

## Decision Tree:-

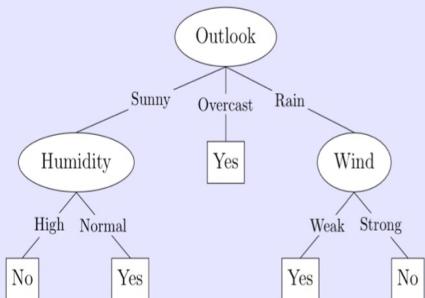
Day	Outlook	Temp	Humidity	Wind	Tennis?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	No
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

= {D}



given a new query you just need  
to pass your data to DT to  
find your output.

→ DT



## Algorithm 1 DECISIONTREETRAIN(*data*, *remaining features*)

```
1: guess  $\leftarrow$  most frequent answer in data // default answer for this data
2: if the labels in data are unambiguous then (only single class in the leaf node)
3:   return LEAF(guess) // base case: no need to split further
4: else if remaining features is empty then (if no features left to split then stop)
5:   return LEAF(guess) // base case: cannot split further
6: else // we need to query more features
7:   { for all f  $\in$  remaining features do
8:     NO  $\leftarrow$  the subset of data on which f=no
9:     YES  $\leftarrow$  the subset of data on which f=yes
10:    score[f]  $\leftarrow$  # of majority vote answers in NO
11:      + # of majority vote answers in YES
12:    } // the accuracy we would get if we only queried on f
13:   end for  $\rightarrow$  Storing score of a feature (Information gain)
14:   f  $\leftarrow$  the feature with maximal score(f)
15:   NO  $\leftarrow$  the subset of data on which f=no
16:   YES  $\leftarrow$  the subset of data on which f=yes
17:   left  $\leftarrow$  DECISIONTREETRAIN(NO, remaining features \ {f})
18:   right  $\leftarrow$  DECISIONTREETRAIN(YES, remaining features \ {f})
19:   return NODE(f, left, right)
20: end if
```

How to calculate these scores?

Entropy:

Let's say  $y = y_1, y_2, \dots, y_K$  is random variable.

$$H(y) = - \sum_{i=1}^K p(y_i) \log_b(p(y_i))$$

for categorical values

$$H(y) = - \int p(y_i) \log_b(p(y_i)) dy$$

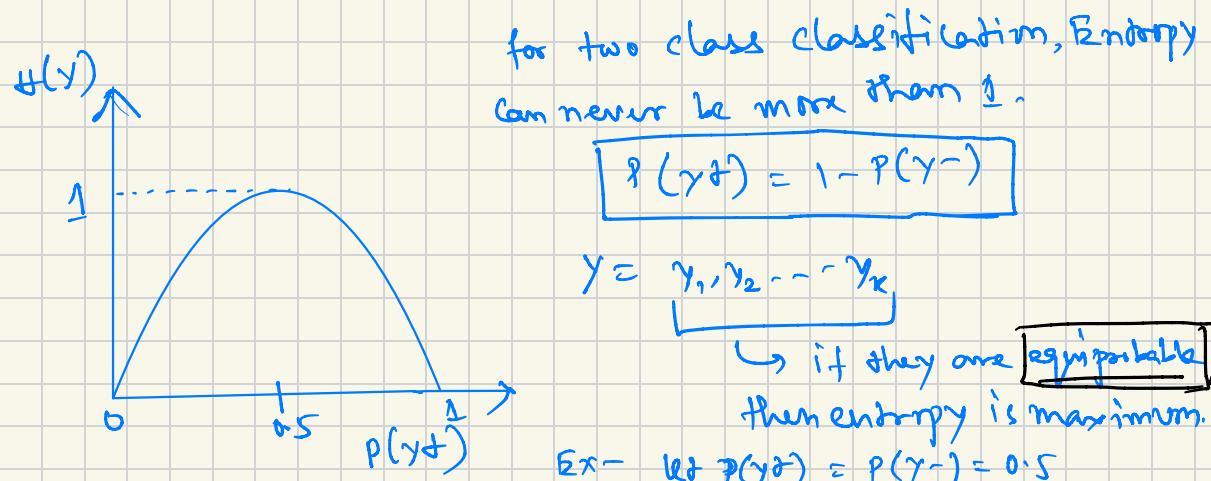
for real values

Here  $b = 2$

$$b = 2.718$$

$$\log_2 = \lg$$

$$\log_e = \ln$$



$$Y = Y_1, Y_2, \dots, Y_k$$

↪ if they are equiprobable then entropy is maximum.

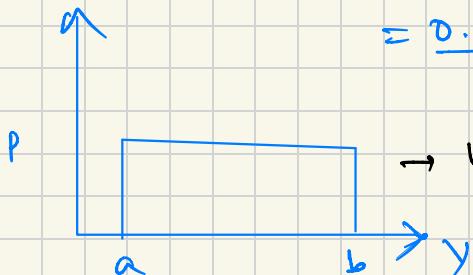
$$\text{Ex- let } P(Y+) = P(Y-) = 0.5$$

$$\begin{aligned} H(Y) &= -P(Y+) \log(P(Y+)) - P(Y-) \log(P(Y-)) \\ &= -0.5 \log 0.5 - 0.5 \log 0.5 \\ &= -(\log 0.5) = -(\log 1 - \log 2) \\ &= \underline{1.} \quad \checkmark \end{aligned}$$

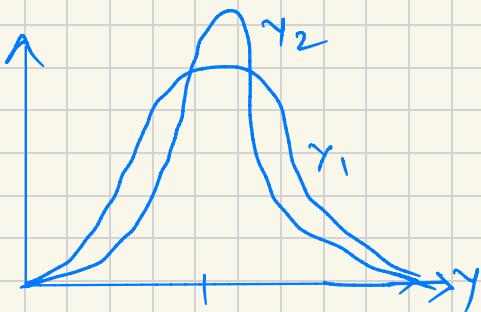
If any one of random variable is most probable & others are 0 then, → entropy is minimal.

$$\text{let } P(Y+) = 1 \quad P(Y-) = 0$$

$$\begin{aligned} H(Y) &= -P(Y+) \log(P(Y+)) - P(Y-) \log(P(Y-)) \\ &= -1 \times 0 - 0 \times \log 0 \\ &= \underline{0.} \end{aligned}$$



→ Uniform distribution, so, if  $y_i$  v  $y_j$  is equiprobable then it has maximum Entropy.



$\gamma_1 \rightarrow$  less peaked  
 $\gamma_2 \rightarrow$  very peaked

So,  $\gamma_1$  has more data variance  
 whereas  $\gamma_2$  has low data variance

So,  $H(\gamma_1) > H(\gamma_2)$

entropy of  $\gamma_1$  always greater than  
 $\gamma_2$ .

## Partition of Data:

initial Entropy

$$H_D(y) = \frac{8}{14} \times \log\left(\frac{8}{14}\right) - \frac{6}{14} \times \log\left(\frac{6}{14}\right)$$

$$= -0.57 \times (-0.8074) - 0.42 \times (-1.2224)$$

$$= 1.46 + .52$$

$$= .98$$

Let,  
 $D \rightarrow D_1$  (5 rows)  $H_{D_1}(y) = 0.97$   
 $D \rightarrow D_2$  (4 rows)  $H_{D_2}(y) = 0$   
 $D \rightarrow D_3$  (5 rows)  $H_{D_3}(y) = 0.97$

## Information gain (IG):

$$\left( \frac{0.97 \times \frac{5}{14}}{1} + 0 \times \frac{4}{14} + 0.97 \times \frac{5}{14} \right) - 0.98$$

weight entropy after  
 partitioning dataset

entropy before  
 partitioning

$$IG(y, var) = H_D(y) - \left[ \sum_{i=1}^k \frac{|D_i|}{|D|} * H_{D_i}(y) \right]$$

◻ Gini Impurity  $\% (I_2)$

$\sim$  similar to entropy

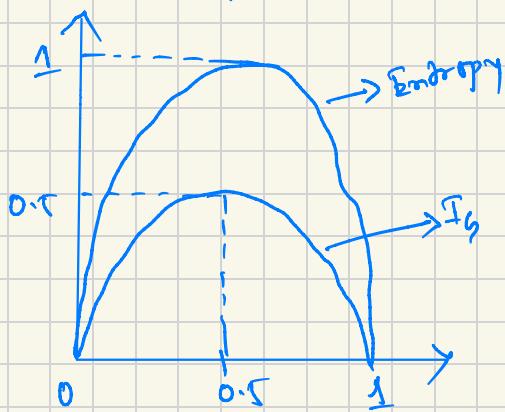
$$I_3(y) \neq I_6(y)$$

for random variable  $X = \gamma_1, \gamma_2 - \gamma_M$

$$I_2(y) = 1 - \sum_{i=1}^k (p(y_i))^2$$

$$\text{Case 1: } P(Y+) = 0.5 \quad I_2(Y) = 1 - \left[ (0.5)^2 + (0.5)^2 \right] \\ P(Y-) = 0.5 \quad = 1 - [0.25 + 0.25]$$

$$\text{Case 2: } P(\gamma_0) = 1 \quad P(\gamma_1) = 0 \quad F_Q(\gamma) = 1 - [1^{\alpha_0} + 0^{\alpha_1}] = 1 - 1 = 0$$



max value for Gini Impurity  
is 0.5.

H(y) is computationally intensive due to log calculation whereas computing square has less computational cost compare to log computation.

So,  $I_5$  is preferable.

## Construct a DT

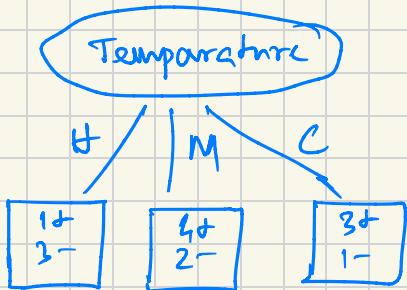
Step 1:  $D = \{y^+, y^-\} = \{8, 6\}$  &  $H(y) = .98$

Day	Outlook	Temp	Humidity	Wind	Tennis?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	No
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Step 2:

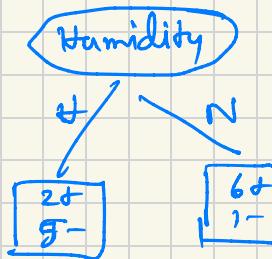
$$\begin{aligned}
 & H(y) = -\frac{5}{14} \left( \frac{2}{5} \log\left(\frac{2}{5}\right) + \frac{3}{5} \log\left(\frac{3}{5}\right) \right) \\
 & - \frac{4}{14} \left( \frac{3}{4} \log\left(\frac{3}{4}\right) + \frac{1}{4} \log\left(\frac{1}{4}\right) \right) \\
 & - \frac{5}{14} \left( \frac{3}{5} \log\left(\frac{3}{5}\right) + \frac{2}{5} \log\left(\frac{2}{5}\right) \right) \Big] - .98 \\
 & = -\frac{5}{2} \left( \frac{2}{5} \log\left(\frac{2}{5}\right) + \frac{3}{5} \log\left(\frac{3}{5}\right) \right) - \frac{4}{14} \left( \frac{3}{4} \log\left(\frac{3}{4}\right) \right. \\
 & \quad \left. + \frac{1}{4} \log\left(\frac{1}{4}\right) \right) \\
 & = -\frac{5}{2}(-.5288 - .4422) + \frac{4}{14}(.3112 - .5) \\
 & = (.3914 + .2312) - .98 \\
 & = (-.6231 - .98) \\
 & = \underline{.35}
 \end{aligned}$$

Step 3:



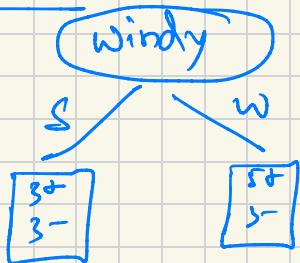
find  $IG$  as above

Step 4:



find  $IG$  as above

Step 5:



find  $IG$  as above

Step 6:

$$IG(y, \text{outlook}) = 1.35$$

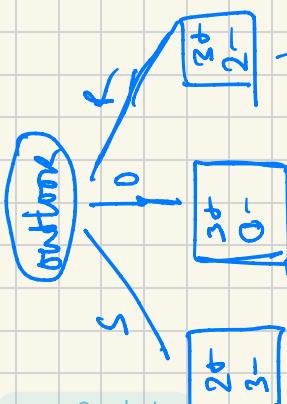
$$IG(y, \text{Temperature}) = -$$

$$IG(y, \text{Humidity}) = -$$

$$IG(y, \text{Windy}) = -$$

We will choose the highest  $IG$  as root.

$$IG(y, f) = b_D(y) - \sum_{i=1}^K \frac{|D_i|}{|D|} * H_{D_i}(y)$$



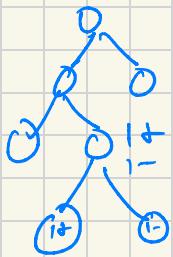
→ no negative points so we can stop here.  
This mode is called pure mode.

We are recursively branching each nodes based on 16.

\*\* three points where to stop growing the tree —

① pure node → stop growing the tree.

② Can't grow the tree anymore due to lack of points



In this case no point of growing the node.  
Suppose, your dataset has 10k points & you have only one points which follows one path. That point may be noisy. (It will overfit)

③ if we are too deep → stop growing the tree,

As, the depth of the tree ↑ → chance of overfitting ↑

depth is small → underfit.

In DT → depth of the tree is a hyperparameter.

↳ we can use  $C_V$  to find it.

Pruning: Remove subtrees for better generalization (decrease variance)

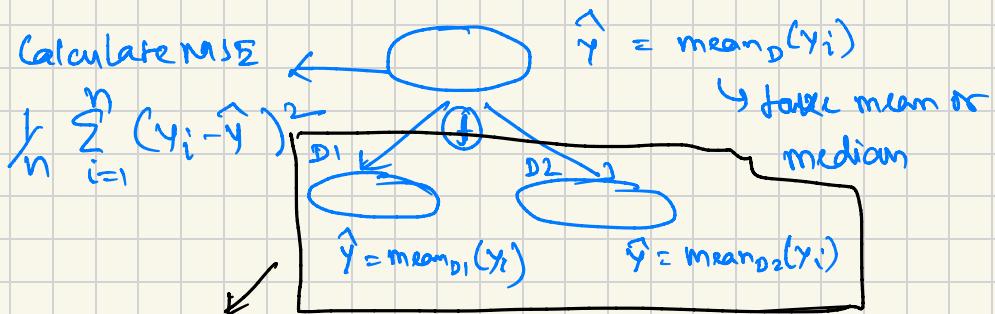
Prepruning: Early stopping

Postpruning: Grow the whole tree then prune subtrees which overfit on the pruning set

Prepruning is faster but postpruning is more accurate.

## DT for Regression :

- IG used for classification
- MSE or MAE is used for regression



find weighted mean square

$$w_1 \text{MSE}_{D_1} + w_2 \text{MSE}_{D_2}$$

Now subtract  $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2 - [w_1 \text{MSE}_{D_1} + w_2 \text{MSE}_{D_2}]$

whichever split minimizes MSE, should consider for split.

# Random Forest :

---

## **Algorithm 1: Pseudo code for the random forest algorithm**

---

To generate  $c$  classifiers:

**for**  $i = 1$  to  $c$  **do**

    Randomly sample the training data  $D$  with replacement to produce  $D_i$

    Create a root node,  $N_i$  containing  $D_i$

    Call BuildTree( $N_i$ )

**end for**

**BuildTree(N):**

if  $N$  contains instances of only one class **then**

**return**

**else**

    Randomly select  $x\%$  of the possible splitting features in  $N$

    Select the feature  $F$  with the highest information gain to split on

    Create  $f$  child nodes of  $N$ ,  $N_1, \dots, N_f$ , where  $F$  has  $f$  possible values ( $F_1, \dots, F_f$ )

**for**  $i = 1$  to  $f$  **do**

        Set the contents of  $N_i$  to  $D_i$ , where  $D_i$  is all instances in  $N$  that match

$F_i$

        Call BuildTree( $N_i$ )

**end for**

**end if**

## AdaBoost Algorithm:

Given  $n$  examples  $(x_i, y_i)$  where  $x \in X \subseteq \mathbb{R}^d$ ,  $y \in \{-1, 1\}$

① Initialize  $D_1(i) = 1/n$

$\downarrow$   $D_1$  is first dataset.

② for  $t$  in  $(1 \dots T)$ :

i) Train weak classifier on Distribution

$$D_t(i), h_t(x) : x \rightarrow \{-1, 1\}$$

ii) Find error with  $h_t(x)$

$$\epsilon_t = \sum_{i=1}^n D_t(i) [y_i \neq h_t(x_i)]$$

if  $\epsilon_t < 0.5$  then stop here

else:

find weights for each record of the dataset

$$\alpha_t = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$$

iii) update  $D_{t+1}(i) = \frac{D_t(i) \exp \{-\alpha_t y_i h_t(x_i)\}}{Z_t}$

$$\text{Where } Z_t = \sum_{i=1}^n D_t(i) \exp \{-\alpha_t y_i h_t(x_i)\}$$

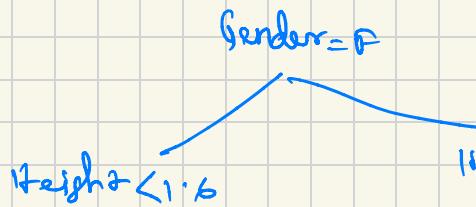
output classifier is

$$h(x) \geq \operatorname{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$$

## Gradient boosting Algorithm :

height	Favorite color	Gender	weight (kg)	1st residual
1.6	Blue	Male	88	16.8
1.6	Green	Female	76	4.8
1.5	Blue	Female	56	-15.2
1.8	Red	Male	73	1.8
1.5	Green	Male	77	5.8
1.4	Blue	Female	57	-14.2

$$\text{avg weight} = \underline{11.2}$$



## Dimensionality Reduction

volume of a d-dim hyperspace

$$V_d(r) = \frac{r^d \pi^{d/2}}{\Gamma(\frac{d}{2} + 1)}$$

$$V_3(r) = \frac{r^3 \pi^{3/2}}{\Gamma(\frac{5}{2})}$$

$$\Gamma(\frac{5}{2}) = \int_0^{\frac{\pi}{2}} e^{-x} x^{\frac{5}{2}-1} dx$$

$$\Gamma(1 + \frac{3}{2})$$

$$V_3(r) = \frac{r^3 \pi^{3/2}}{\frac{3}{4} \pi^{1/2}}$$

$$= \frac{3}{2} \sqrt{\frac{3}{2}}$$

$$V_3 r = \frac{4}{3} \pi r^3 \quad \checkmark$$

$$= \frac{3}{2} \times \frac{1}{2} \sqrt{\pi}$$

$$= \frac{3}{4} \sqrt{\pi}$$

$$\frac{V_d(\text{spare})}{V_d(\text{cube})} = \frac{r^d \pi^{d/2}}{\frac{(2\pi)^d}{2^d} \Gamma(\frac{d}{2})}$$

## Gamma function

$$\Gamma(n) = \int_0^\infty e^{-x} x^{n-1} dx$$

$$\Gamma(1) = 1, \quad \Gamma(\frac{1}{2}) = \sqrt{\pi}$$

$$\Gamma(n+1) = n \Gamma(n)$$

$$\Gamma(n+1) = n!$$

$$\textcircled{a} \quad \Gamma(n) \Gamma(1-n) = \frac{\pi}{\sin n\pi}$$

two way we can solve PCA

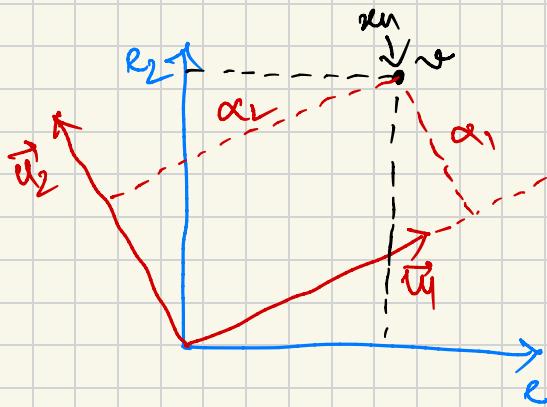
- ① Maximum variance formulation
- ② Minimum Error formulation

### Minimum Error Formulation

PCA can be obtained by minimizing the reconstruction error.

Find a transformation that minimizes?

$$J(x, \tilde{x}) = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|$$



$$v = \alpha_1 u_1 + \alpha_2 u_2 \dots$$

$$x_n = \sum_{i=1}^D \alpha_{ni} u_i$$

$$\alpha_{ni} = x_n^T u_i$$

Therefore, we get the projection onto a new orthonormal basis.

$$x_n = \sum_{i=1}^D (x_n^T u_i) u_i$$

$$u_i^T u_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

For reconstruction, we can use the first M elements of the basis and a shared offset for the remaining dimensions.

$$\tilde{x}_n = \sum_{i=1}^M (x_n^T u_i) u_i + \sum_{i=M+1}^D b_i u_i$$

$$\begin{aligned}
 x_n - \tilde{x}_n &= \sum_{i=1}^D (x_n^T u_i) u_i - \sum_{i=1}^M (x_n^T u_i) u_i - \sum_{i=M+1}^D b_i u_i \\
 &= \sum_{i=M+1}^D (x_n^T u_i) u_i - \sum_{i=M+1}^D b_i u_i \\
 &= \sum_{i=M+1}^D (x_n^T u_i) u_i - b_i u_i
 \end{aligned}$$

$$\begin{aligned}
 J(x, \tilde{x}) &= \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 \\
 &= \frac{1}{N} \sum_{n=1}^N \left\| \sum_{i=M+1}^D (x_n^T u_i) u_i - b_i u_i \right\|^2
 \end{aligned}$$

Now we need to minimize  $u_i$  &  $b_i$ .

$$\begin{aligned}
 \frac{\partial L(x, \tilde{x})}{\partial b_i} &= -\frac{1}{N} \sum_{n=1}^N 2 \left( \sum_{i=M+1}^D (x_n^T u_i) u_i - b_i u_i \right) u_i \geq 0 \\
 &= \frac{2}{N} \sum_{n=1}^N \left( \sum_{i=M+1}^D x_n^T u_i - b_i \right) = 0 \\
 \Rightarrow & \boxed{b_i = \sum_{i=M+1}^D \tilde{x}^T u_i}
 \end{aligned}$$

$$J(x, \tilde{x}) = \frac{1}{N} \sum_{n=1}^N \left\| \sum_{i=M+1}^D (x_n^T u_i - \tilde{x}^T u_i) u_i \right\|^2$$

$$\begin{aligned} \frac{\partial J(x, \tilde{x})}{\partial u_i} &= \frac{1}{N} \sum_{n=1}^N 2 \left[ \sum_{i=M+1}^D (x_n^T u_i - \tilde{x}^T u_i) u_i (x_n^T - \tilde{x}^T) \right] \\ &\Rightarrow \frac{2}{N} \sum_{n=1}^N \sum_{i=M+1}^D (x_n^T u_i^T u_i - \tilde{x}^T u_i^T u_i) (x_n^T - \tilde{x}^T) \\ &= \frac{2}{N} \sum_{n=1}^N \sum_{i=M+1}^D u_i^T (x_n^T - \tilde{x}^T) (x_n^T - \tilde{x}^T) u_i = 0 \end{aligned}$$

$$\Rightarrow \sum_{i=M+1}^D u_i^T S u_i = 0$$

so, if we put  $u_i = 0$  then  $J(x, \tilde{x})$  will be minimized

so, we need to put constraint

$$J(x, \tilde{x}) = \sum_{i=M+1}^D u_i^T S u_i + \lambda_i (1 - u_i^T u_i)$$

we know  $u_i$  is orthonormal.

for  $M=0 \ L=1$

$$\frac{\partial J(x, \tilde{x})}{\partial u_1} = 2(u_1^T S u_1) + \lambda_1 \frac{(1 - u_1^T u_1)}{\partial u_1}$$

$$= u_1^T S - \lambda_1 u_1 = 0$$

$$\Rightarrow u_1^T S = \lambda_1 u_1$$

Here  $u_1$  is the corresponding eigenvector of eigenvalue  $\lambda_1$ .

Our intention is to find top  $(D-M)$  eigenvalues so that  $J$  will be minimized.

# K-Means clustering:

Consider this as a minimization problem.

Data:  $X = \{x_1, x_2, \dots, x_n\} \quad x_n \in \mathbb{R}^D$

Goal: Partition into K clusters by minimizing

$$J = \sum_{i=1}^n \sum_{k=1}^K z_{ik} \|x_n - \mu_k\|^2$$

find cluster assignment (latent var)

$$z_{ik} \in \{0, 1\}$$

and

$$\text{Cluster means } \mu_k \in \mathbb{R}^D$$

$z_{ik} = 1$  iff point n is assigned to cluster k.

Algorithm (EM algorithm)

$$z_{ik} = \begin{bmatrix} z_{i1} \\ z_{i2} \\ \vdots \\ z_{ik} \end{bmatrix} \in \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

1-hot-encoding.

- ① Initialize with a random  $\mu_k \in \mathbb{R}^D$
- ② Repeat until convergence  
→ find the assignment (fixed means)

$$z_{ik} = \begin{cases} 1 & \text{if } k = \arg\min \|x_n - \mu_k\|^2 \\ 0 & \text{otherwise} \end{cases}$$

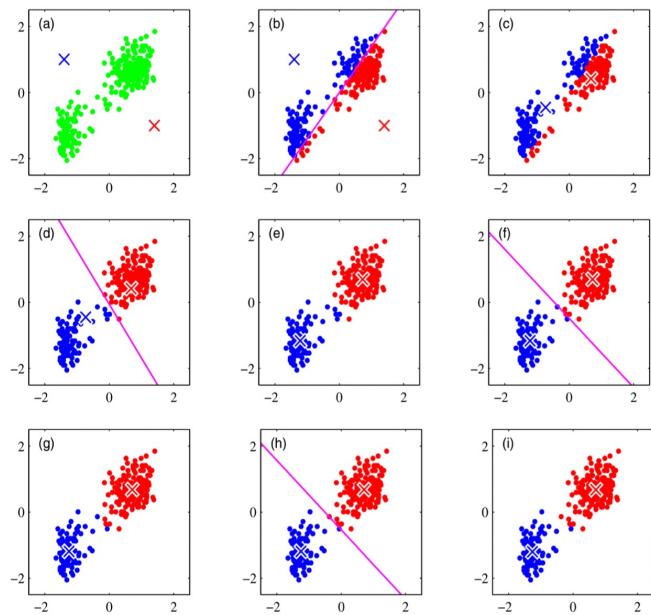
- find the means (fixed assignment)

$$\mu_k = \frac{\sum z_{ik} x_n}{\sum z_{ik}}$$

$\Sigma M \rightarrow$  Expectation

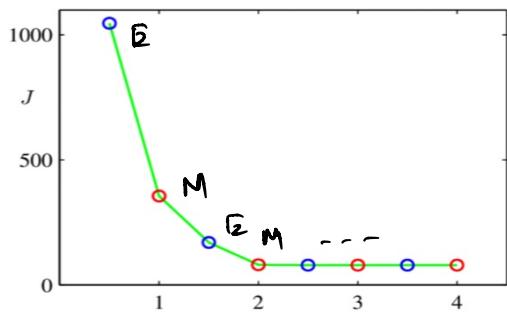
Maximization.

J is non-convex for  
 $\mu_1$  and  $\mu_2$  together  
so K-means converges  
to a local minima.



**Figure 9.1** Illustration of the  $K$ -means algorithm using the re-scaled Old Faithful data set. (a) Green points denote the data set in a two-dimensional Euclidean space. The initial choices for centres  $\mu_1$  and  $\mu_2$  are shown by the red and blue crosses, respectively. (b) In the initial E step, each data point is assigned either to the red cluster or to the blue cluster, according to which cluster centre is nearer. This is equivalent to classifying the points according to which side of the perpendicular bisector of the two cluster centres, shown by the magenta line, they lie on. (c) In the subsequent M step, each cluster centre is re-computed to be the mean of the points assigned to the corresponding cluster. (d)–(i) show successive E and M steps through to final convergence of the algorithm.

**Figure 9.2** Plot of the cost function  $J$  given by (9.1) after each E step (blue points) and M step (red points) of the  $K$ -means algorithm for the example shown in Figure 9.1. The algorithm has converged after the third M step, and the final EM cycle produces no changes in either the assignments or the prototype vectors.



$$J = \sum_{h=1}^n \sum_{k=1}^K z_{hk} \|x_n - \mu_k\|^2$$

It is a convex function for  $\mu_k$  if  $z_{hk}$  is fixed.

$$\text{so, } \frac{\partial J}{\partial \mu_k} = -2 \sum_{n=1}^N z_{hk} (x_n - \mu_k)^T = 0$$

$$\Rightarrow \sum_{n=1}^N z_{nk} x_n - m_k \sum_{n=1}^N z_{nk} = 0$$

$\Rightarrow$

$$m_k = \frac{\sum_{n=1}^N z_{nk} x_n}{\sum_{n=1}^N z_{nk}}$$

Summing up the datapoints within  $k$  class

no of points within  $k$  class.