

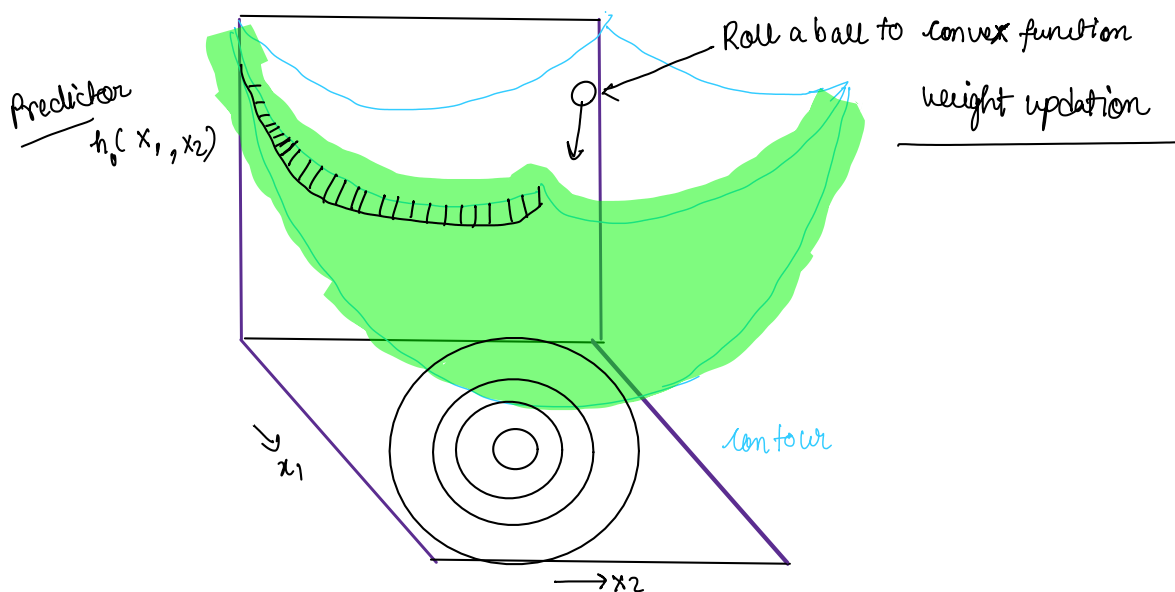
Gradient Descent algorithm as optimizer

Monday, 11 March 2024 10:50 AM

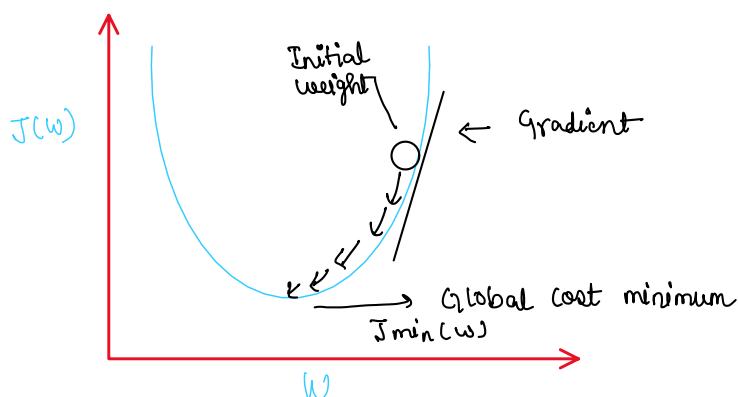
We have studied Normal Equation and its matrix form this is analytical method, which is good when the number of features X are not large

for large number of features, we need the algorithmic tool based on Gradient Descent Algorithm, it is called optimizer.

when we apply on a convex quadratic functions, the global optimum solution is guaranteed



Example of a convex function



$$J(w) = \frac{1}{2m} \sum_{i=1}^m (h_w(x_i^{(i)}) - y^{(i)})^2$$

$$\text{for } i=1 \rightarrow J(w) = \frac{1}{2} (h_w(x_1) - y)^2 \rightarrow \textcircled{1}$$

for m normal distribution

for n no. of features

$$h_w(x)_{j=1,2,\dots,n} = w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$= w_j^T x_j$$

$$w_j = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$x_j = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \rightarrow (2)$$

for optimization, we need to find

$$\frac{\partial J(w)}{\partial w_j}$$

Now let us substitute $S = h_w(x_j) - y \rightarrow (3)$

$$\text{for (2)} \rightarrow S = w_j^T x_j - y$$

$$\boxed{\frac{\partial S}{\partial w_j} = x_j \rightarrow (4)}$$

$$\text{from (1) \& (3)} \rightarrow J(w) = \frac{1}{2} S^2$$

$$\text{Hence, } \frac{\partial J(w)}{\partial S} = \frac{1}{2} * 2 * S = S \rightarrow (5)$$

$$\frac{\partial J(w)}{\partial w_j} = \frac{\partial J(w)}{\partial S} \cdot \frac{\partial S}{\partial w_j} = S \cdot x_j \text{ (from (4) and (5))}$$

$$= (h_w(x_j) - y) \cdot x_j$$

$$\text{for 'm' samples } \boxed{\frac{\partial J(w)}{\partial w_j} = \sum_{i=1}^m (h_w(x_j^{(i)}) - y^{(i)}) \cdot x_j^{(i)}}$$

(i) \leftarrow feature size
(i) = 0, 1, 2, ..., n
(j) \leftarrow no. of samples

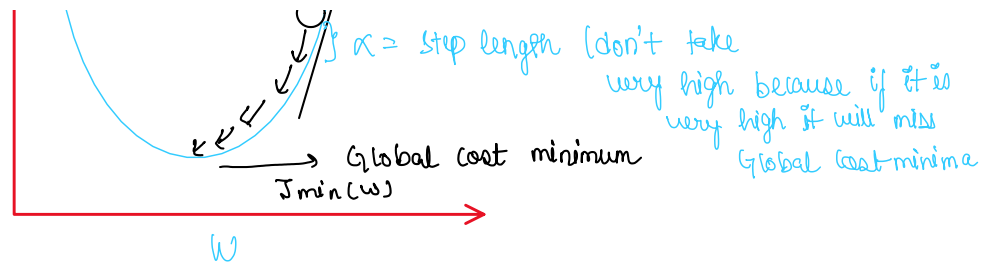
for Gradient Descent Algorithm

$$w_j^{\wedge} = w_j^{\circ} - \alpha \cdot \frac{1}{m} \frac{\partial J(w)}{\partial w_j}$$

hyper parameter
 α = learning rate
plays a important role
how model parameter
will converge

$$w_j^{\wedge} = w_j^{\circ} - \alpha \frac{1}{m} \sum_{i=1}^m (h_w(x_j^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$





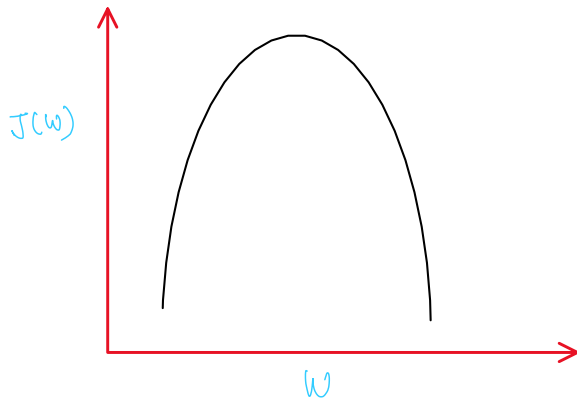
for Gradient **Ascent** Algorithm

$J(w)$ = Profit function ($w_1, w_2, w_3 \dots w_n$)

$$w_j^o = w_j^o + \alpha \cdot \frac{1}{m} \frac{\partial J(w)}{\partial w_j^o}$$

hyper parameter
 α = learning Rate
 plays a Important Role
 how model parameter will converge

$$w_j^o = w_j^o + \alpha \frac{1}{m} \sum_{i=1}^m (h_w(x_j^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$



Batch / vanilla Gradient Descent Algorithm

Repeat {

$$w_j^o = w_j^o - \alpha \frac{1}{m} \sum_{i=1}^m (h_w(x_j^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

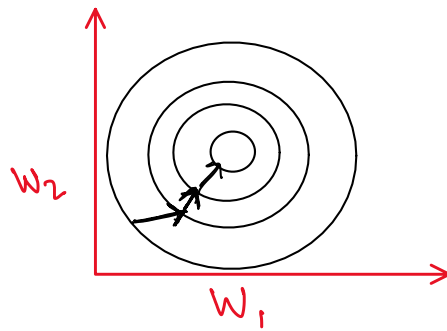
$$w_0 := w_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_w(x_j^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

$$w_1 := w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_w(x_j^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$$

$$\vdots$$

$$w_n := w_n - \alpha \frac{1}{n} \sum_{i=1}^n (h_w(x) - y) \cdot x_n$$

}



contour plot



GDA optimizer

m is very large
 n is also large 1000 or higher

↓ We shift to

Stochastic Gradient Descent Algorithm

- ① Randomly shuffle training samples.
- ② Repeat {

for $i = 1, 2, \dots, m$ {

$$w_j := w_j - \alpha (h_w(x^{(i)}) - \tilde{y}^{(i)}) \cdot x_j^{(i)} \leftarrow \frac{\partial J(w)}{\partial w_j}$$

$$w_0 := w_0 - \alpha (h_w(x^{(i)}) - \tilde{y}^{(i)}) \cdot x_0^{(i)}$$

$$w_1 := w_1 - \alpha (h_w(x^{(i)}) - \tilde{y}^{(i)}) \cdot x_1^{(i)}$$

⋮

$$w_n := w_n - \alpha (h_w(x^{(i)}) - \tilde{y}^{(i)}) \cdot x_n^{(i)}$$

}

}

Mini-Batch Gradient Descent Algorithm

- ① Create a Batch size, b
- ② Randomly shuffle the batches

Repeat {

$$\text{for } i = 1, 2, \dots, b \quad \left\{ \frac{\partial J(w)}{\partial w_j} \right.$$
$$w_j = w_j - \frac{\alpha}{b} \sum_{i=1}^b (h_w(x)^{(i)} - y^{(i)}) \cdot x_j^{(i)}$$

$$w_0 = w_0 - \frac{\alpha}{b} \sum_{i=1}^b (h_w(x)^{(i)} - y^{(i)}) \cdot x_0^{(i)}$$

$$w_1 = w_1 - \frac{\alpha}{b} \sum_{i=1}^b (h_w(x)^{(i)} - y^{(i)}) \cdot x_1^{(i)}$$

\vdots

$$w_n = w_n - \frac{\alpha}{b} \sum_{i=1}^b (h_w(x)^{(i)} - y^{(i)}) \cdot x_n^{(i)}$$

}

b = hyper parameter

How to choose batch size - b ?

b may be 16, 32, 64, 128
 $2^4, 2^5, 2^6, 2^7$

