

# Bayes Classifiers

## Training time

Joint model  $p(y, \underline{x}) = p(y)p(\underline{x}|y)$

$$p(y=k) = \pi_k \approx \frac{\# \text{ } k \text{ labels}}{N}$$

$$p(\underline{x}|y=k) \dots \text{ eg } \mathcal{N}(\underline{x}; \mu^{(k)}, \Sigma^{(k)})$$

Not Bayesian

Assuming we know  
all parameters.

|| Mean & cov of  
 $\underline{x}$ 's in class  $k$

$$\text{Naive Bayes } p(\underline{x}|y=k) = \prod_d p(x_d|y=k)$$

Univariate Gaussian,  
Discrete, ...

## Test time

$$p(y|\underline{x}) \propto p(y, \underline{x}) \quad (\text{Bayes' Rule})$$

$$y_{\text{guess}} = \operatorname{argmax}_k \underbrace{p(y=k, \underline{x})}_{\text{"goodness"}}$$

quadratic decision boundary.

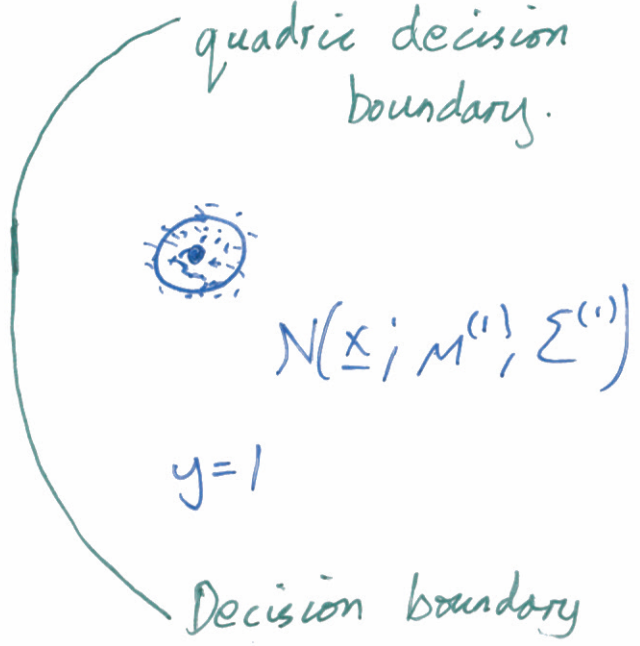


$N(\underline{x}; \mu^{(0)}, \Sigma^{(0)})$   
 $y=0$



$N(\underline{x}; \mu^{(1)}, \Sigma^{(1)})$

$y=1$



Decision boundary

$\pi_0 N(\underline{x}; \mu^{(0)}, \Sigma^{(0)}) =$

$\pi_1 N(\underline{x}; \mu^{(1)}, \Sigma^{(1)})$



$p(y|\underline{x}) = \frac{p(y, \underline{x})}{p(\underline{x})}$

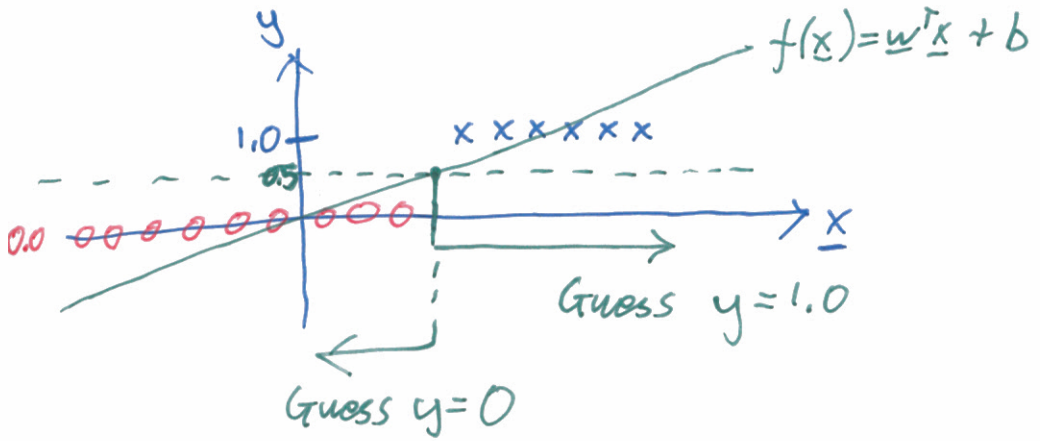
$p(\underline{x}) = \sum_{y_i} p(y_i, \underline{x})$

Discipuli Domini Colini Drummond qui vigesimo-septimo die  
 Februarii MDCCXIX subscripserunt 1719

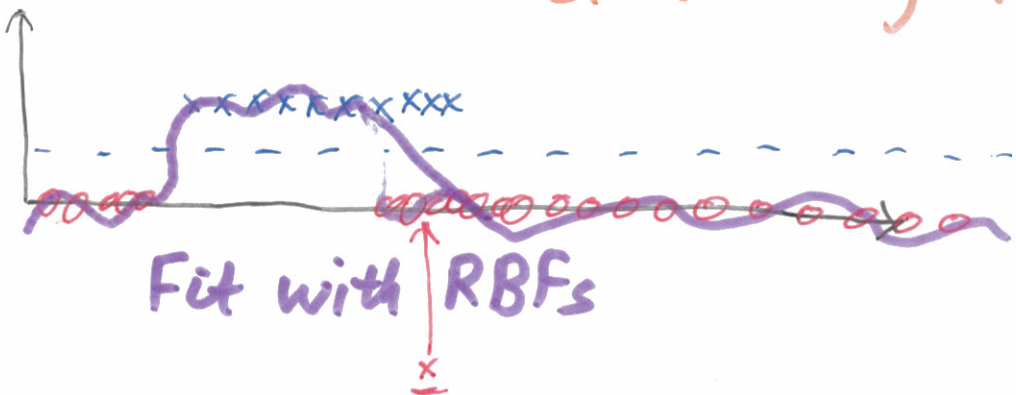
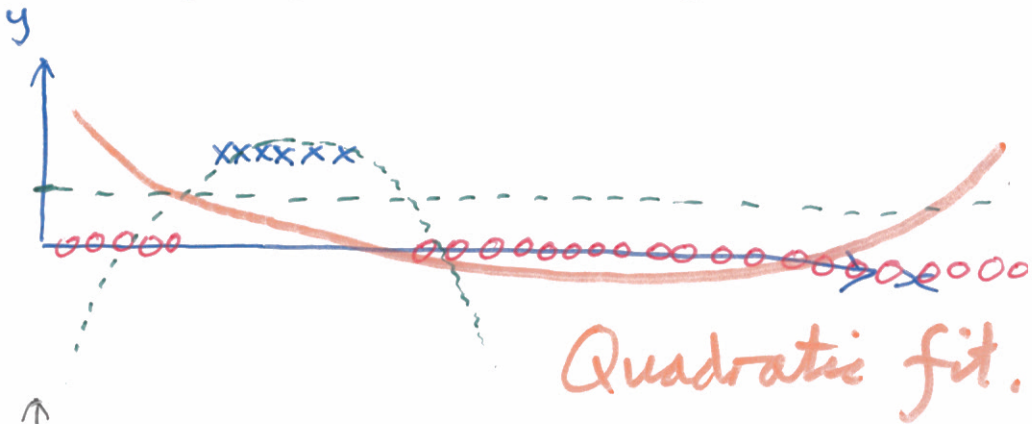
Arch Rennie 4	And: Beard 3	Alexr Crokat 3
Edwa: Aukine 2	Geo: Carruthers 4	David Lindsay 1
Geo: Gordon 2	Geo: Sewar	Geo: Douglass 2
Geo: Wilson 4	Geo: M <sup>r</sup> L <sup>r</sup> 2	Gul: Taylor 1
Gul: Horsburgh 3	Gul: Ramsay 3	Jo: Barclay 2
Hen: Ker 4	Gul: Adell 4	Jo: Nicolson
Joan: Boston 3	John Connell 2	Jo: Howley
Jo: Carruthers 4	J <sup>r</sup> Gulekrist m <sup>r</sup> 2	John: Patoun 1
Joan: Morison 3	Jo: Gilson	John: Rutherford 2
John: Paxton 2	Jo: Shaw 1	Jo: Smith 1
Jo: Full 2	Jo: Jaques 1	Jo: Thomson 2
Mich: Robertsone 2	Ruth: Bell	Isa: Maddox
Pat: Murdock 3	Ro: Dalrymple 1	Rob: Cleland 4
Simon Elliot 2	Row: Dunbar 4	Rob: Douglass 1
	Skinner: Smith	Rob: Risharby 1
Thomas Carmichael 2	Tho: Davidson 2	Th: Bayes
		Tho: Morison 2



# Regression for classification



If  $f(x) > 1/2$ , guess  $y=1$



If minimize square loss?

$$\min_{f(x)} E [(y - f(x))^2] \text{ at some } \underline{x}$$

Cost

$$= p_1 (1-f)^2 + \underbrace{(1-p_1)}_{p(y=0|x)} (0-f)^2$$

$$= p_1 (1 - 2f + f^2) + (1 - p_1) f^2$$

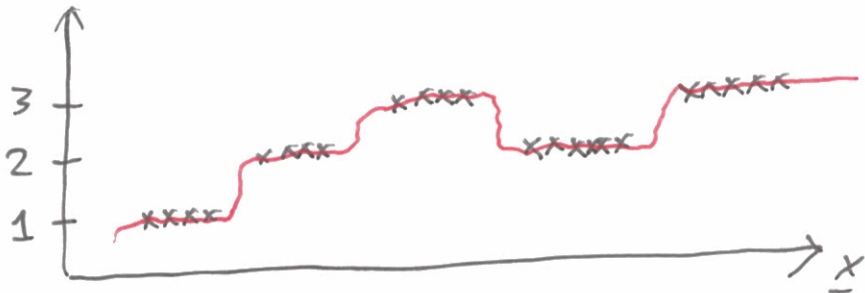
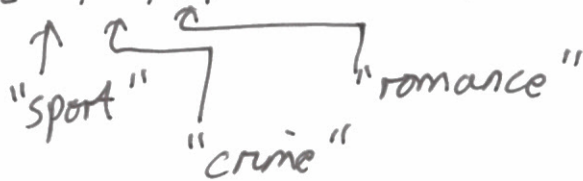
$$= f^2 (p_1 - p_1 + 1) - 2p_1 f + p_1$$

$$\frac{\partial \text{cost}}{\partial f} = 2f - 2p_1 = 0 \text{ at optimum}$$

$$\boxed{f = p_1}$$

# Multiple classes

$$y \in \{1, 2, 3, \dots, 103\}$$



$$f(\underline{x}) = \underline{w}^T \underline{x}$$

[ Replace  $\underline{x}$   
with  $\phi(\underline{x})$   
if you like ]

$$f(\underline{x}^{(1)}) \approx 1 \Rightarrow \text{"sport"}$$

$$f(\underline{x}^{(2)}) \approx 3 \Rightarrow \text{"romance"}$$

$$f\left(\frac{\underline{x}^{(1)} + \underline{x}^{(2)}}{2}\right) \approx 2 \Rightarrow \text{"crime"}$$

# One-hot encoding, One-of-k-encoding

M

Vector output

$$y^{(n)} = [0 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]^T$$

↙ k<sup>th</sup> position

k x 1

for k classes.

If n<sup>th</sup> example

is in class k

Fit k functions, one for each bit  $y_k$

Predict class where  $f_k(\underline{x})$  biggest

Pre-processing also useful for features  $\underline{x}$

$x_d \in \{ \text{"red", "green", "blue"} \}$

↓  $\{ 1, 2, 3 \}$

3 features

red	→	1 0 0	1 star	1 0 0 0
green	→	0 1 0	2 star	1 1 0 0
blue	→	0 0 1	3 star	1 1 1 0

Puzzle: R libraries don't create. Why?