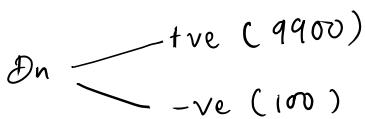


## Classification Algorithm in various situations

### 1. What is Imbalanced and balanced dataset?

Algorithm may get biased if we have a lot of +ve or -ve point in our data set.

ex: lets say we have a lot of positive point, then algorithm will always predict +ve for future unseen data



### 1. Balanced data set (2 class classification)

$$D_n = n_1 + n_2$$

$$\begin{matrix} \downarrow & \downarrow \\ +ve & -ve \end{matrix}$$

if  $n_1 \approx n_2$

{ its okay if  $n_1 + n_2$   
but it roughly similar

### 2. Imbalanced dataset :-

$$\text{if } n_1 \ll n_2 \quad \text{or} \quad n_2 \ll n_1$$

{ there is very much difference in  
# of points

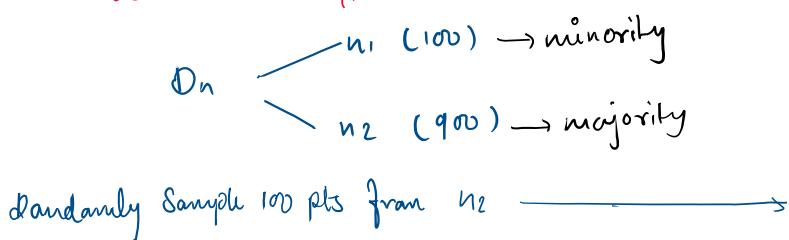
ex: K-NN with Imbalanced data

$n_1 = 50 \quad n_2 = 950$  \* model will be biased towards -ve because of majority  
 +ve                  -ve  
 value  
 \* Dominant class.

Problem: You can get high accuracy with imbalanced data with a dumb model

# How to work around imbalanced dataset?

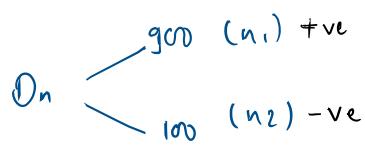
#### 1. undersampling:



Problem: \* we are throwing away 800 pts which should be available  
 \* There are very few points in D'n so model may not even work

# Note: Throwing away data should be avoided.

#### 2. Oversampling



{ 100 -ve pts (n2)  
 ↓  
 repeat 9 times  
 ↓  
 900 -ve pts

## Classification Algorithm in various situations

D<sub>n</sub> : 900 +ve pts  
900 -ve pts

\* placing more points from minority class in the dataset

1. replacing points from minority class
2. Extrapolation (creating artificial or synthetic points)

3. Weightage (not available for KNN)

D<sub>n</sub>   
900 (-ve pts)      w+ve = 9 → more weight to minority class  
100 (+ve pts)      w-ve = 1 → less weight to majority class

It is equivalent to repetition

## 2. Define Multi-class classification?

1. Binary Classification:  $y_i \in \{0, 1\}$

2. Multiclass classification:  $y_i \in \{0, 1, 2, 3, \dots, k\}$

Some algorithms can only do binary classification (ex: logistic regression)

→ we can do one vs rest

## # KNN Given a distance measure.

Sometimes we don't get  $x_i$  as a numerical vector, but we get it as similarity matrix

As K-NN cares mostly about distance or similarity

→ Given a sim-matrix or distance matrix explicitly → K-NN can work

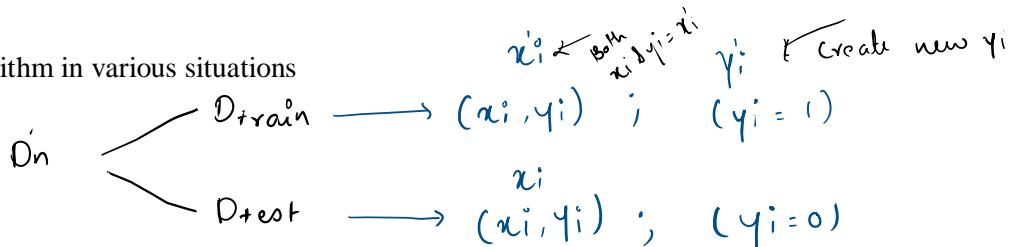
$$\text{distance} = \frac{1}{\text{similarity}}$$

## 3. How to determine if data is changing over time? Or Distribution (Dtrain) != Distribution (Dtest)?

D<sub>n</sub>   
Dtrain ( $x_i, y_i$ )  
Dtest ( $x_i, y_i$ )

① Create a new data set using D<sub>n</sub> to check the distribution

## Classification Algorithm in various situations

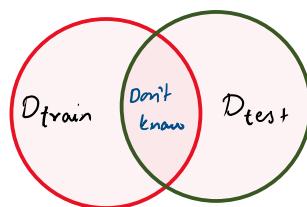


② Create a binary classifier on  $D'_n$

$$f_x \rightarrow 1 \text{ or } 0$$

High error means similar distribution

18



③ if binary classification has accuracy of 70%. i.e 30% missclassified

This mean train and test data cannot be seperated completely

∴ distribution is not very similar

18



→ If is almost overlapping

Binary classification accuracy low

This means  $D_{train}$  &  $D_{test}$  are coming from same distribution

Higher the misclassification mean similar distribution.

## 4. Explain Impact of Outliers?

Outliers can change the model drastically as it can easily affect decision surface

for k-NN → \* if  $k$  is small : outlier can easily impact model

\* if  $k$  is large : it will be less prone to outliers.

## 5. What is LocalOutlier Factor?

④  $k$ -distance:  $k$ -distance ( $x_i$ ) = distance to the  $k^{th}$  nearest neighbour of  $x_i$  from  $x_i$

$N(A) = \text{neighbourhood of } A \dots \text{ - set of all point that belong to } k\text{-NN of } x_i$

### (b) Reachability distance:-

$$\text{Reachability distance } (x_i, x_j) = \max(k\text{-distance}(x_j), \text{distance}(x_i, x_j))$$

If  $x_i$  is in neighbourhood of  $x_j$  ( $N(x_j)$ ) then reachability distance is  $k^{\text{th}}$  distance of  $x_j$  (i.e. farthest distance of nearest neighbour)

else, reachability distance is equal to distance between  $(x_i, x_j)$

### ⑥ Local reachability density $\rightarrow \text{lrd}(A)$

$$\text{lrd}(x_i) = \frac{1}{\sum_{x_j \in N(x_i)} \left\{ \frac{\text{reachability distance}(x_i, x_j)}{|N(x_i)|} \right\}}$$

size of set or  
elements of set  
Generally  $k$  but not necessarily because there can be  
multiple point at same distance.

Denominator: Avg reachability distance of  $x_i$  from its neighbours.

$\text{lrd}(x_i)$ : inverse of average reachability distance of  $x_i$  from its neighbours

$$\text{lrd}(x_i) = \frac{N(x_i)}{\sum_{x_j \in N(x_i)} \frac{\text{reachability distance}(x_i, x_j)}{\text{reachability distance}}}$$

density      # of points  
local      reachability distance

% density =  $\frac{\text{# of point}}{\text{unit area}}$

### ⑦ Local outlier factor(LOF)

$$\text{LOF} = \frac{\sum_{x_j \in N(x_i)} \text{lrd}(x_j)}{|N(x_i)|} \times \frac{1}{\text{lrd}(x_i)}$$

Avg lrd of point in  
neighbourhood of  $x_i$   
lrd of  $x_i$

If  $\text{LOF}(x_i) \rightarrow$  large  $\rightarrow$  outlier  
small  $\rightarrow$  inlier

## Classification Algorithm in various situations

1. for each point ( $x_i$ ) compute  $\text{LOF}(x_i)$
2. pick points with highest  $\text{LOF}$  -  $\rightarrow$  OUTLIER

Disadvantages:-

- \* LOF is very hard to interpretate
- \* use local outlier probability (modified LOF)