

## Logistic Regression

The most common use of logistic regression models is in binary classification problems.

### Examples of classification

- 1) finance company wants to know whether a customer is default or not.
- 2) Predicting an email is spam or not.
- 3) Categorizing email into promotional, personal and official.

A binary classification problem is a type of supervised learning problem where the goal is to classify input data into one of two possible categories or classes.

### Key features

- ① Two classes : the target variable (output) has exactly two unique values often represented as : → 0 and 1 (e.g Yes/No, Positive/Negative)
  - True or False
  - Any two distinct labels  
(Spam not spam)
- ② Objective : to learn a decision boundary or function that can accurately distinguish between two classes based on features in a dataset.

## Simple Boundary Decision Method in LR

It is the process of using a decision boundary to separate data points into two classes based on their predicted probabilities.

Logistic regression models the probability of a binary outcome  $Y$  (eq 0 or 1) as a function of the input features  $X$ .

Steps:

### ① Probability Prediction:

→ Logistic regression uses the sigmoid func to predict the probability  $P(y = 1|x)$ , which lies b/w 0 and 1:

$$P(y = 1|x) = \frac{1}{1 + e^{-(B_0 + B_1 x)}}$$

Here  $B_0$  is the Intercept

$B_1$  is the coefficient for the feature  $x$ .

### ② Decision Threshold:

→ A threshold  $T$  (commonly  $T=0.5$ ) is chosen to classify the outcome

if  $P(y = 1|x) \geq T$ , classify as 1

if  $P(y = 1|x) < T$ , classify as 0

### 3. Decision Boundary

→ If it is the value of  $x$  where  $P(y=1/x) = T$  for  $T = 0.5$ , this corresponds to solving

$$\beta_0 + \beta_1 x = 0$$

In higher dimensions this becomes a hyperplane

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n = 0$$

### Sigmoid Curve

The sigmoid curve is the graphical representation of the sigmoid function, which maps any real number  $Z$  to a value between 0 and 1. It is used in the logistic regression to predict the probabilities

$$\Rightarrow \text{Sigmoid func} = \sigma(Z) = \frac{1}{1+e^{-Z}}$$

Where

$$Z = \beta_0 + \beta_1 x. \quad (\text{a linear combination of input features and model parameters})$$

#### \* Properties of Sigmoid Curve:

##### ① S shape

↳ Sigmoid curve is S-shaped.

It asymptotically approaches 1 as  $Z \rightarrow +\infty$  and 0 as  $Z \rightarrow -\infty$ .

##### ② Range → O/P of the Sigmoid func always b/w 0 and 1

### 3. Interpretation

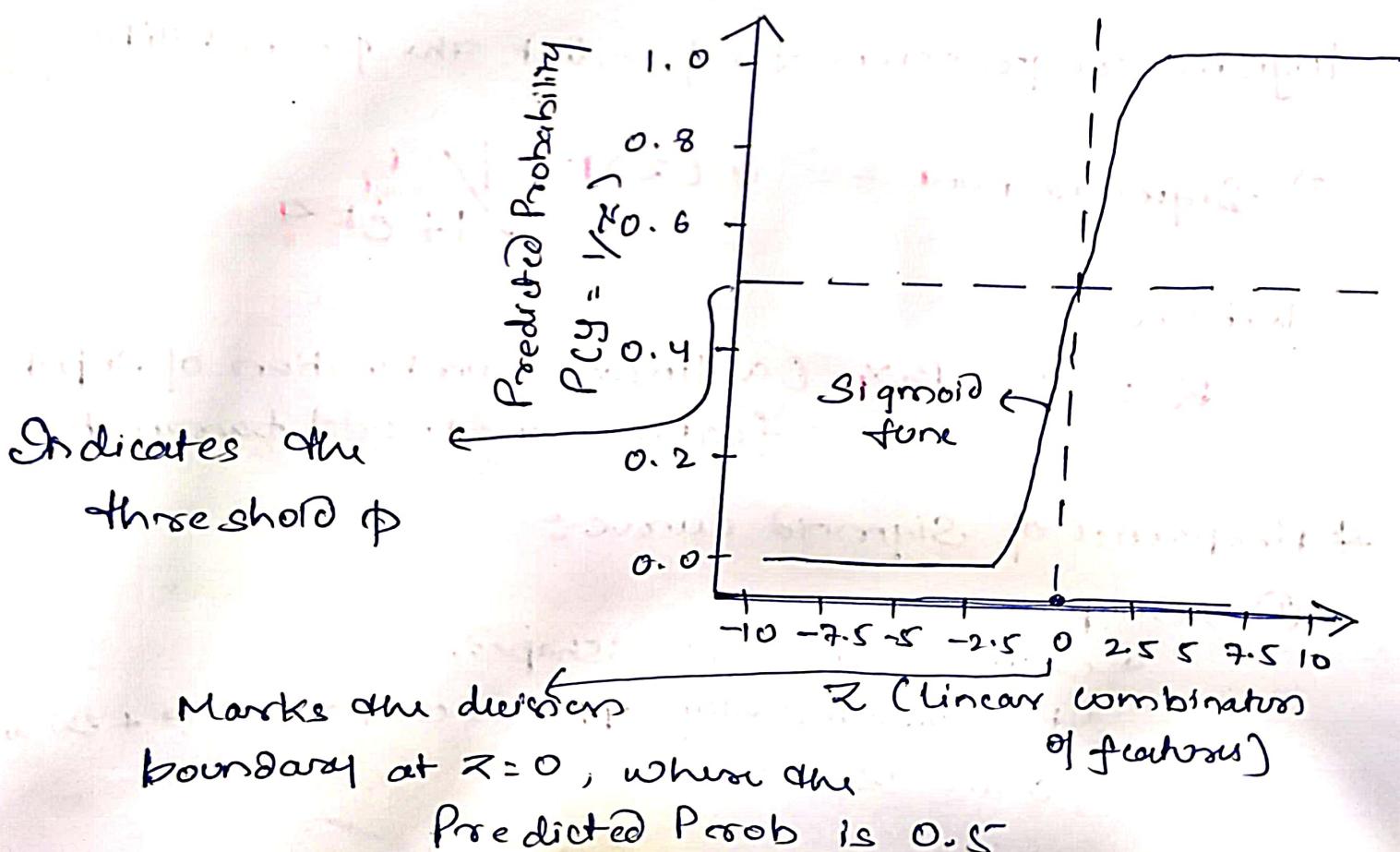
↳ The value  $\sigma(z)$  represents the predicted prob of a positive class

the point where  $\sigma(z) = 0.5$  corresponds to the decision boundary in logistic regression.

#### \* Connection b/w Decision Boundary and Sigmoid Curve

→ the sigmoid curve defines the probability of the Positive class across the feature space.

→ the decision boundary corresponds to the point on the sigmoid curve, where the probability is equal to the threshold  $T = 0.5$ .



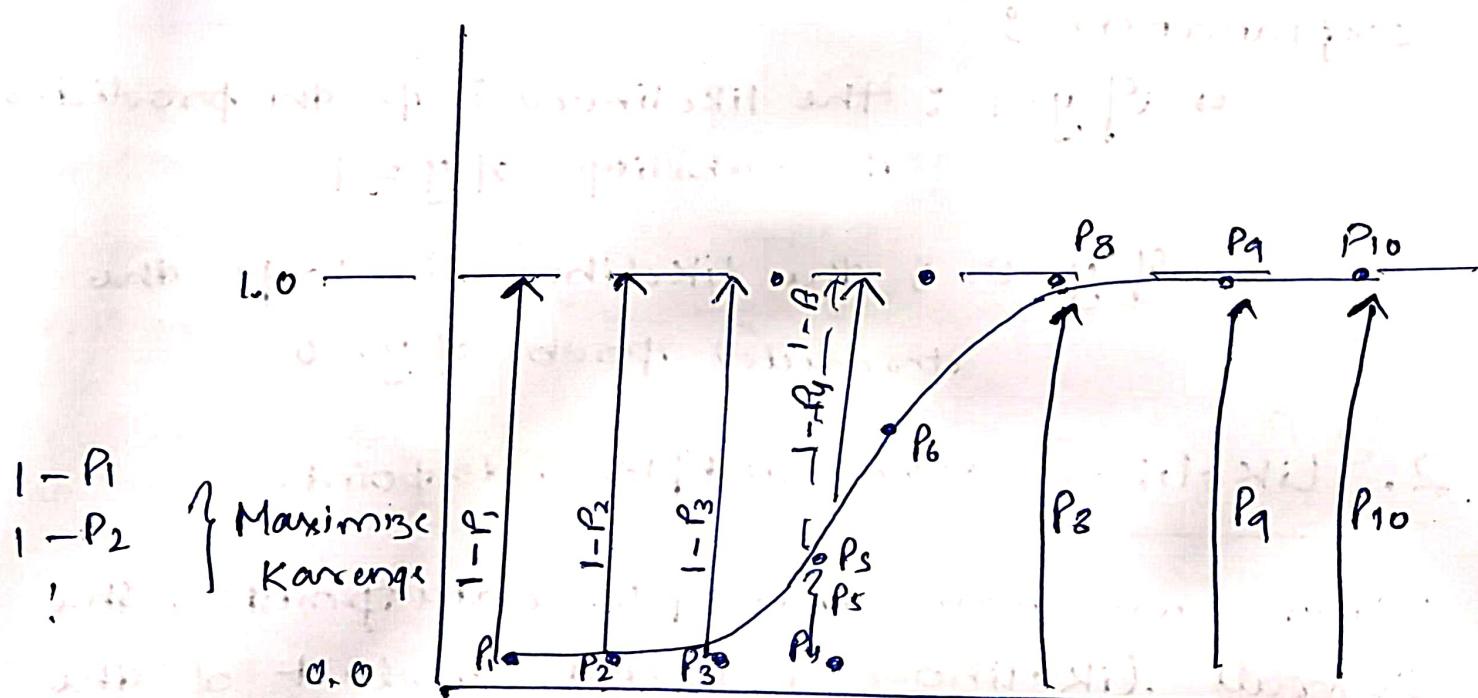
The sigmoid curve even

$$y = P(\text{Diabetes}) = \frac{1}{1 + e^{-(B_0 + B_1 x)}}$$

let's say  
you take  $B_0 = -15$  and  $B_1 = 0.065$  now what  
will be the prob. of diabetes for a patient with  
sugar level 220  $\Rightarrow \frac{1}{1 + e^{-(-15 + 0.065 \times 220)}} = 0.33$

By varying the values of  $B_0$  and  $B_1$ , we get diff  
sigmoid curves, now based on some function that  
we have to minimize or maximize, we get the  
best sigmoid curve.

$\Rightarrow$  Finding the best fit sigmoid Curve



for  
Best  $B_0, B_1$

$$\text{Product} = (1-P_1)(1-P_2)(1-P_3)(1-P_4) \\ (1-P_5)(P_6)(P_7)(P_8)(P_9)(P_{10})$$

Has to be maximize



## What is Likelihood?

Likelihood measures how well the predicted probabilities of a model align with the actual outcomes in the data. It's a core concept in statistical modelling and is particularly useful for evaluating models like logistic regression.

### 1. Likelihood for a single datapoint

↳ For a binary outcome  $y$ , the likelihood depends on whether  $y=1$  or  $y=0$ .

formula :  $L = p^y \cdot (1-p)^{1-y}$

Explanation :

↳ If  $y=1$  : the likelihood is  $p$ , the predicted probability of  $y=1$

↳ If  $y=0$  : the likelihood is  $1-p$ , the predicted prob of  $y=0$

### 2. Likelihood for multiple datapoints

When we have multiple datapoints, the overall likelihood  $L$  is the product of the likelihoods for all individual data points.

$$L = \prod_{i=1}^n (P_i^{y_i} \cdot (1 - P_i)^{1-y_i})$$

Example Table:

Datapoint	Predicted Prob (P)	Actual outcome y	Likelihood contribution
1	0.3	0	$1 - 0.3$
2	0.4	0	$1 - 0.4$
3	0.7	1	$0.7$
4	0.9	1	$0.9$

$$L = (1 - 0.3)(1 - 0.4)(0.7)(0.9) = 0.2268$$

### 3. Log-Likelihood

Computing likelihood directly for many datapoints can result in very small numbers (due to multiplying probabilities), leading to numerical instability, to simplify we use the natural algorithm:

$$l = \sum_{i=1}^n (y_i \cdot \log(P_i) + (1-y_i) \cdot \log(1-P_i))$$

Log likelihood converts the product into sum making it easier to optimize.

## Why is likelihood Important?

- ① Model Evaluation : Likelihood Quantifies how well the model's predicted probabilities align with observed data.
- ② Parameter evaluation : In logistic regression we maximize the likelihood (or equivalently the log-likelihood) to find the best fitting model parameters.

## # What are $B_0$ and $B_1$ ?

these are the coefficients in logistic regression:  
→  $B_0$  is the intercept (shifts the curve up or down)  
→  $B_1$  is the slope (controls how steep the curve is)

We want to find the best values for  $B_0$  and  $B_1$ , so the curve fits the data perfectly, to do this we maximize a likelihood func.

The likelihood func measures how well the logistic curve predicts the actual data. It looks like this

$$L = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}$$

Here  $p_i$  is the predicted probability.

$y_i$  is the actual outcome (1 or 0)

## Logistic Regression Equation

$$P(\text{diabetes}) = \frac{1}{1 + e^{-(B_0 + B_1 x)}}$$

where  
 $B_0 = -13.5$   
 $B_1 = 0.06$

Now we can plugin  $x$  and compute the Prob of a person having diabetes.

## Linearising the Sigmoid Equation

$$P = \frac{1}{1 + e^{-(B_0 + B_1 x)}} \quad \& \quad 1 - P = \frac{e^{-(B_0 + B_1 x)}}{1 + e^{-(B_0 + B_1 x)}}$$

$$\frac{P}{1 - P} = \frac{1}{1 + e^{-(B_0 + B_1 x)}} \stackrel{?}{=} \frac{e^{-(B_0 + B_1 x)}}{1 + e^{-(B_0 + B_1 x)}}$$

$$\Rightarrow \frac{P}{1 - P} = e^{(B_0 + B_1 x)} \quad \# \text{ taking log on both sides}$$

$$\Rightarrow \ln\left(\frac{P}{1 - P}\right) = B_0 + B_1 x$$

Here  $\frac{P}{1 - P}$  is called as odds

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

If  $\frac{P}{1 - P}$  is 4  $\Rightarrow P(\text{Diabetes}) = 4 \times P(\text{No Diabetes})$   
 $\therefore P \rightarrow \text{Prob of diabetes}$

## Relationship b/w ODDS and Probability

Suppose  $P = 0.5 \Rightarrow 1 - P = 0.5 \Rightarrow \text{ODDS} = \frac{P}{1-P} = \frac{0.5}{0.5} = 1$

If  $\text{ODDS} = 2 \Rightarrow \frac{P}{1-P} = 2 \Rightarrow P = 2 - 2P \Rightarrow 3P = 2 \Rightarrow P = \frac{2}{3} = 0.67$

# In logistic regression, ODDS and log ODDS are used to model probabilities in a way that ensures the predicted values remains b/w 0 and 1.

The ODDS measure the likelihood of an event happening compared to it not happening.

Mathematically ODDS are calculated as

$$\text{ODDS} = \frac{P}{1-P}$$

Where  
 $P \Rightarrow$  Prob of event happening  
 $1-P \Rightarrow$  Prob of event not happening,

e.g if  $P = 0.75$  then ODDS are  $\Rightarrow \frac{0.75}{1-0.75} = 3$  meaning the event is 3 times as likely to happen as not

## Log-odds (Logit)

↳ Log odds is the natural logarithm of the odds

$$\text{Log-odds (Logit)} = \ln(P/(1-P))$$

The log-odds transforms the probability range  $[0, 1]$  into a range of  $(-\infty, \infty)$ , this transformation is important because logistic regression models a linear relationship b/w the input variables and the log-odds.

## Why to use Odds and Log-odds in Logistic Regression

① Handling Non-linearity: Logistic regression cannot directly model probabilities because they are bounded b/w 0 and 1, by modelling the log odds the regression becomes a linear problem

⇒ Probabilities are bounded b/w 0 and 1, if we try to directly model probabilities P as a linear func of the Input X.

$$P = B_0 + B_1 X$$

this approach has two major issues:

① Bound violation: The linear equation can produce values outside the range  $[0, 1]$ , which are not valid probabilities.

2. Non linear relationship : The change in Probability for a one unit increase in  $x$  is not consistent across the range of  $P$  for eg Increasing  $x$  might lead to a large change in  $P$  when  $P$  is near 0.5 but only a small change when  $P$  is near 0 or 1.

## # The Role of Log-odds

To overcome these issues , logistic regression models the log-odds instead of directly modeling probabilities.

$$\text{Log-odds} = \ln(P/(1-P))$$

this transformation maps the probability range  $[0, 1]$  to the log-odds range  $(-\infty, \infty)$

unlike probabilities log-odds are not bounded and can take any real value

Modelling the log-odds as a linear func of  $x$

$$\ln(P/(1-P)) = B_0 + B_1 x$$

this even represents a linear relationship b/w log odds and  $x$  , which can be solved using standard regression techniques.

the coefficients  $B_0$  and  $B_1$  can now be interpreted as changes in log-odds per unit change in  $X$ .

## # Mapping Back to Probabilities

To predict the probabilities from the linear model of log-odds, we use the inverse transformation (the sigmoid func)

$$P = \frac{1}{1 + e^{-(B_0 + B_1 X)}}$$

the sigmoid func ensures the predicted values are always b/w 0 and 1

This effectively captures the non linear relationship b/w  $X$  and  $P$ .



## Multivariate Logistic Regression

### (Model Building)

Carry used in case of univariate logistic regression

$$P = \frac{1}{1 + e^{-(B_0 + B_1 X)}} \Rightarrow \text{the above eqn has only one feature variable } X, \text{ for which the coefficient is } B_1$$

If we have multiple features, say  $n$ , you can simply extend the above eqn with  $n$  features variables

$$P = \frac{1}{1 + e^{-(B_0 + B_1 X_1 + B_2 X_2 + B_3 X_3 + \dots + B_n X_n)}}$$

# Multivariate Logistic Regression Telecom churn example:

## Problem statement

- ↳ We will be looking at the telecom churn prediction example → we will be using features to predict whether a particular customer will switch to another telecom provider or not (ie churn or not)
- Here 'Churn' is a binary variable : 1 means the customer has churned and 0 means the customer has not churned.

Datasets that we have : ① churn dataset  
② Internet\_Data dataset  
③ Customer\_Data dataset

\* Dataframe 1 : churn\_data columns

cust\_id, tenure, phoneservice, contract, PaperlessBilling, MonthlyCharges, TotalCharges, churns columns.

\* Dataframe 2 : customer data

cust\_id, gender, Senior Citizen, Partner, Dependents

\* Dataframe 3 : Internet\_Data

cust\_id, multiple lines, internet service, online security, online backup, Device Protection, Tech Support, streaming TV, streaming movies

the columns are same in all the Dataframes, hence we can merge the data into single Dataframe.

## Let's see the Data Cleaning & Preparation

```
# suppress warnings
```

```
> import warnings
```

```
> warnings.filterwarnings('ignore')
```

```
# importing Pandas and numpy
```

```
> import pandas as pd, numpy as np
```

```
# Importing all datasets
```

```
> churn_data = pd.read_csv("churn_data.csv")
```

```
> customers_data = pd.read_csv("customers_data.csv")
```

```
> internet_data = pd.read_csv("Internet_data.csv")
```

```
# combining all data files into one consolidated df
```

```
> df_1 = pd.merge(churn_data, customers_data,  
    how='inner', on='customerID')
```

```
# final dataframe with all the predictor vars
```

```
> telecom = pd.merge(df_1, internet_data, how='inner'  
    on='customerID')
```

```
> telecom.head()
```

```
> telecom.shape
```

```
6 (7043, 21)
```

```
# Statistical Info  
> telecom.describe()  
  
# Let's see the type of each col  
> telecom.info()  
↳ No Null values  
  
# Step 3 : Data Preparation  
  
# Converting some binary variables (Yes/No) to (1/0)  
  
# List of variables to map  
> varlist = ['PhoneService', 'PaperlessBilling', 'Churn',  
             'Partner', 'Dependents']  
  
# Defining the map func  
> def binary_map(x):  
    return x.map({'Yes': 1, 'No': 0})  
  
> telecom[varlist] = telecom[varlist].apply(binary_map)  
  
# Creating a dummy variable for some of the  
categorical variables and dropping the first one  
> dummy1 = pd.get_dummies(telecom[['contract',  
                                    'PaymentMethod', 'gender', 'InternetService'],  
                                    drop_first=True])
```

```
# Adding the results to the master Dataframe  
telecom = pd.concat([telecom, dummy1],  
axis=1)  
# Now we got the telecom data.  
# Now Doing the same for the remaining categorical variables  
  
# Creating dummy variables for the variable 'Multiple Lines'  
  
> ml = pd.get_dummies(telecom['MultipleLines'],  
prefix='MultipleLines')  
# Dropping 'MultipleLines_No phone service' column  
> mli = ml.drop(['MultipleLines_No phone service'], 1)  
# Adding the results to master Dataframe  
  
> telecom = pd.concat([telecom, mli], axis=1)  
# Creating ... for "OnlineSecurity" col  
  
> OS = pd.get_dummies(telecom['Online Security'],  
prefix='OnlineSecurity')  
> OSI = OS.drop(['Online Security-No internet service'], 1)  
> telecom = pd.concat([telecom, OSI], axis=1)  
# Creating dummy for the variable 'Online Backup'  
> Ob = pd.get_dummies(telecom['Online Backup'],  
prefix='Online Backup')  
> OBI = Ob.drop(['Online Backup-No internet service'], 1)  
> telecom = pd.concat([telecom, OBI], axis=1)
```

# Creating dummy variable for the variable  
'DeviceProtection'

```
> dp = pd.get_dummies(telecom['DeviceProtection'],  
                      Prefix = 'DeviceProtection')  
> dpl = dp.drop(['DeviceProtection_No internet service'], 1)  
> telecom = pd.concat([telecom, dpl], axis=1)
```

# Creating dummy variable for the variable  
'TechSupport'

```
> fs = pd.get_dummies(telecom['TechSupport'],  
                      Prefix = 'TechSupport')  
> fs1 = fs.drop(['TechSupport_No internet service'], 1)  
> telecom = pd.concat([telecom, fs1], axis=1)
```

# Creating dummy variables for 'StreamingTV'

```
> st = pd.get_dummies(telecom['StreamingTV'],  
                      Prefix = 'StreamingTV')  
> st1 = st.drop(['StreamingTV_No internet service'], 1)  
> telecom = pd.concat([telecom, st1], axis=1)
```

# Creating dummy variables for the variable  
'StreamingMovies'

```
> sm = pd.get_dummies(telecom['StreamingMovies'],  
                      Prefix = 'StreamingMovies')  
> sm1 = sm.drop(['StreamingMovies_No internet service'], 1)  
> telecom = pd.concat([telecom, sm1], axis=1)
```

# Dropping the repeated variables

```
> telecom = telecom.drop(['Contract', 'PaymentMethod'], 1)
```

# the variable was imported as string we need to convert it to float

```
> telecom['TotalCharges'] = telecom['TotalCharges'].convert_objects(convert_numeric=True)
```

(#) So the process of dummy variable creation was quite familiar, except this time, you manually dropped one of the columns for many dummy variables for eg the column 'MultipleLines' we dropped the level 'MultipleLines\_No phone service' manually, instead of simply using the drop-first=True, the reason we did this is that If you check the variable multiple lines

```
telecom['MultipleLines'].as_type('category').value_counts()
```

No	3390
Yes	2991
No Phone Service	682

Now out of these, It is best to drop 'No Phone Service' since it isn't of any use because it is anyway being indicated by the variable 'Phone Service' already present in dataframe

# Checking for outliers in continuous variables

> num\_telecom = telecom[['Tenure', 'MonthlyCharges', 'SeniorCitizen', 'TotalCharges']]

# checking outliers at 25, 50, 75, 90, 95, 99

> num\_telecom.describe(Percentiles=[.25, .5, .75, .9, .95, .99])

↳ from the distribution shown above, you can see that there are no outliers in the Data, the nos are gradually increasing.

# Checking the missing values and Inputting them

# Adding up the missing values (columnwise)

> telecom.isnull().sum()

↳ total charges → 11 values are missing

# rest all the values are present in the dataset

It means that 11 out of 7043 records have missing values ie 0.1%. → It is best to remove these observations from analysis.

# Checking the %age of Missing Values

> round(100 \* (telecom.isnull().sum() / len(telecom.index))), 2)

# removing NAN total charges rows

> telecom = telecom[~np.isnan(telecom['TotalCharges'])]

```
# checking the % of missing values after removing  
the missing values.
```

```
> round(100 * (telecom.isnull().sum() / len(  
telecom.index)), 2)
```

```
# now we don't have missing values
```

### # Step 4 : Test-Train Split

```
> from sklearn.model_selection import train_test_split
```

```
# Putting feature variables to X
```

```
> X = telecom.drop(['churn', 'customerID'], axis=1)
```

```
# Putting response variable to Y
```

```
> Y = telecom['churn']
```

```
# splitting the data onto train and test
```

```
> X_train, X_test, Y_train, Y_test = train_test_split  
(X, Y, train_size=0.7, test_size=0.3,  
random_state=32)
```

### # Step 5 : feature Scaling

```
> from sklearn.preprocessing import StandardScaler
```

```
> Scaler = StandardScaler()
```

```
> X_train[['tenure', 'MonthlyCharges', 'TotalCharges']] = Scaler.fit_transform
```

```
(X_train[['tenure', 'MonthlyCharges', 'TotalCharges']] )
```

# Step 6: Looking at correlations

```
> import matplotlib.pyplot as plt  
> import Seaborn as sns  
%matplotlib inline
```

# Let's see the correlation matrix

```
> plt.figure(figsize=(20,10))  
> sns.heatmap(telecom.corr(), annot=True)  
> plt.show()
```

# Dropping highly correlated dummy variables

(Here in this first little, business context also helps)

```
> X_test = X_test.drop(['MultipleLines_No', 'OnlineSecurity  
-No', 'OnlineBackup_No', 'DeviceProtection_No',  
'StreamingTV_No', 'StreamingMovie_No'], 1)  
  
> X_train = X_train.drop(['MultipleLines_No', 'Online  
Security-No', 'OnlineBackup-No',  
'DeviceProtection-No',  
'StreamingMovie-No'], 1)
```

# After dropping highly correlated variables, we  
will check the correlations

```
> plt.figure(figsize=(20,10))  
> sns.heatmap(X_train.corr(), annot=True)  
> plt.show()
```

# Running the first training Model

> import statsmodels.api as sm

# Logistic Regression Model

> logml = sm.GLM(y\_train, sm.add\_constant(X\_train),  
> family=sm.families.Binomial())

# this is for the whole binomial  
thing

> logml.fit().summary()

↳ we can see the Pvalues

# the lower the Pvalue, the higher is its  
significance.

### \* Feature elimination using RFE

We built the first model in the previous segment, based on the summary statistics we concluded that many of the features might be Insignificant hence we need to do some feature elimination.

# Step 8 : feature selection using RFE

> from sklearn.linear\_model import LogisticRegression

> logreg = LogisticRegression()

```
> from sklearn.feature_selection import RFE  
> rfe = RFE(LogReg, 15)  
    ↳ 15 top features to select  
> rfe = rfe.fit(X_train, y_train)  
> rfe.support_  
    ↳ It will return true/false if feature  
        is there in the top 15.  
> list(zip(X_train.columns, rfe.support_,  
          rfe.ranking_))  
    ↳ Colnames support rank  
  
## Let's see the included cols  
> col = X_train.columns[rfe.support_]  
> print(col)  
  
## Let's see the excluded cols  
> X_train.columns[~rfe.support_]  
    ↳ Negation operator  
    ↳ Cols which are excluded  
  
## Assessing the models with StatsModels  
X_train_sm = sm.add_constant(X_train[col])  
Logm2 = sm.GLM(y_train, X_train_sm,  
                 family=sm.families.Binomial())  
res = Logm2.fit()  
res.summary()
```

- # All of them have a good p value.
- # Getting the predicted values on the train set
  - >  $y\_train\_pred = \text{reg.predict}(X\_train\_sm)$
  - >  $y\_train\_pred = y\_train\_pred\_values.reshape(-1)$   
↳ reshaping it
- # Creating a Dataframe with the actual churn flag and predicted probabilities
  - >  $y\_train\_pred\_final = \text{Pd.DataFrame}(\{'Churn': y\_train\_values, 'Churn\_prob': y\_train\_pred\})$
  - >  $y\_train\_pred\_final['cust\_id'] = y\_train\_index$
  - >  $y\_train\_pred\_final.head()$
- # Churn Churn\_Prob cust\_id
 

Churn	Churn_Prob	cust_id
0	0.2251	329
0	0.2748	5790
1	0.6921	6498
- # creating a new column 'Predicted' with 1 if the Churn\_Prob > 0.5 else 0
  - >  $y\_train\_pred\_final['predicted'] = y\_train\_pred\_final.Churn\_Prob.map(\lambda x: 1 \text{ if } x > 0.5 \text{ else } 0)$
  - >  $y\_train\_pred\_final.head()$



## Confusion Matrix and Accuracy

We chose a cutoff of 0.5 in order to classify the customers into 'churn' and 'not churn'. Now since we are classifying the customers into two classes, we will obviously have some errors.

The classes of errors that would be there are:

- Churn customers being (incorrectly) classified as Non-churn
- Non churn customers being (incorrectly) classified as 'churn'
- ⇒ To capture these errors and to evaluate how well the model is → we will use confusion matrix

A confusion Matrix would look like

Actual	Predicted	
	No (Non-churn)	Yes (churn)
No (Non-churn)	1406	143
Yes (churn)	263	296

The table shows the comparison of predicted and actual labels; the actual labels are along the vertical axis while the predicted labels are along the horizontal axis.

thus the second row and first column (263) is the nos of customers who have actually churned but the model has predicted them as non churn

Similarly the cell at the second row, second col (298) is the number of customers who are actually churned and also predicted as churned.

# The simplest model evaluation metric for classification model is "accuracy" → It is the percentage of correctly predicted labels

so what would be the correctly predicted labels be

- ↳ Churn customers being actually identified as churn
- ↳ Non churn customers being actually identified as non churn.

$$\text{Accuracy} = \frac{\text{Correctly predicted labels}}{\text{total Number of labels}}$$

$$\Rightarrow \frac{1406 + 298}{1406 + 143 + 263 + 298} = 80.75$$

```
> y-train-pred-final.head()
```

churn	Churn-prob	custID	predicted
0	0.225	879	0
0	0.244	8790	0
1	0.692	6498	1
1	0.504	880	1

```
> from sklearn import metrics
```

```
# Confusion Matrix
```

```
> Confusion = metrics.confusion_matrix(y-train-pred-final.churn, y-train-pred-final.predicted)
```

```
> print(Confusion)
```

```
[ [ 3270  365 ]  
[ 579  208 ] ]
```

```
> print(metrics.accuracy_score(y-train-pred-final.churn, y-train-pred-final.predicted))
```

```
# Manual Feature Elimination
```

→ We used RFE to select 15 features, but still few features having higher correlation, so for better feature elimination we use NIF

$$VIF_i = \frac{1}{1 - R_i^2}$$

where  $i$  refers to the  $i^{th}$  variable which is being represented as a combination of the rest of the independent variables.

Here  $R^2$  means, how well the other features predict this feature.

# Checking the VIF values for the feature variables

> from statsmodels.stats.outliers\_influence import variance\_inflation\_factor

# creating a VIF Dataframe

> vif = pd.DataFrame()

> vif['Features'] = x\_train[col].columns

> vif['VIF'] = [variance\_inflation\_factor(x\_train[col].values, i) for i in range(x\_train[col].shape[1])]

> vif['VIF'] = round(vif['VIF'], 2)

> vif = vif.sort\_values(by='VIF', ascending=False)

> vif

↳ PhoneService → 3.86 VIF (Hence drop it)

> col = col.drop('PhoneService', 1)

> col

- # Let's rerun the model using selected variables
- ```
> x_train_sm = sm.add_constant(x_train[col])
> logm3 = sm.GLM(y_train, x_train_sm,
family = sm.families.Binomial())
> res = logm3.fit()
> res.summary()
> y_train_pred = res.predict(x_train_sm), values,
reshape(-1)
> y_train_pred_final['churn_prob'] = y_train_pred
# Creating new col 'Predicted' with 1 if churn_prob > 0.5
else 0
> y_train_pred_final['Predicted'] = y_train_pred_final.
churn_prob.map(lambda x:
1 if x > 0.5 else 0)
> y_train_pred_final.head()

```
- |   | Churn | Churn_prob | CustID | Predicted |
|---|-------|------------|--------|-----------|
| 0 | 0.254 | 829        | 0      |           |
| 0 | 0.224 | 5790       | 0      |           |
| 1 | 0.693 | 6498       | 1      |           |
| 1 | 0.510 | 380        | 1      |           |

# let's check the overall accuracy

Point( metrics.accuracy\_score(y\_train\_pred\_final,  
churn, y\_train\_pred\_final.predicted))

↳ 0.805 → 80.5%

# Again check for VIF

↳ total charges → VIF ↗

# Now we can drop those columns based on  
business model or business use case.

Better is VIF < 5 for all features.

Q You are building a logistic regression model  
for a media company to predict whether a  
person will like a new show "Sacred Curses"  
variable value of 1 indicates that the person liked  
the show , while 0 means they didn't , the coefficient  
for the 5 shows are as follows :

| variable-name        | coefficientval |
|----------------------|----------------|
| True-detective-liked | 0.47           |
| Modern-family-liked  | -0.45          |
| Mindhunters-liked    | 0.39           |
| Friends-liked        | -0.23          |
| Narcos-liked         | 0.55           |

Now we have the data of 3 consumers, Reetesh, Kshitij, Shruti for these 5 shows indicating whether or not they liked these shows

| Consumer | True-Detective | Modernfamily | MindHunter | Forwards | Narcos |
|----------|----------------|--------------|------------|----------|--------|
| Reetesh  | 1              | 0            | 0          | 0        | 1      |
| Kshitij  | 1              | 1            | 1          | 0        | 1      |
| Shruti   | 0              | 1            | 0          | 1        | 1      |

To find the person who is most likely to like the show we can use log odds.

$$\ln(P/(1-P)) = B_0 + B_1 x_1 + B_2 x_2 + \dots + B_n x_n$$

$$\Rightarrow 0.47 x_1 - 0.45 x_2 + 0.39 x_3 - 0.23 x_4 + 0.55 x_5$$

For Reetesh:  $0.47 \times 1 - 0.45 \times 0 + 0.39 \times 0 - 0.23 \times 0 + 0.55 \times 1 = 1.02$

For Kshitij:  $0.47 \times 1 - 0.45 \times 1 + 0.39 \times 1 - 0.23 \times 0 + 0.55 \times 1 = 0.96$

for Shruti  $\Rightarrow -0.13$

$\Rightarrow$  Reetesh liking the show is the highest

## Metrics Beyond Accuracy : Sensitivity & Specificity

Ques is : Is accuracy enough to assess the goodness of the model?  $\rightarrow$  Big NO!

To understand why accuracy is often not the best metric (Example)

$\rightarrow$  Increasing churn is the most serious issue in the telecom company, company desperately wants to retain customers, Marketing Head decides to roll out discounts/offers to the customers likely to churn. It's imp that the model identifies almost all the churn customers correctly. It's fine if it incorrectly predicts some of the non-churn customers as churn. Since in that case  $\rightarrow$  the worst that will happen is that company will offer discounts to them.

Confusion matrix we got for our model again -  
the actual labels are along the column while the predicted labels are along the rows.

| Actual / Predicted | Not churn | churn |
|--------------------|-----------|-------|
| Not churn          | 3269      | 366   |
| churn              | 595       | 692   |

From the table we can see that there are

$$595 + 692 = 1287 \text{ actual churn customers}$$

So ideally model should predict all of them as churn but only of those 1287, the current model only predicts 692 as churn  $\Rightarrow$  only 692 out of 1287 or only about 53% of churn customers will be predicted by the model as churn  $\Rightarrow$  It's very risky.

So even if the accuracy is 80%, the model only predicts 53% of churn cases correctly.

Here what's happening is that we care more about one class (class = 'churn') than the other. It's a very common situation in classification problem  $\rightarrow$  we may always care more about one class than the other.

On the other hand the accuracy tells you the model's performance on both classes combined  $\rightarrow$  which is fine but not the most imp metric.

Example 2: Suppose we are building a logistic regression model for cancer patients, based on certain features we need to predict whether the patient has cancer or not. Here if we incorrectly predict many diseased patients as not having cancer, it can be very risky.

In such cases it is better that instead of looking at the overall accuracy, you care about predicting the 1st (the disease) correctly.

Similarly, if we are building a model to determine whether you should block (where blocking is 1 and not blocking is 0) a customer's transactions or not based on the past transaction behaviour in order to detect frauds → Here, we would care more about getting the  $\text{O}^{\circ}$  right, Because we might not want to wrongly block a good customer's transactions as it might lead to a very bad customer experience.

\* Hence It's crucial that we consider the overall business problem you are trying to solve to decide the metric you want to maximize or minimize.

This brings us to two of the most commonly used metrics to evaluate a classification model.

① Sensitivity

② Specificity

Sensitivity is defined as :

$$\text{Sensitivity} = \frac{\text{No. of actual Yeses correctly predicted}}{\text{Total No. of Actual Yeses}}$$

Here 'yes' means 'churn' and 'no' means 'non-churn'

| Actual / Predicted | Not churn | churn |
|--------------------|-----------|-------|
| Not churn          | 3069      | 366   |
| churn              | 595       | 692   |

The different elements in the matrix can be labelled as follows

| Actual / Predicted | Not churn       | churn          |
|--------------------|-----------------|----------------|
| Not churn          | True negatives  | False Positive |
| churn              | False Negatives | True Positives |

⇒ the first cell contains the actual 'notchurns' being predicted as 'not churns' hence labelled as True Negatives. (Negatives imply that class is 0 here, Not-churn)

⇒ the second cell contains actual 'not churns' being predicted as 'churn' → Hence it is labelled as False Neg.

⇒ Similarly the third cell contains the actual 'churns' being predicted as 'not churn' → False Neg.

⇒ fourth cell contains actual 'churns' which are predicted as 'churns' hence True Neg.

$$\text{Sensitivity} \Rightarrow \frac{692}{595+692} \rightarrow \text{Actual Yes} \approx 88.76\%$$

So the sensitivity is quite low.

Similarly we have Specificity:

$$\text{Specificity} = \frac{\text{Nos of actual No correctly predicted}}{\text{Total Nos of actual No}}$$

$$= \frac{3269}{3269 + 366}$$

# When evaluating the performance of the classification model, accuracy (the ratio of correctly predicted observations to the total observations) is often insufficient, especially when dealing with Imbalanced Datasets, so to better assess model performance we use metrics like Sensitivity and specificity, these are derived from Confusion Matrix.

Confusion Matrix is a must to understand for both the Sensitivity and Specificity.

## Confusion Matrix

It is a table that summarizes the performance of a classification model with two classes (binary classification)

|                 | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | True Positive      | False Negative     |
| Actual Negative | False Positive     | True Negative      |

# True Positive : correctly predicted positives (TP)

# True Negative (TN) : correctly predicted negatives

# False Positives (FP) : incorrectly predicted as true (Type I error)

# False Negatives (FN) : incorrectly predicted as -ve. (Type II error)

### ① Sensitivity (Recall or True Positive Rate)

→ Sensitivity measures the proportion of actual +ves that are correctly identified by the model.

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

Key points :

Answers the ques : How good is the model at identifying the cases?

High sensitivity means fewer false negatives which is critical in applications like disease diagnosis (e.g. detecting cancer).

e.g.: In a cancer test, sensitivity ensures that most cancer patients are correctly identified as +ve.

## 2. Specificity (True negative Rate)

It measures the proportion of actual -ve that are correctly identified as -ve by the model.

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

Key points:

→ Answers the question: How good is the model at identifying -ve cases?

→ Higher specificity means fewer false +ve, which is critical in applications like fraud detection.

e.g. In a fraud detection system, it ensures that the legitimate transactions are not flagged as fraudulent.

In our model

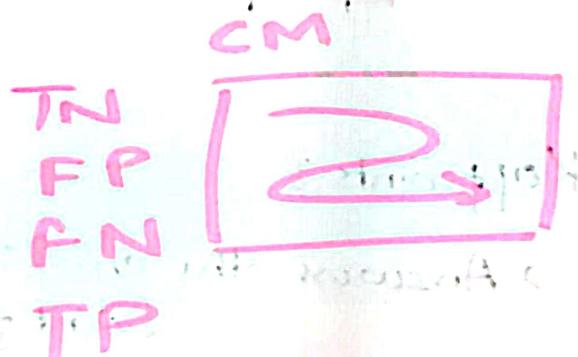
> confusion = metrics.confusion\_matrix(y\_train\_pred\_final.churn, y\_train\_pred\_final.predicted)

> Confusion

↳  $\begin{bmatrix} [3269, 366], \\ [595, 692] \end{bmatrix}$

> metrics.accuracy\_score(y\_train\_pred\_final.churn, y\_train\_pred\_final.predicted)

↳ 80.25



> TP = confusion[1, 1]

> TN = confusion[0, 0]

> FP = confusion[0, 1]

> FN = confusion[1, 0]

> Sens = TP / float(TP+FN)

> Spec = TN / float(TN+FP)

# Calculate false +ve rate → Predicting churn when customer Does not have churned

> FP / float(TN+FP)

## Notes:

A confusion Matrix is a tool often used in the classification problems to evaluate the performance of the model  $\Rightarrow$  It is particularly useful for the binary and multi-class classification problem.

Let's see in detail

### Confusion Matrix Overview

for a binary classification problem, the confusion matrix is structured as follows:

|                 | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | True Positive      | False Negative     |
| Actual Negative | False Positive     | True Negative      |

$\checkmark$  Predicted +ve class &  $\exists$  Actual +ve class  $\Rightarrow$  True Positive (TP)

$\checkmark$  Predicted +ve class &  $\nexists$  Actual +ve class  
 $\Rightarrow$  False Negative (FN)

$\checkmark$  Predicted -ve class &  $\exists$  Actual -ve class  
 $\Rightarrow$  True Negative (TN)

$\checkmark$  Predicted +ve class &  $\nexists$  Actual -ve class  
 $\Rightarrow$  False Positive (FP)

- True Positive (TP) : the model actually predicts a positive case, i.e. correctly Predicts a positive case.
- False Positive (FP) : Model incorrectly predicts a positive case, when it's actually negative.
- True Negative (TN) : Model correctly predicts a -ve case
- False Negative (FN) : The Model Incorrectly predicts a -ve case, when it's actually +ve.

### \* Key Metrics Derived from Confusion Matrix

① **Accuracy** : Measures the overall correctness of the model.

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

e.g. If TP = 50, TN = 30, FP = 10, FN = 10

$$\text{Accuracy} = \frac{50 + 30}{50 + 30 + 10 + 10} = \frac{80}{100} = 0.8$$

② **Sensitivity (Recall or True Positive Rate)**

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

Measures the proportion of actual +ve correctly identified

By using the same numbers

$$\text{Sensitivity} = \frac{50}{50+10} = 0.83$$

### 3. Specificity (True -ve Rate)

$$\text{Specificity} = \frac{TN}{TN + FP}$$

Measures the proportion of actual -ves correctly identified.

### 4. False Positive Rate (FPR)

$$FPR = \frac{FP}{FP + TN}$$

Measures the proportion of actual negatives, incorrectly classified as positive.

$$FPR = 1 - \text{Specificity}$$

$$FPR = \frac{10}{10+30} = \frac{10}{40} = 0.25$$

### 5. Precision (Positive Predicted Value)

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

Measures the proportion of predicted values which are actually +ve.

Q Given the Confusion matrix, to calculate the True Positive Rate and False Positive Rate

$$\text{True Positive Rate} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

|           | Churn | Not Churn |
|-----------|-------|-----------|
| Churn     | 300   | 200       |
| Not Churn | 100   | 400       |

$$\text{False Positive Rate} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

$$\Rightarrow \text{TPR} = \frac{300}{300 + 200} = \frac{60}{500} \times 100 = 60$$

$$\text{FPR} = \frac{100}{100 + 400} = \frac{100}{500} \times 100 = 20$$

Q You have the following table showcasing the actual churn labels and predicted probabilities for 5 customers

| Customer | Churn | Predicted Churn Prob |
|----------|-------|----------------------|
| A1       | 1     | 0.52                 |
| A2       | 0     | 0.56                 |
| A3       | 1     | 0.28                 |
| A4       | 0     | 0.45                 |
| A5       | 0     | 0.22                 |

To calculate true rate and false rate for cutoffs of 0.4 and 0.5, which of these cutoffs gives a better model

Note : the good model is the one in which TPR is high and FPR is low

Ans : To determine which cutoff (0.4 or 0.5) gives a better model we calculate TPR and FPR for both cutoff values.

- True +ve (TP) : Predicted "churn" when actual churn is 1.
- False +ve (FP) : Predicted churn when actual churn is 0.
- False -ve (FN) : Predicted Non-churn when actual churn is 1.
- True -ve (TN) : Predicted Non-churn when actual churn is 0.

$$\text{True Positive Rate (TRate)} = \frac{TP}{TP + FN}$$

$$\text{False Positive Rate (FPR)} = \frac{FP}{FP + TN}$$

1. Calculation for

cutoff = 0.4  $\Rightarrow$  if Predicted Churn Prob  $\geq 0.4$   
Predict Churn (1) else Predict Non-churn (0)

| Customer | Actual Churn | Pred Prob | Predicted Churn |
|----------|--------------|-----------|-----------------|
| A1       | 1            | 0.52      | 1               |
| A2       | 0            | 0.56      | 1               |
| A3       | 1            | 0.78      | 1               |
| A4       | 0            | 0.45      | 1               |
| A5       | 0            | 0.22      | 0               |

## Step 2: Confusion Matrix

$$\rightarrow TP = 2 (A_1, A_3)$$

$$\rightarrow RP = 2 (A_2, A_4)$$

$$\rightarrow FN = 0$$

$$\rightarrow TN = 1 (A_5)$$

$$TPR = \frac{2}{2+0} = 1$$

$$FPR = \frac{2}{2+1} = 0.67$$

2. Calculations for cutoff = 0.5

↳ if Predicted churn Prob  $\geq 0.5 \rightarrow$  Predict churn 1

otherwise predict non-churn 0

| Customer | Actual Churn | Predicted Prob | Predicted Churn |
|----------|--------------|----------------|-----------------|
| A1       | 1            | 0.52           | 1               |
| A2       | 0            | 0.56           | 1               |
| A3       | 1            | 0.78           | 1               |
| A4       | 0            | 0.45           | 0               |
| A5       | 0            | 0.22           | 0               |

$$TP = 2 (A_1, A_3)$$

$$RP = 1 (A_2)$$

$$FN = 0$$

$$TN = 2 (A_4, A_5)$$

$$TPR = \frac{2}{2+0} = 1$$

$$FPR = \frac{1}{1+2} = \frac{1}{3} = 0.33$$

↳ Better



## ROC Curve : Receiver Operating characteristic Curve

the ROC curve is a graphical representation of the trade-off b/w the True Positive Rate (TPR) also called as sensitivity and the False Positive Rate (FPR) at various classification thresholds.

### Why do we need ROC curve?

When building a binary classification model (e.g. predicting churn vs non-churn), the classification decision depends on a threshold applied to predicted probabilities (commonly 0.5 by default). However, this threshold may not be optimal.

- Sensitivity (TPR) : How well the model identifies actual positive.
- Specificity : How well the model identifies actual -ve's.

### The ROC curve helps us:

- ① Visualize the trade-off b/w TPR and FPR for different thresholds
- ② Choose an optimal threshold for classification balancing TPR and FPR
- ③ Evaluate overall performance of the model using AUC (Area Under curve)

## Key Metrics

① True Positive Rate (TPR) =  $\frac{TP}{TP + FN}$

↳ Also called as sensitivity or recall

It measures the proportion of Actual +ve correctly identified.

② False Positive Rate (FPR) =  $\frac{FP}{FP + TN}$

↳ It measures the proportion of actual -ve incorrectly classified as +ve

## How is ROC curve constructed

① calculate TPR and FPR for various thresholds  
e.g. (0.1, 0.2 ... 0.9)

② Plot TPR (y axis) vs FPR (x axis)

→ threshold = 1 : All predictions are -ve  
(TPR=0, FPR=0)

→ threshold = 0 : All Predictions are +ve  
(TPR=1, FPR=1)

As you vary the threshold the curve moves from (0,0) to (1,1)

## Interpreting the ROC Curve

- ① Closer to the top-left corner
    - ↳ Indicates better performance (high TPR and low FPR)
  - ② Diagonal line :
    - ↳ Represents a model with no discriminatory power ( $AUC = 0.5$ )
  - ③ AUC (Area under the curve)
    - $AUC = 1$  (Perfect Model)
    - $AUC = 0.5$  (Random guessing)
    - Higher AUC value indicates better model performance
- # Choosing the optimal threshold
- ① Maximize TPR while minimizing FPR
    - Identify the point on the curve closest to top left corner.
    - This often corresponds to a balanced sensitivity and specificity.
  - ② Business context:
    - Optimal threshold depends on the problem.
    - For e.g.
      - In medical diagnosis, prioritize sensitivity (reduce false negatives)
      - In fraud detection → Prioritize specificity (reduce false positives)

Eq Let's assume the following predicted prob for 5 customers?

| Customer | True Label | Predicted Prob |
|----------|------------|----------------|
| A        | 1          | 0.9            |
| B        | 0          | 0.8            |
| C        | 1          | 0.7            |
| D        | 0          | 0.4            |
| E        | 0          | 0.2            |

- ① Set diff thresholds (eg 0.5, 0.6 etc)
- ② Classify each customer as +ve or -ve
- ③ Calculate TPR and FPR for each threshold.
- ④ Plot TPR vs FPR for each threshold to get ROC curve.

### ROC Curve in Python

Step 1: import necessary libraries

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
from sklearn.metrics import roc_curve, auc
```

```
# True Labels
```

```
y_true = np.array([1, 0, 1, 0, 0])
```

```
# Predicted Prob from your model
```

```
y_scores = np.array([0.9, 0.8, 0.7, 0.4, 0.2])
```

Step 3: Compute TPR, FPR, thresholds

>  $fpr, tpr, \text{thresholds} = \text{roc\_curve}(y_{\text{true}}, y_{\text{scores}})$

Step 4: Calculate AUC

>  $\text{roc\_auc} = \text{auc}(fpr, tpr)$

Step 5: Plot the ROC curve

> plt.figure(figsize=(8, 6))

> plt.plot(fpr, tpr, color="blue") Diagonal line  
on ROC

> plt.plot([0, 1], [0, 1], linestyle="--", color='gray') ↑

> plt.xlabel("False pos rate") ↓

> plt.ylabel("True pos rate") Serves as a  
baseline

> plt.show()

## ④ Finding the optimal threshold for classification Model.

To find the optimal threshold for a classification model, we balance specificity, sensitivity and accuracy.

### I. Understanding Metrics

→ Sensitivity (Recall / True Positive Rate):

↳ Proportion of actual positives correctly identified

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

→ Specificity (True Negative Rate):

↳ Proportion of actual negatives correctly identified.

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

→ Accuracy : Proportion of all correct predictions (both +ve and -ve)

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

## 2. Generate Predictions Across Thresholds

→ For classification models, probabilities are often output instead of binary labels. Vary the threshold from 0 to 1 to compute the corresponding confusion matrix and metrics at each threshold.

## 3. Compute Metrics for each threshold

→ At each threshold

- ① Classify data points as positive or negative based on threshold.
- ② Compute TP, TN, FP, FN from the confusion matrix.
- ③ Calculate specificity, sensitivity and accuracy.

## 4. Choose an Optimization Criterion

Decide what you want to optimize

→ Maximize Accuracy: Choose the threshold where accuracy is highest.

→ Balanced Sensitivity and Specificity:

choose the threshold where sensitivity ≈ specificity

## Eq Python Implementation

```
import numpy as np  
from sklearn.metrics import confusion_matrix,  
    roc_curve, auc
```

```
y_true = np.array([0, 1, 0, 1, 0, 1, 1, 0, 0, 1])
```

```
y_probs = np.array([0.1, 0.4, 0.35, 0.8, 0.2, 0.9,  
    0.7, 0.15, 0.05, 0.35])
```

```
# Generate metrics for the thresholds
```

```
thresholds = np.arange(0, 1, 0.01)
```

```
sensitivities, specificities, accuracies = [], [], []
```

```
for threshold in thresholds:
```

```
y_pred = (y_probs >= threshold).astype(int)
```

```
fn, fp, tn, tp = confusion_matrix(y_true, y_pred,  
    ).ravel()
```

```
sensitivity = tp / (tp + fn)
```

```
specificity = tn / (tn + fp)
```

```
accuracy = (tp + tn) / (tp + tn + fp + fn)
```

```
sensitivities.append(sensitivity)
```

```
specificities.append(specificity)
```

```
accuracies.append(accuracy)
```

# find optimal threshold by maximizing accuracy

Optimal\_idx = np.argmax(accuracies)

Optimal\_threshold = thresholds[Optimal\_idx]

# Plotting Metrics

import matplotlib.pyplot as plt

plt.figure(figsize=(10, 6))

plt.plot(thresholds, sensitivities, label='Sensitivity')

plt.plot(thresholds, specificities, label='Specificity')

plt.plot(thresholds, accuracies, label='Accuracy')

plt.xlabel('Threshold')

plt.ylabel('Metric value')

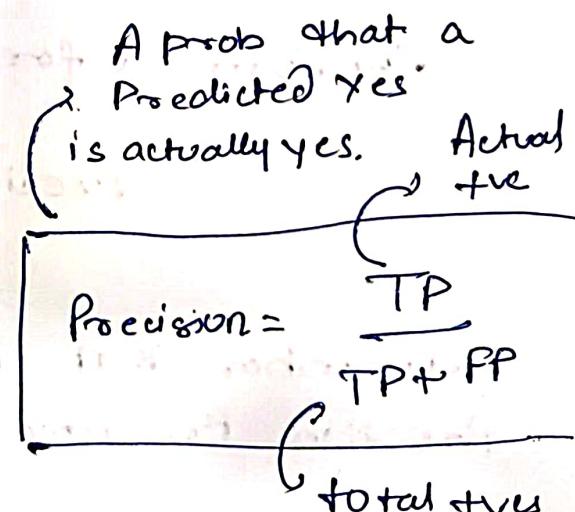
plt.title('Sensitivity, Specificity, Accuracy')

plt.show()

## \* Precision & Recall

confusion Matrix

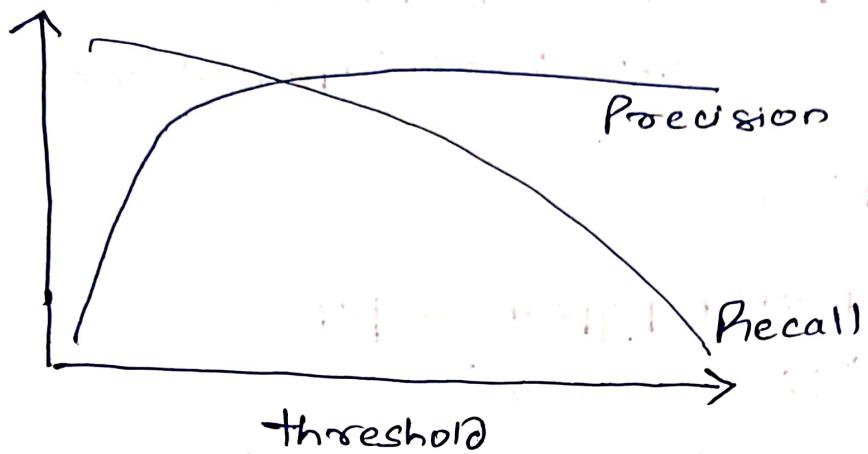
| True/Predicted | No | Yes |
|----------------|----|-----|
| No             | TN | FP  |
| Yes            | FN | TP  |



Recall:  $\frac{TP}{TP+FN}$

→ Prob that a yes case is predicted as such

## Precision and Recall Tradeoffs



As threshold  $\uparrow \rightarrow$  Precision  $\uparrow$

As threshold  $\downarrow \rightarrow$  Recall  $\uparrow$

\* Precision: Probability that a predicted Yes is actually a Yes.

|    |     | True/Predicted |          |
|----|-----|----------------|----------|
|    |     | No             | Yes      |
| No | No  | TN             | FP<br>TP |
|    | Yes | FN             |          |

The formula for precision can be given as:

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

Remember that Precision is the same as the Positive Predicted value.

\* Recall : Probability that our actual Yes is predicted correctly.

| True / Predicted | No | Yes |
|------------------|----|-----|
| No               | TP | FP  |
| Yes              | FN | TP  |

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (\text{Same as Sensitivity})$$

# we can use Sklearn library

- > from sklearn.metrics import precision\_score, recall\_score
- > precision\_score(y\_train\_pred\_final.churn, y\_train\_pred\_final.predicted)
- > recall\_score(y\_train\_pred\_final.churn, y\_train\_pred\_final.predicted)

Q We have the confusion matrix as

| Actual / Predicted | Not churn | churn |
|--------------------|-----------|-------|
| Not churn          | 1294      | 234   |
| churn              | 223       | 359   |

What will be the approximate value of accuracy?

$$\Rightarrow \text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{Total Observations}}$$

$TN \rightarrow 1294$  (not churn predicted as not churn)

$TP = 359$  (churn predicted as churn)

$FN \rightarrow 223$  (wrongly predicted as not churn)

$FP \rightarrow 234$  (wrongly predicted as churn)

$$\text{Total observations} = TN + TP + FN + FP = 1294 + 359 + 223 + 234$$

$$TN + TP = 1294 + 359 = 1653 \quad = 2110$$

$$\Rightarrow \frac{1653}{2110} = 0.78$$

$$\text{Recall} = \frac{TP}{TP + FN} \Rightarrow \frac{359}{359 + 223} = 0.61$$