


Data Mining:

Concepts and Techniques

(3rd ed.)

— Chapter 3 —

Chapter 3: Data Preprocessing

- Data Preprocessing: An Overview
 - Data Quality
 - Major Tasks in Data Preprocessing
- Data Cleaning
- Data Integration
- Data Reduction 
- Data Transformation and Data Discretization
- Summary

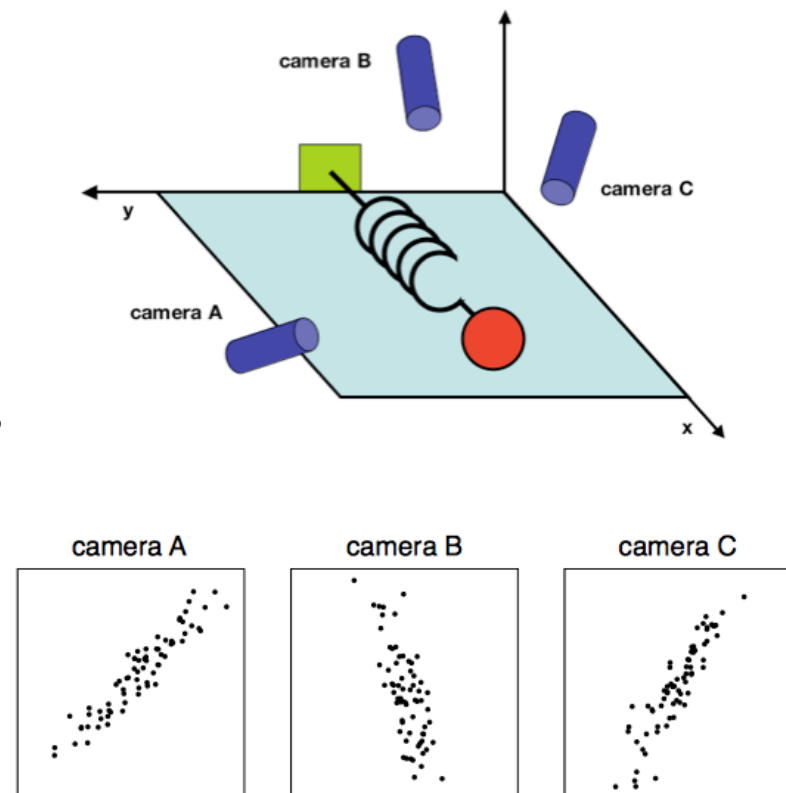
Focused Topic: **Principal Component Analysis**

A dimensionality reduction method

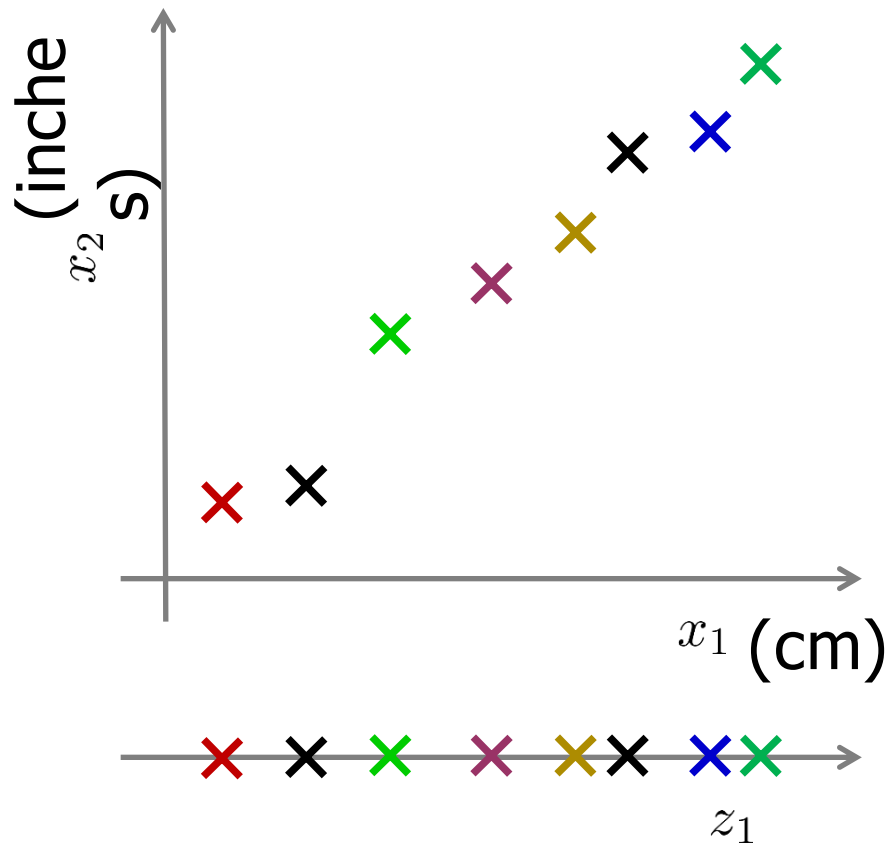
Motivation: Data Compression for Feature Selection and Visualiation

Data Compression Example 1

- Ball of mass m attached to massless, frictionless spring
- Ball moved away from equilibrium results in spring
- oscillating indefinitely along x-axis
- Three cameras ~ three dimensions
- However, all dynamics can be a function of only a single variable x



Data Compression Example 2



Reduce data from
2D to 1D

$$x^{(1)} \rightarrow z^{(1)}$$

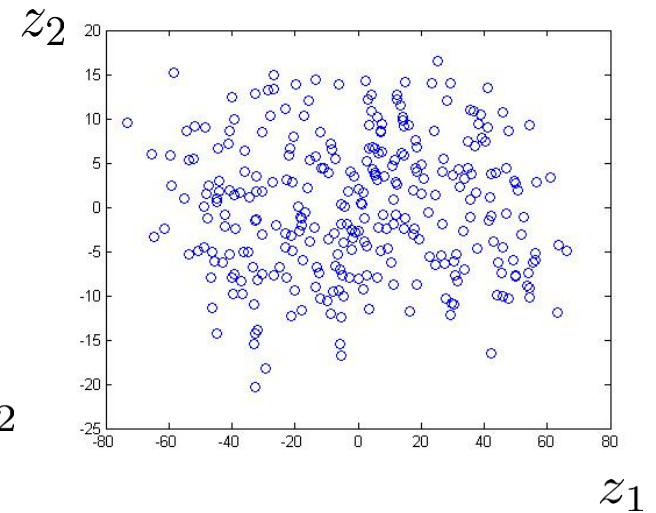
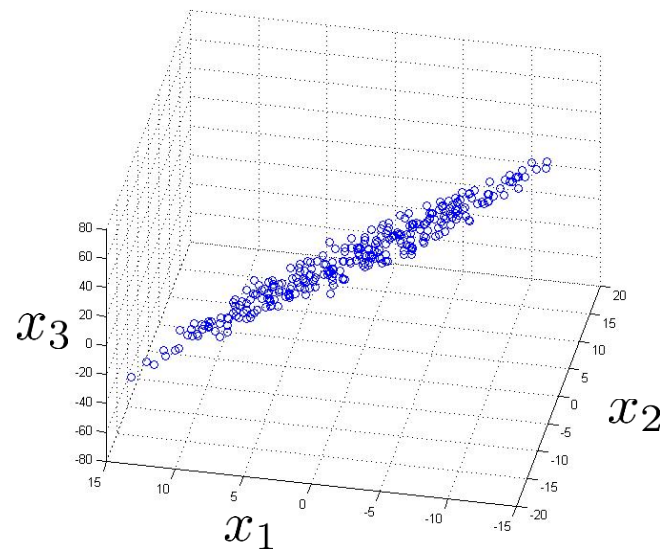
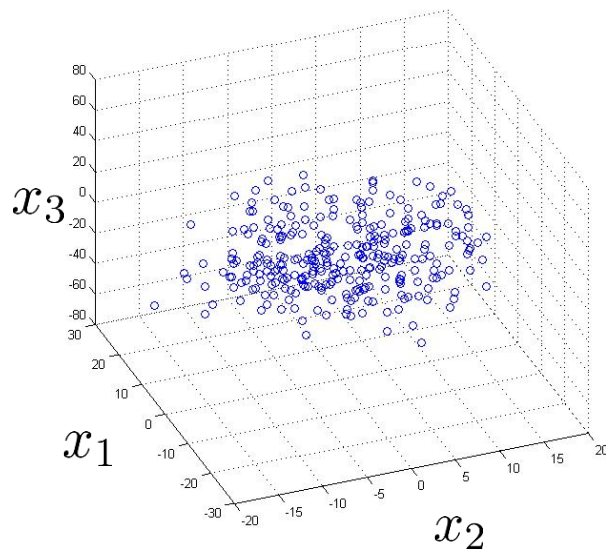
$$x^{(2)} \rightarrow z^{(2)}$$

\vdots

$$x^{(m)} \rightarrow z^{(m)}$$

Data Compression Example 3

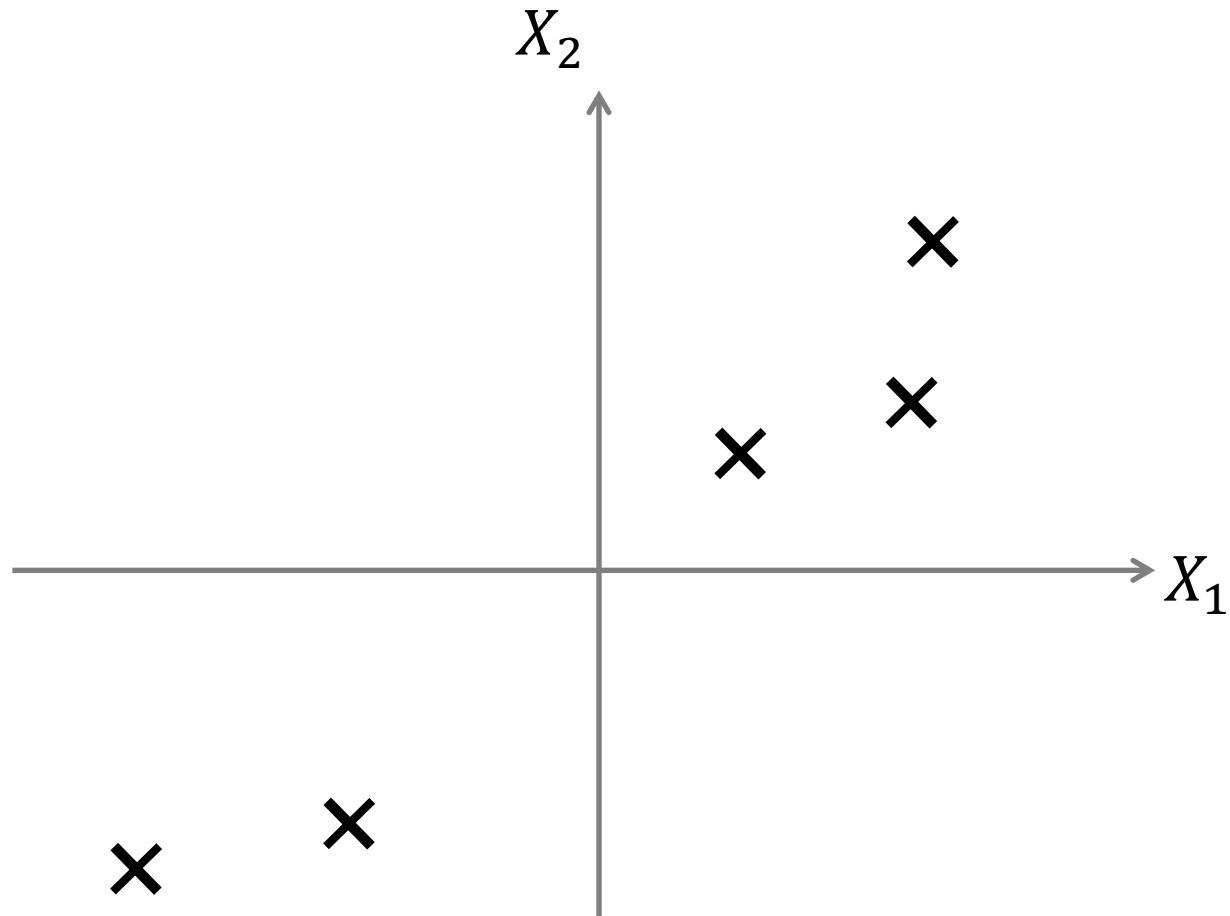
Reduce data from 3D to 2D



PCA PROBLEM FORMULATION

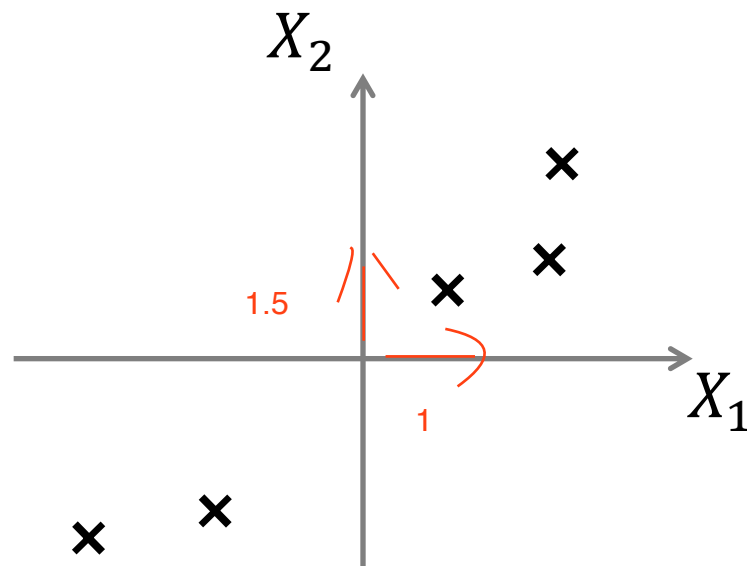
Principal Component Analysis (PCA)

Problem formulation

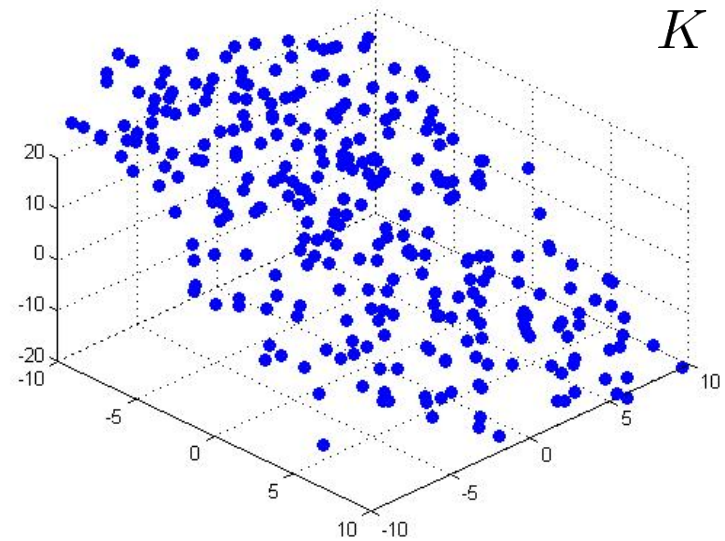


Principal Component Analysis (PCA)

Problem formulation



$3D \rightarrow 2D$
 $K = 2$



Orig Base

Reduce from 2-dimension to 1-dimension: Find a direction (a vector $P_1 \in R^2$) onto which to project the data so as to minimize the projection error.

Reduce from n-dimension to k-dimension: Find k vectors P_1, P_2, \dots, P_k onto which to project the data, so as to minimize the projection error.

In summary, the goal of PCA is...

- Compute the most meaningful basis to re-express a noisy data set
- Hope that this new basis will filter out the noise and reveal hidden structure
- In the toy example:
 - Determine that the dynamics are along the x-axis

PCA ALGORITHM

Naïve Basis: formed directly from the method used to gather data

- At each point in time, record 2 coordinates of ball position in each of the 3 images
- After 10 minutes at 120Hz, we have $10 \times 60 \times 120 = 7200$ 6-dimensional vectors
- These vectors can be represented in arbitrary coordinate systems

6×7200

$$\vec{X} = \begin{bmatrix} x_A \\ y_A \\ x_B \\ y_B \\ x_C \\ y_C \end{bmatrix}$$

PCA: Changing basis to express the data better

- PCA: Is there another basis, which is a linear combination of the original basis, that best re-expresses our data set?
- Assumption: linearity
 - Restricts set of potential bases
 - Simplifies the characterization of a complex system

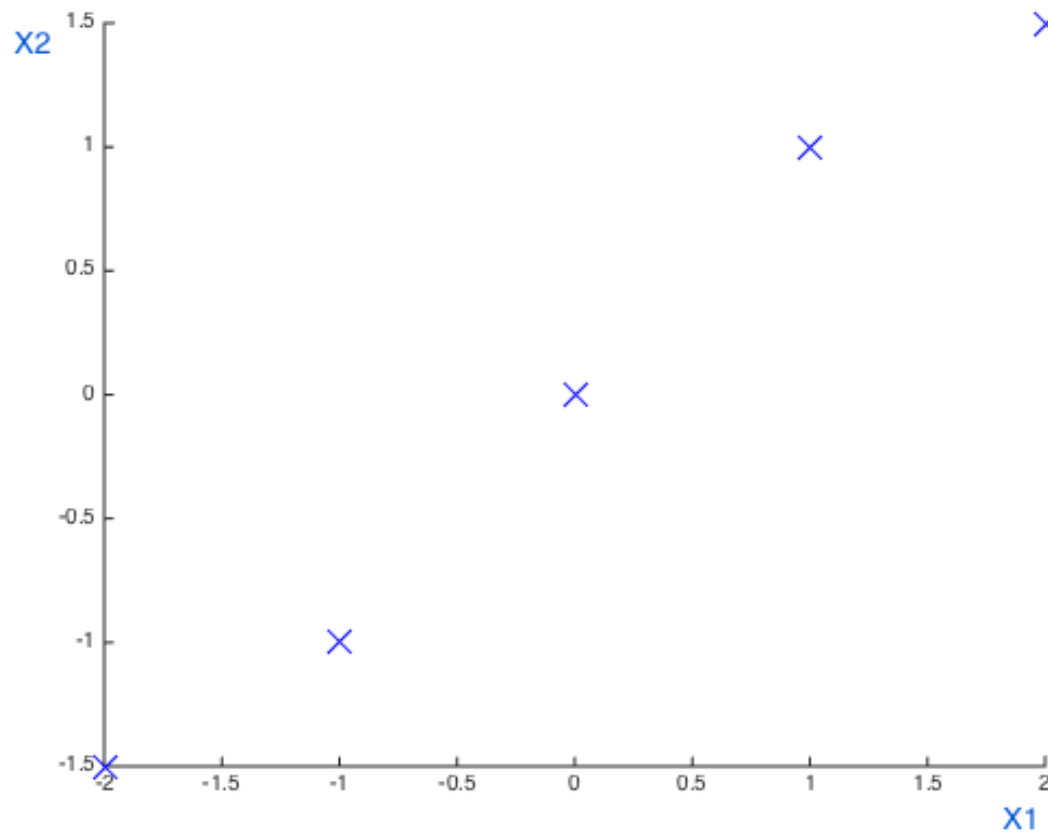
Change of basis expressed in linear algebra

- X is original data ($m \times n$, $m=6$, $n=7200$)
- Let Y be another $m \times n$ matrix such that $Y=PX$ (or, $k \times n$ if first k dimensions).
- P is a matrix that transforms X into Y
 - What is the size of P?
 - Geometrically it is a rotation and stretch
 - The rows of P $\{p_1, \dots, p_m\}$ are the new basis vectors for the columns of X
 - Called principal components of X
 - Each element of y_i is a dot product of x_i with the corresponding row of P (a projection of x_i onto p_j)

$$PX = \begin{bmatrix} p_1 \\ \vdots \\ p_m \end{bmatrix} \begin{bmatrix} x_1 & \cdots & x_n \end{bmatrix}$$
$$Y = \begin{bmatrix} p_1 \cdot x_1 & \cdots & p_1 \cdot x_n \\ \vdots & \ddots & \vdots \\ p_m \cdot x_1 & \cdots & p_m \cdot x_n \end{bmatrix} \quad y_i = \begin{bmatrix} p_1 \cdot x_i \\ \vdots \\ p_m \cdot x_i \end{bmatrix}$$

Running example for two measurements

	x1	x2	x3	x4	x5
X1	-2	-1	0	1	2
X2	-1.5	-1	0	1	1.5



Running example: Changing basis

- Assume that by some way we know the new basis is a 2-d space, i.e., specified by

- Vector $p1=(0.9, 0.45)$
- Vector $p2=(-0.45, 0.9)$

- Matrix P: (P is orthonormal)

$$P = \begin{bmatrix} 0.9 & 0.45 \\ -0.45 & 0.9 \end{bmatrix}$$

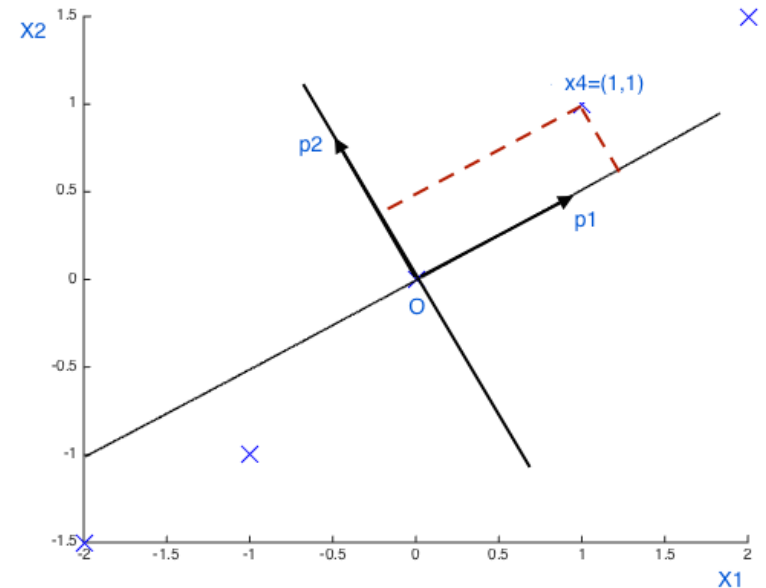
- Change basis for $x4 (1,1)$: two dot products:
 $A \cdot v1$ and $A \cdot v2$ equivalent to two projections

$$x4' = P * x4 = \begin{bmatrix} 0.9 & 0.45 \\ -0.45 & 0.9 \end{bmatrix} * \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1.35 \\ 0.45 \end{bmatrix}$$

- Change basis for entire dataset X:

$$Y = PX = \begin{bmatrix} 0.9 & 0.45 \\ -0.45 & 0.9 \end{bmatrix} * \begin{bmatrix} -2 & -1 & 0 & 1 & 2 \\ -1.5 & -1 & 0 & 1 & 1.5 \end{bmatrix}$$

$$Y = \begin{bmatrix} -2.48 & -1.35 & 0 & 1.35 & 2.48 \\ -0.45 & -0.45 & 0 & 0.45 & 0.45 \end{bmatrix}$$



	x1	x2	x3	x4	x5
X1	-2	-1	0	1	2
X2	-1.5	-1	0	1	1.5

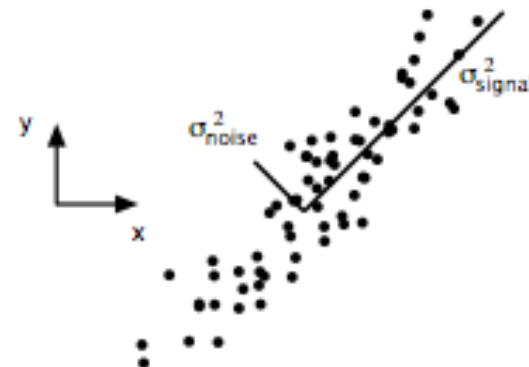
Finding an appropriate change of basis by minimizing noise and redundancy

- Finding change of basis = Finding matrix P
- What is the best way to re-express X ? What features would we like Y to exhibit?
- If we call X “garbled data”, garbling in a linear system can refer to two things:
 - Noise
 - Redundancy

Signal and noise can be measured by variance

- Measurement noise in any reasonable data set should be low.
- In the toy example: ball travels in straight line
 - Any deviation must be noise.
 - The variance due to the signal and noise are indicated by each line in the diagram.
- Assumption: directions with largest variances in our measurement space contain the dynamics of interest.
 - Signal-to-Noise Ratio (SNR)

$$SNR = \frac{\sigma_{signal}^2}{\sigma_{noise}^2}$$



1. → Goal 1: Maximize variances of new dimensions.

Redundancy happens when we have little info about the dimensions of interest when collecting data

- Is it necessary to record 2 variables for the ball-spring system?
- Is it necessary to use 3 cameras?
- Covariance $\text{Cov}(X_1, X_2)$ (or correlation) reveals redundancy between X_1 and X_2 .
- Redundancy should be removed.
 - the new basis can use smaller number of dimensions than the naïve basis does.

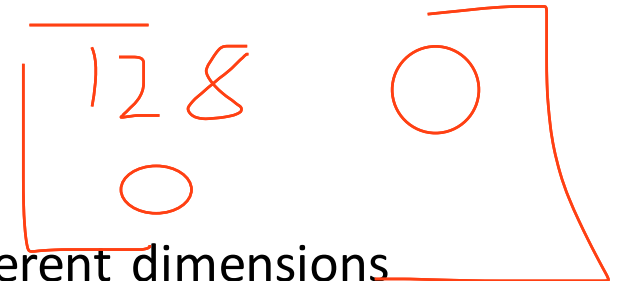


1. → Goal 2: **Minimize co-variance** between new dimensions.

Covariance Matrix tells us about signals and redundancies

$$\Sigma = \begin{bmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & E[(X_1 - \mu_1)(X_2 - \mu_2)] & \cdots & E[(X_1 - \mu_1)(X_n - \mu_n)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[(X_2 - \mu_2)(X_2 - \mu_2)] & \cdots & E[(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - \mu_n)(X_1 - \mu_1)] & E[(X_n - \mu_n)(X_2 - \mu_2)] & \cdots & E[(X_n - \mu_n)(X_n - \mu_n)] \end{bmatrix}$$

- Diagonal elements of C_x : variances of dimensions
 - Large \rightarrow interesting dynamics
 - Small \rightarrow noise
- Off-diagonal elements of C_x : covariances of two different dimensions
 - Large \rightarrow high redundancy
 - Small \rightarrow low redundancy
- **Goal 1+2: Covariance matrix of the new space C_y should be diagonal!**



If all dimensions have zero-mean, we have a simpler formula for covariance matrix

$$\mathbf{C}_X = \begin{bmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & E[(X_1 - \mu_1)(X_2 - \mu_2)] & \cdots & E[(X_1 - \mu_1)(X_n - \mu_n)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[(X_2 - \mu_2)(X_2 - \mu_2)] & \cdots & E[(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - \mu_n)(X_1 - \mu_1)] & E[(X_n - \mu_n)(X_2 - \mu_2)] & \cdots & E[(X_n - \mu_n)(X_n - \mu_n)] \end{bmatrix}$$

- If all dimensions have zero-mean, covariance matrix $\mathbf{C}_X = \frac{1}{n} \mathbf{X} \mathbf{X}^T$
- In the toy example, because both measures are **zero-mean normalized**:

$$C_x = \frac{1}{5} * \begin{bmatrix} -2 & -1 & 0 & 1 & 2 \\ -1.5 & -1 & 0 & 1 & 1.5 \end{bmatrix} * \begin{bmatrix} -2 & -1.5 \\ -1 & -1 \\ 0 & 0 \\ 1 & 1 \\ 2 & 1.5 \end{bmatrix} = \begin{bmatrix} 2.0 & 1.6 \\ 1.6 & 1.3 \end{bmatrix}$$

Signal

Signal

- **First:** Zero-mean normalize each dimension, to simplify computation.

Solving PCA (1): Using eigenvector decomposition

- The objective

Find some orthonormal matrix \mathbf{P} in $\mathbf{Y} = \mathbf{P}\mathbf{X}$ such that $\mathbf{C}_Y \equiv \frac{1}{n}\mathbf{Y}\mathbf{Y}^T$ is a diagonal matrix. The rows of \mathbf{P} are the *principal components* of \mathbf{X} .

- New assumption: \mathbf{P} must be orthonormal
 - There are many possible \mathbf{P} in the search space
 - But if \mathbf{P} is orthonormal, there is an efficient solution to find \mathbf{P} :
 - Eigenvector decomposition

Running example: C_Y is not always diagonal

- In the running example, we chose an orthonormal P :

- p_1 and p_2 are unit vectors: $|p_1| = |p_2| = 1$

- p_1 and p_2 are orthogonal:

$$p_1 \cdot p_2 = (0.9, 0.45) \cdot (-0.45, 0.9) = 0$$

- However, C_Y is not diagonal

$$C_y = \frac{1}{5} * Y * Y^T = \begin{bmatrix} 3.18 & 0.69 \\ 0.69 & 0.16 \end{bmatrix}$$

- How to find P to make C_Y diagonal?

Solving PCA (2): Rewrite \mathbf{C}_Y in terms of \mathbf{P} and \mathbf{C}_X

$$\begin{aligned}\mathbf{C}_Y &= \frac{1}{n} \mathbf{Y} \mathbf{Y}^T \\ &= \frac{1}{n} (\mathbf{P} \mathbf{X}) (\mathbf{P} \mathbf{X})^T \\ &= \frac{1}{n} \mathbf{P} \mathbf{X} \mathbf{X}^T \mathbf{P}^T \\ &= \mathbf{P} \left(\frac{1}{n} \mathbf{X} \mathbf{X}^T \right) \mathbf{P}^T \\ \mathbf{C}_Y &= \mathbf{P} \mathbf{C}_X \mathbf{P}^T\end{aligned}$$

Solving PCA (3): As we want C_Y diagonal, put a diagonal matrix to the formula of C_Y by eigenvector decomposition

Covariance matrix C_X is always symmetric

E: Eigenvectors of C_X

D: Diagonal matrix

(Theorem 4) When C_X is symmetric, $C_X = E D E^T$

$$\begin{aligned} C_Y &= P C_X P^T \\ &= P (E D E^T) P^T \end{aligned}$$

Solving PCA (4): Canceling out anything other than the diagonal matrix in C_Y by selecting P as a matrix where each row is an eigenvector of C_X

Select $E = P^T$

$$P = E^T$$

$$\begin{aligned} C_Y &= P C_X P^T \\ &= P(E D E^T) P^T \\ &= P(P^T D P) P^T \\ &= (P P^T) D (P P^T) \end{aligned}$$

$$E = P^T = P^{-1}$$

(Theorem 1) The inverse of an orthogonal matrix is its transpose. As P is an orthogonal matrix, C_Y turns out to be diagonal

$$\begin{aligned} C_Y &= (P P^{-1}) D (P P^{-1}) \\ &= D \end{aligned}$$

PCA by eigenvalue decomposition in the running example

X is already zero-mean

Eigenvectors of C_X (using Matlab): $[0.627 \ -0.779]$ and $[-0.779 \ -0.627]$

Each row of P is an eigenvector

$$P = \begin{bmatrix} 0.627 & -0.779 \\ -0.779 & -0.627 \end{bmatrix}$$

$$C_X * e = \lambda * e$$

Project X to the new space to obtain Y

X4

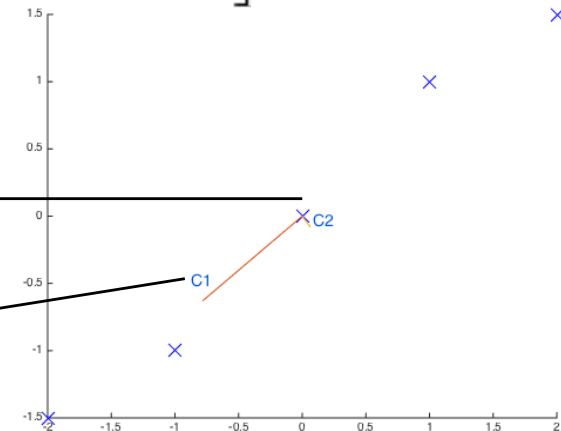
$$Y = PX = \begin{bmatrix} 0.627 & -0.779 \\ -0.779 & -0.627 \end{bmatrix} * \begin{bmatrix} -2 & -1 & 0 & 1 & 2 \\ -1.5 & -1 & 0 & 1 & 1.5 \end{bmatrix}$$

$$Y = \begin{bmatrix} -0.0855 & 0.1520 & 0 & -0.1520 & 0.0855 \\ 2.4985 & 1.4060 & 0 & -1.4060 & -2.4985 \end{bmatrix}$$

Guess what, C_Y is diagonal!

$$C_Y = \frac{1}{5} * Y * Y^T = \begin{bmatrix} 0.012 & 0 \\ 0 & 3.288 \end{bmatrix}$$

Y4



If we want to use only the most important principal component

$$Y = \begin{bmatrix} -0.779 & -0.627 \end{bmatrix} * \begin{bmatrix} -2 & -1 & 0 & 1 & 2 \\ -1.5 & -1 & 0 & 1 & 1.5 \end{bmatrix} = \begin{bmatrix} 2.4985 & 1.4060 & 0 & -1.4060 & -2.4985 \end{bmatrix}$$

Summary:

PCA by Eigenvector Decomposition

- Make \mathbf{X} a mean-normalized data matrix with m dims x n points.
- Find covariance matrix $\mathbf{C}_\mathbf{X} = \frac{1}{n} \mathbf{X}\mathbf{X}^T$, an $m \times m$ symmetric matrix.
- For $\mathbf{C}_\mathbf{X}$: $\mathbf{E} = [\mathbf{e}_1, \dots, \mathbf{e}_m]$ (eigenvectors) and \mathbf{D} ($\mathbf{C}_\mathbf{Y} = \mathbf{D}$, new covar), sorted by variance each -- Use eigenvector decomposition or SVD.
- Find $\mathbf{P} = \begin{bmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_m^T \end{bmatrix}$ as principle components.
- Select \mathbf{P}_k as first k principle components.
- New data $\mathbf{Y} = \mathbf{P}_k \mathbf{X}$.
- What k to choose? Large enough to preserve sufficient covariance.

Singular Value Decomposition (SVD)

- An alternative to finding eigenvectors for the covariance matrix
- More numerically stable than directly finding eigenvectors
- Widely used in practice for PCA

Chapter 3: Data Preprocessing

- Data Preprocessing: An Overview
 - Data Quality
 - Major Tasks in Data Preprocessing
- Data Cleaning
- Data Integration
- Data Reduction
- Data Transformation and Data Discretization
- Summary



Data Transformation

- A function that maps the entire set of values of a given attribute to a new set of replacement values s.t. each old value can be identified with one of the new values
- Methods
 - Smoothing: Remove noise from data
 - Attribute/feature construction
 - New attributes constructed from the given ones
 - Aggregation: Summarization, data cube construction
 - Normalization: Scaled to fall within a smaller, specified range
 - min-max normalization
 - z-score normalization
 - normalization by decimal scaling
 - Discretization: Concept hierarchy climbing

Normalization

- **Min-max normalization:** to $[new_min_A, new_max_A]$

$$v' = \frac{v - min_A}{max_A - min_A} (new_max_A - new_min_A) + new_min_A$$

- Ex. Let income range \$12,000 to \$98,000 normalized to [0.0, 1.0].
Then \$73,000 is mapped to $\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$

- **Z-score normalization** (μ : mean, σ : standard deviation):

$$v' = \frac{v - \mu_A}{\sigma_A}$$

- Ex. Let $\mu = 54,000$, $\sigma = 16,000$. Then $\frac{73,600 - 54,000}{16,000} = 1.225$

Discretization

- Three types of attributes
 - Nominal—values from an unordered set, e.g., color, profession
 - Ordinal—values from an ordered set, e.g., military or academic rank
 - Numeric—real numbers, e.g., integer or real numbers
- Discretization: Divide the range of a continuous attribute into intervals
 - Interval labels can then be used to replace actual data values
 - Reduce data size by discretization
 - Supervised vs. unsupervised
 - Split (top-down) vs. merge (bottom-up)
 - Discretization can be performed recursively on an attribute
 - Prepare for further analysis, e.g., classification

Data Discretization Methods

- Typical methods: All the methods can be applied recursively
 - Binning
 - Top-down split, unsupervised
 - Histogram analysis
 - Top-down split, unsupervised
 - Clustering analysis (unsupervised, top-down split or bottom-up merge)
 - Decision-tree analysis (supervised, top-down split)
 - Correlation (e.g., χ^2) analysis (unsupervised, bottom-up merge)

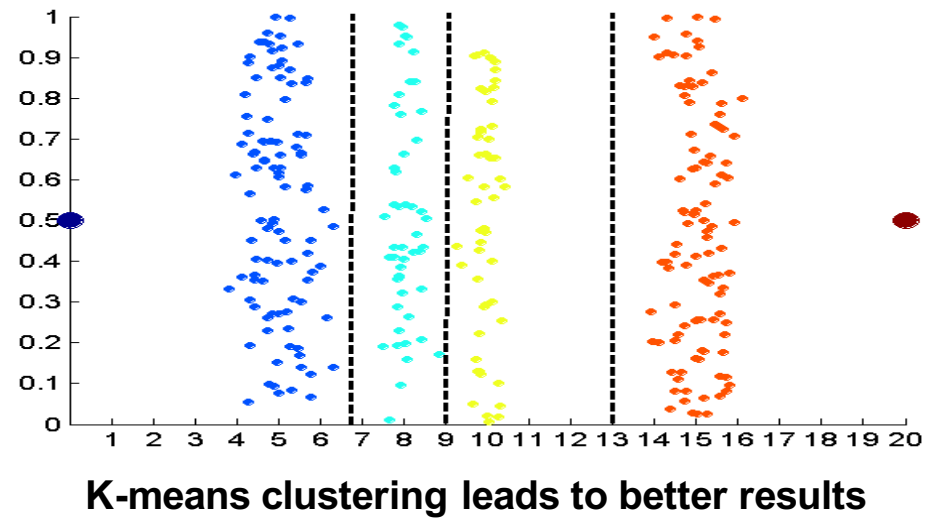
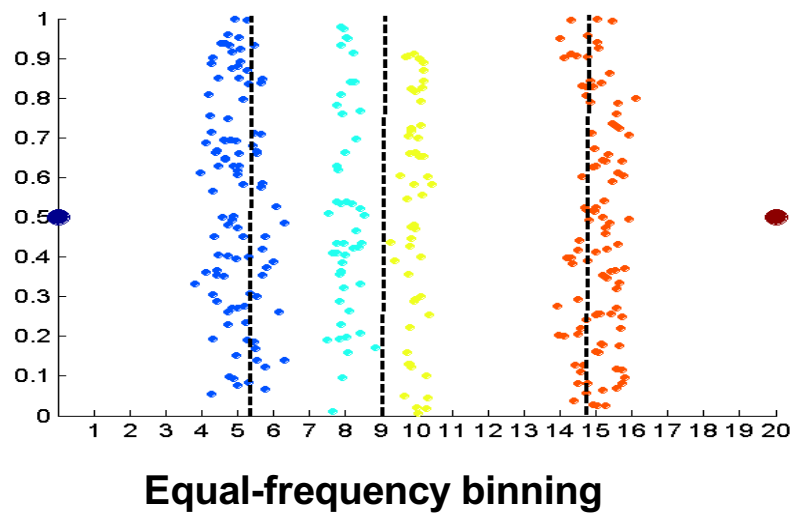
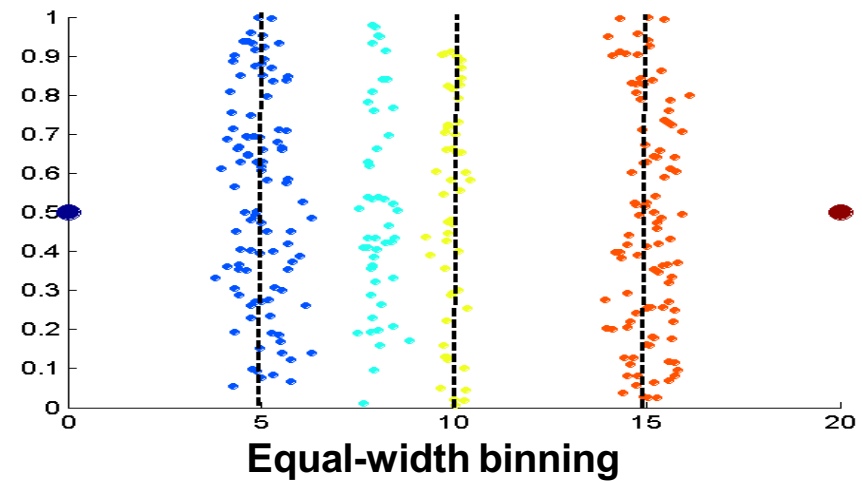
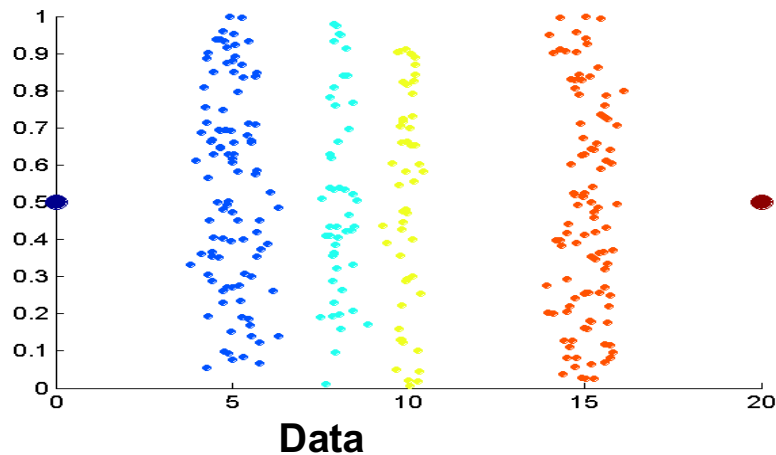
Simple Discretization: Binning

- **Equal-width** (distance) partitioning
 - Divides the range into N intervals of equal size: uniform grid
 - if A and B are the lowest and highest values of the attribute, the width of intervals will be: $W = (B - A)/N$.
 - The most straightforward, but outliers may dominate presentation
 - Skewed data is not handled well
- **Equal-depth** (frequency) partitioning
 - Divides the range into N intervals, each containing approximately same number of samples

Binning Methods for Data Smoothing

- Sorted data for price (in dollars): 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34
- * Partition into equal-frequency (**equi-depth**) bins:
 - Bin 1: 4, 8, 9, 15
 - Bin 2: 21, 21, 24, 25
 - Bin 3: 26, 28, 29, 34
- * Smoothing by **bin means**:
 - Bin 1: 9, 9, 9, 9
 - Bin 2: 23, 23, 23, 23
 - Bin 3: 29, 29, 29, 29
- * Smoothing by **bin boundaries**:
 - Bin 1: 4, 4, 4, 15
 - Bin 2: 21, 21, 25, 25
 - Bin 3: 26, 26, 26, 34

Discretization Without Using Class Labels (Binning vs. Clustering)




Discretization by Classification & Correlation Analysis

- Classification (e.g., decision tree analysis)
 - Supervised: Given class labels, e.g., cancerous vs. benign
 - Using *entropy* to determine split point (discretization point)
 - Top-down, recursive split
 - Details to be covered in Chapter “Classification”
- Correlation analysis (e.g., Chi-merge: χ^2 -based discretization)
 - Supervised: use class information
 - Bottom-up merge: find the best neighboring intervals (those having similar distributions of classes, i.e., low χ^2 values) to merge
 - Merge performed recursively, until a predefined stopping condition

Concept Hierarchy Generation for Nominal Data

- **Concept hierarchy** organizes concepts (i.e., attribute values) hierarchically and is usually associated with each dimension in a data warehouse
- Specification of a partial/total ordering of attributes explicitly at the schema level by users or experts
 - *street < city < state < country*
- Specification of a hierarchy for a set of values by explicit data grouping
 - {Urbana, Champaign, Chicago} < Illinois
- Specification of only a partial set of attributes
 - E.g., only *street < city*, not others
- Automatic generation of hierarchies (or attribute levels) by the analysis of the number of distinct values
 - Country 15 distinct values
 - State 365 distinct values
 - City 3567 distinct values
 - Street 674339 distinct values

Chapter 3: Data Preprocessing

- Data Preprocessing: An Overview
 - Data Quality
 - Major Tasks in Data Preprocessing
- Data Cleaning
- Data Integration
- Data Reduction
- Data Transformation and Data Discretization
- Summary 

Summary

- **Data quality:** accuracy, completeness, consistency, timeliness, believability, interpretability
- **Data cleaning:** e.g. missing/noisy values, outliers
- **Data integration** from multiple sources:
 - Entity identification problem; Remove redundancies; Detect inconsistencies
- **Data reduction**
 - Dimensionality reduction; Numerosity reduction; Data compression
- **Data transformation and data discretization**
 - Normalization; Concept hierarchy generation