



Stochastic Gradient Descent

Mohammad Azodi, Shahid Beheshti University

27/07/2020

Optimization for data science
Teacher: Dr. Bijan Ahmadi

Stochastic Gradient Descent

A custom implementation of Stochastic Gradient Descent for Linear Regression that optimizes the weights $\mathbf{W}(\mathbf{i})$ of each component and the bias \mathbf{b} term.

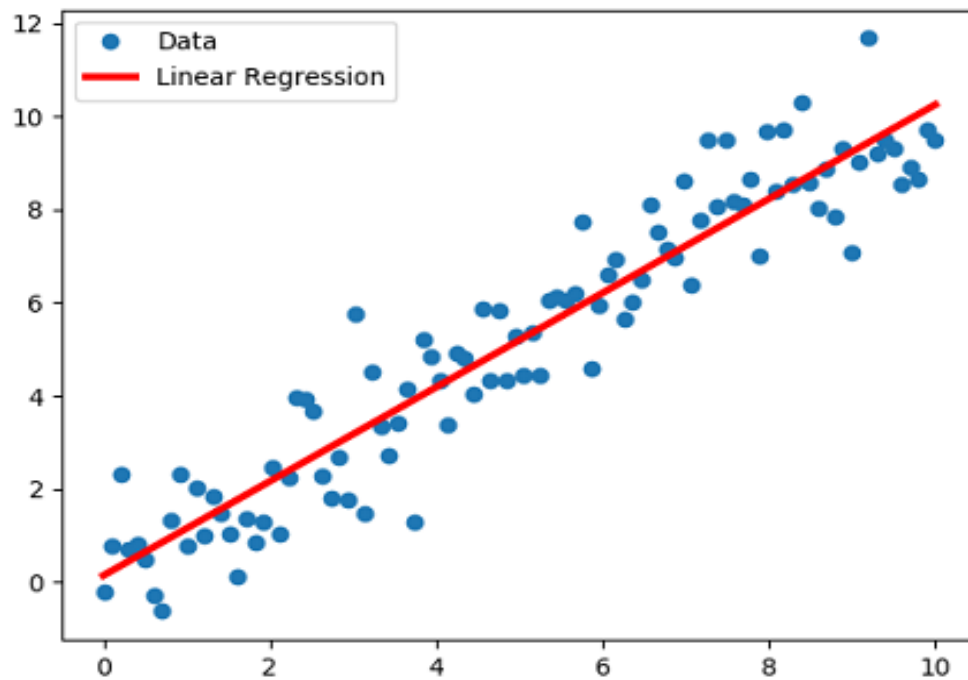
Gradient Descent:

Gradient Descent is an iterative optimization algorithm that can be used to converge to an optimal value easily with the use of modern computational power.

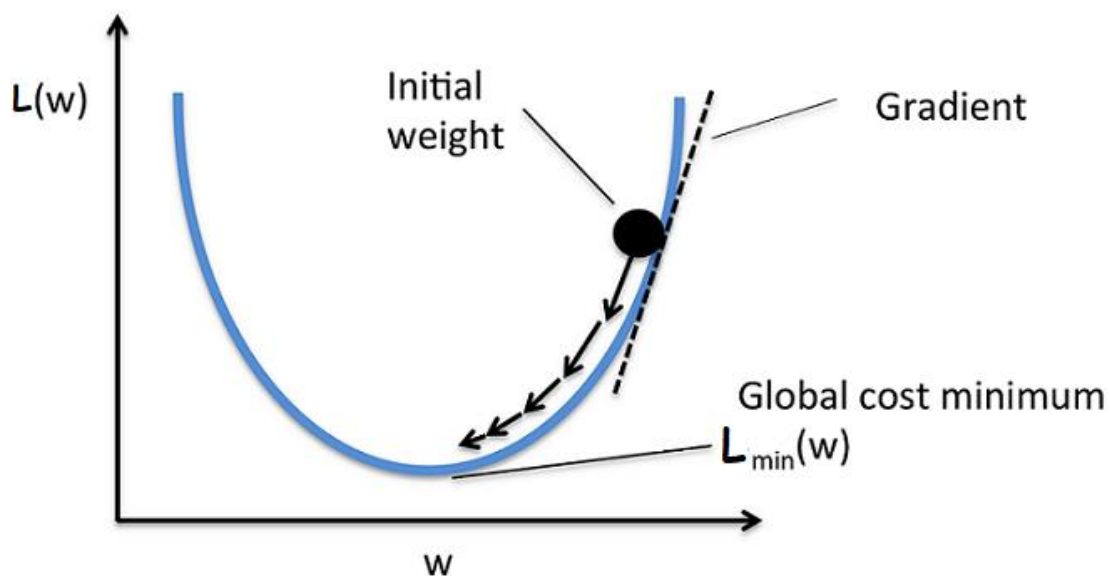
- The update equation for every iteration is, $\mathbf{x}(\mathbf{i}) = \mathbf{x}(\mathbf{i}-1) - \mathbf{r} * [\mathbf{df/dx}]_{\mathbf{x}(\mathbf{i}-1)} ; \mathbf{i} : 1 \rightarrow \mathbf{n}$.
- It mainly depends on the **learning rate** (or) **step size** denoted by " \mathbf{r} ".
- The **learning rate** tells how fast to converge to the optimal value. So giving a right value of " \mathbf{r} " matters.
- If the right value of " \mathbf{r} " ain't given then the updation might jump over optimal value(min) and we'll not be converging at the right solution. Hence need to check with different values of " \mathbf{r} ".

Let us consider a simple linear regression,

- **Objective** : Find the line/plane that best fits the data as,



- The dataset is $D = \langle \mathbf{x}_i, y_i \rangle$; $\mathbf{x}_i \in \mathbb{R}^d$; $y_i \in \mathbb{R}$. Here d is for # of dimensions.
- We can find the line/plane that fits the real values data of form $y_i = \mathbf{W} \cdot \mathbf{T} \cdot \mathbf{x}_i + \mathbf{B}$ for given \mathbf{x}_i . Here \mathbf{B} is for bias.
- Then we've **Mean Squared Error(MSE)** = $\sum ([y_i - (\mathbf{W} \cdot \mathbf{T} \cdot \mathbf{x}_i + \mathbf{B})]^2) / n$; $i: 1 \rightarrow n$.
- The optimal weight vector will be $\mathbf{W}^* = \text{argmin}(\mathbf{W}) \sum ([y_i - (\mathbf{W} \cdot \mathbf{T} \cdot \mathbf{x}_i + \mathbf{B})]^2) / n$ i.e the one which gives minimum sum of squared errors.



We can write the optimization problem $\mathbf{W}^* = \text{argmin}(\mathbf{W}) \sum ([\mathbf{y}_i - (\mathbf{W}.\mathbf{T}^*\mathbf{x}_i + \mathbf{B})]^2) / n$ as,

$$\mathbf{L}(\mathbf{W}) = \sum ([\mathbf{y}_i - (\mathbf{W}.\mathbf{T}^*\mathbf{x}_i + \mathbf{B})]^2) / n ; i : 1 \rightarrow n.$$

Then the vector differentiation or grad of $\mathbf{L}(\mathbf{w})$ is $\nabla_{\mathbf{w}} \mathbf{L} = \sum \{ 2*(\mathbf{y}_i - (\mathbf{W}.\mathbf{T}^*\mathbf{x}_i + \mathbf{B}))(-\mathbf{x}_i) \} / n$

