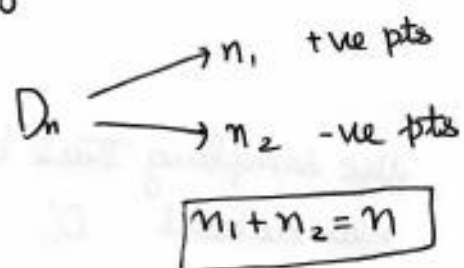


## Classification Algorithm in various situations / cases

- ↳ learning how to apply the techniques to real world data.  
 (KNN, logistic regression, Decision Tree etc)

### Imbalanced vs balanced data

2 class - classification problem:-



- 1) if  $n_1 \approx n_2$   
 $58\% \leftarrow n_1 \approx 580$      $n_1 \neq n_2$   
 $42\% \leftarrow n_2 = 420$

if  $n_1 \approx n_2$  then it is called balanced data set

- 2) if  $n_1 \ll n_2$  (or)  $n_2 \ll n_1 \rightarrow$  Imbalanced dataset  
 $n_1 = 100$     or     $n_2 = 150$   
 $n_2 = 900$         $n_1 = 850$

## K-NN (imbalanced dataset)

$$n_1 = 50 \text{ (+ve)} ; n_2 = 950 \text{ (-ve)}$$

(not always)

↑ majority class

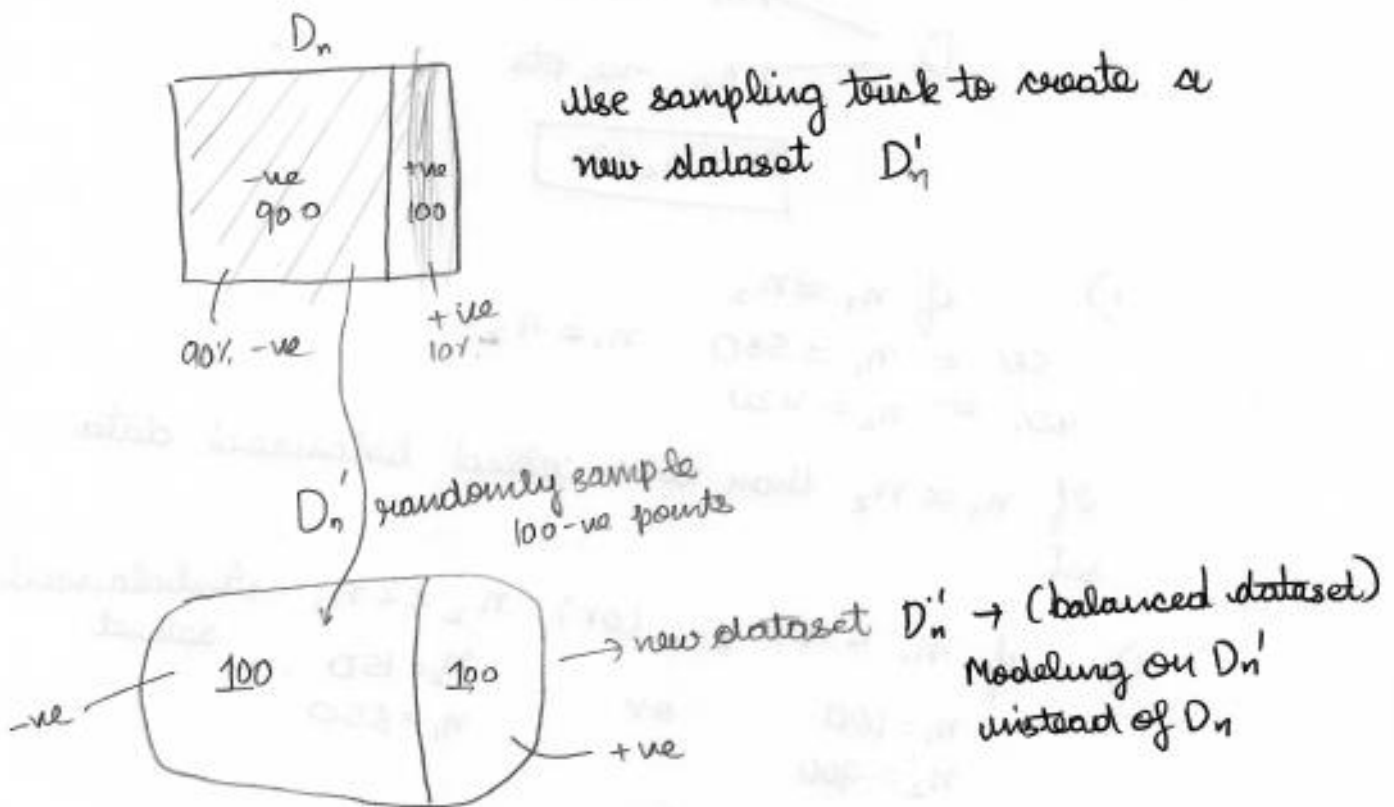
Dominating class

could have some advantage



? How to work around imbalanced dataset problem / issue?

① Undersampling  $\textcircled{D}$   $n_1 = 100 \text{ (+ve)}$   
 $n_2 = 900 \text{ (-ve)}$

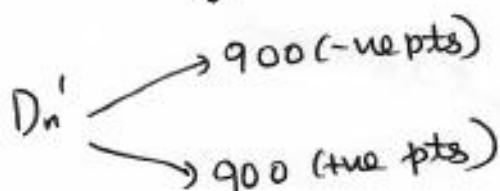
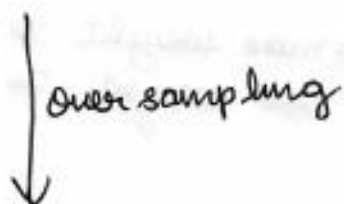
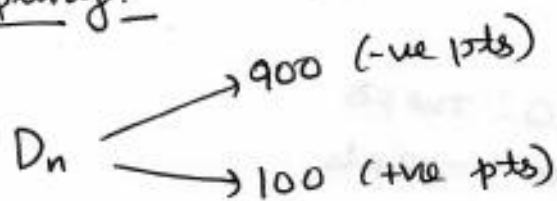


$$|D'_n| < |D_n|$$

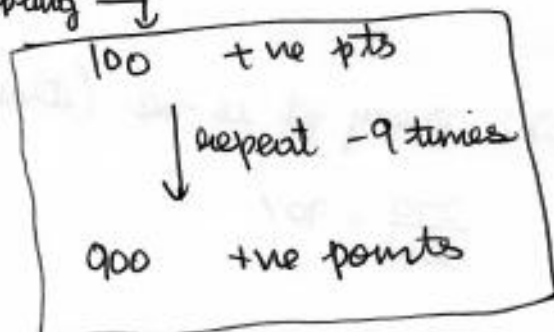
model have smaller amount of data  
it will not work accurately

? Problem with undersampling  
Throwing away data is not good

## ② Over sampling:



Oversampling →



⊗ - +ve pts  
 x - -ve pts  
 • → oversampled repetition  
 +ve pts

☆ Natural / Synthetic pts  $\rightarrow$  (Interpolation)  
add two pts in particular neighborhoods



◇ class-weight  $\div$  100 : +ve pts  
900 : -ve pts

$w_+ = 9 \rightarrow$  more weight to minority class  
 $w_- = 1 \rightarrow$  less weight to majority class



High accuracy  $f(x)$ : every pt is -ve (that is dumb model)

$$\frac{270}{300} = 90\%$$

multi class classification

Binary classifier:  $y_i \in \{0, 1\}$

Multi class classifier:  $y_i \in \{0, 1, 2, 3, \dots, 9\} \rightarrow$  MNIST.

(10 classes)

K-NN

$$D_n = \{(x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \{1, 2, 3, \dots, c\}\}$$

7-NN

$$x_q: \begin{array}{|c|c|c|c|c|} \hline 0 & 6 & 1 & \dots & 0 \\ \hline 1 & 2 & 3 & \dots & c \\ \hline \end{array}$$

Majority vote

$$y_q = \text{class 2}$$

$$y_q \neq 2$$

probabilistic classifier

$$\begin{cases} P(y_q = 2) = \frac{6}{7} \\ P(y_q = 3) = \frac{1}{7} \\ P(y_q = 1) = 0 = P(y_q = 4) = \dots = P(y_q = c) \end{cases}$$

NOTE : Logistic Regression  $\rightarrow$  They are fundamentally binary classifier & cannot do multi-class classification easily as KNN

Q Given a multi class classification problem; can we convert it into a binary classifier prob

$$f(x) \begin{cases} \rightarrow 0 \\ \rightarrow 1 \end{cases} \text{ binary}$$

$$f'(x) \begin{cases} \rightarrow 0 \\ \rightarrow 1 \\ \rightarrow 2 \\ \vdots \\ \rightarrow c \end{cases}$$

$y_i \in \{1, 2, 3, \dots, c\} \rightarrow c \text{ class} \rightarrow \{c \text{ binary classifiers}\}$

①  $D_n$    
  $\rightarrow \{(x_i, y_i) \mid y_i = 1\} \rightarrow \text{+ve pts}$    
  $\rightarrow \{(x_i, y_i) \mid y_i \neq 1\} \rightarrow \text{-ve pts}$    
 (Binary classifier) 1 or not

② Again  $D_n$

$D_n$    
  $\rightarrow y_i = 2$    
  $\rightarrow y_i \neq 2$    
 }  $\rightarrow \text{using this}$    
 (Binary classifier) 2 or not

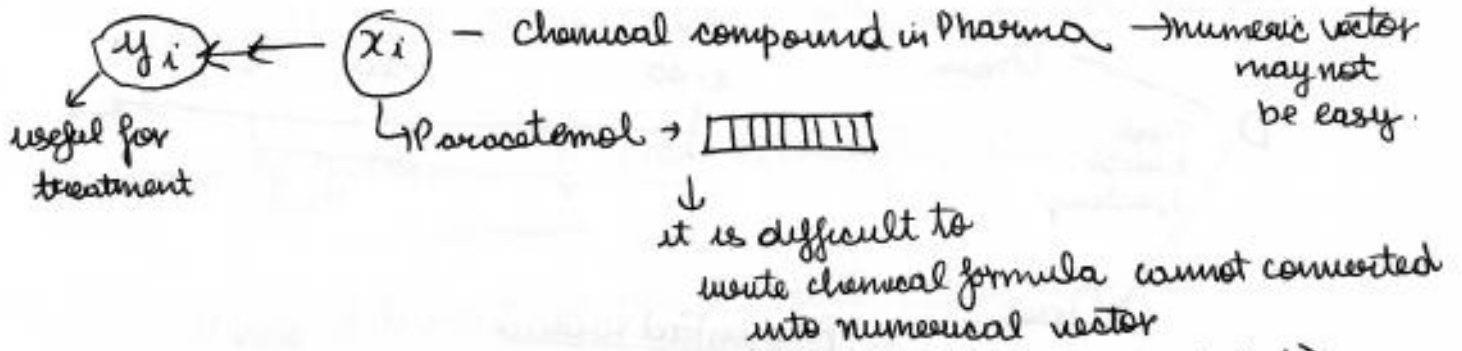
③  $D_n$    
  $\rightarrow y_i = 3$    
  $\rightarrow y_i \neq 3$    
 } (Binary classifier) 3 or not

Multi-class classification problem

$\downarrow$    
  $c$  binary classification problem   
  $\downarrow$    
  $\left\{ \begin{array}{l} f_1(x) \rightarrow \text{class 1 or not} \\ f_2(x) \rightarrow \text{class 2 or not} \\ \vdots \\ f_c(x) \rightarrow \text{class } c \text{ or not} \end{array} \right.$

→ KNN, given a distance measure

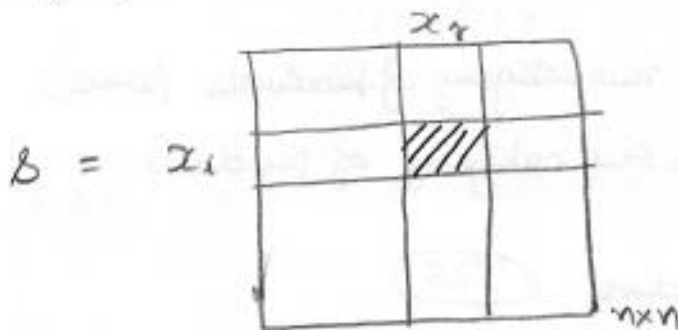
Classification:  $x_i \in \mathbb{R}^d$   
 ↳ vector is not easy



$\text{sim}(x_i, x_j) \rightarrow$  so similarity (distance is calculated)

$$\text{sim}(x_i, x_j) =$$

n datapoints

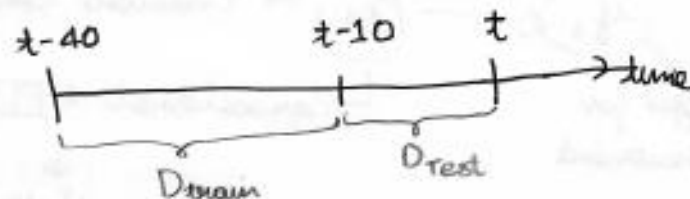
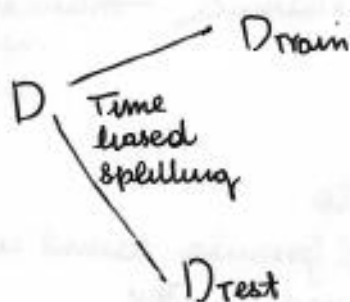


$$S_{ij} = \text{sim}(x_i, x_j)$$

$$\text{Distance: } d_{ij} = \frac{1}{S_{ij}}$$

→ Instead of each  $x_i$  as paracetamol we are given similarity matrix or distance matrix

## Train and test set differences



→ Amazon Food reviews

$D_{train}$  &  $D_{test}$  could be very different

Food reviews

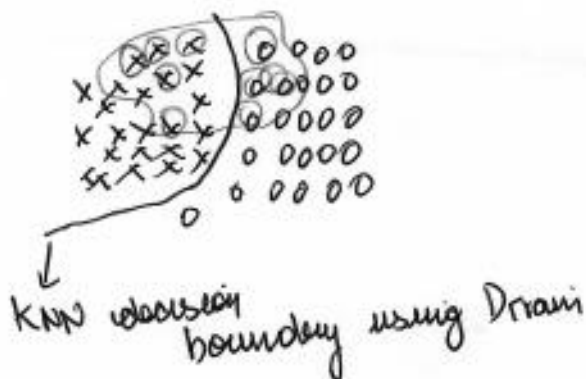
$t-10$  to  $t$  : - new category of products (wine)  
 $t-40$  to  $t-10$  : - old category of products  
 → data changes over time TBS

x: Train dataset -ve

o: Train +ve

⊗: Test -ve

⊙: Test +ve



$D_{train}$  and  $D_{test}$  are fundamentally different



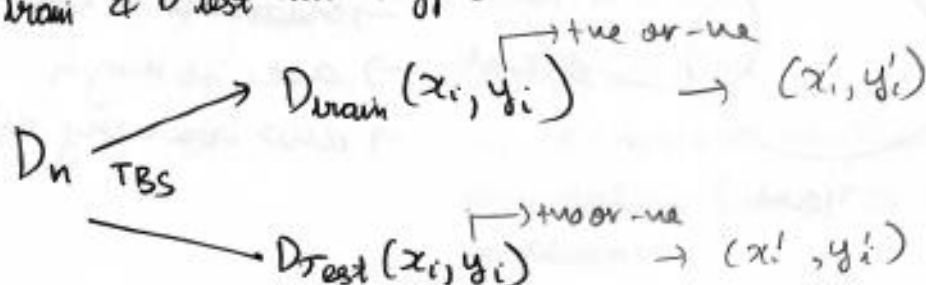
distribution of -ve data has changed from  $D_{train}$  to  $D_{test}$ .

- ① On  $D_{train}$  line perform very well but for  $D_{test}$  error is high.  
We need to check that  $D_{train}$  &  $D_{test}$  cannot do well.

Q how to determine if data is changing over time?  
or

$D_{train}$  &  $D_{test}$  have different distribution?

Sol



To solve problem we create new dataset

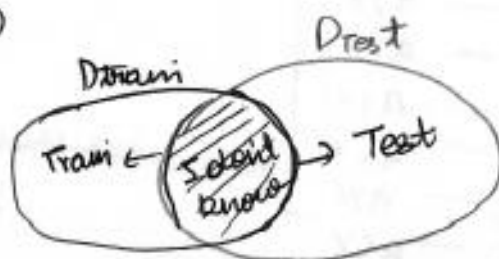
$D_n' \div$   $D_{train} \div y_i' = 1; x_i' = \text{concat}(x_i, y_i)$  in  $D_{train}$ .  
 $D_{test} \div y_i' = 0; x_i' = \text{concat}(x_i, y_i)$  in  $D_{test}$

Build a binary classifier on  $D_n'$

$f(x) \longrightarrow +1 \text{ or } 0$

KNN

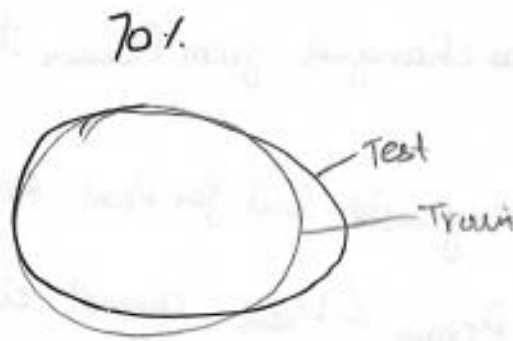
Case I



Binary classifier has accuracy of 70%.

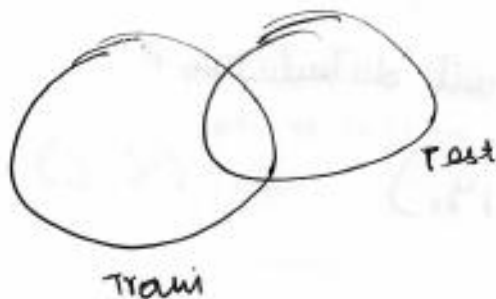
↓  
Train & Test data can be separated partially distribution are different

Case II



- almost overlapping
- binary classification (KNN)
- ↳ accuracy is low
- distributions are very similar

Case III:



- overlap v. low
- acc. is high
- dists are very different

$D_{train}$  &  $D_{test}$

→ same distribution

↳ not from same distribution

↳ This means features are changing with time

## Impact of outliers

For KNN, when  $K$  is small outliers can easily impact your model

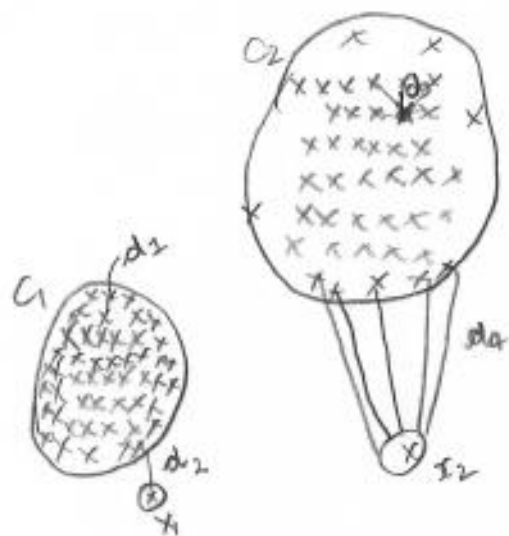
10 fold CV

$K=1$	—	97%
$K=2$	—	97%
$K=3$	—	97%
$K=4$	—	97%
$K=5$	—	97%
$K=6$	—	95%
$K=7$	—	92%

$K=5$  is less prone to outliers than  $K=1$

## Local outlier Factor (LOF)

Detect outliers in data inspired by K-NN.



→  $C_1$ : very dense cluster  
 $C_2$ : sparser cluster.

both  $x_1$  and  $x_2$  are outliers

→ simple solution

$x_i$ : K-nearest neighbours

mean distance from  $x_i$  to its K-nearest neighbours

for 5 nearest neighbours.

$d_1, d_2, d_3, d_4, d_5$  are average distances from point

$d_4$  is the largest value  
 mean distance from outlier is large consider it as an outlier