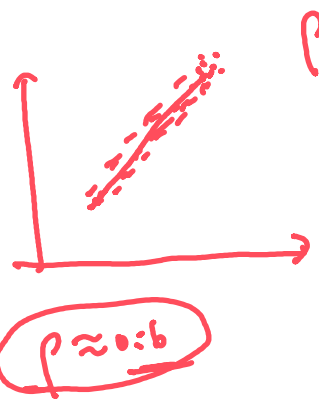$\rho \approx 0$          → $\rho \approx 0$

→ $\rho \approx 0$  linear relar

(3) Sensitivity to outliers

$$\left\{ \rho = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} \right\}$$

[ pearson

spearman corre

robust

sensiti

$\rho$

$\rho \approx 0.6$

correlation

weight

age

$\rho$

$\rho \approx 0$

$\dfrac{SS_{mol}}{SS_{Total}}$    $SS_{Total} = SS_M + SS_{Err}$

⊕ $\boxed{R^2}$ and Model significance

→ [ $\beta = 0$ or $\beta \neq 0$ ]

p-value

Goodness-of-fit measure

prediction power

SS_Error small        SS_Err large
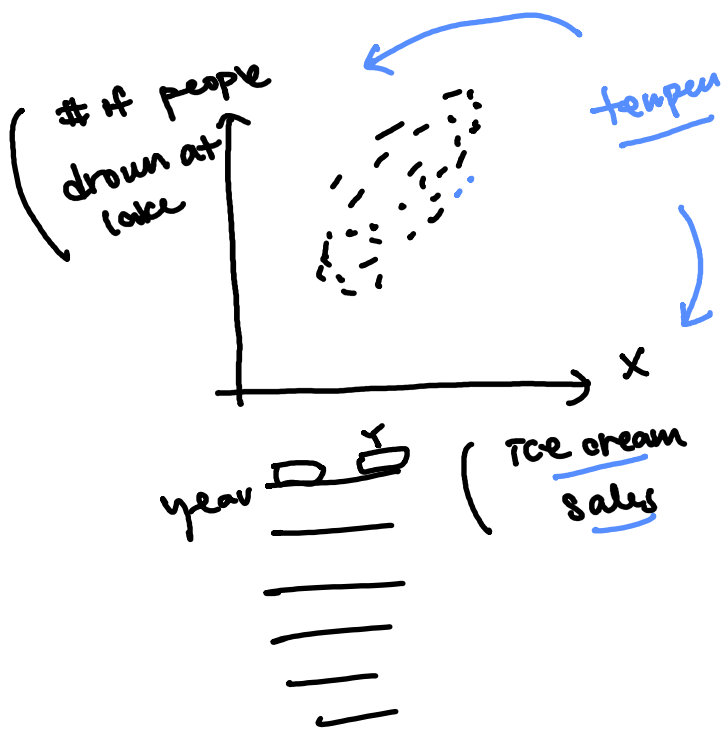


$R^2 \uparrow$        $R^2 \downarrow$

$R^2 = 0.1$

(eg) ( Smoking ~ lung — cancer )

andell $R^2 = 0.3$
$R^2 = 0.8$
mean $R^2 = 0.99$    r
mse

⊠ Causal effect ↔ linear relationship



# of people
drown at
lake        temper

( $H_0: \beta = 0$ vs. $H_a: \beta \neq 0$ )

X

year  $\square$ $\dot{Y}$   ( Ice cream
                              sales

Diagnostics — (1) Model assumption check — Linear X, Y
                                            Normality
                                            Equal var.

        (2) Influence point check ✓



$\left(E(Y) = \beta_0 + \beta_1 x\right)$

$\hat{y} = \hat{\beta_0} + \hat{\beta_1} x$

$\boxed{\text{Residual}} = Y - \hat{Y}$     Expect to see — Normal
                                                          Equal var.

[standardized residual]

Standard
0.1
$N(0,1)$

95%

-2   0   +2

√ standardized residual

$\sqrt{2} \approx 1.5$

Normality check by ( QR plot
                     square-root standardized residual
                              ↓
                     ( below 1.5 )

→ Equal variance check by

$$\left( \frac{Residual}{\sqrt{standardized\ residual}} \right.$$

$\hat{y} = \hat{\beta_0} + \hat{\beta_1} x$

$\sqrt{Stan\ Ro}$

$y - \hat{y}$