

Problem Set 0

*Handed Out: August 31, 2017**Due: Not Mandatory*

- The goal of this assignment is to make sure you have the basic terminology, ideas and techniques needed for CIS700. The questions below do not test this exhaustively but provide a sample that will allow you to evaluate whether you are prepared to take this advanced machine learning class.
- Please write down a brief and clear solution. We will not grade it, but rather give you a solution to study and compare with your own work.
- I highly recommend that you spend time on it. My expectation is that you will spend no more than 3 hours on this and solve correctly at least 5 of the 7 problems. Please record the time you spend.
- It is fine if you need to use additional material to refresh your memory or your understanding of the material. My goal is to make sure that you are familiar with a lot of the material, and have sufficient understanding of the area so that you can navigate your way and solve these problems.
- There are plenty of resources you can consult in case you need help. In particular, you can use the notes/videos from my Machine Learning class: <http://L2R.cs.illinois.edu/danr/Teaching/CS446-17/schedule.html>
- Please make sure you complete this assignment before the next lecture. We will provide a solution then.

1. **[Perceptron - 14 points]** In this question, we will be asking you about Perceptrons and their variants.

Let $D = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$, where the j -th example $\mathbf{x}^{(j)}$ is associated with the label $y^{(j)} \in \{-1, +1\}$. Each example $\mathbf{x}^{(j)}$ is a bit-vector of length n , i.e. $\mathbf{x}^{(j)} \in \{0, 1\}^n$, with the interpretation that the i -th bit of the vector ($x_i^{(j)}$) is 1 if the element described by $\mathbf{x}^{(j)}$ has the i -th attribute on.

- (a) [7 points] Let us first consider a Perceptron where the positive example \mathbf{x} satisfies $\mathbf{w} \cdot \mathbf{x} \geq \theta$, where $\mathbf{w} \in \mathbb{R}^n$, $\theta \in \mathbb{R}$ and \mathbf{x} is some example $\mathbf{x}^{(j)}$ from D .

1. [3 points] Suggest an equivalent representation of this Perceptron in the form of $\underline{\mathbf{w}'} \cdot \mathbf{x}' \geq 0$ given an example $\mathbf{x}^{(j)}$, where $\mathbf{x}' \in \{0, 1\}^{n'}$ for some suitable integer n' .

Define $n' = n + 1$

Define $\mathbf{w}' = \langle \mathbf{w}, -\theta \rangle$

Define $\mathbf{x}' = \langle \mathbf{x}^{(j)}, 1 \rangle$

2. [4 points] In the following table, we describe a specific data set S . Using an initialization of $\mathbf{w}' = \mathbf{0}$, i.e. the zero vector, and a learning rate of $R = 1$, complete the columns under (a) of the table using the Perceptron learning algorithm.

S				(a)		(b)	
j	$\mathbf{x}_1^{(j)}$	$\mathbf{x}_2^{(j)}$	$y^{(j)}$	Mistake? Y/N	Updated \mathbf{w}'	Mistake? Y/N	Updated \mathbf{w}'
Initialization				—	$\mathbf{0}$	—	$\mathbf{0}$
1	1	1	+1	N	$\langle 0, 0, 0 \rangle$	N	$\langle 1, 1, 1 \rangle$
2	1	0	-1	Y	$\langle -1, 0, -1 \rangle$	Y	$\langle 0, 1, 0 \rangle$
3	0	1	+1	Y	$\langle -1, 1, 0 \rangle$	N	$\langle 0, 2, 1 \rangle$

Recommend taking off 1 point per wrong entry as opposed to adding 1 for each correct.

(b) [7 points] Using the same data set used above, we now consider a Perceptron with margin $\gamma > 0$. We can also represent this with $\mathbf{w}' \cdot \mathbf{x}' \geq 0$ like in Perceptron but using a different update rule for the weights.

- [3 points] Let the margin $\gamma > 0$ and learning rate $R > 0$. For a given $(\mathbf{x}^{(j)}, y^{(j)})$, write down the update rule for the Perceptron with margin.

$$\text{If } y^{(j)}(\mathbf{w}' \cdot \mathbf{x}') \leq \gamma \text{ then } \mathbf{w}' = \mathbf{w}' + Ry^{(j)}\mathbf{x}'$$

$$\text{otherwise } \mathbf{w}' = \mathbf{w}'$$

- [4 points] We described a specific data set S in a table earlier. Using an initialization of $\mathbf{w}' = \mathbf{0}$, that is, the zero vector, a learning rate of $R = 1$ and margin $\gamma = 1.5$, complete the columns under (b) of the table using the *Perceptron with margin* learning algorithm.

- (c) [11 points] Suppose we have the same data set S and now we would like to learn a linear separator of the form $\mathbf{w}' \cdot \mathbf{x}' \geq 0$, the canonical representation for any separating hyperplane. This time however, we would like to learn the weights \mathbf{w}' by *minimizing* the error made by the linear separator over S .

We define the error made by \mathbf{w}' over S using the *hinge loss* function, defined as $L(y^{(j)}, \mathbf{x}^{(j)}, \mathbf{w}') = \max(0, 1 - y^{(j)} \mathbf{w}' \cdot \mathbf{x}^{(j)'})$, where $\mathbf{x}^{(j)'}$ is the representation of example $\mathbf{x}^{(j)}$ in the form of \mathbf{x}' in the canonical representation.

Thus the goal of learning is to minimize the following error:

$$\text{Error}(\mathbf{w}', D) = \sum_{j=1}^m L(y^{(j)}, \mathbf{x}^{(j)}, \mathbf{w}') = \sum_{j=1}^m \max(0, 1 - y^{(j)} \mathbf{w}' \cdot \mathbf{x}^{(j)'})$$

One way to do this is to make use of Stochastic Gradient Descent.

1. [9 points] Write the pseudocode for Stochastic Gradient Descent using this hinge loss function with a fixed learning rate of $R > 0$. **Set \mathbf{w}' to a random vector drawn from $[0, 1]^{n+1}$.**
 For $(\mathbf{x}^{(j)}, y^{(j)}) \in S$
 Set $\mathbf{x}^{(j)'}$ = $\langle \mathbf{x}^{(j)}, 1 \rangle$.
 Evaluate $o^{(j)} = \mathbf{w}' \cdot \mathbf{x}^{(j)'}$
 If $y^{(j)} o^{(j)} < 1$ then
 Set $\mathbf{w}' = \mathbf{w}' + R y^{(j)} \mathbf{x}^{(j)'}$
2. [2 points] Suggest a condition on the problem definition that will make the Stochastic Gradient Descent algorithm identical to the Perceptron with Margin algorithm. **Change the hinge loss function to be $\max(0, \gamma - y^{(j)} \mathbf{w}' \cdot \mathbf{x}^{(j)'})$**

2. [Kernels - 16 points]

In this question we will develop a learning algorithm that will take as input a URL and classify it according to whether it is relevant to the topic “Machine Learning” or not. The classifier will only depend on the string of the URL, and not on the web page itself.

In the following we will develop a kernel that will be used to learning how to map a URL string to “relevant” and “irrelevant”.

We are given a collection of m URLs u_1, u_2, \dots, u_m . Each URL consists of characters taken from a vocabulary V of n characters c_1, \dots, c_n . We can assume that V includes *all* ASCII characters.

The *basic* feature vector for each URL u is $F(u)$. $F(u)$ is a binary vector, $F(u) \in \{0, 1\}^n$, where the j th component in $F(u)$ indicates whether the character c_j appears in URL u ($F(u)[j] = 1$) or not ($F(u)[j] = 0$). For example, for the URL $u = \text{www.cnn.com}$, the set of active features in $F(u)$ is $A = \{\mathbf{w}, \mathbf{c}, \mathbf{n}, \mathbf{.}, \mathbf{o}, \mathbf{m}\}$, i.e. components in $F(u)$ that correspond to the indices of the characters of A will be 1, all others will be 0.

Each u is also labeled as relevant ($l = 1$) or irrelevant ($l = 0$).

(a) [17 points] The presence of some *set* of characters can be indicative of “machine learning”, e.g. *ml*, *sv*. So in addition to the basic features in $F(u)$, we want to include features that indicate if a *different* pair of characters c_i and c_j appear anywhere in the URL u . Let us call this new feature space $\phi(u)$.

1. [3 points] What is the total number of features in the expanded feature space of $\phi(u)$ (as a function of the vocabulary size n)?

$$n + \binom{n}{2}$$

2. [3 points] Assume a URL $u = \text{www.a.sg}$

Write down the active features in the expanded feature vector $\phi(u)$.

$$A = \{\mathbf{w}, \mathbf{.}, \mathbf{a}, \mathbf{s}, \mathbf{g}, \mathbf{w.}, \mathbf{wa}, \mathbf{ws}, \mathbf{wg}, \mathbf{.a}, \mathbf{.s}, \mathbf{.g}, \mathbf{as}, \mathbf{ag}, \mathbf{sg}\}$$

3. [3 points] For a URL u of length s , where all the s characters are different, what is the number of active features in $\phi(u)$ (as a function of s)?

$$s + \binom{s}{2}$$

4. [8 points] We want to use the Kernel Perceptron algorithm to learn the function above. Design a kernel $K(u_1, u_2)$ that allows us to compute the value of the dot product $\phi(u_1)\phi(u_2)$ in time linear in n by directly computing it in the $F(u)$ space without expanding to the $\phi(u)$ space (assume that length of u_1 and u_2 is at most n). Give the formula for $K(u_1, u_2)$ and explain why it is true.

$K(u_1, u_2) = \text{same}(u_1, u_2) + \binom{\text{same}(u_1, u_2)}{2}$ where $\text{same}(u_1, u_2)$ is the number of characters that appear both in u_1 and u_2

- (b) [8 points] Now we are going to define a new feature space $\psi(u)$. $\psi(u)$ will include all features from $F(u)$ and also include features that indicate whether a pair of characters appear consecutively in URL u . For example, for the URL $u = \text{www.cnn.com}$, the set of active features in $\psi(u)$ will be $A = \{\mathbf{w}, \mathbf{c}, \mathbf{n}, \mathbf{.}, \mathbf{o}, \mathbf{m}, \mathbf{ww}, \mathbf{w.}, \mathbf{.c}, \mathbf{cn}, \mathbf{nn}, \mathbf{n.}, \mathbf{co}, \mathbf{om}\}$.

1. [2 points] What is the size of the expanded feature space $\psi(u)$ (as a function of the vocabulary size n)?

$$n + n^2$$

2. [3 points] For a URL u of length s where all the s characters are different, what is the number of active features of $\psi(u)$ (as a function of s)?

$$s + s - 1$$

3. [3 points] Write a kernel $K(u_1, u_2)$ that can compute the dot product $\psi(u_1)\psi(u_2)$ in time linear in n (assume that length of u_1 and u_2 is at most n).

Since the length of $\psi(u_1)$ and $\psi(u_2)$ are linear in n , computing the dot product directly will be time linear in n . $K(u_1, u_2) = \psi(u_1) \cdot \psi(u_2)$.

3. [Multiclass - 15 points] Multi-class classification [25 points]

You are a newspaper editor and you notice that your journalism interns repeatedly make mistakes in using appropriate prepositions in the essays they write. You collect 100 example sentences where interns made mistakes and label them with the correct preposition from set $Y = \{\text{with, on, for, at}\}$. There are 25 example sentences for each. You want to build a classifier to predict which preposition should appear in a given sentence. Let us denote each example sentence as a n -dimensional boolean feature vector $\langle x_1, \dots, x_n \rangle, x_i \in \{0, 1\}$ that has a label $y \in Y$.

- (a) [10 points] You decide to use the **One vs. All** scheme, but are not sure whether to use the **Perceptron** algorithm or **naïve Bayes**. Fill in the table below, using the options from Table 3.

Perceptron			Naïve Bayes																																																			
(i) Training protocol (training the model)																																																						
<p>[3 points] What classifiers do you need to learn?</p> <p>Fill the table below using values from Table 3. You don't have to use all the rows below if you don't need them.</p> <table border="1" style="width: 100%; border-collapse: collapse; margin-top: 10px;"> <thead> <tr> <th rowspan="2" style="text-align: center;">Classifier #</th> <th colspan="2" style="text-align: center;">Positive examples</th> <th colspan="2" style="text-align: center;">Negative examples</th> </tr> <tr> <th style="text-align: center;">label(s) of examples</th> <th style="text-align: center;"># of examples</th> <th style="text-align: center;">label(s) of examples</th> <th style="text-align: center;"># of examples</th> </tr> </thead> <tbody> <tr><td style="text-align: center;">1</td><td style="text-align: center;">1</td><td style="text-align: center;">20</td><td style="text-align: center;">14</td><td style="text-align: center;">22</td></tr> <tr><td style="text-align: center;">2</td><td style="text-align: center;">2</td><td style="text-align: center;">20</td><td style="text-align: center;">13</td><td style="text-align: center;">22</td></tr> <tr><td style="text-align: center;">3</td><td style="text-align: center;">3</td><td style="text-align: center;">20</td><td style="text-align: center;">12</td><td style="text-align: center;">22</td></tr> <tr><td style="text-align: center;">4</td><td style="text-align: center;">4</td><td style="text-align: center;">20</td><td style="text-align: center;">11</td><td style="text-align: center;">22</td></tr> <tr><td style="text-align: center;">5</td><td></td><td></td><td></td><td></td></tr> <tr><td style="text-align: center;">6</td><td></td><td></td><td></td><td></td></tr> <tr><td style="text-align: center;">7</td><td></td><td></td><td></td><td></td></tr> <tr><td style="text-align: center;">8</td><td></td><td></td><td></td><td></td></tr> </tbody> </table>			Classifier #	Positive examples		Negative examples		label(s) of examples	# of examples	label(s) of examples	# of examples	1	1	20	14	22	2	2	20	13	22	3	3	20	12	22	4	4	20	11	22	5					6					7					8					<p>[3 points] What parameters do you need to estimate from the data?</p> <p>Choose appropriate options from Table 3 to answer the question. 17,19</p>		
Classifier #	Positive examples			Negative examples																																																		
	label(s) of examples	# of examples	label(s) of examples	# of examples																																																		
1	1	20	14	22																																																		
2	2	20	13	22																																																		
3	3	20	12	22																																																		
4	4	20	11	22																																																		
5																																																						
6																																																						
7																																																						
8																																																						

(1) {with}	(5) {with, on}	(11) {with, on, for}	(16) $\Pr(x_i)$	}	$i \in \{1, \dots, n\}$	
(2) {on}	(6) {with, for}	(12) {with, on, at}	(17) $\Pr(x_i y)$			
(3) {for}	(7) {with, at}	(13) {with, for, at}	(18) $\Pr(y x_i)$			$y \in Y$
(4) {at}	(8) {on, for}	(14) {on, for, at}	(19) $\Pr(y)$			
	(9) {on, at}		(20) 25		(22) 75	
	(10) {for, at}	(15) {with, on, for, at}	(21) 50		(23) 100	

Table 1: Options to choose from. You may choose an option multiple times. In the table, the notation $\{\}$ means all examples belonging to the labels given inside the curly brackets.

Perceptron	Naïve Bayes
(ii) Test protocol (labeling unseen examples)	
<p>[2 points] How will you label an unseen example based on the classifiers you learned above? $\text{argmax}_y f_y(x), y \in \{\text{with, on, for, at}\}$</p>	<p>[2 points] How will you label an unseen example based on the parameters you defined above? $\text{argmax}_y P(\vec{x} y) \times P(y), y \in \{\text{with, on, for, at}\}$</p>

(b) [15 points] After you trained the model, you notice that the interns make another common mistake, while using the preposition **in**. You collect 25 more examples where the correct preposition is **in**, and want to learn a classifier that distinguishes *five* classes instead of four. Fill in the table below.

Perceptron	Naïve Bayes
(iii) New training protocol	
<p>[1 point] Do you have to change the Perceptron classifiers you learned in part (i)? Yes</p> <p>[4 points] Justify your answer. Each classifier must be trained using the new data set consisting of 25 more negative examples. The reason for this is that with these new examples the linear separators may need to change.</p> <p>[2 points] How will the test protocol (part (ii)) change? We now do: $\text{argmax}_y f_y(\vec{x}), y \in \{\text{with, on, for, at, in}\}$, so we need to check one more function.</p>	<p>[2 points] What parameters do you need in order to learn the new classifier? $P(y), P(x_i y), y \in \{\text{with, on, for, at, in}\}$</p> <p>[2 points] What parameters from part (i) can you reuse as-is? $P(x_i y), y \in \{\text{with, on, for, at}\}$</p> <p>[2 points] What are the new values of the parameters from part (i) that you have to re-estimate? $P(y) = \frac{1}{5}, y \in \{\text{with, on, for, at}\}$</p> <p>[2 points] Which parameters will you have to estimate afresh (i.e. the parameters that you did not have in part (i))? $P(y), y \in \{\text{in}\}$ and $P(x_i y), y \in \{\text{in}\}$</p>

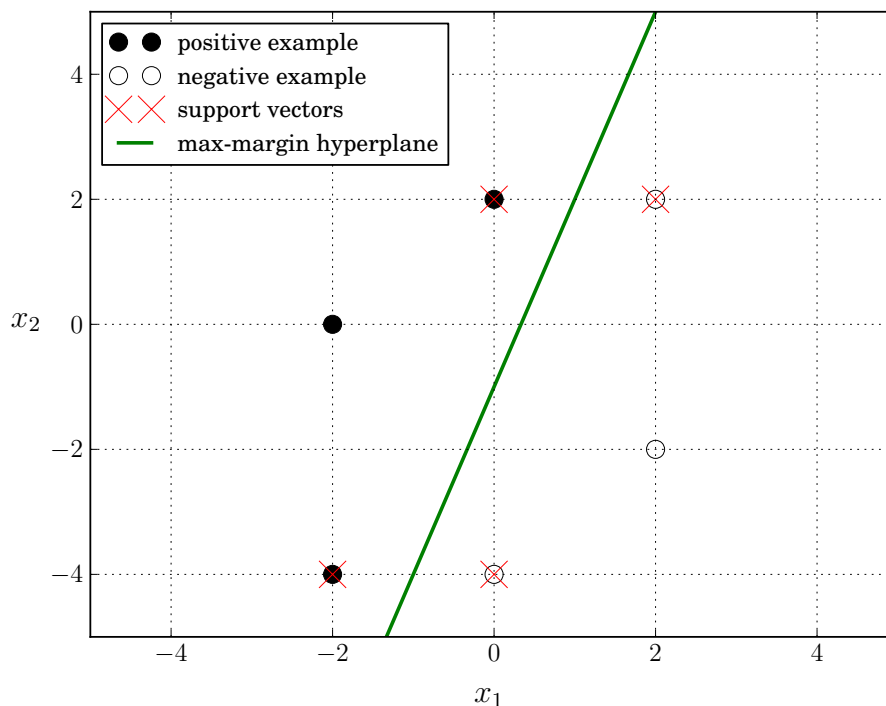
4. **[Support Vector Machines - 15 points]** We are given the following set of training examples $D = \{(x_1^{(i)}, x_2^{(i)}, y^{(i)})\}, i = 1, \dots, m$, where $x_j^{(i)}$ are integer-valued features, and $y^{(i)}$ are binary labels.

x_1	x_2	y
-2	-4	+
-2	0	+
0	2	+
2	2	-
2	-2	-
0	-4	-

Our objective is to learn a hyperplane $w_1x_1 + w_2x_2 + b = 0$ using the hard-SVM objective:

$$\begin{aligned} & \text{minimize} && \frac{1}{2} (w_1^2 + w_2^2) \\ & \text{subject to} && y^{(i)} (w_1x_1^{(i)} + w_2x_2^{(i)} + b) \geq 1, \quad i = 1, \dots, m. \end{aligned}$$

Use the grid below to answer the following questions (you will place a few points and lines on this grid).



(a) [10 points] Finding the hard-SVM hyperplane:

1. [2 points] Place the training examples on the grid, and indicate the support vectors.
See the figure.
2. [3 points] Draw the hyperplane produced by the hard-SVM on the grid.
See the figure.
3. [5 points] Find the values of $w_1, w_2, b \in \mathbb{R}$ that optimize the hard-SVM objective.

The support vectors are at $(-2, -4), (0, 2), (0, -4), (2, 2)$. At the support vectors, the constraints in the optimization of the hard-SVM will be an equality. This gives us the following system of equations:

$$\begin{aligned}+1(-2w_1 - 4w_2 + w_3) &= 1 \\+1(0w_1 + 2w_2 + w_3) &= 1 \\-1(2w_1 + 2w_2 + w_3) &= 1 \\-1(0w_1 - 4w_2 + w_3) &= 1.\end{aligned}$$

By solving this system of equations, we get $w_1 = -1, w_2 = \frac{1}{3}, w_3 = \frac{1}{3}$.

(b) [10 points] Experimental evaluation:

1. [2 points] Provide the classification rule used to classify an example with features x_1, x_2 , using the hyperplane produced by hard-SVM.

$$y = \begin{cases} +1 & \text{if } w_1x_1 + w_2x_2 + b \geq 0 \\ -1 & \text{otherwise.} \end{cases}$$

2. [2 points] What will the error of your classifier be on the training examples D (expressed as the fraction of training examples misclassified)?

0

3. [2 points] Draw on the grid, the hyperplane that will be produced by hard-SVM when you use all training examples except $a = (0, 2, +)$. Using this hyperplane, will you classify a correctly?

The hyperplane will be $x_1 = -1$. No, a will be misclassified.

4. [2 points] Draw on the grid, the hyperplane that will be produced by hard-SVM when you use all training examples except $b = (2, 2, -)$. Using this hyperplane, will you classify b correctly?

The hyperplane will be the same as the one you computed in (a). Yes, b will be classified correctly.

5. [2 points] What will be the average error if you use 6-fold cross validation on the training set D ?

The classifier will make an error when either the support vector at $(0, 2)$, or $(0, -4)$ is left out for validation. In these two cases, the error will be $\frac{1}{6}$, and in all other cases, the error will be 0. The average error will be $\frac{2}{6} \times \frac{1}{6} = \frac{1}{18}$

(c) [5 points] Soft-SVM formulation:

1. [3 points] Write the soft-SVM objective below. Circle either min or max.

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \max \left(0, 1 - y^{(i)}(w_1x_1^{(i)} + w_2x_2^{(i)} + b) \right)$$

2. [2 points] For what range of positive C values, will the hyperplane produced by this soft-SVM objective be most similar to the hyperplane produced by hard-SVM. Circle one of the following.

very small

moderate

very large

5. [Naive Bayes - 14 points]

You would like to study the effects of *irrigation*, *fertilization* and *pesticide* use with respect to the **yield** of a farm. Suppose you are provided with a collection $D = \{D_1, \dots, D_m\}$ of m data points corresponding to m different farms. Each farm has three binary attributes *IsIrrigated* (X_1), *IsFertilized* (X_2) and *UsesPesticide* (X_3), and each has either a high yield ($V = 1$) or a low yield ($V = 0$). The label is **Yield**. A natural model for this is the **multi-variate Bernoulli model**.

Below is a table representing a *specific collection* S of data points for 8 farms to illustrate how a collection might look like.

#	<i>IsIrrigated</i> (X_1)	<i>IsFertilized</i> (X_2)	<i>UsesPesticide</i> (X_3)	Yield (V)
1	No (0)	Yes (1)	No (0)	High (1)
2	Yes (1)	Yes (1)	No (0)	High (1)
3	No (0)	Yes (1)	No (0)	Low (0)
4	No (0)	Yes (1)	No (0)	High (1)
5	No (0)	No (0)	Yes (1)	Low (0)
6	Yes (1)	No (0)	Yes (1)	Low (0)
7	No (0)	No (0)	No (0)	Low (0)
8	No (0)	Yes (1)	No (0)	High (1)

- (a) [6 points] Circle *all* the parameters from the table below that you will need to estimate in order to completely define the model. You may assume that $i \in \{1, 2, 3\}$ for all entries in the table.

(1) $\alpha_i = \Pr(X_i = 1)$	(7) $\beta = \Pr(V = 1)$
(2) $\gamma_i = \Pr(X_i = 0)$	(8) $\varphi = \Pr(V = 0)$
(3) $p_i = \Pr(X_i = 1 \mid V = 1)$	(9) $q_i = \Pr(V = 1 \mid X_i = 1)$
(4) $r_i = \Pr(X_i = 0 \mid V = 1)$	(10) $s_i = \Pr(V = 0 \mid X_i = 1)$
(5) $t_i = \Pr(X_i = 1 \mid V = 0)$	(11) $u_i = \Pr(V = 1 \mid X_i = 0)$
(6) $w_i = \Pr(X_i = 0 \mid V = 0)$	(12) $y_i = \Pr(V = 0 \mid X_i = 0)$

We need (3), (4), (5), (6), (7), (8).

- (b) [3 points] How many **independent** parameters do you have to estimate to learn this model?

Each of X_1 , X_2 and X_3 are binary, so we need to learn $2 - 1 = 1$ independent parameters for each value of $V = v$ for the associated conditional distribution $\Pr(X_i \mid V = v)$. We also need to learn one parameter for $\Pr(V)$ (say $\Pr(V = 1)$) since V is binary also.

This means we need a total of $2 + 2 + 2 + 1 = 7$ independent parameters.

- (c) [5 points] Write explicitly the naïve Bayes classifier for this model as a function of the model parameters selected in (a):

$$\begin{aligned}
 & \Pr(V = v \mid X_1 = x_1, X_2 = x_2, X_3 = x_3) \\
 &= \frac{\Pr(V = v) \prod_{i=1}^3 \Pr(X_i = x_i \mid V = v)}{\Pr(X_1, X_2, X_3)} \\
 &\propto \Pr(V = v) \prod_{i=1}^3 \Pr(X_i = x_i \mid V = v) \\
 &= v \left(\beta \prod_{i=1}^3 p_i^{x_i} r_i^{(1-x_i)} \right) + (1 - v) \left(\varphi \prod_{i=1}^3 t_i^{x_i} w_i^{(1-x_i)} \right)
 \end{aligned}$$

- (d) [5 points] Write the expression for L , the log likelihood of the entire data set D , using the parameters that you have identified in (a).

Since we wrote the naïve Bayes classifier in (b), we can easily write

$$\begin{aligned}
 L = \sum_{j=1}^m & \left(v_j \left[\ln \beta + \sum_{i=1}^3 \left(x_i^{(j)} \ln p_i + (1 - x_i^{(j)}) \ln r_i \right) \right] + \right. \\
 & \left. (1 - v_j) \left[\ln \varphi + \sum_{i=1}^3 \left(x_i^{(j)} \ln t_i + (1 - x_i^{(j)}) \ln w_i \right) \right] \right)
 \end{aligned}$$

where $x_i^{(j)}$ refers to value of X_i for D_j .

- (e) [6 points] We would like to train a Naïve Bayes classifier on S to help us predict the yield on a new farm S_9 .

- [3 points] What is the decision rule for the Naïve Bayes classifier trained on S ?

$$v_{\text{NB}} = \arg \max_{v \in \{0,1\}} \Pr(V = v) \prod_{i=1}^3 \Pr(X_i = x_i \mid V = v)$$

- [3 points] Predict the yield for the following farm using the decision rule written earlier.

#	<i>IsIrrigated</i> (X_1)	<i>IsFertilized</i> (X_2)	<i>UsesPesticide</i> (X_3)	Yield (V)
9	Yes (1)	Yes (1)	Yes (1)	?

We observe that $\Pr(X_3 = 1 \mid V = 1) = 0$, so the value of the objective function for $V = 1$ is 0.

Computing for $V = 0$, we get $\frac{4}{8} \cdot \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{2}{4} > 0$. So by the decision rule we predict the yield of this farm to be Low.

6. **[Probability - 10 points]** You are given the following sample S of data points in order to learn a model. This question will use this data.

Example	A	B	C
1	1	1	0
2	0	1	1
3	1	0	0
4	0	0	0
5	1	1	0
6	0	0	0
7	1	0	1
8	0	1	1
9	1	1	0
10	0	0	0
11	1	1	1
12	0	0	0

- (a) **[3 points]** What would be your estimate for the probability of the following data points, given the sample S , if you were not given any information on a model? (That is, you would estimate the probability directly from the data.)

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

We just count the number of times $A = 1, B = 1$, and $C = 0$ occurs in the dataset. By counting we get $P(A = 1, B = 1, C = 0) = \frac{3}{12} = 0.25$.

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

By counting we get $P(A = 0, B = 1, C = 1) = \frac{2}{12} = \frac{1}{6}$.

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

By counting we get $P(A = 0, B = 0, C = 1) = 0$.

- (b) [10 points] Consider the following graphical model M over three variables A, B , and C .

$$A \rightarrow B \rightarrow C$$

1. [5 points] What are the parameters you need to estimate in order to completely define the model M ? Circle these parameters from Table 2.

(1) $P[A = 1]$	(5) $P[B = 1]$	(9) $P[C = 1]$
(2) $P[A = 1 B = b] \quad b \in \{0, 1\}$	(6) $P[B = 1 C = c] \quad c \in \{0, 1\}$	(10) $P[C = 1 A = a] \quad a \in \{0, 1\}$
(3) $P[A = 1 C = c] \quad c \in \{0, 1\}$	(7) $P[B = 1 A = a] \quad a \in \{0, 1\}$	(11) $P[C = 1 B = b] \quad b \in \{0, 1\}$
(4) $P[A = 1 B, C = b, c]$ $b, c \in \{0, 1\}$	(8) $P[B = 1 A, C = a, c]$ $a, c \in \{0, 1\}$	(12) $P[C = 1 A, B = a, b]$ $a, b \in \{0, 1\}$

Table 2: Options to choose from to explain model M .

2. [5 points] Use the data to estimate the parameters you have circled in (b).1.

The parameters we need to estimate are $P_M[A = 1]$, $P_M[B = 1|A = 1]$, $P_M[B = 1|A = 0]$, $P_M[C = 1|B = 1]$ and $P_M[C = 1|B = 0]$. We estimate these parameters as:

$$P_M[A = 1] = \frac{1}{2} = 0.5$$

$$P_M[B = 1|A = 1] = \frac{4}{6} = \frac{2}{3}$$

$$P_M[B = 1|A = 0] = \frac{2}{6} = \frac{1}{3}$$

$$P_M[C = 1|B = 1] = \frac{3}{6} = 0.5$$

$$P_M[C = 1|B = 0] = \frac{2}{6} = \frac{1}{3}$$

(c) [6 points] Use the parameters chosen in (b).1 to write down expressions for the probabilities of the same data points according to model M and compute these probabilities using the estimated parameters.

1. $P_M(A = 1, B = 1, C = 0)$

Using the conditional independence structure of the graphical model M , we get $P_M(A = 1, B = 1, C = 0) = P_M[A = 1]P_M[B = 1|A = 1]P_M[C = 0|B = 1] = 0.5 * \frac{2}{3} * 0.5 = \frac{1}{6}$.

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

$P_M(A = 0, B = 1, C = 1) = P_M[A = 0]P_M[B = 1|A = 0]P_M[C = 1|B = 1] = 0.5 * \frac{1}{3} * 0.5 = \frac{1}{12}$.

3. $P_M(A = 0, B = 0, C = 1)$

$P_M(A = 0, B = 0, C = 1) = P_M[A = 0]P_M[B = 0|A = 0]P_M[C = 1|B = 0] = 0.5 * \frac{2}{3} * \frac{1}{3} = \frac{1}{9}$.

(d) [6 points] Use the parameters chosen in (b).1 to write down the expressions for the following probabilities for model M and compute these probabilities.

1. $P_M(B = 1)$

$P_M(B = 1) = P_M(B = 1|A = 1)P_M(A = 1) + P_M(B = 1|A = 0)P_M(A = 0) = 0.5\frac{2}{3} + 0.5\frac{1}{3} = 0.5$

2. $P_M(A = 1|B = 0)$

$P_M(A = 1|B = 0) = \frac{P_M(B=0|A=1)P_M(A=1)}{P_M(B=0)} = \frac{0.5 * \frac{1}{3}}{0.5} = \frac{1}{3}$

3. $P_M(A = 0|B = 0, C = 0)$

Since A is independent of C given B , so we have $P_M(A = 1|B = 0, C = 0) = P_M(A = 1|B = 0) = \frac{1}{3}$.
 $P_M(A = 0|B = 0) = 1 - P_M(A = 1|B = 0) = \frac{2}{3}$.

7. [Expectation Maximization - 16 points]

Assume that a set of 3-dimensional points (x, y, z) is generated according to the following probabilistic generative model over Boolean variables $X, Y, Z \in \{0, 1\}$:

$$Y \leftarrow X \rightarrow Z$$

- (a) [4 points] What parameters from Table 3 will you need to estimate in order to completely define the model?

(1) $P(X=1)$	(2) $P(Y=1)$	(3) $P(Z=1)$	
(4) $P(X Y=b)$	(5) $P(X Z=b)$	(6) $P(Y X=b)$	(7) $P(Y Z=b)$
(8) $P(Z X=b)$	(9) $P(Z Y=b)$	(10) $P(X Y=b, Z=c)$	(11) 3

Table 3: Options to choose from. $b, c \in \{0, 1\}$.

Parameters needed are $P(X = 1) = \alpha$, $P(Y = 1|X = 1) = \gamma_{11}$, $P(Y = 1|X = 0) = \gamma_{10}$, $P(Z = 1|X = 1) = \delta_{11}$, $P(Z = 1|X = 0) = \delta_{10}$

- (b) [4 points] You are given a sample of m data points sampled independently at random. However, when the observations are given to you, the value of X is always omitted. Hence, you get to see $\{(y^1, z^1), \dots, (y^m, z^m)\}$. In order to estimate the parameters you identified in part (a), in the course of this question you will derive update rules for them via the EM algorithm for the given model. Express $\Pr(y^j, z^j)$ for an observed sample (y^j, z^j) in terms of the unknown parameters.

$$\begin{aligned}
 P(y^j, z^j) &= P(y^j, z^j, x^j = 0) + P(y^j, z^j, x^j = 1) \\
 &= P(y^j, z^j | x^j = 0)P(x^j = 0) + P(y^j, z^j | x^j = 1)P(x^j = 1) \\
 &= P(y^j | x^j = 0)P(z^j | x^j = 0)P(x^j = 0) + \\
 &\quad P(y^j | x^j = 1)P(z^j | x^j = 1)P(x^j = 1) \\
 &= \gamma_{10}^{y^j} (1 - \gamma_{10})^{(1-y^j)} \delta_{10}^{z^j} (1 - \delta_{10})^{(1-z^j)} (1 - \alpha) + \\
 &\quad \gamma_{11}^{y^j} (1 - \gamma_{11})^{(1-y^j)} \delta_{11}^{z^j} (1 - \delta_{11})^{(1-z^j)} \alpha
 \end{aligned}$$

- (c) [4 points] Let $p_i^j = Pr(X = i | y^j, z^j)$ be the probability that hidden variable X has the value $i \in \{0, 1\}$ for an observation $(y^j, z^j), j \in \{1, \dots, m\}$. Express p_i^j in terms of the unknown parameters.

$$\begin{aligned}
p_0^j &= P(x^j = 0 | y^j, z^j) \\
&= \frac{P(y^j, z^j | x^j = 0) P(x^j = 0)}{P(y^j, z^j)} \\
&= \frac{P(y^j | x^j = 0) P(z^j | x^j = 0) P(x^j = 0)}{P(y^j | x^j = 0) P(z^j | x^j = 0) P(x^j = 0) + P(y^j | x^j = 1) P(z^j | x^j = 1) P(x^j = 1)} \\
&= \frac{(1-\alpha)\gamma_{10}y^j(1-\gamma_{10})^{(1-y^j)}\delta_{10}z^j(1-\delta_{10})^{(1-z^j)}}{(1-\alpha)\gamma_{10}y^j(1-\gamma_{10})^{(1-y^j)}\delta_{10}z^j(1-\delta_{10})^{(1-z^j)} + \alpha\gamma_{11}y^j(1-\gamma_{11})^{(1-y^j)}\delta_{11}z^j(1-\delta_{11})^{(1-z^j)}}
\end{aligned}$$

- (d) [4 points] Let (x^j, y^j, z^j) represent the completed j^{th} example, $j \in \{1, \dots, m\}$. Derive an expression for the expected log likelihood (LL) of the completed data set $\{(x^j, y^j, z^j)\}_{j=1}^m$, given the parameters in (a).

$$\begin{aligned}
E(LL) &= E\left(\sum_{j=1}^m \log P(x^j, y^j, z^j)\right) \\
&= \sum_{j=1}^m p_0^j \log P(x^j = 0, y^j, z^j) + p_1^j \log P(x^j = 1, y^j, z^j) \\
&= \sum_{j=1}^m p_0^j (\log P(y^j | x^j = 0) + \log P(z^j | x^j = 0) + \log P(x^j = 0)) + \\
&\quad \sum_{j=1}^m p_1^j (\log P(y^j | x^j = 1) + \log P(z^j | x^j = 1) + \log P(x^j = 1)) \\
&= \sum_{j=1}^m x_0^j y^j \log \gamma_{10} + x_0^j (1 - y^j) \log(1 - \gamma_{10}) + x_0^j z^j \log \delta_{10} + x_0^j (1 - z^j) \log(1 - \delta_{10}) + \\
&\quad x_1^j y^j \log \gamma_{11} + x_1^j (1 - y^j) \log(1 - \gamma_{11}) + x_1^j z^j \log \delta_{11} + x_1^j (1 - z^j) \log(1 - \delta_{11}) + \\
&\quad p_1^j \log \alpha + p_0^j \log(1 - \alpha)
\end{aligned}$$

- (e) [9 points] Maximize LL , and determine update rules for any two unknown parameters of your choice (from those you identified in part (a)). We will determine α and γ_{10} by differentiating $E(LL)$ by α and γ_{10} .

For α , we get

$$\begin{aligned}\sum_{j=1}^m \frac{p_1^j}{\alpha} - \frac{p_0^j}{1-\alpha} &= 0 \\ \alpha &= \sum_{j=1}^m \frac{p_1^j}{m}\end{aligned}$$

For γ_{10} , we get

$$\begin{aligned}\sum_{j=1}^m p_0^j y^j \frac{1}{\gamma_{10}} - p_0^j (1-y^j) \frac{1}{1-\gamma_{10}} &= 0 \\ \gamma_{10} &= \frac{\sum_{j=1}^m p_0^j y^j}{\sum_{j=1}^m p_0^j}\end{aligned}$$