

[Fall 2025] ECE 5290/7290 and ORIE 5290
Distributed Optimization for Machine Learning and AI
Homework 1
Gradescope Due: September 12th at 5PM

Objective of This Assignment

The objective is to bridge the gap between the mathematical foundations of optimization and their practical application in machine learning. The first two problems are designed to reinforce the essential prerequisite knowledge from linear algebra and multivariate calculus, focusing on the concepts of matrix properties and gradient calculations that form the bedrock of continuous optimization. The subsequent problems will build on this foundation, guiding you to implement a gradient descent algorithm on real loss functions.

Question 1: Linear Algebra Review (15 points)

This problem reviews the concept of positive semidefinite (PSD) matrices, which is fundamental to convex optimization.

- (a) (5 points) Consider the following three matrices:

$$A = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}, \quad B = \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix}, \quad C = \begin{pmatrix} 4 & -2 \\ -2 & 1 \end{pmatrix}$$

For each matrix, determine if it is positive definite (PD), positive semidefinite (PSD), or indefinite.

- (b) (5 points) Briefly explain your reasoning for each matrix. You can justify your answer by computing eigenvalues, checking principal minors, or using the definition $\boldsymbol{\theta}^T M \boldsymbol{\theta}$.
- (c) (5 points) Consider the quadratic function

$$f(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^T M \boldsymbol{\theta} - \mathbf{b}^T \boldsymbol{\theta}.$$

Explain why the convexity of this function depends on the properties of the matrix M . What property must M have for $f(\boldsymbol{\theta})$ to be convex?

Question 2: Calculus Review (40 points)

This problem reviews multivariate calculus, which is essential for gradient-based optimization.

- (a) (5 points) Find the first derivative, $f'(\theta)$, for the following function, which requires the chain rule:

$$f(\theta) = \exp(-(\theta - 2)^2)$$

- (b) (5 points) Consider the function $h(\theta) = (\theta - 5)^2 + 3$. Find the value of θ that minimizes this function. Explain how you can use the derivative to find this minimum.
- (c) (5 points) Let $\boldsymbol{\theta} \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$, and $\mathbf{b} \in \mathbb{R}^m$. Consider the standard least-squares objective function:

$$f(\boldsymbol{\theta}) = \|A\boldsymbol{\theta} - \mathbf{b}\|_2^2$$

Derive the gradient of this function with respect to $\boldsymbol{\theta}$, which is the vector $\nabla f(\boldsymbol{\theta}) \in \mathbb{R}^n$.

- (d) (5 points) Derive the Hessian of $f(\boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$, which is the matrix $\nabla^2 f(\boldsymbol{\theta}) \in \mathbb{R}^{n \times n}$.
- (e) (5 points) Using your result from part (d), what can you say about the convexity of $f(\boldsymbol{\theta})$? (Hint: Think about the properties of the matrix $A^T A$ and your conclusions from Problem 1).
- (f) (5 points) The sigmoid function, $\sigma(z) = (1 + e^{-z})^{-1}$, is a core component of logistic regression. Find its derivative, $\sigma'(z)$, with respect to z .
Hint: Show that the derivative can be simplified to the well-known form: $\sigma(z)(1 - \sigma(z))$.
- (g) (5 points) The Binary Cross-Entropy (BCE) loss for a single example (\mathbf{x}, y) with $y \in \{0, 1\}$ is:

$$L_{BCE}(\boldsymbol{\theta}) = -y \log(\sigma(\boldsymbol{\theta}^T \mathbf{x})) - (1 - y) \log(1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}))$$

A more compact form, often used in optimization literature, is the log-sum-exp form:

$$L_{\text{comp}}(\boldsymbol{\theta}) = -y\boldsymbol{\theta}^T \mathbf{x} + \log(1 + \exp(\boldsymbol{\theta}^T \mathbf{x}))$$

Prove that these two forms are mathematically equivalent.

Hint: Let $z = \boldsymbol{\theta}^T \mathbf{x}$. Start with the BCE form and substitute the definition of $\sigma(z)$. You may find the identity $\log(1 + e^{-z}) = \log(1 + e^z) - z$ useful.

- (h) (5 points) Using the **compact form**, $L_{\text{comp}}(\boldsymbol{\theta})$, derive its gradient with respect to the parameter vector $\boldsymbol{\theta}$, which is $\nabla L_{\text{comp}}(\boldsymbol{\theta})$.
Hint: Your answer should be in a simple form involving the prediction $\sigma(\boldsymbol{\theta}^T \mathbf{x})$ and the true label y .

Question 3: Linear Regression (20 points)

Consider the linear regression problem of finding $\boldsymbol{\theta}$ that minimizes the following least square loss function

$$L(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^N (\boldsymbol{\theta}^T \mathbf{x}_i - y_i)^2 \quad (1)$$

where $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ is the training data set with $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$, $\boldsymbol{\theta} \in \mathbb{R}^d$ is the parameter to be determined, and $\lambda \in \mathbb{R}$ is the regularization coefficient. We plan to implement batch gradient descent and stochastic gradient descent algorithms to solve this problem, with a constant learning rate of α .

- (a) (5 points) Write the batch gradient descent update at the t^{th} iteration, for solving this problem.

i	\mathbf{x}_i	y_i
1	$(1, 1)^\top$	-1
2	$(-1, 1)^\top$	-1
3	$(-1, 0)^\top$	1

Table 1: Dataset for Problem 3

- (b) (5 points) Now consider you are provided with the data set given in Table 1. For this part, let $\alpha = 1$. By taking the initial parameter $\boldsymbol{\theta}^0 = (0, 0)^\top$, find the parameter $\boldsymbol{\theta}^1$ obtained after updating for one iteration using batch gradient descent.
- (c) (10 points) For the batch gradient descent in (b), after updating $\boldsymbol{\theta}^0$ via two batch gradient descent iterations, does the loss $L(\boldsymbol{\theta})$ decrease at $\boldsymbol{\theta}^2$ compared with $\boldsymbol{\theta}^0$? If yes, please show the decrements $L(\boldsymbol{\theta}^2) - L(\boldsymbol{\theta}^0)$; if no, please suggest a way to address the stepsize α that decreases in the loss $L(\boldsymbol{\theta})$.

Question 4: Logistic Regression (30 points)

We consider using logistic regression for a 2-class classification setting. Let $D = \{\mathbf{x}_i, y_i\}_{i=1}^N$ be a dataset for 2-class classification, where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \{0, 1\}$.

- (a) (10 points) The assumption for logistic regression is that the log-odds over the label class is affine i.e.

$$\log \left(\frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} \right) = \boldsymbol{\theta}^\top \mathbf{x}, \quad (2)$$

Assuming additionally the label distribution is even (i.e. $P(y=0) = P(y=1) = 0.5$), derive the posterior $P(y|\mathbf{x})$ for both $y=0$ and $y=1$. Hint: use the Bayes rule.

- (b) (10 points) For binary classification with $y_i \in \{0, 1\}$, the prediction rule using logistic regression is:

$$\begin{aligned} \hat{y}_i &= 1 \text{ if } P(y=1|\mathbf{x}) > 0.5, \\ \hat{y}_i &= 0 \text{ otherwise.} \end{aligned}$$

The loss function of the logistic regression is given by

$$L(\boldsymbol{\theta}) = \sum_{n=1}^N [-y_i \boldsymbol{\theta}^\top \mathbf{x}_i + \log(1 + \exp(\boldsymbol{\theta}^\top \mathbf{x}_i))]. \quad (3)$$

We would like to conduct the Gradient Descent on $\boldsymbol{\theta}$ to minimize the loss function. Can you perform 1-step Gradient Descent update with step size $\eta > 0$ to find $\boldsymbol{\theta}^{t+1}$ from $\boldsymbol{\theta}^t$?

- (c) (10 points) Following all the settings in question (b), additionally assume that at the current iteration t we have $\boldsymbol{\theta}^t$ that successfully classifies all data points, and the data points with both labels exist in the dataset. Compare the norm of $\boldsymbol{\theta}^{t+1}$ to the norm of $\boldsymbol{\theta}^t$. Which is greater?