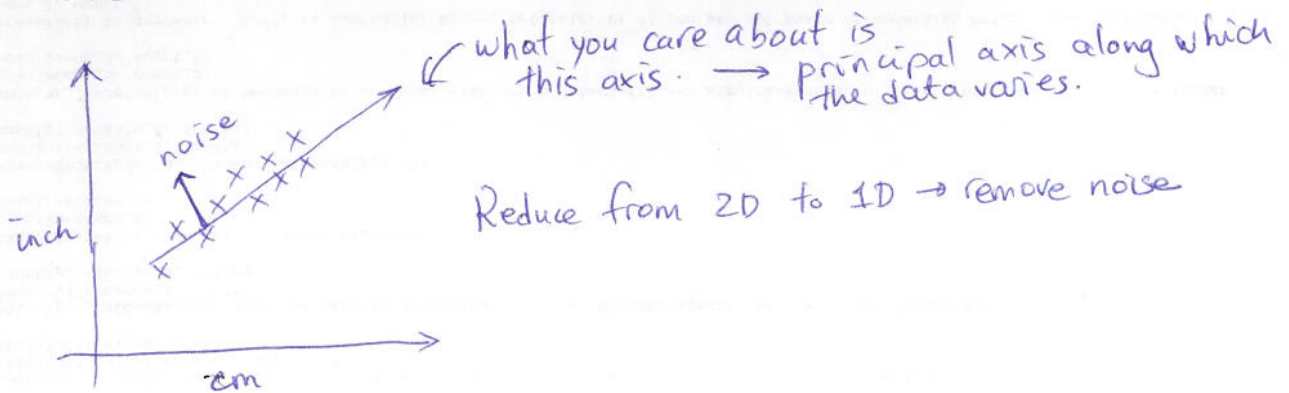


Principal Component Analysis

Given $\{x^{(1)} \dots x^{(m)}\}$ $x^{(i)} \in \mathbb{R}^n$

Goal: reduce to lower dim. data ($k \ll n$)



Pre-processing:

Set $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$

Replace $x^{(i)}$ with $x^{(i)} - \mu$

} zero-out the mean.

Set $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)})^2$

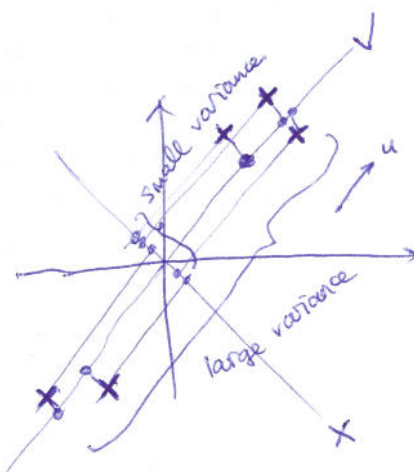
} variance of each feature

Replace $x_j^{(i)}$ with $\frac{x_j^{(i)}}{\sigma_j}$

} Normalize to unit variance

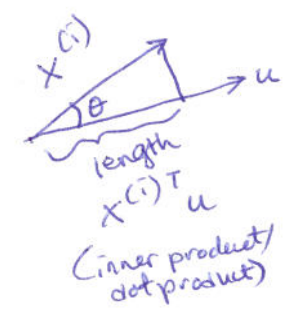
-> each of the feature has now equal variance.

} When features are in diff scales.



find direction u , st when I project the data to it; the projection varies widely.

If $\|u\|=1$, $x^{(i)}$ projected on u has length $x^{(i)T} u$



Choose u to maximize:

$$\text{Max}_{u: \|u\|=1} \frac{1}{m} \sum_{i=1}^m (x^{(i)T} u)^2$$

$$= \frac{1}{m} \sum_{i=1}^m (u^T x^{(i)}) (x^{(i)T} u) = u^T \left[\frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \right] u$$

⇒ u is the principal eigenvector of $\Sigma = \text{cov. matrix}$
 $= \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T}$

Eigenvector.

$Au = \lambda u$
↑ Eigenvector
eigenvalue
largest eigenvalue

derivation.

$\max u^T \Sigma u \text{ st } u^T u = 1$
 $L(u, \lambda) = u^T \Sigma u - \lambda (u^T u - 1)$ (Lagrange)

$\frac{\partial L(u, \lambda)}{\partial u} = \Sigma u - \lambda u \stackrel{\text{set}}{=} 0$
 $\Sigma u = \lambda u$

u = principal eigenvector of Σ .

$\Sigma = n \times n$
 $U = h \times l$
 $x^{(i)} = n \times l$

If we want k-dim subspace,
choose u_1, \dots, u_k to be k top eigenvector of Σ
↳ corresp. to k highest eigenvalues.

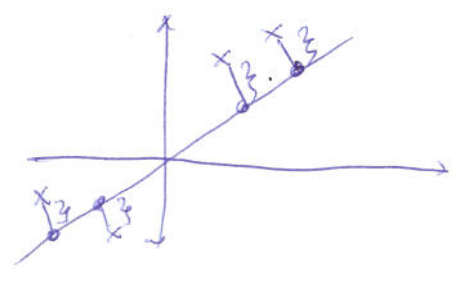
u_1, \dots, u_k : new basis of representing data
at beginning we have $x^{(1)}, \dots, x^{(m)} \in \mathbb{R}^n$

New representation :

$y^{(i)} = (u_1^T x^{(i)}, u_2^T x^{(i)}, \dots, u_k^T x^{(i)}) \cdot y^{(i)} \in \mathbb{R}^k$

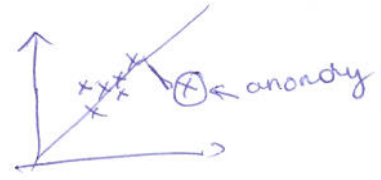
another view of PCA.

minimise sum of squared distance btwn the pt of its projection.



Use of PCA

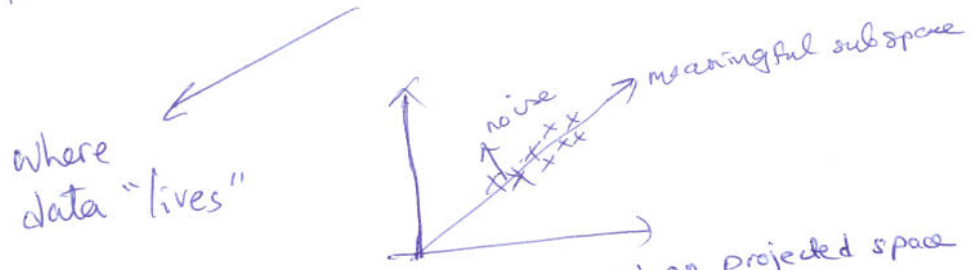
- visualization of high-dim data.
- compression. → ~~feature selection~~ work with lower dim data instead of high-dim data
- learning → more features → more complex → overfitting
use PCA to reduce dimensionality of features
- Anomaly detection → pts far away from your subspace



matching/

- Distance calculation.

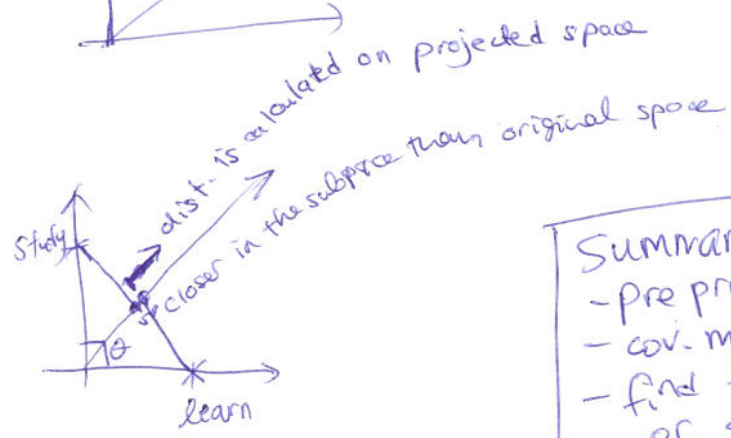
datapts often lie in a subspace low-dim. — the rest of the space maybe noise



In orig space

$$\text{Similarity} = \frac{x^{(i)T} x^{(j)}}{\|x^{(i)}\| \|x^{(j)}\|} = \cos \theta$$

= no words in common.



Summary PCA

- pre process Σ
- cov. matrix computation
- find top k eigenvectors of Σ

Latent Semantic Indexing LSI.

- does not have the preprocessing step

Problem of PCA

- Covariance matrix computation.
when you data

100 x 100 pixel image ; $x^{(i)} \in \mathbb{R}^{10,000}$
 Σ 10,000 x 10,000 = 100 mill. entries

Implement PCA using SVD — singular value decomposition

4

Say we have any $A \in \mathbb{R}^{m \times n}$

$$\text{decompose } A = U D V^T$$

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

$$D = \text{diagonal} = \begin{bmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & \\ & & \ddots & \\ 0 & & & \sigma_n \end{bmatrix}$$

$\sigma_i = \text{singular values of matrix } A$

$$\begin{bmatrix} A \\ m \times n \end{bmatrix} = \begin{bmatrix} U \\ m \times n \end{bmatrix} \begin{bmatrix} D \\ n \times n \end{bmatrix} \begin{bmatrix} V^T \\ n \times n \end{bmatrix}$$

(using svd command in matlab).
 svd = $O(n^3)$ not sure?

U's Column: eigenvector of AA^T

V's columns: eigenvector of $A^T A$

Can be used to compute eigenvector of PCA efficiently.

$$\Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T}$$

given Design matrix

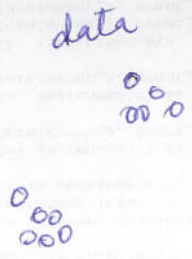
$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix}$$

$$\Sigma = X^T X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \begin{bmatrix} \text{---} x^{(1)} \text{---} \\ \text{---} x^{(2)} \text{---} \\ \vdots \\ \text{---} x^{(m)} \text{---} \end{bmatrix}$$

To get top k eigenvector of Σ ,

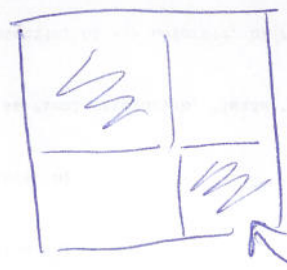
~~XXXX~~ $X = U D V^T$ top k columns of V are the top k eigenvectors of $X^T X = \Sigma$.

Spectral clustering



affinity matrix

$$W_{nm} = e^{-\frac{\|s_n - s_m\|}{\sigma^2}}$$



$$\Rightarrow A = \text{func}(W) \Rightarrow$$

disconnected graph



Smallest ~~top~~ k Eigenvectors of A

Clustering w/ k -means

map back to orig. data



$$d_i = \sum_j W_{ij} = \text{degree of a vertex}$$

$$D = \text{diag}(d_1, d_2, \dots, d_n) \text{ degree matrix}$$

$$|A| = \# \text{ of vertices in } A$$

$$\text{vol}(A) = \sum_{i \in A} d_i$$

$$\text{cut}(A, B) := \sum_{i \in A, j \in B} W_{ij}$$

Balanced min-cut

Ratio cut

Normalized cut

$$\min \text{cut}(A, B) \text{ st } |A| = |B|$$

$$\text{cut}(A, B) \left(\frac{1}{|A|} + \frac{1}{|B|} \right)$$

$$\text{cut}(A, B) \left(\frac{1}{\text{vol}(A)} + \frac{1}{\text{vol}(B)} \right)$$

NP-hard use spectral clustering to approx

$$\min \text{Cut}(A, B) = \min \frac{1}{2} f^T (D - W) f \quad S_i = \begin{cases} 1 & \text{if } x_i \in A \\ -1 & \text{if } x_i \in B \end{cases}$$

$$L = D - W$$

un-normalized Graph Laplacian

$$\min_{f \in \mathbb{R}^n} f^T L f \text{ st. } f^T \mathbf{1} = 0 \quad f^T f = n$$

$$= \min_{f \in \mathbb{R}^n} \frac{f^T L f}{f^T f} \text{ st. } f^T \mathbf{1} = 0 \Rightarrow \lambda = \text{smallest eigenvalue of } L$$

for connected graph; first eigenvector is constant (all 1s)
 disconnected graph; Laplacian is block diagonal
 & first k Laplacian eigenvectors are

⑥



Another paper has used $P = D^{-1}W$ ("normalized" affinity matrix)
 instead of L ; and take the eigenvector corresponding to the largest
 eigenvalue instead of the smallest. The two algorithms, the one using L
 matrix and the one using P matrix can be shown to be similar:

$$\underbrace{(D - W)}_L r = \mu D r$$

$$D^{-1}(D - W)r = \mu r$$

$$\underbrace{D^{-1}(W)}_P r = (1 - \mu)r$$

$$\lambda = 1 - \mu$$

$$v = r$$

equality between using
 L and P matrix: i.e.

$$L v = \lambda v \rightarrow \text{take smallest eigenval}$$

$$P v = (1 - \mu)v \rightarrow \text{take largest eigenval}$$