# Dimensionality Reduction

Sunita Sarawagi

CS 725 Fall 2023

# Dimensional Reduction
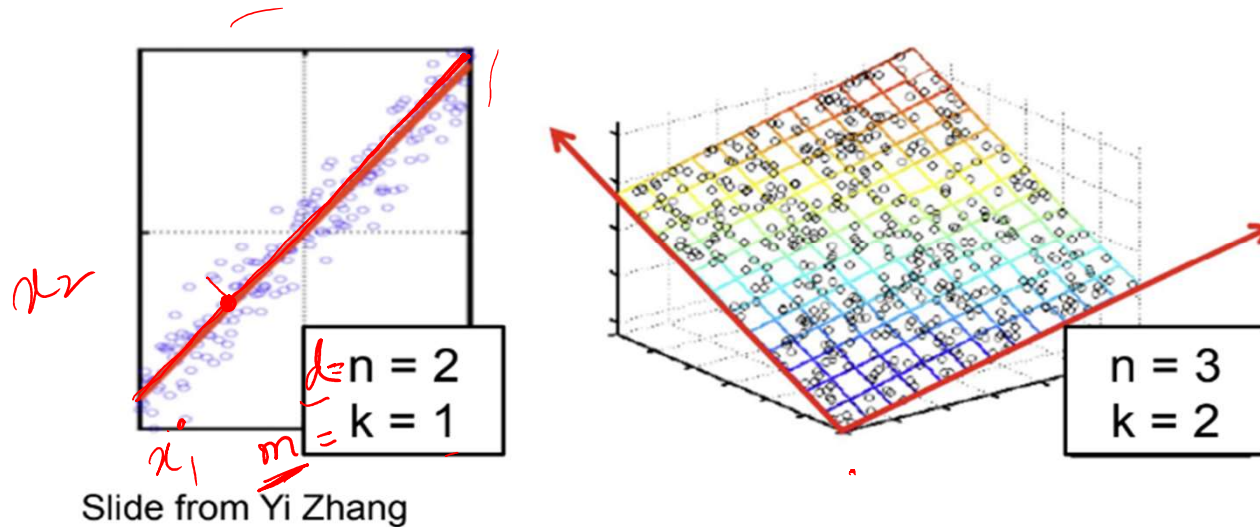
*t-sne*

- A form of unsupervised learning where we learn a mapping from a *high* dimensional space $x \in R^d$ to a lower dimensional space $z \in R^k$, *m* *d is large* where k << d

- Often original data in raw form might have a very high dimensionality (e.g. an image), but the information in each instance can be expressed compactly in smaller dimensions.

- Applications
  - Visualization *continues to be relevant. Eg: visualizing high-dimensional embedding t-sne*
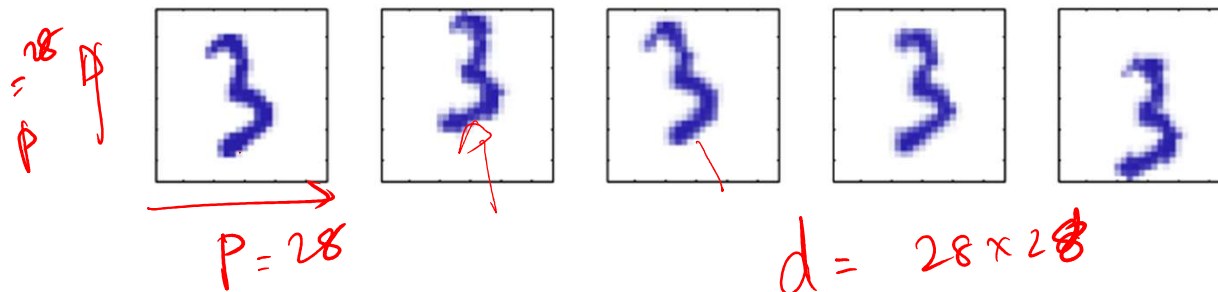  - Efficient learning
  - Data understanding.

# Dimension reduction

- Assumption: data (approximately) lies on a lower dimensional space
- Examples:



$d=n = 2$

$m= k = 1$

$n = 3$

$k = 2$

Slide from Yi Zhang

# Example (from Bishop)

- Suppose we have a dataset of digits ("3") perturbed in various ways:



$P = 28$

$d = 28 \times 28$

- What operations did I perform? What is the data's intrinsic dimensionality?

- Here the underlying manifold is *nonlinear*

Random displacement and rotation

# Linear projections

Given: N data points $x^1, x^2 \, \text{---} \, \text{---} \, x^N$

- Each original high-dimensional $x \in R^d$ is projected to a lower dimensional space $z \in R^m$ using just a linear projection matrix.

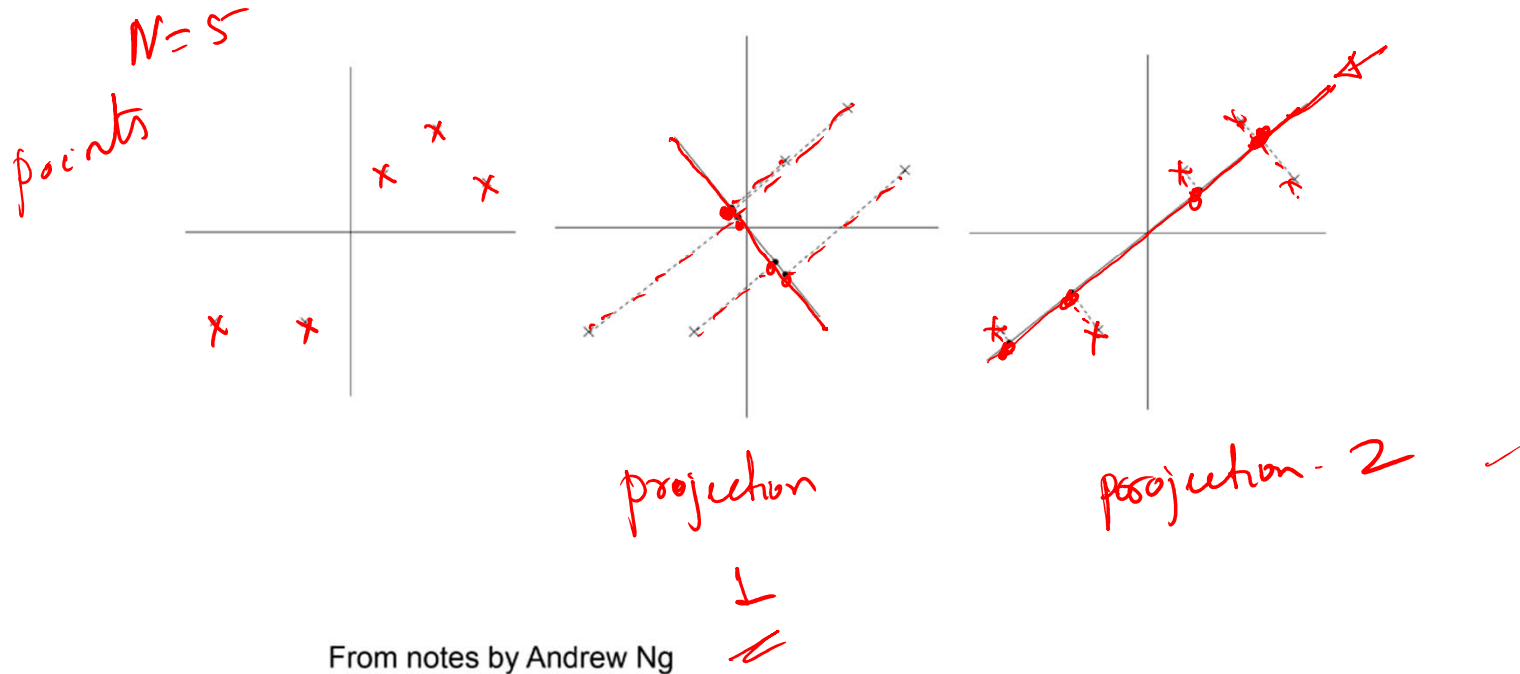$$z = Wx \qquad W = \begin{bmatrix} W_1^T \\ W_2^T \\ \vdots \\ W_m^T \end{bmatrix} \qquad W_j = \begin{bmatrix} w_{j1} \\ \vdots \\ w_{jd} \end{bmatrix} \qquad x = \begin{bmatrix} u_1 \\ \vdots \\ x_d \end{bmatrix} \qquad z = \begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix}$$

$m \times d$

$$z_j = W_j^T x$$

$$= \langle w_j, x \rangle$$

$$W = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{j1} & w_{j2} & \cdots & w_{jd} \\ \vdots & & & \\ w_{m1} & w_{m2} & \cdots & w_{md} \end{bmatrix} \begin{array}{l} \to W_1^T \\ \\ -W_j^T \\ \\ \to W_m^T \end{array}$$
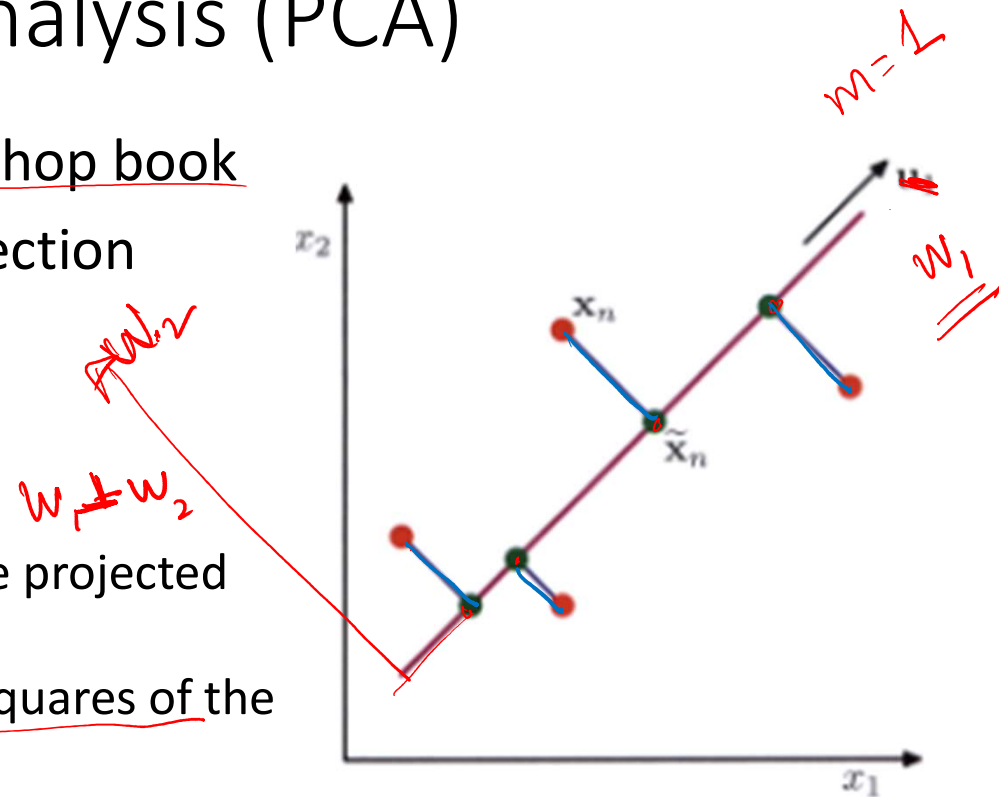
No labelled supervision on the desired z. How do we choose the best z?



From notes by Andrew Ng

# Principal Component Analysis (PCA)

- Reading material: Chapter 12 of Bishop book
- Simple and widely used linear projection
- Projection basis are orthogonal
  - $w_j \perp w_r \quad j, r \leq m$
- Objectives
  - Maximize the variance (spread) of the projected points --- green points
  - Equivalently, minimizing the sum of squares of the projection error --- blue lines.

# Consider single dimensional projection m=1, to ~~minimize~~ *maximize* variance of projected points

Assume $\|W_1\| = 1$

Given $x^1, x^2 \cdots x^N$    $x^i \in R^d$

$z^i = W_1^T x^i$

find $W_1$ s.t

     Variance of $\{z^1, z^2 \cdots z^N\}$ is ~~minimized~~ *maximized*

where $z^i = W_1^T x$

$$W_1 = \arg\max_{W} \sum_{i=1}^{N} (z^i - \bar{z})^2 \quad \text{where}$$

$$z^i = W_1^T x^i$$

$$\bar{z} = \frac{1}{N} \sum_{i=1}^{N} z^i \quad (\text{Average } \bar{z})$$

Two

$N = 5$



$x^3$

$x^4$ $(2, 2)$

$(0.5, 1)$ $x^5$

$x^1$ $x(3, 1)$

$(-2, -2)$ $x$

$(-1, -2)$

$x^1$ $x^2$

$w^1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

$0 = w_b^T x^1$

$\begin{array}{cc} z^1 & 0 \\ z^2 & 1 \\ z^3 & -0.5 \\ z^4 & 0 \\ z^5 & 2 \end{array}$

Variance $\{0, 1, -0.5, 0, 2\}$

$\bar{z} = 0.5$

=

$w_1^1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$z^1 = -4$

$z^2 = -3$

$z^3 = 1.5$

$z^4 = 4$

$z^5 = 4$

$\bar{z} = 0.5$

# Solving for the optimal projection

$$w^1 = \text{argmax} \sum_{i=1}^{N} \left( w_1^T x^i - w_1^T \bar{x} \right)^2$$

It can be shown that

$$\sum_{i=1}^{N} \left( w_1^T x^i - w_1^T \bar{x} \right)^2 \equiv w_1^T S w_1 \quad \text{where} \quad S = \sum_{i=1}^{N} \left( x_i - \bar{x} \right) \left( x_i - \bar{x} \right)^T$$

$$\max_{w} \boxed{w_1^T S w_1} \equiv \text{Eigen value of } w_1$$

$w_1 \equiv \text{Eigen vector}$

where $\|w_1\| = 1$

or $\boxed{w_1^T w_1 = 1}$

covariance ~~correlation~~ between two dimensions of $x$

$S$
$d$ $\begin{bmatrix} S_{jr} \end{bmatrix}$

$d-$

covariance between dimension $j$ & $r$ of N points

$$\max_{w_1} \quad w_1^T S w_1 - \lambda \underbrace{(w_1^T w_1 - 1)}_{} \quad \text{— Lagrangian multiplier.}$$

Equating gradient w.r.t $w_1$ to $0$

$$\boxed{S w_1 = \lambda w_1}$$

$\lambda$ is a scalar $\in R$. S is a square matrix

$\Rightarrow$ $w_1 =$ is an Eigen vector of S

Left multiply by $w_1$

$$w_1^T S w_1 = w_1^T \lambda w_1 = \lambda \boxed{w_1^T w_1} = \underline{\underline{\lambda}} \quad [\because w_1^T w = 1].$$

objective:

$w_1 = $ Eigen-vector corresponding to which Eigen value is maximum.

# In general..

- The best m dimensional linear projection of a d-dimensional dataset for maximizing variance of the projected points are
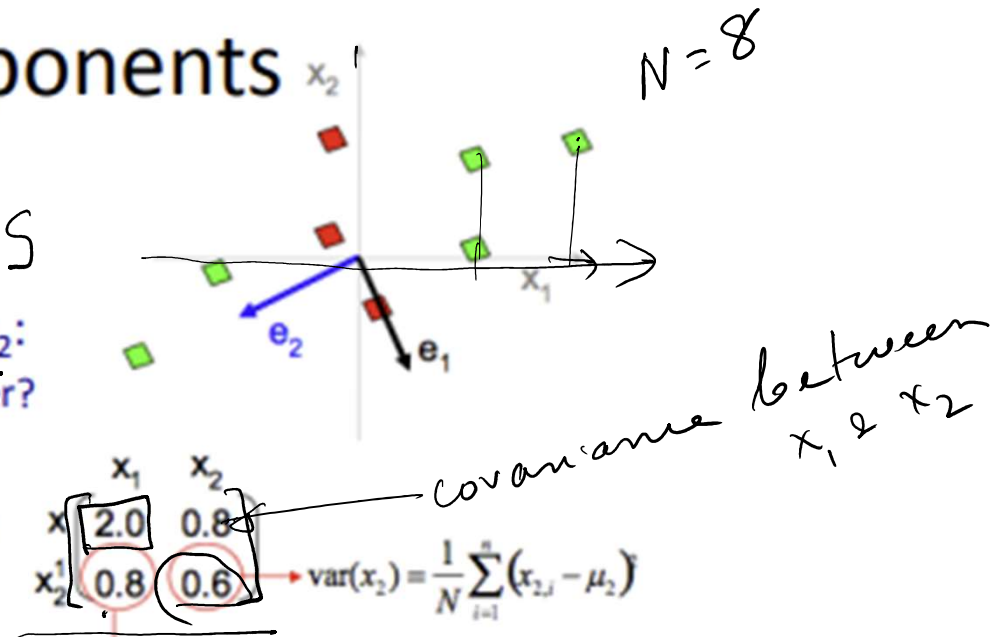    - The first m eigen vectors of the covariance matrix S of the given data points.

# Principal components

$x_2$

$N = 8$

- Compute covariance matrix $\Sigma = S$
  - covariance of dimensions $x_1$ and $x_2$:
    - do $x_1$ and $x_2$ tend to increase together?
    - or does $x_2$ decrease as $x_1$ increases?
  - covariance: measure of variability

|       | $x_1$ | $x_2$ |
|-------|-------|-------|
| $x_1$ | 2.0   | 0.8   |
| $x_2$ | 0.8   | 0.6   |

covariance between $x_1$ & $x_2$

$$\text{var}(x_2) = \frac{1}{N}\sum_{i=1}^{n}(x_{2,i} - \mu_2)^2$$

$$\text{cov}(x_1, x_2) = \frac{1}{N}\sum_{i=1}^{n}(x_{1,i} - \mu_1)(x_{2,i} - \mu_2)$$

- Find the basis of $\Sigma = S$
  - find vectors $e_i$ which aren't turned by $\Sigma$
    - $\Sigma e_i = \lambda_i e_i$: eigenvalue / eigenvector
  - 1$^{st}$ PC: "longest" $e_i$ (has largest $\lambda_i$), 2$^{nd}$ PC: next longest, ...

|             | $e_1$ | $e_2$ |
|-------------|-------|-------|
| $\lambda_1$ 0.26 | x 0.4 | -0.9 |
| $\lambda_2$ 2.42 | y -0.9 | -0.4 |

$w_1$

$\|w_1\| = 1$

$= (0.9)^2 + (0.4)^2$

$= 0.81 + 0.16$

$w_2 \qquad = 0.81 + 0.16$

$\approx 1$

pca.pdf (ed.ac.uk)

# PCA in a nutshell

## 1. correlated hi-d data
("urefu" means "height" in Swahili)

urefu [cm]

height [inches]

## 2. center the points

u

h

want dimension of
highest variance

## 3. compute covariance matrix

$$\begin{array}{cc} & h \quad u \\ h & \begin{pmatrix} 2.0 & \boxed{0.8} \\ 0.8 & 0.6 \end{pmatrix} \end{array} \rightarrow \mathrm{cov}(h,u) = \frac{1}{n}\sum_{i=1}^{n} h_i u_i$$

## 4. eigenvectors + eigenvalues

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix}\begin{pmatrix} e_h \\ e_u \end{pmatrix} = \lambda_e \begin{pmatrix} e_h \\ e_u \end{pmatrix}$$

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix}\begin{pmatrix} f_h \\ f_u \end{pmatrix} = \lambda_f \begin{pmatrix} f_h \\ f_u \end{pmatrix}$$

`eig(cov(data))`

## 5. pick m<d eigenvectors w. highest eigenvalues

u

e

h

f

$m = 1$

$W$

$$z^i = \begin{bmatrix} w_1^T \\ \vdots \\ w_m^T \end{bmatrix} x^i$$

## 6. project data points to those eigenvectors

$$x_e' = x^T e = \sum_{j=1}^{d} x_{ij} e_j$$

## 7. uncorrelated low-d data

e

# PCA example: Eigen Faces

input: dataset of N face images

face: K x K bitmap of pixels

"unfold" each bitmap to $K^2$-dimensional vector

arrange in a matrix
each face = column

$K^2$ x N
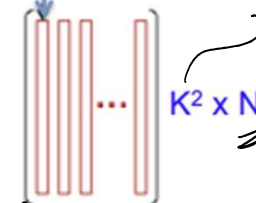
$= d$

PCA

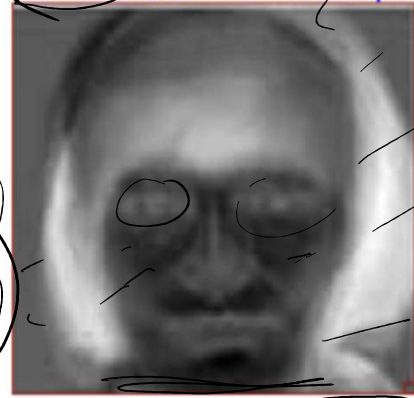"fold" into a K x K bitmap

can visualize eigenvectors: *m* "aspects" of prototypical facial features

$K^2$ x m

set of *m* eigenvectors each is $K^2$-dimensional

Create the matrix.
(Can find PCA without explicitly creating $\Sigma$)

$\Sigma$ or S

# Eigen Faces: Projection



= 0.9 *   - 0.2 *   + 0.4 *   + ...

- Project new face to space of eigen-faces
- Represent vector as a linear combination of principal components
- How many do we need?

m=1  m=2  .            m=6

# (Eigen) Face Recognition

- Face similarity
  - in the reduced space
  - insensitive to lighting expression, orientation
- Projecting new "faces"
  - everything is a face



new face (not in training)

projected to eigenfaces

# Non-linear dimensionality reduction

- Also called manifold learning
- Manifold is a topological space that is locally Euclidean
  - Example: surface of the earth is a curved 2d surface embedded in a high-dimensional space by at each point on the surface, the earth seems flat
- Manifold hypothesis:
  - Most "naturally occurring" high-dimensional dataset lie on a low-dimensional manifold, also called the intrinsic dimensonality of the data

- Many methods exist for learning manifolds: main idea is to preserve local neighborhood of each point in the given dataset.

# Examples

# Stochastic Neighborhood Embedding (SNE)

- Convert high-dimensional Euclidean distances into conditional probabilities that represent similarities