Outline

- Linear regression
- Batch / Stochastic gradient descent
- Normal equation

Supervised Learning
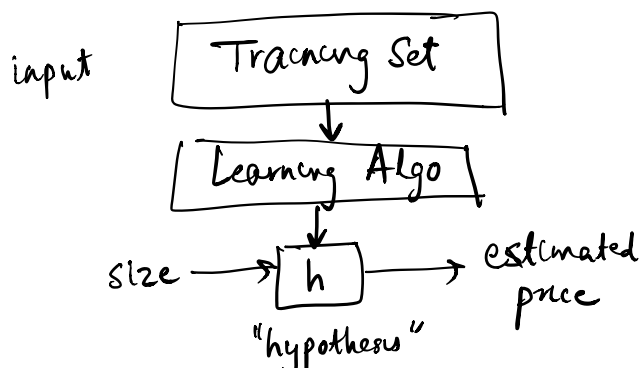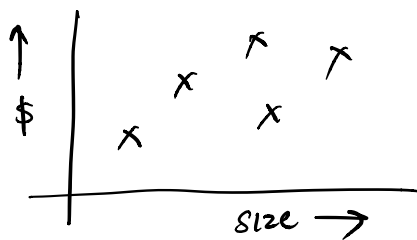
$$X \longrightarrow y$$

picture        steering direction

| Regression | (o/p continuous)

v/s classification

Housing dataset

| Size | Price ($ 1,000s) |
|------|------------------|
| 2104 | 400 |
| 1416 | 232 |
| 1534 | 315 |



input

Training Set

↓

Learning Algo

size → [ h ] → estimated price

"hypothesis"

How to represent $h$?

$h(x) = \theta_0 + \theta_1 x$     (technically affine fn)

More features

|  | Size | # bedrooms | Price |
|---|---|---|---|
| $x^{(1)}$  1 | 2104 | 4 | 400 |
| $x^{(2)}$  1 | 1416 | 3 | 232 |

$x_1^{(1)} = 2104$
$x_1^{(2)} = 1416$

$X_1 = \text{size}, \quad X_2 = \text{\# bedrooms}$

$h(x) = \theta_0 + \theta_1 X_1 + \theta_2 X_2$

$h(x) = \sum_{j=0}^{2} \theta_j X_j$

Define $X_0 = 1$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \qquad X = \begin{bmatrix} X_0 \\ X_1 \\ X_2 \end{bmatrix} \begin{matrix} \text{always } 1 \\ \text{size} \\ \text{\# bedrooms} \end{matrix}$$

parameters.

$n = $ # training examples

$X = $ "inputs" / features.

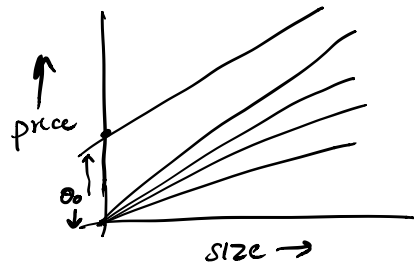$y = $ "output" / target variable.

$(X, y) = $ training example

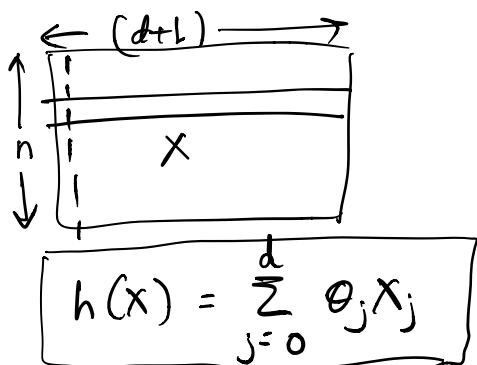$(X^{(i)}, y^{(i)})$ : $i^{th}$ training example

$X_1^{(i)}$ : $i$ runs from $1$ to $n$

$d = $ # features

$(d = 2)$

$X^{(i)}, \theta$ $(d+1)$ dimensional

$$\overset{\longleftarrow (d+1) \longrightarrow}{}$$



$$\boxed{h(x) = \sum_{j=0}^{d} \theta_j X_j}$$

Choose $\theta$ s.t. $h(x) \approx y$

$$h_\theta(x) = h(x)$$

Cost $f^n$  $\quad J(\theta) = \frac{1}{2} \sum_{i=1}^{n} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

$$\underset{\theta}{\min} \; J(\theta)$$

Gradient Descent

Start with $\theta$ (say $\theta = \vec{0}$)

Keep changing $\theta$ to reduce $J(\theta)$



$J(\theta)$     $\theta$    local minimum

local min = global min

Gradient Descent

Start with $\theta$

Repeat until convergence

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \qquad (j = 0, 1, \dots d)$$
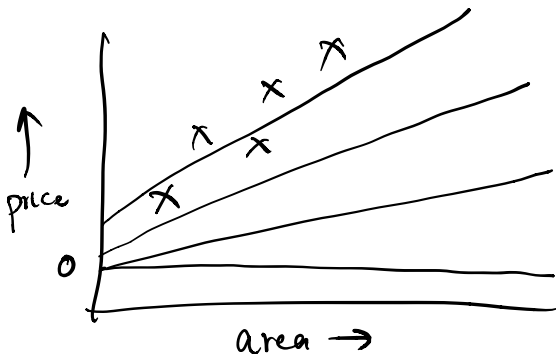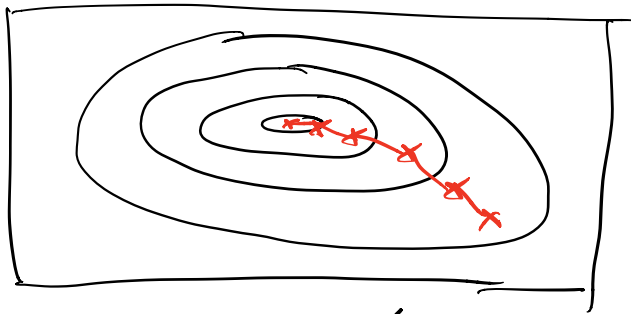
$\uparrow$ learning rate

$a := a + 1$ ✓
$a = a + 1$ ✗

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

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

(chain rule)

$$= (h_\theta(x) - y) \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \theta_1 x_1 \cdots + \theta_d x_d - y)$$
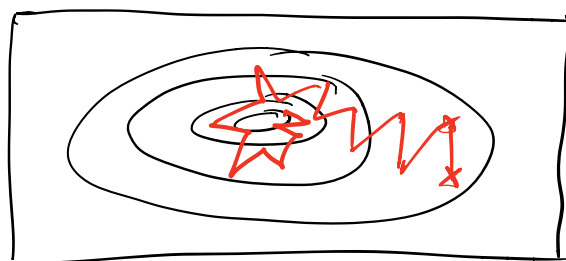
$$= (h_\theta(x) - y) x_j$$

$$\theta_j := \theta_j - \alpha (h_\theta(x) - y) x_j$$

$$\theta_j := \theta_j - \alpha \underbrace{\sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\frac{\partial}{\partial \theta_j} J(\theta)}$$





"Batch" gradient Descent
Stochastic gradient Descent

Repeat {
    For $i = 1$ to $n$ {
        For $j = 0$ to $d$ {
            $\theta_j := \theta_j - \alpha \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)}$
    }



Mini-batch

Normal Equation

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

$A \in \mathbb{R}^{2 \times 2}$      $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$

$f(A)$

$f : \mathbb{R}^{2 \times 2} \longrightarrow \mathbb{R}$

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

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

$$J(\theta) = \frac{1}{2} \sum_{i=1}^{n} \left( h(x^{(i)}) - y^{(i)} \right)^2$$

$$X = \begin{bmatrix} \leftarrow (x^{(1)})^T \longrightarrow \\ \\ \longrightarrow (x^{(n)})^T \end{bmatrix} \quad \text{design matrix}$$

$$X\theta = \begin{bmatrix} \quad \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \end{bmatrix} = \begin{bmatrix} h_\theta(x^{(1)}) \\ \vdots \\ h_\theta(x^{(n)}) \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

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

$$\nabla_\theta J(\theta) = X^T X \theta - X^T y = \vec{0}$$

$$X^T X \theta = X^T y \quad \text{"Normal equation"}$$

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