

# HW 2 Li Chen Liang

1. a)  $x \ y \ z$   $E(x) = (1+4)/2 = 2.5$   $Cov(x, y) = Cov(y, x) = E(xy) - E(x)E(y)$   
 $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$   $E(y) = (2+5)/2 = 3.5$   $= \left( \frac{1 \cdot 2 + 4 \cdot 5}{2} \right) - 2.5(3.5)$   
 $E(z) = (3+6)/2 = 4.5$   $= 2.25$

$Cov(x, z) = Cov(z, x) = E(xz) - E(x)E(z) = 2.25$   
 $Cov(y, z) = Cov(z, y) = E(yz) - E(y)E(z) = 2.25$   
 $Cov(x, x) = E(x^2) - E(x)^2$   
 $= \left( \frac{1^2 + 4^2}{2} \right) - (2.5)^2 = 2.25$

$\begin{matrix} & x & y & z \\ x & \begin{bmatrix} 2.25 & 2.25 & 2.25 \\ 2.25 & 2.25 & 2.25 \\ 2.25 & 2.25 & 2.25 \end{bmatrix} \end{matrix}$  Covariance matrix

$Cov(y, y) = E(y^2) - E(y)^2 = 2.25$

$Cov(z, z) = E(z^2) - E(z)^2 = 2.25$

b) A zero eigenvalue means there is no variance between that eigenvector and the data.

c) If two eigenvalues are the same, then the variance between the corresponding eigenvectors to the data are the same. In this case we have to decide which eigenvalue to choose for the PCA.

2.

a) Yes, greedy algorithm will essentially reach the purest splitting.  
 i.e. There are  $N$  students with  $N$  IDs, then there can be  $N$  branches achieving purest split.

b) Assume  $N$  values with equal possibility  $1/N$  then the entropy is  
 $\underbrace{-1/N \log_2 1/N - 1/N \log_2 1/N \dots - 1/N \log_2 1/N}_N \Rightarrow N \cdot (-1/N \log_2 1/N) = \log_2 N$

- c) Information gain is used to reduce the entropy by splitting into  $k$  partitions
- $$\text{Gain}_{\text{split}} = \text{old entropy} - \text{new entropy} = \text{Entropy}(P) - \left( \sum_{i=1}^k \frac{n_i}{N} \text{Entropy}(i) \right)$$

Assume there are  $N$  values and we split into  $k=N$  partitions and each partition has only 1 value, then  $\text{Gain} = \text{Entropy}(P) - 0$  will be the maximum Gain that can be achieved.

Since SplitInfo is the weighted total of  $-\log \frac{n_i}{N}$ . When we split into  $k=N$  partitions with  $n_i=1$ , then  $-\log \frac{n_i}{N}$  goes to  $\infty$  and  $\text{Gain ratio} = \frac{\text{Gain}_{\text{split}}}{\infty} = 0$

- d) Pre-Pruning: Using certain threshold or condition to stop in order to prevent the tree from fully grown

Post-Pruning: Let decision tree to fully grow then trim from bottom up.

3. a)  $\text{precision} = \frac{5}{5+20} = 0.2$      $\text{recall} = \frac{5}{5+15} = 0.25$

$$\text{F-measure} = \frac{2(5)}{2(5)+15+20} = 0.22$$

- b) ROC curve shows the performance tradeoff between True Positive and False Positive so  $(0,0)$  means we make everything negative so it's not detecting any positives and  $(1,1)$  means we make everything positive and it's detecting all positives. The diagonal line on the ROC curve indicates random classifier.

4. a) Try each  $k$  and record their accuracy. Pick the  $k$  with highest accuracy.  
If  $k$  is too small, it may contain too little of points which might be noise.  
If  $k$  is too large, it may contain points that belong to other classes

- b) Since we have to calculate the distance with each point in the training set, then it will take  $O(n)$  time.

5.

a)  $P(\text{Disease}) = 0.11$   $P(\text{no disease}) = 0.89$

$P(\text{Positive}|\text{disease}) = 0.9$   $P(\text{Positive}|\text{no disease}) = 0.02$

$$P(\text{Disease}|\text{Positive}) = \frac{P(\text{Positive}|\text{Disease})P(\text{Disease})}{P(\text{Positive})}$$

$$\begin{aligned} &= \frac{P(\text{Positive}|\text{Disease})P(\text{Disease}) + P(\text{Positive}|\text{No disease})P(\text{No disease})}{0.9 \times 0.11 + 0.02 \times 0.89} \\ &= \frac{0.9 \times 0.11}{0.9 \times 0.11 + 0.02 \times 0.89} = 0.848 \end{aligned}$$

b) Discretize, Two way Split, Probability density estimation

c) Given the weather condition of a part of sea, such as wind, wave, rain, etc. each day, decide whether should permit a ship to leave the dock.

d) Naive Bayes is simple and efficient for most of the classification tasks. It's robust to noise and irrelevant points. It also handle missing value pretty well. Duplicate data will make the trained model bias toward duplicates.

b. Accuracy  $\approx 80\%$