

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

Parameters : θ_0, θ_1

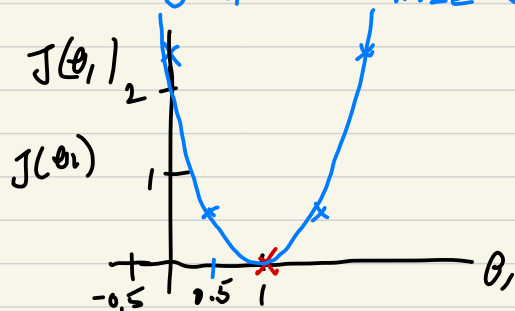
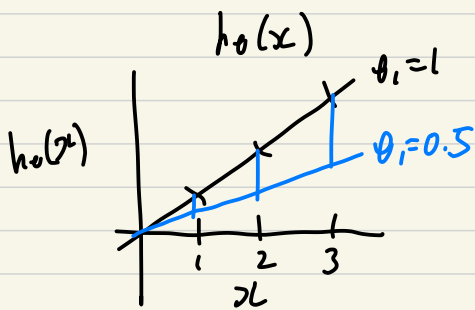
Cost function

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

Goal: minimize θ_0, θ_1 $\underbrace{J(\theta_0, \theta_1)}_{\text{cost-function} = \text{Squared error function}}$

hypothesis $\begin{matrix} h_{\theta}(x) : \text{직접} \\ y : \text{실측} \end{matrix}$

$\theta_0 = 0$ 일때 simplified cost function \rightarrow goal minimize $J(\theta_1)$



$\theta_1 = 1$ 일때 $J(\theta_1) = ?$

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

$$= \frac{1}{2 \cdot 3} (0^2 + 0^2 + 0^2) = 0$$

$$3) J(0) = \frac{1}{2 \cdot 3} [1^2 + 2^2 + 3^2]$$

$$= \frac{1}{6} \cdot 14 = 2.33 \dots$$

* 실측치와 가장 잘 맞는 $\theta_1 = 1$ 에서 $J(\theta_1)$ 이 최소가 되는 것을 확인함.

$$2) J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (3-1.5)^2]$$

$$= \frac{1}{2 \cdot 3} (3.5) = 0.58$$

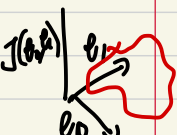
(가설치해보기)

Gradient Descent : $J(\theta_0, \theta_1, \dots, \theta_n)$ 을 minimize 하기 위한 방법
+ 꼭 선형회귀가 아니더라도 이방식 사용가능

• start with some θ_0, θ_1

• keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$

• 만약 꼭대기 ($J(\theta_0, \theta_1)$) 이 높은 리전부터 아래로 내려가기 위한 제일 최단거리로 계속 내려감.



Gradient Descent algorithm

repeat until convergence

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

learning rate (학습률)

: 내려가기 위한 보폭을 어느정도로 설정할 것인가.

$\alpha \uparrow$ - 공격적인 가르기 가능

c.f)

$a := b$ means b 를 a 에 쓰인다.

$a = b$ T or F

for $j=0$ and $j=1$

simultaneously update θ_0 & θ_1
⊕ explain

⊕ explain

θ_0, θ_1 을 구하고 나서 새로운 θ_0, θ_1 을 갖고 다시 계산하는 것.

θ_0 구했다면 θ_1 구할 때 새로운 θ_0 을 갖고 계산해야 함.

temp0 := 공식

temp1 := 공식

$\theta_0 := \text{temp0}$

$\theta_1 := \text{temp1}$

→ 이걸 끝내고

← 그 다음 이걸

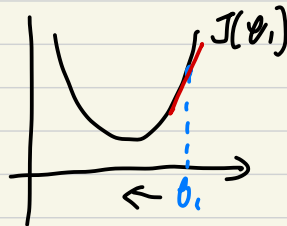
→ simultaneous update

θ_0, θ_1 를 동시에.

미분계수

$$\alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \text{ 가변}$$

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



$$\theta_1 \in \mathbb{R}$$

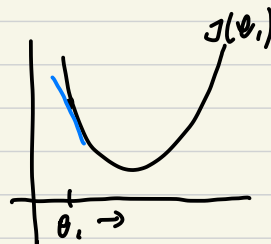
θ_1 바뀌어 가요

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

≈ 0

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

θ_1 이 왼쪽으로 이동: $\frac{d}{d\theta_1} J(\theta_1)$ 이
가까워질



$$\frac{d}{d\theta} J(\theta_1) < 0$$

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

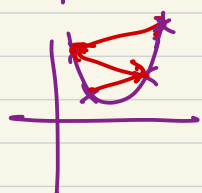
θ_1 이 오른쪽으로 이동하게 됨: $\frac{d}{d\theta} J(\theta_1)$ 이
가까워질

: α 이 작으면



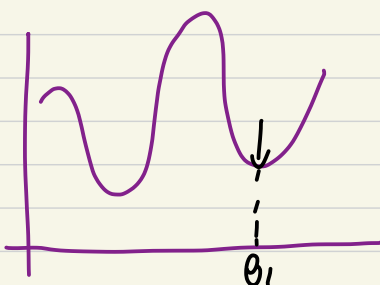
작게 움직이고

α 이 크면



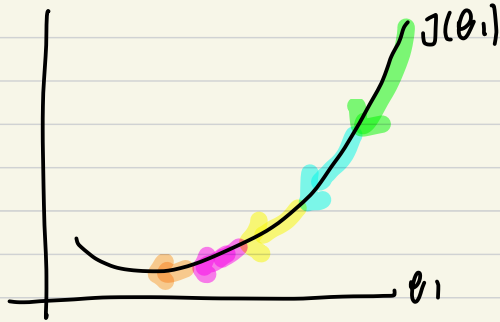
크게 움직이는데
지나칠 수도 있다.

c.f)



θ_1 의 값이 0이면 θ_1 이 local minimum이라 움직이지 X.

· α 가 고정되어 있더라도 계산을 반복할수록 미분계수의 절대값이 작아지기 때문에
최저값에 가까워질수록 이동 거리가 작아진다. $\therefore \alpha$ 를 변경하지 않아도 괜찮다.



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

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

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

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

Gradient Descent의 이용.

"Batch" gradient Descent : 점단 가르기 가능

Vector: $n \times 1$ matrix

vector $y = \begin{bmatrix} \quad \end{bmatrix}_n$ 표시: \mathbb{R}^n

matrix 표시 $\mathbb{R}^{m \times n}$

y_i = i th element

1-indexed VS 0-indexed
index start는 1 이냐 0 이냐

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$y = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \end{bmatrix}$$

행렬 표시 | 다들 다 사용
vector, scalar 표시 | 소문자 사용 \rightarrow 항상 그런건 아니지만 일반적

Multiple features (variables)

n = number of features

몇번이든 $x^{(i)}$ = input (features) of i^{th} training sample
 2번째 $x_j^{(i)}$ = value of feature j in i^{th} training sample
 몇번이든

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

E.g $h_0(x) = 80 + 0.1x_1 + 0.01x_2 + 3x_3 - 2x_4$

쉬운 이해를 위해 $x_0 = 1$ 로 놓자. ($x_0^{(i)} = 1$)

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

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

$$\underbrace{\begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix}}_{\substack{\theta^T \\ (n+1) \times 1 \text{ matrix}}} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

Hypothesis: $h_0(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$
 $= \theta^T x$ (이로 간편히 쓸 수 있다)

Parameters: $\theta_0, \theta_1, \dots, \theta_n$ 이므로 θ 벡터 그 자체라고 놓을 수 있다.

Cost function: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$

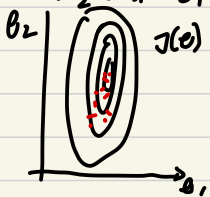
Gradient Descent

$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_j^{(i)} \text{ simultaneously} \\ \text{update } \theta_j \text{ for } j=0, \dots, n \end{array} \right\}$$

Feature Scaling

: make sure features are on a similar scale

E.g $x_1 = \text{size (0-2000 ft}^2\text{)}$
 $x_2 = \# \text{ of bedroom (1-5)}$



→ it takes more time to find $\min J(\theta)$
 큰 값일수록 θ 가 천천히 이동하므로
 범위의

So, $x_1 = \frac{\text{size}}{2000}$, $x_2 = \frac{\# \text{ of bedroom}}{5}$



feature들이 비슷한 범위의 값을 갖도록 만들자. $-3 \text{ to } 3$, $-\frac{1}{3} \text{ to } \frac{1}{3}$ 같이

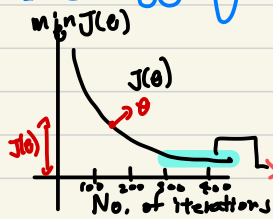
- mean normalization 도 있음

$$\left. \begin{array}{l} x_1 = \frac{\text{size} - 1000}{2000} \\ x_2 = \frac{\# \text{ bedroom} - 2}{5} \end{array} \right\} \text{이렇게.}$$

or $\frac{x_1 - \mu_1}{\sigma_1 (\text{max} - \text{min})}$ 도 되고 standard deviation 도 된다.

여기 있는 것만 다 근사값이므로 엄청 정확할 필요는 X

Debugging: make sure gradient descent working correctly



: 반복할 때마다 $J(\theta)$ 그래프

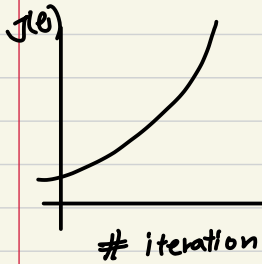
correct: $J(\theta)$ 가 점점 줄어들어야 한다, 수렴해야 함.

이 경우는 거의 수렴, $J(\theta)$ 가 수렴하고 있다.

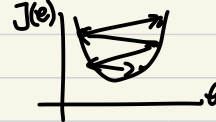
수렴하는지 판단하는 법

: $J(\theta)$ 의 감소량이 10^{-3} 보다 작을 때

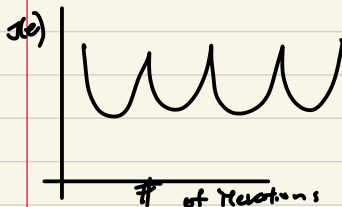
이 값에 도달할 즈음 불명확해서 이 방법을 잘 만들, 그래프 그리는 것이 나을



: 이런 그래프는 α 가 커서



하나라 그런것으로 결과를 작게 하면 된다.



: 이 경우도 위와 같음.

Summary

: if α is too small, slow convergence

: " large: $J(\theta)$ may not decrease

To choose α try: 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1..

가능한 가장 작은 값, 큰 값을 선택하고 작은 값에서 한 단계 큰 값, 큰 값에서 한 단계 작은 값을 선택해가며 α 를 결정한다.

집의 길이, 너비가 있을 때 이 둘을 곱해서 넓이라는 하나의 feature를 가지는 모델을 만들 수 있다.

Normal Equation 은 feature scaling 없이 만들어짐

x_0, x_1, \dots, x_n 값들 $X = \begin{bmatrix} 1 & x_1 & \dots & x_n \end{bmatrix}$ 이 한 vector씩 다 넣는다

$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$ $n \times 1$ vector

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

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

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad \text{design matrix} \quad X = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \text{---} (x^{(n)})^T \text{---} \end{bmatrix}$$

If $x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix}$ $X = \begin{bmatrix} 1 & x_1^{(1)} \\ \vdots & \vdots \\ 1 & x_m^{(1)} \end{bmatrix} \quad m \times (n+1)$ $y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$

ocean: $\text{pinv}(X^T * X) * X^T * y$

Gradient Descent (VS) Normal Equation
- m training examples, n features

- Need to choose α
- Needs many iterations

- Works well even when n is large

이러한 algorithm은 꼭 필요

For large n

$$n > 10^6$$

- No need to choose α
- Don't need to iterate

- Need to compute $(X^T X)^{-1}$

- Slow if n is very large

For small n

$$n \approx 1000 \text{ 정도는 ok}$$

* $X^T X$ is non-invertible 일때

(1) linearly dependent

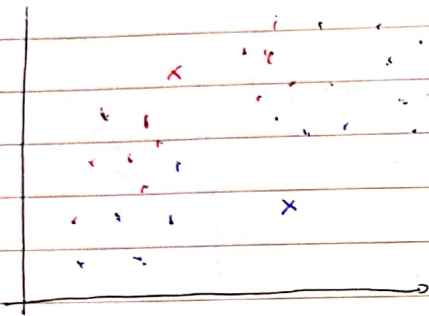
(2) too many features, $(m \leq n)$, 작은 양의 데이터로 파악하기 어렵다.

이럴 땐 불필요한 feature가 있는지 살펴보기

1) Unsupervised learning : y가 없음. unlabeled sample.

- 예시 : ① 고객 세분화 : Market segmentation → 각 그룹에게 물품 팔기
 ② Social network analysis : 이메일을 주고 받은 것을 보고 그룹이 형성
 ③ Organize computing clusters : 데이터 센터 체계화
 ④ Astronomical data analysis

2) K-means algorithm. ① cluster assignment ② move centroid



- 1) randomly initialize two points → cluster centroids
 2개인 이쁜 group 2개로 나눌 수 있나
 2) 각 점들은 제일 가까운 centroid의 색이 될
 3) 같은 색의 점들의 ~~중~~ 평균을 그 그룹의 새로운 점으로 고정
 4) ①부터 다시 반복. 이게 바뀌지 않을 때 끝남.

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)

3) K-means algorithm.

Initialize K cluster ^{centroids} ~~algorithm~~ $\mu_1, \dots, \mu_K \in \mathbb{R}^n$

repeat for $i=1$ to m

$c^{(i)} := \text{index (from 1 to K) of cluster centroid closest to } x^{(i)} \leftarrow \min_k \|x^{(i)} - \mu_k\|$

x 와의 거리가 제일 작은 μ_k 의 k 가 $c^{(i)}$ 가 대입.

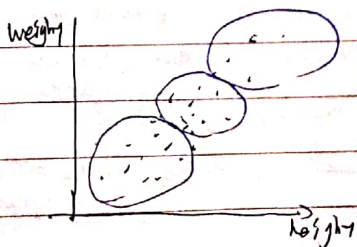
for $k=1$ to K

$\mu_k := \text{average (mean) of points assigned to cluster } k$

ex) $x^{(1)}, x^{(3)}, x^{(4)}, x^{(6)} \rightarrow \mu_2$ group $\Rightarrow c^{(1)}=2, c^{(3)}=2, c^{(4)}=2, c^{(6)}=2$

$$\mu_2 = \frac{x^{(1)} + x^{(3)} + x^{(4)} + x^{(6)}}{4} \in \mathbb{R}^n \quad \text{이제 만약 이걸 평균 구하면 K개 cluster를 나눌 수 있게 됨}$$

4) K-means for non-separated clusters



← 3 market clustering

5) Optimization Objective

① knowing what is the optimization objective of k-means

→ debug the learning algorithm

→ k-means is running correctly.

make sure

→ help k-means find better cuts for this

→ avoid the local minima

$c^{(i)}$ = cluster index

μ_k = cluster centroid k

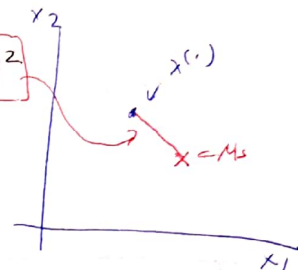
$\mu_{c(i)}$ = $c^{(i)}$ centroid

Optimization objective:

이걸 최소화
 $c^{(i)}, \mu_k$ 를
 찾는다...

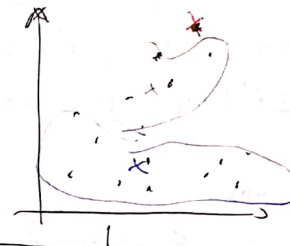
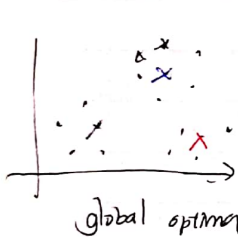
$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

minimize $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$
 $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k$
 Distortion or Cost function
 of the k-means algorithm



6) Random Initialization. (실제 많은 풀, training sample 중에서 centroid를 설정하는 것이 낫다. 앞의 것 2개만 알려주)

Initialization을 어떻게 하느냐에 따라 결과가 달라진다.



→ try multiple Initialization. $k=2 \sim 10$ 일 때 각각, 그 이후는 multiple initialization의 포기는 의미없음
 50 ~ 100

For i = 1 to 100 { Randomly initialize k-means

Run k-means, Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k$

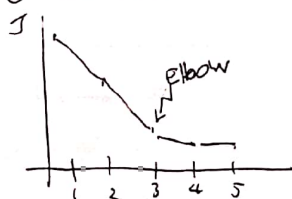
Compute cost function (distortion)

$$\rightarrow J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$$

Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$

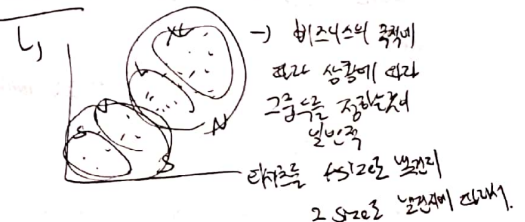
7) Choosing the number of clusters. (실제 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100) 가장 흔한 것.

① Elbow Method



이전보다

이전보다



이전보다



이전보다

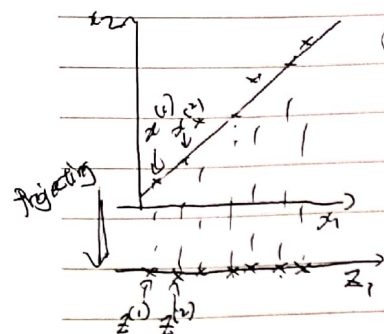
k-means got stuck in a bad local minimum,

You should try re-running k-means with multiple random initialization.

Dimensionality Reduction

1) Data compression: compress the data
 ↳ speed up learning algorithm

Reduce data from 2D to 1D → reduce redundancy



$$x^{(1)} \in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R}$$

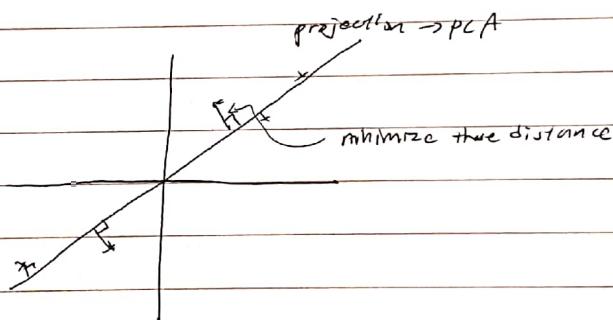
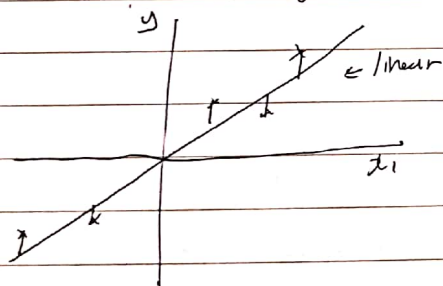
$$x^{(2)} \in \mathbb{R}^2 \rightarrow z^{(2)} \in \mathbb{R}$$

↳ 하나의 숫자만 알고 위치를 알아낼 수 있다 (data量少, 속도 빨라짐)

2) Data visualization

PCA: Try to find lower dimensional surface onto which to project the data.

PCA is not linear regression



→ there is y to predict with x_1

⊗ → there is no special value to predict

all of features treated equally

DATA PROCESSING

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

Preprocessing (feature scaling / mean normalization):

- $\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)}$ Replace each $x_j^{(i)}$ with $x_j^{(i)} - \mu_j$
- scale features to have comparable range of values

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

Reduce data from 2D to 1D

$$x^{(i)} \in \mathbb{R}^2 \rightarrow z^{(i)} \in \mathbb{R}$$

PCA algorithm.

• compute "covariance matrix"

$$\Sigma = \frac{1}{m} \sum_{i=1}^m \begin{bmatrix} x^{(i)} & x^{(i)T} \\ n \times 1 & 1 \times n \end{bmatrix} \quad n \times n$$

• compute eigenvectors of matrix Σ :

$$[U, S, V] = \text{svd}(\Sigma)$$

↳ singular value decomposition

$$U = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(k)} \\ | & | & | & \dots & | \end{bmatrix} \quad \begin{matrix} U \in \mathbb{R}^{n \times n} \\ u^{(1)}, \dots, u^{(k)} \end{matrix}$$

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$$

$$z^{(i)} = \begin{bmatrix} | & | & \dots & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & \dots & | \end{bmatrix}^T x = \underbrace{\begin{bmatrix} -(u^{(1)})^T & \dots & \\ \vdots & & \\ -(u^{(k)})^T & \dots & \end{bmatrix}}_{k \times n} \underbrace{x}_{n \times 1}$$

$n \times k$ matrix

$$z_j = (u^{(j)})^T x$$

$$z = U_{\text{reduce}}^T x$$

$$x = U_{\text{reduce}} \cdot z$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T \leadsto \Sigma = \left(\frac{1}{m}\right) * X' * X$$

Choosing k .

Average squared projection error: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$

Total variation in the data: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

Typically, choose k to be smallest value so that

94% of variance is retained.

90% variance explained 가능.

PCA가 어느정도 잘 되었는지 알 수 있는 것.

Algorithm.

$k=1$ 부터

compute $U_{\text{reduce}}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{\text{approx}}^{(1)}, \dots, x_{\text{approx}}^{(m)}$

check if $\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01$ 이 아니면 $k+1$ 씩 해서 늘려가!

for given k .

SVD 값들이

가장 크다고

이게 더 강한. SVD 1만 하면 보면 돼.

$$1 - \frac{\sum_{i=1}^k s_i}{\sum_{i=1}^n s_i} \leq 0.01 \Rightarrow \frac{\sum_{i=1}^k s_i}{\sum_{i=1}^n s_i} \geq 0.99$$

Applying PCA

Supervised learning speedup.

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

$$x^{(i)} \in \mathbb{R}^{1000}$$

Extract inputs: unlabeled dataset: $x^{(1)}, \dots, x^{(m)} \in \mathbb{R}^{1000}$

$$\downarrow \text{PCA}$$

$$z^{(1)}, \dots, z^{(m)} \in \mathbb{R}^{100}$$

New training set

$$(z^{(1)}, y^{(1)}) \dots (z^{(m)}, y^{(m)})$$

Application of PCA

- Compression

- Reduce memory/disk needed to store data

- Speed up learning algorithm

- Visualization

$$k=2 \text{ or } k=3$$

Bad use of PCA: to prevent overfitting \rightarrow 오히려 오히려 잘못된 방법이 X

Regularization 이 낫다. \leftarrow 이런 y를 한 번 보이기 때문에

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

가치 높은 x와 y가
높아지기 싫어함
PCA보다 낫다.

PCA를 무조건 사용하지 말고, 일단 원래 data 가지고 해보고 속도가 문제거나 원하는 결과가 잘 안 나오면 시도해봐라.