Data Preparation

Lijun Zhang <u>zlj@nju.edu.cn</u> http://cs.nju.edu.cn/zlj





Outline

Introduction

- □ Feature Extraction and Portability
- Data Cleaning
- Data Reduction and Transformation
 - Sampling
 - Feature Subset Selection
 - Dimensionality Reduction with Axis Rotation
 - Dimensionality Reduction with Type Transformation
- **Summary**



Introduction

"Success depends upon previous preparation, and without such preparation there is sure to be failure."—Confucius

"凡事豫(预)则立,不豫(预)则废"——《礼记·中庸》

□ Feature Extraction and Portability

- Raw logs, documents, semistructured data
- Data may contain heterogeneous types
- Data Cleaning
 - Missing, Erroneous, and Inconsistent
- Data Reduction, Selection, and Transformation
 - Efficiency, Effectiveness



Outline

- Introduction
- Feature Extraction and Portability
- Data Cleaning
- Data Reduction and Transformation
 - Sampling
 - Feature Subset Selection
 - Dimensionality Reduction with Axis Rotation
 - Dimensionality Reduction with Type Transformation
- **Summary**



Feature Extraction

Domain	Raw Data	Features
Sensor	Low-level signals	Wavelet or Fourier transforms
Image	Pixels	Color histograms Visual words
Web logs	Text strings	IP address Action
Network traffic	Characteristics of the network packets	Number of bytes transferred Network protocol
Document data	Text strings	Bag-of-words Entity extraction

Feature extraction is an art form that is highly dependent on the skill of the analyst



Data Type Portability (1)

Data is Often Heterogeneous

A demographic data set may contain both numeric and mixed attributes

Possible Solutions

- Designing an algorithm with an arbitrary combination of data types
 - Time-consuming and sometimes impractical
- Converting between various data types
 - Utilize off-the-shelf tools for processing



Data Type Portability (2)

□ Ways of Transforming Data

Source data type	Destination data type	Methods
Numeric	Categorical	Discretization
Categorical	Numeric	Binarization
Text	Numeric	Latent semantic analysis (LSA)
Time series	Discrete sequence	SAX
Time series	Numeric multidimensional	DWT, DFT
Discrete sequence	Numeric multidimensional	DWT, DFT
Spatial	Numeric multidimensional	2-d DWT
Graphs	Numeric multidimensional	MDS, spectral
Any type	Graphs	Similarity graph
		(Restricted applicability)

Table 2.1: Portability of different data types

Numeric to Categorical Data: Discretization (1)



Divides the ranges of the numeric attribute into ϕ ranges



Age Attribute

✓ [0, 10], [21, 20], [21, 30], ...

Salary

× [0, 10000], [10001, 20000], [20001, 30000], ...

Numeric to Categorical Data: Discretization (2)



Equi-width Ranges

- Each range [a, b] is chosen such that b a is a constant
- Equi-log Ranges
 - Each range [a, b] is chosen such that log b – log a is a constant
 - For example, [1, a], $[a, a^2]$, $[a^2, a^3]$, ...
- Equi-depth Ranges
 - Each range has an equal number of records
 - Sorting and Selecting

Categorical to Numeric Data: Binarization



Two categories
0,1 or -1,1

- $\Box \phi$ categories
 - ϕ -dimensional indicator vector
 - The position of 1 indicates the category

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \xrightarrow{\text{1st Category}} 2^{\text{nd Category}} \\ \xrightarrow{\text{3rd Category}} \phi = 3$$



Text to Numeric Data

Tokenization, Stop Word Removal, Stemming, Weighting (TF-IDF)

Tokenization for Chinese sentence is difficult.

"生产鞋子和服装" "今天真热"



Assignment

Text to Numeric Data

 Tokenization, Stop Word Removal, Stemming, Weighting (TF-IDF)
 Document Term Matrix

三国

Document-Term Matrix

将进酒 念奴娇·赤壁怀古 0.3 0 ···· 0 0.5 ····

酒

Dimensionality Reduction
 Latent Semantic Analysis
 Normalization

Time Series to Discrete Sequence Data



- Symbolic Aggregate Approximation (SAX)
 - Window-based averaging
 - Evaluate the average value in each windows
 - Value-based discretization
 - Discretize the average value by equi-depth intervals
- □ How to Ensure Equi-depth?
 - Assume certain distribution, such as Gaussian
 - Estimate the distribution



Time Series to Numeric Data

□ Discrete Wavelet Transform (DWT)

□ Discrete Fourier transform (DFT)

AdvantagesRemove Dependence

Discrete Sequence to Numeric Data



- Discrete sequence to a Set of (binary) Time Series
 - ACACACTGTGACTG (4 Symbols)
 - 10101000001000 (A)
 - 01010100000100 (C)
 - 00000010100010 (T)
 - 0000001010001 (G)
- Map Each of These Time Series into a Multidimensional Vector
- Features from the Different Series are Combined

Any Type to Graphs for Similarity-Based Applications



- □ A Neighborhood Graph for a Set of *n* Points $\mathcal{O} = \{O_1, ..., O_n\}$
 - A Single Node is defined for each O_i
 - An edge exists between O_i and O_j , if

 $d(O_i, O_j) \leq \epsilon$

The weight W_{ij} of edge (i,j) is defined as $W_{ij} = e^{-\frac{d(o_i, o_j)^2}{t^2}}$

Many Variants Exist



Other Transformations

Spatial to Numeric Data
 Similar to Time-series Data

Graphs to Numeric Data
 Multidimensional Scaling (MDS)
 ✓ Edge represents distance
 Spectral Transformations
 ✓ Edge represents similarity



Outline

- Introduction
- □ Feature Extraction and Portability
- Data Cleaning
- Data Reduction and Transformation
 - Sampling
 - Feature Subset Selection
 - Dimensionality Reduction with Axis Rotation
 - Dimensionality Reduction with Type Transformation
- □ Summary



The Reason of Cleaning

Data Collection Technologies are Inaccurate

- Sensors
- Optical character recognition
- Speech-to-text data
- Privacy Reasons
- Manual Errors
- Data Collection is Expensive
 - Medical Test



Handling Missing Entries

- Delete the Data Record Containing missing entries
 - What to do if nothing left?
- Estimate or Impute the Missing Values
 - Additional errors may be introduced
 - Good under certain conditions (e.g., Matrix Completion)
- Designing an Algorithm that Works with Missing Data

Handling Incorrect and Inconsistent Entries



Inconsistency Detection
 E.g., full name and abbreviation
 Domain Knowledge
 Age cannot be 800
 Data-centric Methods





Scaling and Normalization

Features have Different Scales

- Age versus Salary
- Standardization
 - If the *j*-th attribute has mean μ_j and standard derivation σ_i

$$z_i^j = \frac{x_i^j - \mu_j}{\sigma_j}$$

- Map to [0,1]
- Sensitive to noise

$$z_i^j = \frac{x_i^j - min_j}{max_j - min_j}$$



Outline

- Introduction
- Feature Extraction and Portability
- Data Cleaning

Data Reduction and Transformation

- Sampling
- Feature Subset Selection
- Dimensionality Reduction with Axis Rotation
- Dimensionality Reduction with Type Transformation
- **Summary**





□ The Advantages

- Reduce space complexity
- Reduce time complexity

Reduce noise

- Reveal hidden structures
 - E.g., manifold learning

The DisadvantagesInformation loss



Outline

- Introduction
- □ Feature Extraction and Portability
- Data Cleaning
- Data Reduction and Transformation
 - Sampling
 - Feature Subset Selection
 - Dimensionality Reduction with Axis Rotation
 - Dimensionality Reduction with Type Transformation
- **Summary**



Sampling for Static Data

Unbiased (Uniform) Sampling

- Sampling without replacement
- Sampling with replacement
 - Duplicates are possible
- Biased Sampling
 - Some parts of the data are emphasized
 - E.g., Temporal-decay bias

$$p(\overline{X}) \propto e^{-\lambda \, \delta t}$$

- Stratified Sampling
 - Partition data into a set of strata
 - Sample in each of stratum



An Example of Sampling

- There are 10000 people which contain 100 millionaires
- Unbiased Sampling 100 people
 - In expectation, one millionaire will be sampled
 - In practice, maybe no millionaires are sampled
- Stratified Sampling
 - Unbiased Sampling 1 from 100 millionaires
 - Unbiased Sampling 99 from remaining

Reservoir Sampling for Data Streams



□ The Setting

- Data arrive sequentially
- We want sample k of them uniformly
 - There is a reservoir that can hold k data points
- □ The Algorithm
 - The first k data points are kept
 - Insert the n-th data point with probability k/n
 - ✓ If the *n*-th data is inserted, then drop one of the existing *k* data points uniformly

Reservoir Sampling for Data Streams



□ The Setting

- Data arrive sequentially
- We want sample k of them uniformly

data

After *n* stream points have arrived, the probability of any stream point being included in the reservoir is the same, and is equal to k/n.

Insert the n-th data point with probability k/n

✓ If the *n*-th data is inserted, then drop one of the existing *k* data points uniformly



Outline

- Introduction
- Feature Extraction and Portability
- Data Cleaning
- Data Reduction and Transformation
 - Sampling
 - Feature Subset Selection
 - Dimensionality Reduction with Axis Rotation
 - Dimensionality Reduction with Type Transformation
- **Summary**



Feature Subset Selection

- Unsupervised Feature Selection
 - Using the performance of unsupervised learning (e.g, clustering) to guide the selection
- Supervised Feature Selection
 Using the performance of supervised learning (e.g., classification) to guide the selection



Outline

- Introduction
- Feature Extraction and Portability
- Data Cleaning
- Data Reduction and Transformation
 - Sampling
 - Feature Subset Selection
 - Dimensionality Reduction with Axis Rotation
 - Dimensionality Reduction with Type Transformation
- □ Summary

Dimensionality Reduction with Axis Rotation (1)



Motivations (Perfect Case)

Consider the following 3 points in a 2dimensional space



Dimensionality Reduction with Axis Rotation (2)



Motivations (Perfect Case)

What is the new coordinates if we rotate the axis



Dimensionality Reduction with Axis Rotation (2)



Motivations (Perfect Case)

What is the new coordinates if we rotate the axis



Dimensionality Reduction with Axis Rotation (3)



Motivations (Noisy Case)

Consider the following 3 points in a 2dimensional space


Dimensionality Reduction with Axis Rotation (4)



Motivations (Noisy Case)

What is the new coordinates if we rotate the axis



Dimensionality Reduction with Axis Rotation (4)



Motivations (Noisy Case)

What is the new coordinates if we rotate the axis



Dimensionality Reduction with Axis Rotation (5)



- □ When does it Work?
 - Correlations exist among features
- Axis Rotation
 Remove correlations
 Reduce dimensionality
- How to Determine such Axis System?
 Principal component analysis (PCA)
 Singular value decomposition (SVD)

Axis Rotation—Mathematical Formulation (1)



By default, the Original Coordinates are Defined with respect to the Standard Basis

$$\mathbf{x} = \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^d \end{bmatrix} \in \mathbb{R}^d \iff \mathbf{x} = x^1 \mathbf{e}_1 + x^2 \mathbf{e}_2 + \dots + x^d \mathbf{e}_d$$

Axis Rotation—Mathematical Formulation (2)



- □ The New Coordinates with respect to a Orthonormal Basis {**w**₁, **w**₂, ..., **w**_d}
 - $W = [\mathbf{w}_1, \mathbf{w}_2, ..., \mathbf{w}_d]$ is a orthonormal matrix

$$\mathbf{x} = WW^{\mathsf{T}}\mathbf{x} = \left(\sum_{i=1}^{d} \mathbf{w}_{i}\mathbf{w}_{i}^{\mathsf{T}}\right)\mathbf{x} = \sum_{i=1}^{d} \mathbf{w}_{i}(\mathbf{w}_{i}^{\mathsf{T}}\mathbf{x})$$
$$= (\mathbf{w}_{1}^{\mathsf{T}}\mathbf{x})\mathbf{w}_{1} + (\mathbf{w}_{2}^{\mathsf{T}}\mathbf{x})\mathbf{w}_{2} + \dots + (\mathbf{w}_{d}^{\mathsf{T}}\mathbf{x})\mathbf{w}_{d}$$

Thus, the new coordinates are

$$\mathbf{y} = \begin{bmatrix} \mathbf{w}_1^{\mathsf{T}} \mathbf{x} \\ \mathbf{w}_2^{\mathsf{T}} \mathbf{x} \\ \vdots \\ \mathbf{w}_d^{\mathsf{T}} \mathbf{x} \end{bmatrix} \in \mathbb{R}^d$$

Axis Rotation—Mathematical Formulation (2)



The New Coordinates with respect to a Orthonormal Basis $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_d\}$ $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_d]$ is a orthonormal matrix $\mathbf{x} = WW^{\mathsf{T}}\mathbf{x} = \left(\sum_{i=1}^{d} \mathbf{w}_{i}\mathbf{w}_{i}^{\mathsf{T}}\right)\mathbf{x} = \sum_{i=1}^{d} \mathbf{w}_{i}(\mathbf{w}^{\mathsf{T}}\mathbf{x})$ $= (\mathbf{w}_1^{\mathsf{T}}\mathbf{x})\mathbf{w}_1 + (\mathbf{w}_2^{\mathsf{T}}\mathbf{x})$ Dimensionality Thus, the new c reduction is achieved by dropping some of $\mathbf{y} = \begin{bmatrix} \mathbf{w}_1^\mathsf{T} \mathbf{x} \\ \mathbf{w}_2^\mathsf{T} \mathbf{x} \\ \vdots \\ \mathbf{w}_d^\mathsf{T} \mathbf{x} \end{bmatrix} \in \mathbb{R}^d$ the new coordinates.



Terminology

- $\square \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_d\}$ $\blacksquare Basis$
 - Directions
- $\square \mathbf{w}_i^\top \mathbf{x} = \langle \mathbf{w}_i, \mathbf{x} \rangle$
 - New coordinates
 - Projection of \mathbf{x} along the direction \mathbf{w}_i



Principal Component Analysis (PCA)

- Given a set of Data Points $\{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n\}$ where $\mathbf{x}_i \in \mathbb{R}^d$
- □ Finding a set of directions $\{\mathbf{w}_1, \mathbf{w}_2, ..., \mathbf{w}_k\}$ such that the variance of

$$\left\{ \mathbf{y}_{1} = \begin{bmatrix} \mathbf{w}_{1}^{\mathsf{T}} \mathbf{x}_{1} \\ \mathbf{w}_{2}^{\mathsf{T}} \mathbf{x}_{1} \\ \vdots \\ \mathbf{w}_{k}^{\mathsf{T}} \mathbf{x}_{1} \end{bmatrix}, \mathbf{y}_{2} = \begin{bmatrix} \mathbf{w}_{1}^{\mathsf{T}} \mathbf{x}_{2} \\ \mathbf{w}_{2}^{\mathsf{T}} \mathbf{x}_{2} \\ \vdots \\ \mathbf{w}_{k}^{\mathsf{T}} \mathbf{x}_{2} \end{bmatrix}, \cdots, \mathbf{y}_{n} = \begin{bmatrix} \mathbf{w}_{1}^{\mathsf{T}} \mathbf{x}_{n} \\ \mathbf{w}_{2}^{\mathsf{T}} \mathbf{x}_{n} \\ \vdots \\ \mathbf{w}_{k}^{\mathsf{T}} \mathbf{x}_{n} \end{bmatrix} \right\}$$

are maximized



For the purpose of dimensionality reduction, PCA only learns k directions.

of

nt Analysis (PCA)

la Points

for $\mathbf{x}_i \in \mathbb{R}^d$

\Box Finding set of directions { w_1, w_2, \dots, w_k } such that the variance

> PCA uses variances to measure the quality of new coordinates.



are maximized

 $y_1 =$

[⊤]W₁/



PCA—One-dimensional Case (1)

 $\square \text{ New Coordinates of } \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ $\mathbf{w}_1^\mathsf{T} \mathbf{x}_1, \mathbf{w}_1^\mathsf{T} \mathbf{x}_2, \dots, \mathbf{w}_1^\mathsf{T} \mathbf{x}_n$

□ Variance is

$$\frac{1}{n} \sum_{i=1}^{n} (\mathbf{w}_{1}^{\mathsf{T}} \mathbf{x}_{i} - \mu)^{2}$$

where $\mu = \frac{1}{n} \sum_{i=1}^{n} \mathbf{w}_{1}^{\mathsf{T}} \mathbf{x}_{i}$ is the mean of new coordinates



PCA—One-dimensional Case (2)

□ Let $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i$ be the mean vector □ Then, $\mu = \mathbf{w}_1^{\mathsf{T}} \bar{\mathbf{x}}$

$$\frac{1}{n} \sum_{i=1}^{n} (\mathbf{w}_{1}^{\mathsf{T}} \mathbf{x}_{i} - \mu)^{2} = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{w}_{1}^{\mathsf{T}} \mathbf{x}_{i} - \mathbf{w}_{1}^{\mathsf{T}} \bar{\mathbf{x}})^{2}$$
$$= \frac{1}{n} \sum_{i=1}^{n} \left(\mathbf{w}_{1}^{\mathsf{T}} (\mathbf{x}_{i} - \bar{\mathbf{x}}) \right)^{2} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{w}_{1}^{\mathsf{T}} (\mathbf{x}_{i} - \bar{\mathbf{x}}) (\mathbf{x}_{i} - \bar{\mathbf{x}})^{\mathsf{T}} \mathbf{w}_{1}$$
$$= \mathbf{w}_{1}^{\mathsf{T}} \left(\frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_{i} - \bar{\mathbf{x}}) (\mathbf{x}_{i} - \bar{\mathbf{x}})^{\mathsf{T}} \right) \mathbf{w}_{1}$$



PCA—One-dimensional Case (3)

□ The Optimization Problem of PCA

 $\max_{\mathbf{w} \in \mathbb{R}^d} \mathbf{w}^{\mathsf{T}} C \mathbf{w}$ s.t. $\|\mathbf{w}\|_2^2 = 1$ where $C = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^{\mathsf{T}}$ is the covariance matrix **ID** The Solution (Rayleigh Quotient) Lagrangian: $-\mathbf{w}^{\mathsf{T}} C \mathbf{w} + \lambda (\|\mathbf{w}\|_2^2 - \mathbf{1})$ Set the gradient of \mathbf{w} be zero

 $-2C\mathbf{w} + 2\lambda\mathbf{w} = 0 \Leftrightarrow C\mathbf{w} = \lambda\mathbf{w}$



PCA—One-dimensional Case (4)

- (\mathbf{w}, λ) is eigenvector and eigenvalue of *C*
- The objective becomes

$$\mathbf{w}^{\mathsf{T}} C \mathbf{w} = \lambda \mathbf{w}^{\mathsf{T}} \mathbf{w} = \lambda$$

- Thus, we select the largest eigenvector and eigenvalue of C
- □ The Algorithm
 - **1**. Calculate the mean vector $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}$
 - 2. Calculate the covariance matrix $C = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i \overline{\mathbf{x}}) (\mathbf{x}_i \overline{\mathbf{x}})^{\mathsf{T}}$
 - **3**. Calculate the largest eigenvector of *C*



Property of the Covariance Matrix

$$C = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i - \overline{\mathbf{x}}) (\mathbf{x}_i - \overline{\mathbf{x}})^{\mathsf{T}}$$

□ C is symmetric
□ C is positive semidefinite (PSD)
■ All the eigenvalues are non-negative
□ The rank of C is at most n - 1
■ Let $\overline{X} = [\mathbf{x}_1 - \overline{\mathbf{x}}, ..., \mathbf{x}_n - \overline{\mathbf{x}}] \in \mathbb{R}^{d \times n}$ rank(C) = rank($\overline{X}\overline{X}^{\top}$) = rank(\overline{X}) ≤ n - 1

It has at most n - 1 positive eigenvalues



PCA—k-dimensional Case (1)

□ The Optimization Problem of PCA

 $\max_{W \in \mathbb{R}^{d \times k}} \operatorname{trace}(W^{\top}CW)$ s.t. $W^{\top}W = I$ where $C = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_{i} - \overline{\mathbf{x}}) (\mathbf{x}_{i} - \overline{\mathbf{x}})^{\top}$

□ The Solution (Rayleigh Quotient)

• $W = [\mathbf{w}_1, ..., \mathbf{w}_k]$, where $\mathbf{w}_1, ..., \mathbf{w}_k$ are the k largest eigenvectors of C

Section 5.2.2.(6) of [Lütkepohl 1996]

Can also be defined in an incremental fashion



PCA—k-dimensional Case (2)

□ The Algorithm

- **1**. Calculate the mean vector $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}$
- 2. Calculate the covariance matrix $C = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i \overline{\mathbf{x}}) (\mathbf{x}_i \overline{\mathbf{x}})^{\mathsf{T}}$
- 3. Calculate the *k* largest eigenvectors of *C*

Eigenvalue

 \blacksquare λ_i is the variance of the *i*-th coordinate

Measure the quality of PCA
Captured
$$\frac{\sum_{i=1}^{k} \lambda_i}{\sum_{i=1}^{d} \lambda_i}$$
 $\frac{\sum_{i=k+1}^{d} \lambda_i}{\sum_{i=1}^{d} \lambda_i}$ Lost



An Example

Arrythmia data set from the UCI





Discussions of PCA

The Key Operation Eigendecomposition of *C* PCA can also be derived from the perspective of projection error minimization Section 12.1.2 of [Bishop 2007]

PCA is Linear Since

 $\mathbf{x} \in \mathbb{R}^d \to W^{\mathsf{T}} \mathbf{x} \in \mathbb{R}^k$

where $W = [\mathbf{w}_1, ..., \mathbf{w}_k] \in \mathbb{R}^{d \times k}$

PCA is Unsupervised



Singular Value Decomposition (SVD) [©]

□ SVD of $X = [\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ with $d \le n$

$$X = U\Sigma V^{\top} = \sum_{i=1}^{d} \sigma_i \mathbf{u}_i \mathbf{v}_i^{\top}$$

 $\begin{array}{l} \boldsymbol{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d] \in \mathbb{R}^{d \times d}, \ \boldsymbol{U}^\top \boldsymbol{U} = \boldsymbol{U} \boldsymbol{U}^\top = \boldsymbol{I} \\ \boldsymbol{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d] \in \mathbb{R}^{n \times d}, \ \boldsymbol{V}^\top \boldsymbol{V} = \boldsymbol{I} \\ \boldsymbol{\Sigma} = \operatorname{diag}(\sigma_1, \sigma_2, \dots, \sigma_d) \in \mathbb{R}^{d \times d}, \ \sigma_1 \ge \sigma_2 \ge \cdots \ge \sigma_d \ge 0 \end{array}$



Compact SVD

□ SVD of $X = [\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ with rank(r) < min(d, n) $X = U_r \Sigma_r V_r^\top = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ ■ $U_r = [\mathbf{u}_1, \mathbf{u}_2, ..., \mathbf{u}_r] \in \mathbb{R}^{d \times r}, \ U_r^\top U_r = I$ ■ $V_r = [\mathbf{v}_1, \mathbf{v}_2, ..., \mathbf{v}_r] \in \mathbb{R}^{n \times r}, \ V_r^\top V_r = I$

$$\Sigma_r = \operatorname{diag}(\sigma_1, \sigma_2, \dots, \sigma_r) \in \mathbb{R}^{r \times r}, \ \sigma_1 \ge \sigma_2 \ge \dots \ge \\ \sigma_r > 0$$



□ The Algorithm

- 1. Calculate the k largest left singular vectors $\mathbf{u}_1, \mathbf{u}_2, ..., \mathbf{u}_k$ of $X = [\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n]$
- The New Coordinates of x are

$$U_k^{\mathsf{T}} \mathbf{x} = \begin{bmatrix} \mathbf{u}_1^{\mathsf{T}} \mathbf{x} \\ \mathbf{u}_2^{\mathsf{T}} \mathbf{x} \\ \vdots \\ \mathbf{u}_k^{\mathsf{T}} \mathbf{x} \end{bmatrix} \in \mathbb{R}^k$$

■ $U_k = [\mathbf{u}_1, \mathbf{u}_2, ..., \mathbf{u}_k] \in \mathbb{R}^{d \times k}$ □ The New Coordinates of *X* is $U_k^\top X = U_k^\top U_r \Sigma_r V_r^\top = \Sigma_k V_k^\top$

SVD—A Energy-preserving Interpretation



□ The Optimization Problem of SVD

1-dimensional

 $\max_{\mathbf{w}\in\mathbb{R}^d} \quad \mathbf{w}^{\top}(XX^{\top})\mathbf{w}$ s.t. $\|\mathbf{w}\|_2^2 = 1$

k-dimensional

$$\max_{W \in \mathbb{R}^{d \times k}} trace(W^{\top}(XX^{\top})W)$$

s.t.
$$W^{\top}W = I$$

Left (right) singular vectors of X are the eigenvectors of XX^{T} ($X^{T}X$).



PCA by SVD

Old Algorithm

- **1**. Calculate the mean vector $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}$
- 2. Calculate the covariance matrix $C = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i \bar{\mathbf{x}}) (\mathbf{x}_i \bar{\mathbf{x}})^{\mathsf{T}}$
- 3. Calculate the *k*-largest eigenvectors of *C*

New Algorithm

- 1. Calculate the mean vector $\overline{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}$
- 2. Calculate the *k* largest left singular vectors of $\overline{X} = [\mathbf{x}_1 \overline{\mathbf{x}}, ..., \mathbf{x}_n \overline{\mathbf{x}}]$



 $n^{L_{i=1}} \mathbf{x}_{i}$

PCA by SVD

□ Old Algorithm

- **1**. Calculate the mean vector $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}$
- 2. Calculate the covariance matrix $C = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i)$ 3. Calculate the covariance matrix $C = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i)$ PCA is equivalent to SVD if $X = \overline{X}$, that is, if data are zero-mean.

1. Calculate

2. Calculate the pargest left singular vectors of $\overline{X} = [\mathbf{x}_1 - \overline{\mathbf{x}}, ..., \mathbf{x}_n - \overline{\mathbf{x}}]$



Outline

- Introduction
- Feature Extraction and Portability
- Data Cleaning
- Data Reduction and Transformation
 - Sampling
 - Feature Subset Selection
 - Dimensionality Reduction with Axis Rotation
 - Dimensionality Reduction with Type Transformation
- □ Summary

Dimensionality Reduction with Type Transformation



□ Time Series to Multidimensional

- Can also be viewed as a rotation of an axis system
- Haar wavelet transform

Weighted graphs to multidimensional

- Multidimensional Scaling (MDS)
 - Edge represents distance
- Spectral Transformations
 - Edge represents similarity



Haar Wavelet Transform (1)

- A New Basis for time series data
 - Each element basis is a time series (wavelets)
 - Coefficients can be calculated efficiently
 - Coefficients have nice interpretations





Haar Wavelet Transform (2)

Given a Time Series **t** with length *d*

 $\mathbf{t} = \alpha^1 \mathbf{w}_1 + \alpha^2 \mathbf{w}_2 + \dots + \alpha^d \mathbf{w}_d$

where $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_d$ are wavelets, and they are orthogonal to each other

Normalization

$$\mathbf{t} = \alpha^{1} \|\mathbf{w}_{1}\|_{2} \frac{\mathbf{w}_{1}}{\|\mathbf{w}_{1}\|_{2}} + \alpha^{2} \|\mathbf{w}_{2}\|_{2} \frac{\mathbf{w}_{2}}{\|\mathbf{w}_{2}\|_{2}} + \dots + \alpha^{d} \|\mathbf{w}_{d}\|_{2} \frac{\mathbf{w}_{d}}{\|\mathbf{w}_{d}\|_{2}}$$
$$= \frac{\mathbf{w}_{1}}{\|\mathbf{w}_{1}\|_{2}}, \frac{\mathbf{w}_{2}}{\|\mathbf{w}_{2}\|_{2}}, \dots, \frac{\mathbf{w}_{d}}{\|\mathbf{w}_{d}\|_{2}} \text{ are orthonormal to}$$
each other



Haar Wavelet Transform (3)

□ The New Coordinates

$$\mathbf{y} = \begin{bmatrix} \alpha^1 \| \mathbf{w}_1 \|_2 \\ \alpha^2 \| \mathbf{w}_2 \|_2 \\ \vdots \\ \alpha^d \| \mathbf{w}_d \|_2 \end{bmatrix} \in \mathbb{R}^d$$

Dimensionality Reduction

$$Y = \begin{bmatrix} \alpha_1^1 \| \mathbf{w}_1 \|_2 & \alpha_2^1 \| \mathbf{w}_1 \|_2 \cdots & \alpha_n^1 \| \mathbf{w}_1 \|_2 \\ \alpha_1^2 \| \mathbf{w}_2 \|_2 & \alpha_2^2 \| \mathbf{w}_2 \|_2 \cdots & \alpha_n^2 \| \mathbf{w}_2 \|_2 \\ \vdots & \vdots & \vdots & \vdots \\ \alpha_1^d \| \mathbf{w}_d \|_2 & \alpha_2^d \| \mathbf{w}_d \|_2 \cdots & \alpha_n^d \| \mathbf{w}_d \|_2 \end{bmatrix} \in \mathbb{R}^{d \times n}$$

- Feature Selection, PCA, SVD
- Sparse Representation



Multidimensional Scaling (MDS)

Input

- A graph with *n* nodes

Output

A set of coordinates that fits the distance

Metric MDS

$$\min_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^k} \quad \sum_{i, j: i < j} \left(\left\| \mathbf{x}_i - \mathbf{x}_j \right\|_2 - \delta_{ij} \right)^2$$

Assume the specified distance matrix is Euclidean

□ The Algorithm

1. Calculate the dot-product matrix

$$S = -\frac{1}{2} \left(I - \frac{\mathbf{1}\mathbf{1}^{\mathsf{T}}}{n} \right) \Delta \left(I - \frac{\mathbf{1}\mathbf{1}^{\mathsf{T}}}{n} \right)$$

2. Eigen decompose *S*

$$S = U\Lambda U^{\top} = \sum_{i=1}^{n} \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top}$$

3. The new coordinates are

$$U_k \Lambda_k^{-1/2} \in \mathbb{R}^{n \times k}$$

 $U_k = [\mathbf{u}_1, \dots, \mathbf{u}_k] \in \mathbb{R}^{n \times k}, \ \Lambda_k = \operatorname{diag}(\lambda_1, \dots, \lambda_k) \in \mathbb{R}^{k \times k}$

Assume the specified distance matrix is Euclidean

The Algorithm 1. Calculate the dot-product fatrix Metric MDS is 2. Èn equivalent to PCA, if the distance matrix is Euclidean. 3. <u>í</u>re $U_k \Lambda_k^{-1/2} \in \mathbb{R}^{n \times k}$ $U_k = [\mathbf{u}_1, \dots, \mathbf{u}_k] \in \mathbb{R}^{n \times k}, \ \Lambda_k = \operatorname{diag}(\lambda_1, \dots, \lambda_k) \in \mathbb{R}^{k \times k}$



Spectral Transformations (1)

Input

- A graph with *n* nodes
- w_{ij} = w_{ji} be the similarity between nodes i and j

Output

A set of coordinates that preserves the similarity

□ The Objective

$$\min_{\mathbf{x}_{1},\mathbf{x}_{2},...,\mathbf{x}_{n}\in\mathbb{R}^{k}} \quad \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n}w_{ij}\|\mathbf{x}_{i}-\mathbf{x}_{j}\|_{2}^{2}$$



Spectral Transformations (2)

- □ The Optimization Problem
 - $\min_{\substack{Y \in \mathbb{R}^{n \times k} \\ s.t.}} \text{trace}(Y^{\top}LY)$

$$Y = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times k}, \ L = D - W$$

• *D* is a diagonal matrix with $D_{ii} = \sum_{j=1}^{n} w_{ij}$

Generalized Eigenproblem

$$L\mathbf{y} = \lambda D\mathbf{y}$$

- □ The Solution [Luxburg 2007]
 - $Y = [\mathbf{y}_1, ..., \mathbf{y}_k] \in \mathbb{R}^{n \times k}$, where \mathbf{y}_i is the *i*-th smallest eigenvector



Outline

- Introduction
- □ Feature Extraction and Portability
- Data Cleaning
- Data Reduction and Transformation
 - Sampling
 - Feature Subset Selection
 - Dimensionality Reduction with Axis Rotation
 - Dimensionality Reduction with Type Transformation





Summary

- □ Feature Extraction and Portability
- Data Cleaning
- Data Reduction by Sampling
- Dimensionality Reduction with Axis Rotation
 - PCA, SVD
- Dimensionality Reduction with Type Transformation
 - Haar Wavelet Transform, MDS, Spectral Transformation


Reference

Bishop 2007

- Christopher Bishop. Pattern Recognition and Machine Learning. Springer, 2007.
- Lütkepohl 1996
 - H. Lütkepohl. Handbook of Matrices. Wiley, 1996.
- Luxburg 2007
 - Ulrike von Luxburg. A tutorial on spectral clustering. Statistics and Computing, 17(4): 395-416, 2007.