

University of Central Missouri
Department of Computer Science & Cybersecurity

CS5710 Machine Learning
Fall 2025

Home Assignment 3.

Student name: SAI SRAVAN CHINTALA
#700773836

Submission Requirements:

- Once finished your assignment push your source code to your repo (GitHub) and explain the work through the ReadMe file properly. Make sure you add your student info in the ReadMe file.
- Comment your code appropriately *IMPORTANT*.
- Any submission after provided deadline is considered as a late submission.

Part A: Calculation

Q1. Multi-Input Feedforward

Consider a neural network with 3 inputs $x_1=2, x_2=1, x_3=3$, one hidden layer of 2 sigmoid units, and one sigmoid output.

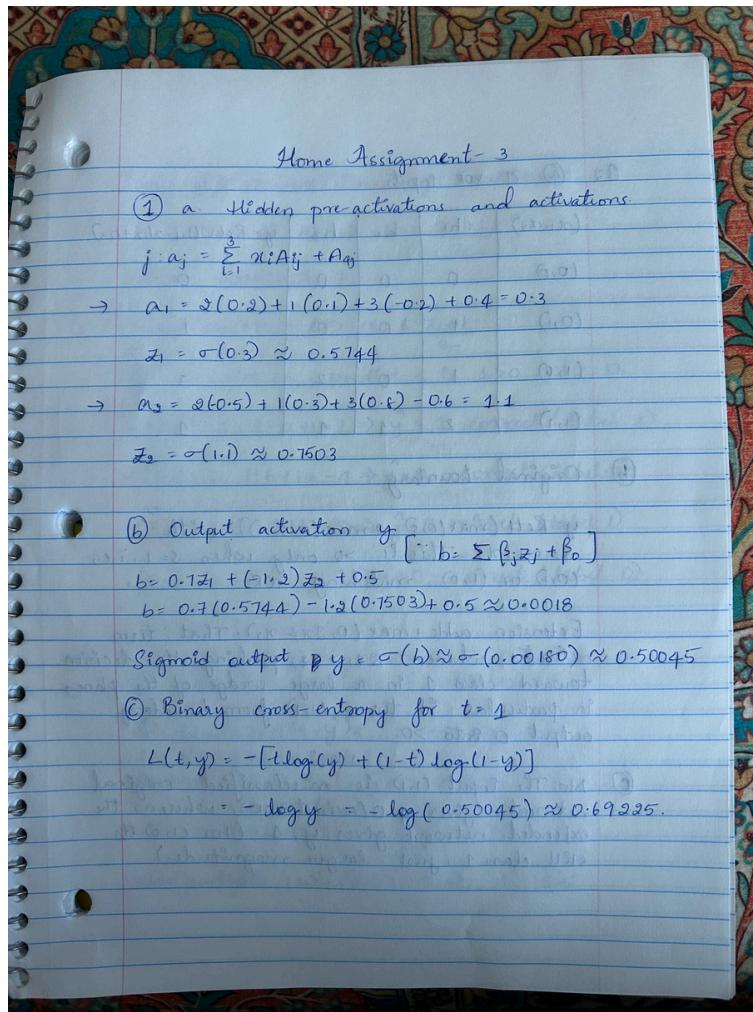
Weights:

$$A = \begin{bmatrix} 0.2 & -0.5 \\ 0.1 & 0.3 \\ -0.2 & 0.8 \\ 0.4 & -0.6 \end{bmatrix}, \quad B = \begin{bmatrix} 0.7 \\ -1.2 \\ 0.5 \end{bmatrix}$$

- A includes biases (last row).
- B includes bias (last element).

Tasks

- Compute hidden pre-activations and activations.
- Compute the output activation y .
- If the true label is $t=1$, calculate the binary cross-entropy loss.



Q2. XOR with ReLU Network

We use the XOR ReLU network from the slides

$$h_1 = \text{ReLU}(x_1 + x_2), \quad h_2 = \text{ReLU}(x_1 + x_2 - 1), \quad y = \text{ReLU}(h_1 - 2h_2)$$

Extend it with an additional hidden unit:

$$h_3 = \text{ReLU}(2x_1 - x_2)$$

and new output:

$$y = \text{ReLU}(h_1 - 2h_2 + h_3)$$

Tasks

- Compute outputs for all four XOR inputs $(0,0), (0,1), (1,0), (1,1)$.
- Compare the decision boundary with the original 2-hidden-unit XOR network.
- Does this extension still compute XOR exactly? If not, which inputs differ?

Q2. (a) 4 XOR inputs

(x_1, x_2)	h_1	h_2	h_3	$y = \text{ReLU}(h_1 - 2h_2 + h_3)$
$(0,0)$	0	0	0	0
$(0,1)$	1	0	0	1
$(1,0)$	1	0	2	3
$(1,1)$	2	1	1	1

(b) Original boundary :

$y = \text{ReLU}(\max(0, s) - 2\max(0, s-1))$ with $s = x_1 + x_2$. This is > 0 only when $s = 1$, i.e., $(0,1)$ or $(1,0)$, matching XOR + AND.

Extension adds $+\max(0, 2x_1 - x_2)$. That term activates whenever $2x_1 > x_2$, pushing the decision toward class 1 in a large wedge of the plane, in particular it turns $(1,1)$ from boundary output 0 into > 0 .

(c) No, The input $(1,1)$ is misclassified: original network gives $y=0$ (correct XOR), whereas the extended network gives $y=1$. (For $(1,0)$ its still class 1, just larger magnitude.)

Q3. Decision Boundary and Misclassification

You train a perceptron with decision rule:

$$y = \begin{cases} 1 & \text{if } w_1x_1 + w_2x_2 + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

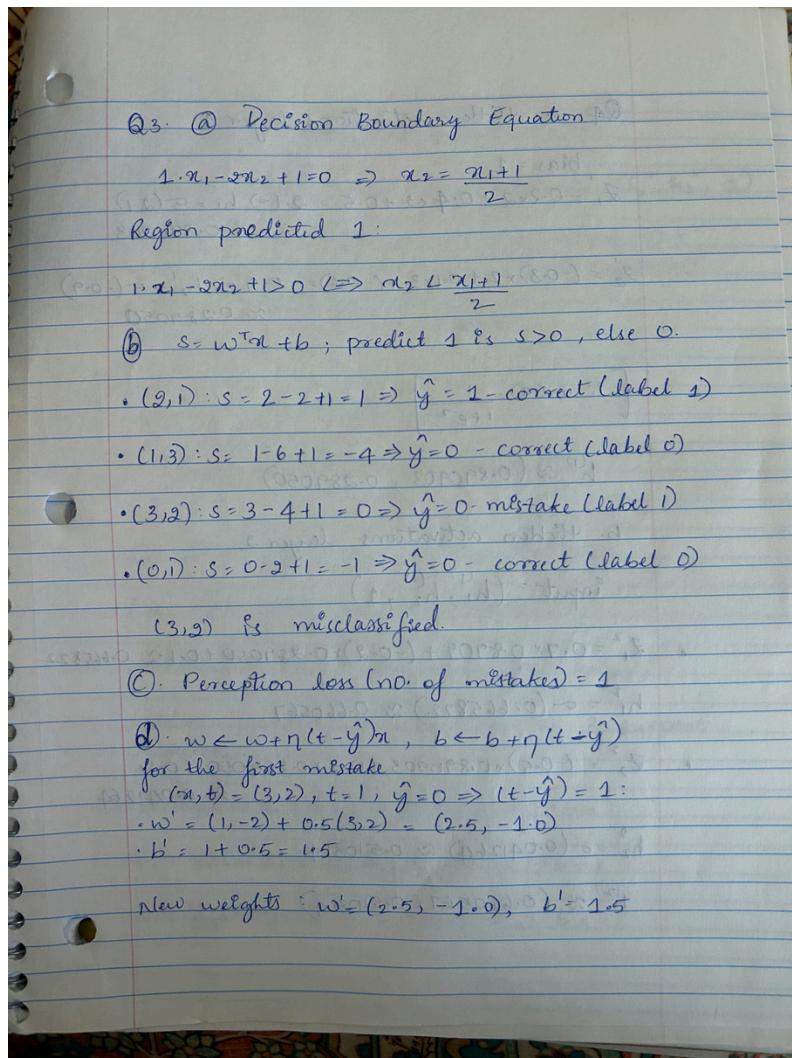
with parameters $w_1=1, w_2=-2, b=1$.

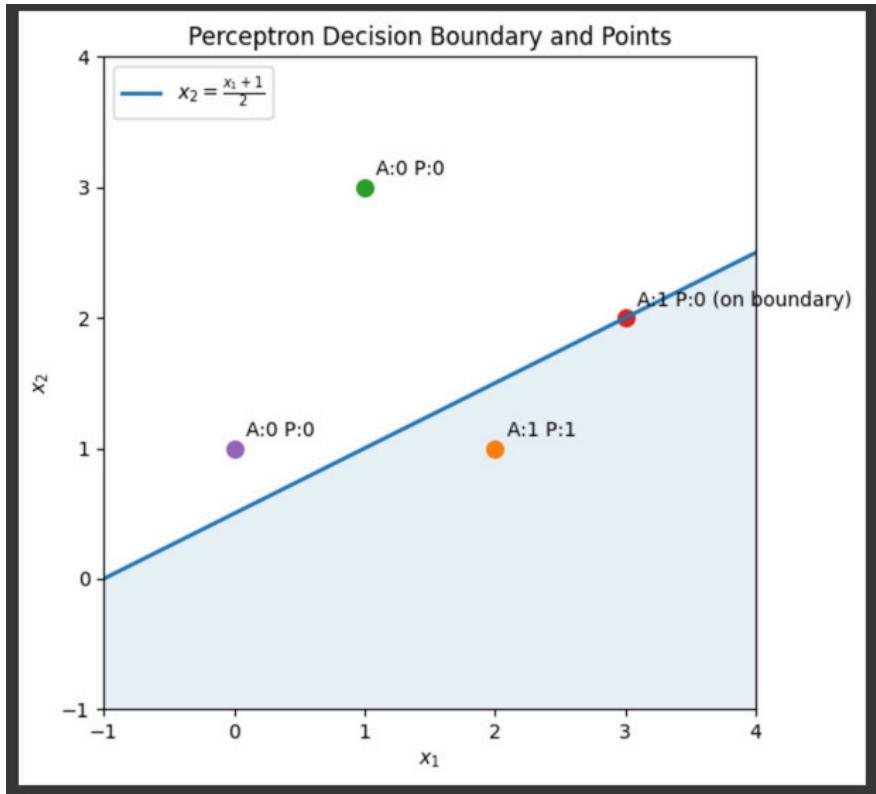
Dataset:

$(2,1) \rightarrow 1, (1,3) \rightarrow 0, (3,2) \rightarrow 1, (0,1) \rightarrow 0$.

Tasks

- Write the decision boundary equation and sketch it.
- Classify each point using the perceptron rule. Which ones are misclassified?
- Compute the perceptron loss (number of mistakes).
- If we apply one weight update with learning rate $\eta=0.5$ on the first mistake, what are the new weights?





Q4. Multi-Layer Forward Pass (Matrix Style)

A network has 2 hidden layers, each with 2 sigmoid units.

- Input $(x_1, x_2) = (2, 3)$
- First layer weights:

$$A^{(1)} = \begin{bmatrix} 0.2 & -0.3 \\ 0.4 & 0.1 \\ 0.5 & -0.6 \end{bmatrix}$$

- Second layer weights:

$$A^{(2)} = \begin{bmatrix} 0.7 & -0.5 \\ -0.2 & 0.3 \\ 0.1 & 0.4 \end{bmatrix}$$

- Output weights: (1.0, -1.2, 0.3)

Tasks:

- Compute hidden activations in layer 1.
- Compute hidden activations in layer 2.
- Compute the final output.

Q4. a) hidden activations layer 1

bias = 1

$$z_1 = 0.2 \times 2 + 0.4 \times 3 + 0.5 = 2.1 \Rightarrow h_1 = \sigma(2.1) \approx 0.890903$$

$$z_2 = (-0.3) \times 2 + 0.1 \times 3 + (-0.6) = -0.9 \Rightarrow h_2 = \sigma(-0.9) \approx 0.289050$$

$$\left[\because \sigma(x) = \frac{1}{1+e^{-x}} \right]$$

$$h^{(1)} \approx (0.890903, 0.289050)$$

b) hidden activations layer 2

$$\text{input} = (h_1^{(1)}, h_2^{(1)}, 1)$$

- $z_1^{(2)} = 0.7 \times 0.8909 + (-0.2) \times 0.289050 + 0.1 \approx 0.665822$
- $h_1^{(2)} = \sigma(0.665822) \approx 0.660567$
- $z_2^{(2)} = (-0.5) \times 0.890903 + 0.3 \times 0.289050 + 0.4 \approx 0.041264$
- $h_2^{(2)} = \sigma(0.041264) \approx 0.510314$

$$h^{(2)} \approx (0.660567, 0.510314)$$

c. Final output (sigmoid)

Output weights $w = (1.0, -1.2, 0.3)$ for $(h_1^T, h_2^T, 1)$

$$z^{(\text{out})} = 1.0 \times 0.660567 + (-1.2) \times 0.510314 + 0.3 \\ \approx 0.348190$$

$$y = \sigma(0.348190) \approx 0.586179.$$

Q5. Linear SVM

We have a dataset in 2D:

- Positive points ($y=+1$): $p_1=(1,3)$, $p_2=(2,2)$
- Negative points ($y=-1$): $n_1=(0,0)$

Supposethese are the support vectors.

Tasks:

1. Augment each point with a bias term.
2. Form the dual constraint equations for $\alpha_1, \alpha_2, \alpha_3$.
3. Solve for the dual variables.
4. Compute the weight vector $\tilde{w} = \sum \alpha_i y_i \tilde{x}_i$.
5. Extract w and b . Write the final separating hyperplane equation.
6. Verify that all support vectors satisfy the margin conditions.

$y = (0.0 + 0.1 + 0) \approx 0.1$ margin
 Q5 (1) $\bar{x} = [x_1, 1]^\top \times 2.0 \approx 1.0$ margin

$$\bar{x}_1 = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix}, \bar{x}_2 = \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix}, \bar{x}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

(2) $\alpha_1 + \alpha_2 - \alpha_3 = 0 \Rightarrow \alpha_3 = \alpha_1 + \alpha_2$

(3) $x_3 = (0, 0)$, $w = \alpha_1 x_1 + \alpha_2 x_2$

Let $w = (w_1, w_2) = (\alpha_1 + 2\alpha_2, 3\alpha_1 + 2\alpha_2)$

Using margin equalities on the 3 support vectors.

- (i) $y_1(w^T x_1 + b) = 1 \Rightarrow (10\alpha_1 + 8\alpha_2) + b = 1$
- (ii) $y_2(w^T x_2 + b) = 1 \Rightarrow (8\alpha_1 + 8\alpha_2) + b = 1$
- (iii) $y_3(w^T x_3 + b) = 1 \Rightarrow 0 + b = -1 \Rightarrow b = -1$



Substitute $b = -1$ into (i) & (ii)

$$10\alpha_1 + 8\alpha_2 = 2 \Rightarrow 5\alpha_1 + 4\alpha_2 = 1$$

$$8\alpha_1 + 8\alpha_2 = 2 \Rightarrow \alpha_1 + \alpha_2 = 0.25$$

$$\alpha_2 = 0.25, \alpha_1 = 0, \alpha_3 = \alpha_1 + \alpha_2 = 0.25$$

$$(4) \bar{w} = \sum_i \alpha_i y_i \bar{x}_i$$

$$\bar{w} = 0 \cdot (+1) \bar{x}_1 + 0.25 \times (+1) \bar{x}_2 + 0.25 (-1) \bar{x}_3 = \begin{bmatrix} 0.5 \\ 0.5 \\ 0 \end{bmatrix}$$

$$(5) w = \begin{bmatrix} 0.5 \\ 0.5 \\ 0 \end{bmatrix}, b = -1$$

Separating hyperplane:

$$0.5x_1 + 0.5x_2 - 1 = 0 \Leftrightarrow x_1 + x_2 = 2$$

$$(6) \cdot (1, 3): 0.5(1) + 0.5(3) - 1 = 1 \Rightarrow y_{+}(\cdot) = +1$$

$$(2, 2): 0.5(2) + 0.5(2) - 1 = 1 \Rightarrow y_{+}(\cdot) = +1$$

$$(0, 0): 0.5(0) + 0.5(0) - 1 = -1 \Rightarrow y_{-}(\cdot) = (-1) \cdot (-1) = +1$$

All three satisfy $y_{\pm}(w^T \bar{x}_i + b) = 1$ on the margin.

Q6. Nonlinear SVM

Consider two support vectors:

- $s_1 = (1, 0), y_1 = -1$
- $s_2 = (2, 1), y_2 = +1$

We map them using a quadratic feature transformation:

$$\phi(x_1, x_2) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

Tasks:

1. Compute $\phi(s_1)$ and $\phi(s_2)$.
2. Write the margin equations in the transformed space using α_1, α_2 .
3. Solve for α_1, α_2 .
4. Compute the separating hyperplane in feature space.
5. Express the decision function $f(x) = w^\top \phi(x) + b$.
6. Using your decision function, classify the new point $q = (1, 1)$.

6Q. (1) $s_1 = (1, 0), y_1 = -1$
 $s_2 = (2, 1), y_2 = +1$

$$\phi(x_1, x_2) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

$$\phi(s_1) = (1, 0, 0)$$

$$\phi(s_2) = (4, 2\sqrt{2}, 1)$$

(2) Dual equality: $-d_1 + d_2 = 0 \Rightarrow d_2 = d_1 = d$

Then $w = d[\phi(s_2) - \phi(s_1)] = d(3, 2\sqrt{2}, 1)$

Margin equalities:

$$-(w^\top \phi(s_1) + b) = 1 \Rightarrow w_1 + b = -1$$

$$(w^\top \phi(s_2) + b) = 1 \Rightarrow 4w_1 + 2\sqrt{2}w_2 + w_3 + b = 1$$

(3) With $w = (3d, 2\sqrt{2}d, d)$:

$$b = -1 - 3d, 12d + 8d\sqrt{2} + d + b = 1 \Rightarrow 21d + (-1 - 3d) = 1$$

$$\Rightarrow d = \frac{1}{9}$$

$$\boxed{d_1 = d_2 = \frac{1}{9}}$$

$$(4) \quad w = \left(\frac{1}{3}, \frac{2\sqrt{2}}{9}, \frac{1}{9} \right), \quad b = -\frac{4}{3}$$

Hyperplane in $z = \phi(x)$ space:

$$w^T z + b = 0 \quad \text{with} \quad z = \phi(x)$$

$$(5) \quad f(x) = w^T \phi(x) + b$$

$$= \frac{1}{3}x_1^2 + \frac{2\sqrt{2}}{9}(2x_1x_2) + \frac{1}{9}x_2^2 - \frac{4}{3}$$

$$= \frac{1}{3}x_1^2 + \frac{4}{9}x_1x_2 + \frac{1}{9}x_2^2 - \frac{4}{3}$$

$$= x_1^2 + 4x_1x_2 + x_2^2 - 12 = 0$$

$$(6) \quad g(1,1)$$

$$f(1,1) = \frac{1}{3} + \frac{4}{9} + \frac{1}{9} - \frac{4}{3}$$

$$= \frac{8}{9} - \frac{12}{9}$$

$$= -\frac{4}{9} < 0$$

$$\Rightarrow y(-1)$$

$$\Rightarrow |y(-1)| = 1$$

Part B – Short-Answer

Q1. From Biological to Artificial (6 pts)

- In one sentence each, define neuron, synapse, and activation function in the biological vs. artificial setting.
- Give two reasons why nonlinearity is required in neural networks.
- Briefly explain why a DAG assumption is used for feed-forward NNs.

Answer:

Q1 (a)	
Biological	Artificial
Neuron: an excitable cell that sends an electrochemical pulse when membrane voltage crosses a threshold	Neuron: a node in a directed acyclic graph that computes a weighted sum of inputs and passes it through a nonlinearity
Synapse: a connection between neurons through which signals are transmitted	Synapse: a learnable weight on an edge b/w artificial neurons.
Activation function: the thresholding behavior that makes a neuron "fire" only when input is sufficiently high	Artificial function: a nonlinear function applied to the weighted sum to enable complex decision functions.

(b) To represent non-linear decision boundaries—a single linear unit cannot solve problems like XOR; hidden nonlinear units make it separable.

To approximate complex functions / learn rich features—with nonlinear activations. NNs are universal function approximators and hidden units can learn nonlinear feature combinations.

(C) A feed-forward NN is defined as a DAG so its computation is a well-defined, acyclic pipeline from inputs to outputs—i.e., a differentiable function you can evaluate forward and optimize by backprop—without feedback loops

Q2. Architecture & Capacity (8 pts)

a) Define depth and width.

b) What are two distinct roles hidden units might learn (e.g., feature selection vs. projection)?

c) Contrast the cases $D < M$ and $D > M$ in a single-hidden-layer network: what representational trade-offs do they imply?

d) State the universal approximation idea for 1-hidden-layer NNs (informally).

Answers:

(a) Depth: the number of hidden layers between input and output (how many learned transformations are stacked). Width: the number of units in a hidden layer.

(b) Lower-dimensional projection / feature selection (compress or pick the most useful features). Nonlinear combinations / higher-dimensional projection (compose features to make data more separable).

(c) $D < M$: acts as a bottleneck/compression → encourages compact representations and noise reduction; risk of underfitting if too small.

$D > M$: expands into a richer/higher-dimensional feature space → more capacity and easier separability; risk of overfitting and higher compute.

(d) A neural network with one hidden layer and suitable nonlinear activations (e.g., sigmoid) can approximate any continuous function on a compact domain arbitrarily well given enough hidden units.

Q3. Perceptron vs. SVM (6 pts)

a) What does the perceptron algorithm optimize?

b) What is the main optimization goal of SVM?

c) Why does SVM typically generalize better than a perceptron? Give one real-world scenario where this difference matters.

Answers:

(a) A perceptron finds any separating hyperplane by making mistake-driven updates that reduce training misclassifications; it does not explicitly maximize a margin.

(c) Because SVM maximizes margin and uses L2 regularization, it controls model capacity and has a unique convex optimum that depends only on a subset of points (support vectors); the

perceptron doesn't optimize margin and is more sensitive to noisy/outlier data. In high-dimensional text/spam classification, SVMs commonly outperform perceptrons thanks to large-margin separation with regularization.

3(b) Maximize the geometric margin b/w classes ; equivalently solve

$$\min_{w, b} \frac{1}{2} \|w\|^2 \text{ s.t. } y_i (w^T x_i + b) \geq 1 \quad (\text{hard margin}),$$

(or)

$$\min_{w, b, \epsilon} \frac{1}{2} \|w\|^2 + C \sum \epsilon_i \text{ s.t. } y_i (w^T x_i + b) \geq 1 - \epsilon_i, \epsilon_i \geq 0 \quad (\text{soft margin}).$$

Q4. Margins and Support Vectors (6 pts)

- a) Define the margin of a dataset.
- b) Explain the role of support vectors in determining the separating hyperplane.
- c) Why does maximizing the margin lead to robustness against noise?

Answers:

- (a) For a separating hyperplane $f(x) = w^T x + b$, the (geometric) margin is the minimum distance from any training point to the hyperplane. Under the SVM's canonical scaling $y_i (w^T x_i + b) \geq 1$, the margin on each side is $1 / \|w\|$ and total margin width is $2 / \|w\|$.
- (b) Support vectors are the points that lie on the margin and satisfy $y_i (w^T x_i + b) = 1$; they alone have non-zero dual weights and thus determine the hyperplane via $w = \sum_i \alpha_i y_i x_i$ and set b . Removing non-SV points does not change the solution; moving an SV does.
- (c) Maximizing the margin increases the perturbation needed to flip a label (larger safety buffer), controls capacity via $\|w\|$ (regularization), and thus reduces overfitting—making the classifier less sensitive to noise and small input variations.

Q5. PCA Concepts (8 pts)

- a) What assumption does PCA make about the structure of the data?
- b) State two equivalent definitions of PCA.
- c) Why must the data be centered before applying PCA?
- d) Explain what the eigenvalue associated with a principal component represents.

Handwritten notes on PCA:

$y_i(w^T x_i + b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0$ (margin)

5Q (a) The data lie (approximately) on a low-dimensional linear subspace of \mathbb{R}^M .
(Assuming the data lies on a low-dimensional linear subspace)

(b) Variance maximization:

$$v_1 = \arg \max_{\|v\|=1} v^T \Sigma v, v_1 \text{ orthogonal to } v_2, \dots, v_{j-1}$$

Minimum reconstruction error:

For orthonormal $V \in \mathbb{R}^{M \times K}$,

$$\min_{\|v\|=1} \frac{1}{N} \sum_{i=1}^N \|x^{(i)} - VV^T x^{(i)}\|^2$$

(c) Centering ensures the covariance uses deviations from the mean, aligning PCA with variance directions:

$$\hat{\mu} = \frac{1}{N} \sum_i \mathbf{x}^{(i)}, \quad \Sigma = \frac{1}{N} \sum_i (\mathbf{x}^{(i)} - \hat{\mu})(\mathbf{x}^{(i)} - \hat{\mu})^T$$

If data aren't centered, the first PC can point toward the mean offset and objectives/eigenvector interpretation are distorted.

(d) If \mathbf{v} is an eigenvector of Σ with eigenvalue λ , then λ is the variance captured along that principal direction (the 'energy' on that axis); the explained-variance ratio is $\lambda/\sum_k \lambda_k$.

Q6. PCA vs. Random Projection (5 pts)

- How does PCA choose directions compared to random projection?
- Give one advantage of PCA over random projection, and one disadvantage.

Answers

: (A) PCA: chooses data-dependent axes—the eigenvectors of the sample covariance — equivalently the directions that maximize variance or minimize reconstruction error.

Random projection: chooses directions at random (e.g., sample a matrix V with i.i.d. Gaussian entries and project $u=Vx$) independent of the data.

(B) Advantage (PCA): it preserves the most variance / least reconstruction error by construction; the eigenvalue on each principal component quantifies the variability captured.

Disadvantage (PCA): it requires eigendecomposition/SVD (potentially costly on large, high-dimensional data), whereas random projection is just sampling V and multiplying.