

NAIVE BAYES CLASSIFIER

$$\underline{X} = (X_1, \dots, X_p)$$

$\underline{y}$  response (class)  $\xrightarrow{P(y=j)}$   $f_j(\underline{x})$

$$P(y=j | \underline{X} = \underline{x}) = \frac{\pi_j f_j(\underline{x})}{\sum_{i=1}^K \pi_i f_i(\underline{x})}$$

① LDA,  $f_j(\underline{x}) \sim N_p(\underline{\mu}_j, \Sigma)$

② QDA,  $f_j(\underline{x}) \sim N_p(\underline{\mu}_j, \Sigma_j), j=1, \dots, K$

③ NB

## NAÏVE BAYES CLASSIFIER

The main assumption is:

$$f_j(\underline{x}) = \prod_{k=1}^p f_{jk}(x_k)$$

$(x_1, x_2, \dots, x_p)$

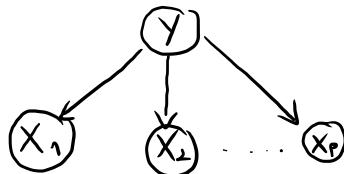
CONDITIONAL  
INDEPENDENCE

$x_1, \dots, x_p$   
DISCRETE

$$P(x_1=x_1, \dots, x_p=x_p | y=j) = f_j(x_1) \cdot f_{j2}(x_2) \cdot \dots \cdot f_{jp}(x_p)$$

$x_1, \dots, x_p$   
CONTINUOUS

$$f(x_1, \dots, x_p | j) = f(x_1 | j) \cdot \dots \cdot f(x_p | j)$$



Predictors are conditionally  
independent given the class  
(naïve, idiots' Bayes, ...)

## (GAUSSIAN) NAIVE BAYES CLASSIFIER

$f_j(x)$  follows a multivariate Gaussian distribution  $\Rightarrow x_1, \dots, x_p$  Gaussian  
 (conditional on class  $j$ )

Recall:

$(X, Y)$  normally distribution:

$$X \perp Y \quad (\text{independent}) \Leftrightarrow \text{Cov}(X, Y) = 0$$

( $\Rightarrow$  always true)

$\Leftrightarrow$  Gaussian

So the matrix  $\Sigma$

$$\Sigma_j = \begin{pmatrix} \sigma_{1j}^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_{pj}^2 \end{pmatrix}$$

$\sigma_{ij}^2 = \text{Var}(x_i | y=j)$

covariance of  $x_1, \dots, x_p$  given  $y=j$

$$\sigma_{ij}^2 = \text{Var}(x_i | y=j)$$

LDA: variances depend on  $j$

QDA: off-diagonal of  $\Sigma_j$  are now zero

## NB DECISION SURFACE

$$f_{jK}(x_K) = \frac{1}{\sqrt{2\pi} \sigma_{Kj}} e^{-\frac{1}{2} \left( \frac{x_K - \mu_{Kj}}{\sigma_{Kj}} \right)^2}$$

So

$$\begin{aligned} P(y=j | x_1, \dots, x_p) &\propto \pi_j f_j(x) = \pi_j \prod_{k=1}^p f_{jk}(x_k) \\ &= \pi_j \prod_{k=1}^p \frac{1}{\sqrt{2\pi} \sigma_{kj}} e^{-\frac{1}{2} \left( \frac{x_k - \mu_{kj}}{\sigma_{kj}} \right)^2} \end{aligned}$$

Find  $j$  which maximizes this probability is equivalent to find the  $j$

$$S_j(x) = \log(\pi_j) - \sum_{k=1}^p \log(\sigma_{kj}) - \frac{1}{2} \sum_{k=1}^p \left( \frac{x_k - \mu_{kj}}{\sigma_{kj}} \right)^2$$

$\downarrow$   
 $(x_1, \dots, x_p)$

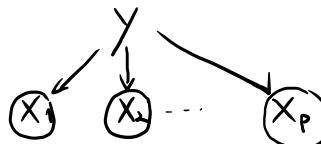
- Quadratic
- Does not involve  $\Sigma$ , so works well for  $p$  very large

## (MULTINOMIAL) NAIVE BAYES CLASSIFIER

$X_1, \dots, X_p$  categorical

The natural assumption is that  $f_{\text{NBR}}$  is multinomial.

How many parameters?



Let us assume that  $Y$  is binary ( $K=2$ )

$X_1, \dots, X_p$  are categorical with  $X_i$  having  $G_i$  categories

$$P(Y=j|X) \propto P(Y=j) P(X_1=x_1|j) \cdots P(X_p=x_p|j)$$

↑

$$\begin{aligned} \theta_{y1} &= P(Y=0) \\ \theta_{y2} &= P(Y=1) \end{aligned}$$

$$\begin{aligned} \theta_{11} &= P(X_1=G_1|Y=0) \\ \theta_{12} &= P(X_1=G_1|Y=1) \\ &\vdots \\ \theta_{1G_1} &= P(X_1=G_1|Y=1) \end{aligned}$$

$$\begin{aligned} \theta_{p1} &= P(X_p=G_p|Y=0) \\ \theta_{p2} &= P(X_p=G_p|Y=1) \\ &\vdots \\ \theta_{pG_p} &= P(X_p=G_p|Y=1) \end{aligned}$$

$X_j$   $(G_j - 1) \times 2$  parameters

$\theta_{ije}$	$i=1, \dots, p$
	$j=1, 2$
	$e=1, \dots, G_e$

Local / Global assumption  
of independence of  
parameters

## PARAMETER ESTIMATION

DATA :  $x_1, \dots, x_n$

→ Think of  $X_1 | Y=0$

independent observations from  $f(x; \underline{\theta})$   
parameters!

### ① FREQUENTIST:

Likelihood       $L(\theta) = \prod_{i=1}^n f(x_i; \theta)$

MLE       $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta)$       maximum likelihood estimate of  $\theta$

### ② BAYESIAN

$\theta$  is treated as a random variable with a prior distribution containing information about  $\theta$  before we see the data       $[\theta \sim U(2, 3)]$

Combine the likelihood (information from data) and the prior knowledge into a posterior distribution

## POSTERIOR DISTRIBUTION

$$g(\theta | \text{data}) = \frac{f(\text{data}|\theta) h(\theta)}{p(\text{data})}$$

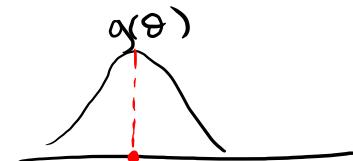
posterior                                  prior

$$= \frac{f(x_1, \theta) \cdot f(x_2, \theta) \cdots f(x_n, \theta) h(\theta)}{p(\text{data})}$$

Likelihood

$$\propto L(\theta) h(\theta)$$

posterior  $\propto$  likelihood  $\times$  prior



$$\hat{\theta} = E[\theta | \text{data}] = \int \theta g(\theta) d\theta$$

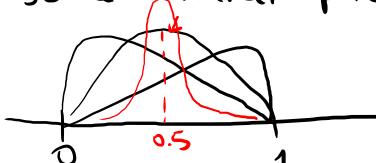
Bayesian estimate of  $\theta$

EXAMPLE (1)    BETA-BINOMIAL    (BINARY PREDICTOR)

$$X \sim \text{Bernoulli}(\theta) \quad f(x; \theta) = \theta^x (1-\theta)^{1-x} \quad x=0,1$$

$\theta$  is between 0 and 1 so a natural prior distribution is

Beta ( $\alpha, \beta$ )



$$h(\theta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \quad 0 \leq \theta \leq 1 \quad \begin{matrix} \text{PRIOR} \\ (\text{conjugate}) \end{matrix}$$

Data:  $x_1, \dots, x_n$

$$L(\theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i} \quad \text{LIKELIHOOD}$$

$$g(\theta | \text{data}) \propto \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$= \theta^{\sum_{i=1}^n x_i + \alpha - 1} (1-\theta)^{n - \sum_{i=1}^n x_i + \beta - 1}$$

$$\text{Beta}\left(\sum_{i=1}^n x_i + \alpha, n - \sum_{i=1}^n x_i + \beta\right)$$

POSTERIOR

updating prior parameters with data

EXAMPLE 2    DIRICHLET - MULTINOMIAL (CATEGORICAL PREDICTOR)

X has G categories

Then the data for n experiments can be summarised as

$$X = (X_1, \dots, X_G)$$

with  $X_i$  = "number of times category i is observed out of n trials"

MULTINOMIAL DISTRIBUTION

$X \sim \text{Multinomial}(\underline{\theta}_1, \dots, \underline{\theta}_G; n)$

$$f(x; \theta) = \frac{n!}{\prod_{i=1}^G x_i!} \prod_{i=1}^G \theta_i^{x_i}$$

$$\begin{aligned} x_i &= 0, \dots, n \\ \sum_{i=1}^G x_i &= n \end{aligned}$$

$G = 2$   
Binomial

## ESTIMATION OF $\theta_1, \dots, \theta_g$

### FREQUENTIST

Likelihood  $L(\underline{\theta}) = \frac{n!}{\prod_{i=1}^g x_i!} \prod_{i=1}^g \theta_i^{x_i}$

$$l(\underline{\theta}) = \log(n!) - \sum_{i=1}^g \log(x_i!) + \sum_{i=1}^g x_i \log(\theta_i)$$

### MLE

$$\max l(\underline{\theta})$$

s.t.  $\sum_{i=1}^g \theta_i = 1$

$$\rightarrow \boxed{\hat{\theta}_i = \frac{x_i}{n}}$$

OBSERVED  
FREQUENCIES

$$\sum_{i=1}^g \hat{\theta}_i = \frac{\sum x_i}{n} = 1$$

	1	2	...	g
$x_1$	$x_2$	-	-	$x_g$

$n = 100$

## BAYESIAN

PRIOR  $h(\theta) \sim \text{DIRICHLET} (\alpha_1, \dots, \alpha_q)$  BETA GENERALIZATION

$$\frac{1}{B(\alpha)} \prod_{i=1}^q \theta_i^{\alpha_{i-1}}$$

with  $B(\alpha) = \frac{\prod_{i=1}^q \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^q \alpha_i)}$

$$\sum_{i=1}^q \theta_i = 1$$

( $\alpha_1 = \alpha$   
 $\alpha_2 = \beta$ )  
 $\theta_2 = 1 - \theta_1$ )

Note :  $\alpha_{i-1}$  plays the role of  $x_i$  in the likelihood, so it can be thought as imaginary counts (before seeing the data)

A common choice is  $h(\theta) \sim D(1, 1, \dots, 1)$

$$E(\theta_i) = \frac{\alpha_i}{\sum_{j=1}^q \alpha_j}$$

PRIOR FREQUENCIES

## POSTERIOR

$$g(\theta | \text{data}) \propto L(\theta) h(\theta)$$

$$= \prod_{i=1}^q \theta_i^{x_i} \prod_{i=1}^q \theta_i^{\alpha_{i-1}} = \prod_{i=1}^q \theta_i^{x_i + \alpha_{i-1}}$$

$x_i + \alpha_{i-1}$

DIRICHLET ( $x_1 + \alpha_1, \dots, x_q + \alpha_q$ )

## PARAMETER ESTIMATES

$$\hat{\theta}_i^{(s)} = \frac{\alpha_i + x_i}{\sum_{j=1}^n \alpha_j + n}$$

BAYESIAN ESTIMATE OF  $\theta_i$

$$\hat{\theta}_i^{(MLE)} = \frac{x_i}{n}$$

MAXIMUM LIKELIHOOD ESTIMATE OF  $\theta_i$   
 n could very small or even zero

Under a uniform prior knowledge:

$$\hat{\theta}_i^{(s)} = \frac{1 + x_i}{G + n}$$

Makes estimation more stable in the  
 case of Naive Bayes classifiers (particularly  
 for large p)