

Zaki

①

## Bayes Classifier

BC directly uses Bayes theorem to predict the class for the new data instance

it estimates the posterior prob for each class  $P(c_i | x)$  & chooses the class w/ the largest prob

$$\hat{y} = \operatorname{argmax} \{ P(c_i | x) \}$$

Bayes theorem allows us to invert the posterior probability in terms of the likelihood & prior probability

$$P(c_i | x) = \frac{P(x | c_i) \cdot P(c_i)}{P(x)}$$

$P(x | c_i)$  is the likelihood defined as the prob of obs  $x$  given  $c_i$  is the true class  
 $P(c_i)$  prior prob of  $c_i$

$P(x)$  Prob of obs  $x$  from any class

$P(x)$  is fixed for a given point, so

$$\hat{y} = \operatorname{argmax} \{ P(x|c_i) P(c_i) \}$$

Predicted class depends on the likelihood of that test class taking into account its prior prob

We need to estimate the likli & prior probs from the training set

### ESTIMATING THE PRIOR PROBS

D<sub>i</sub> Denote subset of D labelled c<sub>i</sub>

$$D_i = \{ \overset{\top}{X_j} \mid X_j \text{ has class } y_j = c_i \}$$

Size D = |D| = n, size of D<sub>i</sub> = |D<sub>i</sub>| = n<sub>i</sub>

Prior Prob has be estimated as:

$$\hat{P}(c_i) = \frac{n_i}{n}$$

## ESTIMATING THE LIKELIHOOD

To estimate the likelihood  $P(x|c_i)$ , we have to estimate the joint probability of  $x$  across all the  $d$  dimensions

That is to say, we have to estimate  
 $P(x = (x_1, x_2, \dots, x_d) | c_i)$

the prob of seeing a data instance with these parameters given the class of  $c_i$

There is a difference in approach if the all the dimensions are numeric

The approach can be parametric or non-parametric. Here we explore parametric

To do this we assume each class  $c_i$  is normally distributed around some mean  $\mu_i$  & a covariance matrix  $\Sigma_i$

And estimated from the train set  $D$ .

4

for any class  $c_i$ , the prob density at  $x$  is thus given as

$$f_i(x) = f(x|\mu_i, \Sigma_i) = \frac{1}{(\sqrt{2\pi})^d \sqrt{|\Sigma_i|}} \exp\left(-\frac{(x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i)}{2}\right)$$

$\downarrow$  Det of Conver mat

This is the probability density function

We use this because the probability of a precise point is 0, hence we need to look @ a range (PDF)

This gives us the relative values of likelihood of being near  $x$

the PDF doesn't give us the probability directly. instead the prob density - prob per unit of  $x$

to get the actual prob we need to consider an interval around  $x$

(5)

the prob of  $X$  falling within that interval is given by the area under the PDF curve w/r to that interval

When we calculate PDF, we are calculating so for class  $c_1$  @ point  $x$

the part of the PDF that make ours class class specific is the inclusion of the classes  $\mu_{\text{mean}}$  &  $\Sigma_{\text{var}}$

recall the mean of a dataset will be a vector of the mean of each dim

Variance  $\Sigma_{\text{var}} = \text{covariance matrix}$   
where Diags are individual attribute variance & rest are covariance relationship

We compute the likelihood by considering a small interval  $\epsilon > 0$  centered @  $x$

6

$$P(x|c_i) = 2\epsilon \cdot f_i(x)$$

we use two because it represent  
two equal area both sides of  $x$

$\epsilon$  is a tiny number

the posterior probability is then given:

~~Bayes~~

	likelihood	prior	proportionality
	$f_i(x)$	$P(c_i)$	$\cancel{2\epsilon}$
$P(c_i x)$	$\frac{2\epsilon \cdot f_i(x) P(c_i)}{\sum 2\epsilon \cdot f_j(x_j) P(c_j)}$	$= \frac{f_i(x) P(c_i)}{\sum f_j(x_j) P(c_j)}$	
	$\underbrace{\quad}_{\text{same for all classes}}$		

Note if we want to calc ~~as~~ the likelihood we need  
the  $2\epsilon$  interval to get the area

However, in Bayes classifier we are interest in proportions  
between a class vs all classes

Hence the interval can be cancelled out

7

the bottom/denominator is the marginal prob  
prob

The prob of observing  $X$  over all classes  
- or regardless of class

This value will remain fix/constant  
for all calculations

that is for calcing the prob that  $X$   
belongs to each class

Because it is constant we can ~~ignore~~  
So, the calculation becomes:

$$\hat{y} = \operatorname{argmax} \{ f_i(x) P(c_i) \}$$

See that we just select the highest value = pred

note, if we want to see the actual probs  
pred and compare we would need the  
full calc w/ the denominator

(8)

Here, the  $2E$  still cancels out

Only need to include  $2E$  if you want to score even further in and get the likelihood

## Calculating Mean & Variance Params

Bayes uses sample mean & sample covariance

$$\text{mean} = \hat{\mu}_i = \frac{1}{n_i} \sum_{x_j \in D_i} x_j$$

$$\text{covar} = \hat{\Sigma} = \frac{1}{n_i} \bar{D}_i^T \bar{D}_i$$

where

$\bar{D}_i$  = centered matrix, calc mean of col  
& subtract from all values

$$\bar{D}_i = D_i - \mathbf{1} \cdot \hat{\mu}_i^T$$

these values are then plugged into the pdf eqs above

(9)

# Bayes Categorical Attributes

Let  $X_j$  be a categorical attribute

It can take  $n_j$  distinct categorical values

$$\text{Dom}(X_j) = \{a_{j1}, a_{j2}, \dots, a_{jn_j}\}$$

This is the set of values the attribute can take with  $a_{j1}$  being a particular value

Each cat attribute  $X_j$  takes on  $n_j$  dist  
 Each categorical attribute is modelled  
 as a  $n_j$ -dimensional Multivariate Bernoulli  
 random variable

Bernoulli = two categories

Multivariate =  $\geq 2$  categories rep as vector

$n_j$  dimensional = length of vector, num of cats

Vector values  $e_{j1}, e_{j2}, \dots, e_{jn_j}$  where  $e_r$   
 $r$  is the  $r^{\text{th}}$  standard Basis Vector in  $\mathbb{R}^{n_j}$

$e_{jR}$  and corresponds to the  $r^{th}$  value  
is  $a_{jR} \in \text{Dom}(x_j)$  (the var categories)

Standard basis vector =  $\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$   
Hot on encoding

where the  $j$  selects the categories

$$e_{jR} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} = a_{jR} = a_{j3}$$

why create this Multivariate Bernoulli Variable?

- it is not just about rewriting in a new format
- original cat simply tells you what category a data point belongs to
- the random variable allows you to express probabilities w.r.t each variables
- Essential for handle uncertainty & prediction
- it also allows to conduct mathematical operations on categorical data & algos

(11)

the entire  $D$ -dimensional dataset  
is modelled as the vector  $\text{RV} = \mathbf{X}$

$$\mathbf{X} = (X_1, X_2, \dots, X_d)^\top$$

This of this as tabular dataset

- $X_n$  = the column
  - each row would be a data instead
  - and each row holds a set of vectors which are the hot one encoded such
- For the columns attributes

~~that~~ Let  $d' = \sum_{j=1}^d M_j$

this is an important step in transforming the dataset into a format suitable for probabilistic modeling w/ categorical attr

it is a way to represent ~~all~~ the total dimension of dataset after all attrs have been hot one encoded

~~work w/ single encoded vector~~

(12)

$d'$  is the sum of the number of categories across all attributes

i.e. total number of columns if every this was not one encoded

it is also the length of a single concatenated vector of all attrs

Example

$X_1$ : Colour {red, green}  $M_1 = 2$

$X_2$ : Shape {Sq, C<sub>1</sub>, tri}  $M_2 = 3$

$X_3$ : Size {Sm, lg}  $M_3 = 2$

$d = 3$  = original Dimension

$$d' = 7 \quad M_1 + M_2 + M_3 = 2 + 3 + 2 = 7$$

(13)

A categorical data point is represented

$$x = (x_1, x_2, \dots, x_d)^T$$

where this is a data point where all dimensions are categorical

Binning as hot one encoding

Therefore the  $d'$ -dimensional Binning vector is:

$$v = \begin{pmatrix} v_1 \\ \vdots \\ v_d \end{pmatrix} \approx \begin{pmatrix} e_{1r_1} \\ \vdots \\ e_{dr_d} \end{pmatrix} \quad v_j = e_{jr_j}$$

tells us the vector is the corresponding standard basis hot-one from within each attr

tells the vector  $v$  is made up of other vectors from hot-one

j

# The Probability of Point $X$ (data point) (PMF)

is obtained from the joint prob mass function:  
for the vector Random variable  $X$ :

$$P(X|c_i) = f(v|c_i) = f(X_1 = e_{1,c_i}, \dots, X_d = e_{d,c_i} | c_i)$$

$P(X|c_i)$  the prob of obs  $X$  data point given that  
it belongs to class  $c_i$

$f(v)$  reps a function (Prob mass function) applied  
to a vector (single hot one encoded)

$f(X_1 = e, \dots, e | c_i)$  represents attributes of a  
data point, its hot one-encoded vector ( $e$ )  
& again belonging to class  $c_i$  w/  
the PMF function applied

the PMF can be directly estimated as

$$\hat{f}(v | c_i) = \frac{n_i(v)}{n_i}$$

$n_i(v)$  count the exact vector occurs in  $n_i$

Note that if the PMF (likelihood) comes out as 0, i.e. no vector counts for the class, then Bayes classify will come out as 0 prob. This is an issue for a number of reasons:

- (1) indicates limited data/missis data/sample
- (2) risks overfitting as model will be certain it wont occur
- (3) Violation of Bayes assumption that data is indep - dependant data causes 0's
- (4) Math issue  $\rightarrow X \cdot 0 = 0 = \text{nullify results}$   
 $\rightarrow \log(0) = \text{undefined}$

A Solution = Pseudo-Count

Assume a +1 count for each value

i.e. to assume  $X$  occurs at least once

$1 + \text{actual observed count of } v^{\text{in}} c_i$

the creates an adjusted probability mass at  $v$  given as:

$$\hat{f}(v|c_i) = \frac{n_i(v) + 1}{n_i + \prod_{j=1}^p M_j}$$

$\prod$  = Product, i.e. multiplication

Product of all numbers  $M_j$  from 1 to  $D$   
i.e. each Attribute  $D$

~~Given~~  $M_1 = 2, M_2 = 3, M_3 = 2$  Unique  
 $= 2 \times 3 \times 2 = 12$  Combinations

unique combos based on all Categorical Attr

PMF needs to then be calculated for each class to extract a prediction of class

(17)

## Challenges with Bayes Class

the main issue is lack of data to estimate the joint prob density (numerical) or mass function (categorical)

for numeric we need to calc  $\Theta(d^2)$  covariances

- As dims/Attrs increase this requires the estimation of too many paras
- $d^2$  = complexity grows exponentially w/ dimensions increasing

for cats we ned to calc joint prob for all possible combinations of attrs

To reduce these issues we can use a reduced set of paras in practice known as Naive Bayes Classifier