

*Machines will be capable, within twenty years, of doing any work that a man can do.*

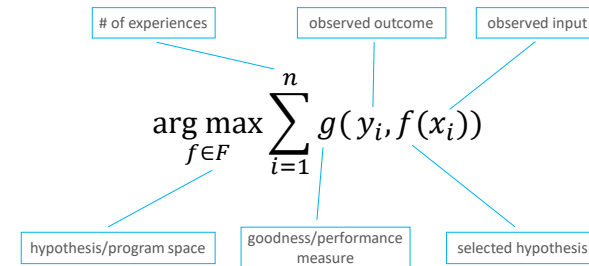
- Herbert Simon, 1965

# CSE552 Machine Learning

Spring 2025  
Dimensionality Reduction  
© 2013-2025 Yakup Genc

1

## Supervised Machine Learning



Spring 2025

CSE552 Machine Learning

2

2

## What is Feature Reduction?

- Feature reduction refers to the mapping of the original high-dimensional data onto a lower-dimensional space
  - Criterion for feature reduction can be different based on different problem settings
    - Unsupervised setting: minimize the information loss
    - Supervised setting: maximize the class discrimination
- Given a set of data points, compute the linear transformation (projection)

$$G \in \mathbb{R}^{d \times k} : x \in \mathbb{R}^d \rightarrow y = G^T x \in \mathbb{R}^k \quad (k \ll d)$$

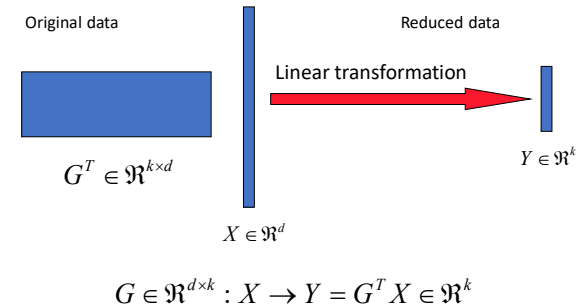
Spring 2025

CSE552 Machine Learning

3

3

## What is Feature Reduction?



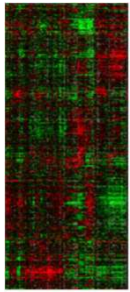
Spring 2025

CSE552 Machine Learning

4

4

## High-dimensional Data



Gene expression



Face images



Handwritten digits

Spring 2025

CSE552 Machine Learning

5

## Why Reduce Dimensionality?

1. Reduces time complexity: Less computation
2. Reduces space complexity: Less parameters (compression)
3. Saves the cost of observing the feature
4. Simpler models are more robust on small datasets
5. More interpretable; simpler explanation
6. Noise removal
7. Data visualization (structure, groups, outliers, etc) if plotted in 2 or 3 dimensions

Spring 2025

CSE552 Machine Learning

6

## Accuracy, Dimensions & Overfitting

- Most machine learning and data mining techniques may not be effective for high-dimensional data
  - Curse of Dimensionality
  - Query accuracy and efficiency degrade rapidly as the dimension increases
- The intrinsic dimension may be small
  - For example, the number of genes responsible for a certain type of disease may be small

Spring 2025

CSE552 Machine Learning

7

## Dimensionality Reduction Algorithms

- Unsupervised
  - Principal Component Analysis (PCA)
  - Latent Semantic Indexing (LSI): truncated SVD
  - Independent Component Analysis (ICA)
  - Canonical Correlation Analysis (CCA)
- Supervised
  - Linear Discriminant Analysis (LDA)
- Semi-supervised
  - Research topic

Spring 2025

CSE552 Machine Learning

8

## Feature Selection vs Extraction

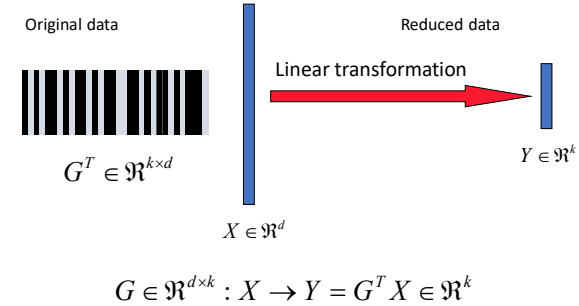
- **Feature selection:** Choosing  $k < d$  important features, ignoring the remaining  $d - k$ 
  - Subset selection algorithms
- **Feature extraction:** Project the original  $x_i, i = 1, \dots, d$  dimensions to new  $k < d$  dimensions,  $z_j, j = 1, \dots, k$ 
  - Principal components analysis (PCA),
  - Linear discriminant analysis (LDA),
  - Factor analysis (FA) ...

Spring 2025

CSE552 Machine Learning

9

## Feature Selection



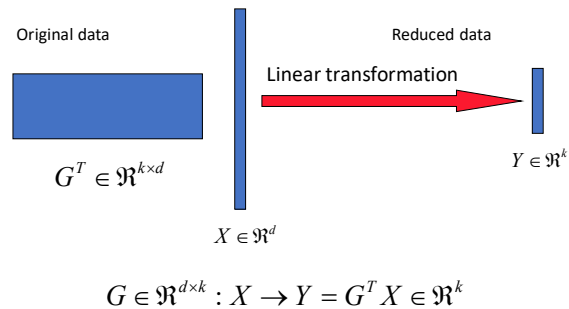
Spring 2025

CSE552 Machine Learning

10

10

## Feature Extraction



Spring 2025

CSE552 Machine Learning

11

11

## Subset Selection

- There are  $2^d$  subsets of  $d$  features
- **Forward search:** Add the best feature at each step
  - Set of features  $F$  initially  $\emptyset$
  - At each iteration, find the best new feature
 
$$j = \arg \min_i E(F \cup x_i)$$
  - Add  $x_j$  to  $F$  if  $E(F \cup x_j) < E(F)$
  - Greedy  $O(d^2)$  algorithm
- **Backward search:** Start with all features and remove one at a time, if possible
- **Floating search:** Add  $k$ , remove  $l$

Spring 2025

CSE552 Machine Learning

12

12

## Key Feature Selection Methods

- Open-loop (filter / front-end / preset bias)
  - Select features for which the reduced data set maximizes between-class separability (by evaluating within-class and between-class covariance matrices)
  - no feedback mechanism from the processing algorithm
- Closed-loop (wrapper/ performance bias)
  - Select features based on the processing algorithm performance (feedback mechanism), which serves as a criterion for feature subset selection

Result: data set with reduced number of features according to a specified optimal criterion

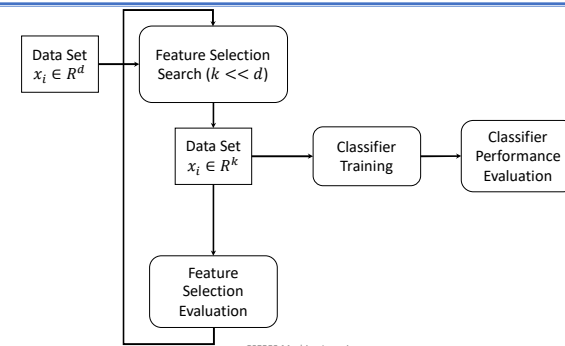
What about decision trees?

Spring 2025

CSE552 Machine Learning

13

## Open Loop



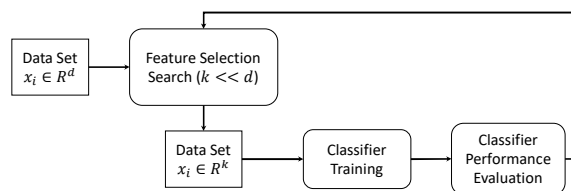
Spring 2025

CSE552 Machine Learning

14

14

## Closed Loop



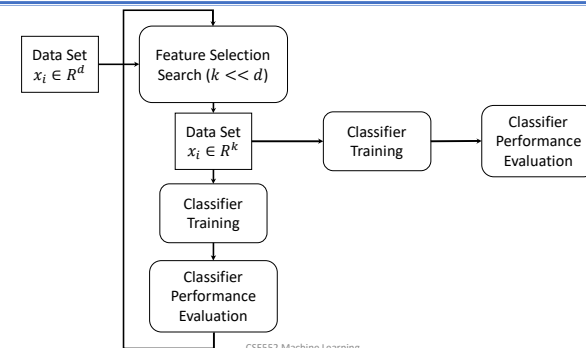
Spring 2025

CSE552 Machine Learning

15

15

## Open Loop



Spring 2025

CSE552 Machine Learning

16

16

## Optimal Feature Selection

- Procedure for optimal FS:
  - Search procedure, to search through candidate subsets of features (given initial step of a search and stop criteria)
  - FS criterion,  $J_i$ , to judge if one subset of features is better than another
- Since feature selection methods are computationally intensive we use heuristic search methods; as a result only sub-optimal solutions can be obtained

Spring 2025

CSE552 Machine Learning

17

## Feature Selection

- FS criteria
  - We use criteria based on maximization, where a better subset of features always gives a bigger value of a criterion
  - and the optimal feature subset gives the maximum value of the criterion
- In practice:
  - For the limited data set and FS criterion based on a classifier performance, removing a feature may improve algorithm's performance (up to a point as it then starts to degrade) – peaking phenomenon

Spring 2025

CSE552 Machine Learning

18

17

18

## FS Paradigms

- Paradigms of optimal FS: minimal representations
- Occam's Razor:
  - The simplest explanation of the observed phenomena in a given domain is the most likely to be a correct one.
- Minimal Description Length (MDL) Principle:
  - Best feature selection can be done by choosing a minimal feature subset that fully describes all classes in a given data set.

Spring 2025

CSE552 Machine Learning

19

19

## MDL Principle

- Can be seen as a formalization of the Occam's razor heuristic
- In short, if a system can be defined in terms of input and the corresponding output data, then in the worst case (longest) it can be described by supplying the entire data set
- On the other hand, if regularities can be discovered, then a much shorter description is possible and can be measured by the MDL principle

Spring 2025

CSE552 Machine Learning

20

20

## Feature Selection

- Criteria
  - A feature selection algorithm uses predefined feature selection criterion (which measures goodness of the subset of features)
- Hope (via MDL principle) is that:
  - by reducing dimensionality we improve generalization ability, up to some max value, but we know that it will start to degrade at some point of reduction

Spring 2025

CSE552 Machine Learning

21

21

## Search

- Goal of SEARCH METHODS: search only through a subset of all possible feature subsets.
- Only sub-optimal subset of features is obtained but at a (much) lower cost.
- Reason
  - The number of possible feature subsets is  $2^n$  where  $n$  – original number of features;
  - search for that number of subsets is computationally very expensive.
- Optimal feature selection is NP-hard thus we need to use sub-optimal feature selection methods

Spring 2025

CSE552 Machine Learning

22

22

## Search Methods

- Exhaustive search
- Branch and Bound
- Individual Feature Ranking
- Sequential Forward and Backward FS
- Stepwise Forward Search
- Stepwise Backward Search
- Probabilistic FS

Spring 2025

CSE552 Machine Learning

23

23

## What is Principal Component Analysis?

- Principal component analysis (PCA)
  - Reduce the dimensionality of a data set by finding a new set of variables, smaller than the original set of variables
  - Retains most of the sample's information
  - Useful for the compression and classification of data
- By information we mean the variation present in the sample, given by the correlations between the original variables
  - The new variables, called principal components (PCs), are uncorrelated, and are ordered by the fraction of the total information each retains

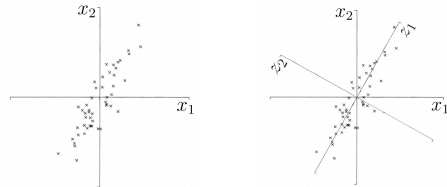
Spring 2025

CSE552 Machine Learning

24

24

## Geometry



- 1<sup>st</sup> principle component (PC)  $z_1$  is a minimum distance fit to a line in  $X$  space
- 2<sup>nd</sup> PC  $z_2$  is a minimum distance fit to a line in the plane perpendicular to the 1<sup>st</sup> PC
- PCs are a series of linear least squares fits to a sample, each orthogonal to all the previous

Spring 2025

CSE552 Machine Learning

25

## Algebraic Definition of PCs

Given a sample of  $n$  observations on a vector of  $d$  variables

$$\{x_1, x_2, \dots, x_n\} \in \mathbb{R}^d$$

define the first principal component of the sample by the linear transformation

$$z_1 = a_1^T x_j = \sum_{i=1}^d a_{i1} x_{ij}, \quad j = 1, 2, \dots, n.$$

where the vector

$$a_1 = (a_{11}, a_{21}, \dots, a_{d1})$$

$$x_j = (x_{1j}, x_{2j}, \dots, x_{dj})$$

is chosen such that  $\text{var}[z_1]$  is maximum.

Spring 2025

CSE552 Machine Learning

26

## Algebraic Definition of PCs

To find  $a_1$  first note that  $\text{var}[z_1] = E((z_1 - \bar{z}_1)^2) = \frac{1}{n} \sum_{i=1}^n (a_1^T x_i - a_1^T \bar{x})^2$

$$= \frac{1}{n} \sum_{i=1}^n a_1^T (x_i - \bar{x})(x_i - \bar{x})^T a_1 = a_1^T S a_1$$

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

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \text{ is the mean.}$$

In the following, we assume the Data is centered.

$$\bar{x} = 0$$

Spring 2025

CSE552 Machine Learning

27

## Algebraic Definition of PCs

Assume  $\bar{x} = 0$

Form the matrix:  $X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{d \times n}$

then  $S = \frac{1}{n} X X^T$

Obtain eigenvectors of  $S$  by computing the SVD of  $X$ :

$$X = U \Sigma V^T$$

Spring 2025

CSE552 Machine Learning

28

## Algebraic Definition of PCs

To find  $a_1$  that maximizes  $\text{var}[z_1]$  subject to  $a_1^T a_1 = 1$

Let  $\lambda$  be a Lagrange multiplier

$$L = a_1^T S a_1 - \lambda(a_1^T a_1 - 1)$$

$$\frac{\partial}{\partial a_1} L = S a_1 - \lambda a_1 = 0$$

$$\Rightarrow (S - \lambda I_p) a_1 = 0$$

therefore  $a_1$  is an eigenvector of  $S$

corresponding to the largest eigenvalue  $\lambda = \lambda_1$ .

Spring 2025

CSE552 Machine Learning

29

29

## Algebraic Definition of PCs

To find the next coefficient vector  $a_2$  maximizing  $\text{var}[z_2]$

subject to  $\text{cov}[z_2, z_1] = 0$

and to  $a_2^T a_2 = 1$

**uncorrelated**

First note that  $\text{cov}[z_2, z_1] = a_1^T S a_2 = \lambda_1 a_1^T a_2$

then let  $\lambda$  and  $\phi$  be Lagrange multipliers, and maximize

$$L = a_2^T S a_2 - \lambda(a_2^T a_2 - 1) - \phi a_1^T a_2$$

Spring 2025

CSE552 Machine Learning

30

30

## Algebraic Definition of PCs

$$L = a_2^T S a_2 - \lambda(a_2^T a_2 - 1) - \phi a_1^T a_2$$

$$\frac{\partial}{\partial a_2} L = S a_2 - \lambda a_2 - \phi a_1 = 0 \Rightarrow \phi = 0$$

$$S a_2 = \lambda a_2 \quad \text{and} \quad \lambda = a_2^T S a_2$$

Spring 2025

CSE552 Machine Learning

31

31

## Algebraic Definition of PCs

We find that  $a_2$  is also an eigenvector of  $S$

whose eigenvalue  $\lambda = \lambda_2$  is the second largest.

In general

$$\text{var}[z_k] = a_k^T S a_k = \lambda_k$$

- The  $k^{\text{th}}$  largest eigenvalue of  $S$  is the variance of the  $k^{\text{th}}$  PC.
- The  $k^{\text{th}}$  PC  $z_k$  retains the  $k^{\text{th}}$  greatest fraction of the variation in the sample.

Spring 2025

CSE552 Machine Learning

32

32



## Algebraic Definition of PCs

- Main steps for computing PCs
  - Form the covariance matrix  $\Sigma$ .
  - Compute its eigenvectors:  $\{a_i\}_{i=1}^d$
  - Use the first  $d$  eigenvectors  $\{a_i\}_{i=1}^d$  to form the  $d$  PCs.
  - The transformation  $G$  is given by  $G \leftarrow [a_1, a_2, \dots, a_k]$

A testpoint  $x \in \mathcal{R}^d \rightarrow G^T x \in \mathcal{R}^k$ .

Spring 2025

CSE552 Machine Learning

33

33

## Principal Components Analysis (PCA)

- Find a low-dimensional space such that when  $\mathbf{x}$  is projected there, information loss is minimized
- The projection of  $\mathbf{x}$  on the direction of  $\mathbf{w}$  is:  $z = \mathbf{w}^T \mathbf{x}$
- Find  $\mathbf{w}$  such that  $\text{Var}(z)$  is maximized

$$\begin{aligned} \text{Var}(z) &= \text{Var}(\mathbf{w}^T \mathbf{x}) = E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})^2] \\ &= E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})] \\ &= E[\mathbf{w}^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{w}] \\ &= \mathbf{w}^T E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \mathbf{w} = \mathbf{w}^T \Sigma \mathbf{w} \end{aligned}$$

where

$$\text{Var}(\mathbf{x}) = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] = \Sigma$$

Spring 2025

CSE552 Machine Learning

34

34

## Principal Components Analysis (PCA)

- Maximize  $\text{Var}(z)$  subject to  $\|\mathbf{w}\| = 1$

$$\max_{\mathbf{w}_1} \mathbf{w}_1^T \Sigma \mathbf{w}_1 - \alpha (\mathbf{w}_1^T \mathbf{w}_1 - 1)$$

$\Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1$  that is,  $\mathbf{w}_1$  is an eigenvector of  $\Sigma$

Choose the one with the largest eigenvalue for  $\text{Var}(z)$  to be max

- Second principal component: Max  $\text{Var}(z_2)$ , s.t.,  $\|\mathbf{w}_2\| = 1$  and orthogonal to  $\mathbf{w}_1$

$$\max_{\mathbf{w}_2} \mathbf{w}_2^T \Sigma \mathbf{w}_2 - \alpha (\mathbf{w}_2^T \mathbf{w}_2 - 1) - \beta (\mathbf{w}_2^T \mathbf{w}_1 - 0)$$

$\Sigma \mathbf{w}_2 = \alpha \mathbf{w}_2$  that is,  $\mathbf{w}_2$  is another eigenvector of  $\Sigma$  and so on.

Spring 2025

CSE552 Machine Learning

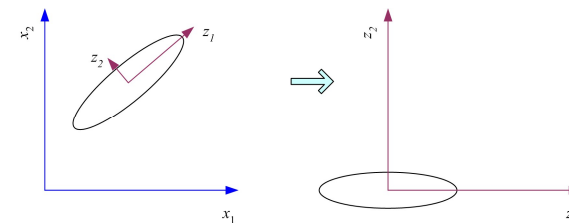
35

35

## What PCA Does

$$\mathbf{z} = \mathbf{W}^T (\mathbf{x} - \mathbf{m})$$

where the columns of  $\mathbf{W}$  are the eigenvectors of  $\Sigma$ , and  $\mathbf{m}$  is sample mean  
Centers the data at the origin and rotates the axes



Spring 2025

CSE552 Machine Learning

36

36

## How to choose k ?

- Proportion of Variance (PoV) explained

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_d}$$

when  $\lambda_i$  are sorted in descending order

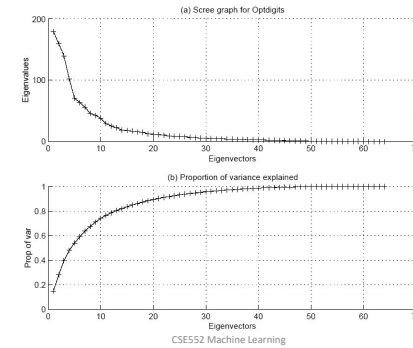
- Typically, stop at PoV > 0.9
- Scree graph plots of PoV vs k, stop at “elbow”

Spring 2025

CSE552 Machine Learning

37

## How to choose k ?

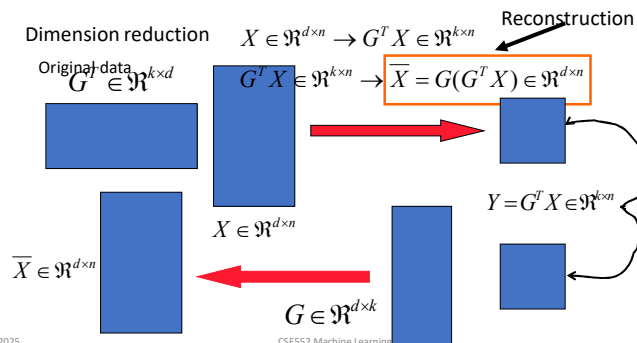


Spring 2025

CSE552 Machine Learning

38

## Optimality Property of PCA



Spring 2025

CSE552 Machine Learning

39

## Optimality Property of PCA

### Main theoretical result:

The matrix  $G$  consisting of the first  $d$  eigenvectors of the covariance matrix  $S$  solves the following min problem:

$$\min_{G \in \mathbb{R}^{d \times k}} \|X - G(G^T X)\|_F^2 \text{ subject to } G^T G = I_k$$

$$\|X - \bar{X}\|_F^2$$

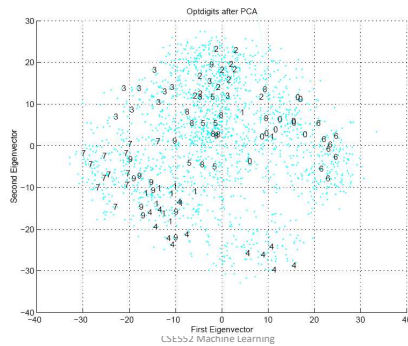
reconstruction error

PCA projection minimizes the reconstruction error among all linear projections of size  $k$ .

Spring 2025

CSE552 Machine Learning

40



Spring 2025

CSE552 Machine Learning

41

## Factor Analysis

- Find a small number of **factors**  $\mathbf{z}$ , which when combined generate  $\mathbf{x}$ :

$$x_i - \mu_i = v_{i1}z_1 + v_{i2}z_2 + \dots + v_{ik}z_k + \varepsilon_i$$

where  $z_j, j=1, \dots, k$  are the **latent factors** with

$$E[z_j] = 0, \text{Var}(z_j) = 1, \text{Cov}(z_i, z_j) = 0, i \neq j,$$

$\varepsilon_i$  are the **noise sources**

$$E[\varepsilon_i] = \psi_i, \text{Cov}(\varepsilon_i, \varepsilon_j) = 0, i \neq j, \text{Cov}(\varepsilon_i, z_j) = 0,$$

and  $v_{ij}$  are the **factor loadings**

Spring 2025

CSE552 Machine Learning

42

## Factor Analysis

- Find a small number of **factors**  $\mathbf{z}$ , which when combined generate  $\mathbf{x}$ :

$$x_i - \mu_i = v_{i1}z_1 + v_{i2}z_2 + \dots + v_{ik}z_k + \varepsilon_i$$

where  $z_j, j=1, \dots, k$  are the **latent factors** with

$$E[z_j] = 0, \text{Var}(z_j) = 1, \text{Cov}(z_i, z_j) = 0, i \neq j,$$

$\varepsilon_i$  are the **noise sources**

$$E[\varepsilon_i] = \psi_i, \text{Cov}(\varepsilon_i, \varepsilon_j) = 0, i \neq j, \text{Cov}(\varepsilon_i, z_j) = 0,$$

and  $v_{ij}$  are the **factor loadings**

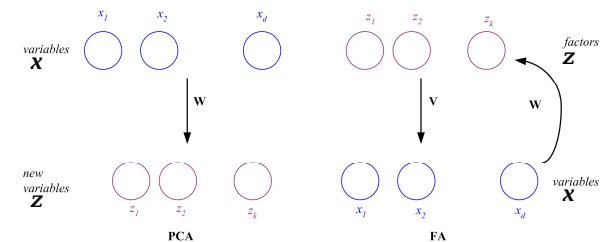
Spring 2025

CSE552 Machine Learning

43

## PCA vs FA

- PCA From  $\mathbf{x}$  to  $\mathbf{z}$   $\mathbf{z} = \mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu})$
- FA From  $\mathbf{z}$  to  $\mathbf{x}$   $\mathbf{x} - \boldsymbol{\mu} = \mathbf{V}\mathbf{z} + \boldsymbol{\varepsilon}$



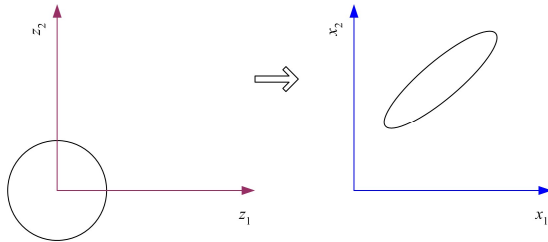
Spring 2025

CSE552 Machine Learning

44

## Factor Analysis

- In FA, factors  $z_j$  are stretched, rotated and translated to generate  $\mathbf{x}$



Spring 2025

CSE552 Machine Learning

45

## Multidimensional Scaling

- Given pairwise distances between  $N$  points,  
 $d_{ij}, i, j = 1, \dots, N$   
 place on a low-dim map s.t. distances are preserved.

- $\mathbf{z} = \mathbf{g}(\mathbf{x} \mid \vartheta)$  Find  $\vartheta$  that min **Sammon stress**

$$E(\theta \mid \mathcal{X}) = \sum_{r,s} \frac{(\|\mathbf{z}^r - \mathbf{z}^s\| - \|\mathbf{x}^r - \mathbf{x}^s\|)^2}{\|\mathbf{x}^r - \mathbf{x}^s\|^2}$$

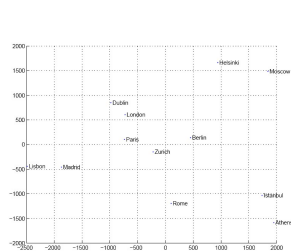
$$= \sum_{r,s} \frac{(\|\mathbf{g}(\mathbf{x}^r \mid \theta) - \mathbf{g}(\mathbf{x}^s \mid \theta)\| - \|\mathbf{x}^r - \mathbf{x}^s\|)^2}{\|\mathbf{x}^r - \mathbf{x}^s\|^2}$$

Spring 2025

CSE552 Machine Learning

46

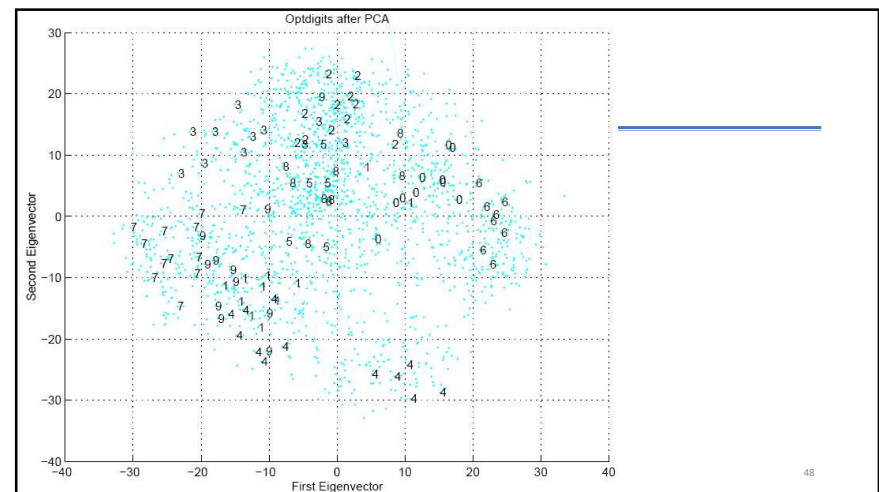
## Map of Europe by MDS

Map from CIA - The World Factbook: <http://www.cia.gov/>

Spring 2025

CSE552 Machine Learning

47



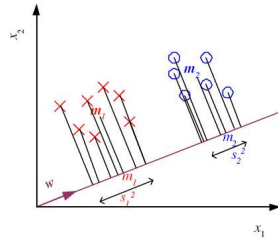
48

## Linear Discriminant Analysis

- Find a low-dimensional space such that when  $\mathbf{x}$  is projected, classes are well-separated.
- Find  $\mathbf{w}$  that maximizes

$$J(\mathbf{w}) = \frac{(\mathbf{m}_1 - \mathbf{m}_2)^2}{s_1^2 + s_2^2}$$

$$\mathbf{m}_1 = \frac{\sum_t \mathbf{w}^T \mathbf{x}^t r^t}{\sum_t r^t} \quad s_1^2 = \sum_t (\mathbf{w}^T \mathbf{x}^t - \mathbf{m}_1)^2 r^t$$



Spring 2025

CSE552 Machine Learning

49

49

## Linear Discriminant Analysis

- Between-class scatter:

$$\begin{aligned} (\mathbf{m}_1 - \mathbf{m}_2)^2 &= (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\ &= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} \\ &= \mathbf{w}^T \mathbf{S}_B \mathbf{w} \quad \text{where } \mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \end{aligned}$$

- Within-class scatter:

$$\begin{aligned} s_1^2 &= \sum_t (\mathbf{w}^T \mathbf{x}^t - \mathbf{m}_1)^2 r^t \\ &= \sum_t \mathbf{w}^T (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} r^t = \mathbf{w}^T \mathbf{S}_1 \mathbf{w} \end{aligned}$$

where  $\mathbf{S}_1 = \sum_t (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T r^t$

$$s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_W \mathbf{w} \quad \text{where } \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

Spring 2025

CSE552 Machine Learning

50

50

## Fisher's Linear Discriminant

- Find  $\mathbf{w}$  that max

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} = \frac{|\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|^2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

- LDA solution:

$$\mathbf{w} = \mathbf{C} \cdot \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

- Parametric solution:

$$\mathbf{w} = \Sigma^{-1} (\mu_1 - \mu_2)$$

when  $p(\mathbf{x} | C_i) \sim \mathcal{N}(\mu_i, \Sigma)$

Spring 2025

CSE552 Machine Learning

51

51

## K>2 Classes

- Within-class scatter:

$$\mathbf{S}_W = \sum_{i=1}^K \mathbf{S}_i \quad \mathbf{S}_i = \sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i) (\mathbf{x}^t - \mathbf{m}_i)^T$$

- Between-class scatter:

$$\mathbf{S}_B = \sum_{i=1}^K N_i (\mathbf{m}_i - \mathbf{m}) (\mathbf{m}_i - \mathbf{m})^T \quad \mathbf{m} = \frac{1}{K} \sum_{i=1}^K \mathbf{m}_i$$

- Find  $\mathbf{W}$  that max

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}$$

The largest eigenvectors of  $\mathbf{S}_W^{-1} \mathbf{S}_B$   
Maximum rank of K-1

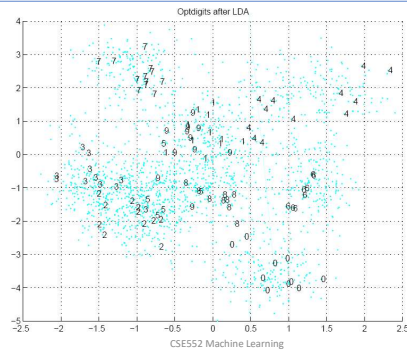
Spring 2025

CSE552 Machine Learning

52

52

## LDA

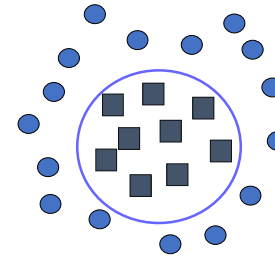


Spring 2025

CSE552 Machine Learning

53

## Motivation



Linear projections will not detect the pattern

Spring 2025

CSE552 Machine Learning

54

## Nonlinear PCA using Kernels

- Traditional PCA applies linear transformation
  - May not be effective for nonlinear data
- Solution: apply nonlinear transformation to potentially very high-dimensional space.

$$\phi: x \rightarrow \phi(x)$$

- Computational efficiency: apply the kernel trick.
  - Require PCA can be rewritten in terms of dot product.

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

Spring 2025

CSE552 Machine Learning

55

## Nonlinear PCA using Kernels

### Rewrite PCA in terms of dot product

Assume the data has been centered, i.e.,  $\sum_i x_i = 0$ .

The covariance matrix  $S$  can be written as  $S = \frac{1}{n} \sum_i x_i x_i^T$

Let  $v$  be The eigenvector of  $S$  corresponding to nonzero eigenvalue

$$Sv = \frac{1}{n} \sum_i x_i x_i^T v = \lambda v \Rightarrow v = \frac{1}{n\lambda} \sum_i (x_i^T v) x_i$$

Eigenvectors of  $S$  lie in the space spanned by all data points.

Spring 2025

CSE552 Machine Learning

56

## Nonlinear PCA using Kernels

$$Sv = \frac{1}{n} \sum_i x_i x_i^T v = \lambda v \Rightarrow v = \frac{1}{n\lambda} \sum_i (x_i^T v) x_i$$

The covariance matrix can be written in matrix form:

$$S = \frac{1}{n} XX^T, \text{ where } X = [x_1, x_2, \dots, x_n].$$

$$v = \sum_i \alpha_i x_i = X\alpha \quad Sv = \frac{1}{n} XX^T X\alpha = \lambda X\alpha$$

$$\frac{1}{n} (X^T X)(X^T X)\alpha = \lambda (X^T X)\alpha$$

$$\frac{1}{n} (X^T X)\alpha = \lambda \alpha$$

Any benefits?

Spring 2025

CSE552 Machine Learning

57

## Nonlinear PCA using Kernels

Next consider the feature space:  $\phi: x \rightarrow \phi(x)$

$$S^\phi = \frac{1}{n} X^\phi (X^\phi)^T, \text{ where } X^\phi = [x_1^\phi, x_2^\phi, \dots, x_n^\phi].$$

$$v = \sum_i \alpha_i \phi(x_i) = X^\phi \alpha \quad \frac{1}{n} (X^\phi)^T X^\phi \alpha = \lambda \alpha$$

The (i,j)-th entry of  $(X^\phi)^T X^\phi$  is  $\phi(x_i) \bullet \phi(x_j)$

Apply the kernel trick:  $K(x_i, x_j) = \phi(x_i) \bullet \phi(x_j)$

K is called the kernel matrix.

$$\frac{1}{n} K \alpha = \lambda \alpha$$

Spring 2025

CSE552 Machine Learning

58

## Nonlinear PCA using Kernels

- Projection of a test point  $x$  onto  $v$ :

$$\begin{aligned} \phi(x) \bullet v &= \phi(x) \bullet \sum_i \alpha_i \phi(x_i) \\ &= \sum_i \alpha_i \phi(x) \bullet \phi(x_i) = \sum_i \alpha_i K(x, x_i) \end{aligned}$$

Explicit mapping is not required here.

Spring 2025

CSE552 Machine Learning

59

## Eigenfaces

(Slides adapted from Lazebnik, Grauman & Lowe)

60

## Face Detection and Recognition



Spring 2025

CSE552 Machine Learning

61

61

## Face Detection and Recognition



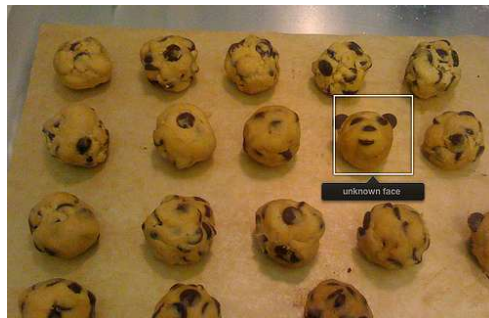
Spring 2025

CSE552 Machine Learning

62

62

## Consumer Application: iPhoto 2009



Spring 2025

CSE552 Machine Learning

63

63

## The Space of All Face Images

- When viewed as vectors of pixel values, face images are extremely high-dimensional
  - 100x100 image = 10,000 dimensions
- However, relatively few 10,000-dimensional vectors correspond to valid face images
- We want to effectively model the subspace of face images



Spring 2025

CSE552 Machine Learning

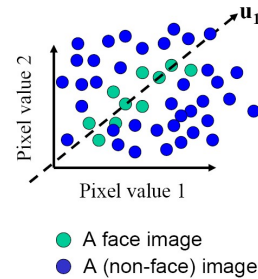
64

64



## The Space of All Face Images

We want to construct a low-dimensional linear subspace that best explains the variation in the set of face images



Spring 2025

CSE552 Machine Learning

65

65

## Principal Component Analysis

- Given:  $N$  data points  $\mathbf{x}_1, \dots, \mathbf{x}_N$  in  $\mathbb{R}^d$
- We want to find a new set of features that are linear combinations of original ones:

$$u(\mathbf{x}_i) = \mathbf{u}^T(\mathbf{x}_i - \boldsymbol{\mu})$$

( $\boldsymbol{\mu}$ : mean of data points)

- What unit vector  $\mathbf{u}$  in  $\mathbb{R}^d$  captures the most variance of the data?

Spring 2025

CSE552 Machine Learning

66

66

## Principal Component Analysis

Direction that maximizes the variance of the projected data:

$$\begin{aligned} \text{var}(u) &= \frac{1}{N} \sum_{i=1}^N \underbrace{\mathbf{u}^T(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{u}^T(\mathbf{x}_i - \boldsymbol{\mu}))^T}_{\text{Projection of data point}} \\ &= \mathbf{u}^T \left[ \underbrace{\sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T}_{\text{Covariance matrix of data}} \right] \mathbf{u} \\ &= \mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u} \end{aligned}$$

The direction that maximizes the variance is the eigenvector associated with the largest eigenvalue of  $\boldsymbol{\Sigma}$

Spring 2025

CSE552 Machine Learning

67

67

## Principal Component Analysis

- The direction that captures the maximum covariance of the data is the eigenvector corresponding to the largest eigenvalue of the data covariance matrix
- Furthermore, the top  $k$  orthogonal directions that capture the most variance of the data are the  $k$  eigenvectors corresponding to the  $k$  largest eigenvalues

Spring 2025

CSE552 Machine Learning

68

68

## Eigenfaces: Key idea

- Assume that most face images lie on a low-dimensional subspace determined by the first  $k$  ( $k < d$ ) directions of maximum variance
- Use PCA to determine the vectors or “eigenfaces”  $u_1, \dots, u_k$  that span that subspace
- Represent all face images in the dataset as linear combinations of eigenfaces

M. Turk and A. Pentland, [Face Recognition using Eigenfaces](#), CVPR 1991

Spring 2025

CSE552 Machine Learning

69

## Eigenfaces (1)

- Calculation of Eigenfaces
  - Calculate average face :  $v$
  - Collect difference between training images and average face in matrix  $A$  ( $M$  by  $N$ ), where  $M$  is the number of pixels and  $N$  is the number of images

$$A = [u_1^1 - v, \dots, u_n^1 - v, \dots, u_1^p - v, \dots, u_n^p - v]$$

- The eigenvectors of covariance matrix  $C$  ( $M$  by  $M$ ) give the eigenfaces
- $M$  is usually big, so this process would be time consuming
- What to do?

$$C = AA^T$$

Spring 2025

CSE552 Machine Learning

70

## Eigenfaces (2)

- Calculation of Eigenvectors of  $C$ 
  - If the number of data points is smaller than the dimension ( $N < M$ ), then there will be only  $N-1$  meaningful eigenvectors.
  - Instead of directly calculating the eigenvectors of  $C$ , we can calculate the eigenvalues and the corresponding eigenvectors of a much smaller matrix  $L$  ( $N$  by  $N$ )

$$L = A^T A$$

- If  $\lambda_i$  are the eigenvalues of  $L$  then  $A \lambda_i$  are the eigenvectors for  $C$ 
  - The eigenvectors are in the descent order of the corresponding eigenvalues

Spring 2025

CSE552 Machine Learning

71

## Eigenfaces (3)

- Representation of Face Images using Eigenfaces
- The training face images and new face images can be represented as linear combination of the eigenfaces.
- When we have a face image  $u$  :

$$u = \sum_i a_i \phi_i$$

- Since the eigenvectors are orthogonal :

$$a_i = u^T \phi_i$$

Spring 2025

CSE552 Machine Learning

72

## Eigenfaces Example

Training images  $\mathbf{x}_1, \dots, \mathbf{x}_N$



Spring 2025

CSE552 Machine Learning

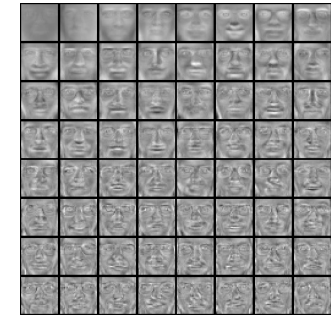
73

## Eigenfaces Example

Mean:  $\mu$



Top eigenvectors:  $\mathbf{u}_1, \dots, \mathbf{u}_k$



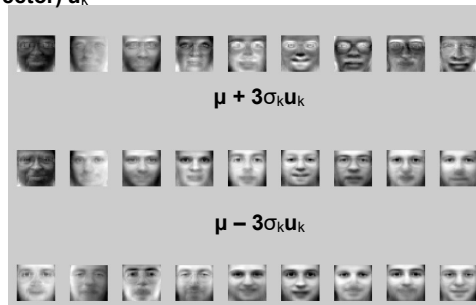
Spring 2025

CSE552 Machine Learning

74

## Eigenfaces Example

Principal component (eigenvector)  $\mathbf{u}_k$



Spring 2025

CSE552 Machine Learning

75

## Eigenfaces Example

Face  $\mathbf{x}$  in "face space" coordinates:



$$\mathbf{x} \rightarrow [\mathbf{u}_1^T(\mathbf{x} - \mu), \dots, \mathbf{u}_k^T(\mathbf{x} - \mu)] \\ = w_1, \dots, w_k$$

Reconstruction:



$$\hat{\mathbf{x}} = \mu + w_1 \mathbf{u}_1 + w_2 \mathbf{u}_2 + w_3 \mathbf{u}_3 + w_4 \mathbf{u}_4 + \dots$$



Spring 2025

CSE552 Machine Learning

76

## Eigenfaces Example



Face  $\mathbf{x}$  in "face space" coordinates:

$$\mathbf{x} \rightarrow [\mathbf{u}_1^T(\mathbf{x} - \boldsymbol{\mu}), \dots, \mathbf{u}_k^T(\mathbf{x} - \boldsymbol{\mu})]$$

$$= w_1, \dots, w_k$$

Spring 2025

CSE552 Machine Learning

77

## Eigenfaces Example



First three eigenfaces

Spring 2025

CSE552 Machine Learning

78

## Recognition with Eigenfaces

- Process labeled training images:
  - Find mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$
  - Find  $k$  principal components (eigenvectors of  $\boldsymbol{\Sigma}$ )  $\mathbf{u}_1, \dots, \mathbf{u}_k$
  - Project each training image  $\mathbf{x}_i$  onto subspace spanned by principal components:  
 $(w_{i1}, \dots, w_{ik}) = (\mathbf{u}_1^T(\mathbf{x}_i - \boldsymbol{\mu}), \dots, \mathbf{u}_k^T(\mathbf{x}_i - \boldsymbol{\mu}))$
- Given novel image  $\mathbf{x}$ :
  - Project onto subspace:  
 $(w_1, \dots, w_k) = (\mathbf{u}_1^T(\mathbf{x} - \boldsymbol{\mu}), \dots, \mathbf{u}_k^T(\mathbf{x} - \boldsymbol{\mu}))$
  - Optional: check reconstruction error  $\mathbf{x} - \hat{\mathbf{x}}$  to determine whether image is really a face
  - Classify as closest training face in  $k$ -dimensional subspace

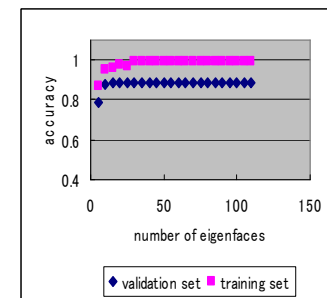
Spring 2025

CSE552 Machine Learning

79

## Classification Using Nearest Neighbor

- Save average coefficients for each person. Classify new face as the person with the closest average
- Recognition accuracy increases with number of eigenfaces till 15.
- Later eigenfaces do not help much with recognition
- Best recognition rates
  - Training set 99%
  - Test set 89%



Spring 2025

CSE552 Machine Learning

80

## Limitations

Global appearance method: not robust to misalignment, background variation



Spring 2025

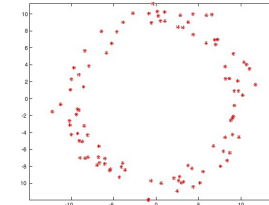
CSE552 Machine Learning

81

81

## Limitations

- PCA assumes that the data has a Gaussian distribution (mean  $\mu$ , covariance matrix  $\Sigma$ )



- The shape of this dataset is not well described by its principal components

Spring 2025

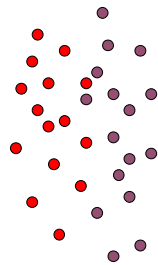
CSE552 Machine Learning

82

82

## Limitations

The direction of maximum variance is not always good for classification



Spring 2025

CSE552 Machine Learning

83

83

Thanks for listening!

84