# Week 2 Practice

# CSCI 567 Machine Learning

## Fall 2025

Instructor: Haipeng Luo

- 1. MULTIPLE-CHOICE QUESTIONS: One or more correct choice(s) for each question.
  - **1.1.** Which one of these is a sign of 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
  - **1.2.** Which of the following can help prevent overfitting?
    - a. Using more training data
    - b. Training until you get the smallest training error
    - c. Including a regularization term in the loss function
    - d. All of the above
  - **1.3.** Let  $\mathbf{X} \in \mathbb{R}^{N \times D}$  be a data matrix with each row corresponding to the feature of an example and  $\mathbf{y} \in \mathbb{R}^N$  be a vector of all the outcomes. The least square solution is  $(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$ . Which of the following is the least square solution if we scale each data point by a factor of 4 (i.e. the new dataset is  $4\mathbf{X}$ )?
    - a.  $4(\mathbf{X}^{\dagger}\mathbf{X})^{-1}\mathbf{X}^{\dagger}\mathbf{y}$
    - b.  $\frac{1}{4}(\mathbf{X}^{\dagger}\mathbf{X})^{-1}\mathbf{X}^{\dagger}\mathbf{y}$

c. 
$$\frac{1}{2}(\mathbf{X}^{\intercal}\mathbf{X})^{-1}\mathbf{X}^{\intercal}\mathbf{y}$$

d. None of the above

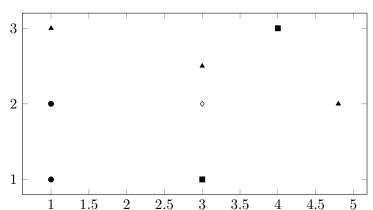
1.4. Consider the following two-dimensional dataset with N=7 training points of three classes (triangle, square, and circle), and additionally one test point denoted by the diamond. Which of the following configuration of the K-nearest neighbor algorithm will predict triangle for the test point?

a. 
$$K = 1$$
, L2 distance

b. 
$$K = 3$$
, L1 distance

c. 
$$K = 3$$
, L2 distance





- **1.5.** Which of the following on linear regression is correct?
  - a. The least square solution has a closed-form formula, even if L2 regularization is applied.
  - b. The covariance matrix  $X^TX$  is not invertible if and only if the number of data points N is smaller than the dimension D.
  - c. When the covariance matrix  $X^TX$  is not invertible, the Residual Sum of Squares (RSS) objective has no minimizers.
  - d. Linear regression is a parametric method.

#### 2. Nearest Neighbor Classification

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

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

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

$$C(\mathbf{x}, \mathbf{x}') = 1 - \frac{\mathbf{x}^{\mathsf{T}} \mathbf{x}'}{||\mathbf{x}||_{2}||\mathbf{x}'||_{2}} = 1 - \frac{\sum_{d=1}^{D} (x_{d} \cdot x'_{d})}{||\mathbf{x}||_{2}||\mathbf{x}'||_{2}}$$
(2)

where the L2 norm of x is defined as

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

Show that, if data is normalized with unit L2 norm, that is,  $||\mathbf{x}||_2 = 1$  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.

### 3. Linear Regression

In the lectures, we have described the least mean square solution for linear regression as

$$\boldsymbol{w}^* = (\tilde{\boldsymbol{X}}^{\mathrm{T}}\tilde{\boldsymbol{X}})^{-1}\tilde{\boldsymbol{X}}^{\mathrm{T}}\boldsymbol{y}$$

where  $\tilde{\boldsymbol{X}}$  is the design matrix (N rows, D+1 columns) and  $\boldsymbol{y}$  is the N-dimensional column vector of the true values in the training data  $\mathcal{D} = \{(\boldsymbol{x}_n, y_n)\}_{n=1}^N$ .

Question 1 We mentioned a practical challenge for linear regression: when  $\tilde{X}^T\tilde{X}$  is not invertible. Please use a concise mathematical statement (in one sentence) to summarize the relationship between the training data  $\tilde{X}$  and the dimensionality of w when this scenario happens. Then use this statement to explain why this scenario must happen when N < D+1.

Question 2 In this problem we use the notation  $w_0 + \mathbf{w}^T \mathbf{x}$  for the linear model, that is, we do not append the constant feature 1 to x. In the lecture we saw that when D = 0, the bias  $w_0^*$  is simply the mean of the sample responses

$$w_0^* = \frac{1}{N} \mathbf{1}_N^{\mathrm{T}} \boldsymbol{y} = \frac{1}{N} \sum_n y_n, \tag{4}$$

where  $\mathbf{1}_N = [1, 1, \dots, 1]^{\mathrm{T}}$  is an N-dimensional column vector whose entries are all ones. Now, we would like you to generalize this to arbitrary D and arrive at a more general condition where Eqn. (4) holds. Please do so by following the three steps below:

- 1) write down the residual sum of squares objective w.r.t. the variable of interest;
- 2) take derivative with respect to  $w_0$  and set it to 0;
- 3) solve the obtained equation and conclude that Eqn. (4) holds if

$$\frac{1}{N} \sum_{n} x_{nd} = 0, \quad \forall d = 1, 2, \dots, D,$$
 (5)

that is, each feature has zero mean.