

## Classification

Given some data, which of the given  $X$  classes does that observation belong to

Why not LR.

Eg → classifying { Epilepsy  
drug overdose  
Stroke

I can cascade the output as 1, 2, 3  
this is because giving 1 then 2  
would imply PO = stroke - Epilepsy  
that is correct so we wont get right and

for 2 classes we can use a dummy variable  $y = \begin{cases} 0 & \text{if stroke} \\ 1 & \text{if drug overdose} \end{cases}$

Why not to use LR

i) won't work for  $> 2$  classes

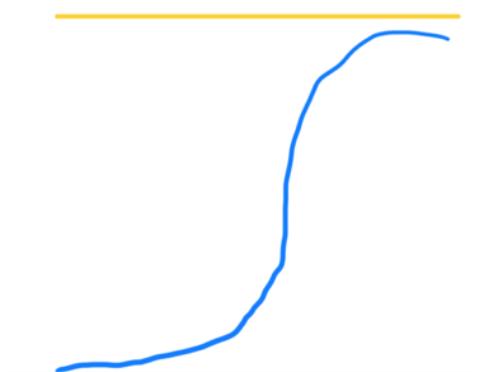
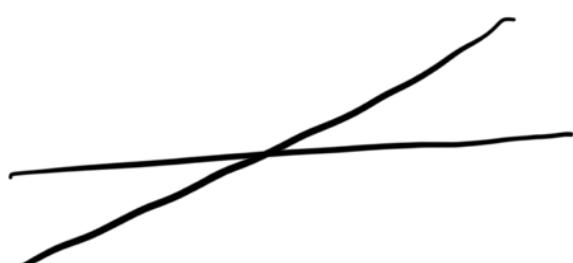
2) Won't provide meaningful estimates.

## Logistic Regression

models the probability that Y - belongs to a particular category

e.g. → Default  $\xrightarrow{\text{balance}}$

Log-R.  $\rightarrow P(\text{default} = \text{Yes} / \text{balance})$   
 $P(\text{balance})$



Using  
Linear Regression

Using logistic  
Regression

if we model with

$$p(x) = \beta_0 + \beta_1 x$$

then we can get a -ve prob.  
& a prob.  $p(x) > 1$

to get output between 0 and 1,

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

We use a method Maximum likelihood  
Indicator.

$$\ln \frac{P(n)}{1 - P(n)} = \beta_0 + \beta_1 x$$



$$\log\left(\frac{p_{(n)}}{1-p_{(n)}}\right) = \beta_0 + \beta_1 x$$

$\hookrightarrow$  Log odds / logit

Estimating the Regression Coefficients

Likelihood function :-

$$l(\beta_0, \beta_1) = \prod_{i:y_i=1} p(u_i) \prod_{i:y_i=0} (1-p(u_i))$$

We have z-statistic to measure significance.

$$H_0 : \beta_1 = 0 \Rightarrow p_{(n)} = \frac{e^{\beta_0}}{1+e^{\beta_0}}$$

Maximum Likelihood Estimation

## Multiple logistic regression

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

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

Using  $P_{21}$  may give a diff co-eff  
in terms of its sign when compared

to  $P_{12}$  in MLR. as other parents  
may also effect it.

{ This is known as confounding }

effect of outcome of one variable is  
mixed up with the effect of another  
variable.

# Multinomial Logistic Regression

e.g. → Stroke, Drug Overdose, Epileptic Seizure.

$$\underline{k=3}$$

$$; \underline{W.L.O.G}$$

we select  $k^{th}$  class

$$P(Y=k | X=n) = \frac{e^{\beta_0 + \dots + \beta_{kp} n_p}}{1 + \sum_{l=1}^{k-1} e^{\beta_{l0} + \beta_{l1} + \dots + \beta_{lp}}}$$

$$\log \left( \frac{P(Y=k | X=n)}{P(Y=l | X=n)} \right) = \beta_{k0} + \beta_{k1} n_1 + \dots + \beta_{kp} n_p$$

Software Coding → alternate for  
Multinomial Regression

we treat all  $K$  classes symmetrically

$$0, \dots, 1, \dots, 1$$

$$\beta_0 + \beta_1 n_1 + \dots + \beta_{kp} n_p$$

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

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

Generative model for Classification

$$P(Y = k | X = x) \rightarrow \text{Bayes' Theorem}$$

$\Rightarrow$  when there is substantial separation  
btw the two classes. The parameter  
estimates for logistic regression are highly  
unstable.

for  $K \geq 2$

$\pi_k$  → prior probability.

$$f_k(x) = P(X = k)$$

density of  $f^m$  of  $x$ , for an obs.  
coming from  $K^{tm}$  class.

Using Bayes theorem

$$P(Y = k | X = x) = \frac{n_k f_k(x)}{\sum_{l=1}^K n_l f_l(x)}$$

$$P_k(n) \doteq \underline{P(Y = k \mid X = n)}$$

$$m_k \text{ is Simple} \rightarrow \frac{\text{no. of obs. in } k^{\text{th}} \text{ class}}{\text{Total Obs.}}$$

finding  $f_k(n)$  is tricky, typically we have to make simplifying assumptions.

A) Lineer DISCRIMINANT

$\Rightarrow$  assume  $P=1$  {only 1 predictor}

$\Rightarrow$  we assume  $f_k(x)$  is Normal

$$f_k(x) = \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{(x-\mu_k)^2}{2\sigma_k^2}}$$

$\Rightarrow$  assume  $\sigma_1^2 = \sigma_2^2 = \sigma_3^2 = \dots = \sigma_k^2 = \sigma^2$   
for simplicity.

$$P_k(x) = \frac{n_k \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu_k)^2}{2\sigma^2}}}{\sum_{i=1}^K n_i \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu_i)^2}{2\sigma^2}}} \quad \rightarrow \textcircled{1}$$

from Stat quest  $\rightarrow$  LDA reduces

the dimensionality and helps project

the data points on a single line

$\mu_k \rightarrow$  Mean of  $k^{\text{th}}$  category.

diff. b/w means

$$\frac{(\mu_1 - \mu_2)^2}{S_1^2 + S_2^2}$$

minimize this

space

minimize this

$\therefore$  Maximize separation and minimize  
the scatter / spread.

this is better than PCA for higher  
dimensions.

taking log of  $c_g^n$  & and re-arranging

$$S_k(n) = n \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(n_k)$$

$c_g$ : for  $k=2$  & if  $n_1 = n_2$

decision  
boundary

$$\rightarrow x = \frac{\mu_1^2 - \mu_2^2}{2(\mu_1 - \mu_2)} = \frac{\mu_1 + \mu_2}{2}$$


---

In reality we plug estimates for  $\mu_1, \mu_K$

and  $\sigma^2$

$$\hat{\mu}_K = \frac{1}{n_K} \sum_{i:y_{ik}} u_i ; \quad \sigma^2 = \frac{1}{n-K} \sum_{k=1}^K \sum_{i \neq i_k} (u_i - \hat{\mu}_k)^2$$

where  $n =$  total number of observations

$$\hat{\mu}_K = n_K / n$$

$$\hat{f}_K(u) = n \cdot \frac{\hat{\mu}_K}{\sigma^2} - \frac{\mu_K^2}{2\sigma^2} + \log(\hat{\pi}_K)$$

discriminant function

assign the observation to the

(class for which eq<sup>n</sup> 2.2 is

the largest

**Major assumption** is that each observation  
comes from a normal distribution

comes from normal dist -

LDA  $\rightarrow$  Linear discriminant Analysis

for  $p > 1$

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

Multivariate normal distribution.

$$X \sim N(\mu, \Sigma)$$

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

$\Delta \text{cov}(X) = \Sigma$  is a  $p \times p$

covariance matrix.

$$f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} [(x-\mu)^T \Sigma^{-1} (x-\mu)]}$$

$$\tilde{S}_K(x) = x^T \Sigma^{-1} \mu_K - \frac{1}{2} \mu_K^T \Sigma^{-1} \mu_K + \log \pi_K$$

Measuring the error → Confusion Matrix  
 for credit default

		True		$\Sigma$
		No	Yes	
Preds	No	9649	252	9896
	Yes	23	81	104
		9667	333	

$$\text{Accuracy} = \frac{\text{correct}}{\text{total}} \times 100 = 97.25\%$$

$$\text{Error} = 2.75\%$$

but a trivial pred if it always says no.  
 It have an error of 3.33%. so it  
 is only slightly worse.

but of 333 actual defaulters it  
 predicted 252 to not default  
 that is an error of 75.7%.

which is unacceptable.

$$FN \approx 75.7\%$$

$$\underline{\text{Precision}} \rightarrow \frac{\text{True } + \text{VR}}{\text{TP} + \text{FP}}$$

$$\underline{\text{Recall / Sensitivity}} \rightarrow \frac{\text{TP}}{\text{TP} + \text{FN}}$$

for the example  $\hat{=} 24.1\%$

$$\text{Specificity} \rightarrow \frac{\text{TN}}{\text{TN} + \text{FP}} \rightarrow \text{False that are detected correctly.}$$

$$\underline{\text{F1 score}} \rightarrow \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Recall} + \text{Precision}}$$

This large disparity is because LDA is rewarded to minimize the incorrect predictions so it optimizes based on data

the one with majority, but a banking company will prefer LDA correctly identifying defaulter by sacrificing some accuracy on non-defaulter.

Modifying it to create a better classifier

$$P(\text{default} = \text{Yes} | X = n) > 0.5$$

we can lower this to 50% to

20%. if we are concerned by incorrectly predicting the default class

$$P(\text{default} = \text{Yes} | X = n) > 0.2$$

taking 20% we get

		True	
		Yes	No
Pred	Yes	195	235

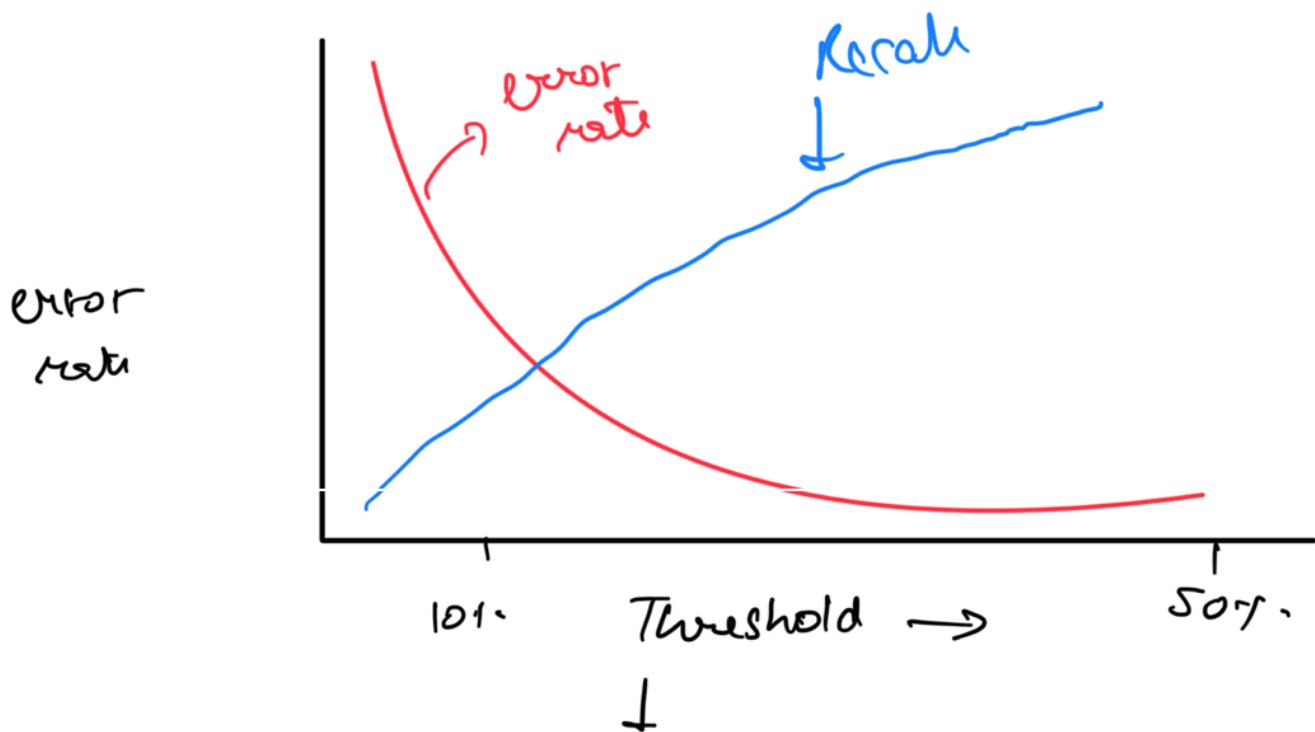
$N_0$	138	9432
-------	-----	------

333      9667

$$\text{Recall} = \frac{138}{333} \rightarrow 41.4\%$$

Error rate

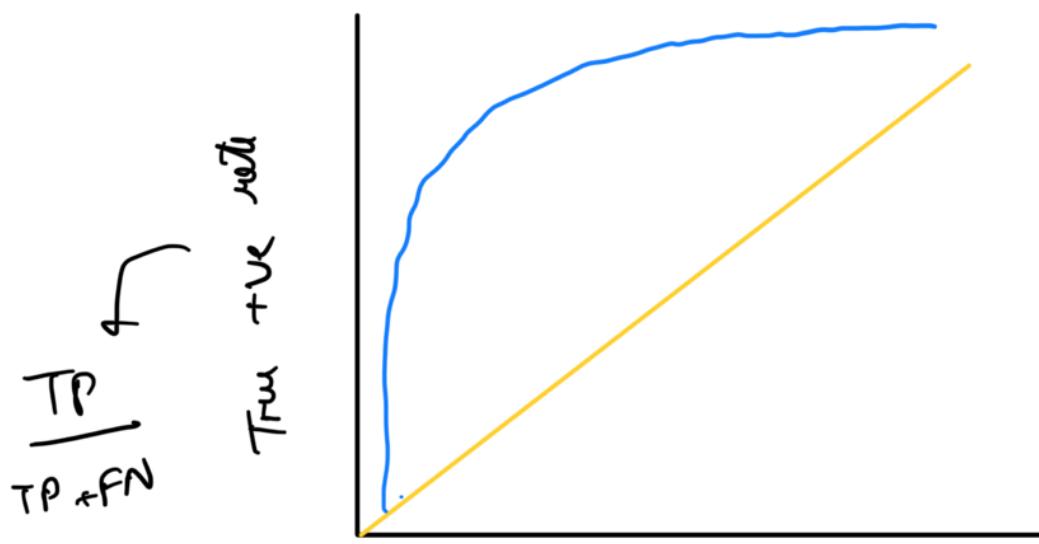
increased to 3.73%



$$P(\text{disease} = \text{Yes} / x = u) \geq \text{Threshold}$$

deciding the threshold value requires  
domain knowledge

ROC → Receiver operating characteristic



False positive rate  
↳  $\frac{FP}{FP + TN}$

area under curve of ROC shows accuracy closer to 1 is best.

by chance classifier will have 0.5

Quadratic discriminant Analysis

- more flexible
- covariance matrix

## Naive Bayes Classifier

$$P_K(n) = P(Y = K \mid X = n)$$

assumption  $\rightarrow$  within the  $K^{\text{th}}$  class,  
the  $p$  predictors are independent

$$f_K(n) = f_{K_1}(n_1) \cdot f_{K_2}(n_2) \cdot \dots \cdot f_{K_p}(n_p)$$

$n_K \rightarrow$  prior of being in class  $K$

$f_{K_p}(n_p) =$  if it was in class  $K$   
then  $p$  of  $n_p$

$$P(Y = K \mid X = n) = n_K \cdot f_{K_1}(n_1) \cdot \dots \cdot f_{K_p}(n_p)$$

$$\frac{\sum_{k=1}^K n_k \cdot f_{k_1}(n_1) \times \dots \times f_{k_p}(n_p)}{\sum_{k=1}^K n_k}$$

## Poisson Regression on count data

$$P(Y=k) = \frac{e^{-\lambda} \lambda^k}{k!} \quad k \in N \cup \{0\}$$

$$E(Y) = \lambda$$

$$\log(\lambda(x_1, \dots, x_p)) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$$\textcircled{3} - \lambda(x_1, x_2, \dots, x_p) = e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}$$

(so that  $\lambda(x_1, \dots, x_p)$  takes non-negative values)

$$\ell(\beta_0, \beta_1, \dots, \beta_p) = \frac{\prod_{i=1}^n e^{-\lambda(n_i)} \lambda(n_i)^{y_i}}{y_i!} \quad \text{(4)}$$

$$\text{where } \lambda(n_i) = e^{\beta_0 + \beta_1 n_{i1} + \dots + \beta_p n_{ip}} \begin{cases} \text{from eq 3} \\ \end{cases}$$

Under poisson regression

... ... ... ...

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

## Generalised Linear Models

all use predictors	$x_1, x_2, \dots, x_p$	<u>distribution</u>
<u>Method</u>		
Linear Regression		Normal
Logistic Regression		Bernoulli
Poisson Regression		Poisson

### Expected Value

$$\text{Lin } E(Y | x_1, \dots, x_p) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$$\text{LOR } E(Y | x_1, \dots, x_p) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}$$

$$\text{PdL } E(Y | x_1, \dots, x_p) = \lambda(x_1, \dots, x_p)$$

$$= e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}$$

In general

$$\eta(E(Y|x_1, \dots, x_p)) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$$\text{Lin } \eta(u) = u$$

$$\text{Log } \eta(u) = \log\left(\frac{u}{1-u}\right)$$

$$\text{Pois } \eta(u) = \log(u)$$

These are GLM Generalised  
Linear Modeling (GLM)