

10-601 Machine Learning  
Fall 2021  
Exam 2 Practice Problems  
October 22, 2021  
Time Limit: N/A

Name:  
Andrew Email:  
Room:  
Seat:  
Exam Number:

---

**Instructions:**

- Fill in your name and Andrew ID above. Be sure to write neatly, or you may not receive credit for your exam.
  - Clearly mark your answers in the allocated space **on the front of each page**. If needed, use the back of a page for scratch space, but you will not get credit for anything written on the back of a page. If you have made a mistake, cross out the invalid parts of your solution, and circle the ones which should be graded.
  - No electronic devices may be used during the exam.
  - Please write all answers in pen.
  - You have N/A to complete the exam. Good luck!
-

## Instructions for Specific Problem Types

For “Select One” questions, please fill in the appropriate bubble completely:

**Select One:** Who taught this course?

- Matt Gormley
- Marie Curie
- Noam Chomsky

If you need to change your answer, you may cross out the previous answer and bubble in the new answer:

**Select One:** Who taught this course?

- Matt Gormley
- Marie Curie
- Noam Chomsky

For “Select all that apply” questions, please fill in all appropriate squares completely:

**Select all that apply:** Which are scientists?

- Stephen Hawking
- Albert Einstein
- Isaac Newton
- I don't know

Again, if you need to change your answer, you may cross out the previous answer(s) and bubble in the new answer(s):

**Select all that apply:** Which are scientists?

- Stephen Hawking
- Albert Einstein
- Isaac Newton
- I don't know

For questions where you must fill in a blank, please make sure your final answer is fully included in the given space. You may cross out answers or parts of answers, but the final answer must still be within the given space.

**Fill in the blank:** What is the course number?

10-601

10-~~7~~601

# 1 Logistic Regression

1. [2 pts] If today I want to predict the probability that a student sleep more than 8 hours on average (SA) given the Course loading (C), I will choose to use linear regression over logistic regression.

Circle one:      True      False

2. Answer the following questions with brief explanations where necessary.

- a) [2 pts] A generalization of logistic regression to a multiclass settings involves expressing the per-class probabilities  $P(y = c|x)$  as the softmax function  $\frac{\exp(w_c^T x)}{\sum_{d \in C} \exp(w_d^T x)}$ , where  $c$  is some class from the set of all classes  $C$ .

Consider a 2-class problem (labels 0 or 1). Rewrite the above expression for this situation, to end up with expressions for  $P(Y = 1|x)$  and  $P(Y = 0|x)$  that we have already come across in class for binary logistic regression.

- b) [3 pts] Given 3 data points  $(1, 1), (1, 0), (0, 0)$  with labels 0, 1, 0 respectively. Consider 2 models, Model 1:  $\sigma(w_1x_1 + w_2x_2)$ , Model 2:  $\sigma(w_0 + w_1x_1 + w_2x_2)$  ( $\sigma(z)$  is the sigmoid function  $\frac{1}{1+e^{-z}}$ ) that compute  $p(y = 1|\mathbf{x})$ . Using the given data, we can learn parameters  $\hat{w}$  by maximizing the conditional log-likelihood.

Suppose we switched  $(0, 0)$  to label 1 instead.

Do the parameters learnt for Model 1 change?

Circle one:      True      False

One-line explanation:

What about Model 2?

Circle one:      True      False

One-line explanation:

- c) [2 pts] For logistic regression, we need to resort to iterative methods such as gradient descent to compute the  $\hat{w}$  that maximizes the conditional log likelihood. Why?
- d) [3 pts] Considering a Gaussian prior, write out the MAP objective function  $J(w)_{MAP}$  in terms of the MLE objective  $J(w)_{MLE}$ . Name the variant of logistic regression this results in.
3. Given a training set  $\{(x_i, y_i), i = 1, \dots, n\}$  where  $x_i \in \mathbb{R}^d$  is a feature vector and  $y_i \in \{0, 1\}$  is a binary label, we want to find the parameters  $\hat{w}$  that maximize the likelihood for the training set, assuming a parametric model of the form

$$p(y = 1|x; w) = \frac{1}{1 + \exp(-w^T x)}.$$

The conditional log likelihood of the training set is

$$\ell(w) = \sum_{i=1}^n y_i \log p(y_i, |x_i; w) + (1 - y_i) \log(1 - p(y_i, |x_i; w)),$$

and the gradient is

$$\nabla \ell(w) = \sum_{i=1}^n (y_i - p(y_i|x_i; w))x_i.$$

- [5 pts.] Is it possible to get a closed form for the parameters  $\hat{w}$  that maximize the conditional log likelihood? How would you compute  $\hat{w}$  in practice?
  - [5 pts.] For a binary logistic regression model, we predict  $y = 1$ , when  $p(y = 1|x) \geq 0.5$ . Show that this is a linear classifier.
  - Consider the case with binary features, i.e,  $x \in \{0, 1\}^d \subset \mathbb{R}^d$ , where feature  $x_1$  is rare and happens to appear in the training set with only label 1. What is  $\hat{w}_1$ ? Is the gradient ever zero for any finite  $w$ ? Why is it important to include a regularization term to control the norm of  $\hat{w}$ ?
4. Given the following dataset,  $\mathcal{D}$ , and a fixed parameter vector,  $\theta$ , write an expression for the binary logistic regression conditional likelihood.

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)} = 0), (\mathbf{x}^{(2)}, y^{(2)} = 0), (\mathbf{x}^{(3)}, y^{(3)} = 1), (\mathbf{x}^{(4)}, y^{(4)} = 1)\}$$

- Write your answer in terms of  $\theta$ ,  $\mathbf{x}^{(1)}$ ,  $\mathbf{x}^{(2)}$ ,  $\mathbf{x}^{(3)}$ , and  $\mathbf{x}^{(4)}$ .
- Do not include  $y^{(1)}$ ,  $y^{(2)}$ ,  $y^{(3)}$ , or  $y^{(4)}$  in your answer.
- Don't try to simplify your expression.

**Conditional likelihood:**

5. Write an expression for the decision boundary of binary logistic regression with a bias term for two-dimensional input features  $x_1 \in \mathbf{R}$  and  $x_2 \in \mathbf{R}$  and parameters  $b$  (the intercept parameter),  $w_1$ , and  $w_2$ . Assume that the decision boundary occurs when  $P(Y = 1 | \mathbf{x}, b, w_1, w_2) = P(Y = 0 | \mathbf{x}, b, w_1, w_2)$ .

- Write your answer in terms of  $x_1$ ,  $x_2$ ,  $b$ ,  $w_1$ , and  $w_2$ .

**Decision boundary equation:**

- What is the geometric shape defined by this equation?

6. We have now feature engineered the two-dimensional input,  $x_1 \in \mathbb{R}$  and  $x_2 \in \mathbb{R}$ , mapping

it to a new input vector:  $\mathbf{x} = \begin{bmatrix} 1 \\ x_1^2 \\ x_2^2 \end{bmatrix}$

- (a) Write an expression for the decision boundary of binary logistic regression with this feature vector  $\mathbf{x}$  and the corresponding parameter vector  $\boldsymbol{\theta} = [b, w_1, w_2]^T$ . Assume that the decision boundary occurs when  $P(Y = 1 | x, \boldsymbol{\theta}) = P(Y = 0 | x, \boldsymbol{\theta})$ . Write your answer in terms of  $x_1$ ,  $x_2$ ,  $b$ ,  $w_1$ , and  $w_2$ .

**Decision boundary expression:**

- (b) What is the geometric shape defined by this equation?

- (c) If we add an L2 regularization on  $[w_1, w_2]^T$ , what happens to **parameters** as we increase the  $\lambda$  that scales this regularization term?

- (d) If we add an L2 regularization on  $[w_1, w_2]^T$ , what happens to the **decision boundary shape** as we increase the  $\lambda$  that scales this regularization term?

## 2 Feature Engineering and Regularization

1. **Model Complexity:** In this question we will consider the effect of increasing the model complexity, while keeping the size of the training set fixed. To be concrete, consider a classification task on the real line  $\mathbb{R}$  with distribution  $D$  and target function  $c^* : \mathbb{R} \rightarrow \{\pm 1\}$  and suppose we have a random sample  $S$  of size  $n$  drawn iid from  $D$ . For each degree  $d$ , let  $\phi_d$  be the feature map given by  $\phi_d(x) = (1, x, x^2, \dots, x^d)$  that maps points on the real line to  $(d + 1)$ -dimensional space.

Now consider the learning algorithm that first applies the feature map  $\phi_d$  to all the training examples and then runs logistic regression as in the previous question. A new example is classified by first applying the feature map  $\phi_d$  and then using the learned classifier.

- a) [4 pts.] For a given dataset  $S$ , is it possible for the training error to increase when we increase the degree  $d$  of the feature map? **Please explain your answer in 1 to 2 sentences.**
- b) [4 pts.] Briefly **explain in 1 to 2 sentences** why the true error first drops and then increases as we increase the degree  $d$ .

### 3 Neural Networks

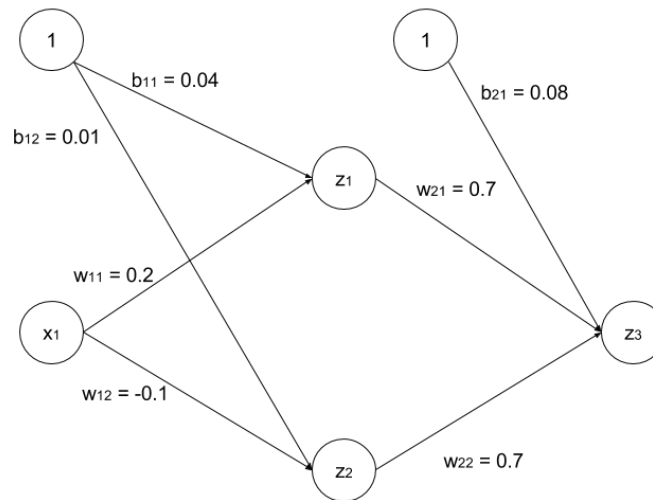


Figure 1: neural network

1. Consider the neural network architecture shown above for a 2-class  $(0, 1)$  classification problem. The values for weights and biases are shown in the figure. We define:

$$a_1 = w_{11}x_1 + b_{11}$$

$$a_2 = w_{12}x_1 + b_{12}$$

$$a_3 = w_{21}z_1 + w_{22}z_2 + b_{21}$$

$$z_1 = \text{relu}(a_1)$$

$$z_2 = \text{relu}(a_2)$$

$$z_3 = \sigma(a_3), \sigma(x) = \frac{1}{1+e^{-x}}$$

Use this information to answer the questions that follow.

- (i) **[6 pts]** For  $x_1 = 0.3$ , compute  $z_3$ , in terms of  $e$ . **Show all work.**

$$z_3 =$$

- (ii) **[2 pts]** To which class does the network predict the given data point ( $x_1 = 0.3$ ), i.e.,  $\hat{y} = ?$  Note that  $\hat{y} = 1$  if  $z_3 > \frac{1}{2}$ , else  $\hat{y} = 0$ .

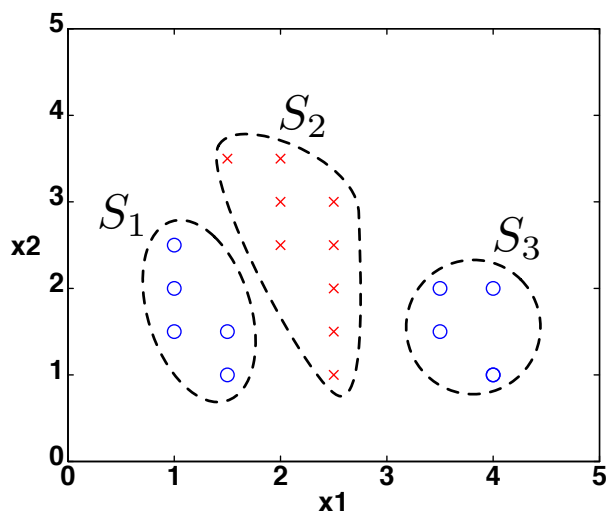
**Circle one:**      0      1

- (iii) **[6 pts]** Perform backpropagation on the bias  $b_{21}$  by deriving the expression for the gradient of the loss function  $L(y, z_3)$  with respect to the bias term  $b_{21}$ ,  $\frac{\partial L}{\partial b_{21}}$ , in

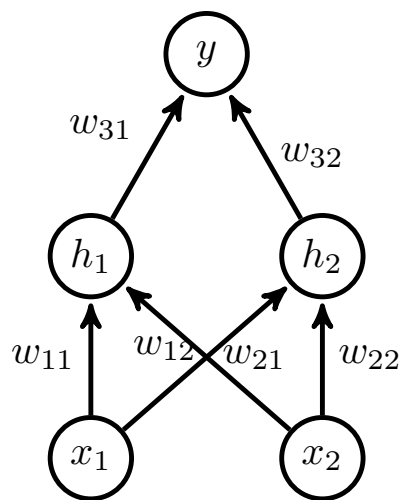
terms of the partial derivatives  $\frac{\partial \alpha}{\partial \beta}$ , where  $\alpha$  and  $\beta$  can be any of  $L, z_i, a_i, b_{ij}, w_{ij}, x_1$  for all valid values of  $i, j$ . Your backpropagation algorithm should be as explicit as possible—that is, make sure each partial derivative  $\frac{\partial \alpha}{\partial \beta}$  cannot be decomposed further into simpler partial derivatives. Do *not* evaluate the partial derivatives.

- (iv) [6 pts] Perform backpropagation on the bias  $b_{12}$  by deriving the expression for the gradient of the loss function  $L(y, z_3)$  with respect to the bias term  $b_{12}$ ,  $\frac{\partial L}{\partial b_{12}}$ , in terms of the partial derivatives  $\frac{\partial \alpha}{\partial \beta}$ , where  $\alpha$  and  $\beta$  can be any of  $L, z_i, a_i, b_{ij}, w_{ij}, x_1$  for all valid values of  $i, j$ . Your backpropagation algorithm should be as explicit as possible—that is, make sure each partial derivative  $\frac{\partial \alpha}{\partial \beta}$  cannot be decomposed further into simpler partial derivatives. Do *not* evaluate the partial derivatives.
2. In this problem we will use a neural network to classify the crosses ( $\times$ ) from the circles ( $\circ$ ) in the simple dataset shown in Figure 2a. Even though the crosses and circles are not linearly separable, we can break the examples into three groups,  $S_1, S_2$ , and  $S_3$  (shown in Figure 2a) so that  $S_1$  is linearly separable from  $S_2$  and  $S_2$  is linearly separable from  $S_3$ . We will exploit this fact to design weights for the neural network shown in Figure 2b in order to correctly classify this training set. For all nodes, we will use the threshold activation function

$$\phi(z) = \begin{cases} 1 & z > 0 \\ 0 & z \leq 0. \end{cases}$$



(a) The dataset with groups  $S_1, S_2$ , and  $S_3$ .



(b) The neural network architecture

Figure 2



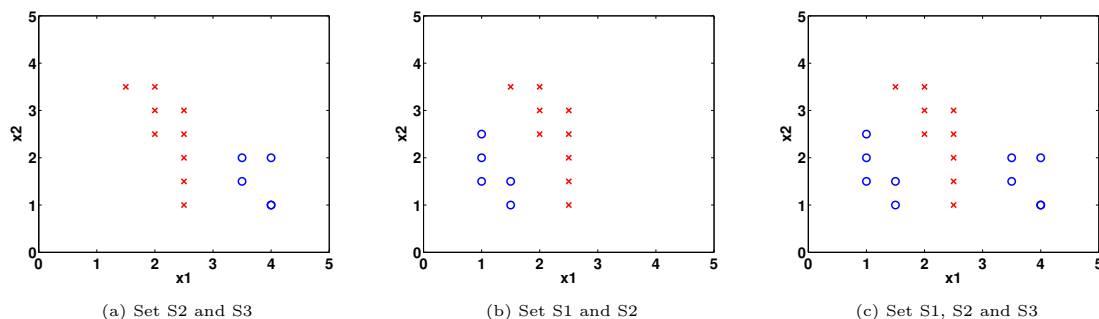


Figure 3: NN classification.

- (i) First we will set the parameters  $w_{11}, w_{12}$  and  $b_1$  of the neuron labeled  $h_1$  so that its output  $h_1(x) = \phi(w_{11}x_1 + w_{12}x_2 + b_1)$  forms a linear separator between the sets  $S_2$  and  $S_3$ .
- (a) [1 pt.] On Fig 3a, draw a linear decision boundary that separates  $S_2$  and  $S_3$ .
- (b) [1 pt.] Write down the corresponding weights  $w_{11}, w_{12}$ , and  $b_1$  so that  $h_1(x) = 0$  for all points in  $S_3$  and  $h_1(x) = 1$  for all points in  $S_2$ . One solution would suffice and the same applies to (ii) and (iii).
- (ii) Next we set the parameters  $w_{21}, w_{22}$  and  $b_2$  of the neuron labeled  $h_2$  so that its output  $h_2(x) = \phi(w_{21}x_1 + w_{22}x_2 + b_2)$  forms a linear separator between the sets  $S_1$  and  $S_2$ .
- (a) [1 pt.] On Fig 3b, draw a linear decision boundary that separates  $S_1$  and  $S_2$ .
- (b) [1 pt.] Write down the corresponding weights  $w_{21}, w_{22}$ , and  $b_2$  so that  $h_2(x) = 0$  for all points in  $S_1$  and  $h_2(x) = 1$  for all points in  $S_2$ .
- (iii) Now we have two classifiers  $h_1$  (to classify  $S_2$  from  $S_3$ ) and  $h_2$  (to classify  $S_1$  from  $S_2$ ). We will set the weights of the final neuron of the neural network based on the results from  $h_1$  and  $h_2$  to classify the crosses from the circles. Let  $h_3(x) = \phi(w_{31}h_1(x) + w_{32}h_2(x) + b_3)$ .
- (a) [1 pt.] Compute  $w_{31}, w_{32}, b_3$  such that  $h_3(x)$  correctly classifies the entire dataset.
- (b) [1 pt.] Draw your decision boundary in Fig 3c.
- (iv) **Back propagation**
- In the above example, we need to learn the weights by according to the data. At first step, we need to get the gradients of the parameters of neural networks.

Suppose there  $m$  data points  $x_i$  with label  $y_i$ , where  $i \in [1, m]$ .  $x_i$  is a  $d \times 1$  vector and  $y_i \in \{0, 1\}$ . We use the data to train a neural network with one hidden layer:

$$h(x) = \sigma(W_1x + b_1)$$

$$p(x) = \sigma(W_2h(x) + b_2),$$

where  $\sigma(x) = \frac{1}{1+\exp(-x)}$  is the sigmoid function,  $W_1$  is a  $n$  by  $d$  matrix and  $b_1$  is a  $n$  by 1 vector,  $W_2$  is a 1 by  $n$  matrix and  $b_2$  is a 1 by 1 vector.

We use cross entropy loss function and minimize the negative log likelihood to train the neural network:

$$l = \frac{1}{m} \sum_i l_i = \frac{1}{m} \sum_i -(y_i \log p_i + (1 - y_i) \log(1 - p_i)),$$

where  $p_i = p(x_i)$ ,  $h_i = h(x_i)$ .

- Describe how you would drive the gradients w.r.t the parameters  $W_1, W_2$  and  $b_1, b_2$ . (No need to write out the detailed mathematical expression.)
  - When  $m$  is large, we typically use a small sample of all the data set to estimate the gradient, this is call stochastic gradient descent (SGD). Explain why we use SGD instead of gradient descent.
  - Work out the following gradient:  $\frac{\partial l}{\partial p_i}, \frac{\partial l}{\partial W_2}, \frac{\partial l}{\partial b_2}, \frac{\partial l}{\partial h_i}, \frac{\partial l}{\partial W_1}, \frac{\partial l}{\partial b_1}$ . When deriving the gradient w.r.t. the parameters in lower layers, you can may assume the gradient in upper layers are available to you (i.e., you can use them in your equation). For example, when calculating  $\frac{\partial l}{\partial W_1}$ , you can assume  $\frac{\partial l}{\partial p_i}, \frac{\partial l}{\partial W_2}, \frac{\partial l}{\partial b_2}, \frac{\partial l}{\partial h_i}$  are known.
3. Consider the following neural network for a 2-D input,  $x_1 \in \mathbb{R}$  and  $x_2 \in \mathbb{R}$  where:

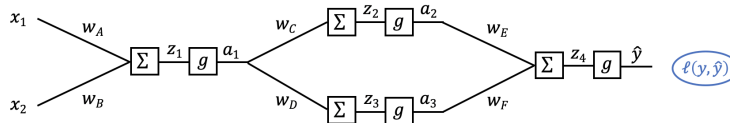


Figure 4: Neural Network

- All  $g$  functions are the same arbitrary non-linear activation function with no parameters
- $\ell(y, \hat{y})$  is an arbitrary loss function with no parameters, and:

$$z_1 = w_A x_1 + w_B x_2 \quad a_1 = g(z_1)$$

$$z_2 = w_C a_1 \quad a_2 = g(z_2)$$

$$z_3 = w_D a_1 \quad a_3 = g(z_3)$$

$$z_4 = w_E a_2 + w_F a_3 \quad \hat{y} = g(z_4)$$

**Note:** There are no bias terms in this network.

- What is the chain of partial derivatives needed to calculate the derivative  $\frac{\partial \ell}{\partial w_E}$ ?

Your answer should be in the form:  $\frac{\partial \ell}{\partial w_E} = \frac{\partial?}{\partial?} \frac{\partial?}{\partial?} \dots$ . Make sure each partial derivative  $\frac{\partial?}{\partial?}$  in your answer cannot be decomposed further into simpler partial derivatives.

**Do not evaluate the derivatives.** Be sure to specify the correct subscripts in your answer.

$$\frac{\partial \ell}{\partial w_E} =$$

(b) The network diagram from above is repeated here for convenience: What is the

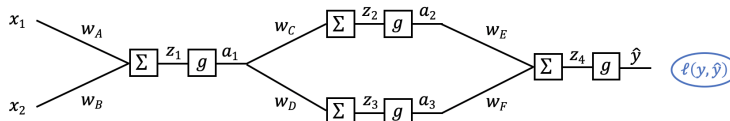


Figure 5: Neural Network

chain of partial derivatives needed to calculate the derivative  $\frac{\partial \ell}{\partial w_C}$ ?  
Your answer should be in the form:

$$\frac{\partial \ell}{\partial w_C} = \frac{\partial ?}{\partial ?} \frac{\partial ?}{\partial ?} \dots$$

Make sure each partial derivative  $\frac{\partial ?}{\partial ?}$  in your answer cannot be decomposed further into simpler partial derivatives. **Do not evaluate the derivatives.** Be sure to specify the correct superscripts in your answer.

$$\frac{\partial \ell}{\partial w_C} =$$

(c) The gradient descent update step for weight  $w_c$  is:

$$w_c \leftarrow w_c - \alpha \frac{\partial Q}{\partial t} = \frac{\partial s}{\partial t}$$

where  $\alpha$  (alpha) is the learning rate (step size).

Now, we want to change our neural network objective function to add an L2 regularization term on the weights. The new objective is:

$$\ell(y, \hat{y}) + \lambda \frac{1}{2} \|w\|_2^2$$

where  $\lambda$  (lambda) is the regularization hyperparameter and  $\mathbf{w}$  is all of the weights in the neural network stacked into a single vector,  $\mathbf{x} = [w_A, w_B, w_C, w_D, w_E, w_F]^T$ . Write the right-hand side of the new gradient descent update step for weight  $w_C$  given this new objective function. You may use  $\frac{\partial \ell}{\partial w_C}$  in your answer.

**Update:**  $w_C \leftarrow \dots$

## 4 MLE/MAP

1. Please circle **True** or **False** for the following questions, providing brief explanations to support your answer.

- (i) [2 pts] Consider the linear regression model  $y = w^T x + \epsilon$ . Assuming  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  and maximizing the conditional log-likelihood is equivalent to minimizing the sum of squared errors  $\|y - w^T x\|_2^2$ .

Circle one:      True      False

One line justification (only if False):

- (ii) [4 pts] Consider  $n$  data points, each with one feature  $x_i$  and an output  $y_i$ . In linear regression, we assume  $y_i \sim \mathcal{N}(wx_i, \sigma^2)$  and compute  $\hat{w}$  through MLE.

Suppose  $y_i \sim \mathcal{N}(\log(wx_i), 1)$  instead. Then the maximum likelihood estimate  $\hat{w}$  is the solution to the following equality:

$$\sum_{i=1}^n x_i y_i = \sum_{i=1}^n x_i \log(wx_i)$$

Circle one:      True      False

Brief explanation:

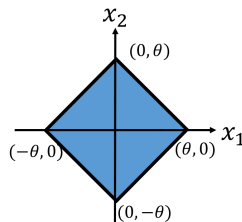
2. **Select all that apply:** Which of the following are correct regarding Gradient Descent (GD). Assume data log-likelihood is  $L(\theta|X)$ , which is a function of the parameter  $\theta$ , and the objective function is negative log-likelihood .

- GD requires that  $L(\theta|X)$  is concave with respect to parameter  $\theta$  in order to converge
- GD requires that  $L(\theta|X)$  is convex with respect to parameter  $\theta$  in order to converge
- GD update rule is  $\theta \leftarrow \theta - \alpha \nabla_{\theta} L(\theta|X)$
- Given a fixed small learning rate (say  $\alpha = 10^{-10}$ ), GD will always reach the optimum after infinite iterations (assume that the objective function satisfies the convergence condition).

3. Let  $X_1, X_2, \dots, X_N$  be i.i.d. data from a uniform distribution over a diamond-shaped area with edge length  $\sqrt{2}\theta$  in  $\mathbb{R}^2$ , where  $\theta \in \mathbb{R}^+$  (see Figure 6). Thus,  $X_i \in \mathbb{R}^2$  and the distribution is

$$p(x|\theta) = \begin{cases} \frac{1}{2\theta^2} & \text{if } \|x\| \leq \theta \\ 0 & \text{otherwise} \end{cases}$$

where  $\|x\| = |x_1| + |x_2|$  is L1 norm. Please find the maximum likelihood estimator of  $\theta$ .

Figure 6: Area of  $\|x\| \leq \theta$ 

4. **Short answer:** Suppose we want to model a 1-dimensional dataset of  $N$  real valued features  $(x^{(i)})$  and targets  $(y^{(i)})$  by:

$$y^{(i)} \sim \mathcal{N}(\exp(wx^{(i)}), 1)$$

Where  $w$  is our unknown (scalar) parameter and  $\mathcal{N}$  is the normal distribution with probability density function:

$$f(a)_{\mathcal{N}(\mu, \sigma^2)} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(a - \mu)^2}{2\sigma^2}\right)$$

Can the maximum conditional negative log likelihood estimator of  $w$  be solved analytically? If so, find the expression for  $w_{\text{MLE}}$ . If not, say so and write down the update rule for  $w$  in gradient descent.

5. Assume we have  $n$  iid random variables  $x_i, i \in [1, n]$  such that each  $x_i$  belongs to a normal distribution with mean  $\mu$  and variance  $\sigma^2$ .

$$p(x_1, x_2, \dots, x_n | \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{(x_i - \mu)^2}{2\sigma^2}\right\}$$

- Write the log likelihood function  $l(x_1, x_2, \dots, x_n | \mu, \sigma^2)$
  - Derive an expression for the Maximum Likelihood Estimate for the variance ( $\sigma^2$ )
6. Assume we have a random variable that is Bernoulli distributed  $X_1, \dots, X_n \sim \text{Bernoulli}(\theta)$ . We are going to derive its MLE. Recall that in a Bernoulli  $X = \{0, 1\}$  and the pdf of a Bernoulli is

$$p(X; \theta) = \theta^x (1 - \theta)^{1-x}$$

- Derive the likelihood,  $L(\theta; X_1, \dots, X_n)$
- Derive the following formula for the log likelihood

$$l(\theta; X_1, \dots, X_n) = \left(\sum_{i=1}^n X_i\right) \log(\theta) + \left(n - \sum_{i=1}^n X_i\right) \log(1 - \theta)$$

- Derive the MLE,  $\hat{\theta}$ , and show that  $\hat{\theta} = \frac{1}{n} \left(\sum_{i=1}^n X_i\right)$

7. Assume we have a random sample that is Bernoulli distributed  $X_1, \dots, X_n \sim \text{Exponential}(\theta)$ . We are going to derive the MLE for  $\theta$ . Recall that a exponential random variable  $X$  has p.d.f:

$$P(X; \theta) = \theta \exp(-\theta X).$$

- a) Derive the likelihood,  $L(\theta; X_1, \dots, X_n)$ .
  - b) Find  $\theta$  that maximizes  $L(\theta; X_1, \dots, X_n)$ .
8. For each question state **True** or **False** and give one line justifications.
- a) **T or F** The value of the Maximum Likelihood Estimate (MLE) is equal to the value of the Maximum A Posteriori (MAP) Estimate with a uniform prior.
  - b) **T or F** The bias of the Maximum Likelihood Estimate (MLE) is typically less than or equal to the bias of the Maximum A Posteriori (MAP) Estimate.
  - c) **T or F** The MAP estimate is always better than the MLE.
  - d) **T or F** In the limit as  $n$  (the number of samples) increases, the MAP and MLE estimates become the same.
  - e) **T or F** Naive Bayes can only be used with MAP estimates, and not MLE estimates.

## 5 Probability, Naive Bayes and MLE

### 5.1 Probability

- For each question, circle the correct option.
  - Which of the following expressions is equivalent to  $p(A|B, C, D)$ ?
    - $\frac{p(A, B, C, D)}{p(C|B, D)p(B|D)p(D)}$
    - $\frac{p(A, B, C, D)}{p(B, C)p(D)}$
    - $\frac{p(A, B, C, D)}{p(B, C|D)p(B)p(C)}$
  - Let  $\mu$  be the mean of some probability distribution.  $p(\mu)$  is always non-zero.
    - True
    - False
- Assume we have a sample space  $\Omega$ . Just state **T** or **F**, no justification needed.
  - If events  $A$ ,  $B$ , and  $C$  are disjoint then they are independent.
  - $P(A|B) \propto \frac{P(A)P(B|A)}{P(A|B)}$ .
  - $P(A \cup B) \leq P(A)$ .
  - $P(A \cap B) \geq P(A)$ .

### 5.2 Naive Bayes

- Consider the following data. It has 4 features  $\mathbf{X} = (x_1, x_2, x_3, x_4)$  and 3 labels  $(+1, 0, -1)$ . Assume that the probabilities  $p(\mathbf{X}|y)$  and  $p(y)$  are both Bernoulli distributions. Answer the questions that follow under the Naive Bayes assumption.

$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	1	0	1	+1
0	1	1	0	+1
1	0	1	1	0
0	1	1	1	0
0	1	0	0	-1
1	0	0	1	-1
0	0	1	1	-1

- Compute the Maximum Likelihood Estimate for  $p(x_i = 1|y), \forall i \in [1, 4], \forall y \in \{+1, 0, -1\}$ .
- Compute the Maximum Likelihood Estimate for the prior probabilities  $p(y = +1), p(y = 0), p(y = -1)$



3. Use the values computed in the above two parts to classify the data point  $(x_1 = 1, x_2 = 1, x_3 = 1, x_4 = 1)$  as either belonging to class  $+1, 0$  or  $-1$
2. You are given a data set of 10,000 students with their sex, height, and hair color. You are trying to build a machine learning classifier to predict the sex of a student, so you randomly split the data into a training set and a testing set. Here are the specifications of the data set:
  - $\text{sex} \in \{\text{male}, \text{female}\}$
  - $\text{height} \in [0, 300]$  centimeters
  - $\text{hair} \in \{\text{brown}, \text{black}, \text{blond}, \text{red}, \text{green}\}$
  - 3240 men in the data set
  - 6760 women in the data set

Under the assumptions necessary for Naive Bayes (not the distributional assumptions you might naturally or intuitively make about the dataset) answer each question with **T** or **F** and a one sentence explanation of your answer:

1. **T or F:** Height is a continuous valued variable. Therefore Naive Bayes is not appropriate since it cannot handle continuous valued variables.
2. **T or F:** Since there is not a similar number of men and women in that dataset Naive Bayes will have high test error.
3. **T or F:**  $p(\text{height}|\text{sex}, \text{hair}) = p(\text{height}|\text{sex})$ .
4. **T or F:**  $p(\text{height}, \text{hair}|\text{sex}) = p(\text{height}|\text{sex}) * p(\text{hair}|\text{sex})$ .

### 5.3 Naive Bayes, Logistic Regression

1. Suppose you wish to learn  $P(Y|X_1, X_2, X_3)$ , where  $Y, X_1, X_2$  and  $X_3$  are all boolean-valued random variables. You consider both Naive Bayes and Logistic Regression as possible approaches.

For each of the following, answer True or False, and give a *one sentence* justification for your answer.

1. T or F: In this case, a good choice for Naive Bayes would be to implement a Gaussian Naive Bayes classifier.
  2. T or F: To learn  $P(Y|X_1, X_2, X_3)$  using Naive Bayes, you must make conditional independence assumptions, including the assumption that  $Y$  is conditionally independent of  $X_1$  given  $X_2$ .
  3. T or F: Logistic regression is certain to be the better choice in this case.
2. Parameter estimation

1. How many parameters must be estimated for your Gaussian Naive Bayes classifier, and what are they (i.e., please list them).
  2. How many parameters must be estimated for your Logistic Regression classifier, and what are they (i.e., please list them).
  3. T or F: We can train Naive Bayes using maximum likelihood estimates for each parameter, but not MAP estimates. Justify your answer *in one sentence*.
  4. T or F: We can train Logistic Regression using maximum likelihood estimates for each parameter, but not MAP estimates. Justify your answer *in one sentence*.
3. Mixing discrete and continuous variables. Suppose we add a numeric, real-valued variable  $X_4$  to our problem. Note we now have a mix of some discrete-valued  $X_i$  and one continuous  $X_j$ .
1. Explain in *two sentences* why we can no longer use Naive Bayes, or if we can, how we would modify our first solution.
  2. Explain in *two sentences* why we can no longer use Logistic Regression, or if we can, how we would modify our first solution.

## 6 PAC Learning

1. **True and Sample Errors:** Consider a classification problem with distribution  $D$  and target function  $c^* : \mathcal{R}^d \mapsto \pm 1$ . For any sample  $S$  drawn from  $D$ , answer whether the following statements are true or false, along with a brief explanation.

a) [4 pts] For a given hypothesis space  $H$ , it is possible to define a sufficient size of  $S$  such that the true error is bounded by the sample error by a margin  $\epsilon$ , for all hypotheses  $h \in H$  with a given probability.

b) [4 pts] The true error of any hypothesis  $h$  is an upper bound on its training error on the sample  $S$ .

2. Let  $X$  be the feature space and there is a distribution  $D$  over  $X$ . We have training samples

$$S : (x_1, c^*(x_1)), \dots, ((x_m, c^*(x_m))),$$

$x_i$  i.i.d from  $D$ . We assume labels  $c^*(x_i) \in \{-1, 1\}$ .

Let  $\mathcal{H}$  be a concept class and let  $h \in \mathcal{H}$  be a concept. In this question we restrict ourselves to  $\mathcal{H}$ . We use

$$err_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbf{I}(h(x_i) \neq c^*(x_i))$$

to denote the training error and

$$err_D(h) = \Pr_{x \sim D}(h(x) \neq c^*(x))$$

to denote the true error. Recall the theorem from class, if the concept class is finite, in the realizable case

$$m \geq \frac{1}{\epsilon} \left[ \ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right]$$

labeled examples are sufficient so that with probability at least  $1 - \delta$ , all  $h \in \mathcal{H}$  with  $err_D(h) \geq \epsilon$  have  $err_S(h) > 0$ ; in the agnostic case,

$$m \geq \frac{1}{2\epsilon^2} \left[ \ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right]$$

labeled examples are sufficient such that with probability at least  $1 - \delta$ , all  $h \in \mathcal{H}$  have  $|err_D(h) - err_S(h)| < \epsilon$ .

a) [2 pts] Briefly explain what is realizable case and what is agnostic case.

b) [4 pts] What is the full name of PAC learning? What is the correspondence between  $\epsilon, \delta$  and the full name?

- c) [4 pts] (True or False) Consider two concept finite classes  $\mathcal{H}_1$  and  $\mathcal{H}_2$  such that  $\mathcal{H}_1 \subset \mathcal{H}_2$ . Let  $h_1 = \arg \min_{h \in \mathcal{H}_1} \text{err}_S(h)$  and  $h_2 = \arg \min_{h \in \mathcal{H}_2} \text{err}_S(h)$ . Thus according to the theorem, because  $|\mathcal{H}_2| \geq |\mathcal{H}_1|$ ,  $\text{err}_D(h_2) \geq \text{err}_D(h_1)$ . Briefly justify your answer.