

COURSERA ML: WEEK 2

Linear Regression w/ Multiple Variables:

eg:

size	# bedrooms	# floors	age	price

↓
1 → m

4 variables/features $x_1, x_2, x_3, x_4 \rightarrow 1 \rightarrow n$

↳ H features = # variables = n (above $n=4$)

↳ $x_j^{(i)}$ = val of feature j in i th training example (row)

eg: $x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$ ← site
 ← # bedrooms
 etc.
 ↓

→ $x^{(2)} \in R^4 (R^n)$

↳ $x_3^{(2)} = 2$

↳ NOW: $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$

$(h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n)$

↳ for convenience, $x_0^{(i)} = 1 \rightarrow x =$

$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in R^{n+1}$
 ↳ $n = H$ features

NOTE: $J(\theta) = \frac{1}{2m} (x\theta - y)^T (x\theta - y)$

good form.

$\delta \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$ $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}^{2 \times 1}$

↳ $h_{\theta}(x) = \theta^T x$ ($= \theta \cdot x$)

eg:
 $x = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} \\ x_0^{(2)} & x_1^{(2)} \\ \vdots & \vdots \\ x_0^{(m)} & x_1^{(m)} \end{bmatrix}^{m \times 2}$
 ↳ $n=1$ here

↳ "multivariate linear regression"

$h_{\theta}(x) = X\theta$

Gradient Descent for Multiple Variables

$$\text{Hypothesis: } h_{\theta}(x) = \theta^T x = \underbrace{\theta_0}_{x_0=1} x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

Parameters: $\theta_0, \theta_1, \theta_2, \dots, \theta_n = \underline{\theta}$ ($n+1$ -dimensional vector)

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

$$\Rightarrow \text{Gradient Descent: } \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\underline{\theta}) \quad \text{for } j = 0 \dots n$$

$$\text{for } n > 1 \rightarrow \frac{\partial}{\partial \theta_j} J(\underline{\theta}) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\hookrightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\underline{\theta}) \quad (\text{w/ simultaneous update!})$$

Gradient Descent in Practice I: Feature Scaling & Mean Normalization

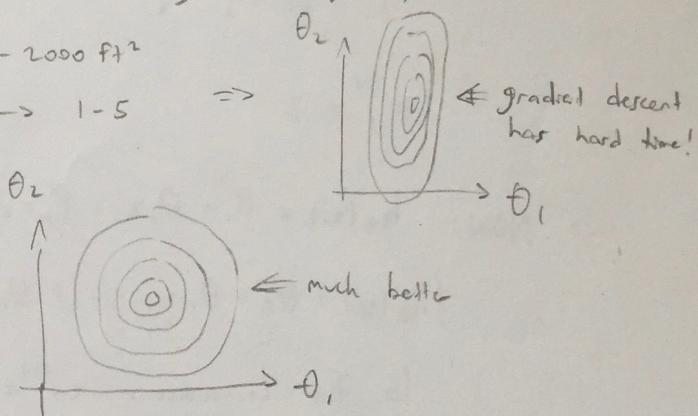
\Rightarrow When features are on similar scale, gradient descent converges more quickly

\hookrightarrow eg 2 features x_1 - size $\rightarrow 0-2000 \text{ ft}^2$
 x_2 - # bedrooms $\rightarrow 1-5 \Rightarrow$

$\left| \begin{array}{l} \text{range} \\ \text{or std. dev.} \end{array} \right| \Rightarrow \text{SCALE FEATURES:}$

$$x_1 - \text{size}/2000 \text{ ft}^2 \rightarrow 0-1$$

$$x_2 - \text{# bedrooms}/5 \rightarrow 0-1$$



\Rightarrow generally, get every feature into a $-1 \leq x_i \leq 1$ range (or close enough)

\hookrightarrow rule of thumb: $-3 \leq x_i \leq 3$ OK, $-\frac{1}{3} \leq x_i \leq \frac{1}{3}$ OK

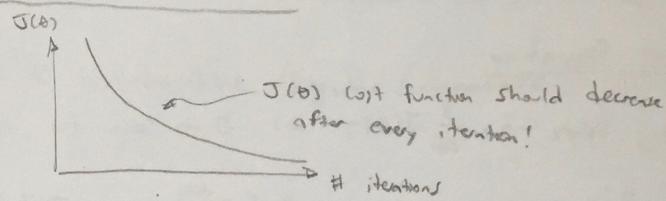
\Rightarrow another strategy: Mean Normalization: $x_i \rightarrow \frac{x_i - \bar{x}_i}{s_i}$ where \bar{x}_i = average

$$\hookrightarrow \text{eg. } x_1 = \frac{\text{size} - 1000}{2000}, \quad x_2 = \frac{\text{bed} - 2}{5} \quad \begin{matrix} \text{mean normalization} \\ \text{feature scaling} \end{matrix}$$

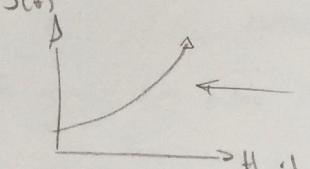
$s_i = \text{max-min or std. dev.}$

Gradient Descent in Practice II: De-bugging & Choosing α

⇒ helpful: plot $J(\theta)$ vs. # iterations



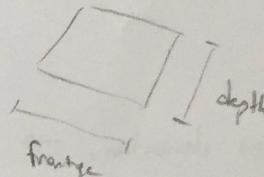
⇒ declare some convergence threshold, e.g. $J(\theta)$ decreases by less than 10^{-3}

⇒ if  likely that α is too big!

↳ for sufficiently small α , $J(\theta)$ should decrease on every iteration (too small = slow)

Features & Polynomial Regression:

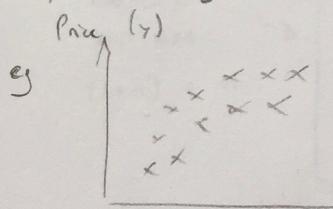
e.g. Housing Price Prediction: $h_\theta(x) = \theta_0 + \theta_1 \cdot \text{frontage} + \theta_2 \cdot \text{depth}$



↳ possibly decide area is real feature!

↳ $h_\theta(x) = \theta_0 + \theta_1 x$ where $x = \text{frontage} \cdot \text{depth}$

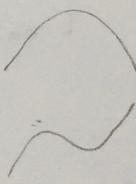
Polynomial Regression:



→ fit quadratic: $\theta_0 + \theta_1 x + \theta_2 x^2$

→ fit cubic: $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$

↳ $x = \text{size}$



⇒ $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$ where $x_1 = \text{size}$

$$x_2 = \text{size}^2$$

$$x_3 = \text{size}^3$$

↳ say size ranges $1 \leq x \leq 1000 \leftarrow \text{FEATURE SCALING V. IMPORTANT}$

↳ another possibility: $h_\theta(x) = \theta_0 + \theta_1 \text{size} + \theta_2 \sqrt{\text{size}}$

Computing Parameters Analytically

Normal Equation

Minimizing before $\rightarrow \frac{d}{d\theta} J(\theta) \rightarrow \text{Set } = 0 \rightarrow \text{solve for real } \# \theta.$

$\hookrightarrow \theta$ is now a vector $\rightarrow \theta \in \mathbb{R}^{n+1}$, $J(\theta_0, \theta_1, \theta_2, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

$\Rightarrow \frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad \text{for each } j \rightarrow \text{solve for } \underbrace{\theta_0, \theta_1, \theta_2, \dots, \theta_n}_{\text{values that minimize } J(\theta)}$

ex] $m = 4$ training examples

x_0	size x_1	# beds x_2	# floors x_3	age (yr) x_4	price y
1	2104	5	1	45	460
:	:	:	:	:	:
1	852	2	1	36	178

$$\Rightarrow X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & \cdot & \cdot & \cdot & \cdot \\ 1 & \cdot & \cdot & \cdot & \cdot \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \quad m \times (n+1)$$

$$y = \begin{bmatrix} 460 \\ \vdots \\ 178 \end{bmatrix} \quad m \times 1$$

\Rightarrow Skips derivation, $\rightarrow \theta_{\min} = (X^T X)^{-1} X^T y$

\Rightarrow generally: m training examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ & n features

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \Rightarrow X = \begin{bmatrix} \cdots & (x^{(1)})^T \\ \cdots & (x^{(2)})^T \\ \vdots & \vdots \\ \cdots & (x^{(m)})^T \end{bmatrix} \quad \begin{array}{l} \text{DESIGN} \\ \text{MATRIX} \\ m \times (n+1) \end{array}$$

* Here, feature scaling not necessary! Not using gradient descent

\Rightarrow w/ m training examples, n features:

Gradient Descent	Normal Eqn
• need to choose α	• no need to choose α
• need many iterations	• no iterations
• works well, even w/ large n	• slow if n is very large $\rightarrow (X^T X)^{-1}$ is $n \times n$, $O(n^3)$
• $O(kn^2)$	\hookrightarrow max $n \approx 10,000$

Normal Equation Non-Invertibility:

$\theta = (X^T X)^{-1} X^T y \rightarrow$ what if $X^T X$ non-invertible / singular / degenerate?

↳ Should happen pretty rarely

↳ octave has pinv & inv \rightarrow using pinv guarantees no problem

↳ causes for non-invertibility:

- redundant features (e.g. x_1 = size in feet, x_2 = size in meters)
- ↳ linearly dependent

• too many features (e.g. $m=10$, $n=100$) $\rightarrow m \leq n$
 ↳ $\notin \mathbb{R}^{10}$

↳ delete linearly dependent features or excess features / regularization