

SABIN

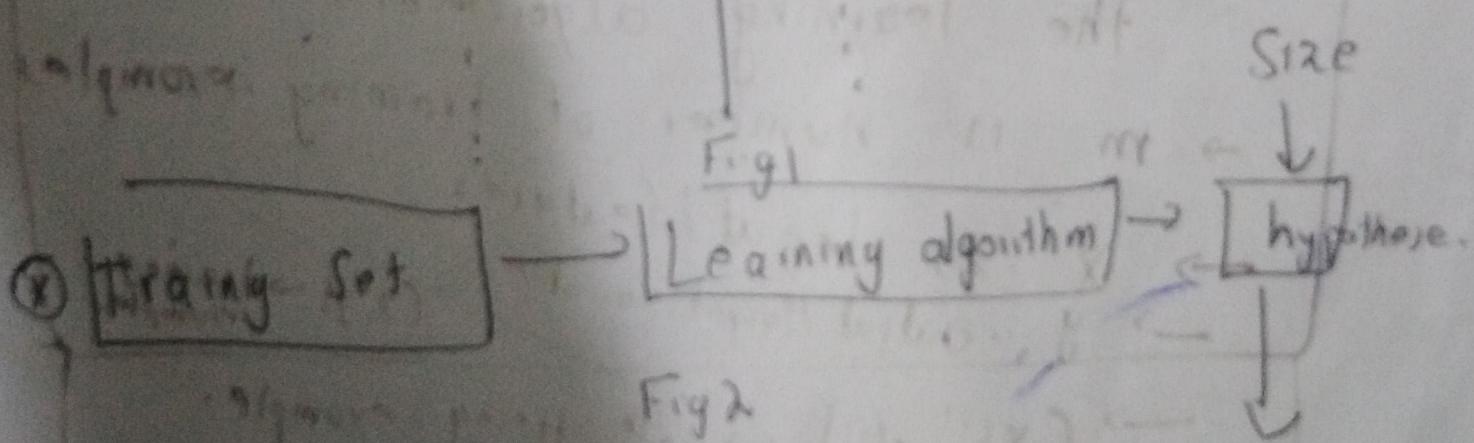
Q G5-229

Latost
Notes
Lecture-2

① What is supervised learning problem?

② Linear Regression:

③ Size (foot)	Price (C\$1000)
2104	900
1416	232
1539	315
852	178



④ How to Represent the Hypothesis

$$\rightarrow h(x) = \theta_0 + \theta_1 x \quad (\text{Linear Regression Hypothesis})$$

(2)

- ④ → There could be additional features as well

⑤ → $h(x) = \sum_{j=0}^n \theta_j x_j$ (Concise form)
where $x_0 = 1$

↳ line

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}, x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

definition

- θ is called the parameters of the learning algorithm
- m is number of training examples
- x is input / feature
- y - output / target variable, $(x, y) \rightarrow$ example
- $(x^{(i)}, y^{(i)}) \rightarrow$ i-th training example
- $x_j^{(i)}$ where $i \rightarrow$ training example
 $j \rightarrow$ feature
- $n \rightarrow$ no. of features (total)

Choose θ s.t $h(\mathbf{x}) \approx y$ for training examples

① $h_{\theta}(\mathbf{x}) \rightarrow$ Hypothesis depends on both \mathbf{x} and parameter θ .

→ In Linear Regression minimize

$$\min_{\theta} \left(h_{\theta}(\mathbf{x}) - y \right)^2 \text{ Eq. 1}$$

$$J(\theta) \equiv \min_{\theta} \frac{1}{2} \sum_{i=1}^m \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2. \text{ Eq. 1(a)}$$

$$\rightarrow \min_{\theta} J(\theta).$$

→ How can we minimize $J(\theta)$?

② Gradient descent - (Baby step to go downhill by taking small steps)

- Repeat ← ① Start with some θ (say $\theta = \vec{0}$).
Until converge → ② Keep changing θ to reduce $J(\theta)$

③ ~~which~~ step Learning rate.

$$\text{④ } \theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (j=1, 2, \dots, n).$$

(4)

→ Derivative of function determine the step size.

↳ Negative sign mean you need to take the step down the hill.

• Assuming one Example:-

$$\rightarrow \frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2.$$

$$= 2 \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y).$$

$$= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n)$$

$$= (h_{\theta}(x) - y) \cdot x_j$$

So Now (Gradient descent equation step ③) on Linear Regression

$$\rightarrow \theta_j := \theta_j - \alpha (h_{\theta}(x) - y) \cdot x_j \quad (\text{Eq. 2})$$

④ This is for one example
therefore for m examples:-

$$\rightarrow \theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \quad (\text{Eq. 3})$$

→ Now ~~for~~^{GD} will take equation 3 and gradient descent would update θ_j for $i = 1, 2, \dots, n$ until convergence

④ The cost $J(\theta)$ is quadratic function
therefore it is a big bowl and
would have a global optima.

⑤ If we draw the contours for $J(\theta)$ (SSE)
then it will have perfect ellipses:

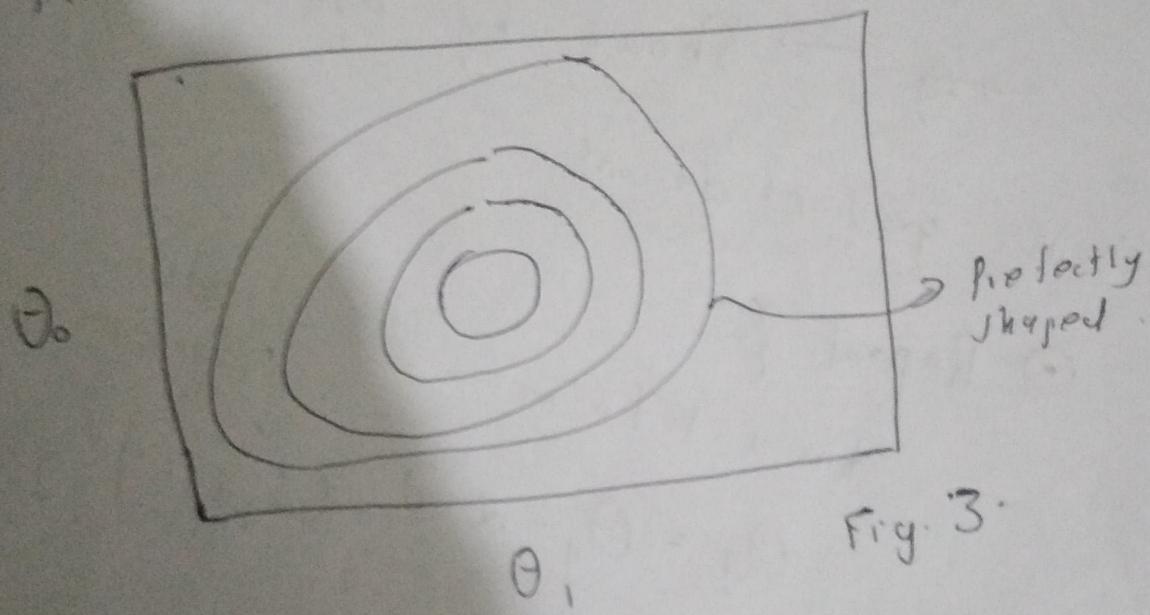


Fig. 3.

⑥ So if you ran gradient descent then
the stepdownhill under this cost would

always be 90 degree (orthogonal)

↳ Algorithm will converge
to global optimia

⑦ $\alpha \rightarrow$ Defines how large will be
step.

(6)

④ Batch gradient descent:

→ Every step of gradient descent will be slow if we take one step of SGD on entire data.

→ frequent {

→ shown by equation 3.

⑤ Stochastic gradient descent:

① Repeat {

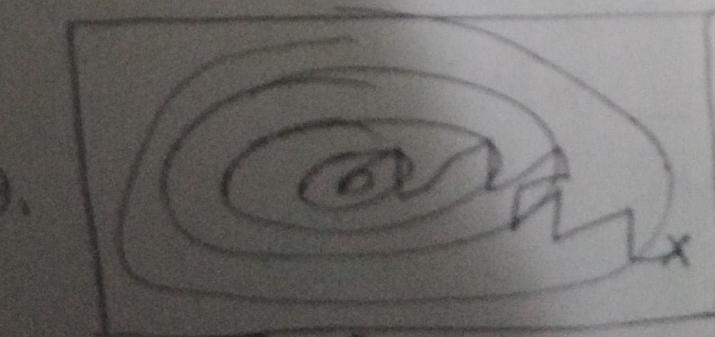
For $i = 1 \text{ to } M$

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot \mathbf{x}_j^{(i)}$$

Update
for every g

} → Eq. 4

② Taking gradient descent step
just by one example at each
time.



→ SGD
will take
a more
noisy step
than

Batch
gradient

Fig. 4

① Even at global minimum SGD will still oscillate

② SGD allows you to learn the hypothesis more fully

* Mini-Batch Gradient Descent

→ Batch size is of some fix size of k examples

* If you have small dataset → use Batch gradient descent.

* Non-iterative Algorithm to solve for parameters

for linear regression:

③ Normal equation (Chaining to show how to derive it)

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix}$$

↑
∈ ℝⁿ⁺¹

* Let's say we have A H.R^{2x2}

$$f(A) = A_{11} + A_{12}^2 \quad (\text{Eq. 5}) \quad f: \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}$$

then

$$f\left(\begin{bmatrix} 5 & 2 \\ 3 & 4 \end{bmatrix}\right) = 5 + 2^2 = 9$$

④ $\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} \\ \frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} \end{bmatrix}$

~~defining
function f
with respect
to matrix A
and computing
partial derivatives
of f with respect
to each element of A~~

⑤ We compute ~~A~~ derivative of f w.r.t A .

⑥ $\nabla_A f(A) = \begin{bmatrix} 1 & 2A_{12} \\ 0 & 0 \end{bmatrix} \rightsquigarrow$ taking the partial function f (See Eq-5)

↓
This show the definition of derivative of a matrix with respect to a specific function.

(x)

⑥ How are we going to derive
normal equation?

$$\nabla_{\theta} J(\theta) = \vec{0}$$

⑦ More derivation to help us in above normal equation derivation.

$\text{Tr}(A) = \sum \text{of diagonal elements}$

Trace:

$$① \text{tr } A = \text{tr } A^T$$

$$③ \text{tr } AB = \text{tr } BA$$

$$② f(A) = \text{tr } AB$$

$$④ \text{tr } ABC = \text{tr } CAB$$

$$\nabla_A f(A) = B^T$$

$$⑤ \nabla_A \text{tr } AA^T C = CA + C^T A$$

⑧ Now returning back to derivation of N.E:

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

$$X = \begin{bmatrix} & (x^{(1)})^T \\ & \cdots \\ \downarrow & x^{(2)^T} \\ \text{design} & \vdots \\ \text{mat.} & x^{(m)^T} \end{bmatrix}$$

$$\rightarrow X\theta = \begin{bmatrix} - & X^{(1)\top} \\ - & X^{(2)\top} \\ \vdots & \\ - & X^{(m)\top} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$= \begin{bmatrix} X^{(1)\top} \theta \\ X^{(2)\top} \theta \\ \vdots \\ X^{(m)\top} \theta \end{bmatrix} = \begin{bmatrix} h_\theta(X^{(1)}) \\ \vdots \\ h_\theta(X^{(m)}) \end{bmatrix}$$

$$\rightarrow y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

\rightarrow Min. Out:

$$J(\theta) = \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

$$X\theta - y = \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}$$

remember $J(\theta) = \sum z^2$ therefore $(X\theta - y)^T (X\theta - y)$
 is sum of square errors.

④ Continuing derivation of N.E

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (\theta^T X - y)^T (\theta^T X - y)$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T X^T y^T) (\theta^T X - y)$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T X^T \theta - \theta^T X^T y - y^T X \theta + y^T y)$$

Similar to $(ax-b)(ax+b)$

$$= a^2 x^2 - axb - bax + b^2$$

z Taking derivatives w.r.t θ

$$= \frac{1}{2} [X^T X \theta + X^T X \theta - X^T y - X^T y]$$

Simplify

$$\# X^T X \theta - X^T y = 0$$

$$\# X^T X \theta = X^T y \text{ (Normal equation).}$$

$$\theta = (X^T X)^{-1} \cdot X^T y \text{ (Optimal value for}$$

θ through
Normal Equation).

→ if X is invertible

then use pseudo inverse.