
Machine Learning, 2024 Fall

Homework 1

Notice

Due 23:59 (CST), Oct 29, 2024

Plagiarizer will get 0 points.

L^AT_EX is highly recommended. Otherwise you should write as legibly as possible.

A Linear Regression with Multiple Variable [30]

In class, we primarily focused on cases where our target variable y is a scalar value, briefly mentioning scenarios involving multi-dimensional predictions. However, in the real world, we are often more interested in high-dimensional cases. We follow the notation in slides:

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_m^T \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y^m \end{bmatrix}$$

Thus for each training example, $y^{(i)}$ is scalar-valued, with p entries. We wish to use a linear model to predict the outputs, as in least squares, by specifying the parameter matrix Θ in

$$y = \Theta^T \mathbf{x}$$

where $\Theta \in \mathbb{R}^{n \times p}$.

(1) Given that the average total error in the one-dimensional case is expressed as:

$$J(\Theta) = \frac{1}{2} \sum_{i=1}^m \left(\Theta^T x^{(i)} - y^{(i)} \right)^2$$

Please provide the error in multiple-dimensional without using any matrix-vector notation.

(2) Given your answer above, please provide the objective function can be expressed as:

$$J(\Theta) = \frac{1}{2} \text{tr} \left((X\Theta - Y)^T (X\Theta - Y) \right)$$

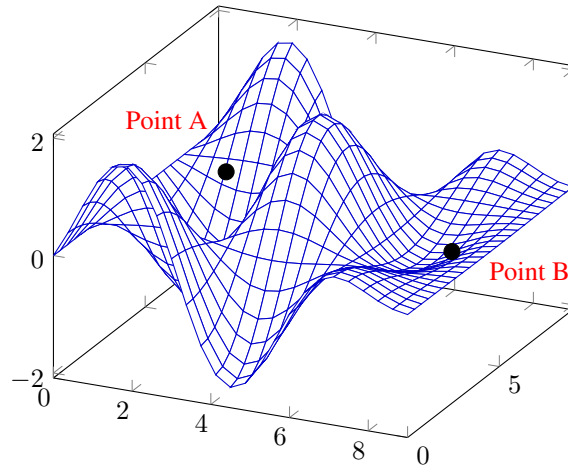
where the trace of a square matrix A , denoted $\text{tr}(A)$, is the sum of the elements on its main diagonal.

(3) Find the closed form solution for Θ which minimizes $J(\Theta)$. You MAY find following equation useful:

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}\mathbf{A}) = \mathbf{A}^T$$

$$\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}^T \mathbf{B} \mathbf{X}) = \mathbf{B} \mathbf{X} + \mathbf{B}^T \mathbf{X}$$

B Gradient Descent [20]



- (a) What is a convex problem? Why does gradient descent work for solving convex problems?
- (b) Why descent direction is the direction that gives the largest decrease in function value only holds for small step size.
- (c) If we visualize the optimization plane, gradient descent is much like a “downhill process,” starting from a point with a high loss (the initial point), selecting an appropriate direction (the negative gradient direction), and taking a small step forward. Please reasonably analyze, based on the above figure, whether initial points A and B can both be optimized to the same final point.

C Maximum Likelihood estimator [10]

We begin by considering a single binary random variable and following a Bernoulli distribution.

$$\text{Bern}(x \mid \mu) = \mu^x (1 - \mu)^{1-x}$$

- (1) Please provide the mean and variance of this distribution and the maximum likelihood estimator.
- (2) Note that the log-likelihood function of bernoulli distribution is also known as cross-entropy cost function. In practical applications, we often first use softmax to predict valid probabilities, which are then input into the cross-entropy function to calculate the loss. Show that the softmax function is equivalent to the sigmoid function in the 2-class case with 0-1 labels.

D Your First learning algorithm-PLA [15]

We introduced a simple learning algorithm called PLA. The weight update rule has the nice interpretation that it moves in the direction of classifying $\mathbf{x}(t)$ correctly. The update rule is

$$\mathbf{w}(t+1) = \mathbf{w}(t) + y(t)\mathbf{x}(t)$$

- (a) Show that $y(t)\mathbf{w}^T(t)\mathbf{x}(t) < 0$. [Hint: $\mathbf{x}(t)$ is misclassified by $\mathbf{w}(t)$.]
- (b) Show that $y(t)\mathbf{w}^T(t+1)\mathbf{x}(t) > y(t)\mathbf{w}^T(t)\mathbf{x}(t)$.
- (c) As far as classifying $\mathbf{x}(t)$ is concerned, argue that the move from $\mathbf{w}(t)$ to $\mathbf{w}(t+1)$ is a move in the right direction.

E Logistic Regression [10]

Answer the following multiple-choice questions, keeping in mind that there may be more than one correct answer. Please provide a reasonable justification; otherwise, no points will be awarded.

(1) Is logistic regression a supervised machine learning algorithm?

- TRUE
- FALSE

(2) Is Logistic regression mainly used for regression?

- TRUE
- FALSE

(3) Is it possible to apply a logistic regression algorithm on a 3-class classification problem?

- TRUE
- FALSE

(4) Which of the following methods do we use to best fit the data in logistic regression?

- Least Square Error
- Jaccard distance
- Maximum Likelihood
- all of them

F Logistic Regression [15]

We are given a sample in which each point has only one feature. Consider a binary classification problem in which sample values $x \in \mathbb{R}$ are drawn randomly from two different class distributions. The first class, with label $y = 0$, has its mean to the left of the mean of the second class, with label $y = 1$. We will use a modified version of logistic regression to classify these data points. We model the posterior probability at a test point $z \in \mathbb{R}$ as

$$P(y = 1|z) = s(z - \alpha),$$

where $\alpha \in \mathbb{R}$ is the sole parameter that we are trying to learn and $s(\gamma) = \frac{1}{1+e^{-\gamma}}$ is the logistic function. The decision boundary is $z = \alpha$ (because $s(z) = \frac{1}{2}$ there).

We will learn the parameter α by performing gradient descent on the logistic loss function. That is, for a data point x with label $y \in \{0, 1\}$, we find the α that minimizes

$$J(\alpha) = -y \ln s(x - \alpha) - (1 - y) \ln(1 - s(x - \alpha)).$$

(a) Derive the stochastic gradient descent update for J with step size $\epsilon > 0$, given a sample value x and a label y . Hint: feel free to use s as an abbreviation for $s(x - \alpha)$.

(b) Is $J(\alpha)$ convex over $\alpha \in \mathbb{R}$? Justify your answer.

(c) Now we consider multiple sample points. Since the feature dimension $d = 1$, we are given an $n \times 1$ design matrix X and a vector $y \in \mathbb{R}^n$ of labels. Consider batch gradient descent on the cost function $\sum_{i=1}^n J(\alpha; X_i, y_i)$. There are circumstances in which this cost function does not have a minimum over $\alpha \in \mathbb{R}$ at all. What is an example of such a circumstance?