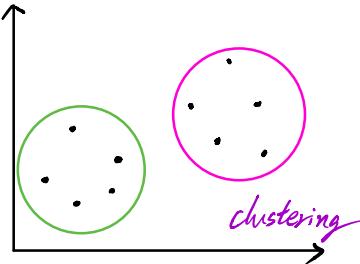


Unsupervised Learning



unlabeled data:

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

K -means : a clustering algorithm

input :- k (no. of clusters)

- training set $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\} \quad x^{(i)} \in \mathbb{R}^n$

optimization objective

$c^{(i)}$ = index of cluster ($1, 2, \dots, K$) to which example $x^{(i)}$ is currently assigned

μ_k = cluster centroid k ($\mu_k \in \mathbb{R}^n$)

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

$$\min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

step: 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 minimize $J(c^{(1)}, \dots, c^{(m)})$ holding μ_1, \dots, μ_K fixed

$\min_k \|x^{(i)} - \mu_k\|^2 \rightarrow c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$

$c_i := k$ [for $k = 1$ to K minimize $J(c^{(1)}, \dots, \mu_k)$

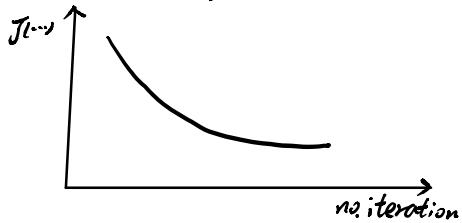
more centroid } $\mu_k :=$ average (means) of points assigned to cluster k

$$\mu_k := \frac{1}{t} \sum_{i=1}^m x^{(i)}$$

$c_{ik} = t =$ no. of $c_i = k$, if $t = 0$, eliminate a cluster

$$\mu_k \in \mathbb{R}^n$$

the process of iteration always decrease $J(\dots)$



random initialization :

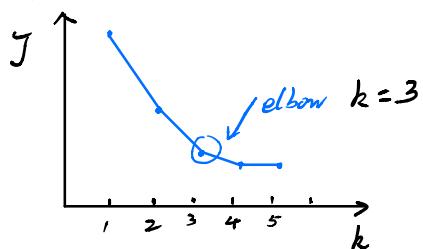
should have $k < m$, randomly pick k training example, set μ_1, \dots, μ_k equal to these k examples, random initialization may reach local optima result, to solve this : using many times of random initialization.

```
for i=1 to 100 {  
    random initialize  $\mu_1, \mu_2, \dots, \mu_k$   
    run k-means  
    compute  $J(\dots)$   
}
```

pick clustering that gave the lowest cost $J(\dots)$

choose the numbers of clusters - k :

elbow method:



plot J - k graph and find elbow point where J decrease speed apparently slow down.

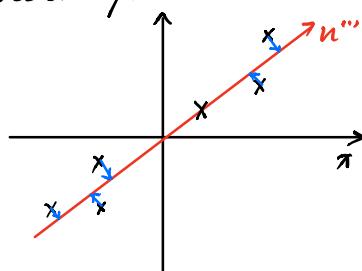
Principal Component Analysis (PCA)

reduce data from high dimension to low dimension.

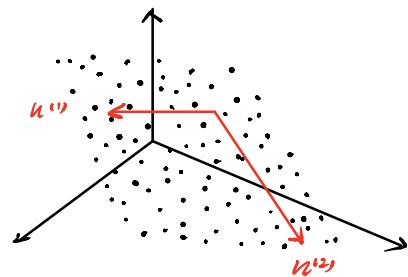
$$X = \begin{bmatrix} X^{(1)} \\ X^{(2)} \\ \vdots \\ X^{(m)} \end{bmatrix} \xrightarrow{\text{PCA}} Z = \begin{bmatrix} Z^{(1)} \\ Z^{(2)} \\ \vdots \\ Z^{(m)} \end{bmatrix}$$

$$X^{(i)} \in \mathbb{R}^n \quad Z^{(i)} \in \mathbb{R}^k$$

PCA problem formulation:



$$\begin{matrix} m=2 \\ k=1 \end{matrix} \quad 2D \rightarrow 1D$$

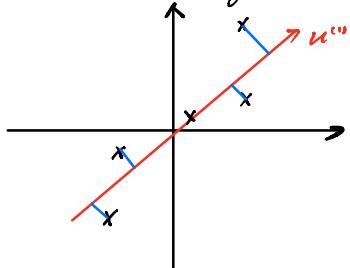


$$\begin{matrix} m=3 \\ k=2 \end{matrix} \quad 3D \rightarrow 2D$$

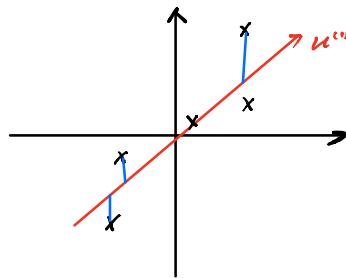
reduce from 2D to 1D: find a direction $u^{(1)} \in \mathbb{R}^n$ onto which to project the data so as to minimize the projection error.

reduce from nD to kD. find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data so as to minimize the projection error.

PCA is not linear regression:



PCA



linear regression

blue line : the error

PCA Algorithm

1. training set: $x^{(1)}, x^{(2)}, \dots, x^{(n)}$,
reduce data from n -dimension to k -dimension
 $(n > k)$

2. preprocessing (feature scaling/mean normalization):

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \text{ensure every feature has zero mean}$$

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j} \quad s_j \text{ can be } \max x_j - \min x_j$$

if different features on different scales, scale features to have comparable range of values

3. compute covariance matrix

$$\Sigma = \frac{1}{m} \sum_{i=1}^n \underbrace{x^{(i)}_{nx1}}_{nx1} \underbrace{(x^{(i)})^T}_{nxn} \quad \Sigma \in \mathbb{R}^{n \times n} \quad \Sigma = \text{sigma}$$

4. compute eigenvectors of matrix

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

$$U = \underbrace{\begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ | & | & | & \dots & | \end{bmatrix}}_{\text{select } n \times n \text{ columns}} \in \mathbb{R}^{n \times n} \quad u^{(i)} \in \mathbb{R}^n$$

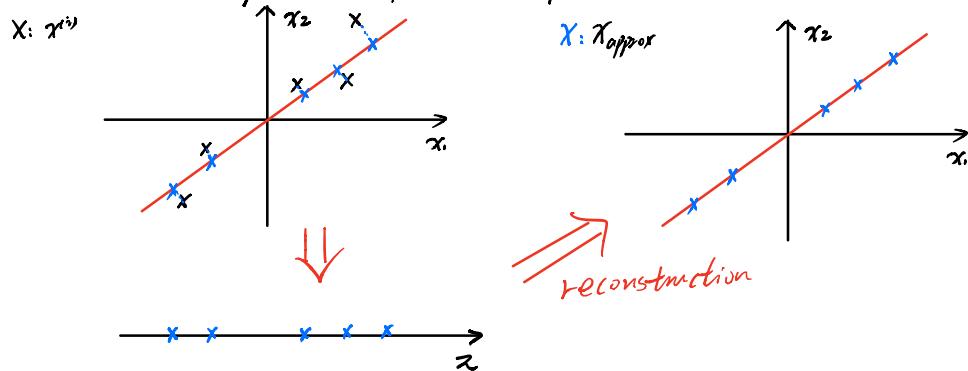
octave: $U_{\text{reduce}} = U(:, 1:k);$

$$U_{\text{reduce}} = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(k)} \\ | & | & | & \dots & | \end{bmatrix} \in \mathbb{R}^{n \times k}$$

octave: $Z = U_{\text{reduce}}^T * X;$

$$Z^{(i)} = (U_{\text{reduce}}^T)_{k \times n}^T x^{(i)}_{nx1} \in \mathbb{R}^k$$

Reconstruction from compressed representation:



$$z = (U_{\text{reduce}})^T x \quad \text{reconstruction} \quad x_{\text{approx}}^{(i)} = U_{\text{reduce}}^{\frac{n \times k}{k \times k}} z^{(i)}$$

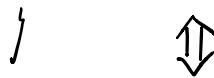
$$z \in \mathbb{R}^k \quad x_{\text{approx}}^{(i)} \in \mathbb{R}^n$$

choose the number of principle components — k

try PCA with $k=1, 2, 3, \dots$ }

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

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



Octave:

$$[U, S, V] = svdsigma$$

$$S = \begin{bmatrix} s_{11} & & & \\ & s_{22} & & \\ & & s_{33} & \dots \\ & & & s_{nn} \end{bmatrix}$$

$$1 - \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \geq 0.01$$

choose the largest k . (99% of variance retained)

advice for applying PCA

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

↓ extract

unlabeled dataset: $x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^n$

↓ PCA

$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^k$

new training set: $\{(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})\}$

Note: map $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA 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 set

application of PCA

- compression
 - reduce memory/disk need to store data
- speed up learning algorithm
- visualization ($k=2, 3$ for plotting graph)

Note: **DO NOT** use pca to reduce features for preventing overfitting, which is a bad idea. This won't work, choose to regularization instead.

Design Machine Learning System:

Note:
think if it's
necessary

- get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- run PCA to reduce features, get $\{(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})\}$
- training hypothesis model on $\{(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})\}$
- test on test set: map $x_{test}^{(i)}$ to $z_{test}^{(i)}$, run $h_0(z)$ on $\{(z_{test}^{(1)}, y_{test}^{(1)}), (z_{test}^{(2)}, y_{test}^{(2)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$