

C02774

Machine Learning

Aug 25, 2021

Last class:-

Linear Regression

Convex Function

$\rightarrow f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex iff

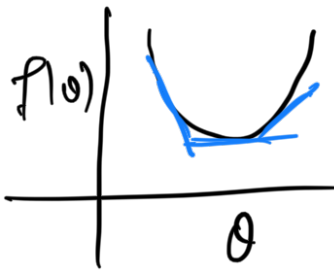
$$f(\alpha \theta^{(1)} + (1-\alpha) \theta^{(2)})$$

$$\leq \alpha f(\theta^{(1)}) + (1-\alpha) f(\theta^{(2)})$$

Minima  
Convex Fun

All local  
minima = Global

Minima



Hessian:-

$(H)$



$$Z^T H Z \geq 0$$

(+ve semi-definite)

$$H_{ijk} = \frac{\partial^2 f(\theta)}{\partial \theta_j \partial \theta_k}$$

strictly

2 assignments:-

① Ass 2 out  
by the summer  
2.5 About 2.5  
Python weeks

② Make-up  
class Sat Aug  
28, 2021

8am-9am

Convex Function



f is concave iff

$$f(\alpha \theta^{(1)} + (1-\alpha) \theta^{(2)})$$

$$\geq \alpha f(\theta^{(1)}) + (1-\alpha) f(\theta^{(2)})$$

Hessian:-

$$Z^T H Z \leq 0$$

-ve semi-definite

Local maxima  
Global maxima

Further Topics:-

① SGD (Stochastic Gradient Descent)

② Newton's Method (2<sup>nd</sup> order)

③ Analytical Solution for Linear reg  

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \theta_0 \ln^{(i)})^2$$

④ Probabilistic interpretation of Linear Regression / Least Squared Regression.  
 $\hookrightarrow \underline{MLE} := \left[ \begin{array}{l} \text{Maximum Likelihood Estimator} \end{array} \right] := \underline{MLE}$

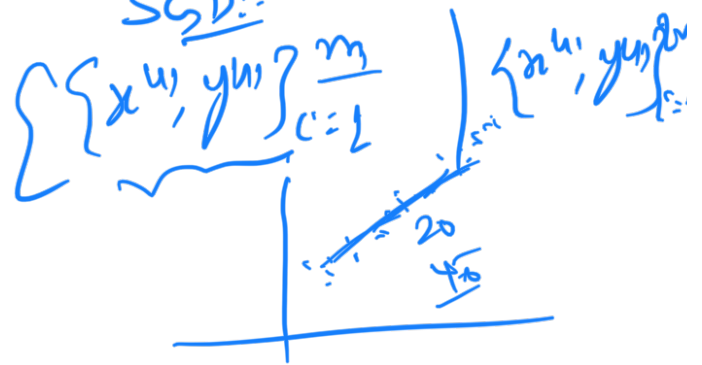
Logistic Regression - Classification

Stochastic Gradient Descent :- (Intuition)

GD :-  
 $t \leftarrow 0;$   
 $\theta(t) \leftarrow \text{init}();$   
 $\text{do } \{$   
 $\leftarrow \theta(t+1) \leftarrow \theta(t) - \nabla J(\theta) \cdot \eta$   
 $t \leftarrow t+1;$   
 $\} \text{ while ! converged}$   
 $O(\underline{m \cdot n})$

$n$  :- # of examples  
 $m$  :- # of features

SGD :-



$\theta^{(i)} = (x^{(i)}, y^{(i)}) \sim \text{Dist. i.i.d. (independently distributed)}$

Binomial Variance

$$\frac{\nabla J(\theta)}{m} \approx \sum_{i=1}^m (x^{(i)}, y^{(i)}) \theta^{(i)}$$

$$J(\theta) \equiv J(\theta) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \theta_0 \ln^{(i)})^2$$

$$\frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \theta_0 \ln^{(i)}) \nabla \theta_0 \ln^{(i)}$$

→ SGD

$t \leftarrow 0;$   
 $\theta(t) \leftarrow \text{init}(\theta);$   
 do  
    $B(t) \leftarrow \text{RandomSample}(\{x^u, y^u\}_{u=1}^m, \delta);$   
    $\theta(t+1) \leftarrow \theta(t) - \eta \nabla_{\theta} f_{B(t)}(\theta) |_{\theta(t)}$   
    $t \leftarrow t+1;$   
 while (not converged)

SGD will converge (in expectation) to local min.

$\{ (x^u, y^u) \}_{u=1}^m$

Cost function:

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

Gradient:

$$\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^m (y^u - h_{\theta}(x^u)) \cdot \nabla_{\theta} h_{\theta}(x^u)$$

$$= \frac{1}{2\delta} \sum_{(x^u, y^u) \in B(t)} (y^u - h_{\theta}(x^u))^2$$

$$\nabla_{\theta} J_{B(t)}(\theta) = \frac{1}{\delta} \sum_{(x^u, y^u) \in B(t)} (y^u - h_{\theta}(x^u)) \nabla_{\theta} h_{\theta}(x^u)$$

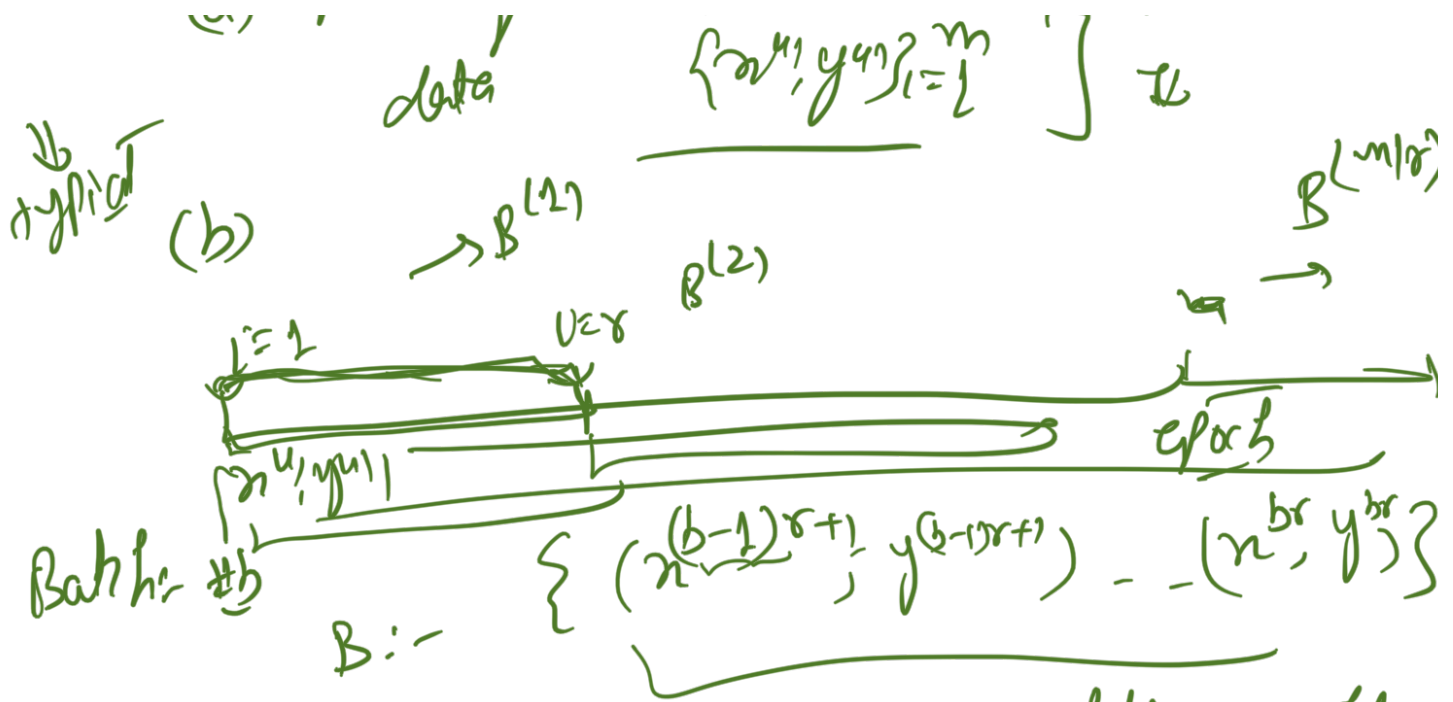
⇒

$$\nabla_{\theta} J(\theta) \xrightarrow{\text{Expectation}} E[\nabla_{\theta} J_{B(t)}(\theta)]$$

H.W.  $B(t) \approx \{x^u, y^u\}_{u=1}^m$

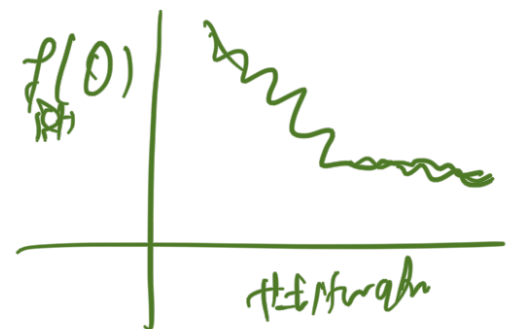
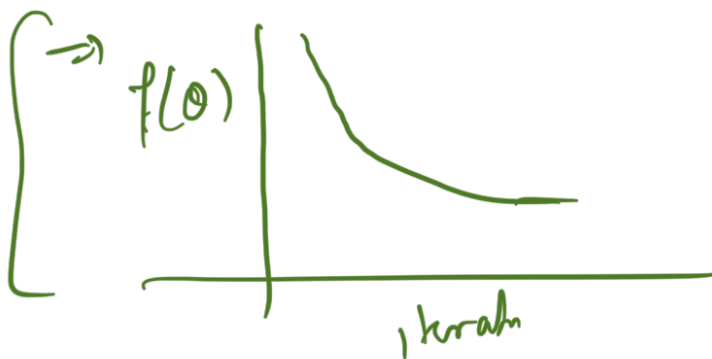
↳ Batch: ...

can Randomly Shuffle entire training



$\Rightarrow \equiv$  when examples are randomly sampled.

$\rightarrow$  GPD:- parallelization:-

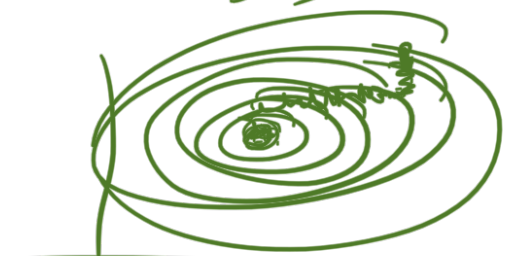


① Why is GGD helpful?

Each parameter update is incremental  $\rightarrow$  time taken for each parameter update?  $O(\frac{1}{\epsilon})$

$\frac{1}{\epsilon} \ll \frac{m}{\epsilon}$   $\rightarrow$  Approximate SGD

$\frac{1}{\epsilon} \ll \frac{m}{\epsilon}$   $\rightarrow$  Approximate SGD



~~\_\_\_\_\_~~

→ SGD: same generated GD

SGD:-

Each update is much faster

:- stochasticity

\_\_\_\_\_