# Day-27.

## K-Nearest Neighbour.

- Classification Technique                                          9
- Classify records with help of euclidean distance.               10

  find closeness of 2 points

## features.                                                       12

- All instances correspond to pts. in n-dimensional Euclidean space
- Classification is delayed till new instance arrives.             2
- Target $f^n$ may be discrete or real valued.                     3
- Classification done by comparing feature vectors of             4
  diff. points.                                                   5
                                                                  6

It is instance Based Learning
→ Based on characteristics of record we classify the O/P.

Euclidean distance →

1st response     2nd response              record 1

$$\sqrt{(x_1 - \mu_1)^2 + (x_2 - \mu_2)^2 + \ldots + (x_p - \mu_p)^2}$$

find closeness or
of 2 records.

$$D(x,y) = \sqrt{\sum_{i=1}^{n}(x_i - y_i)^2}$$

Ex Jay:
Age = 35
Salary = 95k
CD = 3

record 2.

Riya:
Age = 41
Salary = 215k
CD = 2

Ex:

euclidean dis. $= \sqrt{\underbrace{(35-41)^2}_{Age} + \underbrace{(95000 - 215000)^2}_{Salary} - \underbrace{(3-2)^2}_{CD}}$

13 Friday
May 2011

# Note Choosing K: k=5 (default),
choose that value of K which
has lowest error
rate in
Validation data.

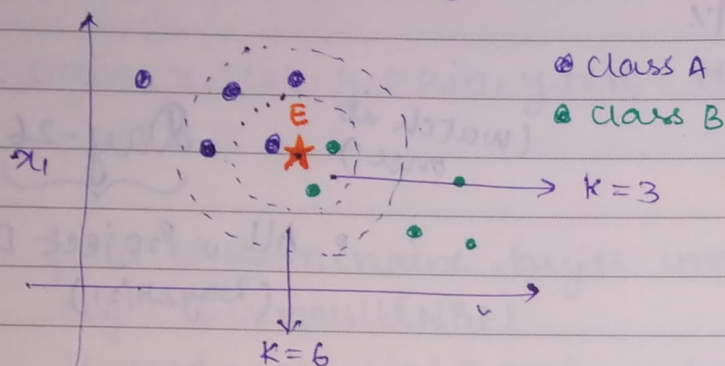May 2011
Wk | Mo Tu We Th Fr Sa Su
17/22 | 30 31                    1
18 |  2  3  4  5  6  7  8
19 |  9 10 11 12 13 14 15
20 | 16 17 18 19 20 21 22
21 | 23 24 25 26 27 28 29

Week 19 ▪ 133-232

## KNN → find how much target variable is closer to predicted variable.



● Class A
● Class B

→ k=3

K=6

E(★) needs to be predicted,
from which class it
belongs, nearest
* for K=3, 3 neighbours are
taken & from majority
from 3, it is predicted.
* for K=6, 6 neighbours
were taken for prediction

. Prediction.
K=3 → Class B
K=6 → Class A

9

10

soln.

Ex.

12  ① If Both are same : 2 class A, 2 Class B then we need to ↑ the
K value. & we need to do iteration.

1

2

3

| Strength | Weakness |
|---|---|
| — Simple to implement & use. | — Need lot of space to store all eg (example) |
| — easy to explain pred$^n$ | — Takes more time to classify a new ex. than with a model. |
| — Robust to noisy data by avg. KNN | |

4

5

6

| Advantage. | Disadvantage. |
|---|---|
| — can be applied to data from any distribution | — Choosing best K maybe difficult |
| — Good classification if no. of sample is large enough. | — Need large no. of samples for accuracy. |
| | — can never fix without assuming parametric distribution. |

DV → categorical
{ IDV → " + cont^n.

2011
June
Wk| Mo Tu We Th Fr Sa Su
1  2  3  4  5
22  6  7  8  9 10 11 12
23 13 14 15 16 17 18 19
24 20 21 22 23 24 25 26
25 27 28 29 30
26

S a t u r d a y

May          2011

14

Week 19  ▪  134-231

## Practical Implementation

Dataset used → Titanic.

• import all essential lib.

```
df = pd.read_csv('train.csv')
df.columns
```

```
le = preprocessing.LabelEncoder()
df['sex'] = le.fit_transform(df['Sex'])
```

```
from sklearn import neighbors
y = df['Pclass']
x = df[['Pclass','PassengerId'], axis=1)
   df.drop
```

• DV → categorical

```
x-train, x-test, y-train, y-test = train-test-split(x, y, test_size = 0.3,
                                                      random_state = 0)
```

```
knn = neighbors.KNeighboursclassifier(n-neighbors = 3)
   knn.fit(x-train, y-train).score(x-test, y-test)
```
↳ 85.39%   (model accuracy).

```
ypred = knn.predict(xtest)
confusion-matrix(ytest, ypred)
```

Traing Pclass

(Here DV is const.
→ × change k)

|   |   | 1 | 2 | 3 |
|---|---|---|---|---|
| Test Pclass | 1 | 60 | 6 | 4 |
|  | 2 | 7 | 27 | 15 |
|  | 3 | 3 | 4 | 141 |

Total = 267
Correct = 228
Incorr = 39

| DV | IDV | K | Score |
|---|---|---|---|
| Pclass | All. | 3 | 85.39%. |
| " | " | 4 | 83.89% |
| " | " | 2 | 85.39%. |

228/267 = 85.39%.

hw. create f^n.

K = 1 to 167.
plt. plot list.

# 16
**Monday**

**May**　　　2011

May
| Wk | Mo | Tu | We | Th | Fr | Sa | Su |
|----|----|----|----|----|----|----|----|
| 17/22 | 30 | 31 | | | | | 2011 |
| 18 | 2 | 3 | 4 | 5 | 6 | 7 | 1 |
| 19 | 9 | 10 | 11 | 12 | 13 | 14 | 8 |
| 20 | 16 | 17 | 18 | 19 | 20 | 21 | 15 |
| 21 | 23 | 24 | 25 | 26 | 27 | 28 | 22 |
| | | | | | | | 29 |

Week 20 ▪ 136-229

## <u>SVM</u> (support vector Machine)
↳ classify record with help of hyperplane.

```
              SVM classification
         → Linear SVM   (HyperPlane)
         → Non-Linear SVM (Kernel Trick)
```

⊙ SVM is used to classify record for over-dimensional data.
　　— Used for both Regression & Classification problems.
　　— mostly used in classification problem.
　　　( we plot each data item as pt. in n-dimensional space),
　　　when n is no. of features.
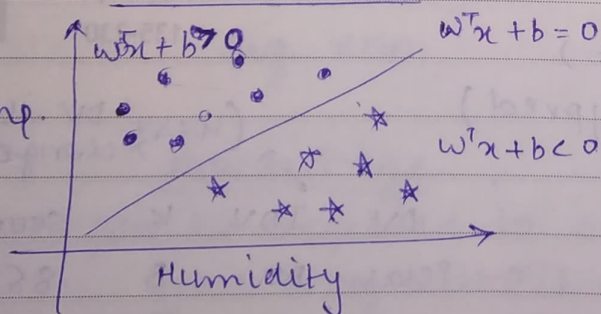


hyperplane is a line which divide the 2 groups.

**Rules for Hyperplane**

It should be equidistant to both groups ↓

① Hyperplane should divide the 2 groups.

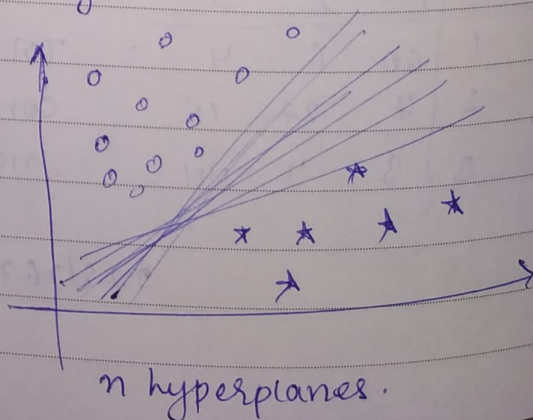② whichever pt. is closer to hyperplane that pt. is called support vector.

\# dist. b/w 2 support vector is called margin.

② **Linear SVM.**



$w^T x + b > 0$

$w^T x + b = 0$

$w^T x + b < 0$

Temp.

Humidity

$$f(x) = sign(w^T x + b)$$

eqⁿ of hyperplane from algebra.



n hyperplanes.

2011

| | | | Th | Fr | Sa | Su |
June
Mo Tu We Th Fr Sa Su
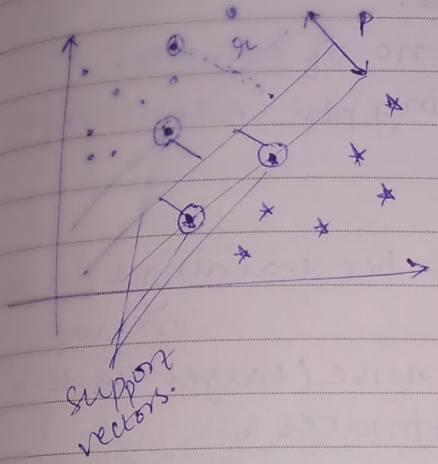| | 1 | 2 | 3 | 4 | 5 |
| 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 |
| 20 | 21 | 22 | 23 | 24 | 25 | 26 |
| 27 | 28 | 29 | 30 | | | |

**T u e s d a y**
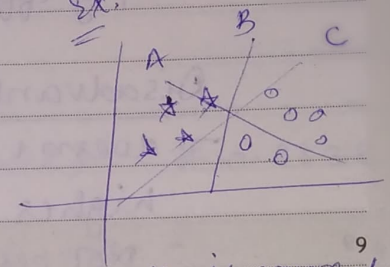
May          2011

17

Week 20  ▪  137-228

we can draw n-no. of hyperplanes which dincle 2 groups.
But AIM is to find such hyperplane that correctly clasify data.



**Good Hyperplane.**

① ⎫ Both rules
② ⎭ should be
          followed.

ex.



B is correct
hyperplane

support
vectors

2×2.



· C is correct hyperplane.
( A & B ko agr lenge to ek to support vector k pas hoga
   lekin dusra vector dur hojayega.)
   margin difference. ( they are not equidistant )

| | 9 |
| 10 |
| 11 |
| 12 |
| 1 |
| 2 |
| 3 |
| 4 |
| 6 |

## Non-Linear.

we convert 2D to multidimensional. is called kernel-trick.



2D

$\phi : x \rightarrow \phi(x)$

3D

# 18

**W e d n e s d a y**

**May**     **2011**

May
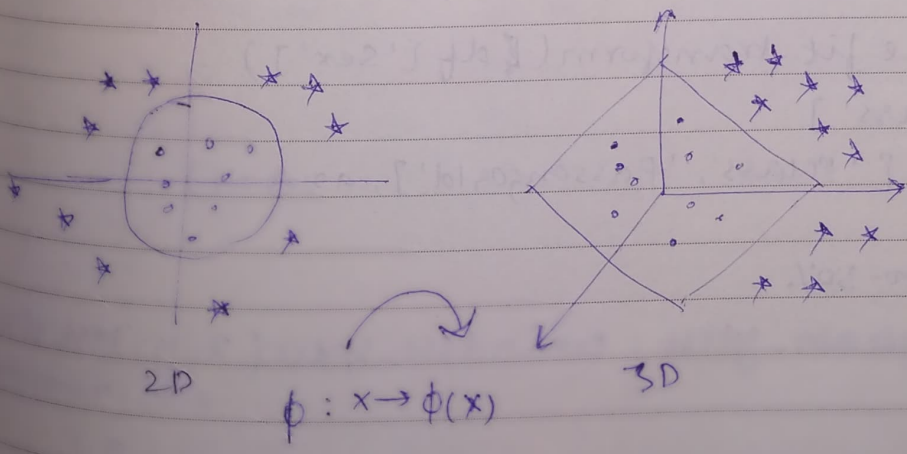| Wk | Mo | Tu | We | Th | Fr | Sa | Su |
|---|---|---|---|---|---|---|---|
| 17/22 | 30 | 31 | | | | | 2011 |
| 18 | 2 | 3 | 4 | 5 | 6 | 7 | 1 |
| 19 | 9 | 10 | 11 | 12 | 13 | 14 | 8 |
| 20 | 16 | 17 | 18 | 19 | 20 | 21 | 15 |
| 21 | 23 | 24 | 25 | 26 | 27 | 28 | 22 |
| | | | | | | | 29 |

Week 20 ▪ 138-227

## Advantage

- works really well with clear margin of separation
- effective in high dimensional spaces.
- effective where no. of dimension > no. of samples.
- M/my efficient (uses subset of training pts. in decision fn).

## Disadvantage

- doesn't perform well, with large dataset bcz training time is higher.
- Not perform well, when data has noise (target classes overlaps.)
- doesn't directly provide prob. estimates, calculated using expensive 5 fold cross-validn.

## Python Implementation

$$dataset \rightarrow Titanic$$

import necessary libraries.

```
from sklearn import svm
df = pd.read-csv('train*csv')

df['sex'] = le.fit-transform(& df['sex'])
y = df['Pclass']
x = df.drop([['Pclass', 'PassengerId'], axis = 1)
```

Split data into 70-30%.

```
x-train, x-test, ytrain, ytest = train-test-split(x, y, test size=
0.3, random-state = 0).
```

2011

June
| Wk | Mo | Tu | We | Th | Fr | Sa | Su |
|----|----|----|----|----|----|----|----|
| | | | | | 1 | 2 | 3 | 4 | 5 |
| 22 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| 23 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
| 24 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |
| 25 | 27 | 28 | 29 | 30 | | | |
| 26 | | | | | | | |

T h u r s d a y

May            2011

**19**

Week 20  ▪  139-226

clf = svm.SVC (gamma = 0.01, C = 100)
↓                      ↳ complexity factor ie
99.9% this record         100 times it perform
is going to execute          iteration
accurately.

clf.fit (xtrain, ytrain)
ypred = clf.predict (x-test)
accuracy-score (ytest, ypred, normalize = True)    # 0.8838
confusion.matrix (ytest, ypred).

Training Pclass

| | 1 | 2 | 3 |
|---|---|---|---|
| 1 | 61 | 6 | 8 |
| 2 | 4 | 35 | 10 |
| 3 | 2 | 6 | 140 |

Total = 267     ⎫  236/287 = 88.38 %
correct = 236    ⎬
Incorrect = 31   ⎭

HW

| DV | IDV | Acc. |
|----|-----|------|
| Pclass | All | ~~88.38%~~ 89.55 |
| Survived | " | 76.11 |
| Gender | " | 75% |
| Embarked | " | 74.62% |
| Parch | " | 82.83% |
| Sibsp | " | 74.25% |

which comb$^n$
is giving
high accuracy?

. Create fn.