

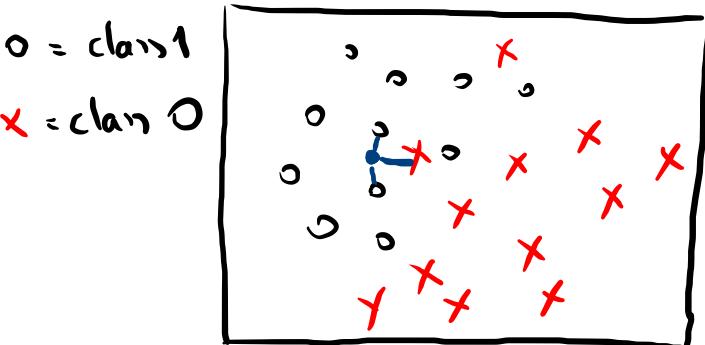
REGRESSION:

$$\textcircled{1} \text{ (parametric)} \quad E[\hat{Y}|X] = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p$$

$$\textcircled{2} \text{ (non-parametric) K-nn} \quad E[\hat{Y}|X] = \text{Ave}(y_i \mid x_i \in N_k(x))$$

CLASSIFICATIONK-NEAREST NEIGHBOUR (K-nn)

K positive integer

 $(x_i, y_i), i=1, \dots, n$ training data x_0 Estimate $P(Y=j \mid X=x_0)$, $j=1, \dots, C$ number of classes according to some metric

$$K=3 \quad \hat{P}(Y=1 \mid X=x_0) = \frac{2}{3}$$

$$P(Y=0 \mid X=x_0) = \frac{1}{3}$$

① Identify the K training points closest to x_0 ② Denote the indices of these points with N_0

$$\hat{P}(Y=j \mid X=x_0) = \frac{1}{K} \sum_{i \in N_0} I(y_i=j), j=1, \dots, C$$

CHOICE OF K

① How complex is this model?

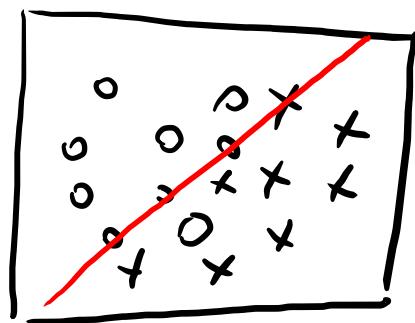
"Effective" number of parameters $\approx \frac{n}{K}$ ($\frac{n}{K}$ non-overlapping regions)

K large \rightarrow simple model

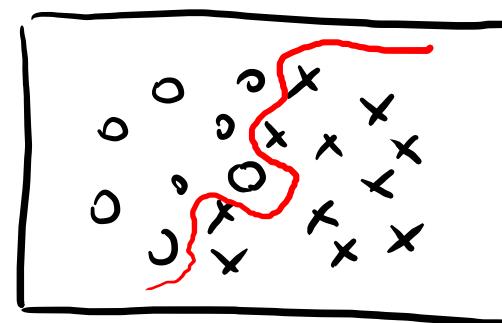
K small \rightarrow complex model

② Complexity of decision surface

$$P(Y=1 | X=x) = 0.5$$

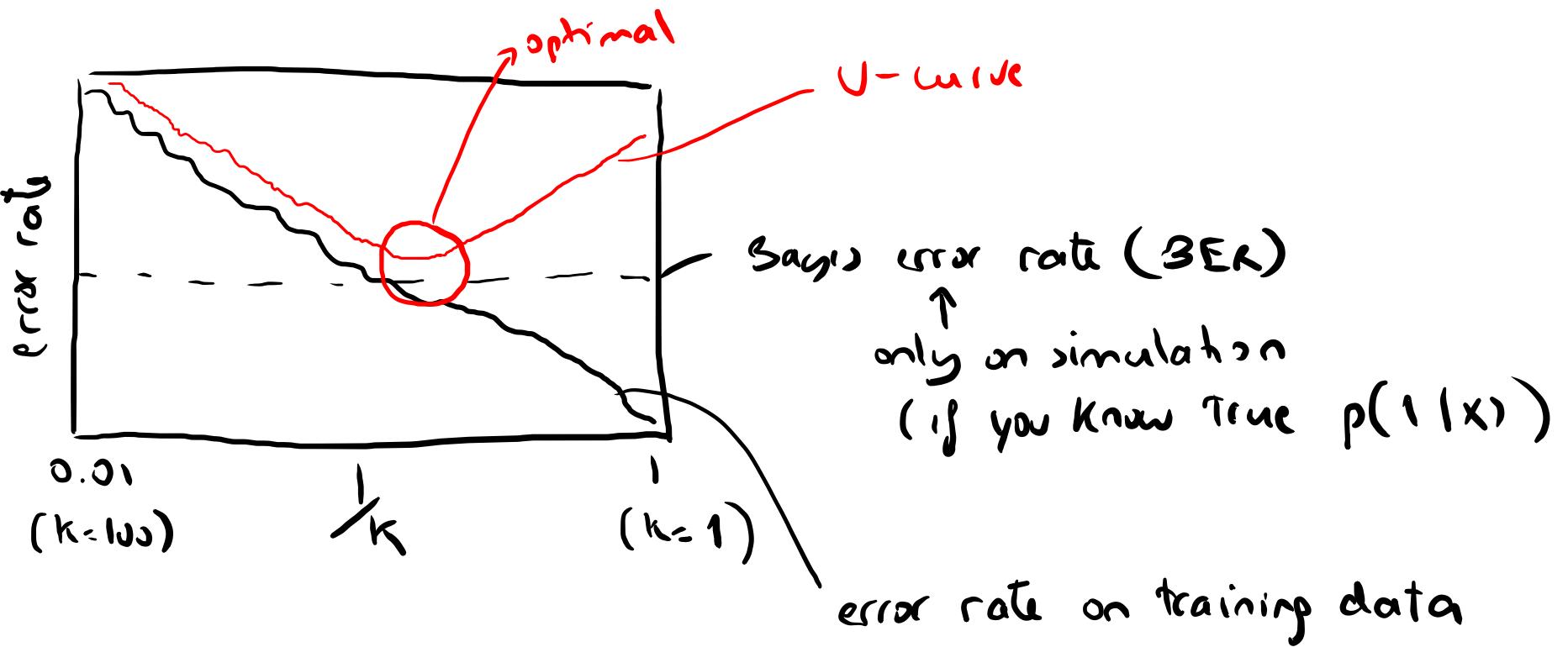


$K=100$



$K=1$

(choose K on some test data (unseen, not used for estimating $\hat{P}(Y=j|X)$)



error rate on training data

LOGISTIC REGRESSION

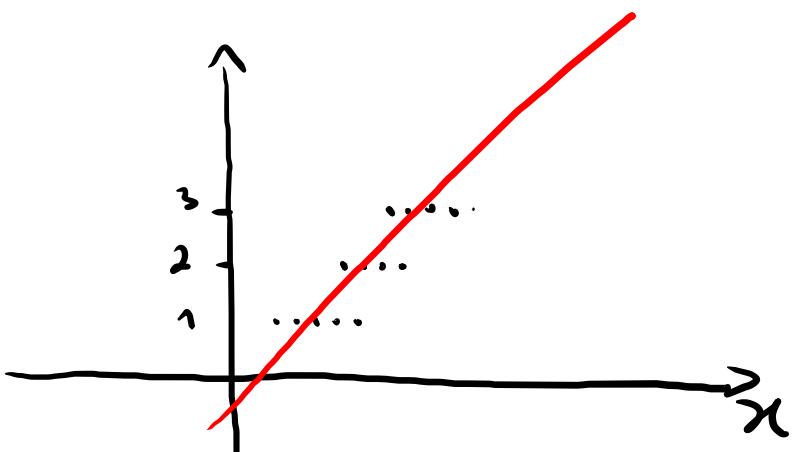
Why not linear regression?

Ex 1

Emergency room

$$y = \begin{cases} 1 & \text{if stroke} \\ 2 & \text{drug overdose} \\ 3 & \text{epileptic seizure} \end{cases}$$

Predict from symptoms



It could even work but

- ① No ordering between the categories, but we impose one
- ② We are imposing equal differences between categories

In most situations, we have no knowledge about ① and ②

Ex 2 Only two classes

Predict only between stroke and drug overdose

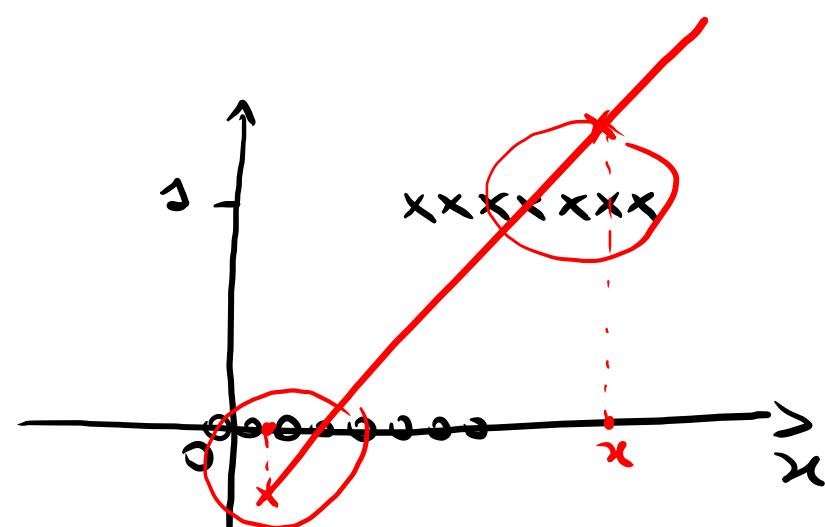
$$y = \begin{cases} 0 & \text{if stroke} \\ 1 & \text{if drug overdose} \end{cases} \quad (\text{dummy variable})$$

Fit linear regression, get \hat{y} and predict to 1 if $\hat{y} > 0.5$
 $E[\hat{y}|x]$

$$E[Y|x] = 0 \cdot P(0|x) + 1 \cdot P(1|x) = P(1|x)$$

As a linear regression model:

$$P(1|x) = \underbrace{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}_{[0, 1]} \quad (-\infty, \infty)$$



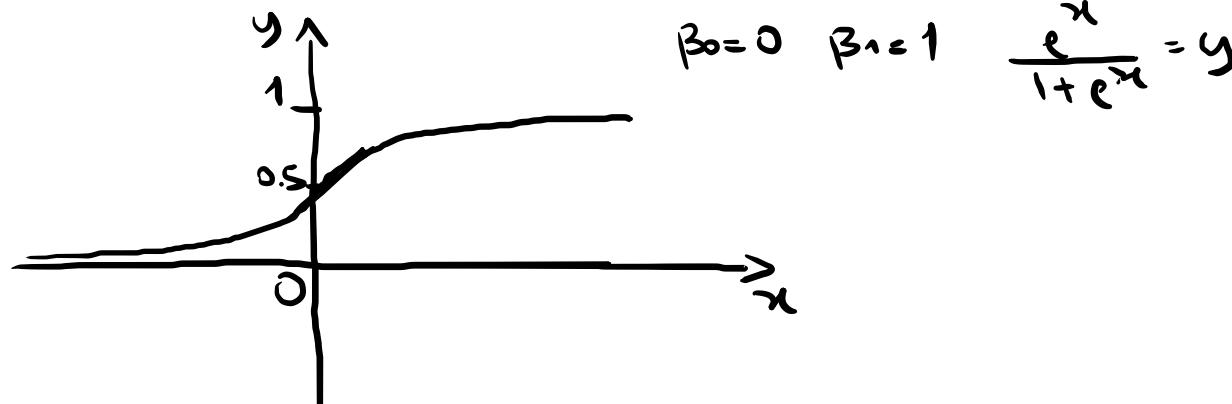
LOGISTIC
FUNCTION
(sigmoid
function)

$$P(1|\underline{x}) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}} = \frac{e^{\beta^t \underline{x}}}{1 + e^{\beta^t \underline{x}}} \quad (\underline{x}_0 = 1)$$

$$P(1|\underline{x}) \in (0, 1)$$

$$\underline{p} = 1$$

$$P(1|\underline{x}) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$



$$\beta_0 = 0 \quad \beta_1 = 1 \quad \frac{e^x}{1 + e^x} = y$$

ODDS

$$\frac{P(1|\underline{x})}{1 - P(1|\underline{x})} = \frac{e^{\beta^t \underline{x}} / (1 + e^{\beta^t \underline{x}})}{1 - \frac{e^{\beta^t \underline{x}}}{1 + e^{\beta^t \underline{x}}}} =$$

$$= \frac{\frac{e^{\beta^t \underline{x}}}{1 + e^{\beta^t \underline{x}}}}{\frac{1}{1 + e^{\beta^t \underline{x}}}} = e^{\beta^t \underline{x}}$$

LOG-
ODDS
(or LOGIT)

$$\log\left(\frac{P(1|\underline{x})}{1 - P(1|\underline{x})}\right) = \beta^t \underline{x} = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

LOGISTIC
REGRESSION
MODEL

INTERPRETATION OF PARAMETERS

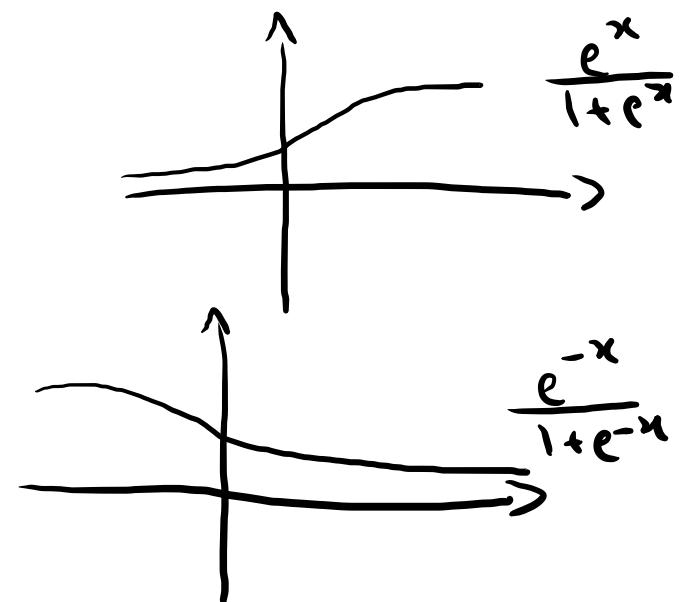
β_j As x_j increases by one unit (while everything else remain constant), the log-odds increase/decrease by β_j units

LOG
ODDS
RATIO

$$\log \left(\frac{p(1|x_0+1)}{1-p(1|x_0+1)} \right) - \log \left(\frac{p(1|x_0)}{1-p(1|x_0)} \right) = \\ = \beta_0 + \beta_1(x_0+1) - \beta_0 - \beta_1 x_0 = \beta_1$$

$\beta_j > 0 \rightarrow p(1|x)$ increases if x_j increases

$\beta_j < 0 \rightarrow p(1|x)$ decreases as x_j increases



PARAMETER ESTIMATION

Typically by maximum likelihood

Training data (\underline{x}_i, y_i) , $i = 1, \dots, n$

y binary response with a Bernoulli distribution:

$$y | x = \begin{cases} 1 & p(1|x) = p(x) \\ 0 & p(0|x) = 1 - p(x) \end{cases}$$

LIKELIHOOD

$$\begin{aligned} L(\beta) &= \prod_{i=1}^n f(\underline{x}_i, y_i; \beta) \\ &= \prod_{i=1}^n p(\underline{x}_i)^{y_i} (1 - p(\underline{x}_i))^{1-y_i} \end{aligned}$$

$$\begin{aligned} y_i = 1 &\rightarrow p(\underline{x}_i) = p(1|\underline{x}_i) \\ y_i = 0 &\rightarrow 1 - p(\underline{x}_i) = p(0|\underline{x}_i) \end{aligned}$$

LOG-LIKELIHOOD

$$L(\beta) = \prod_{i=1}^n p(x_i)^{y_i} (1-p(x_i))^{1-y_i}$$

$$\ell(\beta) = \log(L(\beta)) = \sum_{i=1}^n [y_i \log(p(x_i)) + (1-y_i) \log(1-p(x_i))]$$

$$= \sum_{i=1}^n [y_i \log p(x_i) + \log(1-p(x_i)) - y_i \log(1-p(x_i))]$$

$$= \sum_{i=1}^n \left[y_i \underbrace{\log \frac{p(x_i)}{1-p(x_i)}}_{\text{p}(x_i|x_i)} + \log(1-p(x_i)) \right]$$

$$= \sum_{i=1}^n \left[y_i \underbrace{\beta^t x_i}_{\beta_0 x_{i0} + \beta_1 x_{i1} + \dots + \beta_j x_{ij} + \dots + \beta_p x_{ip}} - \log(1+e^{\beta^t x_i}) \right]$$

FUNCTION
OF $\beta_0, \beta_1, \dots, \beta_p$

$$\begin{aligned} & \log(1 - \frac{e^{\beta^t x_i}}{1+e^{\beta^t x_i}}) \\ &= \log\left(\frac{1}{1+e^{\beta^t x_i}}\right) = -\log(1+e^{\beta^t x_i}) \end{aligned}$$

MAXIMUM LOG-LIKELIHOOD

$$\frac{\partial \ell(\beta)}{\partial \beta_j} = \sum_{i=1}^n \left[y_i x_{ij} - \frac{e^{\beta^t x_i}}{1+e^{\beta^t x_i}} x_{ij} \right] = \sum_{i=1}^n [y_i x_{ij} - p(x_i) x_{ij}] \quad j=0, 1, \dots, p$$

$\underset{p+1 \text{ equations in}}{\underset{p+1 \text{ parameters}}{}}$

Numerical methods such as Newton-Raphson are used to solve these equations, returning $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$

Using these estimates, we calculate probabilities for any x value:

$$\hat{P}(1|x) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p}}$$

> 0.5 classify x to class 1
 < 0.5 classify x to class 0

Or plot a ROC curve on a test data, etc etc

A SIMPLE EXAMPLE where $\hat{P}(1|x)$ can be obtained by hand:

y binary response

x binary

So data can be summarised in a table:

x	0	y	0	1
			a	b
	1		c	d

with $a, b, c, d \neq 0$

	Y
0	a
1	b

x	0	1
0	c	d
1		

$$\log \left(\frac{P(Y=1 | X=1)}{1 - P(Y=1 | X=1)} \right) = \beta_0 + \beta_1 \cdot 1$$

$$\log \left(\frac{P(Y=1 | X=0)}{1 - P(Y=1 | X=0)} \right) = \beta_0$$

So parameters can be easily estimated by :

$$\hat{\beta}_0 = \log \left(\frac{\hat{P}(Y=1 | X=0)}{\hat{P}(Y=0 | X=0)} \right) = \log \left(\frac{\frac{b}{a+b}}{\frac{a}{a+b}} \right) = \log \left(\frac{b}{a} \right)$$

$$\begin{aligned} \hat{\beta}_1 &= \log \left(\frac{\hat{P}(Y=1 | X=1)}{1 - \hat{P}(Y=1 | X=1)} \right) - \hat{\beta}_0 = \log \left(\frac{d}{c} \right) - \log \left(\frac{b}{a} \right) = \\ &= \log \left(\frac{a \cdot d}{b \cdot c} \right) \end{aligned}$$

Final comments

- ① More than 2 classes?

Could fit a model for $p(1|x)$ (versus all the other classes) and another for $p(2|x)$.

$$\text{Then } p(3|x) = 1 - p(1|x) - p(2|x)$$

But logistic regression is not as suited to the case of more than 2 classes as other methods (LDA).

- ② Predictors can take any form: continuous, categorical, ordinal, etc...

- ③ Although the model is linear (in the parameters), it can accommodate rather complex structures, such as "polynomial" terms (x_1^2, x_1^3, \dots) and interaction effects ($x_j \cdot x_k, x_j^2 \cdot x_k, \dots$)