

## Classification

most of the problems would be to classify. be it classifying observations into categories or predicting the category of a new observation.

examples of classification problems:

- 1) set of symptoms of a person  $\rightarrow$  which medical condition
- 2) online banking service based on IP address, transaction history etc  $\rightarrow$  fraud or not?
- 3) DNA sequences with & without diseases  $\rightarrow$  which DNA mutations are disease causing?

Reasons why we can't use linear regression for classification

- 1) regression method can't accommodate a qualitative response with more than 2 classes  
(the order would alter the prediction coefficients.  
we can only set the order if the difference b/w the classes are clear, & equal. such as low, medium, high)
- 2) a regression method will not provide meaningful estimates of  $P(Y|X)$ , even with just 2 classes  
(the probability value won't be bound b/w 0 & 1)

## Logistic model

$p(x) = \beta_0 + \beta_1 x$  - but the value will not be bound to  $[0, 1]$   
to avoid this problem, model  $p(x)$  using a function that gives the outputs between 0 & 1. for all values of  $x$

$$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

the logistic function will always produce an S shaped curve, regardless of the value of  $x$ .

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

- odds ratio  $[0, \infty)$

By taking logs at both sides

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

log odds or logit

increase of one unit of  $X$  changes the log odds by  $\beta_1$ .  
 the amount that  $P(X)$  changes depends on the value of  $X$ .  
 if  $\beta_1$  is pos, increasing  $X$  will increase  $P(X)$   
 vice versa.

maximum likelihood - we seek the estimates for  $\beta_0$  &  $\beta_1$ , such that the predicted probability  $\hat{p}(x_i)$  of a category is as close to the category of individual's observation.

$$\text{likelihood function } l(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i:y_i=0} (1-p(x_i))$$

multiply the prob of all observations where outcome is 1  
 ↓  
 outcome is 0

$H_0: \beta_1 = 0$  (probability doesn't depend on the predictor)

z statistic high  $\rightarrow$  reject the null hypothesis

$$\text{z statistic of } \beta_1 = \frac{\hat{\beta}_1}{\text{SE}(\hat{\beta}_1)}$$

z stat high  $\Rightarrow$  p stat low

$$H_0 \text{ implies that } P(X) = \frac{e^{\beta_0}}{1+e^{\beta_0}}$$

? Why z statistic here, but not t statistic?

For qualitative predictors, use dummy variables like 0,1 etc.

confounding - when the coefficient is of one sign when considered alone vs considered with other predictors (correlated)

e.g. students on avg are more likely to default

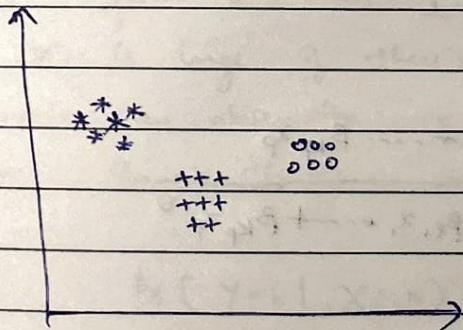
but when only the credit card balance is considered, their prob of default is lower than non students

hence it is important to identify and consider different terms before taking a decision from model

Multiple Logistic regression:-  $\log\left(\frac{P(x)}{1-P(x)}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$

$$P(x) = \frac{e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}$$

Multinomial Logistic regression



Select a single class as baseline  
then calculate the prob of other  
classes with reference to this

$$\Pr(Y=k | X=x) = \frac{e^{\beta_{k0} + \beta_{k1} x_1 + \dots + \beta_{kp} x_p}}{1 + \sum_{l=1}^{K-1} e^{\beta_{l0} + \beta_{l1} x_1 + \dots + \beta_{lp} x_p}}$$

for  $k=1 \dots K-1$

$$\Pr(Y=k | X=x) = \frac{1}{1 + \sum_{l=1}^{K-1} e^{\beta_{l0} + \beta_{l1} x_1 + \dots + \beta_{lp} x_p}}$$

$$\log\left(\frac{\Pr(Y=k | X=x)}{\Pr(Y=l | X=x)}\right) = \beta_{k0} + \beta_{k1} x_1 + \beta_{k2} x_2 + \dots + \beta_{kp} x_p$$

coefficient estimates will differ based on the chosen baseline class, but the prediction log odds ratio will be same.

be careful when interpreting the coefficients.

ex: if we set A to be baseline (A, B, C classes), then we can interpret  $\beta_B$  as the log odds of B vs A.

further one unit increase in  $x_j$  is associated with a  $\beta_{Bj}$  increase in the log odds of B over A.

if  $x_j$  increases by one unit, then

$$\frac{\Pr(Y=B | X=x)}{\Pr(Y=A | X=x)} \text{ increases by } e^{\beta_{Bj}}$$

Softmax coding - rather than selecting a baseline class, we treat all K classes symmetrically and assume for  $k=1, 2, \dots, K$

$$\Pr(Y=k | X=x) = \frac{e^{\beta_{k0} + \beta_{k1}x_1 + \dots + \beta_{kp}x_p}}{\sum_{l=1}^K e^{\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p}}$$

$$= \frac{e^{\beta_{k0} + \beta_{k1}x_1 + \dots + \beta_{kp}x_p}}{\sum_{l=1}^K e^{\beta_{l0} + \beta_{l1}x_1 + \dots + \beta_{lp}x_p}}$$

$$\log \left( \frac{\Pr(Y=k | X=x)}{\Pr(Y=k' | X=x)} \right) = (\beta_{k0} - \beta_{k'0}) + (\beta_{k1} - \beta_{k'1})x_1 + \dots + (\beta_{kp} - \beta_{k'p})x_p$$

Generative models for classification

Logistic Regression involves directly modeling  $\Pr(Y=k | X=x)$ .

In another approach, we first model the distributions of predictors  $X$  separately in each of the response classes (for each value of  $Y$ )

then use Baye's theorem to flip these around into estimates for  $P_3(Y=k | X=x)$

When the distribution of  $X$  within each class is assumed to be normal, it turns out that the model is very similar to logistic regression model.

why do we need another method?

- when there's substantial separation b/w 2 classes, the parameter estimates for the logistic regression model are unstable
- if the distribution of predictors  $X$  is approximately normal in each of the classes & the sample size is small, then these approaches are ill to logistic regression

$\pi_k$  - prior probability that a random observation comes from  $k$ th class

$f_k(x) = P(X=x | Y=k)$  - density function of  $X$  of an observation from  $k$ th class

( $f_k(x)$  is large if there's high prob that  $x$  is in  $k$ th class, there's an observation i.e.  $X=x$ . & vice versa for lower prob)  
then Baye's theorem states that

$$P_3(Y=k | X=x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}$$

$P_k(x) = \Pr(Y=k | X=x)$  - posterior prob that an observation  $X=x$  belongs to class  $k$ .

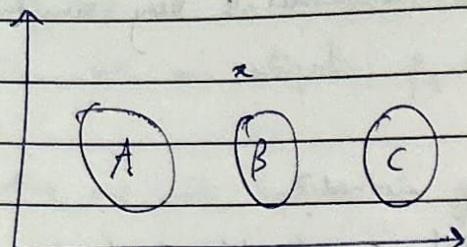
calculating  $\pi_k$  is easy. find out fraction of observations that belong to class  $k$ . (prob of class  $k$ )

but calculating the estimating the density function  $f_k$  requires some assumptions

Some techniques - linear discriminant analysis, quadratic discriminant analysis, naive Bayes

Bayes theory:  $P(\theta \mid \text{data}) \propto P(\text{data} \mid \theta) \times P(\theta)$

posterior                      likelihood              prior



$$P(A|x) \propto p(x|A) p(A)$$

A function  $f(x)$  is a density function if

- $f(x) \geq 0$  for all  $x$

- $\int_{-\infty}^{\infty} f(x) dx = 1$

- $P(a \leq X \leq b) = \int_a^b f(x) dx$

example: if  $X \sim N(0, 1)$

then  $f(0)$  is high  $\rightarrow$  many values near 0

$P(-1 \leq X \leq 1)$  is area under the curve

density function describes how probability is distributed over values of a continuous random variable

Probability mass function

- $p(x) \geq 0$

- $\sum p(x) = 1$

- $P(a \leq X \leq b) = \sum_a^b p(x)$

this is for discrete random variables

example: coin toss

$$P(\text{head}) = 0.5 - P(\text{tail}) = 0.5$$

$$P_r(Y=k | X=x) = \frac{\pi_k f_k(x)}{\sum_{k=1}^K \pi_k f_k(x)} \quad \textcircled{A}$$

Linear discriminant analysis for p=1

assume  $f_k(x)$  is normal or Gaussian,

normal density,  $f_k(x) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{1}{2\sigma_k^2}(x-\mu_k)^2\right) \quad \textcircled{B}$

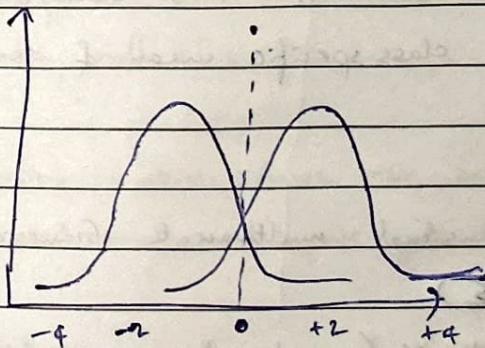
assume  $\sigma_1^2 = \dots = \sigma_K^2$

substituting \textcircled{B} in \textcircled{A} and taking the log

$$\delta_k(x) = \frac{x\mu_k - \frac{\mu_k^2}{2\sigma^2}}{\sigma^2} + \log(\pi_k) \quad \textcircled{C}$$

If  $K=2$  &  $\pi_1 = \pi_2$ , then Bayes classifier assigns an observation to class 1 if  $2x(\mu_1 - \mu_2) > \mu_1^2 - \mu_2^2$ .

else to class 2



$$\mu_1 = -1.25$$

$$\mu_2 = 1.25$$

$$\sigma_1^2 = \sigma_2^2 = 1$$

$$\text{assume } \pi_1 = \pi_2 = 0.5$$

decision boundary - point at which  $\delta_1(x) = \delta_2(x)$  &  
that's equal to  $\frac{\mu_1^2 - \mu_2^2}{2(\mu_1 - \mu_2)} = \frac{\mu_1 + \mu_2}{2}$

then for the above graph  $\frac{-1.25 + 1.25}{2} = 0$ .

if the point is before 0, it is class 1. else class 2.

In reality, to apply Bayesian classifier we need to know

$$\mu_1, \mu_2, \dots, \mu_K, \sigma^2, \pi_1, \pi_2, \dots, \pi_K$$

linear discriminant analysis estimates these values

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{\substack{i: y_i=k}} x_i \quad - \text{avg of all training obs from class } k$$

$$\hat{\sigma}^2 = \frac{1}{n-k} \sum_{k=1}^K \sum_{\substack{i: y_i=k}} (x_i - \hat{\mu}_k)^2 \quad - \begin{array}{l} \text{weighted avg of sample variances} \\ \text{for each of } K \text{ classes} \end{array}$$

$$\hat{\pi}_k = \frac{n_k}{n}$$

plugging the above

$$\hat{S}_k = x \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k}{2\hat{\sigma}^2} + \log(\hat{\pi}_k) \quad \textcircled{2}$$

assign  $X=x$  to  $k$  for which  $\textcircled{2}$  gives the largest value

LDA classifies assumes that observations within each class come from normal distribution, class specific mean & common variance

LDA for  $p > 1$

to indicate that a  $p$  dimensional multivariate Gaussian dist., we write  $X \sim N(\mu, \Sigma)$ .

$E(X)=\mu$  is the mean of  $X$  (vector of  $p$  components)

$Cov(X)=\Sigma$  is the  $p \times p$  covariance matrix of  $X$

multivariate Gaussian density fn.

$$f(x) = \frac{1}{2\pi^{p/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right)$$

in the case of  $p > 1$ , LDA assumes that observations in  $K$ th class are drawn from multivariate Gaussian distribution  $N(\mu_k, \Sigma)$ ,  $\mu_k$  is class specific mean vector,

$\Sigma$  is covariance matrix that's common to all classes.

Performing some algebra, we get-

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

assign to class  $k$  for which the above formula gives largest value.  
above is the matrix equivalent of (2) 2 pages ago.

as before, we need to estimate  $\mu_1, \mu_2, \dots, \mu_k, \pi_1, \pi_2, \dots, \pi_k, \Sigma$

lookout for low error rates: it can be because

- higher the p:u ratio - overfitting
- when the observations in a class itself is low, a useless classifier which predicts wrong will have low error rate  
(only 3% ppl default so a classifier that predicts no one defaults will have error rate 3%)
- aka null classifier

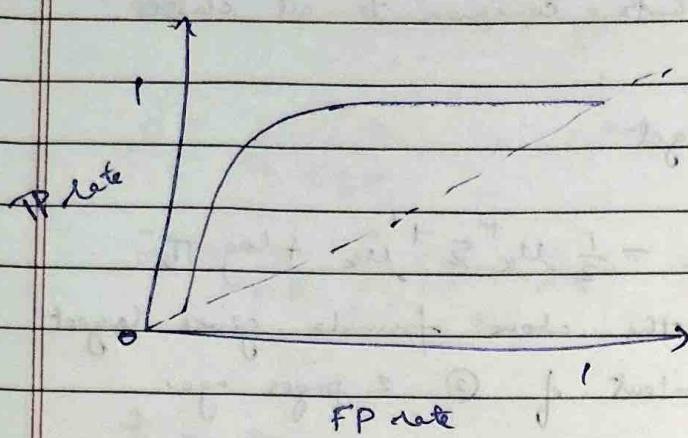
confusion matrix shows the categories of how many were correctly predicted and not; we get to know false pos & false neg

sensitivity - true pos %, all pos :  $TP / TP + FN$

specificity - true neg %, all neg :  $TN / TN + FP$

in a Bayes classifier, we can change the threshold to accommodate the risk to handle the risk that's caused by FP or FN according to business / problem

ROC curve - displays 2 types of errors for all possible thresholds, a good ROC curve will bring the top left corner.  
overall performance of the classifier is given by area under the ROC curve



$\dots \rightarrow$  no information  
classif., useless

		True Class	
		-	+
Predicted class	-	TN	FN
	+	FP	TP
		N	P

FP rate	$FP/N$	Type I error, 1 - specificity
TP rate	$TP/P$	1 - type II error, power, recall, sensitivity
Pos pred value	$TP/P^*$	precision, 1 - false discovery proportion
Neg pred value	$TN/N^*$	

OMG.

### Quadratic discriminant analysis (QDA)

Unlike LDA, here it's assumed that each class has its own covariance matrix  $\Sigma_k$

$$x \sim N(\mu_k, \Sigma_k)$$

$$\delta_k(x) = (\text{Some huge formula})$$

$$= -\frac{1}{2}x^T \sum_k^{-1} x + x^T \sum_k^{-1} \mu_k - \frac{1}{2} \mu_k^T \sum_k^{-1} \mu_k$$

$$-\frac{1}{2} \log |\Sigma_k| + \log P_k$$

Why LDA or QDA?  $\Rightarrow$  Bias-variance tradeoff

bias - error from wrong assumptions in model

high bias = underfitting

straight line fitted on curved line

variance - over sensitive

high variance = overfitting

wiggly line that passes through every data point

- when there are  $p$  predictors, estimating covariance matrix requires  $p(p+1)/2$  parameters

for QDA when  $\Sigma$  is different for each class,  $Kp(p+1)/2$  parameters

- so, if we assume a single covariance matrix in LDA, it becomes linear in  $x$  - only  $Kp$  linear coefficients.

LDA is much less flexible,

- if the assumption that  $K$  classes share same covariance matrix is false, then LDA can suffer from high bias.

- prefers

LDA - relatively fewer training observations

QDA - training set is huge

### Naive Bayes

In LDA & QDA the assumption that  $f_k$  is a multivariate normal random variable distribution, class specific mean  $\mu_k$ , single of class specific covariance matrix ( $\Sigma$  &  $\Sigma_k$ ) were made.

& the problem of estimating  $K$  p-dimensional density functions were replaced into estimating  $K$  p-dimensional vectors of  $1 \times K$   $p \times p$  covariance matrices

in Naive Bayes we make assumption that within  $k^{th}$  class, the  $p$  predictors are independent

$$f_k(x) = f_{k_1}(x_1) \times f_{k_2}(x_2) \times \dots \times f_{k_p}(x_p)$$

estimating a  $p$ -dimensional density  $f_k$  is hard because we should consider marginal distribution of each predictor, also the joint distribution.

But by assuming  $p$  covariates are independent, i.e. each class, we eliminate the need of association b/w predictors.

in reality it's not independent, but it works. so lalalalala

$$\Pr(Y=k | X=x) = \frac{\pi_k \times f_{k_1}(x_1) \times f_{k_2}(x_2) \times f_{k_3}(x_3) \times \dots \times f_{k_p}(x_p)}{\sum_{i=1}^K \pi_i \times f_{k_i}(x_1) \times f_{k_i}(x_2) \times \dots \times f_{k_i}(x_p)}$$

to estimate 1D density function  $f_{kj}$  using training data:

$x_{ij}, x_{j1}, \dots, x_{jk}$ , we have the options

- $x_j | Y=k \sim N(\mu_{jk}, \sigma_{jk}^2)$  (if  $x_j$  is quantitative)

univariate normal distribution

like in QDA, but the covariance matrix is diagonal since the predictors are independent

- non parametric estimate for  $f_{kj}$  (if  $x_j$  is quantitative)
  - plot histograms with same width for each class
  - check which bin ~~has~~ the it falls to, compare & choose the class which has highest density
  - OR kernel density for - smoothed histogram, then evaluate the curve at  $x_j$  to get the density
- count the proportion of training observations for  $j^{th}$  predictor corresponding to each class (qualitative)

Sometimes Naïve Bayes has slightly less overall error rate when  $p$  is larger or  $n$  is smaller, Naïve Bayes will perform better than LDA or QDA.

### Comparison of classification methods

#### Analytical comparison

- LDA is a special case of QDA
- (LDA is restricted version of QDA with  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K$ )
- Bayes classifier with a linear decision boundary is a special case of Naïve Bayes.

Surprisingly this means LDA is a special case of Naïve Bayes

- Naïve Bayes is a special case of LDA with  $\Sigma$  restricted to be a diagonal matrix with  $j^{th}$  diagonal element equal to  $\sigma_j^2$
- Neither QDA nor Naïve Bayes is a special case of others.

Naïve Bayes produces a more flexible fit which is additive.

If the interactions among the predictors are important in discriminating b/w classes QDA is better.

$$\log \left( \frac{P_k(Y=k | X=x)}{P_l(Y=l | X=x)} \right) = a_k + b_{kj} x_j \quad -\text{LDA}$$

like logistic regression,

assumes log odds of posterior prob is linear in  $x$ .

$$= a_k + \sum_{j=1}^p b_{kj} x_j + \sum_{j=1}^p \sum_{l=1}^p c_{kj} x_j x_l \quad -\text{QDA}$$

assumes log odds of posterior prob is quadratic in  $x$

$$= a_k + \sum_{j=1}^p g_{kj} x_j \quad -\text{Naïve Bayes}$$

takes the form of generalised additive model

- choice of method depends on the distribution of predictors in each of  $K$  classes & the bias-variance tradeoff ( $n, p$  numbers size)

- LDA will outperform logistic regression when the normality assumption holds true.
- KNN - non parametric. will perform good when the decision boundary is highly non linear and  $n$  is very large and  $p$  is small.
- When  $n$  is modest &  $p$  is not very small & decision boundary is non linear, QDA  $\Rightarrow$  KNN
- unlike logistic regression, KNN doesn't tell which predictors are important.

### Empirical comparison

Refer textbook for scenarios of the graph of error rates

when the true decision boundaries are linear, LDA & logistic regression will perform well  
 moderately non linear - QDA & Naive Bayes  
 more complicated - non parametric approach such as KNN  
 but the level of smoothness must be chosen carefully

### Generalized linear models

#### \* Linear regression on bike share data

- the output is not necessarily an integer
- it could be negative numbers
- data is heteroscedastic. as the mean increases, variance also increases - thus not suitable for linear regression

$$Y = \sum_{j=1}^p X_j \beta_j + \epsilon$$

here  $\epsilon$  is assumed to be mean zero term,

or<sup>2</sup> variance as constant which is violation of assumptions

of linear model

#### \* Poisson regression

$$P(Y=k) = \frac{e^{-\lambda}}{k!} \lambda^k \quad \text{for } k=1, 2, \dots$$

$Y$  takes on non-negative integers

$$\lambda = E(Y) = \text{Var}(Y)$$

(larger the mean - larger the variance)

Poisson distn is typically used to model counts

if we model  $Y$  as poisson distn with:  ~~$E(Y)$~~   $E(Y) = \lambda = 5$ ,

$$\text{the } P(Y=0) = \frac{e^{-\lambda} \lambda^0}{0!} = 0.0067$$

probability of finding no users(counts) at that hour/cond

we expect the mean  $\lambda = E(Y)$  to vary

$$\log(\lambda(X_1, X_2, \dots, X_p)) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

$$\lambda(X_1, X_2, \dots, X_p) = e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p} \quad \text{---(1)}$$

to estimate the coefficients  $\beta_0, \beta_1, \dots, \beta_p$  use the max likelihood est'n

$$l(\beta_0, \beta_1, \dots, \beta_p) = \prod_{i=1}^n \frac{e^{-\lambda(x_i)} \lambda(x_i)^{y_i}}{y_i!} \quad \text{where } \lambda = \text{---(1)}$$

### Important distinctions

- an increase in  $X_j$  by one unit is associated with a change in  $E(Y) = \lambda$  by a factor of  $e^{\beta_j}$
- unlike the assumptions on linear regression, here the mean & variance relationship are handled
- no negative predictions

[ over dispersion - variance much higher than mean. inflates Z values  
topic to be checked later ]

\* we assume  $Y$  follows

normal (Gaussian) distn in linear regression

Bernoulli distn

Poisson distn

logistic regression

Poisson regression

these 3 distributions come from exponential family.

other dist<sup>n</sup>s in this are Gamma dist<sup>n</sup>, negative binomial dist<sup>n</sup>

We can perform regression by modeling the response  $y$  as coming from the exponential family, then transform the mean so that it is a linear fn of predictors.

This regression approach is called Generalised Linear Model