

SC4000/CZ4041/CE4041: Machine Learning

Lecture 2 Tutorial Question Sets

$$\begin{aligned} P(B=1|A=1) &= P(C=0, B=1|A=1) + P(C=1, B=1|A=1) \\ 0.3 &\equiv P(C=0, B=1|A=1) + 0.1 \\ P(C=0, B=1, A=1) &\leq 0.1 \end{aligned}$$

$$P(A=1) = 1 - P(A=0) = 0.7$$

$$P(A=1|B=0) + P(B=1|A=1) = P(A=1) \quad P(B=1, A=1) = 0.3$$

Question 1: Suppose A , B and C are three variables of binary values (0 or 1). Given the probabilities $P(A = 1, B = 0) = 0.4$, $P(A = 0) = 0.3$, and $P(A = 1, B = 1, C = 1) = 0.1$, compute the following probabilities:

$$1. P(B = 1|A = 1) = \frac{P(C=1, A=1)}{P(A=1)} \text{ or } \frac{P(A=1|B=1) \cdot P(B=1)}{P(A=1)}$$

$$2. P(C = 0|B = 1, A = 1)$$

$$\frac{P(C=0, B=1, A=1)}{P(B=1, A=1)} = \frac{P(C=0, B=1, A=1)}{0.3} = 0.2$$

Question 2: Suppose that if a person has lung cancer, his/her probability of having gene X is 0.9, and if a person does not have lung cancer, his/her probability of having gene X is 0.2. The probability of a person having lung cancer is 0.01. Now, we know that a patient A has gene X . $P(Z) = 1$

marginal probability

1. Use Bayesian decision theory with 0/1 loss to predict whether the patient A has lung cancer or not.

2. Consider that costs of misclassification are different. Assume that the cost for correct decisions is 0, the cost of misclassifying a person who does not have lung cancer to be a patient with lung cancer is 0.007, and the cost of misclassifying a person who has lung cancer to be a healthy person is 1. Please use Bayesian decision theory with the predefined loss to predict whether the patient A has lung cancer or not.

$$Y \stackrel{\text{given}}{\sim} P(Z=1|Y=1) = 0.9 \rightarrow P(Z=0|Y=1) = 0.1$$

$$Z \stackrel{\text{given}}{\sim} P(Z=1|Y=0) = 0.2 \rightarrow P(Z=0|Y=0) = 0.8$$

marginal probability

$$P(Y=1) = 0.01 \rightarrow P(Y=0) = 0.99$$

2.1) predict $Y=1$ if $P(Y=1|Z=1) > P(Y=0|Z=1)$

$$\frac{P(Z=1|Y=1)P(Y=1)}{P(Z=1)} > \frac{P(Z=1|Y=0)P(Y=0)}{P(Z=1)}$$

$$0.9 \cdot 0.01 > 0.2 \cdot 0.99 \\ 0.009 > 0.198 \quad \text{false}$$

Patient A
doesn't have
lung cancer

$$R(a_0|x) = \sum_{C=0}^{C=1} \lambda_C P(Y=C|x)$$

$$\lambda_{00} = 0$$

$$R(a_1|Z=1) = \lambda_{11} \times P(Y=1|Z=1) + \lambda_{10} \times P(Y=0|Z=1)$$

$$\lambda_{11} = 0$$

$$= 0 + 0.007 \times \frac{0.198}{P(Z=1)}$$

$$\lambda_{10} = 0.007$$

$$= \frac{0.0014}{P(Z=1)} < R(a_0|Z=1) = \frac{0.009}{P(Z=1)}$$

$$\lambda_{01} = 1$$

Predicted:
Patient A
has lung
cancer

Question 1

Question 1: Suppose A , B and C are three variables of binary values (0 or 1). Given the probabilities $P(A = 1, B = 0) = 0.4$, $P(A = 0) = 0.3$, and $P(A = 1, B = 1, C = 1) = 0.1$, compute the following probabilities:

1. $P(B = 1|A = 1)$.

2. $P(C = 0|B = 1, A = 1)$.

1.1

- $P(B = 1|A = 1) = \frac{P(B=1,A=1)}{P(A=1)}$
- $P(A = 1) = 1 - P(A = 0) = 1 - 0.3 = 0.7$
- $P(A = 1) = 0.7$
 $= P(A = 1, B = 0) + P(A = 1, B = 1)$
 $= 0.4 + P(A = 1, B = 1)$
- $P(A = 1, B = 1) = 0.7 - 0.4 = 0.3$
- Thus, $P(B = 1|A = 1) = \frac{0.3}{0.7} = \frac{3}{7}$

1.2

- $P(C = 0|B = 1, A = 1) = \frac{P(C=0,B=1,A=1)}{P(B=1,A=1)} = \frac{P(C=0,B=1,A=1)}{0.3}$
- $P(B = 1, A = 1) = 0.3$
 $= P(C = 0, B = 1, A = 1) + P(C = 1, B = 1, A = 1)$
 $= P(C = 0, B = 1, A = 1) + 0.1$
- $P(C = 0, B = 1, A = 1) = 0.3 - 0.1 = 0.2$
- Thus, $P(C = 0|B = 1, A = 1) = \frac{0.2}{0.3} = \frac{2}{3}$

Question 2

Question 2: Suppose that if a person has lung cancer, his/her probability of having gene X is 0.9, and if a person does not have lung cancer, his/her probability of having gene X is 0.2. The probability of a person having lung cancer is 0.01. Now, we know that a patient A has gene X .

1. Use Bayesian decision theory with 0/1 loss to predict whether the patient A has lung cancer or not.
2. Consider that costs of misclassification are different. Assume that the cost for correct decisions is 0, the cost of misclassifying a person who does not have lung cancer to be a patient with lung cancer is 0.007, and the cost of misclassifying a person who has lung cancer to be a healthy person is 1. Please use Bayesian decision theory with the predefined loss to predict whether the patient A has lung cancer or not.

2.1

If a person has lung cancer, his/her probability of having gene X is 0.9

$$P(Z = 1|Y = 1) = 0.9$$



$$P(Z = 0|Y = 1) = 0.1$$

If a person does not have lung cancer, his/her probability of having gene X is 0.2

$$P(Z = 1|Y = 0) = 0.2$$



$$P(Z = 0|Y = 0) = 0.8$$

The probability of a person having lung cancer is 0.01

$$P(Y = 1) = 0.01$$



$$P(Y = 0) = 0.99$$

2.2

λ_{ik} : loss (or cost) of action a_i that predicts $Y = i$ while the true label is k

The cost for correct decisions is 0
 $\lambda_{00} = 0$, and $\lambda_{11} = 0$

Cost of misclassifying a person not having lung cancer to be a patient with lung cancer is 0.007
 $\lambda_{10} = 0.007$

Cost of misclassifying a person having lung cancer to be a healthy person is 1
 $\lambda_{01} = 1$

- We know that patient A has gene X , i.e., $Z = 1$
- To estimate $P(Y = 0|Z = 1)$ v.s. $P(Y = 1|Z = 1)$
- With the 0/1 loss

Predict $Y = 1$ if $P(Y = 1|Z = 1) > P(Y = 0|Z = 1)$

Predict $Y = 0$ otherwise

OR $P(Z = 1|Y = 1)P(Y = 1) > P(Z = 1|Y = 0)P(Y = 0)$

- Take action a^* if $a^* = \arg \min_{a_i} R(a_i|x)$, where
 $R(a_i|x) = \sum_{c=0}^{C-1} \lambda_{ic} P(Y = c|x)$

• Expected risk of predicting A not having lung cancer:

$$\begin{aligned} R(a_0|Z = 1) &= \lambda_{01} \times P(Y = 1|Z = 1) + \lambda_{00} \times P(Y = 0|Z = 1) \\ &= 1 \times \frac{P(Z = 1|Y = 1) \times P(Y = 1)}{P(Z = 1)} + 0 \\ &= 1 \times \frac{0.9 \times 0.01}{P(Z = 1)} = \frac{0.009}{P(Z = 1)} \end{aligned}$$

- We have
 - $P(Z = 1|Y = 1)P(Y = 1) = 0.9 \times 0.01 = 0.009$
 - $P(Z = 1|Y = 0)P(Y = 0) = 0.2 \times 0.99 = 0.198$
 - $P(Z = 1|Y = 0)P(Y = 0) > P(Z = 1|Y = 1)P(Y = 1)$
- $P(Y = 0|Z = 1) > P(Y = 1|Z = 1)$
- Prediction: patient A does not have lung cancer

• Expected risk of predicting A having lung cancer:

$$\begin{aligned} R(a_1|Z = 1) &= \lambda_{11} \times P(Y = 1|Z = 1) + \lambda_{10} \times P(Y = 0|Z = 1) \\ &= 0 + 0.007 \times \frac{P(Z = 1|Y = 0) \times P(Y = 0)}{P(Z = 1)} \\ &= 0.007 \times \frac{0.198}{P(Z = 1)} \\ &= \frac{0.0014}{P(Z = 1)} < R(a_0|Z = 1) = \frac{0.009}{P(Z = 1)} \end{aligned}$$

• Prediction: patient A has lung cancer

Question 1.2

Task: Predict class label for $(A = 1, B = 1, C = 1)$ using NB

$$\text{Let } P(A = 1, B = 1, C = 1) = K$$

$$P(+|A = 1, B = 1, C = 1)$$

$$\begin{aligned} &= \frac{P(A = 1, B = 1, C = 1|+) \times P(+)}{K} \\ &= \frac{P(A = 1|+) \times P(B = 1|+) \times P(C = 1|+) \times P(+)}{K} \\ &= \frac{0.5 \times 0.5 \times 0.4}{K} \\ &= \frac{0.1}{K} \end{aligned}$$

$$P(-|A = 1, B = 1, C = 1)$$

$$\begin{aligned} &= \frac{P(A = 1, B = 1, C = 1|-) \times P(-)}{K} \\ &= \frac{P(A = 1|-) \times P(B = 1|-) \times P(C = 1|-) \times P(-)}{K} \\ &= \frac{0.33 \times 0.33 \times 0.33 \times 0.6}{K} \\ &= \frac{0.0222}{K} < \frac{0.1}{K} \quad P(+|A = 1, B = 1, C = 1) \end{aligned}$$

Class label = “+”

SC4000/CZ4041/CE4041: Machine Learning Lecture 3 Tutorial Question Sets

- 6/10
+ 4/10

Question 1.1

Estimate conditional probabilities

$$P(x_i = k|y = c) = \frac{|(x_i = k) \wedge (y = c)|}{|y = c|}$$

	$P(A = 1 -)$	$P(A = 0 -)$
	2/6 = 0.33	4/6 = 0.67
	$P(B = 1 -)$	$P(B = 0 -)$
	2/6 = 0.33	4/6 = 0.67
	$P(C = 1 -)$	$P(C = 0 -)$
	2/6 = 0.33	4/6 = 0.67
	$P(A = 1 +)$	$P(A = 0 +)$
	2/4 = 0.5	2/4 = 0.5
	$P(B = 1 +)$	$P(B = 0 +)$
	2/4 = 0.5	2/4 = 0.5
	$P(C = 1 +)$	$P(C = 0 +)$
	4/4 = 1	0/4 = 0

Record	A	B	C	Class
1	0	0	0	-
2	0	0	1	-
3	0	1	1	+
4	0	1	1	+
5	0	0	1	-
6	1	0	1	+
7	1	1	0	-
8	1	0	0	-
9	1	0	1	+
10	0	1	0	-

Table 1: Data set for Question 1.

Record	A	B	C	Class
1	0	0	0	-
2	0	0	1	-
3	0	1	1	+
4	0	1	1	+
5	0	0	1	-
6	1	0	1	+
7	1	1	0	-
8	1	0	0	-
9	1	0	1	+
10	0	1	0	-

Table 1: Data set for Question 1.

$$m = 3$$

$p = 1/3$ for all discrete features of class Yes

$p = 2/3$ for all discrete features of class No

$$\begin{aligned} P(X|\text{Class}=\text{No}) &= P(\text{HomO}=\text{No} | \text{Class}=\text{No}) \\ &\times P(\text{Status}=\text{Married} | \text{Class}=\text{No}) \\ &\times P(\text{Income}=120\text{K} | \text{Class}=\text{No}) \\ &= 6/11 \times 6/13 \times 0.0072 = 0.0018 \end{aligned}$$

$$\begin{aligned} P(X|\text{Class}=\text{Yes}) &= P(\text{HomO}=\text{No} | \text{Class}=\text{Yes}) \\ &\times P(\text{Status}=\text{Married} | \text{Class}=\text{Yes}) \\ &\times P(\text{Income}=120\text{K} | \text{Class}=\text{Yes}) \\ &= 4/5 \times 1/6 \times 1.2 \times 10^{-9} = 1.6 \times 10^{-10} \end{aligned}$$

$$\begin{aligned} P(\text{X}|\text{No}) \times P(\text{No}) &> P(\text{X}|\text{Yes}) \times P(\text{Yes}) \\ \text{Thus, } P(\text{No}|\text{X}) &> P(\text{Yes}|\text{X}) \end{aligned}$$

Cheat = No

Tid	Home Owner	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

1. Estimate the conditional probabilities for $P(A = 1|+)$, $P(B = 1|+)$, $P(C = 1|+)$, $P(A = 1|-)$, $P(B = 1|-)$, and $P(C = 1|-)$.

2. Use the estimate of conditional probabilities given in the previous question to predict the class label for a test example $(A = 1, B = 1, C = 1)$ using the naïve Bayes approach.

Question 2: On Page 28 of the lecture notes “Lecture 3”, recalculate the likelihoods using m-estimate. Compare the m-estimate method with the original method shown on Page 25 for estimating probabilities. Which method is better and why?

Home Owner	Marital Status	Taxable Income	Cheat
No	Married	120K	?

$$\begin{aligned} P(\text{HomO}=\text{Yes}|\text{No}) &= (3+2)/(7+3) = 5/10 \\ P(\text{HomO}=\text{No}|\text{No}) &= (4+2)/(7+3) = 6/10 \end{aligned}$$

$$\begin{aligned} P(\text{HomO}=\text{Yes}|\text{Yes}) &= (0+1)/(3+3) = 1/6 \\ P(\text{HomO}=\text{No}|\text{Yes}) &= (3+1)/(3+3) = 4/6 \end{aligned}$$

$$P(\text{Status}=\text{Single}|\text{No}) = (2+2)/(7+3) = 4/10$$

$$P(\text{Status}=\text{Divorced}|\text{No}) = (1+2)/(7+3) = 3/10$$

$$P(\text{Status}=\text{Married}|\text{No}) = (4+2)/(7+3) = 6/10$$

$$P(\text{Status}=\text{Single}|\text{Yes}) = (2+1)/(3+3) = 3/6$$

$$P(\text{Status}=\text{Divorced}|\text{Yes}) = (1+1)/(3+3) = 2/6$$

$$P(\text{Status}=\text{Married}|\text{Yes}) = (0+1)/(3+3) = 1/6$$

For taxable income:

If class=No: sample mean=110

sample variance=2975

If class=Yes: sample mean=90

sample variance=25

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

$$p = 2/3$$

$$P(\text{Class} = \text{No}) = 7/10$$

$$P(\text{Class} = \text{Yes}) = 3/10$$

$$m = 3$$

$$p = 1/3$$

SC4000/CZ4041/CE4041: Machine Learning Lecture 4 Tutorial Question Sets

$$P(D=1 | E=1, A=1, B=1)$$

Question 1: Given the Bayesian Belief Network shown in Figure 1 (the same as on the lecture notes), if the person has high blood pressure, but exercises regularly and eats a healthy diet, to estimate the probability that the person has heart disease.

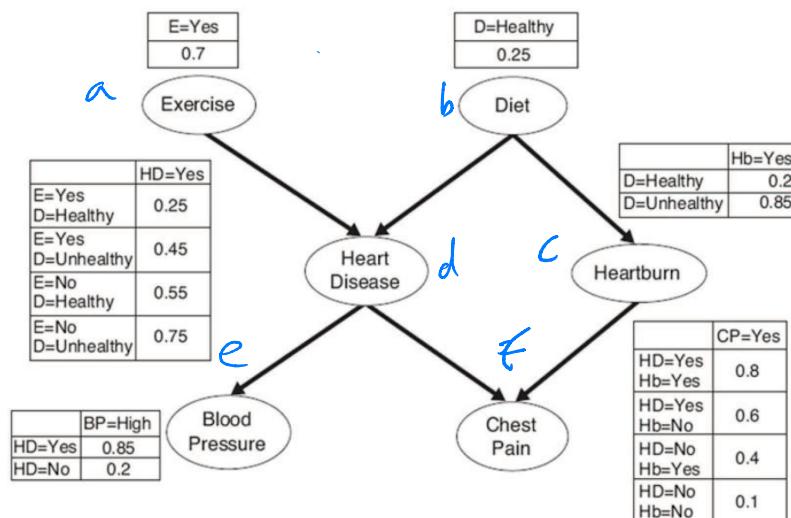


Figure 1: The Bayesian Network for Question 1.

Question 2: Consider the Bayesian Belief Network shown in Figure 2. Given that the outcomes of medical test A and medical test B for a specific patient are positive and negative, respectively. That is, $A = \text{Pos}$ and $B = \text{Neg}$. Predict the probability that the patient has high cholesterol.

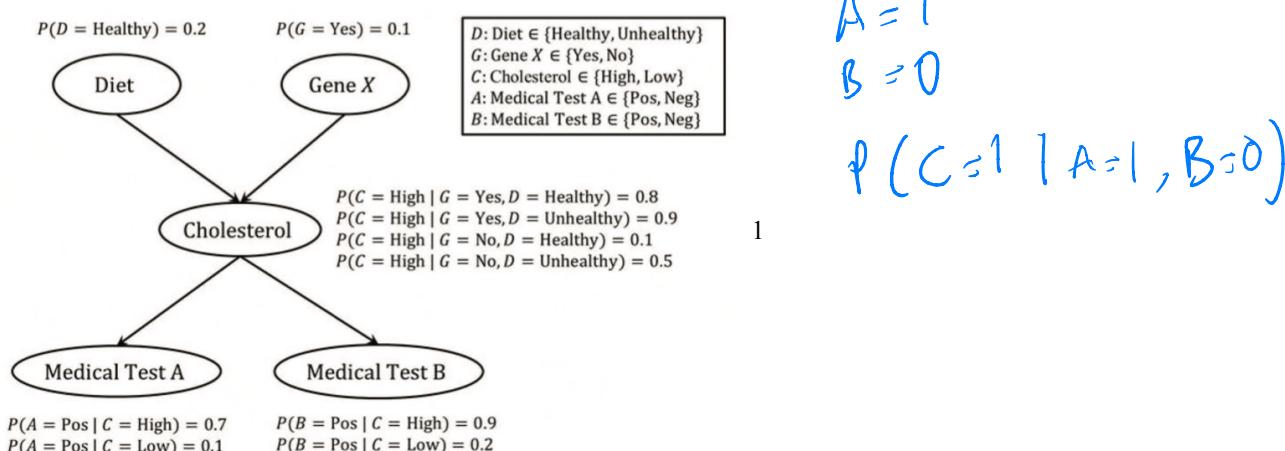


Figure 2: The Bayesian Network for Question 2.

SC4000/CZ4041/CE4041: Machine Learning

Lecture 5 Tutorial Question Sets

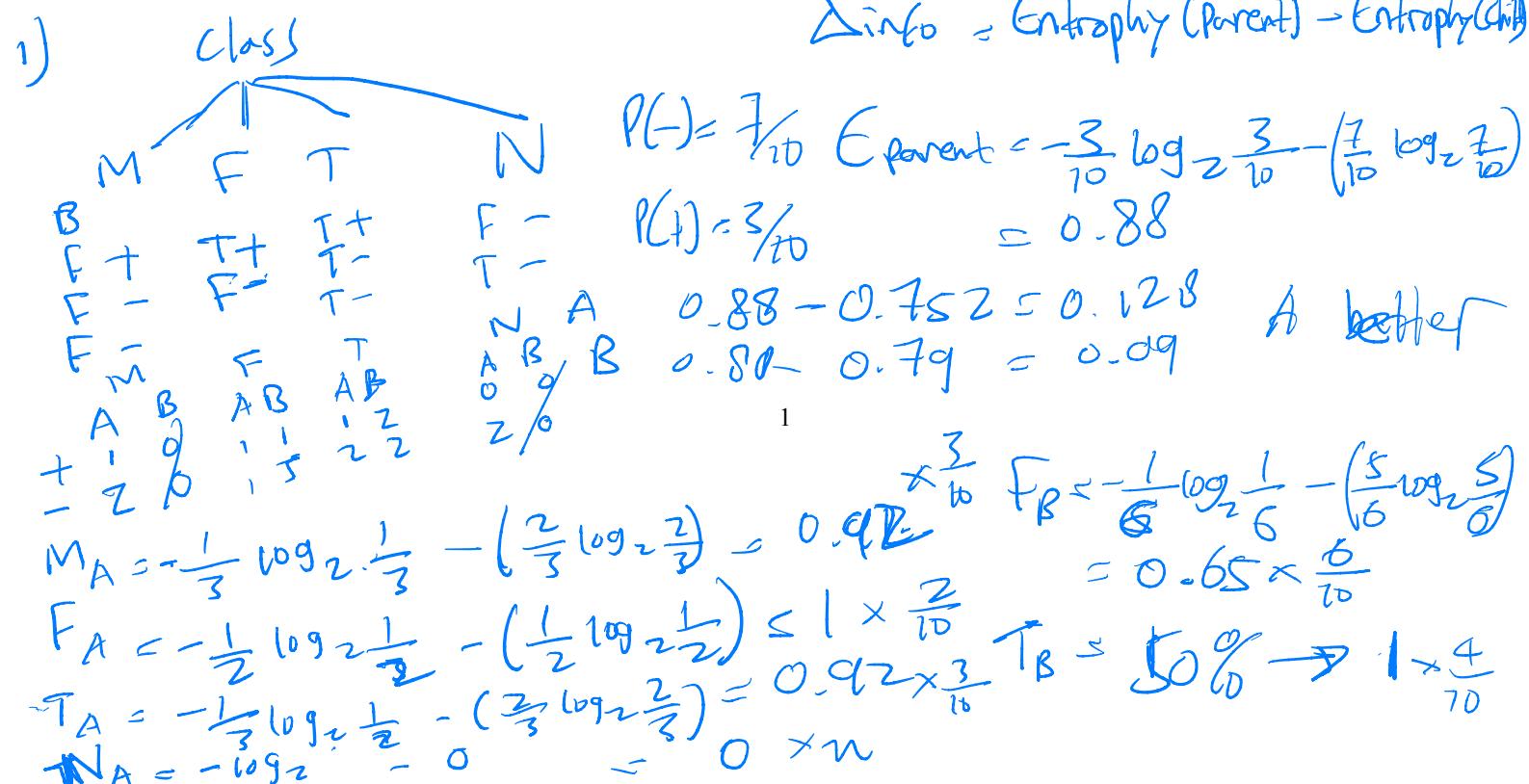
Question 1: Consider the data set shown in Table 1 for a binary classification problem.

$$3) \Delta_{\text{info}} R = \frac{\Delta_{\text{info}}}{\text{split INFO}}$$

Table 1: Data set for Question 1.

A	B	Class Label
M	F	+
F	T	+
T	T	+
M	F	-
M	F	-
F	F	-
N	F	-
N	T	-
T	T	-
T	F	-

1. Calculate the information gain when splitting on A and B (using multi-way split on A). Which feature would the decision tree induction algorithm choose?
 2. Calculate the gain ratio when splitting on A and B (using multi-way split on A). Which feature would the decision tree induction algorithm choose?



Entropy & Information Gain

$$\text{Entropy}(t) = -\sum_c P(y=c; t) \log_2 P(y=c; t)$$

- Suppose a parent node t is split into P partitions (children)

- Information Gain:

$$\Delta_{\text{info}} = \text{Entropy}(t) - \sum_{j=1}^P \frac{n_j}{n} \text{Entropy}(j)$$

Number of examples at child j

Number of examples at node t

- To choose a feature whose test condition maximizes the gain

Question 1.1

Table 1: Data set for Question 1.

A	B	Class Label
M	F	+
F	T	+
T	T	+
M	F	-
M	F	-
F	F	-
N	F	-
N	T	-
T	T	-
T	F	-

Task: calculate information gain when splitting on A (multi-way) and B . Which feature to choose?

Parent
+
-

$A = T$	$A = F$	$A = M$	$A = N$
+	1	1	1
-	2	1	2

Split on A

$B = T$	$B = F$
+	2
-	2

Split on B

$$\text{Entropy}(\text{Parent}) = -\left(\frac{3}{10}\right) \log_2 \left(\frac{3}{10}\right) - \left(\frac{7}{10}\right) \log_2 \left(\frac{7}{10}\right) = 0.8813$$

Parent
+
-

$$\text{Entropy}(A = T) = -\left(\frac{1}{3}\right) \log_2 \left(\frac{1}{3}\right) - \left(\frac{2}{3}\right) \log_2 \left(\frac{2}{3}\right) = 0.9183$$

$$\text{Entropy}(A = F) = -\left(\frac{1}{2}\right) \log_2 \left(\frac{1}{2}\right) - \left(\frac{1}{2}\right) \log_2 \left(\frac{1}{2}\right) = 1$$

$A = T$	$A = F$	$A = M$	$A = N$
+	1	1	1
-	2	1	2

$$\text{Entropy}(A = M) = -\left(\frac{1}{3}\right) \log_2 \left(\frac{1}{3}\right) - \left(\frac{2}{3}\right) \log_2 \left(\frac{2}{3}\right) = 0.9183 \quad \text{Split on } A$$

$$\text{Entropy}(A = N) = -\left(\frac{0}{2}\right) \log_2 \left(\frac{0}{2}\right) - \left(\frac{2}{2}\right) \log_2 \left(\frac{2}{2}\right) = 0$$

$$\text{Entropy}(\text{Split}_A) = \left(\frac{3}{10}\right) \times 0.9183 + \left(\frac{2}{10}\right) \times 1 + \left(\frac{3}{10}\right) \times 0.9183 + \left(\frac{2}{10}\right) \times 0 = 0.7510$$

$$\Delta_{\text{info}}(A) = 0.8813 - 0.7510 = 0.1303$$

$$\text{Entropy}(\text{Parent}) = -\left(\frac{3}{10}\right) \log_2 \left(\frac{3}{10}\right) - \left(\frac{7}{10}\right) \log_2 \left(\frac{7}{10}\right) = 0.8813$$

Parent
+
-

$$\text{Entropy}(B = T) = -\left(\frac{2}{4}\right) \log_2 \left(\frac{2}{4}\right) - \left(\frac{2}{4}\right) \log_2 \left(\frac{2}{4}\right) = 1$$

$$\text{Entropy}(B = F) = -\left(\frac{1}{6}\right) \log_2 \left(\frac{1}{6}\right) - \left(\frac{5}{6}\right) \log_2 \left(\frac{5}{6}\right) = 0.65$$

Split on B

Question 1.2: Calculate Gain Ratio

- Suppose a parent node t is split into P partitions (children)

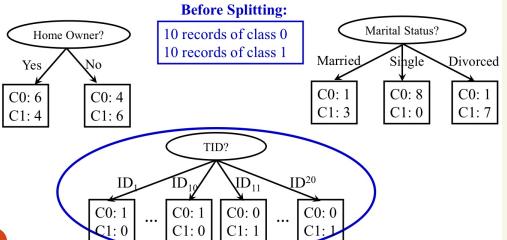
$$\text{Gain Ratio} = \frac{\Delta_{\text{info}}}{\text{SplitINFO}}$$

$$\text{where SplitINFO} = -\sum_{j=1}^P \frac{n_j}{n} \log_2 \frac{n_j}{n}$$

The number of records in partition j

Why Gain Ratio?

- Disadvantage: tends to prefer splits that result in large number of partitions, each being small but pure



$$\text{Entropy}(\text{Split}_B) = \left(\frac{4}{10}\right) \times 1 + \left(\frac{6}{10}\right) \times 0.65 = 0.79$$

$$\Delta_{\text{info}}(B) = 0.8813 - 0.79 = 0.0913 \quad < \quad \Delta_{\text{info}}(A) = 0.1303 \quad \checkmark$$

Question 1.2

$$\Delta_{\text{info}}(A) = 0.1303$$

$$\Delta_{\text{info}}(B) = 0.0913$$

$$\text{Gain Ratio} = \frac{\Delta_{\text{info}}}{\text{SplitINFO}} \quad \text{where SplitINFO} = -\sum_{i=1}^P \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

$A = T$	$A = F$	$A = M$	$A = N$
3	2	3	2

SplitINFO(A)

$$= -\left(\frac{3}{10}\right) \log_2 \left(\frac{3}{10}\right) - \left(\frac{2}{10}\right) \log_2 \left(\frac{2}{10}\right) - \left(\frac{3}{10}\right) \log_2 \left(\frac{3}{10}\right) - \left(\frac{2}{10}\right) \log_2 \left(\frac{2}{10}\right) = 1.9710$$

$$\text{GainRatio}_A = \frac{\Delta_{\text{info}}(A)}{\text{SplitINFO}(A)} = \frac{0.1303}{1.9710} = 0.0661$$

$B = T$	$B = F$
4	6
-	-

$$\text{SplitINFO}(B) = -\left(\frac{4}{10}\right) \log_2 \left(\frac{4}{10}\right) - \left(\frac{6}{10}\right) \log_2 \left(\frac{6}{10}\right) = 0.9710$$

$$\text{GainRatio}_B = \frac{\Delta_{\text{info}}(B)}{\text{SplitINFO}(B)} = \frac{0.0913}{0.9710} = 0.094 \quad \checkmark$$

$$Q) e'(t) \leq e(t) + N \times k = (3+2) + 2 \times 0.5 = 6$$

Case 1
before pruning: $e'(t) = 7 + 1 \times 0.5 = 7.5$, not prune

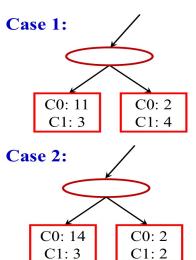
Case 2
before pruning: $e'(t) = (3+2) + 2 \times 0.5 = 6$,
 $e'(t) = 5 + 1 \times 0.5 = 5.5$, prune

Question 1

- Pessimistic error?

$$e'(T) = e(T) + N \times 0.5$$

PRUNE?



00/CZ4041/CE4041: Machine Learning Lecture 6 Tutorial Question Sets

Given the two cases shown in Figure 1, should we perform post-pruning using pessimistic error for each of them?

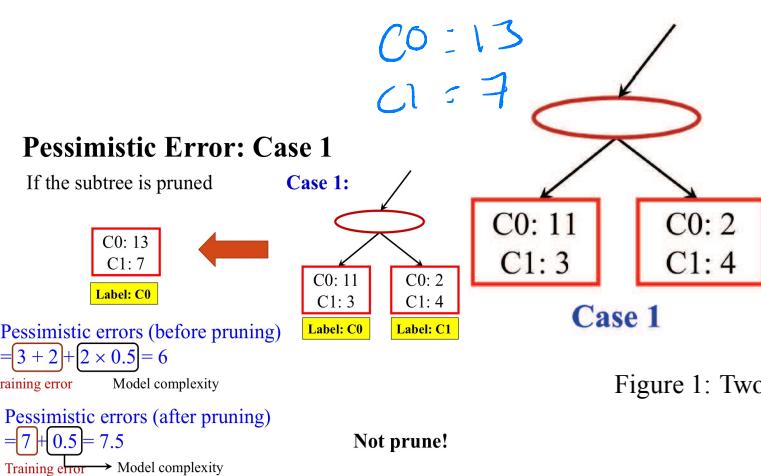


Figure 1: Two cases of decision trees.

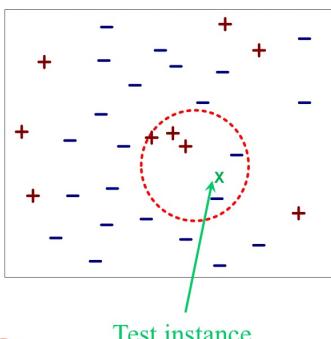
Table 1: The 5 nearest (training) neighbors to a test instance.

Training instance index	Class label	Distance to the test instance
1	+	3
2	+	3.5
3	+	4
4	-	1.5
5	-	2

Retrieved instances →

Question 2

Consider a binary classification problem, and a 5-NN classifier



Training instances	Class label	Distance to test instance
1	+	3
2	+	3.5
3	+	4
4	-	1.5
5	-	2

- Majority voting:
 $+ : 3 > - : 2$
- Distance-weighted voting:



Question 2 (cont.)

Distance-Weighted voting for +:

$$\left(\frac{1}{3}\right)^2 + \left(\frac{1}{3.5}\right)^2 + \left(\frac{1}{4}\right)^2 = 0.2552$$



Distance-Weighted voting for -:

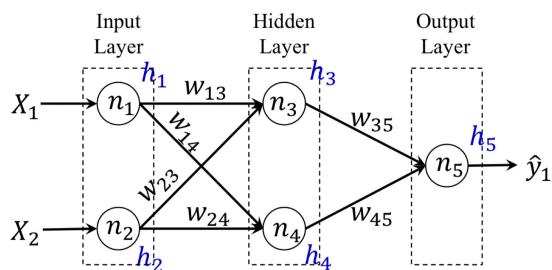
$$\left(\frac{1}{2}\right)^2 + \left(\frac{1}{1.5}\right)^2 = 0.6944$$

Training record	Class label	Distance to test record
1	+	3
2	+	3.5
3	+	4
4	-	1.5
5	-	2



$$\lambda = 0.4$$

$$\theta = 0$$



x_1	x_2	y_1
0	0	-1
0	1	1
1	0	1
1	1	1

$$h_1 = h_2 = 0$$

$$h_3 = \text{Sign}(0) = 1$$

$$h_4 = 1$$

$$h_5 = \text{Sign}(-1+1) \\ \Rightarrow -1$$

SC4000 Machine Learning Tutorial Artificial Neural Networks

Question 1: On Slide 64 of Lecture 7, we have shown how to use backpropagation to update the parameters of ANN with one initialization setting for w . Suppose now we initialize w with another set of values: $w_{13} = -1$, $w_{14} = -1$, $w_{23} = -1$, $w_{24} = -1$, $w_{35} = -1$, and $w_{45} = -1$. Run one epoch (i.e., run through the whole training dataset once), to show how the parameters are updated at each iteration.

Question 2: Consider a 2-dimensional dataset for three-class classification by ANN, as shown in Figure 1. Which ANN model as shown in Figure 2 is proper to solve the classification problem? Why?

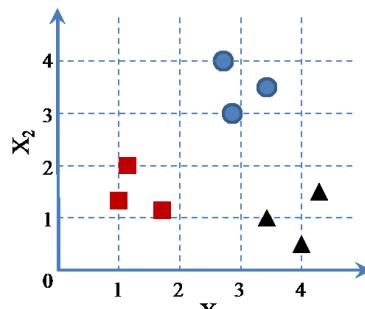


Figure 1: Dataset for Question 2.

X²
2X

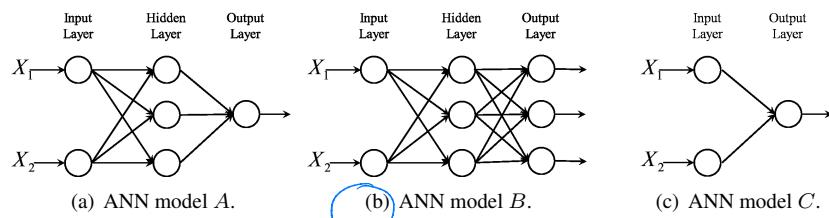


Figure 2: Different ANN structures for Question 2.

Question 3: Compute the derivative of the sigmoid function $f(z) = \frac{1}{1+e^{-z}}$ w.r.t. z .

$$\frac{1}{1+e^{-z}}$$

$$u = 1 + e^{-z}$$

$$du = -e^{-z}$$

$$\sim \frac{1}{u^2} (-e^{-z})$$

$$\frac{dz}{du} = \frac{e^{-z}}{(1+e^{-z})^2}$$

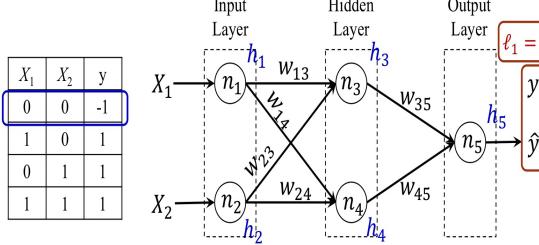
$$\frac{1}{1+e^{-z}} \left(\frac{e^{-z}}{1+e^{-z}} \right)$$

$$f(z) = \frac{1}{1+e^{-z}}$$

$$(f(z))(1-f(z))$$

Question 1

- In Lecture 7, we showed how to use backpropagation to update the parameters of ANN with one initialization setting for w .
- Suppose now we initialize w with another set of values: $w_{13} = -1$, $w_{14} = -1$, $w_{23} = -1$, $w_{24} = -1$, $w_{35} = -1$, and $w_{45} = -1$.
- Run one epoch (i.e., run through the whole training dataset once), to show how the parameters are updated at each iteration.

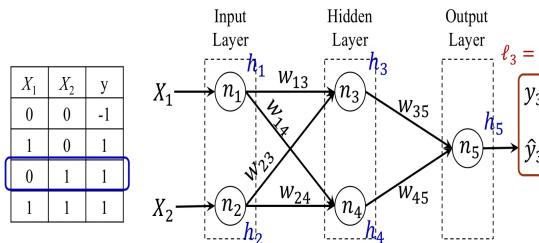


Forward pass:

For the 1st example: $h_1 = 0$ and $h_2 = 0$

$h_3 = \text{sign}(0 \times (-1) + 0 \times (-1)) = 1$ and $h_4 = \text{sign}(0 \times (-1) + 0 \times (-1)) = 1$

Then $\hat{y}_1 = h_5 = \text{sign}(1 \times (-1) + 1 \times (-1)) = -1$



For the 3rd example: $h_1 = 0$ and $h_2 = 1$

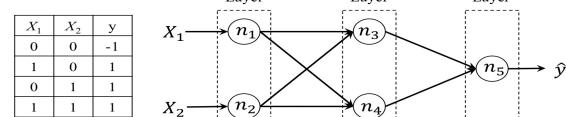
$h_3 = \text{sign}(0 \times (-1) + 1 \times (-1)) = -1$ and $h_4 = \text{sign}(0 \times (-1) + 1 \times (-1)) = -1$

Then $\hat{y}_3 = h_5 = \text{sign}((-1) \times (-1) + (-1) \times (-1)) = 1$

Question 1

$$\lambda = 0.4, \theta = 0$$

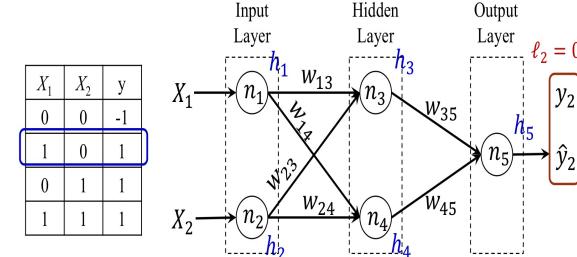
Sign function as $a(\cdot)$



• Initialization 1 (lecture notes): $(w_{13} = 1, w_{14} = 1, w_{23} = 1, w_{24} = 1, w_{35} = 1, w_{45} = 1)$

• Initialization 2:

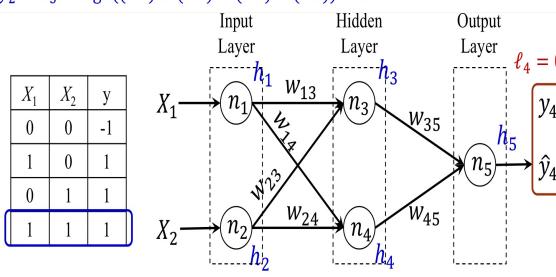
$$(w_{13} = -1, w_{14} = -1, w_{23} = -1, w_{24} = -1, w_{35} = -1, w_{45} = -1)$$



For the 2nd example: $h_1 = 1$ and $h_2 = 0$

$h_3 = \text{sign}(1 \times (-1) + 0 \times (-1)) = -1$ and $h_4 = \text{sign}(1 \times (-1) + 0 \times (-1)) = -1$

Then $\hat{y}_2 = h_5 = \text{sign}((-1) \times (-1) + (-1) \times (-1)) = 1$



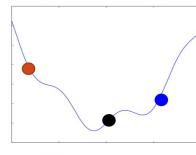
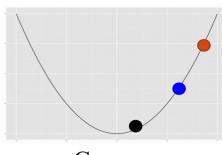
For the 4th example: $h_1 = 1$ and $h_2 = 1$

$h_3 = \text{sign}(1 \times (-1) + 1 \times (-1)) = -1$ and $h_4 = \text{sign}(1 \times (-1) + 1 \times (-1)) = -1$

Then $\hat{y}_4 = h_5 = \text{sign}((-1) \times (-1) + (-1) \times (-1)) = 1$

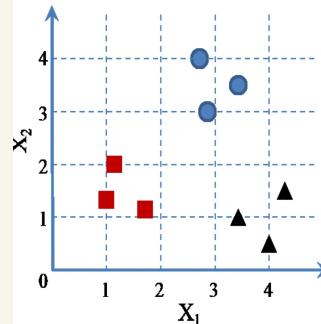
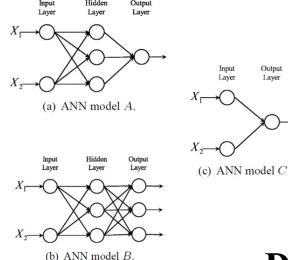
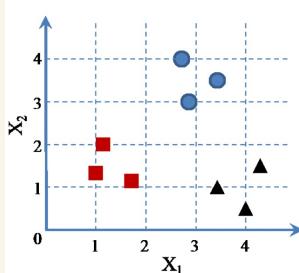
Question 1: Conclusion

- Different initializations may lead to different convergence iterations
- Different initializations may lead to different local optima



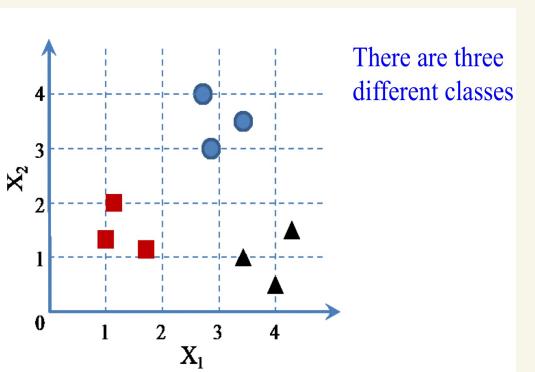
Question 2

- Consider a 2-dimensional dataset for three-class classification by ANN, as shown in Figure 1. Which ANN model as shown in Figure 2 is proper to solve the classification problem? Why?



The instances of different classes are not able to be linearly separated

Therefore, the Perceptron model does not work



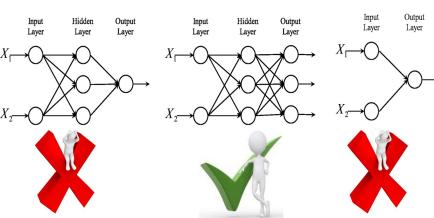
Design Issues for ANN

- The number of nodes in the input layer
 - Assign an input node to each numerical or binary input variable
- The number of nodes in the output layer
 - Binary class problem → single node
 - C -class problem → C output nodes
- We output a one-hot encoding of the class
 - [1, 0, 0] for class 1
 - [0, 1, 0] for class 2

There are three different classes

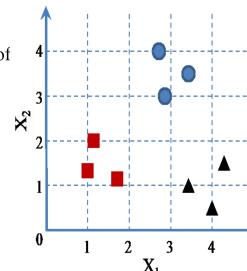
Question 2 (cont.)

Therefore, there should be three output nodes in the output layer



[Optional] Multi-class activation

- In a multi-class classification problem with C classes, we should output a one-hot encoding of the class
 - [1, 0, 0] for class 1
 - [0, 1, 0] for class 2
 - We take the pre-activation outputs, $z_{L,1}, z_{L,2}$, and $z_{L,3}$, and feed them to a softmax function.
- $$\hat{y}_i = \frac{\exp(z_{L,i})}{\sum_j \exp(z_{L,j})}$$
- This results in $1 > \hat{y}_i > 0$ and $\sum_i \hat{y}_i = 1$



Question 3

- Compute the derivative of the sigmoid function w.r.t. z :

$$f(z) = \frac{1}{1 + e^{-z}}$$

- Denote $y = 1 + e^{-z}$, i.e., y is a function of z
- And then f is a function of y : $f = \frac{1}{y}$

- Based on the chain rule:

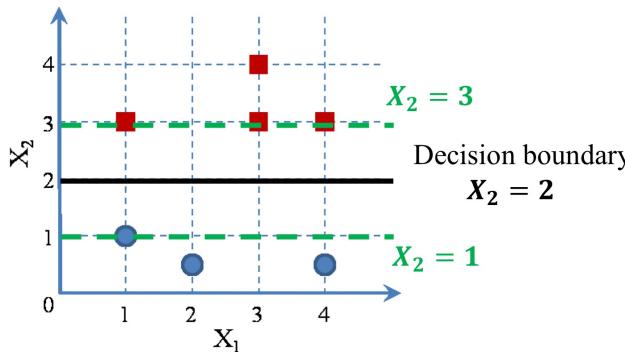
$$\frac{\partial f(z)}{\partial z} = \frac{\partial f}{\partial y} \frac{\partial y}{\partial z} = \left(-\frac{1}{y^2} \right) \left(\frac{\partial 1}{\partial z} + \boxed{\frac{\partial e^{-z}}{\partial z}} \right) = \frac{e^{-z}}{y^2}$$

$$\frac{\partial f(z)}{\partial z} = \frac{\partial f}{\partial y} \frac{\partial y}{\partial z} = \frac{e^{-z}}{y^2}$$

$$y = 1 + e^{-z}$$

$$\begin{aligned} \frac{\partial f(z)}{\partial z} &= \frac{e^{-z}}{(1 + e^{-z})^2} = \frac{1}{1 + e^{-z}} \left(\frac{e^{-z}}{1 + e^{-z}} \right) \\ f(z) &= \frac{1}{1 + e^{-z}} = f(z) \left(1 - \frac{1}{1 + e^{-z}} \right) \\ &= f(z)(1 - f(z)) \end{aligned}$$

Question 1



$$\frac{w^2 x_i^2}{2} - \frac{\|w\|^2}{2} + \frac{\lambda}{2} w^2$$

SC4000 Machine Learning Tutorial

Support Vector Machines and Linear Regression

$$\begin{aligned} & \sum_i^n x_i^T w - x_i y_i + \lambda w \\ & \sum_{i=1}^n x_i (w x_i - y_i) + \lambda w \\ & \sum_{i=1}^n (w x_i) x_i - \sum_{i=1}^n y_i x_i + \lambda w = 0 \end{aligned}$$

Question 1: Consider a 2-dimensional dataset for two-class classification by SVM, as shown in Figure 1, where the red “square” and blue “circle” denote the positive and negative classes respectively. Is this dataset separable by a linear SVM classifier? If no, why? If yes, what is the decision boundary of the linear SVM? And what are the pair of parallel hyperplanes associated with the decision boundary? (No need to provide proofs)

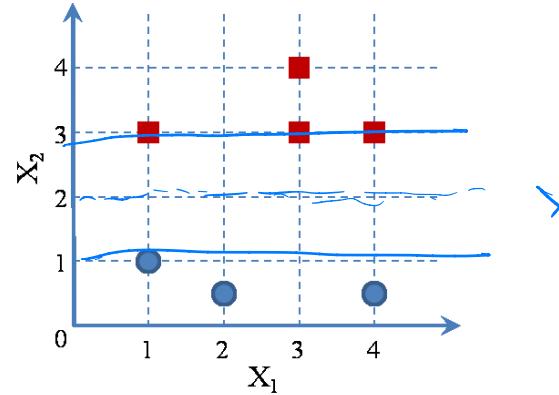


Figure 1: Dataset for Question 1.

Question 2: Please induce why the two parallel hyperplanes of a decision boundary,

$$w \cdot x + b = 0,$$

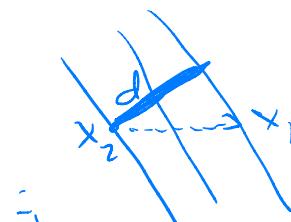
can be written as

$$w \cdot x + b = \bar{k}, \text{ and } w \cdot x + b = -\bar{k},$$

where $\bar{k} > 0$, respectively.

Question 3: In an SVM, we have two support vectors $(2, 3)$ from class 1 and $(-1, 4)$ from class 2. The separating hyperplane can be written as $3x_1 + x_2 + b = 0$. Compute the margin of separation.

$$w \cdot (x_1 - x_2) \quad \|w\| \|x_1 - x_2\| \cos \theta = d$$



$$d = \frac{w \cdot (x_1 - x_2)}{\|w\|} \quad \frac{(3, 1) \cdot (-1, 3)}{\sqrt{10}} = \frac{9 - 1}{\sqrt{10}} = \frac{8}{\sqrt{10}}$$

Question 2

The two parallel hyperplanes passing the closest circle(s) and square(s) can be written as

$$\mathbf{w} \cdot \mathbf{x} + b = k, \text{ where } k > 0$$

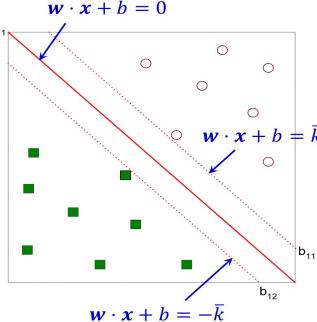
$$\mathbf{w} \cdot \mathbf{x} + b = k', \text{ where } k' < 0$$

It can be shown that, these two parallel hyperplanes can be further rewritten as

$$\mathbf{w} \cdot \mathbf{x} + b = \bar{k}$$

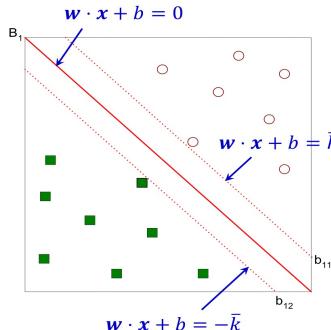
$$\mathbf{w} \cdot \mathbf{x} + b = -\bar{k}$$

$$\text{where } \bar{k} > 0$$



Question 2

- \mathbf{w} determines the orientation (slope) of the decision boundary.
- The support vectors determine how the decision boundary moves in parallel motion.
- Together, they determine b .

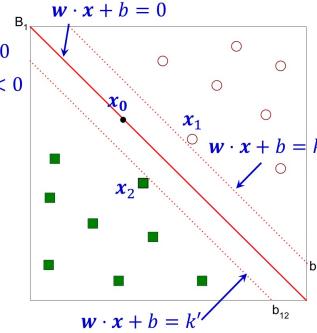


Question 2 (cont.)

$$b_{11}: \mathbf{w} \cdot \mathbf{x}_1 + b = k, \text{ where } k > 0$$

$$b_{12}: \mathbf{w} \cdot \mathbf{x}_2 + b = k', \text{ where } k' < 0$$

Given two support vectors (or two points on b_{11} and b_{12} respectively), I can choose b such that $k = -k'$

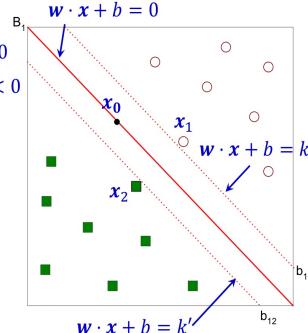


Question 2 (cont.)

$$b_{11}: \mathbf{w} \cdot \mathbf{x}_1 + b = k, \text{ where } k > 0$$

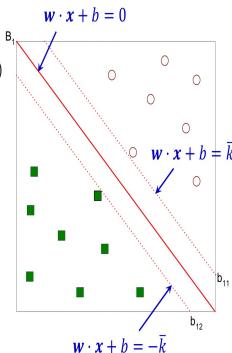
$$b_{12}: \mathbf{w} \cdot \mathbf{x}_2 + b = k', \text{ where } k' < 0$$

$$\begin{aligned} \mathbf{w} \cdot \mathbf{x}_1 + b &= -(\mathbf{w} \cdot \mathbf{x}_2 + b) \\ 2b &= -\mathbf{w} \cdot \mathbf{x}_1 - \mathbf{w} \cdot \mathbf{x}_2 \\ b &= -\frac{1}{2}\mathbf{w} \cdot (\mathbf{x}_1 + \mathbf{x}_2) \end{aligned}$$



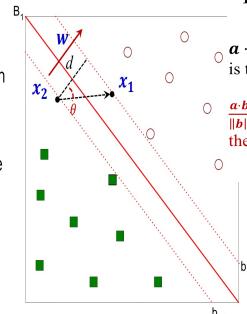
Question 3

- Two support vectors: $(2, 3)$ and $(-1, 4)$ from the two classes, respectively.
- Decision boundary:
 $\mathbf{w} \cdot \mathbf{x} + b = 3x_1 + x_2 + b = 0$
- What is the margin?



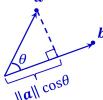
Question 3

- From the lecture, \mathbf{w} is orthogonal to the decision boundary.
- All we need is to find the length of projection of the vector $(x_1 - x_2)$ onto the direction of \mathbf{w} .



Review: Geometry of Inner Products

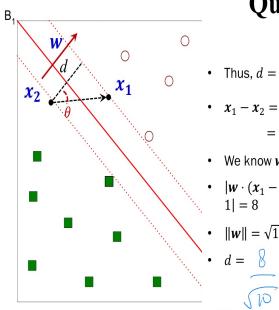
$$\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos\theta, \text{ where } \theta \text{ is the angle between } \mathbf{a} \text{ and } \mathbf{b}$$



$$\frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\|} = \|\mathbf{a}\| \cos\theta \text{ is the length of the projection of } \mathbf{a} \text{ onto } \mathbf{b}$$

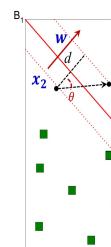
Question 3

- We know that
 $\mathbf{w} \cdot (x_1 - x_2) = \|\mathbf{w}\| \|x_1 - x_2\| \cos(\theta) = \|\mathbf{w}\| d \text{ or } -\|\mathbf{w}\| d$
- Thus, $d = \frac{|\mathbf{w} \cdot (x_1 - x_2)|}{\|\mathbf{w}\|}$
- $x_1 - x_2 = (2, 3) - (-1, 4) = (3, -1)$



Question 3

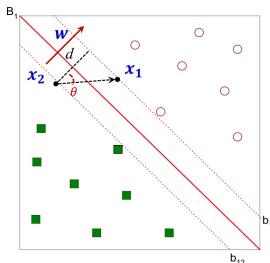
- Thus, $d = \left| \frac{\mathbf{w} \cdot (x_1 - x_2)}{\|\mathbf{w}\|} \right|$
- $x_1 - x_2 = (2, 3) - (-1, 4) = (3, -1)$
- We know $\mathbf{w} = (3, 1)$
- $\|\mathbf{w} \cdot (x_1 - x_2)\| = |3 \cdot 3 + 1 \cdot -1| = 8$
- $\|\mathbf{w}\| = \sqrt{10}$
- $d = \frac{8}{\sqrt{10}} = \frac{4\sqrt{10}}{5}$



Question 3

- Why can't we directly use the equation $\|\mathbf{w}\| d = 2$?
- That equation is only valid when \mathbf{w} is properly rescaled.
- That is, when
 $\mathbf{w} \cdot x_1 + b = 1$
 $\mathbf{w} \cdot x_2 + b = -1$

Can only use if $\frac{1}{\|\mathbf{w}\|} = 1$ and $\frac{-1}{\|\mathbf{w}\|} = -1$



Q4: Regularized Linear Regression

- To solve the unconstrained minimization problem, we can set the derivative of $\mathcal{L}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$ w.r.t. \mathbf{w} to zero

$$\frac{\partial (\mathcal{L}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2)}{\partial \mathbf{w}} = \frac{\partial \left(\frac{1}{2} \sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \mathbf{w} \cdot \mathbf{w} \right)}{\partial \mathbf{w}} = \mathbf{0}$$

- We can obtain a closed-form solution

Some Concepts: Review (cont.)

$$\begin{array}{c} X \in \mathbb{R}^{d \times d} \quad w \in \mathbb{R}^{d \times 1} \quad z = \mathbf{X}w \in \mathbb{R}^{d \times 1} \\ \begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{matrix} \times \begin{matrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \end{matrix} = \begin{matrix} z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \end{matrix} \quad x_1^T w = \sum_{j=1}^d x_{1j} w_j = z_1 \\ \\ X \in \mathbb{R}^{d \times k} \quad w \in \mathbb{R}^{k \times 1} \quad z = \mathbf{X}w \in \mathbb{R}^{d \times 1} \\ \begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{matrix} \times \begin{matrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \end{matrix} = \begin{matrix} z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \end{matrix} \end{array}$$

Question 4

- Denote by $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^T$

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_{N0} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{1d} & \cdots & \mathbf{x}_{Nd} \end{pmatrix}^T = \begin{pmatrix} \mathbf{x}_{10} & \cdots & \mathbf{x}_{1d} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{N0} & \cdots & \mathbf{x}_{Nd} \end{pmatrix}$$

- And by $\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$ How to get this closed-form solution?

- The closed-form solution for \mathbf{w} :

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\begin{array}{c} \mathbf{x}_i \qquad \qquad \qquad \mathbf{x}_i^T \\ \begin{bmatrix} x_{i0} \\ x_{i1} \\ \vdots \\ x_{id} \end{bmatrix} \qquad [x_{i0} \quad x_{i1} \quad \dots \quad x_{id}] \end{array}$$

$$\mathbf{x}_i \mathbf{x}_i^T = \begin{pmatrix} x_{i0} \times x_{i0} & \cdots & x_{i0} \times x_{id} \\ \vdots & \ddots & \vdots \\ x_{id} \times x_{i0} & \cdots & x_{id} \times x_{id} \end{pmatrix}$$

Some Concepts: Review

Transpose of a vector/matrix

If \mathbf{X} is a square matrix, then its numbers of rows and columns are the same

If \mathbf{X} is a symmetric matrix, then it is square and $\mathbf{X}^T = \mathbf{X}$

Some Concepts: Review (cont.)

- For a square matrix \mathbf{X} , if \mathbf{X} is invertible, then

$$\mathbf{X} \mathbf{X}^{-1} = \mathbf{I} \quad \text{Identity matrix}$$

$$\begin{array}{c} \mathbf{X} \qquad \qquad \qquad \mathbf{X}^{-1} \qquad \qquad \qquad \mathbf{I} \\ \times \qquad \qquad \qquad \times \qquad \qquad \qquad = \qquad \qquad \qquad \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \end{array}$$

$$\frac{\partial \left(\frac{1}{2} \sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \mathbf{w} \cdot \mathbf{w} \right)}{\partial \mathbf{w}} = \mathbf{0}$$

$$\begin{aligned} \frac{\partial \frac{1}{2} \sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}_i - y_i)^2}{\partial \mathbf{w}} &= \frac{1}{2} \sum_{i=1}^N 2(\mathbf{w} \cdot \mathbf{x}_i - y_i) \frac{\partial (\mathbf{w} \cdot \mathbf{x}_i - y_i)}{\partial \mathbf{w}} \\ &= \sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}_i - y_i) \mathbf{x}_i \end{aligned}$$

$$\sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}_i - y_i) \mathbf{x}_i + \lambda \mathbf{w} = \mathbf{0}$$

$$\sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}_i) \mathbf{x}_i - \sum_{i=1}^N y_i \mathbf{x}_i + \lambda \mathbf{w} = \mathbf{0}$$

$$\mathbf{x}_i \mathbf{x}_i^T = \begin{pmatrix} x_{i0} \times x_{i0} & \cdots & x_{i0} \times x_{id} \\ \vdots & \ddots & \vdots \\ x_{id} \times x_{i0} & \cdots & x_{id} \times x_{id} \end{pmatrix}$$

$$\sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) = \begin{pmatrix} \sum_{i=1}^N x_{i0} \times x_{i0} & \cdots & \sum_{i=1}^N x_{i0} \times x_{id} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^N x_{id} \times x_{i0} & \cdots & \sum_{i=1}^N x_{id} \times x_{id} \end{pmatrix}$$

$$\text{Therefore } \left(\sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) \right) \mathbf{w} - \sum_{i=1}^N y_i \mathbf{x}_i + \lambda \mathbf{I} \mathbf{w} = \mathbf{0}$$

$$(\mathbf{X}^T \mathbf{X}) \mathbf{w} - \mathbf{X}^T \mathbf{y} + \lambda \mathbf{I} \mathbf{w} = \mathbf{0}$$

Always invertible as long as λ is positive

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{I} \mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Some Concepts: Review

- The transpose of \mathbf{XY} : $(\mathbf{XY})^T = \mathbf{Y}^T \mathbf{X}^T$

$$(\mathbf{XY})^T = \mathbf{Y}^T \mathbf{X}^T$$

- The transpose of \mathbf{Xw} : $(\mathbf{Xw})^T = \mathbf{w}^T \mathbf{X}^T$

$$(\mathbf{Xw})^T = \mathbf{w}^T \mathbf{X}^T$$

- The transpose of $\mathbf{x}^T \mathbf{w}$: $(\mathbf{x}^T \mathbf{w})^T = \mathbf{w}^T \mathbf{x}$

$$(\mathbf{x}^T \mathbf{w})^T = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{x}^T = \mathbf{x}$$

Some Concepts: Review (cont.)

- Any vector (or matrix) \mathbf{x} (or \mathbf{X}) times identity matrix \mathbf{I} equals to the vector (or matrix) itself

$$\mathbf{Ix} = \mathbf{x} (\mathbf{x}^T \mathbf{I} = \mathbf{x}^T) \quad \text{OR} \quad \mathbf{XI} = \mathbf{X} (\mathbf{IX} = \mathbf{X})$$

$$\begin{array}{c} \mathbf{X} \qquad \qquad \qquad \mathbf{I} \qquad \qquad \qquad \mathbf{X} \\ \times \qquad \qquad \qquad \times \qquad \qquad \qquad = \qquad \qquad \qquad \mathbf{X} \\ \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \qquad \qquad \qquad \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \qquad \qquad \qquad \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \\ \mathbf{x} \qquad \qquad \qquad \mathbf{I} \qquad \qquad \qquad = \qquad \qquad \qquad \mathbf{x} \\ \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \qquad \qquad \qquad \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \qquad \qquad \qquad \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \end{array}$$

$$\sum_{i=1}^N (\mathbf{w} \cdot \mathbf{x}_i) \mathbf{x}_i - \sum_{i=1}^N y_i \mathbf{x}_i + \lambda \mathbf{w} = \mathbf{0}$$

$$\begin{aligned} \mathbf{x}_i (\mathbf{w} \cdot \mathbf{x}_i) &= \mathbf{x}_i (\mathbf{x}_i^T \mathbf{w}) \\ &= (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{w} \end{aligned}$$

Identity matrix

$$\left(\sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) \right) \mathbf{w} - \sum_{i=1}^N y_i \mathbf{x}_i + \lambda \mathbf{w} = \mathbf{0}$$

A matrix of $(d+1)$ by $(d+1)$ size, where \mathbf{x}_i is a column vector of $(d+1)$ dimensions

$$\sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) = \begin{pmatrix} \sum_{i=1}^N x_{i0} \times x_{i0} & \cdots & \sum_{i=1}^N x_{i0} \times x_{id} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^N x_{id} \times x_{i0} & \cdots & \sum_{i=1}^N x_{id} \times x_{id} \end{pmatrix}$$

$$\begin{array}{c} \mathbf{x}_1 \qquad \qquad \qquad \mathbf{x}_1^T \\ \begin{bmatrix} x_{10} \\ x_{11} \\ \vdots \\ x_{1d} \end{bmatrix} \qquad [x_{10} \quad x_{11} \quad \dots \quad x_{1d}] \end{array}$$

$$\begin{aligned} \sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) &= \mathbf{X}^T \mathbf{X} \\ \text{Therefore } \left(\sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) \right) \mathbf{w} - \sum_{i=1}^N y_i \mathbf{x}_i + \lambda \mathbf{I} \mathbf{w} &= \mathbf{0} \\ (\mathbf{X}^T \mathbf{X}) \mathbf{w} - \mathbf{X}^T \mathbf{y} + \lambda \mathbf{I} \mathbf{w} &= \mathbf{0} \\ (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ \mathbf{I} \mathbf{w} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ \mathbf{w} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

Question 1 *Q1* *not tested*

- Suppose there are 3 base classifiers, each classifier has error rate, $\varepsilon = 0.65$ or accuracy $acc = 0.35$
- Consider to combine the 3 base classifiers to make a prediction on a test instance using a majority vote
- Therefore, probability that the ensemble classifier makes a wrong prediction is:

$$\sum_{i=2}^3 \binom{3}{i} \varepsilon^i (1 - \varepsilon)^{3-i} = 3 \times 0.65^2 \times 0.35 + 1 \times 0.65^3 \times 1 = 0.71825$$

SC4000 Machine Learning Tutorial Ensemble Learning

Question 1: Why is the condition “the base classifiers should do better than a classifier that performs random guessing” necessary for ensemble learning?

Question 2: Suppose we have trained 5 base binary classifiers: f_1, f_2, f_3, f_4 and f_5 . Their predictions on a validation dataset are shown in Table 1, where the last column denotes the ground-truth class labels. Which base classifiers would you choose to construct an ensemble learner?

The Binomial Distribution

- This is the general expression for the binomial distribution.
- We perform N experiments. In each experiment, the probability of observing a negative outcome is ε . What is the probability for observing exactly M negative outcomes?

$$P(N, M) = \binom{N}{M} \varepsilon^M (1 - \varepsilon)^{N-M}$$

Table 1: Data set for Question 2.

ID	f_1	f_2	f_3	f_4	f_5	Ground Truth
P1	+	+	-	-	+	+
P2	+	+	-	+	-	+
P3	-	-	+	+	-	+
P4	-	-	+	-	+	+
P5	-	-	+	+	-	-
P6	-	-	-	+	+	+
P7	+	+	+	+	-	+
P8	-	+	+	-	+	-
P9	+	+	-	+	+	+
P10	-	-	-	+	-	-

Applying That Formula ...

- Suppose there are N (odd) independent base classifiers, each of which has the same error rate ε
- Therefore, probability that the ensemble classifier makes a wrong prediction is:

$$P(N) = \sum_{i=1}^{N/2} \binom{N}{i} \varepsilon^i (1 - \varepsilon)^{N-i}$$

Question 1 (cont.)

- It can be proved that, with an odd number N
 - If $p > 0.5$ then $P_c(N)$ is monotonically increasing in N , and $P(N) \rightarrow 1$ as $N \rightarrow \infty$
 - If $p = 0.5$ then $P_c(N) = 0.5$ for all N
 - If $p < 0.5$ then $P_c(N)$ is monotonically decreasing in N , and $P(N) \rightarrow 0$ as $N \rightarrow \infty$
- Detailed proof can be found in the paper “Application of Majority Voting to Pattern Recognition: An Analysis of Its Behavior and Performance, 1997”

Another Application

- Alternatively, let p be the probability that a single classifier makes the correct decision. The probability that the ensemble of $2n + 1$ base classifier makes the correct decision is $P_c(2n + 1)$

$$P_c(2n + 1) = \sum_{m=n+1}^{2n+1} \binom{2n+1}{m} p^m (1 - p)^{2n+1-m}$$

Q2

Table 1: Data set for Question 2.

ID	f_1	f_2	f_3	f_4	f_5	Ground Truth
P1	✓	✓	✗	✗	✓	+
P2	✗	✓	✗	✓	✗	+
P3	✗	✗	✓	✓	✗	+
P4	✗	✗	✓	✗	✓	+
P5	✓	✓	✗	✓	✓	-
P6	✗	✗	✗	✓	✓	+
P7	✓	✓	✓	✓	✗	+
P8	✓	✗	✗	✓	✗	-
P9	✓	✓	✗	✓	✓	+
P10	✗	✗	✓	✗	✗	-

0.7 0.6 0.4 0.6 0.6

SC4000 Machine Learning Tutorial

Clustering

Question 1: Given the distance matrix shown in Table 1, use a dendrogram to show how to perform agglomerative hierarchical clustering algorithm with Single Link on the distance matrix.

Table 1: Distance matrix.

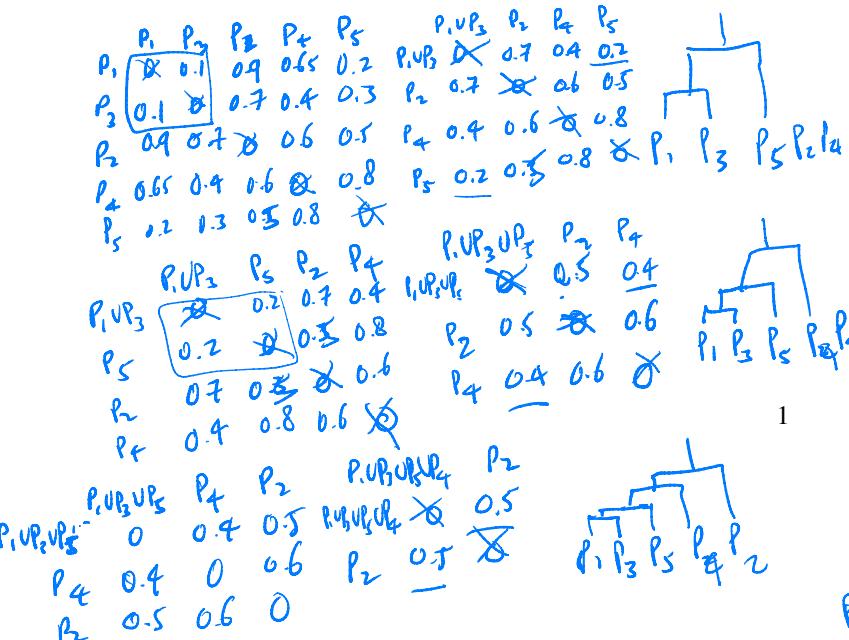
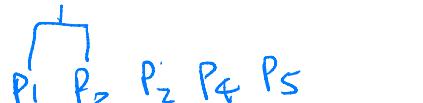
	P1	P2	P3	P4	P5
P1	0	0.9	0.1	0.65	0.2
P2	0.9	0	0.7	0.6	0.5
P3	0.1	0.7	0	0.4	0.3
P4	0.65	0.6	0.4	0	0.8
P5	0.2	0.5	0.3	0.8	0

Question 2: Refer to the clustering problem on Slide 59 of Lecture 10. Use a dendrogram to show how to perform hierarchical clustering with Complete Link on the similarity matrix.

largest similarity = take Max
but min it except 1

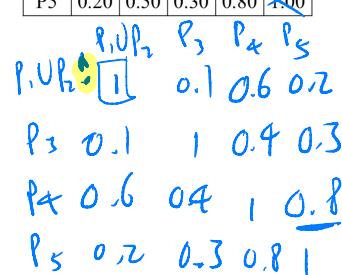
Q1)

	P1	P2	P3	P4	P5
P1	0	0.9	0.1	0.65	0.2
P2	0.9	0	0.7	0.6	0.5
P3	0.1	0.7	0	0.4	0.3
P4	0.65	0.6	0.4	0	0.8
P5	0.2	0.5	0.3	0.8	0



Q2)

	P1	P2	P3	P4	P5
P1	0.90	0.90	0.10	0.65	0.20
P2	0.90	0.70	0.70	0.60	0.50
P3	0.10	0.70	0.70	0.40	0.30
P4	0.65	0.60	0.40	0.40	0.80
P5	0.20	0.50	0.30	0.80	0.00

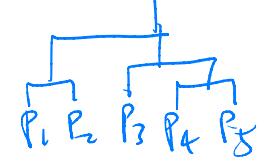


1

	P1	P2	P3	P4	P5
P1	0.90	0.90	0.10	0.65	0.20
P2	0.90	0.70	0.70	0.60	0.50
P3	0.10	0.70	0.70	0.40	0.30
P4	0.65	0.60	0.40	0.40	0.80
P5	0.20	0.50	0.30	0.80	0.00



	P1	P2	P3	P4	P5
P1	0.90	0.90	0.10	0.65	0.20
P2	0.90	0.70	0.70	0.60	0.50
P3	0.10	0.70	0.70	0.40	0.30
P4	0.65	0.60	0.40	0.40	0.80
P5	0.20	0.50	0.30	0.80	0.00



Question 1

- A dataset of five 4-dimensional instances is given in Table 1. Suppose an SVD is performed on the data matrix \mathbf{X} (5-by-4) via $\mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{U}^T$. The matrices \mathbf{V} , \mathbf{D} , and \mathbf{U} are shown in Tables 2-4, respectively. Use principal component analysis to project the 5 data points in Table 1 to 2-dimensional space.

	X_1	X_2	X_3	X_4
	2	4	1	3
	1	2	3	5
	-2	-4	-4	-1
	0	-1	-2	-6
	-1	-1	2	-1

Step 1: Center the data points
s.t. the mean is $\mathbf{0}$

$\hat{\mathbf{u}}$	X_1	X_2	X_3	X_4
	0	0	0	0

Already centered

PCA

- Step 1: Center the data points s.t. the mean is $\mathbf{0}$
- Step 2: Compute sample covariance matrix $\tilde{\Sigma}$
- Step 3: Compute eigenvectors of $\tilde{\Sigma}$, $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d\}$, which are sorted based on their eigenvalues in non-increasing order, i.e., $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$
- Step 4: Select the first k eigenvectors to construct principal components
- Step 5: Project the data instances to k-dimensional space

Dimensionality Reduction

- Question 1:** A dataset of five 4-dimensional instances is given in Table 1. Suppose an SVD is performed on the data matrix \mathbf{X} (5-by-4) via $\mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{U}^T$. The matrices \mathbf{V} , \mathbf{D} , and \mathbf{U} are shown in Tables 2-4, respectively. Use principal component analysis to project the 5 data points in Table 1 to 2-dimensional space.

Table 1: Data set for Question 1.

Data Points	X_1	X_2	X_3	X_4
P1	2	4	1	3
P2	1	2	3	5
P3	-2	-4	-4	-1
P4	0	-1	-2	-6
P5	-1	-1	2	-1

$\mathbf{U} = \text{the Eigenvectors}$

$\mathbf{U} = \text{the Eigenvectors}$

$$\tilde{\Sigma} = \frac{1}{5-1} \mathbf{X}^T \mathbf{X} = \frac{1}{4} (\mathbf{V}\mathbf{D}\mathbf{U}^T)^T \mathbf{V}\mathbf{D}\mathbf{U}^T$$

$$= \frac{1}{4} \mathbf{U} \mathbf{D}^T \boxed{\mathbf{V}^T} \boxed{\mathbf{V}} \mathbf{D} \mathbf{U}^T$$

$$\tilde{\mathbf{D}} = \frac{1}{4} \mathbf{D}^T \mathbf{D} = \mathbf{U} \tilde{\mathbf{D}} \mathbf{U}^T$$

$$\begin{aligned} \tilde{\Sigma} &= \frac{1}{5-1} \mathbf{X}^T \mathbf{X} = \frac{1}{4} (\mathbf{V}\mathbf{D}\mathbf{U}^T)^T \mathbf{V}\mathbf{D}\mathbf{U}^T \\ &= \frac{1}{4} \mathbf{U} \mathbf{D}^T \boxed{\mathbf{V}^T} \boxed{\mathbf{V}} \mathbf{D} \mathbf{U}^T \end{aligned}$$

$$\tilde{\mathbf{D}} = \frac{1}{4} \mathbf{D}^T \mathbf{D} = \mathbf{U} \tilde{\mathbf{D}} \mathbf{U}^T$$

$$\tilde{\Sigma} \mathbf{U} = \mathbf{U} \tilde{\mathbf{D}} \quad \mathbf{U}^T \mathbf{U} = \mathbf{I} \quad \tilde{\mathbf{D}}$$

$$\tilde{\Sigma} \mathbf{u}_i = \tilde{\mathbf{D}}_{ii} \mathbf{u}_i$$

$$\begin{array}{|c c c c|} \hline & 10.9040 & 0 & 0 & 0 \\ \hline & 0 & 4.8385 & 0 & 0 \\ \hline & 0 & 0 & 3.3973 & 0 \\ \hline & 0 & 0 & 0 & 0.3867 \\ \hline & 0 & 0 & 0 & 0 \\ \hline \end{array}$$

Question 1 (cont.)

$$A = U^T \begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_D \end{bmatrix} U$$

$$A = [u_1 \ u_2 \ \dots \ u_D] \begin{bmatrix} \lambda_1 u_1^T \\ \lambda_2 u_2^T \\ \vdots \\ \lambda_D u_D^T \end{bmatrix}$$

$$\alpha A = \alpha \lambda_1 x_1 x_1^T + \alpha \lambda_2 x_2 x_2^T + \dots + \alpha \lambda_D x_D x_D^T$$

$$A = [u_1 \ u_2 \ \dots \ u_D] \begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_D \end{bmatrix} \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_D^T \end{bmatrix}$$

$$A = \lambda_1 u_1 u_1^T + \lambda_2 u_2 u_2^T + \dots + \lambda_D u_D u_D^T$$

$$\alpha A = \alpha \lambda_1 u_1 u_1^T + \alpha \lambda_2 u_2 u_2^T + \dots + \alpha \lambda_D u_D u_D^T$$

$$\alpha A = U^T \begin{bmatrix} \alpha \lambda_1 & 0 & 0 & 0 \\ 0 & \alpha \lambda_2 & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \alpha \lambda_D \end{bmatrix} U$$

Question 1 (cont.)

- Step 2: Compute sample covariance matrix $\tilde{\Sigma}$
- Step 3: Compute eigenvectors of $\tilde{\Sigma}$, $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d\}$, which are sorted based on their eigenvalues in non-increasing order, i.e., $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$



- Perform SVD on the centered data matrix \mathbf{X} (5-by-4) to obtain

$$\mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{U}^\top$$
 (already done in this question)

- Step 4: Select the first k eigenvectors to construct principal components



- Select the first 2 eigenvectors of $\tilde{\Sigma}$, i.e., the first two column vectors of the matrix \mathbf{U} to construct the principle component matrix $\mathbf{U}_2 = [\mathbf{u}_1, \mathbf{u}_2]$

\mathbf{U}_2			
-0.2224	0.3430	-0.3302	-0.8508
-0.4880	0.5756	-0.4046	0.5166
-0.4479	0.2924	0.8400	-0.0911
-0.7154	-0.6823	-0.1473	-0.00309

Table 2: The matrix \mathbf{V} (5-by-5) obtained by SVD ($\mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{U}^\top$)

-0.4577	0.2550	-0.5536	0.4680	0.4472
-0.5612	-0.2150	0.1896	-0.6348	0.4472
0.4497	-0.7183	-0.2749	0.0787	0.4472
0.5206	0.6063	-0.1153	-0.3849	0.4472
0.0486	0.0720	0.7542	0.4730	0.4472

Table 3: The matrix \mathbf{D} (5-by-4) obtained by SVD ($\mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{U}^\top$)

10.9040	0	0	0
0	4.8385	0	0
0	0	3.3973	0
0	0	0	0.3867
0	0	0	0

Table 4: The matrix \mathbf{U} (4-by-4) obtained by SVD ($\mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{U}^\top$)

-0.2224	0.3430	-0.3302	-0.8508
-0.4880	0.5756	-0.4046	0.5166
-0.4479	0.2924	0.8400	-0.0911
-0.7154	-0.6823	-0.1473	-0.0309

Question 1 (cont.)

- Step 5: Project the instances to k-dimensional space



- Compute $\mathbf{X}\mathbf{U}_2$

2

$$\begin{array}{|c|c|c|c|} \hline 2 & 4 & 1 & 3 \\ \hline 1 & 2 & 3 & 5 \\ \hline -2 & -4 & -4 & -1 \\ \hline 0 & -1 & -2 & -6 \\ \hline -1 & -1 & 2 & -1 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline -0.2224 & 0.3430 \\ \hline -0.4880 & 0.5756 \\ \hline -0.4479 & 0.2924 \\ \hline -0.7154 & -0.6823 \\ \hline \end{array} = \begin{array}{|c|c|} \hline -4.9908 & 1.2338 \\ \hline -6.1191 & -1.0402 \\ \hline 4.9038 & -3.4759 \\ \hline 5.6762 & 2.9335 \\ \hline 0.5399 & 0.3486 \\ \hline \end{array}$$

\mathbf{X} \mathbf{U}_2 \mathbf{Z}