

Linear Regression with One Variable

Model Representation:

regression: predict real-valued output \rightarrow predict house prices

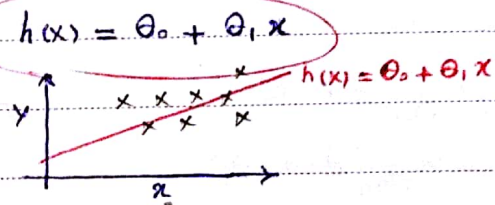
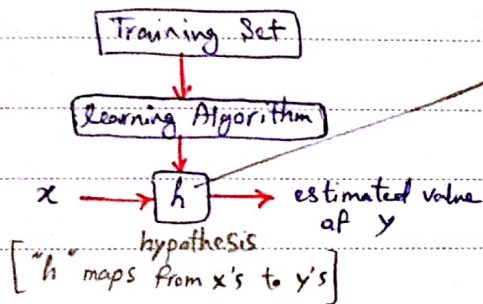
classification: classify into discrete-valued output \rightarrow tumor classification (malignant or benign)

مکروسکوپیک مکتوبیات با نظارت (Supervised) است؟ زیرا ما داده‌ها را به دست می‌آوریم و از آن‌ها یاد می‌گیریم.

training set $\left\{ \begin{array}{l} m = \text{number of examples} \\ x\text{'s} = \text{input variable / features} \\ y\text{'s} = \text{output variable / target} \end{array} \right.$

$(x, y) = \text{one training example}$

$(x^{(i)}, y^{(i)}) = i^{\text{th}} \text{ training example}$



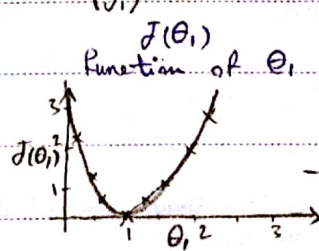
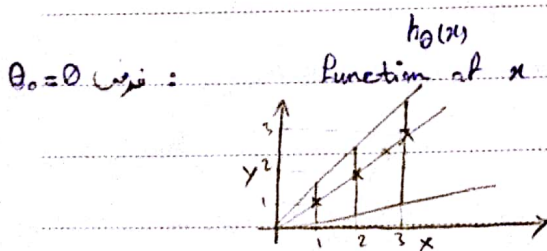
univariate linear regression (linear reg. with 1 variable)

چرا تابع مکتوبیات را خطی فرض کردیم؟ زیرا راحتی کار (در شروع)

Cost Function:

Find best θ_0 and θ_1 that minimize
$$\frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$= J(\theta_0, \theta_1)$$
 Cost Func. (Squared error function)



minimize $J(\theta_1)$

اگر θ_0 را داشته باشیم، آن‌گاه $J(\theta_0, \theta_1)$ به دو متغیر بستگی دارد. پس یک نمودار سه بعدی فراهم داشت: $\theta_0, \theta_1, J(\theta_0, \theta_1)$. این نمودار سه بعدی (Contour) شبیه یک کاسه (bowl-shaped) خواهد بود که تغییر آن جابجایی که تابع هزینه مینی‌موم می‌شود.

باید اکثریت داشته باشیم که بتواند به صورت اتوماتیک θ_0 و θ_1 مشخص را بیابد که با آن‌ها تابع $J(\theta_0, \theta_1)$ مینی‌موم شود.

Gradient Descent:

تابع gradient descent برای پیدا کردن یک تابع به کار می رود و برای این کار هدف کارایی تابع را کم کردن است. Have some function $J(\theta_0, \theta_1)$ want min $J(\theta_0, \theta_1)$ start with some θ_0, θ_1 keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ end up at a minimum $J(\theta_0, \theta_1)$

Gradient descent algorithm:

$$\text{repeat until convergence } \left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \right\} \quad (\text{for } j=0 \text{ and } j=1)$$

learning rate derivative term

نرخ یادگیری (α) مشخص می کند که در هر تکرار مقدار θ_0 و θ_1 چقدر تغییر می کند. اگر α بزرگ باشد مقدار θ_0 و θ_1 به سرعت تغییر می کند (معمولاً بهتر است). اگر α کوچک باشد مقدار θ_0 و θ_1 به آرامی تغییر می کند و ممکن است به حداقل نرسد.

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

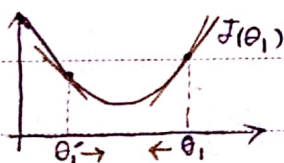
$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

(Simultaneous update)

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

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

Gradient Descent Intuition:



فرض $\theta_0 = 0$ و $\theta_1 \in \mathbb{R}$ به صورت متناوب باشد:

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

شیب منفی است

$$\theta_1 := \theta_1 - \alpha \quad (\text{positive number})$$

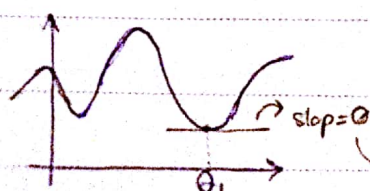
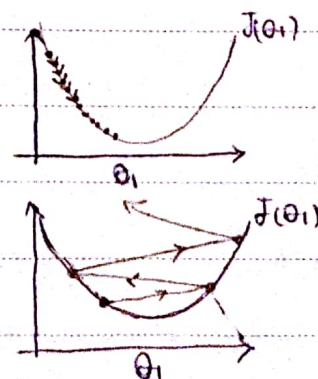
$$\theta_1' := \theta_1 - \alpha \quad (\text{negative number})$$

در نقطه θ_1 شیب منفی است

if α is too small, gradient descent can be slow

نرخ یادگیری (α)

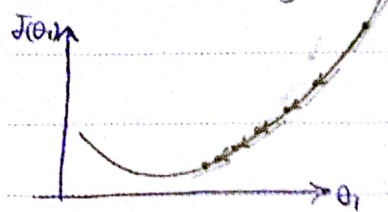
if α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge!



اگر θ_1 در نقطه ای که شیب صفر باشد (مثلاً در نقطه $\theta_1 = 0$) قرار گیرد:

$$\theta_1 := \theta_1 - (\alpha \cdot 0) \rightarrow \theta_1 := \theta_1$$

نکته مهم: در صورت انجام تکرارهای متوالی gradient descent نیازی به تغییر مقدار α نیست! کافایت α مقادیر ثابت باشد.
 هر چه به نقطه مینی موم نزدیک تر می شویم، α باید بزرگ تر می شود، در نتیجه $\alpha \propto \frac{\partial}{\partial \theta_i} J(\theta_i)$ کوچک تر می شود!



Gradient Descent for Linear Regression:

باید ابتدا مشتق جزئی تابع J را نسبت به θ_0 و θ_1 بیابیم:

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\hat{h}_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

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

نکته: تابع هدف برای گرادیان منفی همواره یک تابع Convex Function است و یک استیما بیش تر ندارد. در نتیجه هرگز با gradient descent به استیما محلی بر نمیخوریم.

به الگوریتم مورد بحث batch gradient descent هم گفته می شود؛ چرا که در هر تکرار از تمام دیتاست (training examples) استفاده می کنیم.
 روش های دیگری مانند Normal Equation Method هم برای محاسبه مقدار مینی موم تابع هدف وجود دارند. اما روش gradient descent در دیتاست های بزرگ تر بهتر از روش normal equation عمل می کند.

Linear Regression with Multiple Variables

Multiple Features:

m = number of samples

n = number of features

در گرادیان منفی متغیره، هم رکورد تنها یک عدد نیست بلکه یک ستون از اعداد است که هر عدد بیانگر مقدار هر ویژگی برای آن رکورد است:

$$X^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^n$$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{m \times n}$$

hypothesis تابع: $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Example: $h_{\theta}(x) = 80 + 0.1 x_1 + 0.01 x_2 + 3 x_3 - 2 x_4$
 (base price) (price per square foot) (each floor) (age of the house)

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

$$h_{\theta}(x) = \theta^T \cdot X = [\theta_0 \theta_1 \dots \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

Gradient Descent For Multiple Variables:

parameters: $\theta_0, \theta_1, \dots, \theta_n \rightarrow \theta: (n+1)$ -dimensioned vector

Cost function: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

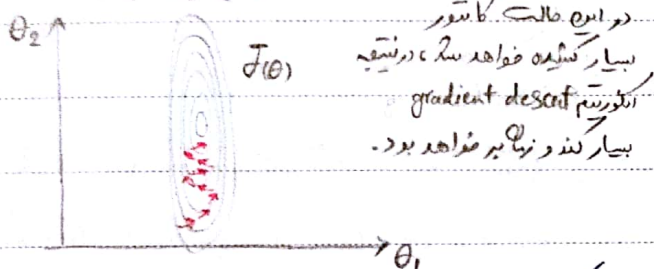
Gradient descent: Repeat $\left\{ \begin{aligned} \theta_j &:= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \\ &= \theta_j - \alpha \left(\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) \end{aligned} \right.$
(Simultaneously update θ_j for $j=0,1,\dots,n$)

Feature Scaling:

Idea: Make sure features are on a similar scale

Example: $0 \leq x_1 \leq 2000$

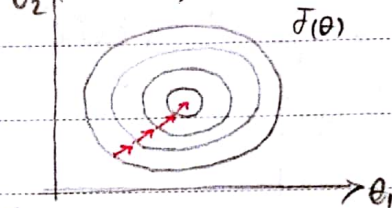
$1 \leq x_2 \leq 5$



Feature scaling

$x_1 = \frac{x_1}{2000} \Rightarrow 0 \leq x_1 \leq 1$

$x_2 = \frac{x_2}{5} \Rightarrow 0 \leq x_2 \leq 1$



اگر دامنه تغییرات هر یک از تغییرها (x) یکسان نباشد اما نزدیک باشد، مشکلی نیست. اما اگر دامنه تغییرات بسیار متفاوت باشد باید سعی کنیم هکلی تغییرها در بازه $-1 \leq x_i \leq 1$ قرار بگیرند.

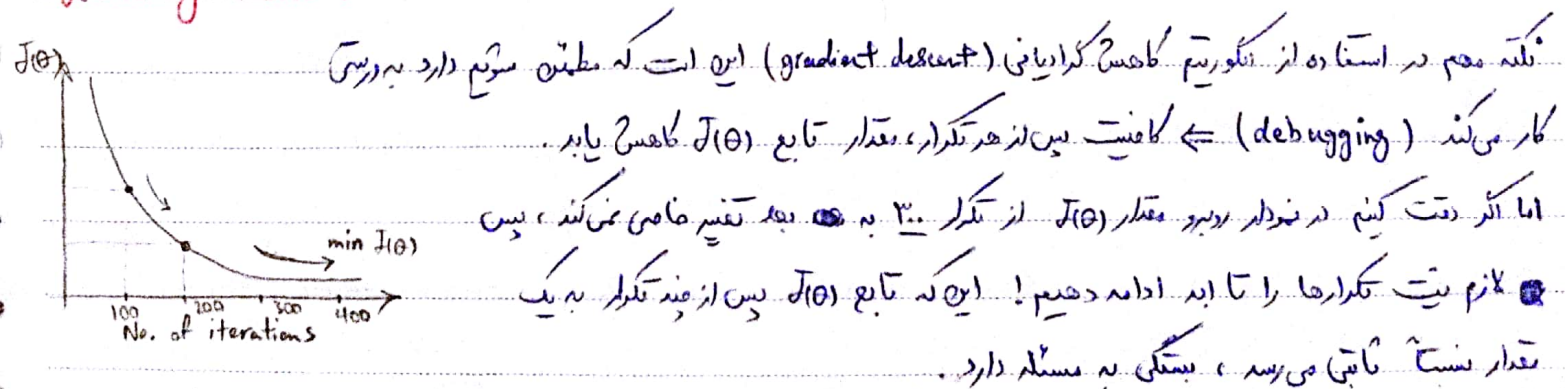
یکی از روش‌های کوچک کردن بازه تغییرات mean normalization است:

Example: $0 \leq x_1 \leq 2000 \rightarrow x_1 = \frac{x_1 - 1000}{2000} \Rightarrow -0.5 \leq x_i \leq 0.5$

$$x_i \leftarrow \frac{x_i - \mu_i}{s_i}$$

 μ_i : میانگین بازه (μ_i)
 s_i : Standard deviation $(\max(x_i) - \min(x_i))$
 average value $(\frac{\max(x_i) + \min(x_i)}{2})$

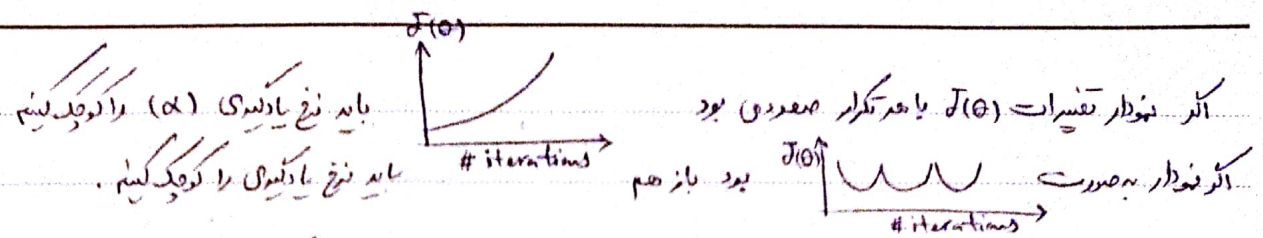
Learning Rate:



Declare convergence if $J(\theta)$ decreases less than $\epsilon = 10^{-3}$ in one iteration

می‌توانیم برای الگوریتم یک convergence (همگرایی) مشخص کنیم:

اما مقدار ϵ چقدر باید باشد، پس تعیین از روی نمودار اولویت دارد.



Summary:

بسیار کوچک: slow convergence

بسیار بزرگ: $J(\theta)$ may not decrease on every iteration; may not converge (slow convergence also possible)

روش انتخاب درست نرخ یادگیری؟ صریح و خطا

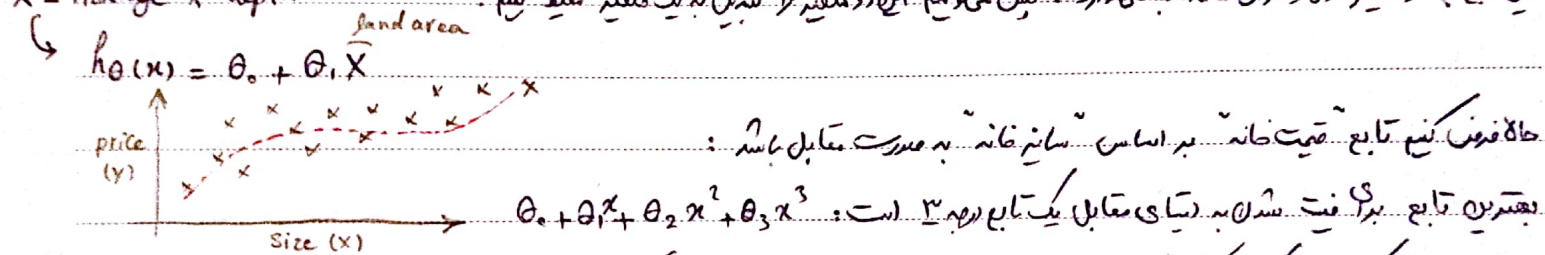
Features and Polynomial Regression:

فرض کنیم تابع مساحت خانه به صورت زیر باشد:

$h_\theta(x) = \theta_0 + \theta_1 \times \text{Frontage} + \theta_2 \times \text{depth}$

$X = \text{Frontage} \times \text{depth}$

این تابع به درجه اول و عرض خانه بستگی دارد، پس می توانیم این درجه اول را تبدیل به یک معین خطی کنیم:



با توجه به آن که تنها یک ویژگی بزرگ داریم (مساحت) می توانیم سه معین بالا را این گونه بسازیم:

$h_\theta(x) = \theta_0 + \theta_1(x_1) + \theta_2(x_2) + \theta_3(x_3)$

در دو مثال بالا دیدیم که می توانیم چند معین را ادغام کنیم یا معین های جدید به اضافه قبلی ها بسازیم!

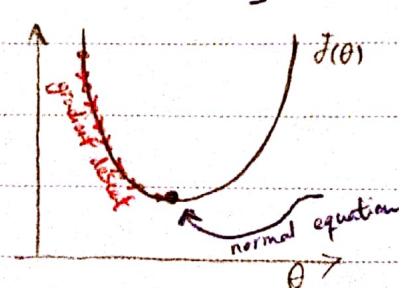
نکته: برای نمودار محدود به تابع درجه ۳ ذکر شده نیستیم؛ ممکن است ترابع دیگری داشته $h_\theta(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\sqrt{\text{size}})$ مناسب باشند.

روش انتخاب بهترین تابع، شکل داده های نمودار است: تابع انتخابی باید به بهترین شکل به دیتا fit شود.

به توابع غیر خطی که در بالا مشاهده کردیم، توابع کلم سیرال (چند جمله ای) (polynomial) گفته می شود.

Normal Equation:

گاهی اوقات بهتر است به جای استفاده از gradient descent، تابع $J(\theta)$ را حل کنیم تا θ ی بهینه را بیابیم.



برای مثال در تابع روبه بالا اگر بخواهیم گامی گامی قدم به قدم مقدار θ را تغییر می دهیم تا به θ ی بهینه برسیم. در حالی که با حل یک معادله ساده می توانیم بهترین θ بزرگترین θ را پیدا کنیم.

$$J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{\partial}{\partial \theta} J(\theta) = \dots = 0 \rightarrow \text{Solve for } \theta!$$

Subject

Date

x_0	x_1	x_2	x_3	x_4	y
	Size	no. of bedrooms	no. of floors	age	price

$m = \dots$
 $n = 4$

design matrix

$$y = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{m \times 1}$$

$$X = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}_{m \times (n+1)}$$

فرض کنیم داده ها ما به صورت روبه رو باشد

آن کاد مقدار بهینه ماتریس θ به صورت زیر به دست می آید :

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

Octave/MATLAB : $\text{pinv}(X' * X) * X' * y$
 $(X^T X)^{-1}$

به این روش حل معادله یافتن θ روش Normal Equation گفته می شود.

Gradient Descent

Normal Equation :

- need to choose α
- needs many iterations
- works well even when n is large.

- No need for α
- No need to iterate
- Need to compute $(X^T X)^{-1}$ $n \times n$ matrix $\rightarrow O(n^3)$
- slow if n is very large

نکته : برای n های کوچک روش normal equation بهتر است ؛ اما اگر n بسیار بزرگ باشد (مثلاً $n > 10000$) روش گرادینت برای یافتن جواب خواهد داد.

Normal Equation Noninvertibility :

در معادله $\theta = (X^T X)^{-1} X^T y$ اگر $X^T X$ معکوس پذیر نباشد ، چه اتفاقی خواهد افتاد ؟
این اتفاق به ندرت رخ می دهد . اما توجه کنیم که Octave/MATLAB دو عملگر معکوس کننده دارد : inv و pinv
دستور pinv که ما از این استفاده می کنیم ، معکوس $X^T X$ را محاسبه می کند ، حتی اگر معکوس پذیر نباشد !

★ در این جا وارد جزئیات ریاضیات دستور pinv و تفاوت آن با inv نمی شویم . اما می خواهیم بدانیم که در چه صورتی ماتریس $X^T X$ معکوس ناپذیر خواهد بود ؟

(۱) اگر Feature ها اضافی داشته باشیم (به جایگزین هم تبلی خطی دارند) مثلاً :
 $\begin{cases} x_1 = \text{Size in } \text{ft}^2 \\ x_2 = \text{Size in } \text{m}^2 \end{cases}$

(۲) اگر Feature ها بسیار زیادی داشته باشیم (مثلاً $m \leq n$) ← آن گاه باید برخی Feature ها را حذف کنیم یا از رگرسیون استفاده کنیم

در آینده راجع به آن خواهم نوشت

پس اگر ماتریس $X^T X$ معکوس ناپذیر بود باید به فکر کاهش تعداد Feature ها باشیم .