

Homework #0

Due: January 31, 2025 at 11:59 PM

Welcome to CS1810! The purpose of this assignment is to help assess your readiness for this course. It will be graded for completeness and effort. **Areas of this assignment that are difficult are an indication of areas in which *you* need to self-study.** If you find you are struggling with many of these questions, it might be prudent to postpone taking this course until after you have mastered the necessary prerequisites. *During the term, the staff will be prioritizing support for new material taught in CS1810 over teaching prerequisites.* If you are unsure about your readiness, please contact the head TFs for advice.

1. Please type your solutions after the corresponding problems using this L^AT_EX template, and start each problem on a new page.
2. Please submit the **writeup PDF to the Gradescope assignment ‘HW0’**. Remember to assign pages for each question.
3. Please submit your L^AT_EX file and code files (i.e., anything ending in .py, .ipynb, or .tex) to the **Gradescope assignment ‘HW0 - Supplemental’**.

Problem 1 (Modeling Linear Trends - Linear Algebra Review)

In this class, we will be exploring the question of “how do we model the trend in a dataset” under different guises. In this problem, we will explore the algebra of modeling a linear trend in data. We call the process of finding a model that capture the trend in the data, “fitting the model.”

Learning Goals: In this problem, you will practice translating machine learning goals (“modeling trends in data”) into mathematical formalism using linear algebra. You will explore how the right mathematical formalization can help us express our modeling ideas unambiguously and provide ways for us to analyze different pathways to meeting our machine learning goals.

Let’s consider a dataset consisting of two points $\mathcal{D} = \{(x_1, y_1), (x_2, y_2)\}$, where x_n, y_n are scalars for $n = 1, 2$. Recall that the equation of a line in 2-dimensions can be written: $y = w_0 + w_1x$.

1. Write a system of linear equations determining the coefficients w_0, w_1 of the line passing through the points in our dataset \mathcal{D} and analytically solve for w_0, w_1 by solving this system of linear equations (i.e., using substitution). Please show your work.
2. Write the above system of linear equations in matrix notation, so that you have a matrix equation of the form $\mathbf{y} = \mathbf{X}\mathbf{w}$, where $\mathbf{y}, \mathbf{w} \in \mathbb{R}^2$ and $\mathbf{X} \in \mathbb{R}^{2 \times 2}$. For full credit, it suffices to write out what \mathbf{X} , \mathbf{y} , and \mathbf{w} should look like in terms of $x_1, x_2, y_1, y_2, w_0, w_1$, and any other necessary constants. Please show your reasoning and supporting intermediate steps.
3. Using properties of matrices, characterize exactly when an unique solution for $\mathbf{w} = (w_0 \ w_1)^T$ exists. In other words, what must be true about your dataset in order for there to be a unique solution for \mathbf{w} ? When the solution for \mathbf{w} exists (and is unique), write out, as a matrix expression, its analytical form (i.e., write \mathbf{w} in terms of \mathbf{X} and \mathbf{y}).

Hint: What special property must our \mathbf{X} matrix possess? What must be true about our data points in \mathcal{D} for this special property to hold?

4. Compute \mathbf{w} by hand via your matrix expression in (3) and compare it with your solution in (1). Do your final answers match? What is one advantage for phrasing the problem of fitting the model in terms of matrix notation?
5. In real-life, we often work with datasets that consist of hundreds, if not millions, of points. In such cases, does our analytical expression for \mathbf{w} that we derived in (3) apply immediately to the case when \mathcal{D} consists of more than two points? Why or why not?

Solution

1. We assume that $x_1 \neq x_2$. Then, the system of linear equations is

$$\begin{aligned} & \begin{cases} y_1 = w_0 + w_1 x_1 \\ y_2 = w_0 + w_1 x_2 \end{cases} \quad (1.1, 1.2) \\ \therefore & \begin{cases} x_2 y_1 = x_2 w_0 + w_1 x_1 x_2 \\ x_1 y_2 = x_1 w_0 + w_1 x_1 x_2 \end{cases} \end{aligned}$$

Subtracting gives

$$x_1 y_2 - x_2 y_1 = w_0 (x_1 - x_2)$$

Assuming that $x_1 \neq x_2$,

$$w_0 = \frac{x_1 y_2 - x_2 y_1}{x_1 - x_2}$$

Subtracting (1.1) from (1.2) gives

$$\begin{aligned} y_2 - y_1 &= w_1 (x_2 - x_1) \\ \therefore & \quad w_1 = \frac{y_2 - y_1}{x_2 - x_1} \end{aligned}$$

2. We have

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \end{bmatrix}; \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}; \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

Note that

$$\mathbf{X} \cdot \mathbf{w} = \begin{bmatrix} w_0 + w_1 x_1 \\ w_0 + w_1 x_2 \end{bmatrix}$$

corresponds to the RHS of (1.1,1.2) and thus

$$\mathbf{y} = \mathbf{X} \mathbf{w} \quad (1.3)$$

3. Note that (1.3) is a matrix equation. Thus, \mathbf{w} has a unique solution iff \mathbf{X} is invertible, which occurs iff our data points are unique. When this is true, we have

$$\mathbf{w} = \mathbf{X}^{-1} \mathbf{y} \quad (1.4)$$

4. Solving (1.4) gives

$$\begin{aligned} \mathbf{w} &= \frac{1}{x_2 - x_1} \begin{bmatrix} x_2 & -x_1 \\ -1 & 1 \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \\ &= \frac{1}{x_2 - x_1} \begin{bmatrix} x_2 y_1 - x_1 y_2 \\ -y_1 + y_2 \end{bmatrix} \\ \therefore \quad \mathbf{w} &= \begin{bmatrix} \frac{x_1 y_2 - x_2 y_1}{x_1 - x_2} \\ \frac{y_2 - y_1}{x_2 - x_1} \end{bmatrix} \end{aligned}$$

matching up with our answer in (1).

5. When we have a very large number of data points, the equation derived in (3) may not hold. This is because \mathbf{X} may have many more rows than it does columns. In such a case—unless all rows are linear combinations of two basis rows—we may have contradictory rows, and therefore no exact solution.

Problem 2 (Optimizing Objectives - Calculus Review)

In this class, we will write real-life goals we want our model to achieve into a mathematical expression and then find the optimal settings of the model that achieves these goals. The formal framework we will employ is that of mathematical optimization. Although the mathematics of optimization can be quite complex and deep, we have all encountered basic optimization problems in our first calculus class!

Learning Goals: In this problem, we will explore how to formalize real-life goals as mathematical optimization problems. We will also investigate under what conditions these optimization problems have solutions.

In her most recent work-from-home shopping spree, Nari decided to buy several house plants. *Her goal is to make them to grow as tall as possible.* After perusing the internet, Nari learns that the height y in mm of her Weeping Fig plant can be directly modeled as a function of the oz of water x she gives it each week:

$$y = -3x^2 + 72x + 70.$$

1. First, plot the height function. What does the plot tell you about the existence and uniqueness of a maximum plant height? Next, support your claim solely based on the form of the function.
2. Use calculus to find how many ounces of water per week Nari should give to her plant in order to maximize its height. With this much water, how tall will her plant grow?

Now suppose that Nari want to optimize both the amount of water x_1 (in oz) *and* the amount of direct sunlight x_2 (in hours) to provide for her plants. After extensive research, she decided that the height y (in mm) of her plants can be modeled as a two variable function:

$$y = f(x_1, x_2) = \exp(-(x_1 - 2)^2 - (x_2 - 1)^2)$$

3. Using `matplotlib`, visualize in 3D the height function as a function of x_1 and x_2 using the `plot_surface` utility for $(x_1, x_2) \in (0, 6) \times (0, 6)$. Then, determine the values of x_1 and x_2 that maximize plant height. Do these yield a global maximum?

Hint: You don't need to take any derivatives here; reasoning about the form of $f(x_1, x_2)$ suffices.

Solution

1. The plot appears to show that there exists a unique maximum plant height at some value of x between 10 and 15. Based on our knowledge of polynomials, we know that for $a = -3 < 0$, we will observe a concave curve with a unique, global maximum.
2. Set the first derivative to 0:

$$\begin{aligned}\frac{dy}{dx} &= -6x + 72 = 0 \\ \therefore \quad &\boxed{x = 12}\end{aligned}$$

We can confirm this is a maximum by taking the second derivative:

$$\frac{d^2y}{dx^2} = -6 < 0$$

Thus, Nari should give her plant 12 ounces of water per week. Also,

$$\begin{aligned}y &= -3(12^2) + 72(12) + 70 \\ \therefore \quad &\boxed{y = 502}\end{aligned}$$

Thus, with 12 ounces of water, Nari's plant will grow 502mm.

3. Re-writing out f we have

$$y = e^{-\{(x_1-2)^2+(x_2-1)^2\}}$$

We can see that $-\{(x_1-2)^2+(x_2-1)^2\} \leq 0$. Thus, y achieves its global maximum when $-\{(x_1-2)^2+(x_2-1)^2\} = 0$, which occurs iff $\boxed{x_1 = 2}$ and $\boxed{x_2 = 1}$. Moreover, in this case, $\boxed{y = e^0 = 1}$.

Problem 3 (Reasoning about Randomness - Probability and Statistics Review)

In this class, one of our main focuses is to model the unexpected variations in real-life phenomena using the formalism of random variables. In this problem, we will use random variables to model how much time it takes an USPS package processing system to process packages that arrive in a day.

Learning Goals: In this problem, you will analyze random variables and their distributions both analytically and computationally. You will also practice drawing connections between said analytical and computational conclusions.

Consider the following model for each package that arrives at the US Postal Service (USPS):

- Every package has a random size S (measured in in^3) and weight W (measured in pounds), with joint distribution

$$(S, W)^T \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \text{ with } \boldsymbol{\mu} = \begin{bmatrix} 120 \\ 4 \end{bmatrix} \text{ and } \boldsymbol{\Sigma} = \begin{bmatrix} 1.5 & 1 \\ 1 & 1.5 \end{bmatrix}.$$

- The size and weight of each package is independent of those of all the other packages.
- Processing time T (in seconds) for each package is given by $T = 60 + 0.6W + 0.2S + \epsilon$, where ϵ is an independent random noise variable with Gaussian distribution $\epsilon \sim \mathcal{N}(0, 5)$.

1. Perform the following tasks:

- (a) Give one reason for why the Gaussian distribution may not be appropriate for modeling the size and weight of packages.
- (b) Empirically estimate the most likely combination of size and weight of a package by sampling 500 times from the joint distribution of S and W and generating a bivariate histogram of your S and W samples. A visual inspection is sufficient – you do not need to be incredibly precise. How close are these empirical values to the theoretical expected size and expected weight of a package, according to the given Bivariate Gaussian distribution?

Hint: For this part, you may find the `multivariate_normal` module from `scipy.stats` especially helpful. You may also find the `seaborn.histplot` function quite helpful.

2. For 1001 evenly-spaced values of W between 0 and 10, plot W versus the joint Bivariate Gaussian PDF $p(W, S)$ with S fixed at $S = 118$. Repeat this procedure for S fixed at $S = 122$. Comparing these two PDF plots, what can you say about the correlation of random variables S and W ?
3. Because T is a linear combination of random variables, it itself is a random variable. Using properties of expectations and variance, please compute $\mathbb{E}(T)$ and $\text{Var}(T)$ analytically.
4. Define N to be the number of packages that arrive today, and suppose that packages that weigh less than 4 pounds are considered fragile. Conditional on $N = n$, what is the name and PMF of the distribution of the number of fragile packages that arrive today?
5. Now suppose that $N = \sum_{h=1}^{24} P_h$, where the P_h are independent and identically distributed as $\text{Pois}(\lambda = 3)$. Then define $T^* = \sum_{i=1}^N T_i$ as the *total* amount of time it takes to process *all* these packages, where T_i follows the distribution of T that we previously defined for each package.
 - (a) Write a function to simulate draws from the distribution of T^* .
 - (b) Using your function, empirically estimate the mean and standard deviation of T^* by generating 1000 samples from the distribution of T^* .

Solution

1.
 - a. The support of a Gaussian distribution is $(-\infty, \infty)$. Thus, a Gaussian model implies a positive probability of $S, W < 0$, which does not make any physical sense, as size and weight cannot be negative.
 - b. A visual inspection of our plot shows high density around the values 4 and 120, corresponding to the expected values of W and S respectively. This shows that the empirical values of the average size and weight of the package correspond well with the theoretical values.
2. This plot shows us that S and W are positively correlated: For a fixed, larger value of $S = 122$ (versus $S = 118$), we see that the conditional distribution of W (given S) is shifted to the right.
3. We have

$$T = 60 + 0.6W + 0.2S + \epsilon$$

Thus,

$$\begin{aligned}\mathbb{E}(T) &= 60 + 0.6\mathbb{E}(W) + 0.2\mathbb{E}(S) + \mathbb{E}(\epsilon) \\ &= 60 + 0.6(4) + 0.2(120) + 0\end{aligned}$$

$$\therefore \boxed{\mathbb{E}(T) = 86.4}$$

and

$$\begin{aligned}\text{Var}(T) &= 0.6^2\text{Var}(W) + 0.2^2\text{Var}(S) + 2(0.6)(0.2)\text{Cov}(W, S) + \text{Var}(\epsilon) \\ &= 0.36(1.5) + 0.04(1.5) + 0.24(1) + 5\end{aligned}$$

$$\therefore \boxed{\text{Var}(T) = 5.84}$$

4. Let F denote the number of fragile packages that arrive today. Note that any given package i has independent probability $P(W_i < 4) = \frac{1}{2}$ of being fragile. Thus, given a fixed number of packages, each of which is fragile with probability $\frac{1}{2}$, we are asking how many of these is fragile. This is the story of the Binomial distribution. Thus,

$$\boxed{F|(N = n) \sim \text{Bin}\left(n, \frac{1}{2}\right)}$$

and the PMF is

$$\boxed{P(F = f) = \binom{n}{f} \left(\frac{1}{2}\right)^n ; \quad f \in \{0, 1, \dots, n\}}$$

5. Empirical estimates:

```
mean = 6235.987638799891
stdev = 726.8731836729792
```


Problem 4 (Implementing a Linear Regression - Coding Review)

In this class, we will bridge theory and practice through implementing the methods that we cover from scratch. In this problem, we follow up on Problem 1 through exploring a more practical version of linear regression (fitting a linear model). Namely, we use ordinary least squares (OLS) to estimate a *line of best fit* rather than a perfect fit to our data. Note that the focus of this problem is on coding rather than math—we will cover the relevant theory in much more depth during the course.

Learning Goals: In this problem, you will gain experience with the procedure of modeling real-world data. You will also get useful practice with debugging and writing clean, efficient code in Python.

Steve is a fictional CS 1810 TF giving a live demo of how to fit a linear regression. However, he quickly realizes that coding live in front of an audience isn't for the faint of heart. As a star student, you will help him with his code. Just like Problem 1, the demo uses a 2-D dataset, so that the goal is to model the relationship between the x and y coordinates. The data are stored in the `data` variable, with the first column corresponding to the x -coordinate and the second corresponding to the y -coordinate.

1. Using the provided data, Steve has defined variables `y` and `x` corresponding to the respective coordinates. What is wrong with his current code? Fix the code and then plot the data. Does there appear to be a linear trend?
2. Steve then defines a new variable `X`, which is meant to resemble \mathbf{X} from Problem 1. Specifically, `X` is supposed to have one column of all ones (recall that this allows us to fit an intercept) and one column which is just `x`, the x -coordinates. However, he realizes that his code yields the wrong shape for `X`. What's going on here? Fix the code and then report what `y.shape` and `X.shape` are. Why is there no second coordinate in the output for `y.shape`?

Hint: check the documentation for `np.hstack`.

3. Steve takes a much-needed break from coding to give the following high level overview of linear regression: given a target (response) \mathbf{y} and features (predictors) \mathbf{X} , the goal of linear regression is to find weights \mathbf{w} such that $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$ closely approximates the true data \mathbf{y} . In OLS, we estimate \mathbf{w} to be

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Steve skips over the derivation of the result but assures you that you will learn it later in the course. What should the shape of $\hat{\mathbf{w}}$ be in Steve's demo?

4. Having walked through the idea of linear regression, Steve then attempts to implement a `LinearRegression` class. He correctly identifies that we need 3 components: a constructor, a `fit` function for computing $\hat{\mathbf{w}}$ from the data, and a `predict` function for computing the estimate $\mathbf{X}\hat{\mathbf{w}}$. However, he realizes that there is something wrong (meaning logic or syntax) with at least one of these components. Please point out the issues, fix them, and include the plot of the fitted line.
5. As his final act for the day, Steve introduces the Mean Squared Error (MSE) loss function:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

This captures how well the outputs of our model, $\hat{\mathbf{y}}$, fit the actual data \mathbf{y} . Steve manages to correctly implement an MSE computation! However, you realize that he can vectorize his code to make it faster, meaning that he can directly compute the MSE from NumPy arrays without using any for loops. Implement the vectorized MSE and write down the corresponding mathematical expression, which should directly be in terms of the vectors \mathbf{y} and $\hat{\mathbf{y}}$ rather than their components.

Solution

1. `data` is an array of (x, y) pairs. Thus, Steve is just taking the first 2 rows of his data.
2. Steve used `np.hstack` when he probably meant to use `np.column_stack`. He wanted to “pair up” the ones with the x values, but ended up just concatenating all the x values onto the end of the ones array.

```
y.shape = (100,)
X.shape = (100, 2)
```

There is no second coordinate for `y.shape` because `y` is a vector and is a 1D ndarray rather than a 2D ndarray (as is `X`).

3. The shape of `X` is 100×2 so that of $\mathbf{X}^T \mathbf{X}$ is 2×2 and $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is 2×100 . Finally, the dimensions of $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ should be 2×1 since those of `y` are 100×1 .
4. See code.
5. Corresponding mathematical expression in terms of vectors is

$$\text{MSE} = \frac{1}{n} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}})$$