

# Logistic Regression

## • Logistic Regression

- Don't be confused by the use of the term "regression" in its name.
- Logistic Regression is a classification algorithm designed to predict categorical target labels.

## Section Overview

- Transforming linear Regression to logistic regression.
- Mathematical theory behind Logistic Regression
- Simple implementation of Logistic Regression for classification problems.
- Interpreting Results:-
  - Odds Ratio and Coefficients
  - Classification Metrics
    - Accuracy
    - Precision
    - Recall
  - ROC curves ( Receiver Operator Characteristic Curves ) ✓
- Multiclass classification with logistic Regression -
  - Logistic Regression project
  - Logistic Regression project solution.

12

- Logistic Regression will allow us to predict a categorical label based on historical features data-

- The categorical target column is two or more discrete class labels.

- Classification algorithm predict a class or categorical label:

- Class 0 : Car Image
- Class 1 : Street Image
- Class 2 : Bridge Image

keep in mind, any continuous target can be converted into categories through discretization.

- Class 0      price - \$0 - 100K
- Class 1      price \$ 100K - 200K
- Class 2      price \$ 200K &

- Classification algorithm also often produce a probability prediction of belonging to a class:

- Class 0 : 10% probability - Car Image
- Class 1 : 85% probability - Street Image
- Class 2 : 5% probability - Bridge Image

◦ Model reports back prediction of Class 1, image is a street.

- Also note our prediction  $y$  will be a category, meaning we ~~would~~ would not be able to calculate a difference based on  $y$ .

- car Images - street Images does not make sense -
- We need to discover a completely different set of error metrics and performance evaluation.

### 3- Logistic Regression

#### - Theory & Intuition part One

logistic regression works by transforming a linear regression into a classification model through the use of logistic function

$$f(x) = \frac{1}{1+e^{-x}}$$

- logistic function = Sigmoid function

- for now we just are discussing logistic function not logistic regression - its theory and past development.

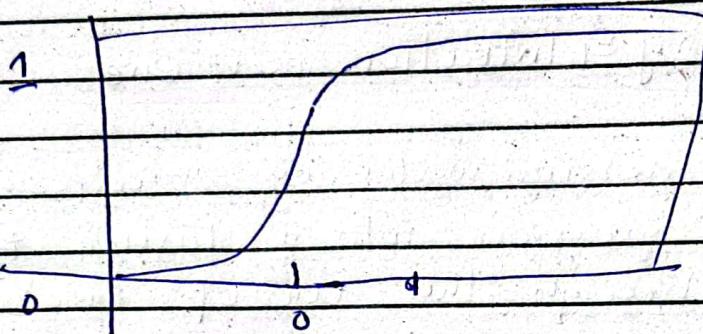
- 1830 - 1850 - Pierre Franoise Verhulst developed logistic function:

$$f(x) = \frac{1}{1+e^{-x}}$$

In 1883 : logistic regression function was independently developed in field chemistry as a model of autocatalysed by ~~by~~ Wilhelm Ostwald.

- 1830-1850 - Developed for the purposes of modelling population growth.

- Why we need for a logistic function versus a logarithmic function?



Note: There is a "family" of logistic Regression.

$$\bullet \operatorname{erf}\left(\frac{\sqrt{2}}{2}x\right) = \bullet \frac{x}{\sqrt{1+x^2}}$$

$$\bullet \tanh(x)$$

$$\bullet \frac{2}{\pi} \operatorname{gd}\left(\frac{\pi}{2}x\right) = \bullet \frac{2}{\pi} \arctan\left(\frac{\pi}{2}x\right)$$

$$\bullet \frac{x}{1+|x|}$$

- Notice the "leveling off" behaviour of the curve.

output range

$$0 - 1$$

- Many natural ~~so~~ world systems have a "carrying capacity" or a natural limiting factor.

• 1940s : Using logistic function for statistical modelling was developed by Joseph Berkson.

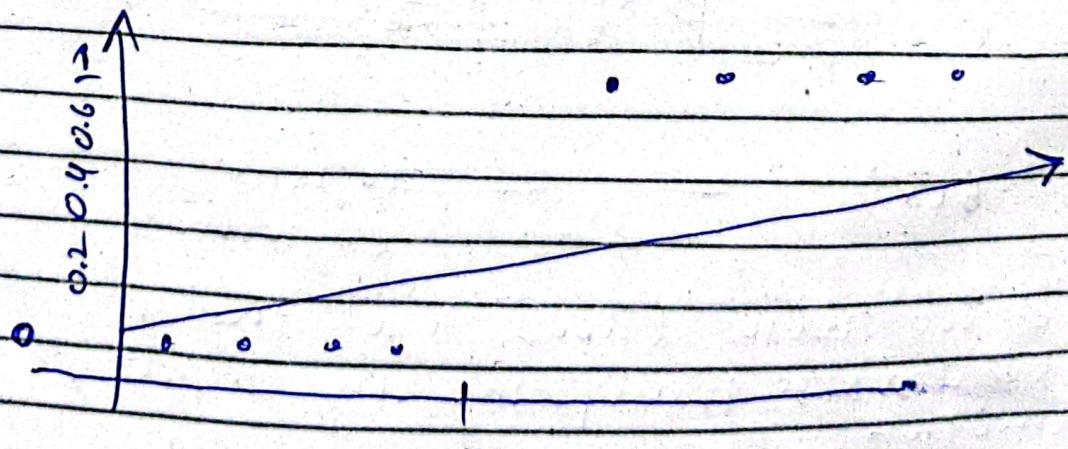
• 1944 : "Application of the logistic function to bio-assay" in the journal of the American Statistical Association

## 004 - Logistic Regression - Theory & Intuition

### part a - Linear to Logistic

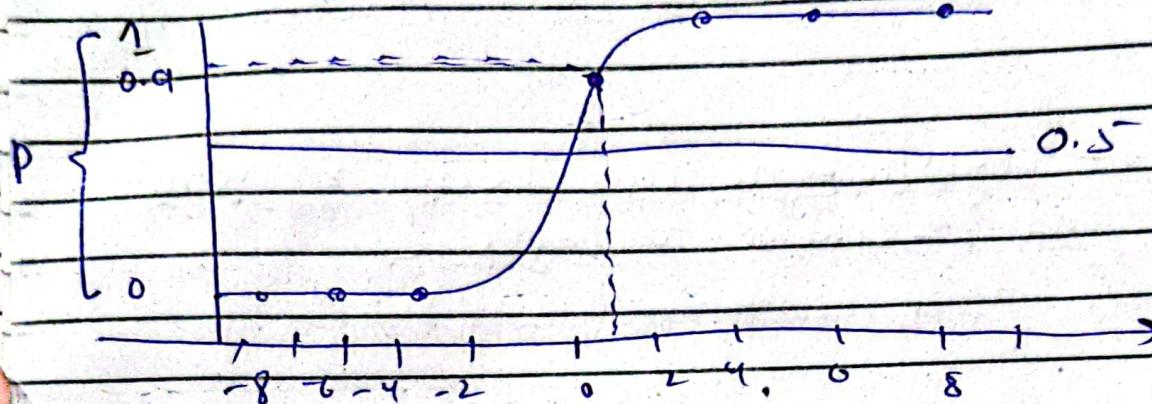
• our dataset

Income	loan paid
-5	0
-4	0
-2	0
-1	0
0	0
2	1
3	1
4	1
5	1
	1



What if we want to predict if a person would be able to

- g) Someone's meant?
- fitting the linear Regression would not work (recalls Anscombe's quartet)
  - o we could make use of a logistic function



the above point will be considered 1 and assumed to have 90% probability based on directionality,

$m(x)$

## 005 Logistic Regression - Theory & Intuition

part two Linear to logistic Math

$$\hat{y} = \beta_0 x_0 + \dots + \beta_n x_n$$

$$\hat{y} = \sum_{i=0}^n \beta_i x_i$$

↳ Linear Regression

$$g(x) = \frac{1}{1+e^{-x}}$$

↳ Logistic regression

and we also know that the logistic regression transforms any input to function be between 0 and 1.

$$\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

$$\hat{y} = \beta_0 + \sum_{i=1}^n \beta_i x_i$$

$$\hat{y} = \frac{1}{1 + e^{-\sum_{i=0}^n \beta_i x_i}}$$

- First we need to understand the terms odds.
- A term you may be familiar with from gambling odds
- In gambling odds are often referred to in the sense of N to 1.
- But where does this actually come from?

$$\begin{array}{c} p \rightarrow \text{probability of happening} \\ 1-p \rightarrow \text{not happening} \end{array}$$

- Solving for log odds

$$\hat{y} = \frac{1}{1 + e^{-\sum_{i=0}^n \beta_i x_i}}$$

$$\hat{y} + \hat{y} e^{-\sum_{i=0}^n \beta_i x_i} = 1$$

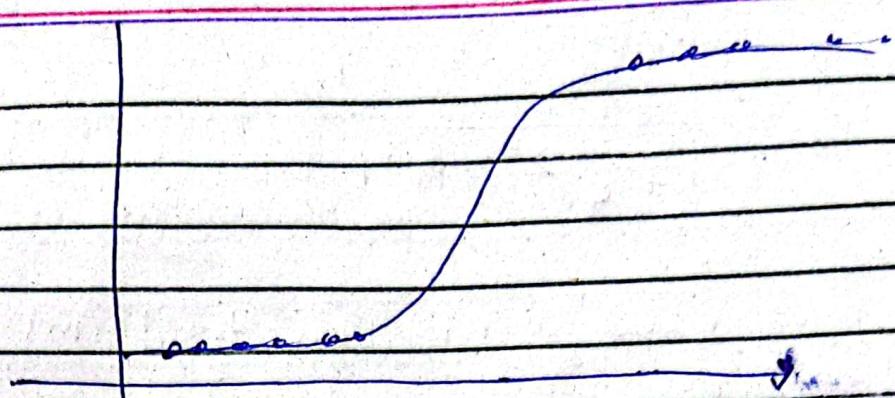
$$e^{\sum_{i=0}^n \beta_i x_i} = 1 - \hat{y}$$

$$\frac{\hat{y}}{1 - \hat{y}} = e^{\sum_{i=0}^n \beta_i x_i}$$

$$\ln\left(\frac{\hat{y}}{1 - \hat{y}}\right) = \sum_{i=0}^n \beta_i x_i$$

What would the function curve look like in terms of log-odds?

$$\ln\left(\frac{0.5}{1-0.5}\right) = 0$$



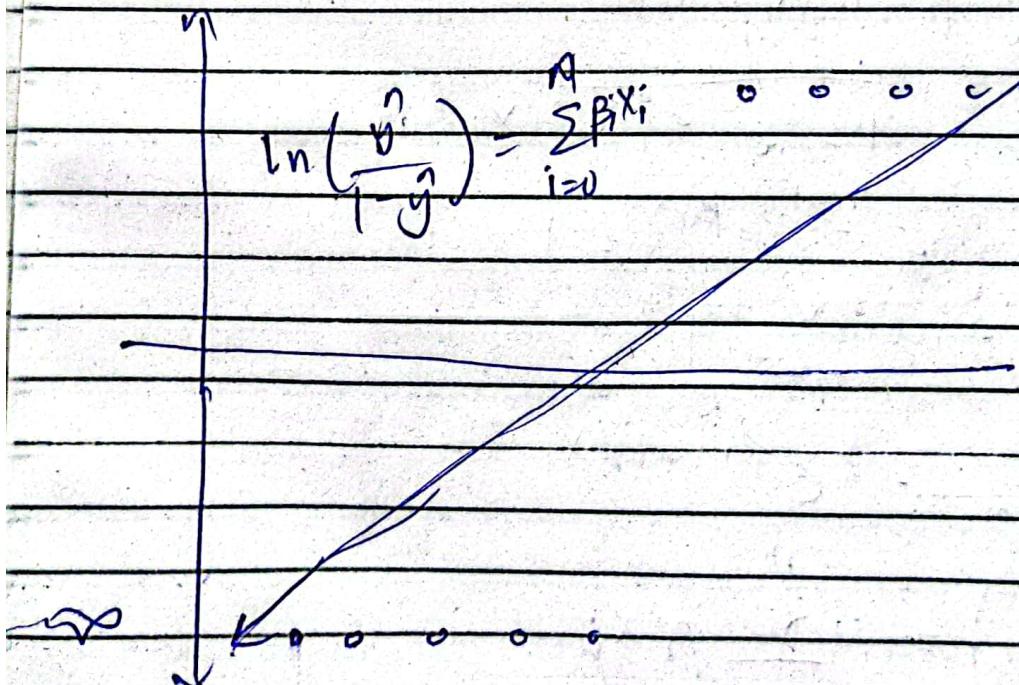
• As  $p$  goes to 1 then log odds become  $\infty$

$$\text{As } \lim_{p \rightarrow 1} \ln\left(\frac{P}{1-P}\right) = \infty$$

• As  $p$  goes to 0 then log odds become  $-\infty$

$$\lim_{p \rightarrow 0} \left( \frac{P}{1-P} \right) = \infty$$

$$\ln\left(\frac{y}{1-y}\right) = \sum_{i=0}^n \beta_i x_i$$



- Since the log odds scale is non-linear a  $\beta$  value can not be directly linked to "one unit increase" as it could in linear regression -

$$\ln\left(\frac{y}{1-y}\right) = \sum_{i=0}^n \beta_i x_i$$

- Sign of coefficient

- Positive  $\beta$  indicate an increase in likelihood of belonging to 1 class with increase in associated X features.

- Negative  $\beta$  indicates an decrease in likelihood of belonging to 1 class with increase in associated X features.

- Magnitude of coefficient

- Harder to directly interpret magnitude of  $\beta$  directly, especially when we could have discrete and continuous X features values.

- We can however begin to use odds ratio, essentially comparing magnitude against each other.

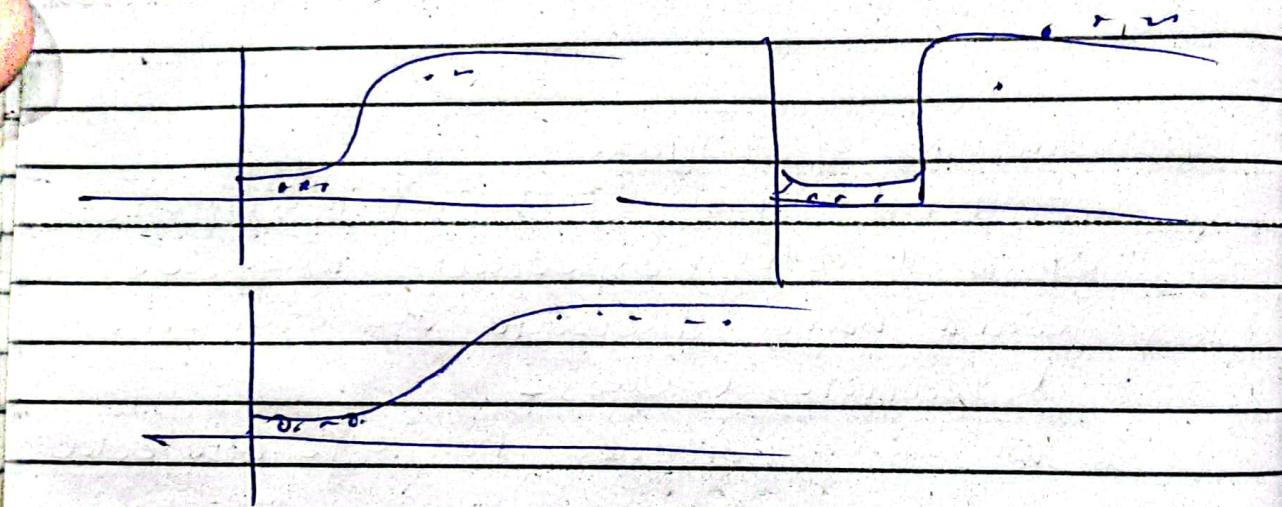
- Comparing magnitudes of coefficients against each other can lead to insight over which features have the strongest effect on prediction output.

- We will discuss the basic of fitting the best curve with maximum likelihood in the next lecture.

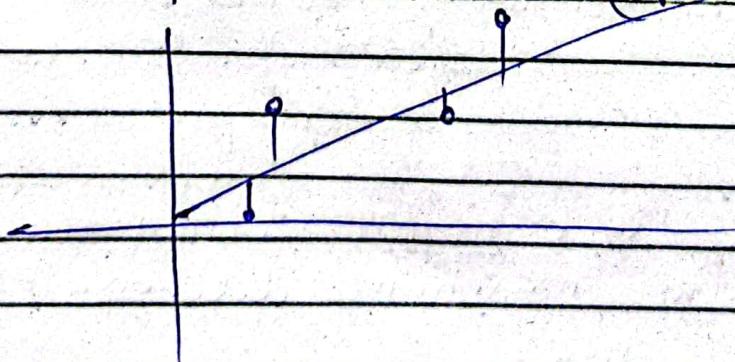
06 Logistic Regression - Theory-Intuition  
part three : finding the best fit

- logistic regression uses maximum likelihood to find the best fitting model.
- This lecture will also ~~cover~~ display cost function and gradient descent that is solved by for by the computer.

Here we see three different  $\beta$  values.  
How do we measure which is best fit?

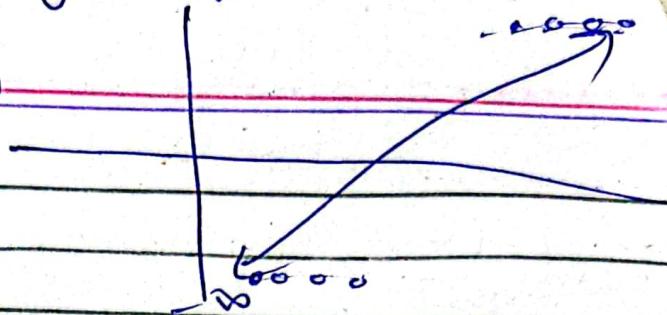


Recall in ~~least square~~ linear regression we seek to minimize the residual sum of ~~errors~~ squares (RSS).



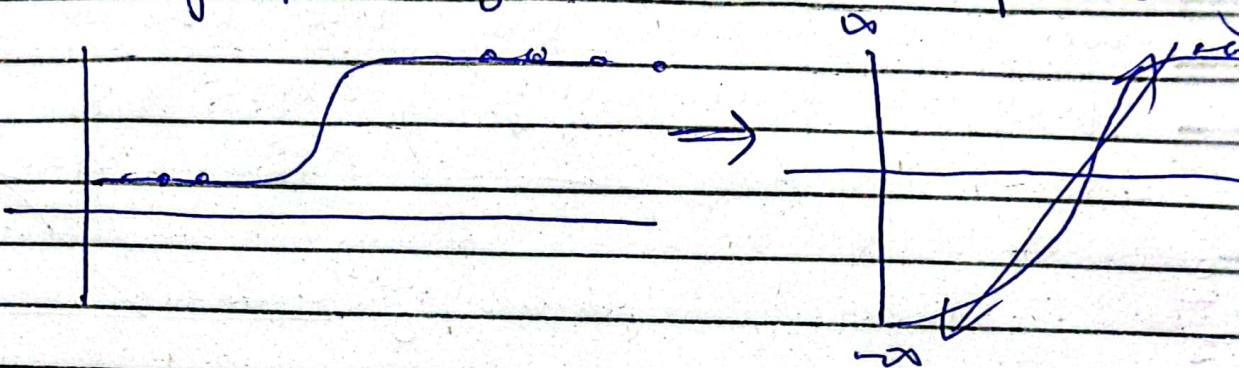
Unfortunately, even in log odds targets are at infinity, making R&S unfeasible.

$$\ln\left(\frac{y}{1-y}\right) = \sum_{i=0}^n p_i x_i$$



Maximum likelihood:-

- The first step is to know that we can go from log odds back to probability



$$\ln\left(\frac{p}{1-p}\right) = \ln(\text{odds})$$

$$\frac{p}{1-p} = e^{\ln(\text{odds})}$$

$$p = (1-p)e^{\ln(\text{odds})}$$

$$p = \frac{e^{\ln(\text{odds})}}{1 + e^{\ln(\text{odds})}} - pe^{\ln(\text{odds})}$$

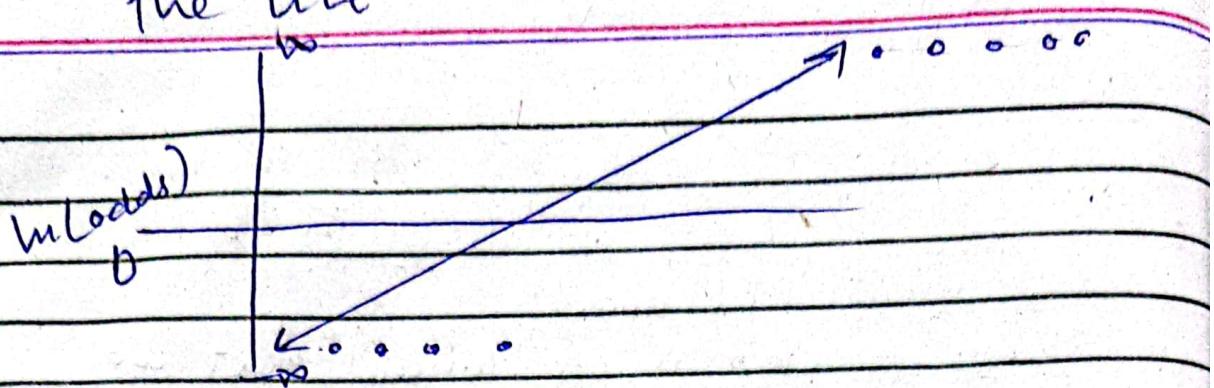
$$p + pe^{\ln(\text{odds})} = e^{\ln(\text{odds})}$$

$$p(1 + e^{\ln(\text{odds})}) = e^{\ln(\text{odds})}$$

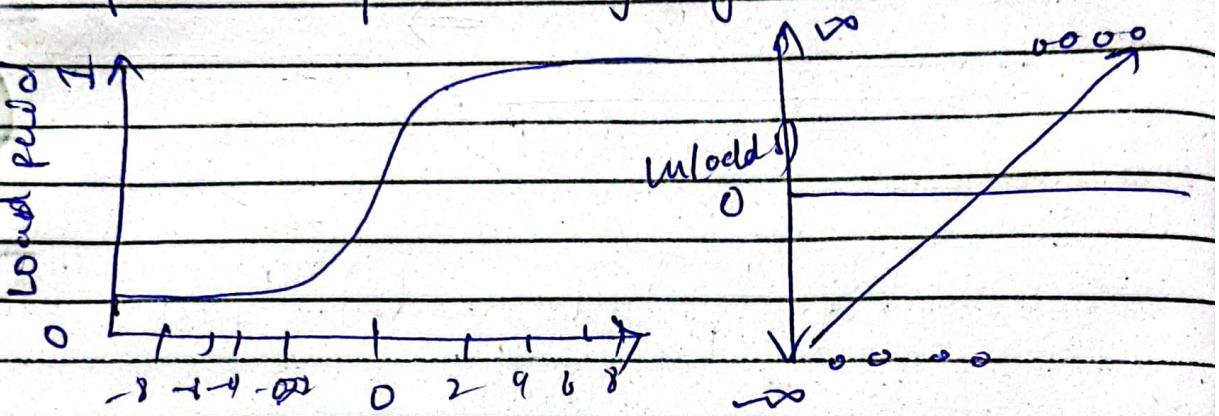
$$p = \frac{e^{\ln(\text{odds})}}{1 + e^{\ln(\text{odds})}}$$

Now we are able to convert  $\ln(\text{odds})$  into probability

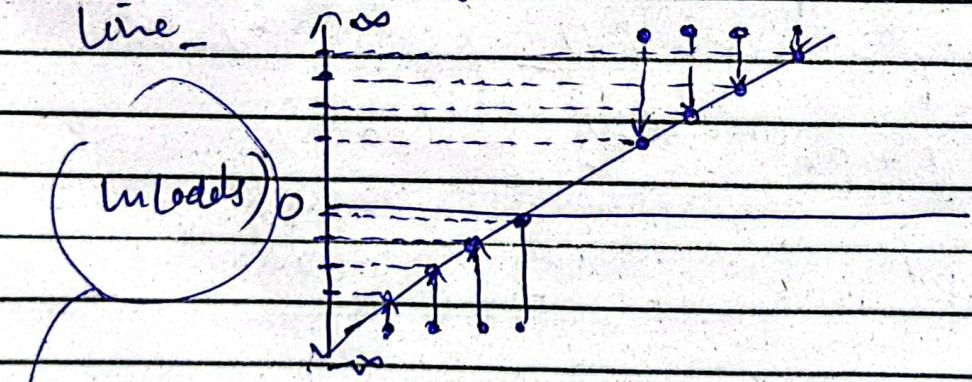
• We choose a line in the log(odd) axis and project the point on the line



• we also know that line has a form in probability y-axis-



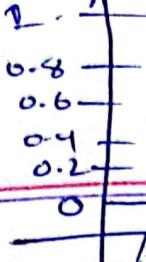
- we chose a line in the log(odd) axis and project the points on the line.



- calculate the logs odds for the projected points on the this line

put that project points

$$f = \frac{e^{ln(Odds)}}{1 + e^{ln(Odds)}}$$



- plot these values as the probability on the Logistic Regression model -
- And we now measure the likelihood of these probabilities -

→ likelihood = product of probabilities of belonging to class 1.

Note

- if it has a perfect fit
- blue points should be all the way at 1
- black should be at 0

$$\text{likelihood} = 0.9 \times 0.8 \times 0.65 \times 0.55 \\ \times (1-0.2) \times (1-0.08) \times (1-0.02)$$

↓  
class 0

$$\text{likelihood} = 0.129$$

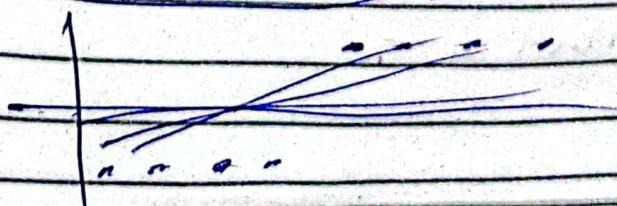
we want to maximize this to  $\phi 1$

Note..

In practical we actually maximize the log of the likelihood  
(e.g.  $\ln(0.9) \times \ln(0.8) \times \dots$ )

Mathematics works better

- We just need some set of coefficients that will maximize these likelihood log



- While we are trying to maximize the likelihood, we still need something to minimize, since the computer's gradient descent methods can only search for minimums -

- In terms of cost function, we seek to minimize the following (logloss):

$$J(x) = -\frac{1}{m} \sum_{j=1}^m y^j \log(\hat{y}^j) + (1-y^j) \log(1-\hat{y}^j)$$

$$J(x) = -\frac{1}{m} \sum_{j=1}^m \left( y^j \log \left( \frac{1}{1+e^{-\sum_{i=0}^n \beta_i x_i^j}} \right) + (1-y^j) \log \left( 1 - \frac{1}{1+e^{-\sum_{i=0}^n \beta_i x_i^j}} \right) \right)$$

- just with linear regression gradient descent can solve this for us!

## 007 Logistic Regression with Scikit-learn

### part one: Exploratory Data Analysis

# code

import numpy as np

import pandas as pd

import seaborn as sns

import matplotlib.pyplot as plt

```
df = pd.read_csv("hearing-test.csv")
```

```
df.head()
```

```
df.describe() -
```

of ['test-result'].value\_counts()

```
out 1 3000  
0 2000
```

sns.countplot(df['test-result'])

sns.boxplot(x="test-result", ~~hue=y="age"~~)

sns.boxplot(x="test-result", y="physical-score")

sns.scatterplot(x='age', y="physical-score")

hue = 'test-result')

will tell us  
about intuition

sns.pairplot(hue = "test-result")

sns.heatmap(df.corr(), annot=True)

sns.scatterplot(x="physical-score", y="test-result")

from mpl\_toolkits.mplot3d import Axes3D

fig = plt.figure()

ax = fig.add\_subplot(111, projection='3d')

ax = scatter(~~df['age']~~, df['physical-score'],  
df['test-result'])

c = 'test-result')

## 008 Logistic Regression with Scikit Learn

### part two : creating and Training a Model

# code

```
X = df.drop("test result", axis=1)  
y = df['test result']
```

```
from sklearn.model_selection import train_test_split  
from sklearn.preprocessing import StandardScaler  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

scaler = StandardScaler()

scaled\_X\_train = scaler.fit\_transform(X\_train)

scaled\_X\_test = scaler.transform(X\_test)

```
from sklearn.linear_model import LogisticRegression
```

```
help(LogisticRegression)
```

```
log_model = LogisticRegression()
```

```
log_model.fit(scaled_X_train, y_train)
```

```
log_model.coef_
```

```
y_pred = log_model.predict(scaled_X_test)
```

```
o predict_proba()
```

# 009 Classification Metrics - Confusion Matrix and Accuracy

## part 1: Confusion Matrix Basis

- You've probably heard of terms such as "false positive" or "false negative". As well as metrics like "accuracy".
- But what does these terms actually mean mathematically?

- Imagine we have developed a test or model to detect presence of a virus infection in a person based on some biological feature
- We could treat this as a Logistic Regression, predicting:
  - 0 - Not infected (Test Negative)
  - 1 - Infected (Test positive)
- It is unlikely our model will perform perfectly. This means there 4 possible outcomes:
  - Infected person test positive
  - Healthy person test negative
  - Infected person test negative
  - (D) Healthy person test positive
- Based on these 4 possibilities, there are many error metrics we can calculate.
- First let start by visualizing these four possibilities as a matrix-

		Actual	
		Infected	Healthy
		True positive	false positive
Infected			
Healthy		False negative	True Negative

100 people - 5 are infected

		Infected	Healthy
predicted	Infected	4	2
Healthy	1	93	

Accuracy :-

How often is the model correct?

$$\text{Acc} = (TP + TN) / \text{Total}$$

$$= \frac{(4+93)}{100} = 97\% \text{ Accuracy}$$

- Is this good value?

The accuracy paradox could be caused by imbalance data label data

\* This is the Accuracy paradox!

- o Any classifier dealing with imbalanced classes has to confront the issue of the accuracy paradox -

o Imbalanced classes will always result in a distorted accuracy reflecting better performance than what is truly warranted.

Φ Imbalance classes are often found in real world datasets:

- o Medical condition can affect small portion of the population -

o Fraud is not common (e.g. Real vs Fraud credit card usage).

• If a class is only a small percentage ( $n\%$ ), then a classifier that always predict the major class will always have an accuracy of  $(1-n)$ .

• In our previous ~~case~~ examples we saw infected were only 5% of the data -

• Allowing the accuracy to be 95%.

- This mean we should rely on accuracy.

## 10 - Classification performance Metrics part 2: Precision and Recall

let's explore three more metrics  
that can help give a clearer

picture of performance:

- Recall (a.k.a sensitivity)
- precision (a.k.a specificity)
- F1-Score

Let begin with Recall

Recall:-

When it is actually is a  
positive case, how often is it  
correct?

$$\frac{(\text{TP})}{\text{Total actual positive}} = \text{Recall}$$

$$\text{Recall} = \frac{4}{5} = 80\%$$

If always positive

$$\text{Recall} = \frac{0}{5} = 0\%$$

$$\frac{\text{TP}}{\text{TP} + \text{FN}} = \text{precision}$$

Now lets explore precision:

precision:-

When prediction is positive; how often it is correct-

$$\frac{\text{TP}}{\text{TP} + \text{FP}} = \text{precision}$$

Total predicted positive

$$\frac{\text{TP}}{\text{TP} + \text{FP}} = \text{precision}$$

$$\frac{4}{6} = 0.666 \%$$

If & always healthy

$$\frac{0}{0+0} = 0$$

- Recall and precision can help illuminate our performance specifically in regards to the relevant or positive cases.

- Depending on the model, there is typically a trade-off between precision and recall, which will explore later on with the ROC curve -

ROC (Receiver operator characteristic curve)

## F1-Score:-

Since precision and recall are related to each other through the numerators (TP), we often also repeat the F1-score, which is harmonic mean of precision and recall.

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

- The harmonic mean (instead of the normal mean) allows the entire harmonic mean to go to zero either precision or recalls ends up being zero.

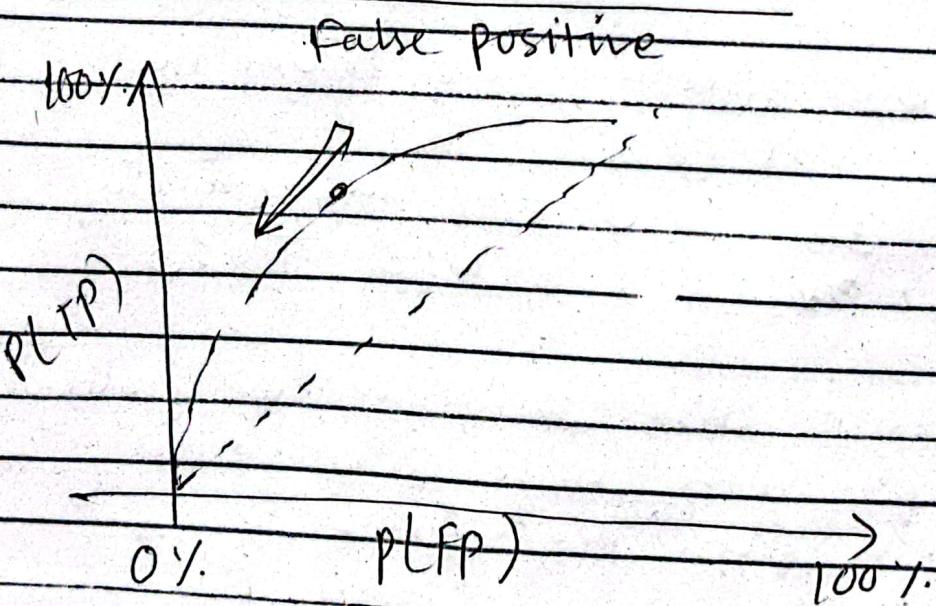
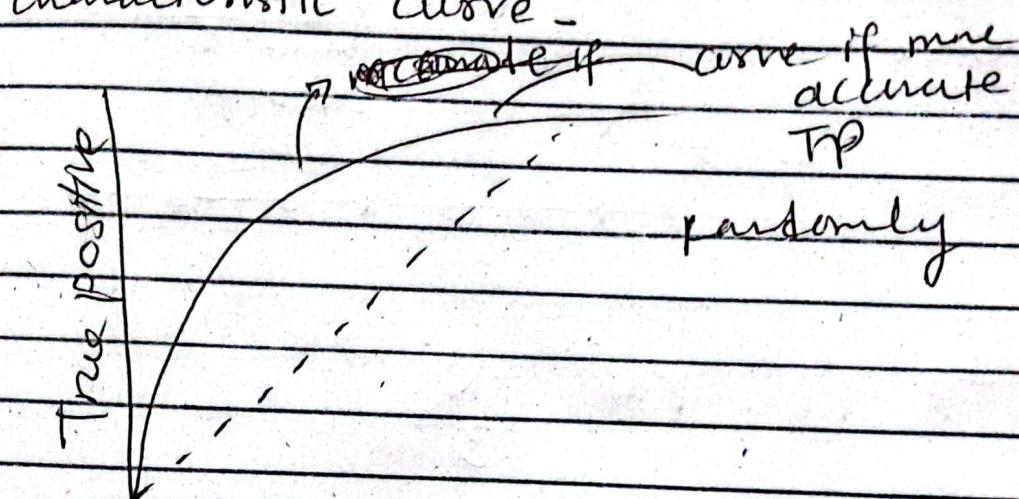
# As a final note there are many more metrics available : can check them online -

## Oct 1 Classification Performance Metrics part 3: ROC curves

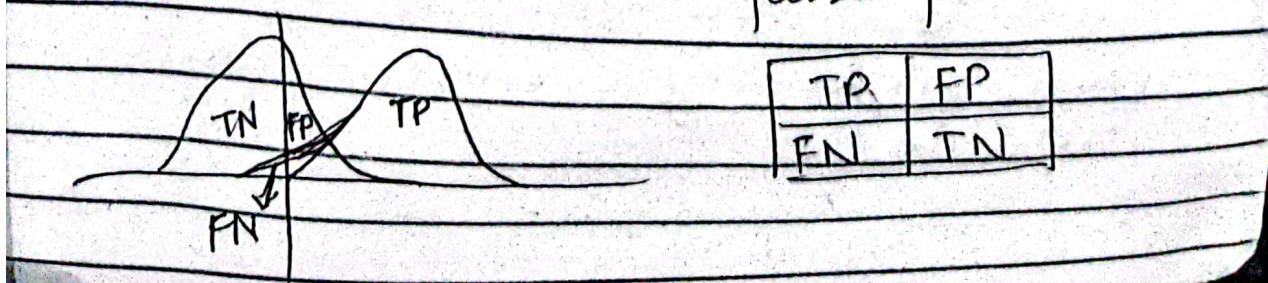
History :-

During WW2 radar was developed now to evaluate radar performance they developed Receiver operator characteristic curve.

Characteristic curve -

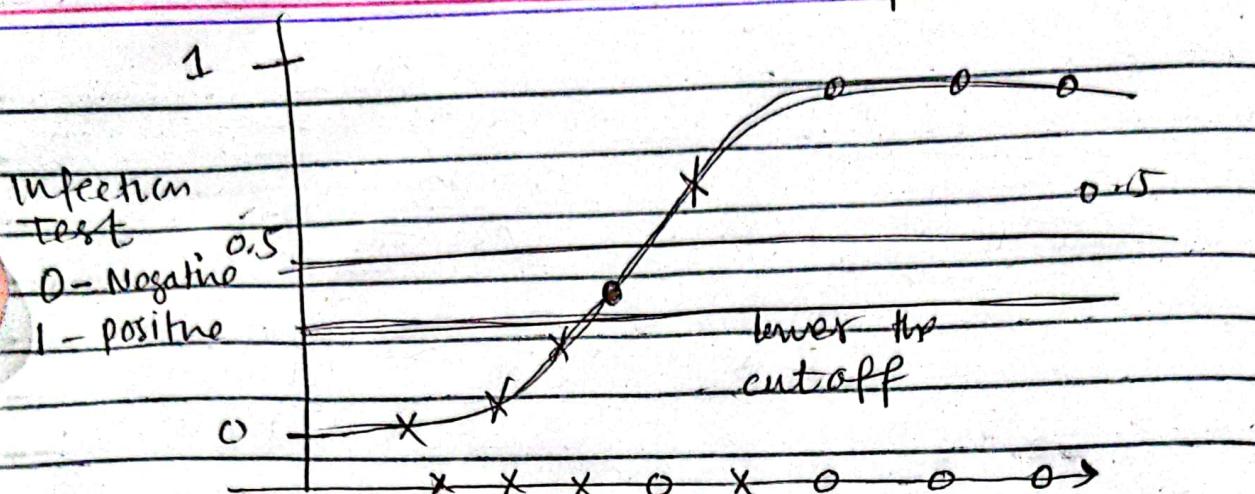


There can be a trade off between true positives and false positives -



Our previous test  
infection

Actual  
X Negative  
O Positive



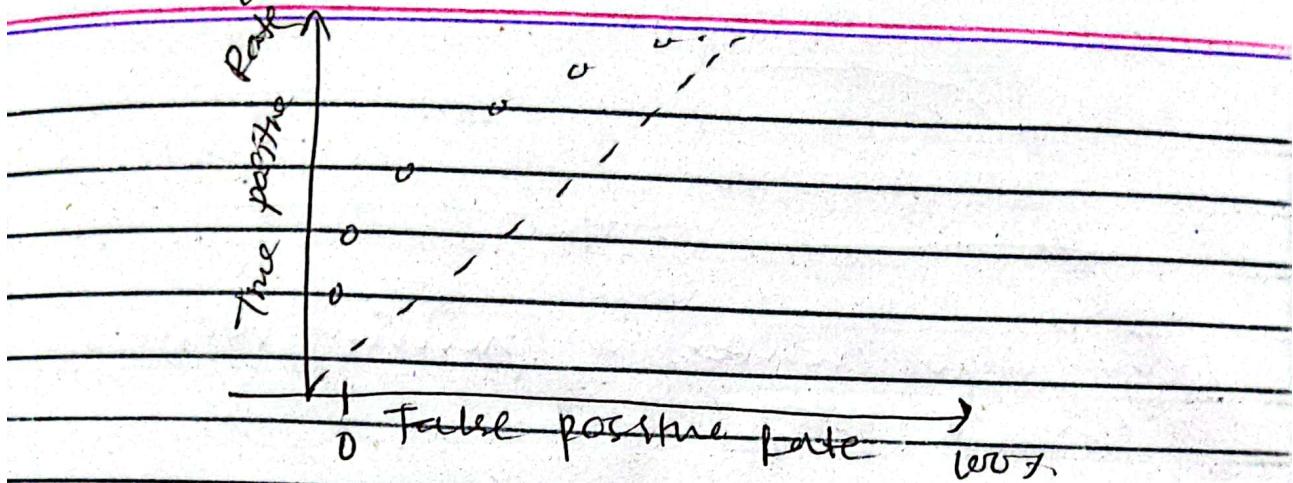
But sometime we lower the  
cutoff

3 O — my understanding  
TP : 3 FP : 1 , FN : 0 TN : 3

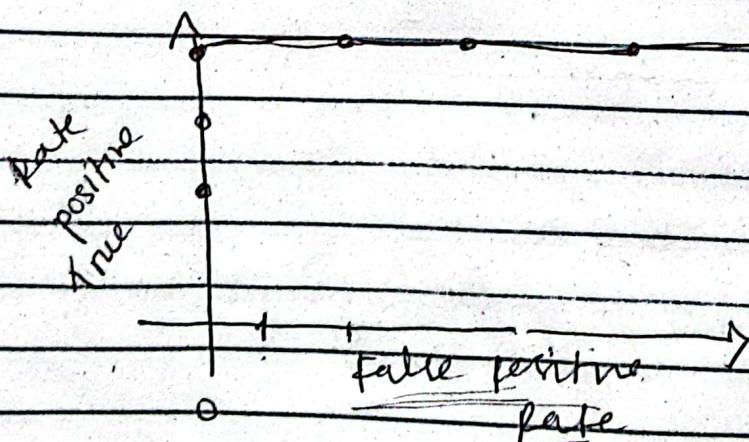
In certain cases/ situation we gladly  
accept more false positive to  
reduce false negatives -

Imagine a dangerous virus test,  
we ~~should~~ would much rather  
produce false positive too and  
later do more stringent examination  
then accidentally release a false  
negative!

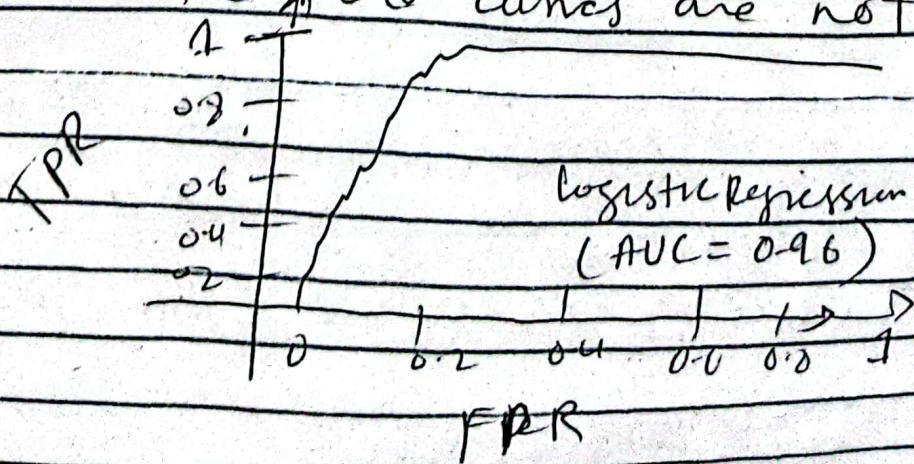
- Chart the True vs. False positives for various cut-offs for the ROC curve.



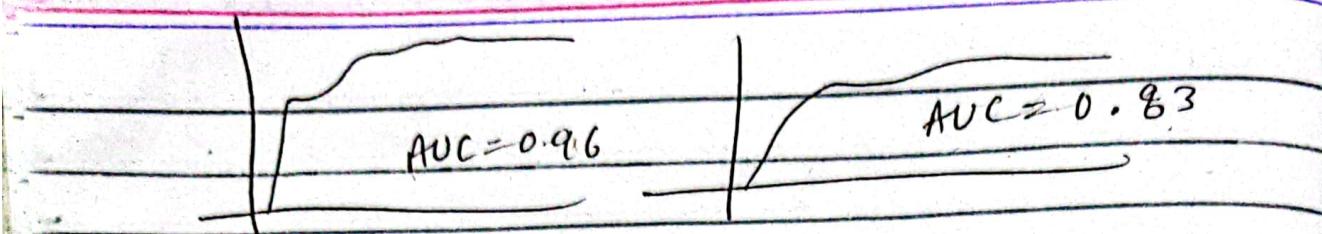
- A perfect model would have a zero FPR.
- Random guessing is the red line.



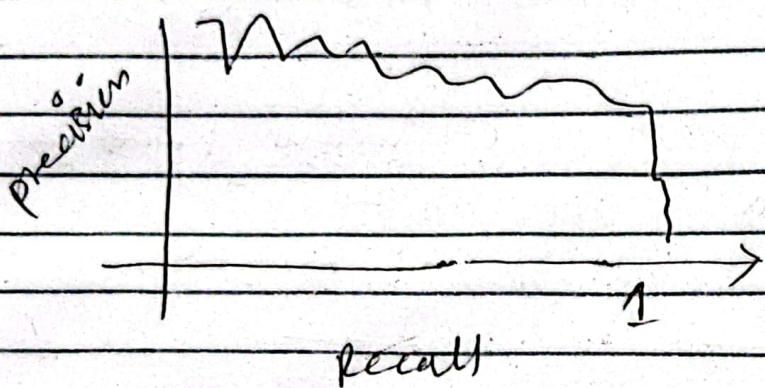
- Realistically with smaller datasets the ROC curves are not as smooth.



- AUC - Area under the curve, allows us to compare ROCs for different models.



can also use a precision versus recall curves:



## 12 Logistic Regression with Scikit-learn part 3 - performance Evaluation.

# code

log\_model.coef\_

df.head()

from sklearn.metrics import accuracy\_score,  
confusion\_matrix,  
classification\_Report

y\_pred = log\_model.predict(X\_Scaled\_test)

accuracy\_score(y\_test, y\_pred)

out 93.0 %

0.93

confusion\_matrix(y-test, y-preds)  
from sklearn.metrics import plot\_confusion\_matrix  
~~plot\_confusion\_matrix(log-model, scaled-X-test, y-test)~~  
normalize = 'all'

len(y-test)  
print  
(classification\_report(y-test, y-preds))  
precision-score()  
recall-score()

From sklearn.metrics import plot\_precision\_recall\_curve, plot\_roc\_curve

plot\_roc\_curve(log-model, scaled-X-test, y-test)  
plot\_precision\_recall(log-model, predict\_proba(scaled-X-test)[0])

log-model, predict\_proba(scaled-X-test)[0]

013

Logistic Regression Multiclass  
problems

part 1 - Data and Model

# Iris Datasets

df = pd.read\_csv("../csv")

df.head()

df.info()

df.describe()

df[["Species"]].value\_counts()

sns.countplot(data=df, x='species')  
sns.scatterplot(x='petal-length', y='petal-width', hue='species')  
~~sns.pairplot(data=df)~~  
sns.heatmap(data.corr(), annot=True)

X = df.drop(['species'], axis=1)  
y = df['species']

from sklearn.model\_selection import train\_test\_split  
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()  
scaler\_X\_test = scaler.transform(X\_train)

Q14)

## Logistic Regression

Multiclass problem

part two training and performance evaluation

from sklearn.model\_selection import GridSearchCV

log\_model = LogisticRegression()  
hyperparameters:-

penalty: { "L1", "L2", "elasticnet", "none" }  
default = 'L2'

The "newton-cg" and "lbfgs" solvers support only 'L2' penalties, "elasticnet" is only supported by "saga" solver. If none (not supported by the liblinear solver), no regularization is applied

L1\_ratio: float if used in elasticnet.

multi\_class: {'auto', 'ovr', 'multinomial'}

|  
One  
versus  
rest  
|  
default "auto"

Solver: {"newton-cg", "lbfgs", "liblinear",  
'sag', 'saga'}  
default = 'lbfgs')

optimization algorithms  
try to minimize cost  
function -

# Code

model = LogisticRegression(solver='saga', multi\_class='ovr', max\_iter=5000)

penalty = ['L1', 'L2', 'elasticnet']

L1\_ratio = np.linspace(0, 10, 10)

param\_grid = [{"penalty":

L1\_ratio: L1\_ratio,  
C: C}])

param\_grid = [{"penalty": penalty,  
L1\_ratio: L1\_ratio,  
C: C}])

grid\_model = GridSearchCV(~~log model~~,  
paramsgrid  
= params\_grid)

~~grid\_model.fit(Scaled\_X\_Scaled, y\_train)~~

from sklearn.metrics import accuracy\_score,  
confusion\_matrix

y\_pred = grid\_model.best\_params\_

{'C': 1.088}

y\_preds = grid\_model.predict(~~Scaled\_X\_test~~)

accuracy\_score(y\_test, y\_preds)

ROC wont work with Multiclass

will need to do this using  
scikit learn code.

11

Code

11