

CLASSIFICATION AND Regression MODELS

K-nearest neighbors

18.1 How "classification" works?

Amazon fire food review

Review \rightarrow Text \rightarrow Vector (Bow/tfidf) w2v)

classification:

we have 364K review (x_i) \rightarrow $\begin{pmatrix} x_i \\ y_i \end{pmatrix}$, +ve/-ve
Vector representation.

\rightarrow given a new review predict if the review is +ve & -ve

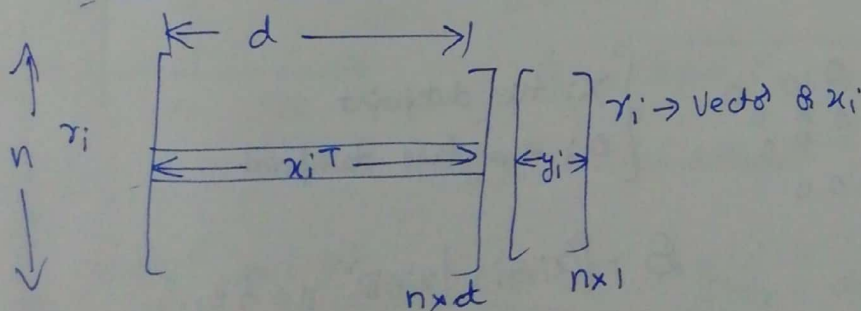
\rightarrow finding a function $y = f(x)$
+ve/-ve \rightarrow review text

classification algo working:

Training \rightarrow $\begin{matrix} \text{algo} \\ \boxed{f} \end{matrix}$ Training stage.
(x_i, y_i)
 $i = 1 \text{ to } 100K$

Testing/Evaluation $(x_q) \rightarrow \boxed{f} \rightarrow (y_q)$

18.2 Data matrix notation:



$\mathcal{D} = \{ (x_i, y_i)_{i=1}^n \}$ $x_i \in \mathbb{R}^d$, $y_i \in \{0, 1\}$

18.3 Classification vs Regression (Examples)

$$\mathcal{D} = \{(x_i, y_i)_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1\}\}$$

Amazon food reviews

There are two classes (2 class classification)
Binary classifier

MNIST:

$y_i \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\} \rightarrow 10$ class/multiclass classification

what if $y_i \in \mathbb{R}$

$\hookrightarrow y_i$ is no more part of a small finite set of classes

Ex:

$i = 1 \rightarrow 10K$

$x_i = \langle \text{weight, age, gender, race} \rangle$

$y_i = \text{height}$ (real number)

$y_i = f(x)$

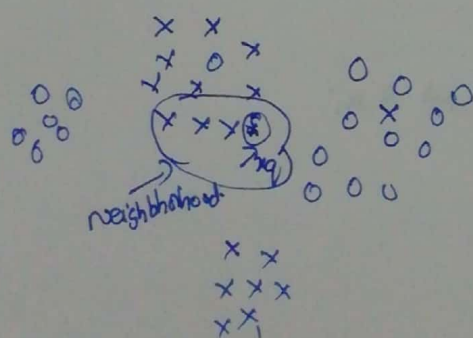
where y_i is a real number then it is called Regression algorithm.

18.4 k nearest neighbours Geometric intuition with a toy example

k-nearest neighbours (knn)

2D toy dataset

Let's take a case of binary classification



$\begin{cases} x: \text{+ve datapoint} \\ o: \text{-ve class datapoint} \end{cases}$

$$\mathcal{D} = \{(x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1\}\}$$

geom + since x_q is close to \bar{x} we can conclude that x_q is positive

steps

① Find K nearest points to x_q is \odot

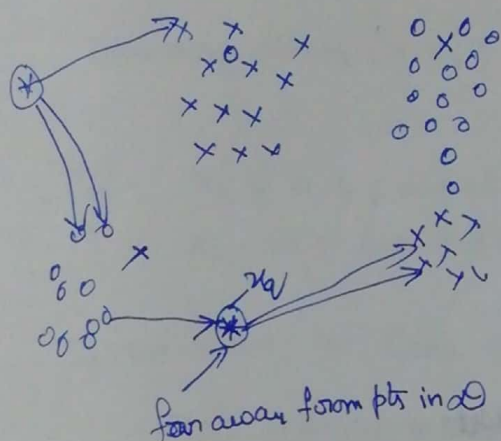
Let $K=3$ 3 nearest points are
 (x_1, x_2, x_3)
 y_1, y_2, y_3

②

② Take all the class labels (y_1, y_2, y_3)

Based on the majority we will declare the class

18.5 Failure cases of KNN



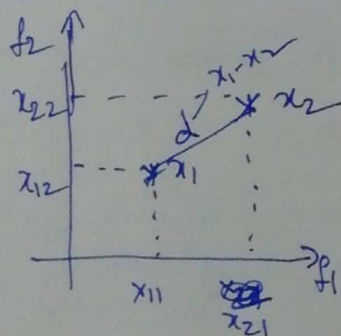
x o x o x o x o x o
 x x o x o x o o
 x o o o x x x o o x x o o
 x x o o x o x o o

↓
 → Jumble of +ve/-ve
 → Randomly spread
 → There is no useful information

→ As the x_q is far away, then it is
 Very good say don't know

→ There is useful information

18.6 Distance measure: Euclidean (L_2), Manhattan (L_1), Minkowski, Hamming.



$$x_1 = \begin{pmatrix} f_1 & f_2 \\ x_{11} & x_{12} \end{pmatrix} \quad x_2 = \begin{pmatrix} f_1 & f_2 \\ x_{21} & x_{22} \end{pmatrix}$$

d = len of shortest line from x_1 to x_2

$$d = \sqrt{(x_{21} - x_{11})^2 + (x_{22} - x_{12})^2} = \underbrace{\|x_1 - x_2\|}_{\text{length}} \quad \text{Pythagorean theorem. Euclidean distance.}$$

$$x_1 \in \mathbb{R}^d, x_2 \in \mathbb{R}^d$$

$$\|x_1 - x_2\| = \sqrt{\left(\sum_{i=1}^d (x_{1i} - x_{2i})^2 \right)}$$

$$\|x_1 - x_2\|_2 = L_2 \text{ norm}$$

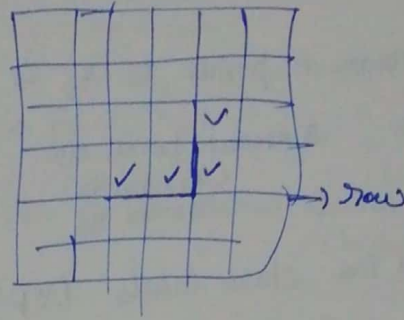
$$\|x_1\|_2 \rightarrow \text{dist of } x_1 \text{ from origin} \\ = \left(\sum_{i=1}^d x_{1i}^2 \right)^{1/2}$$

Manhattan dist:-

$$\sum_{i=1}^d |x_{1i} - x_{2i}|$$

L₁ norm of vector $(x_1 - x_2)$

$$\|x_1 - x_2\|_1$$



L₁ norm of x_1 is

$$\|x_1\|_1 = \sum_{i=1}^d |x_{1i}|$$

abs.

Generalization of L₁ & L₂ are L_p norms

L_p-norm \rightarrow Minkowski dist

$$\|x_1 - x_2\|_p = \left(\sum_{i=1}^d |x_{1i} - x_{2i}|^p \right)^{1/p}$$

If $p=2 \rightarrow$ Minkowski dist \Rightarrow Euclidean dist

If $p=1 \rightarrow$ " \rightarrow Manhattan dist

$$\|x_1\|_p = \left(\sum_{i=1}^d |x_{1i}|^p \right)^{1/p} \quad p \neq 0$$

\rightarrow Distances are always computed b/w two points
and norms are always for a vector.

$$\begin{aligned} \text{Euclid}(x_1, x_2) &= \text{L}_2 \text{ norm of } (x_1 - x_2) \\ &= \|x_1 - x_2\|_2 \end{aligned}$$

$$\text{manhattan}(x_1, x_2) = \|(x_1 - x_2)\|_1$$

$$\text{minkowski}(x_1, x_2) = \|(x_1 - x_2)\|_p$$

Hamming Distance

→ This is used in text processing when we have a boolean vector

$x_1, x_2 \rightarrow$ boolean vector \rightarrow Binary Bow

$$x_1 = [0, 1, 1, 0, 1, 0, 0, \dots]$$
$$x_2 \quad [1, 0, 1, 0, 1, 0, 1, \dots]$$

Hamming distance(x_1, x_2) = # location/dimensions where binary vectors differ

Hamming distance $(x, x_2) = 4$

Ex 2 $x_1 = a b c a d e f f g h i k$
 $x_2 = a c b a e d e g f h i k$

If we have a Gene codes/seq of AGTC character

$$u_1 = AAGTC TCAG -$$

$x_2 = \text{AGATCTCCA} \dots$

18.7 Coline distance & Coline similarity

Similarity vs distance

→ x_1, x_2 distance increases, similarity decreases

→ 11 de Greaser 11 11 Greasy

$$\cos \text{dist}(x_1, x_2) = 1 - \cos \text{sim}(x_1, x_2)$$

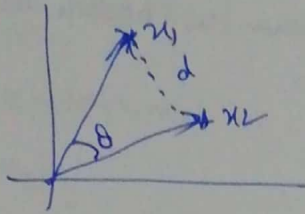
$\cos(\text{angle})$ is b/w $[-1, 1]$

→ Cos siml of $(x_1, x_2) = +1$ when they are very similar

→ " = -1 " very dissimilar

Cosine Similarity

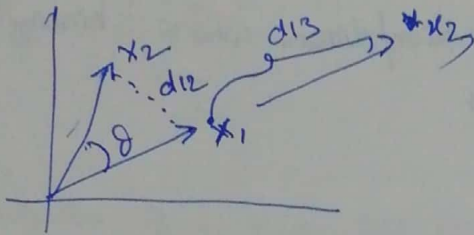
$d = \text{Euc. dist}$



$$\text{Cos-sim} = \cos \theta$$

where θ is angle b/w x_1 & x_2

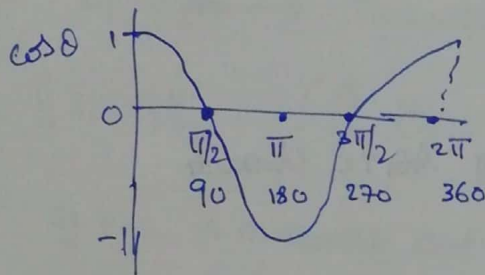
→ Difference b/w Euc. dist & Cos similarity.



$$\text{Cos-sim}(x_1, x_2) = \cos \theta$$

$$\text{Cos-sim}(x_1, x_3) = 1$$

$$\theta_{x_1, x_3} = 0^\circ$$



$$\text{Cos dist}(x_1, x_3) = 1 - \text{Cos-sim}(x_1, x_3)$$

$$= 1 - 1$$

$$= 0$$

Ex

$$\cos \theta = \frac{x_1 \cdot x_2}{\|x_1\|_2 \|x_2\|_2}$$

① If x_1 & x_2 are unit vectors (length=1)

$$\cos \theta = x_1 \cdot x_2$$

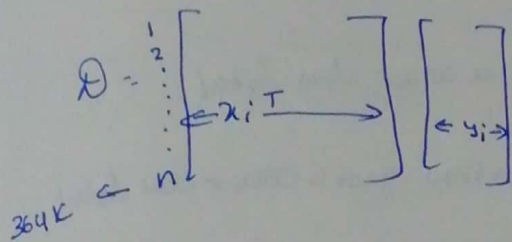
Relation b/w Euc dist vs Cos-sim

If x_1 & x_2 are unit vectors

$$\begin{aligned} [\text{Euc dist}(x_1, x_2)]^2 &= 2(1 - \underbrace{\cos(\theta)}_{\text{Cos dist}}) \\ &= 2(\text{Cos dist}(x_1, x_2)) \end{aligned}$$

18.8 How to Measure The Effectiveness of KNN

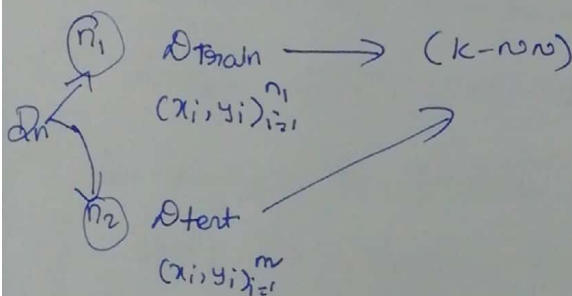
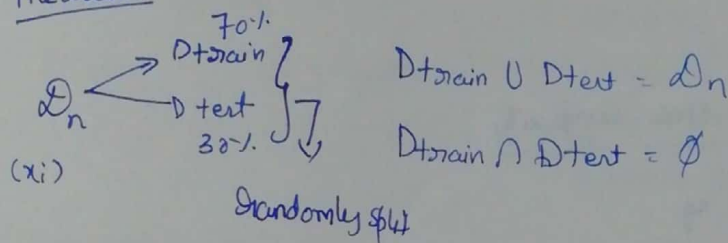
Amazon Finefood review



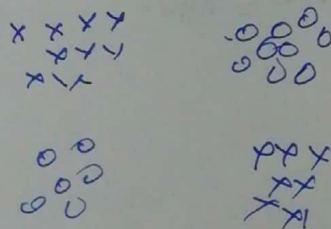
Problem: Given a new food review (x_q) what is its polarity (+ve/-ve)

$$x_q \rightarrow \text{Text}_q \rightarrow (x_q)$$

measure:



what is training in D_{train} ?



Each point x_i in $D_{\text{test}} = (x_i, y_i)_{i=1}^{n_2}$

$$x_q = (x_i)$$

for each point in D_{test} $x_q = \text{pt}$ case D_{train} & know to determine.

$$y \neq y_q \Rightarrow y_{pt}$$

$$Cnt \mp \pm 1$$

End

Count = # of for which $D_{train} + knn$ gave a correct class label

$$\text{Accuracy} = \frac{\text{Count}}{n_2} \rightarrow \# \text{ pts for which } D_{train} + knn \text{ gave a correct class label}$$

→
pts in D_{test}

$$0 \leq Acc \leq 1$$

Acc = 0.91 = 91% of times it is predicting correctly

18.9 Test/Evaluation time & space Complexity

We are given a query point x_q

$$x_q \rightarrow y_q$$

Input = $D_{train}, k, x_q \in \mathbb{R}^d$; output = y_q

$$knnpts = []$$

for each x_i in D_{train} : → $n_{pts} = d\text{-dim}$

- Compute $d(x_i, x_q) \rightarrow d_i$
- Keep the smallest k distances → $knnpts []$ $\rightarrow (x_i, y_i, d_i) \in \text{tuple}$.
- * k is small (5 or 10)

$$\text{Count}_{pts} = 0, \text{Count}_{neg} = 0$$

for each x_i in $knnpts$

if y_i is +ve

$$\text{Count}_{pts} += 1$$

else

$$\text{Count}_{neg} += 1$$

If $\text{Count}_{pts} > \text{Count}_{neg}$

$$\text{return } y_q = 1$$

$$\text{else } y_q = 0$$

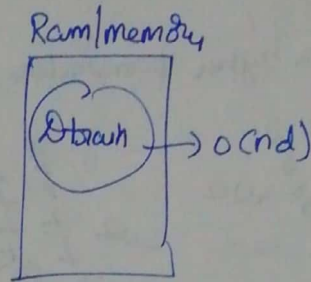
Time Complexity $\div O(nd) + O(1) + O(1)$

Amazon food reviews

$n \approx 364K$

$d \approx 100K$ (Bow)

300 (tfidf)

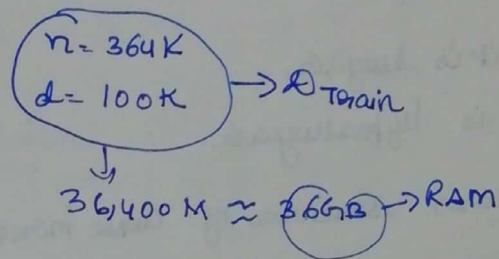


18.10 Known Limitations

Large space Complexity

Time Complexity $\div O(nd)$

Space $\div O(nd)$



Time Complexity

36 Billion Computation.

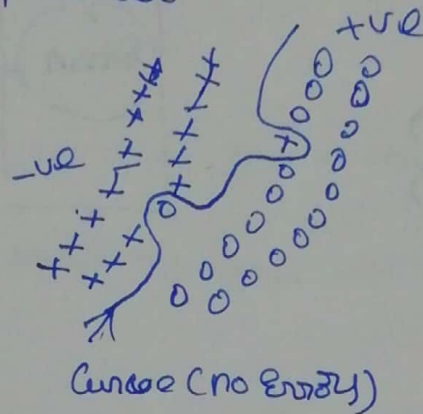
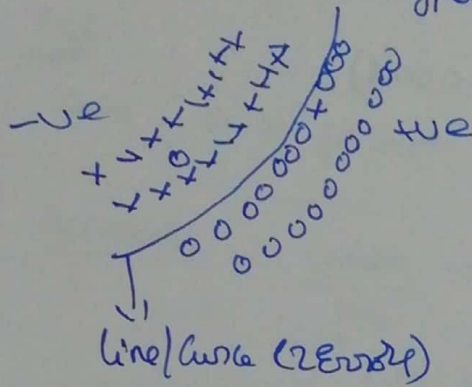
Review $\xrightarrow{1ms}$ +ve/-ve. (Internet application)

Low-latency \div Time it takes to predict y_q given x_q (It should be fast) Ex: Trading, Finance

- Simple Implementation we have known $\rightarrow O(nd)$, $O(nd)$
- Ppl don't use known because of Time & Space Complexity

18.11 Decision Surface for k-NN as k changes

→ Here k is a hyper parameter



→ These Curve +ve form -ve & vice versa are called decision surfaces

→ In 3D it is surface

→ In nD it is hypersurface.

→ As k increases smoothness of Curve increase.

Let's say $k = n$

when $k = 1, 2, \dots, (n)$

n = Total number of points

$n_1 = +ve (600)$

$n_1 + n_2 = n$

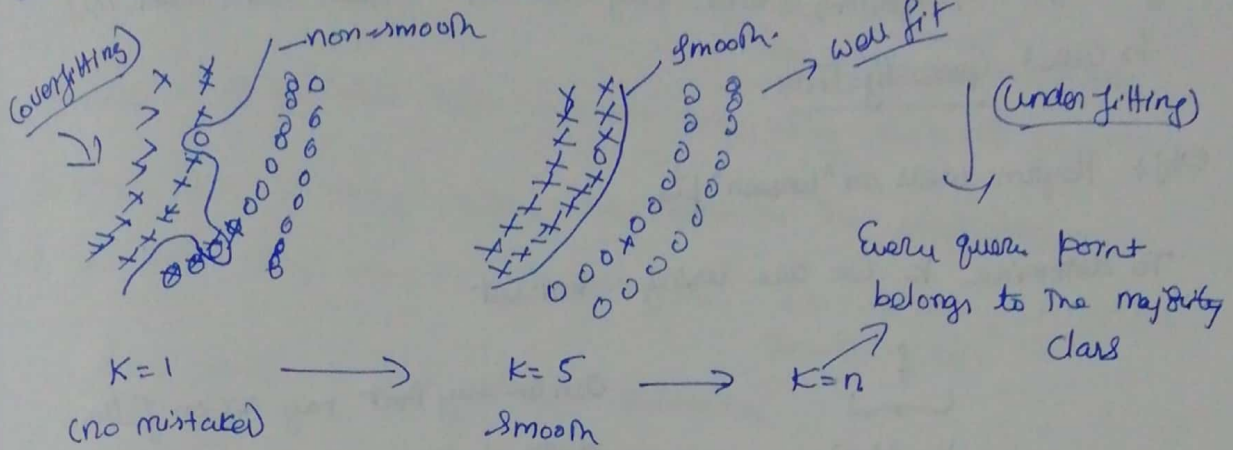
$n_2 = -ve (400)$

Let's say $n_1 > n_2$

If my $k = 1000$, 1000 neighbours ~~are~~ should be considered which means we have to consider all the points.

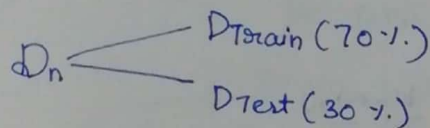
Hence here the majority class is n_1 our model will predict every input as positive only

18.12 Overfitting & Underfitting



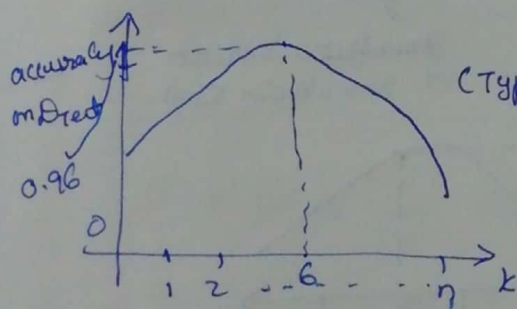
fitting because $f(x_q) = y_q$. The process of finding a function is called fitting.

18.13 Need for Cross Validation



one-idea:

	D_{train}	accuracy on Test
K=1	D_{train}	0.78
K=2	"	0.82
K=3	"	0.85

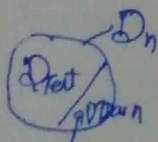


(Typically, need not be the same)

$k=6$, give me the best accuracy when using the D_{train} .

$k=6 \rightarrow (6-n) \rightarrow 96\%$

Using D_{train} & 6-fold on Amazon review dataset It got an acc of 96%.

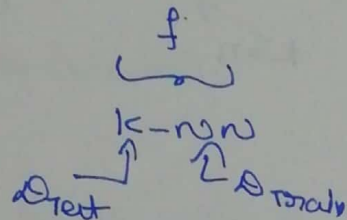


Obj: For a feature vector point x_q how accurately I can predict y_q

→ If a function/algorithm work very well on unseen data then this is called Generalization

Obj: Perform well on "unseen" pts

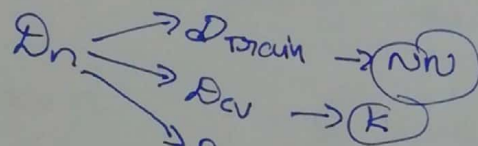
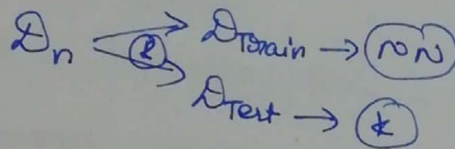
To determine k we are using Test set



Can we say that my acc on future data is 96%.

Ans: we cannot guarantee

To overcome this problem we use Cross Validation.



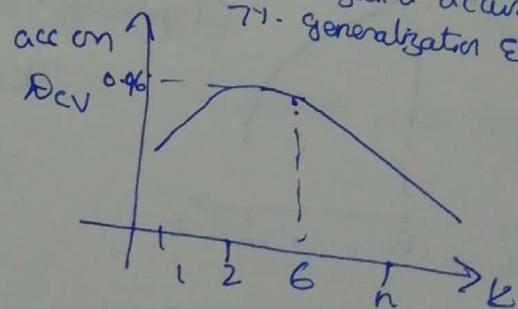
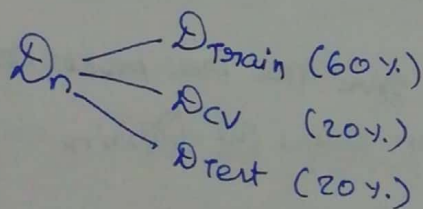
$f: k \rightarrow n$

→ did not use D_{test} to compute acc

→ x_i in D_{test} as x_q
→ $acc: 93\%$

→ Can we say that 6-nn has an acc of 93% on unseen data?

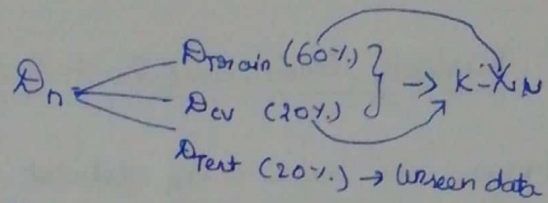
Yes



Generalization's accuracy
71. Generalization Error.

19.14

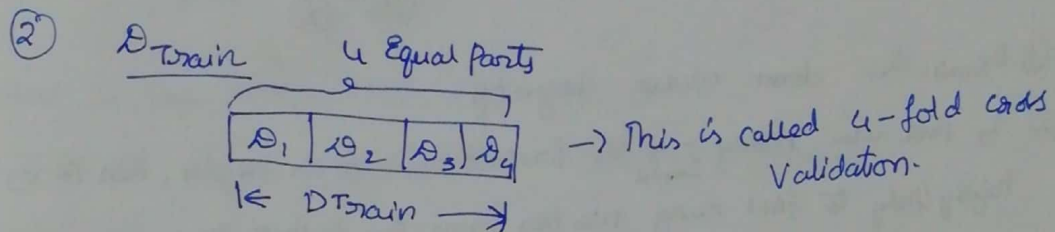
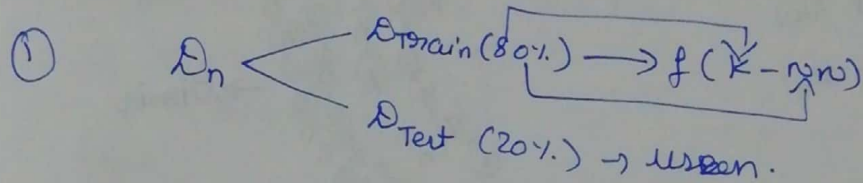
K-Fold Cross Validation



Problem : of the total data we are only using 60% to find run

→ more the Training data, the better is the algorithm.

→ Is there any way to combine the 60% & 20%, it can be done using K-fold cross validation.



		D_{train}	CV	acc. on CV
4 times	$k=1$	$D_1 D_2 D_3$	D_4	a_4
	$k=2$	$D_1 D_2 D_4$	D_3	a_3
	$k=3$	$D_1 D_3 D_4$	D_2	a_2
	$k=4$	$D_2 D_3 D_4$	D_1	a_1
} $avg(a_1, a_2, a_3, a_4) = a_{k=1}$				
4 times	$k=2$			
	\vdots			
4 times	$k=3$			

③ K' fold CV → $D_{train} \xrightarrow{\text{run}} k' \text{ in } k' \text{ times}$

→ what is the right K'

$K'=4$; $K'=10$; $K'=100$

rule of thumb : 10 fold CV

We are repeating multiple time, to compute the best K' in $k' \text{ times}$

Increases k' time when we are doing K' CV

18.15

Visualizing train, validation & test datasets

Imagine we have a big dataset D_n

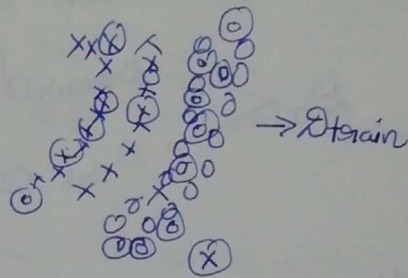
x : -ve class datapoint in D_{train}

o : +ve class datapoint in D_{train}

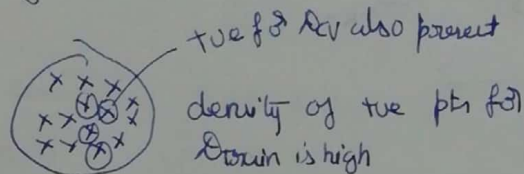
$T_{train}, cv, Test = (x_i, y_i)$
 \uparrow Input data point
 \uparrow class label

\otimes -ve point D_{cv}

\odot +ve point D_{cv}

Tips

- ① D_{train}, D_{cv} don't overlap perfectly
- ② If there are many + pts from D_{train} in a region, then it is highly likely to find many + pts from D_{cv} in that region.
- ③ If there are very few +ve/-ve pts in a region from D_{train} then it is very unlikely to find +ve/-ve from D_{cv} in that region

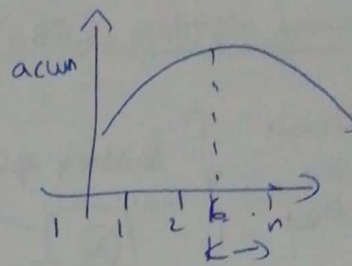
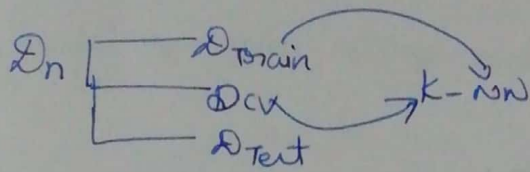
Intuitively18.16 How to determine underfitting & overfitting

k -fold on $D_{cv} \rightarrow$ best $k \rightarrow$ neither overfit
 \rightarrow underfit

⑥ How can we be sure about O.F & U.F

$$\text{accuracy} = \frac{\# \text{ correctly class pts}}{\text{Total \#pts}} = 0.93$$

$$\text{Error} = 1 - \text{accuracy} = 0.07$$



Training Error

what is the train error for 2-NN?

for each x_i in D_{train}

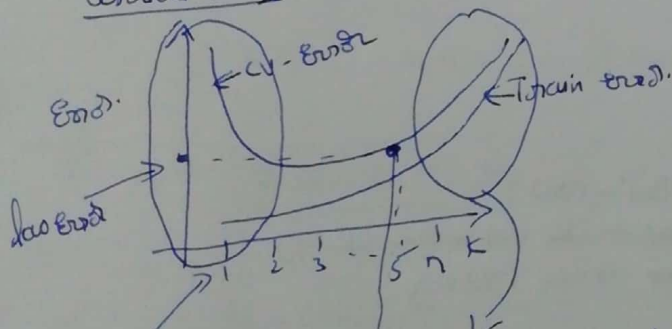
- find 2 nearest neigh to x_i from D_{train}
- majority vote to get the class label
- if $y_i == \text{class label}$
accurate

else
error

what is the ~~train~~ ^{D_{cv}} error for 2-NN?

→ same as above, calculate the error & accuracy

Consider $k=3$



$k=1$ overfit

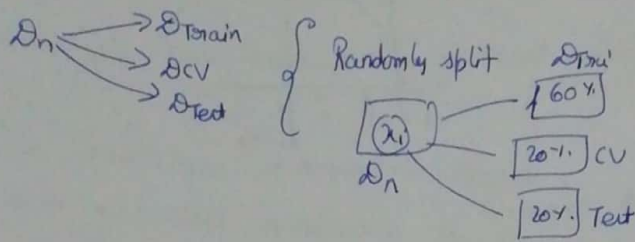
$k=n$ underfit

If k is in this region overfit
Train Error ↓
CV Error ↑

If k is in this region underfit
Train Error ↑
CV Error ↑

✓ Trade off
Train Error & CV Error will be close

18-17 Time based splitting (TBS)

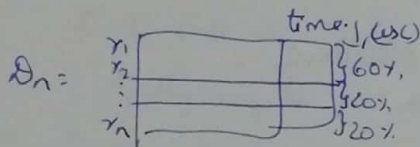


Probability of x_i to be part of
train is 60%.
cv is 20%.
Test is 20%.

For Amazon food reviews, time based splitting is better than Random splitting

$$D = \begin{bmatrix} f_1 & f_2 & \dots & f_d \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

In the actual dataset

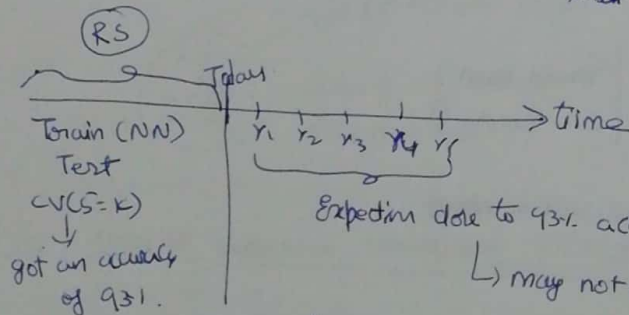


→ Sort the reviews on the time column in the
asc order.

Why TBS?

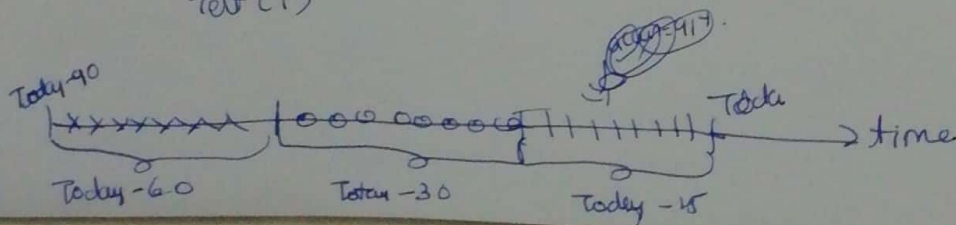
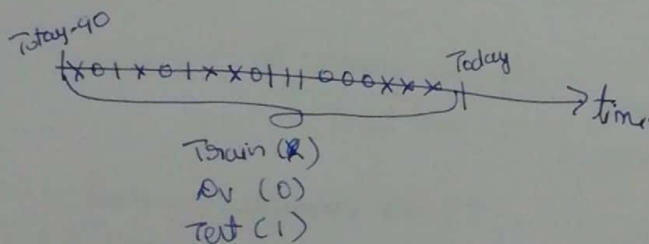
Let's say I used random splitting

$D_{train} = nms$
 $D_{cv} = k$ in $k-nn$ ($k=5$)
 $D_{test} \rightarrow acc$ (93%)



Deployed-like
(D_{train})
($k=5$)

* The new reviews may be different from D_{train}
* with time my products & review change



If I train my model 15 days back and I tested on last 15 days accuracy is 91%.

We can predict that when we predict the next 15 days data we can expect 91% of accuracy

When even time is available & things/behaviour & data change over time then time based splitting is preferable than random splitting

18.18 k-NN for regression

Classification $\mathcal{D} = \{(x_i, y_i)_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \{0, 1\}\}$
 $x_q \rightarrow y_q \rightarrow \text{class label}$

Regression $\mathcal{D} = \{(x_i, y_i)_{i=1}^n \mid x_i \in \mathbb{R}^d, y_i \in \mathbb{R}\}$
 $x_q \rightarrow y_q \rightarrow \text{number}$

① Given x_q , find k -nearest neighbours

$(x_1, y_1) (x_2, y_2) \dots (x_k, y_k)$

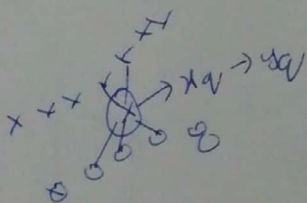
② $y_q \leftarrow \underbrace{y_1, y_2, y_3 \dots y_k}_{\substack{\text{all are real valued numbers.}}}$

$$y_q = \text{mean}(y_i)_{i=1}^k$$

$$y_q = \text{median}(y_i)_{i=1}^k \rightarrow \text{less prone to outliers}$$

* For some algorithm we can extend/modify classification to regression.

18.19 weighted k-NN



5-NN
 $k=5$

2 \rightarrow -ve } \rightarrow -ve
 3 \rightarrow +ve } \rightarrow +ve

$x_q \rightarrow$

$x_1, y_1, d_1 \rightarrow 0.1$	-ve
$x_2, y_2, d_2 \rightarrow 0.2$	-ve
$x_3, y_3, d_3 \rightarrow 1.0$	+ve
$x_4, y_4, d_4 \rightarrow 2.0$	+ve
$x_5, y_5, d_5 \rightarrow 4.0$	+ve

Given my $x_q \rightarrow x_1, -ve -0.1$
 $x_2, -ve -0.2$
 $x_3, -ve -1.0$
 $x_4, -ve -1.0$
 $x_5, -ve \rightarrow 4.0$

y_i distance w_i
 10
 5 } 15
 1
 0.5
 0.26 } 1.75

$w_i = 1/d_i$ (one way to find weights)
 $d_i \uparrow, w_i \downarrow$
 $d_i \downarrow, w_i \uparrow$

$\rightarrow w_i = \frac{1}{d_i}$ Simplest weight function

18.20 Voronoi diagram

\rightarrow Related to k -nn ($k=1$)

\rightarrow It rays

18.21 Binary Search Tree

k -nn + $O(n)$ If d is small \rightarrow Time Complexity
 k is small

$O(n)$ \rightarrow space Complexity

Time: $O(n) \rightarrow O(\lg(n)) \rightarrow$ Kd-trees

$n=1024 \rightarrow \lg(n)=10$

Applied in
 Computer graphics, geometry, Applied math

Binary search tree (BST)

- Data structure

\rightarrow Binary Search

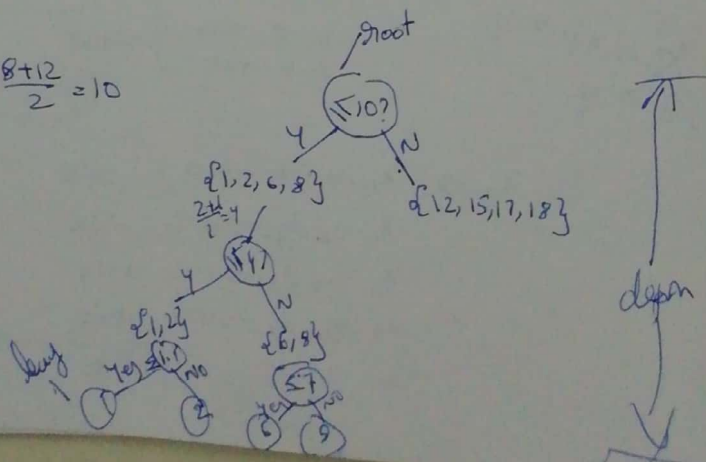
Prob: given a sorted array, find a number is present in array or not

1 2 3 4 5 6 7 8

A: [1, 3, 6, 8, 12, 15, 17, 18]

① Construct a BST

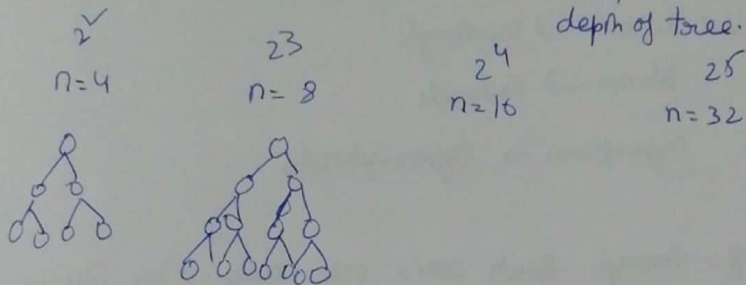
Take median: $\frac{8+12}{2} = 10$



② Find an element in an array
 $q = 15$

To find 15, with comparison we can find 15 exists or not

BST: Time Complex to search is $O(\log(n))$

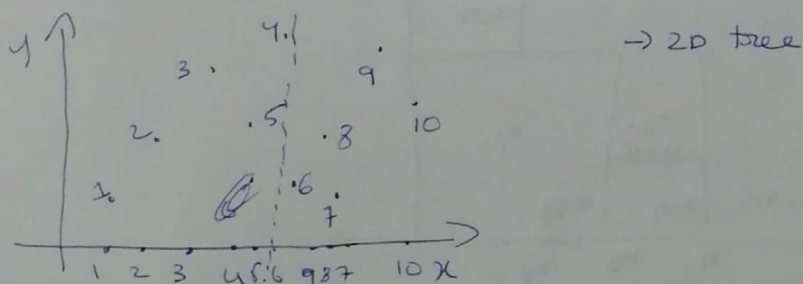


10:29 How to build a kd-tree

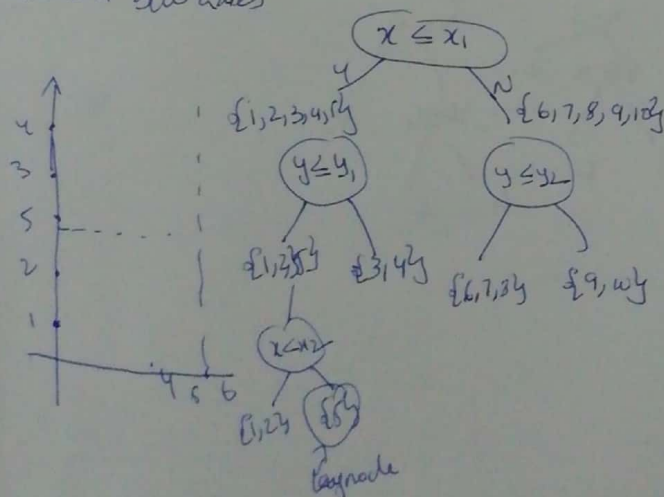
BST: given a list of sorted number \rightarrow Given a tree
⏟
scalar
'1D'

Extend BST to kd tree

\hookrightarrow It will operate in 2D, 3D ... n



- ① Pick x-axis, project points on to x-axis.
- ② Compute the median (will be b/w 5 & 6)
 \hookrightarrow Half the points to the left of median. and the other half to the right
- ③ Alternate b/w axes



18.30 Code sample : Decision Boundary

sklearn.neighbors.KNeighborsClassifier

```
class sklearn.neighbors.KNeighborsClassifier (n_neighbors = 5, weights = 'uniform',  
algorithm = 'auto', leaf_size = 30, p = 2, metric = 'minkowski', metric_params = None,  
n_jobs = 1, **kwargs)
```

n_neighbors : int optional (default=5)

weights : string or callable, optional (default='uniform')

- 'uniform': uniform weights. All points in the neighborhood are weighted equally
- 'distance': weight points by the inverse of their distance. In this case, closer neighbors of a query point will have a greater influence than neighbors which are further away
- [callable]: a user-defined function which accepts an array of distances, and returns an array of the same shape containing the weights

algorithm: {'auto', 'ball-tree', 'kd-tree', 'brute'}, optional

Algorithm used to compute the nearest neighbors:

- * 'balltree' will use BallTree
- * 'kd-tree' will use KDTree
- * 'brute' will use the brute-force search
- * 'auto' will attempt to decide the most appropriate algorithm based on the values passed to fit method

Note: fitting on sparse input will override the setting of this parameter, use brute force

- * Dimensionality is low it will pick KDTree
- * If dataset is small it will pick Brute

leaf_size: int, optional (default=30)

leaf_size passed to BallTree & KDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree.

The optimal value depends on the nature of the problem

p: Integer, optional (default=2)

Power parameter for the Minkowski metric. When $p=1$, this is equivalent to using Manhattan distance (L_1) & Euclidean distance (L_2) for $p=2$. For arbitrary p , Minkowski distance (L_p) is used

metric: string & callable, default 'minkowski'

The distance metric to use for the tree. The default metric is Minkowski and with $p=2$ is equivalent to the standard Euclidean metric.

See the documentation of the distance-metric class for a list of available metrics

metric_params: dict, optional (default=None)

Additional keyword arguments for the metric function

n_jobs: int, optional (default=1)

The number of parallel jobs to run for neighbors search.

If -1 , then the number of jobs is set to no. of CPU cores.

Doesn't affect fit method.

Examples In sklearn documentation

```
X = [[0], [1], [2], [3]]
```

```
Y = [0, 0, 1, 1]
```

```
from sklearn.neighbors import KNeighborsClassifier
```

```
neigh = KNeighborsClassifier(n_neighbors=3)
```

```
neigh.fit(X, Y)
```

```
print(neigh.predict([[1.1]]))
```

```
O/p: [0]
```

```
print(neigh.predict_proba([[0.9]]))
```

```
[[0.666667, 0.333333]]
```

18.31 Code Sample: Cross Validation

```
import sklearn. cross-validation import cross_val_score
```

```
name = ['x', 'y', 'class'] # define class names.
```

```
df = pd.read_csv('Concentration.csv', header=None, name=name)
```

```
df.head()
```

```
# Create design matrix X and target vector Y
```

```
X = np.array(df.iloc[:, 0:4]) # End index is exclusive
```

```
Y = np.array(df['class']) # showing you two ways of indexing a pandas df
```

Simple Cross Validation

```
# Split the dataset into train & test.
```

```
X_train, X_test, y_train, y_test = cross_validation.train_test_split
```

```
(X, y, test_size=0.3, random_state=42)
```

Split the train data set into cross validation train & cross validation test

```
X_tr, X_cv, y_tr, y_cv = cross_validation.train_test_split(X, y,  
test_size=0.2)
```

```
for i in range(1, 30, 2):
```

```
    knn = KNeighborsClassifier(n_neighbors=i)
```

```
    knn.fit(X_tr, y_tr)
```

```
    pred = knn.predict(X_cv) # prediction of cv train
```

```
    # Evaluate cv accuracy
```

```
    acc = accuracy_score(y_cv, pred, normalize=True) * float(100)
```

```
    print('cv accuracy for k = %d is %d%%' % (i, acc))
```

```
knn = KNeighborsClassifier(1)
```

```
knn.fit(X_tr, y_tr)
```

```
pred = knn.predict(X_test)
```

```
acc = accuracy_score(y_test, pred, normalize=True) * float(100)
```

```
print('Test accuracy
```

O/P

k=1 100 %.

3 100 %.

5 100 %.

creating odd list of k for knn

myList = list(range(0, 50))

neighbors = list(filter(lambda x: x % 2 != 0, myList))

Empty list that will hold cv scores.

cv_scores = []

for k in neighbors:

knn = KNeighborsClassifier(n_neighbors=k)

scores = cross_val_score(knn, X_train, Y_train, cv=10, scoring='accuracy')

cv_scores.append(scores.mean())

changing the misclassification error.

MSE = [1 - x for x in cv_scores]

determining best k

optimal_k = neighbors[MSE.index(min(MSE))]

print('The optimal number of neighbor is x.d.', x.optimal_k)

plt.plot(neighbors, MSE)

for xy in zip(neighbors, np.round(MSE, 3)):

plt.annotate('(x.s, y.s)' % xy, xy=xy, textcoords='data')

plt.xlabel('Number of neighbor k')

plt.ylabel('misclassification error')

plt.show()

KFold

```
class sklearn.model_selection.KFold (n_splits, shuffle = False, random_state = None)
```

→ K-Folds Cross Validator

→ provides train/test indices to split data in train/test sets. Split dataset into K consecutive folds (without shuffling by default)

→ Each fold is then used once as a validation when K-1 remaining folds form the training set.

n_splits : int, default=3

Number of folds. Must be at least 2.

Shuffle : boolean, optional

Whether to shuffle data before splitting into batches.

Random-state : int optional.

```
// from sklearn.model_selection import KFold
```

```
X = np.array([[1, 2], [3, 4], [1, 2], [3, 4]])
```

```
y = np.array([1, 2, 3, 4])
```

```
kf = KFold(n_splits = 2)
```

```
kf.get_n_splits(X)
```

```
// 2
```

```
print(kf)
```

```
// KFold(n_splits = 2, random_state = None, shuffle = False)
```

```
for train_index, test_index in kf.split(X):
```

```
    X_train, X_test = X[train_index], X[test_index]
```

```
    y_train, y_test = y[train_index], y[test_index]
```

StratifiedKFold

StratifiedKFold (n_splits=3, shuffle=False, random_state=None)

- Stratified K-Folds cross validation.
- provides train/test indices to split data in train/test sets
- This cross validation object is a variation of KFold that returns stratified folds.
- The folds are made by preserving the percentage of samples for each class.

Notes

All the folds have size $\text{train} (n\text{-samples}/n\text{-splits})$, the last one has the complementary.

```
from sklearn.model_selection import StratifiedKFold
```

```
X = np.array([[1, 2], [3, 4], [1, 2], [3, 4], [5, 6], [7, 8]])
```

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

```
skf = StratifiedKFold(n_splits=2)
```

```
skf.get_n_splits(X, y)
```

```
for train_index, test_index in skf.split(X, y):
```

```
    X_train, X_test = X[train_index], X[test_index]
```

```
    y_train, y_test = y[train_index], y[test_index]
```


Repeated KFold

RepeatedKFold (n_splits=5, n_repeats=10, random_state=None)

→ Repeats K-Fold n times with different randomization in each repetition

→ rkf = RepeatedKFold(n_splits=2, n_repeats=2, random_state=2652124)

for train_index, test_index in rkf.split(X):

X_train, X_test = X[train_index], X[test_index]
y_train, y_test = y[train_index], y[test_index]

Leave One Out Cross Validation (LOOCV)

→ In this approach, we reserve only one data point from the available dataset and train the model on the rest of the data.

→ This process iterates for each data point.

→ This also has its own advantages & disadvantages.

- * we make use of all data points, hence bias will be low

- * we repeat the cross validation process n times (where n is the number of data points) which results in high execution time.

- * This approach leads to higher variation in testing model effectiveness because we test against one data point. So, our estimations get highly influenced by the data point. If the data point turns out to be an outlier it can lead to a higher variation.

//
from sklearn.model_selection import LeaveOneOut

X = np.array([[1, 2], [3, 4]])

y = np.array([1, 2])

loo = LeaveOneOut()

loo.get_n_splits(X)

for train_index, test_index in loo.split(X):

 X_train, X_test = X[train_index], X[test_index]

 y_train, y_test = y[train_index], y[test_index]