

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!*

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

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 ( $x,y$ )	Estimate $f(x)$
$x= 2 , y = 10$	8
$x= 5 , y = 15$	16
$x= 5 , y = 14$	16
$x= 6 , y = 16$	19

MSE = .....

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.
- 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?

NAME:

ID:

**2) 24pt – Entropy, Decision Tree Learning**

Given:

<b>x</b>	<b>log<sub>2</sub> x</b>
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.**

**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;

- decrease
- increase
- remain unchanged

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.

- **T/F** 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).
- State **one of the most important advantages** of using a decision tree classifier.

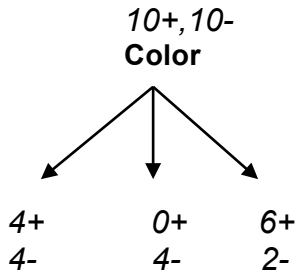
.....

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.**



2pt - Entropy at left leaf: .....

2pt - Entropy at middle leaf: .....

3pt - Entropy at right leaf: .....

3pt - Remaining Entropy = .....

**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.
  
- **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.

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 light 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.*

**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
X=2	0.1	0.4
X=3	0.3	0.1

i)  $P(X=2)=$  ..... 2pt

ii)  $P(Y=Red | X=2) =$  ..... 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).

•  $P(A, B) =$  .....

•  $P(A | B) =$  .....

**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 ☺

•  $P(X, Y | C) =$  .....

•  $P(X | Y, C) =$  .....

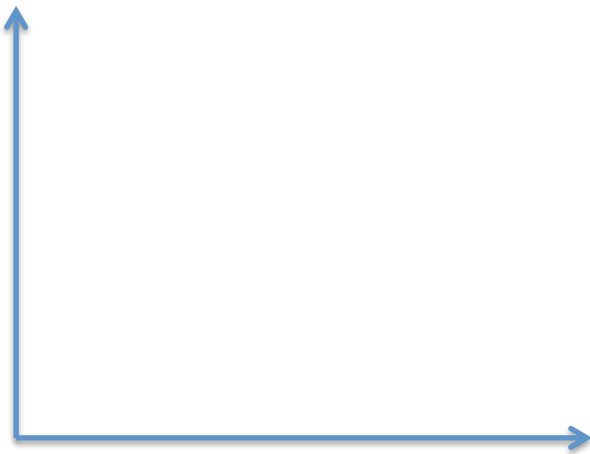
NAME:

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)



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  $\mathbf{x}$  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  $\mathbf{x}$ . *One line formula.*

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

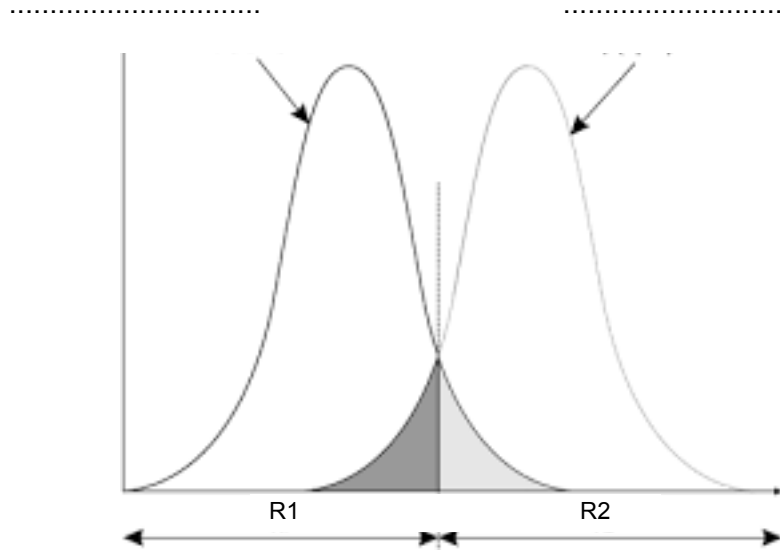
**c) 3pt** – Assume we are given  $\mathbf{x} = [a_1 \ a_2 \ \dots \ a_k]$  where  $a_i$  are the attributes  $C_j$  is the  $j$ th class. How is the term below simplified, if we assume the Naive Bayes classifier. Be careful about details/indices....

$$P(a_1, a_2, \dots, a_d \mid C_j) = \dots\dots\dots$$

NAME:

ID:

d) 6pt – Assume we have the following two distributions for  $x$  for the two classes C1 and C2.



- **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 C2 samples being classified as class C1 and give its formula as an integral.

NAME:

ID:

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

- **4pts** – Use Laplace smoothing to calculate only:

$P(\text{Fever}=\text{High} \mid \text{Flu}) = \dots\dots\dots$

$P(\text{RunnyNose}=\text{No} \mid \text{Flu}) = \dots\dots\dots$

NAME:

ID:

**7) 3pts –**

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?