# برازش خطی چندمتغیره

❖ **برازش خطی چندگانه** (Multiple Linear Regression) برای پیش‌بینی یک متغیر وابسته ($y$) که به صورت خطی با بیش از یک متغیر مستقل ($x_1, x_2, \ldots, x_n$) مرتبط است به کار می‌رود.

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n = \theta^T X$$

$$x_0 = 1$$

$x_j^{(i)}$: مقدار ویژگی $j$ در نمونه $i$

$x^{(i)}$: ویژگی‌های نمونه $i$ که می‌توانیم آن‌ها را به شکل بردار $X = [1, x_1, x_2, \ldots, x_n] \in \mathbb{R}^{n+1}$ تعریف کنیم. $x_0^{(i)} = 1$ است.

$\theta$: بردار ضرایب که به شکل $\theta = [\theta_0, \theta_1, \theta_2, \ldots, \theta_n] \in \mathbb{R}^{n+1}$ آن را تعریف می‌کنیم.

$J(\theta) : MSE \longrightarrow$
$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} = \theta_j - \frac{\alpha}{m} \sum_{i=1}^{m} (h_\theta(x)^{(i)} - y^{(i)}) x_j^{(i)}$$

# Multiple features

Notation:

    n= number of features

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

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

| Size (feet²) | Price ($1000) |
|---|---|
| $x$ | $y$ |
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

|  | $x_1$ | $x_2$ | $x_3$ | $x_4$ |  |
|---|---|---|---|---|---|
|  | Size (feet²) | Number of bedrooms | Number of floors | Age of home (years) | Price ($1000) |
|  | 2104 | 5 | 1 | 45 | 460 |
|  | 1416 | 3 | 2 | 40 | 232 |
|  | 1534 | 3 | 2 | 30 | 315 |
|  | 852 | 2 | 1 | 36 | 178 |
|  | ... | ... | ... | ... | ... |

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

**Multiple features**

$x^{(i)} \longrightarrow y^{(i)}$

Hypothesis

Previously: $h_\theta(x) = \theta_0 + \theta_1 x$

Multivariate linear regression: $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$

$\times x_0$

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

$\theta = [\theta_0, \theta_1, \ldots, \theta_n]^T$

$x = [x_0, x_1, \ldots, x_n]^T$

size $n+1$

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|-------|-------|-------|-------|-----|
| Size (feet²) | Number of bedrooms | Number of floors | Age of home (years) | Price ($1000) |
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| ... | ... | ... | ... | ... |

$X^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$

$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$

# Gradient descent for multiple variables

Hypothesis: $h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$

Parameters: $\theta_0, \theta_1, \ldots, \theta_n$

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

Gradient descent:

Repeat $\{$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \ldots, \theta_n)$$

$\}$        (simultaneously update for every $j = 0, \ldots, n$)

# Gradient Descent

Previously (n=1):

Repeat {

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

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

(simultaneously update $\theta_0, \theta_1$ )

}

New algorithm $(n \geq 1)$:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$$

(simultaneously update $\theta_j$ for $j = 0, \dots, n$)

}

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

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

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

...

[handwritten annotations: $n=1$, $n>1$, $x_0^{(i)} = 1$, $\theta_j = \theta - \alpha \frac{\partial J}{\partial \theta}$]

# Feature Scaling

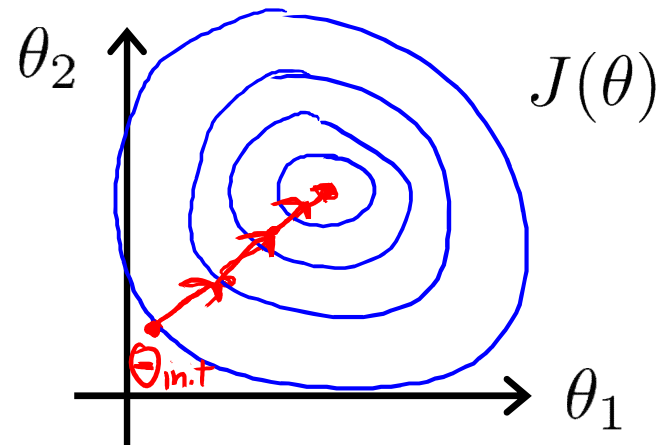Idea: Make sure features are on a similar scale.

E.g. $x_1$ = size (0-2000 feet²)

$x_2$ = number of bedrooms (1-5)

$$x_1 = \frac{\text{size (feet}^2)}{2000}$$

$$x_2 = \frac{\text{number of bedrooms}}{5}$$

$0 \leq x_1 \leq 1$     $0 \leq x_2 \leq 1$

# Mean normalization

Replace $x_i$ with $x_i - \mu_i$ to make features have approximately zero mean
(Do not apply to $x_0 = 1$).

E.g.    $x_1 = \frac{size - 1000}{2000}$

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

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

$$x_1 \leftarrow \frac{x_1 - \boxed{\mu_1}}{\boxed{S_1}}$$

← avg value of $x_1$ in training set

range (max - min)
(or standard deviation)

$$x_2 \leftarrow \frac{x_2 - \mu_1}{S_2}$$

# تغییر مقیاس متغیرها

Gradient descent
without scaling

Gradient descent
after scaling variables

$x_1 \gg x_2$

$w_2$      $J(w)$

$w_1$

$0 \leq x_1 \leq 1$
$0 \leq x_2 \leq 1$

$w_2$      $J(w)$

$w_1$

min-max normalization

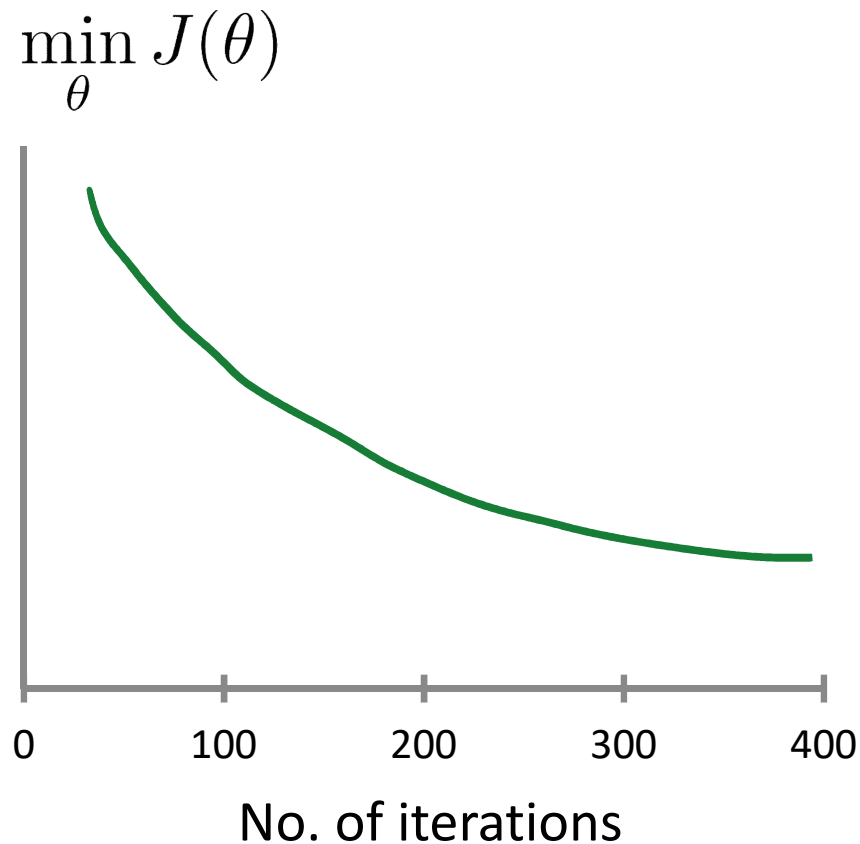$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

Mean normalization

$$x' = \frac{x - mean(x)}{max(x) - min(x)}$$

**Gradient descent**

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

- "Debugging": How to make sure gradient descent is working correctly.
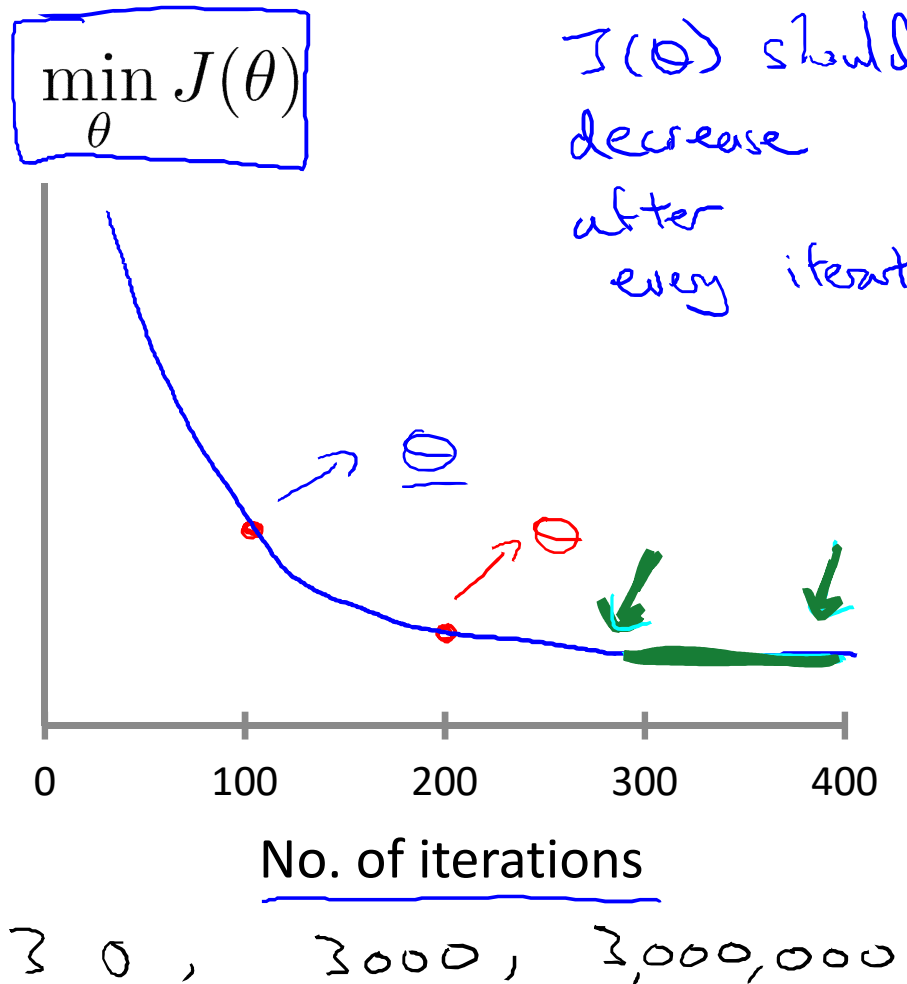
- How to choose learning rate $\alpha$.

**Making sure gradient descent is working correctly.**

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



No. of iterations

Example automatic convergence test:

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

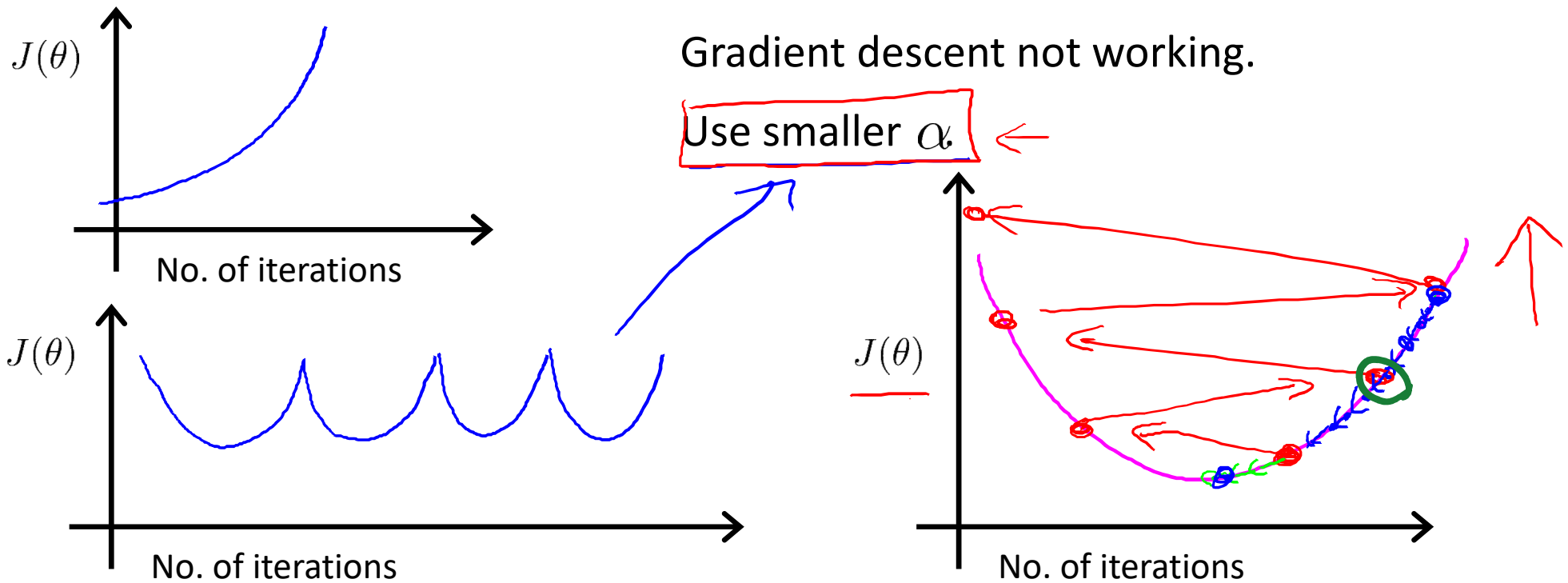**Making sure gradient descent is working correctly.**



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

No. of iterations

$J(\theta)$ should decrease after every iteration.

3  0 ,    3000 ,   3,000,000

→ Example automatic convergence test:

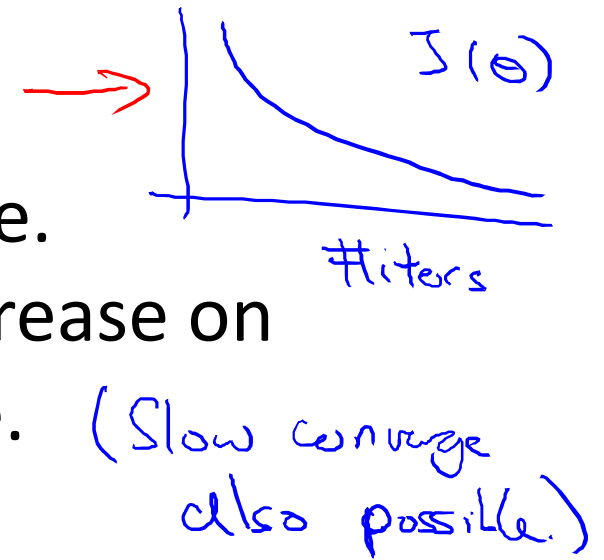→ Declare convergence if $J(\theta)$ decreases by less than $10^{-3}$ in one iteration.

$\varepsilon$

# Making sure gradient descent is working correctly.



Gradient descent not working.

Use smaller $\alpha$.

- For sufficiently small $\alpha$, $J(\theta)$ should decrease on every iteration.
- But if $\alpha$ is too small, gradient descent can be slow to converge.

**Summary:**

- If $\alpha$ is too small: slow convergence.
- If $\alpha$ is too large: $J(\theta)$ may not decrease on every iteration; may not converge.
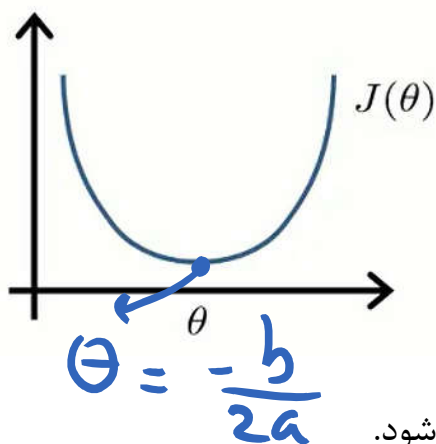
*J(θ)*

*#iters*

*(Slow converge also possible.)*

To choose $\alpha$, try

$$\ldots, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, \ldots$$

*3x    ~3x    3x    ~3x*

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

حالت ساده به‌دست آوردن کمینه یک تابع:

$$\frac{dJ(\theta)}{d\theta} = 0 \longrightarrow \quad \theta = -\frac{b}{2a}$$

$$\theta = -\frac{b}{2a}$$

اما در مسئله ما $\theta \in \mathbb{R}^{n+1}$ است و مشتق باید برای همه $\theta$ ها برابر صفر قرار بگیرد و حل شود.

$$\frac{\partial J(\theta)}{\partial \theta_j} = 0 \ (for\ every\ j) \longrightarrow \quad \theta_0, \cdots, \theta_n$$

$$j = 0 \cdots n$$

# Normal equation

(رگرسیون خطی)

Normal equation: Analytic solution

$\theta \in \mathbb{R}^{n+1} \Rightarrow J(\theta_0, \theta_1, \ldots \theta_n)$

$\frac{\partial J}{\partial \theta_j} = 0 \Rightarrow$

$j = 0 \ldots n$

بهینه $\theta_n \ldots, \theta$

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$X_0 \quad X_1 \quad X_2 \quad X_4 \quad X_6$

$3 \quad 4$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

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

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

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

$(X^T X)^{-1}$ is inverse of matrix $X^T X$

$m$ **training examples,** $n$ **features.**

Gradient Descent

- Need to choose $\alpha$.
- Needs many iterations.
- Works well even when $n$ is large.

$n = 10^6$

Normal Equation

- No need to choose $\alpha$.
- Don't need to iterate.
- Need to compute $(X^T X)^{-1}$ → $O(n^3)$
- Slow if $n$ is very large.

$n = 100$
$n = 1000$

* مقیاس کردن ورژنهای لازم نیت

# What if $X^T X$ is non-invertible? (singular/ degenerate)

$$\begin{bmatrix} x_0 & \dot{x}^{(1)T} & x_1 \\ & x^{(m)T} & \end{bmatrix}$$

- Redundant features (linearly dependent).

  E.g. $x_1$ $\boxed{\begin{array}{l} \text{size in feet}^2 \\ \text{size in m}^2 \end{array}}$ $\Rightarrow$ $X^T X$ singular

  $x_2$

- Too many features (e.g. $m \leq n$).

  $n = 100$  $\boxed{m = 10}$

  - Delete some features, or use regularization.

  $\theta:$

  $10 \not k \times 1$

# Housing prices prediction

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

$x_1$          $x_2$

<u>Area</u>

$x = \underline{frontage * depth}$

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

↖ land area

# برازش چندجمله

❖ **برازش چندجمله‌ای** (Polynomial Regression) برای پیش‌بینی متغیر وابسته‌ای که به صورت چندجمله‌ای با متغیرهای مستقل مرتبط است به کار می‌رود.

| | |
|---|---|
| **مثال** | **برازش چندجمله‌ای:** پیش‌بینی قیمت خانه |



قیمت (y)

متراژ (x)

$$h_\theta(x) = \theta_0 + \theta_1(size)^1 + \theta_2(size)^2 + \theta_3(size)^3$$

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

# Polynomial regression



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

$$x_1 = x \qquad x_2 = x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

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

$$x_1 = (size)$$
$$x_2 = (size)^2$$
$$x_3 = (size)^3$$

Size: $1 - 1000$

Size$^2$: $1 - 1000,000$

Size$^3$: $1 - 10^9$

# Choice of features



Price (y)

Size (x)

$x : size$

$h_\theta(x) = \theta_0 + \theta_1(size) + \theta_2(size)^2$

$h_\theta(x) = \theta_0 + \theta_1(size) + \theta_2\sqrt{(size)}$

$x_1 = x$     $x_2 = \sqrt{x}$

# Stochastic Gradient Descent vs Batch Gradient Descent

## Batch Gradient Descent

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

$$\theta_{t+1} = \theta_t - \alpha J'(\theta)$$

Need to compute for all data,
if you have sample size, m = 1million
Very slow to update θ

## Stochastic Gradient Descent

$$cost\left(\theta, (x^{(i)}, y^{(i)})\right)$$
$$= \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2$$

1. Random Shuffle Data
2. Repeat {
   for i = 1, ... , m {

   $$\theta_{t+1}$$
   $$= \theta_t - \alpha \frac{\partial}{\partial \theta_t} cost\left(\theta, (x^{(i)}, y^{(i)})\right)$$

   }
}

Update θ using only one data point.
Faster

ممکن است حول $\theta_{opt}$ نوسان کند $\leftarrow \alpha$ دهنده باشد.

# Mini-batch Gradient Descent

❖ اگر تعداد نمونه ها (m)، خیلی بزرگ باشد، به چندین **mini-batch** تقسیم می شوند و  در هر مرحله

یکی از  **mini-batch**ها ر به هنگام رسانی پارامترها استفاده می شود.

$$B.S : 32, 64, 128, \dots$$

$$m_B = \frac{m}{B.S.}$$

مثلاً 32نمونه

$X^{\{1\}}$

$X^{\{m_B\}}$

1 epoch

Random shuffle
در هر mini-batch

$$\widehat{\theta_0} = \widehat{b}$$
$$[\theta_1 \cdots \theta_n]' = \underline{W}$$

$$X = \begin{bmatrix} \vdots & \vdots & & \vdots \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ \vdots & \vdots & & \vdots \end{bmatrix}_{n_x \times m}$$

# Mini-batch Gradient Descent

$$\frac{m}{t} = {}^{m}B$$

$t = 1 \rightarrow Stochastic$

$t = m \rightarrow Batch$

$t = 64, 128, 256, 512 \rightarrow Mini - batch$

$for\ every\ t:$

     Optimize based on $X^{(1:t)}$

     Compute cost $J$

     Update θ

Computational resource per epoch

Stochastic      Mini-batch      Number of datapoints

Batch

Epochs required to find good W, b values

$$\left| X^{(1)} \cdots X^{(64)} \right| X^{(65)} \cdots X^{(128)} \right| \cdots$$
$$X^{\{1\}} \qquad X^{\{2\}}$$

$J$ — Batch GD — Iteration

$J$ — mini-batch GD — Iteration

* ممکن است یادگیری $X^{\{2\}}$ از $X^{\{1\}}$ مشکل تر باشد $\Rightarrow$ $J^{\{2\}} > J^{\{1\}}$

$1 < m_B < m$ $\qquad$ $m_B = m \longrightarrow$ stochastic G.D.

$m_B = 1 \longrightarrow$ Batch G.D.

step size 0.100



step size 0.100

$$1 < m_B < m$$

Fastest learning

**Batch Gradient Descent**

**Stochastic Gradient Descent**

**Too long per iteration especially when $m$ is very large**

**Lose speedup from vectorization (Inefficient implementation)**

## Newton Method

$J(\theta)$     $J'(\theta) = 0$

$\theta_{init}$   $\theta_{t+1} = \theta_t - \dfrac{J'(\theta_t)}{J''(\theta_t)}$

نیوتن متغیره تک

$\underline{\theta} = [\theta_0 \ .. \ \theta_n]^T_{n+1 \times 1}$

$J(\underline{\theta})$   $\nabla_\theta J(\underline{\theta}_{opt}) = \left[\dfrac{\partial J}{\partial \theta_0}, .. , \dfrac{\partial J}{\partial \theta_n}\right]^T$

$\underline{\theta}_{init}$   $\underline{\theta} := \underline{\theta} - H^{-1} (\nabla_\theta J(\theta))$

نیوتن راستون

$H_{(n+1)(n+1)} : H_{ij} = \dfrac{\partial J}{\partial \theta_i \theta_j}$

- Faster convergence

$x_{init.}$
$x_{t+1} = x_t - \dfrac{f(x_t)}{f'(x_t)}$

ریشه‌یابی : ریشهٔ $f(x)$

$f(x)$

$(x_n, f(x_n))$  Slope = $f'(x_n)$

$y = f(x)$

y

$(x_{n+1}, 0)$
$x_{n+1}$      $x_n$

$\left[\begin{array}{cc} \dfrac{\partial^2 J}{(\partial \theta_0)^2} & .. \\ & \dfrac{\partial^3}{\partial \theta_n^2} \end{array}\right]$

هر مت همگرای سریع ۔ راصون می تواند سریع تر از GD ها باشد .

ولی $H^{-1}_{(n+1)(n+1)}$ برای $n$ بزرگ هزینه محاساتی زیادی دارد .

وقتی روش نیوتن راصون برای حل مسأله Log. به کارگی رود ←
Reg

Fisher Scoring

Newton Method vs Gradient Descent

# Logistic Regression

## Introduction

- Logistic regression is a widely used discriminative classification model $p(y \mid \mathbf{x}; \boldsymbol{\theta})$.

- where $\mathbf{x} \in \mathbb{R}^D$ is a fixed-dimensional input vector, $y \in \{1, \dots, C\}$ is the class label.

- If $C = 2$, this is known as <span style="color:red">binary</span> logistic regression.

- if $C > 2$, it is known as <span style="color:red">multinomial</span> logistic regression, or alternatively, <span style="color:green">multiclass logistic regression</span>.

**Classification**

Email: Spam / Not Spam?
Online Transactions: Fraudulent (Yes / No)?
Tumor: Malignant / Benign ?

$$y \in \{0, 1\}$$

0: "Negative Class" (e.g., benign tumor)
1: "Positive Class" (e.g., malignant tumor)

Threshold classifier output $h_\theta(x)$ at 0.5:

If $h_\theta(x) \geq 0.5$, predict "y = 1"

If $h_\theta(x) < 0.5$, predict "y = 0"

$h_\theta(x) = \theta^T X$

$-\infty \rightarrow \infty$

Classification:   y = 0 or 1        $h_\theta(X) = \theta^T X$

$h_\theta(x)$ can be > 1 or < 0

**Logistic Regression:**   $\boxed{0 \le h_\theta(x) \le 1}$

$h_\theta(X) = P(y = 1 \mid \underline{X})$

طبقه بندی باینری

# Logistic Regression Model

We want $\quad 0 \le h_\theta(x) \le 1$

$$h_\theta(x) = \sigma(\theta^T x)$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$H_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$\sigma(z)$

**Interpretation of Hypothesis Output** $\{x^{(i)}, y^{(i)}\}_{i=1}^{m}$ $y^{(i)} \in \{0, 1\}$

$h_\theta(x)$ = estimated probability that y = 1 conditioned on input x

Example: If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

$h_\theta(x)_{test} = 0.7 \longrightarrow$ صحيح

$$P(y = 0|x; \theta) + P(y = 1|x; \theta) = 1$$
$$P(y = 0|x; \theta) = 1 - P(y = 1|x; \theta)$$

$h_\theta(x)_{test} = 0.3 \longrightarrow$ حرش ختم

$h_\theta(x^{(i)}) = \sigma(\theta^T x^{(i)})$

$x^{(i)} \longrightarrow y^{(i)}$

$h_\theta(x) = P(y = 1 | X)$

$1 - h_\theta(x) = P(y = 0 | X)$

$x_{test} \longrightarrow \dfrac{1}{1 + e^{-\theta^T x_{test}}} \gtrless^{th} 0.5$ صحيح
حوش حتم

# Logistic regression

$$h_\theta(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1+e^{-z}}$$



Suppose predict "$y = 1$" if $h_\theta(x) \geq 0.5$

predict "$y = 0$" if $h_\theta(x) < 0.5$

th $= 0.5$ مرز تصمیم گیری :

$$\bar{\theta}^T x = 0 \Rightarrow g(\bar{\theta}^T x) = \frac{1}{2}$$

$$\bar{\theta}^T x \geq 0 \longrightarrow y = 1$$

$$\bar{\theta}^T x < 0 \longrightarrow y = 0$$

**Decision Boundary**



$$\theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2$$
$$x_0 = 1$$

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

$(x^{(i)}, y^{(i)})$
$i = 1 \ldots m$
$\Rightarrow \theta_0 = -3 \quad \theta_1 = \theta_2 = 1$

Predict " $y = 1$ " if $-3 + x_1 + x_2 \geq 0 \longrightarrow y = 1$

مرز تصمیم گیری

$\theta^T x > 0 \quad < 0 \longrightarrow y = 0$

$h = 0.5$

# Logistic Unit:

$x_0 = 1$

$x_1$

$x_2$

$\vdots$

$x_n$

$\theta_0$

$\theta_1$

$\theta_2$

$\theta_n$

$\sigma(\theta^T x)$ $- y$

$0 \leq h_\theta(x) \leq 1$

$$P_\theta(y = 1 \mid x) = h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}} \gtrless 0.5$$

$$\theta^T x \gtrless 0$$

$\sigma(\theta^T x)$   $\theta_0, \theta_1, \theta_2$

## Non-linear decision boundaries

$h_\theta(x) = \sigma(P_\theta(\underline{x}))$

$P_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2$

Transform input features in suitable way

$$\phi(x_1, x_2) = [1, x_1^2, x_2^2] \quad x_1 \wedge x_2, x_1, x_2$$

$$w = [-R^2, 1, 1]. \text{ Then } w^\top \phi(x) = \boxed{x_1^2 + x_2^2 - R^2} \geq 0 \rightarrow y = 1$$
$$< 0 \rightarrow y = 0$$

Decision boundary (where f(x) = 0) defines a circle with radius R



$[1, x_1, x_2]$

تَبدیل

$[1, x_1^2, x_2^2]$

با مرز تصمیم گیری حطی تفکیک پذیر نیست.

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \cdots, (x^{(m)}, y^{(m)})\}$

m examples $\quad x \in \begin{bmatrix} x_0 \\ x_1 \\ \cdots \\ x_n \end{bmatrix} \quad x_0 = 1, y \in \{0, 1\}$
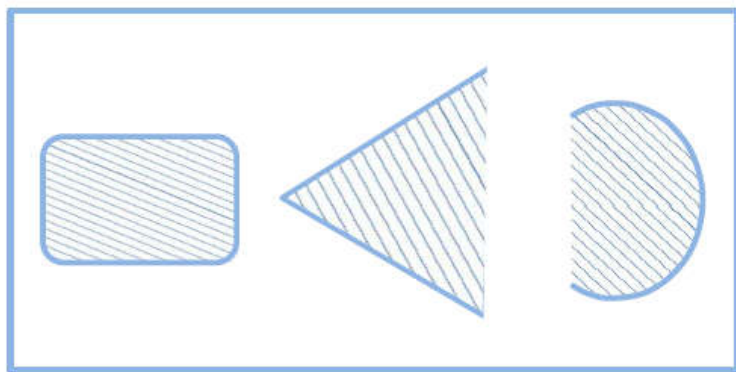
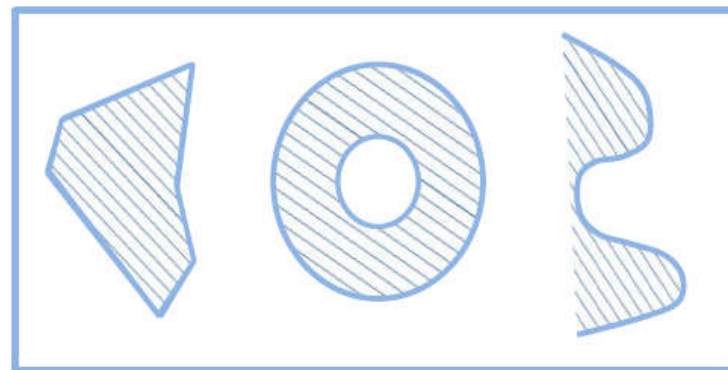$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose parameters $\theta$ ?

## Cost function

Linear regression:  $J(\theta) = \frac{1}{m} \sum\limits_{i=1}^{m} \frac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$   $J(\theta)$

$\sigma(\theta^T x^{(i)})$

$\text{Cost}(h_\theta(x^{(i)}), y^{(i)}) = \frac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$  →  غیر محدب
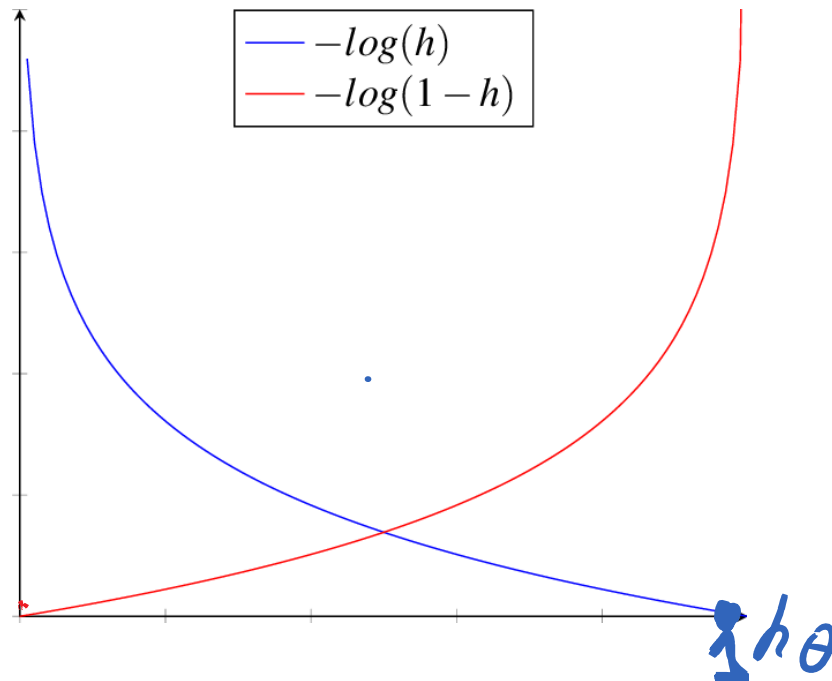
Convex                    Not Convex

**Logistic regression cost function**

$$h_\theta(x) = P(y=1 \mid x)$$

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

Convex



Legend:
- $-log(h)$
- $-log(1-h)$

$h_\theta$

## Logistic regression cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

Note: $y = 0$ or $1$ always

$$\text{Cost}(h_\theta(x), y) = -y\log(h_\theta(x)) - (1-y)\log(1 - h_\theta(x))$$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

## Logistic regression cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \log \underbrace{h_\theta(x^{(i)})}_{P(y=1|X)} + (1 - y^{(i)}) \log \underbrace{(1 - h_\theta(x^{(i)}))}_{P(y=0|X)} \right]$$

# Gradient Descent

$\theta$

$$J(\theta) = -\frac{1}{m}\big[\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)}))\big]$$

Want $\min_\theta J(\theta)$:

Repeat $\{$

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

$\}$  (simultaneously update all $\theta_j$)

$j = 0 \cdots n$

$$\frac{\partial}{\partial \theta_j} \sigma(\theta^T x)?$$

---

$$\frac{\partial}{\partial z} \sigma(z) = \sigma(z)[1 - z]$$

True fact about sigmoid functions

$$\frac{\partial}{\partial \theta_j} \sigma(\theta^T x) = \frac{\partial}{\partial z} \sigma(z) \cdot \frac{\partial z}{\partial \theta_j}$$

Chain rule!

$$\frac{\partial}{\partial \theta_j} \sigma(\theta^T x) = \sigma(\theta^T x)[1 - \sigma(\theta^T x)] x_j$$

Plug and chug

## Gradient Descent

$$J(\theta) = -\frac{1}{m}\left[\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))\right]$$

Want $\min_\theta J(\theta)$:  $\underline{\theta_{init}}$

Repeat $\{$

$$\theta_j := \theta_j - \alpha \sum_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right)x_j^{(i)}$$

$= \sigma'(\theta^T x)$

(simultaneously update all $\theta_j$)

$j = 0 \dots n$

$$\underline{\theta} = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}$$

$\}$

Algorithm looks identical to linear regression!

## logistic regression (Probabilistic view)

$$\underline{w}^T \underline{x} + b \equiv \theta^T x$$

Binary logistic regression often follows the following model

$$\underbrace{\theta_0}_{b} + \underbrace{\theta_1 x_1 + \cdots + \theta_n x_n}_{\underline{w}^T \underline{x}}$$

$$p(y|x; \theta) = \mathrm{Ber}(y|\sigma(w^T x + b))$$

Bernoulli : $p = \sigma(\theta^T x)$

Bernoulli    sigmoid       bias
                        weight

$$f(k; p) = \begin{cases} p & \text{if } k = 1, \\ q = 1 - p & \text{if } k = 0. \end{cases}$$

$$p(y = 1 | x; \theta) = \sigma(a) = \frac{1}{1 + e^{-a}}, \text{ where } a = \log \frac{p}{1 - p}$$

logistic regression (Probabilistic view)

**Finding the Maximum Likelihood (ML) solution is equivalent to minimizing the cross entropy cost function**

$$P(y \mid X; \theta) \qquad P_\theta(Y \mid X)$$

$$\max_\theta \log P_\theta(y \mid X) \equiv \min_\theta J(\theta)$$

$$\text{log\_likelihood} \qquad \qquad \text{cross entropy}$$

$$P_\theta(y = 1 \mid X) = h_\theta(X) = \sigma(\theta^T X)$$

$$P_\theta(y = 0 \mid X) = 1 - h_\theta(X) = 1 - \sigma(\theta^T X)$$

log_likelihood: $\sum_{i=1}^{m}$ ??? $y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))$

$P_\theta(y^{(i)} | x^{(i)}) = (h_\theta(x^{(i)}))^{y^{(i)}} + (1-h_\theta(x^{(i)}))^{(1-y^{(i)})}$

$\log P_\theta(y^{(i)} | x^{(i)}) = \boxed{y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))}$

$i = 1 \ldots m$

$P(\underline{y} | \underline{\underline{X}}) = \prod_{i=1}^{m} P_\theta(y^{(i)} | x^{(i)})$

$\max \log P(\underline{y} | \underline{\underline{X}}) = \max \sum_{i=1}^{m} \log P_\theta(y^{(i)} | x^{(i)})$

## Multiclass classification

Email foldering/tagging: Work, Friends, Family, Hobby

Medical diagrams: Not ill, Cold, Flu

Weather: Sunny, Cloudy, Rain, Snow

$$y \in \{0, 1, 2, 3\}$$

$x_{test} \longrightarrow 2 \text{ س'ط}$
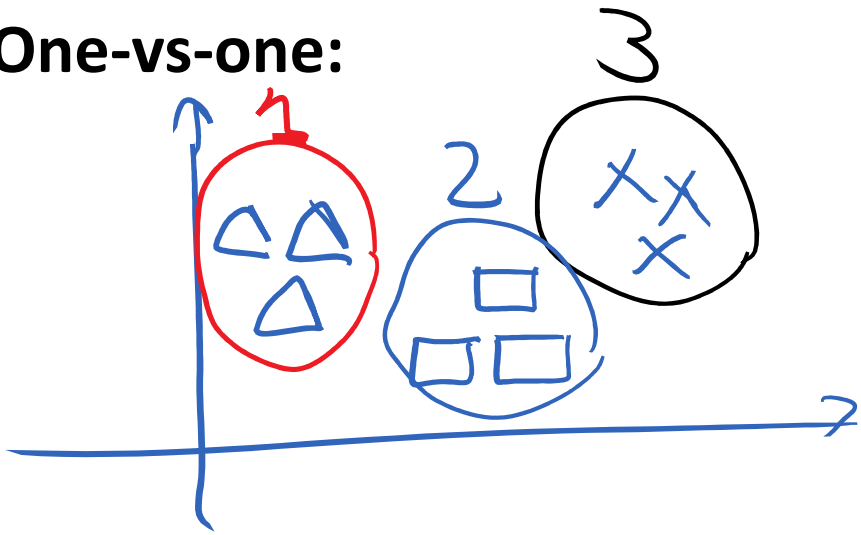
Binary classification:
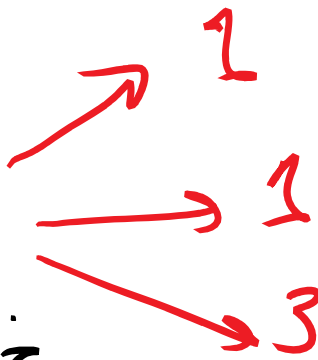
Multi-class classification:

**One-vs-one:**



$3$

$1$  $2$

$k = 3$

$1 - 2$

$2 - 3$

$1 - 3$

$$\binom{k}{2}$$

عداد المصنفات

$$x_{test} \longrightarrow \binom{k}{2}_{خروج}$$

$$\nearrow 1$$
$$\rightarrow 1 \implies y_{test} = 1$$
$$\searrow 3$$

**One-vs-all:**



Class 1: △
Class 2: □
Class 3: ✕

$h_\theta^{(i)}(x) = P(y = i|x;\theta) \qquad (i = 1, 2, 3)$

Softmax regression:

$(x^{(i)}, y^{(i)})$   $i = 1..m$

برای logistic regression اگر $k > 2$ (کلاس) سه تقسیم یافته

$y^{(i)} \in \{1, 2, ..., k\}$

MNIST $K = 10$

$p(y = k \mid \underline{x}) \implies h_\theta(X) = \begin{bmatrix} P(y=1 \mid x; \theta) \\ P(y=2 \mid x; \theta) \\ P(y=k \mid x; \theta) \end{bmatrix} =$

$h_\theta^{(j)}(X) = \dfrac{e^{\theta^{(j)T}X}}{\sum\limits_{i=1} e^{\theta^{(i)T}X}}$

$\sum\limits_{j=1}^{k} h_\theta^{(j)}(X)$

$\underline{\theta} = [\theta^{(1)}, ..., \theta^{(k)}]$ $(n+1) \times k$ : ابعاد سایز :

$\Theta^{(k)} \; ... \; \Theta^{(1)}$
$(n+1) \times 1$

## Cross entropy cost function:

$$\underline{1}\,(\text{صحیح تصادفی}) = 1 \qquad\qquad \underline{1}\,(\text{عبارت درست}) = 0$$

$$J(\Theta) = -\left[\sum_{i=1}^{m}\sum_{k=1}^{\boxed{K}} 1\left\{y^{(i)}=k\right\} \log \frac{\exp\left(\Theta^{(k)T}\underline{x}^{(i)}\right)}{\sum_{j=1}^{k}\exp\left(\Theta^{(j)T}\underline{x}^{(i)}\right)}\right]$$

تعمیم یافته تابع هزینه log Reg.

$$\underline{k=4}$$

$$X_{test} \longrightarrow h_\Theta(X_{test}) = \begin{bmatrix} 0.8 \\ 0.1 \\ \%5 \\ 0/05 \end{bmatrix} \longrightarrow Y_{test} = 1$$

$$\underline{P(y\,|\,x,\,\boxed{\Theta})} \sim \text{Categorical}\,(P_1, P_2, \dots, P_k) \qquad \boxed{P_k = 1 - \sum_{i=1}^{k-1} P_i}$$

# Regularization

# Regularization

A fundamental problem is that the algorithm tries to pick parameters that minimize loss on the training set, but this may not result in a model that has low loss on future data. This is called **overfitting**.

### example

suppose we want to predict the probability of heads when tossing a coin. We toss it $N = 3$ times and observe 3 heads. The MLE is $\theta_{\text{mle}} = N_1/(N_0 + N1) = 3/(3 + 0) = 1$. However, if we use $\text{Ber}(y \mid \theta_{\text{mle}})$ as our model, we will predict that all future coin tosses will also be heads, which seems rather unlikely.

# Regularization

- The core of the problem is that the model has enough parameters to perfectly fit the observed training data, so it can **perfectly match** the empirical distribution.

- However, in most cases the empirical distribution is not the same as the true distribution, so putting all the probability mass on the observed set of $N$ examples will not leave over any probability for novel data in the future. That is, **the model may not generalize**.

# Solution

The main solution to overfitting is to use regularization, which means to add a penalty term to the Cost function. Thus we optimize an objective of the form
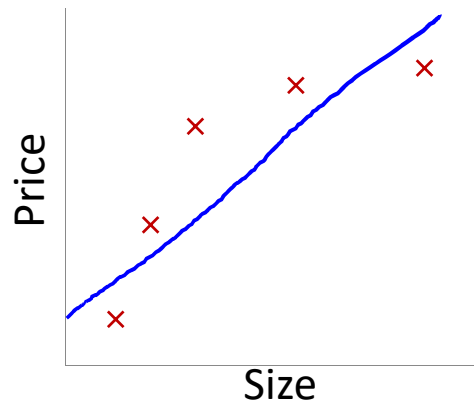
$$\mathcal{L}(\boldsymbol{\theta}; \lambda) = \left[ \frac{1}{M} \sum_{i=1}^{M} \ell(\mathbf{y}_i, \boldsymbol{\theta}; \mathbf{x}_i) \right] + \lambda C(\boldsymbol{\theta})$$

$\lambda \geq 0$ is a tuning parameter and control the relative impact of these two terms on the regression coefficient estimates.

When $\lambda = 0$, the penalty term has no effect
However, as $\lambda \to \infty$, the impact of the shrinkage penalty grows, and the ridge regression coefficient estimates will approach zero.

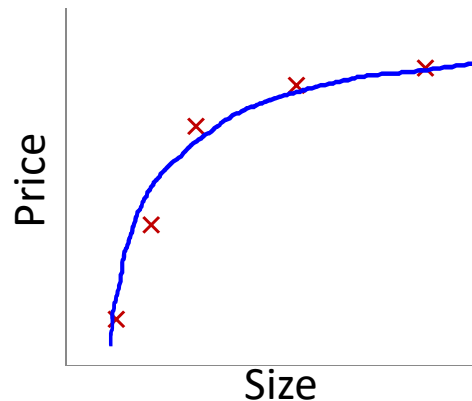# Example: Linear regression (housing prices)



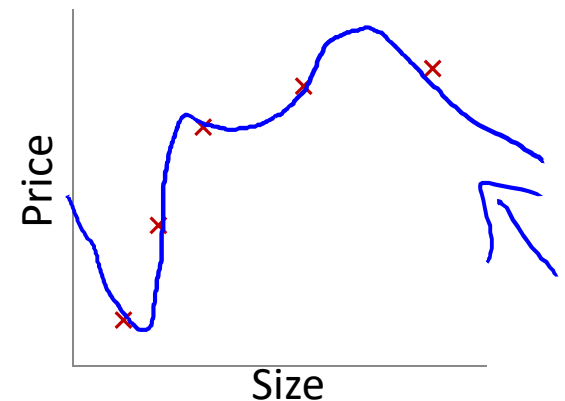$$\rightarrow \theta_0 + \theta_1 x$$

"Underfit"   "High bias"   $\lambda$

$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$$

"Just right"   $\lambda$

$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$
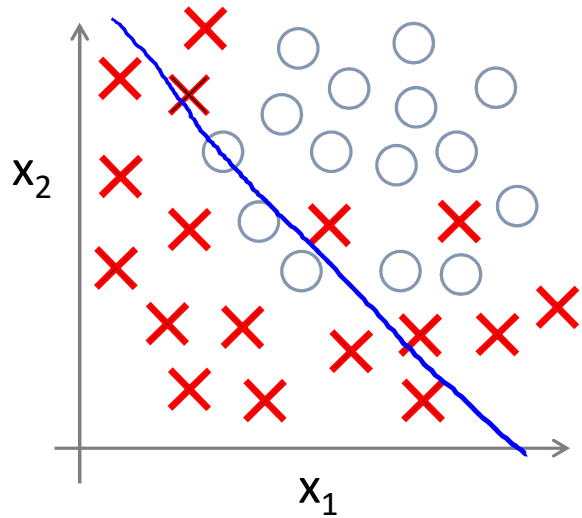
"Overfit"   "High variance"   $\lambda = 0$

**Overfitting:** If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples).
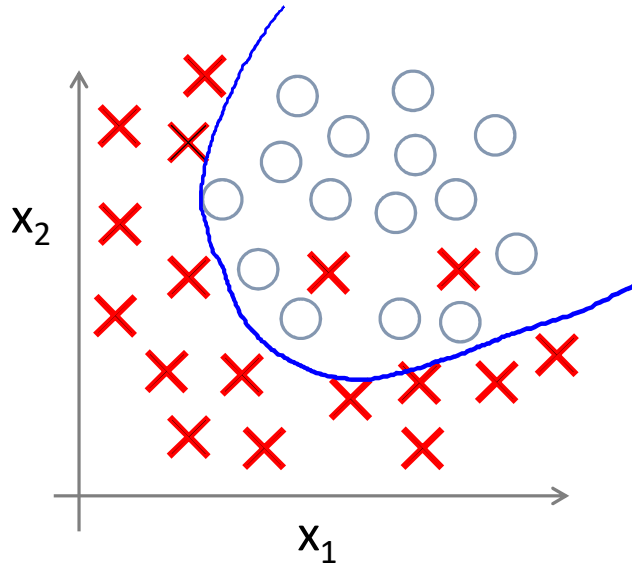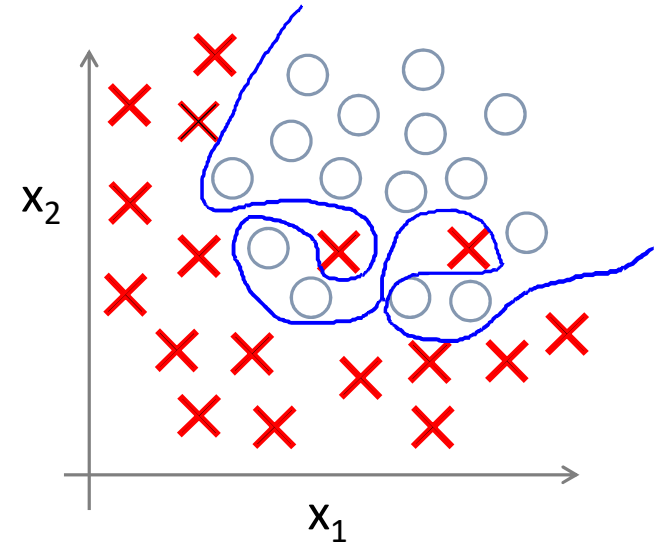
# Example: Logistic regression

(ضقدسنی)



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

( $g$ = sigmoid function)

"Underfit"

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2$
$+\theta_3 x_1^2 + \theta_4 x_2^2$
$+\theta_5 x_1 x_2)$

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$
$+\theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2$
$+\theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \ldots)$

"Overfit"

## Addressing overfitting:

$x_1 = $ size of house
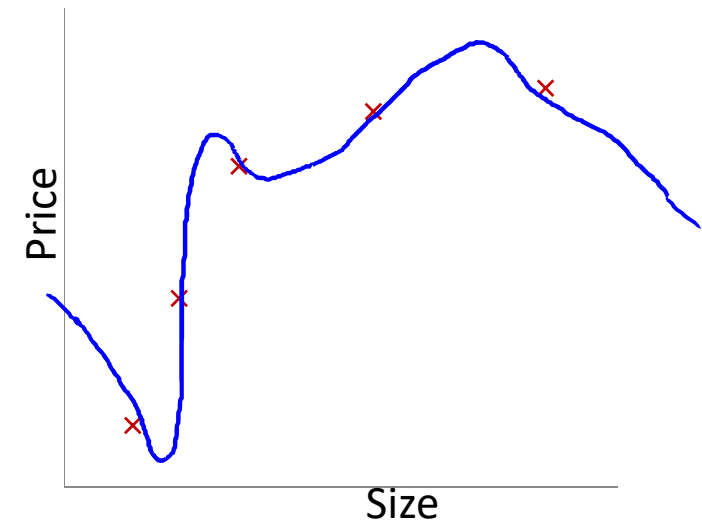$x_2 = $ no. of bedrooms
$x_3 = $ no. of floors
$x_4 = $ age of house
$x_5 = $ average income in neighborhood
$x_6 = $ kitchen size
$\vdots$
$x_{100}$

**Addressing overfitting:**

Options:
1. Reduce number of features
   — Manually select which features to keep.
   — Model selection algorithm (later in course).
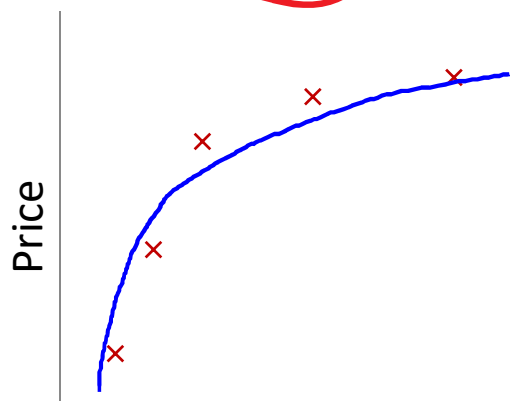2. Regularization
   — Keep all the features, but reduce magnitude/values of parameters $\theta_j$
   — Works well when we have a lot of features, each of which contributes a bit to predicting $y$.
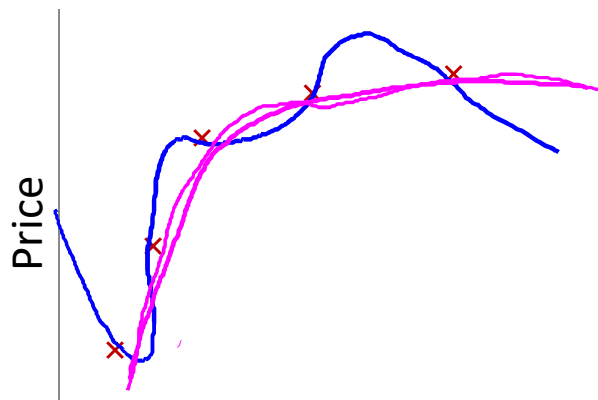
# Regularization

## Cost function

**Intuition**

Price — Size of house

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

Price — Size of house

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we penalize and make $\theta_3, \theta_4$ really small.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000\,\theta_3^2 + 1000\,\theta_4^2$$

$$\theta_3 \approx 0 \qquad \theta_4 \approx 0$$

**Regularization.**

$y = \theta_0$

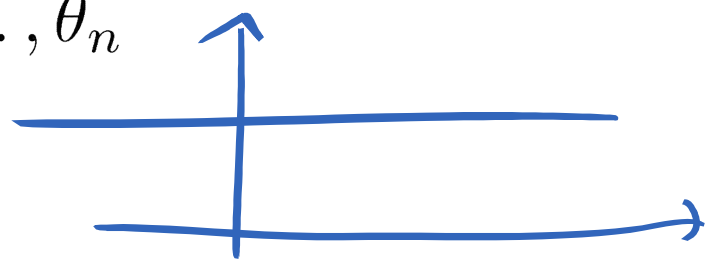$\theta_1 = \theta_2 = \ldots \theta_{100} = 0$

Small values for parameters $\theta_0, \theta_1, \ldots, \theta_n$
- — "Simpler" hypothesis
- — Less prone to overfitting

Housing:
- — Features: $x_1, x_2, \ldots, x_{100}$
- — Parameters: $\theta_0, \theta_1, \theta_2, \ldots, \theta_{100}$

$$J(\theta) = \frac{1}{2m}\left[\sum_{i=1}^{m}(h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{i=1}^{n} \theta_j^2\right]$$
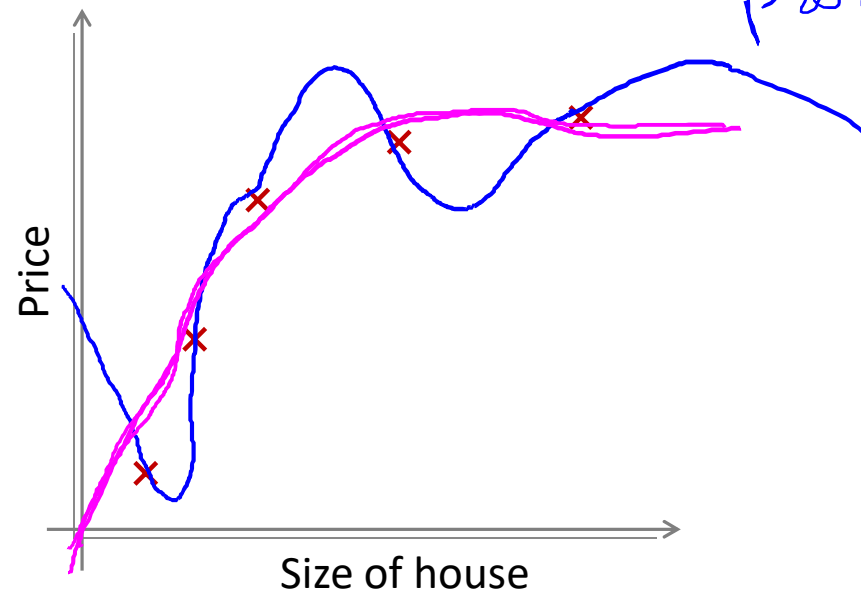
$\theta_1, \theta_2, \theta_3, \ldots, \theta_{100}$

# Regularization.

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

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

regularization parameter



Price

Size of house

$L1\text{-norm}: \ \boxed{\lambda} \cdot \sum_{j=1}^{n} |\theta_j|$

In regularized linear regression, we choose $\theta$ to minimize

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$
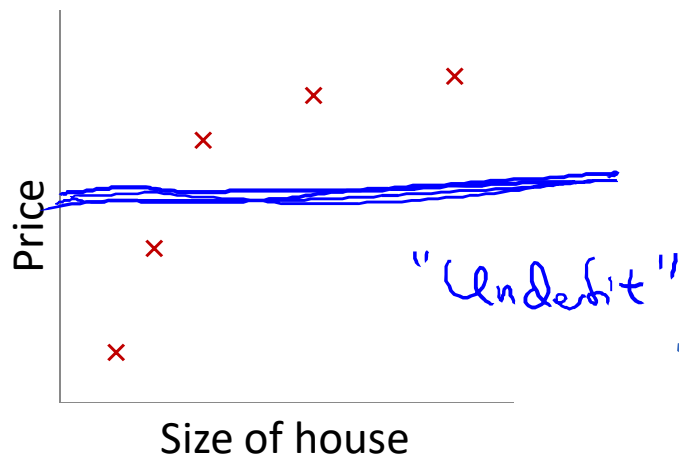
$L_2\text{-norm}$
regularization

What if $\lambda$ is set to an extremely large value (perhaps for too large for our problem, say $\lambda = 10^{10}$)?

- Algorithm works fine; setting $\lambda$ to be very large can't hurt it
- Algortihm fails to eliminate overfitting.
- Algorithm results in underfitting. (Fails to fit even training data well).
- Gradient descent will fail to converge.

In regularized linear regression, we choose $\theta$ to minimize

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

What if $\lambda$ is set to an extremely large value (perhaps for too large for our problem, say $\lambda = 10^{10}$)?



Price

"Underfit"

Size of house

$h_\theta(x)$

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

$\theta_1, \theta_2, \theta_3, \theta_4$

$\theta_1 \approx 0, \theta_2 \approx 0$

$\theta_3 \approx 0, \theta_4 \approx 0$

$h_\theta(x) = \theta_0$

# Regularization

Regularized linear regression

# Regularized linear regression

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

$$\min_\theta J(\theta)$$

**Gradient descent**

$$\boxed{l_2 - norm}$$

Repeat {

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

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)} \quad - \frac{\lambda}{m}\theta_j \right]$$

$$(j = \cancel{0}, 1, 2, 3, \ldots, n)$$

}

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$$

$$\boxed{1 - \alpha \frac{\lambda}{m} < 1}$$

$$0.99 \qquad \theta_j \times 0.99$$

# Normal equation(Regularized) $(L_2\text{-norm})$

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \quad\quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \mathbb{R}^m$$

$m \times (n+1)$

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

$$\frac{\partial}{\partial \theta_j} J(\theta) \overset{\text{set}}{=} 0 \quad \longrightarrow$$

$$\Theta = \left( \boxed{X^T X} + \lambda \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & \cdots & 1 \end{bmatrix} \right)^{-1} X^T y$$

$(n+1) \times (n+1)$

$E.g. \quad n = 2 \quad \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

## Non-invertibility

Suppose $m \leq n,$ $\leftarrow$

    (#examples)  (#features)

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

non-invertible / singular

If $\lambda > 0,$

$$\theta = \left( X^T X + \lambda \begin{bmatrix} 0 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix} \right)^{-1} X^T y$$

invertible

$m << n$

ویژگی‌های وابسته خطی