

CSCI 567 Fall 2018 Midterm Exam
DO NOT OPEN EXAM UNTIL INSTRUCTED TO DO SO
PLEASE TURN OFF ALL CELL PHONES

Problem	1	2	3	4	5	6	Total
Max	16	10	16	42	24	12	120
Points							

Please read the following instructions carefully:

- The exam has a total of **22 pages** (double-sided, including this cover and two blank pages in the end) and **6 questions**. Each question have sub-questions. Once you are permitted to open your exam (and not before), you should check and make sure that you are not missing any pages.
- This is a **closed-book/notes** exam. Consulting any resources is NOT permitted.
- Any kind of cheating will lead to **score 0** for the entire exam and be reported to SJACS.
- Duration of the exam is 2 hours. Questions are not ordered by their difficulty. Budget your time on each question carefully.
- Answers should be **concise** and written down **legibly**. All questions can be done within 5-10 lines.
- You must answer each question on the page provided. You can use the last two blank pages as scratch paper. Raise your hand to ask a proctor for more if needed.
- Select **one and only one answer** for all multiple choice questions.
- You **may not** leave your seat **for any reason** unless you submit your exam at that point.

1 Machine Learning Concepts (16 points)

1.1 Multiple Choice (12 points)

(a) Which of the following is a parametric model?

- (A) Nearest neighbor classifier
- (B) Support vector machines with a Gaussian kernel
- (C) Linear regression with a polynomial kernel
- (D) fully connected feedforward neural nets

Ans: D.

(2 points)

(b) Which of the following phenomenon is called overfitting?

- (A) low training error, low test error
- (B) low training error, high test error
- (C) high training error, low test error
- (D) high training error, high test error

Ans: B.

(2 points)

(c) Suppose the conditional probability of seeing label k given \mathbf{x} is $\mathbb{P}(y = k \mid \mathbf{x})$. Which of the following is the Bayes optimal classifier?

- (A) A randomized classifier that predicts label k with probability $\mathbb{P}(y = k \mid \mathbf{x})$ given \mathbf{x} .
- (B) A deterministic classifier that predicts $\operatorname{argmax}_k \mathbb{P}(y = k \mid \mathbf{x})$ given \mathbf{x} .
- (C) A deterministic classifier that predicts $\operatorname{argmin}_k \mathbb{P}(y = k \mid \mathbf{x})$ given \mathbf{x} .
- (D) Both (A) and (B).

Ans: B. See Lec 1.

(2 points)

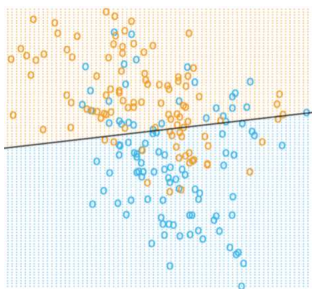
(d) Which of the following cannot be used as regularization to control model complexity?

- (A) $R(\mathbf{w}) = \sum_{d=1}^D |w_d|$
- (B) $R(\mathbf{w}) = \sum_{d=1}^D |w_d|^3$
- (C) $R(\mathbf{w}) = \sum_{d=1}^D w_d^3$
- (D) $R(\mathbf{w}) = \sum_{d=1}^D |w_d|^4$

C. R is unbounded both above and below. As a result, minimizing or maximizing it will push the objective function to $-\infty$ or ∞ .

(2 points)

(e) Which of the these classifiers could have generated this decision boundary?



- (A) SVM
- (B) 1-nearest-neighbor with L2 distance
- (C) Logistic regression
- (D) A & C

Ans: D. Both are linear classifiers. (2 points)

- (f) For a fixed multiclass problem, which of the following multiclass-to-binary reductions has the smallest testing time complexity?

- (A) One-versus-all
- (B) One-versus-one
- (C) Tree reduction
- (D) Both (A) and (C)

Ans: C. (2 points)

1.2 Multiclass reduction (4 points)

Show that one-versus-all can be seen as a special case of error-correcting-output-code (ECOC). Specifically, write down the code matrix \mathbf{M} for ECOC for a problem with C labels so that executing ECOC is the same as doing one-versus-all. (Note: the entry of \mathbf{M} should be either -1 or $+1$.)

\mathbf{M} is a $C \times C$ matrix with $+1$ on the diagonal and -1 on every other entries, or a $C \times C$ matrix with -1 on the diagonal and $+1$ on every other entries. Either one will get full credit.

Only getting the dimension right will get 1 point.

2 Nearest Neighbor Classification (10 points)

2.1 First Problem (4 points)

We mentioned that the Euclidean/L2 distance is often used as the *default* distance for nearest neighbor classification. It is defined as

$$E(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2^2 = \sum_{d=1}^D (x_d - x'_d)^2 \quad (1)$$

In some applications such as information retrieval, the cosine distance is widely used too. It is defined as

$$C(\mathbf{x}, \mathbf{x}') = 1 - \frac{\mathbf{x}^T \mathbf{x}'}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2} = 1 - \frac{\sum_{d=1}^D (x_d \cdot x'_d)}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2}, \quad (2)$$

where the L2 norm of \mathbf{x} is defined as

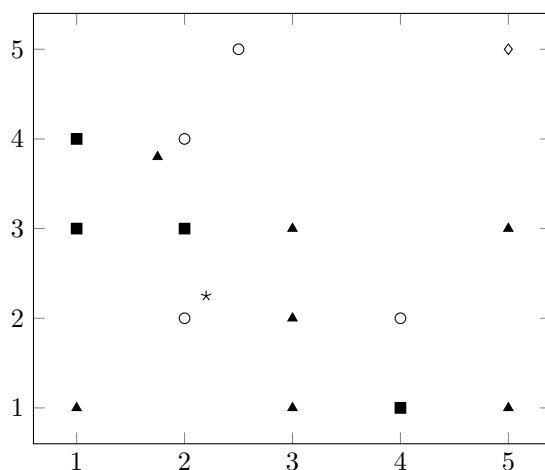
$$\|\mathbf{x}\|_2 = \sqrt{\sum_{d=1}^D x_d^2}. \quad (3)$$

Show that, if data is normalized with unit L2 norm, that is, $\|\mathbf{x}\| = 1$ for all \mathbf{x} in the training and test sets, changing the distance function from the Euclidean distance to the cosine distance will NOT affect the nearest neighbor classification results.

When data is normalized, we have $E(\mathbf{x}, \mathbf{x}') = \sum_{d=1}^D (x_d - x'_d)^2 = \sum_{d=1}^D (x_d^2 + x'^2_d - 2x_d x'_d) = 2(1 - \sum_{d=1}^D (x_d \cdot x'_d)) = 2C(\mathbf{x}, \mathbf{x}')$. Scaling any distance measurement by a factor of two clearly will not affect the results of NNC. (4 points)

2.2 Second Problem (6 points)

For the data given below, squares, triangles, and open circles are three different classes of data in the training set and the diamond (\diamond) and star (*) are test points. We denote the total number of training points as N and consider K-nearest-neighbor (KNN) classifier with L2 distance.



(a) When $K = 1$, how many *circles*(o) will be misclassified (as a validation point) when one performs leave-one-out validation (that is, N -fold cross validation)? (A single number as the answer is enough.)

3

(2 points)

(b) What is the minimum value of K for which the *star*(*) will be classified as a *triangle*? (A single number as the answer is enough.)

4

(2 points)

(c) What is the *diamond* classified as for $K = N$? Explain why.

Triangle

(1 point)

When $K = N$ the prediction is always the majority label of the training set. Triangle is the majority in this example.

(1 point)

3 Linear Regression (16 points)

3.1 First Problem (4 points)

For a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \in \mathbb{R}^D \times \mathbb{R}$, let \mathbf{w}_* be the least square solution with no regularization (assume $\mathbf{X}^T \mathbf{X}$ is invertible where $\mathbf{X} \in \mathbb{R}^{N \times D}$ is the data matrix with each row corresponding to the feature of an example, as used in the class). Find the least square solution (with no regularization again) of the following new training set: $(\mathbf{x}_1, y_1 - \mathbf{w}_*^T \mathbf{x}_1), \dots, (\mathbf{x}_N, y_N - \mathbf{w}_*^T \mathbf{x}_N)$.

Using the formula derived in the class, we know the new least square solution \mathbf{w}'_* is

$$\mathbf{w}'_* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{X} \mathbf{w}_*) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) = \mathbf{0}.$$

Rubrics:

- Write down the correct formula for \mathbf{w}_* . (1 point)
- Use the correct formula to find \mathbf{w}'_* or re-derive in a correct way (e.g. correct gradient). (2 points)
- Get the correct answer finally $\mathbf{w}'_* = \mathbf{0}$. (1 point)

3.2 Second Problem (12 points)

Assume we have a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \in \mathbb{R}^D \times \mathbb{R}$, where each outcome y_n is generated by a probabilistic model $\mathbf{w}_*^T \mathbf{x}_n + \epsilon_n$ with ϵ_n being an independent Gaussian noise with zero-mean and variance σ^2 for some $\sigma > 0$. In other words, the probability of seeing any outcome $y \in \mathbb{R}$ given \mathbf{x}_n is

$$\mathbb{P}(y \mid \mathbf{x}_n; \mathbf{w}_*, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(y - \mathbf{w}_*^T \mathbf{x}_n)^2}{2\sigma^2}\right).$$

- (a) Assume σ is fixed and given, find the maximum likelihood estimation for \mathbf{w}_* . In other words, first write down the probability of seeing the outcomes y_1, \dots, y_N given $\mathbf{x}_1, \dots, \mathbf{x}_N$ as a function of the value of \mathbf{w}_* ; then find the value of \mathbf{w}_* that maximizes this probability. You can assume $\mathbf{X}^T \mathbf{X}$ is invertible where \mathbf{X} is the data matrix as used in Problem 3.1.

The probability of seeing the outcomes y_1, \dots, y_N given $\mathbf{x}_1, \dots, \mathbf{x}_N$ for a linear model \mathbf{w} is

$$\mathcal{P}(\mathbf{w}) = \prod_{i=1}^N \mathbb{P}(y_i \mid \mathbf{x}_i; \mathbf{w}, \sigma) = \prod_{i=1}^N \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2}\right).$$

Taking the negative log, this becomes

$$F(\mathbf{w}) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2 = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2.$$

Maximizing \mathcal{P} is the same as minimizing F , which is clearly the same as just minimizing $\sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2$, the same objective as for least square regression. Therefore the MLE for \mathbf{w}_* is exactly the same as the least square solution:

$$\mathbf{w}_* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Rubrics:

- Write down the correct likelihood function. (2 points)
- Any derivation that reveals this is the same as least square regression. (2 points)
- Arrive at the correct final answer. (2 points)

- (b) Now consider σ as a parameter of the probabilistic model too, that is, the model is specified by both \mathbf{w}_* and σ . Find the maximum likelihood estimation for \mathbf{w}_* and σ .

From previous calculation, the MLE for \mathbf{w}_* and σ is the minimizer of the function

$$F(\mathbf{w}, \sigma) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2.$$

We first fix σ and minimize over \mathbf{w} , which leads to the same MLE for \mathbf{w}_* :

$$\mathbf{w}_* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Next we minimize $F(\mathbf{w}_*, \sigma)$ as function of σ by setting the derivative w.r.t. σ to be 0:

$$\frac{\partial F(\mathbf{w}_*, \sigma)}{\partial \sigma} = \frac{N}{\sigma} - \frac{1}{\sigma^3} \|\mathbf{X}\mathbf{w}_* - \mathbf{y}\|_2^2 = 0.$$

Solving for σ gives the MLE estimate:

$$\sigma = \frac{1}{\sqrt{N}} \|\mathbf{X}\mathbf{w}_* - \mathbf{y}\|_2 = \frac{1}{\sqrt{N}} \|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - \mathbf{y}\|_2.$$

Rubrics:

- Write down the correct likelihood as a function of both \mathbf{w} and σ . (1 point)
- Arrive at the correct solution for \mathbf{w}_* . (2 points)
- Correct derivative with respect to σ . (2 points)
- Arrive at the correct solution for σ . (1 point)

4 Linear Classifiers (42 points)

4.1 Multiple Choice (6 points)

- (a) Which of the following is true?
- (A) When the data is linearly separable, logistic loss (without regularization) does not admit a minimizer.
 - (B) When the data is linearly separable, perceptron converges to a max-margin classifier.
 - (C) Normalizing the output \mathbf{w} of the perceptron algorithm so that $\|\mathbf{w}\|_2 = 1$ changes its test error.
 - (D) Normalizing the output \mathbf{w} of the perceptron algorithm so that $\|\mathbf{w}\|_1 = 1$ changes its test error.

Ans: A. Take the binary case as an example, logistic loss is $\sum_{n=1} \ln(1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n))$. When data is separable, one can find \mathbf{w} such that $y_n \mathbf{w}^T \mathbf{x}_n \geq 0$ for all n . Scaling this \mathbf{w} up will always lead to smaller loss and thus the function does not admit a minimizer. (2 points)

- (b) The perceptron algorithm makes an update $\mathbf{w}' \leftarrow \mathbf{w} + \eta y_n \mathbf{x}_n$ with $\eta = 1$ when \mathbf{w} misclassifies \mathbf{x}_n . Using which of the following different values for η will make sure \mathbf{w}' classifies \mathbf{x}_n correctly?
- (A) $\eta > \frac{y(\mathbf{w}^T \mathbf{x}_n)}{\|\mathbf{x}_n\|_2^2}$ (B) $\eta < \frac{-y(\mathbf{w}^T \mathbf{x}_n)}{\|\mathbf{x}_n\|_2^2 + 1}$ (C) $\eta < \frac{-y(\mathbf{w}^T \mathbf{x}_n)}{\|\mathbf{x}_n\|_2^2}$ (D) $\eta > \frac{-y(\mathbf{w}^T \mathbf{x}_n)}{\|\mathbf{x}_n\|_2^2}$

Ans: D. Solve $y_n \mathbf{w}'^T \mathbf{x}_n = y_n (\mathbf{w} + \eta y_n \mathbf{x}_n)^T \mathbf{x}_n > 0$ for η . (2 points)

- (c) Which of the following about SVM is true?
- (A) Support vectors are training points that are misclassified.
 - (B) Support vectors are training points that are on the learned hyperplane.
 - (C) It is possible that a support vector is on the learned hyperplane.
 - (D) Only misclassified training points could be support vectors, but not all of them are.

Ans: C. (2 points)

4.2 Perceptron (4 points)

The following table shows a binary classification training set and the number of times each point is misclassified during a run of the perceptron algorithm. Write down the final output of the algorithm (that is, the weight vector that represents the hyperplane) assuming we append constant 1 as the first feature to all examples and start with the all-zero weight vector.

\mathbf{x}	y	Times misclassified
(-3, 2)	+1	5
(-1, 1)	-1	5
(5, 2)	+1	3
(2, 2)	-1	4
(1, -2)	+1	3

By the algorithm, the weight is a linear combination of data points weighted by the number of times they are misclassified and their label. That is (2 points)

$$\mathbf{w} = 5 \times (1, -3, 2) - 5 \times (1, -1, 1) + 3 \times (1, 5, 2) - 4 \times (1, 2, 2) + 3 \times (1, 1, -2) = (2, 0, -3). \quad (2 \text{ points})$$

4.3 Kernelized Perceptron (6 points)

Algorithm 1: Perceptron with nonlinear mapping ϕ

```
1 Input: A training set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ 
2 Initialize:  $\mathbf{w} = \mathbf{0}$ 
3 while not converged do
4   randomly pick an example  $(\mathbf{x}_n, y_n)$ , make prediction  $\hat{y} = \text{sign}(\mathbf{w}^T \phi(\mathbf{x}_n))$ 
5   if  $\hat{y} \neq y_n$  then
6      $\mathbf{w} \leftarrow \mathbf{w} + y_n \phi(\mathbf{x}_n)$ 
```

Algorithm 2: Perceptron with kernel function k

```
1 Input: A training set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ 
2 Initialize:  $\alpha_n = 0$  for  $n = 1, \dots, N$  (2 points)
3 while not converged do
4   randomly pick an example  $(\mathbf{x}_n, y_n)$ , make prediction  $\hat{y} = \text{sign}\left(\sum_{m=1}^N \alpha_m k(\mathbf{x}_m, \mathbf{x}_n)\right)$  (2 points)
5   if  $\hat{y} \neq y_n$  then
6      $\alpha_n \leftarrow \alpha_n + y_n$  (2 points)
```

Algorithm 1 is the perceptron algorithm with a nonlinear mapping ϕ . Now given the corresponding kernel function k , please kernelize the algorithm and fill out the missing details in Algorithm 2.

4.4 Multiclass Perceptron (6 points)

Recall that a linear model for a multiclass classification problem with C classes is parameterized by C weight vectors $\mathbf{w}_1, \dots, \mathbf{w}_C \in \mathbb{R}^D$. In the class we derive the multiclass logistic regression by minimizing the multiclass logistic loss. In this problem you need to derive the multiclass perceptron algorithm in a similar way. Specifically, the multiclass perceptron loss on a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \in \mathbb{R}^D \times [C]$ is defined as

$$F(\mathbf{w}_1, \dots, \mathbf{w}_C) = \sum_{n=1}^N \max \left\{ 0, \max_{k \neq y_n} \mathbf{w}_k^T \mathbf{x}_n - \mathbf{w}_{y_n}^T \mathbf{x}_n \right\}.$$

Similarly to the binary case, multiclass perceptron is simply applying SGD with learning rate 1 to minimize the multiclass perceptron loss. Based on the above information, write down the multiclass perceptron algorithm below. (For simplicity, you do not need to worry about the non-differential points of F . In other words, the term $\max_{k \neq y_n} \mathbf{w}_k^T \mathbf{x}_n - \mathbf{w}_{y_n}^T \mathbf{x}_n$ is never 0.)

Solutions:

Algorithm 3: Multiclass Perceptron

```

1 Input: A training set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ 
2 Initialize:  $\mathbf{w}_1 = \dots = \mathbf{w}_C = 0$  (or randomly)
3 while not converged do
4   randomly pick an example  $(\mathbf{x}_n, y_n)$ , make prediction  $\hat{y} = \operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x}_n$  (2 points)
5   if  $\hat{y} \neq y_n$  then
6      $\mathbf{w}_{\hat{y}} \leftarrow \mathbf{w}_{\hat{y}} - \mathbf{x}_n$  (2 points)
7      $\mathbf{w}_{y_n} \leftarrow \mathbf{w}_{y_n} + \mathbf{x}_n$  (2 points)
```

4.5 Perceptron for Non-separable Data (20 points)

Warning: you might find this problem the hardest and want to finish other problems first.

In Homework 1 you are asked to prove a bound on the number of mistakes that perceptron makes for a separable data set. In this problem you need to follow the steps below to prove a bound for the general case with possibly non-separable data. In particular we consider the following variant of perceptron where instead of randomly picking a data point in each iteration, we simply make one pass of the data set. We also assume that the data is normalized so that $\|\mathbf{x}_n\|_2 = 1$ for all $n \in [N]$.

Algorithm 4: One-pass Perceptron

```

1 Input: A training set  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ 
2 Initialize:  $\mathbf{w}_1 = \mathbf{0}$ 
3 for  $n = 1, \dots, N$  do
4   make prediction  $\hat{y} = \text{sign}(\mathbf{w}_n^T \mathbf{x}_n)$ 
5   if  $\hat{y} \neq y_n$  then
6      $\mathbf{w}_{n+1} = \mathbf{w}_n + y_n \mathbf{x}_n$ 

```

(a) Let M be the number of mistakes Algorithm 4 makes. Prove $\|\mathbf{w}_{N+1}\|_2 \leq \sqrt{M}$.

This is exactly the same as the question in Homework 1. By the algorithm we have for any n ,

$$\begin{aligned}
 \|\mathbf{w}_{n+1}\|_2^2 &= \|\mathbf{w}_n + \mathbb{I}[y_n \mathbf{w}_n^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n\|_2^2 \\
 &= \|\mathbf{w}_n\|_2^2 + \mathbb{I}[y_n \mathbf{w}_n^T \mathbf{x}_n \leq 0] (2y_n \mathbf{w}_n^T \mathbf{x}_n + y_n^2 \|\mathbf{x}_n\|_2^2) \\
 &\leq \|\mathbf{w}_n\|_2^2 + \mathbb{I}[y_n \mathbf{w}_n^T \mathbf{x}_n \leq 0]
 \end{aligned}
 \tag{2 points}$$

Applying this inequality recursively we get $\|\mathbf{w}_{N+1}\|_2^2 \leq \sum_{n=1}^N \mathbb{I}[y_n \mathbf{w}_n^T \mathbf{x}_n \leq 0] = M$ and thus $\|\mathbf{w}_{N+1}\|_2 \leq \sqrt{M}$. (1 point)

- (b) Next show that for any \mathbf{w}_* with $\|\mathbf{w}_*\|_2 \leq 1$, we have $\mathbf{w}_{N+1}^T \mathbf{w}_* \geq M - M_*$ where $M_* = \sum_{n=1}^N \max\{0, 1 - y_n \mathbf{w}_*^T \mathbf{x}_n\}$ is the total hinge loss of \mathbf{w}_* .

For any n , we have

$$\begin{aligned} \mathbf{w}_{n+1}^T \mathbf{w}_* &= \mathbf{w}_n^T \mathbf{w}_* + \mathbb{I}[y_n \mathbf{w}_n^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n^T \mathbf{w}_* \\ &= \mathbf{w}_n^T \mathbf{w}_* + \mathbb{I}[y_n \mathbf{w}_n^T \mathbf{x}_n \leq 0] (1 - (1 - y_n \mathbf{x}_n^T \mathbf{w}_*)) \end{aligned} \quad (1 \text{ point})$$

$$\geq \mathbf{w}_n^T \mathbf{w}_* + \mathbb{I}[y_n \mathbf{w}_n^T \mathbf{x}_n \leq 0] (1 - \max\{0, 1 - y_n \mathbf{x}_n^T \mathbf{w}_*\}) \quad (1 \text{ point})$$

$$\geq \mathbf{w}_n^T \mathbf{w}_* + \mathbb{I}[y_n \mathbf{w}_n^T \mathbf{x}_n \leq 0] - \max\{0, 1 - y_n \mathbf{x}_n^T \mathbf{w}_*\} \quad (1 \text{ point})$$

Applying this inequality recursively we get

$$\mathbf{w}_{N+1}^T \mathbf{w}_* \geq \sum_{n=1}^N \mathbb{I}[y_n \mathbf{w}_n^T \mathbf{x}_n \leq 0] - \sum_{n=1}^N \max\{0, 1 - y_n \mathbf{x}_n^T \mathbf{w}_*\} = M - M_*. \quad (1 \text{ point})$$

- (c) Combine the above two results to prove a mistake bound $M \leq 2M_* + 1$. (Hint: as in Homework 1 you need to use the Cauchy-Schwarz inequality $\mathbf{a}^T \mathbf{b} \leq \|\mathbf{a}\|_2 \|\mathbf{b}\|_2$)

Since

$$\mathbf{w}_{N+1}^T \mathbf{w}_* \leq \|\mathbf{w}_{N+1}\|_2 \|\mathbf{w}_*\|_2 \leq \sqrt{M} \quad (1 \text{ point})$$

and

$$\sqrt{M} \leq \frac{1}{2}(M + 1), \quad (1 \text{ point})$$

we have $M - M_* \leq \frac{1}{2}(M + 1)$. Solving for M gives $M \leq 2M_* + 1$. (1 point)

(Directly solving $M - M_* \leq \sqrt{M}$ for M also works too.)

- (d) Assume each data point in the training set S is an i.i.d. sample of a fixed distribution \mathcal{P} , which we denote as $S \sim \mathcal{P}$. Recall that the risk (in terms of 0-1 loss) of a classifier f is defined as $R(f) = \mathbb{E}[\mathbb{I}[f(\mathbf{x}) \neq y]]$ where the expectation is with respect to the random draw of $(\mathbf{x}, y) \sim \mathcal{P}$ and the potential internal randomness of the classifier f . We also define the risk in terms of the hinge loss of a linear classifier \mathbf{w} as $R_{\text{hinge}}(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}}[\max\{0, 1 - y\mathbf{w}^T \mathbf{x}\}]$.

Now consider the following *randomized* classifier f : for any \mathbf{x} , the prediction is $\text{sign}(\mathbf{w}_n^T \mathbf{x})$ where \mathbf{w}_n is drawn uniformly at random from $\mathbf{w}_1, \dots, \mathbf{w}_N$ defined in Algorithm 4. Prove

$$\mathbb{E}_{S \sim \mathcal{P}}[R(f)] \leq 2 \min_{\mathbf{w}: \|\mathbf{w}\|_2 \leq 1} R_{\text{hinge}}(\mathbf{w}) + \frac{1}{N}.$$

The result from the last question can be written as

$$\frac{1}{N} \sum_{n=1}^N \mathbb{I}[\text{sign}(\mathbf{w}_n^T \mathbf{x}_n) \neq y_n] \leq \frac{2}{N} \sum_{n=1}^N \max\{0, 1 - y_n \mathbf{w}^T \mathbf{x}_n\} + \frac{1}{N} \quad (1 \text{ point})$$

for any \mathbf{w} such that $\|\mathbf{w}\|_2 \leq 1$. Taking the expectation on both sides with respect to $S \sim \mathcal{P}$ we arrive at

$$\frac{1}{N} \sum_{n=1}^N \mathbb{E}_{(\mathbf{x}_1, y_1) \sim \mathcal{P}, \dots, (\mathbf{x}_{n-1}, y_{n-1}) \sim \mathcal{P}} \mathbb{E}_{(\mathbf{x}_n, y_n) \sim \mathcal{P}} [\mathbb{I}[\text{sign}(\mathbf{w}_n^T \mathbf{x}_n) \neq y_n]] \leq 2R_{\text{hinge}}(\mathbf{w}) + \frac{1}{N}. \quad (2 \text{ points})$$

Since \mathbf{w}_n is independent of the fresh new example (\mathbf{x}_n, y_n) , we have

$$\mathbb{E}_{(\mathbf{x}_n, y_n) \sim \mathcal{P}} [\mathbb{I}[\text{sign}(\mathbf{w}_n^T \mathbf{x}_n) \neq y_n]] = R(\mathbf{w}_n), \quad (2 \text{ points})$$

the risk of the linear classifier \mathbf{w}_n . Therefore we have

$$\mathbb{E}_{S \sim \mathcal{P}} \left[\frac{1}{N} \sum_{n=1}^N R(\mathbf{w}_n) \right] \leq 2R_{\text{hinge}}(\mathbf{w}) + \frac{1}{N}. \quad (1 \text{ point})$$

By the definition of the randomized classifier f , the left hand side is exactly $\mathbb{E}_{S \sim \mathcal{P}}[R(f)]$, proving the claimed statement since \mathbf{w} is any weight such that $\|\mathbf{w}\|_2 \leq 1$. (1 point)

- (e) Suppose we now make T passes (instead of one pass) of the data to obtain $\mathbf{w}_1, \dots, \mathbf{w}_{TN}$ using the perceptron algorithm. Similar to f defined in the last question, we define a new randomized classifier f' which predicts by randomly picking one of the TN hyperplanes $\mathbf{w}_1, \dots, \mathbf{w}_{TN}$ and then following the prediction of that hyperplane. Does the following hold

$$\mathbb{E}_{S \sim \mathcal{P}}[R(f')] \leq 2 \min_{\mathbf{w}: \|\mathbf{w}\|_2 \leq 1} R_{\text{hinge}}(\mathbf{w}) + \frac{1}{TN}?$$

Prove it if you think so, otherwise briefly explain why this is not true.

This is not true.

(1 point)

Following the previous proof you will find that it breaks if one wants to claim for example

$$\mathbb{E}_{(\mathbf{x}_1, y_1) \sim \mathcal{P}}[\mathbb{I}[\text{sign}(\mathbf{w}_{N+1}^T \mathbf{x}_1) \neq y_1]] = R(\mathbf{w}_{N+1}).$$

This is because \mathbf{w}_{N+1} depends on (\mathbf{x}_1, y_1) .

Another way to see that this should not be true is: if it was true, then by making many passes of a training set of size one, the term $\frac{1}{TN}$ will still go to zero, leading to a very strong risk bound, which is clearly impossible.

Similar explanations along this line will all get 2 points.

5 Neural Networks (24 points)

5.1 Multiple Choice (10 points)

- (a) Overfitting is a major problem for neural networks. Which of the following can help prevent overfitting?
- (A) Retraining on the same data many times.
 - (B) Using a larger learning rate for Backpropagation .
 - (C) Dropping random neurons in each iteration of Backpropagation.
 - (D) Training until you get the smallest training error.

Ans: C.

(2 points)

- (b) Which of the following is true about neural nets?
- (A) A fully connected feedforward neural net without nonlinear activation functions is the same as a linear model.
 - (B) Since we usually apply Backpropagation (i.e. SGD) to learn neural networks, it is a convex problem.
 - (C) A neural net with one hidden layer and a fixed number of neurons can represent any continuous function.
 - (D) A max-pooling layer with a 2×2 filter has 4 parameters to be learned.

Ans: A.

(2 points)

- (c) Suppose a convolution layer takes a 4×6 image with 3 channels as input and outputs a $3 \times 4 \times 6$ volume. Which of the following is a possible configuration of this layer?
- (A) One 2×3 filter with depth 6, stride 1, no zero-padding.
 - (B) Six 2×3 filters with depth 3, stride 1, no zero-padding.
 - (B) Six 3×4 filters with depth 3, stride 2, no zero-padding.
 - (B) Six 3×4 filters with depth 3, stride 1, 1 pixel of zero-padding.

Ans: B.

(2 points)

- (d) How many parameters do we need to learn for the following network structure? An $8 \times 8 \times 3$ image input, followed by a convolution layer with 2 filters of size 2×2 (stride 1, no zero-padding), then another convolution layer with 4 filters of size 3×3 (stride 2, no zero-padding), and finally a max-pooling layer with a 2×2 filter (stride 1, no zero-padding). (Note: the depth of all filters are not explicitly spelled out, and we assume no bias/intercept terms are used.)
- (A) 96
 - (B) 44
 - (C) 100
 - (D) 48

Ans: A. $2 \times (2 \times 2 \times 3) + 4 \times (3 \times 3 \times 2) = 96$.

(2 points)

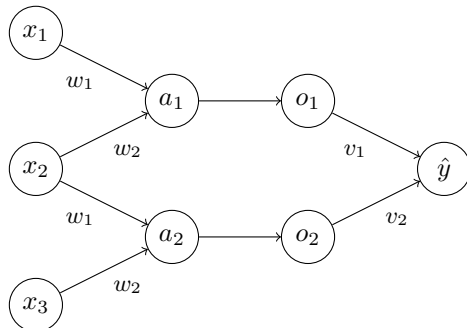
- (e) What is the final output dimension of the last question?
- (A) $2 \times 2 \times 3$
 - (B) $5 \times 5 \times 4$
 - (C) $2 \times 2 \times 1$
 - (D) $2 \times 2 \times 4$

Ans: D.

(2 points)

5.2 Backpropagation for CNN (14 points)

Consider the following mini convolutional neural net, where (x_1, x_2, x_3) is the input, followed by a convolution layer with a filter (w_1, w_2) , a ReLU layer, and a fully connected layer with weight (v_1, v_2) .



More concretely, the computation is specified by

$$\begin{aligned} a_1 &= x_1 w_1 + x_2 w_2 \\ a_2 &= x_2 w_1 + x_3 w_2 \\ o_1 &= \max\{0, a_1\} \\ o_2 &= \max\{0, a_2\} \\ \hat{y} &= o_1 v_1 + o_2 v_2 \end{aligned}$$

For an example $(\mathbf{x}, y) \in \mathbb{R}^3 \times \{-1, +1\}$, we define the logistic loss of the CNN as

$$\ell = \ln(1 + \exp(-y\hat{y}))$$

which is a function of the parameters of the network: w_1, w_2, v_1, v_2 .

- (a) Write down $\frac{\partial \ell}{\partial v_1}$ and $\frac{\partial \ell}{\partial v_2}$. You can use the sigmoid function $\sigma(z) = \frac{1}{1+e^{-z}}$ to simplify your notation.

$$\frac{\partial \ell}{\partial v_1} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v_1} \quad (1 \text{ point})$$

$$= \frac{-ye^{-y\hat{y}}}{1 + e^{-y\hat{y}}} o_1 = -\sigma(-y\hat{y}) y o_1 = (\sigma(y\hat{y}) - 1) y o_1 \quad (1 \text{ point})$$

$$\frac{\partial \ell}{\partial v_2} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v_2} \quad (1 \text{ point})$$

$$= \frac{-ye^{-y\hat{y}}}{1 + e^{-y\hat{y}}} o_2 = -\sigma(-y\hat{y}) y o_2 = (\sigma(y\hat{y}) - 1) y o_2 \quad (1 \text{ point})$$

Either one of the last three expressions is acceptable.

- (b) Write down $\frac{\partial \ell}{\partial w_1}$ and $\frac{\partial \ell}{\partial w_2}$. The derivative of the ReLU function is $H(a) = \mathbb{I}[a > 0]$, which you can use directly in your answer.

$$\frac{\partial \ell}{\partial w_1} = \frac{\partial \ell}{\partial a_1} \frac{\partial a_1}{\partial w_1} + \frac{\partial \ell}{\partial a_2} \frac{\partial a_2}{\partial w_1} \quad (1 \text{ point})$$

$$= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} + \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial a_2} \frac{\partial a_2}{\partial w_1} = (\sigma(y\hat{y}) - 1)y(v_1 H(a_1)x_1 + v_2 H(a_2)x_2) \quad (1 \text{ point})$$

Similarly

$$\frac{\partial \ell}{\partial w_2} = \frac{\partial \ell}{\partial a_1} \frac{\partial a_1}{\partial w_2} + \frac{\partial \ell}{\partial a_2} \frac{\partial a_2}{\partial w_2} \quad (1 \text{ point})$$

$$= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial a_1} \frac{\partial a_1}{\partial w_2} + \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial a_2} \frac{\partial a_2}{\partial w_2} = (\sigma(y\hat{y}) - 1)y(v_1 H(a_1)x_2 + v_2 H(a_2)x_3). \quad (1 \text{ point})$$

Again, other equivalent expressions are acceptable.

- (c) Using the derivations above, fill out the missing details for the Backpropagation algorithm below that is used to train this mini CNN.

Algorithm 5: Backpropagation for the above mini CNN

1	Input: A training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$, learning rate η	
2	Initialize: $w_1 = w_2 = v_1 = v_2 = 0$ (or randomly)	(1 point)
3	while <i>not converged</i> do	
4	randomly pick an example (\mathbf{x}_n, y_n)	
5	Forward propagation: compute	(2 points)
	$a_1 = x_1 w_1 + x_2 w_2, a_2 = x_2 w_1 + x_3 w_2$	
	$o_1 = \max\{0, a_1\}, o_2 = \max\{0, a_2\}, \hat{y} = o_1 v_1 + o_2 v_2$	
6	Backward propagation: update	(3 points)
	$w_1 \leftarrow w_1 - \eta(\sigma(y\hat{y}) - 1)y(v_1 H(a_1)x_1 + v_2 H(a_2)x_2)$	
	$w_2 \leftarrow w_2 - \eta(\sigma(y\hat{y}) - 1)y(v_1 H(a_1)x_2 + v_2 H(a_2)x_3)$	
	$v_1 \leftarrow v_1 - \eta(\sigma(y\hat{y}) - 1)y o_1$	
	$v_2 \leftarrow v_2 - (\sigma(y\hat{y}) - 1)y o_2$	

Deduct 2 points if updating w_1/w_2 with the updated value of v_1/v_2 .

Do not deduct points for using the wrong gradients due to mistakes from previous two questions.

6 Kernel Methods (12 points)

6.1 Multiple Choice (6 points)

(a) Which of the following is true about kernel function?

- (A) If k_1 and k_2 are kernel, then $c_1k_1 + c_2k_2$ is a kernel too for any $c_1, c_2 \in \mathbb{R}$.
- (B) Kernel function must be symmetric, that is, $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$.
- (C) If k is a kernel, then $-k$ is a kernel too.
- (D) If k is a kernel, then $\ln k$ is a kernel too.

Ans: B.

(2 points)

(b) Which of the following is not a kernel function?

- (A) $k(x, x') = (\mathbf{x}^T \mathbf{x}')^2$
- (B) $k(x, x') = (\mathbf{x}^T \mathbf{x}' + 1)^2$
- (C) $k(x, x') = -\|\mathbf{x} - \mathbf{x}'\|_2^2$
- (D) $k(x, x') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2)$

Ans: C.

(2 points)

(c) Vovk's real polynomial kernel $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ is defined as:

$$k(\mathbf{x}, \mathbf{x}') = \frac{1 - (\mathbf{x}^T \mathbf{x}')^p}{1 - (\mathbf{x}^T \mathbf{x}')}$$

where p is a non-negative integer. Which of the following is the corresponding non-linear mapping for this kernel when $D = 2$ and $p = 3$?

- (A) $\phi(\mathbf{x}) = [x_1^2 x_2^2, 2x_1 x_2, x_1, x_2, 1]^T$
- (B) $\phi(\mathbf{x}) = [x_1^2 x_2^2, \sqrt{2}x_1 x_2, x_1, x_2, 1]^T$
- (C) $\phi(\mathbf{x}) = [x_1^2, x_2^2, \sqrt{2}x_1 x_2, x_1, x_2, 1]^T$
- (D) $\phi(\mathbf{x}) = [x_1^2, x_2^2, 2x_1 x_2, x_1, x_2, 1]^T$

Ans: C.

(2 points)

$$\begin{aligned} \frac{1 - (\mathbf{x}^T \mathbf{x}')^3}{1 - (\mathbf{x}^T \mathbf{x}')} &= \frac{(1 - (\mathbf{x}^T \mathbf{x}'))(1 + \mathbf{x}^T \mathbf{x}' + (\mathbf{x}^T \mathbf{x}')^2)}{1 - (\mathbf{x}^T \mathbf{x}')} \\ &= x_1 x'_1 + x_2 x'_2 + (x_1 x'_1 + x_2 x'_2)^2 + 1 \\ &= x_1^2 x_1'^2 + x_2^2 x_2'^2 + 2x_1 x'_1 x_2 x'_2 + x_1 x'_1 + x_2 x'_2 + 1 \end{aligned}$$

Thus, $\phi(\mathbf{x}) = [x_1^2, x_2^2, \sqrt{2}x_1 x_2, x_1, x_2, 1]^T$.

6.2 Kernel Composition (6 points)

Prove that if $k_1, k_2 : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ are both kernel functions, then $k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$ is a kernel function too. Specifically, suppose ϕ_1 and ϕ_2 are the corresponding mappings for k_1 and k_2 respectively. Construct the mapping ϕ that certifies k being a kernel function.

Observe that

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') = \left(\sum_{i=1}^D \phi_1(\mathbf{x})_i \phi_1(\mathbf{x}')_i \right) \left(\sum_{j=1}^D \phi_2(\mathbf{x})_j \phi_2(\mathbf{x}')_j \right) \quad (2 \text{ points})$$

$$= \sum_{i=1}^D \sum_{j=1}^D \phi_1(\mathbf{x})_i \phi_1(\mathbf{x}')_i \phi_2(\mathbf{x})_j \phi_2(\mathbf{x}')_j = \sum_{i=1}^D \sum_{j=1}^D (\phi_1(\mathbf{x})_i \phi_2(\mathbf{x})_j) (\phi_1(\mathbf{x}')_i \phi_2(\mathbf{x}')_j) \quad (2 \text{ points})$$

Therefore, $\phi(\mathbf{x}) = \phi_1(\mathbf{x})\phi_2(\mathbf{x})^T$. Writing it as a D^2 -dimensional vector is acceptable too. (2 points)

This blank page can be used as scratch paper.

This blank page can be used as scratch paper.