

## K-Nearest Neighbors

$$X = \begin{pmatrix} \text{wt} & \text{ht} & \text{age} \\ x_{11} & x_{12} & x_{1p} \\ \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & x_{np} \end{pmatrix} \xrightarrow{\text{standardize}} \begin{pmatrix} \frac{x_{11} - \bar{x}_1}{\text{sd}(x_1)} & \frac{x_{12} - \bar{x}_2}{\text{sd}(x_2)} & \frac{x_{1p} - \bar{x}_p}{\text{sd}(x_p)} \\ \vdots & \vdots & \vdots \end{pmatrix}$$

1. Because not parametric

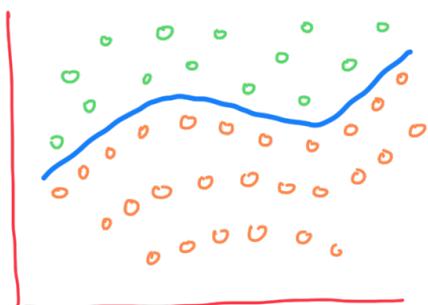
$$2. \text{ Distance } (x_1, x_2) = \sqrt{(x'_{11} - x'_{21})^2 + \dots + (x'_{1p} - x'_{2p})^2}$$

Model complexity  $\propto 1/k$

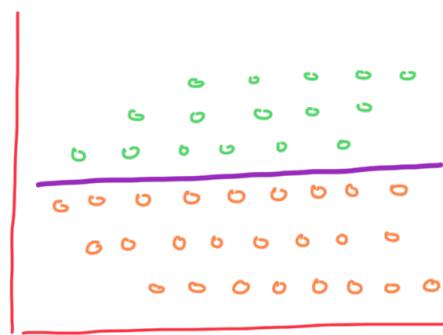
$k=1$	over-fitting, wiggly, flexible, variance $\uparrow$
$k$ large	inflexible bias $\uparrow$

Disadvantage:

1. CURSE OF DIMENSIONALITY: Fail at larger dimensions
2. nearest neighbors tend to be spatially large, bias  $\uparrow$
3. reducing neighbor by  $k \downarrow$  variance  $\uparrow$



Low value of  $k$  is good



Large value of  $k$  is good

Q.8. Bias variation derivation for KNN

# STATS 503

	Generative	Discriminative
	$P(Y=c_k x) \propto P(x Y=c_k)$	$P(Y=c_k x)$
Parametric	LDA    QDA NB	Logistic Regression NN, SVM, DT, RF
Non-parametric	—	kNN

## Parametric Methods

- ① Make some assumption of functional form  
i.e. come up with some model  
eg. Gaussian distribution for  $p(x|Y=c_k)$  OR  
Linear model for logit of  $p(Y=c_k|x)$

② Use training data to fit model or unknown parameters

**Non-Parametric Methods** : No explicit assumptions  
Advantages Disadvantages

1. Flexible
  2. Wider range of shapes

1. Overfitting
  2. Large # observations required

## Baye's Theorem

distribution of  $x$   
given class label

prior probability

$$P(Y=c_k | X) = \frac{p(x|Y=c_k) \cdot P(Y=c_k)}{p(x)}$$

$$\arg\max_k p(y=c_k | x) = \arg\max_k \left\{ \frac{p(x|y=c_k) \cdot p(y=c_k)}{p(x) \rightarrow \sum_{k'=1}^K \pi_{k'} p(x|y=c_{k'})} \right\}$$

$$\underset{\text{independent of } k}{\arg \max_k} \left\{ p(x|y=c_k) \cdot p(y=c_k) \right\} = \arg \max_k \pi_k p_k(x)$$

Generative Methods of classification

$$\text{argmax}_k (Y = C_k | X = x) \propto \pi_k P(x|Y = C_k) = \pi_k p_k(x)$$

specify a specific direction for  $p_k(x)$  (with parameters)

use training data

$$(x_1, Y_1), (x_2, Y_2), \dots, (x_n, Y_n)$$

$$\Rightarrow \hat{\pi}_k \hat{p}_k(x) \Rightarrow \text{argmax}_k [\hat{\pi}_k \cdot \hat{p}_k(x)] \Rightarrow \hat{C}(x) \text{ fitted classifier}$$

Assumptions

1) QDA  $p_k(x) \sim MVN(\mu_k, \Sigma_k)$

2) LDA  $p_k(x) \sim MVN(\mu_k, \Sigma_k)$   $\Sigma_k = \Sigma$  same covariance

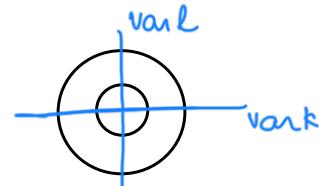
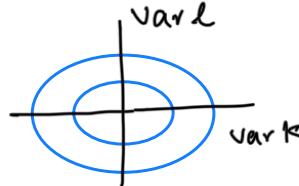
3) NB  $p_k(x) \sim MVN(\mu_k, \Sigma_k)$   $\sigma_{ij} = 0; i \neq j$

$\begin{bmatrix} \sigma_{11} & & \phi \\ & \sigma_{22} & \\ \phi & \ddots & \sigma_{kk} \end{bmatrix}$   $\Sigma_k$ : diagonal matrix no correlation b/w variables

$\Sigma_i$  denotes the dispersion of the  $i^{\text{th}}$  class

$$\Sigma_i = \begin{pmatrix} \text{Var}(X_1) & \dots & \text{Cov}(X_p, X_1) \\ \vdots & \ddots & \vdots \\ \text{Var}(X_2) & & \vdots \\ \vdots & & \vdots \\ \text{Cov}(X_1, X_p) & \dots & \text{Var}(X_p) \end{pmatrix}_{p \times p}$$

LDA: same dispersion for all the  $k$  classes



NB: when correlation b/w variables = 0  
⇒ same dispersion for all classes

Bias: QDA < LDA < NB  
flexible rigid

Variance: QDA > LDA > NB

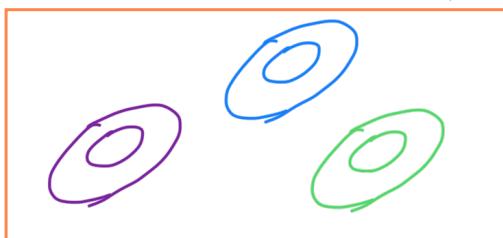
parameters:  $O(kp^2) > O(p^3) > O(kp)$   $\because p > k$

Use Cases

QDA: variances are different among classes and  $n > p$   
enough observations

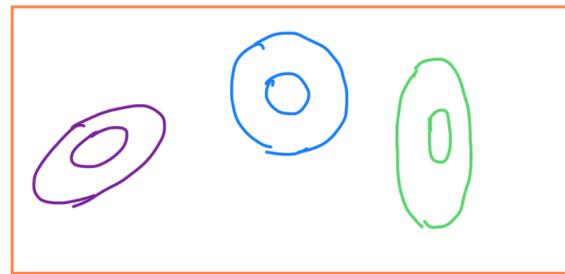
LDA: variances are similar or  $n > p$   
NOT enough observations to get variance

NB: when  $p$  (no of dimensions) is large  $> n$  (sample size)  
e.g. count of words in word vectors / language models



$p = \text{dimension} = 2$   
 $k = \text{classes} = 3$

LDA distribution better



QDA distribution better

L.5. How to plot LDA, QDA, NB curves?

## Logistic Regression

$$\log \frac{p(x)}{1-p(x)} = x^\top \beta \quad x = (1, x_1, \dots, x_p) \quad \beta = (\beta_0, \beta_1, \dots, \beta_p)$$

$p(x) \in [0, 1] \quad x^\top \beta \in [-\infty, \infty]$

$$P(Y=1|x) \left[ p(x) = \frac{e^{x^\top \beta}}{1+e^{x^\top \beta}} \right]$$

$$P(Y=0|x) \left[ 1-p(x) = \frac{1}{1+e^{x^\top \beta}} \right]$$

Classification Boundary

$$x^\top \beta = 0 \quad \text{if} \quad x^\top \beta > 0 \Rightarrow c_1$$

$$x^\top \beta < 0 \Rightarrow c_2$$

Q.1. what is outcome of L.R?

**DEFINITION** Likelihood:

$$Q.2. \text{ when do we use } Y_i=1 \Rightarrow \frac{e^{x_i^\top \beta}}{1+e^{x_i^\top \beta}}$$

Training data  $(x_1, y_1), \dots, (x_n, y_n)$

$$Y_i=0 \Rightarrow \frac{1}{1+e^{x_i^\top \beta}}$$

$$\operatorname{argmax}_{\beta} \prod_{i=1}^n p(x_i, y_i) = \operatorname{argmax}_{\beta} \prod_{i=1}^n [p(x_i) p(y_i|x_i)]$$

$$= \operatorname{argmax}_{\beta} \prod [p(y_i|x_i)] = \operatorname{argmax}_{\beta} \ln [\prod p(y_i|x_i)]$$

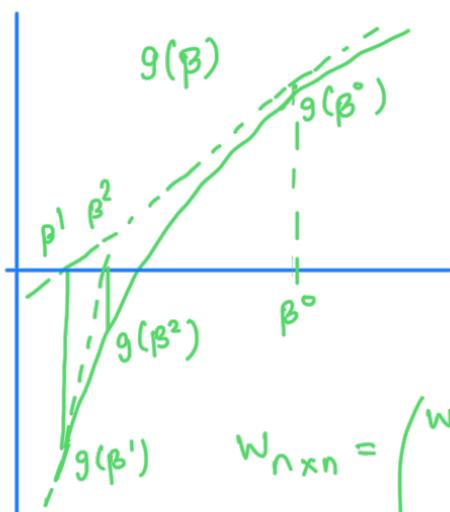
independent of  $\beta$  one-one monotonic function

$$= \operatorname{argmax}_{\beta} \sum_{i=1}^n \ln p(y_i|x_i) = \operatorname{argmax}_{\beta} l(\beta) \leftarrow$$

some function of  $\beta$

How to maximize  $\beta_0$ : Iterative Newton Raphson method

$$\frac{\partial l(\beta)}{\partial \beta} = 0 \Rightarrow \text{For } p+1 \text{ dimension vector} \quad \frac{\partial l(\beta)}{\partial \beta_0} = 0 \quad \frac{\partial l(\beta)}{\partial \beta_1} = 0 \quad \dots \quad \frac{\partial l(\beta)}{\partial \beta_p} = 0$$



$l(\beta)$  is a scalar function

$$\beta_{\text{new}} = \beta_{\text{old}} - g(\beta_{\text{old}})/g'(\beta_{\text{old}})$$

$$g(\beta) = \frac{\partial l(\beta)}{\partial \beta} = 0$$

$$\beta_{\text{new}} = (X^\top W X)^{-1} X^\top W z \quad \text{repeat until convergence}$$

$$W_{n \times n} = \begin{pmatrix} w_{11} & & \Phi \\ & w_{22} & \ddots \\ \Phi & & w_{nn} \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$$

$$W = \operatorname{diag} [p(x_i; \beta_{\text{old}}) \{1 - p(x_i; \beta_{\text{old}})\}]$$

$$z = X\beta_{\text{old}} + W^{-1}(y - p)$$

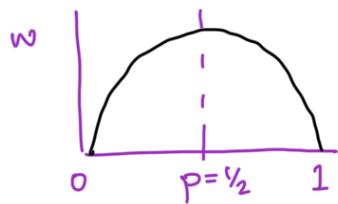
Comparison with Ordinary least squares  
 $\min_{\beta} \sum (z_i - \beta^T x_i)^2 \Rightarrow \hat{\beta} = (X^T X)^{-1} X^T z$

Comparison with weighted least squares  
 $\min_{\beta} \sum w_i (z_i - \beta^T x_i)^2 \Rightarrow \hat{\beta} = (X^T W X)^{-1} X^T W z$

Q.3. What is value of  $w$

### Interpretation of logistic Regression

$$w_i = p(x_i) (1-p(x_i))$$



$p \approx 0$  or  $p \approx 1$   $w$ : small

$p \approx 1/2$   $w$ : large

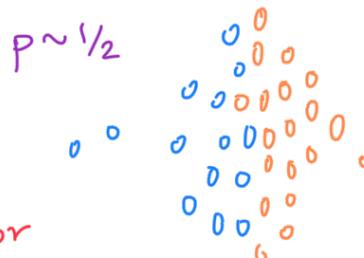
easy for classification

not so easy classification

### Regression Error

$$\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i - \dots - \hat{\beta}_d x_d)^2$$



### Classification Error

$$\text{np.mean } (\hat{y} \neq y_{\text{true}})$$

Q.4. accuracy =  $\uparrow$

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}(\hat{c}(x_i) \neq y_i)$$



$$\hat{c}(x) = \begin{cases} 1 & \hat{p}(x) > 0.7 \\ 0 & \hat{p}(x) \leq 0.7 \end{cases}$$

let  $C = 0.7$   
 $f_{\text{pr}}, t_{\text{pr}}, \text{thresholds} \leftarrow$  from ROC-curve

Plot ROC and get AUC  
 upper left points and larger AUC are better results

### LDA

- Linearity by construction
- More general
- Robust to outliers
- weight depends on distance from decision boundary

- Linearity by consequence
- assumes Gaussian density
- easier to compute, when  $\mathcal{N}(\mu, \Sigma)$
- Uses all data points

## Multinomial logistic regression

$$P(Y = c_1 | x) = \frac{e^{x^\top \beta_1}}{1 + e^{x^\top \beta_1} + e^{x^\top \beta_2} + \dots + e^{x^\top \beta_{k-1}}}$$

$$P(Y = c_2 | x) = \frac{e^{x^\top \beta_2}}{1 + e^{x^\top \beta_1} + e^{x^\top \beta_2} + \dots + e^{x^\top \beta_{k-1}}}$$

$$P(Y = c_k | x) = \frac{1}{1 + e^{x^\top \beta_1} + e^{x^\top \beta_2} + \dots + e^{x^\top \beta_{k-1}}}$$

Use training data  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$$\Rightarrow \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_{k-1}$$

## Bias-Variance Trade off

Model complexity ↑  
kNN as  $k \downarrow$

QDA > LDA > NB  
logistic Regression > LDA  
smooth spline  $\lambda \downarrow$

LOOCV > k-fold CV as most training data is used

computationally intense, fit each model  $n$  times ( $\times$ )  
5/10 fold CV do not suffer high Bias & high Variance

$$V = E[(z_1 - \bar{z})(z_2 - \bar{z}) \dots (z_k - \bar{z})] \quad \text{Var}(\bar{z}) = \text{Var}\left(\frac{1}{k} \sum z_i\right) = \frac{1}{k^2} \text{Var}\left(\sum z_i\right) = \frac{1}{k^2} \left( \sum \text{Var}(z_i) + \sum \text{Cov}(z_j, z_k) \right)$$

KC-VE.  $\bar{z} = \frac{1}{K} \sum z_i$   $z_j$  and  $z_k$  are highly correlated for LOOCV

$\therefore \text{Var}(\bar{z})$  LOOCV is high

$$(y - \hat{y})^2 = \underbrace{(y - E[y])^2}_{a} + \underbrace{E(y) - E(\hat{y})}_{b}^2 + \underbrace{E(\hat{y}) - \hat{y}}_{c}^2 + \underbrace{\text{Var}(\hat{y})}_{o}$$

$$\begin{aligned} E[ab] &= (E[y] - E[\hat{y}]) E[y - E[y]] & E[y] &= \text{constant} \\ &= (E[y] - E[\hat{y}]) (E[y] - E[y]) = 0 & E[c] &= c \end{aligned}$$

$$E[(y_o - \hat{y}_o)^2] = E[(y_o - E[y_o])^2] + E[(E[y_o] - E[\hat{y}_o])^2] + E[(E[\hat{y}_o] - \hat{y}_o)^2]$$

$$\begin{aligned} y &= f(x) + \varepsilon \\ E[y_o] &= E[f(x_o)] + E[\varepsilon] \\ f(x_o) &= \text{constant wrt } E \\ E[\varepsilon] &= 0 ; \text{ assumption} \\ \therefore E[y_o] &= f(x_o) \\ y_o - E[y_o] &= y_o - f(x_o) = \varepsilon \end{aligned}$$

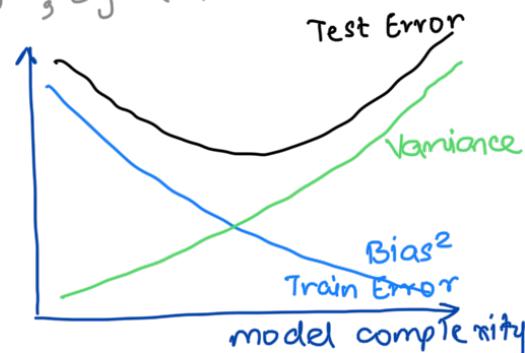
$$E[\varepsilon^2] = \text{Var}(\varepsilon) = \sigma^2 ; \text{ assumption}$$

$$\text{Expected Test error } T.E(x_o) = E[(y_o - \hat{y}_o)^2]$$

$$\therefore T.E(x_o) = \sigma^2 + \text{bias}(\hat{y}_o) + \text{var}(\hat{y}_o)$$

irreducible error

$$\begin{aligned} E[(E[y_o] - E[\hat{y}_o])^2] &= (f(x_o) - E[\hat{y}_o])^2 ; \text{ expectation vanishes} \\ &= \text{Bias}(\hat{y}_o) ; \text{ by definition} \end{aligned}$$



model complexity ↑ bias<sup>2</sup> ↓ variance ↑

Assume: Data points  $(x_1, y_1) \dots (x_n, y_n)$  are iids (R.V.)  
from a joint distribution  $P(x, y)$

Given:

$x$ : Feature vector

$y$ : class Label  
Binary  
Multiclass

Predictors  $\rightarrow$  continuous / discrete

To find: a model / classifier  $\gamma$  as function of  $x$

Goal: minimize error for prediction on unseen "test" data

0-1 Loss:

$$L(y, \hat{y}) = \begin{cases} 1 & \text{if } \hat{y} \neq y \\ 0 & \text{if } \hat{y} = y \end{cases}$$

Symmetric loss, all errors are equal  
sometimes not the best

CROSS-VALIDATION:

$$1 \left[ \frac{1}{20} \right] \text{mean}(y_{\text{true}} \neq y_{\text{hat}})$$

$$2 \left[ \frac{1}{20} \right] \text{mean}(y_{\text{true}} \neq y_{\text{hat}})$$

:

$$10 \left[ \frac{1}{20} \right] \text{mean}(y_{\text{true}} \neq y_{\text{hat}})$$

np.mean()

10-fold cross-validation

Data\_size = 200

Folders\_size = 20

No of observations  $\uparrow$

No of dimensions  $\downarrow$

relation complexity  $\uparrow$

Model complexity  $\uparrow$

Splines



$$\text{Cubic: } ax^3 + bx^2 + cx + d = 0 \quad 4(k+1) \text{ equations}$$

$$\text{continuity: } s_k(\xi_k) = s_{k+1}(\xi_k) \quad k$$

$$\text{First-order cty: } s'_k(\xi_k) = s'_{k+1}(\xi_k) \quad k$$

$$\text{Second order cty: } s''_k(\xi_k) = s''_{k+1}(\xi_k) \quad k$$

---

$k+4$  equations

$$\text{Natural cubic: } s''(\xi_0) = s''(\xi_{k+1}) = 0 \quad 2$$

$$s'''(\xi_0) = s'''(\xi_{k+1}) = 0 \quad 2$$

---

$k$  equations

L.6 verify from graph  
is this first & second order?

$$\text{check } TE = \text{var} + (\text{bias})^2$$

smoothing splines

To predict a function  $g(x)$  for the given data:

$$\min_{g(x)} \underbrace{\sum (y_i - g(x_i))^2}_{\text{RSS}} + \lambda \underbrace{\int g''(t)^2 dt}_{\text{penalty}}$$

1.  $g''(t)^2$ : roughness; second order gives curvature

2. squared to get positive outcome

3.  $\lambda \geq 0$ : tuning parameter

$\lambda = 0$   $\hat{g}(x)$  curve fits all points; over-fitting

$\lambda \rightarrow \infty$   $\hat{g}(x)$  becomes smoothest

$\hat{g}''(x) \therefore$  linear regression model

## Generalised Additive Models

### Advantages

$$y_i = \beta_0 + f_1(wt_i) + f_2(ht_i) + \dots + f_p(age_i) + \epsilon_i$$

### Disadvantages

suits for splines with more than one "x" predictors  
non-linear fit but still additive NOT  $f(wt, ht, age^2) \dots$   
Easy interpretation of individual predictors

$$\ln\left(\frac{p(r_i=c_1|x_i)}{1-p(r_i=c_1|x_i)}\right) = \beta_0 + \sum_{j=1}^n f_j(x_{ij})$$

