

PRINCIPAL COMPONENT ANALYSIS

Dimensionality Reduction

In dimensionality reduction problems we observe

$$x_1, \dots, x_n \in \mathbb{R}^d$$

The goal is to transform these inputs to new variables

$$x_i \mapsto \theta_i \in \mathbb{R}^k$$

where $k < d$, in such a way that information loss is minimized.

Methods for dimensionality reduction may be categorized according to the following issues:

- 1) How is "information loss" quantified?
- 2) Supervised or unsupervised? If labels y_1, \dots, y_n are available, are they used?

3) Is the map $x \mapsto \theta$ linear or nonlinear?

4) Feature selection

(A)

$$\theta = \begin{bmatrix} \\ \end{bmatrix}$$

versus feature extraction

$$\theta = \begin{bmatrix} \\ \\ \\ \end{bmatrix}$$

PCA is described by

1)

2)

3)

4)

Linear Spans and Projections

Let $a_1, \dots, a_k \in \mathbb{R}^d$ be linearly independent column vectors.

Definition The linear span of a_1, \dots, a_k is the set

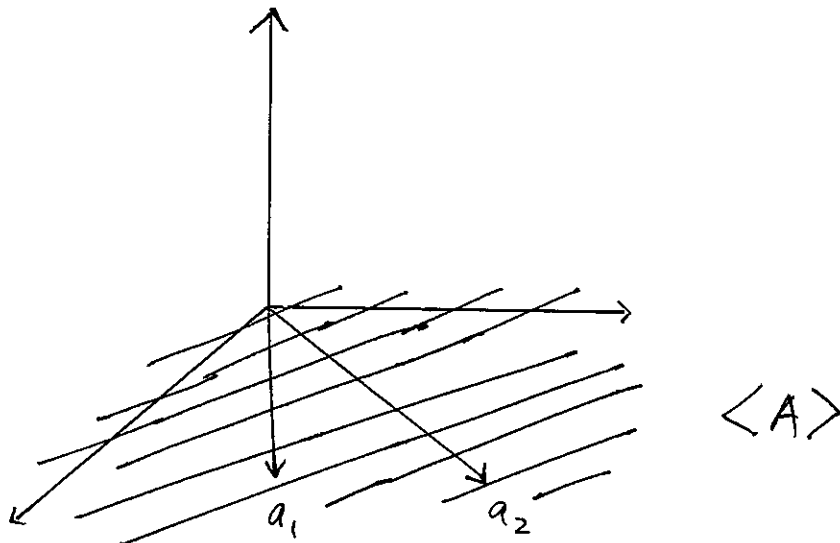
$$\left\{ x = \sum_{j=1}^k \theta^{(j)} a_j : \theta = [\theta^{(1)} \dots \theta^{(k)}]^T \in \mathbb{R}^k \right\}$$

Equivalently, if we set

$$A = \begin{bmatrix} a_1 & \dots & a_k \end{bmatrix} \in \mathbb{R}^{d \times k}$$

then the linear span is the image or column span of A , and will be denoted $\langle A \rangle$.

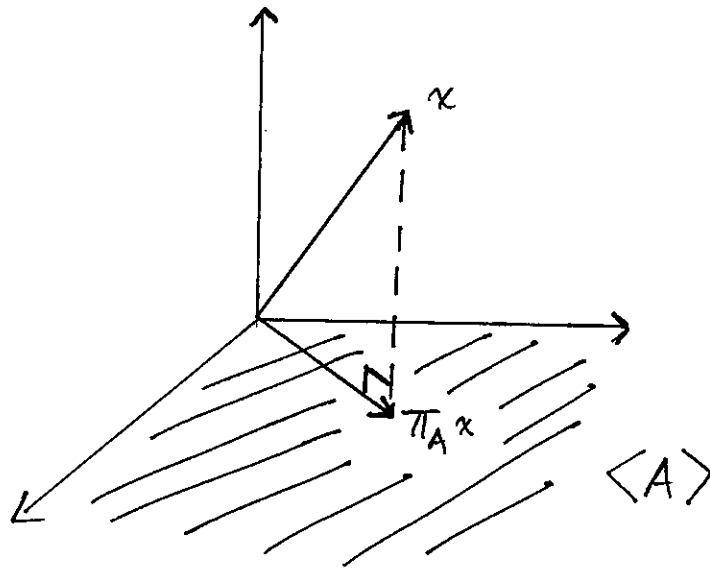
$$\begin{aligned} d &= 3 \\ k &= 2 \end{aligned}$$



Definition | The mapping $\mathbb{R}^d \rightarrow \langle A \rangle$ given by

$$\pi_A : x \mapsto \text{closest point to } x \text{ in } \langle A \rangle$$

is called the projection onto $\langle A \rangle$.



Every point in $\langle A \rangle$ is equal to $A\theta$ for some

$\theta \in \mathbb{R}^k$. Therefore $\pi_A x = A\hat{\theta}$ where

(B) $\hat{\theta} =$

Therefore

$$\pi_A x = A (A^T A)^{-1} A^T x$$

If a_1, \dots, a_k are orthonormal, then

(c) $\pi_A x =$

By the orthogonality principle, $x - \pi_A x$ is orthogonal to every vector in $\langle A \rangle$.

Exercise 1 Verify this result.

Solution | Let $A\theta$ be an arbitrary vector in $\langle A \rangle$.

We wish to show $\langle A\theta, x - \pi_A x \rangle = 0$. Now

$$\langle A\theta, x - \pi_A x \rangle = (A\theta)^T (x - A(A^T A)^{-1} A^T x)$$

(D)

$$= \theta^T A^T (x - A(A^T A)^{-1} A^T x)$$

$$= \theta^T A^T x - \theta^T A^T A (A^T A)^{-1} A^T x$$

$$=$$

Projection matrices are also idempotent, meaning

$$\pi_A^2 = \pi_A \cdot \pi_A$$

(E)

$$=$$

Interpretation :

Principal Component Analysis

The idea behind PCA is to approximate

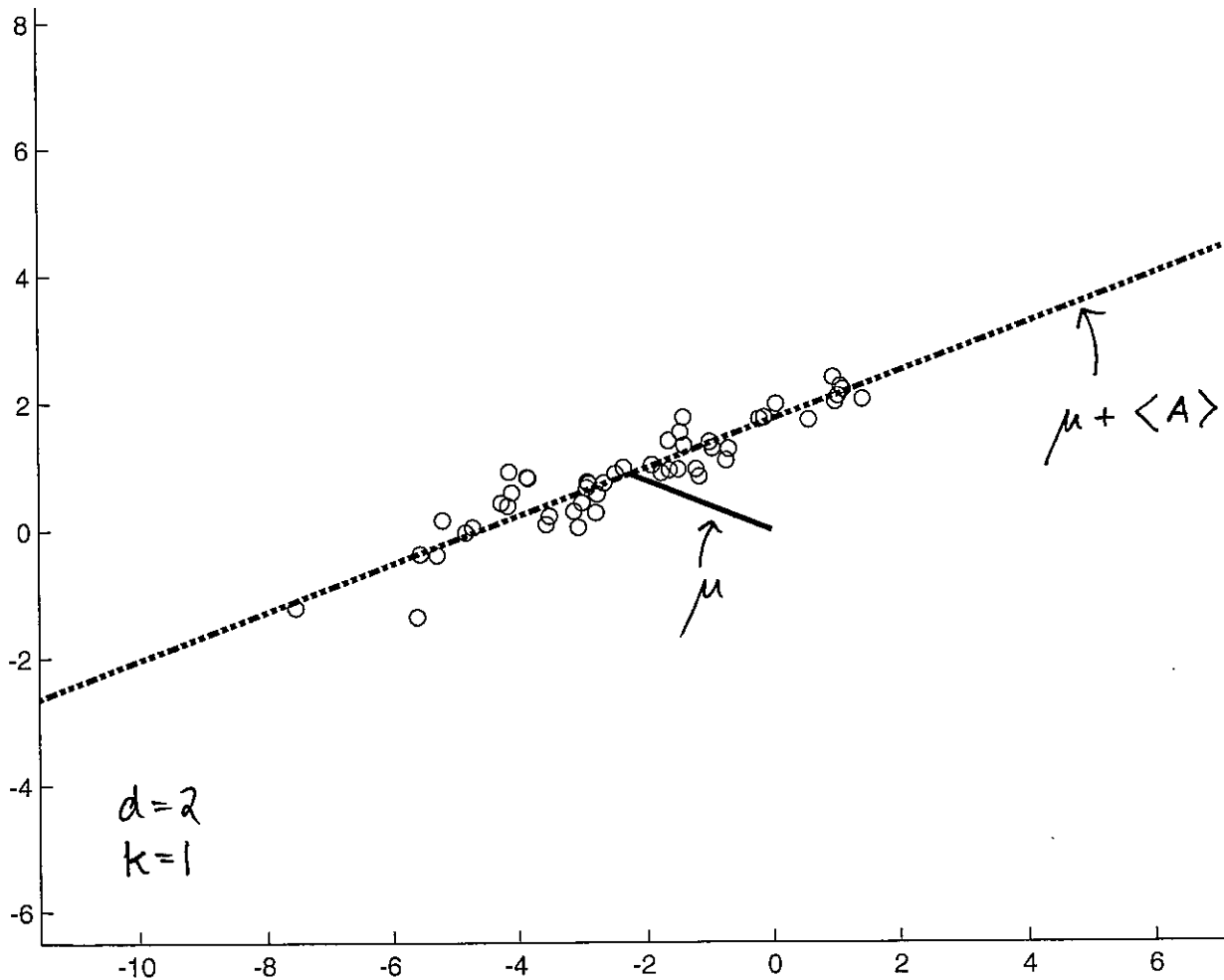
$$x_i \approx \mu + A \theta_i$$

where

$$\mu \in \mathbb{R}^d$$

$$A \in \mathbb{R}^{d \times k}, \text{ orthonormal columns}$$

$$\theta_i \in \mathbb{R}^k$$



Mathematically, we define μ , A , and $\theta_1, \dots, \theta_n$ to be the solution of

$$\min_{\mu, A, \{\theta_i\}} \sum_{i=1}^n \|x_i - \mu - A\theta_i\|^2$$

The solution is given in terms of the spectral representation of

(F)

$$S =$$

In particular, write

$$S = U\Lambda U^T$$

where

$$U = \begin{bmatrix} u_1 & \dots & u_d \end{bmatrix}, \quad \Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \dots & \\ & & & \lambda_d \end{bmatrix}$$

with $U^T U = U U^T = I_{d \times d}$ and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$

Note: $S u_j = \lambda_j u_j \quad \forall j \dots$

We will show that the optimal rank k linear approximation of the data is given by

⑥ $\mu =$

$A =$

$\theta_i =$

Is A unique?

Terminology

- principal component transformation

$x \mapsto \theta =$

- j th principal component

$\theta^{(j)} =$

- j th principal eigenvector \rightarrow

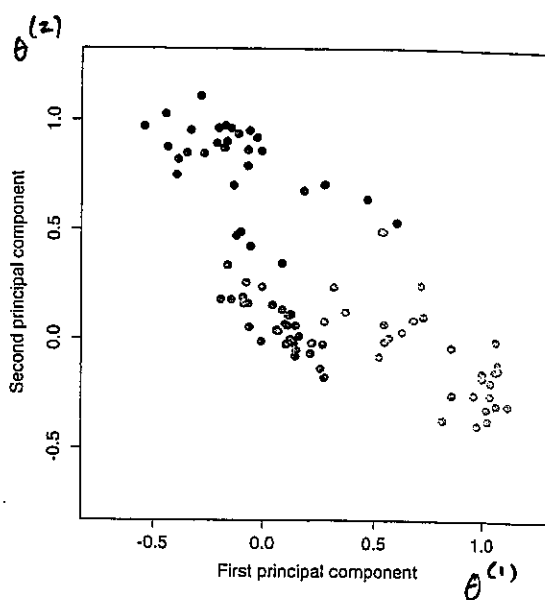
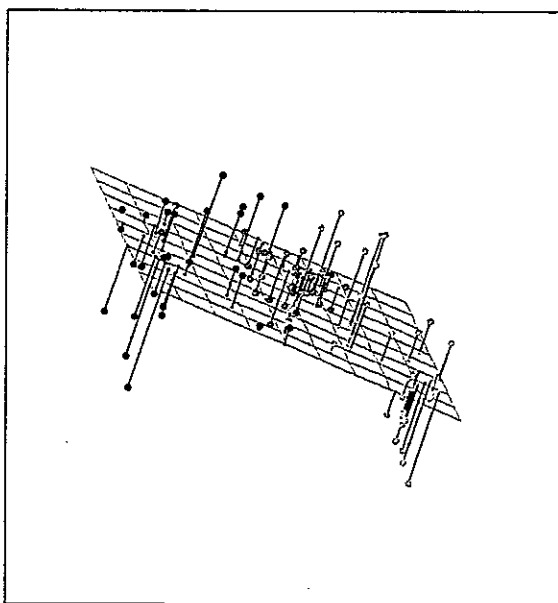


FIGURE 14.21. The best rank-two linear approximation to the half-sphere data. The right panel shows the projected points with coordinates given by ~~the~~, the first two principal components of the data.

Hastie, Tibshirani, & Friedman

Proof: | We wish to minimize

$$\sum_{i=1}^n \|x_i - \mu - A \theta_i\|^2$$

w.r.t $\mu, A, \{\theta_i\}$

Step 1: Eliminate $\{\theta_i\}$

Suppose A, μ are fixed. Then to minimize the expression we take

$$\textcircled{H} \quad \theta_i =$$
$$=$$

Step 2: Eliminate μ

Holding A fixed, we wish to minimize

$$\sum_{i=1}^n \|x_i - \mu - AA^T(x_i - \mu)\|^2$$

$$= \sum_i \| (I - AA^T)(x_i - \mu) \|^2$$

$$= \sum (x_i - \mu)^T \underbrace{(I - AA^T)^T (I - AA^T)}_B (x_i - \mu)$$

$$\begin{aligned}\frac{\partial}{\partial \mu} &= -2 \sum_i B (x_i - \mu) \\ &= -2B \cdot \sum_{i=1}^n (x_i - \mu) = 0\end{aligned}$$

This is satisfied when

$$\textcircled{I} \quad \mu =$$

Step 3 Optimize A

It remains to minimize

$$\sum_{i=1}^n \|x_i - \bar{x} - A A^T (x_i - \bar{x})\|^2$$

w.r.t. A

For convenience, assume $\bar{x} = 0$. (otherwise we could

substitute $\tilde{x}_i = x_i - \bar{x}$)

Observe

$$\begin{aligned}\sum_i \|x_i - AA^T x_i\|^2 &= \sum_i (x_i - AA^T x_i)^T (x_i - AA^T x_i) \\ &= \sum_i x_i^T x_i - 2x_i^T AA^T x_i + x_i^T AA^T AA^T x_i\end{aligned}$$

(J)

Therefore we can focus on maximizing

Introduce

$$b_j := u^T a_j$$

$$y_i := u^T x_i$$

} change of basis

Our goal is to show

$$(K) \quad \langle A \rangle = \langle u_1, \dots, u_k \rangle \iff b_j =$$



Now

$$\sum_i x_i^T A_k A_k^T x_i = \sum_i \|A^T x_i\|^2$$

$$= \sum_{i=1}^n \sum_{j=1}^k (a_j^T x_i)^2$$

$$= \sum_{i=1}^n \sum_{j=1}^k (a_j^T U U^T x_i)^2$$

$$= \sum_{i=1}^n \sum_{j=1}^k (b_j^T y_i)^2$$

$$= \sum_{i=1}^n \sum_{j=1}^k (b_j^T y_i) (y_i^T b_j)$$

$$= \sum_{j=1}^k b_j^T \left[\sum_{i=1}^n y_i y_i^T \right] b_j$$

Exercise

Simplify



Solution

$$\begin{aligned}\sum_{i=1}^n y_i y_i^T &= \sum_{i=1}^n u^T x_i x_i^T u \\ &= u^T \left[\sum_i x_i x_i^T \right] u \\ &= u^T [nS] u \\ &= n \cdot u^T \cdot u \wedge u^T \cdot u \\ &= n \wedge\end{aligned}$$

■

Thus, we need to maximize

$$\begin{aligned}\sum_{j=1}^k b_j^T \wedge b_j &= \sum_{j=1}^k \sum_{l=1}^d (b_j^{(l)})^2 \lambda_l \\ &= \sum_{l=1}^d \left[\sum_{j=1}^k (b_j^{(l)})^2 \right] \lambda_l \\ &= \sum_{l=1}^d h_l \lambda_l \quad \text{where } h_l = \sum_{j=1}^k (b_j^{(l)})^2\end{aligned}$$

Lemma: $0 \leq h_l \leq 1$ for $l=1, \dots, d$

and $\sum_{l=1}^d h_l = k$ (you'll show this on the homework)

Since $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$, we maximize

$$\sum_{l=1}^d h_l \lambda_l$$

by letting

$$h_l =$$

(L)

which implies

$$b_j^{(e)} = 0 \quad \text{for } l > k$$

Therefore

$$a_j = \sum b_j \in \langle u_1, \dots, u_k \rangle$$



Maximum Variance Projections

PCA can be derived from a second perspective.

Suppose the data has zero mean:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = 0.$$

What is the normal vector a_{\perp} ($\|a_{\perp}\| = 1$)
for which

$$\theta^{(1)} =$$

has maximal variance?

The variance of $\theta^{(1)}$ is defined to be

$$\text{Var}(\theta^{(1)}) =$$

Theorem Assume $\lambda_1 > \lambda_2$.

If

$$\theta^{(1)} = a_{\perp}^T x$$

has maximal variance among all a_{\perp} with $\|a_{\perp}\| = 1$,
then

$$a_{\perp} = \text{1st principal eigenvector} = u_1$$

$$\theta^{(1)} = \text{1st principal component}$$

The variance of $u_i^T x$ is

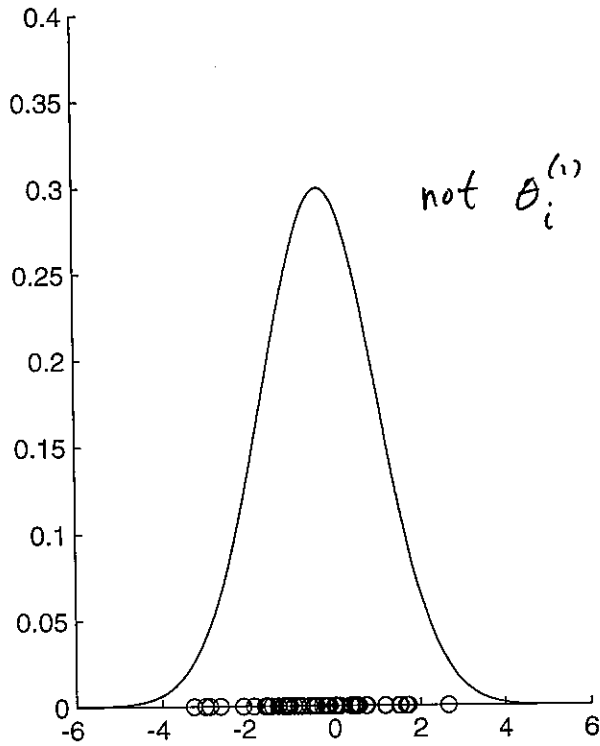
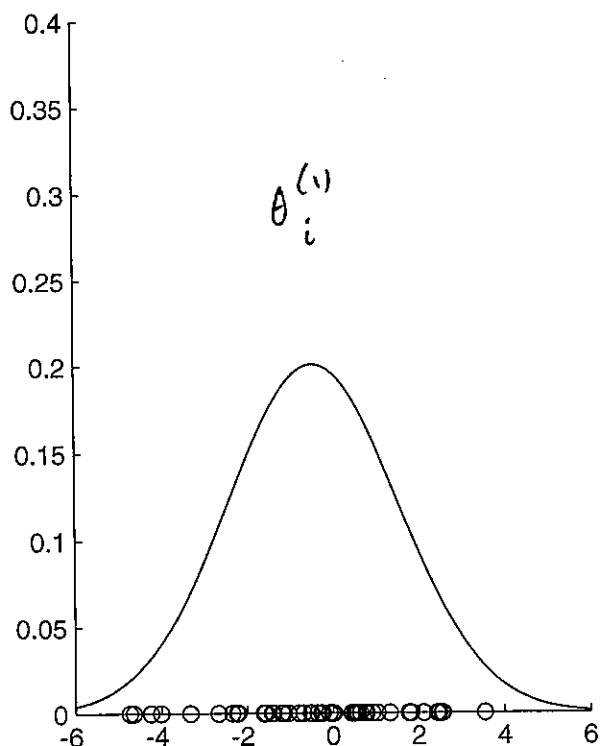
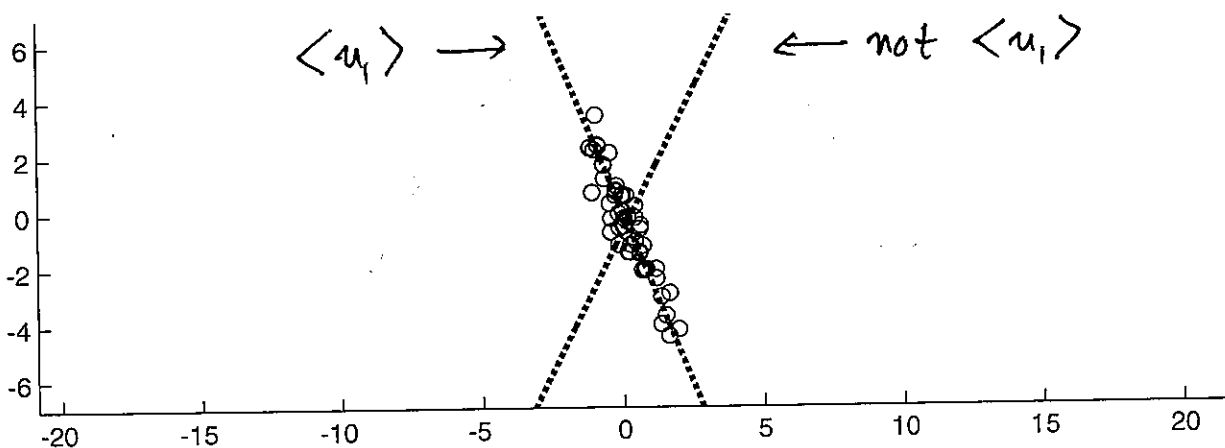
(N)

$$\text{Var}(\theta^{(1)}) =$$

$$=$$

$$=$$

$$=$$



The other principal eigenvectors/components satisfy a similar property:

Theorem 1 Assume $\lambda_k > \lambda_{k+1}$.

If

$$\theta^{(k)} = a_k^T X$$

has maximal variance among a_k such that

- $\|a_k\| = 1$
- $a_k \perp u_1, \dots, u_{k-1}$

then

$a_k = k^{\text{th}}$ principal eigenvector = u_k

$\theta^{(k)} = k^{\text{th}}$ principal component

The variance of $\theta^{(k)}$ is _____.

Geometric Interpretation

If the density of a random variable X is

$$\phi(x; \mu, \Sigma) = (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}, x \in \mathbb{R}^d$$

then we say X is a multivariate Gaussian/normal random variable with parameters $\mu \in \mathbb{R}^d$, $\Sigma \in \mathbb{R}^{d \times d}$, $\Sigma > 0$, and write $X \sim N(\mu, \Sigma)$.

If μ, Σ are unknown, we can estimate them from data $x_1, \dots, x_n \stackrel{iid}{\sim} N(\mu, \Sigma)$.

In particular, the Maximum Likelihood Estimate (MLE) of μ, Σ is obtained by maximizing the likelihood

$$l(\mu, \Sigma) := \prod_{i=1}^n \phi(x_i; \mu, \Sigma)$$

The maximum likelihood estimate of a Gaussian density is given by

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\Sigma} = S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

A contour of $\phi(x; \bar{x}, S)$ consists of all points x such that

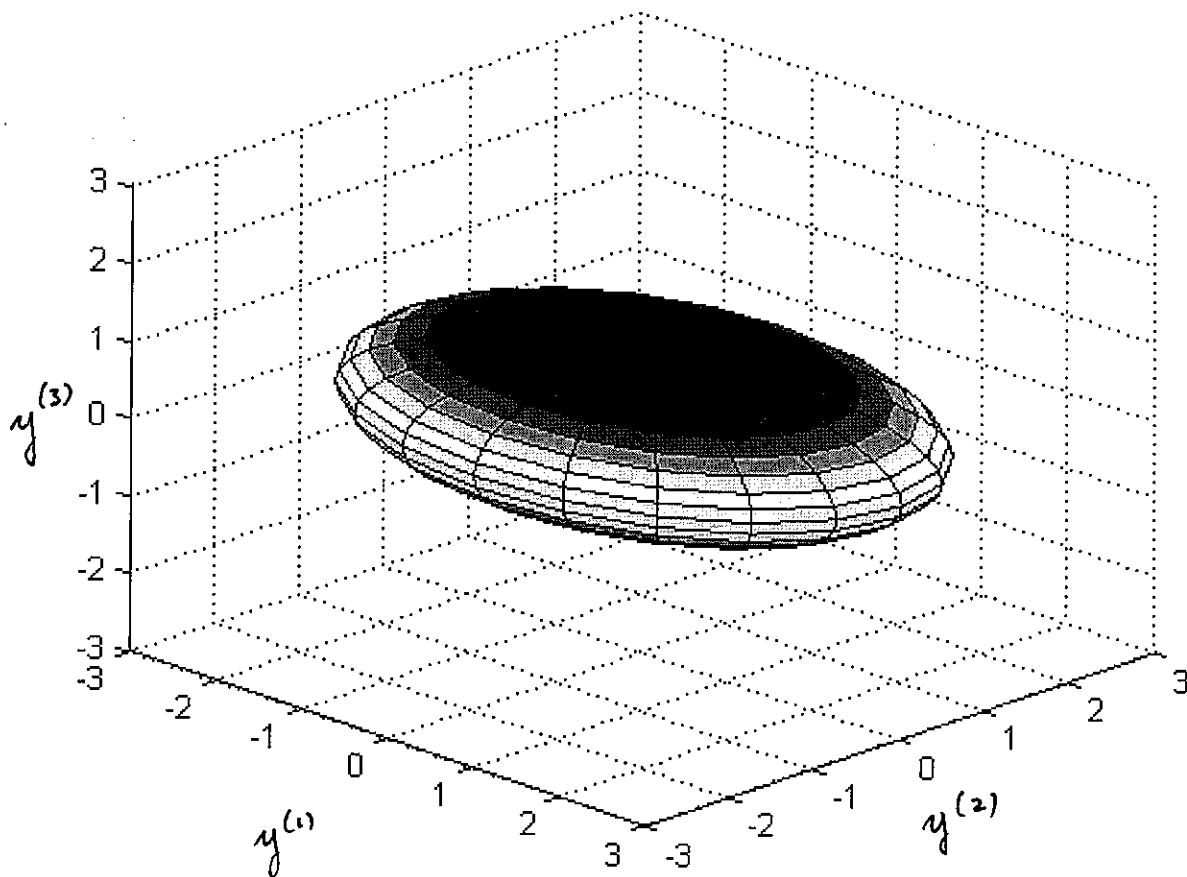
$$= c$$

Let $\bar{y} = U^T x$, $\bar{y} = U^T \bar{x}$. We may write the contour-defining equation as

$$c =$$

In 3-d, the contour is an ellipsoid:

$$\frac{(y^{(1)} - \bar{y}^{(1)})^2}{\lambda_1} + \frac{(y^{(2)} - \bar{y}^{(2)})^2}{\lambda_2} + \frac{(y^{(3)} - \bar{y}^{(3)})^2}{\lambda_3} = c$$



$$\lambda_1 > \lambda_2 > \lambda_3$$

In general, PCA fits an elliptical distribution to the data, and projects onto the first k major axes of the resulting ellipsoid.

Practical Matters

Preprocessing

It is customary to center and scale a data set so that it has 0 mean and unit variance along each feature

For $j = 1$ to d

$$1) \quad \bar{x}^{(j)} = \frac{1}{n} \sum_{i=1}^n x_i^{(j)}$$

$$x_i^{(j)} \leftarrow x_i^{(j)} - \bar{x}^{(j)}, \quad i = 1, \dots, n$$

$$2) \quad \sigma^{(j)} = \left(\frac{1}{n} \sum_i (x_i^{(j)})^2 \right)^{\frac{1}{2}}$$

$$x_i^{(j)} \leftarrow x_i^{(j)} / \sigma^{(j)}, \quad i = 1, \dots, n$$

End

This puts all features on an "equal playing field."

The steps may be omitted when

1) The data are known to be zero mean

2) The data are known to have comparable units of measurement.

Selecting k

It can be shown that

$$\min_{\mu, A_k, \{\theta_i\}} \sum_{i=1}^n \|\chi_i - \mu - A_k \theta_i\|^2 = n (\lambda_{k+1} + \dots + \lambda_d)$$

When $k=d$ this becomes

$$\min_{\mu} \sum_{i=1}^n \|\chi_i - \mu\|^2 = n (\lambda_1 + \dots + \lambda_d),$$

which we call the total variation of the data.

Therefore, the % of total variation captured by k -dimensional PCA is

① $\%TV =$

One heuristic for choosing k is to take the smallest k such that

$$\%TV \geq .95$$

↑
or .99, or some other number.

When to use PCA

- When the data form a single "point cloud" in space
- When the data are approximately Gaussian, or some other "elliptical" distribution.
- When low rank linear subspaces capture a majority of the variation

When not to use PCA

Key

A. Feature selection

Feature extraction

$$\theta = \begin{bmatrix} x^{(2)} \\ x^{(5)} \\ x^{(11)} \\ \vdots \end{bmatrix}$$

$$\theta = \begin{bmatrix} \phi^{(1)}(x) \\ \phi^{(2)}(x) \\ \vdots \end{bmatrix}$$

e.g., $\phi^{(1)}(x) = (x^{(1)})^2 x^{(2)} + e^{-x^{(5)}}$

PCA: sum of squared errors, unsupervised, linear, feature extraction

B. $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \|x - A\theta\|^2 = (A^T A)^{-1} A^T x$

C. $\pi_A x = A \cdot A^T x$

D.
$$\begin{aligned} &= \theta^T A^T (x - A(A^T A)^{-1} A^T x) \\ &= \theta^T (A^T x - A^T A (A^T A)^{-1} A^T x) \\ &= \theta^T (A^T x - A^T x) = 0 \end{aligned}$$

E.
$$\begin{aligned} \pi_A^2 &= \pi_A \cdot \pi_A = A^T (A^T A)^{-1} A^T A (A^T A)^{-1} A^T \\ &= A (A^T A)^{-1} A^T = \pi_A \end{aligned}$$

\Rightarrow second projection has no effect

F. $S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$, the sample covariance matrix

G. $\mu = \bar{x}$, $A = [u_1 \dots u_k]$, $\theta_i = A^T (x_i - \bar{x})$

$x \mapsto \theta = A^T (x - \bar{x})$, $\theta^{(j)} = u_j^T (x - \bar{x})$, u_j

$$H. (A^T A)^{-1} A^T (x_i - \mu) = A^T (x_i - \mu) \quad \text{I. } \mu = \bar{x}$$

$$J. \sum_i x_i^T x_i - x_i^T A A^T x_i, \quad \sum_i x_i^T A A^T x_i$$

$$K. b_j = \left[\begin{array}{c|c} \text{whatever} & 0 \dots 0 \\ \hline & k+1 \leftrightarrow d \end{array} \right]^T$$

$$L. h_l = \begin{cases} 1 & \text{if } l \leq k \\ 0 & \text{if } l > k \end{cases}$$

$$M. \theta^{(1)} = a_1^T x, \quad \text{Var}(\theta^{(1)}) := \frac{1}{n} \sum_{i=1}^n (a_1^T x_i)^2$$

$$N. \text{Var}(\theta^{(1)}) = \frac{1}{n} \sum (u_1^T x_i)^2 = \frac{1}{n} \sum (u_1^T x_i)(x_i^T u_1)$$

$$= u_1^T \cdot \frac{1}{n} \sum x_i x_i^T \cdot u_1$$

$$= u_1^T \cdot U \Lambda U^T \cdot u_1$$

$$= [1 \ 0 \ \dots \ 0] \cdot \Lambda \cdot \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \lambda_1$$

$$O. c = (x - \bar{x})^T S^{-1} (x - \bar{x}) = (u_y - u_{\bar{y}})^T S^{-1} (u_y - u_{\bar{y}}) \\ = (y - \bar{y})^T u^T \cdot u \Lambda^{-1} u^T \cdot u (y - \bar{y}) = (y - \bar{y})^T \Lambda^{-1} (y - \bar{y}) \\ = \sum_{j=1}^d \frac{(y^{(j)} - \bar{y}^{(j)})^2}{\lambda_j}$$

$$P. \%TV = \frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_d}$$