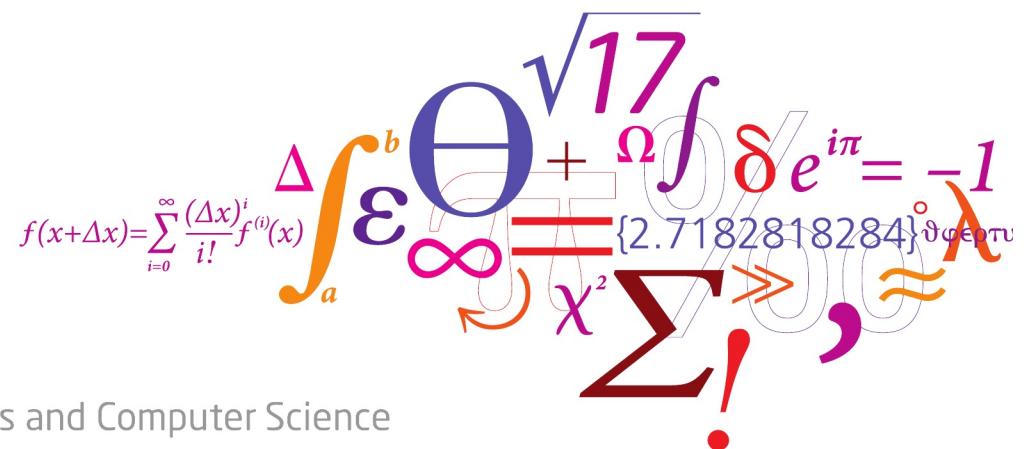


02450: Introduction to Machine Learning and Data Mining

Data, feature extraction and PCA

Tue Herlau

DTU Compute, Technical University of Denmark (DTU)

$$f(x+\Delta x) = \sum_{i=0}^{\infty} \frac{(\Delta x)^i}{i!} f^{(i)}(x)$$


DTU Compute

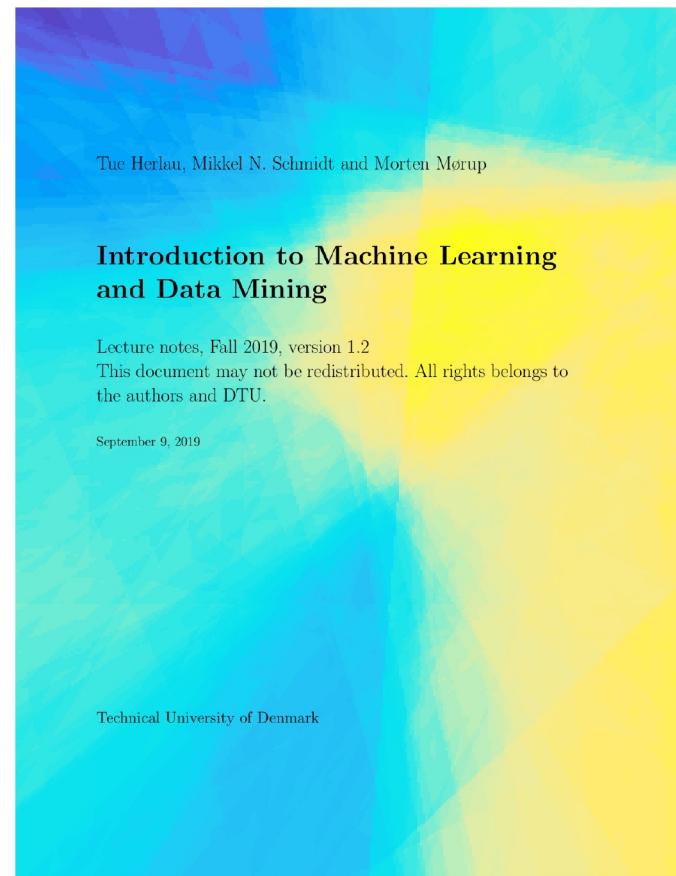
Department of Applied Mathematics and Computer Science

Today

Feedback Groups of the day:

Nicolai André Stæhr Kruhøffer, Chanyu Yang,
Qing Zheng, Sameer Agarwal, Adithya Iyer,
Thomas Phelps, Sayyada Shaama Hasan Kazmi,
Xin Li, Jorge Bertomeu Genis, Joachim
Lyngholm-Kjærby, Frederik Bonde Zilstorff, Doth
Angellina Pernille Andreasen, Ayub Abdi
Mohamed Nur, Emil Ørup Kristensen, Zehaib
Hussain, Peixuan Wang, Søren Meyer Nielsen, Li
Zixuan, Hanlu He, Magnus Fredslund, Jakob
Boëtius Andersen, Jean Ababii, Ali Waleed
Abbas, Abdullahi Abdirahman Mohamed,
Christian Skovmøller Agger, Christian Anberg,
Mads Christian Berggrein Andersen, Hjalte
Bækgaard Andersen, Pernille Lysgaard Andersen,
Frederik Lyhne Andersen, August Semrau
Andersen, Sebastian Sindlev Andersen, Branavan
Annalingam, Joan Maria Arenas Gomez, Shobhit
Arora, Nicklas Oliver Askjær, Olga
Athanasopoulou, Guðrun Anna Atladottir, Nicolo
Aurisano, Luka Avbreht, Anton Baht, Alexandra
Weronika Balicki, Jan-Georges Jersild Balin,
Bosse Bandowski, Javiera Bartolomé, Christoph
Bätz, Katrine Bay, Oliver Repholtz Behrens, Ian
Beissmann, Mikolaj Cyprian Bejster, Elias
Benameur, Colin Rolf Benker, Christian
Alexander Bertram, Sascha Peter Bilert, Asger
Birfau, Anton Birkedal, Rebekka Overgaard

Reading material: Chapter 2, Chapter 3



Lecture Schedule

① Introduction

3 September: C1

Data: Feature extraction, and visualization

② Data, feature extraction and PCA

10 September: C2, C3

③ Measures of similarity, summary statistics and probabilities

17 September: C4, C5

④ Probability densities and data Visualization

24 September: C6, C7

Supervised learning: Classification and regression

⑤ Decision trees and linear regression

1 October: C8, C9 (Project 1 due before 13:00)

⑥ Overfitting, cross-validation and Nearest Neighbor (Note: Tentative)

8 October: TBA

⑦ Bayes, Naive Bayes and performance evaluation (Note: Tentative)

22 October: TBA

Piazza online help: <https://piazza.com/dtu.dk/fall2019/02450>

Videos/streaming of lectures: <https://video.dtu.dk>

⑧ Artificial Neural Networks and Bias/Variance

29 October: C14, C15

⑨ AUC and ensemble methods

5 November: C16, C17

Unsupervised learning: Clustering and density estimation

⑩ K-means and hierarchical clustering

12 November: C18 (Project 2 due before 13:00)

⑪ Mixture models and density estimation

19 November: C19, C20

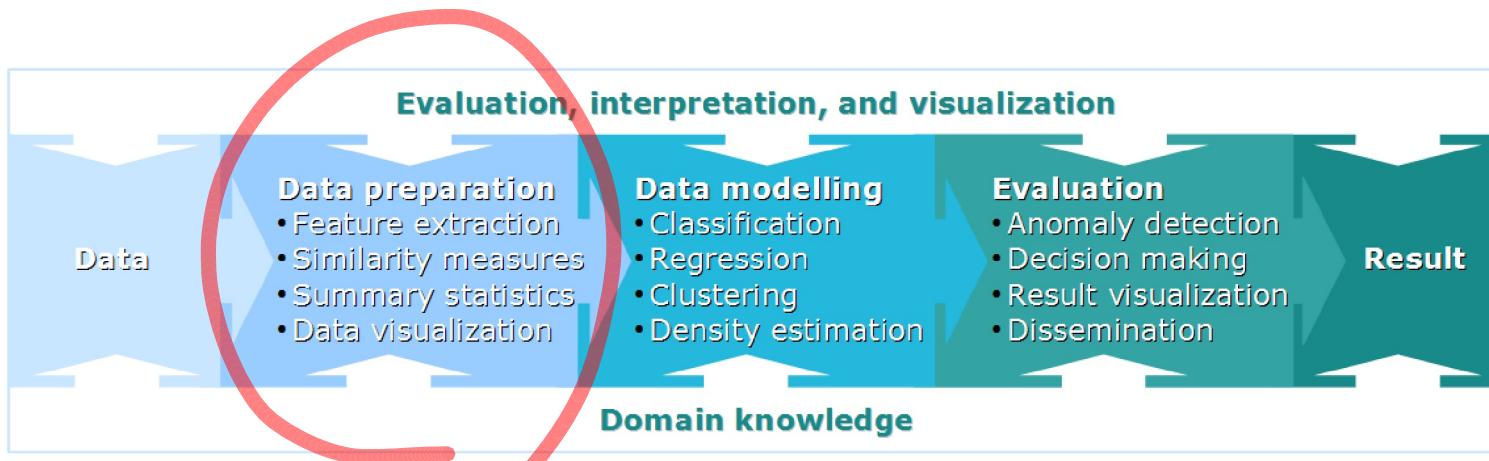
⑫ Association mining

26 November: C21

Recap

⑬ Recap and discussion of the exam

3 December: C1-C21 (Project 3 due before 13:00)



Learning Objectives

- Understand the types of data, their attributes and data issues
- Be able to apply principal component analysis for data visualization and feature extraction

What is data?

- Collection of data objects and their attributes
- An attribute is a property or characteristic of an object
 - Also known as variable, field, characteristic, or feature
- Collection of attributes describe an object
 - Also known as record, point, case, sample, entity, or instance

Attributes				
ID	Age	Gender	Name	
1	31	F	Alex	
2	24	M	Ben	
3	52	F	Cindy	
4	35	M	Dan	
5	58	M	Eric	
6	46	F	Fay	
7	42	M	George	

Discrete / continuous attributes

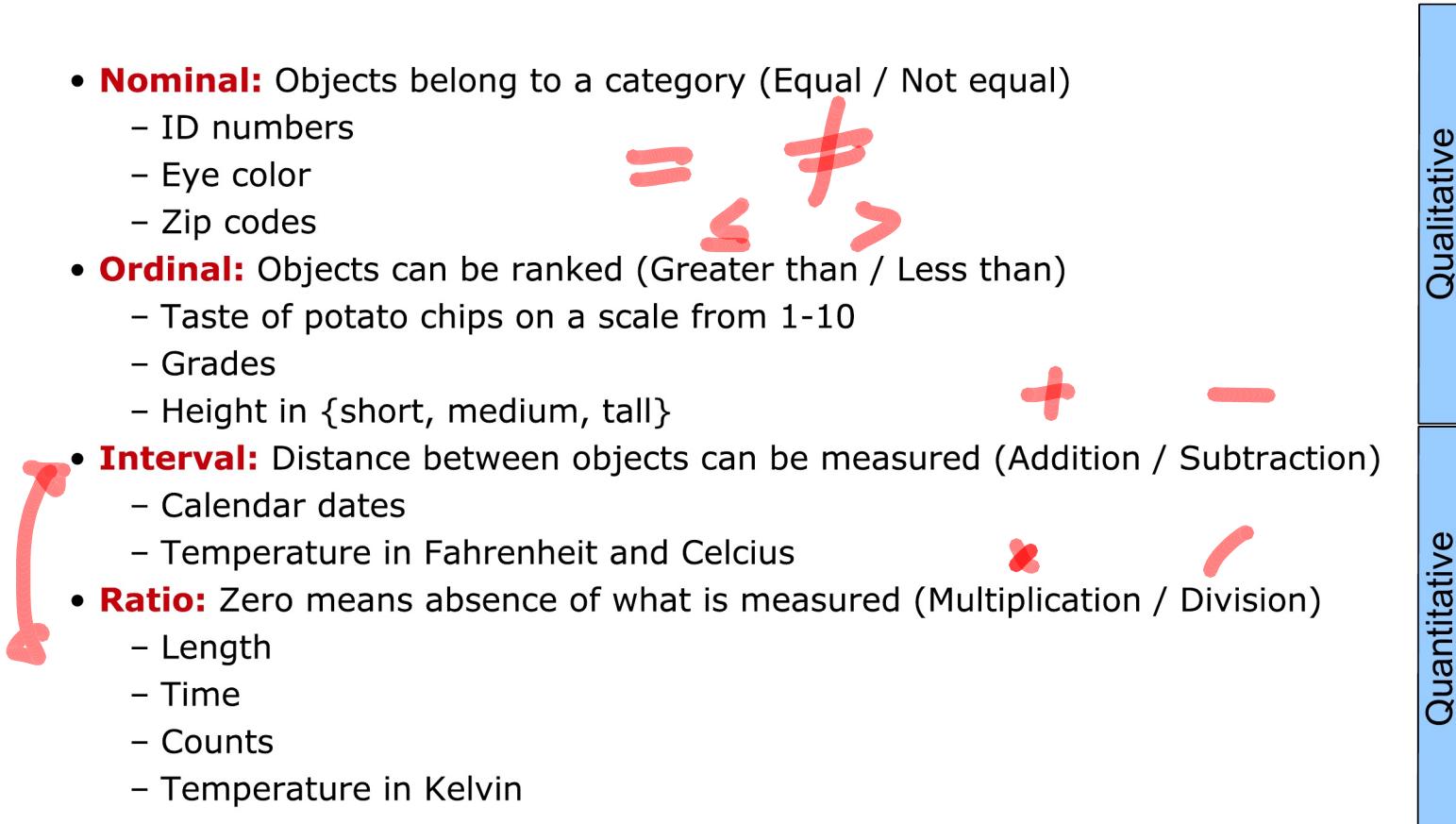
- **Discrete**

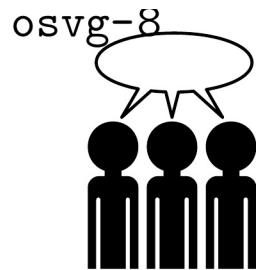
- Finite (or countably infinite) set of values
- Examples:
 - Zip codes
 - Counts
 - Set of words in a collection of documents
- Often represented as integer variables

- **Continuous**

- Has real numbers as attribute values
- Examples:
 - Temperature
 - Height
 - Weight.
- Often represented as floating point variables

Types of attributes

- **Nominal:** Objects belong to a category (Equal / Not equal)
 - ID numbers
 - Eye color
 - Zip codes
 - **Ordinal:** Objects can be ranked (Greater than / Less than)
 - Taste of potato chips on a scale from 1-10
 - Grades
 - Height in {short, medium, tall}
 - **Interval:** Distance between objects can be measured (Addition / Subtraction)
 - Calendar dates
 - Temperature in Fahrenheit and Celcius
 - **Ratio:** Zero means absence of what is measured (Multiplication / Division)
 - Length
 - Time
 - Counts
 - Temperature in Kelvin
- 
- The diagram illustrates the classification of attributes into Qualitative and Quantitative types. A vertical blue bar is divided into two sections: 'Qualitative' at the top and 'Quantitative' at the bottom. To the left of the bar, there is a red curved arrow pointing downwards, starting from the 'Qualitative' section and ending near the 'Ratio' category. To the right of the bar, there are four red symbols: an equals sign (=), a plus sign (+), a minus sign (-), and a multiplication sign (×). These symbols correspond to the operations that can be performed on each type of attribute: equality for Nominal, addition and subtraction for Interval, and multiplication and division for Ratio.



Discussion

- **Classify the following attributes**
 - a) Military rank
 - b) Angles measured in degrees
 - c) A persons year of birth
 - d) A persons age in years
 - e) Coat check number
 - f) Distance from center of campus
 - g) Number of patients in a hospital



- **Discrete**
 - Finite (or countably infinite) set of values
- **Continuous**
 - Real number
- **Nominal** (Equal / Not equal)
 - Objects belong to a category
- **Ordinal** (Greater than / Less than)
 - Objects can be ranked
- **Interval** (Addition / Subtraction)
 - Distance between objects can be measured
- **Ratio** (Multiplication / Division)
 - Zero means absence of what is measured

Quiz 1: Attribute types (Spring 2012)

No.	Attribute description	Abbrev.
x_1	Type (0 = served cold, 1 = served hot)	TYPE
x_2	Calories per serving	CAL
x_3	Grams of protein	PROT
x_4	Grams of fat	FAT
x_5	Milligrams of sodium	SOD
x_6	Grams of dietary fiber	FIB
x_7	Grams of complex carbohydrates	CARB
x_8	Grams of sugars	SUG
x_9	Milligrams of potassium	POT
x_{10}	Vitamins and minerals in 0%, 25%, or 100% of FDA recommendations	VIT
x_{11}	Shelf position (1, 2, or 3, counting from the floor)	SHELF
x_{12}	Weight in ounces of one serving	WEIGHT
x_{13}	Number of cups in one serving	CUPS
x_{14}	Name of cereal brand	NAME
y	Average rating of the cereal (from 0 to 100)	RAT

Table 1: Attributes in a study of cereals (i.e. breakfast products, taken from <http://lib.stat.cmu.edu/DASL/Datafiles/Cereals.html>).

In a study of healthy breakfast habits 77 cereal brands were investigated. The attributes of the data are given in Table 1. There are a total of 14 attributes denoted x_1-x_{14} and one output variable y which defines the average rating of the cereal products by the consumers.

Which statement about the attributes in the data set is *incorrect*?

- A. NAME is discrete and nominal.
- B. PROT, FAT and SOD are all continuous and ratio.
- C. TYPE and VIT are both discrete and ordinal.
- D. An attribute that is ratio will also be interval.
- E. Don't know.

= nominal
 < ordinal
 + interval
 ✗ ratio

Solution:

There are a finite set of brands thus NAME is discrete and as the only operators that can be applied to NAME is equal or not equal NAME is nominal. PROT, FAT and SOD are all continuous and since they have that zero means absence they are ratio. TYPE and VIT are both discrete, however, TYPE is not ordinal, i.e. Hot is not better than Cold, thus

TYPE must be considered nominal, VIT on the other hand is ordinal as 0% is less than 25% which in turn is less than 100%. An attribute that is ratio will also be both interval, ordinal and nominal, i.e. we can apply all the operations $=, \neq, >, <, +, -, *, /$ to a ratio attribute.

Types of data sets

- **Record data**
 - Collection of data objects and their attributes
 - Representation: Table
- **Relational data**
 - Collection of data objects and their relation
 - Representation: Graph
- **Ordered data**
 - Ordered collection of data objects
 - Representation: Sequence

Record data example: Market basket data

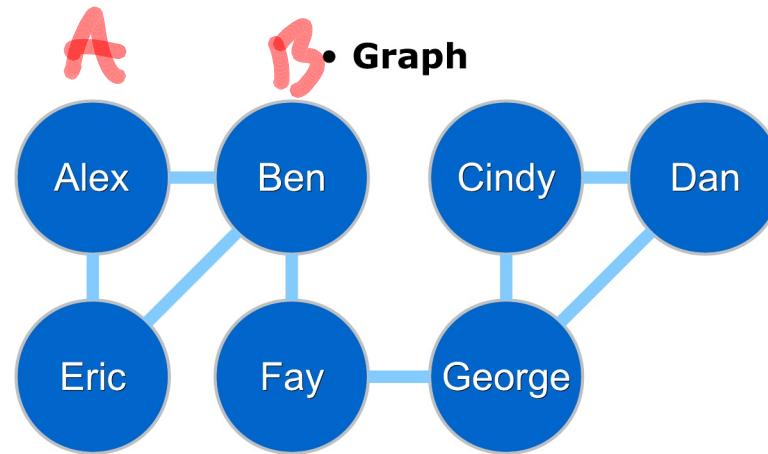
• Transaction data table

ID	Items
1	Bread, Soda, Milk
2	Beer, Bread
3	Beer, Soda, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Soda, Diaper, Milk

• Matrix

ID	Bread	Soda	Milk	Beer	Diaper
1	1	1	1	0	0
2	1	0	0	1	0
3	0	1	1	1	1
4	1	0	1	1	1
5	0	1	1	0	1

Relational data example: Who knows who?

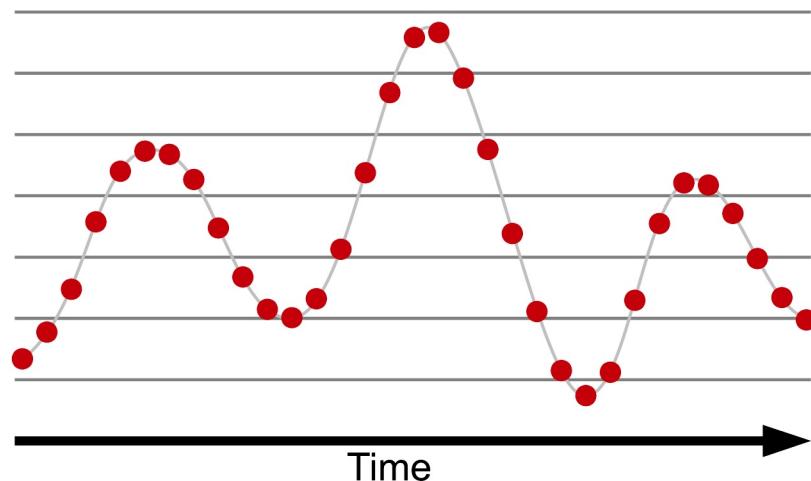


• Matrix

	A	B	C	D	E	F	G
A	0	1	0	0	1	0	0
B	1	0	0	0	1	1	0
C	0	0	0	1	0	0	1
D	0	0	1	0	0	0	1
E	1	1	0	0	0	0	0
F	0	1	0	0	0	0	1
G	0	0	1	1	0	1	0

Ordered data example: Time series

- Sequence



- Matrix

Time	Value
0	1.3
1	1.8
2	2.5
3	3.6
4	4.4
5	4.7
6	4.6
7	4.3
8	2.4
9	2.1
10	2.0
11	2.3
12	3.1

Data quality

- **Data is of high quality if they**
 - Are fit for their intended use
 - Correctly represent the phenomena they correspond to
- **Examples of quality problems**
 - Noise
 - Outliers
 - Missing values



Noise

- **Definition**

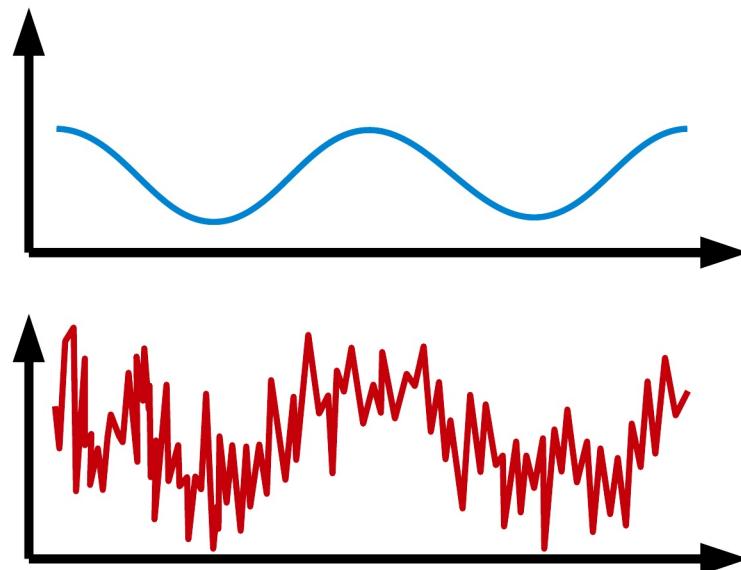
- Unwanted perturbation to a signal
- Unwanted data

- **Reasons for noise**

- Limits in measurement accuracy
- Interference from other signals
- Measurement of attributes not related to the data modeling task

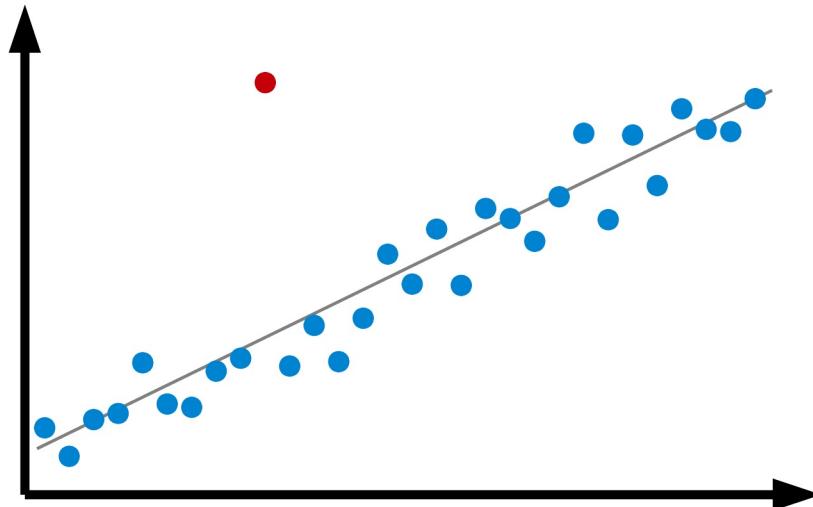
- **Handling noise**

- (– Exclude noisy attributes)
- Remove noise by filtering
 - Include a model of the noise



Outliers

- **Definition**
 - Data objects which are significantly different from most others
- **Reasons for outliers**
 - Measurement error
 - Natural property of data
- **Handling outliers**
 - Identify & exclude outliers
 - Model the outliers



Missing values

- **Definition**

- No value is stored for an attribute in a data object

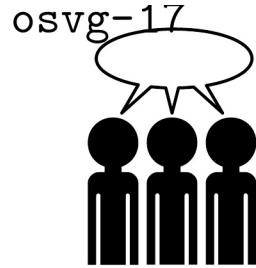
- **Reasons for missing values**

- Information is not collected
 - People decline to give their age
- Attribute is not applicable
 - Annual income is not applicable to children

- **Handling missing values**

- Eliminate data objects
- Estimate missing values (e.g. an average)
- Ignore the missing value in analysis

ID	Age	Gender	Name
1	31	F	Alex
2	(?)	M	Ben
3	52	F	Cindy
4	35	(?)	Dan
5	(?)	M	Eric
6	(?)	F	Fay
7	42	M	(?)



Discussion

- A group of people were asked to write how many children they have
 - Their response was this

3 1 ~~NONE~~ 2 7 3 15 0 1 3 2 zero ~~1~~ 0 1

- A research assistant typed the results into a table
 - His table looked like this

Children	3	1	0	2	7	5	15	0	1	3	-2	0	0	0	1
----------	---	---	---	---	---	---	----	---	---	---	----	---	---	---	---

- Are there any data quality issues?
 - Noise?
 - Outliers?
 - Missing values?
- Why have these issues occurred, and how should they be handled?

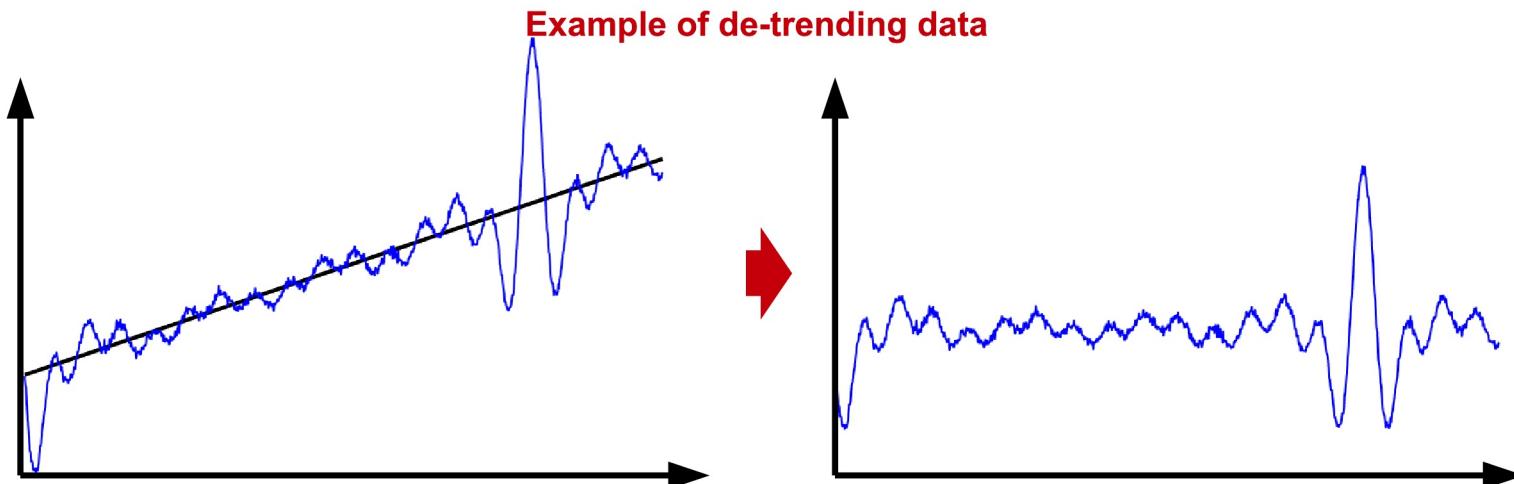
Dataset manipulations

- **Sampling**
 - Selecting a representative subset of data points
- **Feature subset selection**
 - Choose a subset of attributes
- **Feature extraction/transformation**
 - Create new features from existing attributes
 - Discretization and binarization
 - Apply a fixed transformation to an attribute
 - Aggregation several attributes into a single attribute
- **Dimensionality reduction**
 - Project data to a low-dimensional subspace

PCA

Feature processing

- Eliminating, suppressing, or attenuating certain aspects of the data
 - Noise removal in audio signals
 - Elimination of common words in text documents
 - Removal of background in images
 - Removal of examples which are corrupted
 - De-trending data (if it is not stationary)





Common feature transformations

ID	MPG	Cylinders	Horsepower	Weight	Year	Safety	Acceleration	Origin
1	18	8	150	3436	70	4	11	France
2	28	4	79	2625	82	4	18.6	USA
3	26	4	79	2255	76	3	17.7	USA
3	29	4	70	1937	76	1	14.2	Germany
4	NaN	8	175	3850	70	2	11	USA
5	24	4	90	2430	70	3	14.5	Germany
6	17.5	6	95	3193	76	4	17.8	USA
7	25	4	87	2672	70	-100	17.5	France
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
142	15	8	198	4341	70	2	10	USA

$$\mathbf{X} = \begin{bmatrix} 18 & 8 & 150 & 3436 & 70 & 4 & 11 & 3 \\ 28 & 4 & 79 & 2625 & 82 & 4 & 18.6 & 1 \\ \vdots & \vdots \\ 15 & 8 & 198 & 4341 & 70 & 2 & 10 & 1 \end{bmatrix}$$

Standardize:

$$\mathbf{X} = \begin{bmatrix} \cdots & (X_{1j} - \hat{\mu}_j)/\hat{\sigma}_j & \cdots \\ \cdots & (X_{2j} - \hat{\mu}_j)/\hat{\sigma}_j & \cdots \\ & \vdots & \\ \cdots & (X_{Nj} - \hat{\mu}_j)/\hat{\sigma}_j & \cdots \end{bmatrix}$$

$$\hat{\mu}_j = \frac{1}{N} \sum_{i=1}^N X_{ij}, \quad \hat{\sigma}_j = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (X_{ij} - \mu_j)^2}$$

Binarize/threshold:

$$\mathbf{X} = \begin{bmatrix} \cdots & 1_{[\theta, \infty[}(x_{1j}) & \cdots \\ \cdots & 1_{[\theta, \infty[}(x_{2j}) & \cdots \\ & \vdots & \\ \cdots & 1_{[\theta, \infty[}(x_{Nj}) & \cdots \end{bmatrix}$$

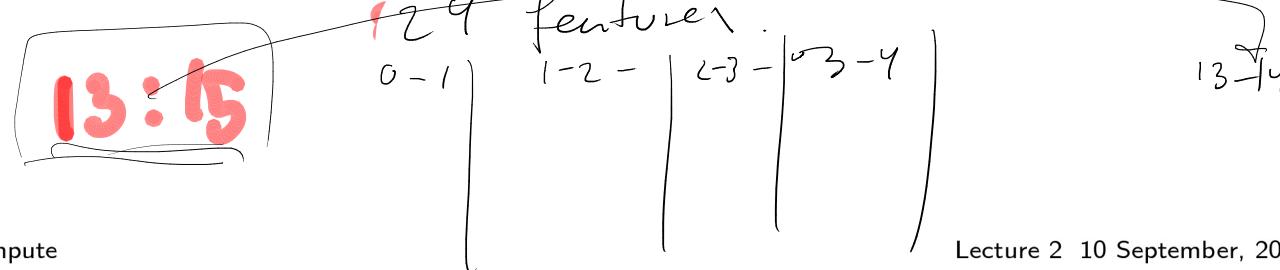
$$1_{[\theta, \infty[}(x) = 1 \text{ if } x \geq \theta \text{ otherwise } 0$$

One-out-of K encoding

Age	Height	Weight	Nationality
-0.2248	-0.4762	-0.2097	'Sweden'
-0.5890	0.8620	0.6252	'Sweden'
-0.2938	-1.3617	0.1832	'Sweden'
-0.8479	0.4550	-1.0298	'Sweden'
-1.1201	-0.8487	0.9492	'Norway'
2.5260	-0.3349	0.3071	'Norway'
1.6555	0.5528	0.1352	'Norway'
0.3075	1.0391	0.5152	'Norway'
X= -1.2571	-1.1176	0.2614	'Norway'
-0.8655	1.2607	-0.9415	'Sweden'
-0.1765	0.6601	-0.1623	'Norway'
0.7914	-0.0679	-0.1461	'Denmark'
-1.3320	-0.1952	-0.5320	'Denmark'
-2.3299	-0.2176	1.6821	'Sweden'
-1.4491	-0.3031	-0.8757	'Sweden'
0.3335	0.0230	-0.4838	'Sweden'
0.3914	0.0513	-0.7120	'Denmark'
0.4517	0.8261	-1.1742	'Sweden'
-0.1303	1.5270	-0.1922	'Norway'
0.1837	0.4669	-0.2741	'Denmark'

One-out-of-K coding

Age	Height	Weight			
-0.2248	-0.4762	-0.2097	0	0	1
-0.5890	0.8620	0.6252	0	0	1
-0.2938	-1.3617	0.1832	0	0	1
-0.8479	0.4550	-1.0298	0	0	1
-1.1201	-0.8487	0.9492	0	1	0
2.5260	-0.3349	0.3071	0	1	0
1.6555	0.5528	0.1352	0	1	0
0.3075	1.0391	0.5152	0	1	0
-1.2571	-1.1176	0.2614	0	1	0
-0.8655	1.2607	-0.9415	0	0	1
-0.1765	0.6601	-0.1623	0	1	0
0.7914	-0.0679	-0.1461	1	0	0
-1.3320	-0.1952	-0.5320	1	0	0
-2.3299	-0.2176	1.6821	0	0	1
-1.4491	-0.3031	-0.8757	0	0	1
0.3335	0.0230	-0.4838	0	0	1
0.3914	0.0513	-0.7120	1	0	0
0.4517	0.8261	-1.1742	0	0	1
-0.1303	1.5270	-0.1922	0	1	0
0.1837	0.4669	-0.2741	1	0	0



Bag of words representation

- First three sentences on [wikipedia.org](https://en.wikipedia.org)
 - 1 – The bag-of-words model is a simplifying assumption used in natural language processing and information retrieval
 - 2 – In this model, a text (such as a sentence or a document) is represented as an unordered collection of words, disregarding grammar and even word order
 - 3 – The bag-of-words model is used in some methods of document classification



(Image source: <https://pixabay.com/p-297223/>)

Bag of words representation

- First three sentences on [wikipedia.org](https://en.wikipedia.org)
 - The bag-of-words model is a simplifying assumption used in natural language processing and information retrieval
 - In this model, a text (such as a sentence or a document) is represented as an unordered collection of words, disregarding grammar and even word order
 - The bag-of-words model is used in some methods of document classification

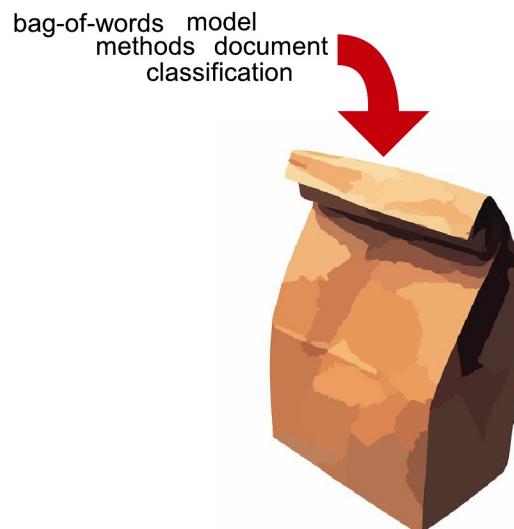


- We will treat **this text** as a data set and create a bag-of-words model of it



Bag of words representation

- Elimination of common words (so-called stop words)
 - The bag-of-words model is a simplifying assumption used in natural language processing and information retrieval
 - In this model, a text (such as a sentence or a document) is represented as an unordered collection of words, disregarding grammar and even word order
 - The bag-of-words model is used in some methods of document classification



Bag of words representation

- Representation as matrix

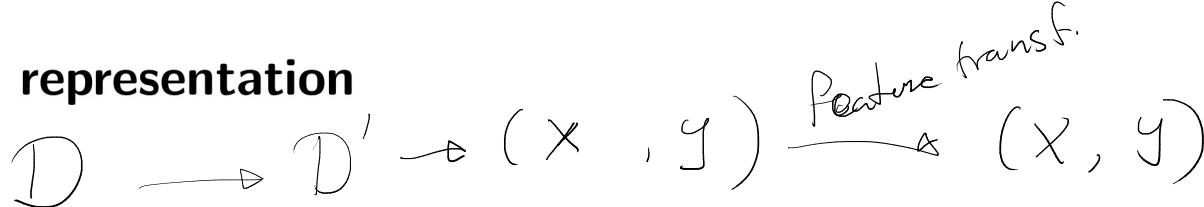
Word	Sentence		
	1	2	3
bag-of-words	1		1
model	1	1	1
simplifying	1		
assumption	1		
natural	1		
language	1		
processing	1		
information	1		
retrieval	1		
text		1	
sentence		1	
document		1	1
represented		1	
unordered		1	
collection		1	
words		1	
disregarding		1	
grammar		1	
word		1	
order		1	
methods			1
classification			1

Bag of words representation

- Stemming

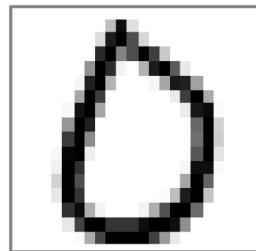
Word	Sentence		
	1	2	3
bag-of-word*	1		1
model*	1	1	1
simplif*	1		
assum*	1		
natural*	1		
languag*	1		
process*	1		
information*	1		
retriev*	1		
text*		1	
sentence*		1	
document*	1		1
represent*		1	
unorder*		1	
collect*		1	
word*		2	
disregard*		1	
grammar*		1	
order*		1	
method*			1
classif*			1

Image representation



- Example: Handwritten digits

- Preprocessing
 - Digitalization
 - Centering
 - Rotation
 - Scaling



28×28

- Vectorization

1×784

$$M_0 = \begin{bmatrix} 0 & \dots & 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & \dots & 0.3 & 1 & 0.2 & 0 & \dots & 0 \\ \vdots & & & & & & & \vdots \\ 0 & \dots & 0 & 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

$$x_0 = [0 \dots 0 \underbrace{0.3 \ 1 \ 0.2 \ 0}_{M=28^2} \dots 0]^T \quad 28^2 \times 1.$$

- Matrix representation of data set

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_N \end{bmatrix}$$

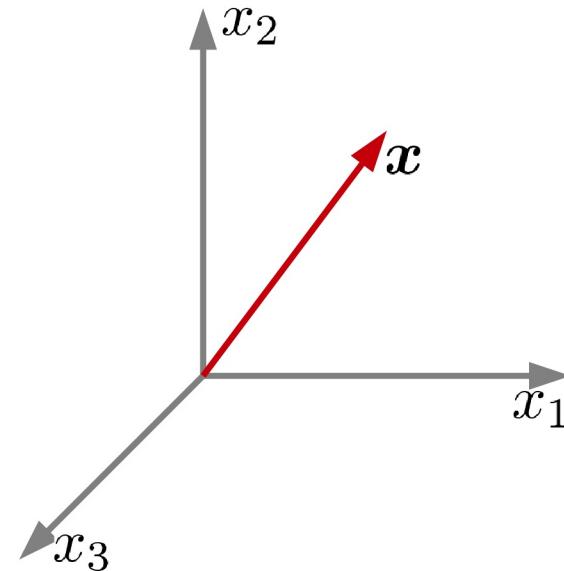


If each image is 28×28 pixels
then \mathbf{X} is a $N \times 784$ matrix. September, 2019

Vector space representation

- All these data objects have a vector space representation

$$\boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_M \end{bmatrix}$$



Plan for the rest of today:

- Linear algebra recap (subspaces and projections)
- The **goal** of Principal Component Analysis (PCA)
- Derivation of PCA
- Singular Value Decomposition used to implement PCA
- Use of PCA for data visualization

Vectors and matrices

- Common matrix notation

$$\mathbf{A}, A, \overline{\mathbf{A}}$$

$$\mathbf{A} = \begin{bmatrix} a_{1,1} & \cdots & a_{1,M} \\ \vdots & & \vdots \\ a_{N,1} & \cdots & a_{N,M} \end{bmatrix} \in \mathbb{R}^{N \times M}$$

- Common vector notation

$$\mathbf{x}, x, \overline{x}, \vec{x}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_M \end{bmatrix} \in \mathbb{R}^M$$

Matrix multiplication

- Two matrices can be multiplied $\mathbf{AB} = \mathbf{C}$
 - if the number of columns in the first equals the number of rows in the second

$$\begin{array}{c} \text{A} \times \text{B} = \text{C} \\[10pt] L \times M \quad M \times N \quad L \times N \\[10pt] \begin{matrix} 3 \times 4 \text{ matrix} \\ \left[\begin{matrix} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ 1 & 2 & 3 & 4 \end{matrix} \right] \end{matrix} \begin{matrix} 4 \times 5 \text{ matrix} \\ \left[\begin{matrix} \cdot & \cdot & \cdot & a & \cdot \\ \cdot & \cdot & \cdot & b & \cdot \\ \cdot & \cdot & \cdot & c & \cdot \\ \cdot & \cdot & \cdot & d & \cdot \end{matrix} \right] \end{matrix} = \begin{matrix} 3 \times 5 \text{ matrix} \\ \left[\begin{matrix} \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & x_{3,4} & \cdot \end{matrix} \right] \end{matrix} \end{array}$$

$$x_{3,4} = 1 \cdot a + 2 \cdot b + 3 \cdot c + 4 \cdot d$$

Example:

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 19 & 22 \\ 43 & 50 \end{bmatrix}$$

Matrix transpose

- The transpose of a matrix

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 6 \\ 7 & 8 & 9 \end{bmatrix} \quad A^\top = \begin{bmatrix} 1 & 3 & 7 \\ 2 & 4 & 8 \\ 3 & 6 & 9 \end{bmatrix}$$

The diagram shows a 3x3 matrix A . The element 7 at position $(3,1)$ is circled in red. Arrows point from the circled 7 to its position in the transpose matrix A^\top at position $(1,3)$, illustrating that the transpose of a matrix swaps rows and columns.

- Transpose of a sum

$$(A + B)^\top = A^\top + B^\top$$

- Transpose of a product

$$(AB)^\top = B^\top A^\top$$

$$(Ax)^\top y = x^\top A^\top y = x^\top (A^\top y)$$

$$\left[\quad \right]^\top \quad \left[\quad \right]$$

Diagram showing two vectors enclosed in brackets, with a superscript \top indicating their transpose.

The identity matrix

- Ones on the diagonal and zeros everywhere else

$$\mathbf{I} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} \quad \mathbf{I}^\top = \mathbf{I}$$

- Multiplying by the identity does not change anything

$$\begin{aligned} \mathbf{IA} &= \mathbf{A} \\ \mathbf{I}_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{A} &= \begin{bmatrix} a & b \\ c & d \end{bmatrix} \\ \mathbf{I}_2 \mathbf{A} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \end{aligned}$$

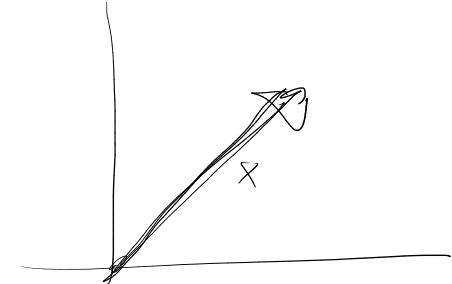
- For a square matrix, the inverse satisfies

$$\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$$

Norms

- The (Euclidian) norm of a vector measures it's length (magnitude):

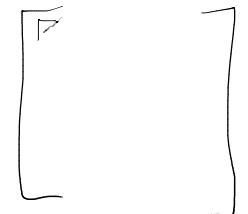
$$\|\mathbf{x}\|_2 = \sqrt{\sum_{n=1}^N |x_n|^2} = \sqrt{\mathbf{x}^\top \mathbf{x}}$$



- The Frobenius norm of a matrix measures it's magnitude:

$$\|\mathbf{X}\|_F^2 = \sum_{i,j} x_{i,j}^2 = \text{trace}(\mathbf{X} \mathbf{X}^T) = \text{trace}(\mathbf{X}^T \mathbf{X})$$

$$\|\mathbf{X}\|_F = \sqrt{\sum_{i,j} x_{i,j}^2}$$



Where trace takes the sum of the diagonal elements, i.e.

$$\text{trace}(\mathbf{A}) = \sum_{i=1}^N a_{i,i}$$

$$\text{trace} \begin{bmatrix} a & b & c \\ 0 & d & e \\ 0 & 0 & f \end{bmatrix} = a + d + f.$$

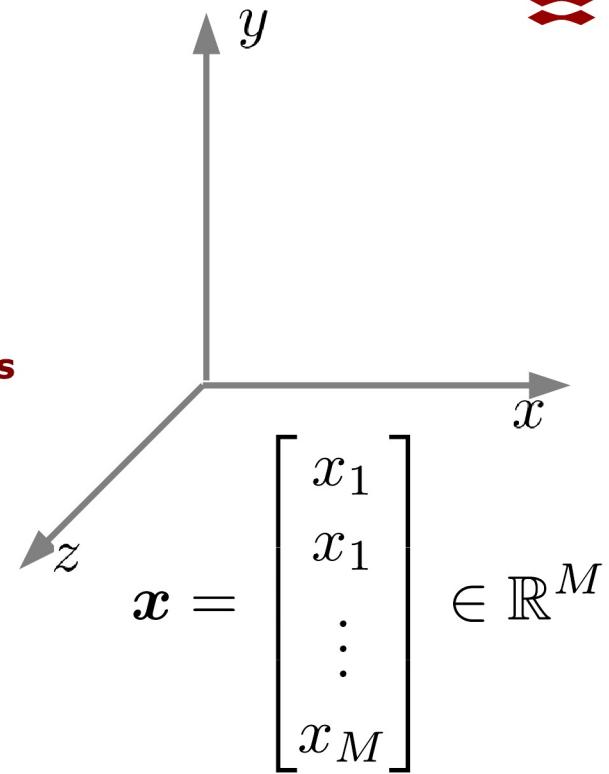
Vector spaces

- A M-dimensional vector space is just \mathbb{R}^M
- This is the set of all M-dimensional vectors
- A vector space is closed under **linear combinations**

$$a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \cdots + a_n \mathbf{x}_n$$

$\mathbf{x}_1, \dots, \mathbf{x}_n$ Vectors

a_1, \dots, a_n Numbers



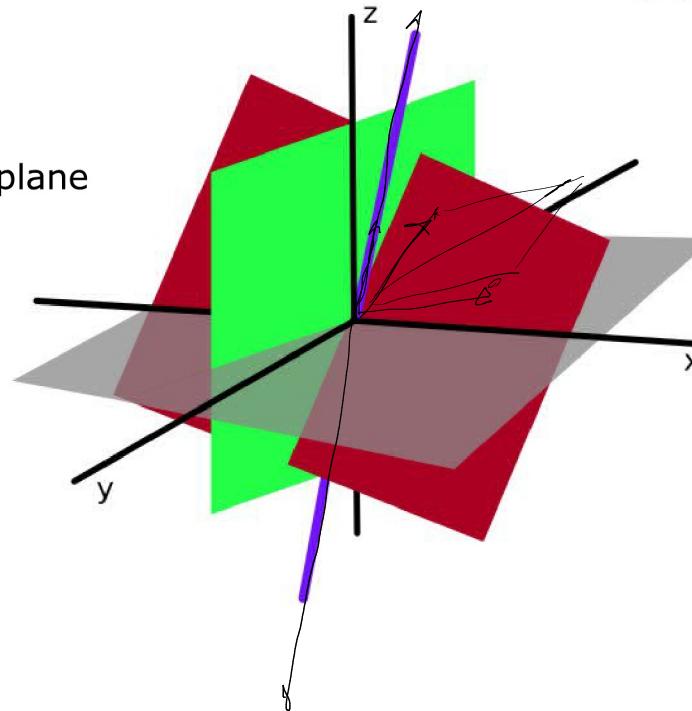
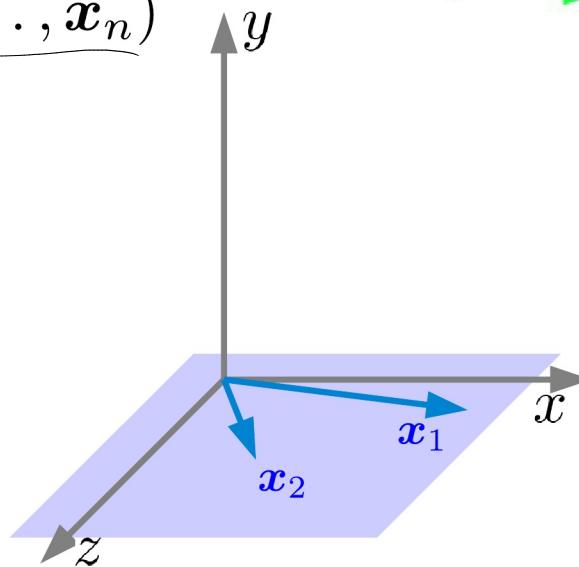
Subspaces

- A **subspace** generalizes the concept of a line/plane
- If we consider n vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$
the **span** is then all linear combinations

$$\mathbf{z} = a_1\mathbf{x}_1 + a_2\mathbf{x}_2 + \cdots + a_n\mathbf{x}_n$$

and it is said to be a **subspace**

$$\underline{V = \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_n)}$$



Basis of a (sub)space

- Vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ are said to be **linearly independent** if

$$\mathbf{0} = a_1\mathbf{x}_1 + a_2\mathbf{x}_2 + \dots + a_n\mathbf{x}_n$$

implies $a_1 = a_2 = \dots = a_n = 0$

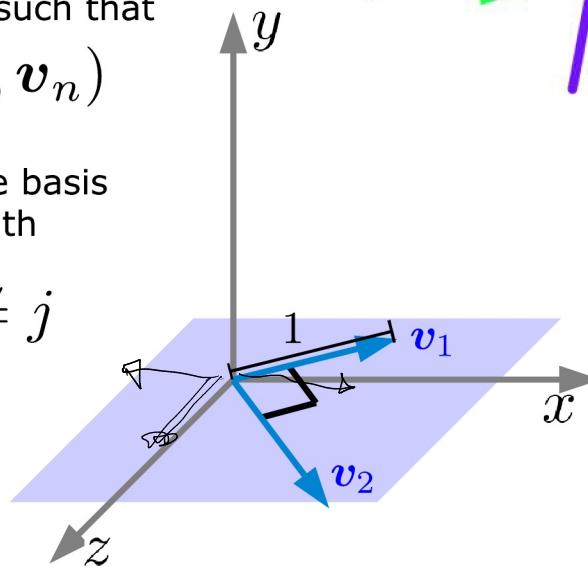
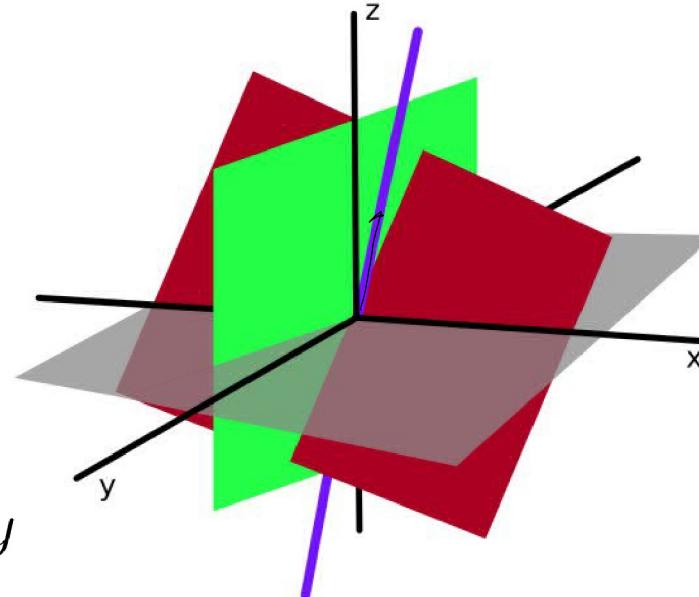
- A **basis** of a vector space V are n linearly independent vectors such that

$$V = \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_n)$$

- A basis is **orthonormal** if the basis is orthogonal and of unit length

$$\mathbf{v}_i^T \mathbf{v}_j = 0 \text{ for } i \neq j$$

$$\|\mathbf{v}_i\| = 1$$



Basis of a (sub)space

- A **basis** of a vector space V are n linearly independent vectors such that

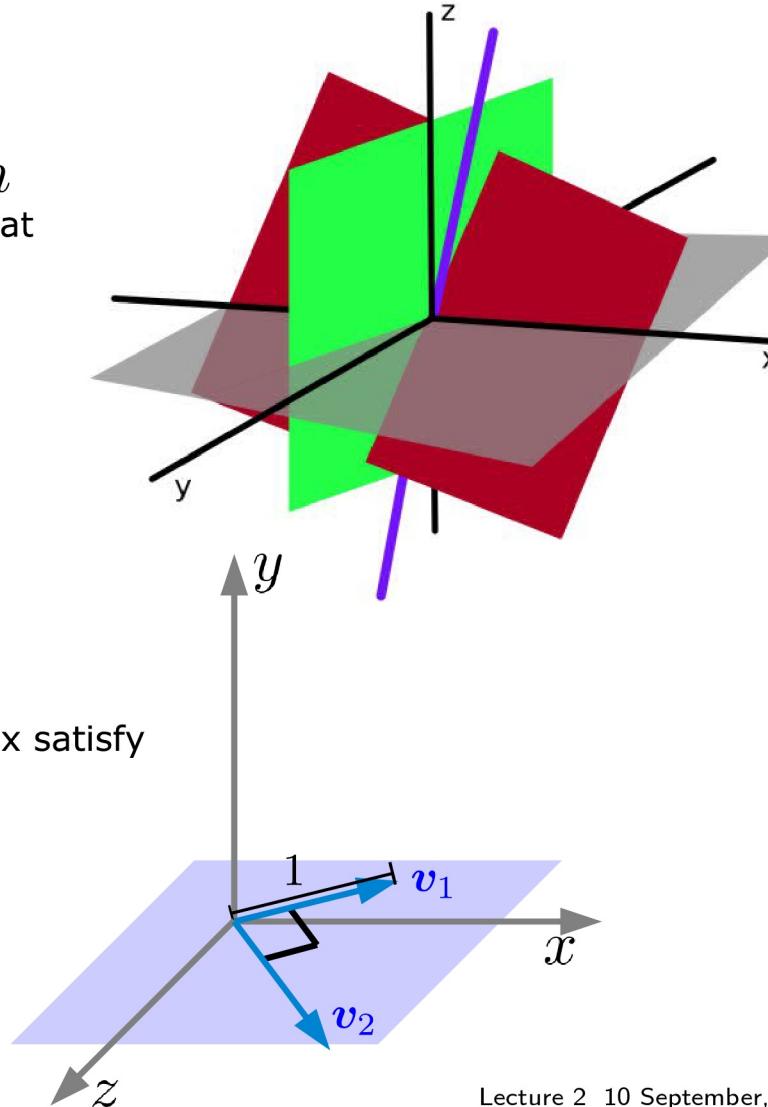
$$V = \text{span}(\mathbf{v}_1, \dots, \mathbf{v}_n)$$

- We collect the basis into a matrix

$$\mathbf{V} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_N \end{bmatrix}$$

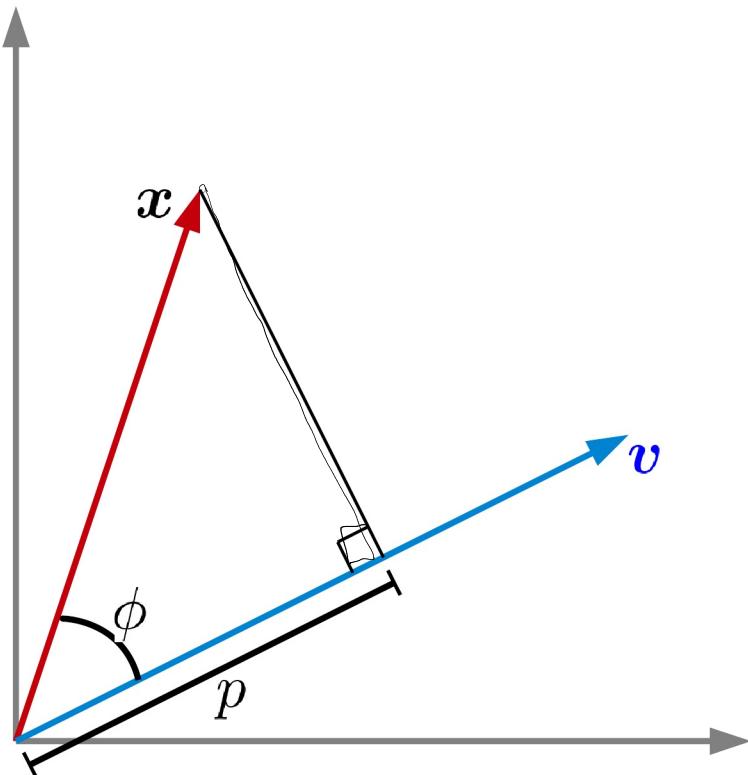
- If the basis is orthonormal the matrix satisfy

$$\mathbf{V}^\top \mathbf{V} = \mathbf{I}, \quad \mathbf{V}^\top = \mathbf{V}^{-1}$$



Projection

- Projection onto a vector



- Angle between vectors

$$\cos(\phi) = \frac{\mathbf{v}^\top \mathbf{x}}{\|\mathbf{x}\|_2 \|\mathbf{v}\|_2}$$

- Length of projection

$$p = \|\mathbf{x}\|_2 \cos(\phi) = \frac{\mathbf{v}^\top \mathbf{x}}{\|\mathbf{v}\|_2}$$

- Projection onto unit vector

$$p = \mathbf{v}^\top \mathbf{x}$$

Projection onto a subspace

DTU
 $b \in \mathbb{R}^3 \times V = \mathbb{R}^{3 \times 2}$

- **Projection onto a subspace**

- Subspace of dimension n defined by a orthonormal basis matrix V
- Projection of \mathbf{x} (M dimensional) onto V given by

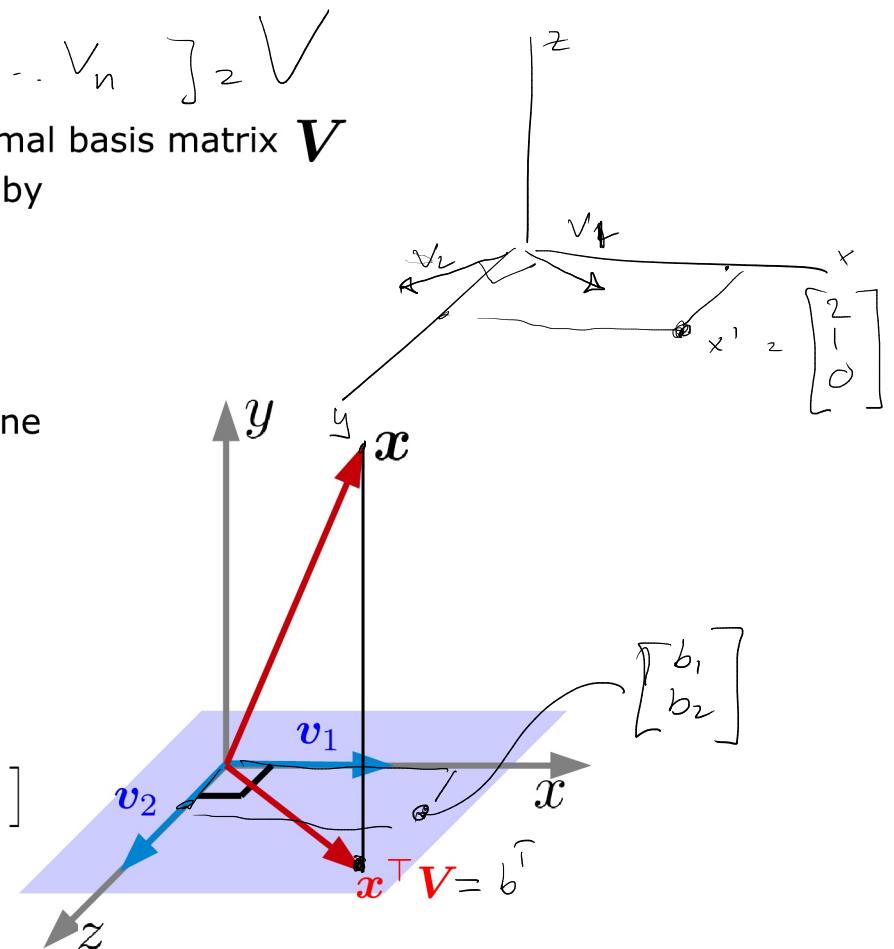
$$\mathbf{b}^T = \mathbf{x}^T V$$

- 'Reconstruction' can be found as: $\mathbf{x}' = V\mathbf{b}$

Example: Projection of 3-D vector onto the (x,z) plane

$$V = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

$$\mathbf{x}^T V = [x \ y \ z] \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix} = [x \ z]$$



Projection onto a subspace

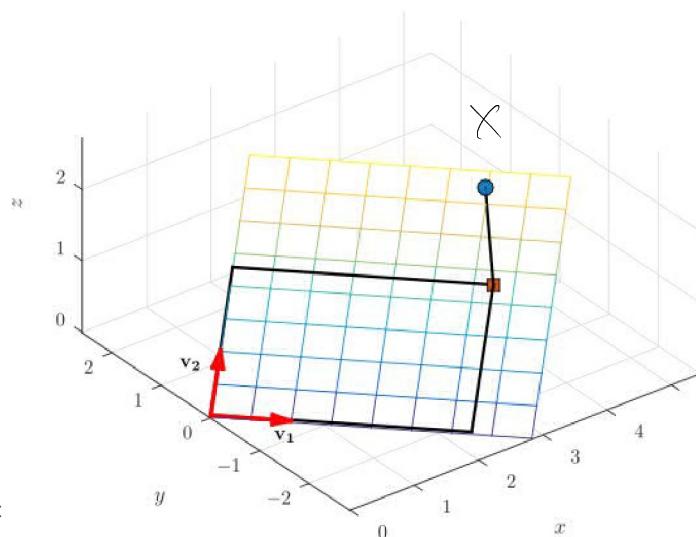
- **Projection onto a subspace**

- Subspace of dimension n defined by a orthonormal basis matrix V
- Projection of \mathbf{x} (M dimensional) onto V given by

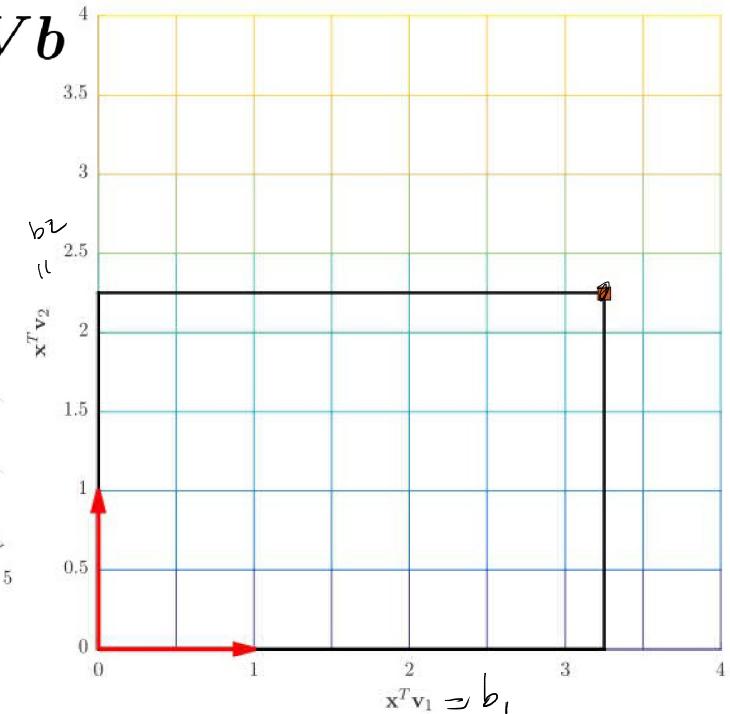
$$\mathbf{b}^T = \mathbf{x}^T V$$

- 'Reconstruction' can be found as: $\mathbf{x}' = V\mathbf{b}$

Example 2:

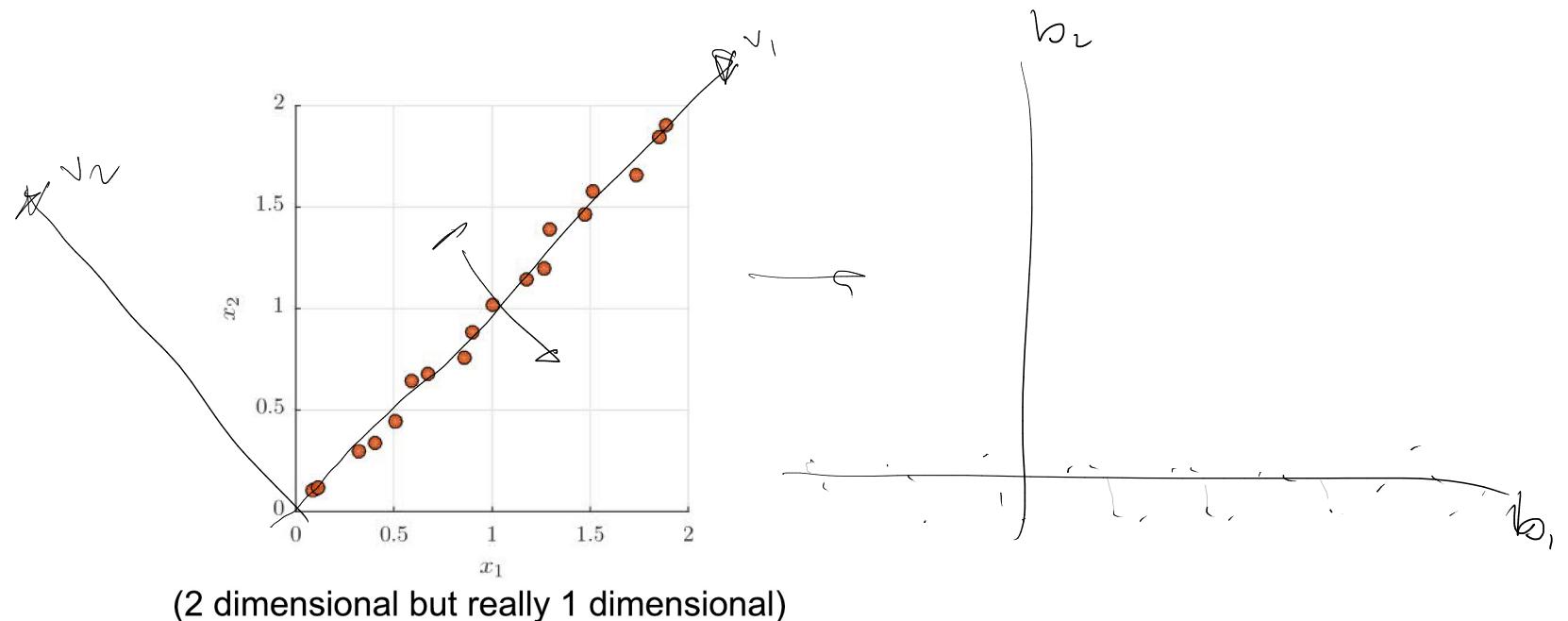


43 DTU Cc



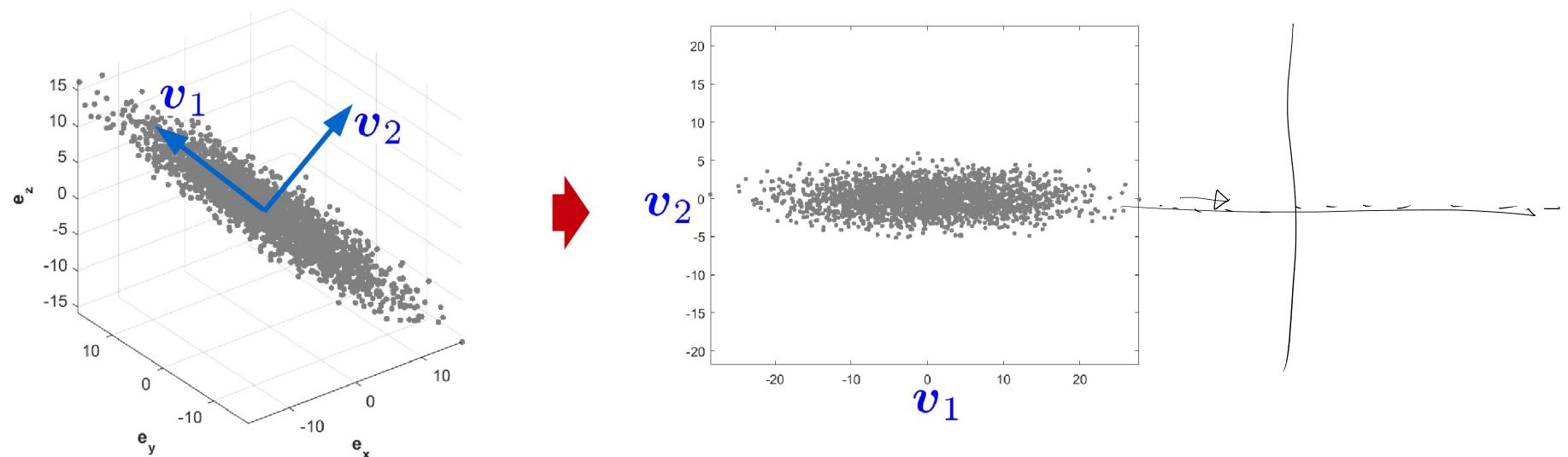
PCA for high-dimensional data

- Much data is high-dimensional
- We want to find a **lower**-dimensional representation of the **high**-dimensional data



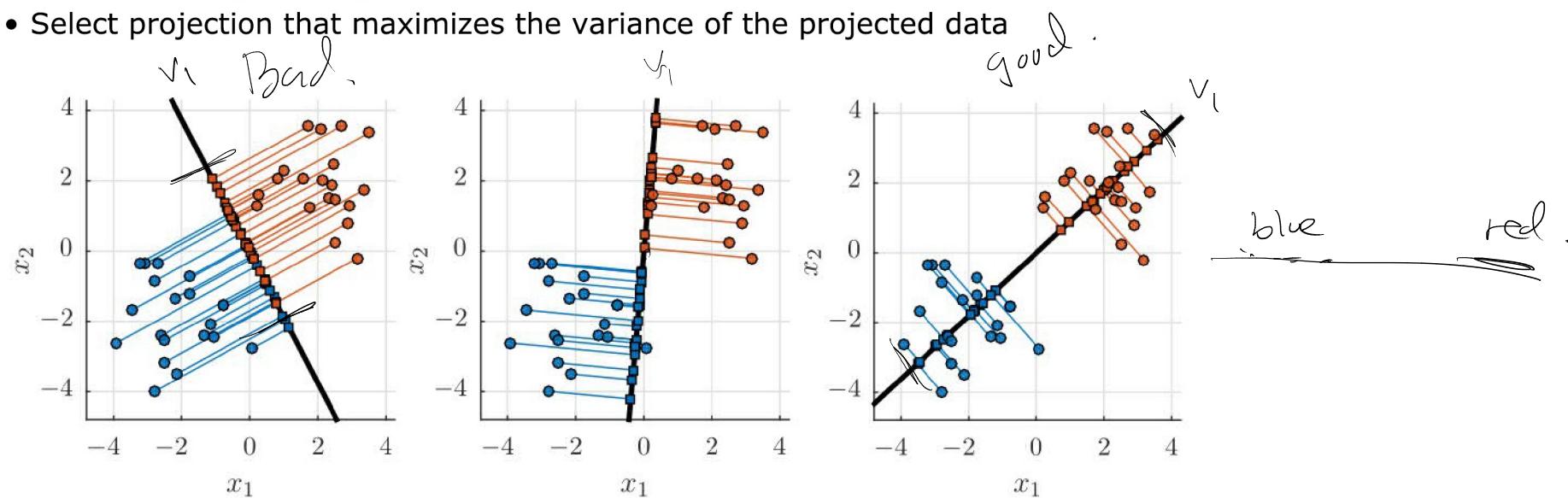
PCA for high-dimensional data

- Much data is high-dimensional
- We can **project high** dimensional data to a **lower** dimensional **subspace**
- But what is a good projection?



PCA for high-dimensional data

- Much data is high-dimensional
- We can project high dimensional data to a lower dimensional subspace
- But what is a good projection?
- Select projection that maximizes the variance of the projected data



osvg-49

$$\text{Var}[x] = \frac{1}{n-1} [\mathbb{E}[x] - \bar{x}]^2$$

PCA derivation

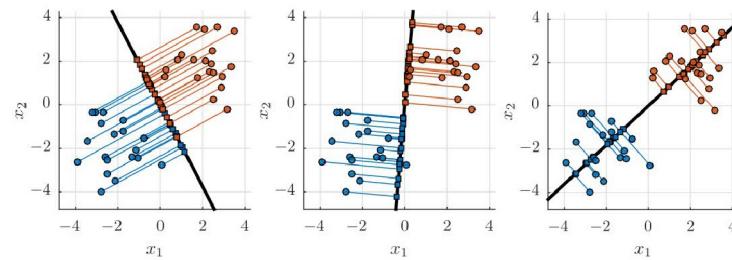
Projection of x_i onto unit vector v : $b_i = \underline{\underline{x}_i v}$

$$\begin{aligned}\text{Var}[b] &= \frac{1}{N-1} \sum_{i=1}^N \left[b_i - \frac{1}{N} \sum_{j=1}^N b_j \right]^2 \\ &= \frac{1}{N-1} \sum_{i=1}^N \left[\underline{x}_i^\top v - \frac{1}{N} \sum_{j=1}^N \underline{x}_j^\top v \right]^2 \\ &= \frac{1}{N-1} \sum_{i=1}^N \left[\left(\underline{x}_i^\top - \frac{1}{N} \sum_{j=1}^N \underline{x}_j^\top \right) v \right]^2 \\ &= \frac{1}{N-1} \sum_{i=1}^N (\tilde{x}_i^\top v)^2 \quad \boxed{\tilde{x}_i = x_i - \mathbf{m}}\end{aligned}$$

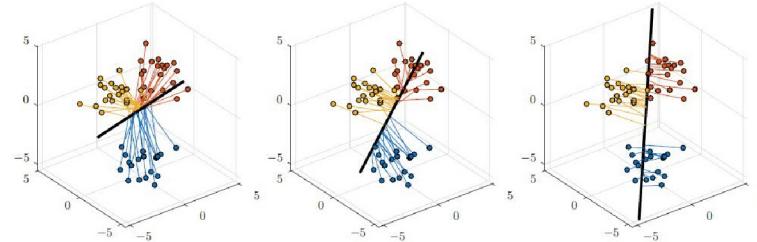
$$= \frac{1}{N-1} \sum_{i=1}^N \tilde{x}_i^\top \tilde{x}_i$$

$$\tilde{x} = \begin{bmatrix} \tilde{x}_1^\top \\ \tilde{x}_2^\top \\ \vdots \\ \tilde{x}_N^\top \end{bmatrix}$$

2D example



3D example



$$A\mathbf{v} = \lambda\mathbf{v}, \quad A \text{ is a } N \times N \text{ matrix}$$

We say \mathbf{v} is an eigenvector with eigenvalue λ

$$\begin{aligned}&= \frac{1}{N-1} \mathbf{v}^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{v} \\ &= \frac{1}{N-1} \sigma_i^2 \mathbf{v}^\top \mathbf{v} \\ &= \frac{1}{N-1} \sigma_i^2\end{aligned}$$

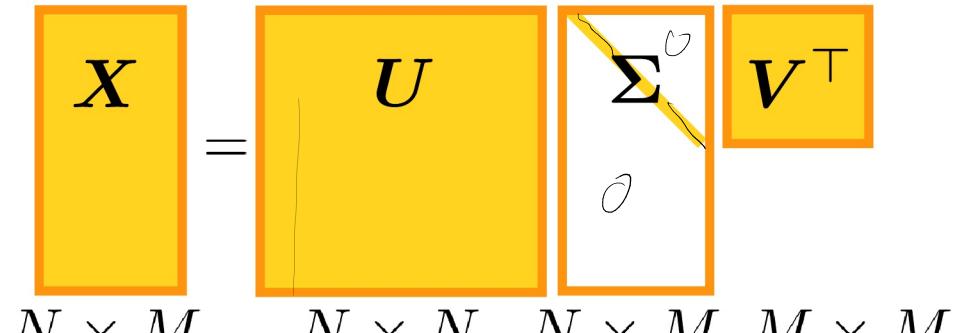
DTU

$$\mathbf{m} = \frac{1}{N} \sum_{j=1}^N x_j.$$

The Singular Value Decomposition (SVD)

(Eugino Beltrami & Camille Jordan, independently, 1873-1874)

Any $N \times M$ matrix can be decomposed as follows:

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


X
 $N \times M$ U
 $N \times N$ Σ
 $N \times M$ V^\top
↑ Orthonormal ↑ Diagonal ↑ Orthonormal

$$U = \begin{bmatrix} \vdots & \ddots & \vdots \\ u_1 & \cdots & u_N \end{bmatrix} \quad V = \begin{bmatrix} \vdots & \ddots & \vdots \\ v_1 & \cdots & v_M \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_M \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix}$$

$\sigma_1, \dots, \sigma_M$
is known as the singular values

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_M \geq 0$$

$$\text{if } i \neq j: \Sigma_{i,j} = 0, \quad U^\top U = I_{N \times N}, \quad V^\top V = I_{M \times M}$$

The Singular Value Decomposition (SVD)

(Eugino Beltrami & Camille Jordan, independently, 1873-1874)

Any $N \times M$ matrix can be decomposed as follows:

$$\tilde{X} = U \Sigma V^\top$$

$$\tilde{X} = \begin{matrix} U \\ N \times N \end{matrix} \Sigma \begin{matrix} V^\top \\ M \times M \end{matrix}$$

$$\begin{aligned} (\tilde{X}^\top \tilde{X})v_i &= (U \Sigma V^\top)^\top U \Sigma V^\top v_i \\ &= (V \Sigma^\top U^\top U \Sigma V^\top) v_i \\ &= V \Sigma^\top \Sigma v_i = \sigma_i^2 v_i \end{aligned}$$

$$A\mathbf{v} = \lambda\mathbf{v}, \quad A \text{ is a } N \times N \text{ matrix}$$

We say \mathbf{v} is an eigenvector with eigenvalue λ

$$V = \begin{bmatrix} | & | & | \\ v_1 & v_2 & \dots & v_M \\ | & | & | \end{bmatrix} \quad \sigma_i = \sqrt{\lambda_i} \quad v_i = \frac{1}{\sqrt{\lambda_i}} v_i$$

$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_M \geq 0$

is known as the singular values

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_M \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}$$

osvg-51

Principal component analysis (PCA)

(Karl Pearson, 1901)

- 1) Subtract the mean from each observation $\tilde{x}_i = x_i - \mathbf{m}$
- 2) Apply singular value decomposition (SVD) $\tilde{\mathbf{X}} = \mathbf{U} \Sigma \mathbf{V}^\top$

$$\tilde{\mathbf{X}} = \mathbf{U} \Sigma \mathbf{V}^\top$$

$N \times M \quad N \times N \quad N \times M \quad M \times M$

$$\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix}, \mathbf{m} = \frac{1}{N} \sum_{i=1}^N x_i$$

DTU

$$\tilde{\mathbf{X}} = \begin{bmatrix} \tilde{x}_1 \\ \vdots \\ \tilde{x}_N \end{bmatrix}$$

$$\mathbf{V}_{(K)} = \begin{bmatrix} v_1 & \dots & v_K \end{bmatrix}, \mathbf{V}_2 = \begin{bmatrix} v_1 & \dots & v_M \end{bmatrix}$$

- 3) Select first K columns of \mathbf{V} (the PCA projection operation) and first K columns of Σ .

$$\hat{\mathbf{X}} = \mathbf{U} \Sigma_{(K)}$$

$N \times K$
(PCA components or PCA projection of the data)

$$\mathbf{V}_{(K)} = \begin{bmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_K \end{bmatrix}$$

$$\mathbf{V}_{(K)}$$

$M \times K$
(PCA loadings)

Principal component analysis (PCA)

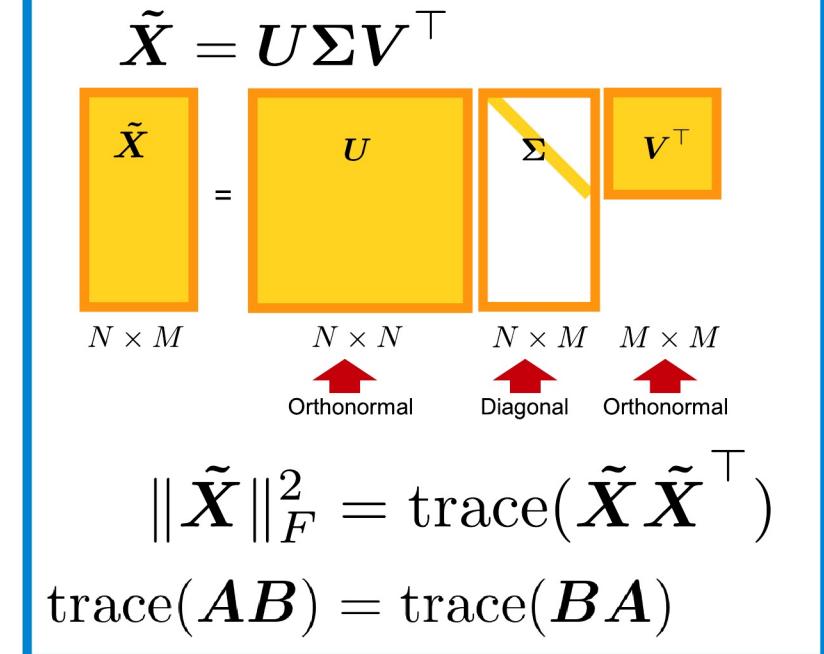
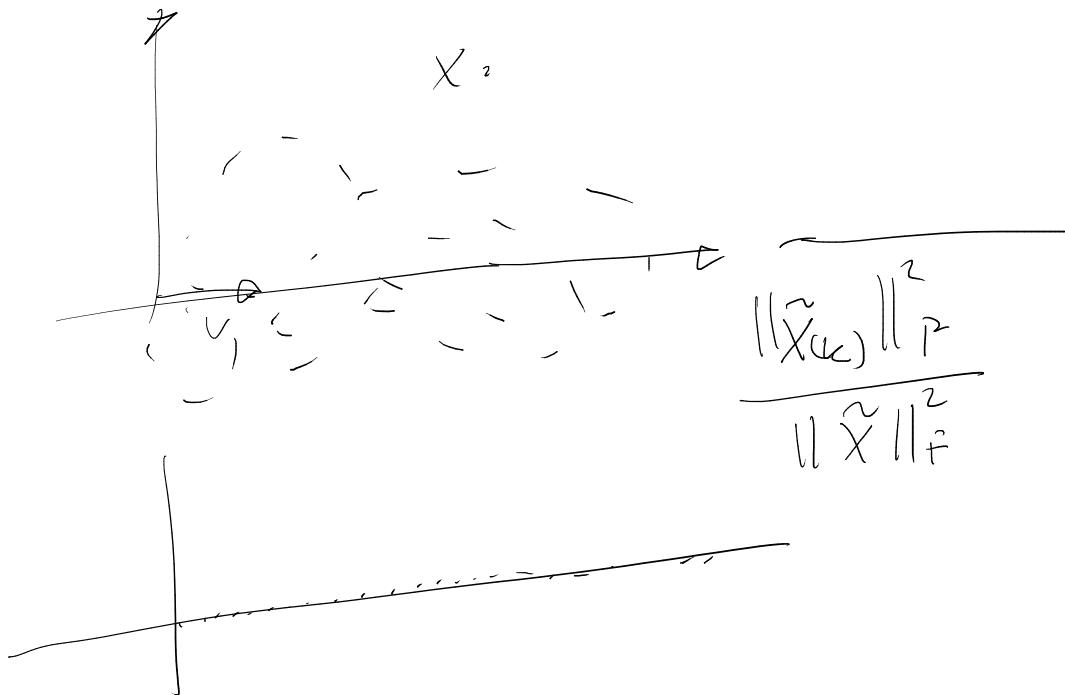
$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$$

- Entries in the diagonal matrix Σ are called **singular values**
 - They are sorted (largest first)
 - Indicate how much variability is explained by the corresponding component
 - 1st component explains most of the variability
 - 2nd component explains most of the remaining variability
 - Etc.

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_N \end{bmatrix} \quad \sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_N \geq 0$$

Explained Variance

- We will show that $\|\tilde{\mathbf{X}}\|_F^2 = \sum_i \sigma_i^2$
where $\sigma_i = \Sigma_{i,i}$



Fraction of the variation in the data explained by the i^{th} principal component is given by:

$$\frac{\sigma_i^2}{\sum_j \sigma_j^2}$$

And by the first K principal components

$$\frac{\sum_{i=1}^K \sigma_i^2}{\sum_j \sigma_j^2}$$

Quiz 2: PCA (Fall 2012)

No.	Attribute description	Abbrev.
x_1	Age (in years)	AGE
x_2	Gender (Female=0, Male=1)	GDR
x_3	Total Bilirubin	TB
x_4	Direct Bilirubin	DB
x_5	Alkaline Phosphotase	AP
x_6	Alamine Aminotransferase	AlA
x_7	Aspartate Aminotransferase	AsA
x_8	Total Proteins	TP
x_9	Albumin	AB
x_{10}	Albumin to Globulin ratio	A/G
y	0=No liver disease, 1=Liver disease	LD

Table 1: Attributes in a study on liver disease among Indians living in the north eastern part of Andhra Pradesh, India. (taken from <http://archive.ics.uci.edu/ml/datasets/ILPD> +%28Indian+Liver+Patient+Dataset%29). The data has 10 input attributes x_1-x_{10} and one output variable y which defines whether the subject considered has a liver disease ($y = 1$) or not ($y = 0$). x_3-x_9 are non-negative measurements giving the concentrations of various quantities measured in a blood test. x_{10} gives the ratio of Albumin to Globulin in the blood.

$$\tilde{X} = U \Sigma V^T$$

$$\Sigma = \begin{bmatrix} 40.1 & & & & 6 \\ & 34.2 & & & \\ & & 28.1 & & \\ & & & 24.8 & \end{bmatrix}$$



A PCA analysis is applied to the standardized data based on the attributes x_1-x_{10} . The squared Frobenius norm of the standardized data matrix \tilde{X} is given by $\|\tilde{X}\|_F^2 = 5780.0$. The first four singular values are $\sigma_1 = 40.1$, $\sigma_2 = 34.2$, $\sigma_3 = 28.1$, and $\sigma_4 = 24.8$. Which of the following statements is correct?

$$\frac{\sigma_i^2}{\sum_{j=1}^n \sigma_j^2} = \frac{\sigma_i^2}{\|\tilde{X}\|_F^2}$$

- A. The first PCA component accounts for more than 35 % of the variation.
- B. The second PCA component accounts for more than 30 % of the variation.
- C. The first three PCA components account for less than 70 % of the variation in the data.
- D. The fourth PCA component accounts for less than 10 % of the variation in the data.
- E. Don't know.

Solution:

The i^{th} principal component accounts for $\frac{\sigma_i^2}{\sum_j \sigma_j^2} = \frac{\sigma_i^2}{\|X\|_F^2}$. We therefore have that the first PCA component accounts for $\frac{40.1^2}{5780.0} = 27.8\%$, the second $\frac{34.2^2}{5780.0} = 20.2\%$, and the first three principal components ac-

count for $\frac{40.1^2 + 34.2^2 + 28.1^2}{5780.0} = 61.7\%$ of the variation whereas the fourth principal component accounts for $\frac{24.8^2}{5780.0} = 10.6\%$. Thus, the first three PCA components account for less than 70% of the variation in the data.

Fishers Iris Data

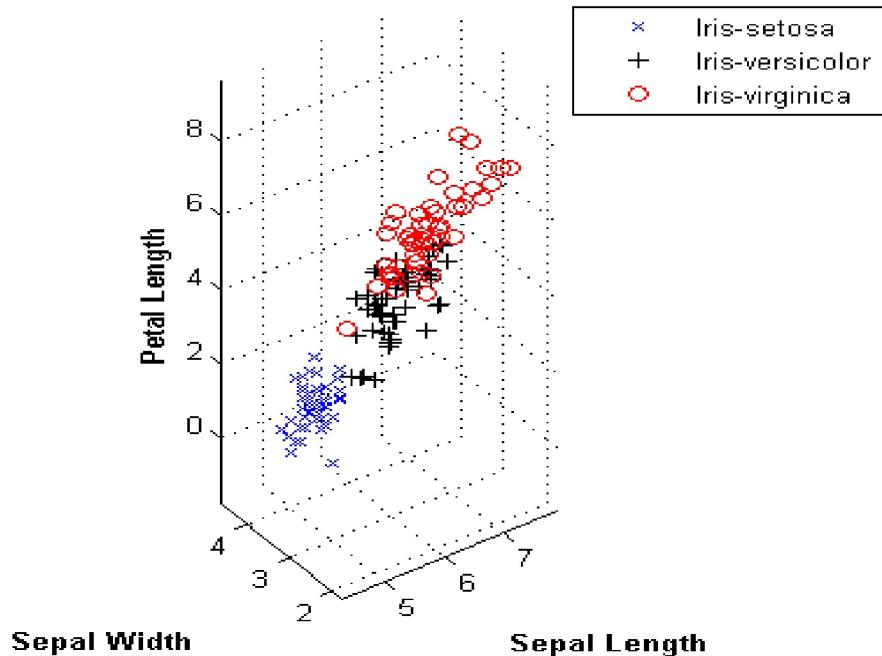


**Three types of flowers:
Iris Setosa, Iris Versicolor, Iris Virginica**

Flower ID	Attribute				Petal Width
	Sepal Length	Sepal Width	Petal Length	Petal Width	
1	5.1	3.5	1.4	0.2	
2	4.9	3.0	1.4	0.2	
3	4.7	3.2	1.3	0.2	
4	4.6	3.1	1.5	0.2	
.	
.	
150	5.9	3.0	5.1	1.8	

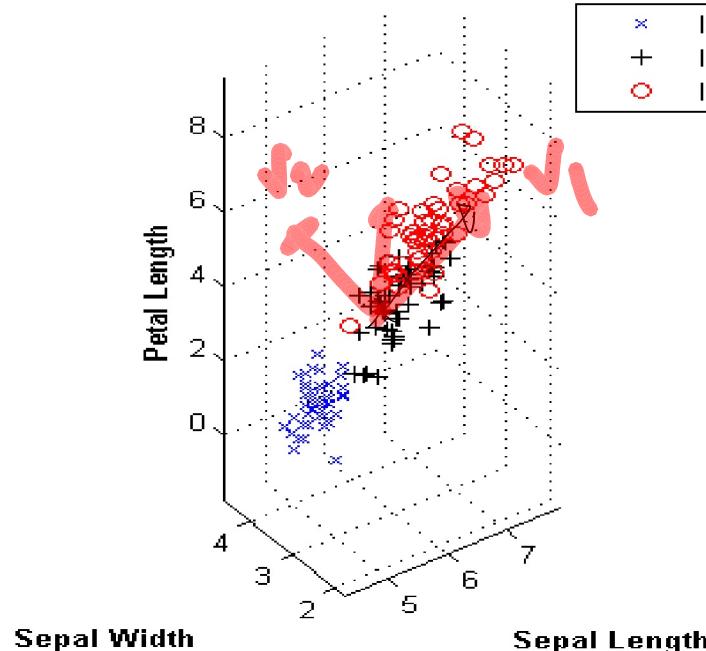
We will presently consider the first 3 attributes, i.e. Sepal length, Sepal Width and Petal Length.

3D scatter plot of Iris Data



What fraction of the total variation in the data will the first principal component account for?

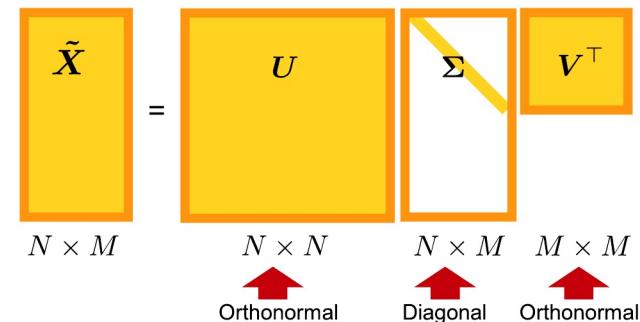
3D scatter plot of Iris Data



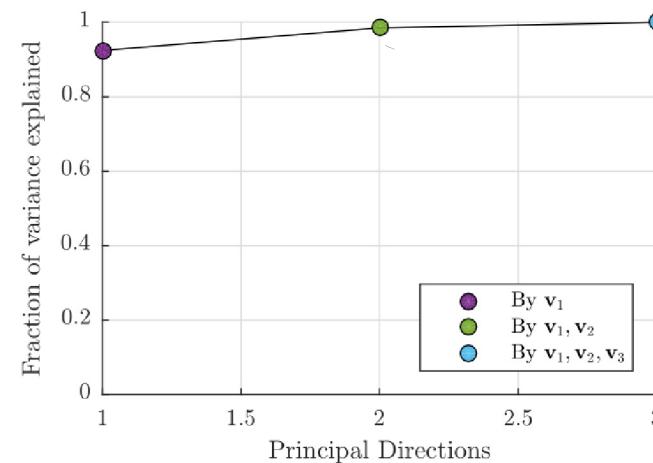
What fraction of the total variation in the data will the first principal component account for?

- 1) Subtract the mean
- 2) Apply singular value decomposition (SVD)

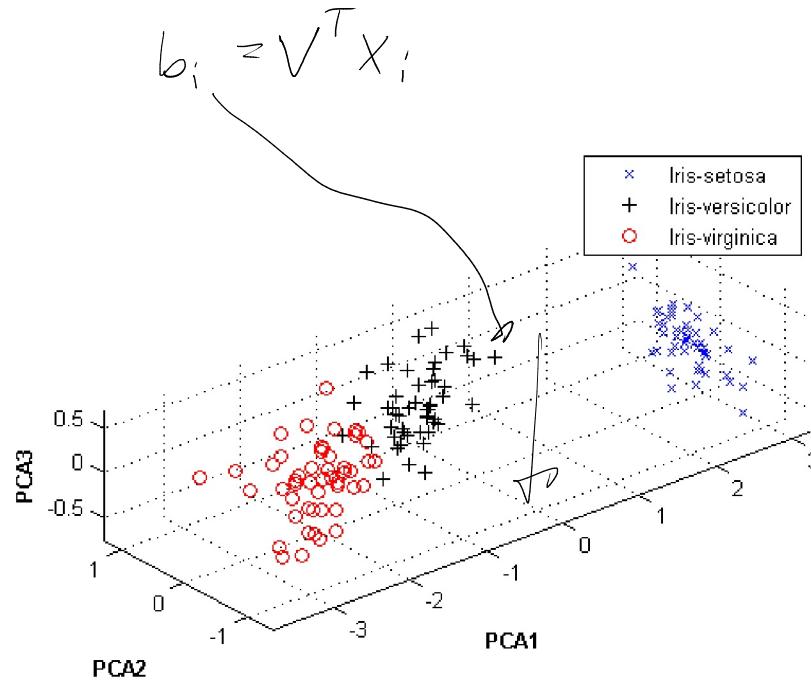
$$\tilde{\mathbf{X}} = \mathbf{U} \Sigma \mathbf{V}^\top$$



Evaluate the singular values to determine how much of the dynamics is lost when reducing the dimensionality



Visualization of the PCA projections of the data



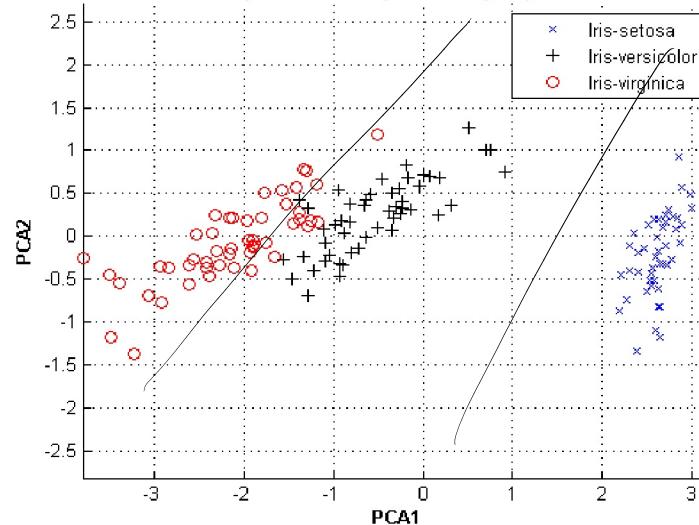
$$\tilde{X} = U\Sigma V^\top$$

\tilde{X}	U	Σ	V^\top
$N \times M$	$N \times N$	$N \times M$	$M \times M$
Orthonormal	Diagonal	Orthonormal	

$$PCA1: b_1 = \tilde{X}v_1 = u_1\sigma_1$$

$$PCA2: b_2 = \tilde{X}v_2 = u_2\sigma_2$$

$$PCA3: b_3 = \tilde{X}v_3 = u_3\sigma_3$$



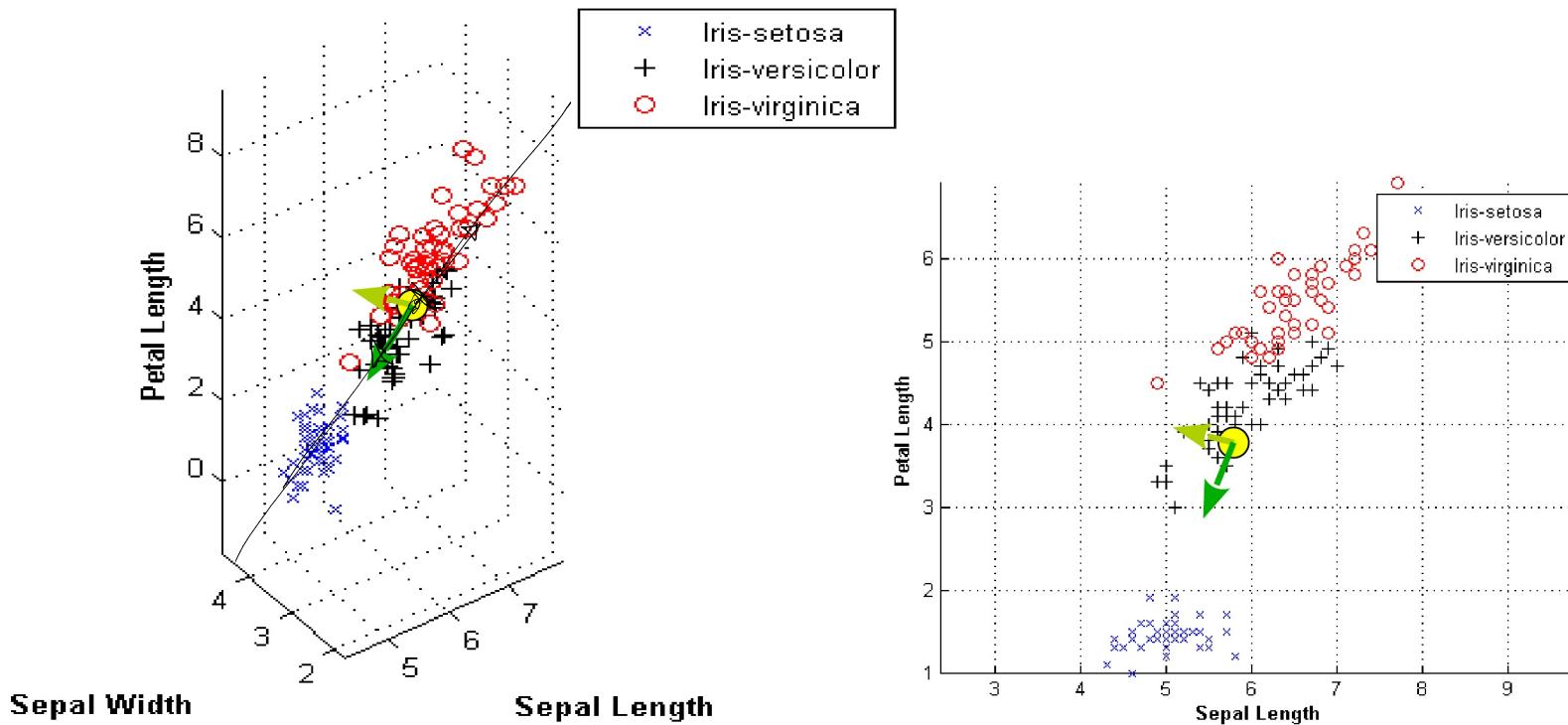
The principal directions V

Sepal Length
Sepal Width
Petal Length

$$V = \begin{bmatrix} -0.39 & -0.64 & -0.66 \\ 0.09 & -0.74 & 0.66 \\ -0.92 & 0.20 & 0.35 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 5.8 \\ 3.1 \\ 3.8 \end{bmatrix}, \quad \vec{v}_1 = \begin{bmatrix} -0.39 \\ 0.09 \\ -0.92 \end{bmatrix}, \quad \vec{v}_2 = \begin{bmatrix} -0.64 \\ -0.74 \\ 0.20 \end{bmatrix}$$

Sepal Length
 Sepal Width
 Petal Length



Quiz 3: PCA Cont. (Fall 2012)

No.	Attribute description	Abbrev.
x_1	Age (in years)	AGE
x_2	Gender (Female=0, Male=1)	GDR
x_3	Total Bilirubin	TB
x_4	Direct Bilirubin	DB
x_5	Alkaline Phosphatase	AP
x_6	Alamine Aminotransferase	AlA
x_7	Aspartate Aminotransferase	AsA
x_8	Total Proteins	TP
x_9	Albumin	AB
x_{10}	Albumin to Globulin ratio	A/G
y	0=No liver disease, 1=Liver disease	LD

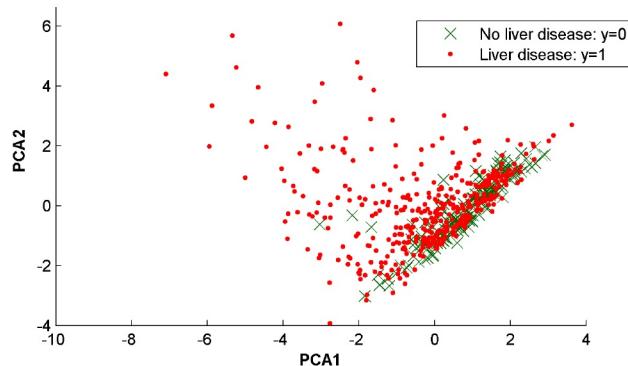


Figure 1: Principal component 1 (PCA1) plotted against principal component 2 (PCA2).

Question 1. The first and second principal compo-

nent directions of the liver-dataset are

$$\mathbf{v}_1 = \begin{bmatrix} -0.1404 \\ -0.1090 \\ -0.4115 \\ -0.4179 \\ -0.2468 \\ -0.2682 \\ -0.3009 \\ 0.2781 \\ 0.4375 \\ 0.3638 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} -0.2859 \\ 0.0130 \\ 0.2510 \\ 0.2622 \\ 0.0525 \\ 0.4162 \\ 0.3927 \\ 0.4197 \\ 0.4323 \\ 0.3052 \end{bmatrix}.$$

In the figure, the data projected onto the first two principal components is plotted, and the colors indicate the presence of liver disease. Which of the following statements is *correct*?

- A. Relatively high values of AGE, GDR, TB, DB, AP, AlA, and AsA and low values of TP, AB, and A/G will result in a positive projection onto the first principal component.
- B. Relatively low values of the projection onto PCA1 and high values of the projection onto PCA2 indicates the subject does not have a liver disease.
- C. PCA2 mainly discriminate between old subjects with low measurements of TB, DB, AlA, AsA, TP, AB, and A/G from young subjects with high values of TB, DB, AlA, AsA, TP, AB, and A/G.
- D. The principal component directions are not guaranteed to be orthogonal to each other since the data has been standardized.

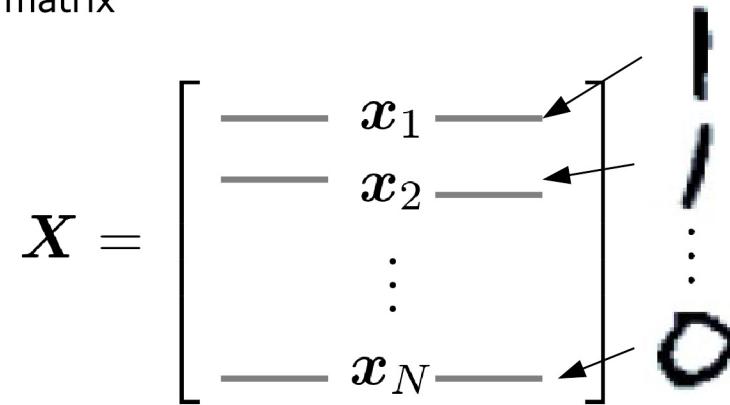
Solution:

AGE, GDR, TB, DB, AP, AlA, and AsA have negative coefficients of PCA1 whereas TP, AB, and A/G have positive coefficients resulting in a negative projection onto the first principal component, thus this is correct. From the figure we observe that observations with low values of PCA1 and high values of PCA2 in general have a red dot meaning they have a liver disease. For PCA2 we observe that GDR has a negative value whereas the remaining entities

have positive values while GDR and AP have small amplitudes. As a result PCA2 mainly discriminate between young subjects with high measurements of TD, DB, AlA, AsA, TP, AB, and A/G from old subjects with low values of TD, DB, AlA, AsA, TP, AB, and A/G hence this is correct. The principal component directions are always orthogonal to each other irrespective of the data preprocessing.

Visualization of hand written digits

- Data matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_N \end{bmatrix}$$


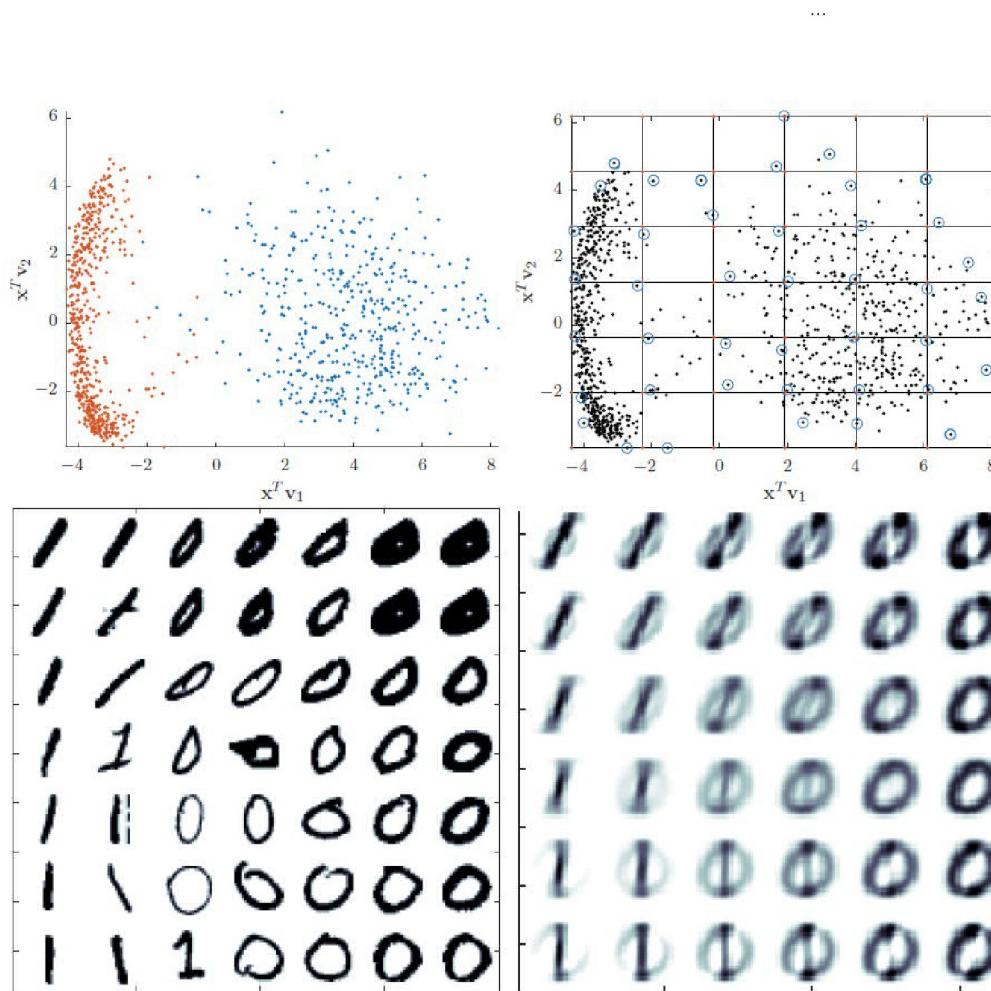
If each image is 28×28 pixels then \mathbf{X} is a $N \times 784$ matrix

- Principal component analysis

$$\tilde{\mathbf{X}} = \mathbf{U} \Sigma \mathbf{V}^\top$$

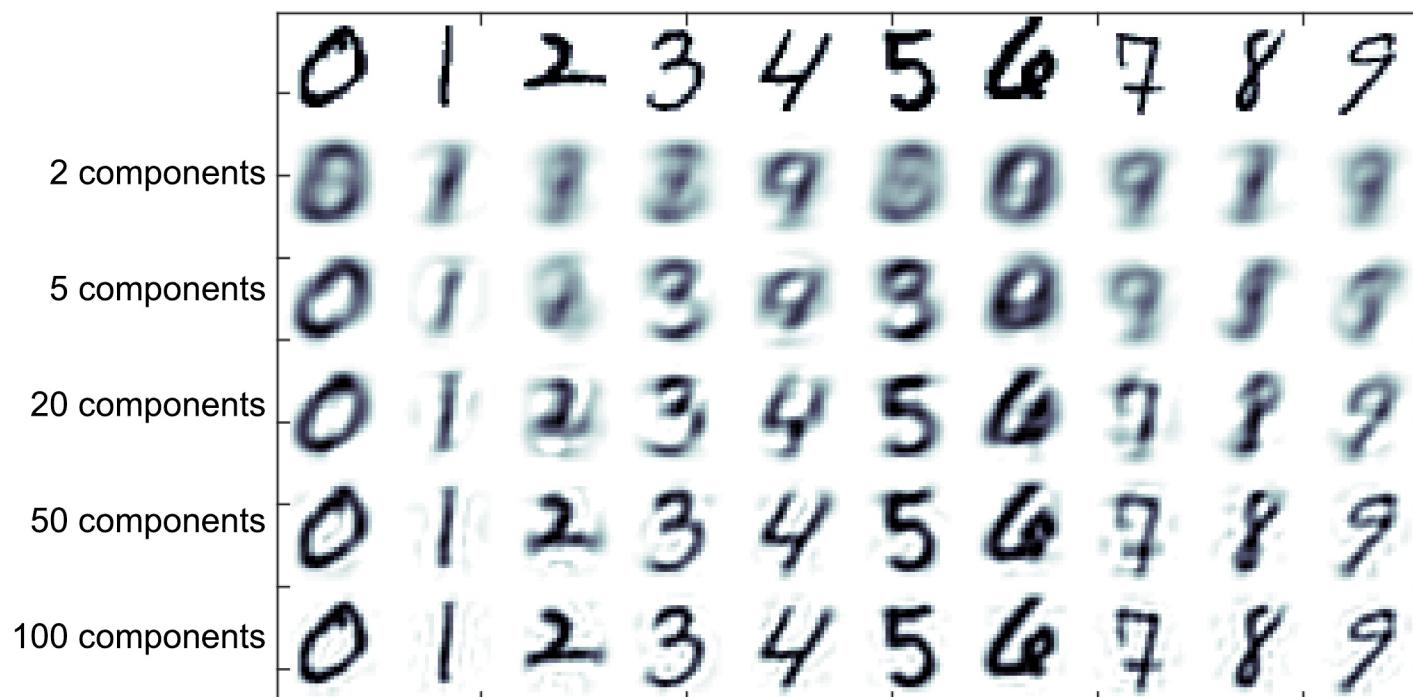
$\tilde{\mathbf{X}}$	$=$	\mathbf{U}	Σ	\mathbf{V}^\top
$N \times M$		$N \times N$	Diagonal	$M \times M$
		Orthonormal		Orthonormal

Visualization of hand written digits



PCA as compression

Only include a few components: $\hat{x}_i = Vb + m$ n=2,5,20,50,100



Data and domain driven feature extraction

PCA is an example of a data driven approach for feature extraction

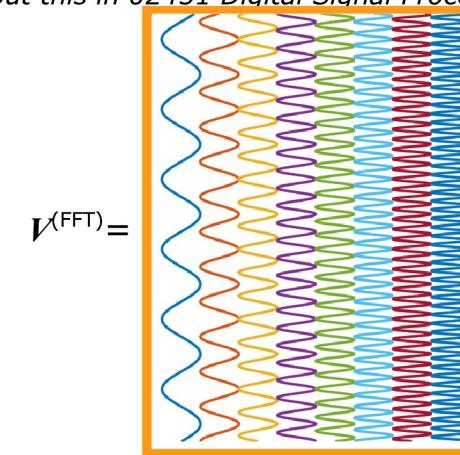
i.e., we define from data the features extracted in terms of the projections $V^{(PCA)}$ that preserve most of the variance in the data

$$\tilde{X} = U \Sigma V^\top$$

\tilde{X} U Σ V^\top
 $N \times M$ $N \times N$ $N \times M$ $M \times M$

The fourier transform is an example of a domain driven approach for feature extraction

i.e., in the analysis of sound good features are often to use spectral representations. These can be found by projecting the data using the so-called fourier transform matrix $V^{(FFT)}$ where the components are defined as specific frequencies such that the projection of the data onto these frequencies defines the extend to which these frequencies are present in the data. (you can learn much more about this in 02451 Digital Signal Processing)



Resources

<http://www2.imm.dtu.dk> Our online PCA demo which highlights key concepts of PCA such as the effect of normalization, variance explained, and much more (<http://www2.imm.dtu.dk/courses/02450/DemoPCA.html>)

<https://arxiv.org> A great and more in-depth tutorial on PCA
(<https://arxiv.org/abs/1404.1100>)

<https://www.3blue1brown.com> An great, animated recap of linear algebra
(<https://www.3blue1brown.com/essence-of-linear-algebra-page/>)