

**CSCI 567 Fall 2022 Practice Problems for Quiz 1**  
Instructor: Vatsal Sharan

# 1 Multiple choice questions

(30 points)

**IMPORTANT:** Select ALL options among  $\{A,B,C,D\}$  that you think are correct. You get 0.5 point for selecting each correct option and similarly 0.5 point for not selecting each incorrect option. You get 1 additional point for selecting all four options correctly.

- (1) Which of the following are true statements about machine learning?
- (A) Cross-validation is often used to tune the hyper-parameters of a machine learning algorithm.
  - (B) The goal of a machine learning algorithm is to achieve zero error on a training set.
  - (C) Regularization is a common way to prevent overfitting.
  - (D) Classification and regression are two common tasks in machine learning.

Answer: ACD

- (2) Which of the following phenomenon is called overfitting?
- (A) low training error, low test error
  - (B) high training error, high test error
  - (C) high training error, low test error
  - (D) low training error, high test error

Answer: D

- (3) Which of the following are true statements about linear regression?
- (A) The least square solution has a closed-form formula, even if  $\ell_2$  regularization is applied.
  - (B) The covariance matrix  $\mathbf{X}^T \mathbf{X}$  is not invertible if and only if the number of data points  $n$  is smaller than the dimension  $d$ .
  - (C) When the covariance matrix  $\mathbf{X}^T \mathbf{X}$  is not invertible, the Residual Sum of Squares (RSS) objective has no minimizers.
  - (D) Linear regression with kernels is a non-parametric method.

Answer: AD

- (4) Consider a training set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  and a probabilistic model  $\mathbb{P}(y_i | \mathbf{x}_i; \mathbf{w})$  which specifies for each  $i$  the probability of seeing outcome  $y_i$  given feature  $\mathbf{x}_i$  and parameter  $\mathbf{w}$ . Which of the following is the Maximum Likelihood Estimation (MLE) for  $\mathbf{w}$ ?
- (A)  $\operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i; \mathbf{w})$
  - (B)  $\operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i; \mathbf{w})$
  - (C)  $\operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^n \ln \mathbb{P}(y_i | \mathbf{x}_i; \mathbf{w})$
  - (D)  $\operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n \ln \mathbb{P}(y_i | \mathbf{x}_i; \mathbf{w})$

Answer: BC

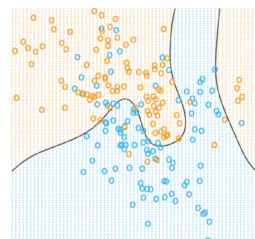
(5) Which of the following are true statements about kernels?

- (A) A machine learning algorithm can be kernelized if it does not need explicit access to the feature vectors and instead only requires access to inner products of the feature vectors.
- (B) A Gram/kernel matrix must be positive semidefinite.
- (C) The product of two kernel functions is still a kernel function.
- (D) The sum of two kernel functions is always a kernel function.

Answer: ABCD

(6) Which of the following ML models could have generated the decision boundary in the binary classification problem on the right?

- (A) SVM without kernels
- (B) 1-nearest-neighbor.
- (C) Logistic regression without kernels
- (D) None of the above.



Answer: D

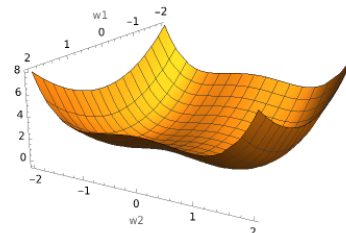
(7) A machine learning objective function is usually of the form  $F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w})$  for some loss function  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$  corresponding to the  $i$ -th training point. We know that  $\nabla f_i(\mathbf{w})$  for  $i$  chosen uniformly at random is a stochastic gradient of this objective since  $\mathbb{E}[\nabla f_i(\mathbf{w})] = \nabla F(\mathbf{w})$ . Which of the following is also a stochastic gradient?

- (A)  $\sum_{i \in S} \nabla f_i(\mathbf{w})$  for a subset  $S \subset \{1, \dots, n\}$  (of some fixed size) chosen uniformly at random.
- (B)  $\frac{1}{|S|} \sum_{i \in S} \nabla f_i(\mathbf{w})$  for a subset  $S \subset \{1, \dots, n\}$  (of some fixed size) chosen uniformly at random.
- (C)  $\nabla f_i(\mathbf{w}) - \nabla f_i(\mathbf{w}_0) + \nabla F(\mathbf{w}_0)$  for  $i$  chosen uniformly at random and some fixed point  $\mathbf{w}_0$ .
- (D)  $\frac{\partial F(\mathbf{w})}{\partial w_i} \cdot d\mathbf{e}_i$  for a coordinate  $i$  chosen uniformly at random ( $\mathbf{e}_i \in \mathbb{R}^d$  is the standard basis vector with 1 in the  $i$ -th coordinate and 0 in all other coordinates).

Answer: BCD

(8) Consider a two-dimensional function  $F(\mathbf{w}) = w_1^2 - w_2^2 + \frac{1}{2}w_2^4$  (a plot is provided below). Which of the following statement is correct?

- (A)  $F$  is convex.
- (B)  $F$  has three stationary points:  $(0, 0)$ ,  $(0, -1)$ , and  $(0, 1)$ .
- (C)  $(0, 0)$  is a saddle point of  $F$ .
- (D) Gradient Descent always converges to  $(0, -1)$ , regardless of initialization.



Answer: BC

(9) Which of the following is true?

- (A) Normalizing the output  $\mathbf{w}$  of the perceptron algorithm so that  $\|\mathbf{w}\|_2 = 1$  changes its test error.
- (B) Normalizing the output  $\mathbf{w}$  of the perceptron algorithm so that  $\|\mathbf{w}\|_1 = 1$  changes its test error.
- (C) Minimizing 0-1 loss is NP-hard for every dataset, even if it is linearly separable.
- (D) Perceptron algorithm may fail to converge on linearly-separable data.

Answer: C

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

(A)  $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}')^2$

(C)  $k(\mathbf{x}, \mathbf{x}') = -\|\mathbf{x} - \mathbf{x}'\|_2^2$

(B)  $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^2$

(D)  $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2)$

Answer: C

## 2 Short answer questions

### 2.1 Linear Classifiers with squared hinge loss (6 points)

In class, we have seen the hinge loss  $\ell(z) = \max\{0, 1 - z\}$ . In this question, we will consider a modification of the hinge loss. Given a binary dataset  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \{-1, 1\}$ , we define the following new loss function for a linear model  $\mathbf{w} \in \mathbb{R}^d$ :

$$F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}), \quad \text{where } f_i(\mathbf{w}) = (\max\{0, 1 - y_i \mathbf{w}^T \mathbf{x}_i\})^2. \quad (1)$$

For a fixed  $i$ , write down the gradient  $\nabla f_i(\mathbf{w})$  (show your derivation), then fill in the missing details in the repeat-loop of the algorithm below which applies SGD to minimize  $F$ .

---

**Algorithm 1:** SGD for minimizing Eq. (1)

---

- 1 **Input:** A training set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \{-1, 1\}$ , learning rate  $\eta > 0$
  - 2 **Initialization:**  $\mathbf{w} = \mathbf{0}$
  - 3 **Repeat:**
  - 4     randomly pick an example  $(\mathbf{x}_i, y_i)$
  - 5     compute  $z = y_i \mathbf{w}^T \mathbf{x}_i$  and update  $\mathbf{w} \leftarrow \mathbf{w} + 2\eta y_i \max\{0, 1 - z\} \mathbf{x}_i$
- 

When  $1 - y_i \mathbf{w}^T \mathbf{x}_i < 0$ , the function is a constant and thus the gradient is  $\mathbf{0}$ . On the other hand, when  $1 - y_i \mathbf{w}^T \mathbf{x}_i > 0$ , the function is simply  $f_i(\mathbf{w}) = (1 - y_i \mathbf{w}^T \mathbf{x}_i)^2$ , and thus by chain rule we have  $\nabla f_i(\mathbf{w}) = -2y_i(1 - y_i \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i$ , which is also  $\mathbf{0}$  when  $y_i \mathbf{w}^T \mathbf{x}_i$  approaches 1. Therefore, we have

$$\nabla f_i(\mathbf{w}) = -2y_i \max\{0, 1 - y_i \mathbf{w}^T \mathbf{x}_i\} \mathbf{x}_i.$$

Rubrics:

- 3 points for the derivation of the gradient. Technically we need some discussion for the case  $y_i \mathbf{w}^T \mathbf{x}_i = 1$  as the solution above shows, but for simplicity it is okay if such discussion is missing as long as the final gradient is correct.
- 3 points for the SGD implementation. Do not deduct more points for using a wrong gradient from earlier wrong derivation.

## 2.2 Nearest Neighbor Classification

(4 points)

In Homework 1, we used the Euclidean distance ( $\ell_2$  distance) and the cosine distance to find the nearest neighbors for the  $k$ -NN algorithm. In this question we will relate the two notions. Let us recall the definitions of Euclidean distance and cosine distance. The Euclidean distance  $E(\mathbf{x}, \mathbf{x}')$  is defined as

$$E(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2 = \sqrt{\sum_{j=1}^d (x_j - x'_j)^2}. \quad (2)$$

The cosine distance  $C(\mathbf{x}, \mathbf{x}')$  is defined as (assuming  $\|\mathbf{x}\| \neq 0$  and  $\|\mathbf{x}'\| \neq 0$ ),

$$C(\mathbf{x}, \mathbf{x}') = 1 - \frac{\mathbf{x}^T \mathbf{x}'}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2} = 1 - \frac{\sum_{j=1}^d (x_j \cdot x'_j)}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2}. \quad (3)$$

Show that if all the datapoints are normalized to have unit  $\ell_2$  norm (that is,  $\|\mathbf{x}\| = 1$  for all  $\mathbf{x}$  in the training and test sets) then changing the distance function from the Euclidean distance to the cosine distance will NOT affect the nearest neighbor classification results.

The goal is to construct a function  $C(\mathbf{x}, \mathbf{x}') = f(E(\mathbf{x}, \mathbf{x}'))$  where  $f$  is monotonic when  $E(\mathbf{x}, \mathbf{x}') \geq 0$ . Once we have a monotonic  $f$ , for any point  $\mathbf{x}$ ,  $E(\mathbf{x}, \mathbf{x}_1) \geq E(\mathbf{x}, \mathbf{x}_2) \iff C(\mathbf{x}, \mathbf{x}_1) \geq C(\mathbf{x}, \mathbf{x}_2)$ ,  $k$ -NN will retrieve the same neighbors with two distance functions.

Given  $\|\mathbf{x}\| = 1$  for all  $\mathbf{x}$ , we have  $C(\mathbf{x}, \mathbf{x}') = 1 - \frac{\mathbf{x}^T \mathbf{x}'}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2} = 1 - \mathbf{x}^T \mathbf{x}'$ . Try to construct a inner product term from  $E(\mathbf{x}, \mathbf{x}')$ .  $E(\mathbf{x}, \mathbf{x}')^2 = \|\mathbf{x}\|_2^2 + \|\mathbf{x}'\|_2^2 - 2\mathbf{x}^T \mathbf{x}' = 2 - 2\mathbf{x}^T \mathbf{x}' = 2C(\mathbf{x}, \mathbf{x}')$ .

Rubrics:

- 3 points for constructing the relationship between two distance functions.
- 1 point for the monotonicity explanation, other reasonable explanation also works.

### 2.3 ML diagnostics: Gradient descent

(5 points)

Let  $F(\mathbf{w}) = w_1^2 + 2w_2^2$ .

(1) Write down the gradient descent update rule for minimizing  $F(\mathbf{w})$ .

$$\frac{\partial F(\mathbf{w})}{\partial \mathbf{w}} = \begin{bmatrix} \partial F(\mathbf{w})/\partial w_1 \\ \partial F(\mathbf{w})/\partial w_2 \end{bmatrix} = \begin{bmatrix} 2w_1 \\ 4w_2 \end{bmatrix}$$
$$\mathbf{w} = \mathbf{w} - \eta \frac{\partial F(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{w} - \eta \begin{bmatrix} 2w_1 \\ 4w_2 \end{bmatrix}, \text{ where } \eta \text{ is the learning rate.}$$

The figure below depicts a run of gradient descent to minimize  $F(\mathbf{w})$  (initializing at  $w_1 = w_2 = 40$ , and running for 20 iterations). The value of  $\mathbf{w}$  after every update is shown as a red circle. The contour lines of  $F(\mathbf{w})$  are also depicted in the figure.

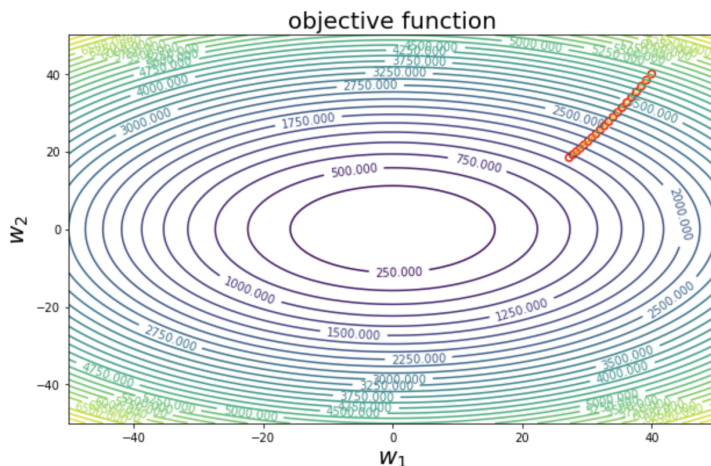


Figure 1: Gradient descent on  $F(\mathbf{w}) = w_1^2 + 2w_2^2$ .

(2a) If we continued to run gradient descent for a sufficiently large number of iterations, to what point would it converge to?

(0, 0)

(2b) Does it seem possible to make gradient descent converge faster than in the figure (with the same initialization and for the function  $F(\mathbf{w})$ )? Explain.

Increase learning rate.

Rubrics:

- 3 points for (1). Correct gradient of  $F(\mathbf{w})$  (2 points), the update function (1 point). Writing two update functions for  $w_1$  and  $w_2$  separately is also fine.

- 1 point for (2a).
- 1 point for (2b).

### 3 Linear Regression

(8 points)

In class, we discussed that if we use Newton's method to solve the least square optimization problem then it only takes one step to converge. We will prove this statement in this problem. Let

$$F(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

where  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ . Recall that Newton's method updates the parameters as follow:

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \mathbf{H}_t^{-1} \nabla F(\mathbf{w}^{(t)})$$

where  $\mathbf{H}_t = \nabla^2 F(\mathbf{w}^{(t)}) \in \mathbb{R}^{d \times d}$  is the Hessian matrix of the objective function evaluated at  $\mathbf{w}^{(t)}$ , i.e. for every index  $u, v \in \{1, \dots, d\}$  the  $(u, v)$ -th entry of  $\mathbf{H}_t$  is  $H_t(u, v) = \frac{\partial^2}{\partial w_u \partial w_v} F(\mathbf{w}) \big|_{\mathbf{w}=\mathbf{w}^{(t)}}$ .

- (1) Find the Hessian of the least square objective function  $F(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$ .

$$\begin{aligned} \frac{\partial}{\partial w_u} F(\mathbf{w}) &= \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i) x_{iu} \\ \frac{\partial^2}{\partial w_u \partial w_v} F(\mathbf{w}) &= \sum_{i=1}^n x_{iu} x_{iv} = (\mathbf{X}^T \mathbf{X})_{uv} \end{aligned}$$

Therefore,  $\mathbf{H} = \nabla^2 F(\mathbf{w}) = \mathbf{X}^T \mathbf{X}$ .

- (2) Show that given any initialization  $\mathbf{w}^{(0)}$ , after one iteration of Newton's method we obtain the optimal  $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ . (Here  $\mathbf{X}$  is the  $n \times d$  data matrix with one row per data point, and  $\mathbf{y}$  is the  $n$ -vector of their labels.)

$$\begin{aligned} \mathbf{w}^{(1)} &= \mathbf{w}^{(0)} - \mathbf{H}^{-1} \nabla L(\mathbf{w}^{(0)}) \\ &= \mathbf{w}^{(0)} - \mathbf{H}^{-1} \mathbf{X}^T (\mathbf{X} \mathbf{w}^{(0)} - \mathbf{y}) \\ &= \mathbf{w}^{(0)} - \mathbf{w}^{(0)} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

Rubrics:

- 4 points for (1).
- 4 points for (2).

Partial credit for incomplete solutions.