

7/01/2019

## Machine Learning

1

### → Books

- Understanding Machine Learning: From Theory to Algorithms  
— By Shalev Shwartz and Shai Ben-David
- Pattern Recognition and Machine Learning  
— By Christopher M. Bishop.

Machine Learning by Tom Mitchell

Introduction to Machine Learning by Alpaydin

### → Learning:

A child learning alphabets

Generalisation

— comes from experience.

Monkey and Banana

Extrapolating

pigeon. Superstition.

Bait Synts. Shyness.

## Machine Learning

→ Learning is the process of counting experience to knowledge.

→ Machine learning is about a program that they can learn from the available inputs to them and to convert it into knowledge and perform the same task.

→ Unseen data - The data not used for training, there is

prior knowledge or Productive Bias.

(2)

- The knowledge of features of human and animal to differentiate between entities
- features or attributes/ - define the entity
- characteristics/ properties.
- Algorithm must be able to learn distinguishing ability
- adaptive.

### Types of learning

⇒ Supervised learning. (guidance exists) labels, classes

— Unsupervised learning.

$$y = mx + c$$

$$y = \sin(2\pi x)$$

can be a vector -  $(x_1, x_2, \dots, x_m)$

where  $m$  is the number of features

14/01/2020

Euclidean - to measure similarity / dissimilarity

Point

Manhattan.

similar

Pythagorean -  $\sqrt{x^2 + y^2}$  and  $d(x, x) = 0$ ,  $d(x, y) = d(y, x)$

Metric: A measure which satisfies

These three properties is called a metric

Euclidean distance is  $\sqrt{\sum (x_i - y_i)^2}$

→ If the points are similar then similarity will be high. Or metric is low.

→ While trying to group the objects we need to consider similarity and dissimilarity. (2)

→ Mapping  $y = f(x)$

Choosing correct function.  $f_1, f_2, f_3$  among these.

A formal model:

domain Set - {Input you expect}

An arbitrary set  $X = \{x_1, x_2, x_3, \dots, x_n\}$ .

label Set - {Target that we need to find}.

$y$  is label set  $y = \{0, 1\}$  for two class problems

$\{+1, -1\}$  useful, think (A) & (B) & (C)

Training data -  $\{(x_1, y_1), (x_2, y_2), \dots, (x_p, y_p)\}$ .

You shouldn't use 100% for training.

We have to have some data for testing of the trained model.

controlled Environment - Already have knowledge for which  $x$  we have which  $y$ .

→ learner's output:

Training Testing

The learner is also called predictor (learning provides hypotheses or a classifier). (mapping of  $x$  value)

The learner is expected to produce  $y_i$  from a production rule. (data and see  $y$  value)

$h: X \rightarrow Y$ .  $y = f(x)$

error  $y_i - \hat{y}_i$

→ We use  $A(S)$  to denote the hypothesis that a learning algorithm  $A$ , returns upon receiving the training sequence  $S$ . (4)

→ Measure of Success:

We define the error of a classifier to be the probability that it does predict the correct label on a random data point guaranteed by the underlying distribution.

Formally, given a domain subset,  $A \subset X$ , the probability distribution  $D_A$  assigns a number,  $D(A)$ , which determines how likely it is to observe a point  $x \in A$ . We refer to  $A$  as an event and express it using a function  $\Pi : X \rightarrow \{0, 1\}$ , namely,  $A = \{x \in X : \Pi(x) = 1\}$ .

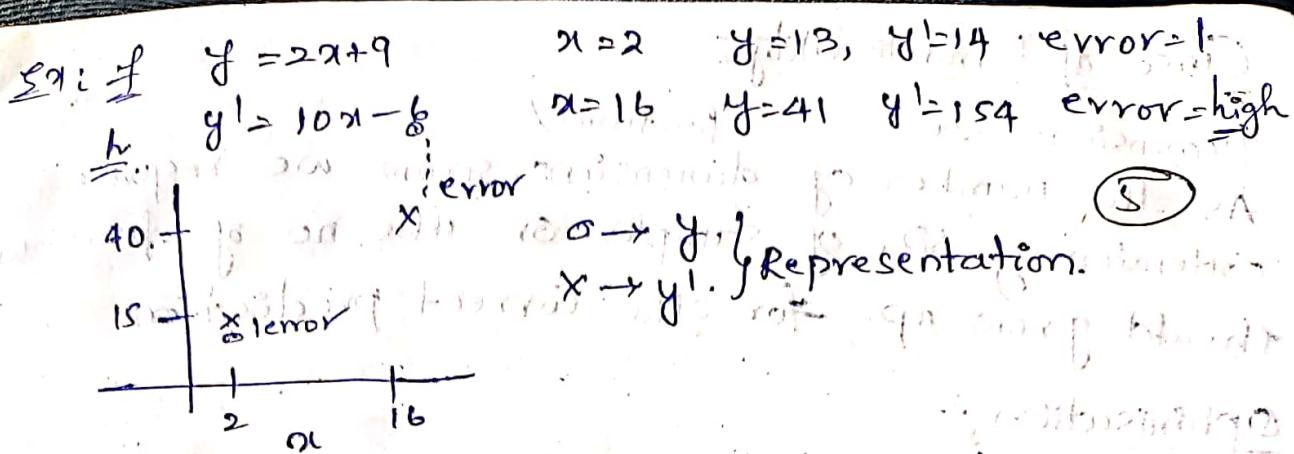
(from TB Shalev and Shai Benyamin)

// If  $x \in A$  to the class. then  $\Pi(x)=1$  or else the belongs  $\Pi(x)=0$ .

→ We define the error of prediction rule.

$h : X \rightarrow Y$  to be

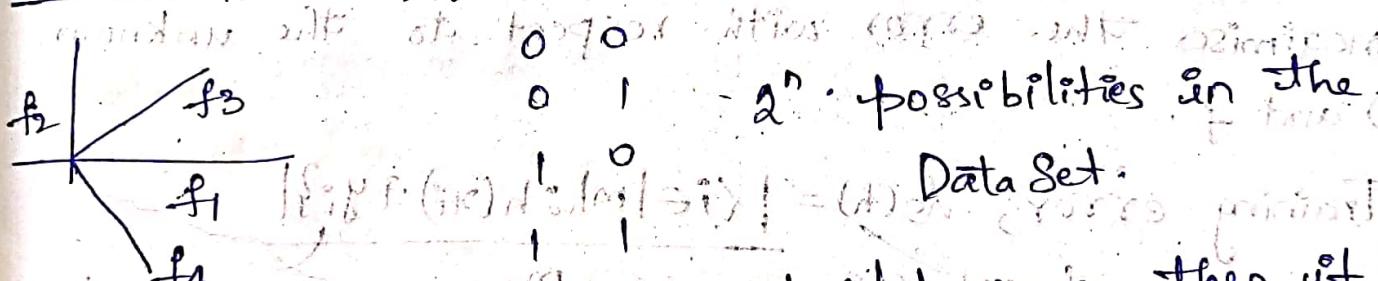
$\text{err}_D, f^{(h)} = P[h(x) \neq f(x)] = D\{x : h(x) \neq f(x)\}$ .  
we derive - the difference b/w them is the error.



→  $(x_D, y_D)$  → has many synonyms, such as

- Generalisation error
- The risk
- The true error of  $\hat{y}$ .

Feature Space: Basic f1, f2, f3, f4, ...  $2^n$  possibilities in the



If data generated from restricted space, then it learns about all the similar data but the disadvantage here is the model cannot learn about dissimilar data. (6)

If we want to represent all the real world data, we have to choose the data/samples such that it not only has the restricted space.

If generalisation not happening correctly because of taking only the restricted area, then it may leads to generalisation error.

## Curse of Dimensionality

(6)

Dimension-size (features)  
As the number of dimensions grow, we require extension of size of samples or the no of features should grow up. for the correct prediction.

## Optimisation:

Machine Learning Algorithms are Optimisation Algorithms.

optimisation in this case is  $(y - \hat{y})^2$ , to make it positive.

## Empirical Risk Minimisation; (ERM):

The goal of the algorithm is to find  $h$  that minimises the error with respect to the unknown  $D$  and  $f$ .

Training error,  $ds(h) = \frac{1}{m} \sum_{i=1}^m |h(x_i) - y_i|$

Hypothesis

where  $[m] = \{1, 2, 3, \dots, m\}$ .  
Expected output

The learning algorithm that comes up with a predictor  $h$  that minimises  $ds(h)$  is called Empirical risk minimisation or ERM.

## Overfitting:

If we are given training data, instead of learning data generalises (Remembering)  
(child learning tables).  $Ls(h)$  if value not there is error huge.

$$\Sigma y = y = \sin(2\pi x) \quad | \quad y = \sum w_i x^i$$

(T.B. Bishop)

$$y = w_0 + w_1 x + w_2 x^2 + \dots$$

exponent value = 2 (Quadratic)

If  $y = [w_0 + w_1 x]$ , only this taken then it becomes linear.

(There by error large.)

If value = 8 (Almost fit)

If exponent value increases more then there will be many oscillations [Because it tries to remember].

16/01/2020

$\alpha$  Def (h) Generalisation Error.

$L_g(h)$  Empirical Risk.



minimum Empirical Risk is achieved, around

→ k-fold Cross-validation:

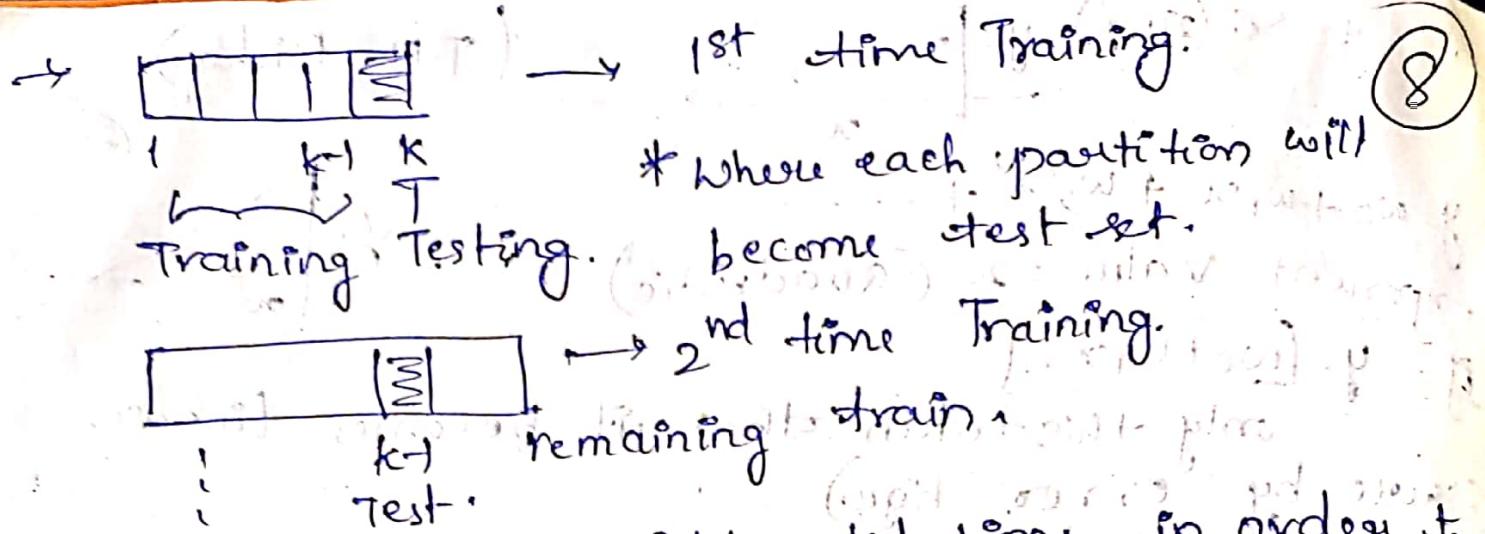
Create k-sets or partition

k-sets and Randomly pick the training set and the testing set.

1	2	3	4	5	6	7	8	9	10
1	2	3	4	5	6	7	8	9	10

randomly pick from this k-divided set.

some partitions for training and some partitions for testing



\* We do the Cross fold validations in order to decrease the error.

For unbiased learning.

More amount of time for training.

→ Prior knowledge; set of features to say the problem requires that.

more training some times. Greater confusion.

In k-fold → 1 part will be k times trained.

### Sample Selection: (Sampling methods)

Random selection of samples from every area.

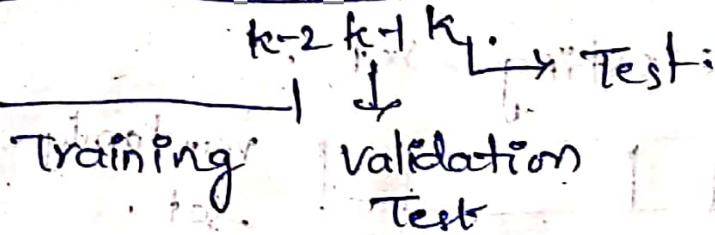
Random Sampling - Good enough for sampling.

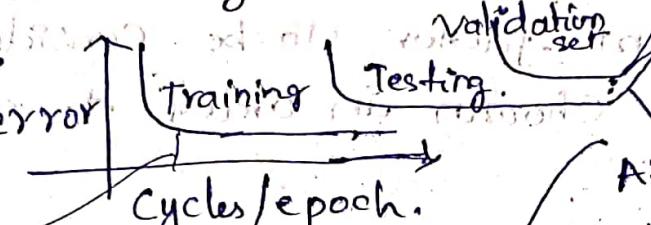
→ Instead of k-fold Cross validation one Validation we can use



Set Method

→ Test can only be tested after the learning completed.



→ So, while the training is going on, the error  $J_{\text{S}}(h) = \underline{0}$ .  
 The hypothesis is tested on the validation set.  
 So, In order to decrease the generalisation error  
 (This set is used for testing generalisation and if error is huge, while training itself we get to know).  


Learned in training set and test that on test data.

After this point training to be stopped since model overlearns. Therefore, only remember to and pass for this purpose. Validation

set is used.

(At that point error is checked and we stop the training there).  
 → Load Learning (child learning tables). (Just Remember (Rote))

→ To minimize the Empirical Risk

we go for Inductive Biasing.

→ A set of predictors called class of Hypothesis is denoted by  $\mathcal{H}$ . If  $h_{\text{eff}}$  is a function mapping  $X \rightarrow y$ .

we have to pick one in order to minimize the Empirical Risk.

$$\underset{\mathcal{H}}{\text{ERM}}(S) \in \arg \min_{h_{\text{eff}}} J_{\text{S}}(h)$$

$$\sin(2\pi x)$$

$$f = \{h_1, h_2, h_3, \dots, h_n\}$$

$$h_1 = ax + b$$

$$h_2 = cx^2 + dx + e$$

$$h_3 = f + g_1 + h_1^2 + \dots + h_n^2$$

\* which gives the minimum error to be chosen

→ Q Another important factor to be considered which hypothesis to be chosen in order to avoid overfitting.

→ how many hypotheses to be considered, how many can a data set handle? Can there be any upperbound?

Assumptions: In statistical theory,

The I.I.D. assumption. (Independently, Identically Distributed according to Distribution D).

21/01/2020

→ Non Representative Set: Take samples from Distribution D. Some Missing Regions: How many number of samples are sufficient? (Probabilistic and approximation) Samples (Independent and Identical)

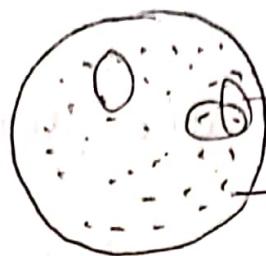
→ The Training Samples are Independently and Identically distributed according to a Distribution D.

$S^n D^m$  where  $m'$  is the size of  $S$ , and  $D^m$  denotes the probability of  $m$ -tuples indeed by applying

'D' to pick each element of the sample independently.  
of the other members of the tuple.

→ Q. What is  $m$ ?

(11)



sets of bad instances.

A set consisting of  $m$  number of samples

→ Confidence Parameter } 2 parameters to estimate

→ Accuracy

What is the confidence that, that hypothesis provides certain accuracy.

→ The Probability of getting a non representative sample is  $\delta$  and call  $(1-\delta)$  the confidence parameter of our prediction.

Since, we cannot guarantee a perfect label, we choose another parameter for the quality of prediction, the accuracy parameter,  $\epsilon$ .

→  $L(D, f) < \epsilon$  for the correct predictor

True Error

$f: X \rightarrow Y$

$g: S \rightarrow Y$

Training

→ Let  $S_{\text{tr}} = \{x_1, x_2, \dots, x_m\}$  is the instances  $m$ , the training set; if there are bad samples and if we find the upperbound for probability, then we get  $L(D, f)$  high leads to  $L(D, g) = 0$ , we can avoid them

(12)

$$\rightarrow D^m \left( \{S|_x : d_{(D,f)}^{(hs)} > \epsilon \} \right)$$

Let  $H_B$  be the set of bad hypothesis, that is

$$H_B = \{h \in H : d_{(D,f)}^{(hs)} > \epsilon \}$$

Let  $M = \{S|_x : \exists h \in H_B, l_S(h) = 0\}$  be the set of misleading samples.

Training error = 0

leads to overfitting.

(ERM) - Called as

If this happens True error is  
- see

$$S|_x : h_{(D,f)}^{(hs)} > \epsilon \} \subseteq M$$

(Subset of M)

$M$  = Union of Bad set of examples.

$$M = \bigcup_{h \in H_B} \{S|_x : l_S(h) = 0\}$$

$$\text{hence, } D^m \left( \{S|_x : d_{(D,f)}^{(hs)} > \epsilon \} \right) \leq D^m(M) = D^m \left( \bigcup_{h \in H_B} \{S|_x : l_S(h) = 0\} \right)$$

$\rightarrow$  Lemma: For any two sets A and B and a distribution

$$D' = D(A \cup B) \leq D(A) + D(B)$$

$$D^m \left( \{S|_x : d_{(D,f)}^{(hs)} > \epsilon \} \right) \leq D^m(M) = D^m \left( \bigcup_{h \in H_B} \{S|_x : l_S(h) = 0\} \right)$$

$$= \sum_{h \in H_B} D^m(\{S|_x : l_S(h) = 0\})$$

Bad hypothesis.

$$\begin{aligned} \rightarrow D^m \left( \{S|_x : l_S(h) = 0\} \right) &= D^m \left( \{S|_x : \forall i, h(x_i) = f(x_i)\} \right) \\ &= \prod_{i=1}^m D(x_i : h(x_i) = f(x_i)) \end{aligned}$$

$$\rightarrow \sum_{h \in H_B} D^m(\{S_i : h(S_i) = 0\}) = 1 - \frac{e^{-m}}{(D_f)} \leq 1 - e$$

(13)

$$D(x_i : h(x_i) = y) = 1 - \frac{e^{-m}}{(D_f)} \leq 1 - e$$

$$\therefore \sum_{h \in H_B} D^m(\{S_i : h(S_i) = 0\}) = (1 - \frac{e^{-m}}{(D_f)})^m \leq (1 - e)^m.$$

$$\Rightarrow D^m(\{S_i : h(S_i) = 0\}) \leq (1 - e)^m \leq e^{-em}. \quad -(2.7)$$

$$\rightarrow D^m(\{S_i : \frac{h_S}{(D_f)} > \epsilon\}) \leq \sum_{h \in H_B} D^m(\{S_i : h(S_i) = 0\}) \quad -(2.9),$$

$$\leq |H_B| e^{-em} \rightarrow \text{If log applied on } b.s. \downarrow \text{were}$$

$\rightarrow$  Corollary: Let  $H$  be a finite hypothesis class, let  $\delta \in (0, 1)$  and  $\epsilon > 0$  and let  $m$  be an integer that satisfies  $\frac{m \log((|H|)/\delta)}{\epsilon}$

$$\frac{m \log((|H|)/\delta)}{\epsilon}$$

$\rightarrow$  Suppose, the hypothesis space  $H$  is the set of conjunctions of literals of  $n$  Boolean variables. In this 3 states for Boolean variable:

1. Negated

2. Unnegated

3. It does not appear  
no of boolean variables

$$|H| = 3^n + 1$$

3 states adding conjunctions, sent falsely,

which is (the conjunction of any atom and its negation).

$\Rightarrow \epsilon = 5\%$  confidence, as  $99\%(\delta)$

Error less than

5% - accuracy tells

$\rightarrow$  If we want to guarantee at most 5% Error and 99% of the time

$$n=30$$

$$\epsilon = 5\% = \frac{1}{20}$$

$$\delta = 99\% = \frac{1}{100}$$

The Sample Complexity

$$= 20 \times (n \ln 3 + \ln(100))$$

$$\approx 752$$

Using no of samples  $\frac{752}{2} = 376$  we can generate error less than 5% and with confidence 99%.

$$752 \leq 2^{30} = 1,073,741,824$$

$\rightarrow$  And, the number of hypothesis

$$= 2^{30} + 1 = 205,891,132,094,650$$

Hypothesis Space =  $2^n$ .

$$(m = \frac{1}{\epsilon} (2^n \ln 2 + \ln(\frac{1}{\delta})))$$

$$m = \frac{1}{\epsilon} (2^n \ln 2 + \ln(\frac{1}{\delta}))$$

$$n=20$$

$$\epsilon = \frac{1}{20}$$

$$\delta = \frac{1}{100}$$

$$m = 20 \times (2^{30} \ln 2 + \ln 100) \approx 14,885,222,452$$

Samples.

14 Billion samples

$\rightarrow$  To choose the hypothesis to do estimation we require prior knowledge (Inductive Bias)

$$m \geq \frac{\log(\frac{1}{\delta})}{\epsilon}$$

14

Biased variance - trade off

(15)

→ For a finite hypothesis class, if the ERM rule with respect to that class is applied on a sufficiently large sample (whose size is independent of underlying distribution or labeling function then the output hypothesis will be (probably) approximately correct.)

PAC Learning model: Probably Approximately Correct learning model: A hypothesis class  $H$  is PAC learnable if there exists a function  $m_H : (0,1)^2 \rightarrow N$  and a learning algorithm with the learning property: for every  $\epsilon, \delta \in (0,1)$ , for every distribution  $D$  over  $X$ , and for every labeling function  $f: X \rightarrow \{0,1\}$ , (belong to) for every realization assumption holds with respect to  $H, D, f$ , then.

When running the learning algorithm on  $n$  i.i.d examples generated by  $D$  and  $m \geq m_H(\epsilon, \delta)$  the algorithm returns a hypothesis labelled by  $g$ , such that, the probability of atleast  $(1-\delta)$  (over the choice of examples),  $d(D, g) \leq \epsilon$

22/01/2020.

## IT Lab - 2 ML

(16)

→ Numpy python Matplotlib pandas python object

← Introduction to Numpy: (Numerical python) c-structure in python

np.linalg → arrays homogeneous (same type)

→ gives documentation of numpy methods  
Understanding data-types in python. obj-ref

np.array([1, 2, 3, 4]) import array.  
↳ list(range(10)) Obj-type

np.array([1, 2, 3, 4]), array.array(i, 4). Obj-size, Obj-digit

fancy indexing → Creating arrays from scratch

masking

Broadcasting.

np.empty(3)

np.zeros(  
dtype=)

range(0, 20, 2) → 0, 2 → 18

linspace(0, 1, 5) → 0, 0.125, 0.25, 0.375, 0.5, 0.625, 0.75, 0.875, 1.0

\* Basic of Numpy Arrays

$$\text{sigmoid} = \frac{1}{1 + e^{-x}}$$

1/2  
2 dimensions

2  
1 dimension

b  
w - 1x2

$$z = w^T x + b$$

VV 2x1

$$\begin{vmatrix} 1 & \\ & 1 & \\ & 0 & \\ & 1 & \\ & 1 & \\ & 1 & \end{vmatrix}$$



23/01/2020  $\rightarrow h^* \in H$  which provides:

(17)

$$\mathbb{E}_D(h^*) = 0$$

Realizability

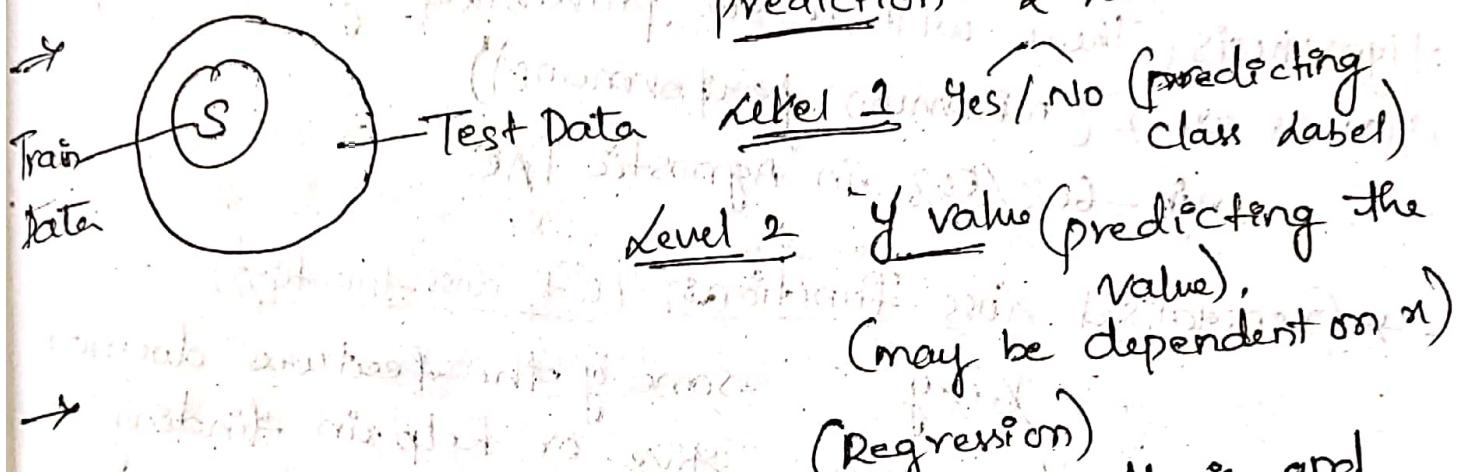
$$m \geq \frac{(-\ln(\delta))}{\epsilon^2}$$

Sample Complexity

Come up with a hypothesis that probably:

$\forall x \rightarrow$  dimension - we need to measure:

Prediction - 2 levels:



Realizability - Considering the finite hypothesis and choosing the correct hypothesis such that the generalisation error is minimised.

(Depends on prior knowledge) too much leads to overfitting.

→ Relaxing Realizability Condition - If finite hypothesis will be called agnostic (no assumption of  $H$ ) PAC learning. Instead of  $x \rightarrow y$  will be mapping to the Distribution  $x \times y \rightarrow$  Distribution.

D is considered as Probability Distribution over  $x \times y$ .

→ The Empirical and True Error

$$L_D(h) \equiv P(\text{probability, not a number}) \quad (17)$$

$$L_D(h) = P : h(x) \neq y \quad (x, y) \in D$$

$$L_S(h) = \underbrace{\{x \in [m] : h(x_i) \neq y_i\}}_{m}$$

Empirical error same  
but the true  
error differ

→ There must be some benchmark to take the hypothesis (There will be 50% probability of success if classes are 2 (minimum performance))

with min - 60%/50% in Agnostic PAC

→ Generalised Loss Functions: (0-1 loss function)

$X \rightarrow Y$  some of the features does not give or help in finding

$$x_i = (x_{i1}, \dots, x_{id})$$

$$y,$$

$$L_D(h) = L_D(h, (x, y)) = \begin{cases} 1 & \text{if } h(x) \neq y \\ 0 & \text{if } h(x) = y \end{cases}$$

This is called 0-1 loss function.

Square loss:  $L_D(h, (x, y)) = (h(x) - y)^2$

Def: A hypothesis class  $H$  is agnostic PAC learnable with respect to a set  $Z$  and a loss function  $l: H \times Z \rightarrow \mathbb{R}_+$ , if there exist a function  $m_H: (0, 1)^2 \rightarrow H$  and a learning algorithm with the following property:

for every  $\epsilon, \delta \in E^{(0,1)}$  and for distribution  $D$  over  $Z$ , when running the learning algorithm on  $m \geq m_H(\epsilon, \delta)$  i.i.d examples generated by  $D$ , the algorithm return  $h \in H$  such that with probability of atleast  $1 - \delta$ ,

$$D \leftarrow \min_{h \in H} L_D(h) + \epsilon \quad L_D(h) \geq 0.$$

28/01/2020.

Model Complexity:  $\{h_1, h_2, h_3, \dots\}$  -  $\text{hng}$  set of hypothesis  
 $\{D_1, D_2, D_3, \dots, D_n\}$  - Data sets

$h_i$  can be used for solving  $D_i$  Data set

- There is no Universal learner.

Theorem: No free lunch Theorem.  
 Let 'A' be any learning algorithm for the task of binary classification with respect to the 0-1 loss over a domain  $X$ . Let 'm' be the any member smaller than  $\frac{|X|}{2}$ , representing a training set size. Then there exists distribution  $D$  over  $X \times \{0,1\}$  such that:

1. There exists a function  $f: X \rightarrow \{0,1\}$  with  $L_D(f) = 0$ .
2. With the probability of atleast  $\frac{1}{f}$  over the choice of  $S \subseteq D^m$  we have that  $L_D(A(S)) \geq \frac{1}{8}$ .

Characteristics that we need to choose for the selection of model is crucial.

## Bias Variance Trade Off (19)

Bias: [comes from the fact that choosing a model.]

[Too much prior knowledge OR

restricting of choosing the hypothesis doesn't benefit]

→ 2 factors that influence Data set

the model

- Bias

- variance.

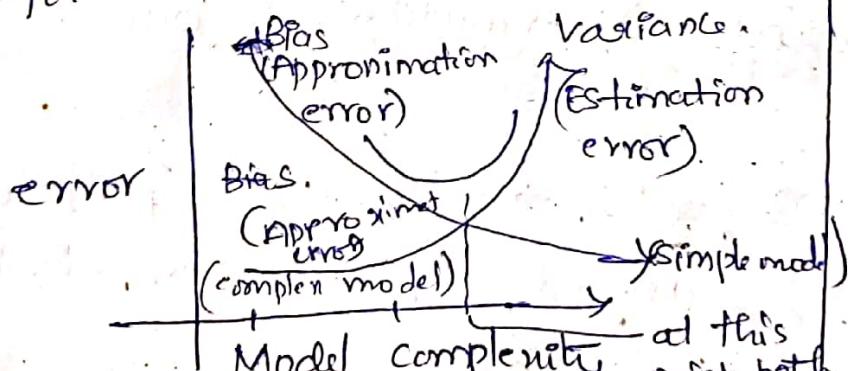
$$E = \text{Approximation Error} + \text{Estimation Error}$$

Approximation Error is Bias.

Estimation Error is Variance.

- what model

- how much it is able to fit in the model.



In Simple model Bias (decrease) that is approximation increases Error increases and In Complex model Variance increases.

Lot of error

$\times \times \times \times \times \times$

$$\theta = \theta_0 + \theta_1 x$$

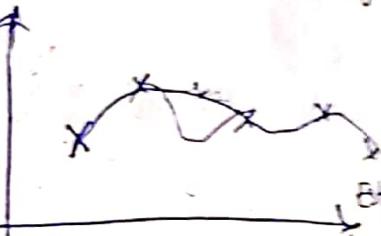
Underfitting happens.

Model - Simple

- Rigid.

(19)

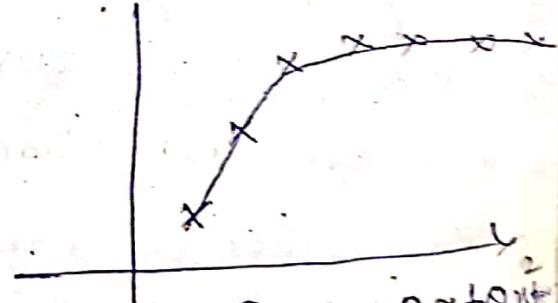
$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$



(In this Approximation error)

$$\theta = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

In this case Empirical Error decreases but the Generalisation or true error Increases.



$$\theta = \theta_0 + \theta_1 x + \theta_2 x^2$$

In this case fit is better seen overfitting and underfitting

(In this

→ NP - Non Deterministic polynomial.  
Problem: Verify the solution in the polynomial time problem solution and verification non polynomial in polynomial time 20

NP complete — Ans yes/no (Decidable problem)

Optimisation → Decision problem.

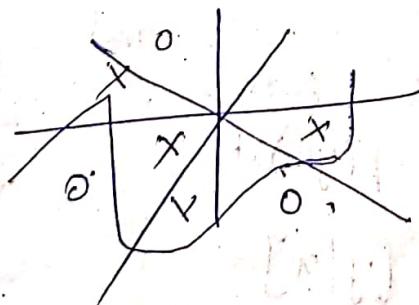
{ So, optimisation Based on constraint cannot be NP-complete they are NP-hard. }

→  $x \rightarrow y$  mapping model such that true error or least error = 0.

→ Bias:

Bias are the simplifying assumptions made by a model to make the target function easier to learn.  
Low Bias - Suggests less assumptions about the form of the target function. Ex: Decision tree, KNN, SVM.  
High Bias - Suggests more assumptions about the form of target function. Ex: Linear Regression, Linear Discriminant analysis, Logistic Regression, etc.

→ Bias



→ Variance:

Variance is the amount that the estimate of a target function will change if different training.

data is used.

Low Variance: Suggests small changes to the estimate of the target function with changes to the training data. Ex: Linear Regression, LDA, Logistic Regression. (21)

High Variance: Suggests large changes to the estimate of the target function with changes to the training data set. Ex: Decision Tree, KNN, Support Vector Machine.

→ Loss Function: (L)

Loss Function (L) for Regression.

$$E(L) = \iint_L L(y, h(x)) p(x, y) dx dy.$$

Actual predicted  
Target Target.

Sum Squared Error:

$$L(y, h(x)) = (y - h(x))^2$$

$$E(L) = \iint (y - h(x))^2 p(x, y) dx dy.$$

→  $\frac{\delta E(L)}{\delta h(x)} = 2 \int (h(x) - y) p(x, y) dy = 0$  (for finding minimum value differentiate and equate it to 0)

Solving for  $y(x)$

$$y(x) = \underbrace{\int y p(x, y) dy}_{p(x)} = \int y p(y|x) dy,$$

$$= E[y|x]$$

Expected y given  $x$ .

$$(h(x) - y)^2 = \{h(x) - E[y|x] + E[y|x] - y\}^2$$

(22)

add and subtract.

$$\rightarrow \{h(x) - E[y|x]\}^2 + \{E[y|x] - y\}^2 + 2(h(x) - E[y|x])(E[y|x] - y).$$

$$\rightarrow E[L] = \underbrace{\int \{h(x) - E[y|x]\}^2 p(x) dx}_{\text{approximation Error}} +$$

$$\underbrace{\int \{E[y|x] - y\}^2 p(x) dx}_{\text{Training data set it had defines } y \text{, not } h(x)} \rightarrow \text{Estimation Error}$$

Training data set it had defines  $y$ , not  $\underline{h(x)}$

$y = \sin(2\pi(x))$  when tried for serial expansion, we try to control the oscillations using the term called Regularisation term.

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_K x^K + \frac{1}{C} \text{, Regularisation Term.}$$

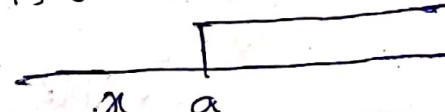
30/01/2020 Learnable Tactile problems

Theory proposed by Vapnik - Chervonenkis (VC):-  
The PAC agnostic PACS etc has finite hypothesis. But, this theory suggested that even though there are infinite hypothesis the problem is learnable.

Theory VC)

Even if the hypothesis class is infinite, if the VC dimension of the class is finite, then that particular model is learnable.

Shattering: There may be threshold function (ex: step function below which is 0 above which is 1)



$h_a$  is a step function.

above  $a'$  is 1, below  $a$  is 0.

2 points Class  $C = \{0, 1\}$ : 2 classes

If one point is used, either 0 or 1 can be assigned as the target.

If two points are  $(0, 0), (0, 1), (1, 0), (1, 1)$

$\begin{matrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{matrix}$  4 types of class labels

If the hypothesis is able to get such type of class labels then it is good to take such hypothesis.

$\rightarrow h_a$  is the  $a$  is the threshold.

Two points  $Q \leq Q_2$ ,  $a = C_1 + 1$

$$h_a(C_1) = 1$$

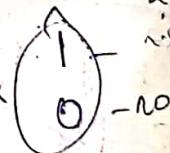
$Q \leq a + 1$  true.

$$\begin{cases} h_a(Q) = 1 & \text{if } Q \leq a \\ h_a(Q) = 0 & \text{if } Q > a \end{cases}$$

$a = Q - 1$  — The next threshold function

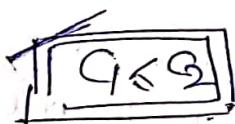
$h_a(C_1) = 0$  since  $Q \leq Q - 1$

$$C_1 \leq a$$



Consider another point  $C_2$

$$a = Q + 1$$



$$Q > a$$

a

$$h_a(C_2) = 1$$

$$h_a(C_1) = 1$$

possible

possible

No threshold can be chosen in this.

$$h_a(C_1) = 0$$

Not possible we cannot choose such threshold.

$$Q > a$$

$$h_a(C_1) = 0$$

$$h_a(C_2) = 0$$

possible

$$C_1, C_2 \approx a, Q \approx a$$

$$h_a(Q) = 1, h_a(C_1) = 1$$

$$h_a(Q) = 1, h_a(C_2) = 1$$

$$h_a(C_2) = 0, C_2 \neq 0$$

possible

possible

Def: Let  $\mathcal{H}$  be a class function from  $\mathcal{X}$  to  $\{0,1\}$ .  
 and let  $C = \{c_1, c_2, \dots, c_n\} \subseteq \mathcal{X}$  The restriction of  $\mathcal{H}$  to  $C$  is the set of functions from  $C$  to  $\{0,1\}$  that can be derived from  $\mathcal{H}$ , that is

(24)

$$\mathcal{H}_C = \{h(c_1), h(c_2), \dots, h(c_n) \mid h \in \mathcal{H}\}$$

where we represent each function from  $C$  to  $\{0,1\}$ .  
 as a vector in  $\{0,1\}^{|C|}$ .

Def (Shattering): A hypothesis class  $\mathcal{H}$  shatters a finite set  $C \subseteq \mathcal{X}$  if the restriction of  $\mathcal{H}$  to  $C$  is the set of all functions from  $C$  to  $\{0,1\}$ .

That is:  $|\mathcal{H}_C| = 2^{|C|}$

Def: A dictionary of a set  $S$  is a partition of  $S$  into two disjoint subsets.

$$S = \{x_1, x_2, \dots, x_m\}$$

$$S^+ = \{x_1, x_2, x_3, \dots, x_8\}$$

$$S^- = \{x_2, x_3, x_4, x_5\}$$

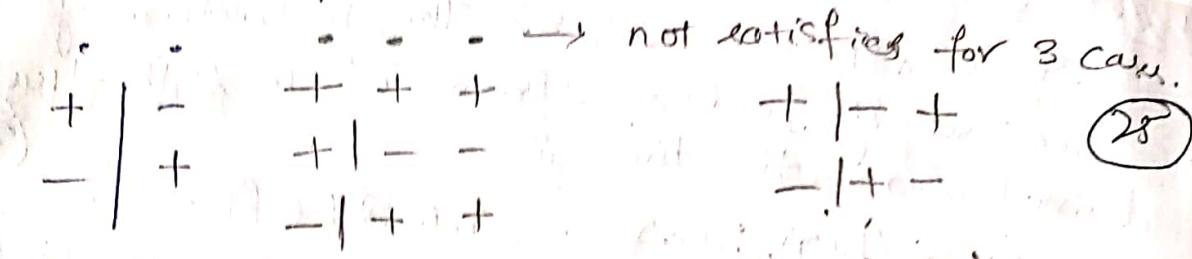
Def: A set of instances  $S$  is shattered by hypothesis space  $\mathcal{H}$  if and only if for every dichotomy of  $S$ , there exists some hypothesis in  $\mathcal{H}$  consistent with this dichotomy.

If  $g \in \mathcal{H}$ , classifies all  $S^+$  as positive  
 and all  $S^-$  as negative.

Ex: 1-dimensional

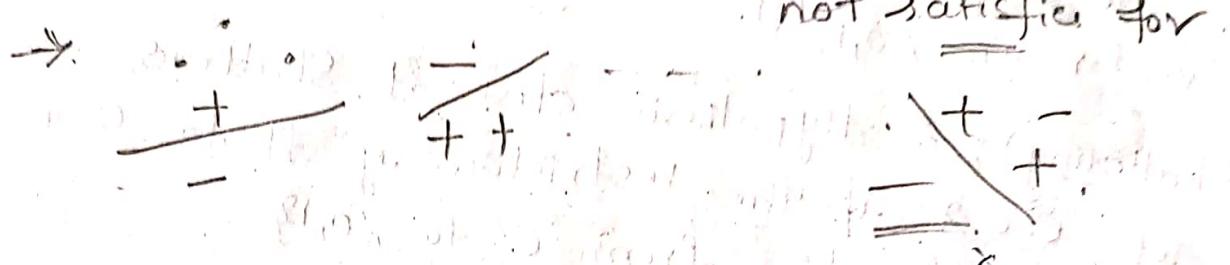
A linear classifier in 1-dimensional

A linear classifier is 1 dimensional



(28)

→ A linear classifier for 2 dimensional (2D)



→ The VC dimension of a hypothesis class  $H$  denoted  $VC\dim(H)$  is the maximal size of a set  $H \subset \mathcal{X}$  that can be shattered by  $H$ . (dimension is  $d$  we can shatter  $d+1$ )

4/02/2020

→ If we take a  $\boxed{++}$   $+ - *$  4 points  
 $2^4 = 16$  labels

$\begin{matrix} + & + \\ - & \square \\ + & + \end{matrix}$   $2^5 = 32$  combinations - Should be able to shatter those 32 Combinations

Linear Predictors, Linear Regression:

→ Set of points, and if there exists linear relationship b/w them then they can be classified as linear regression.

Variable as one

$$y = mx + c$$

Multi variable as

$$y = w_0 + w_1(x_1) + w_2(x_2) + \dots$$

$$y = w_0 + w_1 \phi_1(x) + w_2 \phi_2(x) + \dots + w_{n-1} \phi_{n-1}(x)$$

→ linear combination of non-linear functions.

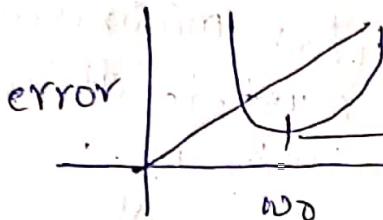
$\phi_i(x)$  can be non linear.

(26)

$$\phi_0(x_0) = 1$$

$$y = \sum w_i \phi_i(x_i)$$

Fitting data is like optimisation problem.

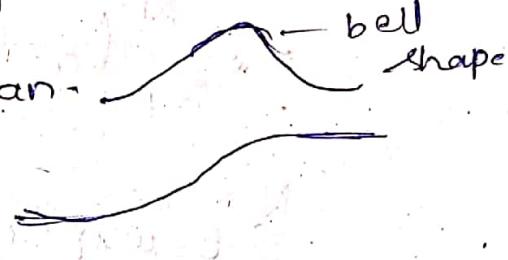


We have to choose such that we can reduce error w.r.t. weights so, it is optimisation problem.

$$\rightarrow \phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\} \text{ - Gaussian}$$

$$\phi_j(x) = \sigma(x - \mu_j) \text{ - sigmoid}$$

$$\sigma(a) = \frac{1}{1 + \exp(-a)} \rightarrow \text{logistic}$$



→ Since, we know the  $w_i$ 's values, if new data comes it would be able to predict the  $y$  value.

$$y = w_0 + w_1 x$$

Functional form error:

$$X = \{x_1, \dots, x_N\}$$

$$x_i = \{x_{i1}, x_{i2}, \dots, x_{iM}\}$$

Sample ( $x_1$ ) -  $\begin{array}{|c|c|c|c|} \hline & & \dots & \\ \hline \end{array}$  M

Sample ( $x_2$ ) -  $\begin{array}{|c|c|c|c|} \hline & & \dots & \\ \hline \end{array}$  M

1 to M attribute

Sample set

$$S = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}$$

Class labels

$$t = \{t_1, t_2, \dots, t_N\}$$

t = target value.

y = Predicted value.

$$y = wx + \epsilon$$

Least Mean Square  $y = wx$  we need to find  $w$ ' value.

(mean square error)  $E_D(w) = \sum_{i=1}^N (t_i - y_i)^2$  — Error.

to find appropriate  $w$  value.

$$\frac{\partial E_D}{\partial w} = \sum_{i=1}^N (t_i - wx_i)^2 \quad (\text{For minimisation we must do the differentiation})$$

$$= 2 \sum_{i=1}^N (t_i - wx_i) \frac{\partial}{\partial w} (wx_i) = -2 \sum_{i=1}^N (t_i - wx_i) x_i = 0$$

$$\sum_{i=1}^N t_i x_i = 0 \quad \text{OR} \quad \sum_{i=1}^N t_i = 0$$

$$w = \frac{t}{x}$$

$$\sum_{i=1}^N (t_i - wx_i) x_i = 0$$

$$tx_i - wx_i^2 = 0$$

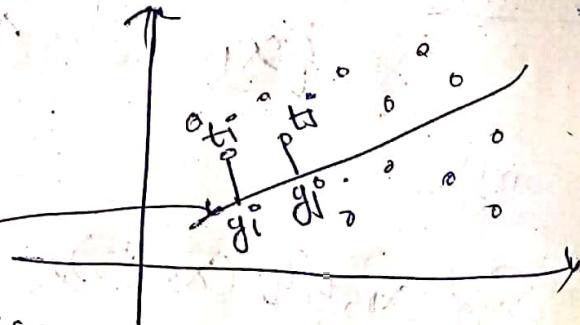
$$w = \frac{\sum t x_i}{\sum x_i^2} = \frac{\text{Covariance } (x_i \text{ and } t)}{\text{Variance } (x)}$$

→ We wanted to get the  $w$  value in order to minimise the error.

$$y_{\text{new}} = wx_{\text{new}}$$

$$y_{\text{new}} = \left( \frac{x_i \cdot t}{x_i^2} \right) x_{\text{new}}$$

for drawing this line



→ Least Means Square Method.

→ Gradient Descent Method.

$$E_D(\omega) = \sum_{n=1}^N (t_n - y_n)^2$$
$$= \sum_{n=1}^N (t_n - (\omega_0 \phi_j(x_n) + \dots + \omega_m \phi_m(x_n)))^2$$

(28)

$$y = \sum_{j=0}^{m-1} \omega_j x_j = \sum_j \omega_j \phi_j(x_j)$$

$\omega^T \phi_j(x_j)$  weight vector  
Input vector

$$\omega = (\omega_0, \omega_1, \dots, \omega_{m-1})^T$$

$$\frac{\partial E_D(\omega)}{\partial \omega_j} = \frac{d}{d \omega_j} \sum_{n=1}^N (t_n - \sum_j \omega_j \phi_j(x_n))^2$$

(In order to do matrix multiplication we have to do transpose)

$$= -2 \sum_{n=1}^N (t_n - \sum_j \omega_j \phi_j(x_n)) \phi_j(x_n)^T$$

$$\Rightarrow = -2 \sum_{n=1}^N (t_n \phi_j(x_n)^T - \sum_j \omega_j \phi_j(x_n)^T) = 0$$

After Rearranging the terms

$$\omega_{ML} = (\phi^T \phi)^{-1} \phi^T \cdot t$$

$\phi$  - matrix

Maximum

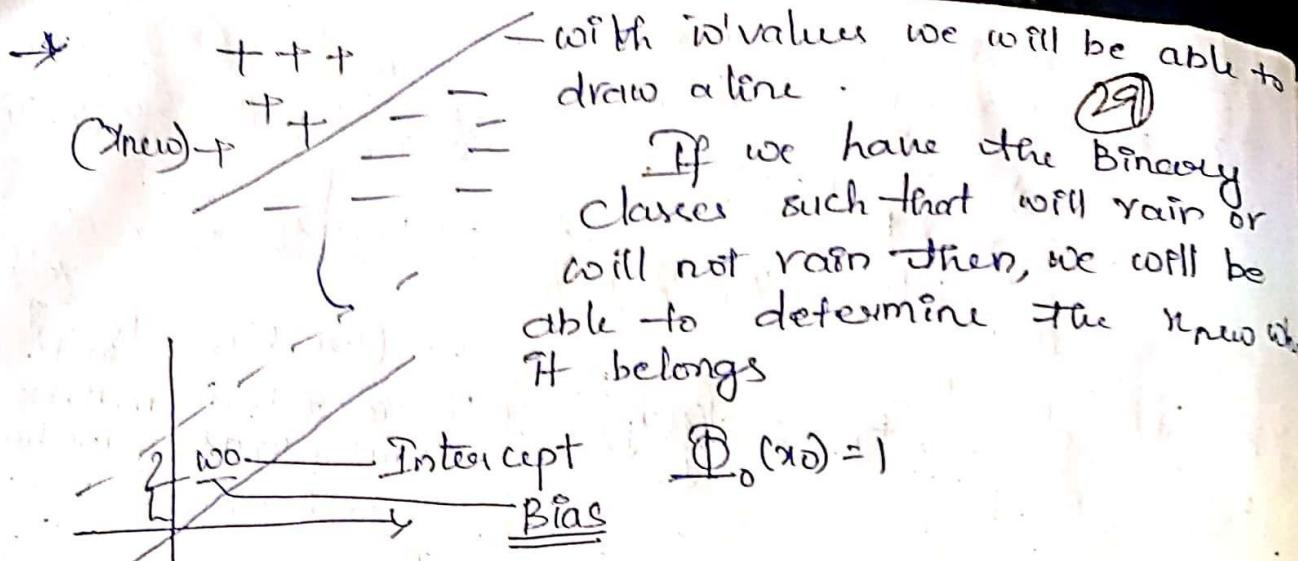
likelihood. These are known as the normal equations for the least square problem.

$\phi$  is an  $N \times M$  matrix called the design matrix whose elements are given by  $\phi_{ij} = \phi_j(x_i)$

$$\underline{\phi} = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_M) & \phi_1(x_M) & \dots & \phi_{M-1}(x_M) \end{pmatrix}$$

Purpose

→  $\phi^+ = (\phi^T \phi)^{-1} \phi^T$  is known as Moore-Penrose Pseudo Inverse Matrix.



With  $w$ ' values we will be able to draw a line. 29  
If we have the binary classes such that will rain or will not rain then, we will be able to determine the new whether it belongs.

### Gradient Descent for linear Regression:

(slope)

Since Least Mean Square involves inverse calculation we don't want to do such complex calculations we choose Gradient method.

$$y_j = \sum w_j \phi(x_j)$$

$$E_p(w) = \sum (t_n - y_j)^2$$

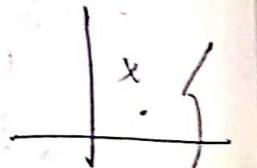
$$\frac{\partial}{\partial w} E_p(w) = \frac{\partial}{\partial w_j} \sum (t_n - y_j)^2$$

$$= -2 \left( \sum t_n - y_j \right) \frac{\partial y_j}{\partial w_j} = -2 \sum (t - y_j) \phi_j'(x_j)$$

Batch learning

Sequential or online learning

as when samples come make changes to the system



where grad  
is highest  
we move at  
that path.  
assuming that we  
may reach our  
target soon.

factor

Types of learning.

→ Gradient Descent Algorithm for linear Regression  
Initialise weights to 0  $w(0) = 0$  (80)

for  $i=1$  to  $\dots$  until convergence (Epochs)

→ predict for each sample  $x_j$  using  $y_k = \sum_{j=0}^N w_j \phi_j(x_k)$

→ compute the gradient descent

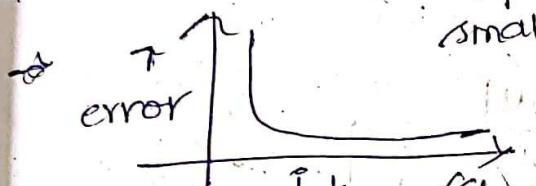
of the loss as  $2 \sum (t_j - y_j) \phi_j(x_j)$  not a matrix simple calculation

→ This is a vector  $g$ .

→ update  $w^{(t+1)} = w^{(t)} - \frac{\eta}{T} g$

small quantity (Learning Parameter)

0 to 1



↑  
Iterations  
(Should make sure not  
to overfit)

→ Regularized Least Squares: (In order to reduce the overfitting effect)

$$E_D(w) + \lambda E_D^{(L)}$$

Total error  $\frac{1}{2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2 + \lambda E_D(w)$   
 $\gamma$  is the Regularisation coefficient

Simplest form of Regularisation:

is  $E_D(w) = \frac{1}{2} w^T w$  - weight decay.

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^T \phi_n(x_n))^2 + \frac{1}{2} w^T w$$

$$w = (I + \gamma \phi^T \phi)^{-1} \phi^T t$$

Bernoulli

Gaussian

1.3 Bernoulli

$$P(x; \phi) = \phi^x (1-\phi)^{1-x}$$

DL

31

6/02/2020

2.4 Gaussian

$$P(x; \mu, \sigma) = \sqrt{\frac{1}{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\rightarrow P(X, \vec{w}; \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2} [X - \mu]^T \Sigma^{-1} [X - \mu]}$$

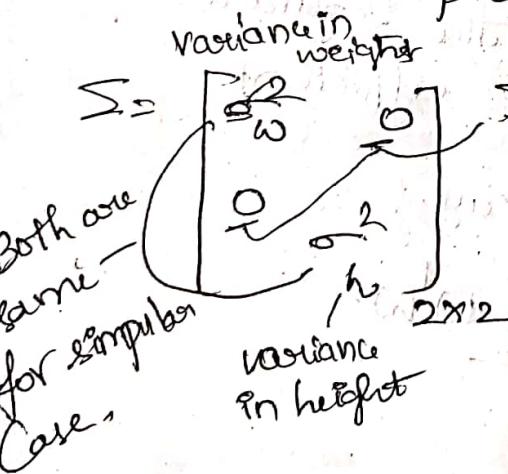
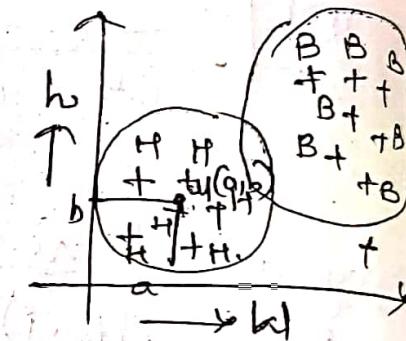
3.4 Exponential

$$P(x; \lambda) = \lambda^x e^{-\lambda x} \quad \text{2x2 matrix}$$

$$4.3 \text{ Laplace. } P(x; \mu, \Sigma) = \frac{1}{\sqrt{2\pi n}} e^{-\frac{1}{2} \frac{|x-\mu|}{\sigma}}$$

$\mu = a, b$   $\Sigma$ -function of  $x$   
in Gaussian

When  $x = a, b$  it gives  
highest or maximum  
probability

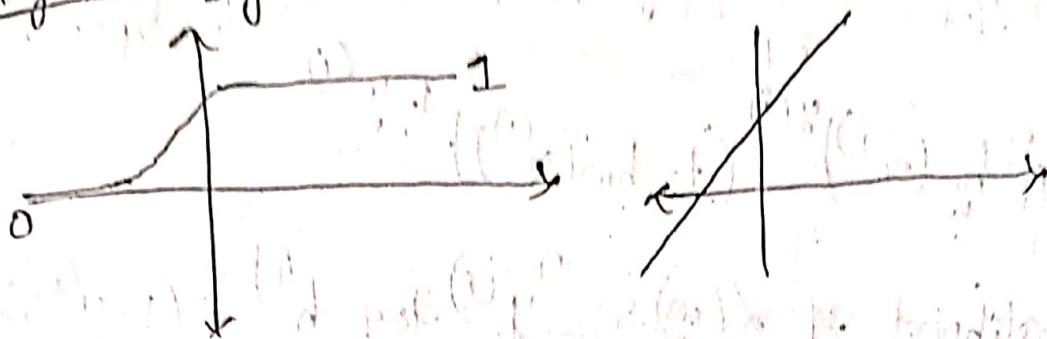


Joint variance in height and weight called as Covariance

$\sigma_w^2$  and  $\sigma_h^2$  are independence, which means Covariance is 0

6/02/2020 → assumption: decision boundary is not a straight line (32)

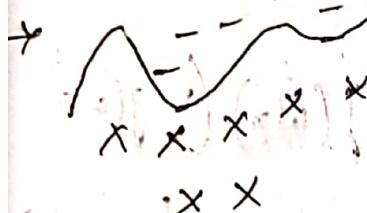
### Logistic Regression:



In logistic function, the value will be b/w 0 and 1.

$$\text{Logistic function} = \frac{1}{1 + e^{-x}}$$

$$f: x \rightarrow (0, 1)$$



How do you learn to differentiate b/w those points?

→ If  $g(z) = \frac{1}{1 + e^{-z}}$  is the logistic function or sigmoid

function.

$$\begin{aligned} g'(z) &= \frac{d}{dz} \left( \frac{1}{1 + e^{-z}} \right) = \frac{1}{1 + e^{-z}} \cdot \frac{(1 + e^{-z})}{e^{-z}} = \frac{1}{1 + e^{-z}} \left( 1 - \frac{1}{1 + e^{-z}} \right) \\ &= g(z)(1 - g(z)) \end{aligned}$$

→ If we assume  $P(y=1|x, \omega) = h_\omega(x)$

$$P(y=0|x, \omega) = 1 - h_\omega(x)$$

This can be written as  $(hw(x))^y (1-hw(x))^{1-y}$ .

(Binomial distribution)

Assuming  $n$  samples are i.i.d., we can write

33

$$l(\omega) = P(Y/x \in \omega) = \prod_{i=1}^n P(Y^{(i)}/x^{(i)}, \omega) \text{ mutually independent}$$
$$= \prod_{i=1}^n h_w(x^{(i)})^{y^{(i)}} (1-h_w(x^{(i)}))^{1-y^{(i)}}$$

$$\text{Log likelihood of } l(\omega) = \sum_{i=1}^n y^{(i)} \log h_w(x^{(i)}) + (1-y^{(i)}) \log (1-h_w(x^{(i)}))$$

Maximisation  $y = \sum_i x_i t_i$

$$h_w(x) = g(w^T x) = \frac{1}{1 + \exp(-x)}$$

→ To maximise likelihood, use stochastic gradient rule:

$$\frac{\partial}{\partial w_j} l(\omega) = \left( y \frac{1}{g(w^T x)} - (1-y) \frac{1}{1-g(w^T x)} \right) \frac{\partial}{\partial w_j} g(w^T x)$$

$t$  = target  
 $x$  = feature  
 $w$  = weight  
 $y$  = target value

$$= \left( y \frac{1}{g(w^T x)} - (1-y) \frac{1}{1-g(w^T x)} \right) g'(w^T x) (1-g(w^T x)) \frac{\partial}{\partial w_j} (w^T x)$$

$$= y (1-g(w^T x)) - (1-y) g(w^T x) x_j$$

$$= (y - h_w(x)) x_j$$

(same like  
 $\hat{t} - y$ )

Classification: To decide whether unseen sample belongs to which particular class based on the learning that has happened on the training data.

(24)

may be binary or multi-class classification.

(1)

$C_1, C_2, \dots, C_k$

(2)

$k$  number of classes for example:

How multiclass classification happens?

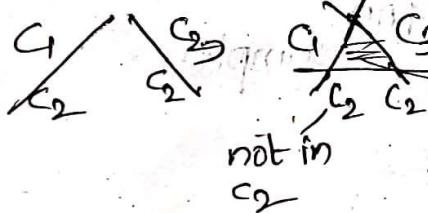
Two methods—

1. One Versus All (making Binary) — (1)  $C_1 / C_2 \dots / C_k$   
Disadv: what is common b/w  $C_2 \dots / C_k$  may not be good enough.

2. Pair wise learning of multiple classes

$(C_1, C_2), (C_1, C_3) \dots, (C_1, C_k)$

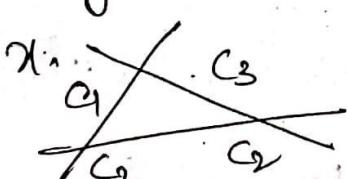
not common keep in other class



this type of gray area may not be classified by any of the classifiers.

By using the binary classifier we may get contradiction.

If  $x$  belongs to  $C_2$  and  $C_2$  and  $C_2$  again we get a contradiction.



11/02/2020 — If this is a Binary classifier  
→ L Class 1      X = class 2

(35)

→ By adding extra features for separation of other classes.

1d 10-samples If we add extra feature density also increases (we want density to be seen)

2d  $10^2$  —

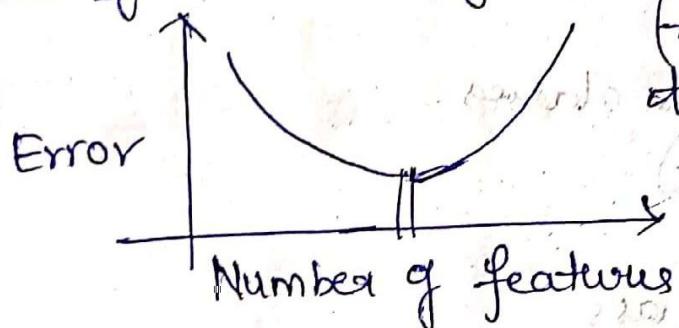
When other features grow the samples

3d 10<sup>3</sup> samples Shouldn't grow ( $10^3$ ) and it's not good

Kind of features

(Most of time we use Euclidean distance for similarity or dissimilarity)

### Curse of Dimensionality:



(As the Dimensionality increases the samples may grow)

so, we may need  $10^{400}$  samples.

NX400

$$P_1 = \dots = 400$$

$$P_2 = \dots = 400$$

$$P_{10} = \dots = 400$$

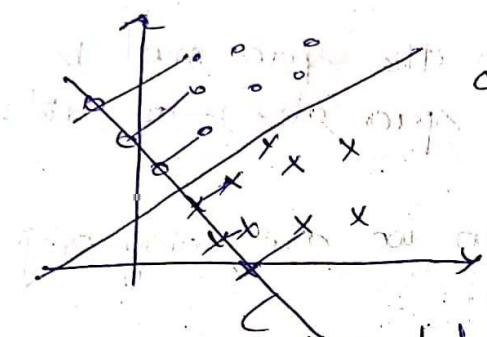
Feature Selection — one way of reduction of dimensions since we select only the particular features

Dimensionality Reduction: PCA, FDA

35  
there  
by require  
also  
be same  
samples  
with  
instances  
easier,  
etc

choose such an axis for which the variation  
in the samples is highest  
(less no of features to project on to the  
line).

No clear way of doing that



In PCA  
can reduce into any  
number of features  
(can fail in some  
cases)

36

may be this helps. (Hence, we can differentiate  
b/w these classes)  
we need to get the direction of line such that it  
will be helpful to differentiate b/w those classes.  
What is the slope/bias needed to be identified this  
is FADA (Feature Dis Fisher's Discriminant Analysis)

Fisher's Linear Discriminant; What is projection you need  
kd (k-dimensions)

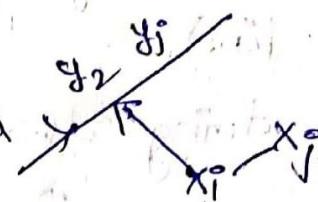
Main Idea: Find projection to a line such that samples  
from different classes are well separated.

Suppose we have 2 classes and  $d$  number of dimensions  
where  $n_1$  samples are from class 1 and  $n_2$  samples are  
from class 2  $n_1 + n_2 = N$   $X = \{x_1, x_2, \dots, x_N\}$

Let the line lie in the direction given by the unit vector  $w$ .

Scalar  $w^t x_i$  is the distance of the projection from the origin.

$$\text{Let } y_i = w^t x_i$$

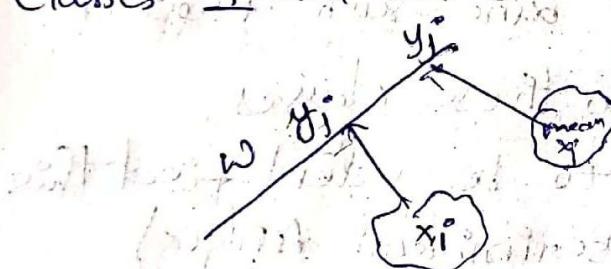


(37)

(Whoever is my neighbours in the space must be also my neighbours in the projected space also) - as said by

Thus we need to take care when we are doing projection in the lower dimensions also.

To measure the separations between projections of classes 1 and 2 (we can make use of mean)



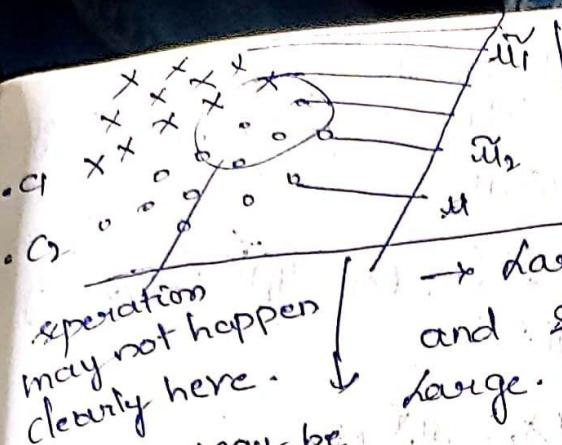
Does it require mean only, or something else?  
let  $\mu_1$  and  $\mu_2$  be the means of classes  $C_1$  and  $C_2$  and  $\tilde{\mu}_1$  and  $\tilde{\mu}_2$  are the means of the projections

If they are well separated then points are also well separated (assumption)

$$\tilde{\mu}_1 = \frac{1}{n_1} \sum_{x_i \in C_1} w^t x_i = w^t \sum_{x_i \in C_1} x_i = w^t \mu_1$$

$$\tilde{\mu}_2 = w^t \mu_2 = \frac{1}{n_2} \sum_{x_i \in C_2} w^t x_i = w^t \sum_{x_i \in C_2} x_i$$

Projection along which we wanted to do projection



So, separation can happen clearly by choosing the correct projection line.

(27) Separation may not happen clearly here. → Large Variations in the data and separation  $\bar{w}_1$  and  $\bar{w}_2$  is always large.

May be this line can be chosen (so, choosing the line is Linear Discriminant Analysis).

Not only mean we should also consider the variances. The  $S_1$  and  $S_2$ . FDA.

The Fisher's Criterion

$$J(k) = \frac{(\bar{w}_1 - \bar{w}_2)^2}{S_1^2 + S_2^2}$$

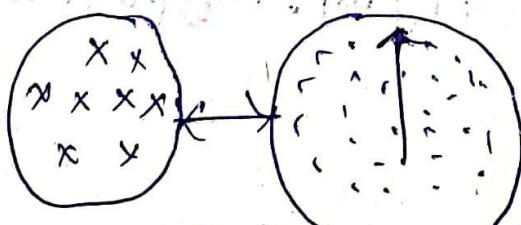
$$\hat{S}_1^2 = \sum_{y_i \in C_1} (y_i - \bar{w}_1)^2$$

$$\hat{S}_2^2 = \sum_{y_i \in C_2} (y_i - \bar{w}_2)^2$$

→  $\hat{S}_1, \hat{S}_2$  - scatter/variance

Scatter in class - must be small.

Scatter across the class - must be as large as possible



$S_W$  = Scatter within the class

$S_B$  = Scatter between the classes,

$S_w = S_1 + S_2$  (Scatter within the class is sum of scatter of class 1 and class 2)

Using  $y_i = w^T x_i$   $\hat{y}_i = w^T \mu_i$

$$\tilde{S}_1^2 = \sum_{y_i \in C_1} (w^T x_i - w^T \mu_1)^2 = \sum_{y_i \in C_1} (w^T (x_i - \mu_1))^2 \quad (39)$$

Var in  
space projected

$$\Rightarrow \sum_{y_i \in C_1} w^T ((x_i - \mu_1)(x_i - \mu_1)^T) w = (x_i - \mu_1)^T w =$$

$$= \sum_{y_i \in C_1} w^T (x_i - \mu_1) \underbrace{(x_i - \mu_1)^T}_{S_1}$$

$$= w^T S_1 w.$$

$$\tilde{S}_2^2 = w^T S_2 w.$$

$$S_1 = \sum_{x_i \in C_1} (x_i - \mu_1)(x_i - \mu_1)^T \quad \text{scatter}$$

$$S_2 = \sum_{x_i \in C_2} (x_i - \mu_2)(x_i - \mu_2)^T$$

$$\tilde{S}_1^2 + \tilde{S}_2^2 = w^T S_1 w + w^T S_2 w = w^T (S_1 + S_2) w = w^T S_w w.$$

Define  $S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$

$$\begin{aligned} \hat{\mu}_1^2 - \hat{\mu}_2^2 &= (w^T \mu_1 - w^T \mu_2)^2 = w^T (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T w \\ &= w^T S_B w. \end{aligned}$$

$$J(w) = \frac{w^T S_w w}{w^T S_B w}$$

$$\frac{d}{dw} J(w) = \frac{d}{dw} (\omega^T S_B w) \omega^T S_B w - \frac{d}{dw} (\omega^T S_B w) \omega^T S_B w$$

$$= \frac{(\omega^T S_B w)^2}{(\omega^T S_B w)^2} \quad (90)$$

$$= \frac{(2S_B w) \omega^T S_B w - (2S_B w) \omega^T S_B w}{(\omega^T S_B w)^2} = 0$$

$$\omega^T S_B w (S_B w) - \omega^T S_B w (S_B w) = 0$$

$$\frac{\omega^T S_B w (S_B w)}{\omega^T S_B w} - \frac{\omega^T S_B w (S_B w)}{\omega^T S_B w} = 0 \Rightarrow S_B w - 1 S_B w = 0$$

$$S_B w = 1 S_B w$$

$$S^T S_B w = 1 w$$

Converting it into Eigen value problem.

Eigen value problem:

$$S^T S_B w = 1 w$$

$S_B X$  for any vector  $X$ , points in the same direction as  $(\lambda_1 - \lambda_2)$

$$S_B X = (\lambda_1 - \lambda_2) (\lambda_1 - \lambda_2)^T X = (\lambda_1 - \lambda_2) (\lambda_1 - \lambda_2) X$$

$$= \lambda (\lambda_1 - \lambda_2)$$

$$S^T S_B w = \lambda w$$

$$S^T S_B (S_B w)$$

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^t$$

$$S_W = \sum_{n \in Q} (x_n - \mu_1)(x_n - \mu_1)^t + \sum_{n \in S_2} (x_n - \mu_2)(x_n - \mu_2)^t$$

$$J(\omega) = \frac{\omega^t S_B \omega}{\omega^t S_W \omega}$$

(4)

$$\underbrace{(\omega^t S_B \omega)}_{\text{Scalar term}} S_W \omega = \underbrace{(\omega^t S_W \omega)}_{\text{Scalar term}} S_B \omega$$

$$S_W \omega = S_B \omega$$

$$(S_W^{-1} S_W) \omega = S_W^{-1} S_B \omega$$

$$\omega \propto S_W^{-1} (\mu_1 - \mu_2)$$

$$S_W^{-1} S_B \omega = \lambda \omega$$

$$S_W^{-1} S_B (S_W^{-1} (\mu_1 - \mu_2)) = S_W^{-1} \propto (\mu_1 - \mu_2) \quad (\text{the projection})$$

$$= \underline{\alpha} \left( \underline{S_W^{-1} (\mu_1 - \mu_2)} \right)$$

Aim: To find out the direction of the line required for

line required for

Eigen value problem  $\rightarrow A\underline{x} = \lambda \underline{x}$  form. this project all the points

$$Q = [Q_1, Q_2] (2, 3) (3, 3) (4, 3) (5, 3)$$

$$Q_2 = [Q_1, Q_2] (2, 1) (3, 1) (3, 2) (5, 3) (6, 3)$$

$$G_1 = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 3 & 3 \\ 4 & 5 \\ 5 & 5 \end{bmatrix}$$

$$G_2 = \begin{bmatrix} 1 & 0 \\ 2 & 1 \\ 3 & 1 \\ 4 & 1 \\ 5 & 1 \end{bmatrix}$$

(42)

①  $\mu_1$  = mean  $G_1 = [3 \quad 3.6]$

$\mu_2$  = mean  $G_2 = [3.3 \quad 2]$

② Compute the Scatter Matrices

$$S_1 = 4 \times \text{cov}(G_1) = \begin{bmatrix} 10 & 8.0 \\ 8.0 & 7.2 \end{bmatrix}$$

$$S_2 = 5 \times \text{cov}(G_2) = \begin{bmatrix} 17.3 & 16 \\ 16 & 16 \end{bmatrix}$$

③ Within the class scatter

$$S_W = S_1 + S_2 = \begin{bmatrix} 27.3 & 24 \\ 24 & 23.2 \end{bmatrix}$$

④ The inverse of  $S_W = S_W^{-1} = \begin{bmatrix} 0.39 & -0.04 \\ -0.04 & 0.47 \end{bmatrix}$

⑤ Optimal line direction  $w = S_W^{-1}(\mu_1 - \mu_2) = \begin{bmatrix} -0.79 \\ 0.89 \end{bmatrix}$

6. Final step is the actual projection of  $x_i$  onto line

$$y_1 = \mathbf{w}^T x_1 = \mathbf{w}_1^T = [-0.79 \quad 0.89] \begin{bmatrix} 1 & 5 \\ 2 & 5 \end{bmatrix} = [0.81, -0.4]$$

$$y_2 = \mathbf{w}^T x_2 = [-0.79 \quad 0.89] \begin{bmatrix} 1 & 6 \\ 2 & 5 \end{bmatrix} = [-0.65, -0.25] \quad (43)$$