# Classification Algorithms in various situations
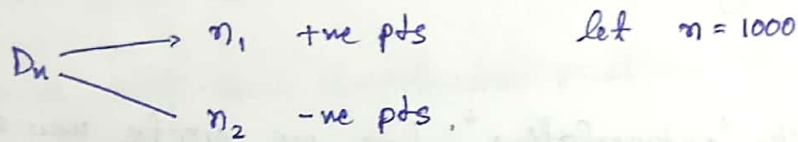
① **Imbalanced Vs Balanced. dataset**

consider example of 2-class classification.

$$D_n \begin{cases} \longrightarrow n_1 \quad \text{+ve pts} \\ \longrightarrow n_2 \quad \text{-ve pts} \end{cases} \qquad \text{let} \quad n = 1000$$

if $n_1 \approx n_2 \longrightarrow$ balanced dataset

ex→   $n_1 = 580$       $n_1 \neq n_2 \quad n_1 \approx n_2$.
$n_2 = 420$

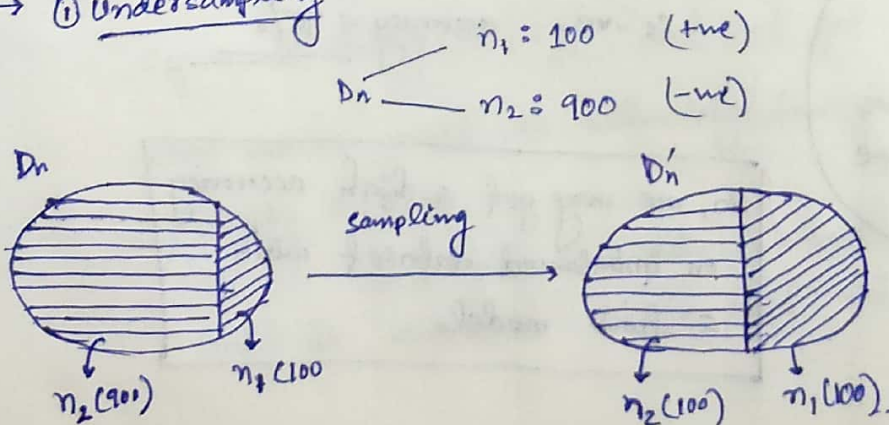if $n_1 << n_2$   or   $n_2 << n_1 \longrightarrow$ imbalanced dataset.

Ex→   $n_1 = 100$          $n_1 = 850$
$n_2 = 900$          $n_2 = 150$.

Imbalanced dataset may create problem.

✗ Just because there are more pts. of the majority class, it may have some advantage. (more biased towards it).

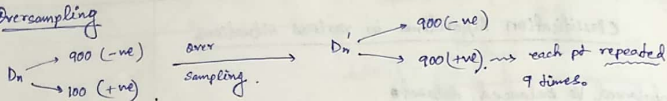So, how to work around the issue of imbalanced dataset?

→ ① **Undersampling**

$$D_n \begin{cases} n_1 : 100 \quad (\text{+ve}) \\ n_2 : 900 \quad (\text{-ve}) \end{cases}$$



$D_n$           sampling           $D_n'$

$n_2 (900)$       $n_1 (100)$           $n_2 (100)$   $n_1 (100)$.

So, basically we sample $n_1$ -ve pts 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

$$D_n \begin{cases} \to 900\,(-ve) \\ \to 100\,(+ve) \end{cases} \xrightarrow[\text{sampling}]{\text{over}} D_n' \begin{cases} \to 900\,(-ve) \\ \to 900\,(+ve) \to \text{each pt repeated} \\ \qquad\qquad\qquad 9 \text{ times.} \end{cases}$$
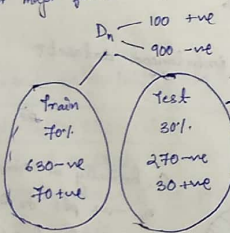
Simply repeating the existing pts is one of the very simplest techniques.

* There are other ways like "extrapolation", where we create new synthetic pts in the region bounded by the original set of pts of that class.

* Another way of implementing repitition of pts instead of manually repeating is by assigning class weights.

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

* Another major problem of imbalanced dataset. →

$$D_n \begin{cases} 100 \; +ve \\ 900 \; -ve \end{cases}$$

Train 70%          Test 30%
630 -ve            270 -ve
70 +ve             30 +ve

→ Even if our model says every pt is -ve, accuracy = 90%.

┌─────────────────────────────┐
│ So, we may get a high accuracy │
│ on imbalanced dataset with    │
│ a dumb model.                 │
└─────────────────────────────┘

Imbalanced datasets are very much prevalent in real-world scenarios.

ex→ medical :- 10% cancer
              90% non-cancer.

e-commerce:- 10% buy
             90% don't buy.

Ⓘ Multi-class classification

There are classifiers like K-NN which can be easily extended to do multi-class classification, there are others like logistic regression that can't not do multi-class classification easily.

So, given a multi-class classification problem, can we convert it to a binary classification problem.?

↳ let $y_i \in \{1,2,3,\ldots c\}$.

① $D_n \begin{cases} \to \{(x_i,y_i)\,|\; y_i=1\} \to (+)ve \; \left(\begin{array}{l}\text{this now can be handled by}\\ \text{a simple binary classifier.}\end{array}\right), \text{ or not.} \\ \to \{(x_i,y_i)\,|\; y_i \neq 1\} \to (-)ve \end{cases}$

Similarly for each of the c-classes, we can build c-classifiers (binary)

So, multiclass classification problem → 'c' binary classification problems

This technique is called 1 Vs rest

Ⓘ Given similarity matrix

There will be cases when we can get only the similarity matrix $S_{n \times n}$
$S_{ij} = $ similarity $(x_i, x_j)$, but can't get vector featurization of data points.

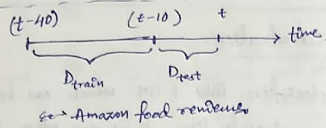KNN works great as ultimately from the given $D_n$, we care about $S_{nn}$ only in KNN.

We'll see how other models like logistic regression, NB etc could be extended to be used with S or not.

(iv) **Train and Test set Differences**

$D_n$ → $D_{train}$ , $D_{test}$ — random splitting
$(x_i, y_i)$.

So random splitting, each pt has equal probability for going to $D_{train}$ & $D_{test}$.
So, both $D_{train}$ & $D_{test}$ will be similar.

$(t-40)$  $(t-10)$  $t$ → time.

$D_n$ → $D_{train}$ , $D_{test}$ — Time based splitting.

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

Ex → Amazon food reviews.

In TBS, $D_{train}$ & $D_{test}$ could get very different.
Ex → Suppose some old category of products get removed & new category get introduced.

**Graphical representation**



x :— train (-)ve
▲ :— train (+)ve.
⊗ :— test (-)ve
Ⓐ :— test (+)ve.

※ $D_{train}$ & $D_{test}$ being fundamentally different causes the distribution of <u>(+)ve</u> data change from $D_{train}$ to $D_{test}$.

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

CV error – low
Test-error– high ⎰ → We need to check whether $D_{train}$ & $D_{test}$ come from same distribution. or not.

How to determine this?
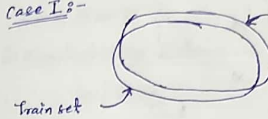→ ① Time split the data to test & train sets. ($D_{test}$ & $D_{train}$).
② $D_{train'}$ → $x_i = (x_i, y_i)$ , $y_i = 1$.
   $D_{test'}$ → $x_i = (x_i, y_i)$, $y_i = 0$. ⎱ → $D_n'$

That is, basically level all the training set pts as 1 class & test set pts as another. This gives $D_n'$.

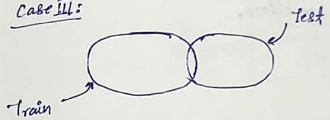③ Now train a binary classifier on $D_n'$.

**Case I :—**



Train set   Test set

Almost overlapping.
So, accuracy of binary classifier is <u>low</u>.
⇓
distributions similar.

**Case II :**



Test
Train

less overlap.
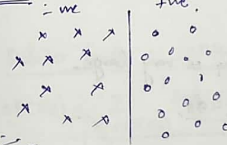Medium accuracy of classifier
⇓
distribut" not very similar.

**Case III :**



Test
Train

very little overlap.
High accuracy of classifier
⇓
distribut" very different.

(v) **Impact of outliers**

KNN



-ve region.
+ve region

(K=1) small

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

|     | accuracy |
|-----|----------|
| k=1 → | 97% |
| k=2 → | 97% |
| k=3 → | 97% |
| k=4 → | 97% → less prone to outliers |
| k=5 → | 95% |
| k=6 → | 92% |

For other models like logistic regression and naïve bayes, we'll see the impact later.

Next we'll know detecting outliers & remove them.

**Detecting Outliers** — consider only pts ~~while~~ corresponding to a single class while doing outlier detection.

Trying a simple solution —



So, we are currently only considering (-)ve class pts.

$C_1$: dense cluster.      $x_1$ & $x_2$ are outliers.

$C_2$: sparse cluster.

Simple sol$^n$ algo:—

① for every pt $x_i$, compute its $k$(say 5) NN.

② calculate avg. dist of $x_i$ fm its K-NN.

③ sort $x_i$'s by the avg. dist.

④ if avg-dist is high $\Rightarrow$ the pt is outlier.

using this, $x_2$ will be declared outlier as $d_4$ is very large.

But $x_1$ won't be declared outlier as $d_2 \approx d_3$.

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

So, we need a better algo.

---

Some terminologies

① K-distance $(x_i)$ → distance of $x_i$ to its $k^{th}$ nearest neighbour.

② Neighbourhood of $x_i$ [$N(x_i)$]?
→ set of all pts that belong to k-NN of $x_i$, which can be more than 'k' in case of tie.

③ Reachability distance →

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

$$\begin{cases} \text{if } x_i \in N(x_j) \\ \quad \text{then } \text{reach-dist}(x_i, x_j) = k\text{-dist}(x_j). \\ \text{else} \\ \quad \text{reach-dist}(x_i, x_j) = \text{dist}(x_i, x_j). \end{cases}$$

④ Local reachability density :- lrd $(x_i)$.

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

→ avg reachability distance of $x_i$ fm its neighbours.

$\text{lrd}(x_i)$ = inverse of avg reachability distance of $x_i$ fm its neighbours.

→ Its calculating some sort of density of the surrounding around that point.

---

Local Outlier factor :- LOF $(x_i)$ — inspired fm KNN.

$$\text{LOF}(x_i) = \underbrace{\cfrac{\displaystyle\sum_{x_j \in N(x_i)} \text{lrd}(x_j)}{|N(x_i)|}}_{\substack{\text{avg lrd of pts in} \\ \text{neighbourhood of } x_i}} * \cfrac{1}{\text{lrd}(x_i)}$$

LOF $(x_i)$ will be large if ① lrd $(x_i)$ is small

② avg. lrd of neighbours of $x_i$ is large.

i.e. density of pts in neighbourhood is high but density of pt itself is low.

So, $LOF(x_i) \longrightarrow$ large $\longrightarrow$ outlier $\longrightarrow$ density of $x_i$ smaller
$\qquad\qquad\qquad$ compared to its neighbours.
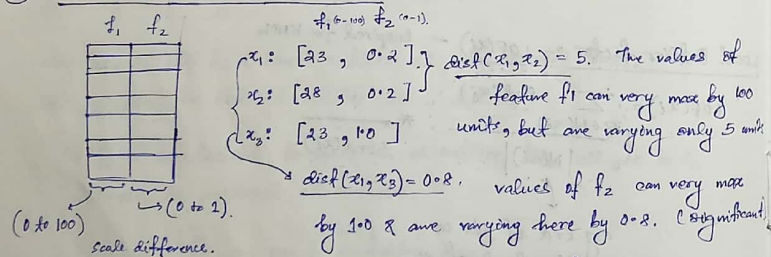$\qquad\qquad\searrow$ small $\longrightarrow$ in-lier.

So, considering the 1st figure,
density of $x_3$ seems same as its neighbours.
Same for $x_4$. Hence, both are inliers.

But density of $x_1$ & $x_2$ are very less as compared to their neighbours.
Hence, outliers.

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

* But if we know fm domain knowledge, that 5% of pts are outliers
then we can sort & remove top 5% pts.

* $LOF(x_i)$ can be as large as infinite. So, there is less interpretability.
It would have been better if $0 \leq LOF(x_i) \leq 1$ which could possibly
be interpreted as — probability of a point $x_i$ to be an outlier.

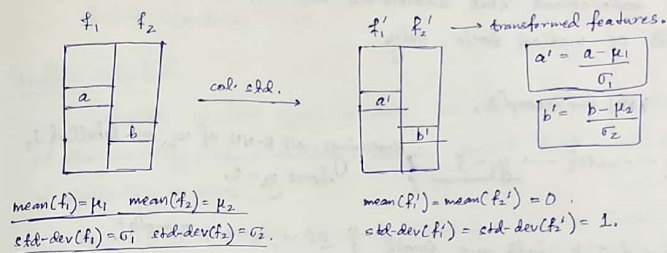There are many other outlier detection techniques apart fm LOF.

---

(VI) Scale and Column Standardization

$f_1 \quad f_2$

$f_1 (0-100) \quad f_2 (0-1)$

$x_1: [23, 0.2]$
$x_2: [28, 0.2]$  $\Big\}$ dist$(x_1, x_2) = 5$. The values of
feature f1 can very max by 100
$x_3: [23, 1.0]$  units, but are varying only 5 units

dist$(x_1, x_3) = 0.8$. values of $f_2$ can very max
$(0\ to\ 100) \quad \rightarrow (0\ to\ 1)$. by 1.0 & are varying here by 0.8 (significant)

scale difference.
So, logically $d(x_1, x_3) \gg d(x_1, x_2)$

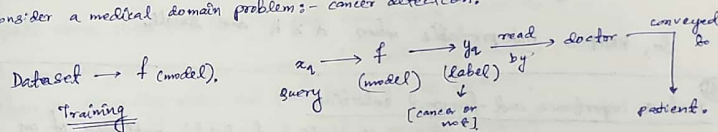But due to scale difference, we are getting contradicting results.

---

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

$f_1 \quad f_2$ $\qquad\qquad$ col. std. $\longrightarrow$ $\qquad$ $f_1' \quad f_2' \longrightarrow$ transformed features.

$\boxed{a' = \dfrac{a - \mu_1}{\sigma_1}}$

$\boxed{b' = \dfrac{b - \mu_2}{\sigma_2}}$

mean$(f_1) = \mu_1$ $\quad$ mean$(f_2) = \mu_2$ $\qquad$ mean$(f_1') = $ mean$(f_2') = 0$.
std-dev$(f_1) = \sigma_1$ $\quad$ std-dev$(f_2) = \sigma_2$. $\qquad$ std-dev$(f_1') = $ std-dev$(f_2') = 1$.

* If we are using KNN with euclidean distance then
1st do col-std. Then apply KNN on transformed features.

* There are models like decision trees that are scale independent also.

---

(VII) Interpretability

Consider a medical domain problem :- cancer detection.

Dataset $\longrightarrow$ f (model). $\qquad\qquad$ $x_q \longrightarrow f \longrightarrow y_q \xrightarrow[by]{read} $ doctor $\xrightarrow{conveyed\ to}$
$\qquad$ Training $\qquad\qquad\qquad$ query (model) (label)
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ [cancer or $\qquad\qquad\qquad\qquad\qquad$ patient.
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ not]

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 blackbox model.

But if the model along with giving some class labels, 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.

ex $\rightarrow$ say $x_q$ is of dimension b, where each feature corresponds
to result of a medical test.

o/p by interpretable model $\longrightarrow$
$\qquad\qquad\qquad$ $f_5 \longrightarrow$ test 5 $\longrightarrow$ value is very high $\Big\}$ Hence patient has
$\qquad\qquad$ & $\quad f_8 \longrightarrow$ test 8 $\longrightarrow$ value low $\qquad\qquad$ cancer.

Doctors understand this additional information as these are results of medical tests only.

Let's take KNN for example.

$$x_q \rightarrow model \rightarrow \underline{y_q = 1}$$

Reasoning: all K-NN of $x_q$ are labelled 1, hence $y_q = 1$.

Let's say $d$ & $k$, both are small $\{ d = 10, k = 7, \text{ example} \}$.

The model (KNN) can also give the 7 nearest neighbours.

So, basically the doctor gets to know about 7 other patient's 'medical test' result for 10 sets, tests, which are similar to query patient.

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

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

---

## Feature importance and forward selection

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

We are given input features — weight, hair_color, hair length, skin color, gender, country, etc.

Clearly, all the features aren't 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'll be useful for us in understanding the model better, hence increasing interpretability.

* KNN in its native form doesn't give feature importance.

---

* In models like logistic regression and decision tree, we can get feature importance

## Feature Selection

Suppose we have 1000-dim 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 whereas PCA & t-SNE preserve distances. 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 & split to $D_{train}$ & $D_{test}$. $\{ D_n \rightarrow D_{train} + D_{test} \}$.
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 new feature selection.

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

→ say $f_{10}$ & $f_6$ are 1st & 2nd most accuracy features. But 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 features, which new feature adds the most value.

e.g. → Suppose $f_1, f_2 \ldots f_{10}$ — 10 features total.
→ After going through each feature individually; say we get $f_{10}$ with highest accuracy.
* Next we'll train the model with features $f_i, f_{10}$ : $i \in [1,10]$ & $i \neq 10$.
* Say we found $f_6, f_{10}$ having highest accuracy. Similarly we'll proceed.
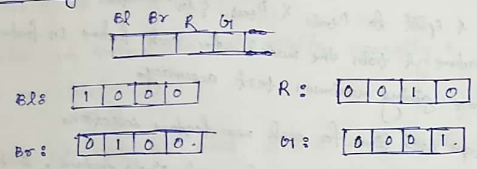
(IX) Categorical & ordinal features

Real valued data $\longrightarrow$ readily used ✓
Text $\longrightarrow$ BoW, tf-idf, word2vec.
categorical data $\longrightarrow$ ??

Ex → hair color → {black, brown, red, grey}.

① Give a number
{black, brown, red, grey}   But numbers have an inherent order, i.e.,
   1     2      3     4
                              1 < 3
         But black < red is absurd.

② One-hot encoding

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

Bl: [1 | 0 | 0 | 0]      R: [0 | 0 | 1 | 0]

Br: [0 | 1 | 0 | 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.

Disadvantage → #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 o/p value for
   that category.

---

this is very much problem specific.

ⓐ domain knowledge.
   let's say fm domain knowledge we know that as distance
   fm equator incr, height ↑.

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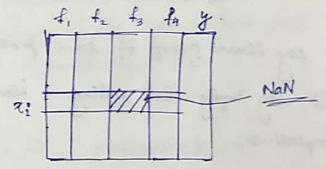numeric features.

ordinal feature

feature:- how do you like Indian food?
   → v.good.
      good.        } logical ordering. So, we can assign them numbers.
      avg.
      bad.
      v.bad.

ⓐ [5,4,3,2,1]   or  ⓑ [10, 6, 4, 3, 1]  –   Both are OK ✓ as long as order
                                              is maintained.

---

Missing value: imputation
   ↳ very common in real world.
   due to ⟶ data corruption
           ⟶ collect^n.



x: ▨                     NaN

① Imputation
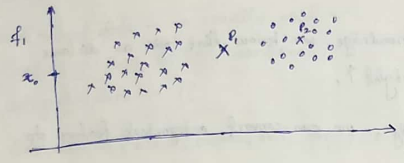   ↳ mean/median/mode.
   Take all the non-missing value in that col. & find its mean/med/mode
& substitute this to all missing values.

      sklearn. preprocessing. Imputer

② class based imputation
   Take of all the non-missing value in a col corresponding to a particular
class & find their mean/median/mode & substitute it to missing values
   of that class.

Geometrical intuition



$\times$ — $y=0$
$\circ$ — $y=1$

Let's say for $x_0$, $f_2$ is missing. $y=1$ for $x_0$.

* If we take mean of ~~all pts~~ $f_2$ of all pts, then $x_0$ will be at $P_1$ & if we only consider pts with $y=1$ then $x_0$ will be at $P_2$

* Intuitively, $P_2$ → makes more sense.

Ⓑ New missing value feature
for this we treat missing values as source of information

Ex — say blood group of some persons is missing.
It may imply that getting a blood test in that locality is expensive.

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.

Ex

| $f_1$ | $f_2$ | $f_3$ | $f_4$ | $f_1'$ | $f_2'$ | $f_3'$ | $f_4'$ |
|---|---|---|---|---|---|---|---|
| $x_1$ ✓ | ✓ | ~~NaN~~ | ✓ | 0 | 0 | 1 | 0 | ✓

④ model based imputation
└ usually KNN used

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

→ say for a bunch of pts, $f_3$ missing.
→ treat $f_3$ as o/p label (y).
→ $D_n$ —$f_3$ present→ $D_{train}$   Now train a model & predict values.
        |
        └ —$f_3$ missing→ $D_{test}$.

---

Curse of dimensionality → (wiki article).
↳ weird things that happen when $d$ be very high.
   that are not common-sensical.

for binary features, #features $=3$ → # data pts $=2^3$.
                  #features $=10$ → # data pts $=2^{10}$

So, no of datapts inc↑ exponentially with dim.
And to perform well, the model requires data pts fm all over the place.
Hence, with inc↑ dimensionality, #datapts required to perform good classification/regression increases exponentially.
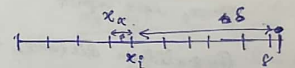
① Hughe's phenomenon
   For fixed ~~no~~ no. of data pts, performance↓ with dim↑.

② Distance functions (euclidean dist).
        The intuition that we have of distances in 2-D, 3-D is not valid for high dimensional spaces.

consider 1-D space → n-random pts

consider $x_i$.   dist_min($x_i$) = dist fm closest pt to $x_i$ to $x_i$ = B.
        similarly, dist_max($x_i$) = S.

So, intuitively

$$\frac{dist\_max(x_i) - dist\_min(x_i)}{dist\_min(x_i)} > 0 \quad \text{when } 1D, 2D \text{ or } 3D.$$

But mathematically it can be proven that as

$$\left[ \lim_{d \to \infty} \frac{dist\_max(x_i) - dist\_min(x_i)}{dist\_min(x_i)} \longrightarrow 0 \right]$$

$$\Downarrow$$

$$dist\_max(x_j) \approx dist\_min(x_i).$$

So, in higher dimensions, every pair of pts are roughly equally distant.

$$d(x_i, x_j) \approx d(x_i', x_j').$$

KNN $\xrightarrow[on]{relies}$ euclidean distance. $\xrightarrow[dim]{in\ high}$ makes no sense

Sol$^n$:- use other distances like cosine-similarity. Its also affected with increasing dim, but not ~~so~~ much as compared to euclidean.

* If data has high dim and dense $\longrightarrow$ curse of dim high↑.
* If data has high dim and sparse $\longrightarrow$ impact of dim lower↓.

The above limit has been done considering
1. uniform distribution of pts randomly. In sparse data, most pts ~~alo~~ will have no component in many directions. Hence, less impact
2. euclidean distance. - If other distance considered, less impact

③ Overfitting and underfitting
$$d↑ \longrightarrow \text{overfitting}↑.$$

---

Sol$^n$ of high dimensions.—

① forward feature selection:- classification oriented.- use class labels.

② dim-reduction:- PCA, t-SNE — don't use class labels.
   └ not classification oriented.

③ using cosine similarity instead of euclidean distance.

④ sparse representation instead of dense representation.

---

**Bias variance trade-off**

Generalization error = Bias$^2$ + variance + irreducible error.
(error on future unseen data)
— (Proof in linear regression lesson).
error that you can't reduce further for a given model.

Bias error $\longrightarrow$ error occured due to simplifying assumptions.



curve makes sense.
If we simplify the assumption to a line/plane → introduces error
↓
bias error.

Bias error is also called underfitting.

In KNN, if k=N, we simplify the assumption to assigning dominant class label to '$x_q$'.

Variance error $\longrightarrow$ how much a model changes as training data changes.



We split dataset randomly twice.
Assume k=1

① $D_n \xrightarrow{\quad} D_{train1}$ {Random $D_{test1}$ split -1}

In $D_{train1}$, $x_1, x_2, x_7$ —
$x_3, x_1, x_5, x_6$ ⊗ — In $D_{train2}$

$\rightarrow$ these are decision boundary for $D_{train1}$.

{Geometrically a model is nothing but decision surface}.

x3, x4, x5, x6 — $D_{train2}$

So, just by randomly splitting 2 times, we get 2 very different decision surfaces (model). Hence high variance for $k=1$, overfit model

---

Ultimate goal → low generalization error.

Gen. error = $Bias^2 \downarrow$ + var. $\downarrow$ + irreducible error.

→ no underfit.   → no overfit.

$k=1 \longrightarrow$ bias $\downarrow$   var $\uparrow$

$k=\infty \longrightarrow$ bias $\uparrow$   var $\downarrow$   } —— Balance needed.

$k=1$ ;   $10 + 100 + 3 \longrightarrow 113$.

$k=5$ ;   $12 + 10 + 3 \longrightarrow \underline{25}$ ~→ right balance in between.

$k=\infty$ ;   $100 + 2 + 3 \longrightarrow 105$

---

Train error = $diff(y_i, \hat{y_i})$ on $D_{train}$.

Test error = $diff(y_i, \hat{y_i})$ on $D_{test}$.

① high-bias

Train error $\uparrow$ ⇒ high bias.

Train error $\downarrow$ ⇒ low variance.

② high variance

Ⓐ Train error $\downarrow$ & Test error $\uparrow$.
   99.99 acc.        90% acc.
   0.01 error       10% error

Ⓑ training data changes slightly & causes significant model change.

---

Best and worst cases (KNN).

① dimensionality
  — not-large ✓
  — large — curse of dim (esp. euclidean dist).
       interpretability ↓.
       runtime of kd-tree/LSH ↑.

② low latency — quick results.
   KNN shouldn't be used → high runtime complexity

③ find right distance measure, then KNN works well.
   for text data, cosine-distance ✓
   If we are given similarity/distance matrix, then use KNN.

* For every ML technique, knowing where it works best and worst is important; which will help in applying them better in future.

---

Performance Measurement of Models

Some classification metrics

Accuracy = $\dfrac{\text{\#pts correctly classified}}{\text{Total \#pts in } D_{test}}$.

Accuracy $\in [0, 1]$.
(bad)   (good)

easy to understand measure.   Measured on test set usually.

100pts: { 60 +ve → 53 +ve, 7 -ve.
          40 -ve → 35 +ve, 5 -ve.

error = 12.   correctly = 88.

Accuracy = $\dfrac{88}{100} = 0.88$.

Problems with accuracy

① Doesn't work well with unbalanced data.