

Loss & Cost Function & Gradient Descent Algorithm

28 July 2024 21:22

What is a cost function? How is it different from loss function?

Loss function

* A loss function, also known as error function, measures the difference between predicted value and actual value.

- Loss function is used when we refer the error for a single training example (a single row)
- Instance level - calculated for a single training example.

Regression: MSE: Mean Squared error:

$$L(y, \hat{y}_i) = (y_i - \hat{y}_i)^2 \text{ for some } i=10$$

Classification: Cross-Entropy Loss

$$y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

predicted probabilities from logistic regression.
actual for some $i=10$

Cost Function

A cost function aggregates the losses over the entire training dataset or batch of training examples
training datapoints

- Dataset level: summarizes the performance of the model over the entire training dataset

For Regression:

Mean Squared Error (MSE) Formula: $J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$h_{\theta}(x^{(i)})$: Predicted value: \hat{y}
 $y^{(i)}$: Actual value: y

Mean Absolute Error (MAE) Formula: $J(\theta) = \frac{1}{m} \sum_{i=1}^m |h_{\theta}(x^{(i)}) - y^{(i)}|$ cost functions for the entire dataset

Mean Absolute Error (MAE) Formula: $J(\theta) = \frac{1}{m} \sum_{i=1}^m |h_{\theta}(x^{(i)}) - y^{(i)}|$ } cost functions for the entire dataset

$|x|$: modulus function

$|0| \rightarrow 0$

For classification

$$y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left\{ y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right\} : \text{Binary Classification}$$

Task why \ominus ve sign in the above equation?

Gradient Descent Algorithm (GDA)

slope

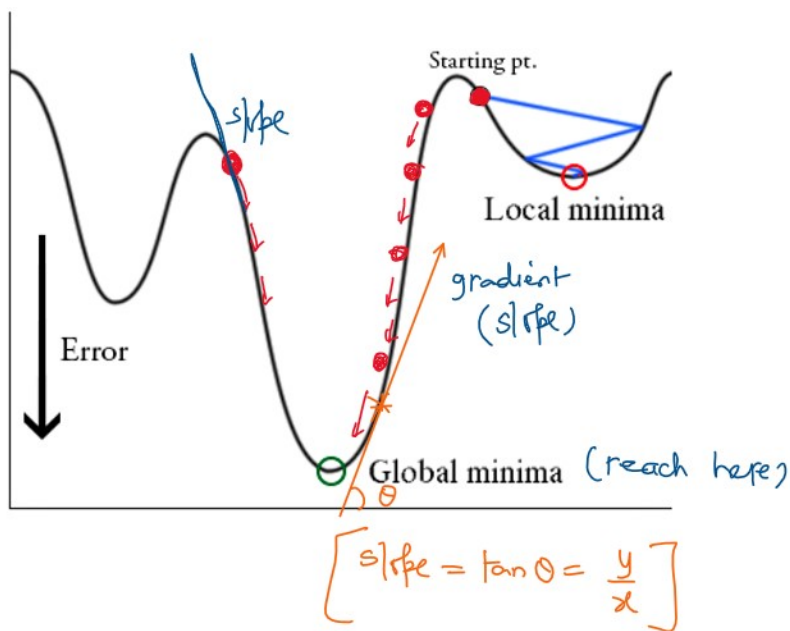
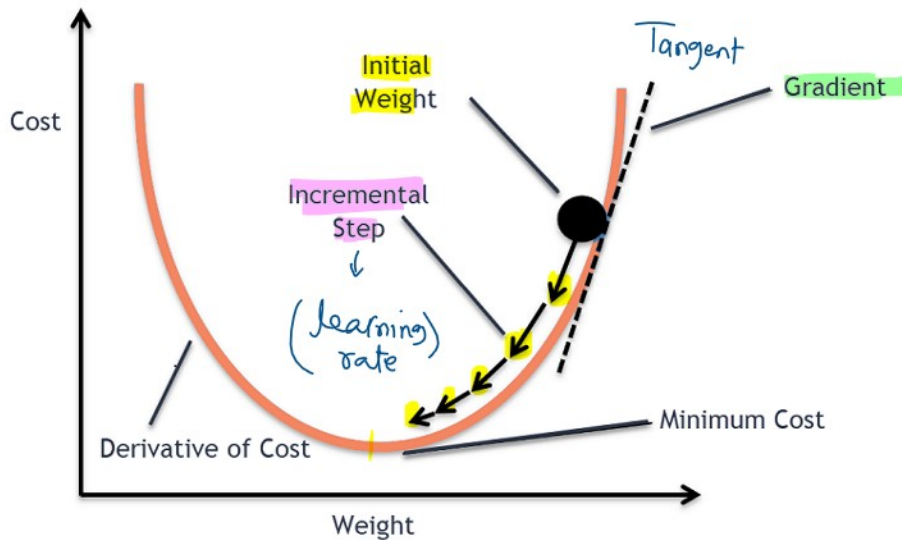
fall down

⊖ op



Follow the path of the steepest descent, taking steps in the direction that slope (tangent)

in the direction ^D that ^{slope (tangent)}
 decreases the slope and brings ^{gradient}
 you closer to the land | plains.



QDA helps the model to find the optimal set of
 parameters (weights and biases) by iteratively
 adjusting them in the opposite direction of the gradient.

GDA

repeat until convergence {

$$\theta_j := \underbrace{\theta_j}_{\substack{\text{initial} \\ \text{value}}} - \underbrace{\alpha}_{\substack{\text{learning rate}}} \underbrace{\frac{\partial}{\partial \theta_j} [J(\theta_0, \theta_1)]}_{\substack{\text{gradient/slope} \\ \text{for } (j=0 \text{ and } j=1)}}$$

cost function

}

For $j=0$

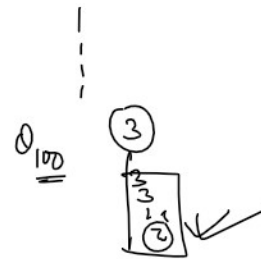
$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} [J(\theta_0, \theta_1)]$$

$$\theta_0 = 10$$

$$\theta_1 = 9$$

$$\theta_2 = 8$$

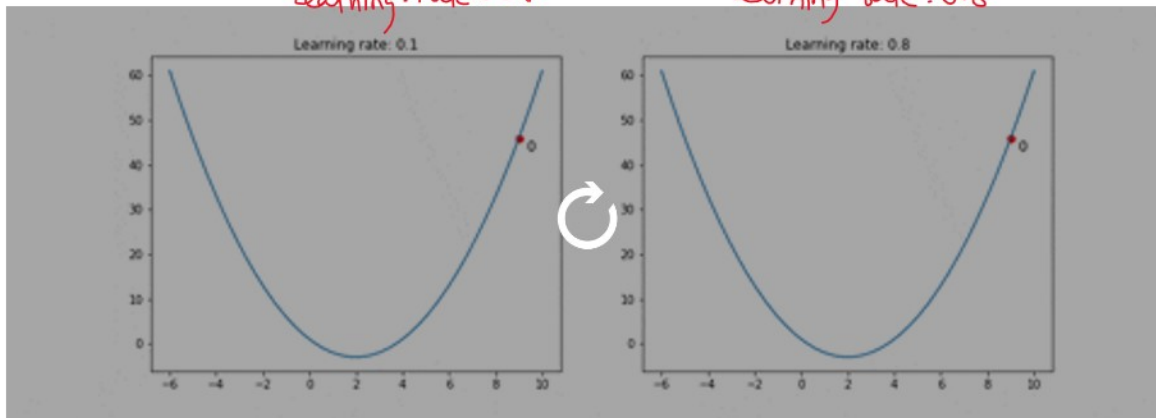
$$\theta_3 = 7.5$$

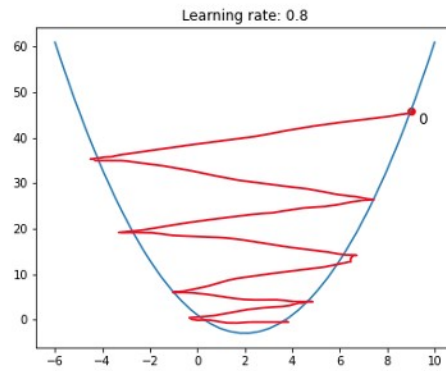
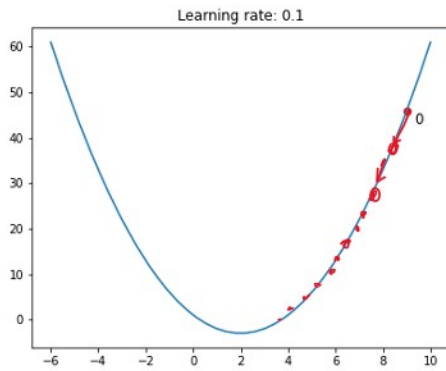


[Animated GIF - Find & Share on GIPHY](#)

learning rate: 0.1

learning rate: 0.8





$$\theta_0)_{\text{updated}} \triangleq \theta_0)_{\text{last}} \rightarrow \text{GPA Converge}$$

Note: If the learning rate is too slow ($\alpha = 0.001$), model takes quite a lot of time to converge (training time is unrealistically high.)

If the learning rate is too high (say $\alpha = 0.8$), model might overshoot the minima and keep bouncing without reaching the minima

$\alpha = 0.1$ to 0.2 : as per the best practices