

# CS4780 Midterm

Fall 2018

NAME:	
Net ID:	
Email:	

I promise to abide by Cornell's Code of Academic Integrity.

Signature: \_\_\_\_\_

# 1 [??] General Machine Learning

Please identify if these statements are either True or False. Please justify your answer **if false**. Correct “True” questions yield 1 point. Correct “False” questions yield 2 points, one for the answer and one for the justification.

1. (T/F) As  $n \rightarrow \infty$ , the 1-NN error is no more than twice the error of the Bayes Optimal classifier.  
T
2. (T/F) MLE can overfit the data if  $n$  (the number of training samples) is small. It tends to work well when  $n$  is large.  
T.
3. (T/F) Both, Gradient descent and Newton's method use only a 1st order approximation of the function to be minimized.  
F. Newton's method uses 2nd order approximation of the function
4. (T/F) If a data set is linearly separable, the Perceptron guarantees that you find a hyperplane but the SVM finds the maximum margin separating hyperplane.  
T
5. (T/F) The best machine learning algorithm make no assumptions about the data.  
F. ML algorithms always make assumptions about the data.
6. (T/F) The k-NN classifier is not a linear classifier. T

7. (T/F) The k-NN algorithm can be used for classification, but not regression.  
F, k-NN can be used for regression by averaging the labels of the k nearest neighbors.
  
8. (T/F) The order of the training points can affect the training time of the Perceptron algorithm. T
  
9. (T/F) Even on non-linearly-separable datasets, the Perceptron algorithm is guaranteed to converge in finite time.  
F. For datasets that are not linearly separable, the Perceptron algorithm can never finish: it runs forever.
  
10. (T/F) In MAP, we find the maximizer of the posterior, so we need to find an expression for the posterior.  
F. Because,  $\arg \max_{\theta} P(\theta|D) = \arg \max_{\theta} \frac{P(D|\theta)P(\theta)}{P(D)} = \arg \max_{\theta} P(D|\theta)P(\theta)$
  
11. (T/F) If you were to use the “true” Bayesian way of machine learning you would put a prior over the possible models and draw several models randomly during training.  
T
  
12. (T/F) If the features are probabilistically dependent on each other, then the naive Bayes assumption cannot hold. F. the features could be **conditionally** independent, given the label.

13. (**T/F**) Logistic regression is a generative model. **F. It's a discriminative model, not a generative model.**
14. (**T/F**) The order of the training points can affect the convergence of the gradient descent algorithm.  
**F, gradient descent just depends on a sum across the training examples, and sums are independent of order.**
15. (**T/F**) For gradient descent, higher learning rates guarantee faster convergence times.  
**F, higher learning rates can lead to divergence.**
16. (**T/F**) For Adagrad, we use the same learning rate for all features.  
**F, Adagrad uses different automatically-chosen learning rates for each feature.**

## 2 [16] K-NN

1. (2 pts) Imagine you apply the kNN classifier with Euclidean distance. Describe what happens if you scale one dimension of the input features by a large positive constant across all examples?

This feature will dominate the distance metric and nearest neighbors will eventually (as the constant is very large) be only measured in that dimension.

2. (2 pts) What is the modeling assumption of kNN?

Points that are close have similar label.

3. (4 pts) Consider points, sampled uniformly at random, within a finite volume in  $d$  dimensions. How do the pairwise distances change as the dimensionality  $d$  increases? How is the distance to a random hyper-plane affected? (No calculation required - just describe what happens.)

Pairwise distances increase with  $d$  and concentrate sharply. Distances between points and a hyperplane stay unaffected by the increased dimensionality.

4. (8 pts) In class, we have learned that the  $k$ -NN algorithm is distance-based.

Suppose we are provided with the following 2D dataset:

- Class +1 (red):  $\{(2, 6)\}$
- Class -1 (green):  $\{(2, 2), (4, 4)\}$

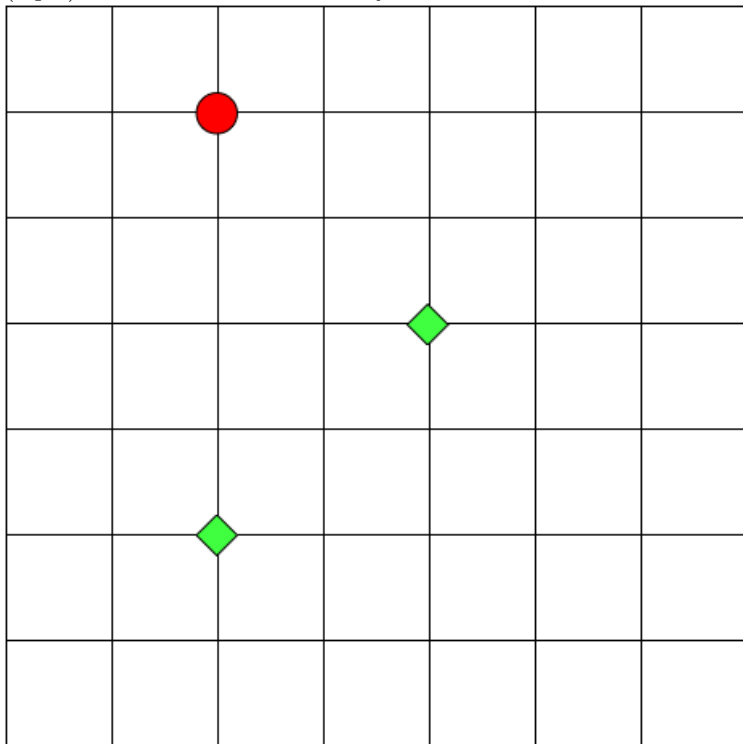
In this problem, we will study the difference between the  $l_2$  and  $l_1$  distances. For two points  $\mathbf{x} = (x_1, x_2)$  and  $\mathbf{z} = (z_1, z_2)$ , the  $l_1$  distance is defined as

$$d_1(\mathbf{x}, \mathbf{z}) = |x_1 - z_1| + |x_2 - z_2|, \quad (1)$$

and the  $l_2$  distance is defined as

$$d_2(\mathbf{x}, \mathbf{z}) = \sqrt{(x_1 - z_1)^2 + (x_2 - z_2)^2}. \quad (2)$$

- (a) (4 pts) Draw the decision boundary for the 1-NN classifier with  $l_2$  distance.



- (b) (4 pts) Show that when  $x_1 > 4$  and  $x_2 > 6$ , 1-NN classifier can't classify the point  $\mathbf{x}^* = (x_1, x_2)$  with  $l_1$  distance. (Hint: show that the closest distances from  $\mathbf{x}^*$  to two classes' points are the same.)

$$d_1(\mathbf{x}^*, (2, 6)) = |x_1 - 2| + |x_2 - 6| = x_1 - 2 + x_2 - 6 = x_1 + x_2 - 8,$$

$$d_1(\mathbf{x}^*, (4, 4)) = |x_1 - 4| + |x_2 - 4| = x_1 - 4 + x_2 - 4 = x_1 + x_2 - 8,$$

$$d_1(\mathbf{x}^*, (2, 2)) = |x_1 - 2| + |x_2 - 2| = x_1 - 2 + x_2 - 2 = x_1 + x_2 - 4.$$

$$d_1(\mathbf{x}^*, (2, 6)) = d_1(\mathbf{x}^*, (4, 4)) < d_1(\mathbf{x}^*, (2, 2)), \text{ which means } \arg \min_{(\mathbf{x}, y) \in \mathcal{D}} d_1(\mathbf{x}^*, \mathbf{x}) = \{((2, 6), +1), ((4, 4), -1)\}.$$

Thus 1-NN can't classify the point  $\mathbf{x}^* = (x_1, x_2)$  with  $l_1$  distance.

### 3 [17] MLE and MAP

- (5 pts) Recall the coin example in the class. A natural assumption about a coin toss is that the distribution of the observed outcomes is a binomial distribution. If a coin was tossed  $n = n_H + n_T$  times and its probability of coming up heads is  $\theta$ , the probability that we would observe exactly  $n_H$  heads and  $n_T$  tails is

$$P(\mathcal{D}|\theta) = \binom{n_H}{n_H + n_T} \theta^{n_H} \cdot (1 - \theta)^{n_T}.$$

In this model,  $\hat{\theta}_{MLE} = \frac{n_H}{n_H + n_T}$ . Furthermore, if  $\theta$  has a prior distribution

$$P(\theta) = \frac{\theta^{\alpha-1}(1-\theta)^{\beta-1}}{B(\alpha, \beta)},$$

where  $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$  is the normalization constant,  $\hat{\theta}_{MAP} = \frac{n_H + \alpha - 1}{n_H + n_T + \beta + \alpha - 2}$ . Please answer following questions:

- (1 pts) What will happen if the prior of  $\theta$  is very wrong and sample number  $n$  is very small?  
MAP can be very wrong.

- (2 pts) Explain why as  $n \rightarrow \infty$ ,  $\hat{\theta}_{MAP} \rightarrow \hat{\theta}_{MLE}$ . (No formal proof required.)  
 $\alpha - 1$  and  $\beta - 1$  become irrelevant compared to very large  $n_H, n_T$ .

- (2 pts) What will happen if the prior of  $\theta$  is very wrong but sample number  $n$  is very large?  
MAP can be as good as MLE without the bad influence from prior.

- (12 pts) Let  $x_1, \dots, x_n$  be iid random data sampled by the poisson distribution

$$P(x|\theta) = e^{-\theta} \frac{\theta^x}{x!} (x \in \mathbb{N} \cup \{0\}). \quad (3)$$

Here  $\theta > 0$  is a hyperparameter of this distribution.

- (3 pts) Write the log likelihood function  $\log P(x_1, x_2, \dots, x_n|\theta)$ .  

$$P(x_1, x_2, \dots, x_n|\theta) = \prod_{i=1}^n P(x_i|\theta) = \prod_{i=1}^n \left( e^{-\theta} \frac{\theta^{x_i}}{x_i!} \right) = e^{-n\theta} \theta^{\sum_{i=1}^n x_i} \prod_{i=1}^n \frac{1}{x_i!}$$

$$\log P(x_1, x_2, \dots, x_n|\theta) = \left( \sum_{i=1}^n x_i \right) \log \theta - n\theta + \sum_{i=1}^n \log \frac{1}{x_i!}$$

- (4 pts) Show that the maximum likelihood estimation  $\hat{\theta}_{MLE} = \frac{\sum_{i=1}^n x_i}{n}$ .

$$\frac{\delta \log P(x_1, x_2, \dots, x_n|\theta)}{\delta \theta} = \frac{\sum_{i=1}^n x_i}{\theta} - n$$



Because  $(\sum_{i=1}^n x_i) > 0$ ,  $\frac{\delta P(x_1, x_2, \dots, x_n | \theta)}{\delta \theta} \geq 0$  if and only if  $\theta \leq \frac{\sum_{i=1}^n x_i}{n}$ .  
Thus  $\hat{\theta}_{MLE} = \arg \max_{\theta} \log P(x_1, x_2, \dots, x_n | \theta) = \frac{\sum_{i=1}^n x_i}{n}$ .

(c) (5 pts) If we have a prior distribution for  $\theta$

$$P(\theta) = e^{-\theta} (\theta > 0) \quad (4)$$

Compute the  $\hat{\theta}_{MAP} = \arg \max_{\theta} \log P(\theta | x_1, x_2, \dots, x_n)$ . Will  $\hat{\theta}_{MAP} \leq \hat{\theta}_{MLE}$  always hold?

$$\begin{aligned} \arg \max_{\theta} \log P(\theta | x_1, x_2, \dots, x_n) &= \arg \max_{\theta} (\log P(x_1, x_2, \dots, x_n | \theta) + \log P(\theta)) \\ &= \arg \max_{\theta} \left( \left( \sum_{i=1}^n x_i \right) \log \theta - n\theta + \sum_{i=1}^n \log \frac{1}{x_i!} - \theta \right) \\ &= \arg \max_{\theta} \left( \left( \sum_{i=1}^n x_i \right) \log \theta - (n+1)\theta + \sum_{i=1}^n \log \frac{1}{x_i!} \right) \end{aligned}$$

By the similar argument as (b), we can have  $\hat{\theta}_{MAP} = \frac{\sum_{i=1}^n x_i}{n+1}$  (4 pts)  
 $\hat{\theta}_{MAP} \leq \frac{\sum_{i=1}^n x_i}{n} = \hat{\theta}_{MLE}$ . (1 pts)

## 4 [18] Naive Bayes

- (2 pts) Write the assumption of Naive Bayes about data.

$$P(\mathbf{x}|y) = \prod_{\alpha=1}^d P(x_{\alpha}|y),$$

where  $x_{\alpha} = [\mathbf{x}]_{\alpha}$  is the value for feature  $\alpha$ , i.e., feature values are independent given the label.

- (16 pts) Ronnie is playing a game named Flippin' Extravaganza, he asks if someone can help him win. Can you help him beat the house? This game is easy to undertake, there is a red hat and a blue hat with a weighted penny respectively. The operator secretly flips one of them and asks you to guess under which hat it is. we made the assumption that you somehow know the coins' weights. In fact, in the Coin Flippin' Extravaganza, the operator never reveals the weights. Instead, you've spent the whole day watching people play the game, recording the results below:

game	penny	nickel	dime	hat
1	T	H	T	Red
2	T	T	H	Blue
3	T	H	T	Blue
4	H	H	H	Red
5	H	H	T	Red
6	T	T	H	Blue
7	H	H	T	Red
8	T	T	H	Blue
9	T	H	H	Blue
10	H	H	H	Red
11	T	T	H	Blue
12	T	H	H	Red
13	H	H	T	Red
14	T	T	H	Blue
15	T	H	H	Blue
16	T	T	H	Blue
17	H	T	H	Red
18	H	T	H	Blue

- (4 pts) This model could be formalized by a Naive Bayes with Bernoulli distributed features. To see this, define the feature space  $\mathcal{X}$  and the label space  $\mathcal{Y}$ . Is the Naive Bayes assumption valid for this problem? Why?

Feature space  $\mathcal{X} \in \{H, T\}^3$ . Label space  $\mathcal{Y} \in \{Red, Blue\}$ . The Naive Bayes assumption is valid because any two coin flips are independent of each other given a label.

- (2 pts) Compute  $P(hat = Red)$  and  $P(hat = Blue)$   
 $P(hat = Red) = \frac{4}{9}$ ,  $P(hat = Blue) = \frac{5}{9}$

- (c) (6 pts) Estimate the following probabilities with +1 smoothing

hat	$P(\text{penny} = H \text{hat})$	$P(\text{nickel} = T \text{hat})$	$P(\text{dime} = T \text{hat})$
Red			
Blue			

Estimation with +1 additive smoothing with 2 categories is as follows:

$$[\hat{\theta}_{jc}]_{\alpha} = \frac{\sum_{i=1}^{18} I(y_i = c)I(x_{ia} = j) + 1}{\sum_{i=1}^{18} I(y_i = c) + 2}$$

Now recalculating the probabilities we get:

hat	$P(\text{penny} = H \text{hat})$	$P(\text{nickel} = T \text{hat})$	$P(\text{dime} = T \text{hat})$
Red	$\frac{7}{10}$	$\frac{2}{10}$	$\frac{5}{10}$
Blue	$\frac{2}{12}$	$\frac{8}{12}$	$\frac{2}{12}$

- (d) (4 pts) if the coins come up  $[H, T, T]$ , compute the probability that they came from the blue hat, i.e.  $P(H, T, T|\text{hat} = \text{Blue})$  by Bayes Rule. Please write the computation formula.

$$\begin{aligned}
 P(y = B|\mathbf{x} = [H, T, T]) &= \frac{P(\mathbf{x} = [H, T, T]|y = B)P(y = B)}{P(\mathbf{x} = [H, T, T])} \\
 &= \frac{P(\mathbf{x} = [H, T, T]|y = B)P(y = B)}{P(\mathbf{x} = [H, T, T]|y = R)P(y = R) + P(\mathbf{x} = [H, T, T]|y = B)P(y = B)} \\
 &= \frac{\prod_{\alpha=1}^k [\hat{\theta}_{jB}]_{\alpha} P(y = B)}{\prod_{\alpha=1}^k [\hat{\theta}_{jR}]_{\alpha} P(y = R) + \prod_{\alpha=1}^k [\hat{\theta}_{jB}]_{\alpha} P(y = B)}, \text{ where } j \in \mathbf{x} = [H, T, T].
 \end{aligned}$$

Then we could use the result in (b) and (c) to compute

$$P(y = R|\mathbf{x} = [H, T, T]) = \frac{\frac{2}{12} \times \frac{8}{12} \times \frac{2}{12} \times \frac{5}{9}}{\frac{7}{10} \times \frac{2}{10} \times \frac{5}{10} \times \frac{4}{9} + \frac{2}{12} \times \frac{8}{12} \times \frac{2}{12} \times \frac{5}{9}} = \frac{125}{503}.$$

## 5 [14] Gradient Descent

You wish to use gradient descent to minimize

$$l(w) = (w - 2)^2 \quad (5)$$

with learning rate  $\alpha > 0$ . Solving these following problems will help you choose a appropriate  $\alpha$  for this particular loss function.

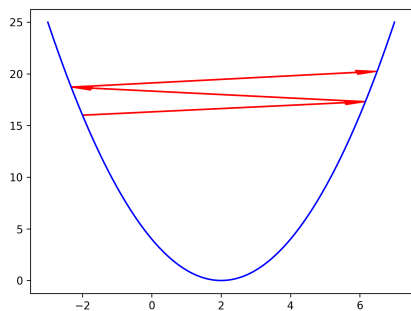
- (3 pts) Suppose at time  $t$ , we have  $w_t$ . Write the  $l'(w)$  and update formula for  $w_{t+1}$  using Gradient Descent.

$$l'(w) = 2(w - 2). \text{ (1 pts)}$$

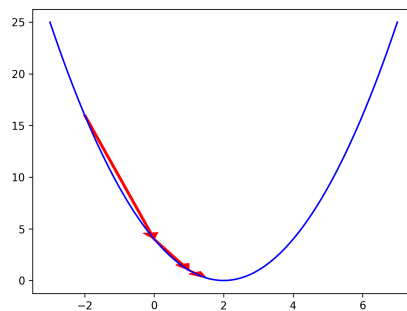
$$w_{t+1} = w_t - \alpha \cdot 2(w_t - 2) \text{ (2 pts)}$$

- (4 pts) Starting with  $w_0 = -2$ , we use gradient descent to update  $w$ . The following three figures are the first 3 updates for different learning rate  $\alpha = 0.25, 0.75, 1.02$  and  $\frac{1}{t+1}$  (when we update  $w_t$  to  $w_{t+1}$ ). Please match different figures to their corresponding learning rates.

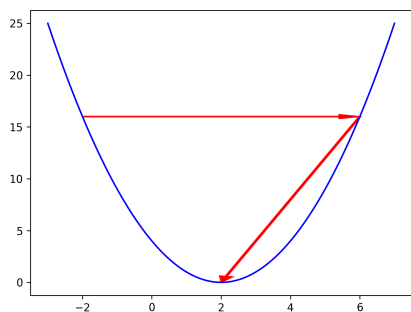
a:  $\alpha = 1.02$ ; b:  $\alpha = 0.25$ ; c:  $\frac{1}{t+1}$ ; d:  $0.75$ .



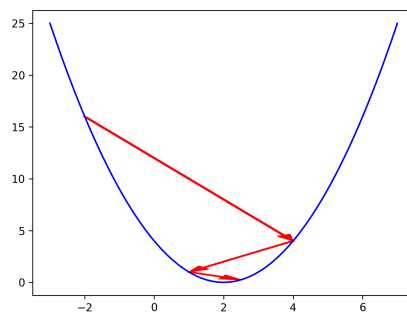
(a)



(b)



(c)



(d)

- (4 pts) Show that  $\forall t, l(w_{t+1}) < l(w_t)$  if and only if  $0 < \alpha < 1$ .

$$w_{t+1} = w_t - \alpha \cdot 2(w_t - 2) \Rightarrow w_{t+1} - 2 = (1 - 2\alpha)(w_t - 2) \Rightarrow l(w_{t+1}) = (1 - 2\alpha)^2 l(w_t).$$

Thus  $l(w_{t+1}) < l(w_t) \Leftrightarrow (1 - 2\alpha)^2 < 1 \Leftrightarrow 0 < \alpha < 1$ .

4. (3 pts) Show that Newton's method will help  $l(w)$  converge to minimum in one step update with arbitrary  $w_0$ . Hessian Matrix here is a scalar 2. Then  $w_1 = w_0 - 2^{-1}l'(w_0) = w_0 - (w_0 - 2) = 2$ .  $l(w_1) = 0$  achieves the minimum.

## 6 [13] Linear Classifiers

1. (10 pts) Consider the following 2D dataset  $\mathcal{D}$ :

- Class +1:  $\{(1, 3), (3, 3)\}$
- Class -1:  $\{(5, 1)\}$

(a) (4 pts) To find a  $\mathbf{w}$  and  $b$  s.t.  $\forall (\mathbf{x}, y) \in \mathcal{D}, y(\mathbf{w}^\top \mathbf{x} + b) > 0$ , please first transform  $\mathcal{D}$  to  $\mathcal{D}' = \{([\mathbf{x}, 1], y) | (\mathbf{x}, y) \in \mathcal{D}\}$  and consider a new  $\mathbf{w}' = [\mathbf{w}, b]$ . Next, apply perception algorithm in this new dataset. Write down the sequence of each updates in the perception algorithm  $[\mathbf{w}'_0, \mathbf{w}'_1, \dots]$  starting with  $\mathbf{w}'_0 = (0, 0, 0)$ . Notice that different sequence of points will lead to different sequence of updates, please consider the points in this fixed order  $[(3, 3), (1, 3), (5, 0)]$  in each iteration of perception.

After append additional 1 for each  $\mathbf{x}$ , the dataset becomes:

- Class +1:  $\{(1, 3, 1), (3, 3, 1)\}$
- Class -1:  $\{(5, 1, 1)\}$

$(0, 0, 0), (3, 3, 1), (-2, -2, 0), (1, 5, 1), (-4, 4, 0), (-1, 7, 1), (-6, 6, 0), (-3, 9, 1)$

(b) (2 pts) If one more point  $(2, 3)$  labeled  $-1$  is added to  $\mathcal{D}$ , what will happen to the perception algorithm?  
This dataset can't be separate after adding  $(2, 3)$ . So the perception algorithm will not converge finally.

(c) (4 pts) In the lecture we learned that SVM will help us find the maximum margin separating hyperplane. Let's guess what this hyperplane is without the proof of it's "maximum". First, let's delete the point  $(1, 3)$  from our dataset. What is the maximum margin separating hyperplane for the remaining two points dataset? After adding the point  $(1, 3)$  back to the dataset, will the maximum margin separating hyperplane change? If we have one more point  $(3, 2)$  labeled  $+1$  to the dataset, will the maximum margin separating hyperplane change?

One example of the maximum margin separating hyperplane is  $x_1 - x_2 - 2 = 0$ , i.e.  $\mathbf{w} = [1, -1]$  and  $b = -2$ . The correct solutions could be  $\alpha \mathbf{w}$  and  $\alpha b \forall \alpha \in \mathcal{R}$ . (2 pts)

After adding the point  $(1, 3)$  back to the dataset, the maximum margin separating hyperplane will keep the same. (1 pts)

If we have one more point  $(3, 2)$  labeled  $+1$  to the dataset, the maximum margin separating hyperplane will change. (1 pts)

2. (2 pts) Write two commonly used binary classification loss functions  $l(h_{\mathbf{w}}(\mathbf{x}_i, y_i))$  (For example zero-one loss or exponential loss).

Hinge-Loss:  $\max[1 - h_{\mathbf{w}}(\mathbf{x}_i)y_i, 0]^p$

Log-Loss:  $\log \left( 1 + e^{-h_{\mathbf{w}}(\mathbf{x}_i)y_i} \right)$

Exponential Loss:  $e^{-h_{\mathbf{w}}(\mathbf{x}_i)y_i}$

Zero-One Loss :  $\delta(\text{sign}(h_{\mathbf{w}}(\mathbf{x}_i)) \neq y_i)$

3. (1 pts) Write one commonly used regression loss function  $l(h_w(\mathbf{x}_i, y_i))$  (For example squared loss or absolute loss).

Squared Loss:  $(h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$

Absolute Loss:  $|h_{\mathbf{w}}(\mathbf{x}_i) - y_i|$

True/False	
kNN	
MLE and MAP	
NB	
Gradient Descent	
Linear Classifiers	
<b>TOTAL</b>	