

Practice Questions for the ECE 661 Spring 2025 Midterm Exam

Introduction to Machine Learning for Engineers
Prof. Gauri Joshi and Prof. Carlee Joe-Wong

The number of practice questions included in this document may not be representative of the length of the exam. These are solely intended for you to gain more experience in answering questions similar to what you will see on the exam.

1 True or False

Problem 1: The computational complexity of KNN scales linearly with the dimension of the data D but quadratically with the number of samples N

☐ True

☐ False

Solution: False, the computational complexity of KNN is $O(ND)$, so it scales linearly with both

Problem 2: Both KNN and Logistic Regression (without applying non-linear transformations to the input features) will result in a linear decision boundary

☐ True

☐ False

Solution: False, KNN does not necessarily result in a linear decision boundary.

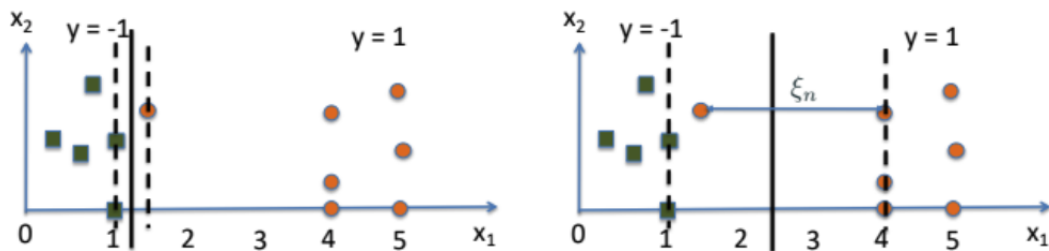
Problem 3: In multi-class classification, taking the One-versus-One approach results in high variance, whereas taking the One-versus-All approach results in high bias.

☐ True

☐ False

Solution: True

Problem 4: Consider the following plot. The graph on the left demonstrates the classifier that would be learned by using soft-margin SVM, while the graph on the right demonstrates the classifier that would be learned using hard-margin SVM.



- ☐ True
- ☐ False

Solution: False, the graph on the left would be hard-margin SVM while the graph on the right would be soft-margin SVM. Soft-margin will allow for the misclassification of outliers

Problem 5: In multinomial logistic regression, our goal is to classify points (\mathbf{x}_n, y_n) for $n = 1, \dots, N$, with each point belonging to one of K classes. We learn K **independent** parameter vectors $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$ such that the probability of the n -th point belonging to the k -th class is equal to

$$\Pr(y = k) = \frac{e^{\mathbf{w}_k^T \mathbf{x}_n}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_n}}.$$

- ☐ True
- ☐ False

Solution: False. The K parameter vectors are not independent since all the probabilities need to sum up to 1. Therefore we only have $K - 1$ parameter vectors to learn.

Problem 6: By using kernel functions to implicitly map input data into a higher-dimensional space, linear algorithms (e.g., linear SVMs) can effectively learn non-linear decision boundaries in the original input space.

- ☐ True
- ☐ False

Solution: True.

2 Multiple Choice

Problem 7: In Naïve Bayes with Laplacian smoothing, if we increase the parameter α (the number of times we add the count of each unique item's occurrences for each label), the bias and the variance of the estimator change as follows:

- ☐ bias increases, variance decreases
- ☐ bias decreases, variance increases
- ☐ bias decreases, variance decreases
- ☐ bias increases, variance increases

Solution: A. bias increases, variance decreases

Problem 8: Answer the following question about the sigmoid function defined as $\sigma(x) = \frac{1}{1+e^{-x}}$

a. Which of the following statements is true -

- ☐ $\sigma(x) = \sigma(-x)$
- ☐ $\sigma(x) = 1 - \sigma(x)$
- ☐ $\sigma(x) = 1 - \sigma(-x)$
- ☐ $\sigma(x) = 1 + \sigma(-x)$

b. Which of the following statements is true -

- ☐ $\frac{d\sigma(x)}{dx} = \sigma(-x)$
- ☐ $\frac{d\sigma(x)}{dx} = (\sigma(x))^2$
- ☐ $\frac{d\sigma(x)}{dx} = (\sigma(-x))^2$
- ☐ $\frac{d\sigma(x)}{dx} = \sigma(x) \cdot \sigma(-x)$

Solution: C, D

Problem 9: Figure 1 shows the plots of loss vs. number of iterations for Batch Gradient Descent (GD), Stochastic Gradient Descent (SGD) and Mini-Batch Stochastic Gradient Descent (mini-batch SGD), which are run on a linear regression dataset with n data points to obtain the optimal parameters. Which of the following statements correctly identifies the plots and gives a correct comparison between the methods?

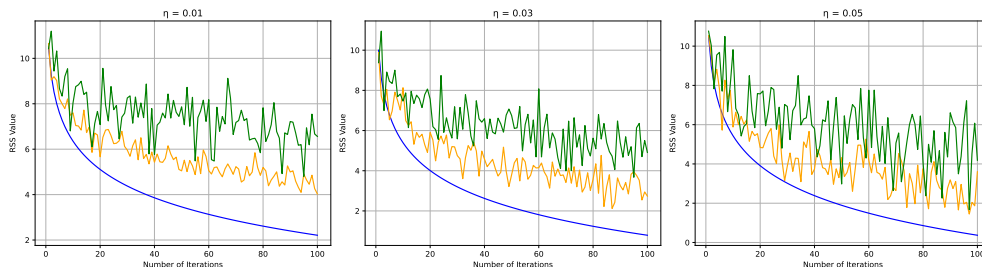


Figure 1: GD, SGD and Mini-batch SGD plots for different η (learning rate)

- ☐ blue: GD, yellow: mini-batch SGD, green: SGD.
Running K iterations of SGD returns the same parameter vector as 1 iteration of GD.

- ☐ blue: GD, yellow: SGD, green: mini-batch SGD.
GD requires more computation per iteration than SGD and mini-batch SGD.
- ☐ blue: GD, yellow: mini-batch SGD, green: SGD.
Mini-batch SGD offers a compromise between GD and SGD as it has faster convergence than SGD but more computation and storage requirements than SGD.
- ☐ blue: mini-batch SGD, yellow: SGD, green: GD.
SGD requires less storage than mini-batch SGD.

Solution: C.

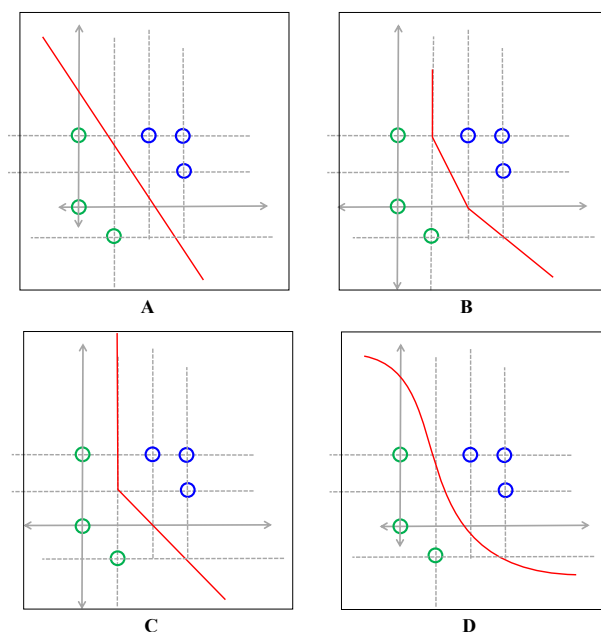
A is false since n iterations of SGD is not equivalent to 1 iteration of GD due to randomness in selection of data points in each iteration in SGD. B is false since yellow is mini-batch SGD and green is SGD. C is true since the plots are correctly identified and the tradeoff explained for mini-batch SGD is correct. D is false since plots are wrongly identified.

Problem 10: Consider a Naive Bayes classifier f and a ham/spam dataset \mathcal{D} consisting of both spam and ham documents, where each document is a sequence of words. We will use the bag-of-words method to calculate the required parameters for the Naive Bayes Classifier and we will not use any smoothing. Which of the below changes to the dataset would also change these Naive Bayes parameters?

- ☐ Duplicating each document in the dataset exactly one time.
- ☐ Changing the order of the words in each document.
- ☐ Adding the word 'start' to the beginning of every document when 'start' already existed in the vocabulary.
- ☐ Duplicating every word in each document exactly two times.

Solution: C

Problem 11: Which of the following is the 1-nearest neighbors decision boundary for the dataset that is shown in the figure? We use the Euclidean distance metric to measure the distance between any two points.



- ☐ A
- ☐ B
- ☐ C
- ☐ D

Solution: C

Problem 12: Which of the following correctly compares the concept of support vectors in hard-margin vs. soft-margin SVM?

- ☐ In a hard-margin SVM, only points on the margin can be support vectors; in a soft-margin SVM, any point (inside or beyond the margin) can also become a support vector.
- ☐ In both hard-margin and soft-margin SVMs, every training point is a support vector.
- ☐ In a hard-margin SVM, points on the correct side but more than one margin-width away from the decision boundary are also considered support vectors.
- ☐ In a soft-margin SVM, no points that are misclassified can be support vectors.

Solution: A

Problem 13: Consider the Maximum-likelihood estimator (MLE) and the Maximum a posteriori (MAP) estimator in an arbitrary parameter estimation problem. Select the correct option from the following.

- ☐ It is always possible to find a closed-form expression for the MLE.
- ☐ If the MAP estimate is unique, then the MLE is unique.
- ☐ MLE is a special case of MAP estimation assuming a uniform prior distribution.
- ☐ MLE maximizes the probability of a set of parameters given the observed data, while MAP maximizes the joint probability of observed data conditioned on a set of parameters.

Solution: C

Problem 14: Consider a multi-class classification problem with $K = 5$ classes. How many binary classifiers you need to train respectively using either “one-versus-one” or “one-versus-rest” approaches?

- ☐ 20 and 5
- ☐ 10 and 5
- ☐ 5 and 10
- ☐ 10 and 10

Solution: B

3 Descriptive

Problem 15: For a course assignment on spam classification using logistic regression, Alice and Bob are given a dataset of emails containing two categories “spam” and “ham”. Alice argues that since the main motive is to figure out “spam” email, one should label the spam emails as $y = 1$ and the “ham” as $y = 0$. Bob argues that class labels are not important in logistic regression and labelled the spam emails as $y = 0$ and the “ham” as $y = 1$.

Recall that the cross-entropy loss function is:

$$\mathcal{E}(\mathbf{w}) = - \sum_{n=1}^N \{y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log[1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)]\} \quad (1)$$

where \mathbf{w} denote the weight vector (including the bias) and \mathbf{x} is the augmented feature vector (including a 1 corresponding to the bias). The sigmoid function is $\sigma(a) = 1/(1 + e^{-a})$. The loss function can be simplified to the following expression:

$$\mathcal{E}(\mathbf{w}) = \sum_{n=1}^N \{(1 - y_n) \mathbf{w}^\top \mathbf{x}_n + \log(1 + e^{-\mathbf{w}^\top \mathbf{x}_n})\} \quad (2)$$

- a. Using the simplified version of the cross-entropy function shown in (2), express Bob’s cross-entropy loss function $\mathcal{E}_{Bob}(\mathbf{w})$ in terms of Alice’s cross-entropy loss function $\mathcal{E}_{Alice}(\mathbf{w})$ and show that

$$\mathcal{E}_{Bob}(\mathbf{w}) = \mathcal{E}_{Alice}(-\mathbf{w}) \quad (3)$$

- b. Suppose that Alice and Bob start with the same initialization $\mathbf{w}_{Alice}^{(0)} = \mathbf{w}_{Bob}^{(0)} = \mathbf{0}$, and run t iterations of gradient descent on their own loss functions using the same fixed learning rate η . Prove that for any $t \in \mathbb{N}$,

$$\mathbf{w}_{Alice}^{(t)} + \mathbf{w}_{Bob}^{(t)} = \mathbf{0}$$

Hint: Use your answer from part (a) and prove by induction.

- c. Compare the equations of the decision boundaries learnt by Alice and Bob and comment on their similarity and/or differences.

Solution:

- a. The cross-entropy function can be simplified as follows:

Let $\mathcal{E}_{Alice}(\mathbf{w})$ be the loss function calculated by Alice. It is equal to the simplified cross-entropy error function:

$$\mathcal{E}_{Alice}(\mathbf{w}) = \sum_{n=1}^N \{(1 - y_n)(\mathbf{w}^\top \mathbf{x}_n) + \log(1 + e^{-\mathbf{w}^\top \mathbf{x}_n})\} \quad (4)$$

Bob’s loss function will replace y_n by $1 - y_n$ everywhere to get

$$\mathcal{E}_{Bob}(\mathbf{w}) = \sum_{n=1}^N \{y_n(\mathbf{w}^\top \mathbf{x}_n) + \log(1 + e^{-\mathbf{w}^\top \mathbf{x}_n})\} \quad (5)$$

$$= \sum_{n=1}^N \{y_n(\mathbf{w}^\top \mathbf{x}_n) + \log(1 + e^{\mathbf{w}^\top \mathbf{x}_n}) - \mathbf{w}^\top \mathbf{x}_n\} \quad (6)$$

$$= \sum_{i=1}^n \left((1 - y_i) (-\mathbf{w}^\top \mathbf{x}_i) + \log(1 + e^{\mathbf{w}^\top \mathbf{x}_i}) \right) \quad (7)$$

$$= \mathcal{E}_{Alice}(-\mathbf{w}) \quad (8)$$

Hence,

$$\nabla_{\mathbf{w}} \mathcal{E}_{Bob}(\mathbf{w}) = -\nabla_{\mathbf{w}} \mathcal{E}_{Alice}(-\mathbf{w}) \quad (9)$$

- b. **Proof by induction:** Clearly the results holds for $k = 0$ as all the weights are zero. Assume it holds for the k -th iteration, i.e., $\mathbf{w}_{Alice}^{(k)} + \mathbf{w}_{Bob}^{(k)} = \mathbf{0}$. The next weight vector after a gradient descent step is given by

$$\mathbf{w}_{Alice}^{(k+1)} = \mathbf{w}_{Alice}^{(k)} - \eta \nabla_{\mathbf{w}} \mathcal{E}_{Alice}(\mathbf{w}_{Alice}^{(k)}). \quad (10)$$

Similarly,

$$\mathbf{w}_{Bob}^{(k+1)} = \mathbf{w}_{Bob}^{(k)} - \eta \nabla_{\mathbf{w}} \mathcal{E}_{Bob}(\mathbf{w}_{Bob}^{(k)}) \quad (11)$$

$$= \mathbf{w}_{Bob}^{(k)} + \eta \nabla_{\mathbf{w}} \mathcal{E}_{Alice}(-\mathbf{w}_{Bob}^{(k)}) \quad (12)$$

$$= -\mathbf{w}_{Alice}^{(k)} + \eta \nabla_{\mathbf{w}} \mathcal{E}_{Alice}(\mathbf{w}_{Alice}^{(k)}) \quad (13)$$

$$= -\mathbf{w}_{Alice}^{(k+1)} \quad (14)$$

This completes the proof.

- c. The Decision boundary is given by the equation $\mathbf{w}^\top x = 0$ and it is invariant of sign. Hence both Alice and Bob will end up with the same decision boundary.

Problem 16: Suppose you are given a training dataset $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ where each (\mathbf{x}_i, y_i) pair consists of a d -dimensional feature vector \mathbf{x}_i and label y_i , $i = 1, 2, \dots, n$.

- a. Suppose that $d = 2$. For each training data point i , the second feature's value is equal to the label y_i : thus, $\mathbf{x}_i = \begin{bmatrix} \mathbf{x}_{i,1} \\ y_i \end{bmatrix}$. You now solve for the linear regression coefficients \mathbf{w}^* that minimize the squared-error $\sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$.

Is the optimal \mathbf{w}^* unique? If so, find a closed-form expression for \mathbf{w}^* . If not, explain why not and suggest a way to modify the linear regression formulation so that it has a unique solution \mathbf{w}^* .

Solution: By inspection, the optimal $\mathbf{w}^* = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$. Note that we can also use the linear regression formula $\mathbf{w}^* = \left(\begin{bmatrix} \mathbf{x}^T \\ \mathbf{y}^T \end{bmatrix} \begin{bmatrix} \mathbf{x} & \mathbf{y} \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{x}^T \\ \mathbf{y}^T \end{bmatrix} \mathbf{y} = \begin{bmatrix} \mathbf{x}^T \mathbf{x} & \mathbf{x}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{x} & \mathbf{y}^T \mathbf{y} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x}^T \mathbf{y} \\ \mathbf{y}^T \mathbf{y} \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$. This solution is unique.

- b. Now suppose you add another new feature to each input vector from part (b), so that $d = 3$. This new feature is equal to the sum of the first input feature and the label, so that each feature vector

$$\mathbf{x}_i = \begin{bmatrix} \mathbf{x}_{i,1} \\ y_i \\ \mathbf{x}_{i,1} + y_i \end{bmatrix}.$$

Is the optimal \mathbf{w}^* that minimizes $\sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$ unique? If so, find a closed-form expression for \mathbf{w}^* . If not, explain why not and suggest a way to modify the linear regression formulation so that it has a unique solution \mathbf{w}^* .

Solution: The features are no longer linear independent, so \mathbf{w}^* is no longer unique. We can use ridge regression to make the solution unique.

- c. Now suppose you apply ridge regression to the training dataset from part (c), i.e., you solve for \mathbf{w}_{ridge}^* that minimizes $\sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_2^2$ for some $\lambda > 0$. Is the optimal \mathbf{w}_{ridge}^* unique? Explain why or why not in one or two sentences.

Solution: Yes, \mathbf{w}_{ridge}^* is unique. The ridge regression solution is always unique, since it is equal to $(2\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$, which is a unique solution.