

Machine learning (Stanford) – Andrew Ng

기석, 2018. 4. 21 updated

<Class material>

강의 : <https://www.coursera.org/learn/machine-learning>

Setting up programming assignment environment

- set up : http://wiki.octave.org/Octave_for_Microsoft_Windows

<Contents MAP>

I. 기본

0. 선형대수학 복습 및 프로그래밍적 표기

1. Octave Coding

II. Machine learning

0. Machine learning Introduction - ML이 무엇이냐! 종류!

(1) Supervised learning :

- regression problem(회귀 문제)

- classification problem(분류 문제) : 6. Classification → logistic regression

(2) Unsupervised learning : Cocktail party algorithm

가장 많이 쓰이는 3가지를 꼽자면, (1) logistic regression (2) neural networks (3) SVM

1. Supervised learning :

기본

1.1. Supervised learning 이해

Solution (구하는 방법)

1.2. Gradient descent

1.3. Normal equation

종류

1.4. Linear regression → regression problem(회귀 문제) : 변수값이 continuous한 경우

1.5. logistic regression → classification problem(분류 문제) : 변수값이 discrete한 경우

1.6. Neural Networks

1.7. Supported vector machine(SVM) - An alternative for logistic regression

2. Unsupervised learning ← unlabeled data set ex) Cocktail party algorithm

2.1. Clustering algorithm : K-means algorithm (the most popular)

2.2. Dimensionality reduction

2.3. (중요) Principle component analysis (PCA)

2.3.1. 문제정의

2.3.2. 해결 PCA 알고리즘

2.3.3. Application

2.3.4. 주의사항

2.4. Anomaly detection

III. Special topics of Machine learning

3. Anomaly detection ← unsupervised로 분류되기도 하나 labeled data도 쓰인다.

4. Recommender systems

IV. Application details of Machine Learning

6. What ML should you apply / diagnostic

7. Machine Learning system design :

 7.1. error analysis

 7.2. large scale machine learning (week 10)

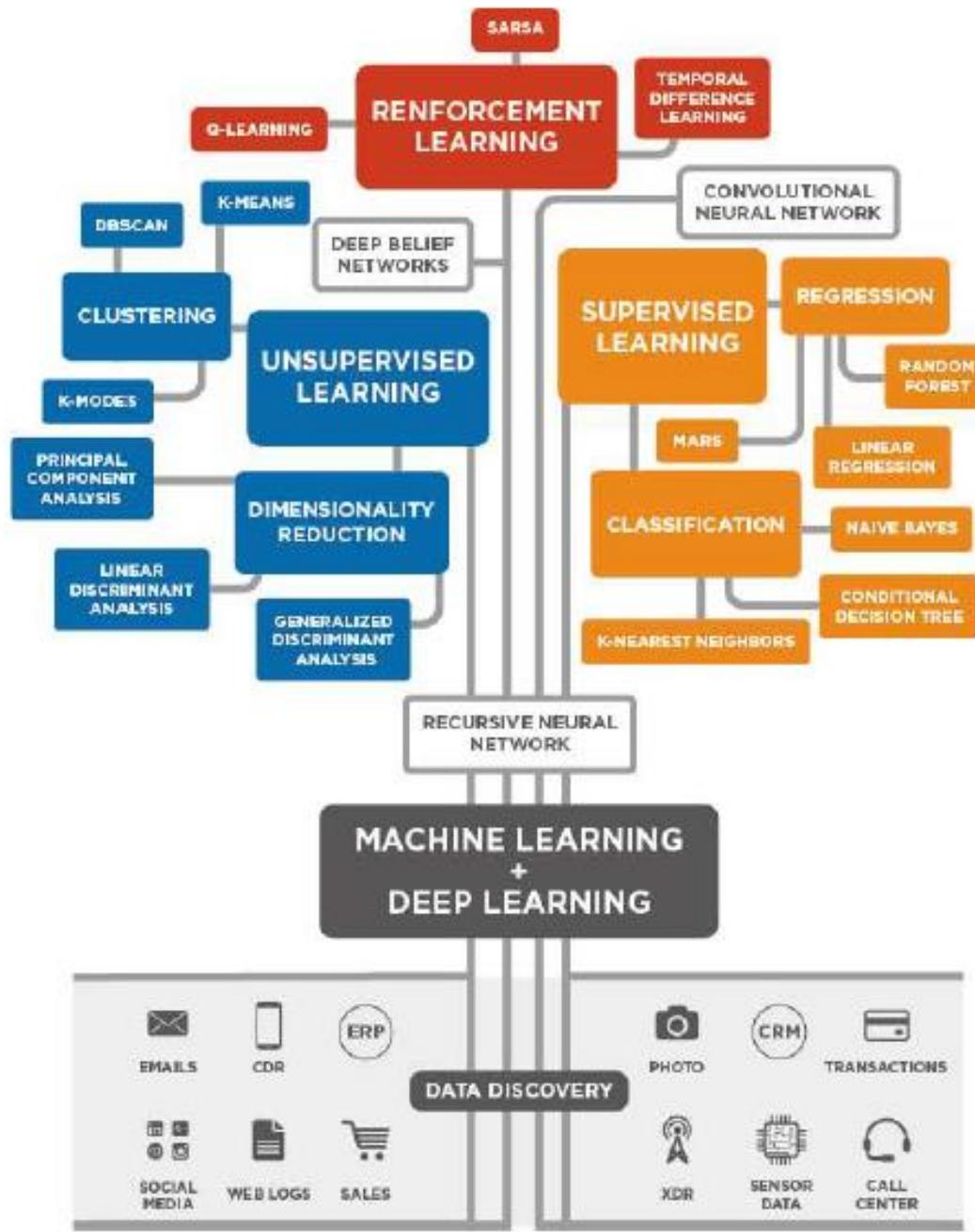
 7.3. Generating large data (artificial data synthesis)

 7.4. Online learning

8. Machine learning pipelines : Photo OCR problem (week 11)

Summary: Main topics

- Supervised Learning $(x^{(i)}, y^{(i)})$
 - Linear regression, logistic regression, neural networks, SVMs
- Unsupervised Learning $x^{(i)}$
 - K-means, PCA, Anomaly detection
- Special applications/special topics
 - Recommender systems, large scale machine learning.
- Advice on building a machine learning system
 - Bias/variance, regularization; deciding what to work on next: evaluation of learning algorithms, learning curves, error analysis, ceiling analysis.



0. 선형대수학 복습

0.1. Matrix and vector 복습

- Matrix :

Matrix Elements (entries of matrix)

$$A = \begin{bmatrix} 1402 & 191 \\ 1371 & 821 \\ 949 & 1437 \\ 147 & 1448 \end{bmatrix}$$

A_{ij} = "i, j entry" in the i^{th} row, j^{th} column.

Matrices are 2-dimensional arrays:

```
 $$\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \\ j & k & l \end{bmatrix}$$
```

The above matrix has four rows and three columns, so it is a 4×3 matrix.

- $\begin{bmatrix} 1 & 2 \\ 4 & 0 \\ 0 & 1 \end{bmatrix}$ is a 3×2 matrix.
- $\begin{bmatrix} 0 & 1 & 4 & 2 \\ 3 & 4 & 0 & 9 \end{bmatrix}$ is a 4×2 matrix.
- $\begin{bmatrix} 0 & 4 & 2 \\ 3 & 4 & 9 \\ 5 & -1 & 0 \end{bmatrix}$ is a 3×3 matrix.
- $\begin{bmatrix} 1 & 2 \end{bmatrix}$ is a 1×2 matrix.

Let A be a matrix shown below. A_{32} is one of the elements of this matrix.

$A = \begin{bmatrix} 85 & 76 & 66 & 5 \\ 94 & 75 & 18 & 28 \\ 68 & 40 & 71 & 5 \end{bmatrix}$

What is the value of A_{32} ?

→ 40이다.

- Vector : vector is a special case of matrix

A vector is a matrix with one column and many rows:

So vectors are a subset of matrices. The above vector is a 4×1 matrix.

Vector: An $n \times 1$ matrix.

$$\underline{y} = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \quad \begin{matrix} \uparrow & \uparrow \\ n=4 \end{matrix} \quad \leftarrow \text{4-dimensional vector.}$$

$y_i = \underline{i^{th} \text{ element}}$

$$y_1 = 460$$

$$y_2 = 232$$

$$y_3 = 35$$

1-indexed vs 0-indexed:

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \leftarrow \quad y = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{bmatrix} \leftarrow$$

1-indexed 0-indexed

- computer notation :

Notation and terms:

- $\$A_{ij}$ refers to the element in the i th row and j th column of matrix A .
 - A vector with ' n ' rows is referred to as an ' n '-dimensional vector.
 - $\$v_i$ refers to the element in the i th row of the vector.
 - In general, all our vectors and matrices will be 1-indexed. Note that for some programming languages, the arrays are 0-indexed.
 - Matrices are usually denoted by uppercase names while vectors are lowercase.
 - "Scalar" means that an object is a single value, not a vector or matrix.
 - \mathbb{R} refers to the set of scalar real numbers.
 - \mathbb{R}^n refers to the set of n -dimensional vectors of real numbers.

Run the cell below to get familiar with the commands in Octave/Matlab. Feel free to create matrices and vectors and try out different things.

```
1 % The ; denotes we are going back to a new row.
2 A = [1, 2, 3; 4, 5, 6; 7, 8, 9; 10, 11, 12]
3
4 % Initialize a vector
5 v = [1;2;3]
6
7 % Get the dimension of the matrix A where m = rows and n = columns
8 [m,n] = size(A)
9
10 % You could also store it this way
11 dim_A = size(A)
12
13 % Get the dimension of the vector v
14 dim_v = size(v)
15
16 % Now let's index into the 2nd row 3rd column of matrix A
17 A_23 = A(2,3)
18
```

```
A =
```

```
1 2 3  
4 5 6  
7 8 9  
10 11 12
```

```
v =
```

```
1  
2  
3
```

```
m = 4  
n = 3  
dim_A =
```

```
4 3
```

```
dim_v =
```

```
3 1
```

```
A_23 = 6
```

0.2. Matrix addition & scalar multiplication

→ matrix dimension이 같아야 한다.

```
1 % Initialize matrix A and B  
2 A = [1, 2, 4; 5, 3, 2]  
3 B = [1, 3, 4; 1, 1, 1]  
4  
5 % Initialize constant s  
6 s = 2  
7  
8 % See how element-wise addition works  
9 add_AB = A + B  
10  
11 % See how element-wise subtraction works  
12 sub_AB = A - B  
13  
14 % See how scalar multiplication works  
15 mult_As = A * s  
16  
17 % Divide A by s  
18 div_As = A / s  
19  
20 % What happens if we have a Matrix + scalar?  
21 add_As = A + s  
22
```

Run

Reset

```

A =
1 2 4
5 3 2

B =
1 3 4
1 1 1

s = 2
add_AB =
2 5 8
6 4 3

sub_AB =
0 -1 0
4 2 1

mult_As =
2 4 8
10 6 4

div_As =
0.50000 1.00000 2.00000
2.50000 1.50000 1.00000

add_As =
3 4 6
7 5 4

```

0.3. Matrix vector multiplication

$$\begin{array}{c}
\underline{A} \quad \times \quad \underline{x} \quad = \quad \underline{y} \\
\left[\begin{array}{cccc|c} & & & & \\ & & & & \\ & & & & \\ \hline 1 & \uparrow & \uparrow & \dots & \uparrow \\ m \times n & \text{matrix} & n \times 1 & \text{matrix} & \text{m-dimensional} \\ & (\text{m rows}, & (\text{n-dimensional} & \text{vector}) & \text{vector} \\ & n \text{ columns}) & \text{vector}) & & \end{array} \right] \times \left[\begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \right] = \left[\begin{array}{c} y_1 \\ y_2 \\ \vdots \\ y_m \end{array} \right]
\end{array}$$

To get y_i , multiply A 's i^{th} row with elements of vector x , and add them up.

- unit 2에서 한 거에 적용

House sizes:

- 2104
- 1416
- 1534
- 852

Matrix $\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}$

4×2

$h_{\theta}(x) = -40 + 0.25x$

$h_{\theta}(x)$

2×1 Vector $\begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$

4×1 matrix $\begin{bmatrix} -40 \times 1 + 0.25 \times 2104 \\ -40 \times 1 + 0.25 \times 1416 \\ \vdots \\ -40 \times 1 + 0.25 \times 852 \end{bmatrix}$

$h_{\theta}(2104)$

$h_{\theta}(1416)$

$\boxed{\text{prediction}} = \boxed{\text{Data Matrix}} \times \boxed{\text{Parameters}}$

4×1

n

for $i = 1:1000$,
prediction(i) = ...

→ 1000개의 데이터가 있으면 for loop(오른쪽)보다는 matrix로 한 줄(왼쪽)로 처리하는 것이 좋다. 이 경우에는 h 를 matrix-vector multiplication으로 처리했다.

0.4. Matrix-matrix multiplication

- 이것은 matrix-vector multiplication의 확장판이다.

$$A \times B = C$$

A : $m \times n$ matrix (m rows, n columns)

B : $n \times o$ matrix (n rows, o columns)

C : $m \times o$ matrix

The i^{th} column of the matrix C is obtained by multiplying A with the i^{th} column of B . (for $i = 1, 2, \dots, o$)

- 적용 : 여러 hypothesis를 동시에 할 수 있다.

House sizes:

$$\left\{ \begin{array}{l} 2104 \\ 1416 \\ 1534 \\ \hline 852 \end{array} \right.$$

Have 3 competing hypotheses:

$$1. h_{\theta}(x) = -40 + 0.25x$$

$$2. h_{\theta}(x) = 200 + 0.1x$$

$$3. h_{\theta}(x) = -150 + 0.4x$$

Matrix

$$\begin{bmatrix} 1 & 2104 \\ 1 & \underline{1416} \\ 1 & \underline{1534} \\ 1 & \underline{852} \end{bmatrix}$$

Matrix

$$\begin{array}{|c|c|} \hline -40 & 200 \\ \hline 0.25 & 0.1 \\ \hline \end{array}$$

-150
0.4

486	410	692
314	342	416
344	353	464
173	285	191

Prediction
of first
he

Predictions of 2nd he

- 특징(Properties)

(1) not commutative : $A * B \neq B * A \rightarrow$ 교환법칙 X

(2) Associative : $A * (B * C) = (A * B) * C \rightarrow$ 결합법칙

(3) Identity matrix

For any matrix A ,

$$A \cdot I = I \cdot A = A$$

$I_{n \times n}$

$m \times n$ $n \times n$ $m \times m$ $m \times n$ $m \times n$

```
1 % Initialize random matrices A and B
2 A = [1,2;4,5]
3 B = [1,1;0,2]
4
5 % Initialize a 2 by 2 identity matrix
6 I = eye(2)
7
8 % The above notation is the same as I = [1,0;0,1]
9
10 % What happens when we multiply I*A ?
11 IA = I*A
12
13 % How about A*I ?
14 AI = A*I
15
16 % Compute A*B
17 AB = A*B
18
19 % Is it equal to B*A?
20 BA = B*A
21
22 % Note that IA = AI but AB != BA
```

(4) Inverse matrix

Not all numbers have an inverse.

$\text{O}(\underline{0^{-1}})$ undefined

Matrix inverse: $\xrightarrow{\text{square matrix}} \xrightarrow{\text{#rows = #columns}} A^{-1}$

If A is an $m \times m$ matrix, and if it has an inverse,

$$\rightarrow \underline{A}(\underline{A^{-1}}) = \underline{A^{-1}} \underline{A} = \underline{I}.$$

Matrices that don't have an inverse are "singular" or "degenerate"

- 코딩

$\text{pinv}(A)$ 로 하면 바로 된다.

$\$A^{-1}\$$ 도 된다.

(5) Transpose

Matrix Transpose

Example:

$$\underline{A} = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 5 & 9 \end{bmatrix} \quad 2 \times 3$$

$$\underline{B} = \underline{A^T} = \begin{bmatrix} 1 & 3 \\ 2 & 5 \\ 0 & 9 \end{bmatrix} \quad 3 \times 2$$

Let A be an $m \times n$ matrix, and let $B = A^T$.

Then B is an $n \times m$ matrix, and

$$B_{ij} = A_{ji}.$$

$$B_{12} = A_{21} = 2$$

$$B_{32} = 9 \quad A_{23}$$

What is $\$ \begin{bmatrix} 0 & 3 \\ 1 & 4 \end{bmatrix}^T \$$?

$\rightarrow {}^T$ 이거로 하네!!

$\$ \$A = \begin{bmatrix} a & b \\ c & d \\ e & f \end{bmatrix} \$ \$$

$\$ \$A^T = \begin{bmatrix} a & c & e \\ b & d & f \end{bmatrix} \$ \$$

1. Octave Coding

- Prototyping language이다!!!

1.1. Basic operations

1.1.1. 사칙연산

5+6

3-2

1/2

2^6 ← python이랑 다른 점

1 == 2

1 ~= 2 ← python이랑 다른 점 !=

&& ← and

|| ← or

xor(1, 0) ← or

PS1('>>');

The default value of the primary prompt string is "octave:> ". To change it, use a command like

PS1 ("WWWu@WWW> ")

1.1.2. assigning

a = 3; ← semicolon은 suppress output

>> c = (3>=1);

>> c

c = 1 ← 신기하다. True라고 나오는 것임.

1.1.3. Displaying

(1) >> a

해도 되고

(2) >> disp (a) ; ← display해도 되구

(3) C 스타일

>> disp (sprintf('2 decimals : %0.2f', a)) ← string print, 0.2f는 무슨 의미인지 알겠지 ㅋㅋㅋ python이랑 똑같잖아.

>> disp (sprintf('6 decimals : %0.6f', a))

1.1.4. Matrix

(1) Matrix

>> A = [1 2; 3 4; 5 6]

A =

1 2

3 4

5 6

(1-1) indexing

```
>> A (3,2)
```

```
>> A(3,2)
```

```
ans = 6
```

- 2열을 다 불러오려면 → : 는 전부를 의미한다.

```
>> A( :, 2)
```

```
ans =
```

```
2
```

```
4
```

```
6
```

- 2행을 다 불러오려면

```
>> A(2,:)
```

```
ans =
```

```
3 4
```

- 1행이랑 3행의 애들을 전부 가져와라

```
>> A ( [1,3] , : )
```

```
ans =
```

```
1 2
```

```
5 6
```

- 만약에 행렬의 모든 원소를 열벡터로 만들고 싶다.

```
A ( : )
```

```
ans =
```

```
1
```

```
3
```

```
5
```

```
200
```

```
10
```

```
11
```

```
12
```

```
201
```

```
202
```

```
101
```

(1-2) Assigning

- 그 자리를 다른 숫자로 **assign**하고 싶다.

$A(:, 2) = [10 ; 11 ; 12]$ \leftarrow 2열 전체가 10, 11, 12로 바뀌게 된다.

$A =$

```
1   10  
3   11  
5   12
```

(1-3) append하고 싶다

$A = [A, [100 ; 101 ; 102]]$ \rightarrow column을 옆에다가 붙인거다!!! ([]로 표시하면 된다.)

\Rightarrow , (쉼표) 랑 ''(띄어쓰기)랑 같은 것이다. 나란히 붙이기!!!

$A =$

```
1   10   100  
3   11   101  
5   12   102
```

- 만약에 row를 append하고 싶으면

```
>> A = [A ; [200 201 202]]
```

$A =$

```
1   10   100  
3   11   101  
5   12   102  
200  201  202
```

(1-4) 두 개의 matrix를 붙이고 싶다.

```
>> A = [1 2; 3 4; 5 6];
```

```
>> B = [11 12; 13 14; 15 16]
```

```
>> C = [A B]
```

$C =$

```
1   2   11   12  
3   4   13   14  
5   6   15   16
```

```
>> C = [A ; B]
```

$C =$

```
1   2
```

```
3   4  
5   6  
11  12  
13  14  
15  16
```

(2) 행렬 만드는 다른 방식들 (**ones**, **zeros**, **eye**, **rand**, **randn**, **magic**)

```
>> v = 1 : 0.1 : 2    ← 1부터 0.1 step으로 2까지 (2포함) : range랑 비슷함
```

→ 열벡터를 만들어준다. 1 1.1 1.2 ... 1.9 2.0

```
>> v = 1:6
```

```
v =
```

```
1   2   3   4   5   6
```

```
>> ones(2,3) → 1로 된 2x3의 matrix를 만들어준다.
```

```
ans =
```

```
1   1   1  
1   1   1
```

```
>> 2*ones(2,3)
```

```
ans =
```

```
2   2   2  
2   2   2
```

```
>> zeros(2,3)
```

```
ans =
```

```
0   0   0  
0   0   0
```

```
>> eye (3,3)
```

```
ans =
```

Diagonal Matrix

```
1   0   0  
0   1   0  
0   0   1
```

```
>> eye(3)
```

```
ans =
```

Diagonal Matrix

```
1 0 0  
0 1 0  
0 0 1
```

```
>> rand(3,3) ← 0부터 1까지 랜덤한 숫자 3x3로!!
```

```
ans =
```

```
0.95038 0.69605 0.61569  
0.70962 0.30824 0.75617  
0.89219 0.58523 0.13924
```

```
>> randn(2,3) ← 가우씨안 distribution에서 2x3 랜덤 matrix (mean=0, std=1)
```

```
ans =
```

```
-0.91661 0.16785 0.72714  
0.57808 0.47086 0.80971
```

```
>> A = magic (3) → 3 by 3 matrix인데 가로 세로 합이 일정한 행렬을 만들어준다.
```

```
A =
```

```
8 1 6  
3 5 7  
4 9 2
```

cf) **help**

```
help eye
```

```
help rand
```

```
help help
```

1.2. How to move data around

```
load
```

```
save the result
```

(1) Matrix size

(1-1) size 함수

```
A = [1 2 ; 3 4; 5 6]
```

```
size(A)
```

```
sz = size(A)
>> sz = size (A)
sz =
```

3 2

→ sz도 함수가 되었다.

size(A, 1) → size 함수의 결과값의 1번째 index 가져와라 → row의 개수
>> size(A, 1)
ans = 3
size(A, 2) → size 함수 결과의 2번째 index 가져와라 → A 행렬의 column의 개수

(1-2) length 함수 → 가장 긴 dimension을 가져다 준다.

```
>> v=[1 2 3 4]
v =
```

1 2 3 4

```
>> length(v)
ans = 4
```

(2) 데이터 저장 및 불러오기

(2-1) pwd

어디에 있는지

(2-2) cd

change directory

(2-3) ls ← 리눅스랑 같네

그 directory에 있는 거 다 나오게

(2-4) data loading

load featuresX.dat

아니면

```
load ('featuresX.dat')
```

다음에

```
>> featuresX
>> size(featuresX)
```

없애려면

```
>> clear featuresX
>> clear → 모든 변수를 다 없앤다.
```

앞의 10개의 data만 변수로 가져오려면

```
v = featuresX(1:10)
```

(2-5) who

who 하면 current data에서 어떤 변수들이 정의됐는지 볼 수 있다.

```
>> who
```

Variables in the current scope:

```
A ans sz v w
```

```
>> whos ← 대박!!!! 어떤 data인지, matrix인지
```

(2-6) 저장 방법

```
>> save hello.mat v;
```

```
>> save hello.txt v -ascii % save as text(ASCII)
```

1.3. Matrix computing

(1) Matrix multiplication

```
>> A * C
```

```
ans =
```

```
5 5  
11 11  
17 17
```

(2) 각 원소별로 연산하기

>> A . * B → 각 원소끼리 곱하는 것이다. 엄청 특이하네. A의 1,1을 B의 1,1과 곱하고, A의 1,2를 B의 1,2와 곱하는 것이다.

```
>> A . ^2 → element wise squaring
```

>> 역수를 취하고 싶을 때

```
v = [ 1 ; 2 ; 3 ]
```

```
>> 1./v ← 1을 v의 각 원소 별로 나눠라.
```

```
ans =
```

```
1.00000
```

0.50000

0.33333

>> log(A) ← 각 원소를!!

ans =

0.00000 0.69315

1.09861 1.38629

1.60944 1.79176

>> exp(v) ← e를 base로!!

ans =

2.7183

7.3891

20.0855

>> abs(v) ← 다 절대값으로

>> -v

ans =

-1

-2

-3

>> 행렬 각 원소에 1을 더하고 싶을때

>> v + ones(length(v), 1) ← v만큼의 열을 만들어서 1을 더해야 하니까 length!!!!

ans =

2

3

4

>> v + 1 ← 같은 거임!!!

(3) Transpose, flip, inverse

(3-1) Transpose

>> A'

ans =

1 3 5

2 4 6

(3-2) flip up down

```
>> flipud(eye(9)) ← 바꿀 수 있음  
ans =
```

Permutation Matrix

```
0 0 1  
0 1 0  
1 0 0
```

(3-3) Inverse

- inverse하고 싶으면
pinv(A) ← pseudo inverse에 해당한다.

(4) 행렬 내의 값 분석

```
a = [1 15 2 0.5]
```

(4-1) MAX 찾기 : 기본이 column단위

```
>> val = max(a) ← 최대값
```

```
val = 15
```

- index까지 알고 싶으면

```
>> [val, ind] = max(a)  
val = 15  
ind = 2
```

>> max (A)하면 → column wise maximum을 한다.

```
>> max (rand(3))  
ans =
```

```
0.60958 0.93724 0.70481
```

→ 열 단위로!!

>> max (rand(3), rand(3)) → rand(3)와 rand(3) 두 개의 같은 index의 원소를 비교하여 3x3의 matrix를 보여줌
ans =

```
0.86301 0.80544 0.97582  
0.90647 0.99723 0.33394  
0.78273 0.23280 0.70618
```

>> max (A) → column wise maximum

```
>> max (A, [ ], 1) → column wise maximum
```

```
>> max (A, [], 2) → per row maximum
```

- 행렬 전체에서 가장 큰 원소를 찾고 싶으면

```
>> max(max(A))
```

아니면

```
>> max ( A(:)) ← vector로 펼쳐놓은 다음에!!
```

(4-2) 부등호 비교

- 값 비교를 하고 싶으면

```
>> a < 3
```

```
>> a
```

```
a =
```

```
1.00000 15.00000 2.00000 0.50000
```

```
>> a < 3
```

```
ans =
```

```
1 0 1 1 ← 2번째 원소는 아니구나
```

```
>> find (a < 3) ← 3보다 작은 원소의 index를 알려준다.
```

```
>> find(a<3)
```

```
ans =
```

```
1 3 4
```

```
>> [r, c] = find ( A >= 7) → 이렇게 해야 원소 index를 알 수 있음. 그냥 vector처럼 1자리면 안됨
```

```
r =
```

```
1
```

```
3
```

```
2
```

```
c =
```

```
1
```

```
2
```

```
3
```

(4-3) 합 : 기본이 column별

```
>> sum (a)
```

```
ans = 18.500
```

- column 별로 sum하려면

```
>> sum(A, 1)
```

- row 별로 sum하려면

```
>> sum(A, 2)
```

- 만약에 행렬의 대각선을 곱하고 싶으면

```
A = magic(9)
```

```
A .* eye(9) ← 원소별로 곱하니까 E 대각행렬의 1과 곱해진 부분만 남게 된다.
```

```
ans =
```

47	0	0	0	0	0	0	0	0
0	68	0	0	0	0	0	0	0
0	0	8	0	0	0	0	0	0
0	0	0	20	0	0	0	0	0
0	0	0	0	41	0	0	0	0
0	0	0	0	0	62	0	0	0
0	0	0	0	0	0	74	0	0
0	0	0	0	0	0	0	14	0
0	0	0	0	0	0	0	0	35

다 더하고 싶은 것인니까

```
>> sum(sum(A.*eye(9)))
```

(4-4) 곱

```
>> prod (a) → 다 곱해!!
```

(4-5) 올림 내림

```
>> ceil (a) → 올림
```

```
>> floor (a) → 내림
```

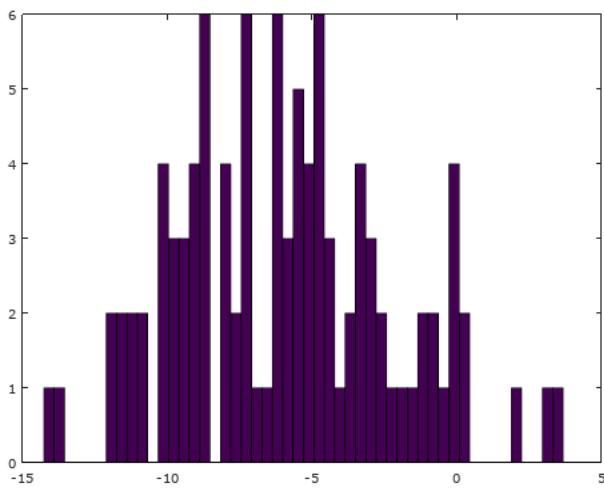
1.4. Plotting

(0) 히스토그램 만들기

```
>> w = -6 + sqrt(10)*randn(1,100);
```

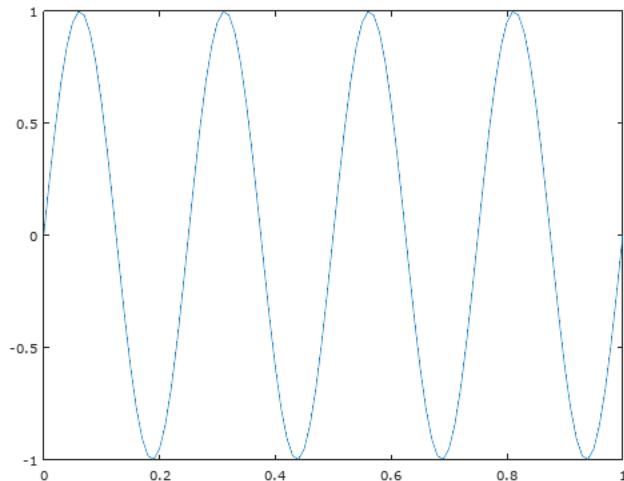
```
>> hist(w)
```

```
>> hist(w,50) ← 50 구간을 만들자.
```



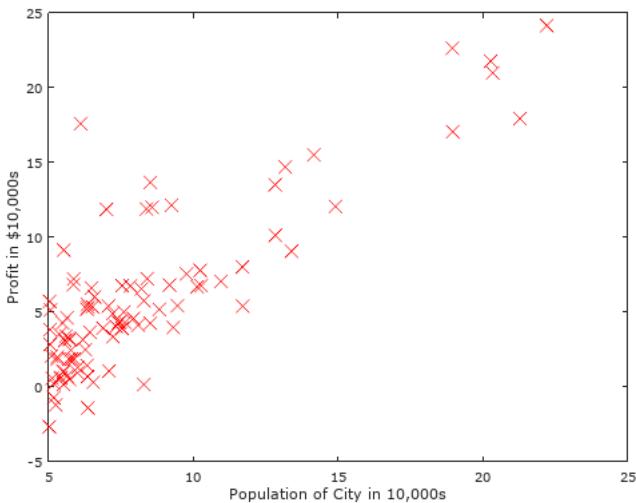
(1) Plot 함수 : (x 축, y 축)으로 하면 된다.

```
>> t = [0 : 0.01 : 1];
>> y1 = sin(2*pi*4*t);
>> plot(t, y1)
```

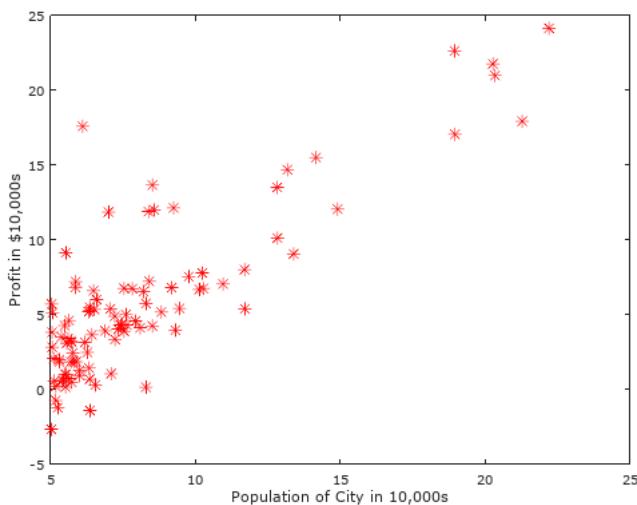


```
>> plot(x, y, 'rx', 'MarkerSize', 10);
```

← rx라는 것은 red x 표시를 하겠다는 뜻이다. MarkerSize는 x의 크기를 의미한다.

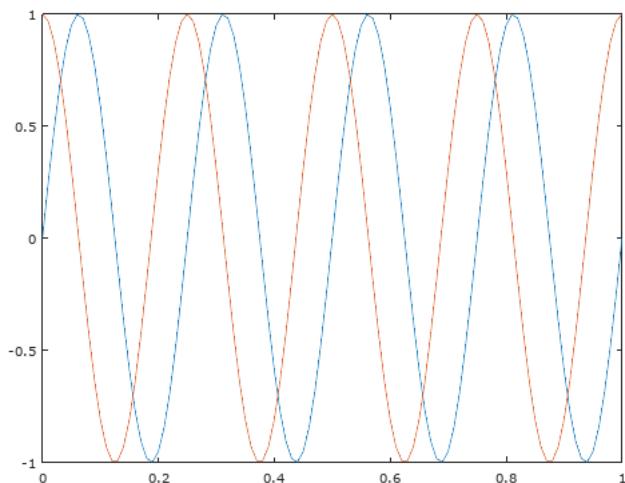


cf) r^* 를 하면 이렇게 된다.



- 위에 고정시키는 방법

```
>> hold on;
>> y2 = cos(2*pi*4*t);
>> plot (t,y2)
```



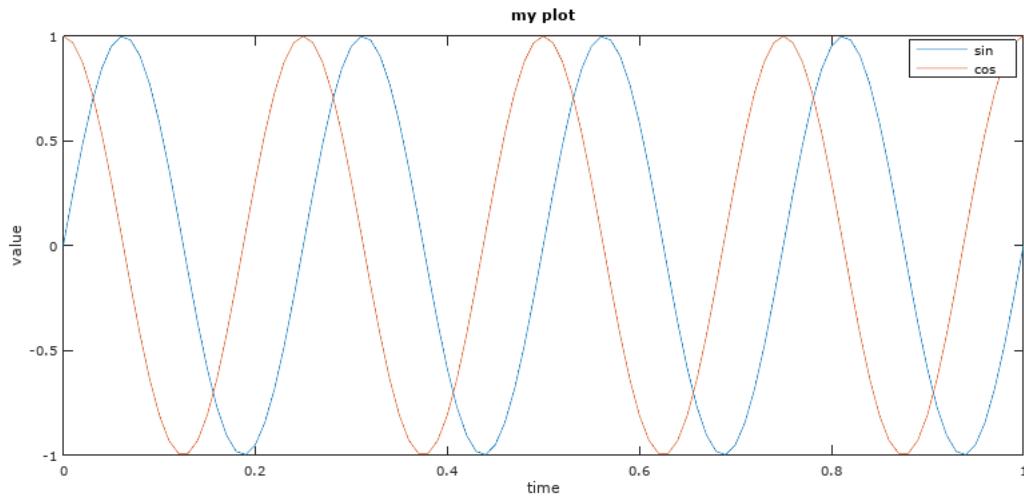
- plot하는 방법 심화 (2주차 octave 과제)

(2) labeling!!!

```
>> xlabel('time')
>> ylabel('value')
>> legend ('sin', 'cos')
>> title ('my plot')
- range 바꾸기
```

>> axis ([0.5 1 -1 1]) ➔ 처음 2개는 x축이 어디서부터 어디까지, 뒤에 두 개는 y축이 어디서부터 어디까지

>> clf ; ➔ clear figure



(3) 저장하려면

```
>> cd 'C:\Users\owner\Desktop'; print -dpng 'myPlot.png'
```

```
>> close
```

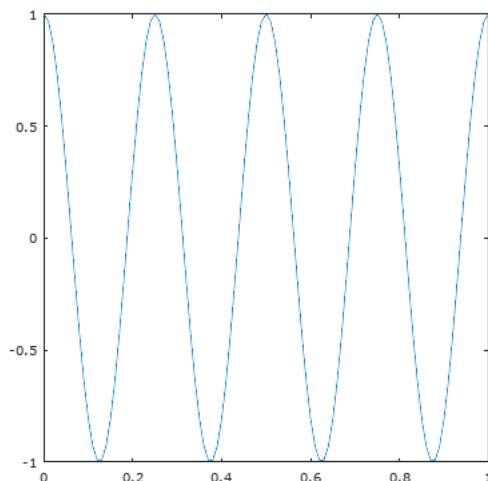
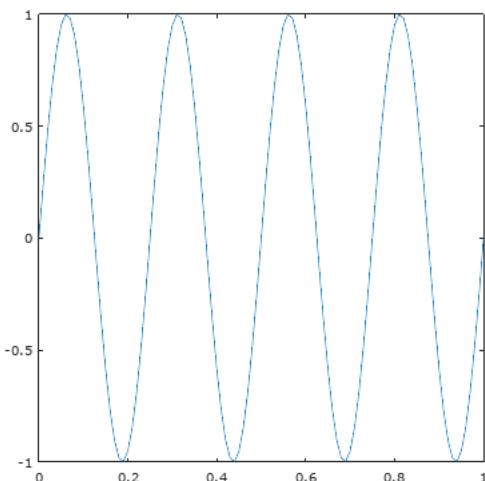
figure 없어짐

(4) 따로따로 띄우려면

```
>> figure (1); plot (t, y1);
>> figure (2); plot (t, y2);
```

- 1개의 창에 같이 띄우려면

```
>> subplot(1,2, 1)      → 1 x 2 grid를 만들고 첫번째 element에 access
>> plot(t, y1);
>> subplot( 1, 2, 2 );
>> plot (t, y2);
```



(5) heat map plotting하기

```
>> A = magic(5)
```

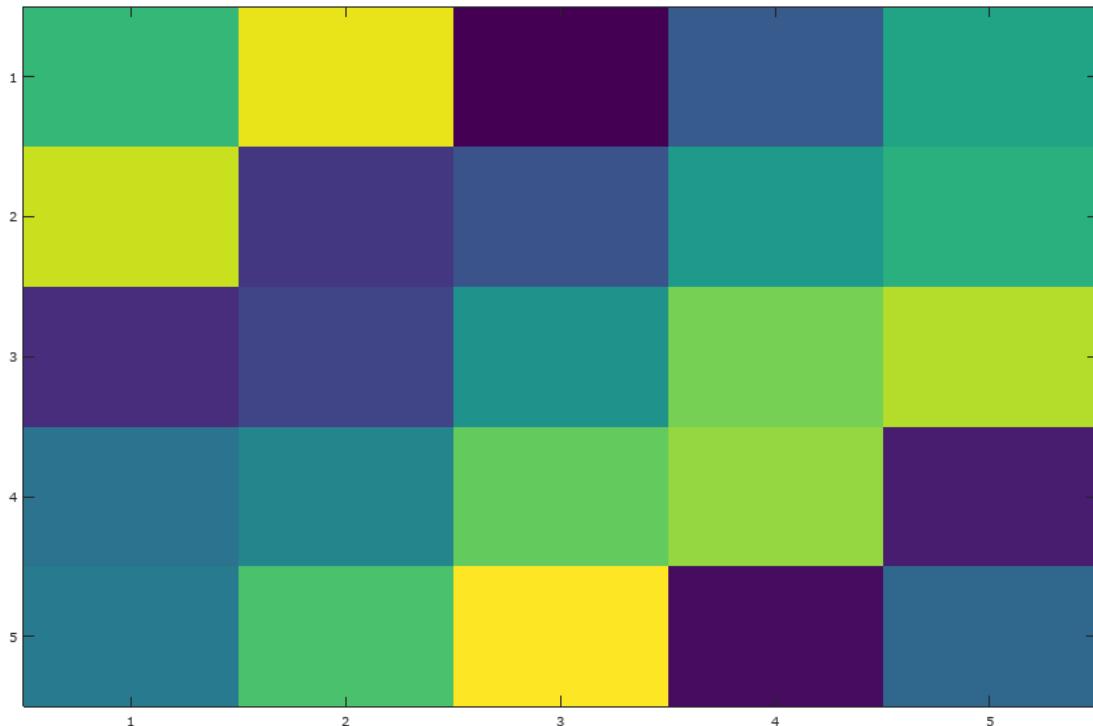
A =

```

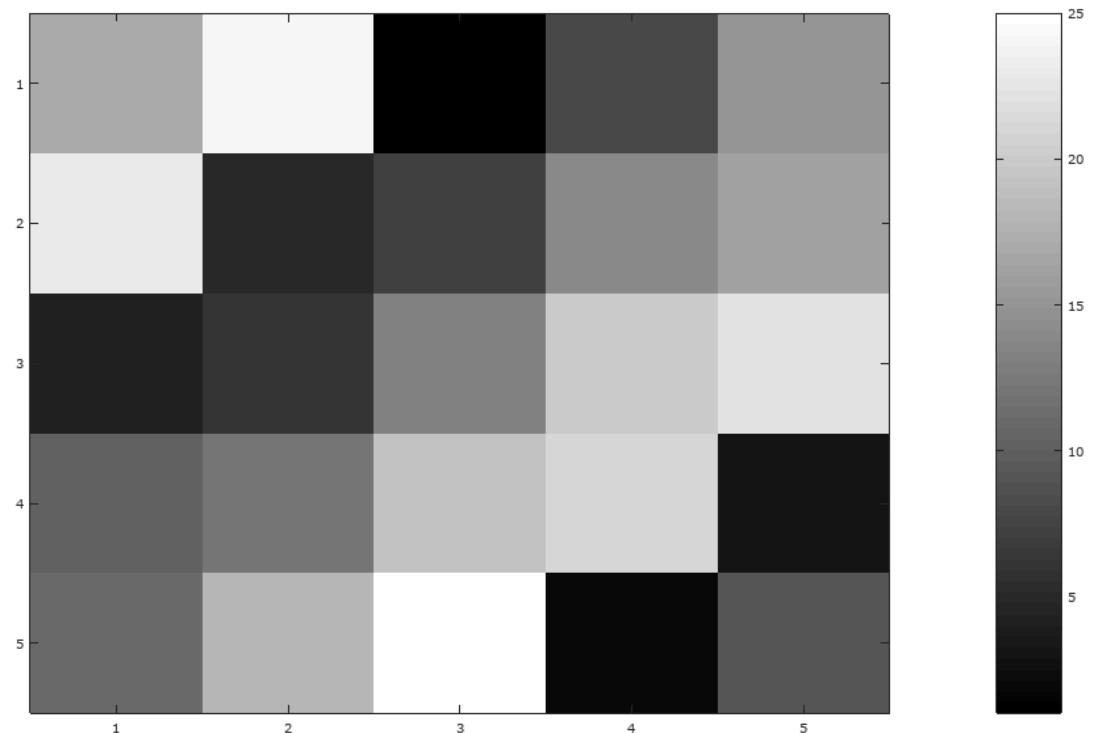
17  24   1   8   15
23   5   7  14  16
  4   6  13  20  22
10  12  19  21   3
11  18  25   2   9

```

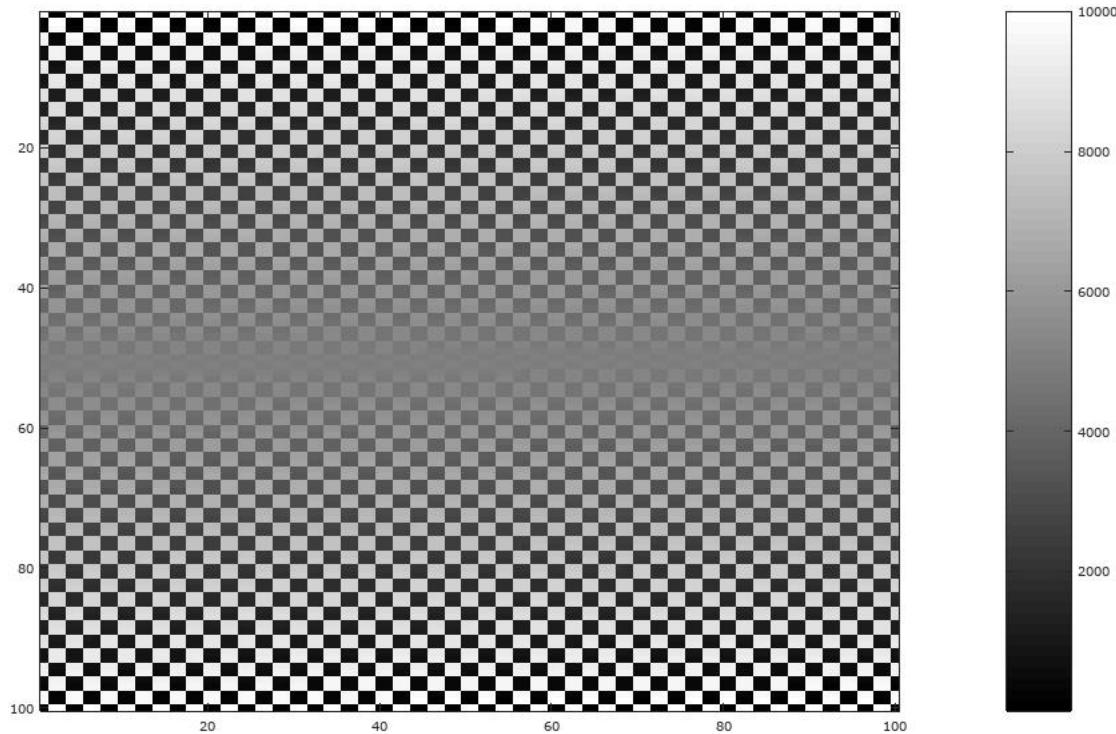
>> **imagesc(A)**



>> **imagesc(A), colorbar, colormap gray;** ➔ comma chaining of commands (연결해서 쓰는 거다)



>> **imagesc(magic(100)), colorbar, colormap gray;**



1.5. Control statements

(1) $v(i) =$ ➔ 이걸로 처리하면 된다!!

- 빈공간 만들기

>> **v = zeros(10, 1)**

>> for i=1:10, ← 1부터 10까지

>> for i=1:10,

$v(i) = 2^i;$

end;

>> v

v =

2

4

8

16

32

64

128

256

512

1024

>> i = 1;

```
>> while i <= 5,  
v(i) = 100;  
i=i+1;  
end;  
>> v  
v =
```

```
100  
100  
100  
100  
100  
64  
128  
256  
512  
1024
```

- break 사용하는 방법

```
i = 1;  
while true,  
    v(i) = 999;  
    i = i+1;  
    if i == 6,  
        break;  
    end;  
end;
```

처음 5개는 999로 바뀌게 된다.

```
>> v(1)  
ans = 999
```

- if statement

```
if v(1) == 1,  
    disp('the value is one');  
elseif v(1) ==2,  
    disp('the value is two');  
else,  
    disp('haha');  
end;
```

→ elseif만 다른 듯. 끝에 항상 ;을 넣는다는 거랑.

- function을 사용해보자.

pwd해서 있는 directory에

function 이름.m

으로 저장하면 된다!!

그 다음에 바로 squareThisNumber(5) 하면 적용된다. → .m 으로 저장한다는 것이 핵심포인트.

- search path 사용하기

>> **addpath('C:\Users\Owner\Desktop')** 이렇게 해놓으면 다른 directory에 있어도 여기있는 함수를 찾아보게 된다.

- octave는 많은 value를 return할 수 있다.

```
function [y1, y2] = squareAndCubeThisNumber(x)
y1 = x ^2;
y2 = x ^3;
```

```
>> [a, b] = squareAndCubeThisNumber(x)
```

```
>>a
```

```
>>b
```

다른 게 나오게 된다.

- 종합

→ Cost function

```
>> X = [1 1; 1 2; 1 3]
X =
    1   1
    1   2
    1   3
>> y = [1; 2; 3]
y =
    1
    2
    3
>> theta = [0;1];
```

```

function J = costFunctionJ(X, y, theta)
% X is the "design matrix" containing our training examples.
% y is the class labels

m = size(X, 1); % number of training examples
predictions = X*theta; % predictions of hypothesis on all m examples
sqrErrors = (predictions-y).^2; % squared errors

J = 1/(2*m) * sum(sqrErrors);

```

1.6. Vectorization

- for loop 보다는 dot product를 하면 한 줄 만에 계산할 수 있다!!

Vectorization example.

$$\rightarrow h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

$$= \theta^T x$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Unvectorized implementation

```

→ prediction = 0.0;
→ for j = 1:n+1,
    prediction = prediction + theta(j) * x(j)
end;

```

Vectorized implementation

```

→ prediction = theta' * x;

```

- Gradient descent 도 vectorized implementation을 할 수 있다. 식을 간단히 써서 $\theta := \theta - \alpha \delta$. 이거랑 같은 것이다. For loop은 가급적 지양하자. 굉장히 오래 걸릴 수 있다!!

$$\begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ \theta_2 &:= \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)} \\ (n = 2) \end{aligned}$$

Vectorized implementation:

$$\Theta := \Theta - \frac{\lambda}{m} \delta \quad \text{where } \delta = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

$\delta = \begin{bmatrix} \delta_0 \\ \delta_1 \\ \delta_2 \end{bmatrix}$

$\delta_0 = \frac{1}{m} \sum_{i=1}^m (h_\theta(\delta) - y^{(i)}) x_0^{(i)}$

$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \end{bmatrix}$

$h_\theta(\delta) = 2\delta_0 + 5\delta_1$

$u(j) = 2v(j) + 5w(j) \quad (\text{for all } j)$

$u = 2v + 5w$

$(h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$

$+ (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$

\dots

- 실전 적용 (Octave)

Gradient descent 하는 법

```
<computeCost.m>
function J = computeCost(X, y, theta)
J = 0;
J = 1/(2*m) * sum(((X * theta) - y).^2)
end
```

```
X = [ones(m, 1), data(:,1)]; % Add a column of ones to x
theta = zeros(2, 1); % initialize fitting parameters
```

```
% Some gradient descent settings
iterations = 1500;
alpha = 0.01;
```

```
fprintf('Testing the cost function ...')
% compute and display initial cost
J = computeCost(X, y, theta);
fprintf('With theta = [0 ; 0] Cost computed = %f', J);
```

- 미분을 하는 방법

디렉토리에서 sympy 다운로드하기. (<https://github.com/cbm755/octsympy>)

* 버그 해결 : <https://www.nntp.perl.org/group/perl.modules/2010/01/msg69392.html>

Download the symbolic-win-py-bundle-2.5.0.zip file from releases.

Start Octave

At the Octave prompt, type pkg install symbolic-win-py-bundle-2.5.0.zip.

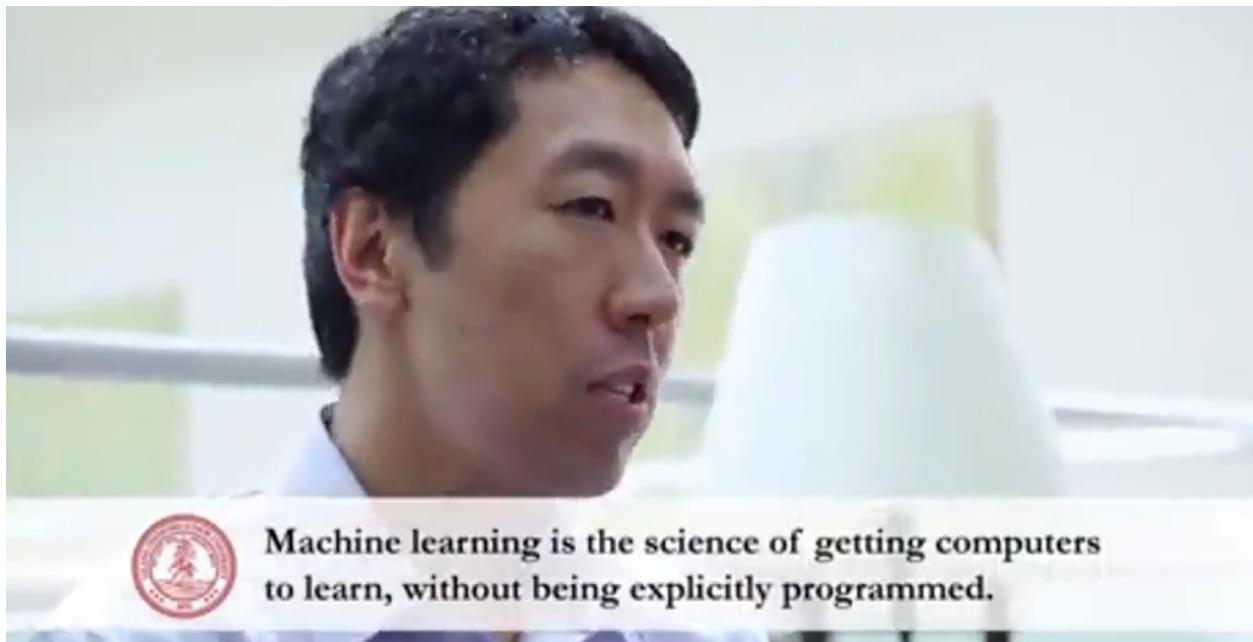
At the Octave prompt, type pkg load symbolic.

At the Octave prompt, type syms x, then $f = (\sin(x/2))^3$, diff(f, x), etc.

→ 버그는 해결했는데 syms x 가 안되네 여전히.. (18. 3. 11)

II. Machine learning

0. What is machine learning!



0.1. 왜 필요한지? (동기부여 뿐뿐)

Examples:

- Database mining
 - Large datasets from growth of automation/web.
E.g., Web click data, medical records, biology, engineering
- Applications can't program by hand.
 - E.g., Autonomous helicopter, handwriting recognition, most of Natural Language Processing (NLP), Computer Vision.
- Self-customizing programs
 - E.g., Amazon, Netflix product recommendations
- Understanding human learning (brain, real AI).

0.2. Machine learning의 정의

- Field of study that gives computers the ability to learn without being explicitly programmed. (Arthur Samuel(1959))
- Well posed learning problem : A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.
→ P는 승률을 의미할 수도
- 예시 :

"A computer program is said to *learn* from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E."

Suppose your email program watches which emails you do or do not mark as spam, and based on that learns how to better filter spam. What is the task T in this setting?

- Classifying emails as spam or not spam. T
- Watching you label emails as spam or not spam. E
- The number (or fraction) of emails correctly classified as spam/not spam. P
- None of the above—this is not a machine learning problem. P

0.3. 머신러닝의 종류는 무엇인가?

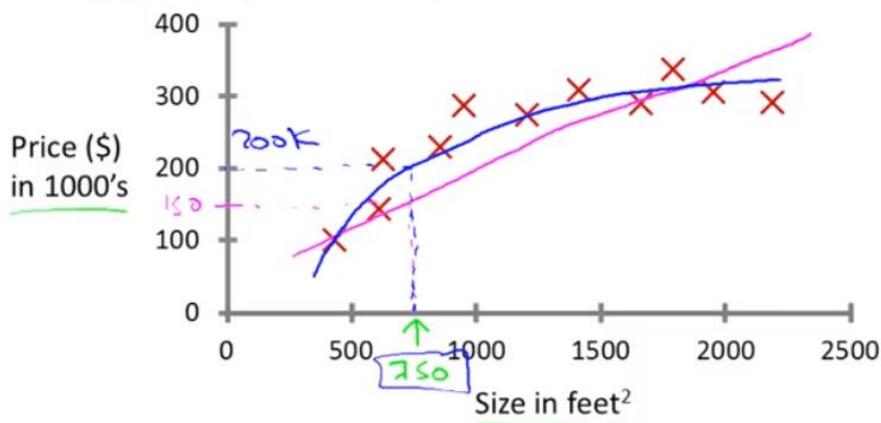
(1) Supervised learning

- we gave the "right answers" → It can be a **regression problem**(continuous) or a **classification problem**(discrete).

→ 즉, 데이터 셋을 먼저 주는 것이 핵심이다. 내가 알고 있었던 machine learning의 종류인 것 같다!

(1-1) regression problem :

Housing price prediction.



Supervised Learning
"right answers" given

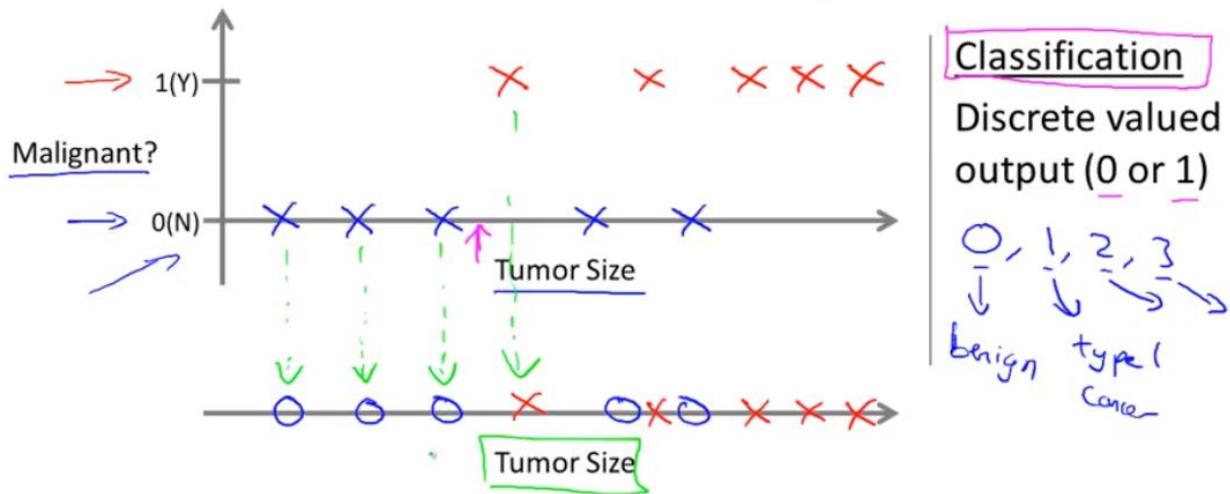
Regression: Predict continuous valued output (price)

→ machine learning question is "how much is the 750 square feet house sold?"

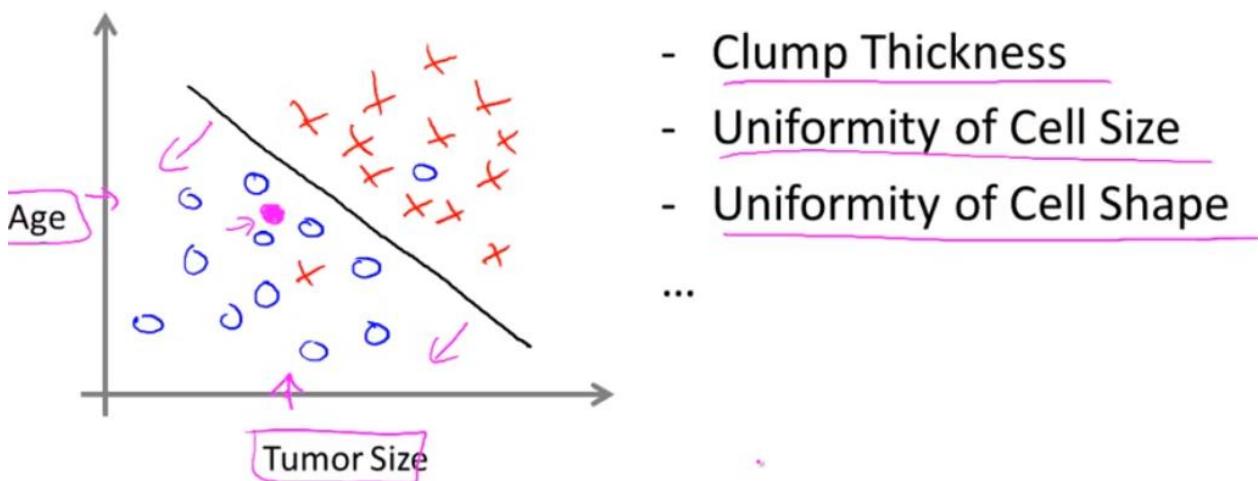
- example 2 : Regression - Given a picture of a person, we have to predict their age on the basis of the given picture

(1-2) classification problem : 여기서의 classification은 종속변수가 discrete하다. 0, 1 뿐만 아니라, 0, 1, 2, 3이어도 괜찮다. → logistic regression

Breast cancer (malignant, benign)



→ 0, 1을 다른 기호로 처리해서 1차원으로 표시할 수도 있다. 그렇게 하면 tumor size 말고 다른 독립변수도 plotting할 수 있게 된다. 2차원으로 하면 Machine learning의 목적은 age, tumor size를 알 때, 과연 그 사람의 tumor가 benign한지, malignant한지를 알아내는 것이다. 그것을 위해서는 그 경계에 선을 그릴 수 있어야 한다. 따라서 그것이 바로 머신러닝의 작업이 된다. 만약에 변수가 2개가 아니라 무한대로 늘어나게 된다면? 그것도 처리할 수 있다.



- example 2 :

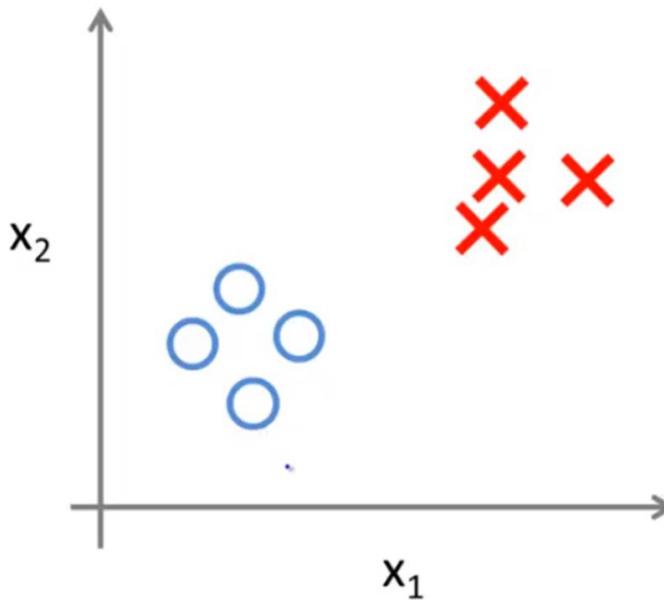
Given data about the size of houses on the real estate market, try to predict their price. Price as a function of size is a continuous output, so this is a regression problem.

We could turn this example into a classification problem by instead making our output about whether the house "sells for more or less than the asking price." Here we are classifying the houses based on price into two discrete categories.

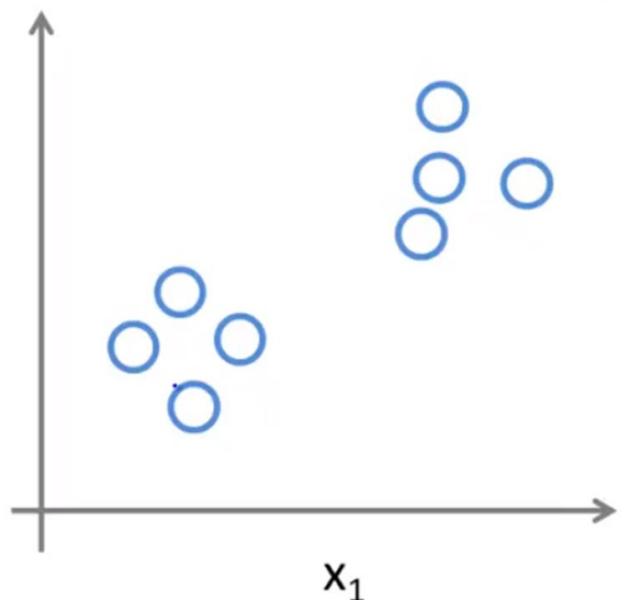
(2) Unsupervised learning

- 정의 : data의 label이 다르다!! : unsupervised learning은 answer를 주지 않은 채 structure를 찾아보라고 하는 것이다. Not labeled. 즉, we have no or little idea what our results should look like.
→ "We can derive structure from data where we don't necessarily know the effect of the variables."
- 목적 : Data set을 받고 find some structure in this data set.

Supervised Learning



Unsupervised Learning

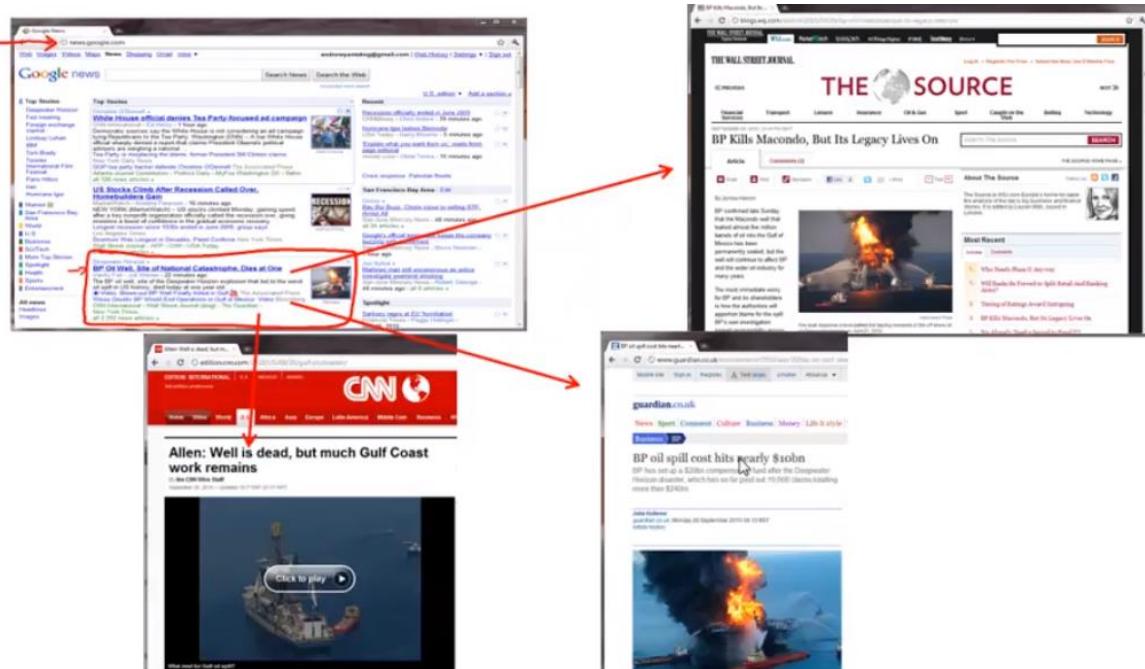


- 차이점에 대한 심화된 설명 : **When you already have the answer or the answer is already given into the data set then it's Supervised learning.** The input data has correct output. Like House price, you already know the price of some area(sq.ft). So your input data (data set you have) already have correct output. Or like if the cancer is malignant or not? You already have a data set with correct answer.

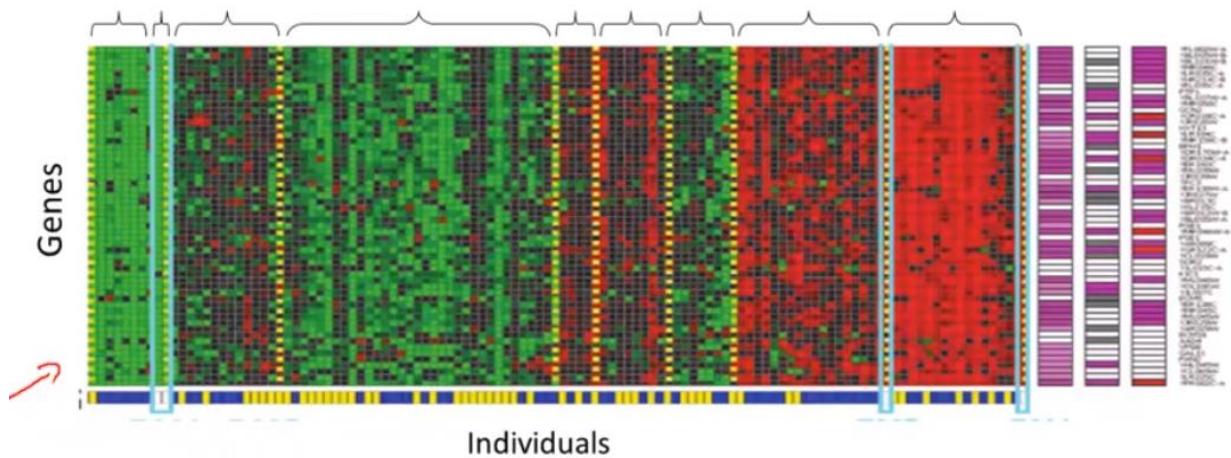
When you have data set and don't know the answer then it's Unsupervised learning. Like you've bunch of photos of several people but you don't know which is which. So you categorise them with which ones are alike. You don't have the correct answer here. You categorise them based on similarities. I hope you get the difference. Cheers!

(2-1) Clustering algorithm : 두 cluster로 나눌 수 있다.

- example 1 : News를 clustering한다. 이게 뭔지는 알려주지 않았는데도..



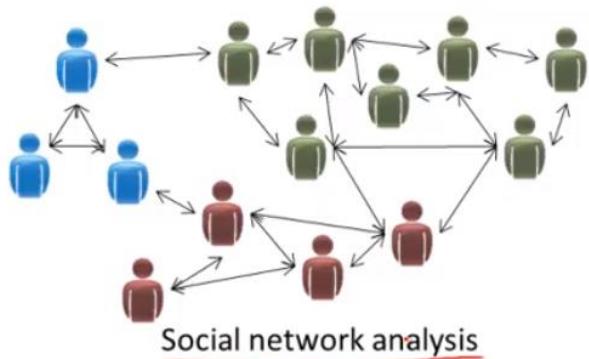
- example 2 : 딱히 먼저 분류해놓은 것은 아닌데, 유전자 발현 여부의 패턴에 따라서 분류를 해놨다. cluster로!



- 다른 예 :



Organize computing clusters



Social network analysis



Market segmentation



Astronomical data analysis

- 예시 문제 :

Of the following examples, which would you address using an unsupervised learning algorithm? (Check all that apply.)

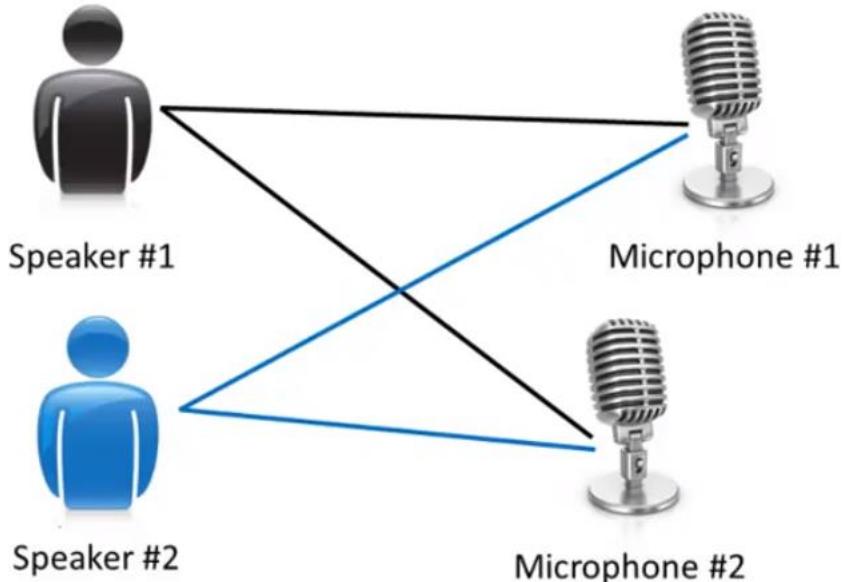
- Given email labeled as spam/not spam, learn a spam filter.
- Given a set of news articles found on the web, group them into set of articles about the same story.
- Given a database of customer data, automatically discover market segments and group customers into different market segments.
- Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not.

(2-2) Non-clustering algorithm : "Cocktail Party Algorithm"

- Non-clustering: The "Cocktail Party Algorithm", allows you to find structure in a chaotic environment. (i.e. identifying individual voices and music from a mesh of sounds at a cocktail party)

→ Cocktail party problem : 파티에 가면 잘 안 들려!

Cocktail party problem



→ 와! 두 개의 목소리를 나눌 수 있다. 거리로 구분을 했었다. 거리에 따라서 한 쪽에서는 크게, 한쪽에서는 작게 들릴 테니.

```
[W,s,v] = svd((repmat(sum(x.*x,1),size(x,1),1).*x)*x');
```

→ 이 것은 octave로 한 줄로 이 알고리즘을 짤 수 있다. python으로 하면 여러 줄로 짜야 한다. Octave은 prototyping을 할 수 있다. 이걸 하면 더 빨리 이해할 수 있고, 더 잘 구동시킬 수 있다. 많은 회사들에서도 octave를 통해서 prototype를 만들어보고 그 다음에 잘 되면 python, java로 이행한다.

1. Supervised learning - Model representation

1.1. Supervised learning의 이해

1.1.1. Notation :

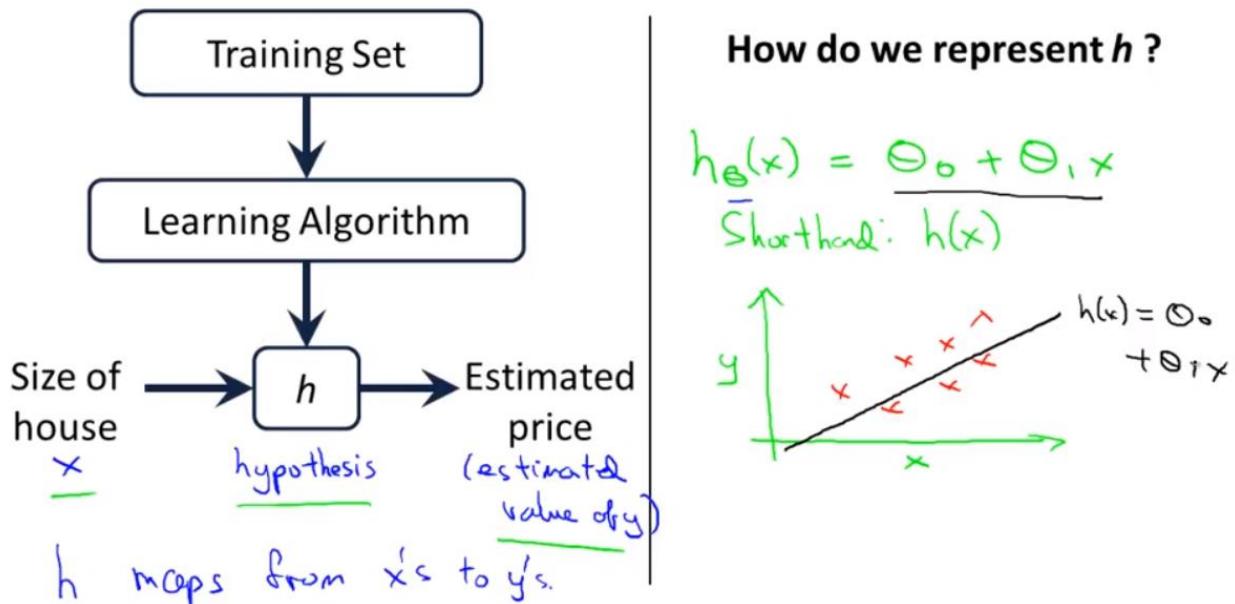
<u>Training set of housing prices (Portland, OR)</u>	<u>Size in feet²(x)</u>	<u>Price (\$ in 1000's)(y)</u>
	→ 2104	460
	1416	232
	1534	315
	852	178

Notation:

- m = Number of training examples
 - x 's = "input" variable / features
 - y 's = "output" variable / "target" variable
- (x, y) - one training example
 $(x^{(i)}, y^{(i)})$ - i^{th} training example

1.1.2. h (hypothesis)

: hypothesis는 독립, 종속변수를 있는 함수이다! To describe the supervised learning problem slightly more formally, our goal is, given a training set, to learn a function $h : X \rightarrow Y$ so that $h(x)$ is a "good" predictor for the corresponding value of y . For historical reasons, this function h is called a hypothesis.



→ 이 경우에는 univariate linear regression이다. 변수가 하나인 경우!!

1.1.3. $J(\theta_0, \theta_1)$ (cost function)

Training Set	Size in feet ² (x)	Price (\$) in 1000's (y)
	2104	460
	1416	232
	1534	315
	852	178

$\left. \begin{array}{l} \\ \\ \\ \\ \end{array} \right\} m=47$

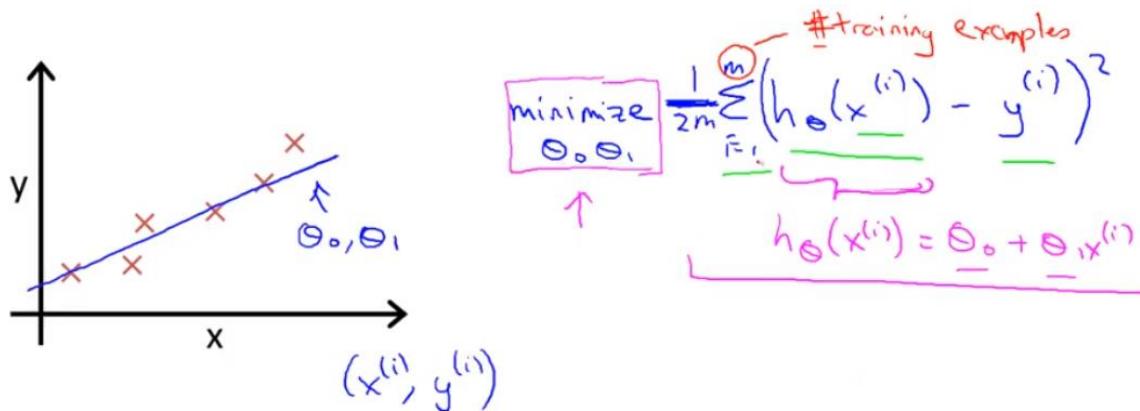
Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

θ_i 's: Parameters

How to choose θ_i 's?

- 목적 : How do we choose "hypothesis"? We use cost function (square error function)

- 방법 : Choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training examples (x, y)



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

minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1 Cost function
Squared error function

- 언어 (이해 못하고 있음 → 이해함!!) :

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

→ To break it apart, it is $\frac{1}{2} \bar{x}^2$ where \bar{x} is the mean of the squares of $h_{\theta}(x^{(i)}) - y^{(i)}$, or the difference between the predicted value and the actual value.

→ 2로 다시 나누는 이유 : The mean is halved $\frac{1}{2}$ as a convenience for the computation of the gradient descent, as the derivative term of the square function will cancel out the $\frac{1}{2}$ term.

→ 제곱을 편미분하면 2가 내려오니까 2로 나눠버리면 편하잖아.

$$f(x)^2 = 2*f(x)*f'(x)$$

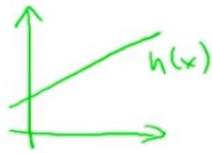
1.1.4. h , J 와 목표(Goal)식 :

Hypothesis:

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

Parameters:

$$\underline{\theta_0, \theta_1}$$



Cost Function:

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

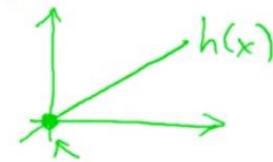
Goal: $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

Simplified

$$h_{\theta}(x) = \underline{\theta_1 x}$$

$$\underline{\theta_0 = 0}$$

$$\underline{\theta_1}$$

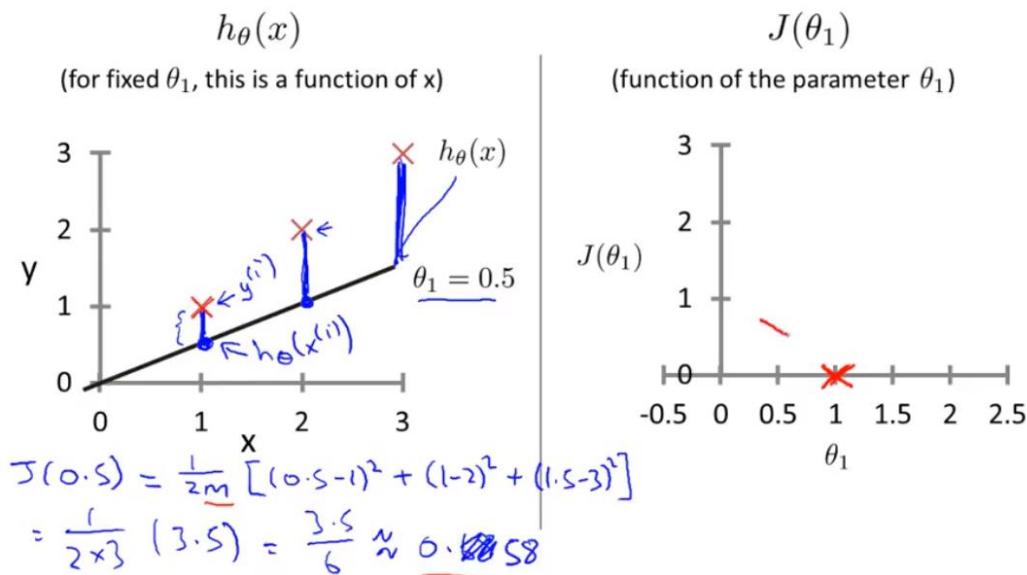


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\underset{\theta_1}{\text{minimize}} J(\theta_1) \quad \underline{\theta_1 x^{(i)}}$$

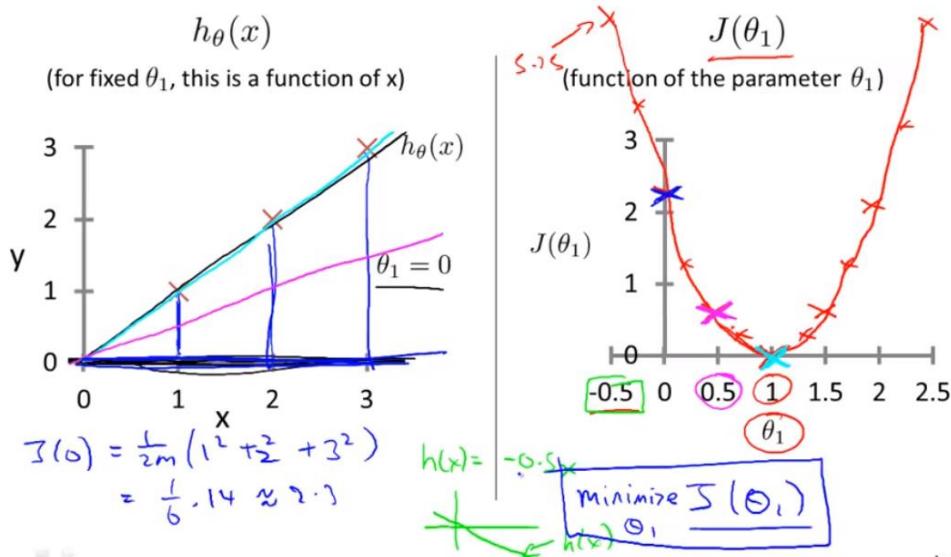
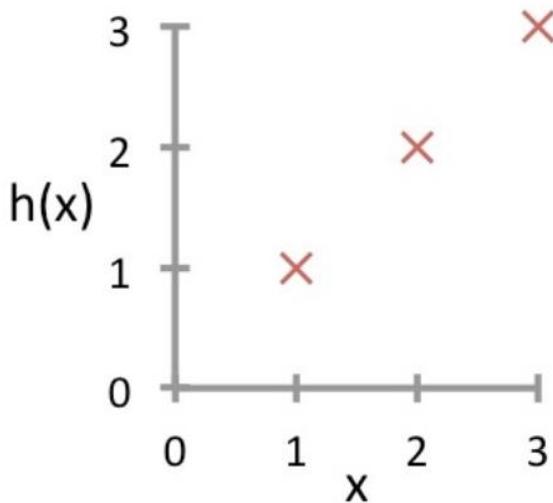
(1) Simplified version 실제 계산해보기 : $\theta_0=0$ 이라고 가정

$$(1) J(0.5) = 0.58$$



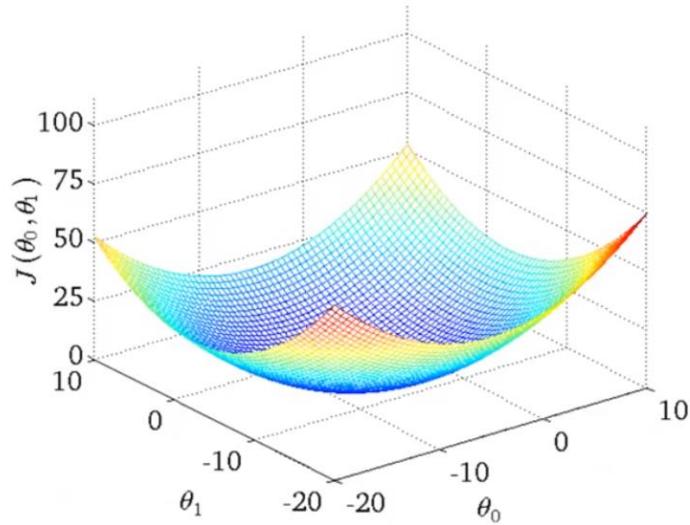
$$(2) J(0) = 14/6$$

Suppose we have a training set with $m=3$ examples, plotted below. Our hypothesis representation is $h_{\theta}(x) = \theta_0 + \theta_1 x$, with parameter θ_1 . The cost function $J(\theta_1)$ is $J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$. What is $J(0)$?

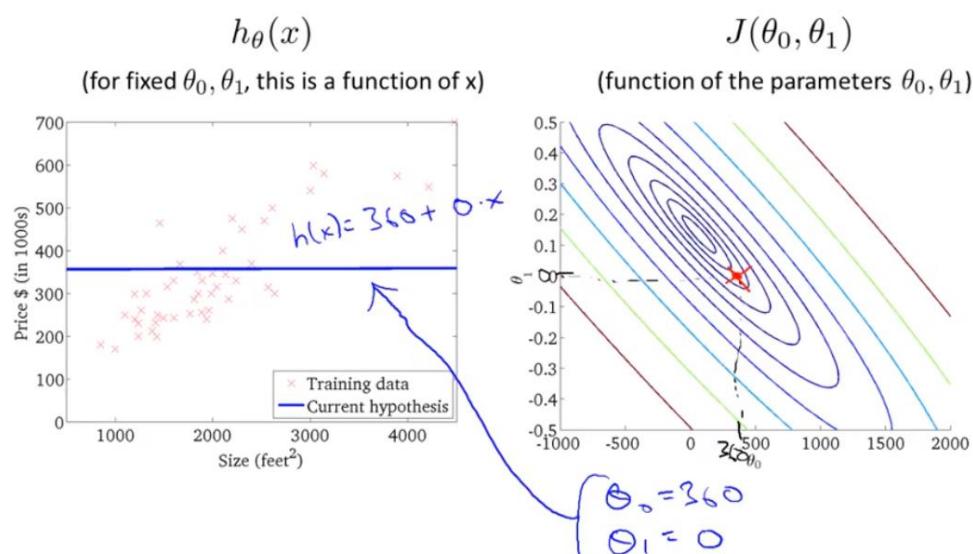
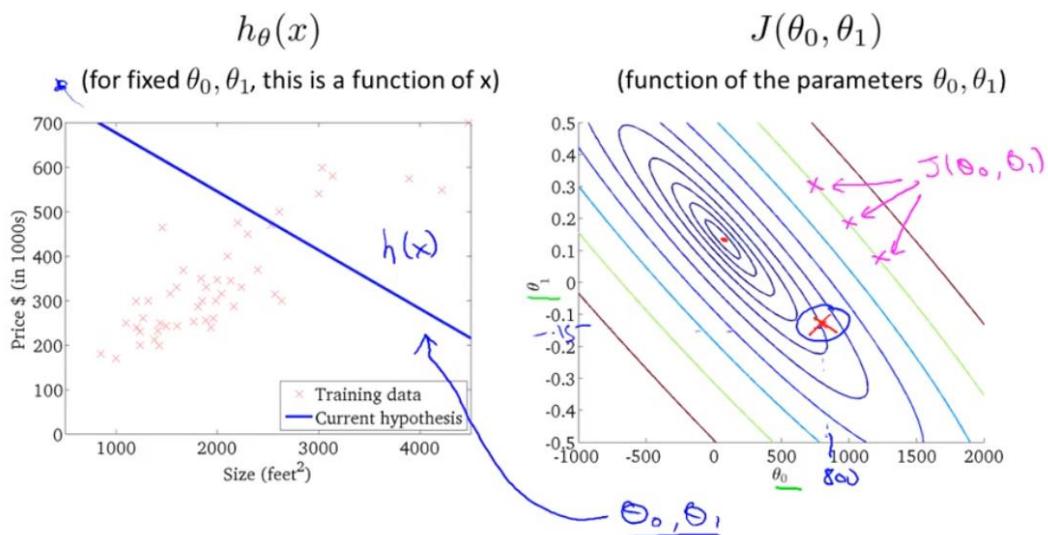


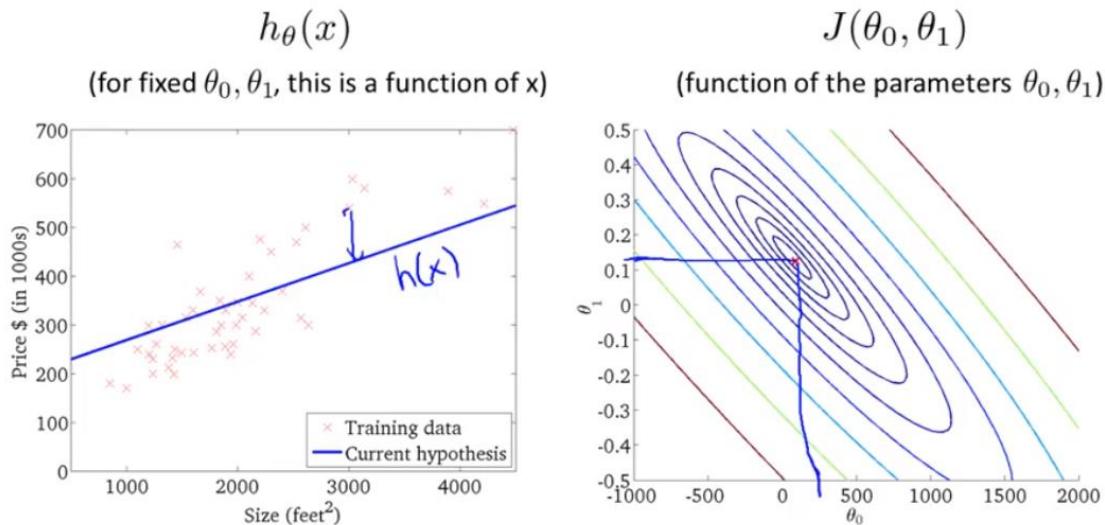
→ 아하! 이런 식으로 $J(0)$, $J(1)$, $J(0.5)$ 를 해가다 보면 포물선이 그려지고, 결국에는 $J(1)$ 이 최선이라는 것을 알 수 있구나!. θ_1 이 1일 때가 Minimize J 를 할 수 있는 것이다.

(2) θ_1 과 θ_0 이 다 살아있을 때의 cost function



→ Contour figure





→ 나중에는 θ_0, θ_1 이 아니라 더 많아져서 θ_2, \dots 등등 이 생기게 되면, 이것은 visualization이 힘들어지게 된다.

Supervised machine learning : Cost function minimize method

→ J (Cost function) minimize method : Gradient descent, Normal equation

1.2. J 를 minimize하는 알고리즘-1 : gradient descent

1.2.1. Gradient descent 정의 : minimize를 하는 알고리즘. 사실 J 가 어떤 것이든 폭넓게 적용할 수 있다.

$$\text{Have some function } J(\theta_0, \theta_1) \quad J(\theta_0, \theta_1, \dots, \theta_n)$$

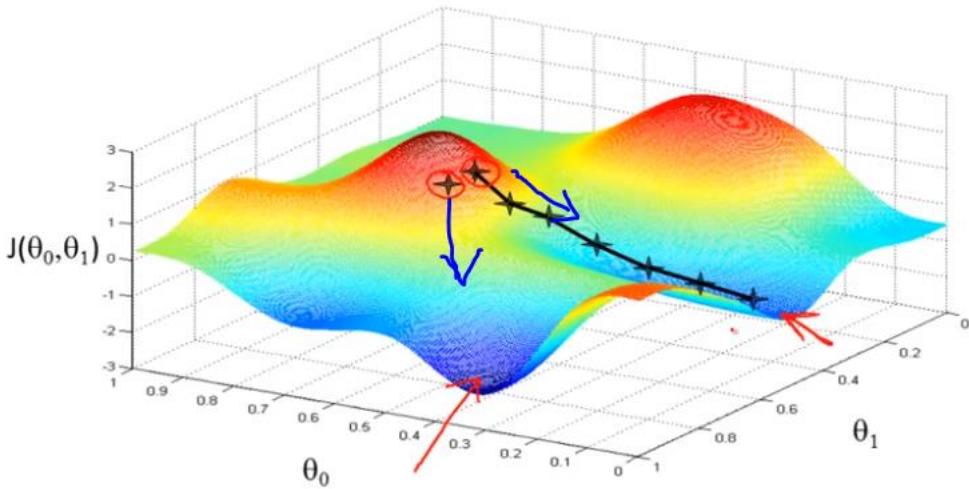
$$\text{Want } \min_{\theta_0, \theta_1} J(\theta_0, \theta_1) \quad \min_{\theta_0, \dots, \theta_n} J(\theta_0, \dots, \theta_n)$$

Outline:

- Start with some $\underline{\theta_0, \theta_1}$ (say $\theta_0 = 0, \theta_1 = 0$)
- Keep changing $\underline{\theta_0, \theta_1}$ to reduce $\underline{J(\theta_0, \theta_1)}$
until we hopefully end up at a minimum

- 직관 : 한 점에서 시작한다. → 360도를 둘러보고, 어느 방향으로 babystep하면 더 빨리 내려갈 수 있을까?를 결정하고 나아간다. → babystep을 간 지점에서 다시 그렇게 한다.

- 특징 : 초기점이 조금만 달라져도 도달하는 local minimum이 달라질 수 있다.



1.2.2. 알고리즘 : We make steps down the cost function in the direction with the steepest descent.

```

repeat until convergence {
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$    (for  $j = 0$  and  $j = 1$ )
}
    learning rate
  
```

Repeat until convergence:

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

:=라는 것은 assignment을 의미하는 것이다.

=라는 것은 truth assertion을 의미한다.

α : learning rate \rightarrow 얼마나 이동하는가. α 가 크면 babystep이 굉장히 큰 것이다.

중요 포인트(for $j=0$ and $j=1$) : simultaneous update이다. θ_0, θ_1 이 동시에 바뀌어야 한다. 만약에 non-simultaneous update는 gradient descent라고 하지는 않는다. 그것은 잘 작동할 수도 있지만, gradient descent라고 부르지는 않고 다른 파생 알고리즘이다.

Gradient descent algorithm

```

repeat until convergence {
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$    (for  $j = 0$  and  $j = 1$ )
}
    learning rate
  
```

θ_0, θ_1 <u>Assignment</u> $a := b$ $a := a + 1$	<u>Truth assertio</u> $a = b$ \leftarrow $a = a + 1$ \times
-----------------------------------------------------------------------	-----------------------------------------------------------------------

Correct: Simultaneous update

```

 $\rightarrow \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ 
 $\rightarrow \text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ 
 $\rightarrow \theta_0 := \text{temp0}$ 
 $\rightarrow \theta_1 := \text{temp1}$ 
  
```

Incorrect:

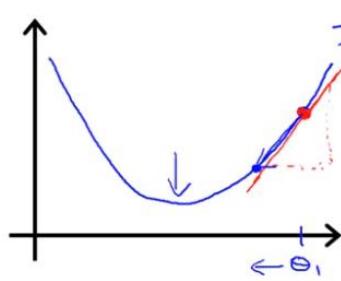
```

 $\rightarrow \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ 
 $\rightarrow \theta_0 := \text{temp0}$ 
 $\rightarrow \text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ 
 $\rightarrow \theta_1 := \text{temp1}$ 
  
```

partial derivative term : 기울기를 의미한다. (방향)

- 이 알고리즘에 대한 직관 :

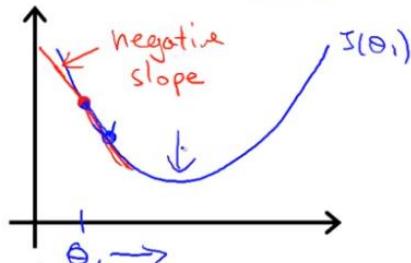
- Derivative term : $J(\theta_1)$ 의 경우를 생각해보자.



$$J(\theta_1) \quad (\theta_1 \in \mathbb{R})$$

$$\theta_1 := \theta_1 - \frac{\partial}{\partial \theta_1} J(\theta_1) \geq 0$$

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



$$\frac{\partial}{\partial \theta_1} J(\theta_1) \leq 0$$

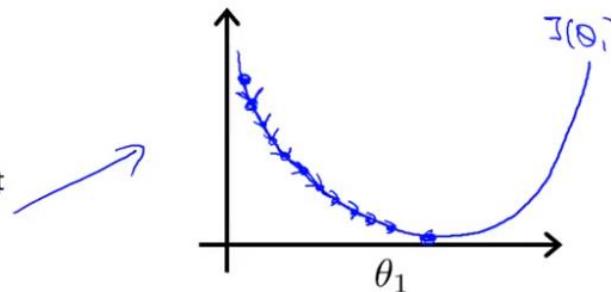
$$\theta_1 := \theta_1 - \alpha \cdot (\text{negative number})$$

→ 가운데로 θ 가 이동하고 있다.

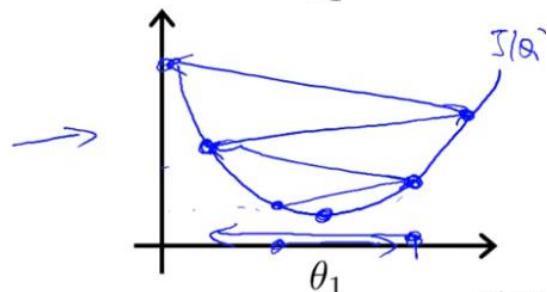
- α (learning rate)가 너무 크면 어떻게 될까?

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

If α is too small, gradient descent can be slow.



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



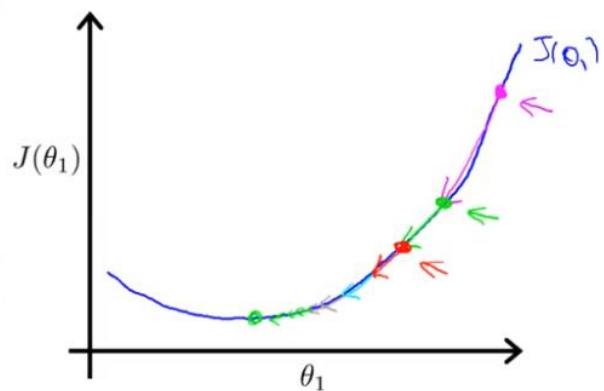
- 만약에 local minimum에 있게 되면 1 step of gradient descent로 인해서 θ 가 변하지 않는다. 미분값이 0이 되기 때문이다.

- 1 step of gradient descent의 특징 : local minimum에 가까워질수록 작은 step을 가게 된다. (α 가 불변이더라도)

Gradient descent can converge to a local minimum, even with the learning rate α fixed.

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

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.



1.2.3. Coding Gradient descent for linear regression (실제 코딩)

- linear regression model에 어떻게 gradient descent를 적용시킬 수 있을까?
- "Batch" gradient descent
 - Batch means each step of gradient descent uses all the training examples. 다 사용!.

Gradient descent algorithm

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

Linear Regression Model

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

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

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

- 미분을 계산해보자.

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{2}{2m} \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{2}{2m} \cdot \frac{1}{m} \sum_{i=1}^m (\underline{\theta_0 + \theta_1 x^{(i)}} - y^{(i)})^2 \end{aligned}$$

$$\theta_0, j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

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

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\
 &= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y) \\
 &= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x_i - y \right) \\
 &= (h_{\theta}(x) - y) x_j
 \end{aligned}$$

Gradient descent algorithm

repeat until convergence {

$$\begin{aligned}
 \theta_0 &:= \theta_0 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \right] \\
 \theta_1 &:= \theta_1 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \right]
 \end{aligned}
 \quad \left. \begin{array}{l} \text{update} \\ \theta_0 \text{ and } \theta_1 \\ \text{simultaneously} \end{array} \right\}$$

}

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

→ $x_0 = 1$ 로 설정하면 vectorization해서 한 줄로 표현할 수 있다.

$$h(x) = \theta_0 * x_0 + \theta_1 * x_1 \rightarrow \theta * x$$

Gradient Descent

Previously ($n=1$):

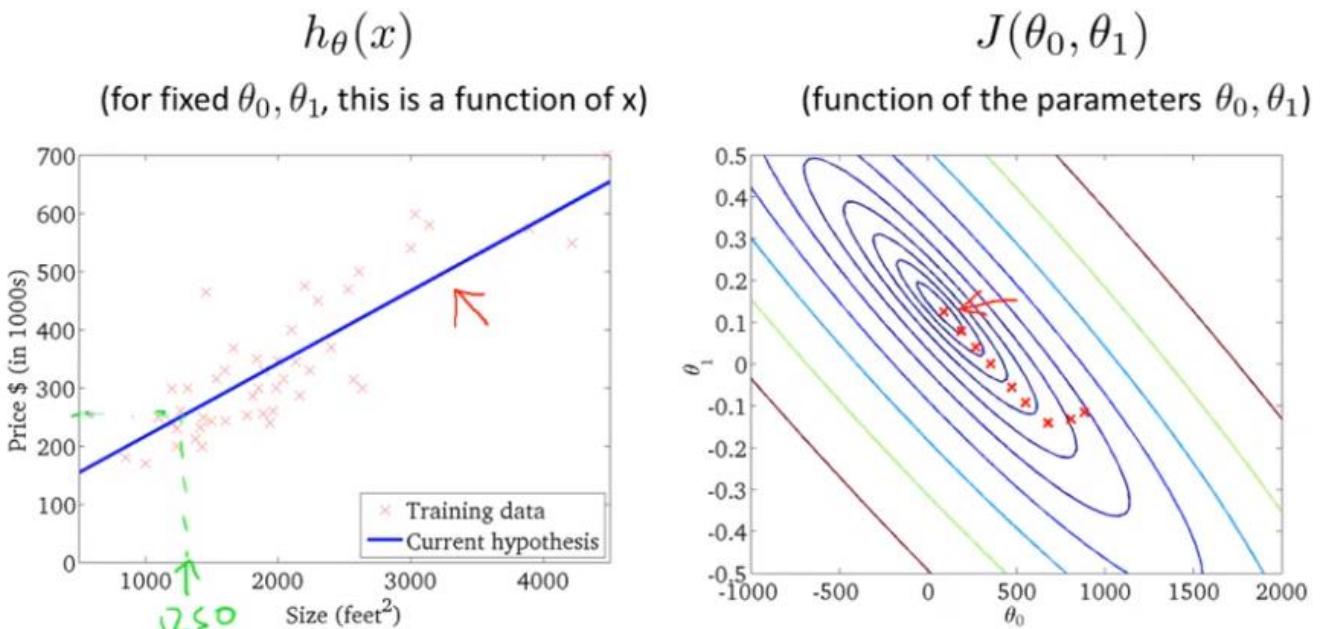
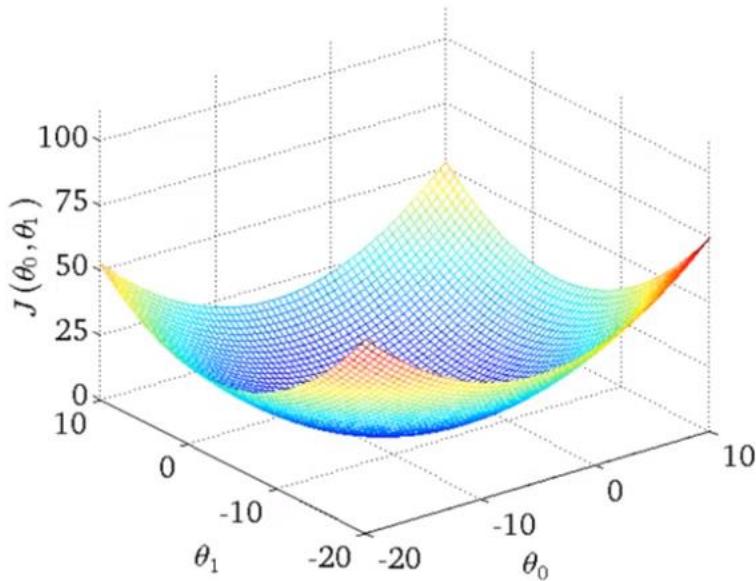
$$\begin{aligned}
 \text{Repeat} \{ & \\
 \rightarrow \theta_0 &:= \theta_0 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \right] \\
 &\quad \boxed{\frac{\partial}{\partial \theta_0} J(\theta)} \\
 \rightarrow \theta_1 &:= \theta_1 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \right] \\
 &\quad \text{(simultaneously update } \theta_0, \theta_1 \text{)}
 \end{aligned}$$

}

New algorithm ($n \geq 1$):

$$\begin{aligned}
 \text{Repeat} \{ & \\
 \rightarrow \theta_j &:= \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right] \\
 &\quad \text{(simultaneously update } \theta_j \text{ for } j = 0, \dots, n \text{)} \\
 &\quad \boxed{x_0^{(i)} = 1} \\
 \rightarrow \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\
 \rightarrow \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} \\
 \rightarrow \theta_2 &:= \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}
 \end{aligned}$$

- linear regression의 cost function은 convex function으로 bowl shape이다. 따라서 무조건 local minimum이 global minimum이 된다.



- 부가 설명 : linear regression을 할 때 gradient descent과 같은 iterative algorithm 말고 선형대수학에서 배운 최소화 방법이 있다. "Normal equation method"가 그것이다. 그러나 larger data sets에서 전자의 효과가 더 좋은 것으로 밝혀졌다.

1.2.4. Gradient descent in Multivariate linear regression (Week 2)

→ Gradient descent for multiple variables

- 여러 개의 feature(특징)을 통해서 집값을 예측하고 싶을 때의 선형 회귀

Size (feet ²)	Price (\$1000)
x	y
2104	460
1416	232
1534	315
852	178
...	...

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

Multiple features (variables).

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

Notation:

- n = number of features $n=4$
- $x^{(i)}$ = input (features) of i^{th} training example.
- $x_j^{(i)}$ = value of feature j in i^{th} training example.

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

→ univariate과 다르게 multivariate은 여러 변수로 hypothesis를 만든다.

$$\rightarrow h_{\theta}(x) = \underline{\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n}$$

For convenience of notation, define $x_0 = 1$. ($x_0^{(i)} = 1$)

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \left[\begin{array}{c} [\theta_0, \theta_1, \dots, \theta_n] \\ \hline \theta^T \end{array} \right]$$

$$h_{\theta}(x) = \underline{\theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n}$$

$$= \underline{\theta^T x}.$$

→ x_0 이 1이라고 가정하면 된다! 이를 bias unit이라고도 한다. (neural network에서 사용하는 용어)

→ 내적, 즉 dot product로 표현할 수 있다.

- 이전과 같으나, x 와 θ 를 $n+1$ 차원의 벡터로 생각하면 된다.

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

Parameters: $\underline{\theta_0, \theta_1, \dots, \theta_n}$ Θ $n+1$ -dimensional vector

Cost function:

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

Gradient descent:

Repeat {

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

}

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

- cost function 코딩 방법

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

For linear regression, which of the following are also equivalent and correct definitions of $J(\theta)$?

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

Correct

$\$ \$ J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2 \$ \$$ (Inner sum starts at 0)

This should be selected

- update rule : 미분값을 계산해 넣으면 x_0 가 1로 놓았기 때문에 일반화가 가능하다. 이를 bias unit이라고도 한다.(neural network에서 사용하는 용어)

Gradient Descent

Previously ($n=1$):

Repeat {

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

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \quad \text{(simultaneously update } \theta_0, \theta_1)$$

}

New algorithm ($n \geq 1$):

Repeat {

$$\downarrow \frac{\partial}{\partial \theta_j} J(\theta)$$

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

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

}

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

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

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

1.2.5. Gradient descent in practice-1 : feature scaling

- 변수의 단위가 같아야 잘 된다. 너무 길쭉하면 gradient descent가 너무 오래 걸린다. → similar ranges of values / 균일화해서 어디서 시작해도 무방하게 하려고 하는 것일 수도 있다.(혁동님) / 만약에 $x_0=0, 1000, 2000$ 이라면 $\alpha * 1/m * 합 * x_0$ 이 들쭉날쭉 될 테니 그것을 방지하려고 하는 것 아닐까?(기석)

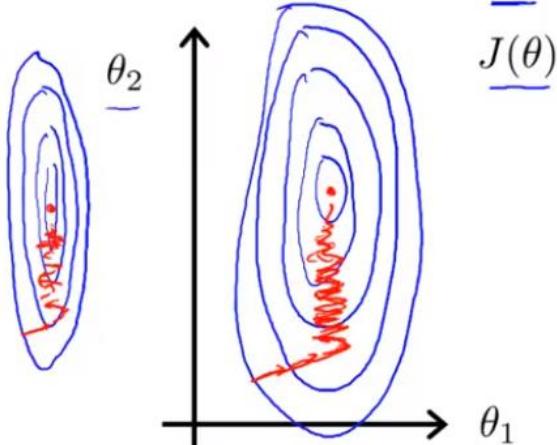
(방법 1) feature scaling을 하면 된다. $-1 < x_i < 1$ range가 되면 좋다. $0 < x < 3$ 도 좋다. $-2 < x < 1$ 도 좋다. $-0.0001 < x < 0.0001$ 은 안 좋다. 결론적으로 $-3 < x < 3$ 이면 괜찮다.

Feature Scaling

Idea: Make sure features are on a similar scale.

E.g. $x_1 = \text{size (0-2000 feet}^2)$

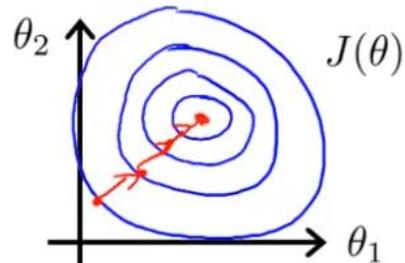
$x_2 = \text{number of bedrooms (1-5)}$



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

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

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



(방법 2) mean normalization

도 좋다.(표준화!!)

- x_0 에는 하지 말 것

- 나눌 때, range로 해도 되고(max-min), standard deviation으로 해도 된다. 엄청 엄밀할 필요는 없다.

Mean normalization

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

E.g. $x_1 = \frac{\text{size} - 1000}{2000}$ Average size = 100
 $x_2 = \frac{\#bedrooms - 2}{5}$ 1-5 bedrooms
 $-0.5 \leq x_1 \leq 0.5$ $-0.5 \leq x_2 \leq 0.5$

$$x_1 \leftarrow \frac{x_1 - \mu_1}{s_1} \quad \begin{array}{l} \text{avg value} \\ \text{of } x_i \\ \text{in training set} \end{array}$$
$$x_2 \leftarrow \frac{x_2 - \mu_2}{s_2}$$

range ($\max - \min$)
(or standard deviation)

- 코딩

```
$$x_i := \frac{x_i - \mu_i}{s_i}$$
```

μ_i is the average of all the values for feature (i)

s_i is the range of values (max - min), or s_i is the standard deviation

1.2.6. Gradient descent in practice-2 : learning rate (α) and "Debugging"

- 목표 :

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

- "Debugging": How to make sure gradient descent is working correctly.

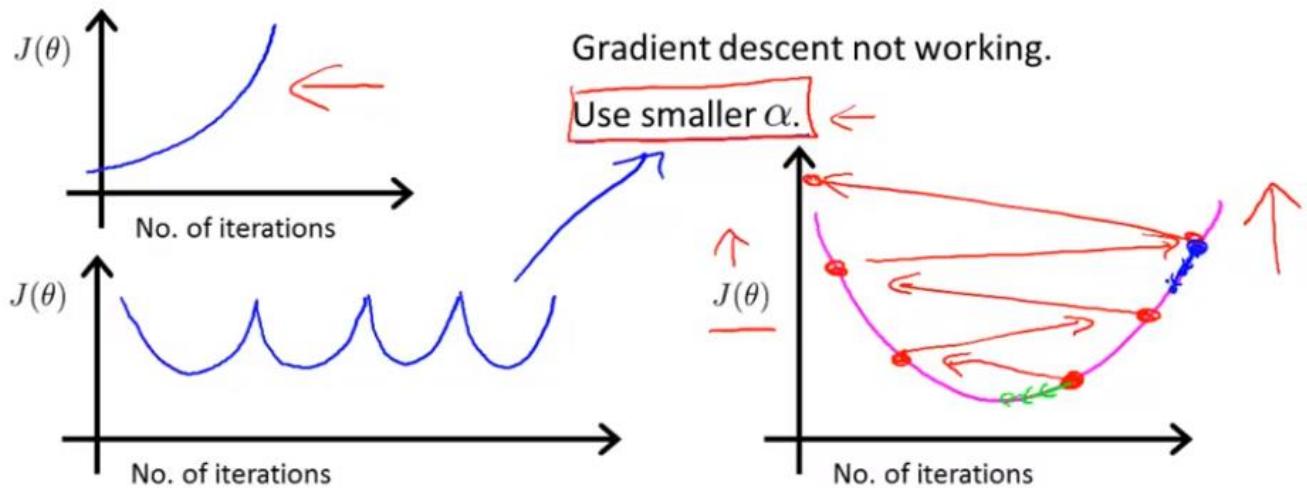
- How to choose learning rate α .

- Make sure gradient descent is working correctly

: $J(\theta)$ should decrease after every iteration \rightarrow it converges

cf) automatic convergence test도 있다.

Making sure gradient descent is working correctly.



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

→ 이렇게 J - iteration 횟수 그래프를 그리면 잘 되는 것이다.

- 잘 안 되는 경우 :

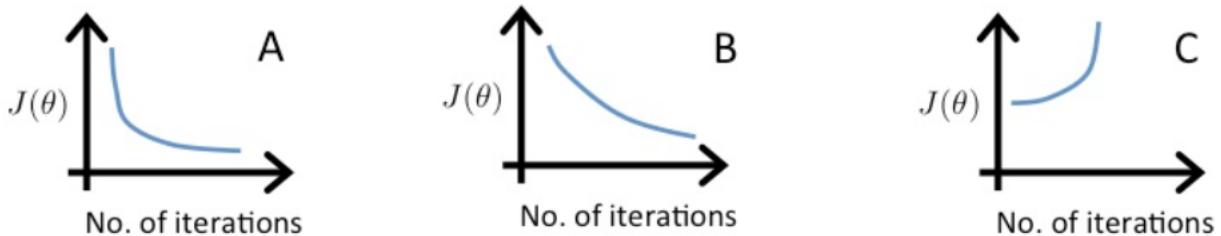
- (1) 만약 수렴하지 않을 때에는 α (learning rate)를 줄인다.
- (2) 파도 칠 때에도 마찬가지 솔루션을 가진다.
- (3) 수렴이 너무 천천히 되면 α 를 증가시킨다 → 증가의 속도는 3배, 10배 정도가 괜찮다.

To choose α , try

$$\dots, \underbrace{0.001}_{\approx x}, \underbrace{0.003}_{\approx 3x}, \underbrace{0.01}_{\approx x}, \underbrace{0.03}_{\approx 3x}, \underbrace{0.1}_{\approx x}, \underbrace{0.3}_{\approx 3x}, 1, \dots$$

- 문제의 예 :

Suppose a friend ran gradient descent three times, with $\alpha = 0.01$, $\alpha = 0.1$, and $\alpha = 1$, and got the following three plots (labeled A, B, and C):



Which plots corresponds to which values of α ?

- A is $\alpha = 0.01$, B is $\alpha = 0.1$, C is $\alpha = 1$.
- A is $\alpha = 0.1$, B is $\alpha = 0.01$, C is $\alpha = 1$.

1.2.7. Gradient descent in practice-3 : Choosing features and polynomial regression

(1) 두 개를 쓰는 것보다 곱한 한 개의 feature를 쓰는 게 나을 수 있다.

Housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

Area

$$\times = \underline{\text{frontage}} * \underline{\text{depth}}$$

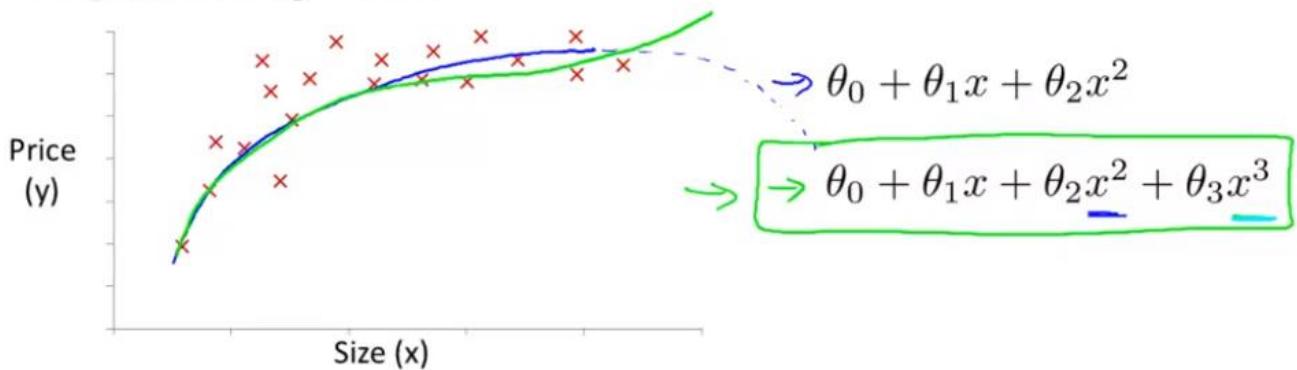


$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{Area}}_{\text{land area}}$$

(2) Polynomial regression : quadratic model

- 방법 : 제곱, 세제곱한 것을 feature로 정하면 된다.

Polynomial regression



$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3 \end{aligned}$$

$$x_1 = (\text{size})$$

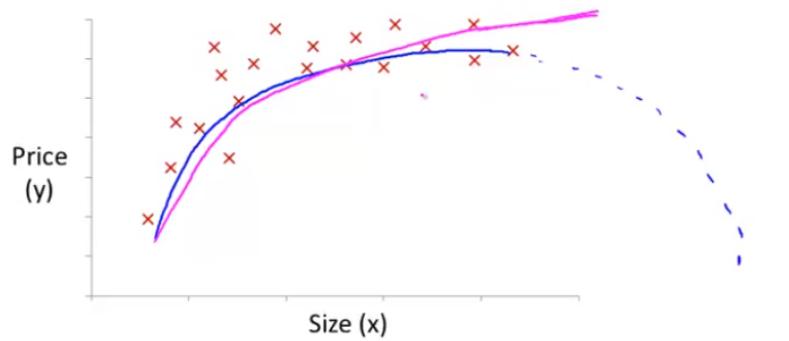
$$x_2 = (\text{size})^2$$

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

- 주의 : scaling을 꼭 한다. range로 나누면 된다.

- cubic model로 안 가도록 하고 싶을 때! square root를 사용한다!!! (2차 함수는 밑으로 포물선을 그리기 때문이다.)

Choice of features



$$\rightarrow h_\theta(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

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

K

1.2.8. Gradient descent in practice-4 : Regularized Gradient descent

- 문제 정의 : polynomial 혹은 many feature를 쓰게 되면 overfitting 문제가 일어날 수 있다.
- 해결 : 이를 해결하기 위해서 regularize시킨다. θ 를 줄이는 cost function 만들기.

Gradient descent

$$\begin{matrix} \theta_0 \\ \vdots \\ \theta_1, \theta_2, \dots, \theta_n \end{matrix}$$

Repeat {

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

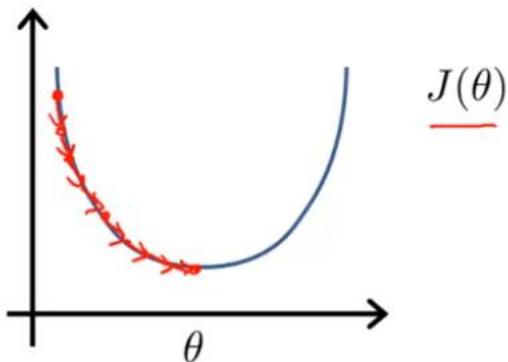
$$\frac{\partial}{\partial \theta_0} J(\theta)$$

$$\begin{aligned} \rightarrow \theta_j &:= \theta_j - \boxed{\alpha} \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \\ &\quad \} \\ \rightarrow \theta_j &:= \underline{\theta_j (1 - \alpha \frac{\lambda}{m})} - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \end{aligned}$$

$\Rightarrow J$ cost function을 편미분하면 $\lambda/m * \theta$ 가 나온다. θ^2 를 미분하면 2가 내려오기 때문이다. 따라서 괄호 안에 $(1 - \alpha * \lambda / m)$ 이라고 쓸 수 있다. $(1 - \alpha * \lambda / m) < 1$ 이다. 거의 0.99.

1.3. Normal equation \leftarrow Supervised machine learning - regression - solution2

Gradient Descent



Normal equation: Method to solve for θ analytically.

→ iterative하지 않게 바로 갈 수 있다.

1.3.1. 방법

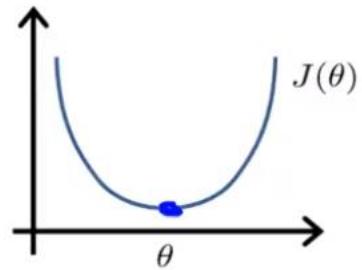
- 직관 : 미분해서 0으로 하면 된다.

Intuition: If 1D ($\theta \in \mathbb{R}$)

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

$$\frac{\partial}{\partial \theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0$$

Solve for θ



$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots \stackrel{\text{set}}{=} 0 \quad (\text{for every } j)$$

Solve for $\theta_0, \theta_1, \dots, \theta_n$

And

- 예시 :

(1) 먼저 data set로 행렬을 만든다. 행렬 x와 벡터 y로 만든다.

Examples: $m = 4$.

x_0	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$ $m \times (n+1)$

$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$ m -dimensional vector

- 행렬 X 를 구하는 방법

m examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$; n features.

$$\underline{x^{(i)}} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad X = \begin{bmatrix} (x^{(1)})^\top \\ (x^{(2)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix} \quad m \times (n+1)$$

(design matrix)

→ m 이 데이터의 개수이고, n 은 변수의 개수이다.

(2) θ 를 구하면 끝난다.

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

- octave에서

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

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

$$\text{Set } A = \underline{X^T X}$$

$$\boxed{(X^T X)^{-1}} = A^{-1}$$

Octave: $\text{pinv}(\boxed{X' * X}) * X' * y$

×

$$\text{pinv}(\underline{X^T * X})$$

- 지금 이 normal equation method에서는 feature scaling을 안 해도 된다.

왜냐하면 표준화의 과정이 포함되어 있기 때문이다.

<굉장히 중요!!!> 왜 $X^T \theta = y$ 가 왜 최소화하는 해를 갖는 것인지!!!

$$- X \theta = y$$

$$\theta = X^{-1} * y$$

이렇게 하면 되지만 역행렬이 없는 경우가 있을 수 있잖아!

$$- X \theta = y$$

$$X^T * X * \theta = X^T * y$$

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

$$\rightarrow X * \theta = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \dots \approx y \quad (y \text{랑 완전히 같아지지는 않는다.!!! 그래도!!!})$$

1.3.2. 장/단점

- n이 크면 gradient descent를 하면 된다. (n=10,000정도 이상이면 크지.. 역행렬 계산이 너무 expensive)
- n이 1000보다 작으면 normal descent가 좋다.
- classification 알고리즘은 normal equation이 작동 안 된다. (정말??)

m training examples, n features.

Gradient Descent

- • Need to choose α .
- • Needs many iterations.
- Works well even when n is large.

$$\nearrow n = 10^6$$



Normal Equation

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

$$\begin{aligned} n &= 100 \\ n &= 1000 \\ \cdots \cdots & n = 10000 \end{aligned}$$

1.3.3. Normal equation non-invertibility

- 만약에 $X^T X$ is non-invertible(singular / degenerate)하면 어떻게 하냐?
→ Octave에서는 $\text{pinv}(X^T * X) * X^T * y$ 해도 잘 된다. (pseudo inverse) → 저번 시간에 했던 것이다! 역행렬이 없는 경우에 바로 역행렬을 구하지
- invertible이 안 되는 경우! $X^T * X$ 가 역행렬이 없는 경우
 - (1) linearly dependent한 경우 : redundant하면 역행렬이 없다.
→ 다중 공선성 문제 : 똑같은 변수가 여러 개 집어 넣으면!! → 비슷한 변수를 빼야 한다. 교차상관성을 다 계산해야 한다. 굉장히 중요하다!!
 - (2) feature가 너무 많은 경우 : 너무 feature(n)가 많으면 안 될 수도... 해결하려면 regularization이 필요하다.

나중에 설명???? 바로 뒤에 나옴. 1.3.4.

→ feature가 sample보다 더 많은 경우, degree of freedom 개념 : 샘플 개수보다 특징 수가 많으면 특징들의 상관관계를 알 수가 없다. 부족하다.

What if $X^T X$ is non-invertible?

- Redundant features (linearly dependent).

E.g. $\underline{x}_1 = \text{size in feet}^2$ $l_m = 3.28 \text{ feet}$

$\underline{x}_2 = \text{size in m}^2$

$$\underline{x}_1 = (3.28)^2 \underline{x}_2$$

$$\rightarrow m = 10$$

$$\rightarrow n = 100$$

$$\Theta \in \mathbb{R}^{101}$$

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

- Delete some features, or use regularization.

cf)

Cross-covariance

$$\gamma_{XY}(\tau) = E[(X_t - \mu_X)(Y_{t+\tau} - \mu_Y)],$$

Cross correlation coefficient

$$\rho_{XY}(\tau) = \frac{1}{\sigma_X \sigma_Y} E[(X_t - \mu_X)(Y_{t+\tau} - \mu_Y)] = \frac{1}{\sigma_X \sigma_Y} \gamma_{XY}(\tau)$$

1.3.4. Normal equation in practice : regularized normal equation

Normal equation

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

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

$$\rightarrow \min_{\theta} J(\theta) \quad \rightarrow \quad \rightarrow \Theta = (X^T X + \lambda \begin{bmatrix} 0 & 1 & \dots & 1 \end{bmatrix})^{-1} X^T y$$

→ derive가 가능하다.

- 일석이조!! : regularization을 하면 non-invertible한 $X^T X$ 도 invertible해진다. (증명 가능)

Non-invertibility (optional/advanced).

Suppose $m \leq n$, $\begin{matrix} \leftarrow \\ (\# \text{examples}) \quad (\# \text{features}) \end{matrix}$

$$\theta = \underbrace{(X^T X)^{-1}}_{\text{non-invertible / singular}} X^T y \quad \begin{matrix} \text{pinv} \\ \underline{\quad} \end{matrix} \quad \begin{matrix} \text{inv} \\ \pi \end{matrix}$$

If $\lambda > 0$,

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

invertible.

→ $m < n$ 하면 해가 정해지지 않기 때문에 non-invertible하다. positive definite matrix (Determinant 0이 무조건 0보다 커진다)

Supervised learning 종류 : regression problem / classification problem

Week 3.

1.4. regression problem : 위 참조

1.5. Classification

1.5. Classification → logistic regression

1. Why logistic regression?

- 목적 : logistic regression을 위해서!! discrete value를 가졌을 경우! continuous가 아닌 경우
- 종류 : two class, multiclass classification

Classification

→ Email: Spam / Not Spam?

→ Online Transactions: Fraudulent (Yes / No)?

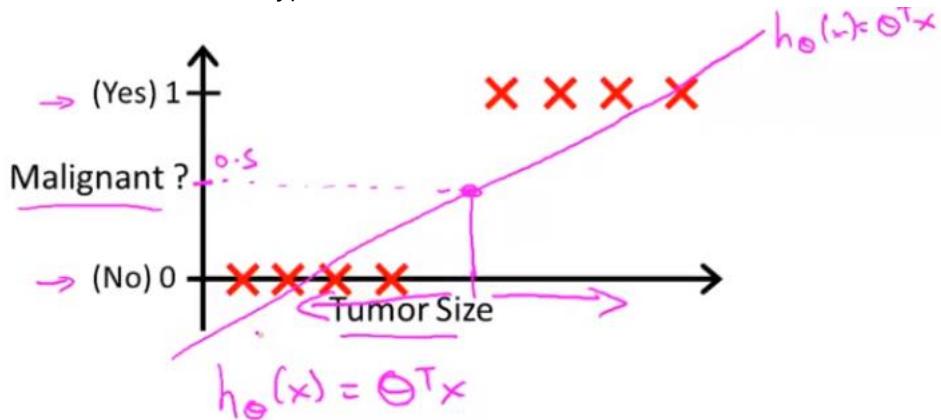
→ Tumor: Malignant / Benign ?

→ $y \in \{0, 1\}$ 0: "Negative Class" (e.g., benign tumor)
1: "Positive Class" (e.g., malignant tumor)

→ $y \in \{0, 1, 2, 3\}$

- 기존 linear regression의 한계 : linear regression으로 classification을 하는 경우의 문제점

(1) 이상값이 존재하게 되면 맞지 않는 hypothesis, threshold를 갖게 되는 문제점

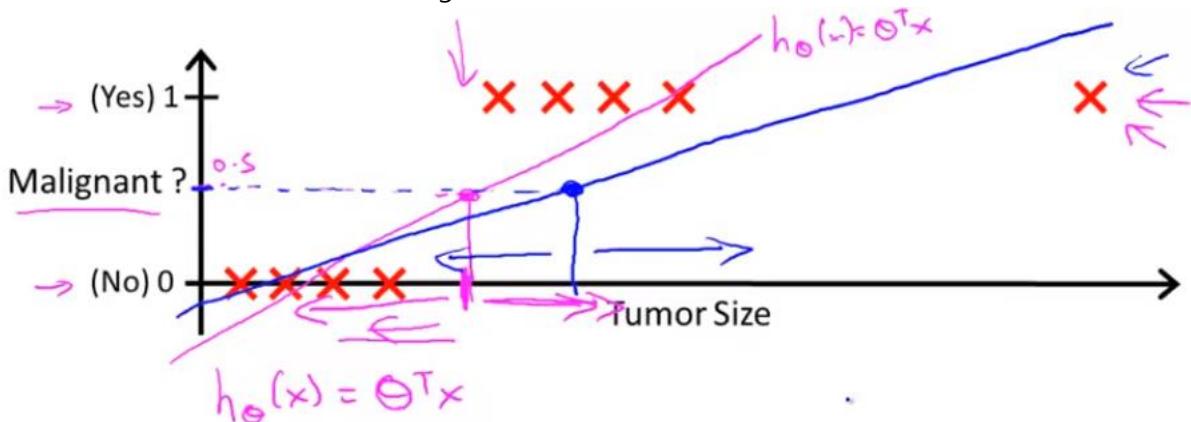


→ Threshold classifier output $h_\theta(x)$ at 0.5:

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

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

→ 그런데 아웃라이어가 있게 되면, linear regression이 부정확한 threshold를 갖게 한다.



→ 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"

(2) hypothesis $h_\theta(x)$ 가 1보다 크거나, 0보다 작은 경우가 생긴다. 하지만 classification에서는 $y=0, 1$ 이 전부이다.

- 문제점 해결 : **logistic regression**을 하게 되면 $0 \leq h_\theta(x) \leq 1$ 이 된다. (regression이라고 칭하기는 하지만, classification model이니까 헷갈려 하지 말 것.)

2. What is Logistic regression model

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

→ 위 아래로 0, 1에 수렴하기 때문에 0과 1을 벗어날 리 없다.

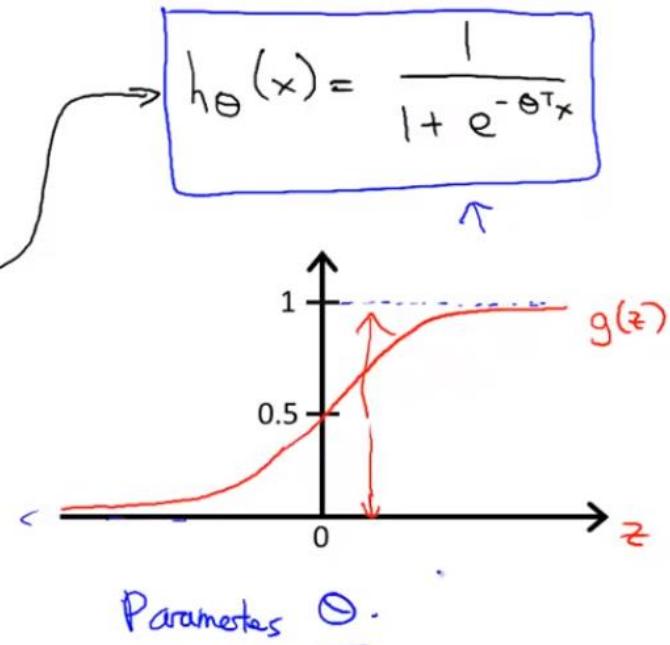
Logistic Regression Model

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

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

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

Sigmoid function
Logistic function



(2) 의미 : $h_\theta(x)$ = estimated probability that $y=1$ on input x

ex) 환자 1명의 tumor size x 가 주어졌을 때, 악성 종양($y=1$)일 가능성이 몇 퍼센트이냐?

$h_\theta(x) = P(y=1 | x; \theta)$ ← probability that $y=1$, given x , parameterized by θ

Interpretation of Hypothesis Output

$$h_\theta(x)$$

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

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

$$h_\theta(x) = 0.7 \quad y=1$$

Tell patient that 70% chance of tumor being malignant

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

"probability that $y=1$, given x , parameterized by θ "

$$y = 0 \text{ or } 1$$

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

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

(3) Decision boundary

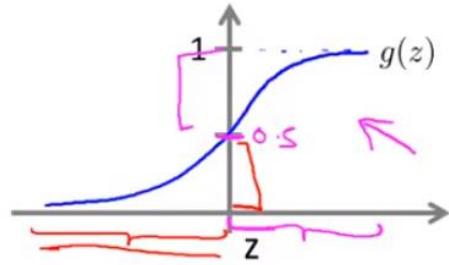
- 0.5가 기준이 될 수 있다. $h_\theta(x) \geq 0.5$ 이면 $y=1$ 로 predict하고, $h_\theta(x) < 0.5$ 이면 $y=0$ 로 predict하면 된다.
다시 말하면, $h_\theta(x) = g(z) = g(\theta^T x) = 0.5$ 가 되는 $\theta^T x$ 가 0이 기준이다. $\theta^T x$ 가 0보다 크면 $y=1$ 로, $\theta^T x$ 가 0보다 작으면 $y=0$ 로 prediction하면 된다.

Logistic regression

$$\rightarrow h_{\theta}(x) = g(\theta^T x) = P(y=1|x; \theta)$$

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

Suppose predict " $y = 1$ " if $h_{\theta}(x) \geq 0.5$
 $\theta^T x \geq 0$



$$g(z) \geq 0.5 \quad \text{when } z \geq 0$$

$$h_{\theta}(x) = g(\theta^T x) \geq 0.5 \quad \text{when } \theta^T x \geq 0$$

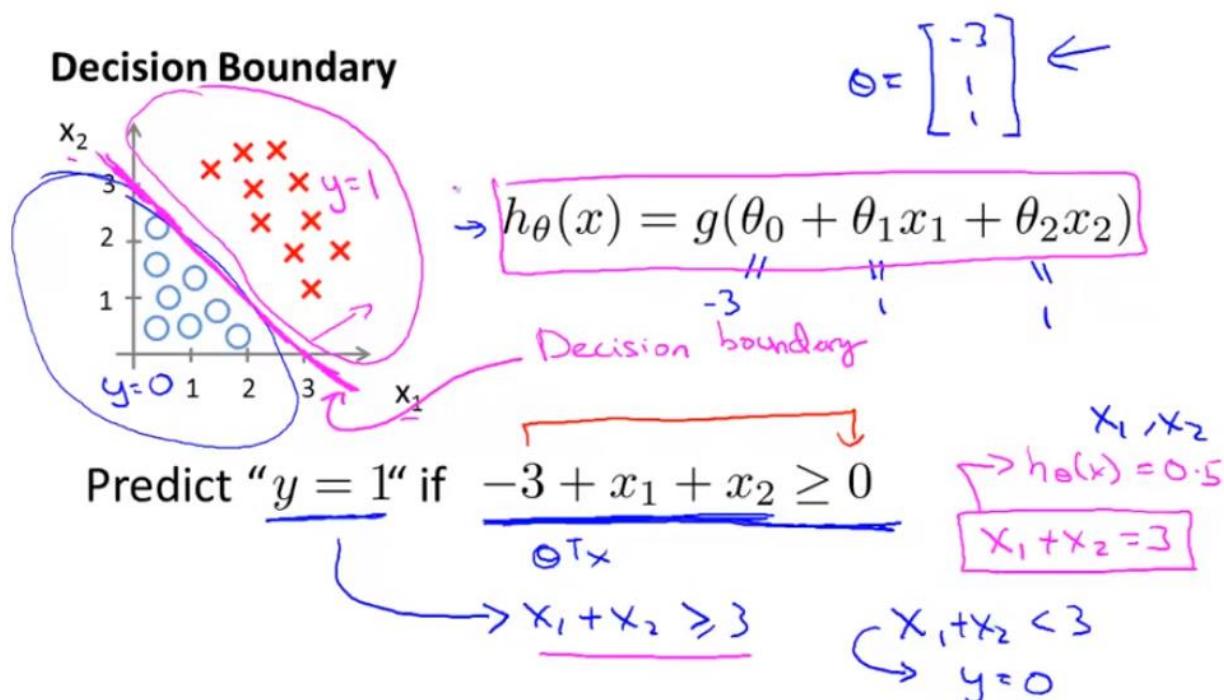
predict " $y = 0$ " if $h_{\theta}(x) < 0.5$
 $h_{\theta}(x) = g(\theta^T x)$
 $\theta^T x < 0$

So if our input to g is $\theta^T X$, then that means:

$$h_{\theta}(x) = g(\theta^T x) \geq 0.5$$

when $\theta^T x \geq 0$

- 예를 들면, $\theta^T x$ 가 ≥ 0 이면, 즉, $g(z) = h_{\theta}(x)$ 가 0.5보다 크면 $y=1$ 이라고 추정할 수 있다. 이 경계를 decision boundary라고 하며, 이는 data set에서 파생되는 속성이 아니라, $h_{\theta}(x)$, 즉 hypothesis에서 알 수 있는 속성이다.

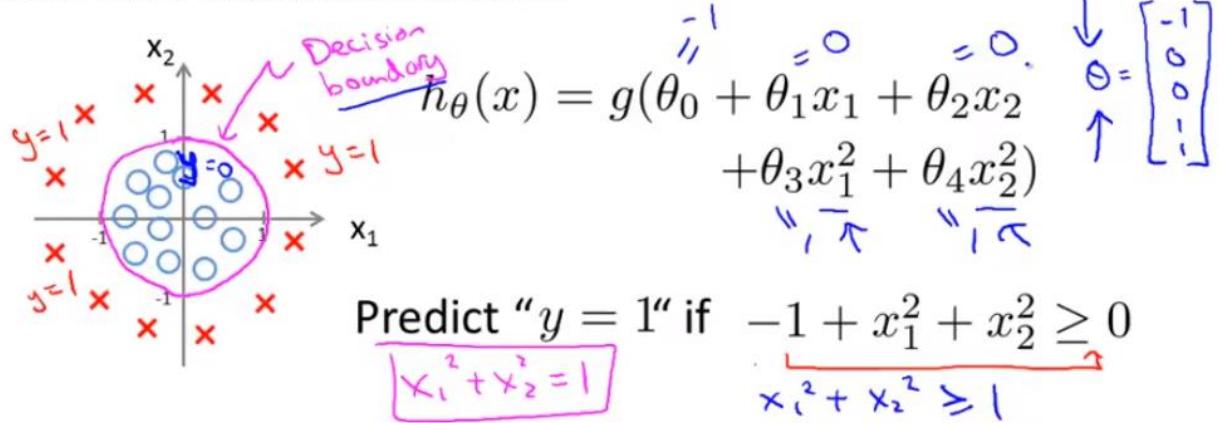


- non-linear decision boundaries

- polynomial term을 더하면 이것도 표현할 수 있다.

- 원으로도 표현 가능

Non-linear decision boundaries



3. How to fit the parameter

- 목표 : cost function 을 가장 minimize 시키자!
- 문제 정의 :

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

m examples

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \subset \mathbb{R}^{n+1} \quad x_0 = 1, y \in \{0, 1\}$$

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

How to choose parameters θ ?

- linear regression 때처럼, $h(x) - y$ 의 제곱의 합으로 할 때의 문제점 \rightarrow non convex 해지기 때문에 gradient descent 를 하면 하나의 global minimum 에 수렴하지 않는다. 따라서 cost function이 convex 해지는 방법을 찾아야 한다.

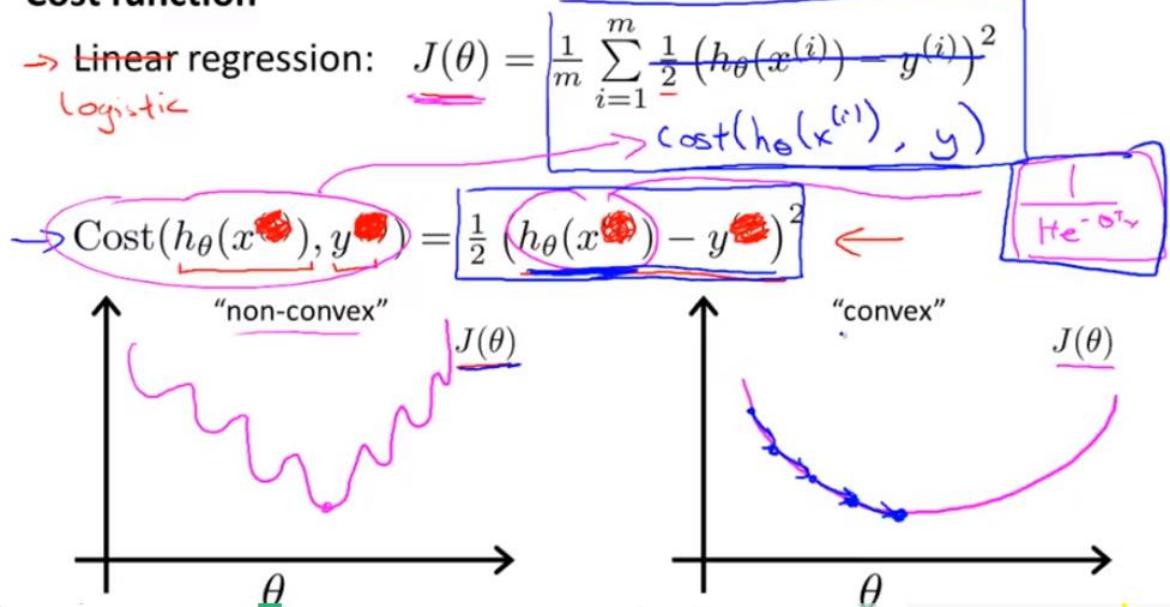
Cost function

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

cost($h_\theta(x^{(i)})$, y)

$$\text{Cost}(h_\theta(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

Cost function



질문!!!

이거 왜 non convex해지나??

- 해결 : log함수를 사용해서

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

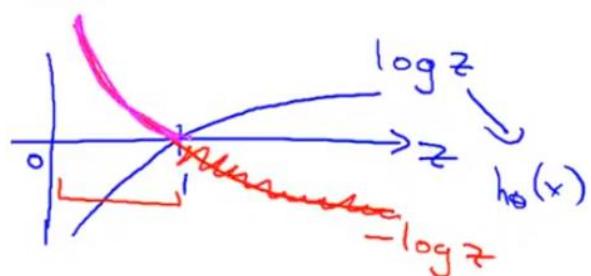
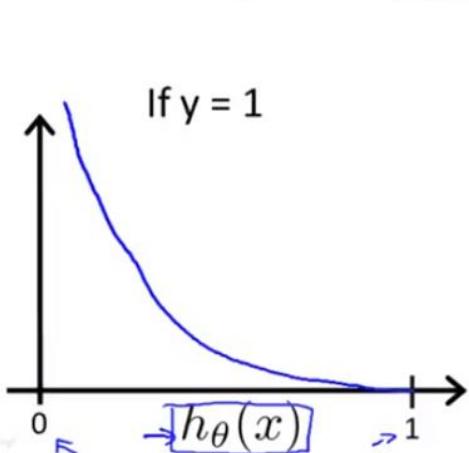
$$\text{Cost}(h_\theta(x), y) = -\log(h_\theta(x)) \quad \text{if } y = 1$$

$$\text{Cost}(h_\theta(x), y) = -\log(1 - h_\theta(x)) \quad \text{if } y = 0$$

(1) $y = 1$ 일 때,

Logistic regression cost function

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



→ 직관적인 의미 : $y=1$ 임에도 불구하고, $h_\theta(x)$ 라고 대답을 해버리면(tumor가 malignant일 확률이 0이라고) cost가 굉장히 크게 되게끔 설정한 것이다.

Andrew I

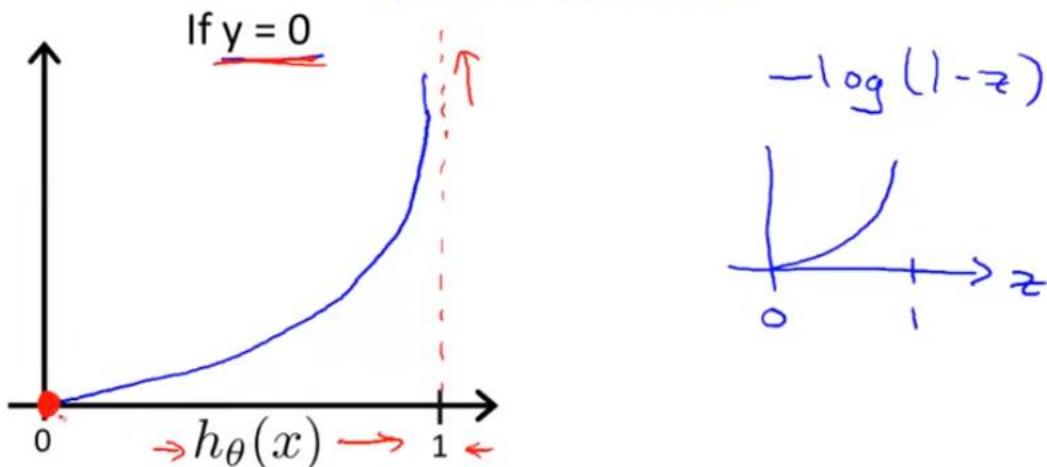
\rightarrow Cost = 0 if $y = 1, h_\theta(x) = 1$
 But as $h_\theta(x) \rightarrow 0$
 $Cost \rightarrow \infty$

\rightarrow Captures intuition that if $h_\theta(x) = 0$,
 (predict $P(y = 1|x; \theta) = 0$), but $y = 1$,
 we'll penalize learning algorithm by a very
 large cost.

(2) $y = 0$ 일 때

Logistic regression cost function

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



$$\text{Cost}(h_\theta(x), y) = 0 \text{ if } h_\theta(x) = y$$

$$\text{Cost}(h_\theta(x), y) \rightarrow \infty \text{ if } y = 0 \text{ and } h_\theta(x) \rightarrow 1$$

$$\text{Cost}(h_\theta(x), y) \rightarrow \infty \text{ if } y = 1 \text{ and } h_\theta(x) \rightarrow 0$$

- 표기 : simplified cost function

\rightarrow 두 줄로 쓴 식을 1줄로 쓸 수 있다.

Logistic regression cost function

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

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

Note: $y = 0$ or 1 always

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

Logistic regression cost function

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

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

→ 왜 이것을 사용하나? maximum likelihood estimation으로도 유도할 수 있다.

To fit parameters θ :

$$\min_{\theta} J(\theta) \quad \text{Get } \underline{\theta}$$

To make a prediction given new x :

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

$p(y=1 | x; \theta)$

- fitting 방법 : gradient descent

- 방법 :

이거 어떻게 저렇게 예쁘게 나오는 거지??

Gradient Descent

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

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

Repeat {

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

}

(simultaneously update all θ_j)

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

→ 유도해봤는데, 매우 비슷하게 나온다. 완전히 성공하지는 못함. 부호가 어딘가 (-)가 들어옴

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

Repeat {

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

(simultaneously update all θ_j)

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

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

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

- α 를 맞게 정하는 방법 :

Plot

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

as a function of the number of iterations and

make sure $J(\theta)$ is decreasing on every iteration.

→ J 가 작아졌다가 커지거나 하면 alpha를 다시 설정해야 한다.

- Vectorization

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

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

⋮

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

We would like a vectorized implementation of the form
 $\theta := \theta - \alpha \delta$ (for some vector $\delta \in \mathbb{R}^{n+1}$).

What should the vectorized implementation be?

$$\Theta \quad \theta := \theta - \alpha \frac{1}{m} \sum_{i=1}^m [(h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}]$$

- feature scaling도 할 수 있다.

- 나중에 7.2.1.에서 mini-batch gradient descent랑 stochastic gradient descent를 배우게 된다.(computational cost를 줄여주기 위한 것이다.)

4. Advanced Optimization for linear / logistic regression

(1) 방법의 종류

1) gradient descent : 위에서 했다.

Optimization algorithm

Cost function $J(\theta)$. Want $\min_{\theta} J(\theta)$.

Given θ , we have code that can compute

$$\begin{aligned}\rightarrow & - J(\theta) \\ \rightarrow & - \frac{\partial}{\partial \theta_j} J(\theta) \quad (\text{for } j = 0, 1, \dots, n)\end{aligned}$$

Gradient descent:

Repeat {

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

}

- Advanced methods

2) conjugate gradient

3) BFGS

4) L-BFGS

- 장점 : no need to manually pick α (learning rate)

(2) Octave에서의 방법 :

(2-1) linear regression에서 어떻게 하는지 한번 해보자.

Example:

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

$$\rightarrow \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \quad \theta_1 = 5, \theta_2 = 5.$$

$$\rightarrow J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\rightarrow \frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$

$$\rightarrow \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

```
function [jVal, gradient] = costFunction(theta)
    jVal = (theta(1)-5)^2 + ...
            (theta(2)-5)^2;
    gradient = zeros(2,1);
    gradient(1) = 2*(theta(1)-5);
    gradient(2) = 2*(theta(2)-5);
```

```
options = optimset('GradObj', 'on', 'MaxIter', 100);
initialTheta = zeros(2,1);
[optTheta, functionVal, exitFlag] ... = fminunc(@costFunction, initialTheta, options);
```

\rightarrow options = optimset('GradObj', 'on', 'MaxIter', 100) \rightarrow 100번 iteration을 하라는 거구나.. 이거 혹시 library인가?

- cost function \leftarrow 따로 함수 설정

A screenshot of Microsoft WordPad showing a MATLAB script named costFunction.m. The script calculates the cost function and its gradient for a given theta vector. A yellow circle highlights the variable `jVal` in the first line of code.

```
function [jVal, gradient] = costFunction(theta)
jVal = (theta(1)-5)^2 + (theta(2)-5)^2;
gradient = zeros(2,1);
gradient(1) = 2*(theta(1)-5);
gradient(2) = 2*(theta(2)-5);
```

A screenshot of an Octave terminal window showing the execution of the costFunction.m script. The user runs PS1('>>'), changes the directory to 'C:\Users\ang\Desktop', sets optimization options, initializes theta to zeros, and then calls fminunc. The output shows the optimized theta values and the function value.

```
octave-3.2.4.exe:1> PS1('>> ')
>> cd 'C:\Users\ang\Desktop'
>>
>> options = optimset('GradObj','on', 'MaxIter', '100');
>> initialTheta = zeros(2,1)
initialTheta =
    0
    0
>> [optTheta, functionVal, exitFlag] = fminunc(@costFunction, initial
```

→ 이렇게 치면 결과가 빨리 나온다.

A screenshot of an Octave terminal window showing the results of the optimization. The user runs fminunc with the costFunction script, initial theta values, and optimization options. The output shows the optimized theta values (5.0000, 5.0000), the function value (1.5777e-030), and the exit flag (1).

```
<ag> = fminunc(@costFunction, initialTheta, options)
optTheta =
    5.0000
    5.0000
functionVal = 1.5777e-030
exitFlag = 1
>>
```

→ Theta가 5, 5가 optimum이구나.

(2-2) logistic regression을 하는 방법 (Gradient descent가 아닌 advanced algorithm)

$$\underline{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \begin{array}{l} \text{theta(1)} \\ \text{theta(2)} \\ \vdots \\ \text{theta(n+1)} \end{array}$$

```
function [jVal, gradient] = costFunction(theta)
    jVal = [ code to compute  $J(\theta)$  ];
    gradient(1) = [ code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$  ];
    gradient(2) = [ code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$  ];
    :
    gradient(n+1) = [ code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$  ];
```

→ indexing이 1부터 된다는 것을 조심.

- 정리 :

We first need to provide a function that evaluates the following two functions for a given input value θ :

$$J(\theta) \\ \frac{\partial}{\partial \theta_j} J(\theta)$$

We can write a single function that returns both of these:

```
1 function [jVal, gradient] = costFunction(theta)
2     jVal = [...code to compute  $J(\theta)$ ...];
3     gradient = [...code to compute derivative of  $J(\theta)$ ...];
4 end
```

Then we can use octave's "fminunc()" optimization algorithm along with the "optimset()" function that creates an object containing the options we want to send to "fminunc()". (Note: the value for MaxIter should be an integer, not a character string - errata in the video at 7:30)

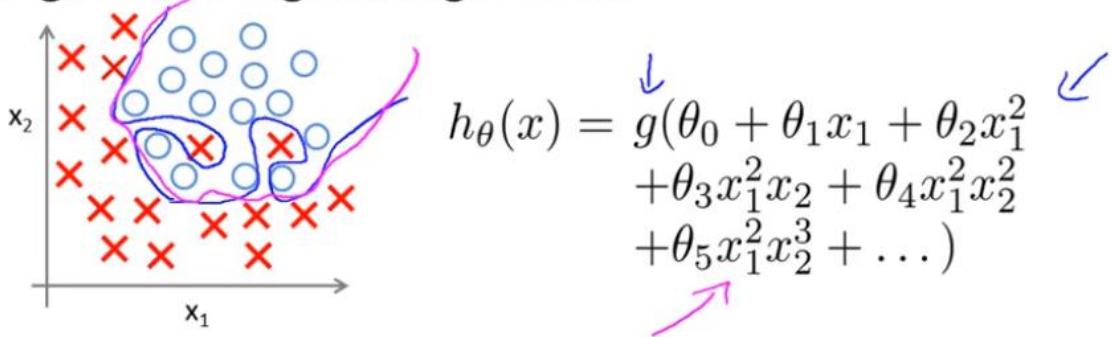
```
1 options = optimset('GradObj', 'on', 'MaxIter', 100);
2 initialTheta = zeros(2,1);
3 [optTheta, functionVal, exitFlag] = fminunc(@costFunction, initialTheta,
    options);
```

We give to the function "fminunc()" our cost function, our initial vector of theta values, and the "options" object that we created beforehand.

5. Application of logistic regression-1 : Regularized logistic regression

- 원리 : θ 의 합을 도입하여 overfitting 문제를 완화시킨다.

Regularized logistic regression.



Cost function:

$$\rightarrow J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad | \boxed{\theta_1, \theta_2, \dots, \theta_n}$$

- 알고리즘 : hypothesis가 다르지만 외관상 linear regression과 비슷하다. 주의 : $j=0$ 을 빼고 하면 된다.

Gradient descent

Repeat {

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

$$\rightarrow \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad | \underbrace{(j=1, 2, 3, \dots, n)}_{\theta_1, \dots, \theta_n}$$

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

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

- regularized gradient descent가 잘 되는지 확인하는 방법

- Plot $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))\right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$ as a function of the number of iterations and make sure it's decreasing.

Correct

→ j cost function이 감소하는지 봐야지.

- Advanced optimization에서 해보기

Advanced optimization

$\underset{\theta}{\text{fminunc}} \text{(@costFunction)}$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \begin{array}{l} \theta_0 \leftarrow \text{theta}(1) \\ \theta_1 \leftarrow \text{theta}(2) \\ \vdots \\ \theta_n \leftarrow \text{theta}(n+1) \end{array}$$

```

function [jVal, gradient] = costFunction(theta)
    jVal = [ code to compute  $J(\theta)$  ];
     $\Rightarrow J(\theta) = \left[ -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$ 
     $\Rightarrow \text{gradient}(1) = [\text{code to compute } \frac{\partial}{\partial \theta_0} J(\theta)];$ 
     $\Rightarrow \text{gradient}(2) = [\text{code to compute } \frac{\partial}{\partial \theta_1} J(\theta)];$ 
     $\vdots$ 
     $\text{gradient}(n+1) = [\text{code to compute } \frac{\partial}{\partial \theta_n} J(\theta)];$ 

```

$\rightarrow \text{fminunc} (\text{@costFunction})$ 하면 j cost function이 minimize된다.

f minimization unconstrained

```

1 options = optimset('GradObj', 'on', 'MaxIter', 100);
2 initialTheta = zeros(2,1);
3 [optTheta, functionVal, exitFlag] = fminunc(@costFunction, initialTheta,
    options);
4

```

We give to the function "fminunc()" our cost function, our initial vector of theta values, and the "options" object that we created beforehand.

6. Application of logistic regression-2 : Multiclass classification in logistic regression

- 예시 :

Multiclass classification

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

$y=1$ $y=2$ $y=3$ $y=4$

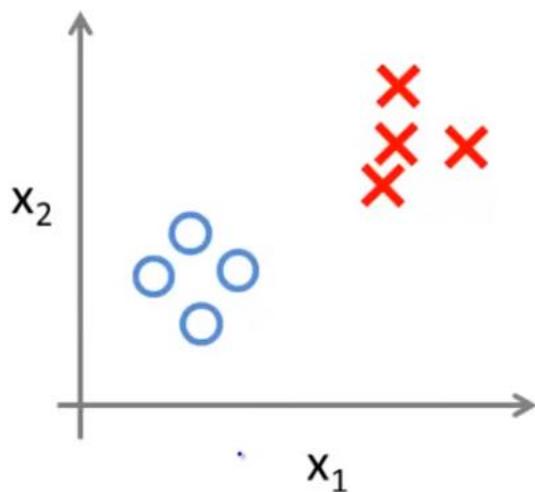
Medical diagrams: Not ill, Cold, Flu

$y=1$ 2 3

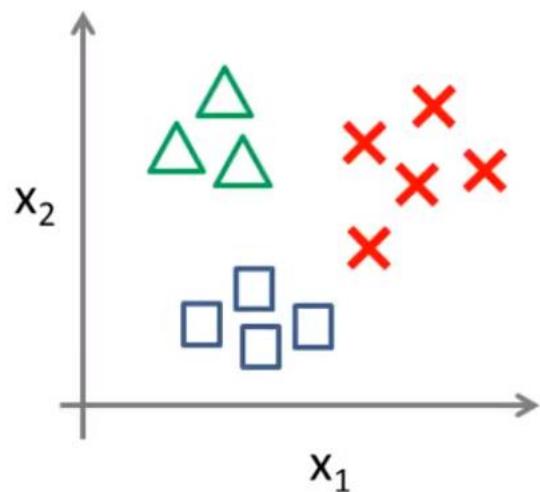
Weather: Sunny, Cloudy, Rain, Snow

$y=1$ 2 3 4

Binary classification:

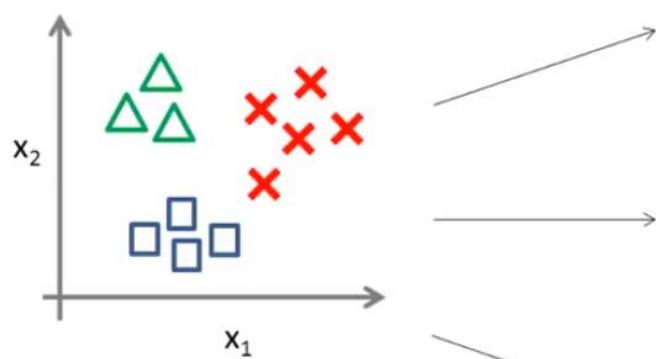


Multi-class classification:



- one versus all classification (one-vs-rest)
- hypothesis를 class 개수 만큼 만들면 된다.

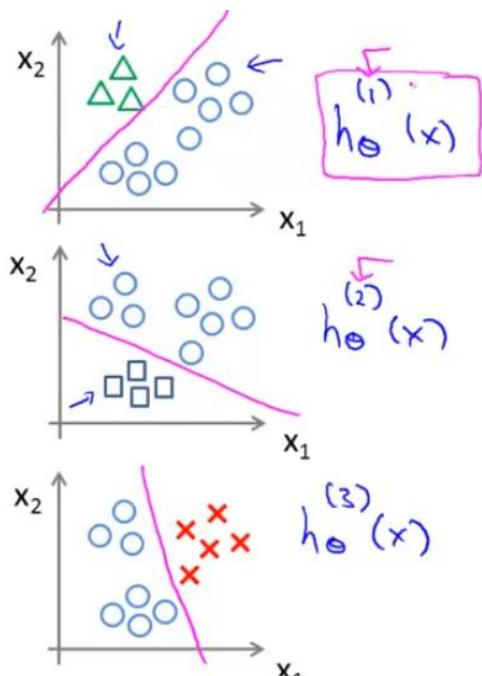
One-vs-all (one-vs-rest):



Class 1: \triangle ←
 Class 2: \square ←
 Class 3: \times ←

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

→ 각각의 hypothesis는 세모를, 네모를, x를 y=1로 취급하는 것이다.



- 어떤 x가 어떤 class에 속하는지 보기 위해서는 $h(x)$ 가 가장 크게 (1에 가깝게 만드는) hypothesis를 선택하면 된다. 그 hypothesis를 사용한 class를 선택하면 되는 것이다.

One-vs-all

Train a logistic regression classifier $h_{\theta}^{(i)}(x)$ for each class i to predict the probability that $y = i$.

On a new input x , to make a prediction, pick the class i that maximizes

$$\max_i h_{\theta}^{(i)}(x)$$

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

$$h_{\theta}^{(0)}(x) = P(y = 0|x; \theta)$$

$$h_{\theta}^{(1)}(x) = P(y = 1|x; \theta)$$

...

$$h_{\theta}^{(n)}(x) = P(y = n|x; \theta)$$

$$\text{prediction} = \max_i(h_{\theta}^{(i)}(x))$$

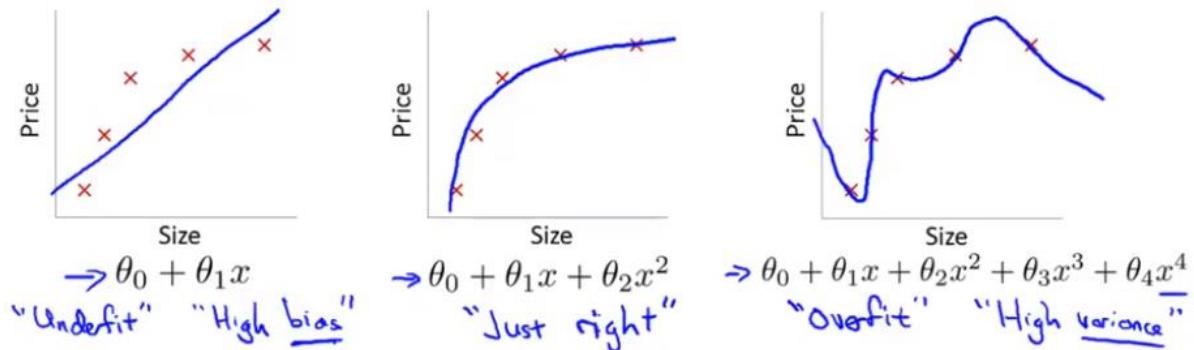
7. Application-3 : The problem of overfitting

7.1. 문제 :

예시1) linear regression

- linear regression의 경우에 \rightarrow under-fitting = high bias
- quadratic \rightarrow 적당함
- high order polynomial \rightarrow overfitting = high variance

Example: Linear regression (housing prices)



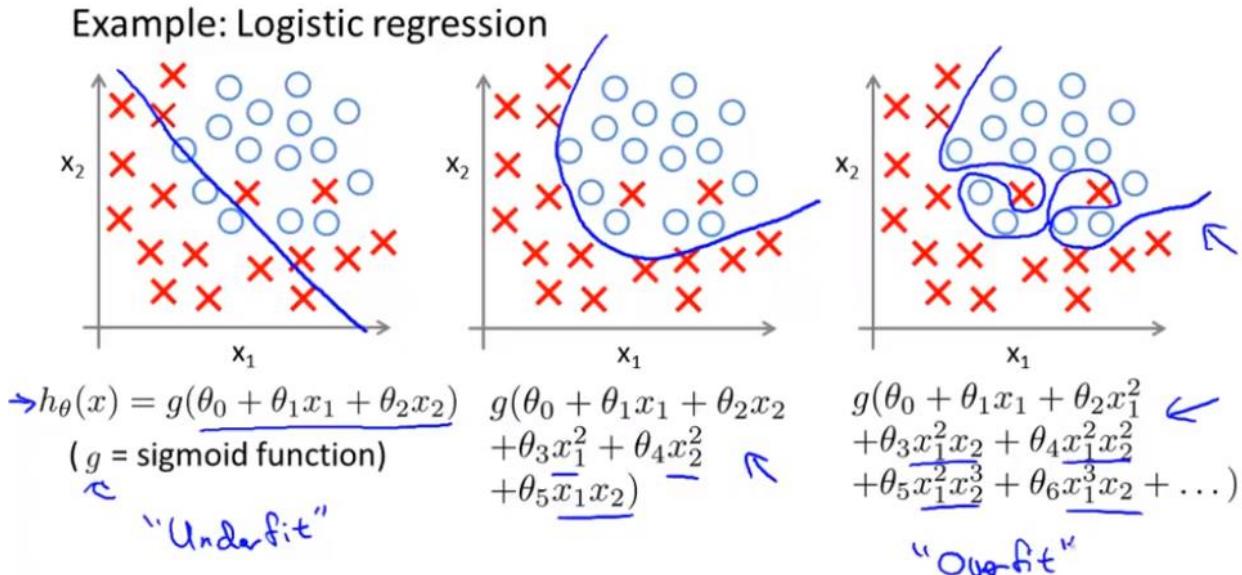
Overfitting: If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples).

\rightarrow 이처럼 hypothesis가 data를 지나치게 fit할 때, 새로운 값들을 예측하는 게 어려울 수 있다. 그것이 overfitting

문제이다.

- 이 문제처럼 물론 degree of polynomial이 커질 때도(3번째 경우) 문제이지만, feature가 많아지는 경우에도 만약에 data set의 수가 적으면 문제가 발생할 수 있다. 구불구불한 것이 많아진다.

예시2) logistic regression



7.2. 해결 : regularization

첫번째는 feature를 다 쓰지 않고 빼는 방법이 있고, 두번째는 feature는 다 쓰면서 θ 를 줄이는 방법(regularization)이 있다.

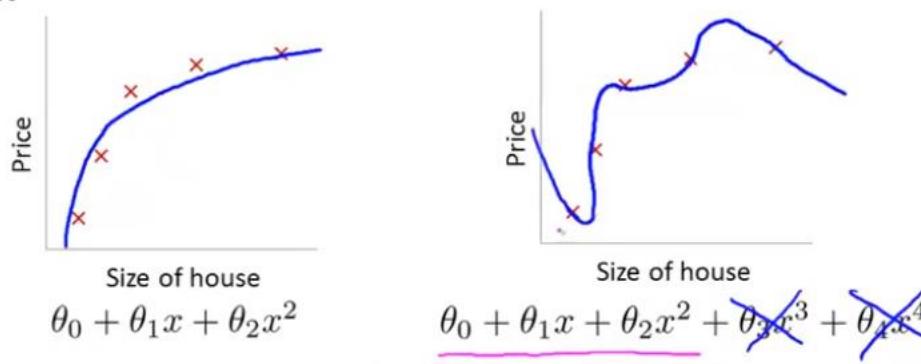
1. Reduce number of features.
 - — Manually select which features to keep.
 - — Model selection algorithm (later in course).
2. Regularization.
 - — Keep all the features, but reduce magnitude/values of parameters θ_j .
 - Works well when we have a lot of features, each of which contributes a bit to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

(1) Reduce the number of features (with algorithm)

(2) Regularization

- 예시 : θ_3 와 θ_4 를 줄임으로써 x^3, x^4 항을 없애는 효과를 낼 수 있다.

Intuition



Suppose we penalize and make θ_3, θ_4 really small.

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

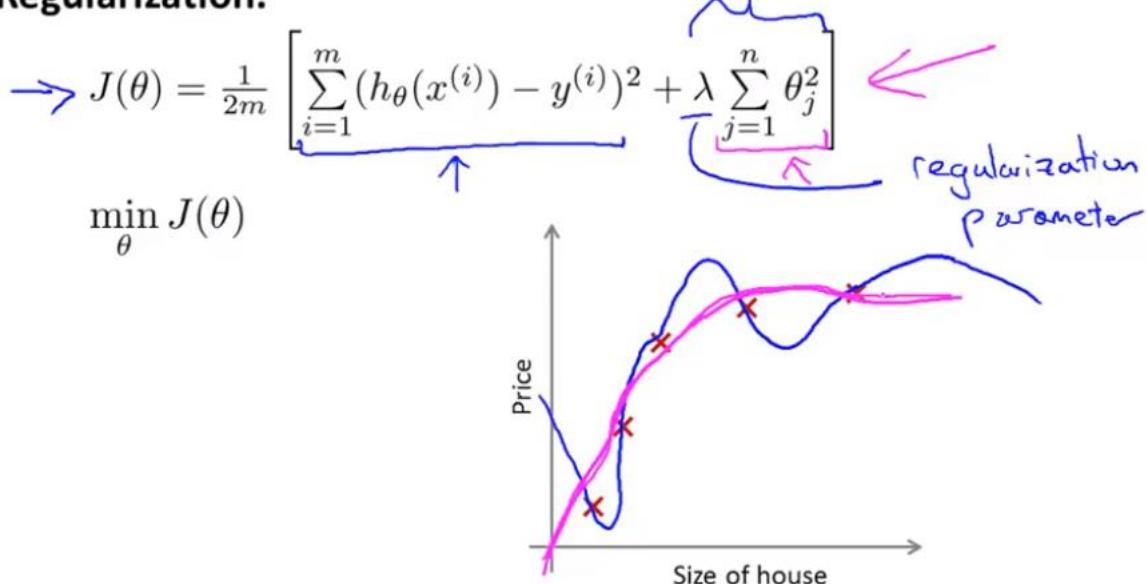
$\theta_3 \approx 0 \quad \theta_4 \approx 0$

→ small values for parameters (θ) leads to "simpler" "Smoothen" hypothesis, less prone to overfitting

질문!!! 왜 저거를 굳이 더해버리는 거지? → minimize하려고

- 원리 :

Regularization.



→ θ 를 모두 더한다. λ 는 regularization parameter이고 θ 의 합에 곱해서 regularization objective에 곱하는 것이다. 이 때, J 가 두 항의 합으로 나타내지게 되는데, 앞의 항은 hypothesis랑 y 의 차이를 줄여주는 효과를 내고, 뒤의 항은 parameter를 줄이는 효과를 내므로, 그것의 균형을 찾을 수 있게 된다.

왜 θ 를 줄이면 이와 같이 굴곡이 사라지는가에 대한 직관은 직접 해봐야지 얻을 수 있다.

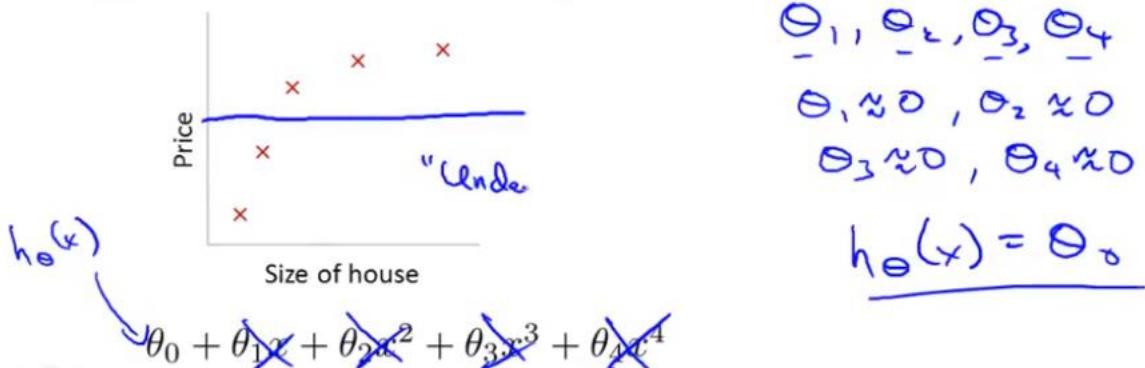
- λ (regularization parameter)의 조절 : λ 를 높인다는 것은 θ 가 높아지는 것에 대한 엄청난 패널티를 가한다는 것이다. 따라서 λ 를 높이게 되면 θ 가 0으로 가게 되고, 그렇게 되면 지나치게 parameter가 축소되어 모델이 단순화되는 문제점이 발생한다.

만약에 $\lambda = 0$ 이 되게 되면 반대로, 다시 overfitting의 문제가 살아날 것이다.

In regularized linear regression, we choose θ 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 λ is set to an extremely large value (perhaps far too large for our problem, say $\lambda = 10^{10}$)?



- 미래 : non-linear classifiers

1.6. Neural Networks

1.6.1. Neural network를 언제 쓰는가.

- Neural network가 필요한 이유 : logistic regression 의 non-linear classification으로는 부족하기 때문이다

(1) feature가 너무 많아지고, overfitting이 일어날 수 있고

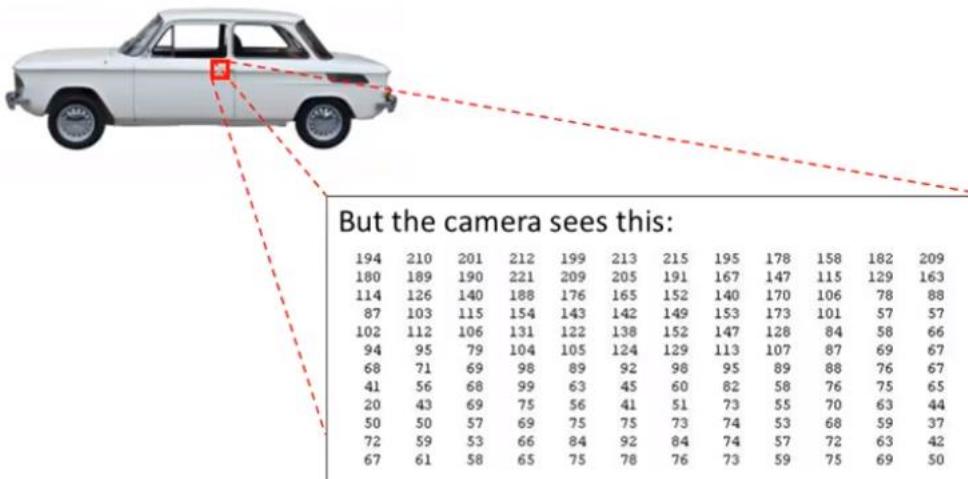
(2) computationally expensive하다.

→ 더 fit가 잘 되는 것을 찾다보면 x_1^2, x_2^2, x_3^2 , 혹은 $x_1^3, x_1^2 \cdot x_2$ 등을 하다 보면 feature가 엄청 많아진다.

➔ why do we need neural network? : non-linear model을 더 잘 represent할 수 있기 때문이다.

- 예시 : Computer vision

You see this:



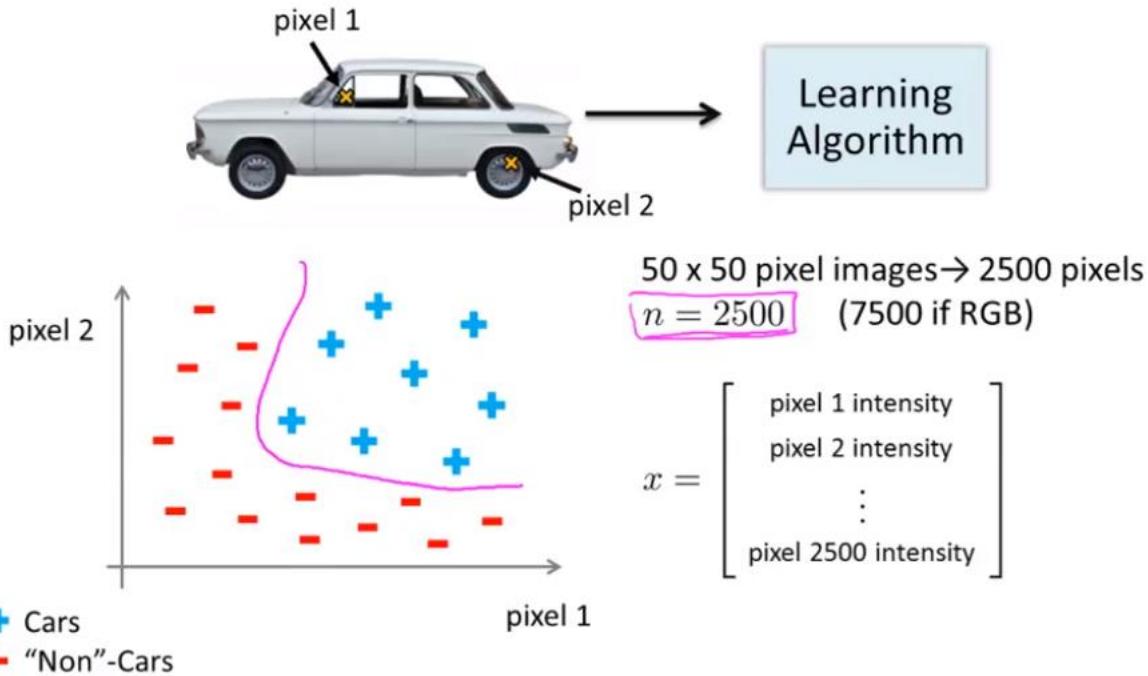
Computer Vision: Car detection



Testing:



- non-linear classification으로 해보려는 경우!



→ n이 크게 되면 변수가 너무 많아진다. 변수 2개를 곱한 형태는 $n(n-1)/2$ 개이다.

- Neural network

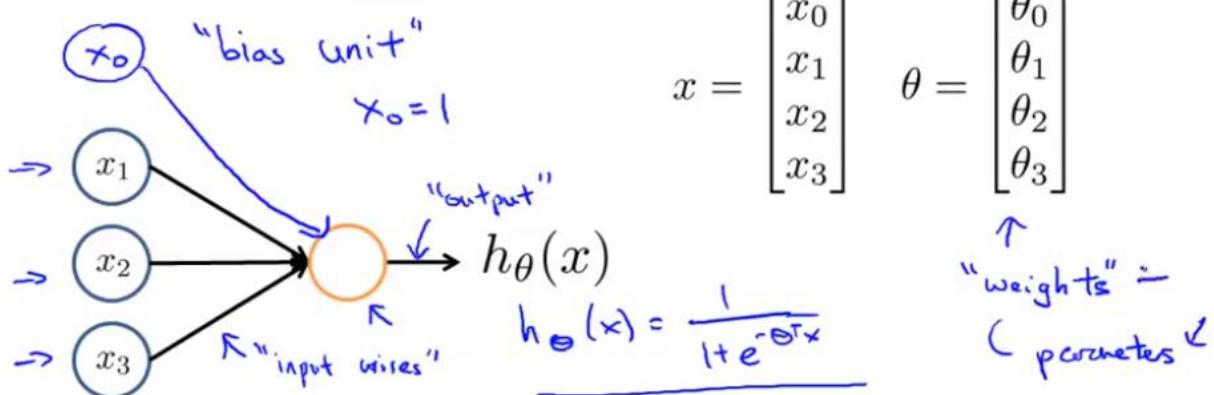
- mimicking the brain : 뉴런을 연결하면 뇌가 data를 배워서 어떻게 해석해야 하는지 알게 된다.

1.6.2. Neural network의 구조와 알고리즘

- Model representation :

- input(dendrite) → output through Exon (activation의 과정)

Neuron model: Logistic unit



Sigmoid (logistic) activation function.

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

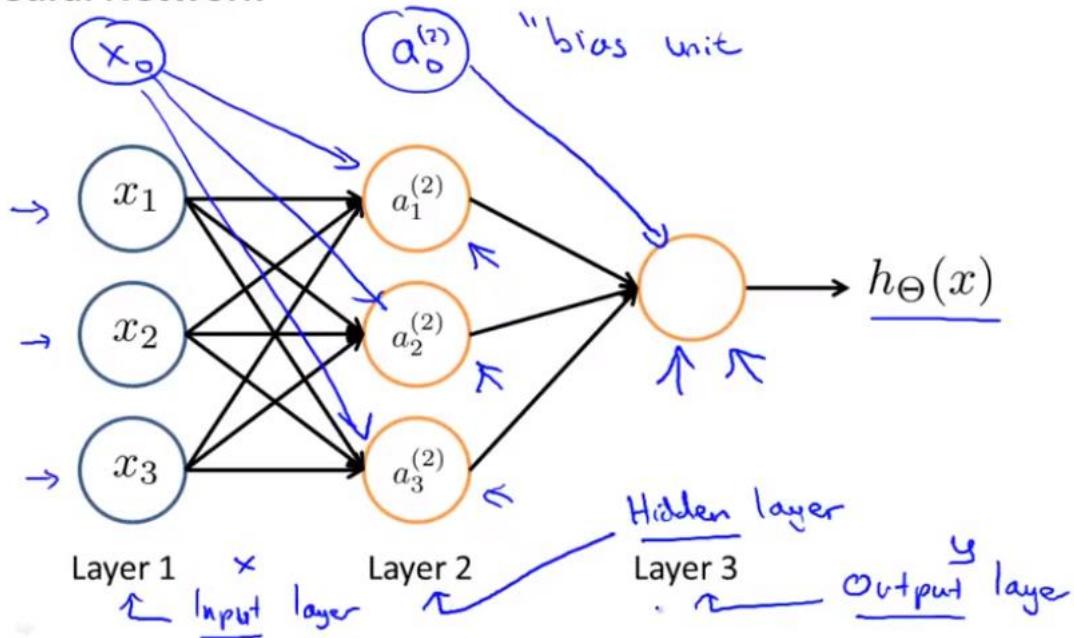
→ parameter가 weight라고 불린다.

$a_i^{(j)}$ = "activation" of unit i in layer j

$\Theta^{(j)}$ = matrix of weights controlling function mapping from layer j to layer $j+1$

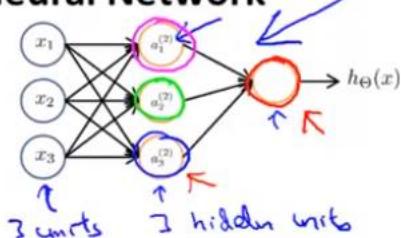
- layer의 개수가 더 많아질 수 있다. : input과 output이 아닌 layer는 hidden layer라고 부른다. 1개 이상일 수도 있다. In this model our x_0 input node is sometimes called the "bias unit." It is always equal to 1.

Neural Network



- Neural network의 알고리즘

Neural Network



$\rightarrow a_i^{(j)}$ = "activation" of unit i in layer j

$\rightarrow \Theta^{(j)}$ = matrix of weights controlling function mapping from layer j to layer $j+1$

$$h_{\Theta}(x)$$

$$\rightarrow a_1^{(2)} = g(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3)$$

$$\rightarrow a_2^{(2)} = g(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3)$$

$$\rightarrow a_3^{(2)} = g(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3)$$

$$\rightarrow h_{\Theta}(x) = a_1^{(3)} = g(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)})$$

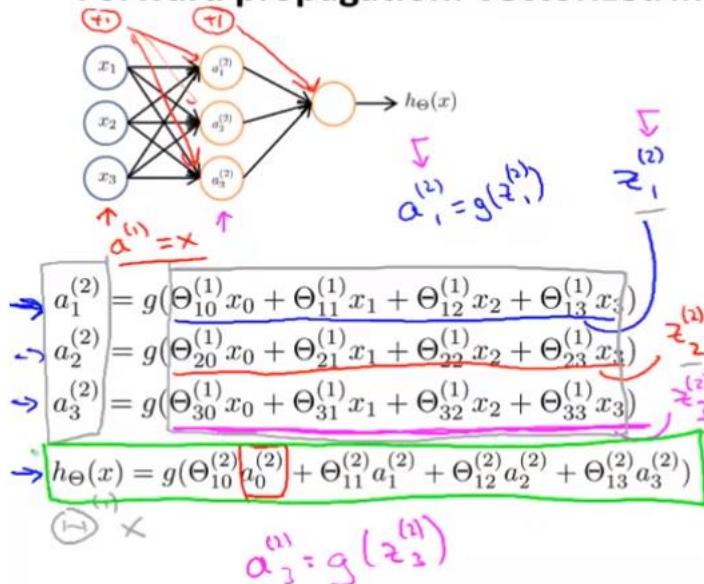
\Rightarrow If network has s_j units in layer j , s_{j+1} units in layer $j+1$, then $\Theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$.

$$s_{j+1} \times (s_j + 1)$$

$\rightarrow \Theta$ 행렬이 x 를 a 로 만들어준다. Θ 행렬의 dimension을 보게 되면, 일단 행의 개수는 다음 단계($j+1$ 단계)의 변수의 개수여야 한다. 열의 개수는 그 전 단계(j 단계) + 1이어야 한다. 왜냐하면 bias unit $\rightarrow 1$ 이 있기 때문이다.

- Vectorized implementation : vector로 만들어보자!!

Forward propagation: Vectorized implementation



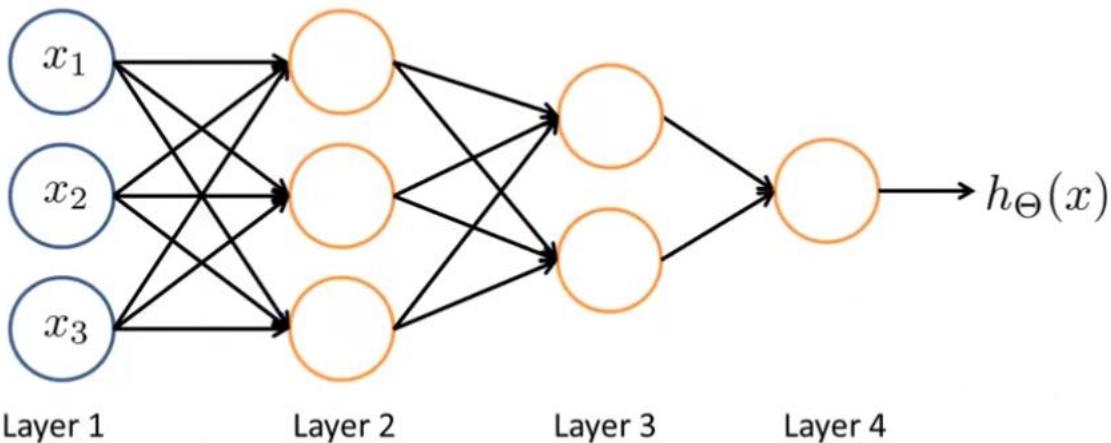
$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$\begin{aligned} z^{(2)} &= \Theta^{(1)} \times a^{(1)} \\ a^{(2)} &= g(z^{(2)}) \\ \text{Add } a_0^{(2)} &= 1. \rightarrow a^{(2)} \in \mathbb{R}^4 \\ z^{(3)} &= \Theta^{(2)} a^{(2)} \\ h_\Theta(x) &= a^{(3)} = g(z^{(3)}) \end{aligned}$$

$$z_k^{(2)} = \Theta_{k,0}^{(1)}x_0 + \Theta_{k,1}^{(1)}x_1 + \dots + \Theta_{k,n}^{(1)}x_n$$

$$\begin{aligned} a_1^{(2)} &= g(z_1^{(2)}) \\ a_2^{(2)} &= g(z_2^{(2)}) \\ \rightarrow a_3^{(2)} &= g(z_3^{(2)}) \end{aligned}$$

Other network architectures



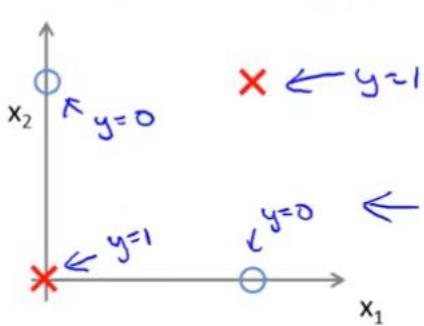
- x_1, x_2, x_3 가 θ 행렬로 가공되어 a_1, a_2, a_3 가 나온다. \rightarrow between layer j and layer $j+1$, we are doing exactly the same thing as we did in logistic regression.

$$\circledcirc z^{(2)} = \Theta^{(1)} a^{(1)}; a^{(2)} = g(z^{(2)})$$

- Example :

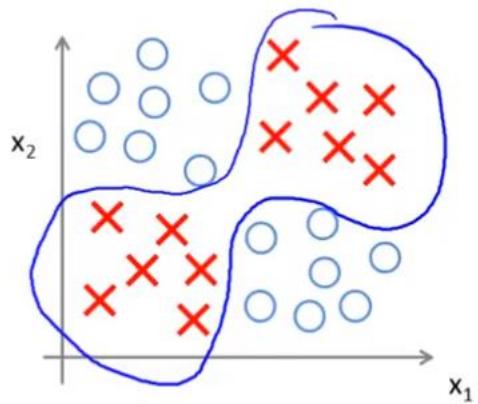
Non-linear classification example: XOR/XNOR

$\rightarrow x_1, x_2$ are binary (0 or 1).



$$y = \underline{x_1 \text{ XOR } x_2}$$

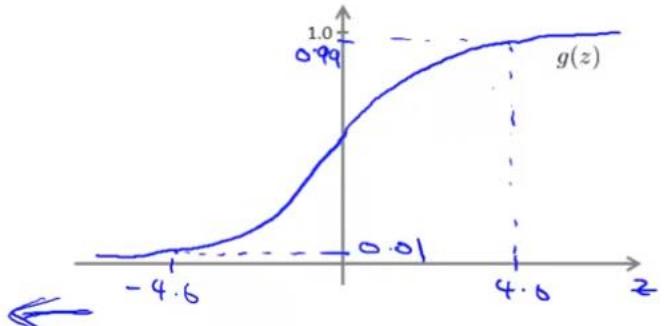
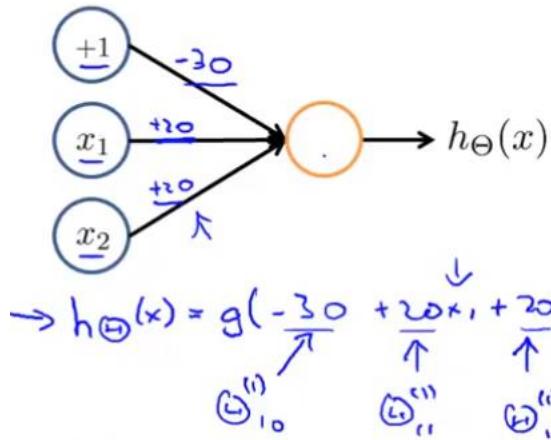
$$\overbrace{\underline{x_1 \text{ XNOR } x_2}}^{\leftarrow} \quad \overbrace{\underline{\text{NOT } (x_1 \text{ XOR } x_2)}}^{\leftarrow}$$



Simple example: AND

$\rightarrow x_1, x_2 \in \{0, 1\}$

$\rightarrow y = x_1 \text{ AND } x_2$

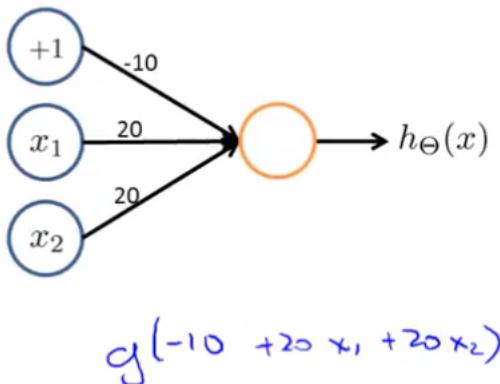


x_1	x_2	$h_{\Theta}(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

$h_{\Theta}(x) \approx x_1 \text{ AND } x_2$

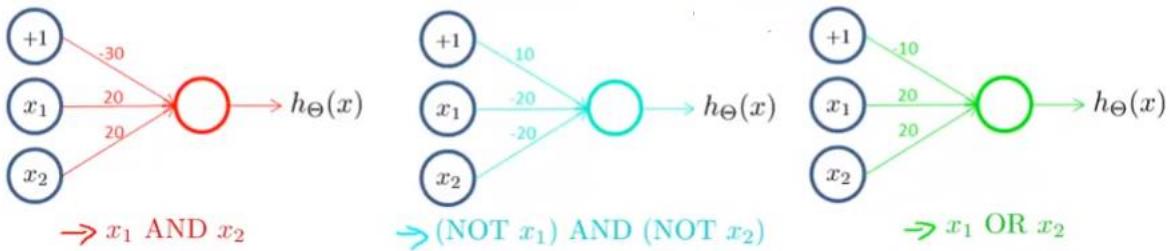
- Or function

Example: OR function

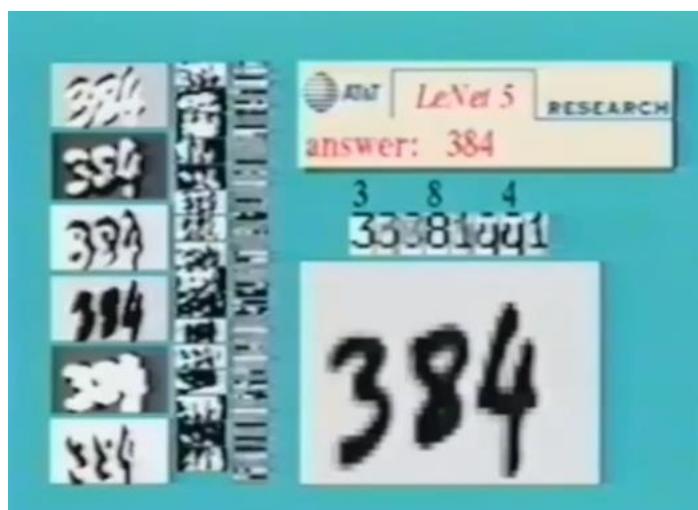
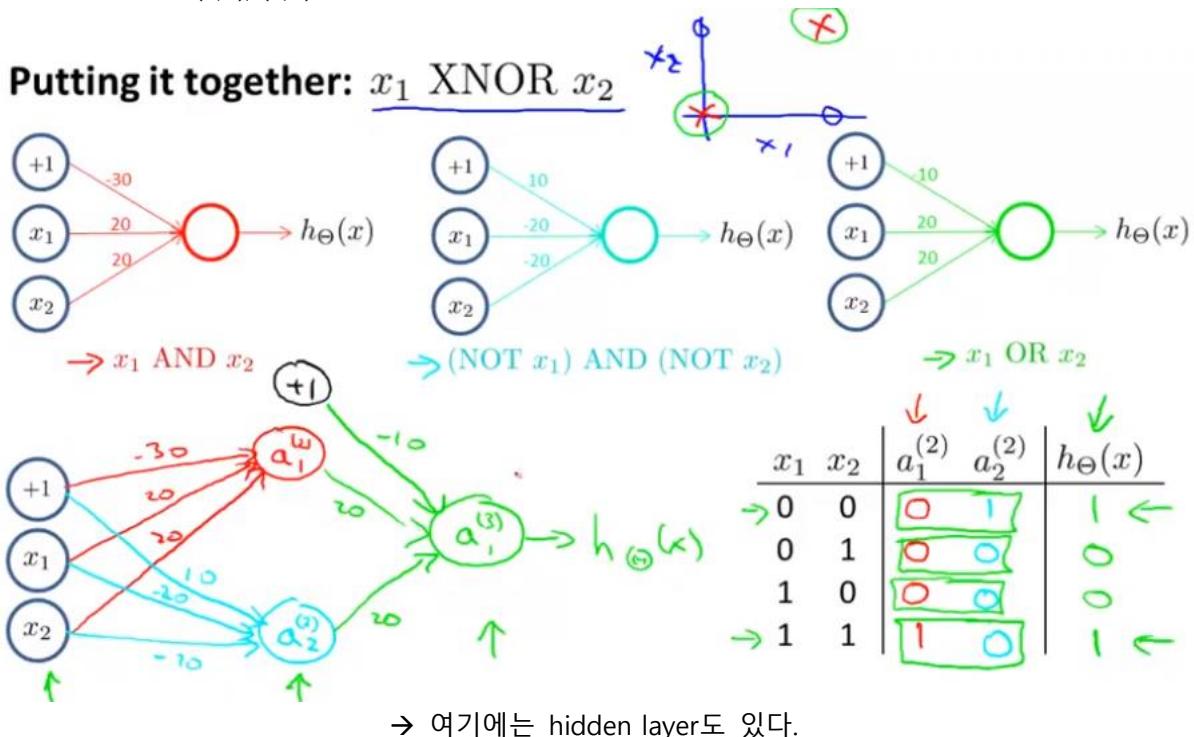


x_1	x_2	$h_{\Theta}(x)$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	≈ 1
1	1	≈ 1

- Other example



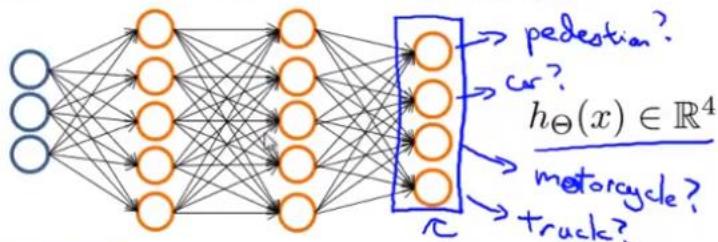
- XNOR function을 만들어보자! 둘다 0이거나, 둘다 1이거나. (exclusive nor이다)
- XOR은 exclusive or → (1,0), (0,1) 만 true



→ 이게 어떻게 가능한거지?? → multiclass classification
→ 예를 들어서 숫자를 삼등분해서 ...

- Multiclass classification
 - 아까는 0아니면 1이 $h(x)$ 로의 전부였으나, 이번에는 결과 output이 정말 많다.

Multiple output units: One-vs-all.



Want $h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, etc.
 when pedestrian when car when motorcycle

Want $h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, etc.
 when pedestrian when car when motorcycle

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

→ $y^{(i)}$ one of $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$
 pedestrian car motorcycle truck

~~Previously~~ $y \in \{1, 2, 3, 4\}$
 $h_{\Theta}(x^{(i)}) \approx y^{(i)}$

- 구조 : → hidden layer를 저절로 만들어준다.

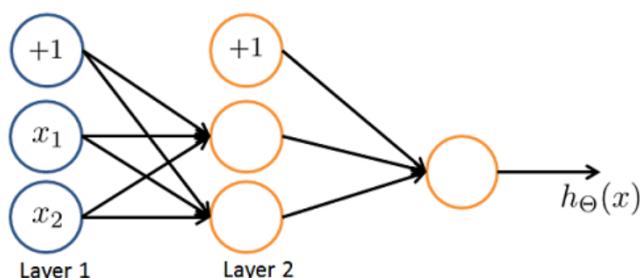
$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \\ \dots \end{bmatrix} \rightarrow \begin{bmatrix} a_0^{(3)} \\ a_1^{(3)} \\ a_2^{(3)} \\ \dots \end{bmatrix} \rightarrow \dots \rightarrow \begin{bmatrix} h_{\Theta}(x)_1 \\ h_{\Theta}(x)_2 \\ h_{\Theta}(x)_3 \\ h_{\Theta}(x)_4 \end{bmatrix}$$

- 질문 :



5. You are using the neural network pictured below and have learned the parameters
 $\Theta^{(1)} = \begin{bmatrix} 1 & -1.5 & 3.7 \\ 1 & 5.1 & 2.3 \end{bmatrix}$ (used to compute $a^{(2)}$) and $\Theta^{(2)} = [1 \quad 0.6 \quad -0.8]$ (used to compute $a^{(3)}$) as a function of $a^{(2)}$). Suppose you swap the parameters for the first hidden layer between its two units so $\Theta^{(1)} = \begin{bmatrix} 1 & 5.1 & 2.3 \\ 1 & -1.5 & 3.7 \end{bmatrix}$ and also swap the output layer so $\Theta^{(2)} = [1 \quad -0.8 \quad 0.6]$. How will this change the value of the output $h_{\Theta}(x)$?

1 / 1 points



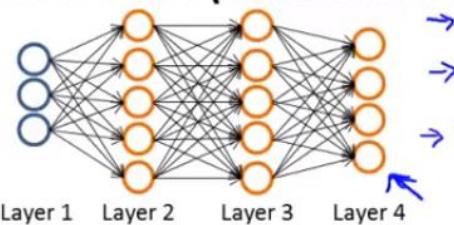
→ 왜 이게 답이 바뀌지 않는다면? 순서가 반대가 되니까 당연히 같지!!

1.6.3. Neural Network의 설계 방법 : How to minimize cost function

→ Backpropagation /

- Multiclass classification

Neural Network (Classification)



- $\rightarrow \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- $\rightarrow L = \text{total no. of layers in network}$ $L=4$
- $\rightarrow s_l = \text{no. of units (not counting bias unit) in layer } l$ $s_1=3, s_2=5, s_3=s_4=L=4$

Binary classification

$$y = 0 \text{ or } 1$$

$$h_{\Theta}(x)$$

$$1 \text{ output unit}$$

$$h_{\Theta}(x) \in \mathbb{R}$$

$$s_L = 1, K=1$$

Let's first define a few variables that we will need to use:

- $L = \text{total number of layers in the network}$
- $s_l = \text{number of units (not counting bias unit) in layer } l$
- $K = \text{number of output units/classes}$

→ logistic regression에서 sample x 의 training set 데이터가 m 개 있으니까 1부터 m 까지 시그마로 하는거다. 다만 이제 Neural Network로 하게 되면 y 값이 0,1의 하나의 값이 아니라 벡터로 존재하다 보니까 cost function도 $i = 1$ 부터 k 까지 더해야 한다.

Cost function

Logistic regression:

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

Neural network:

$$\rightarrow h_{\Theta}(x) \in \mathbb{R}^K \quad (h_{\Theta}(x))_i = i^{\text{th}} \text{ output}$$

$$\begin{aligned} J(\Theta) &= -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\Theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_k) \right] \\ &\quad + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2 \end{aligned}$$

→ regularization term에서는 x_0 부분(bias term)을 하지 않는다.

- regularization term 관련

Note:

- the double sum simply adds up the logistic regression costs calculated for each cell in the output layer
- the triple sum simply adds up the squares of all the individual Θ s in the entire network.
- the i in the triple sum does **not** refer to training example i

- How to minimize cost function of Neural network

1.6.3.1. Backpropagation

Gradient computation

$$\begin{aligned} \Rightarrow J(\Theta) &= -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log h_\theta(x^{(i)})_k + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)})_k) \right] \\ &\quad + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_j^{(l)})^2 \\ \Rightarrow \min_{\Theta} J(\Theta) \end{aligned}$$

Need code to compute:

- $J(\Theta)$
- $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

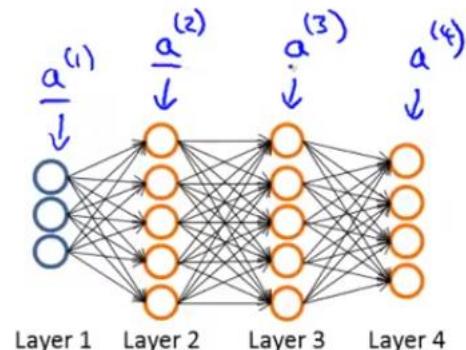
$\rightarrow l$ 은 layer다.

Gradient computation

Given one training example (x, y):

Forward propagation:

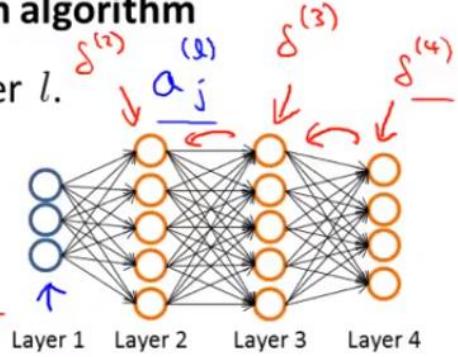
$$\begin{aligned} a^{(1)} &= x \\ \rightarrow z^{(2)} &= \Theta^{(1)} a^{(1)} \\ \rightarrow a^{(2)} &= g(z^{(2)}) \quad (\text{add } a_0^{(2)}) \\ \rightarrow z^{(3)} &= \Theta^{(2)} a^{(2)} \\ \rightarrow a^{(3)} &= g(z^{(3)}) \quad (\text{add } a_0^{(3)}) \\ \rightarrow z^{(4)} &= \Theta^{(3)} a^{(3)} \\ \rightarrow a^{(4)} &= h_\Theta(x) = g(z^{(4)}) \end{aligned}$$



- cost function을 편미분하려면 backpropagation을 하는 방법이 있다. 델타4를 구하고, 대입하여 델타3를 구하고, 델타2를 구하게 되면 거기다가 각각 activation을 곱하면 cost function을 편미분한 값이 된다.
이거 미분이 왜 저렇게 되는거지?

Gradient computation: Backpropagation algorithm

Intuition: $\delta_j^{(l)}$ = "error" of node j in layer l .



For each output unit (layer $L = 4$)

$$\delta_j^{(4)} = \underline{a_j^{(4)}} - \underline{y_j} \quad (\underline{h_\theta(x)_j})_j \quad \underline{\delta^{(4)}} = \underline{a} - \underline{y}$$

$$\rightarrow \delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)})$$

$$\rightarrow \delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \cdot g'(z^{(2)})$$

$$(\text{No } \delta^{(1)}) \quad \frac{\partial}{\partial \Theta_{ij}^{(2)}} J(\Theta) = a_j^{(2)} \delta_i^{(3)} \quad (\text{ignoring } \lambda; \text{ if } \lambda = 0)$$

빨간색은 vectorized implementation

Backpropagation algorithm

→ Training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j). (use to compute $\frac{\partial}{\partial \Theta_{ij}^{(2)}} J(\Theta)$)

For $i = 1$ to m $\leftarrow (x^{(i)}, y^{(i)})$.

Set $a^{(1)} = x^{(i)}$

→ Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \dots, L$

→ Using $y^{(i)}$, compute $\delta^{(L)} = \underline{a^{(L)}} - \underline{y^{(i)}}$

→ Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ ~~$\delta^{(1)}$~~

→ $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ $\Delta^{(2)} := \Delta^{(1)} + \delta^{(2)} (a^{(1)})^T$

→ $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$ if $j \neq 0$

→ $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$ if $j = 0$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$$

- 정리

Back propagation Algorithm

Given training set $\{(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})\}$

- Set $\Delta_{ij}^{(l)} := 0$ for all (l, i, j) , (hence you end up having a matrix full of zeros)

For training example $t = 1$ to m :

1. Set $a^{(1)} := x^{(t)}$

2. Perform forward propagation to compute $a^{(l)}$ for $l=2,3,\dots,L$

3. Using $y^{(t)}$, compute $\delta^{(L)} = a^{(L)} - y^{(t)}$

Where L is our total number of layers and $a^{(L)}$ is the vector of outputs of the activation units for the last layer. So our "error values" for the last layer are simply the differences of our actual results in the last layer and the correct outputs in y. To get the delta values of the layers before the last layer, we can use an equation that steps us back from right to left:

4. Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ using $\delta^{(l)} = ((\Theta^{(l)})^T \delta^{(l+1)}) .* a^{(l)} .* (1 - a^{(l)})$

The delta values of layer l are calculated by multiplying the delta values in the next layer with the theta matrix of layer l. We then element-wise multiply that with a function called g', or g-prime, which is the derivative of the activation function g evaluated with the input values given by $z^{(l)}$.

The g-prime derivative terms can also be written out as:

$$g'(z^{(l)}) = a^{(l)} .* (1 - a^{(l)})$$

5. $\Delta_{i,j}^{(l)} := \Delta_{i,j}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ or with vectorization, $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

Hence we update our new Δ matrix.

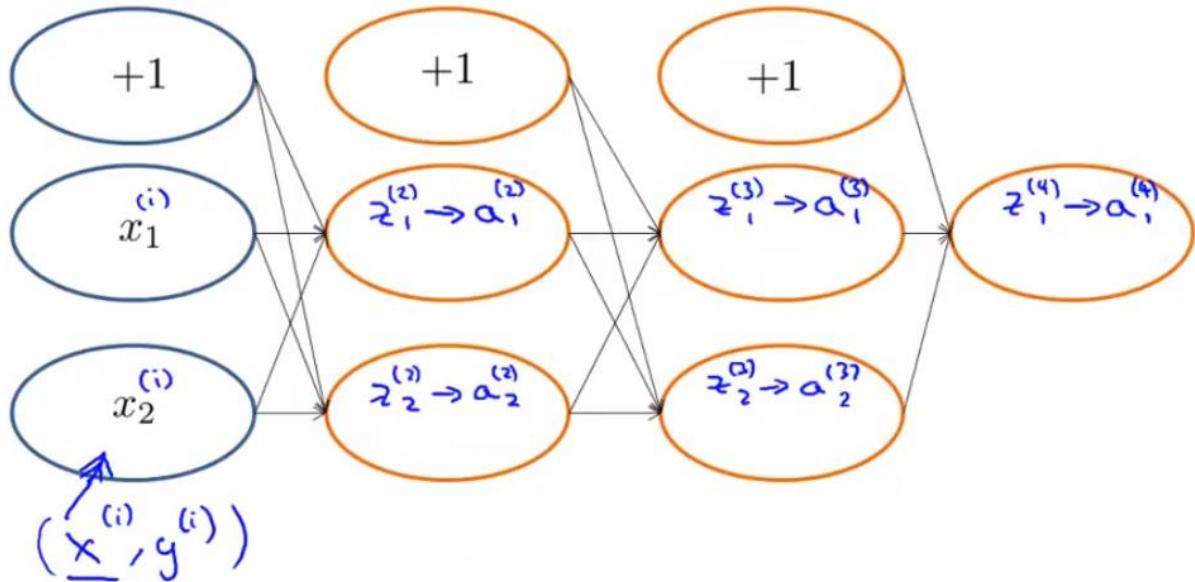
- $D_{i,j}^{(l)} := \frac{1}{m} (\Delta_{i,j}^{(l)} + \lambda \Theta_{i,j}^{(l)})$, if $j \neq 0$.
- $D_{i,j}^{(l)} := \frac{1}{m} \Delta_{i,j}^{(l)}$ if $j = 0$

The capital-delta matrix D is used as an "accumulator" to add up our values as we go along and eventually compute our partial derivative. Thus we get $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$

- 이 과정이 이해가 안되더라도 괜찮다. 어려운 건 누구에게나 마찬가지다. 블랙박스 같지만, 순차적으로 차근차근 이해하도록 해보자.

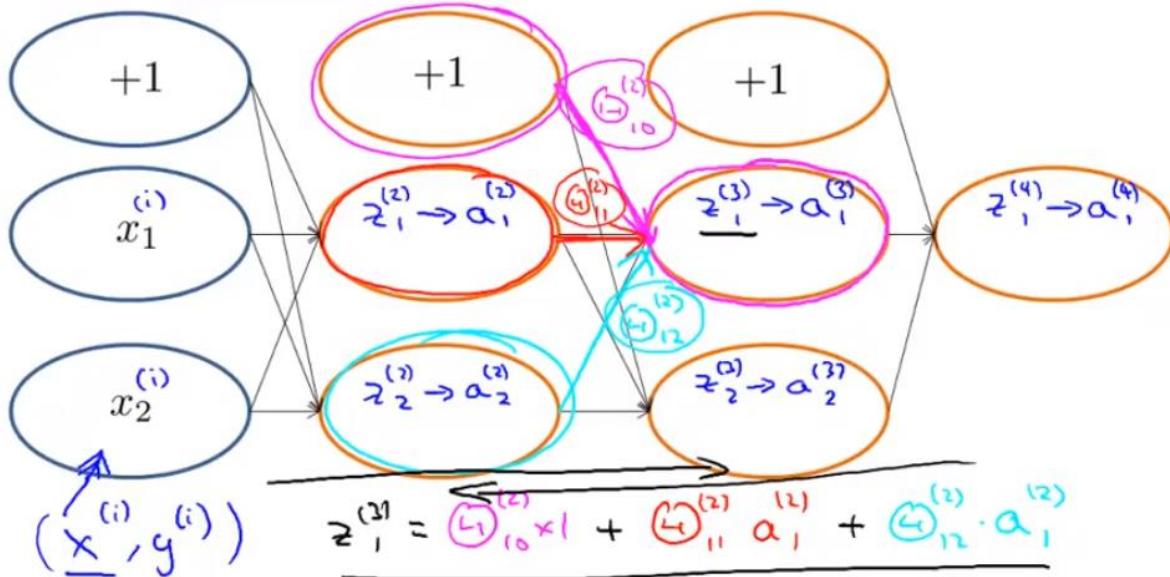
- Forward Propagation의 과정 :

Forward Propagation



→ 이렇게 흘러간다. 좀 더 자세히 $z_1^{(3)}$ 이 어떻게 만들어지는지 살펴보자.

Forward Propagation



- Backpropagation

What is backpropagation doing?

$$J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_\Theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\Theta(x^{(i)})) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$

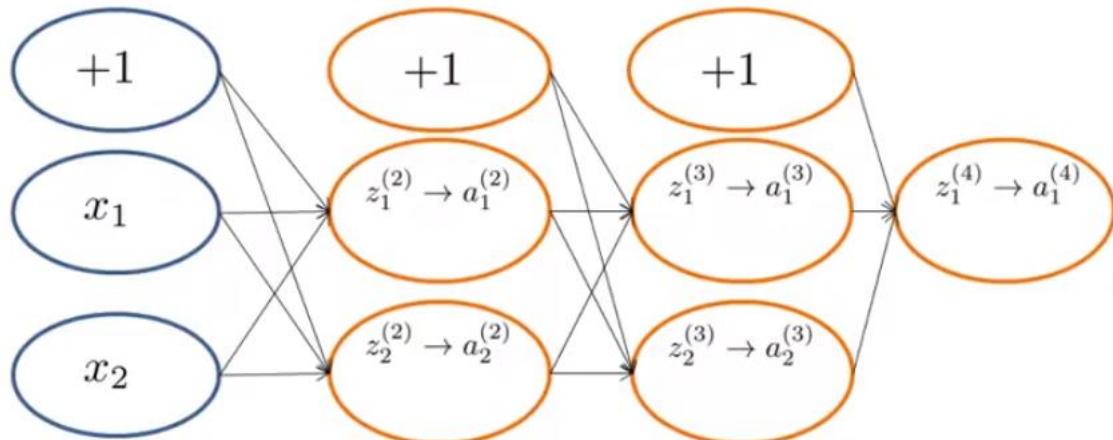
Focusing on a single example $x^{(i)}$, $y^{(i)}$, the case of 1 output unit, and ignoring regularization ($\lambda = 0$),

$$\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log h_\Theta(x^{(i)})$$

(Think of $\text{cost}(i) \approx (h_\Theta(x^{(i)}) - y^{(i)})^2$)

I.e. how well is the network doing on example i?

Forward Propagation

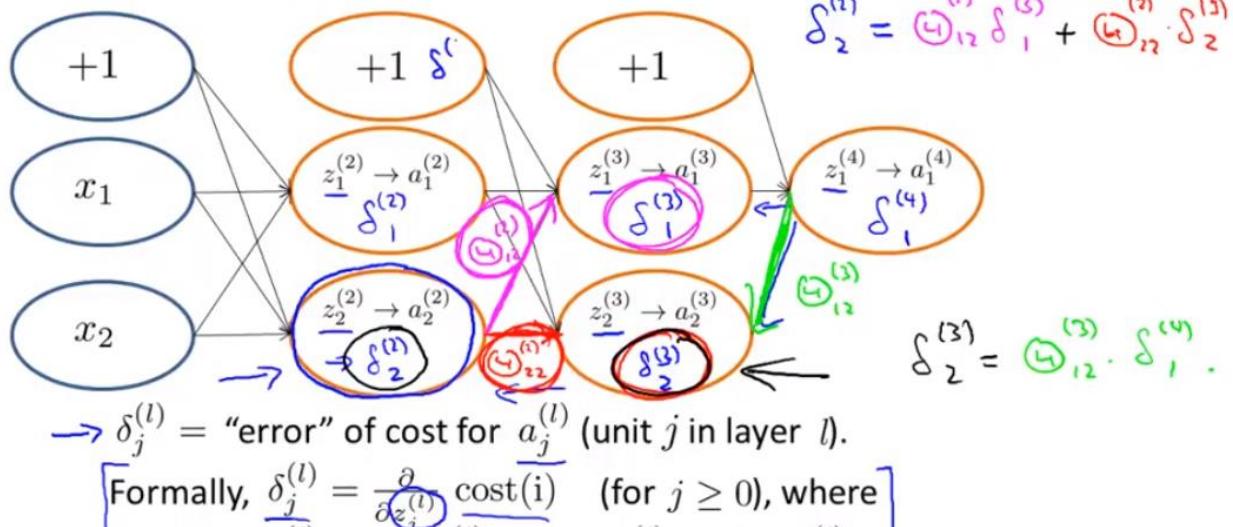


$\rightarrow \delta_j^{(l)}$ = "error" of cost for $a_j^{(l)}$ (unit j in layer l).

Formally, $\delta_j^{(l)} = \frac{\partial \text{cost}(i)}{\partial z_j^{(l)}}$ (for $j \geq 0$), where
 $\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log h_\Theta(x^{(i)})$

→ 이 편미분식은 이해가 안 된다. 어쨌든 간에 $\delta(4)$ 를 구하고 $\delta(3)$ 를 구하고, 거꾸로 구하는 방법이다.

Forward Propagation



If we consider simple non-multiclass classification ($k = 1$) and disregard regularization, the cost is computed with:

$$\text{cost}(t) = y^{(t)} \log(h_\Theta(x^{(t)})) + (1 - y^{(t)}) \log(1 - h_\Theta(x^{(t)}))$$

$\rightarrow \delta(2) = \theta * \delta(3) + \theta * \delta(3)$ 라는 게 어떻게 나온거지? 왜 이렇게 계산을 해야만 하는거지? 굉장히 신기하다.

1.6.4. Backpropagation과 octave implementation - Octave에서 어떻게 해야하는가!!

1.6.4.1. Unrolling parameters

- 목적 : optimization
- 이제는 theta가 vector가 아니다. matrices이다.

Advanced optimization

```
function [jVal, gradient] = costFunction(theta)
    % theta: R^(n+1) (vector)
    ...
    optTheta = fminunc(@costFunction, initialTheta, options)
```

Neural Network ($L=4$):

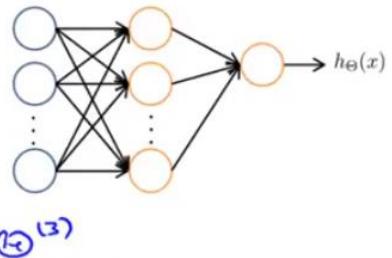
- $\rightarrow \Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices (Theta1, Theta2, Theta3)
- $\rightarrow D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D1, D2, D3)
- "Unroll" into vectors

$\rightarrow 0|$ matrix data를 어떻게 unroll해야지 빨간색으로 있는 gradient나 theta에 vector 형태로 집어넣을 수 있을까?

- $(:)$ 를 쓰고 다시 행렬로 바꿀 때는 **reshape** command를 쓰면 된다.

Example

$s_1 = 10, s_2 = 10, s_3 = 1$
 $\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$
 $D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$
 $\rightarrow \text{thetaVec} = [\Theta^{(1)}; \Theta^{(2)}; \Theta^{(3)}];$
 $\rightarrow \text{DVec} = [D^{(1)}; D^{(2)}; D^{(3)}];$
 $\Theta^{(1)} = \text{reshape}(\text{thetaVec}(1:110), 10, 11);$
 $\Theta^{(2)} = \text{reshape}(\text{thetaVec}(111:220), 10, 11);$
 $\Theta^{(3)} = \text{reshape}(\text{thetaVec}(221:231), 1, 11);$



→ ; 로 연결하니까 열로 쭉쭉 길게 unroll할 수 있다.

- 이것을 활용해보자. 이것은 gradient descent가 아니라 advanced optimization method이다.

Learning Algorithm

\rightarrow Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.
 \rightarrow Unroll to get `initialTheta` to pass to
 \rightarrow `fminunc(@costFunction, initialTheta, options)`

```

function [jval, gradientVec] = costFunction(thetaVec)
    → From thetaVec, get  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$  reshape
    → Use forward prop/back prop to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\Theta)$ 
    Unroll  $D^{(1)}, D^{(2)}, D^{(3)}$  to get gradientVec.

```

→ initial theta를 만들 때는 vector로 길게 해서 집어넣고, reshape한 다음에 D 계산하고(reshape해서 행렬로 만들어야 Forward, Backpropagation이 편하다), D를 다시 vector로 해서 gradientVec에 assign한다.

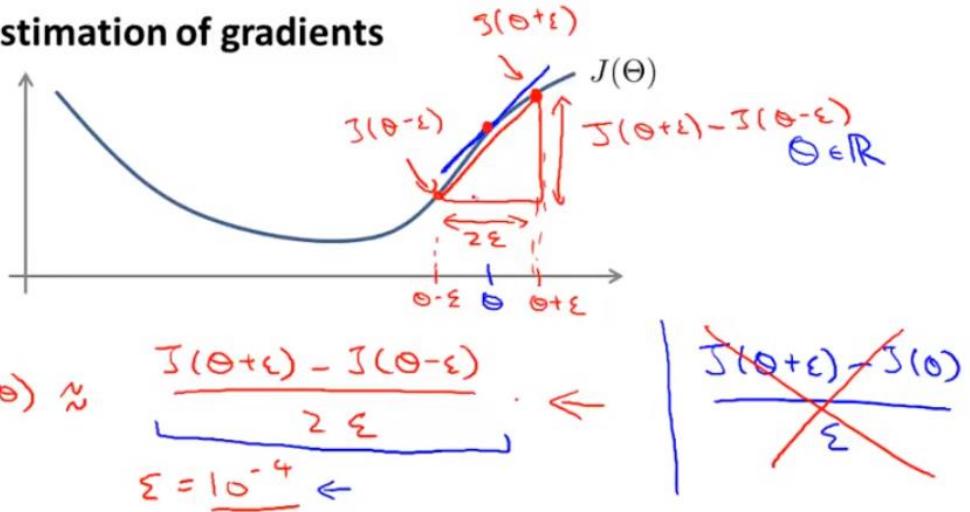
- 문제 : backpropagation이 잘 되는 것 같아 보여도 잘 못 된 것이 있을 수 있다!

→ Gradient checking

- Cost function J 한점의 gradient를 구하기 위해서 함수 위의 두 점을 이어서 기울기를 구한다.

- 계산 방법 :

Numerical estimation of gradients



Implement: gradApprox = $\frac{(J(\text{theta} + \text{EPSILON}) - J(\text{theta} - \text{EPSILON}))}{(2 * \text{EPSILON})}$

→ 이를 다시 쓰면 다음과 같다. epsilon은 10^{-4} 가 적당하다.

Parameter vector θ

$\rightarrow \theta \in \mathbb{R}^n$ (E.g. θ is "unrolled" version of $\underline{\Theta^{(1)}}, \underline{\Theta^{(2)}}, \underline{\Theta^{(3)}}$)

$\rightarrow \theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$

$\rightarrow \frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$

$\rightarrow \frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$

\vdots

$\rightarrow \frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_n - \epsilon)}{2\epsilon}$

- 이를 backpropagation의 D와 얼마나 같은지 비교하면 된다.

```
for i = 1:n,
    thetaPlus = theta;
    thetaPlus(i) = thetaPlus(i) + EPSILON;
    thetaMinus = theta;
    thetaMinus(i) = thetaMinus(i) - EPSILON;
    gradApprox(i) = (J(thetaPlus) - J(thetaMinus)) / (2 * EPSILON);
end;
```

$\frac{\partial}{\partial \theta_i} J(\theta)$

$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i + \epsilon \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \theta_i + \epsilon$

Check that gradApprox ≈ DVec ←
 ↑
 From back prop.

```

1  epsilon = 1e-4;
2  for i = 1:n,
3      thetaPlus = theta;
4      thetaPlus(i) += epsilon;
5      thetaMinus = theta;
6      thetaMinus(i) -= epsilon;
7      gradApprox(i) = (J(thetaPlus) - J(thetaMinus))/(2*epsilon)
8  end;
9

```

- 실행 방법 : gradApprox는 computationally expensive하기 때문에, 맞는지만 확인하고 꺼야 한다.

Implementation Note:

- - Implement backprop to compute D_{vec} (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$).
- - Implement numerical gradient check to compute gradApprox.
- - Make sure they give similar values.
- - Turn off gradient checking. Using backprop code for learning.
↳ D_{vec}

Important:

- - Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of `costFunction(...)`) your code will be very slow.

1.6.4.2. Random Initialization - 어떻게 initial Theta를 정할 것인가?

Initial value of Θ

For gradient descent and advanced optimization method, need initial value for Θ .

```
optTheta = fminunc(@costFunction,
                    initialTheta, options)
```

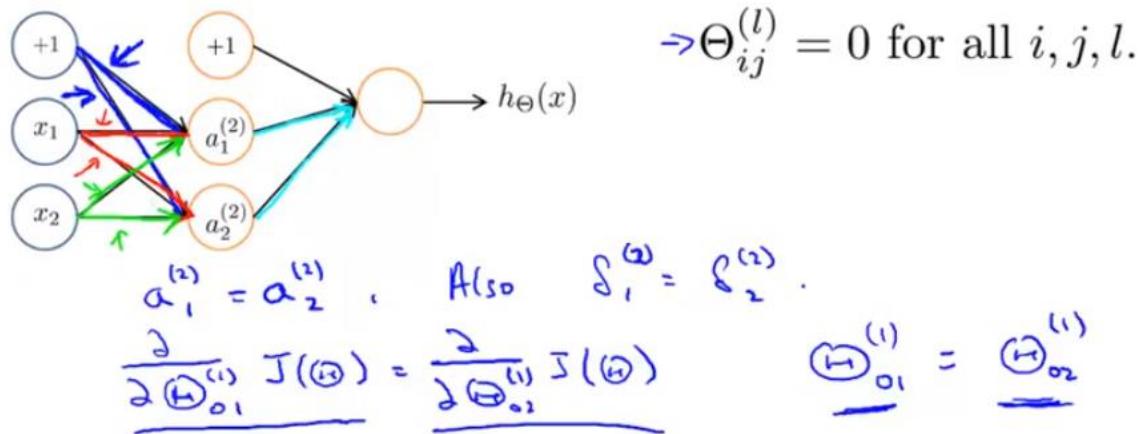
Consider gradient descent

Set initialTheta = zeros(n, 1) ?

- theta를 모두 zero로 하면 안 되는 이유:

→ hidden unit a가 다 똑같이 나와서 ($a_1=a_2$), theta를 업데이트하려고 해도 잘 안된다!! Backpropagation을 한다고 가정을 하면, a가 같게 되면 δ 도 같게 된다. δ 가 같게 되면 cost function도 같게 된다. 따라서 cost function을 미분한 값이 같게 되어서 거기에 learning rate를 곱한 값으로 gradient descent를 해도 theta값이 다르지 않게 된다.

Zero initialization



After each update, parameters corresponding to inputs going into each of two hidden units are identical.

→ 질문!! theta가 다 0이라면서? 왜 꼭 2개씩 짹을 이루어서 같아야 하는 것인가?

- 해결 방법 : symmetry breaking

Random initialization: Symmetry breaking

→ Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$
 (i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)

E.g.

Random 10×11 matrix (betw. 0 and 1)

→ Theta1 = $\boxed{\text{rand}(10, 11) * (2 * \text{INIT_EPSILON})}$
 $\quad - \underline{\text{INIT_EPSILON}}$; $[-\epsilon, \epsilon]$

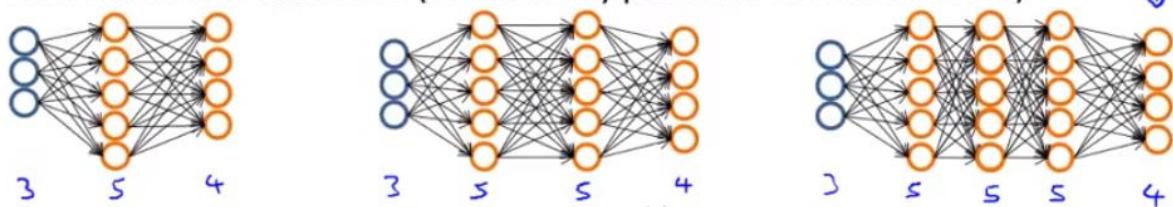
→ Theta2 = $\boxed{\text{rand}(1, 11) * (2 * \text{INIT_EPSILON})}$
 $\quad - \underline{\text{INIT_EPSILON}}$;

1.6.5. Neural network 종합

Step1) Pick a neural network architecture!

Training a neural network

Pick a network architecture (connectivity pattern between neurons)



→ No. of input units: Dimension of features $x^{(i)}$

→ No. output units: Number of classes

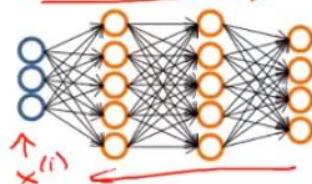
→ hidden layer를 1개 넣는게 가장 좋다. 2개 이상이면, 각 hidden layer에 unit 수가 layer끼리 같으면 좋다.

hidden layer의 unit 수가 많으면 좋다.

Step 2) Training a neural network

Training a neural network

- 1. Randomly initialize weights
- 2. Implement forward propagation to get $h_{\Theta}(x^{(i)})$ for any $x^{(i)}$
- 3. Implement code to compute cost function $J(\Theta)$
- 4. Implement backprop to compute partial derivatives $\frac{\partial}{\partial \Theta_j^{(l)}} J(\Theta)$
- **for** $i = 1:m$ { $(x^{(1)}, y^{(1)})$ $(x^{(2)}, y^{(2)})$, ... , $(x^{(m)}, y^{(m)})$ }
 - Perform forward propagation and backpropagation using example $(x^{(i)}, y^{(i)})$
 - (Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l = 2, \dots, L$).
 - $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l)} (\alpha^{(l)})^T$
 - ...
compute $\frac{\partial}{\partial \Theta_j^{(l)}} J(\Theta)$.



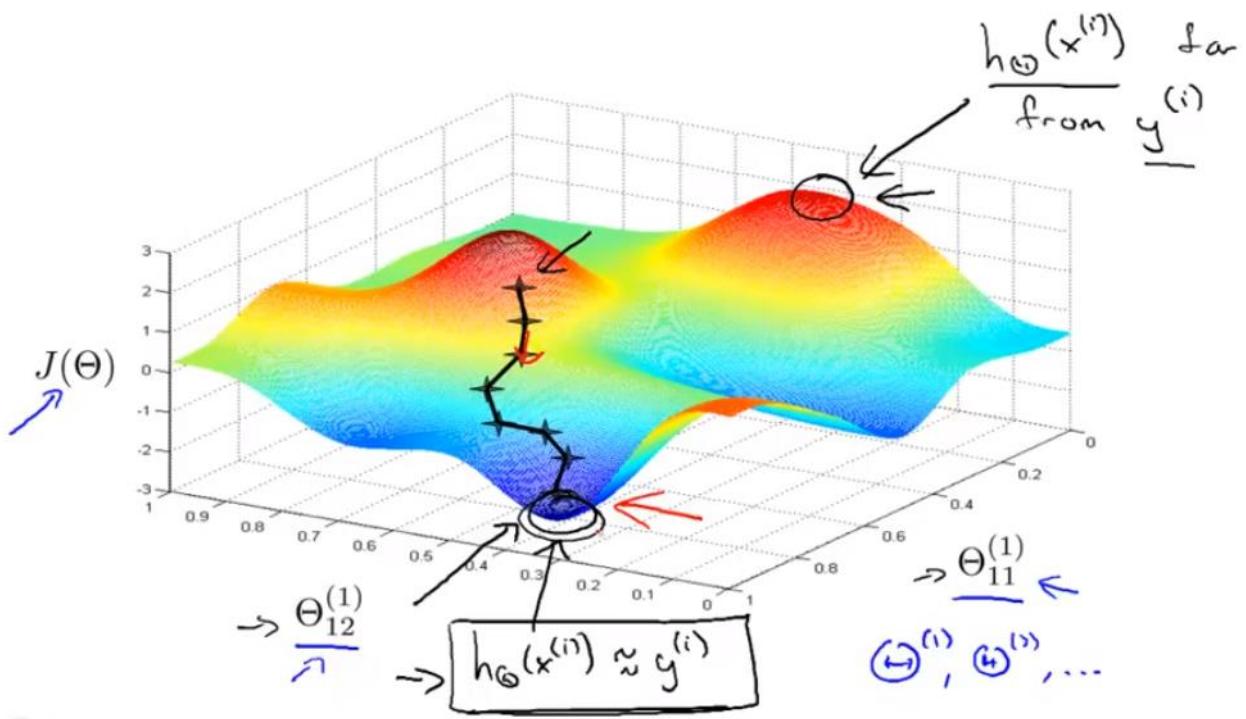
Step3) Gradient checking (검토) / Optimization method

Training a neural network

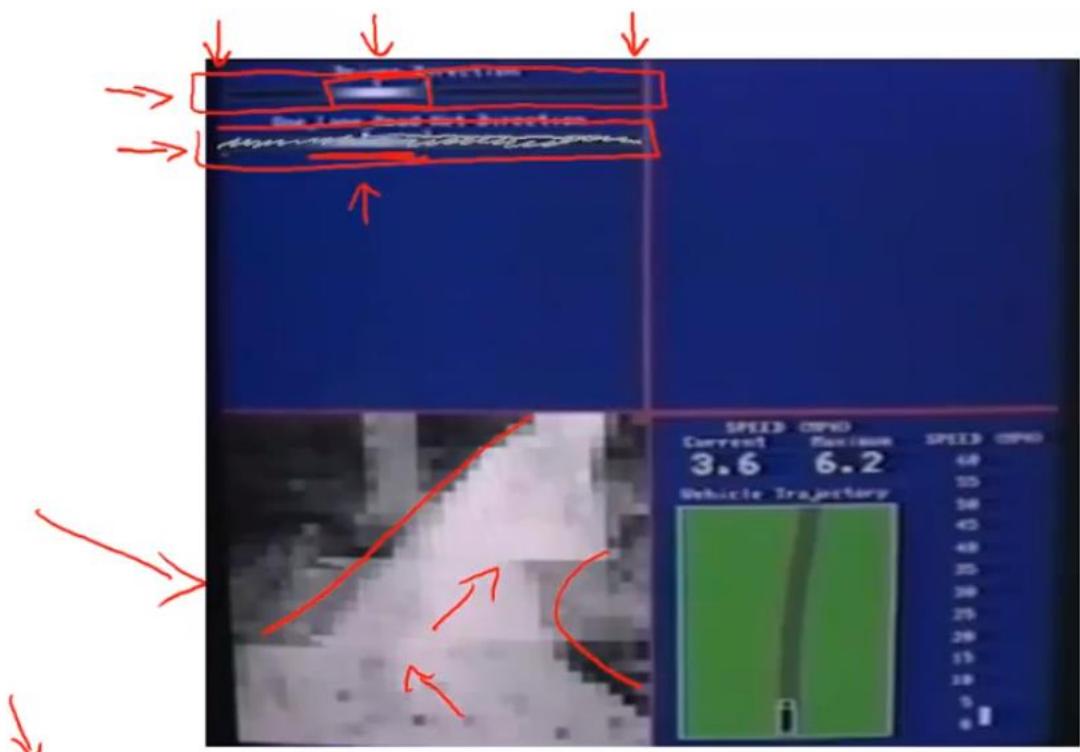
- 5. Use gradient checking to compare $\frac{\partial}{\partial \Theta_j^{(l)}} J(\Theta)$ computed using backpropagation vs. using numerical estimate of gradient of $J(\Theta)$.
Then disable gradient checking code.
- 6. Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\Theta)$ as a function of parameters Θ

→ Gradient descent가 좋다.

ex) gradient descent의 직관



1.6.6. 자율 주행에의 응용



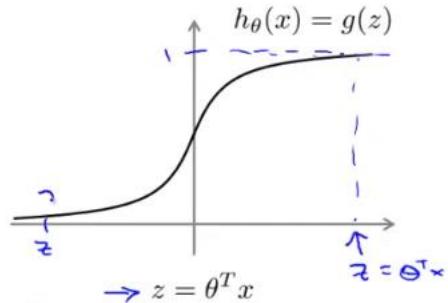
→ 왼쪽 위의 위치가 steering wheel의 위치이다. 위의 bar는 사람이 운전대를 어떻게 잡는지. 주행을 할 때, 사람이 하는 것을 기록한다. 이미지와 함께.

1.7. Supported vector machine - An alternative for logistic regression ← supervised learning

- SVM 개괄 : supervised learning에서 logistic regression의 응용 버전이라고 생각하면 편하다.
- Logistical regression의 복습

Alternative view of logistic regression

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



If $y = 1$, we want $h_{\theta}(x) \approx 1$, $\theta^T x \gg 0$
 If $y = 0$, we want $h_{\theta}(x) \approx 0$, $\theta^T x \ll 0$

→ 여기서 $y=1, y=0$ 일 때, cost는 이렇게 된다.

1.7.1. What is SVM?

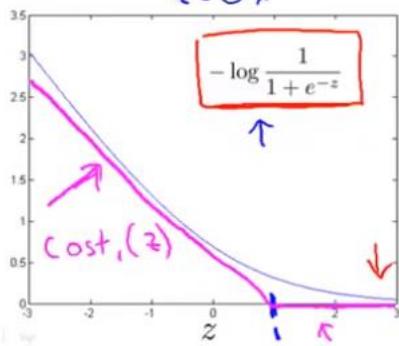
- 목적 : cost function을 간략히 직선화해서 computational cost 줄이는 효과

Alternative view of logistic regression

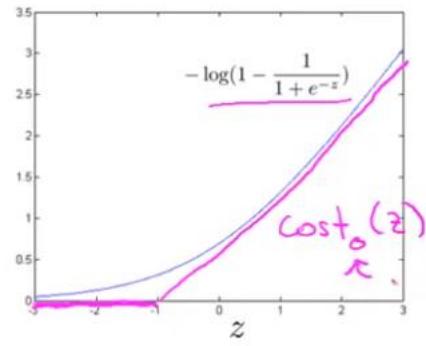
Cost of example: $-(y \log h_{\theta}(x) + (1-y) \log(1 - h_{\theta}(x)))$ ←

$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right) \leftarrow$$

If $y = 1$ (want $\theta^T x \gg 0$):



If $y = 0$ (want $\theta^T x \ll 0$):



- logistic regression의 변형 : cost function을 직선 2개로 표현하게끔 바꾼 것이다. Cost function의 약간의 변형이 있다.

변형-(1) log 함수를 직선의 함수로 바꾼다.

변형-(2) min 함수에 어떤 수를 곱해도 그것을 최소화하는 θ 는 달라지지 않기 때문에 → 관습적으로 regularization term의 lambda를 없애고 cost function 앞에 C라는 상수를 도입한다. (따라서, $C = 1/\lambda$ 이므로 C를 크게 하는 것은 λ 를 작게 하는 것과 같은 효과이다.)

Support vector machine

Logistic regression:

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \underbrace{\left(-\log h_{\theta}(x^{(i)}) \right)}_{cost_1(\theta^T x^{(i)})} + (1 - y^{(i)}) \underbrace{\left(-\log(1 - h_{\theta}(x^{(i)})) \right)}_{cost_0(\theta^T x^{(i)})} \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Support vector machine:

$$\min_{\theta} \cancel{\frac{1}{m}} C \sum_{i=1}^m y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\min_u \frac{(u - 5)^2 + 1}{10} \rightarrow u = 5 \quad \left| \begin{array}{l} A + \lambda B \leftarrow \\ C = \frac{1}{\lambda} \end{array} \right.$$

$$\min_u 10(u - 5)^2 + 10 \rightarrow u = 5 \quad \left| \begin{array}{l} A + B \leftarrow \\ C = A + B \leftarrow \end{array} \right.$$

$$\Rightarrow \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^n \theta_j^2$$

SVM hypothesis

$$\Rightarrow \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^n \theta_j^2$$

Hypothesis:

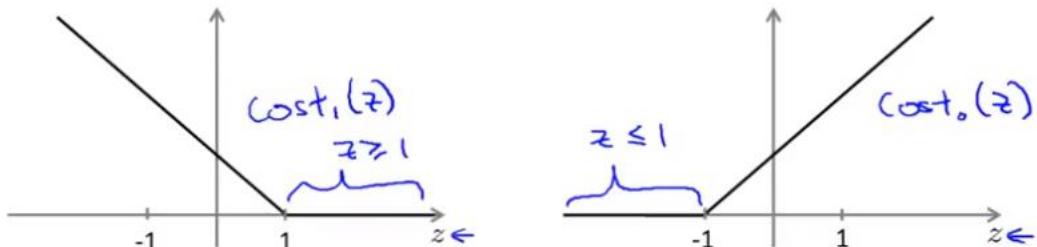
$$h_{\theta}(x) \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

1.7.2. Why do we use SVM? - 이것을 왜 쓰는가? Large margin classifier이기 때문에!!

- implements extra safety margin factor $\rightarrow \theta x$ 가 확실하게 1을 넘기를 원한다.

Support Vector Machine

$$\Rightarrow \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \underbrace{cost_1(\theta^T x^{(i)})}_{z \geq 1} + (1 - y^{(i)}) \underbrace{cost_0(\theta^T x^{(i)})}_{z \leq -1} \right] + \frac{1}{2} \sum_{i=1}^n \theta_j^2$$



\rightarrow If $y = 1$, we want $\theta^T x \geq 1$ (not just ≥ 0)

\rightarrow If $y = 0$, we want $\theta^T x \leq -1$ (not just < 0)

$\theta^T x \geq 1$

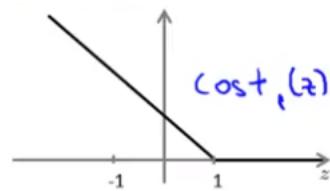
$\theta^T x \leq -1$

- C가 엄청 크면 cost function이 0에 가까워서 minimize될 것이다.

SVM Decision Boundary

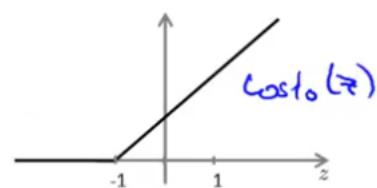
$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Whenever $y^{(i)} = 1$: $\theta^T x^{(i)} \geq 1$



Whenever $y^{(i)} = 0$:

$$\theta^T x^{(i)} \leq 0$$

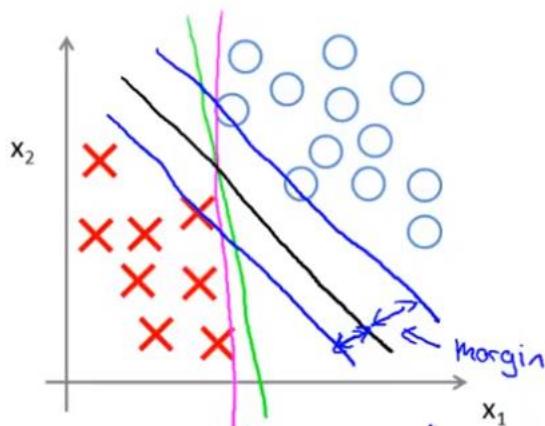


이 식을 최소화하게 되면 cost function의 앞부분은 완전히 0이 되어 버린다. 식을 다시 써보자.

$$\begin{aligned} \min & \quad \cancel{C \sum_{i=1}^m} + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \\ \text{s.t. } & \quad \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1 \\ & \quad \theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0. \end{aligned}$$

- Consequence : 그렇게 되면 large margin을 갖게 된다. (이 부분이 명확한 설명이 없다.) margin이란 decision boundary 선으로부터의 data point의 거리

SVM Decision Boundary: Linearly separable case

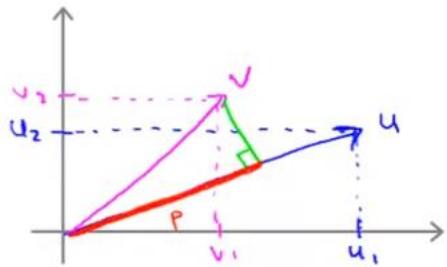


Large margin classifier

- margin을 극대화하는 수학적 증명 :

- 1) inner product가 정사된 길이 곱하기 V 벡터의 길이라는 것을 증명하기.

Vector Inner Product



$$\Rightarrow u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad \Rightarrow v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

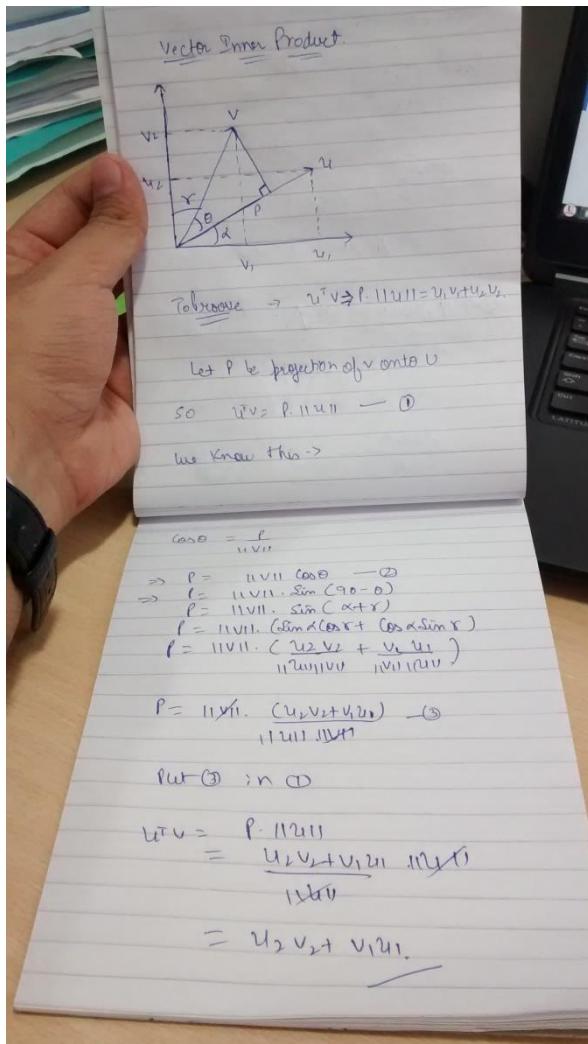
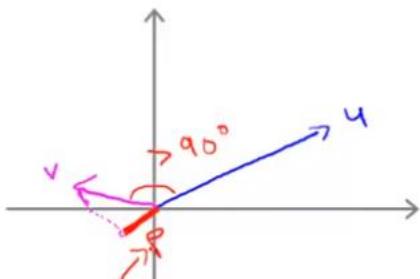
$$u^T v = ? \quad [u_1 \ u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\|u\| = \text{length of vector } u \\ = \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$$

$p = \text{length of projection of } v \text{ onto } u.$

$$\text{Signed } u^T v = \frac{p \cdot \|u\|}{\|u\|} \leftarrow = v^T u \\ = u_1 v_1 + u_2 v_2 \leftarrow p \in \mathbb{R}$$

$$u^T v = p \cdot \|u\| \\ p < 0$$



2) Cost function을 minimize하는 식을 다시 쓰기

$$\omega = (\sum \theta_j)^2$$

SVM Decision Boundary

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} \left(\sqrt{\theta_1^2 + \theta_2^2} \right)^2 = \frac{1}{2} \|\theta\|^2$$

s.t. $\theta^T x^{(i)} \geq 1$ if $y^{(i)} = 1$
 $\theta^T x^{(i)} \leq -1$ if $y^{(i)} = 0$

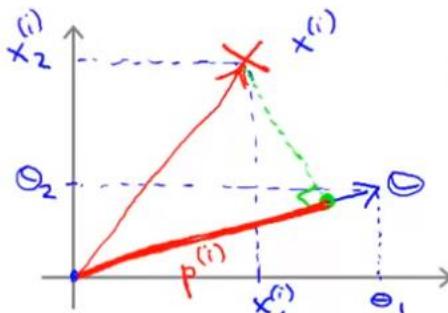
Simplification: $\theta_0 = 0$, $n=2$

$$= \|\theta\|$$

$$\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}, \theta_0 = 0$$

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

\uparrow \uparrow
 $U^T V$



$$\theta^T x^{(i)} = \boxed{p \cdot \|\theta\|} \leftarrow$$

$$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \leftarrow$$

→ subjected to 조건을 1)에서 증명한 식으로 다시 써보자.

3) P가 커지는 곳에 decision boundary가 그려지기 때문에 → margin이 큰 곳에 결정된다.

SVM Decision Boundary

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2$$

s.t. $p^{(i)} \cdot \|\theta\| \geq 1$ if $y^{(i)} = 1$
 $p^{(i)} \cdot \|\theta\| \leq -1$ if $y^{(i)} = -1$

where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector θ .

Simplification: $\theta_0 = 0$

* decision boundary와 theta가 직각인 이유

Q3) How do we know that the theta is perpendicular to the decision boundary?

(thanks to Mentor Chirag for this derivation)

We know that $\theta' * x = 0$ for any x on the boundary since the boundary is where sigmoid = 0.5.

$$\frac{1}{1+e^{-z}} = \frac{1}{2} \Rightarrow e^{-z} = 1 \Rightarrow -z = 0 \Rightarrow \theta' * x = 0$$

So pick two random points on the boundary a and b, then

$$\theta' * a = 0 \text{ and } \theta' * b = 0$$

$$\Rightarrow \theta' * (a - b) = 0 \Rightarrow \theta' \cdot (a - b) = 0$$

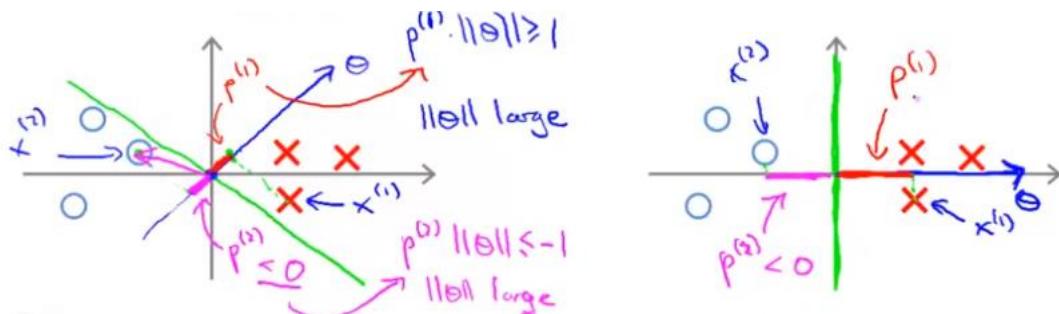
and we know that when the dot product of two vectors is 0, the angle between them is 90 degrees. And since the vector $(a-b)$ is on the decision boundary, θ is perpendicular to the decision boundary.

- decision boundary를 두 가지 경우로 그려서 나눠서 생각해보자.

- decision boundary와 theta는 수직이다.

- x에서 theta에 수직으로 내려서 projection한 벡터의 norm이 p의 값이다.

- 첫번째 경우(왼쪽)에는 p 가 작아진다. 그런데 목적함수를 최소화하려면 $1/2 \|\theta\|^2$ 를 최소화해야 하는데 그러면 $p * |\theta|$ 가 1보다 커지지 않게 될 수 있다. 따라서 제한 조건을 만족시키지 못하게 된다. 두번째 경우처럼 margin을 최대화해야 목적함수의 최소화가 이루어진다. 즉, p 가 큰 곳에 decision boundary(연두색선)이 그려져야 하는 것이다.

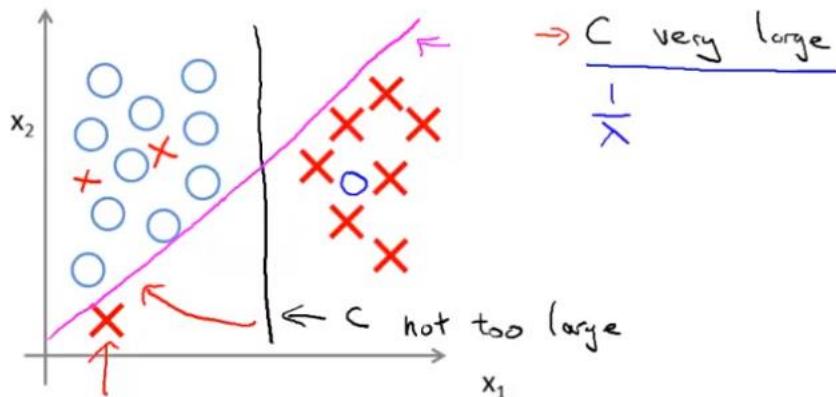


- $\theta_0 = 0$ 이 아니더라도 똑같이 large margin을 찾게 될 것이다.

1.7.3. How do we use SVM? → C조절 (regularization parameter)

- fitting을 할 때, 이상값을 넣으면 안되기 때문에 C값을 낮추면 된다.

Large margin classifier in presence of outliers



→ 여기서 C가 높을 때에는 magenta line처럼 된다. C를 낮추면 검은 색 선처럼 decision boundary가 생긴다.

- 내적을 하게 되면 projection * 길이

- 목적은 theta를 구하는 것이다.

- $p(i) * |\theta|$

- 임의의 theta로 시작을 하고 projection을 다 시킨다. p 가 클수록 좋은 것이다.

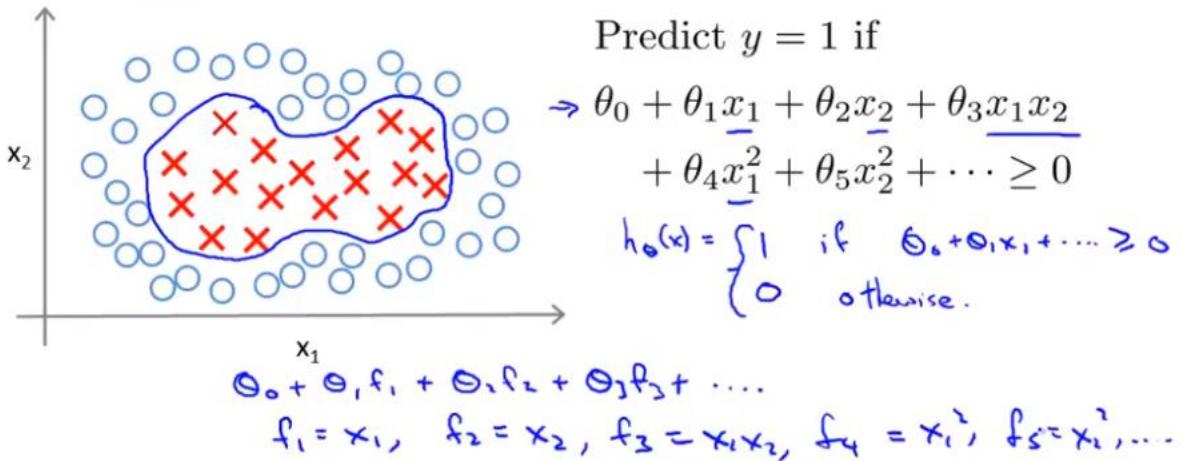
- theta를 이동시켜서 p 값을 점점 크게 만든다.

Decision boundary는 package로 써라.

1.7.4. Using SVM in complex nonlinear classifiers : kernel

- 문제의 정의 : nonlinear classifier를 만들 때, 차수가 너무 많아서 computational cost가 너무 높게 된다. 따라서 여기에 SVM을 적용하면 비용을 낮출 수 있다.

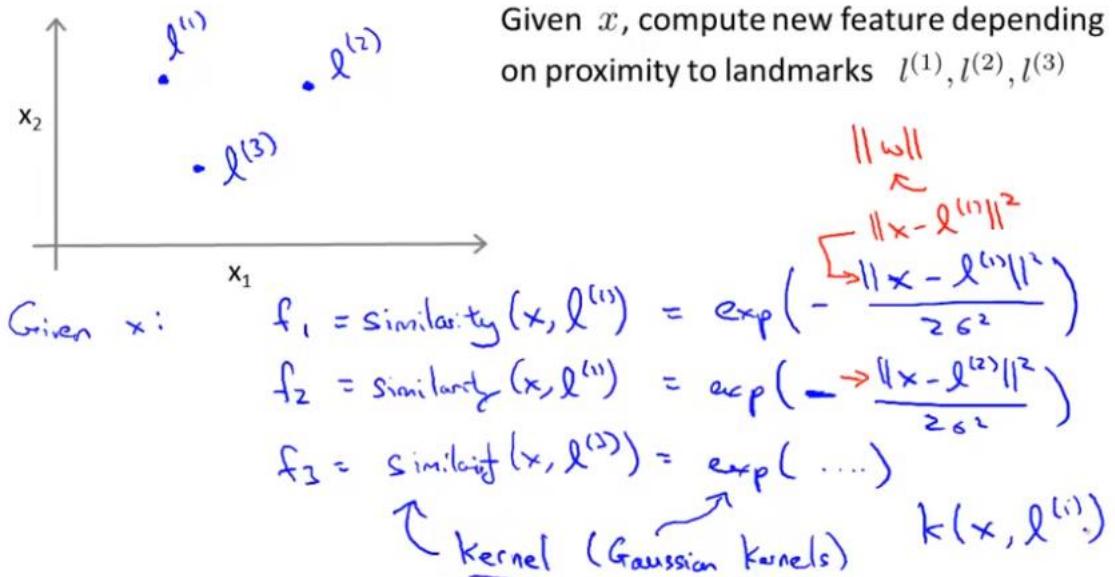
Non-linear Decision Boundary



Is there a different / better choice of the features f_1, f_2, f_3, \dots ?

- 해결 방법 : landmark와의 유사성을 f_1 으로 표시를 하게 되면 어떨까? similarity function은 kernel을 의미한다. kernel 중 하나가 Gaussian kernel이 있다. 다른 kernel 종류도 있다.

Kernel



- 저 식이 작동하는 방식 : 만약에 x 가 landmark에 가까우면 f_1 은 1에 가까워진다. 그러나 landmark에서 멀어지면 f_1 이 0이 된다.

Kernels and Similarity

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{2\sigma^2}\right)$$

If $x \approx l^{(1)}$:

$$f_1 \underset{\uparrow}{\approx} \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$$

$$\begin{aligned} l^{(1)} &\rightarrow f_1 \\ l^{(2)} &\rightarrow f_2 \\ l^{(3)} &\rightarrow f_3 \end{aligned}$$

If x if far from $l^{(1)}$:

$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0.$$

- 예를 들어보자 : f_1 을 y 축으로 하면 (x_1, x_2) 가 3,5에 가까워질수록 1이 되므로 산 모양이 된다. 시그마 제곱이 달라지면 어떻게 될까? 만약에 1에서 0.5가 되면 y 가 1에서 0이 되는 것이 더 빨라지므로 경사가 급격해지고, 만약에 시그마 제곱이 3이 되면 y 가 1에서 0이 되는 것이 더 느려지므로(분모가 커졌으므로 음의 절대값이 더 느리게 커진다) 경사가 완만해진다.

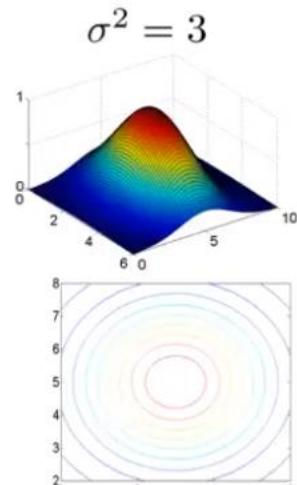
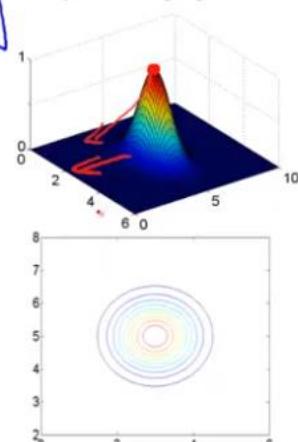
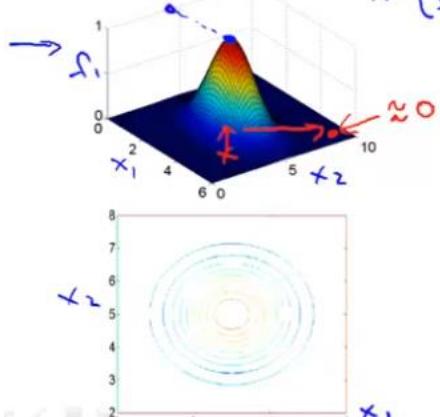
Example:

$$\rightarrow l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

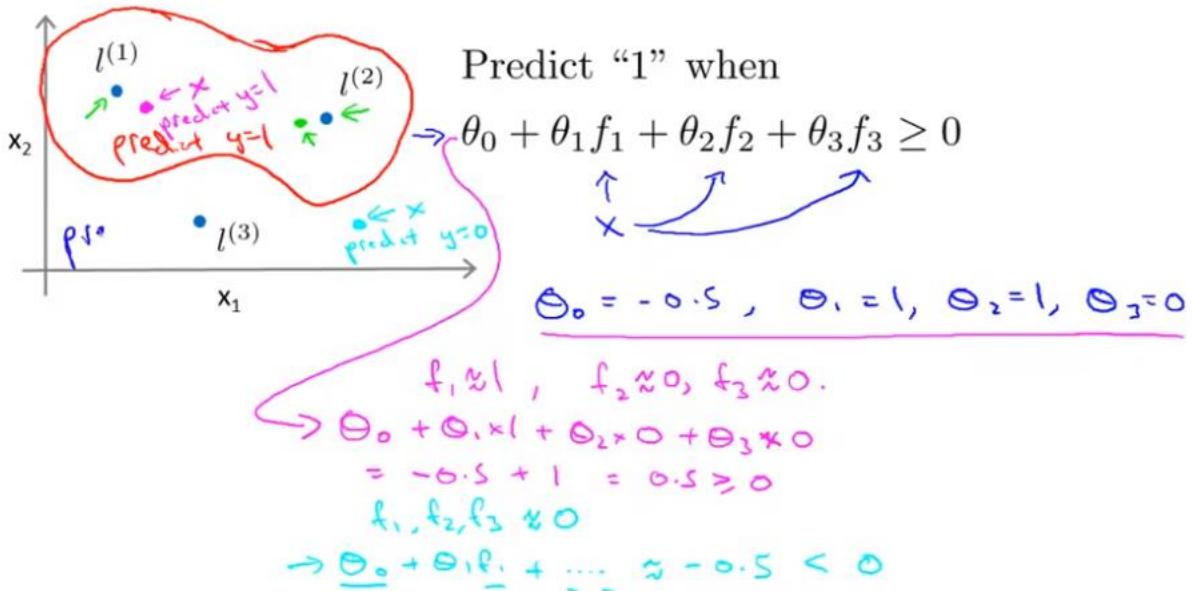
$$\rightarrow \sigma^2 = 1$$

$$x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

$$\sigma^2 = 0.5$$



- 적용 : 만약에 핑크색 점의 data가 있다고 해보자. 이는 landmark 1과 가깝기 때문에 $f_1=1, f_2=0, f_3=0$ 일 것이다. 따라서 이를 parameter를 넣고서 계산해보면 0보다 크므로 $y=1$ 을 predict하게 된다.

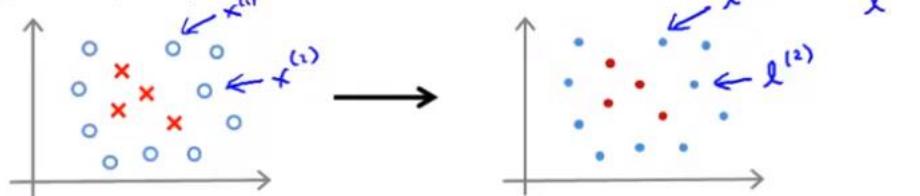


- kernel의 설계 :

- 어떻게 landmark를 설정할 것인가? : m개의 training set을 그래도 landmark로 한다.

Predict $y = 1$ if $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$ ←

Where to get $l^{(1)}, l^{(2)}, l^{(3)}, \dots$?



- 벡터로 만들어버린다.

SVM with Kernels

→ Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$,

→ choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$.

Given example x :

$$\begin{aligned} \rightarrow f_1 &= \text{similarity}(x, l^{(1)}) \\ \rightarrow f_2 &= \text{similarity}(x, l^{(2)}) \\ \dots & \end{aligned}$$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1$$

For training example $(x^{(i)}, y^{(i)})$:

$$\begin{aligned} x^{(i)} \rightarrow f_1^{(i)} &= \text{similarity}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} &= \text{similarity}(x^{(i)}, l^{(2)}) \\ \vdots & \\ f_m^{(i)} &= \text{similarity}(x^{(i)}, l^{(m)}) \end{aligned}$$

$$\begin{aligned} x^{(i)} &\in \mathbb{R}^{n+1} \quad (\text{or } \mathbb{R}^n) \\ f^{(i)} &= \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \\ f_0^{(i)} &= 1 \end{aligned}$$

Andrew Ng

- min 목적식을 만들어보자.

SVM with Kernels

Hypothesis: Given x , compute features $f \in \mathbb{R}^{m+1}$

$$\rightarrow \text{Predict "y=1" if } \theta^T f \geq 0$$

$\theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$

Training:

$$\rightarrow \min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$n = m$
 $\cancel{\theta_0} = m$
 $\rightarrow \theta_0$

- regularization term이 조금은 변형된다.(computationally expensive한 것을 방지하기 위해)

$$\sum_j \theta_j^2 = \theta^T \theta \quad \leftarrow \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix} \quad (\text{ignoring } \theta_0)$$

$$\underline{\theta^T M \theta} \quad \| \theta \|^2$$

- parameter를 어떻게 결정할 것인가?

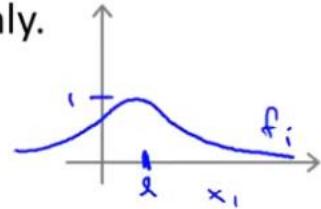
- C의 선택과 시그마 제곱의 선택

SVM parameters:

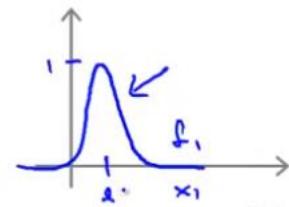
C ($= \frac{1}{\lambda}$). \rightarrow Large C : Lower bias, high variance. (small λ)
 \rightarrow Small C : Higher bias, low variance. (large λ)

σ^2 $\underline{\text{Large } \sigma^2}$: Features f_i vary more smoothly.
 \rightarrow Higher bias, lower variance.

$$\exp\left(-\frac{\|x - \mu^{(i)}\|^2}{2\sigma^2}\right)$$



$\underline{\text{Small } \sigma^2}$: Features f_i vary less smoothly.
Lower bias, higher variance.



Suppose you train an SVM and find it overfits your training data. Which of these would be a reasonable next step? Check all that apply.

Increase C

Un-selected is correct

Decrease C

Correct

Increase σ^2

\rightarrow overfitting이 일어나게 되면 C 를 줄이고, 시그마 제곱을 높게 하면 어느 정도 해결된다.

1.7.5. 실제 SVM의 활용 방법(Using an SVM)

- 실제로 이를 코딩하려고 하지 말고, library를 써라.
- linear kernel(kernel을 사용하지 않는 경우) : feature가 많고, dataset 개수(m)는 적은 경우
- kernel : feature 개수(n)가 작고, data set 개수(m)은 큰 경우가 좋다.

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters θ .



Need to specify:

→ Choice of parameter C.

Choice of kernel (similarity function):

E.g. No kernel ("linear kernel")

Predict "y = 1" if $\underline{\theta^T x} \geq 0$

$$\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0$$

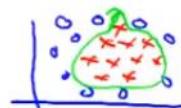
$\rightarrow n$ large, m small $x \in \mathbb{R}^{n+1}$

Gaussian kernel:

$$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right), \text{ where } l^{(i)} = x^{(i)}$$

Need to choose σ^2 .

$x \in \mathbb{R}^n$, n small
and/or m large



- Gaussian kernel을 사용하는 방법

- 1) 함수를 직접 넣어줘야 하는 경우가 많다.
- 2) feature scaling을 해주는 게 좋다.

Kernel (similarity) functions:

function $f = \text{kernel}(x_1, x_2)$

$$f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

return

$x \rightarrow$
 f_1
 f_2
 \vdots
 f_m

→ Note: Do perform feature scaling before using the Gaussian kernel.

$$\begin{aligned} & \boxed{\|x - l\|^2} \\ & \|v\| = x - l \\ & \|v\|^2 = v_1^2 + v_2^2 + \dots + v_n^2 \\ & = (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2 \\ & \cdot \underbrace{1000 \text{ feet}^2}_{1-5 \text{ bedrooms}} \end{aligned}$$

- 3) kernel식이 Mercer's Theorem을 만족해야 한다.

Other choices of kernel

Note: Not all similarity functions $\text{similarity}(x, l)$ make valid kernels.

→ (Need to satisfy technical condition called “Mercer’s Theorem” to make sure SVM packages’ optimizations run correctly, and do not diverge).

Many off-the-shelf kernels available:

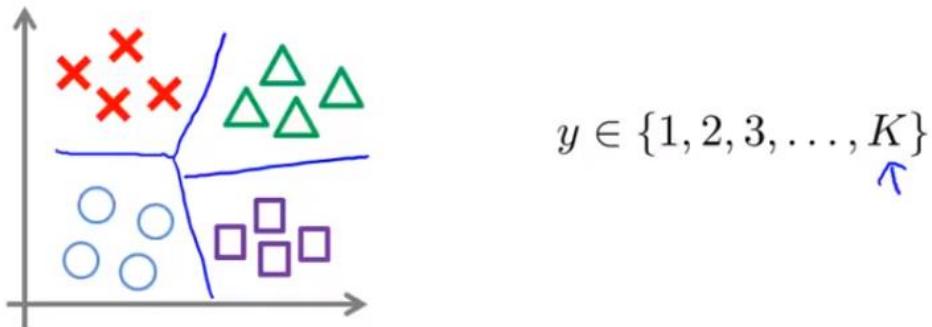
- Polynomial kernel: $k(x, l) = \langle x^T l + \text{constant} \rangle^{\text{degree}}$

$$\langle x^T l \rangle^2, \quad \langle x^T l + 1 \rangle^3, \quad \langle x^T l + 5 \rangle^4$$

- More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...

4) multi-class classification 사용법

Multi-class classification



Many SVM packages already have built-in multi-class classification functionality.

→ Otherwise, use one-vs.-all method. (Train K SVMs, one to distinguish $y = i$ from the rest, for $i = 1, 2, \dots, K$), get $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$
 Pick class i with largest $\underline{(\theta^{(i)})^T x}$

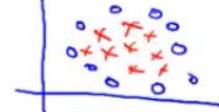
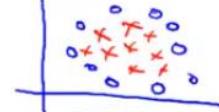
5) SVM vs logistic regression

- kernel SVM은 feature가 작고 dataset이 보통일 때 많이 쓰인다.
- feature가 많고 dataset이 적거나, feature가 적고 dataset이 엄청 많으면 linear kernel 혹은 logistic regression을 쓴다.
 - neural network도 적용이 가능하다. 그런데 neural network가 training할 때 더 느릴 수 있다.
 - SVM은 convex function minimization이기 때문에 global minimum을 찾아준다. local optima 걱정을 할 필요는 없다.

Logistic regression vs. SVMs

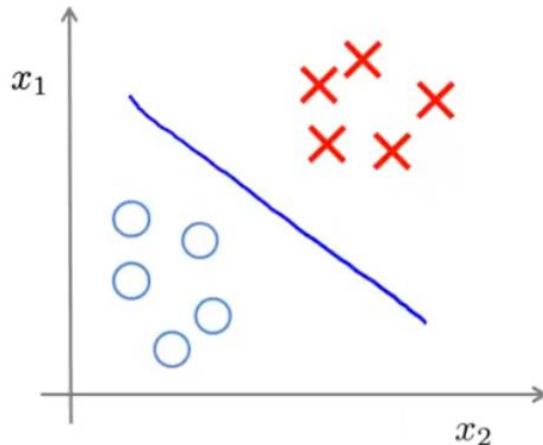
n = number of features ($x \in \mathbb{R}^{n+1}$), m = number of training examples

- If n is large (relative to m): (e.g. $n \geq m$, $n = 10,000$, $m = 10 \dots 1000$)
- Use logistic regression, or SVM without a kernel ("linear kernel")
- If n is small, m is intermediate: ($n = 1-1000$, $m = 10 - 10,000$) ←
 - Use SVM with Gaussian kernel
- If n is small, m is large: ($n = 1-1000$, $m = 50,000+$)
 - Create/add more features, then use logistic regression or SVM without a kernel
- Neural network likely to work well for most of these settings, but may be slower to train.



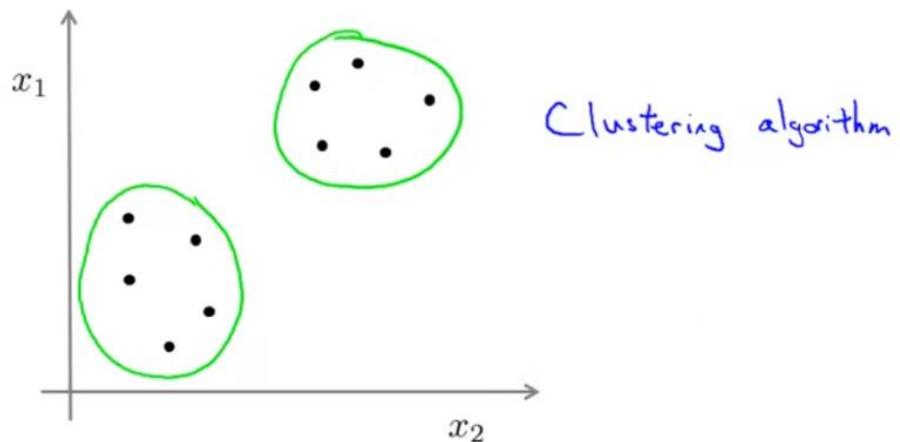
Week 8.**2. Unsupervised learning** ← unlabeled data set

- 복습 : supervised learning

Supervised learning

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$.

- 차이점 : algorithm으로 find some structure for us!

Unsupervised learning

Training set: $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$ ← .

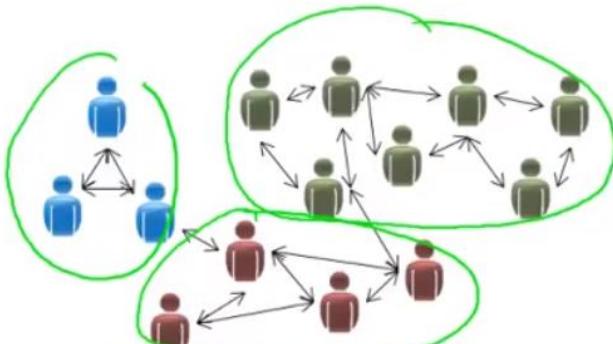
→ clustering하는 방법도 unsupervised 중 하나이다.

- 예

Applications of clustering



→ Market segmentation



→ Social network analysis



Organize computing clusters



Astronomical data analysis

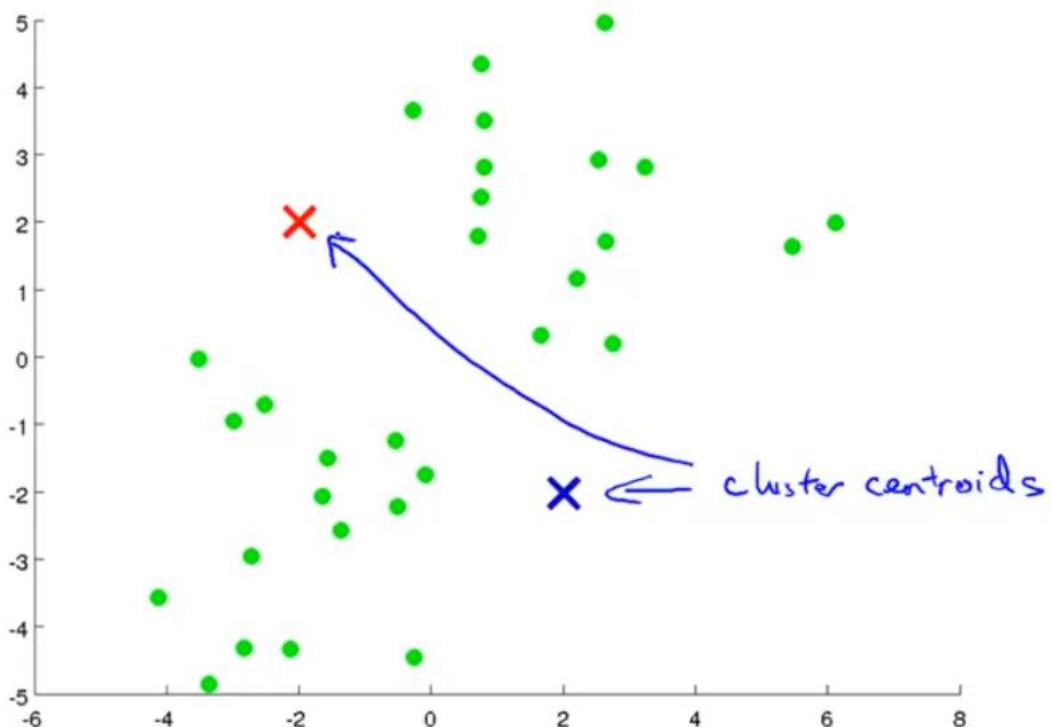
2.1. Unsupervised learning 종류

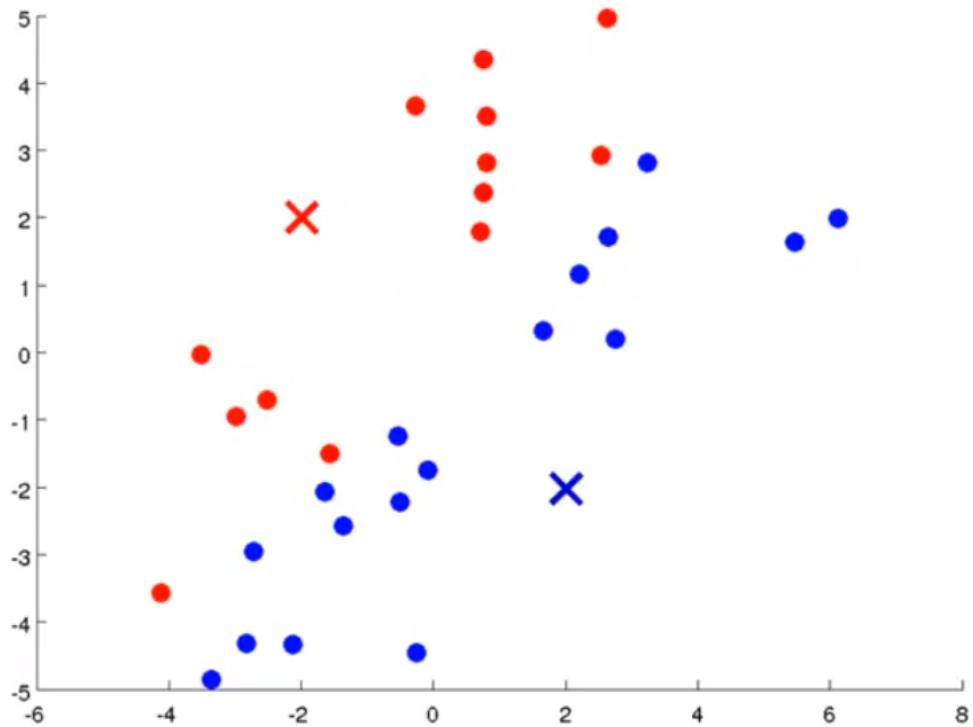
2.1.1. Clustering algorithm : K-means algorithm (the most popular)

(1) K-means Algorithm visualization : 예시로 보는 warm-up

(1단계) cluster assignment step

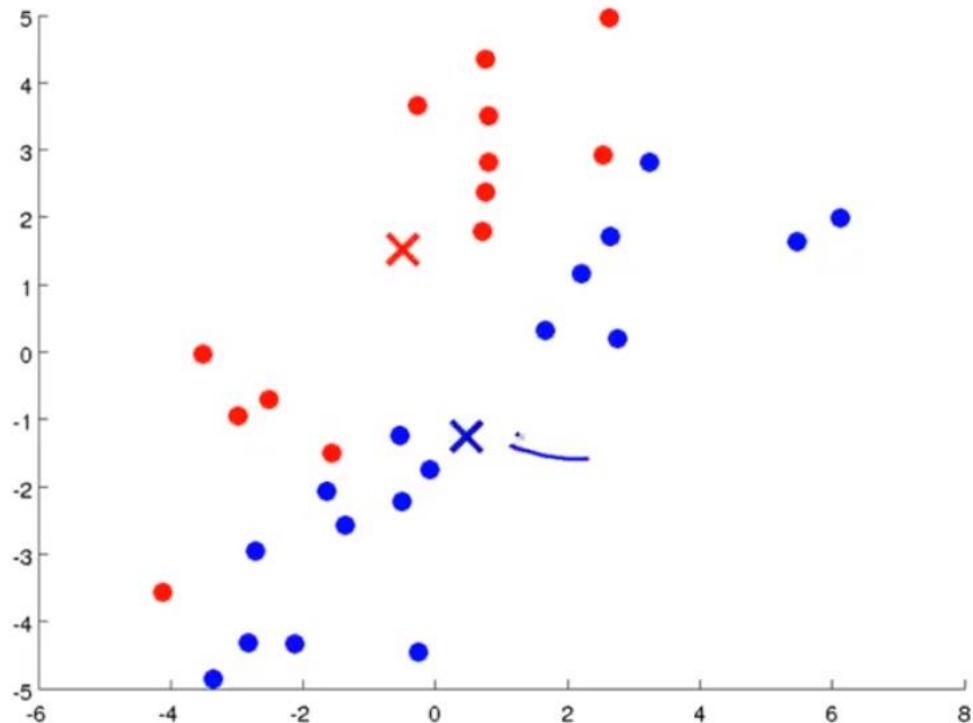
→ 거리 더 가까운 것을 각각 다른 색깔로





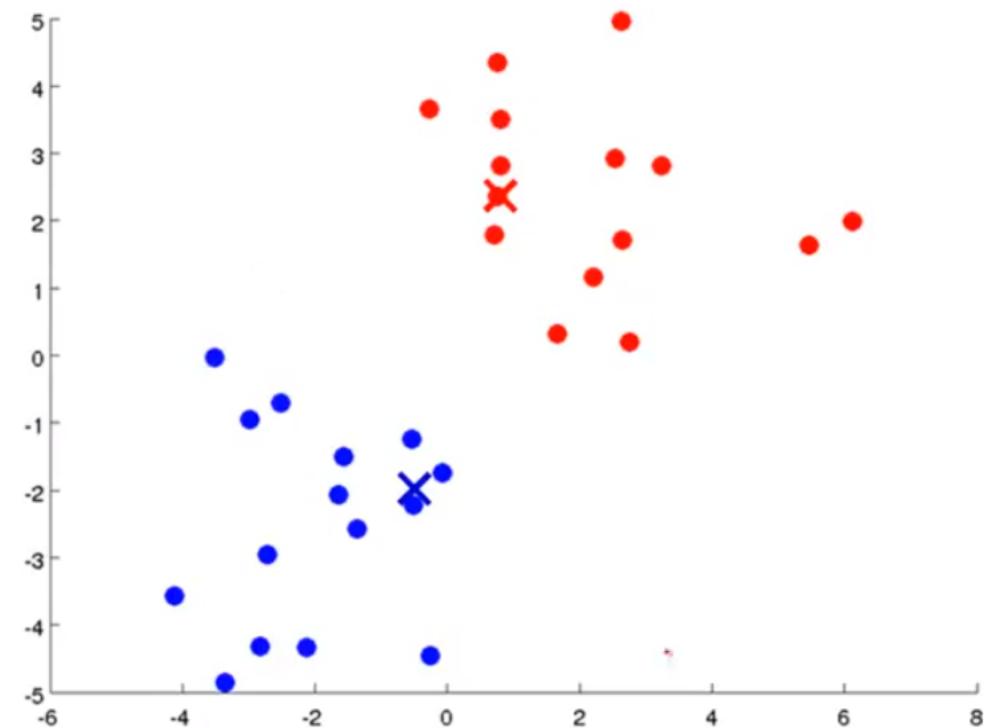
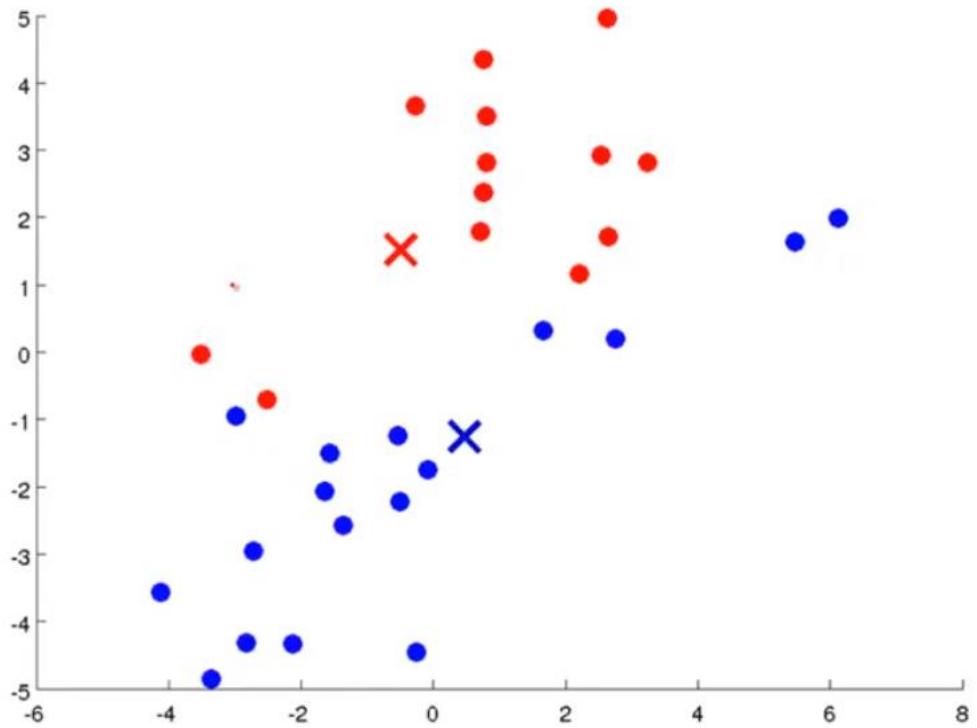
(2단계) move centroids step

- 모든 파란 점을 보고 평균을 계산하여 그 자리로 파란 x를 이동시킨다. 빨간 점도 x, y축을 모두 더하여 n으로 나누고 그 쪽으로 빨간 x를 이동시킨다.

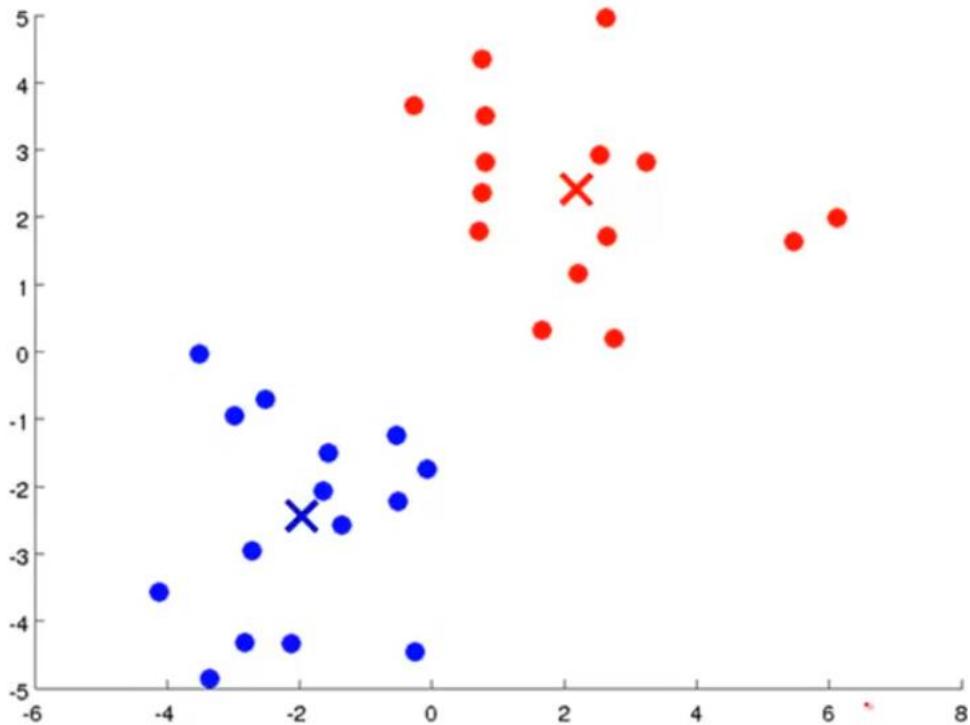


(다시 1단계)

cluster centroid에 거리가 가까운 것을 빨간, 파란색으로 표시한다.



(다시 2단계) move centroid step



→ 완성!!

(2) K-mean algorithm

- 몇 개의 cluster를 설정할 것인가? → k를 정해 (k를 automatically 정하는 알고리즘도 있다.)

Step1. Input : K를 정하는 중요한 단계!

K-means algorithm

Input:

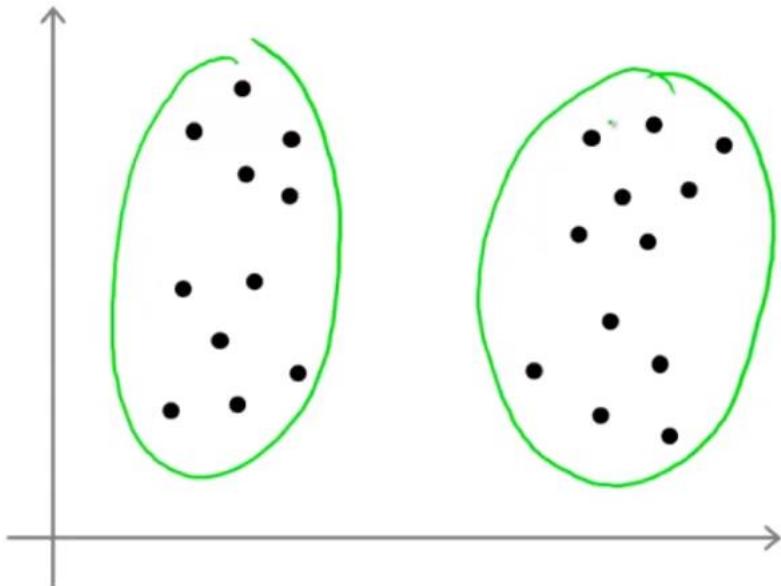
- K (number of clusters) ←
- Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ ←

$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

→ label 0이 없어졌으므로 n개 항으로 이루어진 vector (그 전에는 n+1이었다.)

- capital K (cluster의 개수)를 정하는 방법

(방법 1) visualization을 보고 manually 하는 게 많이 사용된다. → ambiguous 할 수 있다.

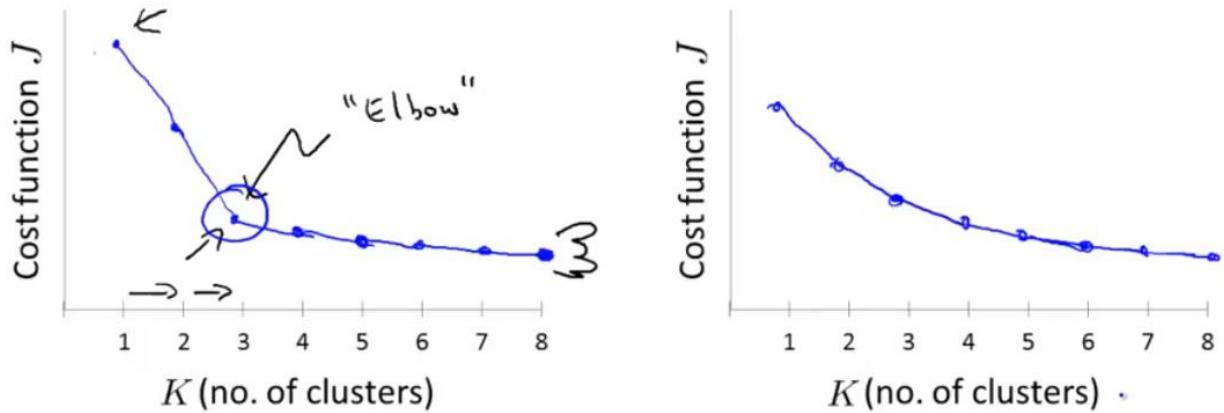


(방법 2) Elbow method도 있다. J cost function이랑 number of cluster를 plotting하는 방법도 있다.

: elbow처럼 꽉 꺾이는 점에서 K 를 정한다. 가파르게 cost가 줄어드는 지점에서! 그러나, 실제로는 꽉 꺾이는 지점이 나타나지 않는 경우가 흔하기 때문에 elbow method가 사용하기가 애매하다.

Choosing the value of K

Elbow method:



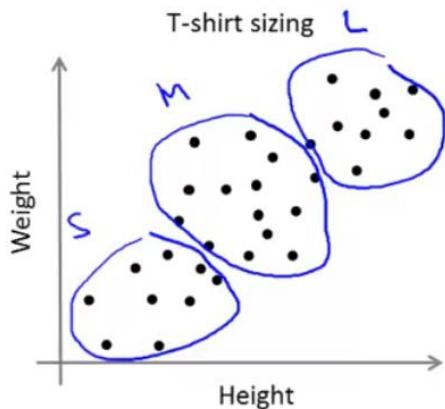
(방법 3) downstream(later) purpose에 부합하는 K 를 고른다. ex) T-shirt segmentation

Choosing the value of K

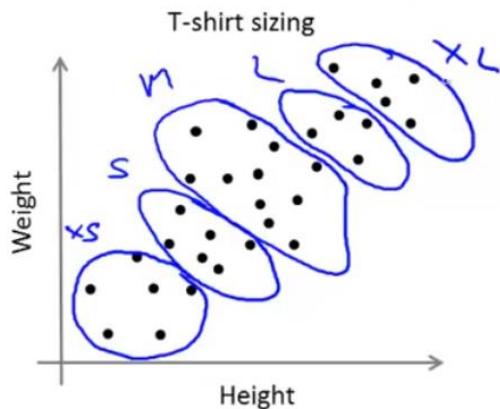
Sometimes, you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.

$K=3 \quad S, M, L$

E.g.



$K=5 \quad XS, S, M, L, XL$



→ 어떤 T-shirt size cluster로 해야지 좋을까? 소비자에게 딱 맞는 shirt를 만드는 게 중요한가, 아니면 규모의 경제를 실현하도록 size option을 줄이는 게 좋을까?

- image compression을 할 때에도, file size를 줄이는 게 중요한가, 아니면 image가 잘 보이는 게 중요한가.

Step2. 다음 algorithm

K-means algorithm

$$\mu_1 \quad \mu_2$$

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

\nwarrow

Cluster assignment step

$\left[\begin{array}{l} \text{for } i = 1 \text{ to } m \\ \quad c^{(i)} := \text{index (from 1 to } K \text{) of cluster centroid} \\ \quad \text{closest to } x^{(i)} \\ \quad \min_{k} \ x^{(i)} - \mu_k\ \\ \text{for } k = 1 \text{ to } K \\ \quad \mu_k := \text{average (mean) of points assigned to cluster } k \end{array} \right]$	$\mu_1 \quad \mu_2$
-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	---------------------

}

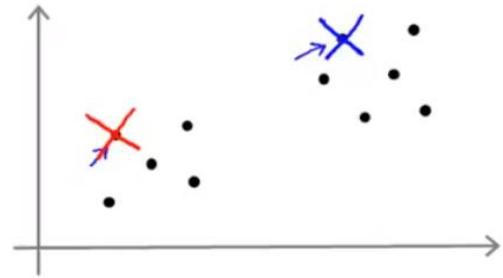
(2-1) How to randomly initialize K cluster centroid → to avoid local optima

- x data set 중에서 random하게 정하는 방법을 강력 추천한다.

Random initialization

Should have $K < m$

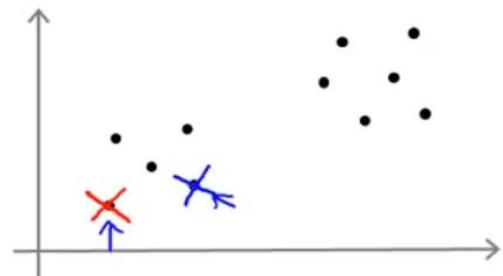
$K=2$



Randomly pick K training examples.

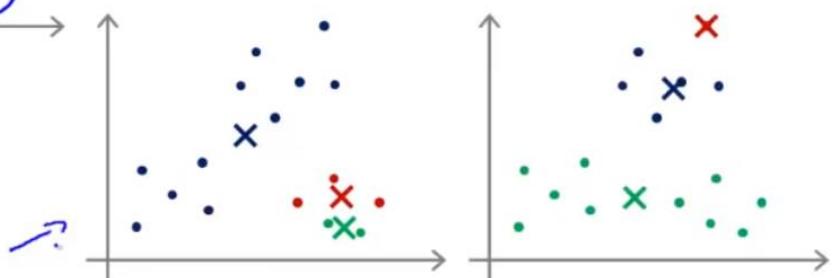
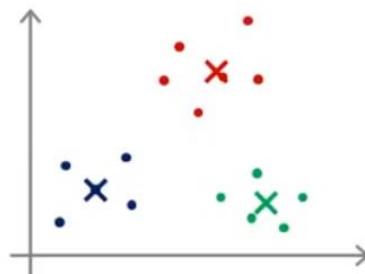
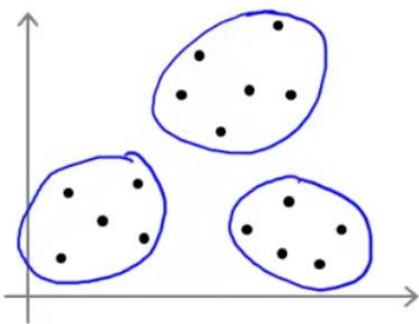
Set μ_1, \dots, μ_K equal to these K examples.

$$\begin{aligned}\mu_1 &= x^{(i)} \\ \mu_2 &= x^{(j)} \\ &\vdots\end{aligned}$$



- local optima of J (cost function)의 함정과 해결

Local optima



→ multiple random initialization

→ k-means를 많이 해서 (100번 이상) 가장 best의 clusterization을 고른다. i 는 iteration의 i 이다. 따라서 몇 번을 찍을 것인가를 의미한다.

Random initialization

```

For i = 1 to 100 { 50 - 1000
    → Randomly initialize K-means.
    Run K-means. Get  $c^{(1)}, \dots, c^{(m)}$ ,  $\mu_1, \dots, \mu_K$ .
    Compute cost function (distortion)
    →  $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$ 
}

```

Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

$K=2-10$ ↑

→ k (cluster의 개수)가 2~10이면 random initialization이 효과가 좋으나, cluster가 100개 이면 사실상 random initialization이 큰 개선 효과가 없다.

(2-2) cluster assignment step : data set x 를 $i = 1$ to m 까지 min 거리하는 μ_k 의 k 로 indexing을 한다.

(2-3) move centroid step :

K-means algorithm

$$\mu_1 \quad \mu_2$$

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

Cluster assignment step

```

for i = 1 to m
     $c^{(i)}$  := index (from 1 to  $K$ ) of cluster centroid
    closest to  $x^{(i)}$ 
        ↪  $\min_{c^{(i)}} \|x^{(i)} - \mu_k\|^2$ 
    ↪  $c^{(1)} = 2, c^{(2)} = 2, c^{(3)} = 2, c^{(4)} = 2$ 
for k = 1 to  $K$ 
    →  $\mu_k$  := average (mean) of points assigned to cluster  $k$ 
        ↪  $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$ 
    }
```

$$\mu_2 = \frac{1}{4} \left[\underline{x^{(1)}} + \underline{x^{(2)}} + \underline{x^{(3)}} + \underline{x^{(4)}} \right] \in \mathbb{R}^n$$

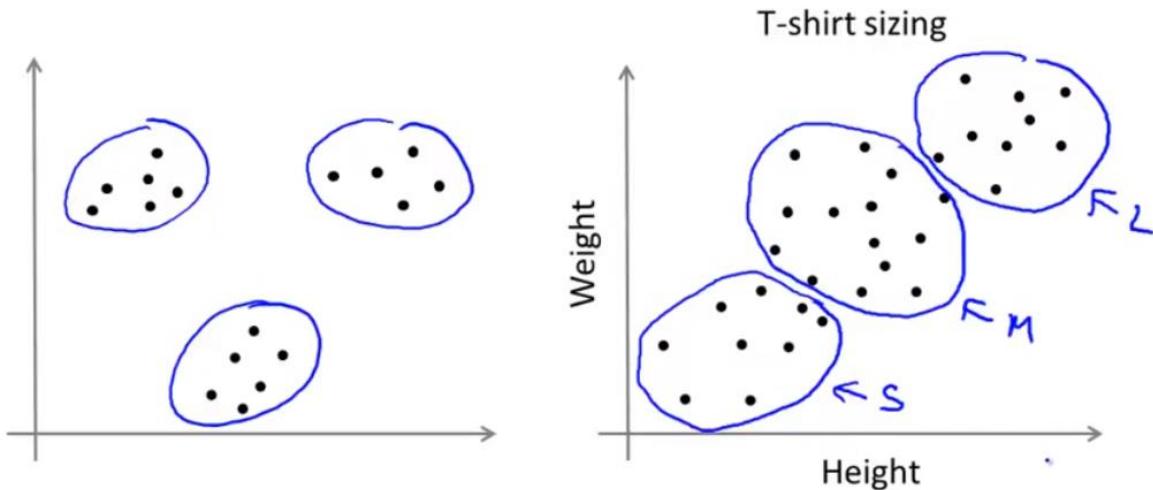
→ cf) 만약에 k 번째 cluster에 아무 data set도 배정되지 않으면 eliminate시키거나, 다시 random하게 μ_k 를 설정해주는 방법이 있다. 그러나 전자가 더 흔하다.

- non well separated cluster도 잘된다는 점.

- T-shirt sizing을 하려는데 clustering이 잘 안 되는 경우, S, M, L을 잘 구분할 수가 없다. 그래도 k-mean을 하면 나름 구분을 할 수 있다. 좋네.

K-means for non-separated clusters

S, M, L



2.1.3. optimization objective

- 목적 : debugging을 할 수 있다. 더 중요하게는 local optimum을 피할 수 있다.

- 변수 :

$$\begin{aligned} \rightarrow c^{(i)} &= \text{index of cluster } (1, 2, \dots, K) \text{ to which example } x^{(i)} \text{ is currently assigned} \\ \rightarrow \mu_k &= \text{cluster centroid } k \quad (\mu_k \in \mathbb{R}^n) \quad \leftarrow \quad k \in \{1, 2, \dots, K\} \\ \mu_{c^{(i)}} &= \text{cluster centroid of cluster to which example } x^{(i)} \text{ has been assigned} \quad x^{(i)} \rightarrow 5 \quad c^{(i)} = 5 \quad \mu_{c^{(i)}} = \mu_5 \end{aligned}$$

- objective

Optimization objective:

$$\rightarrow J(\underbrace{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K}_{}) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\min_{\substack{c^{(1)}, \dots, c^{(m)}, \\ \mu_1, \dots, \mu_K}} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

\rightarrow x와 x가 배정된 cluster의 centroid와의 거리 차이를 minimize하면 된다.

- 아까의 algorithm에서 위의 minimize를 하는 것이다. 즉, cluster assignment step에서는 $c(i)$ 에 대하여 cost function을 minimize하고, move centroid step에서는 μ_k 에 대하여 cost function을 minimize하는 것이다. Cost function은 distortion function이라고도 한다.

K-means algorithm

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

for $i = 1$ to m

$c^{(i)} :=$ index (from 1 to K) of cluster centroid
closest to $x^{(i)}$

*move
centroid*

for $k = 1$ to K

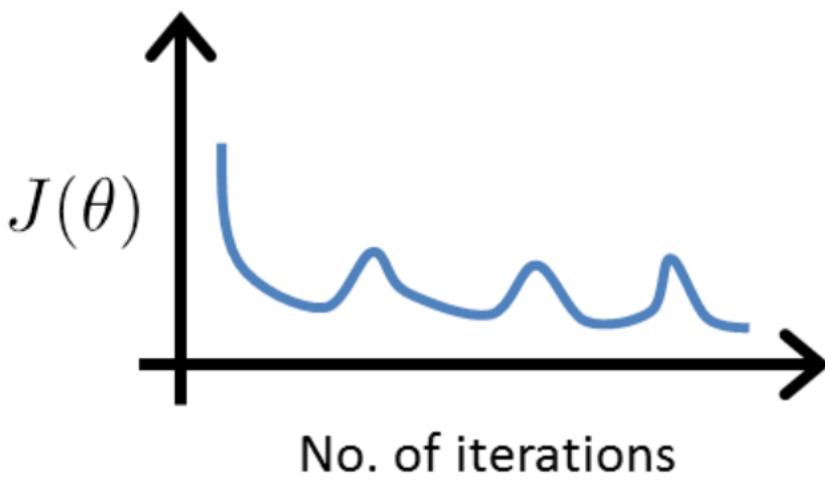
$\mu_k :=$ average (mean) of points assigned to cluster k

}

minimize $J(\dots)$ wpt μ_1, \dots, μ_K

- 만약에 이러면?

Suppose you have implemented k-means and to check that it is running correctly, you plot the cost function $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$ as a function of the number of iterations. Your plot looks like this:



- It is not possible for the cost function to sometimes increase.
There must be a bug in the code.

Correct



Unsupervised learning 종류-2

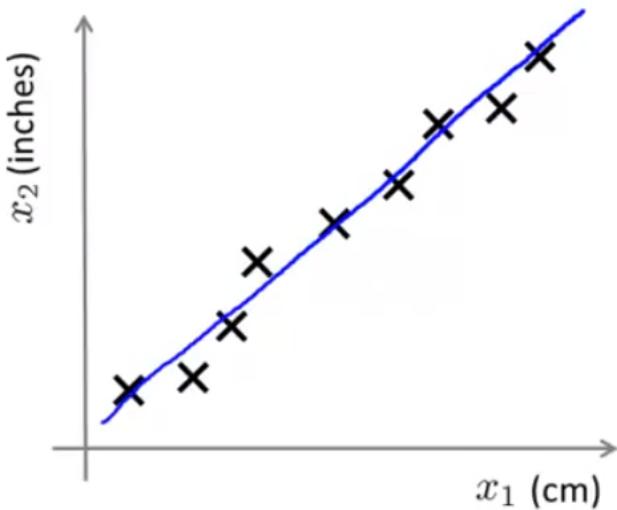
2.2. Dimensionality reduction

- 목적 :

- (1) data compression → speed up algorithms / reduce data storing space
 - (2) Data visualization
- 실제 알고리즘 : Principal Component Analysis

2.2.1. Data compression 원리 (reduce the dimension)

Data Compression

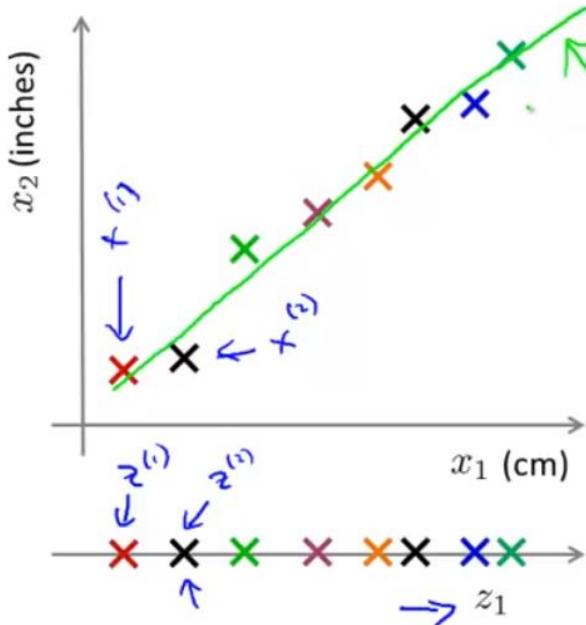


Reduce data from
2D to 1D

→ round off error가 생겨서 왔다갔다 한다.

- 2개의 dimension을 1개의 dimension으로 줄이는 예 : new feature z_1 을 도입하여 x_1, x_2 를 대변할 수 있게 된다.

Data Compression



Reduce data from
2D to 1D

$$x^{(1)} \in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R}$$

$$x^{(2)} \in \mathbb{R}^2 \rightarrow z^{(2)} \in \mathbb{R}$$

$$\vdots$$

$$x^{(m)} \rightarrow z^{(m)}$$

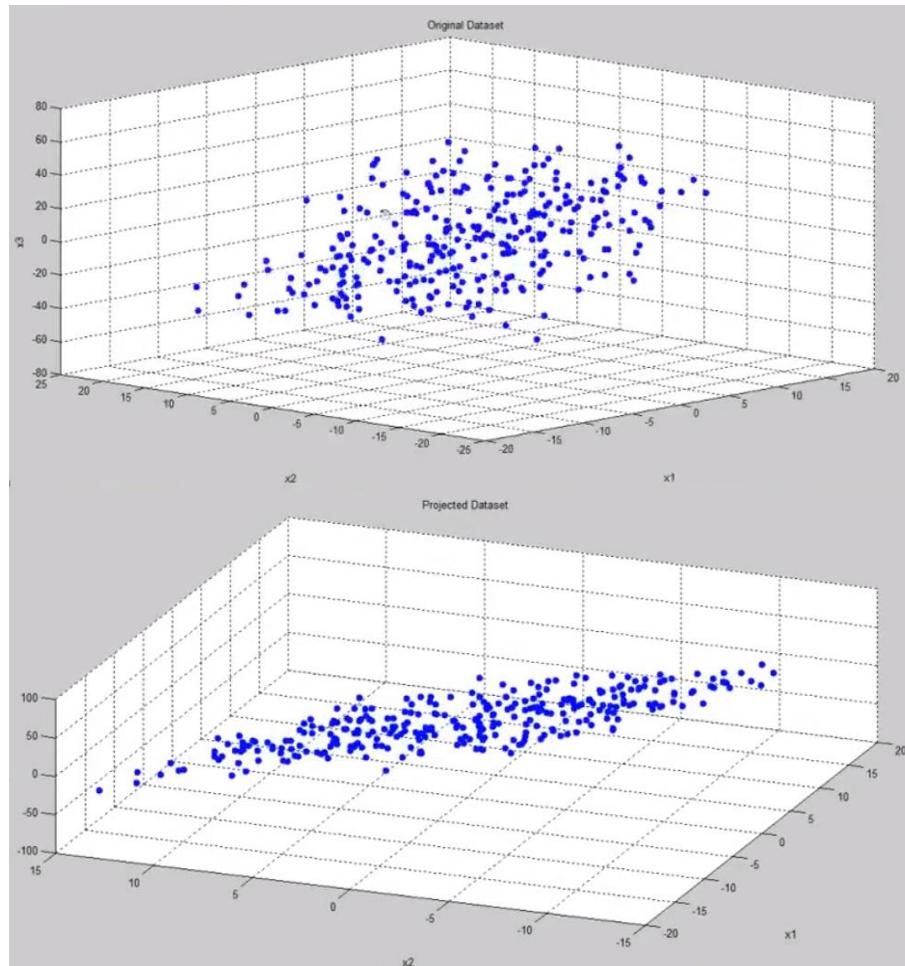
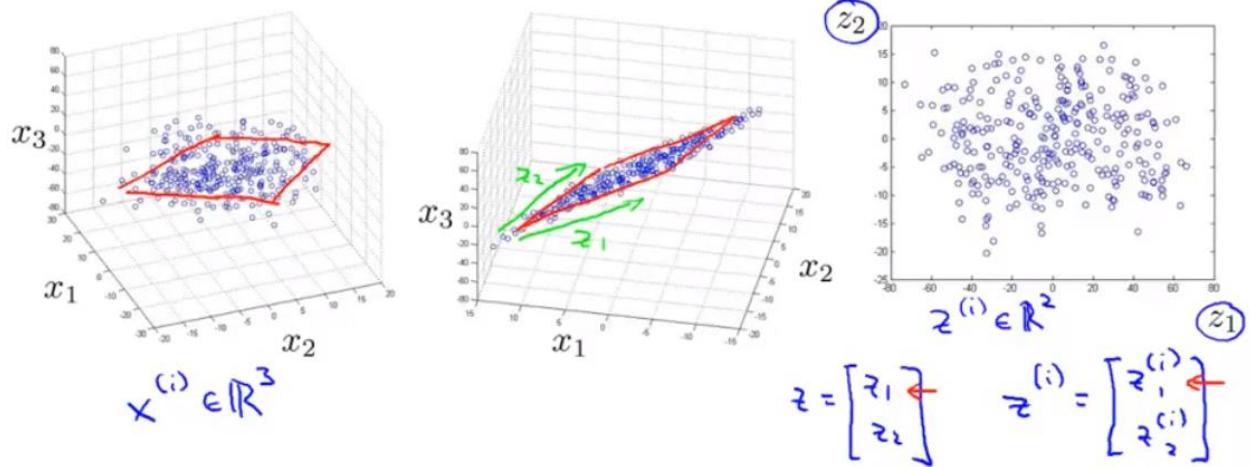
- 3D에서 2D로 reduce하는 예

: project x_3 to the x_1-x_2 plane

Data Compression

$10000 \rightarrow 1000$

Reduce data from 3D to 2D



- 차원을 줄이는 것이다! data set의 개수를 줄이는 것이 아니라!

Suppose we apply dimensionality reduction to a dataset of m examples $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$, where $x^{(i)} \in \mathbb{R}^n$. As a result of this, we will get out:

- ➊ A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(k)}\}$ of k examples where $k \leq n$.
- ➋ A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(k)}\}$ of k examples where $k > n$.
- ➌ A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$ of m examples where $z^{(i)} \in \mathbb{R}^k$ for some value of k and $k \leq n$.
- ➍ A lower dimensional dataset $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$ of m examples where $z^{(i)} \in \mathbb{R}^k$ for some value of k and $k > n$.

2.2.2. Data visualization \leftarrow dimension reduction의 목적-2

- 원래 데이터 (x vector의 차원은 50)

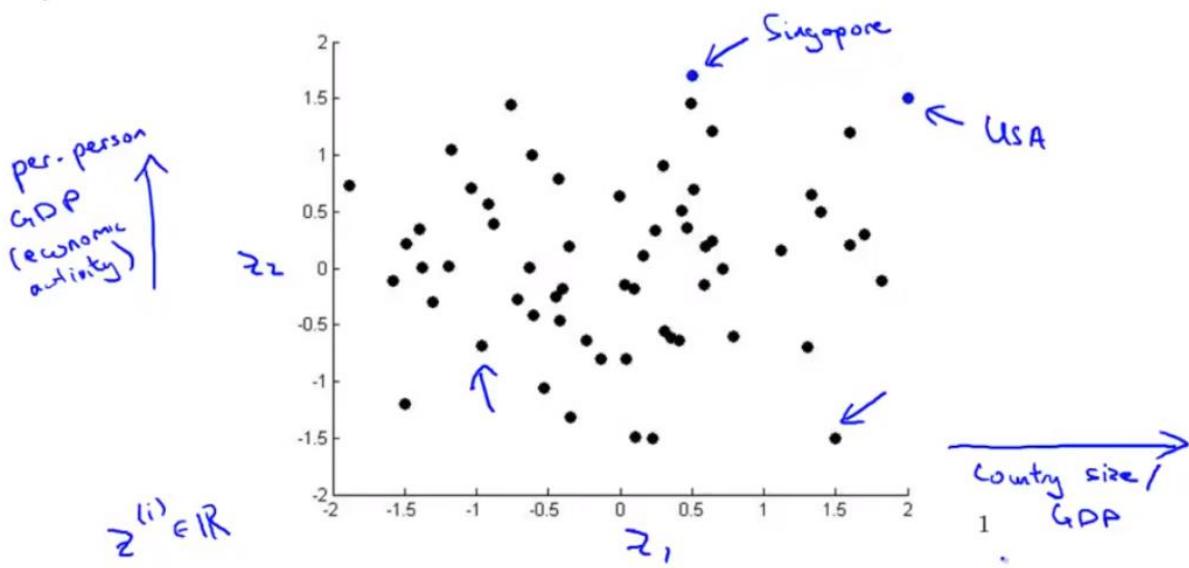
Data Visualization

Country	x_1 GDP (trillions of US\$)	x_2 Per capita GDP (thousands of intl. \$)	x_3 Human Development Index	x_4 Life expectancy	x_5 Poverty Index (Gini as percentage)	x_6 Mean household income (thousands of US\$)	...
	$x \in \mathbb{R}^{50}$	$x^{(i)} \in \mathbb{R}^{50}$					
Canada	1.577	39.17	0.908	80.7	32.6	67.293	...
China	5.878	7.54	0.687	73	46.9	10.22	...
India	1.632	3.41	0.547	64.7	36.8	0.735	...
Russia	1.48	19.84	0.755	65.5	39.9	0.72	...
Singapore	0.223	56.69	0.866	80	42.5	67.1	...
USA	14.527	46.86	0.91	78.3	40.8	84.3	...
...

- 가공한 data (z vector의 차원은 2)

Data Visualization

Country	z_1	z_2	$z^{(i)} \in \mathbb{R}^2$
Canada	1.6	1.2	
China	1.7	0.3	Reduce data
India	1.6	0.2	from 500
Russia	1.4	0.5	to 2D
Singapore	0.5	1.7	
USA	2	1.5	
...	



→ 질문!! 그런데 어떻게 2개로 줄이지????? PCA 대박!!!

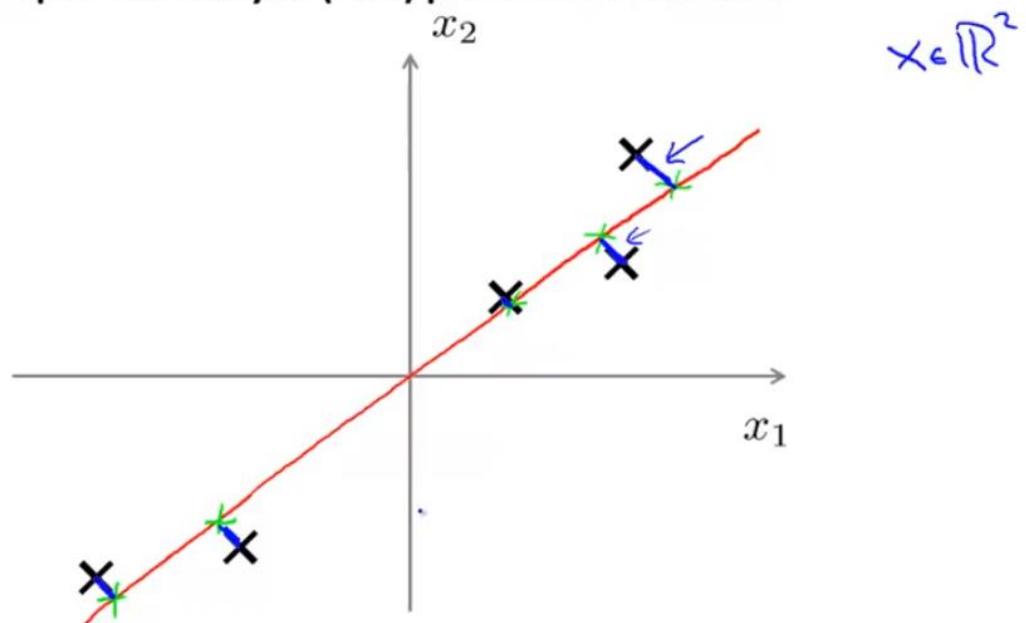
2.3. (중요) Principle component analysis (PCA)

2.3.1. 문제정의 : PCA Problem formulation

→ find lower dimensional surface by projecting data

- 2D 예시 : 2차원 데이터가 있으면, lower dimensional surface(line) = projection surface / projection line을 찾아내서 dimension을 줄인다. 그러면 projection error (square projection error)를 minimize한다. 우선은, 이를 하기 전에 feature scaling / mean normalization을 한다.(평균을 1, range를 -3, 3 사이로 맞춘다.)

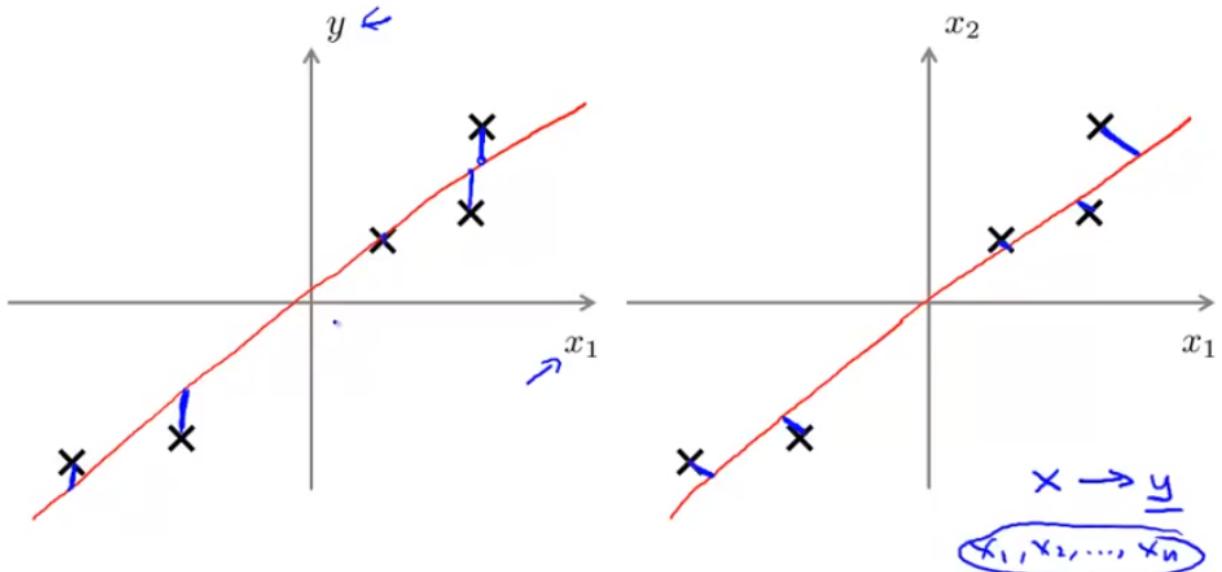
Principal Component Analysis (PCA) problem formulation



Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error.

cf) linear regression과의 차이점 : linear regression은 y 축으로 오차(y 를 기준으로 차이)를 구하지만, PCA는 projection을 한다. 즉, 특별한 y 가 없다. 이것이 바로 unsupervised와의 가장 큰 차이점이다.

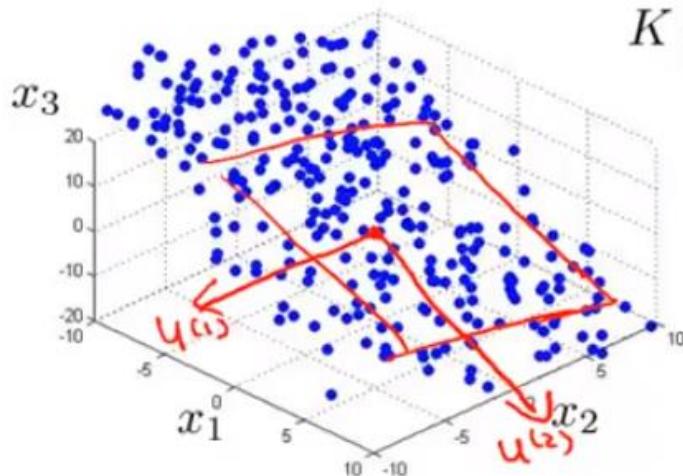
PCA is not linear regression



→ normalization을 했기 때문에 PCA는 (0,0)을 지닌다. mean이 0이기 때문이다.

$$3D \rightarrow 2D$$

$$K = 2$$



Reduce from n-dimension to k-dimension: Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

→ k 개의 vector를 구한다.

2.3.2. 해결 알고리즘 : PCA algorithm

(1) Data preprocessing step : mean normalization하고 feature scaling한다.

Data preprocessing

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)} \leftarrow$

Preprocessing (feature scaling/mean normalization):

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

Replace each $x_j^{(i)}$ with $x_j - \mu_j$.

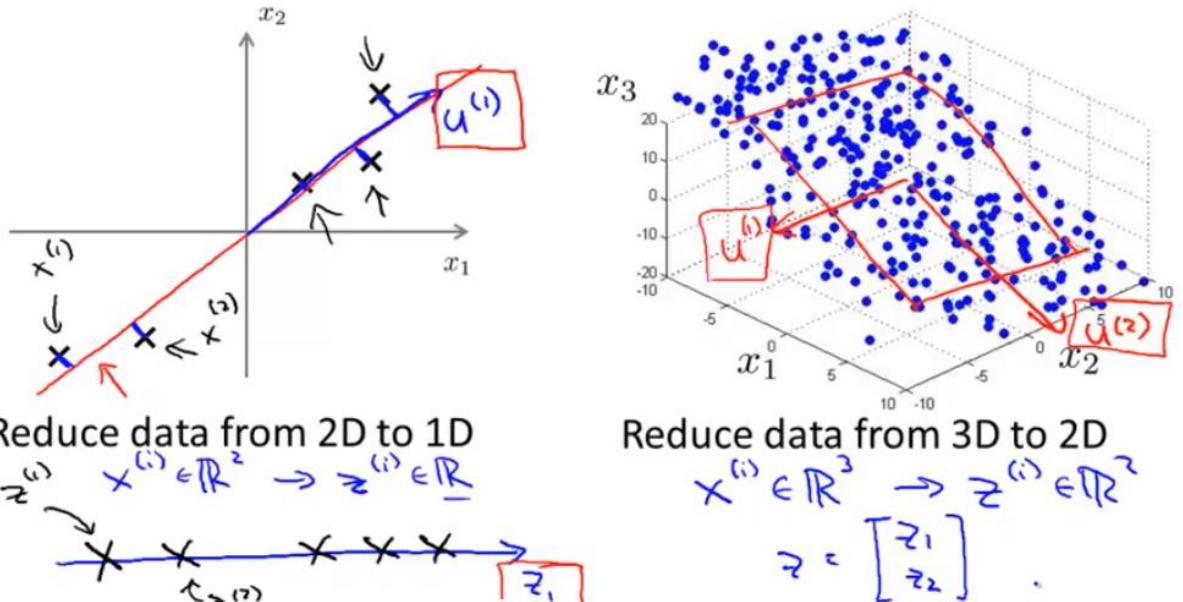
If different features on different scales (e.g., x_1 = size of house, x_2 = number of bedrooms), scale features to have comparable range of values.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

(2) reduce dimension

- 먼저 : $u(1)$ vector를 찾아야지. 3D data면 $u(1), u(2)$ vector를 찾고
- 다음 : z_1 을 찾아야지

Principal Component Analysis (PCA) algorithm



- 수학적 proof : Best low rank approximation을 해준다.

Solving the linear least-squares problem with the SVD

- ▶ $A = U\Sigma V^T$ where $U \in \mathbb{R}^{m \times r}$, $\Sigma \in \mathbb{R}^{r \times r}$, $V \in \mathbb{R}^{n \times r}$.
- ▶ Assume A has rank n (i.e., $r = n$)
- ▶ Solve $Ax \approx b$:

$$\begin{aligned} x &= (A^T A)^{-1} A^T b \\ &= ((U\Sigma V^T)^T (U\Sigma V^T))^{-1} (U\Sigma V^T)^T b \\ &= (V \Sigma^T U^T U \Sigma V^T)^{-1} V \Sigma^T U^T b \\ &= (V \Sigma \Sigma V^T)^{-1} V \Sigma U^T b \\ &= ((V^T)^{-1} (\Sigma \Sigma)^{-1} V^{-1}) V \Sigma U^T b \\ &= V^T \Sigma^{-1} \Sigma^{-1} \Sigma U^T b \\ &= \underline{V^T \Sigma^{-1} U^T b} \end{aligned}$$

- 알고리즘 : n dimension을 k dimension으로 줄이고 싶을 때

<PCA에서 K를 설정하는 방법>

- PCA 분석에서의 K 설정 (number of principal components)
- rule of thumb

Choosing k (number of principal components)

Average squared projection error: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2$

Total variation in the data: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$

Typically, choose k to be smallest value so that

$$\rightarrow \frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

“99% of variance is retained”

→ 이 표현의 뜻은 1%(0.01) 보다 비율이 작다는 것이다. 95% (0.05), 90% (0.1) of variance retained도 용인된다.

- 알고리즘

- 비효율적인 알고리즘 : $k=1$ 부터 $k=2$ 로 놀여가면서 ratio가 0.01보다 작은지 본다. 가장 마지막 k 를 채택하는 방법이 있다.

Algorithm:

Try PCA with $k = 1$ ~~$k=2$~~ ~~$k=3$~~ $k=4$

Compute $U_{reduce}, z^{(1)}, z^{(2)}, \dots, z^{(m)}$, $x_{approx}^{(1)}, \dots, x_{approx}^{(m)}$

Check if

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01?$$

$k = 17$

- 한방에 계산하는 알고리즘 : $[U, S, V] = svd(\Sigma)$

질문!!! 왜 같은 거지???

Choosing k (number of principal components)

Algorithm:

Try PCA with $k = 1$ ~~$k=2$~~ ~~$k=3$~~ ~~$k=4$~~

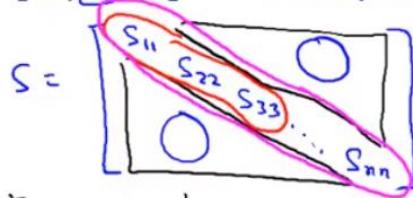
Compute $U_{reduce}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{approx}^{(1)}, \dots, x_{approx}^{(m)}$

Check if

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01?$$

$\underline{k=17}$

$\rightarrow [U, S, V] = svd(Sigma)$



For given k

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \leq 0.01$$

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

$\underline{k=3}$

→ 이렇게 하면 k 만 바꿔가면서 해보면 된다. svd는 한번만 소환하면 되기 때문에 편하다.

Choosing k (number of principal components)

$\rightarrow [U, S, V] = svd(Sigma)$

Pick smallest value of k for which

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

$\underline{k=100}$

(99% of variance retained)

Principal Component Analysis (PCA) algorithm

Reduce data from n -dimensions to k -dimensions

Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)}) (x^{(i)})^T$$

Compute "eigenvectors" of matrix Σ :

$$[U, S, V] = svd(Sigma);$$

→ (혁동님) covariance matrix이지 않을까! x

→ sigma라는 변수에 시그마 행렬을 assgin한다. Singular value decomposition 혹은 eig(Sigma)와 거의 같다.

svd(Sigma)가 더 stable하니까 이것을 쓰도록. 두 값이 같은 이유는 symmetric positive definite 때문이다. Sigma 행렬은 $n \times n$ matrix이다.

$$U = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(m)} \\ | & | & | & \dots & | \end{bmatrix}$$

$U \in \mathbb{R}^{n \times n}$
 $u^{(1)}, \dots, u^{(k)}$

→ U matrix는 $n \times n$ 인데, 그 중에서 k 개의 column vector만 선택하면 된다. → U reduce
 Z 행렬을 새롭게 정의할 수도 있다.

From $[U, S, V] = svd(Sigma)$, we get:

$$\rightarrow U = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ | & | & | & \dots & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$$

$$z = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(k)} \\ | & | & | & \dots & | \end{bmatrix}^T$$

$z \in \mathbb{R}^k$ U_{reduce}

$$x = \begin{bmatrix} -(u^{(1)})^T \\ \vdots \\ -(u^{(k)})^T \end{bmatrix} \quad \begin{matrix} \checkmark \\ \times \\ \sim \\ n \times 1 \end{matrix}$$

$k \times n$ $k \times 1$

- 정리 :

Principal Component Analysis (PCA) algorithm summary

→ After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

$$\text{Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T$$

→ $[U, S, V] = svd(Sigma)$;

→ $U_{\text{reduce}} = U(:, 1:k)$;

→ $z = U_{\text{reduce}}' * x$;

$x \in \mathbb{R}^n$ $\cancel{x_0 = 1}$

→ 질문!!! 왜 $x(i)$ 랑 $x(i)T$ 를 곱한 것이 Sigma가 되는 걸까?

→ 결론적으로 얻는 것은 z vector 들을 모은 Z 행렬이다.

④ $z_j = (u^{(j)})^T x$

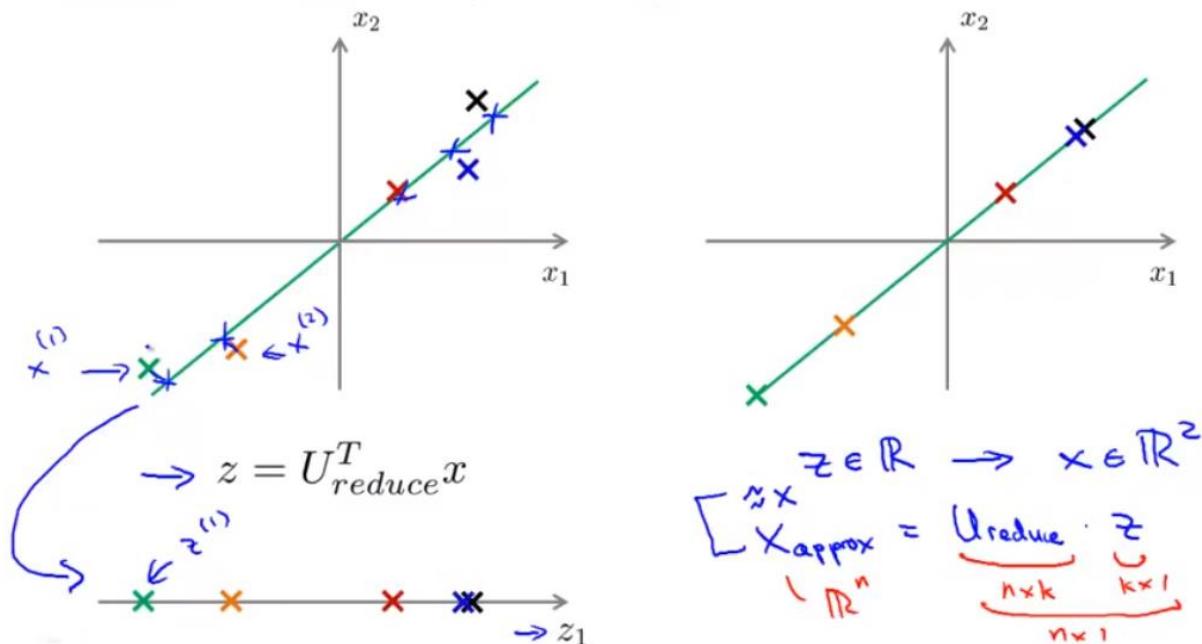
Correct

2.3.3. PCA Application

(1) Reconstruction from Compressed representation

- 문제 정의 : z가 1차원인데, 다시 x (2차원)으로 되돌아 갈 수 있을까?

Reconstruction from compressed representation



- 문제 해결 : $X_{\text{approx.}} = U_{\text{reduce}} * Z$

(2) (가장 많이 사용) PCA application-2 : Supervised learning speedup

→ 원래는 unsupervised에 사용되었는데!! supervised에 적용이 된다!!! 대단!!!!

: extract unlabeled dataset (x vectors)

→ PCA를 통해서 reduced unlabeled dataset (z vectors)

→ new training set로!!! 만든다. (labeled)

Supervised learning speedup

$\rightarrow (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

Extract inputs:

Unlabeled dataset: $x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{10000} \leftarrow$
 $\downarrow \text{PCA}$

$x^{(1)} \in \mathbb{R}^{10000}$
 $\boxed{100 \times 100}$

$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^{1000} \leftarrow$

New training set:

$(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})$

→ dimension이 굉장히 줄어든 z data set을 얻을 수 있다.

- 핵심 : $x \rightarrow z$ (by U_{reduce})로 mapping이 핵심이다. (PCA를 통한) 단, training set을 통해서만! 새로운 example은 U_{reduce} 를 통해서 z로 변환한 다음에 넣으면 된다.

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA $\xrightarrow{x \rightarrow z}$
only on the training set. This mapping can be applied as well to
the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test sets.

→ 예를 들면 classification을 하면 PCA로 data set의 dimension을 줄인다고 하더라도 크게 정확도가 떨어지지는 않는다.

- 목적에 따른 K의 설정

Application of PCA

- Compression

- Reduce memory/disk needed to store data
- Speed up learning algorithm ←

Choose k by % of variance retain

- Visualization

$k=2$ or $k=3$

2.3.4. 주의 사항 (Bad application)

(1) Don't use PCA to reduce overfitting

- 만약에 PCA를 하면 y 에 있는 valuable information이 없어질 가능성이 있다. 물론 99% 이상의 variance가 obtained 되면 큰 문제는 없겠지만 99%를 넘으면 regularization을 쓰는게 overfitting 문제를 막는데 더욱 부작용이 적다.

Bad use of PCA: To prevent overfitting

→ Use $\underline{z^{(i)}}$ instead of $\underline{x^{(i)}}$ to reduce the number of features to $\underline{k} < \underline{n}$.

Thus, fewer features, less likely to overfit.

Bad!

This might work OK, but isn't a good way to address overfitting. Use regularization instead.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \boxed{\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2}$$

(2) Always try to do without PCA first.

- PCA 없이 해보고, 만약에 PCA가 필요하면(data compression / speed up, visualization) 하는 게 낫다.

PCA is sometimes used where it shouldn't be

Design of ML system:

- - Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- - ~~Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$~~
- - Train logistic regression on $\{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\}$
- - Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_{\theta}(z)$ on $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

→ How about doing the whole thing without using PCA?

→ Before implementing PCA, first try running whatever you want to do with the original/raw data $\underline{x^{(i)}}$. Only if that doesn't do what you want, then implement PCA and consider using $z^{(i)}$.

III. Special topics of Machine learning

3. Anomaly Detection

3.1. What is Anomaly detection?

- 특징 : unsupervised, supervised의 두 가지 특징을 모두 가지고 있다.

- Anomaly Detection이 무엇인가?

(예시1) Fraud Detection

Anomaly detection example

→ Fraud detection:

→ $x^{(i)}$ = features of user i 's activities

→ Model $p(x)$ from data.

x_1
 x_2
 x_3
...

Identify unusual users by checking which have $p(x) < \epsilon$

→ flag users that are behaving unusually.

(예시2) Manufacturing

1) 문제 : 새로운 엔진을 만들었을 때, 그것을 소비자에게 보내도 될 만한 것인가?

Anomaly detection example

Aircraft engine features:

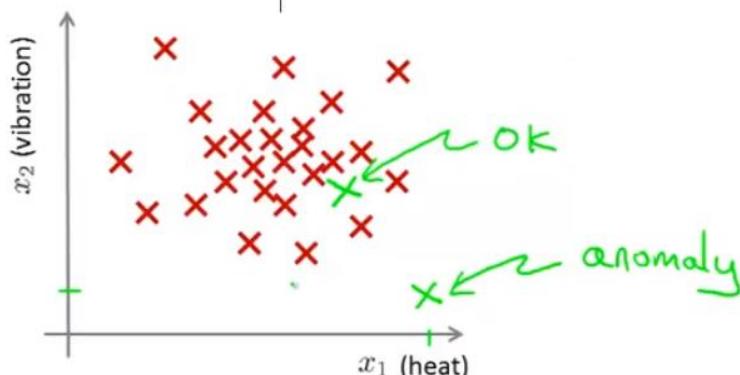
→ x_1 = heat generated

→ x_2 = vibration intensity

...

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

New engine: x_{test}



2) 해결 : Model $P(x)$ 를 새로 만든다! $P(x) < \epsilon$ 이면 anomaly라고 판단하자!

(예시3) Monitoring computers in a data center

→ Monitoring computers in a data center.

→ $x^{(i)}$ = features of machine i

x_1 = memory use, x_2 = number of disk accesses/sec,

x_3 = CPU load, x_4 = CPU load/network traffic.

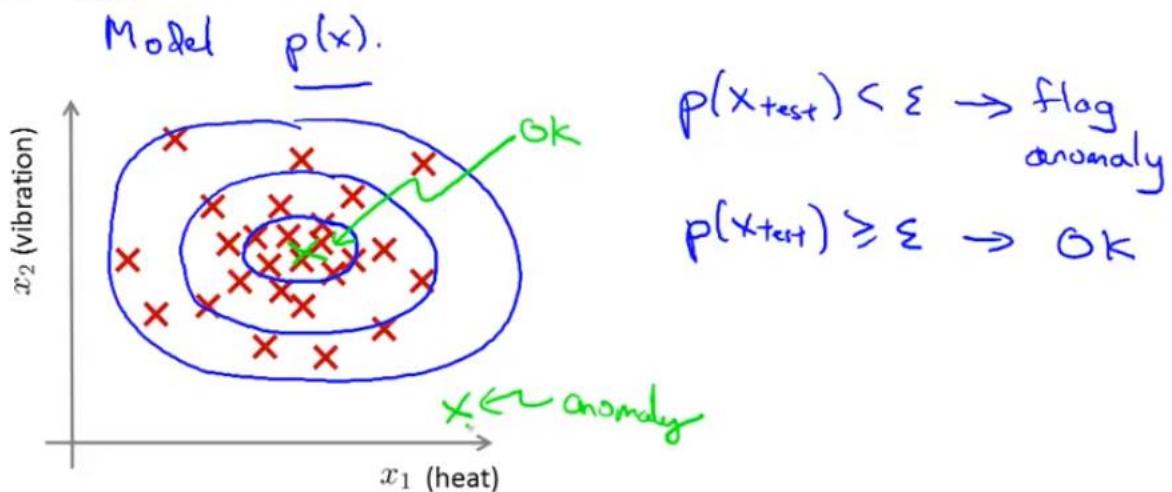
...

$p(x) < \epsilon$

Density estimation

→ Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

→ Is x_{test} anomalous?



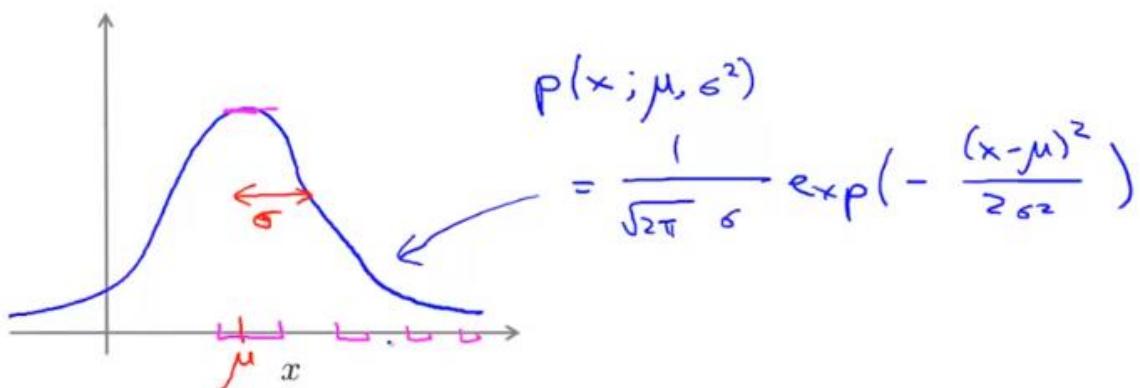
* Gaussian distribution

Gaussian (Normal) distribution

Say $x \in \mathbb{R}$. If x is a distributed Gaussian with mean μ , variance σ^2 .

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

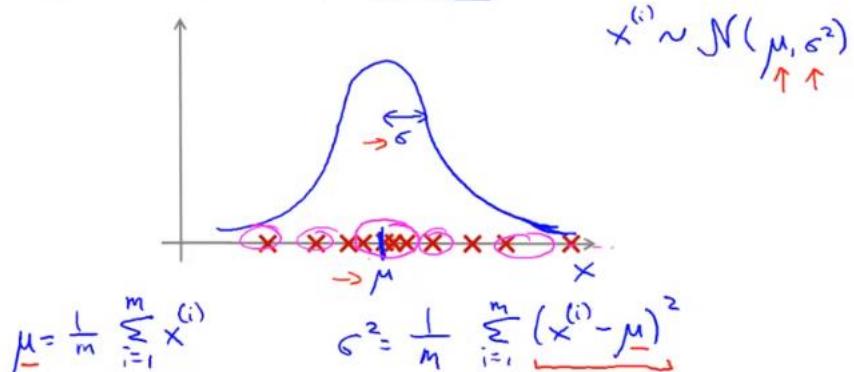
↖ "distributed as"



- Parameter estimation을 해보자 : dataset이 주어지면 어떻게 될 것인가? 이것이 정규분포에서 추출된 것이라고 가정해보고, 정규분포의 평균과 분산을 계산해보면 된다.

Parameter estimation

→ Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ $x^{(i)} \in \mathbb{R}$



→ maximum likelihood estimation에 해당한다.

- Density estimation : 확률 분포를 정규분포를 예측해보자.

Training set : $\{x(1), \dots, x(m)\}$

Each example is $x \in \mathbb{R}^n$

independent assumption

하면 곱해서 확률을 계산할 수 있다.

$$P(x) = P(x_1; \mu_1, \sigma^2) \times P_2 \times P_3 \times P_4$$

Density estimation

→ Training set: $\{x^{(1)}, \dots, x^{(m)}\}$

Each example is $x \in \mathbb{R}^n$

$$\begin{aligned} p(x) &= p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) p(x_3; \mu_3, \sigma_3^2) \dots p(x_n; \mu_n, \sigma_n^2) \\ &= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) \end{aligned}$$

$x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$
 $x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$
 $x_3 \sim \mathcal{N}(\mu_3, \sigma_3^2)$

$\sum_{i=1}^n i = 1+2+3+\dots+n$
 $\prod_{i=1}^n i = 1 \times 2 \times 3 \times \dots \times n$

- 계산 방식 :

(1) 일단 주어진 test set으로 평균과 표준편차를 각 독립변수별로 계산한다.

(2) 새로운 것이 anomaly인지 계산한다. $P = P_1 * P_2 < 0.01$ 로 판단하면 된다.

Anomaly detection algorithm

→ 1. Choose features x_i that you think might be indicative of anomalous examples. $\{x^{(1)}, \dots, x^{(m)}\}$

→ 2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\begin{aligned} \rightarrow \mu_j &= \frac{1}{m} \sum_{i=1}^m x_j^{(i)} & p(x_j; \mu_j, \sigma_j^2) & \mu_1, \mu_2, \dots, \mu_n \\ \rightarrow \sigma_j^2 &= \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2 & \uparrow \quad \uparrow & \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m x^{(i)} \end{aligned}$$

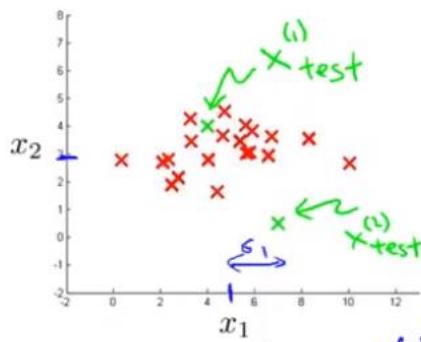
→ 3. Given new example x , compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if $p(x) < \varepsilon$

- 예시 : 성질 변수 x 가 2개인 경우의 예를 보자.

Anomaly detection example

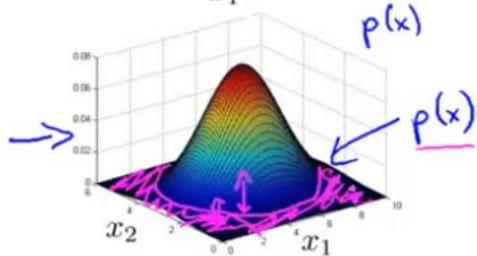


$$\begin{aligned} \mu_1 &= 5, \sigma_1^2 = 4 \\ \mu_2 &= 3, \sigma_2^2 = 1 \end{aligned}$$

$$\rightarrow p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2)$$

$$p(x_1; \mu_1, \sigma_1^2)$$

$$p(x_2; \mu_2, \sigma_2^2)$$



$$\begin{aligned} \varepsilon &= 0.02 \\ p(x_{test}^{(1)}) &= 0.0426 \geq \varepsilon \\ p(x_{test}^{(2)}) &= 0.0021 < \varepsilon \end{aligned}$$

→ 이 핑크색으로 표시한 영역에 example이 있게 되면 anomaly가 된다.

- Developing anomaly detection system을 하려면 evaluate를 해서 더 나아지고 있는지 판단할 수 있어야 한다.

The importance of real-number evaluation

When developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our learning algorithm.

- Assume we have some labeled data, of anomalous and non-anomalous examples. ($y = 0$ if normal, $y = 1$ if anomalous).
- Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples/not anomalous)
- Cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$
- Test set: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

→ 특이사항 : training set에 anomalous한 값들이 조금은 있어도 된다. 그대로 대부분은 normal한 예여야 한다.

- Training set / cross validation set / test set의 구분

Aircraft engines motivating example

- 10000 good (normal) engines
- 20 flawed engines (anomalous) $\frac{2}{10000} = \frac{1}{500}$
- Training set: 6000 good engines ($y=0$) $p(x) = p(x_1; \mu_1, \sigma_1^2) \dots p(x_n; \mu_n, \sigma_n^2)$
CV: 2000 good engines ($y=0$), 10 anomalous ($y=1$)
Test: 2000 good engines ($y=0$), 10 anomalous ($y=1$)

Alternative:

Training set: 6000 good engines

CV: 4000 good engines ($y=0$), 10 anomalous ($y=1$)

Test: 4000 good engines ($y=0$), 10 anomalous ($y=1$)

→ CV랑 test set을 같은 set으로 하는 것은 문제가 있을 수 있다.

- Train, Test, Validation set을 사용하는 방법 (<https://www.youtube.com/watch?v=GtLe9Z2No28>)

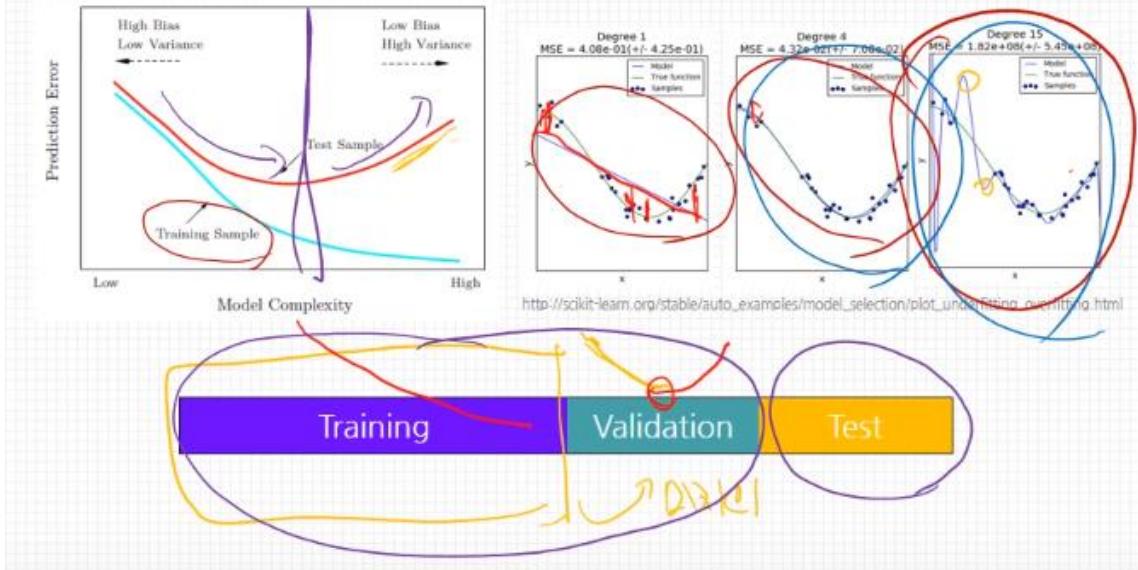
- 목적 : Overfitting을 피할 수 있다.

- 내용 : 학습은 training set으로만 하고, training set이 점점 늘어날수록 fitting error는 줄어들 것이다.

그러나 validation error가 줄다가 저점을 찍고 올라오는 지점이 있을 것이다. 거기서 stop을 해서 training set과 validation set의 황금비율을 찾으면 된다.

- 선정 : test set이 validation set을 잘 대표할 수 있어야 한다.

005 Train/Test/Validation Set



- anomaly detection 과정 예시 :

(1) training set으로 $p(x)$ 를 만들고 (by 평균, 표준편차 계산) → 그러면 각 x_i 별로 y_i 가 나온다.

Algorithm evaluation

- Fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(m)}\}$
- On a cross validation/test example x , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (normal)} \end{cases}$$

(2) cross validation, test set으로 y 를 prediction한다. 얼마나 맞는지를 위의 training set으로 만든 모델과 비교한다.

Possible evaluation metrics:

- True positive, false positive, false negative, true negative
- Precision/Recall
- F_1 -score

→ 엡실론을 정해본다. F_1 score가 뭐지?

- 주의 : skewed data를 쓰기 때문에 classification은 좋은 검증방법이 아닐 수 있다.

Suppose you have fit a model $p(x)$. When evaluating on the cross validation set or test set, your algorithm predicts:

$$y = \begin{cases} 1 & \text{if } p(x) \leq \varepsilon \\ 0 & \text{if } p(x) > \varepsilon \end{cases}$$

Is classification accuracy a good way to measure the algorithm's performance?

- Yes, because we have labels in the cross validation / test sets.
- No, because we do not have labels in the cross validation / test sets.
- No, because of skewed classes (so an algorithm that always predicts $y = 0$ will have high accuracy).

3.2. 언제 Anomaly detection을 해야 할까? vs Supervised learning(Logistic regression, neural network)

- 대상 : $y = 0, y = 1$ 로 labeled data를 얻었는데 supervised learning 이 Anomaly detection보다 낫지 않을까?
- 비교표 (언제 anomaly detection을 쓰는가)
 - sample에서 negative example($y=0$)이 많다. → Anomaly detection
 - anomaly의 종류가 예측 불허 할 때 → Anomaly detection
 - cf) spam의 경우에는 anomaly의 종류가 많긴 하지만, 그것에 대한 data set이 이미 다 있기 때문에 supervised를 하는 것이 낫다.

Anomaly detection

- Very small number of positive examples ($y = 1$). (0-20 is common).
- Large number of negative ($y = 0$) examples. $p(x)$
- Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like;
- future anomalies may look nothing like any of the anomalous examples we’ve seen so far.

vs.

Supervised learning

Large number of positive and negative examples. ←

Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set. ←

Spam ←

- 예시 :

Anomaly detection

- • Fraud detection $y=1$
- • Manufacturing (e.g. aircraft engines)
- • Monitoring machines in a data center

vs.

Supervised learning

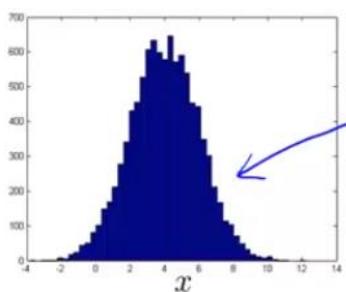
- Email spam classification ←
- Weather prediction ← (sunny/rainy/etc).
- Cancer classification ←

3.3. 어떤 feature를 사용해야 할까?

- non-gaussian feature도 사용할 수 있을까? (정규분포가 아니라면??)

(1) 가공을 해서 변수 x 를 더 정규분포로 만들 수 있다.

Non-gaussian features



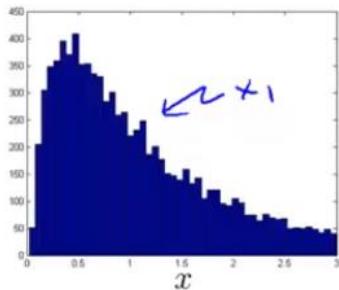
$$p(x_i; \mu_i, \sigma^2_i)$$

$$x_1 \leftarrow \log(x_i)$$

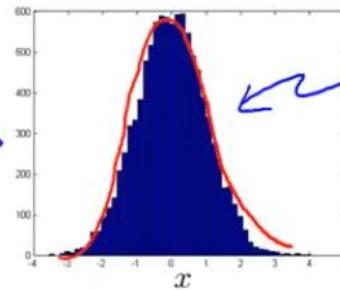
$$x_2 \leftarrow \log(x_2 + 1)$$

$$x_3 \leftarrow \sqrt{x_3} = x_3^{\frac{1}{2}}$$

$$x_4 \leftarrow x_4^{\frac{1}{3}}$$



$$\log(x)$$



```

>> hist(x.^0.05, 50)
>> xNew = x.^0.05;
>> hist(log(x), 50);
>> xNew = log(x);

```

(2) How do you come up with new features of Anomaly detection algorithm

→ Error analysis procedure

(2-1) make new features

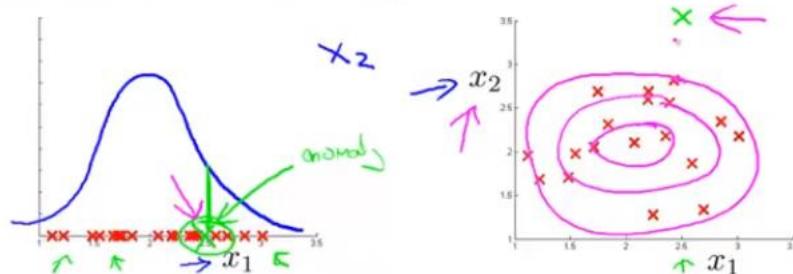
- 문제 상황 : x_1 의 경우에 anomalous example의 $p(x)$ 가 크면, x_1 대신에 x_2 를 찾아본다. 그래서 x_1, x_2 를 동시에 봐서 anomalous example의 $p(x)$ 가 작게 끔 한다.

→ Error analysis for anomaly detection

[Want $p(x)$ large for normal examples x .
 $p(x)$ small for anomalous examples x .]

Most common problem:

[$p(x)$ is comparable (say, both large) for normal
and anomalous examples]



(2-2) 예시 : 변수 간의 나눗셈을 하면 큰 anomaly를 잡아낼 수 있다.

→ Monitoring computers in a data center

→ Choose features that might take on unusually large or small values in the event of an anomaly.

→ x_1 = memory use of computer

→ x_2 = number of disk accesses/sec

→ x_3 = CPU load ←

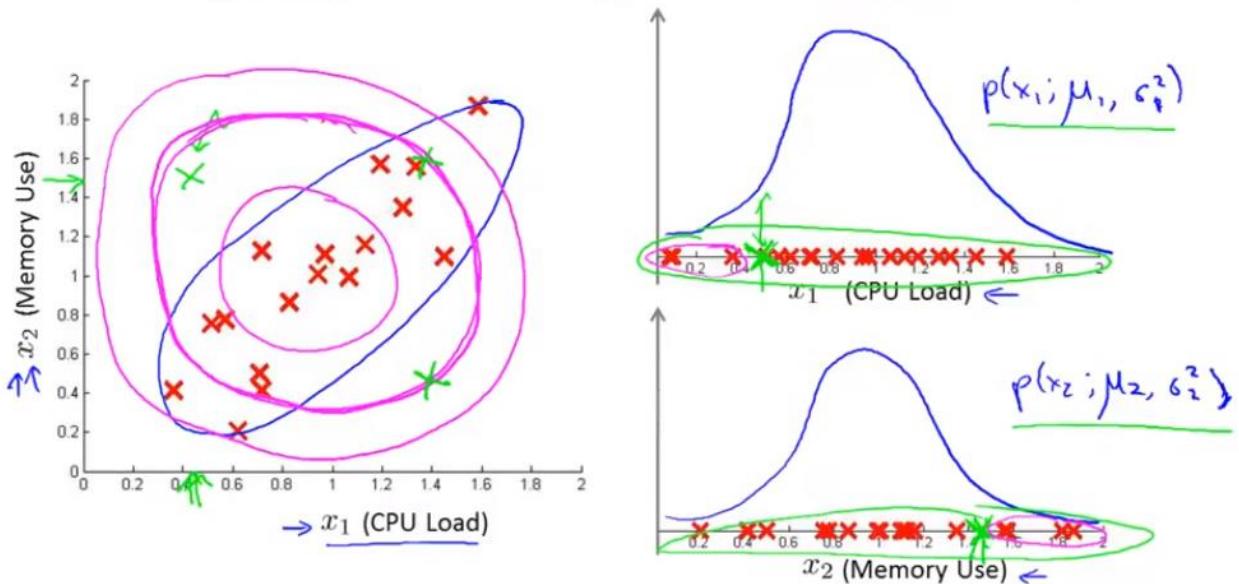
→ x_4 = network traffic ←

$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$

$$x_6 = \frac{(\text{CPU load})^2}{\text{network traffic}}$$

3.4. (심화학습) Multivariate Gaussian distribution

Motivating example: Monitoring machines in a data center



→ $P(x)$ 를 x_1, x_2 별로 하면 높다. 하지만 x_1, x_2 를 동시에 생각하면 이상값이다.

(혁동님) 결합을 해야 anomaly를 detect할 수 있다!!

블록 비유 : 1, 4, 6, 4, 1으로 대각선으로 쌓여있을 때, 위에서 보면 detection이 가능하다. anomaly를!!

- Multivariate gaussian distribution

Multivariate Gaussian (Normal) distribution

→ $x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2), \dots$, etc. separately.
 Model $p(x)$ all in one go.
 Parameters: $\mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$|\Sigma| = \text{determinant of } \Sigma \quad | \det(\text{Sigma})|$

→ 엄청 이상하네 ㅋㅋㅋ 식이.

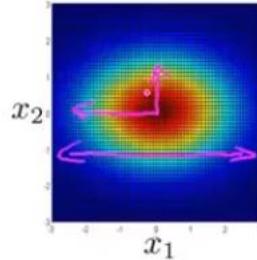
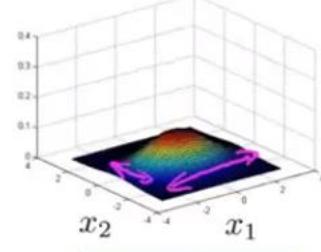
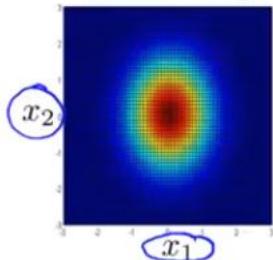
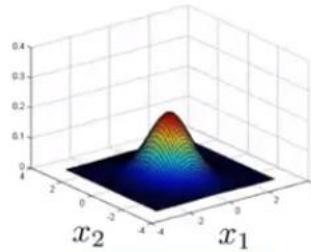
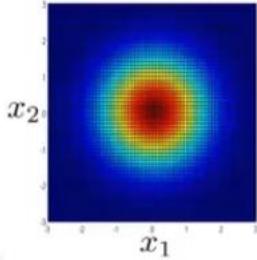
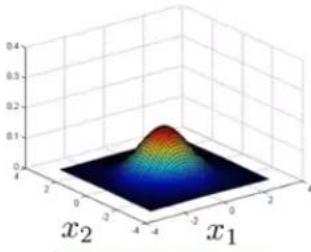
- variance를 바꾸면 어떻게 편할까?

Multivariate Gaussian (Normal) examples

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$



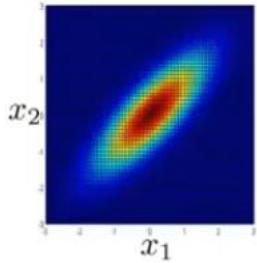
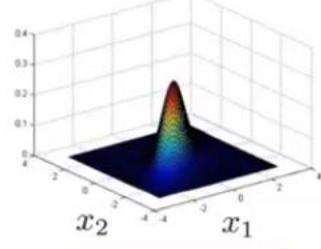
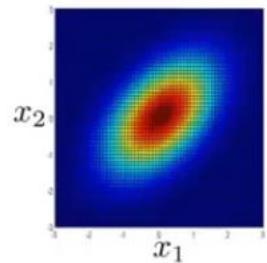
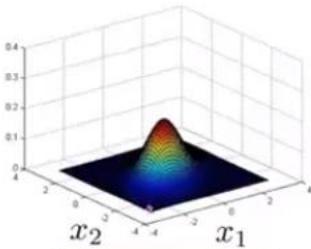
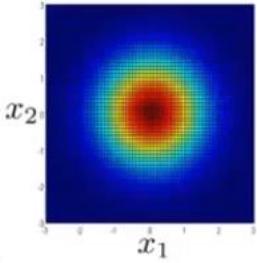
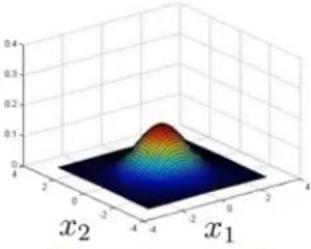
- Correlation이 달라지는 경우. : diagonal side에 있는 걸로!

Multivariate Gaussian (Normal) examples

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

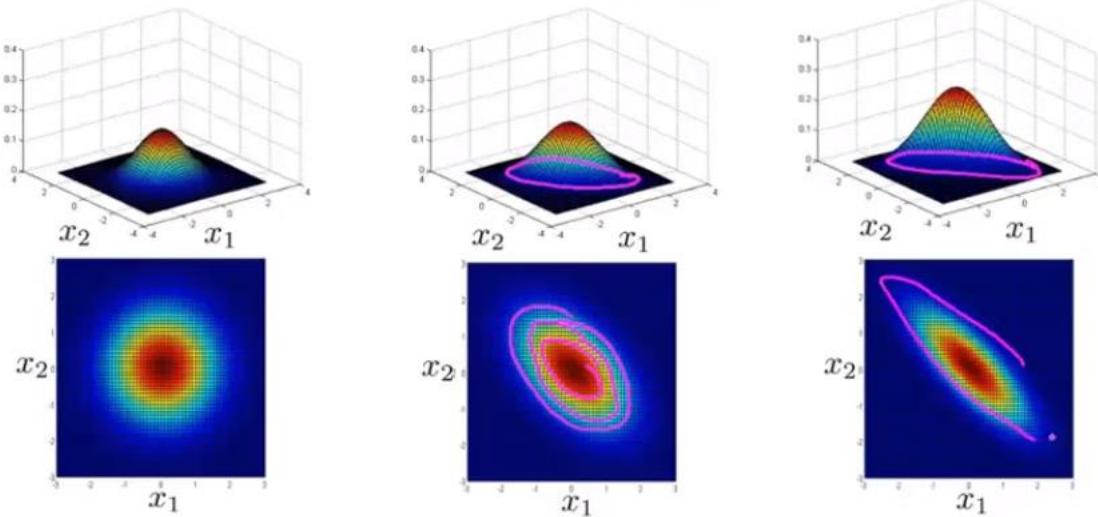
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



- negative correlation

Multivariate Gaussian (Normal) examples

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix} \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$



- 실제로 anomaly detection을 하는 방법

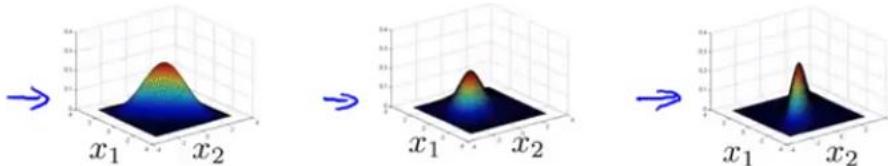
(1) 먼저, 평균과 표준편차 행렬(covariance matrix)을 구해서 정규분포를 가정하여 $P(x)$ 를 fitting한다.

Multivariate Gaussian (Normal) distribution

Parameters $\underline{\mu, \Sigma}$

$$\mu \in \mathbb{R}^n \quad \Sigma \in \mathbb{R}^{n \times n}$$

$$\rightarrow p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$



Parameter fitting:

Given training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \leftarrow x \in \mathbb{R}^n$

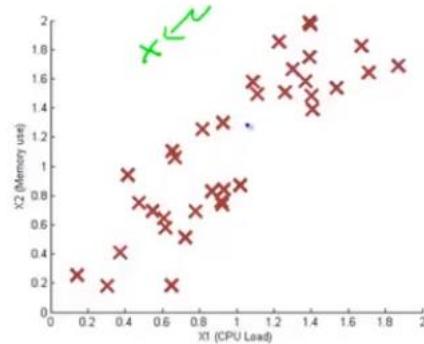
$$\rightarrow \boxed{\mu} = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \rightarrow \boxed{\Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

(2) 새로운 sample x 의 확률을 구해서 anomaly인지 아닌지 판단한다.

Anomaly detection with the multivariate Gaussian

1. Fit model $p(x)$ by setting

$$\begin{cases} \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \\ \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T \end{cases}$$



2. Given a new example x , compute

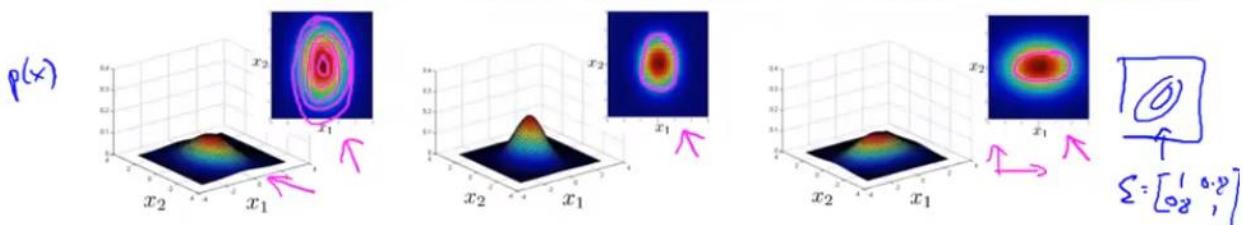
$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

Flag an anomaly if $\underline{p(x) < \varepsilon}$

- multivariate Gaussian model과 original model과의 차이점 : original model은 등고선이 대각 방향으로 향해있을 수가 없다.

Relationship to original model

Original model: $\underline{p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)}$



Corresponds to multivariate Gaussian

$$\rightarrow p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

where

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \dots & \dots \\ \vdots & \ddots & \vdots \\ \dots & \dots & \sigma_n^2 \end{bmatrix}$$

Andrew Ng

- 언제 original model, 언제 multivariate Gaussian model을 사용해야 할까?
- original model은 anomaly를 잡기 위해서 새로운 변수 x_3 를 만들어야 한다. 보통 더 많이 쓴다.
- multivariate Gaussian model은 역행렬을 구한다는 것이 좀 부담스럽다. 그런데, 변수 개수보다 데이터 양이 훨씬 더 많으면 original 방법보다 편하다. original 방법은 새로운 변수를 더 많이 만들어야 하기 때문이다. 변수가 redundant하면 알고 있는 변수가 줄어드는 효과가 있기 때문에 original을 사용하는 것이 더 좋을 수 있다.

→ Original model

$$p(x_1; \mu_1, \sigma_1^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$$

Manually create features to capture anomalies where x_1, x_2 take unusual combinations of values.

$$\rightarrow X_3 = \frac{x_1}{x_2} = \frac{\text{CPU load}}{\text{memory}}$$

→ Computationally cheaper (alternatively, scales better to large n) $n=10,000, n=100,000$

OK even if m (training set size) is small

vs. → Multivariate Gaussian

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

→ Automatically captures correlations between features

$$\Sigma \in \mathbb{R}^{n \times n} \quad \Sigma^{-1}$$

Computationally more expensive

Must have $m > n$ or else Σ is non-invertible. $m \geq 10n$

- 변수 2개만 써도 $n(n+1)/2$ 개가 되니까 너무 많아진다.

4. Recommender system

- 형태 : 영화 추천 시스템

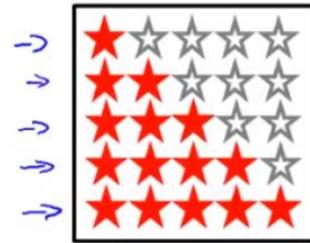
Example: Predicting movie ratings

→ User rates movies using one to five stars zero

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	5	0	6
Romance forever	5	?	0	0
Cute puppies of love	?	5	0	?
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	?

$$n_u = 4$$

$$n_m = 5$$



$\rightarrow n_u = \text{no. users}$
 $\rightarrow n_m = \text{no. movies}$
 $r(i, j) = 1$ if user j has rated movie i
 $y^{(i,j)}$ = rating given by user j to movie i (defined only if $r(i, j) = 1$)

4.1. feature를 먼저 알고 user들의 parameter vector를 구하는 방법

- feature → content based recommendations : 영화의 속성을 넣을 수 있다. 영화의 속성이 평가자의 평가 점수에 미치는 영향을 learning 해보자.

→ 사용자 각각 theta를 정의할 수 있다.

Content-based recommender systems

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$n_u = 4, n_m = 5$	$x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$
$x^{(1)}$ Love at last 1	5	5	0	0	$\downarrow x_1$ (romance) $\rightarrow 0.9 \rightarrow 0$	
$x^{(2)}$ Romance forever 2	5	?	?	0	$\rightarrow 1.0 \rightarrow 0.01$	
$x^{(3)}$ Cute puppies of love 3	?	4	0	?	$\rightarrow 0.99 \rightarrow 0$	
$x^{(4)}$ Nonstop car chases 4	0	0	5	4	$\rightarrow 0.1 \rightarrow 1.0$	
$x^{(5)}$ Swords vs. karate 5	0	0	5	?	$\rightarrow 0 \rightarrow 0.9$	$n=2$

→ For each user j , learn a parameter $\theta^{(j)} \in \mathbb{R}^3$. Predict user j as rating movie i with $(\theta^{(j)})^T x^{(i)}$ stars.

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \Leftrightarrow \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad (\theta^{(1)})^T x^{(3)} = 5 \times 0.99 = 4.95$$

- linear regression과 유사하다!

Problem formulation

→ $r(i, j) = 1$ if user j has rated movie i (0 otherwise)

→ $y^{(i,j)}$ = rating by user j on movie i (if defined)

→ $\theta^{(j)}$ = parameter vector for user j

→ $x^{(i)}$ = feature vector for movie i

→ For user j , movie i , predicted rating: $\underline{(\theta^{(j)})^T (x^{(i)})}$

→ $m^{(j)}$ = no. of movies rated by user j

To learn $\underline{\theta^{(j)}}$:

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i : r(i,j)=1} ((\theta^{(j)})^T (x^{(i)}) - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (\theta_k^{(j)})^2$$

→ regularization term도 있다.

- 방법 정리 :

Optimization algorithm:

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$J(\theta^{(1)}, \dots, \theta^{(n_u)})$

Gradient descent update:

$$\begin{aligned} \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k = 0) \\ \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0) \end{aligned}$$

→ 1/m 은 optimization식에서 없앴기 때문에 없는거다.

4.2. Collaborative filtering → 어떤 feature를 쓸 지 스스로 배운다.

- 목적 : feature가 주어지지 않으면?? User의 parameter vector를 안 후에 feature 값을 찾아낸다.

- 예시 :

(1) 초기 조건

Problem motivation

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	x_1 (romance)	x_2 (action)
Love at last	5	5	0	0	?	?
Romance forever	5	?	?	0	?	?
Cute puppies of love	?	4	0	?	?	?
Nonstop car chases	0	0	5	4	?	?
Swords vs. karate	0	0	5	?	?	?

(2) user들이 취향을 말해줬다.

$$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(2)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}, \theta^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

(3) 이를 통해서 x_1, x_2 를 알아낼 수 있다. 그들이 한 (평가 점수) = $\theta^T * \text{feature}$

- 알고리즘 :

- 아까는 θ 으로 minimize를 했는데, 이제는 x_i 를 통해 식을 minimize한다.

Optimization algorithm

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(i)}$:

$$\rightarrow \min_{x^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(1)}, \dots, x^{(n_m)}$:

$$\min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

$$\odot x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

→ θ 가 주어졌는지, x 가 주어졌는지에 따라서 방법이 달라진다.

- 두 가지를 번갈아 가면서 하면 된다.

Collaborative filtering

[Given $x^{(1)}, \dots, x^{(n_m)}$ (and movie ratings),
can estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$]

$r^{(i,j)}$
 $y^{(i,j)}$

[Given $\theta^{(1)}, \dots, \theta^{(n_u)}$,
can estimate $x^{(1)}, \dots, x^{(n_m)}$]

Guess $\Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \dots$

- Collaborative filtering 알고리즘

- theta랑 x를 번갈아하지 않아도 된다.

→ cost function을 합치면 된다.

Collaborative filtering optimization objective $(i,j) : r(i,j) \neq 1$

→ Given $x^{(1)}, \dots, x^{(n_m)}$, estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$:

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

→ Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, estimate $x^{(1)}, \dots, x^{(n_m)}$:

$$\min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

Minimizing $x^{(1)}, \dots, x^{(n_m)}$ and $\theta^{(1)}, \dots, \theta^{(n_u)}$ simultaneously:

$$J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$$\min_{\substack{x^{(1)}, \dots, x^{(n_m)} \\ \theta^{(1)}, \dots, \theta^{(n_u)}}} J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$$

→ 하늘색, 빨간색 부분이 각각 θ, x 를 기준으로 할 때, 다른 것이 상수가 되므로 합쳐도 무방하다.

→ 주의사항 : x, θ 에는 0행이 포함되지 않기 때문에 항개수가 $n+1$ 이 아니라, n 이다.

- 정리 :

Collaborative filtering algorithm

→ 1. Initialize $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ to small random values.

→ 2. Minimize $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ using gradient descent (or an advanced optimization algorithm). E.g. for every $j = 1, \dots, n_u, i = 1, \dots, n_m$:

$$x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right)$$

3. For a user with parameters $\underline{\theta}$ and a movie with (learned) features \underline{x} , predict a star rating of $\underline{\theta}^T \underline{x}$.

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

→ 1. 처음에 small random value로 initialize하는 이유 : this serves as a symmetry breaking and ensures the algorithm learns features that are different from each other.

- Vectorization : Low rank matrix factorization :

- predicted rating을 하는 방법 : matrix를 만들면 된다.

Collaborative filtering

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

$(\Theta^{(i)})^T(x^{(j)})$

Predicted ratings: $\uparrow_{(i,j)}$

$$\left[\begin{array}{cccc} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \dots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \dots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \dots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{array} \right]$$

Collaborative filtering

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

$(\Theta^{(i)})^T(x^{(j)})$

Predicted ratings: $\uparrow_{(i,j)}$

$$\left[\begin{array}{cccc} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \dots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \dots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \dots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{array} \right]$$

$$\rightarrow X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(n_m)})^T \end{bmatrix} \quad \rightarrow \Theta = \begin{bmatrix} -(\Theta^{(1)})^T \\ -(\Theta^{(2)})^T \\ \vdots \\ -(\Theta^{(n_u)})^T \end{bmatrix}$$

Low rank matrix factorization

$$\text{Let } X = \begin{bmatrix} - & (x^{(1)})^T & - \\ & \vdots & \\ - & (x^{(n_m)})^T & - \end{bmatrix}, \quad \Theta = \begin{bmatrix} - & (\theta^{(1)})^T & - \\ & \vdots & \\ - & (\theta^{(n_u)})^T & - \end{bmatrix}.$$

What is another way of writing the following:

$$\left[\begin{array}{ccc} (x^{(1)})^T(\theta^{(1)}) & \dots & (x^{(1)})^T(\theta^{(n_u)}) \\ \vdots & \ddots & \vdots \\ (x^{(n_m)})^T(\theta^{(1)}) & \dots & (x^{(n_m)})^T(\theta^{(n_u)}) \end{array} \right]$$

- low rank라는 것은 linearly independent한 vector의 수이다. 그것을 줄이면 연산이 편해진다.
(혁동님 자료)

Two types of sparsity for matrices $M \in \mathbb{R}^{n \times p}$

II - Through a factorization of $M = UV^\top$

- Matrix $M = UV^\top$, $U \in \mathbb{R}^{n \times k}$ and $V \in \mathbb{R}^{p \times k}$

- Low rank:** m small

$$M = U V^\top$$

- Sparse decomposition:** U sparse

$$M = \text{Sparse } U V^\top \quad 6)$$

$$\begin{array}{c} \text{Item} \\ \begin{array}{cccc} W & X & Y & Z \end{array} \\ \begin{array}{l} \text{User} \\ \begin{array}{|c|c|c|c|} \hline & 4.5 & 2.0 & \\ \hline A & 4.0 & & 3.5 \\ \hline B & & 5.0 & 2.0 \\ \hline C & 3.5 & 4.0 & 1.0 \\ \hline D & & & \\ \hline \end{array} \end{array} = \begin{array}{l} \text{Rating Matrix} \\ \begin{array}{|c|c|c|} \hline & 1.2 & 0.8 \\ \hline A & 1.4 & 0.9 \\ \hline B & 1.5 & 1.0 \\ \hline C & 1.2 & 0.8 \\ \hline D & & \\ \hline \end{array} \end{array} \times \begin{array}{l} \text{User Matrix} \\ \begin{array}{|c|c|c|c|} \hline & 1.5 & 1.2 & 1.0 & 0.8 \\ \hline & 1.7 & 0.6 & 1.1 & 0.4 \\ \hline \end{array} \end{array} \times \begin{array}{l} \text{Item Matrix} \\ \begin{array}{|c|c|c|c|} \hline & 1.5 & 1.2 & 1.0 & 0.8 \\ \hline & 1.7 & 0.6 & 1.1 & 0.4 \\ \hline \end{array} \end{array}$$

$$\begin{array}{c} \text{Item} \\ \begin{array}{cccc} W & X & Y & Z \end{array} \\ \begin{array}{l} \text{User} \\ \begin{array}{|c|c|c|c|} \hline & 4.5 & 2.0 & \\ \hline A & 4.0 & & 3.5 \\ \hline B & & 5.0 & 2.0 \\ \hline C & 3.5 & 4.0 & 1.0 \\ \hline D & & & \\ \hline \end{array} \end{array} = \begin{array}{l} \text{Rating Matrix} \\ \begin{array}{|c|c|c|c|} \hline & 1.2 & 0.8 & 0.8 \\ \hline A & 1.4 & 0.9 & 0.9 \\ \hline B & 1.5 & 1.0 & 1.0 \\ \hline C & 1.2 & 0.8 & 0.8 \\ \hline D & & & \\ \hline \end{array} \end{array} \times \begin{array}{l} \text{User Matrix} \\ \begin{array}{|c|c|c|c|} \hline & 1.5 & 1.2 & 1.0 & 0.8 \\ \hline & 1.7 & 0.6 & 1.1 & 0.4 \\ \hline & 1.7 & 0.6 & 1.1 & 0.4 \\ \hline \end{array} \end{array} \times \begin{array}{l} \text{Item Matrix} \\ \begin{array}{|c|c|c|c|} \hline & 1.5 & 1.2 & 1.0 & 0.8 \\ \hline & 1.7 & 0.6 & 1.1 & 0.4 \\ \hline & 1.7 & 0.6 & 1.1 & 0.4 \\ \hline \end{array} \end{array}$$

→ 더 적은 양이 연산으로!!

Finding related movies

For each product i , we learn a feature vector $\underline{x^{(i)}} \in \mathbb{R}^n$.

$x_1 = \text{romance}$, $x_2 = \text{action}$, $x_3 = \text{comedy}$, $x_4 = \dots$

- how to find related movies??

How to find movies j related to movie i ?

small $\|x^{(i)} - x^{(j)}\| \rightarrow$ movie j and i are "similar"

5 most similar movies to movie i :

→ Find the 5 movies j with the smallest $\|x^{(i)} - x^{(j)}\|$.

→ 거리로 알 수 있다.

- mean normalization :

- 문제상황 : 만약에 어떤 사람이 영화평가를 1개도 안 했다고 하면, 다 0으로 rating할 것으로 예측이 되어 버린다. 그 이유는 cost function의 첫 항은 정보가 없어서 무용지물이고, regularization term에서 θ 를 각각 다 0으로 해야지 cost function이 최소화되기 때문이다.

Users who have not rated any movies

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)
Love at last	5	5	0	0	?
Romance forever	5	?	?	0	?
Cute puppies of love	?	4	0	?	?
Nonstop car chases	0	0	5	4	?
Swords vs. karate	0	0	5	?	?

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

$$\min_{\substack{x^{(1)}, \dots, x^{(n_m)} \\ \theta^{(1)}, \dots, \theta^{(n_u)}}} \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$n=2 \quad \Theta^{(s)} \in \mathbb{R}^2 \quad \Theta^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \frac{\lambda}{2} [(\Theta_1^{(s)})^2 + (\Theta_2^{(s)})^2] \leftarrow$
 $(\Theta^{(s)})^T x^{(i)} = 0$

- Mean normalization

- average를 따로 저장해놓는다.

- 평가를 한 거에서 평균을 빼버린다. → 평균이 0이 되게 만든다.

- predict를 하고 나서 평균을 더하면 최종 예측값이 되게 한다.

Mean Normalization:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? & 2.5 \\ 5 & ? & ? & 0 & ? & 2.5 \\ ? & 4 & 0 & ? & ? & 2 \\ 0 & 0 & 5 & 4 & ? & ? \\ 0 & 0 & 5 & 0 & ? & ? \end{bmatrix}$$

$$\mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

For user j , on movie i predict:

$$\rightarrow (\underline{\Theta^{(s)}})^\top (\underline{x^{(i)}}) + \underline{\mu_i}$$

\downarrow
learn $\underline{\Theta^{(s)}}, \underline{x^{(i)}}$

User 5 (Eve):

$$\underline{\Theta^{(s)}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\underbrace{(\underline{\Theta^{(s)}})^\top (\underline{x^{(i)}})}_{\approx 0} + \boxed{\underline{\mu_i}}$$

IV. What Machine Learning should I apply?

- 이 단원이 굉장히 중요한 단원이 될 것이다. Debugging 과정과 어디서 머신 러닝 개선을 이룰 수 있는지를 알 수 있게 해준다.

- How to improve the learning algorithm

1. get more training examples

- 문제점 : 2배를 더 모으면 개선될까? 아닐 수 있다. 시간 낭비일 수 있다.

2. feature가 너무 많을 때는 overfitting이 일어날 수 있으므로 feature를 몇 개만 선정한다.

3. feature를 더 늘린다.

4. polynomial feature ($x^2, x_1 \cdot x_2$)

5. decrease/increase "lambda" regularization parameter

- 위의 option 중 어떤 것을 하는 게 좋을지 정하는 방법이 있다.

→ **Machine learning diagnostic** : a test that you can run to gain insight what is / isn't working with a learning algorithm.

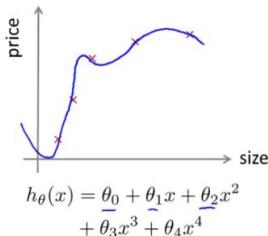
6. Machine Learning diagnostic

6.1. Test set을 배정하여 test error를 보면 overfitting인지 아닌지 볼 수 있다.

- How to evaluate a hypothesis : **training error** & test error를 둘 다 측정한다!!

- training error가 low한지는 좋은 지표가 아니다. → overfitting의 문제

Overfitting example



$$h_\theta(x) = \theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3 + \theta_4x^4$$

- training set(70%) + test set(30%)를 해서 확인할 수 있다.

- test set error가 high하면 overfitting한 것이다.

Dataset:

Size	Price	
2104	400	
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	
<hr/>		
1427	199	
30% 1380	212	
1494	243	

↑ 70%

↑ 30%

Training set

Test set

→

$(x^{(1)}, y^{(1)})$
 $(x^{(2)}, y^{(2)})$
 \vdots
 $(x^{(m)}, y^{(m)})$

→

$(x_{test}^{(1)}, y_{test}^{(1)})$
 $(x_{test}^{(2)}, y_{test}^{(2)})$
 \vdots
 $(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

(1) linear regression

Training/testing procedure for linear regression

- > - Learn parameter θ from training data (minimizing training error $J(\theta)$) $\geq 70\%$

- Compute test set error:

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

(2) logistic regression 도 test set을 사용한다!

- test error 도 좋고, Misclassification error도 계산한다.

- Compute test set error:

$$\rightarrow J_{test}(\theta) = -\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} y_{test}^{(i)} \log h_{\theta}(x_{test}^{(i)}) + (1 - y_{test}^{(i)}) \log h_{\theta}(x_{test}^{(i)})$$

- Misclassification error (0/1 misclassification error):

$$err(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq 0.5, y = 0 \\ 0 & \text{otherwise} \end{cases}$$

or if $h_{\theta}(x) < 0.5, y = 1$ } error

$$\text{Test error} = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} err(h_{\theta}(x_{test}^{(i)}), y_{test}^{(i)}).$$

6.2. Validation set을 배정하여 polynomial model 중에 어떤 것이 진정으로 나은지 선택할 수 있다! (polynomial degree를 fit하는데 test set을 이미 써버렸다!)

- model selection problem

Model selection

1. $h_{\theta}(x) = \theta_0 + \theta_1 x$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3$
- ⋮
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}$

- degree of polynomial / choosing regularization parameter

- 어떤 model을 골라야 할까?

Model selection

$\rightarrow d = \text{degree of polynomial}$

- $d=1$ 1. $h_\theta(x) = \theta_0 + \theta_1 x \rightarrow \Theta^{(1)} \rightarrow J_{test}(\Theta^{(1)})$
- $d=2$ 2. $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \Theta^{(2)} \rightarrow J_{test}(\Theta^{(2)})$
- $d=3$ 3. $h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \Theta^{(3)} \rightarrow J_{test}(\Theta^{(3)})$
- \vdots \vdots
- $d=10$ 10. $h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \Theta^{(10)} \rightarrow J_{test}(\Theta^{(10)})$

Choose $\boxed{\theta_0 + \dots + \theta_5 x^5}$



How well does the model generalize? Report test set error $J_{test}(\theta^{(5)})$.

$\Theta^{(5)}$

Problem: $J_{test}(\theta^{(5)})$ is likely to be an optimistic estimate of generalization error. I.e. our extra parameter ($d = \text{degree of polynomial}$) is fit to test set.

→ 이 예에서는 test set을 가지고 degree of polynomial을 결정했기 때문에 $d=5$ 일 때가 가장 좋은 예측 모델이라고 하기가 힘들다는 것이다 → 그러면 training set을 가지고 하면 되는 것 아니냐? 아 Theta를 training set으로 유도하는 것이다!!! 그 다음에 test set으로 10개 모델 중에서 어떤 게 가장 나은지를 본 것이다. 근데 $d=5$ 가 나온 걸 또 검증해야 하는데.... 이 포인트에서 validation set이 나온다!

- 해결 : 그러면 data set을 1) training set (60%) 2) (cross) validation set (cv) (20%) 3) test set (20%)로 세가지로 나누면 되지 않느냐?

Evaluating your hypothesis

Dataset:

Size	Price	
2104	400	60%
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	20%
1427	199	
1380	212	20%
1494	243	

Training set

Cross validation set (cv)

test set

$M_{cv} = \text{no. of cv examples}$

M_{test}

$(x^{(1)}, y^{(1)})$
 $(x^{(2)}, y^{(2)})$
 \vdots
 $(x^{(m)}, y^{(m)})$

$(x_{cv}^{(1)}, y_{cv}^{(1)})$
 $(x_{cv}^{(2)}, y_{cv}^{(2)})$
 \vdots
 $(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

$(x_{test}^{(1)}, y_{test}^{(1)})$
 $(x_{test}^{(2)}, y_{test}^{(2)})$
 \vdots
 $(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

→ training error, cross validation error, test error 세 가지의 error가 나오게 된다.

Training error:

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

Cross Validation error:

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Test error:

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_\theta(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

→ 여기서 validation set을 가지고 d를 선택하고, 그 다음에 남겨진 test set을 가지고 d=5를 검증하는데 사용하면 되지 않느냐.

- 정리 :

We can now calculate three separate error values for the three different sets using the following method:

1. Optimize the parameters in Θ using the training set for each polynomial degree.
2. Find the polynomial degree d with the least error using the cross validation set.
3. Estimate the generalization error using the test set with $J_{test}(\Theta^{(d)})$, (d = theta from polynomial with lower error);

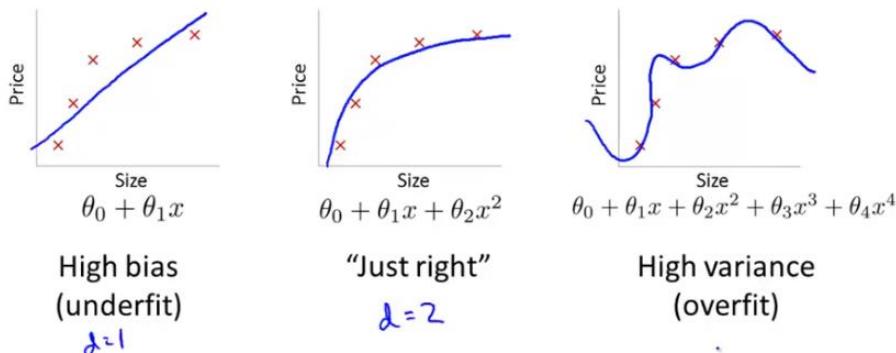
→ 여기서는 validation set을 가지고 polynomial degree d를 fitting했다.

6.3. feature가 많아서 문제냐, 적어서 문제냐! Diagnosing bias vs. variance!!

- 문제 : feature가 너무 많은 경우, feature가 부족한 경우를 어떻게 알까?

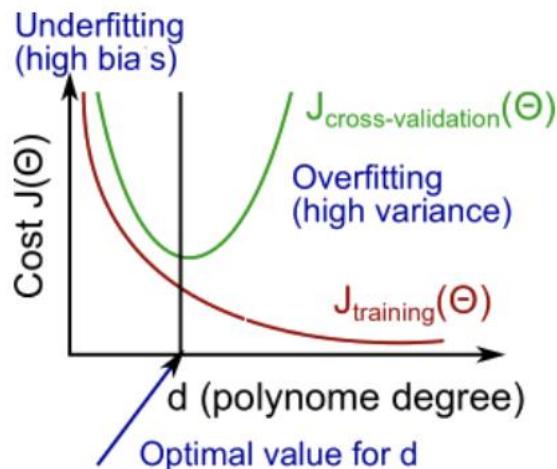
- high bias problem → underfitting problem (bias는 하나의 일직선으로 생각하자)

high variance problem → overfitting problem (variance는 구불구불한 선으로 생각하자)

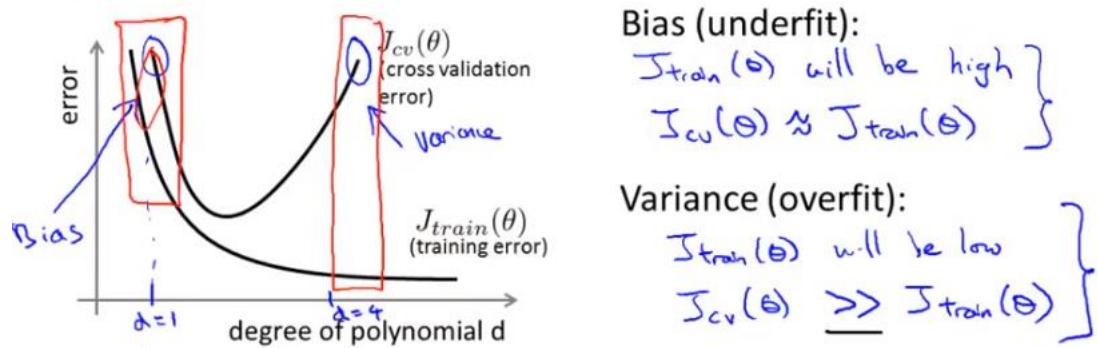


→ d가 오른쪽으로 갈수록 커진다.

- 해결 : training set error가 작으면 이건 overfitting(high variance)이다. cross validation error는 over나 under나 둘 다 높다. 왜냐하면 최적의 degree of polynomial이 있기 때문이다.



Suppose your learning algorithm is performing less well than you were hoping. ($J_{cv}(\theta)$ or $J_{test}(\theta)$ is high.) Is it a bias problem or a variance problem?



Bias (underfit):
 $J_{train}(\theta)$ will be high }
 $J_{cv}(\theta) \approx J_{train}(\theta)$ }

Variance (overfit):
 $J_{train}(\theta)$ will be low }
 $J_{cv}(\theta) \gg J_{train}(\theta)$ }

(1) Bias (underfit) : training error가 높고 validation set error가 높아서 둘이 비슷한 경우.

(2) Variance (overfit) : training error가 낮고, validation set error가 높은 경우. 따라서 cross validation error가 training set error보다 훨씬 큰 경우

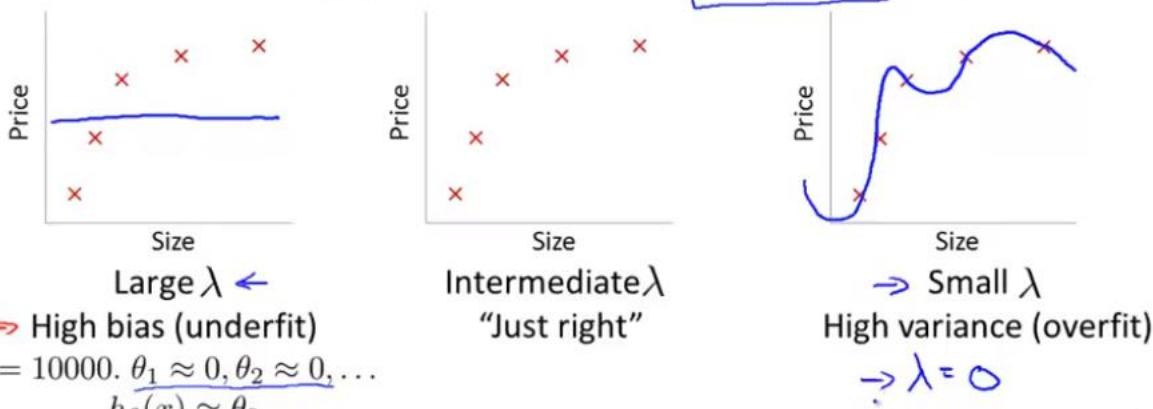
6.4. Lambda (regularization parameter)을 정할 때는 마찬가지로 cross validation set을 사용한다!!

- regularization이 어떻게 bias이랑 variance에 영향을 미치는지?
 - 이 얘기는 예전에 이미 했던 것이다.
 - Lambda (regularization parameter)가 크면 theta가 penalize되기 때문에 Theta 0만이 남게 된다. 따라서 high bias가 된다.
 - 만약 Lambda 가 작게 되면 Theta가 penalize가 적게 되니까, overfitting이 된다.

Linear regression with regularization

Model:
$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \leftarrow$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2 \quad \leftarrow$$



- 어떻게 Lambda를 정할 것인가?

Choosing the regularization parameter λ

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \leftarrow$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2 \quad \leftarrow$$

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

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

→ training error, validation error, test error에는 regularization term이 없다.
regularization term에서 m이 아니라 n이다.

- Lambda를 정하기 위한 model selection 방법

- (1) 일단 cost function J 를 최소화하는 θ 를 구한다. 이 식에는 물론 lambda가 포함되어 있다.
- (2) 다음, cross validation set을 통해서 어떤 Theta vector가 가장 cost function이 낮은지 본다. 여기에는 regularization term이 없다.
- (3) test set으로 마지막에 다시 점검한다.

Model: $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

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

- | | |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ol style="list-style-type: none"> 1. Try $\lambda = 0$ 2. Try $\lambda = 0.01$ 3. Try $\lambda = 0.02$ 4. Try $\lambda = 0.04$ 5. Try $\lambda = 0.08$ ⋮ 12. Try $\lambda = 10$ | $\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
$\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
$\rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
\vdots
$\rightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$ |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
- Pick (say) $\theta^{(5)}$. Test error:

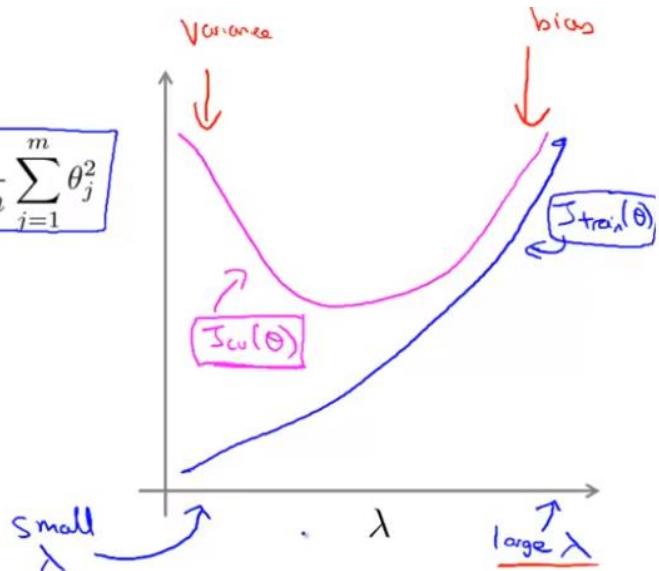
- Training error와 cross validation error를 lambda에 대해 그래프를 그려보자.

- (1) Training error는 lambda가 커지면 bias가 커지게 되기 때문에(underfitting) training error가 커진다.
- (2) Validation error는 lambda가 특정값(optimum)일 때를 제외하고는 양쪽으로 커질 것이다.

$$\rightarrow J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \boxed{\frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2}$$

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

$$\rightarrow \boxed{J_{cv}(\theta)} = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$



- Lambda 고르기 과정 정리

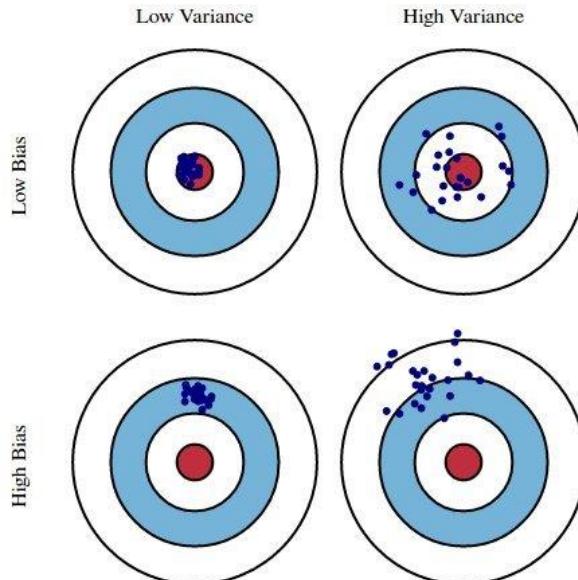
1. Create a list of lambdas (i.e. $\lambda \in \{0, 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24\}$);
2. Create a set of models with different degrees or any other variants.
3. Iterate through the λ s and for each λ go through all the models to learn some Θ .
4. Compute the cross validation error using the learned Θ (computed with λ) on the $J_{CV}(\Theta)$ without regularization or $\lambda = 0$.
5. Select the best combo that produces the lowest error on the cross validation set.
6. Using the best combo Θ and λ , apply it on $J_{test}(\Theta)$ to see if it has a good generalization of the problem.

6.5. training set을 더 늘리는 것이 능사인지 알려면! Learning curve를 그려야 한다.

- Learning curve : code/알고리즘이 잘 작동하는지, bias나 variance가 있는지를 전반적으로 보는 도구이다.
- x축은 training set 표본의 개수.

(1) Training error : training set m이 작으면 fitting이 쉬우니까, training error가 작을 것이다. 그러나 m이 커질수록 training error가 커질 것이다.

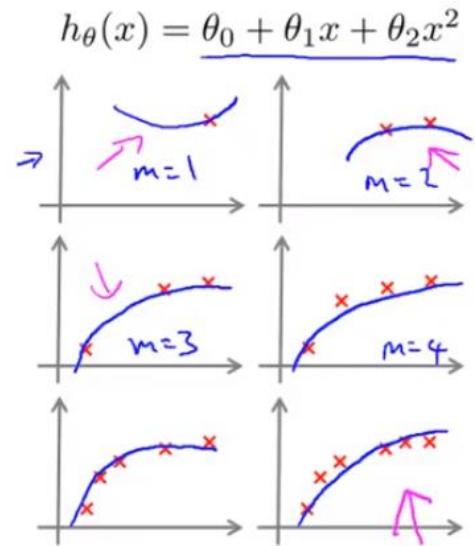
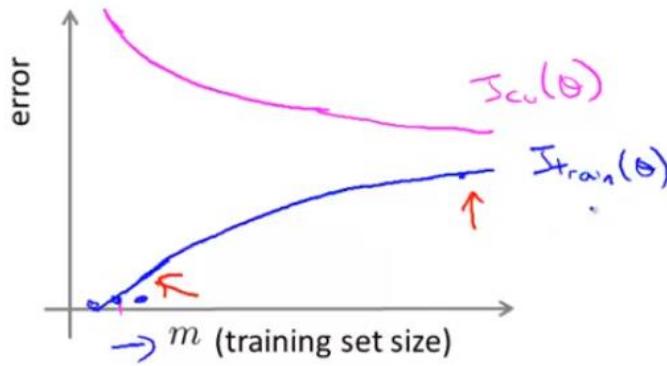
(2) Cross validation error : m이 작으면 hypothesis가 좋지 않다. 그러나 m이 커지면 hypothesis가 generalize가 잘 된다. training이 잘 된다.



Learning curves

$$\rightarrow \underline{J_{train}(\theta)} = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\Rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$



6.5.1. high bias의 경우

(1) Validation error : training set 개수 m 이 커진다고 해서 회귀직선이 크게 개선되지는 않는다. m 이 커진다고 해도 직선일 뿐이다.

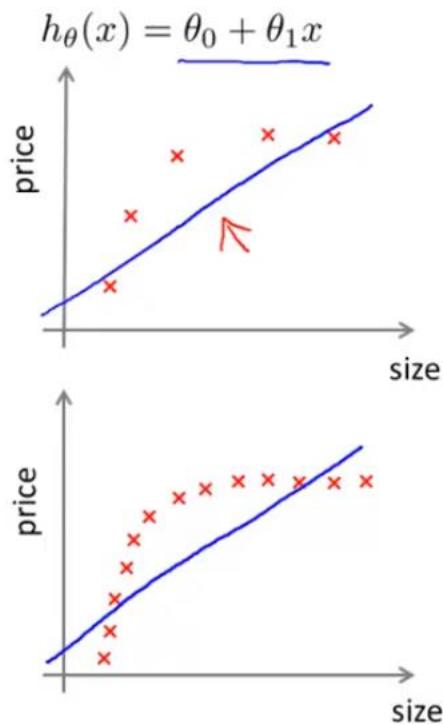
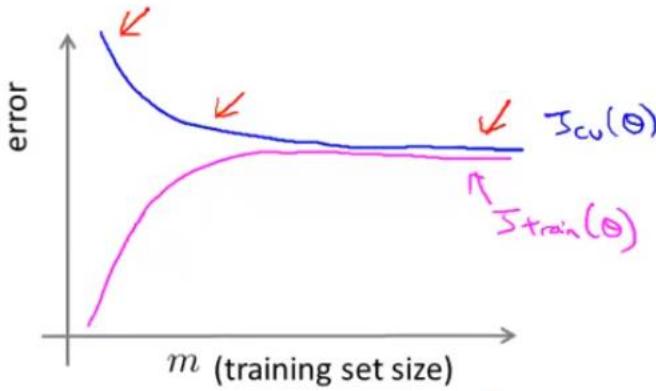
→ 흐.. 이 validation error를 어떻게 정의하는 것인가? 무엇을 validation하는 것인가? 아 training set의 수에 따라서 그것을 validation하는구나.

(2) Training error : m 이 작을 때는 fitting이 잘 되므로 error가 작지만 m 이 커질수록 error가 커지게 되는데 계속 커지지는 않는다. 아 training error와 validation error가 수렵하는데 이건 당연한 현상인건가???

→ 질문하기.

(3) 결론 : High bias일 때는 Validation error랑 Training error가 둘 다 크다. Bias가 크기 때문에(parameter가 2개 밖에 없어) 그리고 training set을 더 늘린다고 해서 딱히 error가 개선되지 않는다.

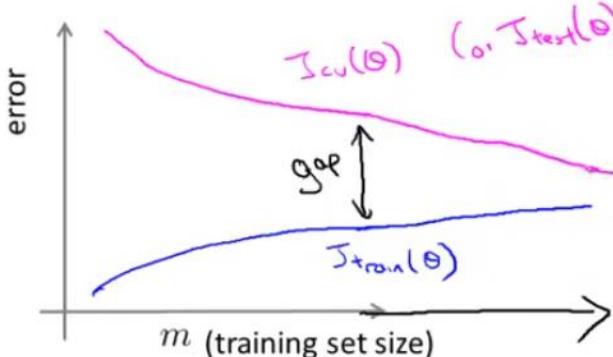
High bias



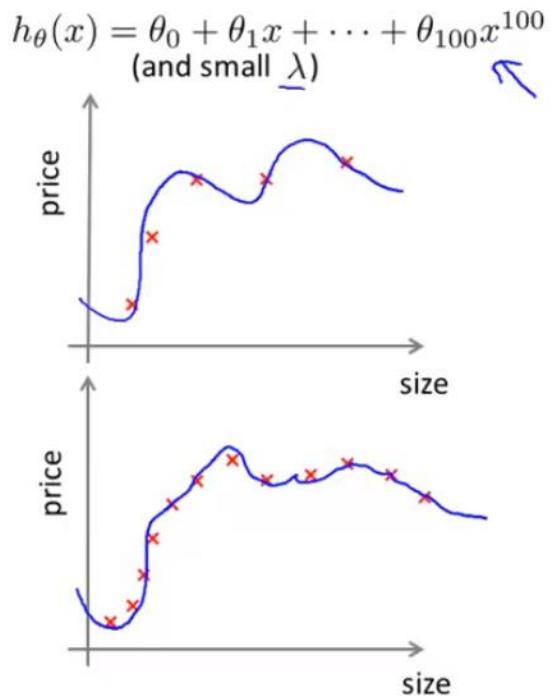
6.5.2 High variance의 경우

- (1) Training error : 표본이 많아지면, overfitting을 하긴 하겠지만 fitting을 모든 점에 하기 힘들어지기 때문에 training error가 조금 커진다. 그래도 낮은 편
- (2) Validation error : overfitting이니까 일단 error는 높은 편. m 이 커지면 그나마 조금 error가 줄어든다.
- (3) 결론 : High variance의 경우에는 gap이 있다. training example을 늘리면 cross validation error를 줄일 수 있겠구나! 도움이 된다.

High variance



If a learning algorithm is suffering from high variance, getting more training data is likely to help.



- 정리 :

Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

- Get more training examples → fixes high variance
- Try smaller sets of features → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features ($x_1^2, x_2^2, x_1 x_2$, etc) → fixes high bias.
- Try decreasing λ → fixes high bias
- Try increasing λ → fixes high variance

5.6. Neural Network의 경우에 고려사항

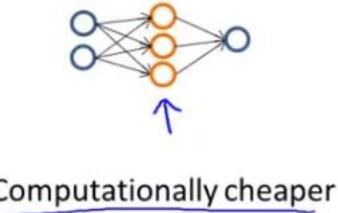
- (1) hidden layer의 개수 :

- 너무 많으면 overfitting이 된다. → regularization 사용
 - A neural network with fewer parameters is prone to underfitting. It is also computationally cheaper.
 - A large neural network with more parameters is prone to overfitting. It is also computationally expensive. In this case you can use regularization (increase λ) to address the overfitting.
- 질문!! 왜 layer가 많으면 parameter가 많은 효과가 나는 것일까?

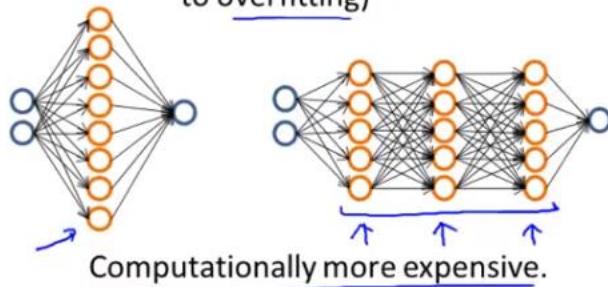
→ parameter가 많아지니까. 행렬로서!

Neural networks and overfitting

→ “Small” neural network
(fewer parameters; more prone to underfitting)



→ “Large” neural network
(more parameters; more prone to overfitting)



Use regularization (λ) to address overfitting.



- hidden layer를 1개 말고 다른 거를 정하고 싶을 때! / hidden layer 1개가 default이다!

: hidden layer 2개, 3개, 4개 등을 해보고 cross validation set, test set으로 어떤 것이 가장 잘 되는지 확인해서 정해본다.

7. Machine learning system design : Building a spam classifier

- 핵심 : 어떻게 learning algorithm을 평가하고 개선할 것인가?

- 어떻게 접근해야 할까? Prioritizing what to work on ← gut feeling으로 이 중에서 하나를 고르지 말 것.

How to spend your time to make it have low error?

- Collect lots of data
 - E.g. “honeypot” project.
- Develop sophisticated features based on email routing information (from email header).
- Develop sophisticated features for message body, e.g. should “discount” and “discounts” be treated as the same word? How about “deal” and “Dealer”? Features about punctuation?
- Develop sophisticated algorithm to detect misspellings (e.g. mortgage, medicine, watches.)

So how could you spend your time to improve the accuracy of this classifier?

- Collect lots of data (for example “honeypot” project but doesn't always work)
- Develop sophisticated features (for example: using email header data in spam emails)
- Develop algorithms to process your input in different ways (recognizing misspellings in spam).

→ 할 수 있는 가능성은 모두 열어두면 좋을 것 같다!

- 해결방법 : Recommended approach

(1) start a simple algorithm that you can implement quickly : 일단 빨리 만들어서 해보자!!

(2) **Learning curve** 그래서 high variance인지 high bias인지 살피기 → more data helpful? more features helpful?

(3) Error analysis : Manually examine the examples that your algorithm made errors on.

→ 1) what type of email it is → 수동 작업으로 쭉 보는 것이 중요하다.

Error Analysis

$m_{CV} = 500$ examples in cross validation set

Algorithm misclassifies 100 emails.

Manually examine the 100 errors, and categorize them based on:

→ (i) What type of email it is *pharma, replica, steal passwords, ...*

→ (ii) What cues (features) you think would have helped the algorithm classify them correctly.

Pharma: 12

→ Deliberate misspellings:

Replica/fake: 4

(m0rgage, med1cine, etc.)

→ Steal passwords: 53

→ Unusual email routing:

Other: 31

→ Unusual (spamming) punctuation:

7.1. Error analysis

7.1.1. Numerical evaluation (Error matrix)

→ 하나의 single real number로 평가하는 방법 ex) error rate

Should discount/discounts/discounted/discounting be treated as the same word?

Can use “stemming” software (E.g. “Porter stemmer”) universe/university.

Error analysis may not be helpful for deciding if this is likely to improve performance. Only solution is to try it and see if it works.

→ stemming을 한 것과 stemming을 안 한 것을 cross validation set로 error rate를 본다.

Need numerical evaluation (e.g., cross validation error) of algorithm's performance with and without stemming.

Without stemming: 5% error With stemming: 3% error

Distinguish upper vs. lower case (Mom/mom): 3.2%

7.1.2. Evaluation for Skewed classes (classification)

- 정의 : skewed classes에서는 위의 single real number error matrix를 쓸 수가 없다. 그 이유는 암 발병률과 같이 원래 확률이 낮은 경우에 learning algorithm 자체의 error rate와 비슷할 수 있기 때문이다. $y=0$ 으로 고정시킨 algorithm에서도 error rate가 0.5%이다. 하지만 이 알고리즘이 좋은 알고리즘은 아니다!! 따라서 다른 알고리즘을 사용할 필요가 있다.

Cancer classification example

Train logistic regression model $h_\theta(x)$. ($y = 1$ if cancer, $y = 0$ otherwise)

Find that you got 1% error on test set.
(99% correct diagnoses)

Only 0.50% of patients have cancer.

→ skewed classes.

```
function y = predictCancer(x)
    → y = 0; %ignore x!
    return
```

→ 0.5% error
→ 99.2% accy (0.8% error)
→ 99.5% accy (0.5% error)

→ test set의 error가 1%인데, 발병률이 0.5%이면 확인하기가 어렵지 않겠느냐. test set에서 error가 나는 이유는 대표성이 저하될 수 있기 때문이다.

- 해결법 : Precision/Recall

- 아까와는 달리 암에 걸린 것을 사건 1로 둔다.
- Precision의 정의는 predicted positive 중에서 true positive의 수를 세는 것이다.
- Recall의 정의는 진짜 암에 걸린 사람들 중에서 암으로 정확하게 진단한 경우(true positive)를 세는 것이다.

Precision/Recall

$y = 1$ in presence of rare class that we want to detect

		Actual class		→ Precision	
		1	0	(Of all patients where we predicted $y = 1$, what fraction actually has cancer?)	
Predicted class	1	True positive	False positive	$\frac{\text{True positives}}{\#\text{predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$	
	0	False negative	True negative	→ Recall (Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)	

→ recall을 사용하게 되면 위에서 $y=0$ 만을 predict하는 알고리즘을 걸러낼 수 있다. 왜냐하면 진짜 암환자 중에서 암이 있다고 진단을 한 비율이 0이기 때문이다.

		Actual class	
		1	0
Predicted class	1	True Positive	False Positive
	0	False Negative	True Negative

$$\text{Precision} = \frac{\text{True positives}}{\# \text{ predicted as positive}} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

$$\text{Recall} = \frac{\text{True positives}}{\# \text{ actual positives}} = \frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$$

- precision and recall의 조절 방법

- modify the classification algorithm's threshold to gain higher precision or recall
- $y = 1, y = 0$ 을 정하는 기준값을 변경하면 이를 가능하게 할 수 있다.

Logistic regression: $0 \leq h_\theta(x) \leq 1$

Predict 1 if $h_\theta(x) \geq 0.5$ ~~or 0.9~~

Predict 0 if $h_\theta(x) < 0.5$ ~~or 0.9~~

(1) Precision을 높이는 방법 :

threshold를 높인다. 그러면 분모가 줄어든다. 암환자로 진단된 환자가 줄어드니까. 정확도가 높아지는 것이다.

Suppose we want to predict $y = 1$ (cancer)
only if very confident.

→ Higher precision, lower recall.

Suppose we want to avoid missing too many cases of cancer (avoid false negatives).

→ Higher recall, lower precision.

(2) Recall을 높이는 방법 :

threshold를 낮춘다. 그러면 분자 true positive는 높일 수 있다. false negative를 낮출 수 있다.

- threshold를 정하는 방법

- precision, recall의 tradeoff가 있는데 어떤 값을 가진 algorithm을 사용하는 게 나을까? single number evaluation matrix로 만들면 바로 판단할 수 있다.

F₁ Score (F score)

How to compare precision/recall numbers?

	Precision(P)	Recall (R)
→ Algorithm 1	0.5	0.4
→ Algorithm 2	0.7	0.1
Algorithm 3	0.02	1.0

- Threshold 고르는 방법

- Sol1 : 평균을 구해보자. 평균이 좋은 single number matrix가 될 수 있을까? 아니다. $y=1$ 만 예측하는 알고리즘을 해보면 precision이 낮고, recall이 높은 algorithm 3과 같은 결과가 나오는데 이 결과가 average가 가장 높은 것을 보면 별로 좋지 못한 방법이라는 것을 알 수 있다.

- Sol2 : F1 Score ($2 * PR / P + R$) → 이 방법을 사용하면 0에 가까우면 가중치를 높게 주기 때문에 P나 R이 0에 가까우면 F score도 0에 가까워진다. 따라서 P나 R이 0이나 1로 극단적인 경우를 걸러낼 수 있다.

→ cross validation set을 통해서 F1 score를 maximization 할 수 있는 threshold를 고르면 된다.

F₁ Score (F score)

How to compare precision/recall numbers?

	Precision(P)	Recall (R)	Average	F ₁ Score
→ Algorithm 1	0.5	0.4	0.45	0.444 ↙
→ Algorithm 2	0.7	0.1	0.4	0.175 ↙
Algorithm 3	0.02	1.0	0.51	0.0392 ↙

Average: $\frac{P+R}{2}$

$F_1 \text{ Score: } 2 \frac{PR}{P+R}$

$P=0 \text{ or } R=0 \Rightarrow F_{\text{score}} = 0$

$P=1 \text{ and } R=1 \Rightarrow F_{\text{score}} = 1$

Predict $y=1$ all the time

7.2. Large data for machine learning

how much data to train on.

- 어떤 조건에서 많은 데이터가 도움이 될까?

- 예시 : logistic regression을 했을 때, training set size가 큰 것이 좋다라는 결론을 낸 연구

Designing a high accuracy learning system

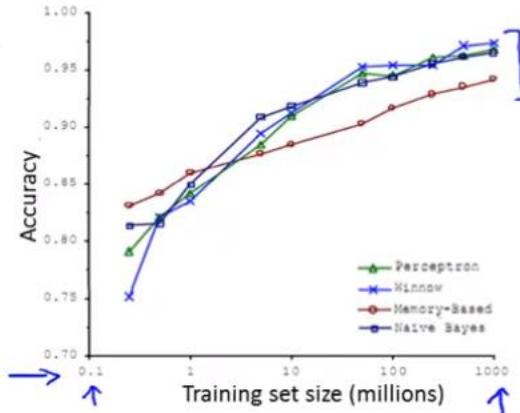
E.g. Classify between confusable words.

{to, two, too} {then, than}

→ For breakfast I ate two eggs.

Algorithms

- - Perceptron (Logistic regression)
- - Winnow
- - Memory-based
- - Naïve Bayes



"It's not who has the best algorithm that wins.

It's who has the most data."

- Large data가 효과가 있을 조건

(1) Many parameter!!!! feature $x \in \mathbb{R}^{n+1}$ 가 y 를 예측하기 위한 충분한 정보가 있다는 전제가 필요하다.

Example: For breakfast I ate two eggs. ←

Counterexample: Predict housing price from only size ← (feet²) and no other features.

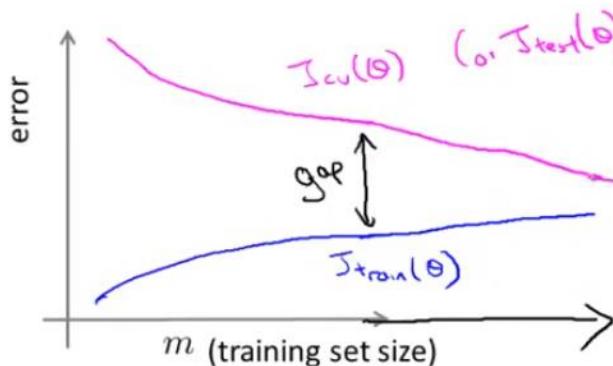
→ 뭐가 data가 충분한 것인지 어떻게 아는가? : human expert한테 물어봤을 때, y 를 잘 예측할 수 있으면 정보가 충분한 것이다.

- 정보가 많다는 것은 무엇일까? parameter가 많다는 뜻이다.

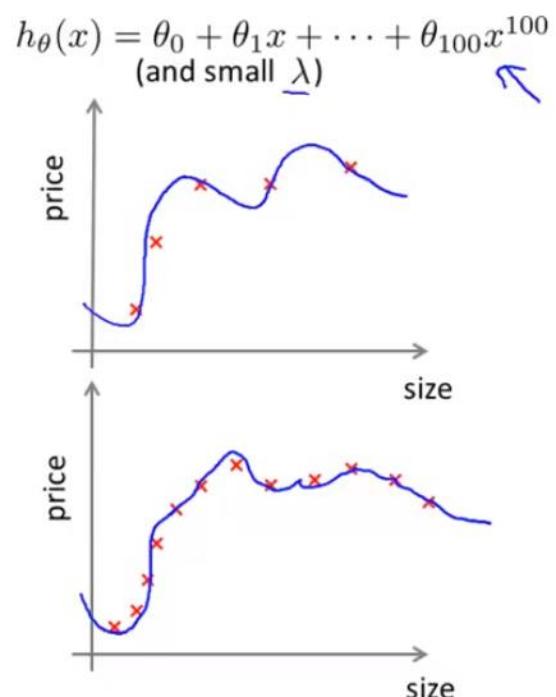
Use a learning algorithm with many parameters (e.g. logistic regression/linear regression with many features; neural network with many hidden units). low bias algorithms. ←

<reminder>

High variance

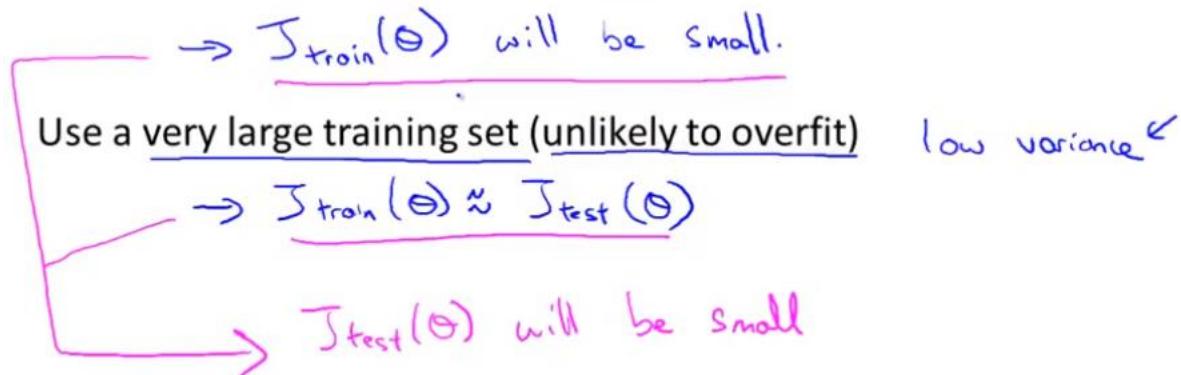


If a learning algorithm is suffering from high variance, getting more training data is likely to help.



→ 그러나 overfitting 문제가 발생할 가능성도 있다. 따라서 (2)번 조건이 필요한 것이다.

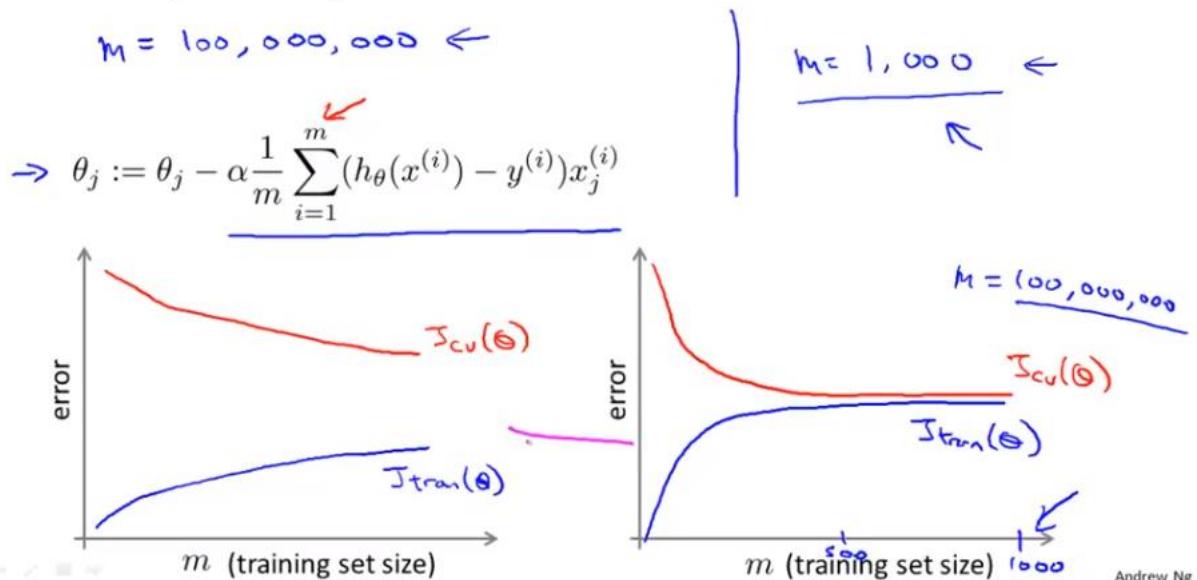
(2) Use very large training set → unlikely to overfit



(1) (2) 조건이 다 갖춰지면 large data를 가졌을 때 high performance learning algorithm에 도달할 수 있다.

- 정리 : training data set(m) = 1억개가 있다고 하더라도 1000까지 해본다. error curve를 그려보고 high variance를 보인다면(왼쪽 그래프) data set을 많이 하는 것이 좋겠다는 결론을 내릴 수 있게 된다. 만약 오른쪽 그래프처럼 나온다면 feature를 더하거나, neural network의 hidden layer를 더하거나 해서 그래프가 왼쪽과 같이 나오도록 할 수 있다.

Learning with large datasets



7.2.1. large data set으로 training을 할 때, 알고리즘을 효율화할 수 있는 3가지 방법

- (1) Stochastic gradient descent
- (2) mini-batch gradient descent
- (3) Map reduce and data parallelism

(1) Stochastic gradient descent

- 사용 목적 : 기존에 우리가 알고리즘을 만드는 방법은 우선 minimize를 하고 싶은 cost function을 찾고 그것을 minimize할 수 있는 parameter를 찾는 것이었다. 그 때, gradient descent를 썼었다. 하지만 이는 computationally expensive해질 수 있다.

- 예시 : linear regression을 예로 사용하지만 logistic regression, neural network에서도 가능하다.

Linear regression with gradient descent

$$\rightarrow h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

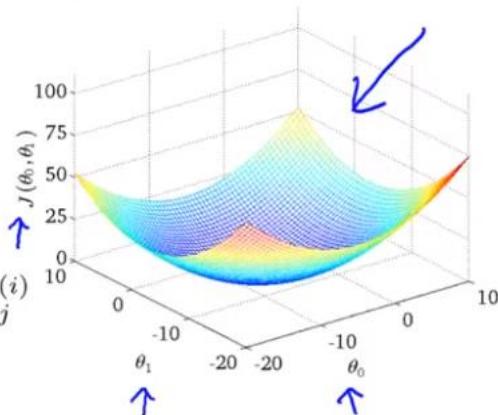
$$\rightarrow J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

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

(for every $j = 0, \dots, n$)

}



$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

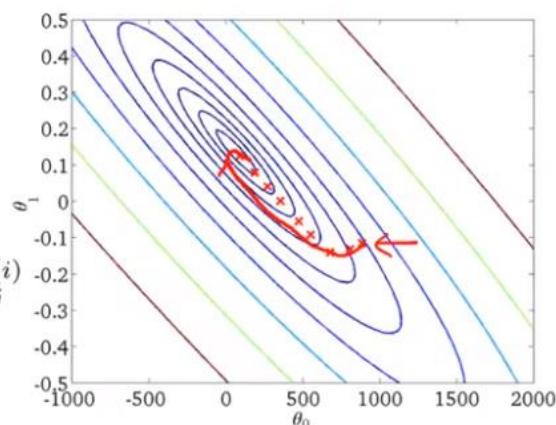
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

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

(for every $j = 0, \dots, n$)

}



- 문제 설정 : batch gradient descent이라고 부른다.(기준에 우리가 배운 것) 근데 $m = 3\text{억}$ 정도 되면 3억 을 메모리에 넣고 한번 gradient descent를 하고 다시 3억 개를 메모리에 넣고 더해서 gradient descent 한 발짝가고 해야 한다. 이는 비효율적이다.

- 문제 해결 : m 개 data set을 다 더해서 빼서 parameter를 update하는게 아니라, 만약에 data set 1개씩해서 update를 하게 되면 어떻게 될까?

Batch gradient descent

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

Repeat {

$$\rightarrow \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m \underbrace{(h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\frac{\partial}{\partial \theta_j} J_{train}(\theta)}$$

(for every $j = 0, \dots, n$)

}

Stochastic gradient descent

$$\rightarrow \underbrace{cost(\theta, (x^{(i)}, y^{(i)}))}_{\frac{\partial}{\partial \theta_j} J_{train}(\theta)} = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\rightarrow J_{train}(\theta) = \frac{1}{m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset. ←

2. Repeat {

$$\left[\begin{array}{l} \text{for } i=1, \dots, m \{ \\ \quad \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{j=0}^n (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \\ \quad \text{for } j=0, \dots, n \} \\ \} \end{array} \right]$$

$$\rightarrow \underbrace{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots}_{\frac{\partial}{\partial \theta_j} cost(\theta, (x^{(i)}, y^{(i)}))}$$

Andrew Ng

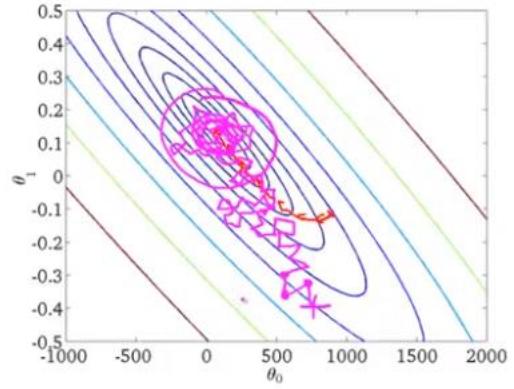
- 정리 : 1개씩 하더라도 batch gradient descent로 했을 때와 가까운 parameter를 얻을 수 있게 된다. m 이 엄청 크다면 1번 ~ 10번 정도만 하면 된다.

→ 질문! 정말 그런가?? 너무 적은거 아닌가? 왜 m 이 큰 경우에는 이게 가능한 것인가?

Stochastic gradient descent

- 1. Randomly shuffle (reorder) training examples

- 2. Repeat {
 - for $i := 1, \dots, m$ {
 - $\theta_j := \theta_j - \alpha(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$
 - (for every $j = 0, \dots, n$)



- 주의 사항

1) Checking for convergence : J cost function의 값이 converging하는지를 보아야 한다. stochastic을 할 때는 그렇게 computing cost가 비싸지 않으므로 gradient descent를 하면서 바로 할 수 있다. cost를 plotting할 때 몇 개씩의 평균을 plotting할지를 정하고 보면 된다.

Checking for convergence

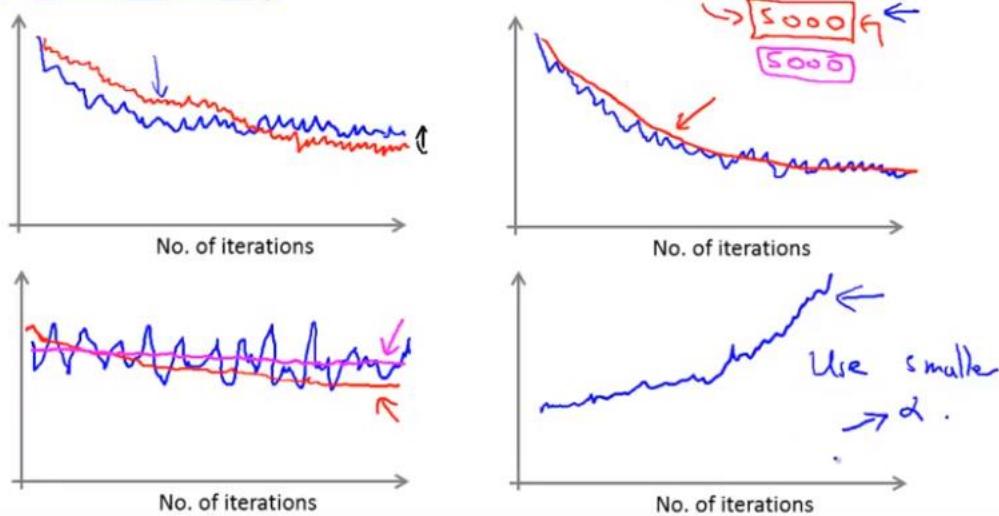
- Batch gradient descent:
 - Plot $J_{train}(\theta)$ as a function of the number of iterations of gradient descent.
 - $$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$
 $M = 300, 500, 500$

- Stochastic gradient descent:
 - $cost(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2$ $\rightarrow (x^{(i)}, y^{(i)}), (x^{(i+1)}, y^{(i+1)}), \dots$
 - During learning, compute $\frac{cost(\theta, (x^{(i)}, y^{(i)}))}{\uparrow}$ before updating θ using $(x^{(i)}, y^{(i)})$.
 - Every 1000 iterations (say), plot $cost(\theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1000 examples processed by algorithm.

- cost를 구할 때 몇 개씩 평균해서 보여주는가에 따라서 learning curve가 달라진다. 만약에 1000개씩 평균하다가 5000개씩 평균하면 smoother graph가 나타난다.

Checking for convergence

Plot $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$, averaged over the last 1000 (say) examples



2) learning rate의 조절

- learning rate(alpha)를 무엇을 쓰는가에 따라서도 graph 양상이 조금 달라질 수 있다.(learning rate이 작으면 cost가 떨어지는 것이 좀 더 느리다.)
- 만약에 diverging하면 alpha를 줄여보자.
- iteration이 증가할 때 learning rate를 줄이면 stochastic gradient descent가 마지막에 날뛰는 것을 줄일 수 있다.

Stochastic gradient descent

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

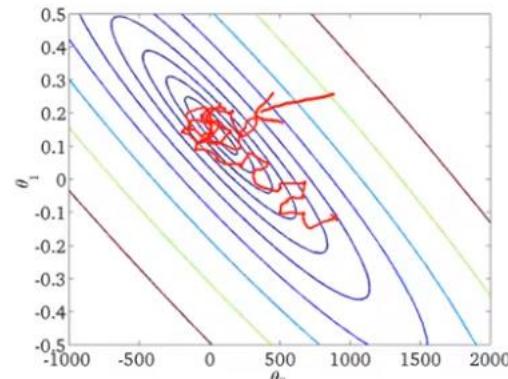
$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.

2. Repeat {

```

    for i := 1, ..., m
        {
             $\theta_j := \theta_j - \alpha(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$ 
            (for j = 0, ..., n)
        }
    }
}
```



Learning rate α is typically held constant. Can slowly decrease α over time if we want θ to converge. (E.g. $\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}$)

Andrew Ng

(2) mini batch gradient descent ← stochastic와 batch 사이 ($b=1$, $b=m$ 사이)

- 목적 : iteration을 한번 할 때마다 1개만(stochastic)하는 게 아니라, batch를 2-100개씩 하는 것이다. 왜냐하면 vectorization하게 되면 훨씬 효율적으로 빠르게 계산할 수 있게 되기 때문이다. parallel하게 10개씩 할 수 있다. $b = 10$ 개로 잡는 것이 좋다.(Andrew 추천)

Mini-batch gradient descent

→ Batch gradient descent: Use all m examples in each iteration

→ Stochastic gradient descent: Use 1 example in each iteration

Mini-batch gradient descent: Use b examples in each iteration

$b = \text{mini-batch size}$. $b = 10$. $\frac{m}{b} = \frac{1000}{10} = 100$

Get $b = 10$ examples $(x^{(i)}, y^{(i)}), \dots (x^{(i+9)}, y^{(i+9)})$

$\rightarrow \theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) \cdot x_j^{(k)}$

$i := i + 10$

- 목적식 :

Mini-batch gradient descent

Say $b = 10, m = 1000$.

Repeat {

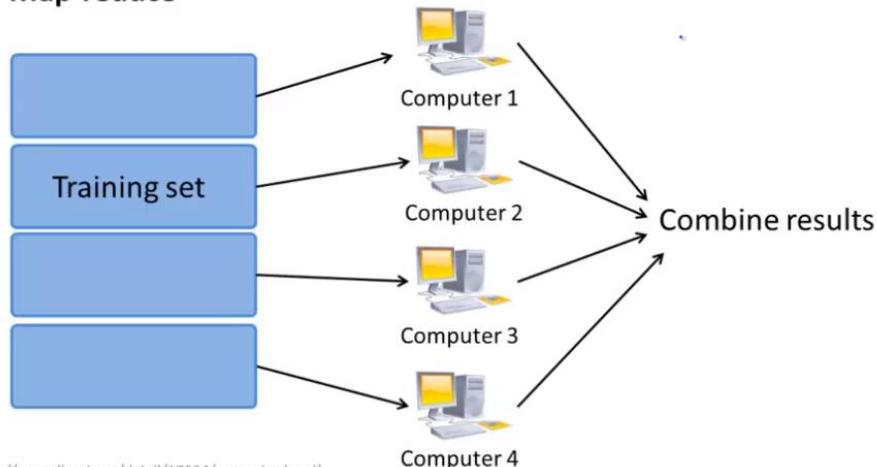
→ for $i = 1, 11, 21, 31, \dots, 991$ {
→ $\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$
(for every $j = 0, \dots, n$)
}
}

(3) Map reduce

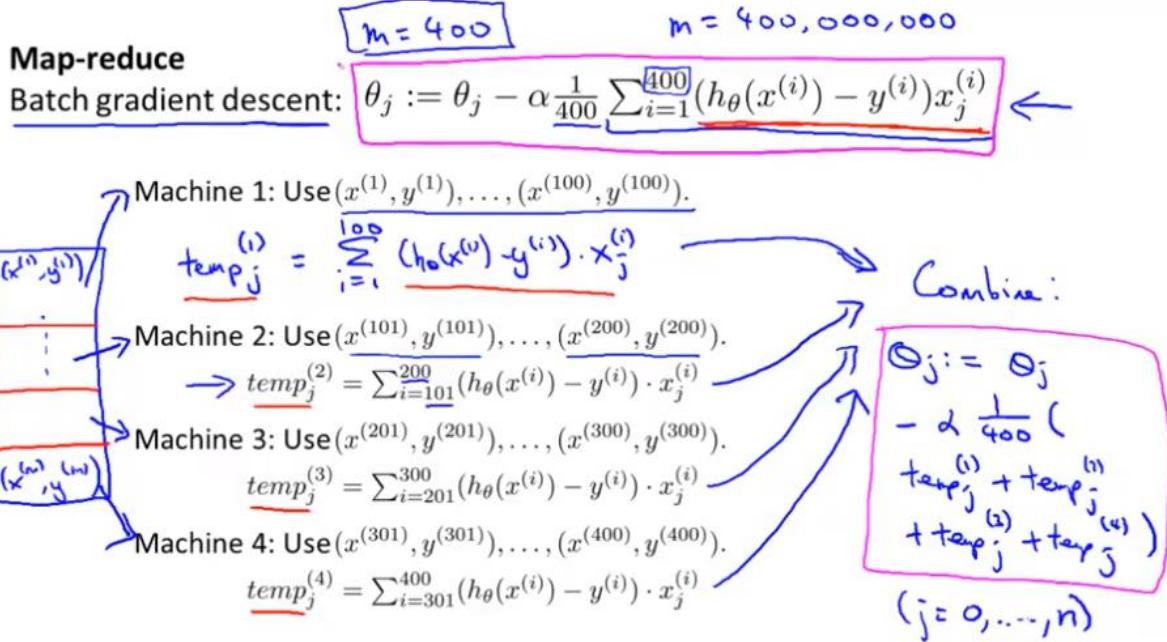
- gradient descent 만큼이나 효율화에 중요하다.

- 400개의 sum을 계산해야 하면 100, 100, 100, 100개로 나눠서 4개의 computer에서 나눠서 계산하면 더 빨라진다.

Map-reduce



- 예시 :



[Jeffrey Dean and Sanjay Ghemawat]

Andrew Ng

- 한 개의 컴퓨터 : core가 여러 개인 computer에서 컴퓨터가 여러 개가 없어도 가능하다. data를 수신, 전송하는 시간도 적게 든다.(network latency) numerical linear library에서는 이것을 한 컴퓨터에서 저절로 해준다.

7.3. Large data generation (artificial data synthesis) → 없으면 만들면 되잖아!!!!

Discussion on getting more data

1. Make sure you have a low bias classifier before expending the effort. (Plot learning curves). E.g. keep increasing the number of features/number of hidden units in neural network until you have a low bias classifier.
2. “How much work would it be to get 10x as much data as we currently have?”
 - Artificial data synthesis
 - Collect/label it yourself
 - “Crowd source” (E.g. Amazon Mechanical Turk)

- **How much work would it take to get 10x much data?!!!** 이게 정말 좋은 질문이다!!!

- 직접 collect, label해도 그렇게 오래 걸리지 않는다면!
- artificial data synthesis → unlimited data for training data
 - font를 바꾸기

Artificial data synthesis for photo OCR



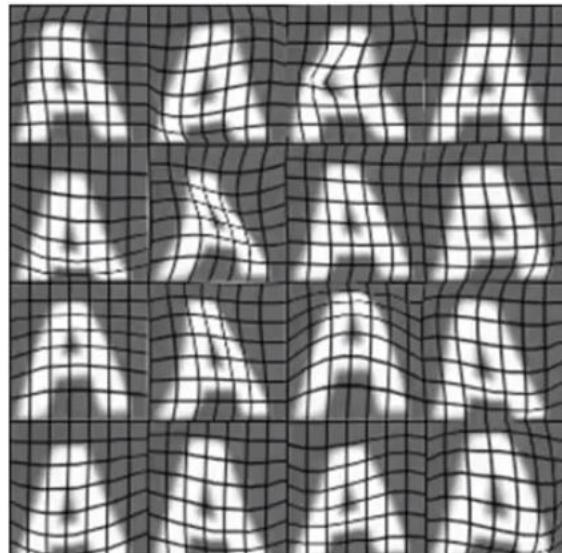
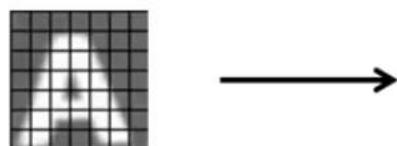
Real data



Synthetic data

- artificial warping

Synthesizing data by introducing distortions



- audio distortion

Synthesizing data by introducing distortions: Speech recognition

Original audio:

Audio on bad cellphone connection

Noisy background: Crowd

Noisy background: Machinery

7.4. Online learning

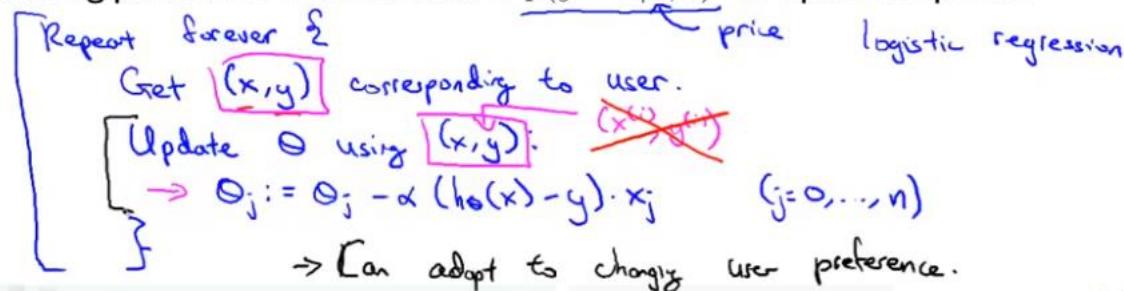
- 목적 : online 상에서 데이터가 계속 들어오는 경우라면 learning을 계속 하게 하면 된다. 데이터를 얻고 그것을 통해서 learning하고, 그 data set을 버리면 된다. Small number of users라면 그것을 모아서 fixed learning set를 사용하면 된다.

만약에 경제상황이 바뀌거나(탄력성이 바뀌거나) 하면 adapt to changing user preference를 할 수 있다.

Online learning

Shipping service website where user comes, specifies origin and destination, you offer to ship their package for some asking price, and users sometimes choose to use your shipping service ($y = 1$), sometimes not ($y = 0$).

Features x capture properties of user, of origin/destination and asking price. We want to learn $p(y = 1|x; \theta)$ to optimize price.



- 예시 2 : user search query (predicted CTR click through rate)

Other online learning example:

Product search (learning to search)

User searches for "Android phone 1080p camera" \leftarrow

Have 100 phones in store. Will return 10 results.

$\rightarrow x$ = features of phone, how many words in user query match name of phone, how many words in query match description of phone, etc.

$\rightarrow y = 1$ if user clicks on link. $y = 0$ otherwise.

\rightarrow Learn $p(y = 1|x; \theta)$. \leftarrow predicted CTR

\rightarrow Use to show user the 10 phones they're most likely to click on.

Other examples: Choosing special offers to show user; customized selection of news articles; product recommendation; ...

8. Machine learning pipelines → Photo OCR(Optical character recognition) example

- 개요 : machine learning이 pipeline의 각 step에서 쓰이는 것이다.

Photo OCR pipeline

1. Text detection



2. Character segmentation



3. Character classification

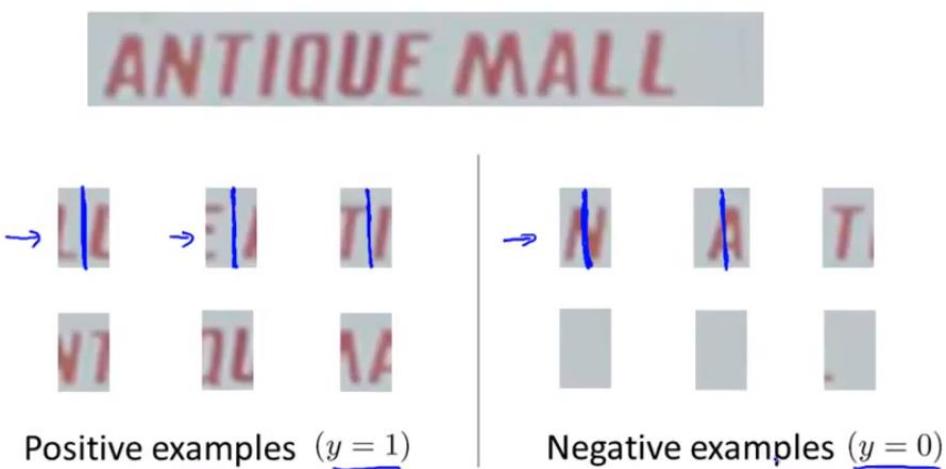


1) text detection



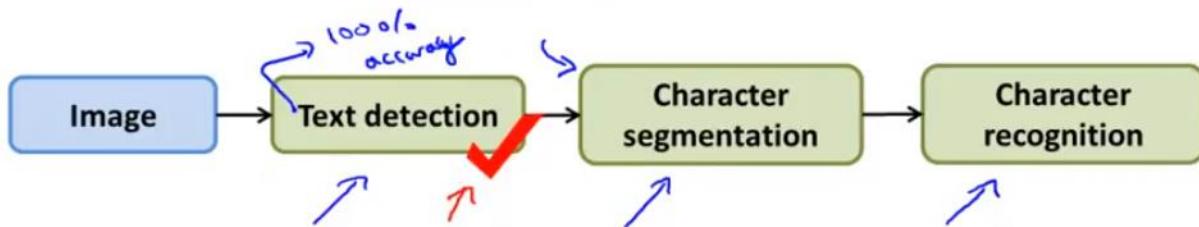
2) segmentation

1D Sliding window for character segmentation



8.1. ceiling analysis : 어떤 부분의 개발에 시간을 할당해야 하는지를 알려주는 분석

Estimating the errors due to each component (ceiling analysis)



What part of the pipeline should you spend the most time trying to improve?

Component	Accuracy
Overall system	72%
Text detection	89%

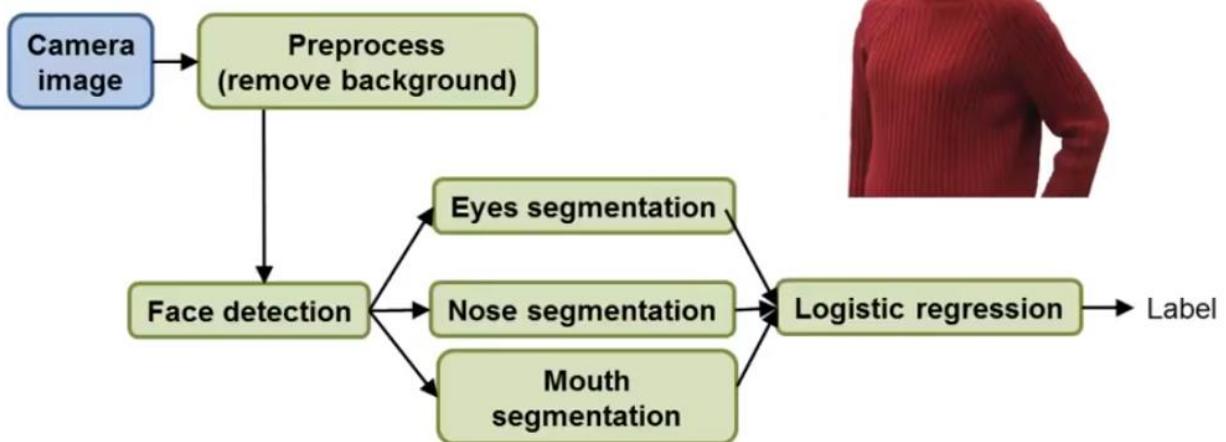
→ 여기서 체크된 부분을 인간이 스스로 해본다. 100% 정확할 때, accuracy가 어떻게 되는가? 72%에서 89%로 증가하는구나.

Component	Accuracy
Overall system	72%
Text detection	89%
Character segmentation	90%
Character recognition	100%

→ 아하! text detection 부분을 열심히 개선하면 되겠구나. 여기에 좀 더 engineer를 많이 투입하자.

Another ceiling analysis example

Face recognition from images
(Artificial example)



Component	Accuracy
Overall system	85%
Preprocess (remove background)	85.1% 5.1%
Face detection	91% 5.9%
Eyes segmentation	95% 4.0%
Nose segmentation	96% 1.0%
Mouth segmentation	97% 1.0%
Logistic regression	100% 3.0%

→ 여기서 중요한 것은 이 논리 구조이다. 이렇게 만약에 과학을 구획화할 수 있으면 어디에 집중할지를 좀 명확히 할 수 있지 않을까? mapping이 역시 중요한 것이다.

- don't trust your gut feeling!!!
- systematic approach is better

Thank you.
- Andrew Ng