

NAME:

ID:

CS 412/512 Machine Learning Midterm 1

100pt

Nov. 2, 2017

- Allocated space should be enough for your answer. Give **brief & clear explanations** for full credits. Points will be taken off for irrelevant/rambling information given within an answer.
- Please write legibly and **circle your final answer**.
- You **may make additional assumptions** if you think it is necessary, but if you do so, **clearly state them**. Your grade will depend on whether the assumption was necessary.
- **Show your work for full credit!**
- **No Internet, no cell phones, no calculators!**

Question		Score	Max Score
1	Basic Concepts		20
2	Decision Trees		24
3	Probability		16
5	Pdfs		12
6	Bayesian Classifiers		25
7	MLE estimate		3
TOTAL			100

NAME:

ID:

1) 20pt – Basic Concepts

- a) 4pts – In a regression problem, you have trained a system to approximate a mapping from x to y . What is the **mean square error** of the estimate ($f'(x)$), over the given test set? *Show your work.*

	Labelled Test $\underline{(x,y)}$	Estimate $f'(x)$	
1	$x = 2, y = 10$	8	y
2	$x = 5, y = 15$	16	x
3	$x = 5, y = 14$	16	y
4	$x = 6, y = 16$	19	y

MSE = ... 4.5

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\frac{1}{4} \cdot ((2)^2 + (1)^2 + (2)^2 + (3)^2)$$

$$4 + 1 + 4 + 9 = 18/4 = 9/2$$

- b) 10pts - You have a training set, a test set and a learning algorithm (for instance a decision tree). Answer as true/False (Note: statements without a qualifier (generally, often etc) claims to hold in general; so choose T if it does indeed.). 2pt each answer. -1 each wrong guess.

- • T / F A zero training set error indicates good generalization performance. ✗
- ✓ • T / F A system that has higher test set error compared to its training set error has overfit to the training set. ✓
- ✓ • T / F As the number of features increases, the risk of overfitting generally increases. ✓
- ✗ • T / F As the number of training samples increases, the risk of overfitting generally decreases. ✓
- ✓ • T / F More complex models with larger number of parameters may fit the training data well, but they are more likely to overfit compared to smaller models. ✓

- c) 6pts – Consider a classification problem with two possible output labels (classes C1 and C2) such that one class (C1) has a 0.9 prior probability.

- 3pts - What is the **base error rate** for this problem, indicate as a percentage? Hint: ZeroR from Weka.

Majority $\frac{9}{10}$ so $\frac{1}{10} = \underline{\underline{\% 10 \text{ base error rate}}}$

- 3pts – What would be the expected error rate if you pick a label randomly (you select C1 and C2 each with a probability of 0.5) for a given x ; indicate as a percentage?

$\frac{5}{10} \cdot \frac{9}{10} + \frac{1}{10} \cdot \frac{5}{10} = \frac{45+5}{100} = \underline{\underline{\% 50}}$

class 1 error
C1 selection & class 2 error
C2 selection & class 1 error

NAME:

ID:

2) 24pt – Entropy, Decision Tree Learning

Given:

shuttle

x	$\log_2 x$
0.25	-2
0.33	-1.6
0.5	-1
0.66	-0.6
0.75	-0.4
1	0

a) 3pt – What is the entropy of a random variable dice that represents the output of a 4-sided (possible outputs are 1,2,3,4) fair dice? Show your work.

1 2 3 4

$$-\frac{1}{4} \log_2 2^{-2} - \frac{1}{4} \log_2 2^{-2} - \frac{1}{4} \log_2 2^{-2} - \frac{1}{4} \log_2 2^{-2}$$

$$\underbrace{-\frac{1}{4} \log_2 2^{-2}}_{1/2} = \frac{1}{2} \cdot 4 = 2$$

b) 3pt – How would the entropy in a) change if the dice was biased (for example, the probability of having a 1 is higher than 2,3,4)? It would:

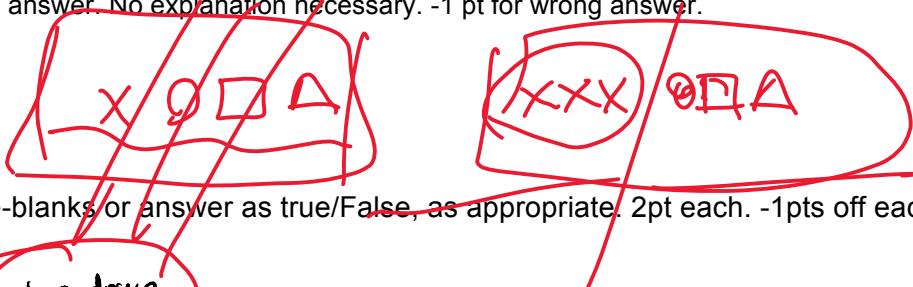
1/4 1/4 1/4

- decrease
- increase
- remain unchanged

$$-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} = +0.3$$

so decrease

Circle the appropriate answer. No explanation necessary. -1 pt for wrong answer.



c) 4 pts – Fill-in-the-blanks or answer as true/False, as appropriate. 2pt each. -1pts off each false guess.

- TF The greedy decision tree learning algorithm ID3 that we saw in class is optimal in the sense that it always generates the smallest tree (least number of nodes). not always
- State one of the most important advantages of using a decision tree classifier.

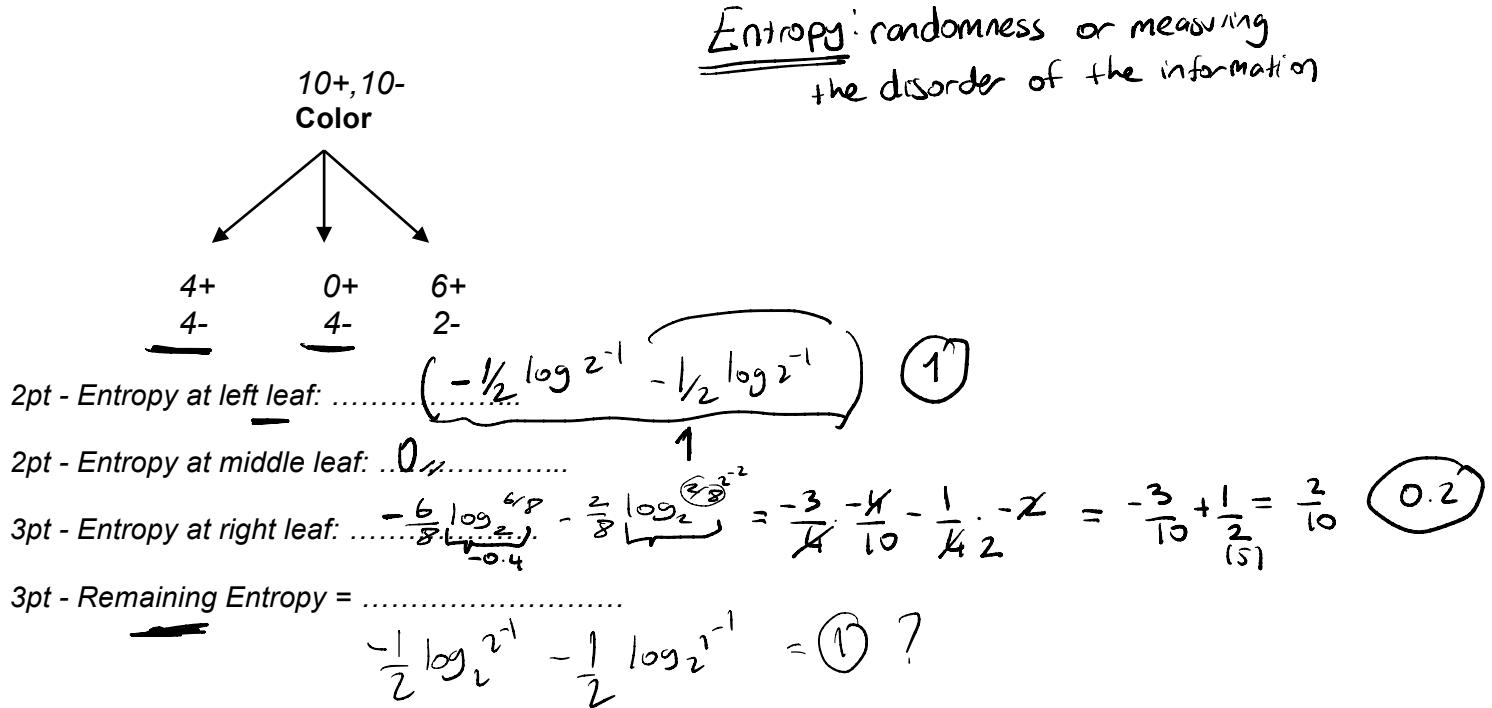
Can be used for both classification and regression problems ✓
 Fast and efficient ✓
 Normalization is not required ✓
Non-parametric less data preparation needed ✓

NAME:

ID:

- d) 10pt – What is the **remaining entropy of class labels** after the tree is split according to the feature “Color”. The +s represent one class, and -s represent another class. There are a total of 20 samples (10+, 10-) at the root of the tree.

NOTE: if the answer is very simple (e.g. entropy is 0 or 1), you can just write that down without writing the whole formula.



- e) 4pts – Your boss gave you a regression problem and asked you to use decision trees with the ID3 algorithm, but the training set size is small.

- 2pts - How would you evaluate different trees, using a validation set or cross-validation? Give a one-line explanation for your reason.
→ In ID3 alg. iteratively divide features. Use entropy and information gain as metrics
→ Two params characterize dt. those we learn by fitting the tree and those we set before training
cv → how to set the hyper params to increase performance of the resulting tree as much as possible
- 2pts - After you are done with training different models and measuring error on validation data, what would you then give to your boss as the **finished system**? Give a one-line answer.
→ Present decision tree models with the comparison of error rates with diff. models.

NAME:

ID:

3) 16pt – Probability Theory

a) 4pt – A pedestrian can be hit by a car with a low probability ($p=0.05$) when crossing the road when the light is green for pedestrians. The probability of being hit by a car while the light is red to pedestrians is expectedly high ($p=0.6$). There is no yellow light in this problem.

What is the total probability of being hit? You should assume that most persons will be reasonable and only get tempted to cross the road at red light with a low probability ($p=0.1$). If you must make an assumption, clearly state it.

$$\frac{5}{100} \cdot \frac{9}{10} + \frac{1}{10} \cdot \frac{6}{10} = \frac{45}{1000} + \frac{60}{1000} = \frac{105}{1000} = 0.105$$

0.105

b) 4pts – Answer the following based on the joint probability table involving two random variables X and Y, given below. Show your work (do not just give a single number).

	Y=Red	Y=Green	
X=1	0.1	0.0	0.1
X=2	0.1	0.4	0.5
X=3	0.3	0.1	0.4

$$P(A \cap B) = P(B|A)P(A)$$

i) $P(X=2) = \underline{0.5} = \underline{1/2}$ 2pt

ii) $P(Y=\text{Red} | X=2) = \frac{P(X=2, Y=\text{Red})}{P(X=2)} = \frac{\underline{0.1}}{\underline{0.5}} = \underline{1/5}$ 2pt

c) 4pts – Assume that two random variables A and B are independent. Simplify the following probabilities (probability terms should be simpler probabilities, involving fewer terms).

• $P(A, B) = \underline{P(A) \cdot P(B)}$ dependent ~~independent~~ $P(A) \cdot P(B|A)$

• $P(A | B) = \frac{\underline{P(B|A) \cdot P(A)}}{\underline{P(B)}}$ For independent event
 $P(A | B) = P(A)$

d) 4pts – Assume that two random variables X, Y are conditionally independent given C. Simplify the following probabilities using the conditional independence information. Hint: Above question ☺

• $P(X, Y | C) = \underline{P(X, Y)} = \underline{P(X) \cdot P(Y)}$

• $P(X | Y, C) = \underline{P(X)}$
 $\frac{\underline{P(Y | C) \cdot P(X)}}{\underline{P(Y | C)}}$ $P(x) =$
 $\underline{P(Y | C) \cdot P(x)}$

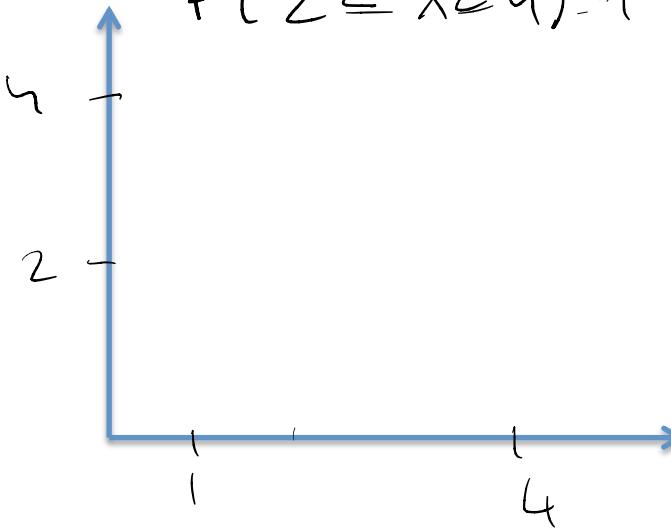
NAME:

ID:

5) 12pt - PDFs

Assume $p(x,y)$ is distributed uniformly in the rectangular area between x in $[1-4]$ and y in $[2-4]$ and 0 elsewhere.

a) 4pts - Draw the pdf $p(x,y)$ making sure to label axes (similar to what we did in our homework)



$$P(2 \leq X \leq 4) = 1 = \left(\frac{1}{2}\right)2$$

$$\begin{aligned} &\text{Find area} = 3, 2 = 6 \\ &\iint_{1}^{4} p(x,y) dx dy \\ &= \int_{2}^{4} \int_{1}^{4} \frac{1}{6} dx dy = 1 \end{aligned}$$

uniform distribution
 $p(x,y) > 1$ conform the rule
 integral = 1

$$\begin{aligned} P &= \int_A p(x,y) dx dy \in [0,1] \\ &\text{for some area subset } A \subseteq [1,4] \times [2,4] \\ p_2(x,y) &= 1 \text{ defined on } [1,2] \times [2,3] \end{aligned}$$

b) 4pts - What is the value of the density ($p(x,y)$) for (x,y) inside the rectangular region?

c) 4pts - What is the marginal probability of $P(2 \leq x \leq 3)$?

NAME:

ID:

6) 25pt – Bayesian Decision Theory

Consider a **classification problem** with input \underline{x} and k possible classes C_i for the questions a)-d).

- a) 3pt - State the Bayes formula that relates **prior**, **posterior** and **conditional probabilities** of a class C_i given some input x . One line formula.

$$P(C_i | \underline{x}) = \frac{P(\underline{x} | C_i) P(C_i)}{P(\underline{x})}$$

likelihood prior

posterior = likelihood * prior

- b) 3pt – What is the **Bayesian decision criterion** that minimizes misclassification error (i.e. to which class do you assign a given x)? One line formula, but do not skip details in the formula.

To be able to minimize misclassification error, we need to choose the one with has higher posterior probability

$$1 - E(\max_i P(c_i | \underline{x}))$$

$$E_{\text{bayes}} = 1 - \sum_{k=1}^K \sum_{c_k} P(c_k) \cdot P(\underline{x} | c_k) dx$$

- c) 3pt – Assume we are given $\underline{x} = [a_1 \ a_2 \ \dots \ a_k]$ where a_i are the attributes C_j is the j th class. How is the term below simplified, if we assume the **Naive Bayes classifier**. Be careful about details/indices....

$$P(a_1, a_2, \dots, a_d | C_j) = \dots \underset{NB}{=} \arg \max_{C_j \in \mathcal{C}} P(C_j)$$

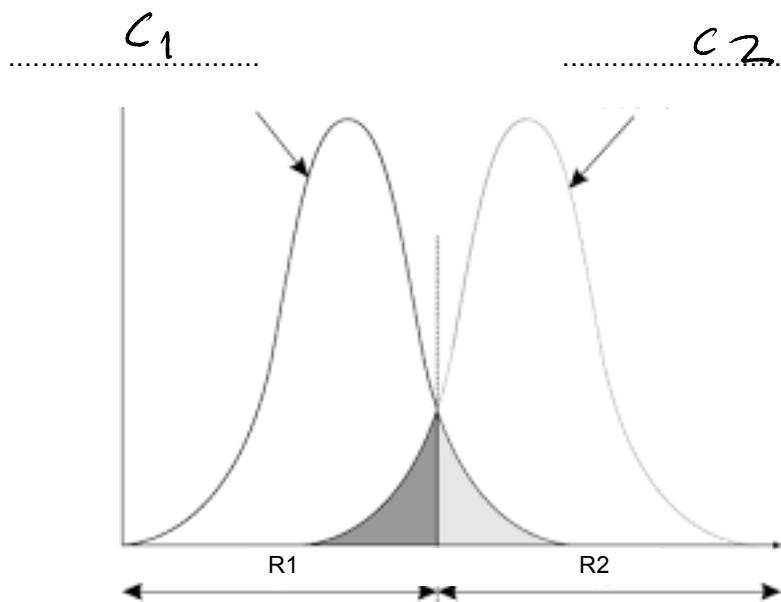
$$\arg \max_{a_i \in a} P(C_j)$$

NAME:

ID:

(d)

- 6pt – Assume we have the following two distributions for x for the two classes C_1 and C_2 .



Graph of
2 normal
distributions

?

- 3pts -- Write the appropriate labels on each of the two distributions, so that the given decision boundary is optimal (minimizes misclassification error). **Be careful, no partial.**
- 3pts – Indicate the area that corresponds to the probability of error corresponding to C_2 samples being classified as class C_1 and give its formula as an integral.

flaw 2

flaw

3

NAME:

ID:

$$\frac{6}{10} \times \frac{1}{2} \times \frac{4}{6} \times \frac{3}{5} \times \frac{5}{6} =$$

f) 10pt -

- 6pts - Consider a naive Bayes classifier trained on the dataset given below. A new patient comes who has $x=[\text{High Fever and Body Ache, but NO runny nose, and NO throat pain}]$. Calculate only the posterior probability of having Flu given these symptoms without considering the denominator ($P(\text{symptoms})$). Do not use smoothing for this example.

You should just leave as a product of probabilities, without doing the final arithmetic.

	Fever	Body Ache	Runny Nose	Throat Pain	Disease
1	High	Yes	No	Yes	Flu
2	High	Yes	No	No	Flu
3	High	No	Yes	No	Flu
4	Medium	Yes	No	No	Flu
5	Medium	No	No	No	Flu
6	High	Yes	No	Yes	Flu
7	Low	No	Yes	Yes	Common cold
8	Low	No	Yes	Yes	Common cold
9	Low	Yes	No	No	Common cold
10	Medium	No	Yes	Yes	Common cold

$$\frac{b_1 b_2 b_3 b_4 b_5}{b_1 b_2 b_3 b_4 b_5} = \frac{1}{10}$$

$$0.14$$

$$0.14$$

$$P(\text{Flu} | \text{High, Yes, No, No}) = \frac{P(\text{High, Yes, No, No} | \text{Flu}) \cdot P(\text{Flu})}{P(\text{High, Yes, No, No})}$$

Use Fever bc

it has higher probability

(1)

- 4pts - Use Laplace smoothing to calculate only:

$$P(\text{Fever=High} | \text{Flu}) = \frac{3+1}{6+1}$$

$$P(\text{Flu} | \text{High}) \cdot P(\text{high})$$

$$P(\text{RunnyNose=No} | \text{Flu}) = \frac{4+1}{6+1}$$

$$P(\text{Flu}) = \frac{6}{10}$$

NAME:

Ixtapa $p(\text{Head})$ ~~$p(\text{Tails})$~~

ID:

phile

7) 3pts -

1

Assume you have observed N coin tosses and 8 of them are Heads and 2 are tails.

- What is the ML estimate of the probability p of observing a head with this coin?

$$\begin{array}{c} \text{Head} \\ \frac{8}{N} \\ \hline N \text{ coin toss} \end{array} \quad \frac{2}{N}$$

8 Head
2 tail

$$\frac{8}{N} = p(\text{head})$$