- Non-linear model
- Neural network

## Linear regression review

$$h_\theta(x) = \theta^T x + b$$

$$J(\theta) = \sum_{i=1}^{n} (y^{(i)} - h_\theta(x^{(i)}))^2 = \sum_{i=1}^{n} (y^{(i)} - \theta^T x^{(i)} - b)^2$$

Run GD or SGD

## Non-linear model

- Kernel $\quad h_\theta(x) = \theta^T \phi(x) \quad$ <span style="color:red">linear in $\theta$, but not $x$</span>

- non-linear in both $\theta$ and $x$

$$\text{ex: } h_\theta(x) = \sqrt{\theta_1^3 x_2 + \sqrt{\theta_3} x_4}$$

Assume we have a dataset $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

$$x^{(i)} \in \mathbb{R}^d, \quad y^{(i)} \in \mathbb{R}$$

$$h_\theta \; \mathbb{R}^d \to \mathbb{R}$$

- Cost func for example $i$

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

- Cost func for dataset

$$J(\theta) = \frac{1}{n} \sum_{i=1}^{n} J^{(i)}(\theta)$$

Note: the minimizer will remain the same w/o $(\frac{1}{n})$

### Optimize

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

GD: $\quad \theta := \theta - \alpha \nabla J(\theta)$

SGD:

for $i = 1$ to $n_{iter}$:
    sample $j$ from $\{1, \dots, n\}$

$$\theta := \theta - \alpha \nabla J'(\theta)$$

Computing $\underline{B}$ gradients $\nabla J^{(j_1)}(\theta), ... \nabla J^{(j_B)}(\theta)$ simultaneously is faster
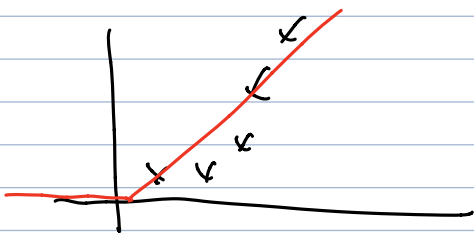(batch size)

for $i = 1$ to $n_{iter}$:

Sample $B$ examples $j_1, ... j_B$ without replacement

$$\theta := \theta - \alpha \frac{1}{B} \sum_{k=1}^{B} \nabla J^{(j_k)}(\theta)$$

## Key Points

① How to define $h_\theta(x)$?    Neural Network

② How to compute $\nabla J^{(i)}(\theta)$?    Back propagation



In deep learning

$$Relu = \max \{t, 0\}$$

$$h_\theta(x) = \underline{relu(wt + b)}$$
activation    $\curvearrowleft$ neural network with one neuron

High dimensional input $x \in \mathbb{R}^d$, $y \in \mathbb{R}$

$$h_\theta(x) = relu(w^T x + b)$$

$$w \in \mathbb{R}^d \qquad x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^d, \qquad b \in \mathbb{R}$$
weigh vector                                          bias

## Stacking neurons

Ex:    $x \in \mathbb{R}^4$.    $x_1, \quad x_2, \quad x_3, \quad x_4$
$\uparrow \qquad \uparrow \qquad \uparrow$
size    # of bed    zip code

intermediate variables:
family size    $a_1$
walkable    $a_2$
school area    $a_3$

family size    $a_1 = relu(\theta_1 x_1 + \theta_2 x_2 + \theta_3)$
$\llcorner$ bias
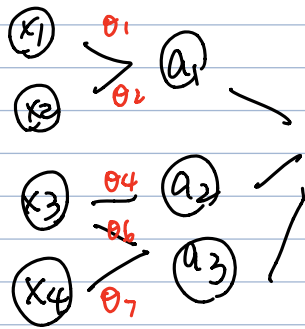
walkable    $a_2 = relu(\theta_4 x_3 + \theta_5)$

$a_3 = relu(\theta_6 x_3 + \theta_7 x_4 + \theta_8)$

$$h_\theta(x) = \text{relu} [\theta_9 a_1 + \theta_{10} a_2 + \theta_{11} a_3 + \theta_{12}]$$

(Generally we don't apply relu at the end)

## Diagram



Intermediate values → **hidden units**

$$a_j = \text{relu}\left(w_j^{[1]T} x + b^{[1]}\right) \quad \forall j = 1 \dots m \qquad - \text{ first layer}$$

$$h_\theta(x) = w^{[2]T} x + b^{[2]} \qquad - \text{ second layer}$$

## Vectorization

$$W^{[1]} = \begin{bmatrix} - & w_1^{[1]T} & - \\ - & w_2^{[1]T} & - \\ & \vdots & \\ - & w_m^{[1]T} & - \end{bmatrix} \in \mathbb{R}^{m \times d}$$

$$W^{[1]} + \begin{bmatrix} b^{[1]} \\ \vdots \\ b^{[m]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]}x + b^{[1]} \\ \vdots \\ w_m^{[1]}x + b^{[m]} \end{bmatrix} \triangleq \begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix}$$

$\mathbb{R}^m$

$$z = w^{[1]}x + b^{[1]} \in \mathbb{R}^m$$

$$a = \begin{bmatrix} \text{relu}(z_1) \\ \vdots \\ \text{relu}(z_m) \end{bmatrix} \triangleq \text{relu}(z)$$

$$W^{[2]} = \left[ w^{[2]T} \right] \in \mathbb{R}^{1 \times m}, \quad b^{[2]} \in \mathbb{R}$$

$$h_\theta(x) = W^{[2]} a + b$$

$$a^{[1]} = relu(w^{[1]} x + b^{[1]})$$

$$a^{[2]} = relu(w^{[2]} x + b^{[2]})$$

$$\vdots$$

$$a^{[r-1]} = relu(w^{[r-1]} x + b^{[r-1]})$$

$$h\theta = w^{[r]} a^{[k-1]} + b^{[k]}$$

$$dim(a_k) = m_k \qquad x \in \mathbb{R}^d, \quad w^{[2]} \in \mathbb{R}^{m_1 \times d}$$

$$w^{[1]} x \in \mathbb{R}^{m_1} \qquad b^{[2]} \in \mathbb{R}^{m_1}$$

$$w^{[2]} \in \mathbb{R}^{m_2 \times m_1}$$

$$w^{[k]} \in \mathbb{R}^{m_k \times m_{k-1}} \qquad b^{[k]} \in \mathbb{R}^{m_x}$$

$$width = max\{m_1, \dots m_{r-1}\}$$