Week 13. Logistic Reg
14. Thanksgiving
15. Review
16. Final

$$\text{odds} \triangleq \frac{P(Y=1)}{1-P(Y=1)} = \begin{array}{c} 1 \\ > 1 \\ < 1 \end{array} \qquad P(Y=1) = \frac{1}{2}$$

$\boxed{Y=1}$

$$\boxed{\text{Odds Ratio}} = \frac{\text{odds }(Y=1 \mid \boxed{\text{Female}})}{\text{odds }(Y=1 \mid \boxed{\text{Male}})} = \begin{array}{c} 1 \\ > 1 \\ < 1 \end{array}$$

$x$
gender etc

F > M

F < M

$0 \leq \boxed{\hat{P}(Y=1 \mid x=\hat{x})} \leq 1$

**Linear Reg.**

conti

$\boxed{Y}$  $X_1, X_2 \cdots X_p$

$\boxed{Y} = \beta_0 + \beta_1 \dot{x_1} + \cdots + \beta_p x_p \boxed{-\varepsilon}$

$Y \sim \boxed{N}(\beta_0 + \beta_1 x_1 + \cdots \beta_p x_p, \sigma^2)$

$\hat{Y} = \hat{\beta_0} + \hat{\beta_1} x_1 - \hat{\beta_p} x_p$

$E(Y)$

**Logistic Reg**

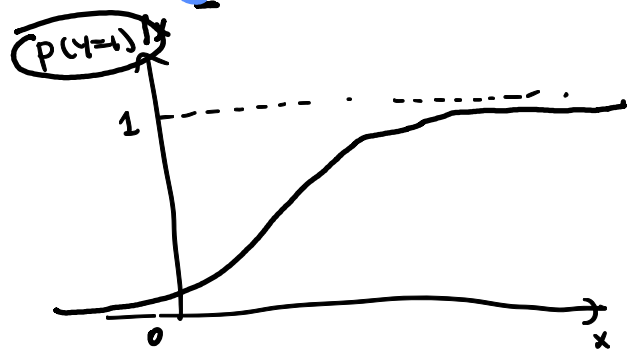$\boxed{Y}$  $X_1, X_2 \cdots X_p$

link

$$\log\left[\frac{P(Y=1)}{1-P(Y=1)}\right] = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p$$

$-\infty < \cdots < \infty$

$$P(Y=1 \mid X) = \frac{\exp(\beta_0 + \beta_1 x + \cdots \beta_p x)}{1 + \exp(\beta_0 + \beta_1 x + \cdots \beta_p x)}$$

$P(Y=1 \mid x) \doteq 0.8$

$Y \sim \text{Bernoulli}(P(Y=1 \mid x))$



$\beta_0 + \beta_1 x_1$

$Y$

$x=x$   $x=x$



$P(Y=1)$

1

0   $x$

Interpretation of $\beta_1$

$$\log\left(\frac{P(Y=1)}{1-P(Y=1)}\right) = \beta_0 + \underline{\beta_1 x}$$

- $x$ is categorical 

  Ref $\quad 0 \leftarrow$ male

  Comp. $\quad 1 \checkmark$ fems

  Female

$$\frac{Odds(Y=1 \mid comp)}{odds(Y=1 \mid ref)} = OR(comp \text{ vs. } ref) = e^{\beta_1}$$

fem $\qquad$ male

$\boxed{\beta_1 = 0}$

$\cdot \beta_1 = 0 \quad e^{\beta_1} = 1$

$\underline{\beta_1 > 0} \quad e^{\beta_1} > 1$

$\beta_1 < 0 \quad < 1$

- $x$ is continuous $\quad x \quad x \to x+1$

$$\frac{odds(Y=1 \mid fib = x+1)}{odds(Y=1 \mid fib = x)} = OR(fib = x+1 \text{ vs. } fib = x)$$

$$= e^{\hat{\beta_1}} = e^{1.91} = > 1 \qquad \text{fib} \uparrow \ P \uparrow$$

$$\left.\frac{odds(Y=1 \mid gaw = x+1)}{odds(Y=1 \mid gaw = x)}\right\} = OR(g = x+1 \text{ vs. } g = x)$$

$$= e^{\hat{\beta_1}} = e^{0.155} > 1 \qquad gaw \uparrow \ P \uparrow.$$

| Model Selection |

$$X_1 \cdots X_p$$

$$\begin{bmatrix} \text{stepwise} \\ \text{forward} \\ \text{backward} \end{bmatrix} \text{selection} + \begin{matrix} \text{AIC} \\ \text{BIC} \end{matrix}$$

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

$R^2$

$R^2 \approx 0.1$

| Model Fitting |

— significance of $x_j$

— Interpretation of $\hat{\beta_j}$

— Goodness-of-fit $\Big[$ pseudo $R^2$

$\hat{P}(Y=1/x)$

$H_0:$ model fits the data well
$H_a: \quad '' \quad$ does

std
$\left( \boxed{Y} - \boxed{\hat{P}(Y=1/x)} \right)$

$\searrow N(0,1)$

| Model Diagnostics |

deviance
pearson

Residual plot $\Big\langle$

$\pm 2$

Cook's D

| x | Y |
|---|---|
| / | 0 |
| / | 0 |
| / | 0 |

| prediction |

$\boxed{\hat{P}(Y=1|x)}$

$\hat{\beta_0} \; \hat{\beta_1} \cdots \hat{\beta_p}$

$$\hat{P}(Y=1|) = \frac{\exp( \cdots )}{1 + \exp(\hat{\beta_0} + \hat{\beta_1} \cdots)}$$

$$[X_1 \quad X_2 \quad X_3 \quad X_4]$$

Backward Selection + AIC / BIC

AIC: the smalle $(2*P$ + models w penalt) error

① Full Model $glm(Y \sim X_1 + \cdots + X_4)$     AIC: 500

② (i) w/o $X_1$     (ii) w/o $X_2$     (iii) w/o $X_3$     (iv) w/o $X_4$

$Y \sim X_2 + X_3 + X_4$    $Y \sim X_1 + X_3 + X_4$    $Y \sim X_1 + X_2 + X_4$    $Y \sim X_1 + X_2 + X_3$

AIC = 450      AIC = 470      AIC = 499      AIC = 480

Remove $X_1$

$$[X_2 \quad X_3 \quad X_4]$$

③ (i) w/o $X_1 . X_2$     (ii) w/o $X_1 . X_3$     (iii) w/o $X_1 . X_4$

AIC $Y \sim X_3 + X_4$    $Y \sim X_2 + X_4$    $Y \sim X_2 + X_3$

AIC = 620      AIC = 510      AIC = 505

455

$$[Y \sim X_2 + X_3 + X_4]$$

Forward selection + AIC

① (i) $Y \sim X_1$
AIC = 700

(ii) $Y \sim X_2$
AIC = 800

(iii) $Y \sim X_3$
AIC = 900

(iv) $Y \sim X_4$
AIC = 750

Add $X_1$

$(X_2 \quad X_3 \quad X_4)$

② (i) $Y \sim X_1 + X_2$
AIC $\approx$ 650

(ii) $Y \sim X_1 + X_3$
AIC = 600

(iii) $Y \sim X_1 + X_4$
AIC = 700

Add $X_3$

$\sim$ AIC
$\sim$ BIC

error $\boxed{|||}$ + $2p$

$\log 300$

$n$

$\log n \cdot p$

$p \uparrow$ : comp$\uparrow$

Amputation Data    significance level    $\alpha = 0.1$

amputation ~ illness-severity + diabetes + ulcers

⊛ who will have the highest chance?

$\boxed{\hat{P}(Y=1 \mid X=x)}$

$OR(\text{L vs. M})$

L vs. 0

- Illness-severity ( L    M    (H) )  ←ref

$$OR(L \text{ vs. } H) = \frac{odds(Y=1 \mid L)}{odds(Y=1 \mid H)} = e^{-2.19} = 0.11 < 1 \quad \begin{bmatrix} \underline{L < H} & ✓ \\ \underline{M < H} & ✓ \end{bmatrix}$$

$$OR(M \text{ vs. } H) = \frac{odds(Y=1 \mid M)}{odds(Y=1 \mid H)} = e^{-0.67} < 1$$

- Diabetes    ( Uncontrolled vs. Controlled )

$$OR(un \text{ vs. } con) = \frac{odds(Y=1 \mid un)}{odds(Y=1 \mid con)} = e^{1.83} > 1$$

$$\underline{Uncon > Con} ↙$$

ref

- Ulcers    ( 1 vs. 0 )

$$OR(1 \text{ vs. } 0) = e^{2.18} \qquad \underline{1 > 0} ↙$$

$\alpha = 0.1$

$\left( \underline{High, \; Uncon, \; \_\_\_} \right)$

if using significance level $\alpha = 0.05$

$L < H$

$M < H$

$\rightarrow$ sig.

$\rightarrow$ insy.

$H_0 : \beta_{zM} = 0$

$\downarrow$

$M = H$

(High or Mode) . uncont   $u/w = 1$

$Y$ ~ $\xi$illness severity

$\downarrow$

prediction   (plasma data)

cut off 0.35

$\log\left[\dfrac{\hat{P}(Y=1)}{1-\hat{P}(Y=1)}\right] = -1.84 + 1.83 \times \text{Fib}$

| Fib | esr | $\hat{P}(esr=1)$ | $\widehat{esr}$ | $\widehat{esr}$ |
|-----|-----|------------------|-----------------|-----------------|
| 2.52 | 0 | 0.3 | 0 | 0 |
| 2.56 | 0 | 0.2 | 0 | 0 |
| 2.19 | 0 | 0.6 | 1 | 1 |
| 2.18 | 0 | 0.4 | 0 | 1 |
| 3.41 | 0 | 0.3 | 0 | 0 |
| ⋮ | | | | |

$\hat{P}(Y=1) = \dfrac{\exp(-1.84 + 1.83 \times \text{Fib})}{1 + \exp(-1.84 + 1.83 \times \text{Fib})}$

if $\hat{P} \leq 1$

$\dfrac{\uparrow}{\downarrow} \quad 0.5 \quad \begin{array}{l}\to 1 \\ \to 0\end{array}$

Y=0
1

$n = 50,000$

$\begin{cases} Y=0 \quad 49,990 \\ Y=1 \quad \boxed{10} \end{cases}$

$\hat{P}$

① $\boxed{\text{0.5 cut off}}$

misclassification

$\dfrac{Y. \quad \hat{Y}}{}$

$\dfrac{1}{5} \quad 20\%$

② prop (Y=1)

sample = 0.35

misclass

$\dfrac{2}{5}$

|  | $\hat{Y}=0$ | $\hat{Y}=1$ |
|--|-------------|-------------|
| Y=1 | → False Negative | |
| Y=0 | | → False positive |

− :
+ :

Cut-off 0.5                    cut-off : sample proportion

rate
50,000

$Y$    $\hat{Y}$

misc  10/50,000

0.001

False Negative = 1

$P(\hat{Y}=0 \mid Y=1)$

$= \frac{10}{10} = 1$

$Y=1$    10

X ✓

1: disease
0: No-disease

$Y$    $\hat{Y}$

False Positive

$P(\hat{Y}=1 \mid Y=0)$

How to choose the optimal cut-off?

(1) ROC curve

(2) Cross-validation          10%

n=500

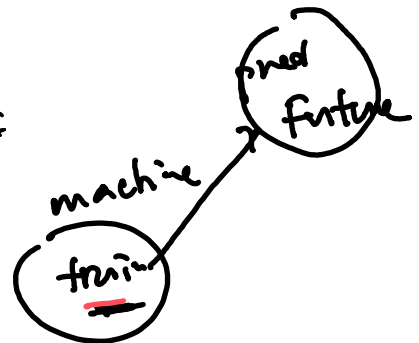90%

450

test

fit by

$\hat{\beta}_0 - \hat{\beta}_p \left[ \hat{P}(Y=1 \mid 0) \right]$

cut-off

machine → pred  future

train

cut-off tuning