

COL 7 #4

Aug 28, 2021

Yesterday:-

Linear Regression

optimization / Finding the local optima.

Gradient Descent :-

$$f(\theta) = (2-3)^2 + 1$$

Generally Applicable
→ Does not make any assumptions
→ Any function (diff.iable)

$$t = 0$$

$$\theta^{(0)} \leftarrow \text{init}();$$

$$\text{do } \{$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \underset{\text{learning rate}}{\eta} \cdot \underset{\text{gradient}}{\nabla_{\theta} f(\theta)}$$

$$t \leftarrow t + 1;$$

} while ! converged;

"overfitting" ←
Early stopping

Converged:-

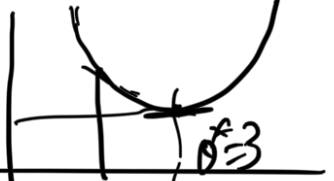
$$1) |\theta^{(t+1)} - \theta^{(t)}| < \epsilon$$

$$\text{or: } |\theta_j^{(t+1)} - \theta_j^{(t)}| < \epsilon_j$$

$$2) \left| \frac{df(\theta)}{d\theta_j} \right| < \epsilon_j$$

$$\left| \nabla_{\theta} f(\theta) \right| < \epsilon$$

$$3) |f(\theta^{(t+1)}) - f(\theta^{(t)})| < \epsilon$$



$$\theta^{(0)} = 2$$

$$\theta^{(1)} = 2.2$$

$$\theta^{(2)} = 2.36$$

When do we stop?
How can we verify this?

threshold

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \nabla_{\theta} f(\theta) \cdot \eta$$

$$\epsilon = \epsilon' \cdot \eta$$

Gradient Descent

$$f(\theta)$$

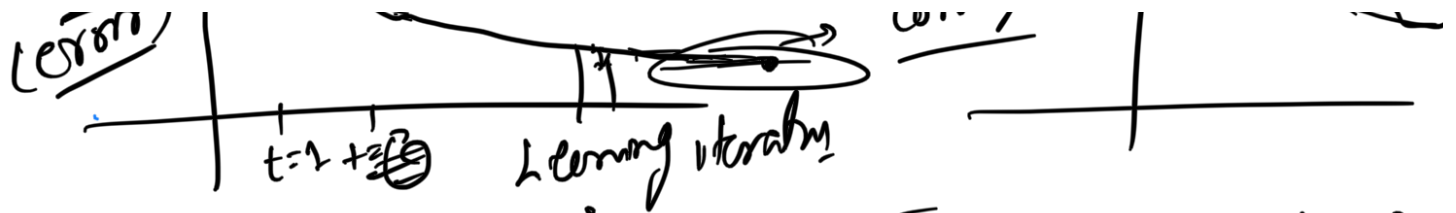
$$\theta^{(t)} = 2$$

$$\theta^{(t+1)} = 2.36$$

SGD

Stochastic Gradient Descent

"converged"

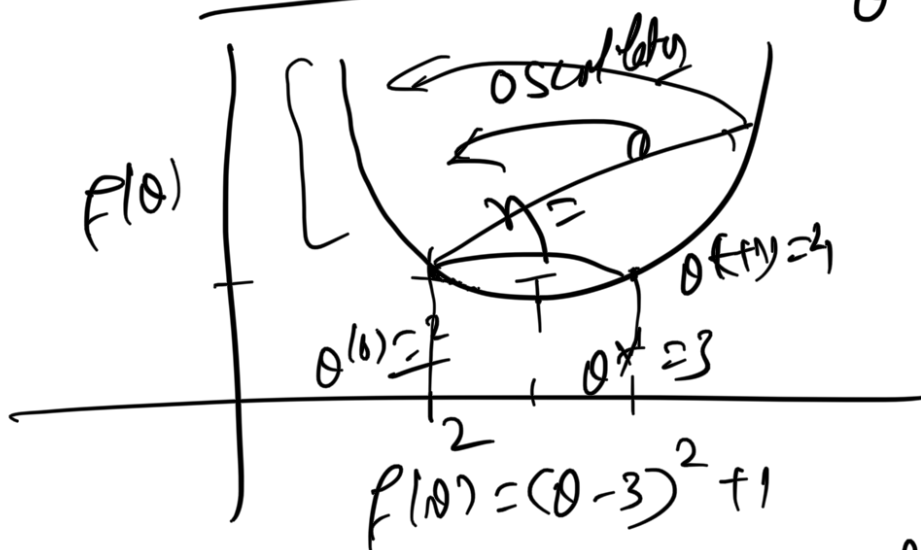


\Rightarrow often!:- \vec{F} is run for a ["certain" no. of (loop enough) iterations] "epoch"

$$J(\theta) = \left[\frac{1}{2m} \sum_{i=1}^m [y^{(i)} - h_{\theta}(x^{(i)})]^2 \right]$$

\Rightarrow Can be expanded if m is large.

choosing "right" η :-



$$\theta(t+1) \leftarrow \theta(t) - \eta \cdot \nabla_{\theta} f(\theta)$$

$$\nabla_{\theta} f(\theta) = 2(\theta - 3)$$

$$\eta = 0.1$$

$$\theta(0) = 2$$

$$\theta(1) = 2.2$$

$$\theta(2) = 2.36$$

\Rightarrow what happens if

η :-

(a) too large

$t=0$

$$\eta = 1$$

$$\theta(0) = 2$$

$$\theta(t+1) \leftarrow 2 - 1 \cdot 2(2-3)$$

$$= 2 - (-2)$$

Too fast

(

$$= 4$$

$$t=1$$

$$\theta^{(t+1)} \leftarrow 4 - 1.2(4-3)$$

$$[2(0-3)]_{\theta=4} \text{ gradient} = 4-2=2$$

$$\downarrow \quad \eta \geq 1$$

\Rightarrow (b) η :- too small.

Too slow

$$t=20 \quad \eta = .01$$

$$\begin{aligned} \theta^{(t+1)} &\leftarrow 2 - .01 * 2(2-3) \\ &= 2 - (-).02 = 2.02 \end{aligned}$$

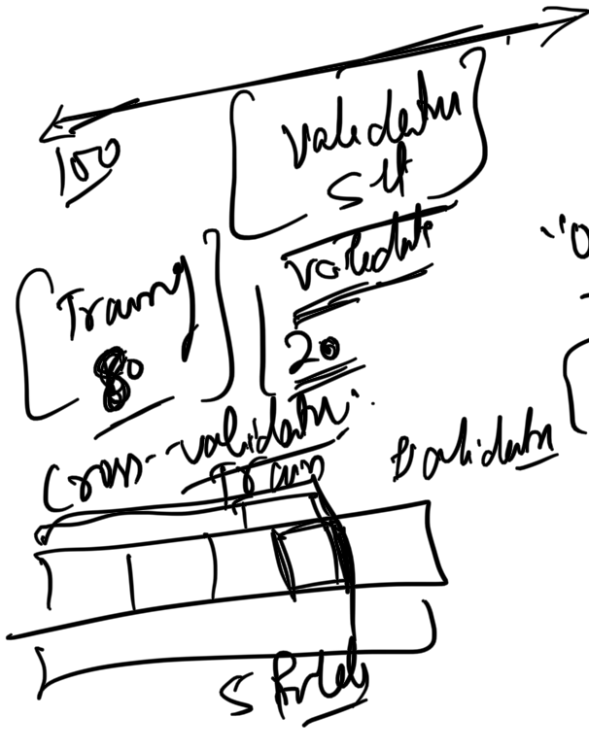
$$t=1$$

$$2 \quad \theta^{(t+1)} \leftarrow 2.02 - .01 * 2(2.02 - 3)$$

$$= 2.02 - (-1)(.96) * .01$$

$$= 2.02 + .0196$$

$$= 2.0396$$



"overfitting"

$$(x^i, y^i)_{i=1}^n$$

$$(x^u, y^u) \text{ wish}$$

computing partial derivatives

$$n^2 + 10x^2 + 20x + 20x^2$$

$$f(\theta) = \theta_1 + \theta_2 + \theta_3 + \theta_4$$

$$\nabla_{\theta} f(\theta) = \begin{bmatrix} 2\theta_2 + 3\theta_1 + 2 \\ 2\theta_1 + 3\theta_2 \end{bmatrix}$$

$$= \begin{bmatrix} 7 \\ 5 \end{bmatrix}$$

$$\leftarrow \begin{array}{l} \frac{\partial f(\theta)}{\partial \theta_2} \\ \frac{\partial f(\theta)}{\partial \theta_1} \end{array} \quad \left| \quad \theta_2 = 1, \theta_1 = 1 \right.$$

Linear Regression:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [y^{(i)} - h_{\theta}(x^{(i)})]^2$$

$$\nabla_{\theta} J(\theta) = \downarrow = \frac{1}{2m} \sum_{i=1}^m \underbrace{[y^{(i)} - \theta^T x^{(i)}]^2}_{\text{quadratic function}}$$

\Rightarrow unique
local minima of $J(\theta)$
(convex fn.)