

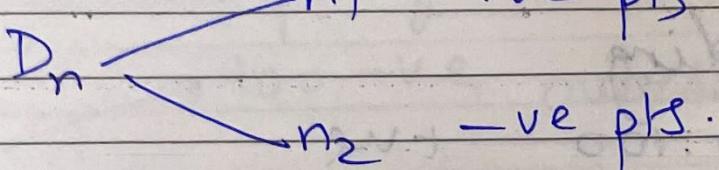
Ch 4: Classification algorithms in various situations

~~Review~~

Q: Imbalanced vs Balanced dataset

g)

Consider eg of 2-class classification



Let $n = 1000$

If $n_1 \approx n_2 \rightarrow$ Balanced dataset

eg $n_1 = 580$ $n_1 \neq n_2$, still $n_1 \approx n_2$
 $n_2 = 420$ significantly less than

eg if $n_1 \ll n_2$ or $n_2 \ll n_1$,
it is imbalanced dataset.

eg $n_1 = 100$
 $n_2 = 900$

$n_1 = 850$
 $n_2 = 150$

Imbalanced dataset may create problem.

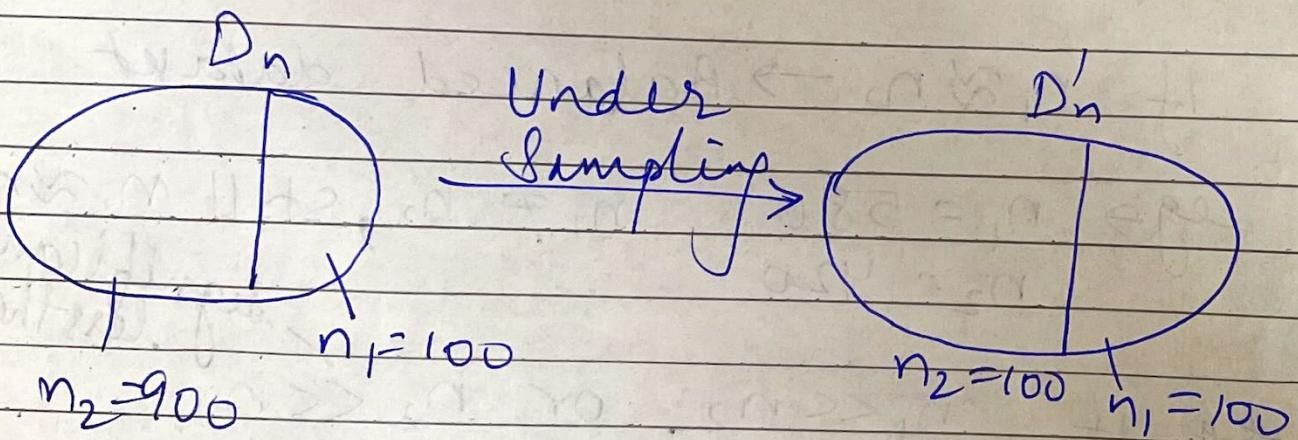
✓

Just because there are more points of the majority class, it may be more biased towards it.

So, how to work around the issue of imbalanced dataset.

1) Undersampling

$$D_n \begin{cases} n_1: 100 & +ve \\ n_2: 900 & -ve \end{cases}$$



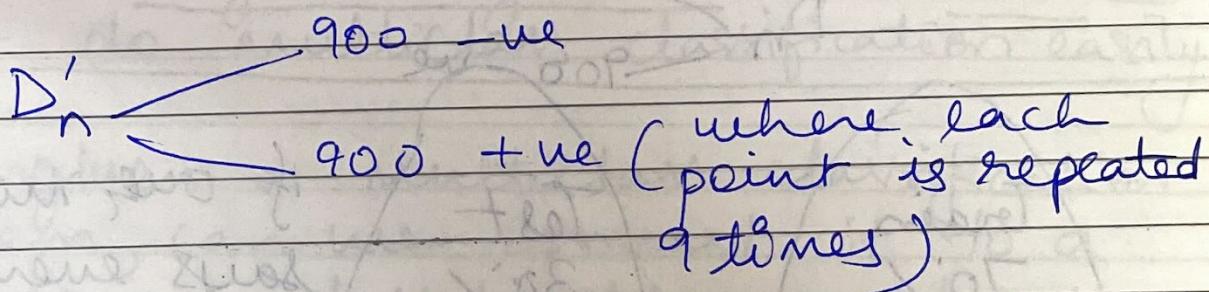
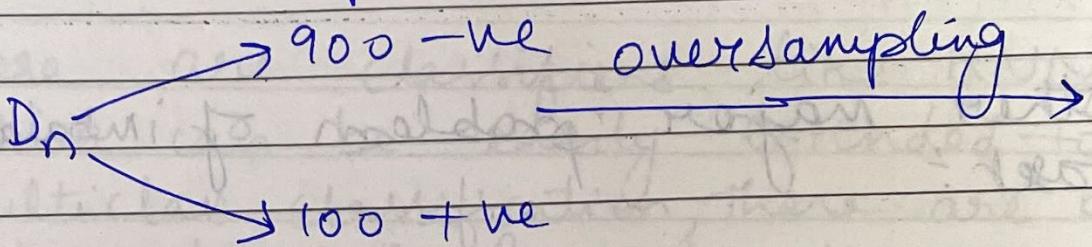
So basically we undersample
the -ve pts ~~equivalent to n_1~~ to D'_n to D_n

Now dataset D'_n is balanced.
Further modelling done on D'_n .

Problem with Undersampling

$|D'_n| \ll |D_n|$, the model may not work that well as we are throwing away a lot of data.

2) Oversampling



1) Simply repeating the existing pts. is one of the very simplest technique.

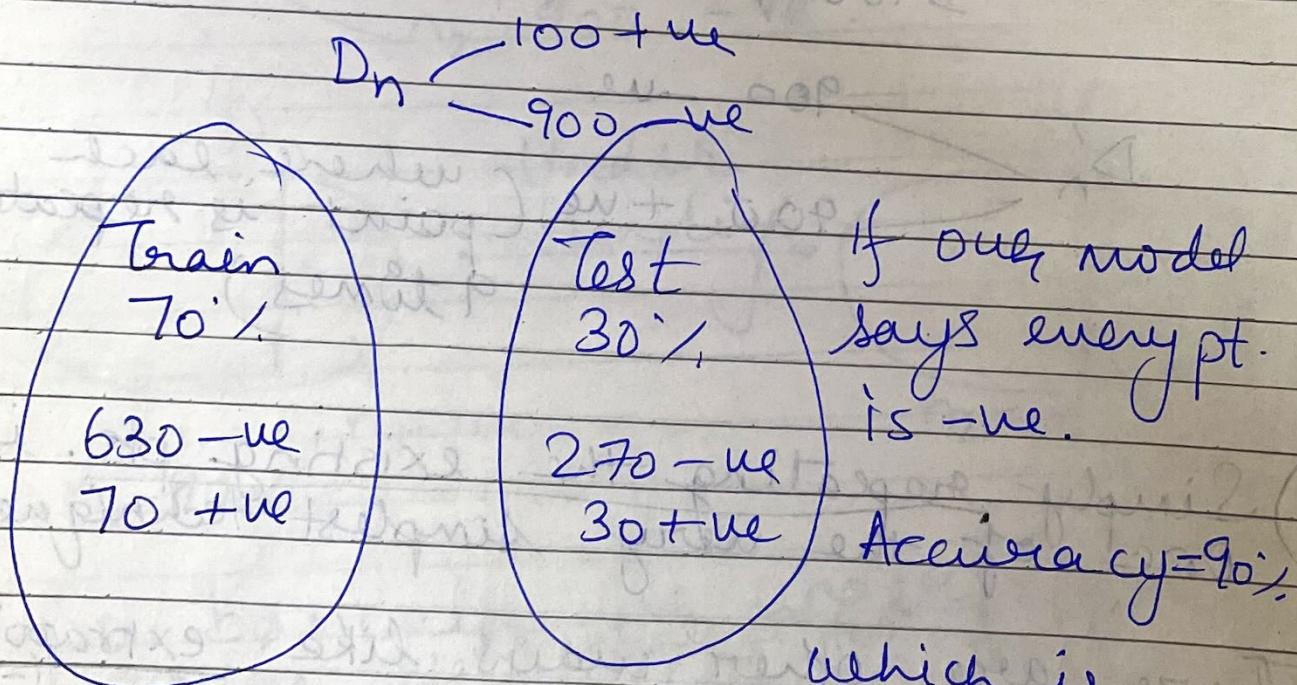
2) There are other ways like extrapolation where we create new synthetic points in the region bounded by the original set of the pts. of that class.

B)

Another way of implementing repetition of pts instead of manually repeating is by assigning class weights (In scikit-learn this method is not implemented for KNN)

In the above example, class weight of + : 9 (more weight to minority class), class weight of - : 1.

Another major problem of imbalanced dataset:



So accuracy is not the right measure for imbalanced data.

Imbalanced datasets are very much prevalent in real world scenarios
eg - medical : 10% cancer
90% non cancer.
ecommerce 10% buy
90% don't buy.

Video 3 Multi-class classification

There are classifiers like KNN which can be easily extended to do multiclass classification, there are other like Logistic Regression that can not do multiclass classification easily.

So, given a multiclass classification problem, can we convert it to a binary classification problem?

Let $y_i \in \{1, 2, 3, \dots, c\}$

This now can be handled by a simple binary classifier.

Similarly, for each of the classes we can build a ~~st~~ binary classifier like $2 \text{ or } \neq 2$ $\text{v} \text{y } 3 \text{ or } \neq 3$ $\text{v} \text{y } 4 \text{ or } \neq 4$ $\text{v} \text{y } 5 \text{ or } \neq 5$, etc.

$\text{v} \text{y } \text{c} \neq c$

So multiclass classification problem
= C binary classification problem.

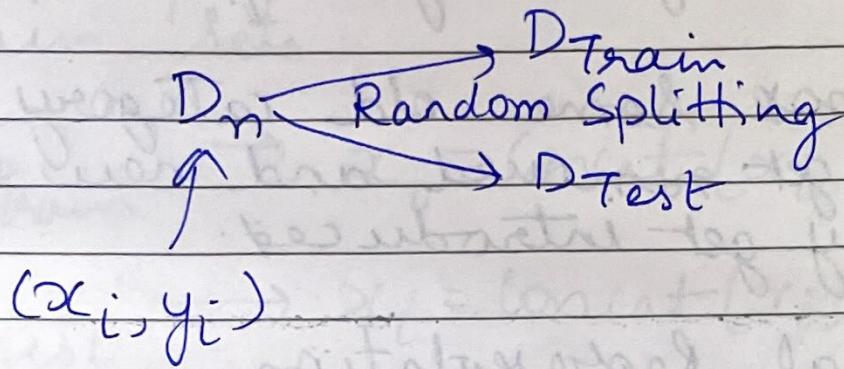
This technique is called One vs Rest

V4: KNN, given a distance or Similarity Matrix (logistic Reg. won't work when dist/sim matrix are given.)
There will be cases when we can get only the similarity Matrix S_{nxn}

$S_{ij} = \text{Similarity}(x_i, x_j)$, but can't get vector featurization of data points.
eg $x_i^i = \text{Paracetamol}$ can't be converted to vector.
But pair of chemical compounds can be written similarity.
We can convert similarity matrix to distance matrix as they are reciprocal of each other.

KNN works great as ultimately from the given D_n , we care about S_{nxn} only in KNN.

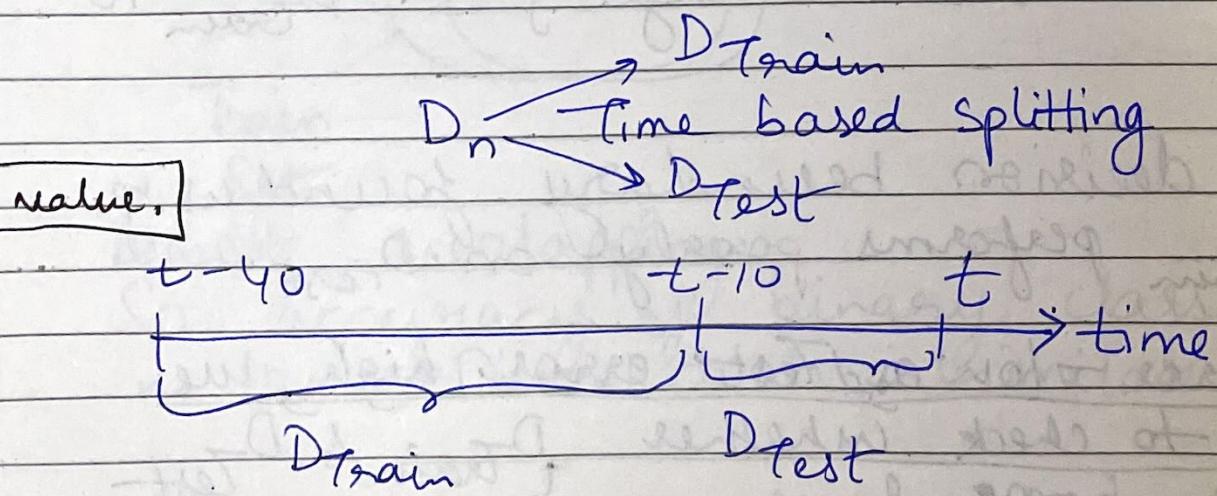
V5: Train and Test Set Differences



- In random splitting each point
- has equal probability for going to D_{Train} & D_{Test}

So both D_{Train} & D_{Test} will be similar.

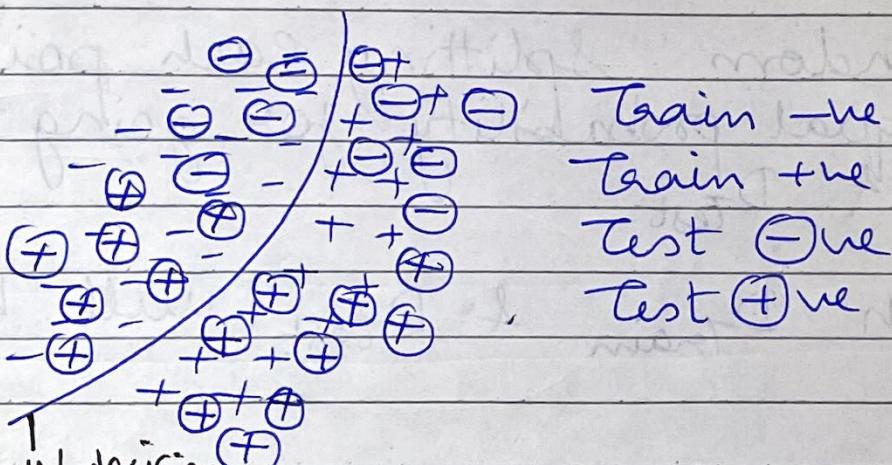
But in Time based splitting it is not the case.



In Time based splitting D_{train} & D_{test} could get very different.

e.g. → Suppose some old category of products get removed and new category get introduced.

Graphical Representation



D_{train} and D_{test} being fundamentally different causes the distribution of -ve and +ve. data change from D_{train} to D_{test} .

The decision boundary found using D_{train} performs poorly on D_{test} .

If train error - low and test error high, we need to check whether D_{train} & D_{test} come from same distribution or not.

How to determine this?

1) Time split the data to test and train sets.

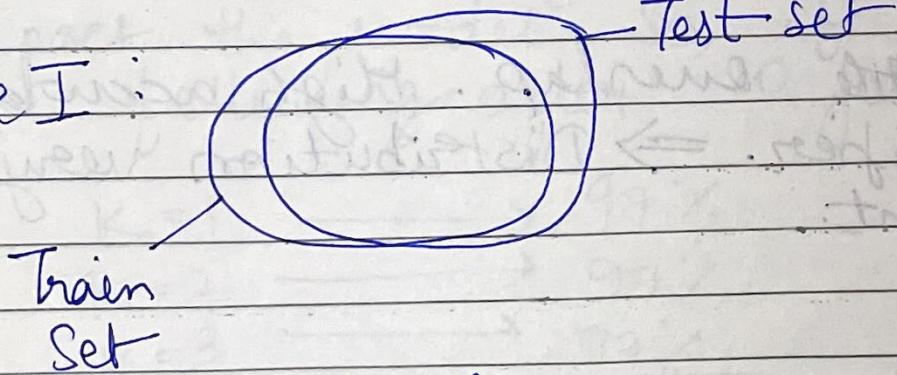
2) $\{D'_n\}_{\text{train}} \rightarrow x_i = \text{concat}(x_i, y_i), y_i = 1$
 $\{D'_n\}_{\text{test}} \rightarrow x_i = \text{concat}(x_i, y_i), y_i = 0$

This is to basically level all the training set points as 1 class and test set points as another. This gives D'_n .

Now train a binary classifier on

D'_n

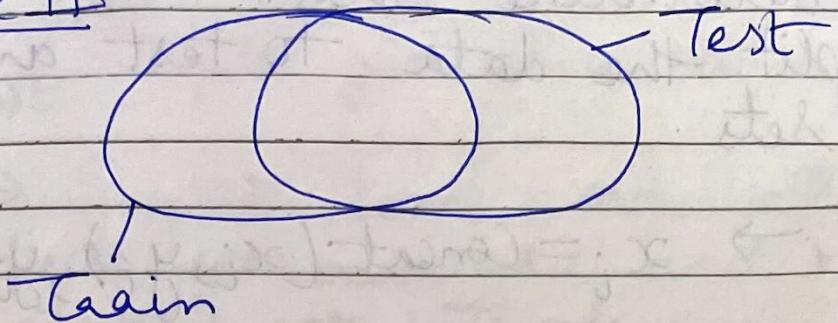
Case I :



Almost overlapping.

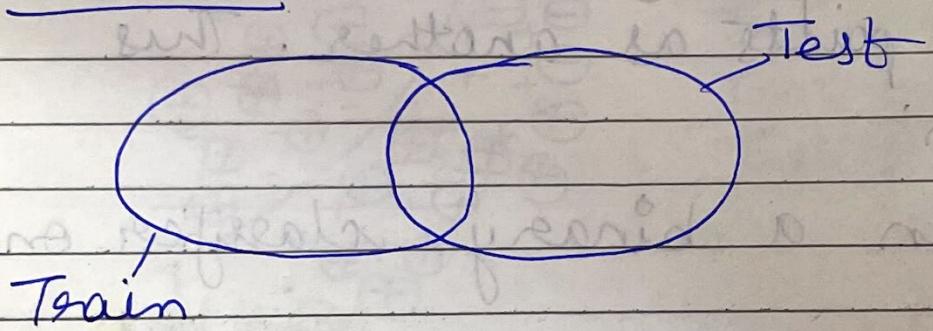
So accuracy of binary classifier is low \Rightarrow Distribution is similar.

Case II



Less overlap. Medium accuracy
of classifier \Rightarrow Distribution not
very similar.

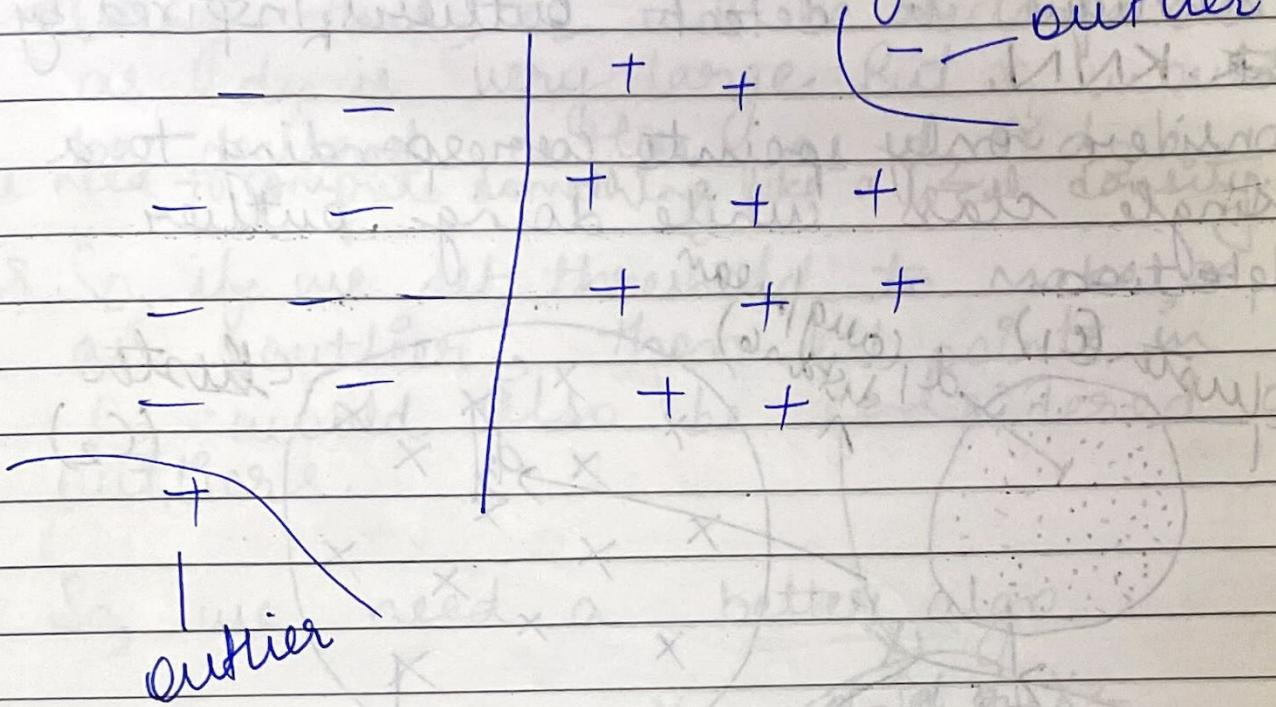
Case III



Very little overlap. High accuracy
of classifier. \Rightarrow Distribution very
different.

V6 Impact of Outliers

Let K be small say $K=1$.



When K is small, outliers can easily impact the model. So, if we get the following after 10 fold CV,

| | | |
|-------|-------------------|---------------------------------|
| $K=1$ | \longrightarrow | 97% |
| $K=2$ | \longrightarrow | 97% |
| $K=3$ | \longrightarrow | 97% |
| $K=4$ | \longrightarrow | 97% — less prone to outliers |
| $K=5$ | \longrightarrow | 95% — to outliers |
| $K=6$ | \longrightarrow | 92% |

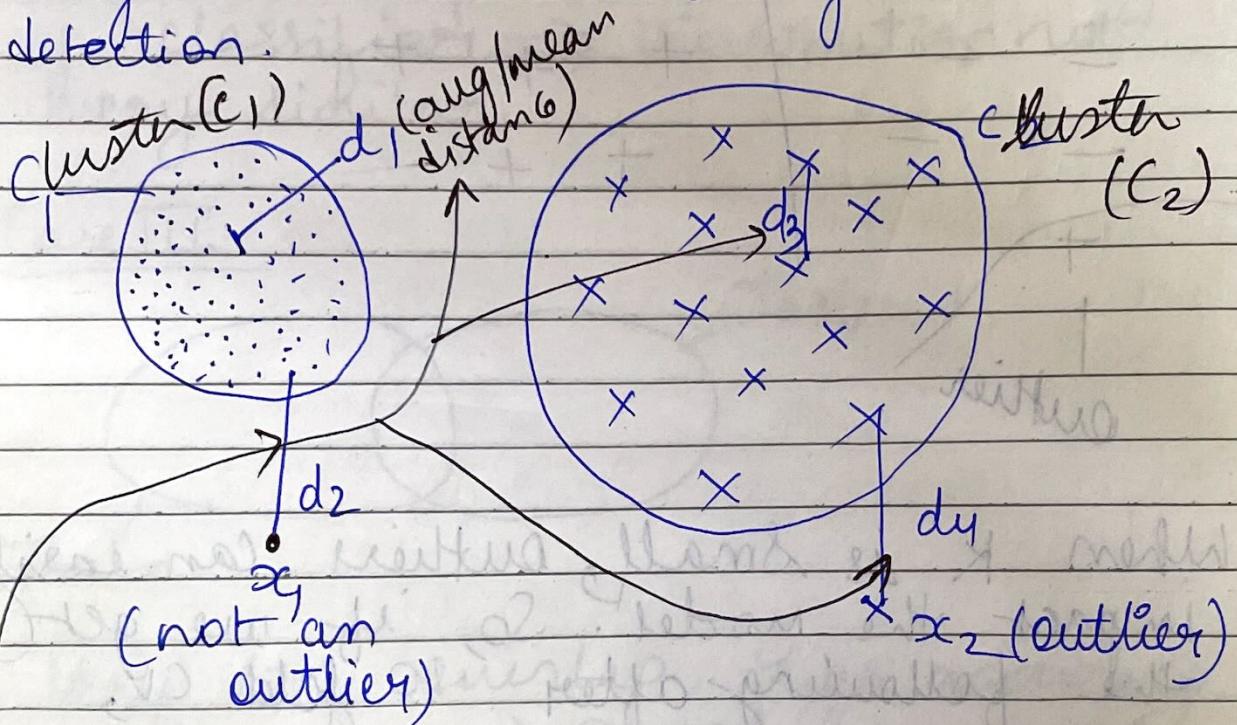
Always prefer to choose large K with high accuracy.

Methods to detect outliers:-

V7: Local Outlier Factor (Simple Solution:
Mean distance to KNN):

Technique to detect outliers, Inspired by
~~KNN~~. KNN.

Consider only points corresponding to a single class while doing outlier detection.



C_1 - Very dense cluster

C_2 - Sparser cluster.

x_1 & x_2 are outliers

Simple Solutions :-

- 1) For every pt x_i compute its KNN.
(mean distance from x_i to its KNN)
- 2) Calculate avg distance of x_i from its KNN
- 3) Sort x_i 's by the avg. distance

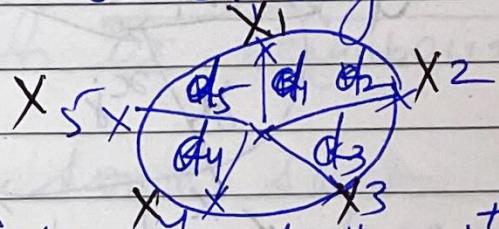
4) If average distance is high \rightarrow The point is outlier.

Using this, x_2 will be declared outlier as d_4 is very large. But x_1 wouldn't be declared outlier as $d_2 \approx d_3$. We need to compute something like a local density.

So, if we set threshold to make x_1 as outlier, then all points in C_2 would also be considered as outliers.

So, we need a better algo.

V8: k-distance



K-distance (x_i) = distance of the k^{th} nearest neighbour of x_i from x_i .

5th distance (x_i) = d_5 ; 1st distance (x_i) = d_1 .

Neighbourhood of x_i [$N_k(x_i)$] = Set of all the points that belong to KNN of x_i which can be more than 'k' in cases of tie.

$$\text{N}_{k=5}(x_i) = \{x_1, x_2, x_3, x_4, x_5\}$$

$$\text{N}_{k=3}(x_i) = \{x_1, x_2, x_3\}$$

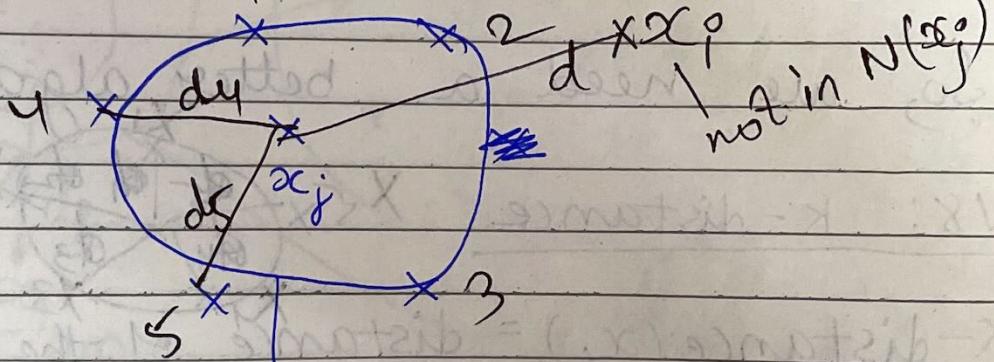
V9: Reachability-Distance (A, B)

reachability distance $(x_i, x_j) =$

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

distance of
 k^{th} NN of x_j
from x_j

actual
distance b/w
 x_i and x_j



Neighbourhood of x_j ($N(x_j)$)
 $K=5$

Neighbourhood of x_j

if $x_i \in N(x_j)$

then $\text{reach-dist}(x_i, x_j) = k\text{-distance}(x_j)$

else if $x_j \notin N(x_j)$

$\text{reach-distance}(x_i, x_j) = \text{distance}(x_i, x_j)$

$$\rightarrow \max(d_5, d_4) = d_5$$

$$\max(d_5, d) = d$$

using
neighbourhood

using
Reachability
distance

calculating
some
sort of
density

V10: Local Reachability-density (A)

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

denominator → no. of neighbours of x_i
 (And it duplicates same distance)

average reachability distance of x_i from its neighbour

$lrd(x_i)$ = inverse of average reachability distance of x_i from its neighbour.

It is calculating some sort of density of the surrounding around that point.

OR.

$$lrd(x_i) = \frac{|N(x_i)|}{\sum_{x_j \in N(x_i)} (\text{reach-dist}(x_i, x_j))}$$

points → distance

VII: Local Outlier Factor (A)

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

↓
avg lrd of
pts in the
neighbourhood of
 x_i

↑ only
lrd of x_i

$\text{LoF}(x_i)$ will be large if 1) $\text{lrd}(x_i)$ is small.
2) avg lrd of neighbours of x_i is large.

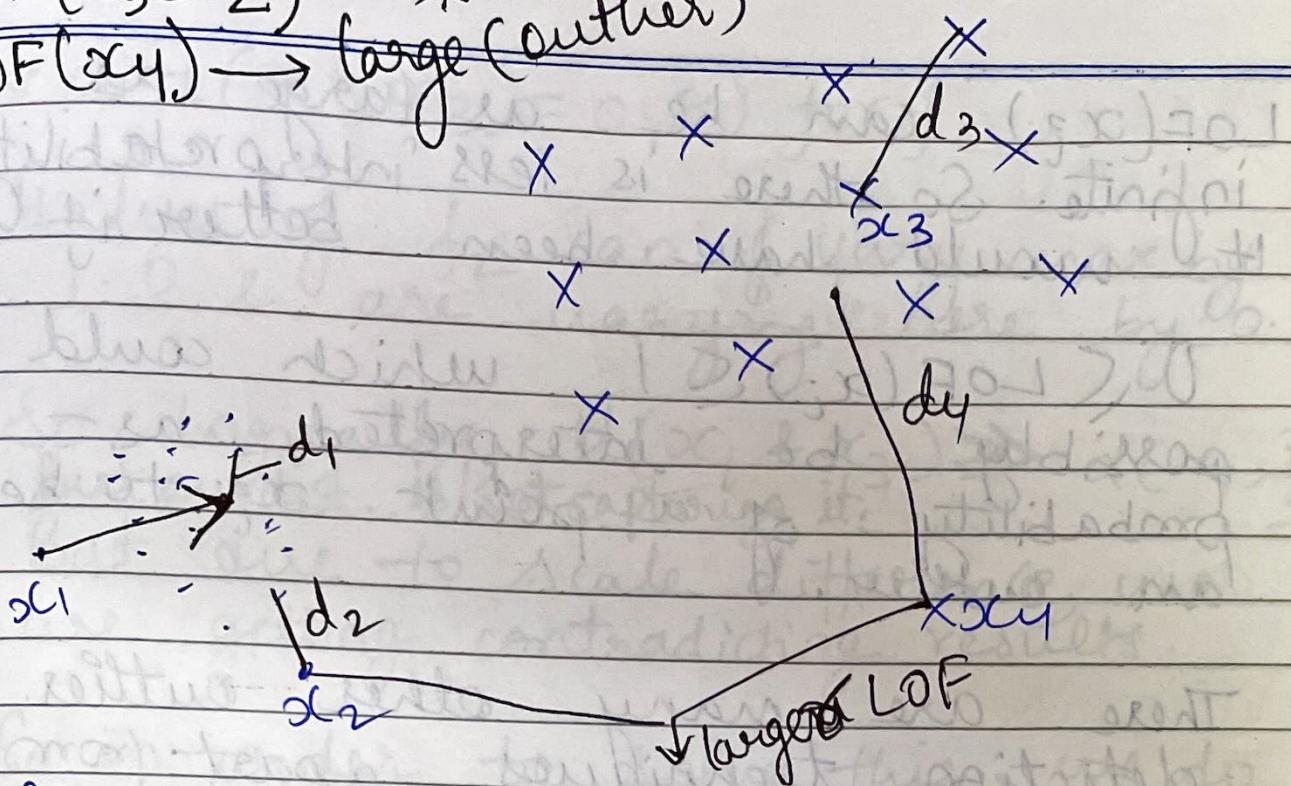
when $\text{LoF}(x_i)$ is ↑ → outliers
↓ → inliers

This means density of points in neighbourhood is high but density of points itself is low.



outlier

$\text{LOF}(x_1 \in C_1) \rightarrow \text{Small}$
 $\text{LOF}(x_2) \rightarrow \text{large (outlier)}$
 $\text{LOF}(x_3 \in C_2) \rightarrow \text{Small}$
 $\text{LOF}(x_4) \rightarrow \text{large (outlier)}$



So considering the figure, density of x_3 seems same as its neighbour same for x_4 . hence both are inliers.

But $\text{reach-dist} = d_1, d_2, d_3, d_4$. Also, $d_1 < d_2 < d_3 < d_4$ but density of x_1 and x_2 are very less as compared to their neighbours.

Hence outliers. $\text{ld}(x_1) = \frac{1}{d_1} > \text{ld}(x_2) = \frac{1}{d_2} > \text{ld}(x_3) = \frac{1}{d_3} > \text{ld}(x_4) \approx \text{ld}(x \in N(x_4))$, $\text{ld}(x_2) < \text{ld}(x \in N(x_2))$

There is no hard and fast rule to set threshold for LOF value beyond which there will be outlier

But if we know from domain knowledge that 5% of points are outliers, then we can sort and remove top 5% points.

for each pt x_i compute $\text{LOF}(x_i)$
pick points with highest $\text{LoF} \rightarrow$ outlier

$\text{LOF}(x_i)$ can be as large as infinite. So there is less interpretability. It would have been better if

$0 \leq \text{LOF}(x_i) \leq 1$ which could possibly be interpreted as probability of a point x_i to be an outlier.

There are many other outlier detection techniques apart from LOF.

V12: Impact of Scale and Column Standardization

| | f_1 | f_2 |
|--|----------|--------|
| | | |
| | 0 to 100 | 0 to 1 |
| | | |

Scale difference.

e.g. $f_1(0-100)$, $f_2(0-1)$

$$\begin{aligned} x_1 &: [23, 0.2] \\ x_2 &: [28, 0.2] \\ x_3 &: [23, 1.0] \end{aligned}$$

$$\text{eu. dist}(x_1, x_2) = \sqrt{5} \text{ for } f_1$$

The values of feature f_1 can vary max by 100 units, but are varying only 5 units.

$\text{eu-dist}(x_1, x_3) = 0.8$) for f_2

Values of f_2 can vary max by 1.0 & are varying here by 0.2

So logically, $d(x_1, x_3) \gg d(x_1, x_2)$
but eu-dist. is not capturing it.

But due to scale difference, we are getting contradicting results.

So in order to solve this, we do column standardization.

| f_1 | f_2 |
|-------|-------|
| a | |
| b | |

col std.

| f'_1 | f'_2 |
|--------|--------|
| a' | |
| b' | |

$$\text{mean } f_1 = \mu_1$$

$$\text{mean } f_2 = \mu_2$$

$$\text{std-dev } f_1 = \sigma_1$$

$$\text{std-dev } f_2 = \sigma_2$$

$$\text{mean } f'_1 = \text{mean } f'_2$$

$$\text{std dev}(f'_1) = 0$$

$$\text{std dev}(f'_2) = 1$$

$$a' = \frac{a - \mu_1}{\sigma_1}$$

$$b' = \frac{b - \mu_2}{\sigma_2}$$

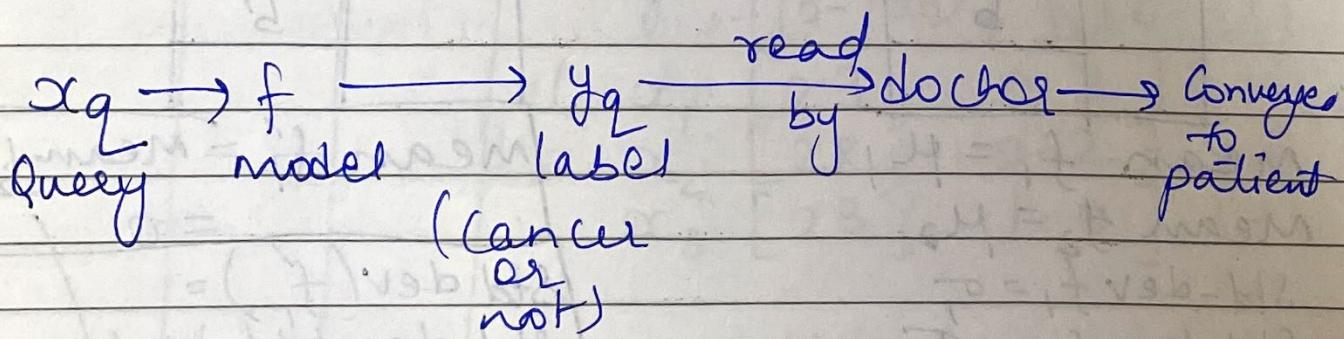
→ transformed feature

- 1) If we are using KNN with euclidean distance then 1st do col-std. Then apply KNN on transformed features.
- 2) There are models like decision tree that are scale independent also.

V18: Interpretability

Consider a medical domain problem be cancer detection.

Dataset $\rightarrow f$ (model)
(Training)



If the model just gives the class label, then it might not help the doctor much as to why the model choose to do this.

This kind of models are called black box model.

But if the model along with giving some class label, also gives some reasoning as to what made the model take the decision, then this reasoning is very useful for the doctor, such models are called interpretable models.

e.g. - Say x_9 is of dimension 10, where each feat. corresponds to result of a medical test.

Output of interpretable model →

$f_5 \rightarrow$ test 5 → value is v. high]

$f_8 \rightarrow$ " 8 → value low]

Hence
Patient
has
Cancer.

Doctor understands this additional info. as these are results of medical tests only.

let's take KNN for eg.

$x_q \rightarrow$ model $\rightarrow y_q = 1$.

Reasoning: All KNN of x_q are labelled 1. Hence $y_q = 1$.

Let's say d & k both are small.

{ $d=10$, $k=7$ - eg}

The model KNN ~~can't~~ ^{can} also give the 7-NN.

So, basically the doc. gets to know about 7 other patients' medical test results for 10 tests, which are similar to query patient.

Hence, doctor can compare the test results manually & get an indication as to why label = 1.

So, KNN is interpretable when d & k are small.

V14: Feature Importance and Forward Feature Selection

Suppose we are solving the problem of height prediction of an individual

We are given input features - weight, hair colour, hair length, skin colour, gender, country.

clearly all the features are equally important.

From domain knowledge, we have idea that weight has a significant impact on height. Similarly, gender, country also play deterministic roles.

If the ML model also gives feature importance it will be useful for us in understanding the model better. Hence increasing interpretability.

KNN doesn't give feature importance.

Logistic Regression & DT can give FImp.

Forward Feature Selection

Suppose we have $1k$ -dimension data. This can cause us to face curse of dimensionality. Where euclidean distance doesn't hold much meaning of similarity anymore.

So we need to reduce dimensions. Our task here is classification (assumed). PCA & t-SNE preserve distance. They don't care about classification.

So we make use of technique called Forward Feature Selection, which works irrespective of classification model.

- 1) Take the data and split to D_{train} & D_{test} . $\{D_n \rightarrow D_{\text{train}} + D_{\text{test}}\}$
- 2) Take single feature & train the model for each feature in feature set.
- 3) Select the feature giving maximum test accuracy.

Repeat above steps each time for each ~~one feature selection~~ feature selected and the new feature. This is called forward selection.

→ Don't care which model you get (advantage) Time complexity (D^2)

Why can't we just calculate accuracy for each single feature, sort them in descending order and then take 1st d' feature?

But its say f_{10} & f_6 are 1st and 2nd most accuracy feature.

But ~~its~~ it's possible that f_{10} & f_6 are very much co-related. Hence, adding it won't help much improvement.

Instead we want - Given that we ~~have~~ already have some feature which new feature adds the most value?

e.g. - Suppose $f_1, f_2, \dots, f_{10} \rightarrow 10$ features total.

After going through each feature say we get f_6 with highest accuracy.

Next we'll train the model with features f_1, f_{10} , i.e. $[1, 10]$ & it's 10.

Say we found f_8, f_9 have highest accuracy. features with lowest accuracy

Backward Selection - removing lowest accuracy.

Time Complexity \rightarrow each iteration :- Training and testing of models 1) d models 2) $d-1$ 3) $d-2$

V15: Handling Categorical and Numerical Features

1) Given a number

{ black, brown, red, gray }
1 2 3 4

But numbers have an inherent order.
i.e., $1 < 3$.

But black < red is absurd.

2) One-hot Encoding

| Bl | Br | R | G | ... |
|----|----|---|---|-----|
| | | | | |

Bl: [1|0|0|0|...]

Br: [0|1|0|0|...]

R: [0|0|1|0|...]

G: [0|0|0|1|...]

Binary vectors of size of number
of distinct elements.

These vectors can't be compared &
that's how these hair colors
should be.

Dis-advantage :-

$$\cancel{\text{#}} \text{ country} = 200.$$

So, we'll have to use a 200 dim vector just for a single feature which is the problem we faced for BOW encoding.

3) Suppose the task is predicting height of person.

So, in place of country, what if we replace it by avg. height of people in that country.

The idea is to replace categorical value with avg. output value for that category.

This is very much problem specific.

4) Domain knowledge

Let's say from domain knowledge we know that a distance from equator ↑, height ↑.

So using this knowledge, we can convert categorical feature to numeric feature.

features that
can be
ordered

5) Ordinal features :-

feature :- How do you like
Indian Food?

→ V. good . }
good . } logical ordering so, we
avg . } can assign them
bad . } numbers.
V. bad } numeric feature.

V16: Handling Missing values by Imputation

It is very common in real world.

It may be due to data corruption,
also called as ^{or} collection corruption.
or missing data,

1) Imputation (Replacement)

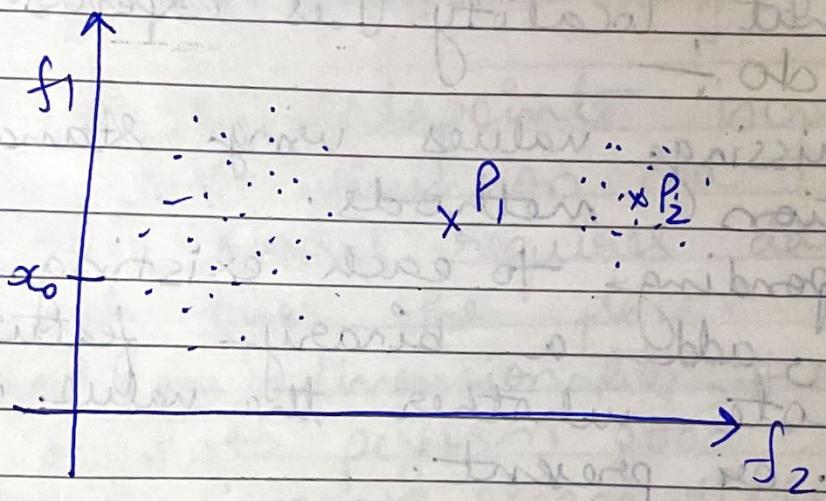
Take all the non-missing value in
that column & find its mean,
median or mode-^{most frequent value} and substitute those
to all missing values.

Sklearn.preprocessing.Imputer

2) Class based imputation

Take all the non-missing value in a column corresponding to a particular class and find their mean/median/mode and substitute it to missing values of that class.

Geometric Intuition



Let's say for x_0 , f_2 is missing. $y=1$ for x_0 . If we take mean of f_2 of all points then x_0 will be at P_1 and if we only consider pts. with $y=1$ then x_0 will be at P_2 .

Intuitively, P_2 makes more sense.

Creating a

3) New missing value feature

In this we treat missing values as source of information.

e.g. - Say blood group of some persons is missing.

It may imply that getting in blood group test in that locality is expensive.

Ways to do :-

① Fill missing values using standard imputation methods.

② Corresponding to each existing feature, add a binary feature to denote whether the value was missing or present.

4) Model based Imputation (Take significantly more time)

Usually KNN is used because of neighborhood concept.

We convert the missing values problem to standard classification and regression.

- ① Say for a bunch of points, f_3 missing
- ② Treat f_3 as output label (y).
- ③ $D_n \xrightarrow{\text{present}} D_{\text{train}}$ Now train a model & predict values.
 $\begin{cases} f_3 \\ \text{missing} \end{cases} \xrightarrow{} D_{\text{Test}}$

~~revised~~

VII: Curse of dimensionality

There are

It's a weird things that happen when d is very high that are not common sense.

For binary features, # of features = 3, then
of data points = $2^3 = 8$

of features = 10, then # of data points = 2^{10}
= 1024

So. no. of datapoints increase exponentially with dimension. And to perform well, the model requires data points from all over the place. Hence with increase in dimensionality, # of datapoints required to perform good classification/regression increase exponentially.

This is unlike Real valued (Regression) features.

a) Hughes phenomenon

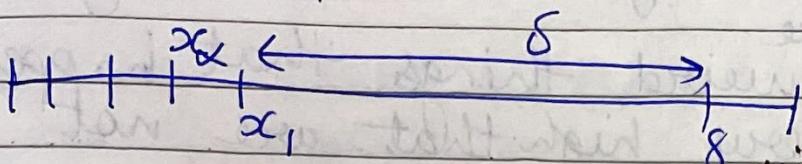
For fixed no. of data points, performance decreases with dimensionality & increases.

b) Distance Functions (Euclidean distance)

The intuition that we have of distances in 2-D, 3-D is not valid for high dimensional spaces.

$$\text{dist-min} = \min_{(x_j)} \left\{ \text{dist}(x_i^0, x_j) \right\}$$

Consider 1-D space \rightarrow n-random points



Consider x_i , $\text{dist-min}(x_i) = \delta_{\text{minimum}}$
from closest point to x_i to $x_i = \beta$

Similarly, distance $\text{max}(x_i) = \delta \max_{x_j \neq x_i} \{\text{dist}(x_i, x_j)\}$

So intuitively in 1D, 2D or 3D,

$\Rightarrow (\text{distance-max}(x_i) - \text{distance-min}(x_i)) \underset{\text{when } 1D, 2D, 3D}{\underset{\text{is greater than } 0}{\approx}} \text{distance-min}(x_i)$

But mathematically it can be proven that as

$$\lim_{d \rightarrow \infty} \frac{\text{distance max}(x_i) - \text{distance min}(x_i)}{\text{distance-min}(x_i)}$$

which implies
that

$\xrightarrow{\text{tends to}} 0$

$$\Rightarrow \text{distance-max}(x_i) \approx \text{distance-min}(x_i)$$

So in higher dimension, every pair of points are roughly equally distant for n random pts if

$$\text{distmax}(x_i) \approx \text{distmin}(x_i)$$

$$d(x_i, x_j) \approx d(x'_i, x'_j)$$

KNN ~~relies~~
on Euclidean
distance

in high
dim

makes no
sense

KNN doesn't
work well.

Soln

Use other distances like
Cosine-similarity. It is also affected
with increasing dimension, but not
much as compared to Euclidean.

If data has high dimension and
dense, curse of dimensionality is
high. And if data has high
dimensionality and sparse, then
impact of dimensionality lowers.

The above limit has been done
considering.

- 1) Uniform distribution of points
randomly. In sparse data, most
points will have no component
in many directions. Hence, less
impact.
- 2) Euclidean distance - If other distances
considered, it gives less impact.

3) Overfitting and Underfitting

As $d \rightarrow \infty \rightarrow$ overfitting \uparrow

Solution of high dimensions :-

- 1) Forward feature Selection :- classification oriented, use class labels (^{model}_{dependent})
- 2) Dimensionality reduction - PCA, TSNE -
Don't use class labels (^{model}_{Independent})
Not classification oriented.
- 3) Using cosine similarity instead of euclidean distance. for eg. in KNN.
- 4) Sparse representation instead of dense representation. for eg. in KNN.

V18: Bias Variance Tradeoff

It was studied in Theory of Machine Learning.

Generalization error = Error to be made on future unseen data.

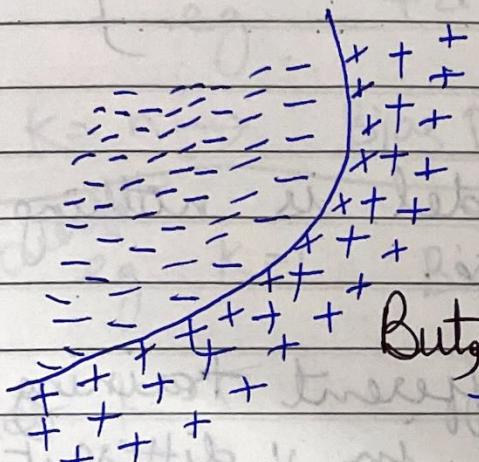
$$= \underbrace{\text{Bias}^2 + \text{Variance} + \text{irreducible error}}_{\text{Error that you can't reduce further for a given model.}}$$

Goal of ML is Gen. error ↓.

1) Bias error = Errors occurred due to simplifying assumptions.

High bias = Underfitting

Simplifying assumption



Modelling assumption
= "Plane" to separate +ve from -ve pts.

But, Curve makes sense in this data.

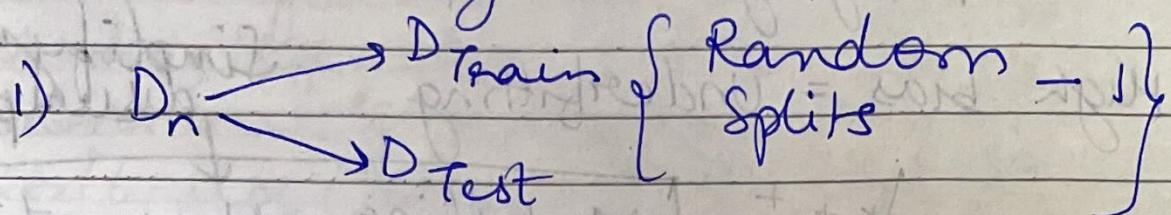
And there is no way to separate my points using plane. Hence bias error

→ In KNN, if $K = n$, we make simplifying assumption to assigning dominant class label to x_q .

→ High bias means underfitting which is you are making simplifying assumptions.

2) Variance error =

Variance means → how much a model changes as training data changes.



Geometrically a model is nothing but decision surface.

When you generate different training datasets, which result in v. different models or decision surface.

↓
Causes large changes in the model

↓
High variance model

↓
Overfit model.

→ As $K \uparrow$ variance \downarrow

↓
K in KNN

$$\text{Gen error} = \text{Bias}^2 \downarrow + \text{Var} \downarrow + \text{For. error}$$

↓ ↓
 no underfit no overfit

$K=1 \rightarrow$ Bias \downarrow Var \uparrow
 $K \uparrow \rightarrow$ Bias \uparrow slightly Var \downarrow drastically
 eg +2 -20

$K=n \rightarrow$ Bias \uparrow very high Var \downarrow

~~eg overfit~~ eg $K=1$ Bias = 10 + Var = 100 + ~~err~~ 113 $\leftarrow z=3$

let $K=5$ 12 + 10 + 3 \leftarrow
 $= 25$

$K=n$ 100 + 2 + 3 = 105
 Under fit

V19: Intuitive Understanding of bias-variance

Train error = difference (y_i, \hat{y}_i) on D_{train}

Test error = difference (y_i, \hat{y}_i) on D_{test}

1) High bias / ~~overfit~~ Underfitting

→ If Train error ↑ = bias ↑

kNN, $k=n$

→ If Train error ↓ = bias ↓

2) High Variance / Overfitting

→ If Train error ↑ & Test error ↑

= Overfit

= ↑ High Variance

→ If Training data changes slightly, should the model change this causes Variance ↑.

e.g kNN, $k=1$, High Variance.

V21: Best and Worst case of algorithm of KNN

1) when dimensionality is small (< 10) not large \rightarrow KNN is good.

When dimensionality is high we go into curse of dimensionality, interpretability \downarrow , runtime of kd tree/LSH \uparrow . KNN - not good
 \nearrow delay

2) Low latency System.
→ gives quick result. as it has high runtime complexity. KNN is not good

3) If you can find right distance measure,
~~so~~ then KNN works well.

→ For text data - cosine-distance.

→ for genomic data for eg,

If we are given similarity/distance matrix then use KNN.

Every ML technique, knowing whether it works best and worst is important, which will help in applying them better in future.