

Binomial distribution $\rightarrow m/e \quad \hat{\pi} = \frac{y}{n}$
 s.e. $\sqrt{\hat{\pi}(1-\hat{\pi})/n}$
 $w = 2 \left[e(\hat{\pi}_1, \hat{\pi}_2) - e(\hat{\pi}, \hat{\pi}) \right]$
 since $w \sim \chi^2_{p=1}$,
 p-value $P[X^2 > w]$, $X^2 \sim \chi^2_1$

influence of x on π

\rightarrow age influence the satisfaction (old customers tends to have higher chances to be satisfied)

\hookrightarrow age is continuous: how can we see the influence of a continuous variable on the response (probability of being among satisfied customers)

π as a function of x
 \downarrow
 logistic regression $\quad \hat{\pi} = \frac{e^{\eta(x)}}{1 + e^{\eta(x)}}$ $\quad \eta(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$
 \uparrow logistic function $\quad \uparrow$ linear predictor
 $\pi \in [0; 1]$, $\eta(x) \in (-\infty; +\infty)$

in the figure $\eta(x) = \beta_0 + \beta_1 x$

\rightarrow logistic regression is a special case of GLM (Generalized Linear Models)

\rightarrow logistic \rightarrow binomial distribution

$g(E[Y|x]) = \eta(x)$
 \uparrow link function

Other distributions
 - gamma
 - poisson
 - exponential family

$\hat{\pi} = \frac{e^{\eta(x)}}{1 + e^{\eta(x)}}$

$\hat{\pi} + e^{\eta(x)} \hat{\pi} = e^{\eta(x)}$
 $e^{\eta(x)} - e^{\eta(x)} \hat{\pi} = \hat{\pi}$

$e^{\eta(x)} (1 - \hat{\pi}) = \hat{\pi}$

$e^{\eta(x)} = \frac{\hat{\pi}}{1 - \hat{\pi}}$

$\frac{\hat{\pi}}{1 - \hat{\pi}} = \text{odds}$

$\hat{\pi} = \frac{\text{odds}}{1 + \text{odds}}$

$\eta(x) = \log\left(\frac{\hat{\pi}}{1 - \hat{\pi}}\right)$

link function: logit

$g(x)$ for binomial response

Example

$$x \in \{0, 1\}$$

old
young

$$\eta(x) = \beta_0 + \beta_1 x$$

$$= 1.5299 + (-0.5446)x$$

$$= 1.5299 - 0.5446x$$

old $\eta(x) = 1.5299 - 0.5446 \times 0 = \underline{1.5299}$ ← baseline

$$\eta(x) = \log\left(\frac{\pi}{1-\pi}\right) \leftarrow \text{log-odds}$$

$$\pi = \frac{e^{\eta(x)}}{1 + e^{\eta(x)}} \rightarrow \hat{\pi}_{\text{old}} = \frac{e^{\hat{\beta}_0}}{1 + e^{\hat{\beta}_0}} = 0.822$$

young: $\eta(x) = 1.5299 - 0.5446 \times 1 = 0.9853$

$$\rightarrow \hat{\pi}_{\text{young}} = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1}} = 0.728$$

reverse:

$$\hat{\beta}_0 = \log\left(\frac{\hat{\pi}_{\text{old}}}{1 - \hat{\pi}_{\text{old}}}\right)$$

$$\hat{\beta}_0 + \hat{\beta}_1 = \log\left(\frac{\hat{\pi}_{\text{young}}}{1 - \hat{\pi}_{\text{young}}}\right)$$

$$\hat{\beta}_1 = \log\left(\frac{\hat{\pi}_{\text{young}}}{1 - \hat{\pi}_{\text{young}}}\right) - \hat{\beta}_0$$

$$= \log\left(\frac{\hat{\pi}_{\text{young}}}{1 - \hat{\pi}_{\text{young}}}\right) - \log\left(\frac{\hat{\pi}_{\text{old}}}{1 - \hat{\pi}_{\text{old}}}\right) \quad \beta_1$$

$$= \log\left(\frac{\hat{\pi}_{\text{young}}}{1 - \hat{\pi}_{\text{young}}} \bigg/ \frac{\hat{\pi}_{\text{old}}}{1 - \hat{\pi}_{\text{old}}}\right) \leftarrow \text{log odds ratio}$$

x continuous

$$W = 2 \left(\ell(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2) - \ell(\hat{\beta}_0, \hat{\beta}_1) \right)$$

$$\left| D(\hat{\beta}) = -2 \ell(\hat{\beta}) \right.$$

$$= D(\hat{\beta}_0, \hat{\beta}_1) - D(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2) \sim \chi^2_1$$

$$\hat{\beta}_1 = \hat{\beta}_0 + \hat{\beta}_1(x+1) - \hat{\beta}_0 - \hat{\beta}_1 x$$

expected increase in
the log odds

$$\hat{\beta}_1 = \log \frac{\hat{\pi}_{x+1}}{1 - \hat{\pi}_{x+1}} - \log \frac{\hat{\pi}_x}{1 - \hat{\pi}_x} = \log \left(\frac{\frac{\hat{\pi}_{x+1}}{1 - \hat{\pi}_{x+1}}}{\frac{\hat{\pi}_x}{1 - \hat{\pi}_x}} \right)$$

$$e^{\hat{\beta}_1} = \frac{\hat{\pi}_{x+1}}{1 - \hat{\pi}_{x+1}} - \frac{\hat{\pi}_x}{1 - \hat{\pi}_x}$$

$e^{\hat{\beta}_1}$ = expected
increase in the odds

In this lecture

$n=30$ (y_i, x_i) $i=1, \dots, n$ from $y = f(x) + \epsilon$
 yesterday data

estimate $f(x)$
 $\hat{f}(x)$

GOAL: predict the value of y when new x become available:

increasing p degree of the polynomial

smaller $D(\hat{f}(x)) = \|y - \hat{y}\|^2$
 larger $R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$

par $D(\hat{f}(x)) = 0$
 $R^2 = 1$

Our goal is not to perfectly fit the data that we have, but to predict the outcome/response for new observations

→ a polynomial of really high degree fit really well the data used to construct it, but it is useless to predict new observation

→ evaluation is done on new data (prediction)

→ "residual deviance" (computed on the new data) first decrease (the model is getting better), until a certain point, then increase (the model is getting worse) as a function of p

→ " R^2 " first increase, then decrease, same about the model

→ why we use p on the x-axis → measure (because we are using a polynomial) of the model complexity

$$\hat{y}_0 = \hat{f}(x_0)$$

y_0 = new response
 x_0 = new input/covariate

$$E \left[\left(\hat{f}(x_0) - f(x_0) \right)^2 \right]$$

$\hat{f}(x_0) = \hat{y}_0$ predicted value
 $f(x_0)$ true value

our aim is to minimize the expected (squared) error

$$\hookrightarrow E \left[\left| \hat{f}(x_0) - f(x_0) \right| \right]$$

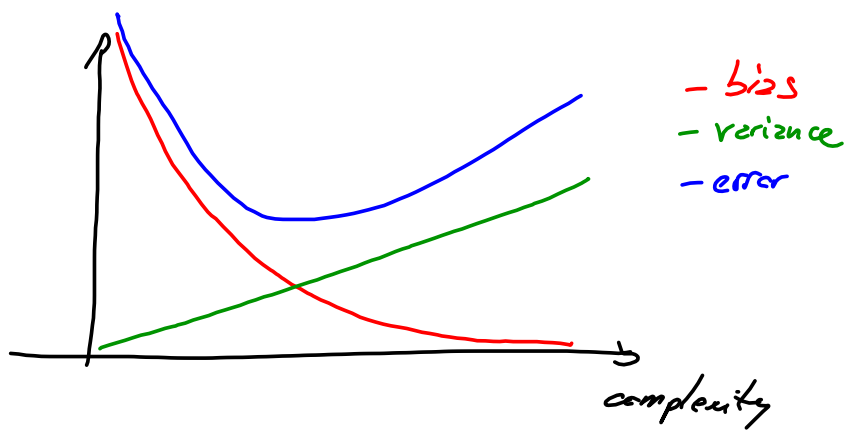
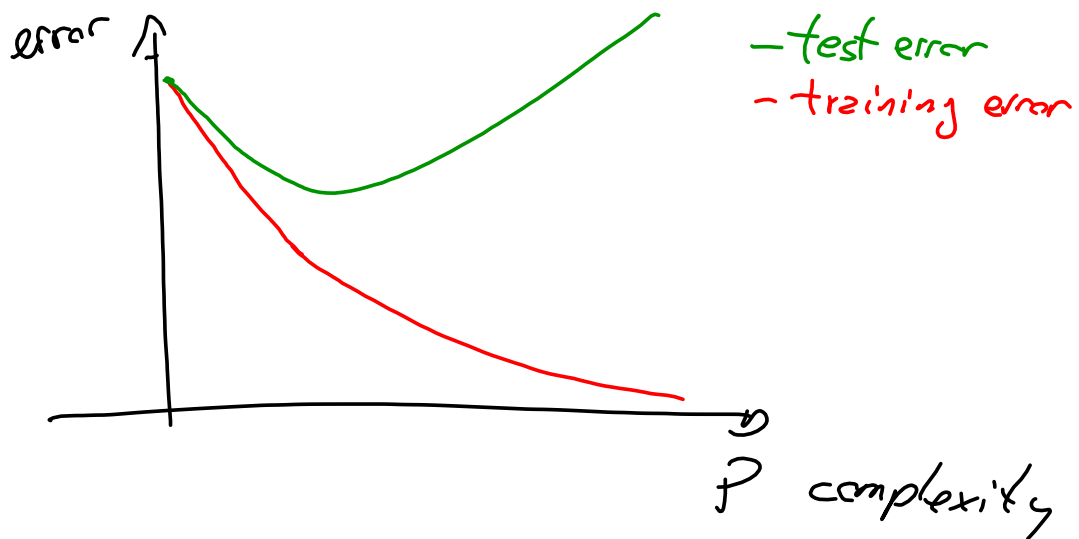
$$E \left[\left(\hat{f}(x_0) - f(x_0) \right)^2 \right] = \underbrace{\left(E \left[\hat{f}(x_0) \right] - f(x_0) \right)^2}_{\text{bias}^2} + \underbrace{\text{Var} \left[\hat{f}(x_0) \right]}_{\text{variance}}$$

→ do exercises 3.1 and 3.2 to derive this decomposition

$$E \left[\left(\hat{f}(x_0) - y_0 \right)^2 \right] = \underbrace{\left(E \left[\hat{f}(x_0) \right] - f(x_0) \right)^2}_{\text{bias}^2} + \underbrace{\text{Var} \left[\hat{f}(x_0) \right]}_{\text{variance}} + \underbrace{\epsilon}_{\text{irreducible error}}$$

$$y_0 = \underline{\underline{f(x_0) + \epsilon}}$$

mean squared error



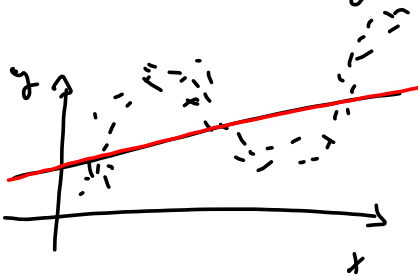
A model:

low complexity \Rightarrow high bias
low variance

the model does not capture well the systematic pattern of the data, but it is insensitive to small changes

high complexity \Rightarrow low bias
high variance

the model fits well the data, but small changes will lead to big variations



complexity low
high bias
low variance



complexity high
low bias
high variance

OVERFITTING

$$f(x) + \frac{\epsilon}{n \cdot x}$$

the estimated function $\hat{f}(x)$ follows random fluctuation of the data \Rightarrow dramatic increase in variance, negligible gain in bias

BIAS - VARIANCE TRADE-OFF

- bias and variance are "contrasting" (one cannot decrease both simultaneously, decreasing one means an increasing other)
- find the best balance between bias and variance to minimize the error

- OLS \rightarrow estimator with the smallest variance among the unbiased
- statistical learning \rightarrow we accept to increase the bias to have a decrease in the variance to reduce the expected prediction error

e.g.: $\hat{\beta}_{ridge} \quad \hat{\beta} (X^T X + \lambda I)^{-1} X^T y$

\rightarrow increase the bias, reduce the variance

(over)optimism: when we evaluate the error in the data that we used to construct the model, we are underestimating the error

