# 1. Loss Functions for Linear Regression

Part 1:

a.  This function could be a good choice for regression. However, it is more sensitive to any outliers compared to the standard squared loss function. Sub-gradient descent is applicable here since it's an absolute value function.

b.  This function would not be a good choice for regression. Since its cubic, it would not be symmetric at 0 which could lead to great error. Gradient descent could be used here.

c.  This function would also not be a good choice. This is an exponential function which means it will continue to grow exponentially. Because of this, the optimization of this function would be unstable. We can use gradient descent here as well.

d.  This function is also not a good choice. This function is used for hinge loss which deals with classification, not regression. We can use sub-gradient descent here as well since the max function is not differential at 0.

Part 2:

a.
$$\partial L_i/\partial w = 3|f(w,b,x_i) - y_i|^2 * \text{sign}(f(w,b,x_i) - y_i) * x_i$$
$$\partial L_i/\partial b = 3|f(w,b,x_i) - y_i|^2 * \text{sign}(f(w,b,x_i) - y_i)$$

b.
$$\partial L_i/\partial w = 3(f(w,b,x_i) - y_i^2 ) * x_i$$
$$\partial L_i/\partial b = 3(f(w,b,x_i) - y_i^2)$$

c.
$$\partial L_i/\partial w = \exp(f(w,b,x_i) - y_i) * x_i$$
$$\partial L_i/\partial b = \exp(f(w,b,x_i) - y_i)$$

d.
When $1 - y_i * f(w,b,x_i) > 0$:

$$\partial L_i/\partial w = -y_i x_i$$
$$\partial L_i/\partial b = -y_i$$

When $1 - y_i * f(w,b,x_i) < 0$:

$$\partial L_i / \partial w = 0$$
$$\partial L_i / \partial b = 0$$

## 2. <u>Loss Functions for Classification</u>

<u>Part 1:</u>

    a. This loss function is a good choice for classification. This function is better than the original hinge loss function as it will penalize incorrect classifications more strongly. Gradient descent can be used when the function is not equal to 0.

    b. This loss function is not a good choice for classification. This function works best for regression as it will penalize correct classifications whose prediction is too far from the true value. Gradient descent can be used here as the function is differential.

    c. This loss function is not a good choice for classification since it is an exponential function. Since its exponential, it can lead to unstable training. We can use gradient descent here since its also differential.

    d. This loss function is a good choice for classification. We can use gradient descent here as well.

<u>Part 2:</u>

a.

If $1 - y_i * f(w,b,x_i) <= 0$, $L_i(w,b) = 0$ with the gradient being 0 as well

If $1 - y_i * f(w,b,x_i) > 0$, $L_i(w,b) = ( 1 - y_i * f(w,b,x_i))^2$

$\partial L_i / \partial w = 2(1 - y_i * f(w,b,x_i))( -y_i x_i) = -2/n \ (y_i x_i(1 - y_i * f(w,b,x_i)))$

$\partial L_i / \partial b = 2(1 - y_i * f(w,b,x_i))( -y_i) = -2/(n y_i(1 - y_i * f(w,b,x_i)))$

b.

$\partial L_i / \partial w = 4(y_i - f(w,b,x_i))^3(x_i) = 4/n \ (x_i(y_i - f(w,b,x_i))^3)$

$\partial L_i / \partial b = 4(y_i - f(w,b,x_i))^3(-1) = 4/n \ ((y_i - f(w,b,x_i))^3)$

c.

$\partial L_i / \partial w = \exp(f(w,b,x_i) - y_i) * x_i$

$\partial L_i / \partial b = \exp(f(w,b,x_i) - y_i)$

d.

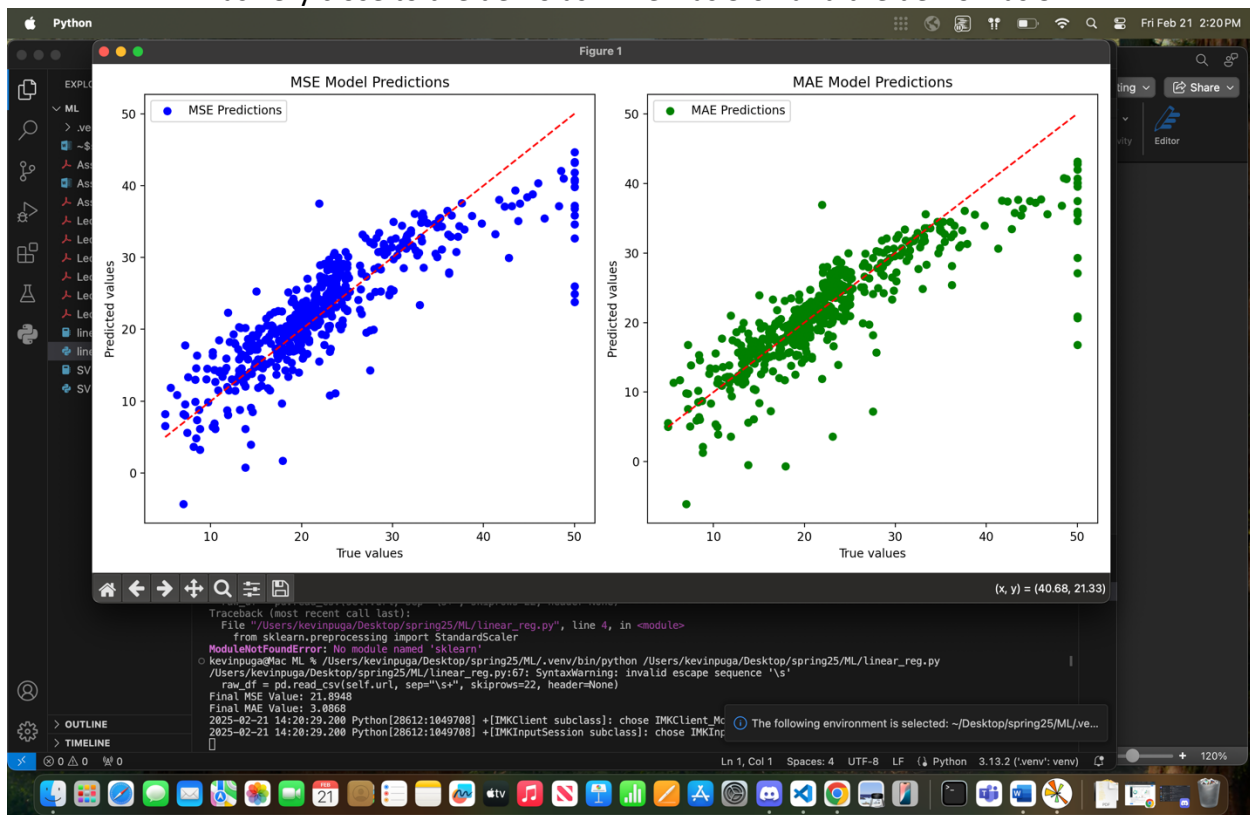$\partial L_i / \partial w = \exp(-y_i * f(w,b,x_i))(y_i x_i) = -y_i x_i * \exp(-y_i * f(w,b,x_i))$

$\partial L_i / \partial b = \exp(-y_i * f(w,b,x_i))(-y_i) = -y_i * \exp(-y_i * f(w,b,x_i))$

## 3. Polynomial and Higher Order features

a. The problem with trying to fit a linear regression model without polynomial features is that it will lead to a poor fit and inaccurate predictions, making the model underfitting.

b. The polynomial features I would use are all powers currently present withing the function ( cubic, square, etc.) In this case, the optimal weight vector would be the set of coefficients that minimizes the error between the predicted set and the actual set of values.

## 4. Implement Linear Regression by Scratch

a. When comparing my values to the demo, my mse value was lower than that of the demo. I got 21.9 when the demo had a value of 24.3 or so. In terms of mae, I was very close to the demo as mine was 3.07 and the demo was 3.2.



## 5. Implement SVM and Perceptrons from Scratch

a. For the perceptron loss graph, the margin that divides both plot points is going through the yellow data points. For the SVM loss graph, the margin is higher and doesn't go through any data points.