Week 2 Discussion

CSCI 567 Fall 2025

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

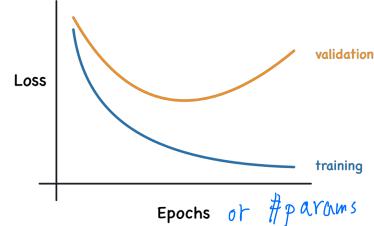
c. High training error, low test error

d. High training error, high test error

- 1. 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

This is like the "definition" of overfitting

The Learning Curves



Credit: Overfitting and Underfitting (Kaggle)

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

- **1.2.** 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

Regularization discourages overly complex model

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. $4(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$

b. $\frac{1}{4}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$

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

d. None 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.
$$4(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$$

b.
$$\frac{1}{4}(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y}$$

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

d. None of the above

$$= \left(\left(\frac{4x}{4x} \right)^{T} \right)^{T}$$

$$= \left(\frac{16x^{T}}{4x} \right)^{T}$$

$$(cX)^T = c \cdot X^T$$
$$(cX)^{-1} = \frac{1}{c} \cdot X^{-1}$$

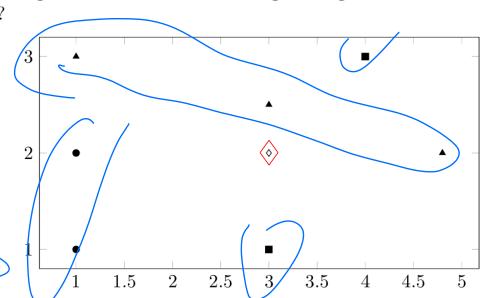
1.4. Consider the following two-dimensional dataset with N=7 training points of three classes (triangle, square, and circle), and additionally one test point denoted by the diamond. Which of the following configuration of the K-nearest neighbor algorithm will predict triangle for the test point?





c. K = 3, L2 distance

d. K = 7, any distance



$$\overline{\chi} \in \mathbb{R}$$
 $\overline{\chi} : \overline{\chi} : \overline{\chi}$

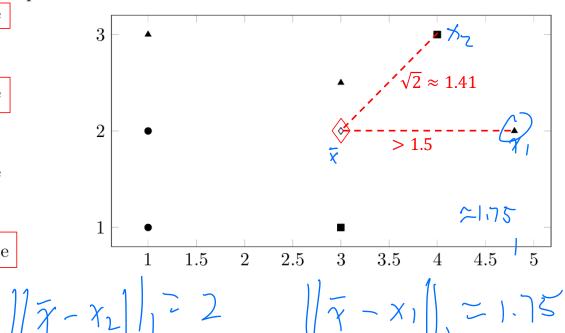
1.4. Consider the following two-dimensional dataset with N=7 training points of three classes (triangle, square, and circle), and additionally one test point denoted by the diamond. Which of the following configuration of the K-nearest neighbor algorithm will predict triangle for the test point?



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

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

d.
$$K = 7$$
, any 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.

- **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.

 Even if $N \ge D$, X^TX can still be **not** invertible, e.g., due to "collinearity"
 - c. When the covariance matrix X^TX is not invertible, the Residual Sum of Squares (RSS) objective has no minimizers.

 Infinitely many solutions

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

$$x^T x = \|x\|_2^2$$

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

When
$$||x|| = ||x'|| = 1$$
, $C(x, x') = 1 - x^T x'$

$$E(x, x')^2 = (x - x')^T (x - x') = ||x||^2 + ||x'||^2 - 2x^T x' = 2C(x, x')$$

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

Question 3.1 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.

Why. If N<D+1, then XX not invertible

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

Question 3.1 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.

$$N < D + 1 \Rightarrow \operatorname{rank}(\tilde{X}^T \tilde{X}) < D + 1$$

 $\Leftrightarrow \tilde{X}^T \tilde{X} \text{ not invertible}$

Why?

- $\tilde{X}^T \tilde{X}$ is $(D+1) \times (D+1)$ matrix
- $\operatorname{rank}(\tilde{X}^T \tilde{X}) = \operatorname{rank}(\tilde{X}) \le \min\{N, D+1\}$

In general, invertible \Leftrightarrow full rank

Question 3.2 In this problem we use the notation $w_0 + \mathbf{w}^T \mathbf{x}$ for the linear model, that is, Design matrix (without bias) we do not append the constant feature 1 to x. In the lecture we saw that when D=0, the the $X = \begin{pmatrix} -x_1^T - \\ -x_2^T - \\ \vdots \\ -x_N^T - \end{pmatrix} \in \mathbb{R}^{N \times D}$ (4)

bias
$$w_0^*$$
 is simply the mean of the sample responses

FIX ONY WY,

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

$$w_0 = \overline{N} \mathbf{1}_N^T \mathbf{y} = \overline{N} \sum_n y_n,$$
 (4) where $\mathbf{1}_N = [1, 1, \dots, 1]^T$ is an N-dimensional column vector whose entries are all ones. Now, Prediction we would like you to generalize this to arbitrary D and arrive at a more general condition $\hat{y}_i = \mathbf{1}_N \mathbf{1}_N^T \mathbf{y} = \mathbf{1}_N^T \mathbf{y} \mathbf{y}_n$

we would like you to generalize this to arbitrary
$$D$$
 and arrive at a more general con where Eqn. (4) holds.

Objective wind mixed (Wo) WIN m

E/x ony
$$w^*$$
, $w^*_b = a \cdot g \cdot w^*_0 + ||y - ||x^*_0 - ||y^*_0||^2$
2) take derivative with respect to w_0 and set it to 0;

derivative with respect to
$$w_0$$
 and set it to 0;

$$0 = -2 \int_{N}^{\infty} \left(y - \chi_{W} - W_{0}^{*} \int_{N} \right) = 1 \quad \text{if } 1 = 1 \quad \text$$

$$($$
 $($ $)$

$$(y-\chi_w-w_01_N)=//v_w$$

the variable of interest;
$$|C_{ij}(x)| \leq |C_{ij}(x)|$$

"Residual" $e V V \circ v$ $\epsilon_i = y_i - \hat{y}_i \qquad e = V - V$

$$= y_i - x_i^T w - w_0$$

$$= \sqrt{\frac{1}{n}} \sqrt{\frac{1$$

$$= y_i - y_$$

$$\frac{2}{i} = \int$$

$$\|\epsilon\|_{2}^{2} = \epsilon^{T}$$

 $\hat{y}_i = x_i^T w + w_0 \qquad y = \chi_0 + w_0$

$$\begin{cases}
 e_i = ||e||_2 = e_i \\
 i = ||e||_2 = e_i \\
 ||e||_2 = e_i = 1
 \end{cases}$$

$$\begin{cases}
 ||e||_2 = e_i = 1 \\
 ||e||_2 = e_i = 1
 \end{cases}$$

3) solve the obtained equation and conclude that Eqn. (4) holds if

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

$$\begin{cases} \langle \zeta \rangle \Rightarrow & \langle \chi \rangle \\ \langle \zeta \rangle \Rightarrow & \langle \chi \rangle \end{cases} \Rightarrow \begin{cases} \langle \chi \rangle \\ \langle \chi \rangle \\ \langle \chi \rangle \end{cases} \Rightarrow \begin{cases} \langle \chi \rangle \\ \langle \chi \rangle \\ \langle \chi \rangle \end{cases} \Rightarrow \begin{cases} \langle \chi \rangle \\ \langle \chi \rangle \\ \langle \chi \rangle \\ \langle \chi \rangle \end{cases} \Rightarrow \begin{cases} \langle \chi \rangle \\ \langle \chi \rangle \\ \langle \chi \rangle \\ \langle \chi \rangle \\ \langle \chi \rangle \end{cases} \Rightarrow \begin{cases} \langle \chi \rangle \\ \langle$$

that is, each feature has zero mean. Wo = 1/1/1/ = 1/2 /2.