

COL 774

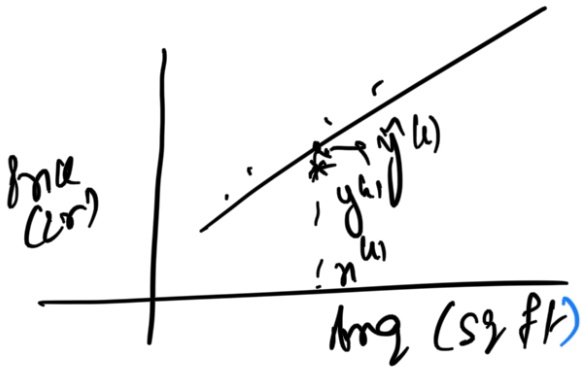
Machine Learning

Sat Aug 21, 2021

Note: -
Sunday Aug 22, 2021
8 am - 9 am

Last class: -

Linear Regression: -



$$\hat{y} = h_{\theta}(x) = \theta_0 + x_1 \theta_1$$

① Hypothesis space?

↳ linear

$$h_{\theta}(x) = \theta^T x$$

$$\text{Hypothesis } \sum_{j=0}^n \theta_j x_j$$

② A good fit?

$$x_0 = 1$$

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

argument $J(\theta)$

③ optimization

how to we compute
argument $J(\theta)$

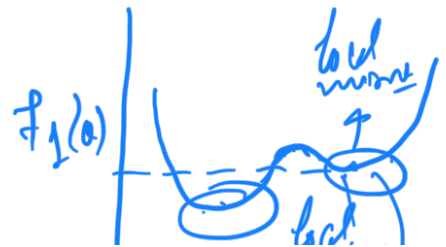
↳ Defn: -

multivariate
con



local minimum
global minimum

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





with
6 or 8 cells
function

⇒ / Conver function :- Local names
≡ Global names

1. "Gravimetric Descent" :-

"Ground Acid"

$$f(0) = (0-3)^2 + 1$$



$f'(0) = 2(0-3) \mid 0=2$
 $= -2$ convex f''

At local minimum:

$f'(0) = 0$ at any value for f

$f'(0) = 2(0-3) + 0$
↳ Finding zeros of gradient

$$\Rightarrow p(1010) \rightarrow [0 \leq 8]$$

$$\left[\frac{d^2 f(\theta)}{d\theta^2} \right]_{\theta=0} = 2 > 0 \Rightarrow \theta=0$$

$\Rightarrow 0=3$ is
of local max.

→ Intertial:- To find the points of local maxima of ϕ .

$f \wedge g \Rightarrow$

$$\frac{2f(0)}{20} = 0$$

may not
always be
possible

Gradient Descent:-

analytically

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \left. \frac{\partial f(\theta)}{\partial \theta} \right|_{\theta^{(t)}}$$

eta learning rate

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

$$\eta = 0.1$$

$$\begin{aligned} \theta^{(1)} &= 2 - 0.1 \left(\frac{\partial (2\theta - 3)}{\partial \theta} \right) \bigg|_{\theta=2} \\ &= 2 - 0.1 \times 2(-1) \\ &= 2 + 0.2 = \underline{2.2} \end{aligned}$$

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

$$\begin{aligned} \theta^{(1)} &= 2.2 - 0.1 \left(\frac{\partial (2\theta - 3)}{\partial \theta} \right) \bigg|_{\theta=2.2} \\ &= 2.2 - 0.1 \times 2(2.2 - 3) \\ &= 2.2 + 0.16 = \underline{2.36} \end{aligned}$$

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

$t \leftarrow t+1$ $t=2$

$$\theta^{(2)} = \theta^{(1)} - \eta \left. \frac{\partial f(\theta)}{\partial \theta} \right|_{\theta^{(1)}}$$

$$= 2.36 - 0.1 \times [2(\theta - 3) \big|_{\theta=2.36}]$$

$$= 2.36 - 0.1 \times 2 \times [-0.36] = 2.36 + 0.072 = 2.432$$

$$= 2.36 + 0.128$$

$$= 2.488 \rightarrow$$

\Rightarrow Using this process, you can get arbitrarily close to local minima if you choose sufficiently small η

Gradient Descent:-

1D Univariate Cost

$$t \leftarrow 0;$$

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

do {

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \cdot \frac{df(\theta)}{d\theta}$$

$$t \leftarrow t+1;$$

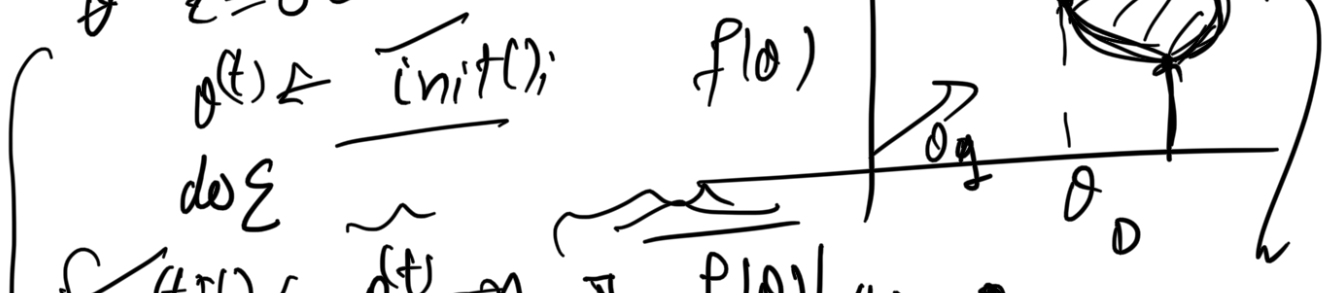
} while (! converged);

Multivariate Cost:-

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

$f(\theta) =$ - nonlinear
Gradient Descent Algorithm:-

$$\text{Sub-gradient } h_0(x) = \theta_0 x + \theta_1$$



$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \cdot \nabla_{\theta} f(\theta) |_{\theta^{(t)}}$
 $\nabla_{\theta} f(\theta) = -$ vector
 $\nabla_{\theta} f(\theta)$ while (unmerged)



$$\begin{bmatrix} \frac{\partial f(\theta)}{\partial \theta_1} \\ \frac{\partial f(\theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial f(\theta)}{\partial \theta_n} \end{bmatrix}_{n+1}$$

$\theta \in \mathbb{R}^{n+1}$

center O_0 speed
"Newton's Method"

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \cdot \nabla_{\theta} f(\theta) |_{\theta^{(t)}} \Rightarrow \theta_j^{(t+1)} \leftarrow \theta_j^{(t)} - \eta \frac{\partial f(\theta)}{\partial \theta_j} |_{\theta^{(t)}}$$

$$\theta_j^{(t+1)} \leftarrow \theta_j^{(t)} - \eta \cdot \frac{\partial f(\theta)}{\partial \theta_j} |_{\theta^{(t)}}$$

\Rightarrow :- ① Impact of η
 ② How do we decide convergence?

① Adaptive η

η :- depends of learning rate ~~no.~~

$$\eta \propto \frac{1}{t} \text{ or } \frac{1}{\sqrt{t}}$$

② Different η for different converge

↳ "Large Diff." \Rightarrow

$$\eta = \frac{y_0}{t} \sim \frac{v_0}{\sqrt{t}}$$