

New York University Tandon School of Engineering
Computer Science and Engineering

CS-GY 6923: Written Homework 1.

Due Thursday, February 10th, 2022, 11:59pm.

Discussion with other students is allowed for this problem set, but solutions must be written-up individually.

For just this first problem set, 10% extra credit will be given if solutions are typewritten (using LaTeX, Markdown, or some other mathematical formatting program).

Problem 1: Practice Minimizing a Loss Function (10pts)

Consider a linear model of the form:

$$f_{\beta}(x) = \beta x,$$

which is the same as the linear model we saw in class, but with the intercept forced to zero. Such models are used when we want to force the predicted value $f_{\beta}(x) = 0$ when $x = 0$. For example, if we are modeling $y =$ output power of a motor vs. $x =$ the input power, we would expect $x = 0 \Rightarrow y = 0$.

- Given data $(x_1, y_1), \dots, (x_n, y_n)$, write the equation for a loss function which measures prediction accuracy using the sum-of-squared distances between the predicted values and target values.
- Derive an expression for the β that minimizes this loss function. Do you get the same expression that we got for β_1 in the full linear model?

Problem 2: Machine Learning Does Averages (15pts)

Suppose we have data $y_1, \dots, y_n \in \mathbb{R}$ and we want to choose a single value $m \in \mathbb{R}$ which is “most representative” of our dataset. This is sometimes called the “central tendency” problem in statistics. A machine learning approach to this problem would measure how representative m is of the data using a loss function.

- Consider the loss function $L(m) = \sum_{i=1}^n (y_i - m)^2$. Show that $L(m)$ is minimized by setting $m = \bar{y}$, where $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ is the **mean** of our data.
- Consider the loss function $L(m) = \max_i |y_i - m|$. What value of m minimizes this loss? **Hint:** Using derivatives will not help here – try just thinking about the minimization problem directly.
- Consider the loss function $L(m) = \sum_{i=1}^n |y_i - m|$. Prove that $L(m)$ is minimized by setting m to the **median** of the data. **Hint:** This question is harder than the previous two and takes some creativity! Again derivatives might not be helpful.
- In a few short sentences, discuss when you might prefer each of the three losses above. Is the median typically considered a more “robust” measure of central tendency than the mean? Why?

Problem 3: Practice with Non-linear Transformations. (8 pts)

A medical researcher wants to model, $f(t)$, the concentration of some chemical in the blood over time. She believes the concentration should decay exponentially in that

$$f(t) = z_0 e^{-\alpha t}, \tag{1}$$

for some parameters z_0 and α . To confirm this model, and to estimate the parameters z_0, α , she collects a large number of time-stamped samples (t_i, c_i) , $i = 1, \dots, n$, where c_i is the measured concentration at time t_i . Unfortunately, the model (1) is non-linear, so she can’t directly apply the linear regression formula to estimate z_0 and α .

- Taking logarithms, show that we can transform our training data so that the conjectured relationship between predictor and target variables is in fact linear.
- Write pseudocode (or actual Python) for how you might estimate z_0 and α using this transformation.

Problem 4: Practice With Gradients (8pts)

For $\mathbf{X} \in \mathbb{R}^{n \times d}$ and target vector $\mathbf{y} \in \mathbb{R}^n$, consider fitting a linear model of the form:

$$f_{\boldsymbol{\beta}}(\mathbf{x}) = \mathbf{X}\boldsymbol{\beta}$$

under the so-called ℓ_p loss: $L_p(\boldsymbol{\beta}) = \|\mathbf{y} - f_{\boldsymbol{\beta}}(\mathbf{x})\|_p^p$. Here $\|\cdot\|_p^p$ denotes the ℓ_p norm raised to the p power. I.e. for any even integer $p = 2, 4, 6, \dots$

$$\|\mathbf{z}\|_p^p = \sum_{i=1}^n z_i^p$$

- Derive an expression for $\nabla g(\mathbf{z})$ where $g(\mathbf{z}) = \|\mathbf{z}\|_p^p$.
- Derive an expression for $\nabla L_p(\boldsymbol{\beta})$.

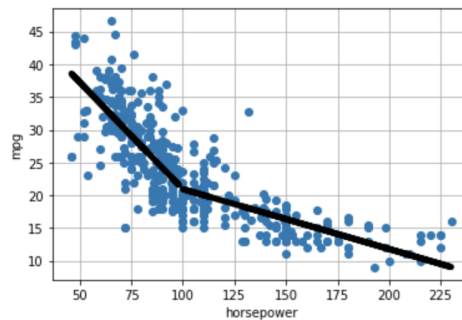
Problem 5: Piecewise Linear Regression via Feature Transformations (15pts)

Your goal is to fit a *piecewise* linear model to a single variate dataset of the form $(x_1, y_1), \dots, (x_n, y_n)$ where all values are scalars. We will only use two pieces. In other words, for some known value λ ,

$$f(x_i) = \begin{cases} a_1 + s_1 x_i & \text{for } x_i < \lambda \\ a_2 + s_2 x_i & \text{for } x_i \geq \lambda \end{cases}$$

with the additional **constraint** that $a_1 + s_1 \lambda = a_2 + s_2 \lambda$. This constraint ensures that our two linear models actually “meet” at $x = \lambda$, which means we get a continuous prediction function.

For example, when $\lambda = 100$, a piecewise linear fit for our MPG data might look like:



- Show that this model is equivalent to the following **unconstrained** model:

$$f(x_i) = \begin{cases} a_1 + s_1 x_i & \text{for } x_i < \lambda \\ a_1 + s_1 \lambda - s_2 \lambda + s_2 x_i & \text{for } x_i \geq \lambda \end{cases}$$

- Show how to fit an optimal f under the squared loss using an algorithm for multiple linear regression. In particular, your approach should:
 - Transform the input data to form a data matrix \mathbf{X} with multiple columns.
 - Use a multiple regression algorithm to find the $\boldsymbol{\beta}$ which minimizes $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$.
 - Extract from the optimal $\boldsymbol{\beta}$ optimal values for a_1, s_1, s_2 .

You need to describe 1) a correct data transformation and 2) a correct mapping from $\boldsymbol{\beta}$ to a_1, s_1, s_2 . **Note that in our model λ is known. It is not a model parameter which needs to be optimized.**

- Implement your algorithm in Python and apply it to the dataset from `demo_auto_mpg.ipynb`. Produce a piecewise linear fit for MPG as a function of Horsepower using the value $\lambda = 100$. Plot the result. You can attach a Jupyter notebook to your submission, or simply include the printed code and plot.

- (d) **(3pts bonus)** Modify your approach to handle the case when λ is unknown. Again obtain a fit for MPG vs. horsepower. What value of λ gives the optimal fit? Include any modified code and a plot of your result.