

## Linear regression with one variable

### 2.1) Model representation

⇒ Data set is called training set.

Notation:

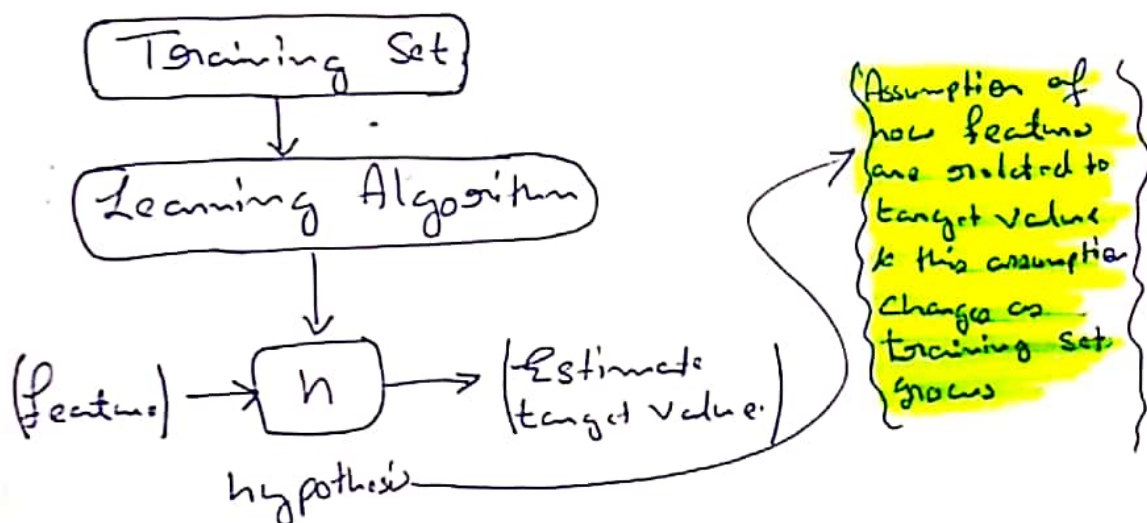
$M \Rightarrow$  Number of training example

$x$ 's  $\Rightarrow$  "input" variables / features

$y$ 's  $\Rightarrow$  "Output" variable / "target" variable

$(x, y) \Rightarrow$  To denote single training example.

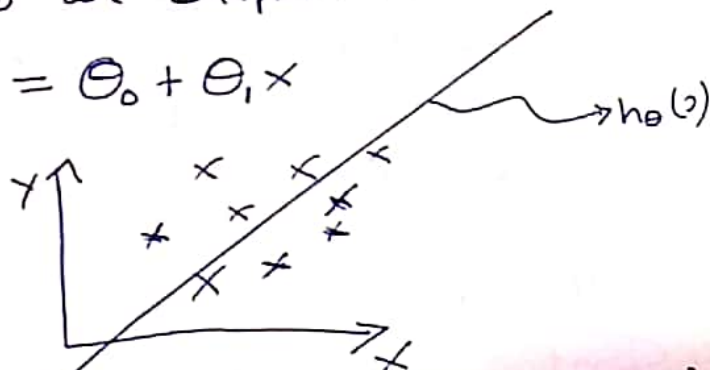
$(x^{(i)}, y^{(i)}) \Rightarrow$  To denote  $i^{\text{th}}$  training example.



⇒  $h$  is a function which maps  $x$ 's to  $y$ 's

⇒ How do we represent  $h$ ?

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



⇒ Linear regression with one variable.  
(i.e. Univariate linear regression)

## 2.2 > Cost function

Hypothesis:  $h_{\theta}(x) = \theta_0 + \theta_1 x$

$\{\theta_0, \theta_1 \Rightarrow \text{Parameters}\}$

Objective: Choose  $\theta_0, \theta_1$  so that  $h_{\theta}(x)$  is close to  $y$  for our training examples  $(x, y)$ .

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Cost function

Squared error cost function

Goal: minimize  $J(\theta_0, \theta_1)$   
 $\theta_0, \theta_1$

Minimize  $J(\theta_0, \theta_1)$  to obtain value of  $\theta_0$  &  $\theta_1$

Contour line  $\Rightarrow$  A contour line of a function of two variables is a curve along which the function has a constant value, so that the curve joins points of equal value.

$\rightarrow$  Different colors are used to represent contour lines at different height.

## 2.5 Gradient descent

→ An algorithm for minimizing the cost function  $J$



$$\min_{\theta_0, \dots, \theta_n} J(\theta_0, \dots, \theta_n)$$

→ First start from a random value

### Gradient descent algorithm

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\forall j=0 \text{ and } j=1)$$

}

⇓ {Correct simulation update}

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

→ {Assignment operator}

$\alpha \Rightarrow$  learning rate  
(i.e. how big the  
Step is)

$\{ \Rightarrow \text{Teach assertion} \}$

Issue: It can be susceptible to local optima.

⇒ The algorithm is also called ("Batch") Gradient Descent.

Each step of gradient descent uses all the training examples

Two extensions

Exact method

Larger number of features

$x_1, x_2, \dots, x_n$   
can be different features

