


Multiple Features:

May 13, 2021

Notation:

$n = \# \text{ of features}$

$x^{(i)}$ = input features of the i^{th} training example

$x_j^{(i)}$ = value of feature j in the i^{th} training example

e.g.

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \rightarrow \begin{array}{l} \text{size} \\ \# \text{ bedrooms} \\ \# \text{ of floors} \\ \text{age of home} \end{array}$$

$$x_2^{(2)} = 3$$

* Now it starts at 0

Now; $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \dots \theta_n x_n$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

∴ $h_{\theta}(x) = \theta^T x$

* Multivariate linear regression

Gradient Descent for Multiple Variables:

Cost Function:

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

Gradient Descent:

repeat {

$$\theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}}_{\frac{\partial}{\partial \theta_0} J(\theta)}$$

Simultaneous update

$$\theta_1 := \theta_1 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}}_{\frac{\partial}{\partial \theta_1} J(\theta)}$$

$$\theta_2 := \theta_2 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}}_{\vdots}$$

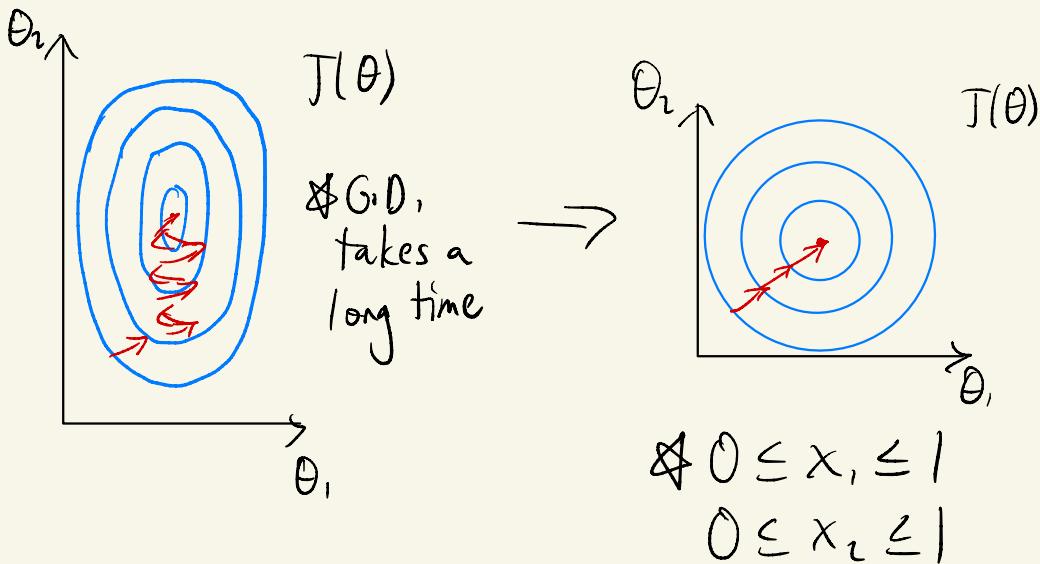
}

repeat { $\theta_j := \theta_j - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\text{Simultaneous update } \theta_j \text{ for } j=0, 1 \dots n}$ }

Simultaneous update θ_j for $j=0, 1 \dots n$

Feature Scaling:

e.g. $x_1 = \text{size } (0 - 2000 \text{ ft}^2) \rightarrow x_1 = \frac{\text{size}}{2000}$
 $x_2 = \# \text{ bedrooms } (1 - 5) \rightarrow x_2 = \frac{\# \text{ bedrooms}}{5}$



- ★ Like $\div 255$ for img pixels
- ★ Not the same as normalization!
 - ↳ changes shape of distribution
- ★ Feature scaling allows for the accuracy to converge faster
 - ↳ usually you want features to be in $-1 \leq x_i \leq 1$ range
 - ★ Some others are fine, too

Mean Normalization:

Mean normalization involves replacing x_i w/ $x_i - \mu$ to make features have \approx zero mean.

* Do not apply to $x_0 = 1$

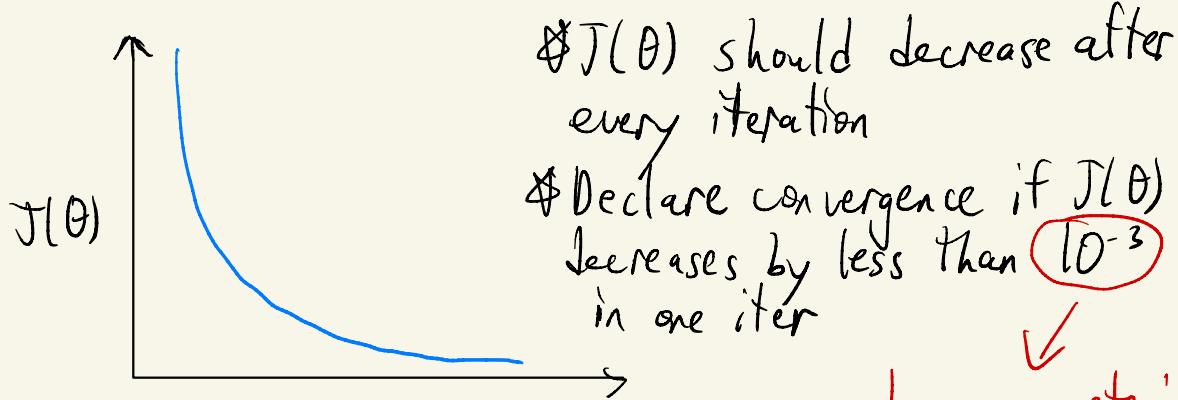
* Can be used w/ feature scaling

e.g. $x_1 = \frac{\text{size} - 1000}{2000} \rightarrow$ avg of x_1 in training set
 \rightarrow range (max-min)

$$x_2 = \frac{\# \text{bedrooms} - 2}{5}$$

$-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$

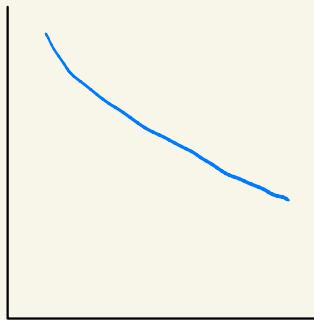
Is Gradient Descent Working?



hyperparameter,
chosen manually

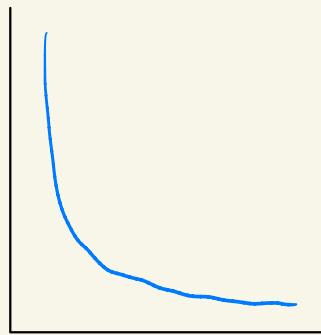
* For sufficiently small α , $J(\theta)$ should decrease on every iter

Learning Rate Cases:

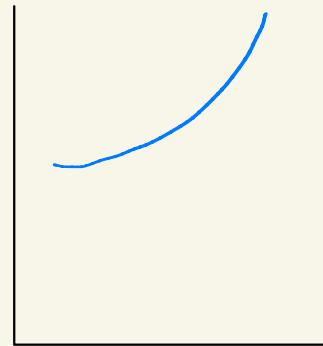


α too low

↳ slow convergence



α just right



α too high

↳ doesn't converge

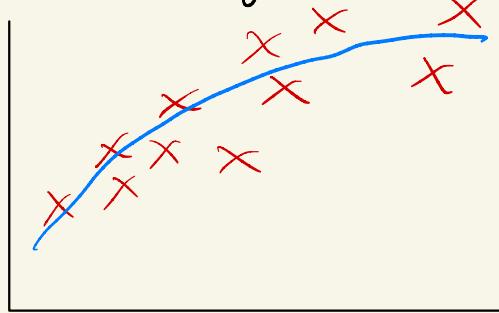
Try $\alpha = 0.001, 0.01, \dots, 1$ (factor of 10)

or $\alpha = 0.001, 0.003, 0.01, \dots, 1$ (factor of ~3)

Polynomial Regression:

e.g.

Price
(y)



Size(x)

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

Must Feature Scale!
or else;

size $(1-1000)$

size 2 $(1-10^6)$

size 3 $(1-10^9)$

$$\text{e.g. } h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 \sqrt{x}$$

let x be size (1-1000)

Feature Scaling!

$$\therefore h_{\theta}(x) = \theta_0 + \theta_1 \frac{x}{1000} + \theta_2 \frac{\sqrt{x}}{\sqrt{1000}}$$