

Motivations

Why?
"Necessity is the mother of invention"

Explosive Growth of Data
Pressing need for the automated analysis of massive data

Emerged in the late 1980s
Major developments in the mid 1990s

Several names over the years
Based on: Statistics, Machine Learning, Database technology

Evolution of Technology

- 1960s: data collection, database creation, & network DBMS
- 1970s: relational data model, relational DBMS implementation
- 1980s: RDBMS, advanced data models (extended-relational, OO, deductive, etc.); application-oriented DBMS (spatial, scientific, engineering, etc.)
- 1990s: data mining, data warehousing, multimedia databases, and Web databases
- 2000s: stream data management and mining, web technology (XML, data integration), global information systems
- 2010s: social networks, NoSQL, unstructured data, etc.
- 2020s: ...

Examples

- **Customer attrition**
 - Given customer information for the past months
 - Predict who is likely to attrite next month, or estimate customer value
- **Credit assessment**
 - Given a loan application
 - Predict whether the bank should approve the loan
- **Customer segmentation**
 - Given several information about the customers
 - Identify interesting groups among them
- **Community detection**
 - Who is discussing what?
- **Sentiment Analysis**
 - ...



Big Data?

Massive datasets

Gigantic quantities of data are being collected in every domain imaginable: science, agriculture, healthcare, transport, sport ...

Exponential increases in computing power has led to ever decreasing costs of instrumentalization and storage

Big Data

The quantity of data or quantity is so large that

- (1) it facilitates predictions not previously possible, and
- (2) it necessitates special processes to deal with

Quantity could be total storage, arrival rate, ...

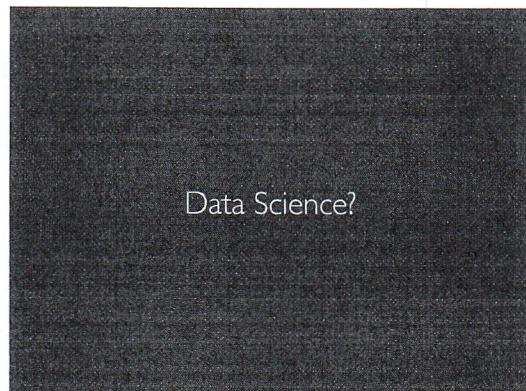
Machine Learning

"A computer program is said to learn from experience E with respect to some class of task T and a performance measure P, if its performance at tasks in T, as measured by P, improves because of experience E."

Machine Learning Paradigms 12

- Suppose we have the experience E encoded as a dataset,
 $D = x_1, x_2, x_3, \dots, x_N$
- Supervised Learning**
 - Given the desired outputs $t_1, t_2, t_3, \dots, t_N$ learns to produce the correct output given a new set of input
- Unsupervised learning**
 - Exploits regularities in D to build a representation to be used for reasoning or prediction
- Reinforcement learning**
 - Producing actions $a_1, a_2, a_3, \dots, a_N$ which affect the environment, and receiving rewards $r_1, r_2, r_3, \dots, r_N$ learn to act in order to maximize rewards in the long term

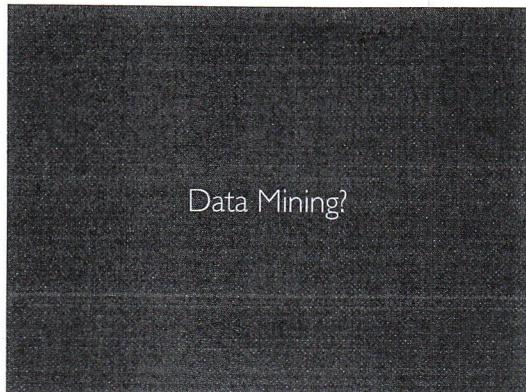
Prof. Pierluca Lanzi POLITECNICO DI MILANO



What Is data science?

Missing data
Statistics
Interpret
Leverage

<http://wikibon.org/blog/role-of-the-data-scientist/>



Data Mining

The non-trivial process of identifying
 (1) valid, (2) novel, (3) potentially useful,
 and (4) understandable patterns in data.

features

target variable

goal: given the features
 assign the right target

An Example Using Contact Lens Data

Age	Spectacle prescription	Astigmatism	Tear production rate	Recommended lenses
Young	Myope	No	Reduced	None
Young	Myope	No	Normal	Soft
Young	Myope	Yes	Reduced	None
Young	Hypermetropic	No	Normal	Hard
Young	Hypermetropic	No	Reduced	None
Young	Hypermetropic	No	Normal	Soft
Young	Hypermetropic	Yes	Reduced	None
Young	Hypermetropic	Yes	Normal	Hard
Pre-presbyopic	Myope	No	Normal	None
Pre-presbyopic	Myope	No	Normal	Soft
Pre-presbyopic	Myope	Yes	Reduced	None
Pre-presbyopic	Myope	Yes	Normal	Hard
Pre-presbyopic	Hypermetropic	No	Reduced	None
Pre-presbyopic	Hypermetropic	No	Normal	Soft
Pre-presbyopic	Hypermetropic	Yes	Reduced	None
Pre-presbyopic	Hypermetropic	Yes	Normal	None
Presbyopic	Myope	No	Reduced	None
Presbyopic	Myope	No	Normal	None
Presbyopic	Myope	Yes	Reduced	None
Presbyopic	Myope	Yes	Normal	Hard
Presbyopic	Hypermetropic	No	Reduced	None
Presbyopic	Hypermetropic	No	Normal	Soft
Presbyopic	Hypermetropic	Yes	Reduced	None
Presbyopic	Hypermetropic	Yes	Normal	None

An example of possible pattern

if astigmatism = yes
 and tear production rate = normal
 and spectacle prescription = myope
 then recommendation = hard

Is it a "Good" Pattern?

- Is it valid?
 - The pattern has to be valid with respect to a certainty level (rule true for the 86%)
- Is it novel?
 - Is the relation between astigmatism and hard contact lenses already well-known?
- Is it useful? Is it actionable?
 - The pattern should provide information useful to the bank for assessing credit risk
- Is it understandable?

What is the General Idea?

- Build computer programs that navigate through databases automatically, seeking regularities or patterns
- However ...
 - Most patterns are uninteresting
 - Most patterns are spurious, inexact, or contingent on accidental coincidences in the particular dataset used
 - Real data is imperfect: Some parts will be garbled, and some will be missing
- Algorithms need to be robust enough to cope with imperfect data and to extract regularities that are inexact but useful

Descriptive vs. Predictive

Are the models built for gaining insight?
(about what already happened)

→ descriptive
data mining task

Or are they built for accurate prediction?
(about what might happen)

→ predictive
data mining task

I want to analyze who are the customers that shop at my store more than twice a week

→ descriptive

I need a model to help me investing in the stock market!

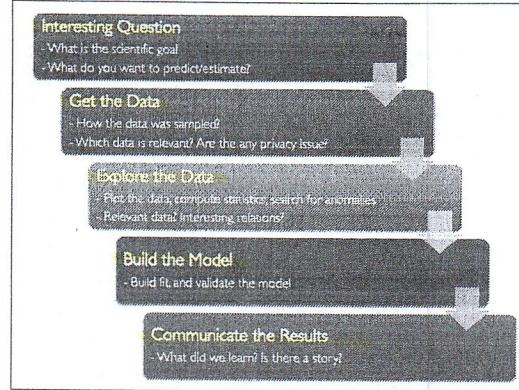
→ predictive

There is a 3rd term :

"Prescriptive"

use descriptive and predictive mining to recommend a course of action
~~using descriptive and predictive data mining techniques we have to give some recommendations on what to do next~~

The Process

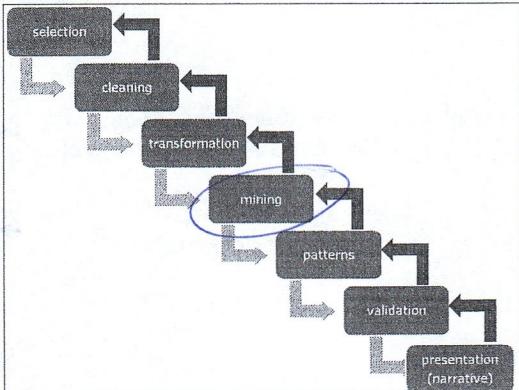


! When we present the results:

What did we learn?

Is there a story?

Can we build a narrative?



What Are the Main Steps?

28

- **Selection**
 - What are data we actually need to answer the posed question?
- **Cleaning**
 - Are there any errors or inconsistencies in the data we need to eliminate?
- **Transformation**
 - Some variables might be eliminated because equivalent to others
 - Some variables might be elaborated to create new variables (e.g. birthday to age, daily measures into weekly/monthly measures, log?)
- **Mining**
 - Select the mining approach: classification, regression, association, etc.
 - Choose and apply the mining algorithm(s)
- **Validation**
 - Are the patterns we discovered sound? According to what criteria?
 - Are the criteria sound? Can we explain the result?
- **Presentation & Narrative**
 - What did we learn? Is there a story to tell? A take-home message?

Prof. Perugini and ...

POLITECNICO DI MILANO

Data Mining Tasks

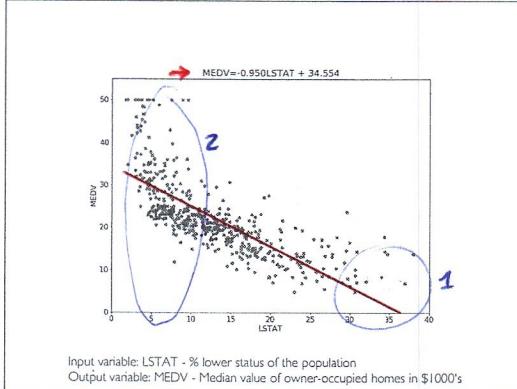
Prediction & Regression

Boston Housing Dataset

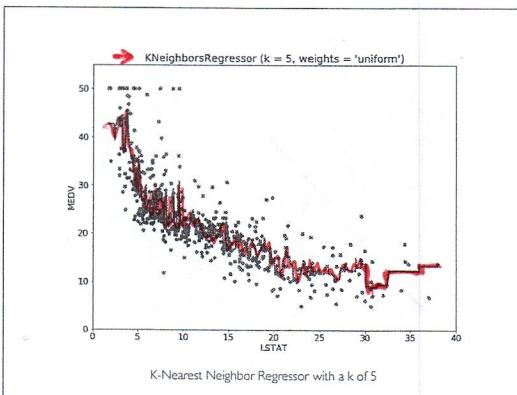
Goal: how much a house was valued in Boston

- Information concerning housing in the area of Boston (MA)
 - collected by U.S Census Service
 - <https://archive.ics.uci.edu/ml/machine-learning-databases/housing/>
- 506 cases and 14 variables:
 - CRIM - per capita crime rate by town;
 - CHAS - Charles River dummy variable (1 if tract bounds river; 0 otherwise);
 - TAX - full-value property-tax rate per \$10,000;
 - PTRATIO - pupil-teacher ratio by town;
 - LSTAT - % lower status of the population;
 - ...
 - MEDV - Median value of owner-occupied homes in \$1000's
- MEDV is the target variable

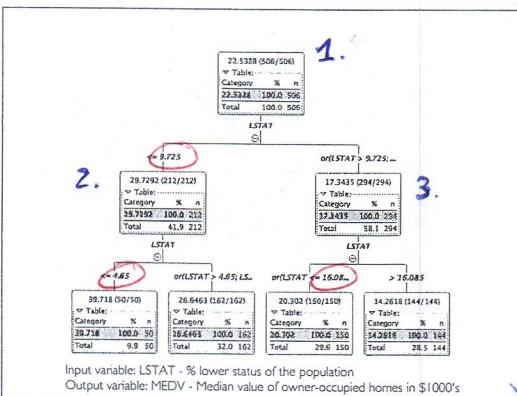
Prof. Pierluca Lanzi POLITECNICO DI MILANO



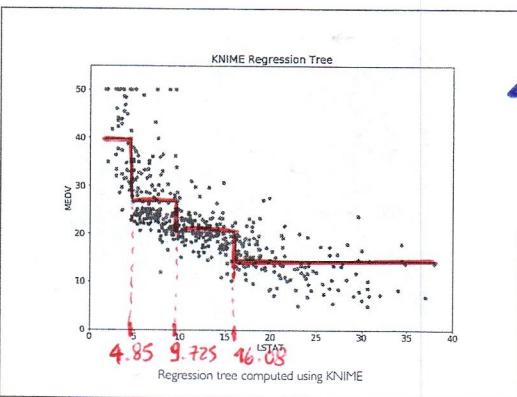
What about the narrative of this solution?
LSTAT is the percentage (%) of the lower status of the population. We can say that if there are several poor people (1) the value of the house is very low. If the percentage of poor people is very low (2) we're in a rich part of the town and the houses have high values.



Another approach is to determine the price of a house basing on the houses in the neighborhood. Here the prediction is build as an average of the 5 nearest points to where we are. If this approach works better we can say that similarities is an important aspect of the problem domain: "if two houses are similar we have more informations (better informations) than just the linear relation".



Another approach: regression tree. It's a tree for which at each node we make a decision. At the beginning there is a root (1). The average value is 22.5328. How can we have a better prediction? We introduce a threshold: if the LSTAT is below 9.725 then the average value is 28.7292, otherwise the average value is 17.3435. By splitting we find subspaces where the results differ a lot. After the first split we can go on splitting. We can split 2. and 3. in two each.



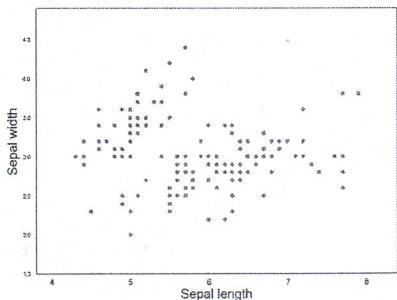
These splits generate a model (piecewise linear model)

What if this works better than the others? We can say that there's a linear relation but it's not as powerful as understanding what are the subspaces in which we have to divide data.

Classification

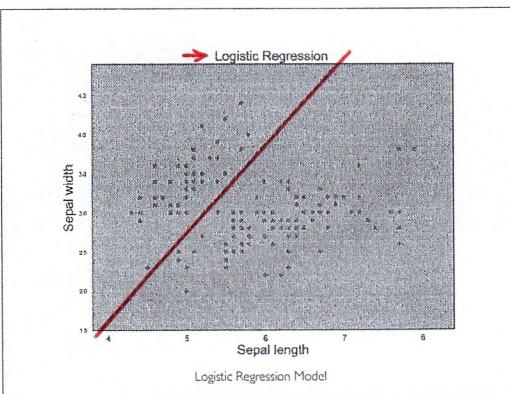
: if we don't have to predict a number but a label

We have to predict a color:
some of the points are green and
some are blue. We need to build
a model that allows to say: "If
your datum is here then it'll be blue/green".
We want a model that is accurate
in predicting what we know and
accurate in predicting what we don't
know (new points).

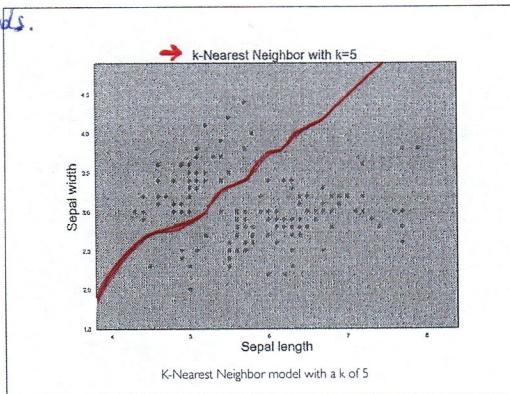


Iris dataset plotted using only two variables and two classes.

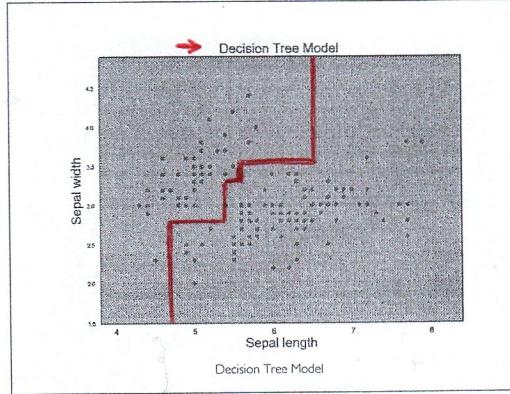
Easiest model: linear model.



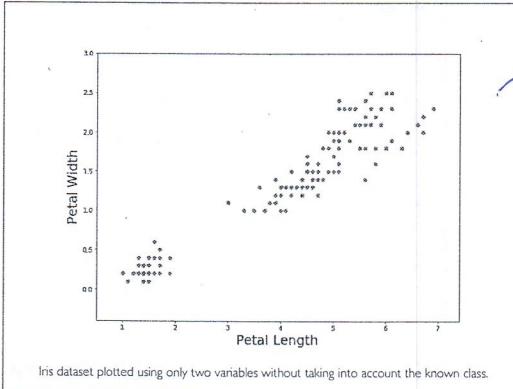
Another option: k nearest neighborhoods.
To assign a class to a point we look
for the k nearest point and the
assigned class will be the most
frequent among the k points.



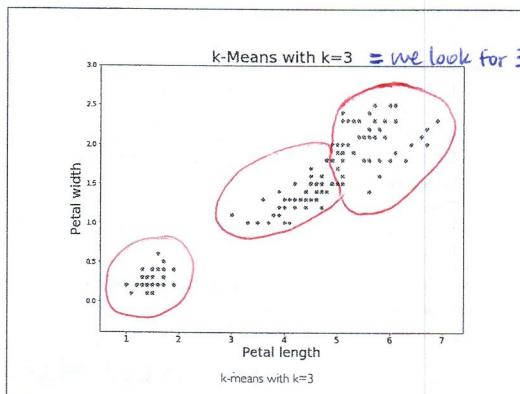
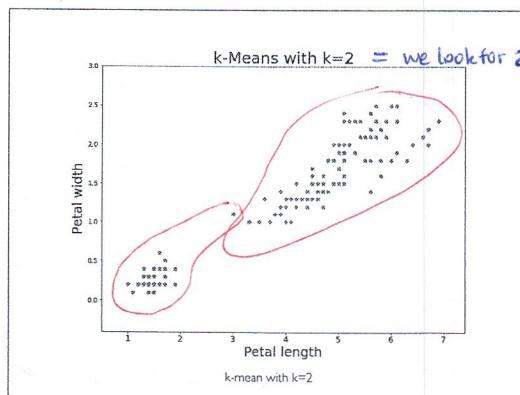
Another option: decision tree model.
The decision tree model considers
hyperplanes.



Clustering : we look at the data and we try to find patterns and analogies



The assumptions we make if the points are nearby then they're very similar. In the PROBLEM SPACE the distance between points is a similarity index.



Associations

Data Mining Tasks: Associations 46

Receives after buying stuff:

- {citrus fruit, semi-finished bread, margarine, ready soups}
- {tropical fruit, yogurt, coffee}
- {whole milk}
- {pip fruit, yogurt, cream cheese, meat spreads}
- {other vegetables, whole milk, condensed milk, long life bakery product}
- {whole milk, butter, yogurt, rice, abrasive cleaner}
- {other vegetables, UHT-milk, rolls/buns, bottled beer, liquor (appetizer)}
- {whole milk, cereals}
- ...

Are there interesting products associations?

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Frequent Itemsets (Minimum Support 0.1%) 47

items	support	count
[1] {hair spray}	0.00118454	11
[2] {flower soil/fertilizer}	0.001931876	19
[3] {rubbing alcohol}	0.001016777	10
[4] {frozen fruits}	0.001220132	12
[5] {prosecco}	0.002033554	20
[6] {honey}	0.001525165	15
[7] {cream}	0.001321810	13
[8] {decalcifier}	0.001525165	15
[9] {organic products}	0.001626843	16
[10] {soap}	0.002643620	26
...		

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Association Rules (Support, Confidence, Lift) 48

```
{rice, sugar} => {whole milk}
Support 0.001220132
Confidence 1
Lift 3.913649

{canned fish, hygiene articles} => {whole milk}
Support 0.001118454
Confidence 1
Lift 3.913649

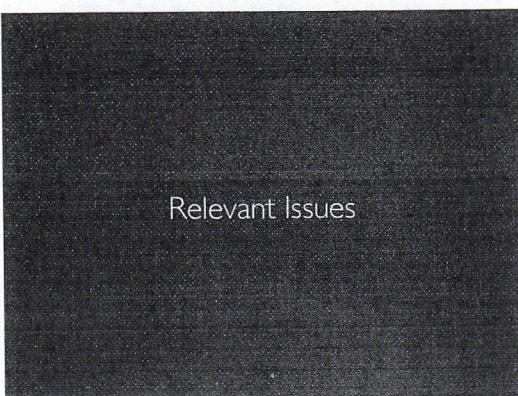
{root vegetables, butter, rice} => {whole milk}
Support 0.001016777
Confidence 1
Lift 3.913649

...

```

Prof. Pierluca Lanzi POLITECNICO DI MILANO

- Other Tasks 49
- **Outlier analysis**
 - An outlier is a data object that does not comply with the general behavior of the data
 - It can be considered as noise or exception but is quite useful in fraud detection, rare events analysis
 - **Trend and evolution analysis**
 - Trend and deviation, regression analysis
 - Sequential pattern mining, periodicity analysis
 - Similarity-based analysis
 - **Text Mining, Topic Modeling, Graph Mining, Data Streams**
 - **Sentiment Analysis, Opinion Mining, etc.**
 - **Other pattern-directed or statistical analyses**
- Prof. Pierluca Lanzi POLITECNICO DI MILANO



Are all the "Discovered" Patterns Interesting?

51

- Data Mining may generate thousands of patterns, but typically not all of them are interesting.
- **Interestingness measures**
 - A pattern is interesting if it is easily understood by humans, valid on new or test data with some degree of certainty, potentially useful, novel, or validates some hypothesis that a user seeks to confirm
- **Objective vs. subjective interestingness measures**
 - Objective measures are based on statistics and structures of patterns
 - Subjective measures are based on user's belief in the data, e.g., unexpectedness, novelty, etc.

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Can We Find All and Only Interesting Patterns?

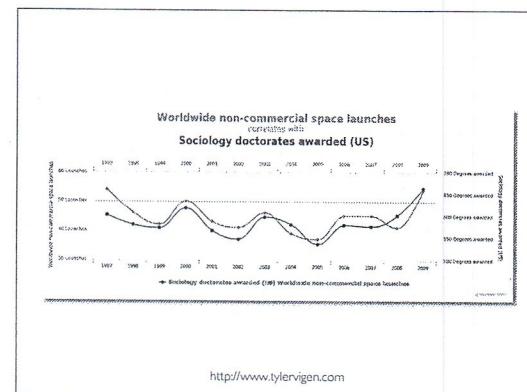
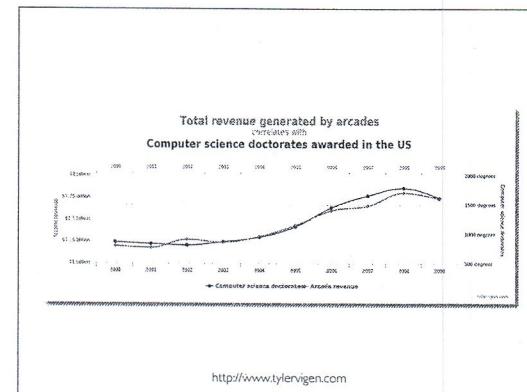
52

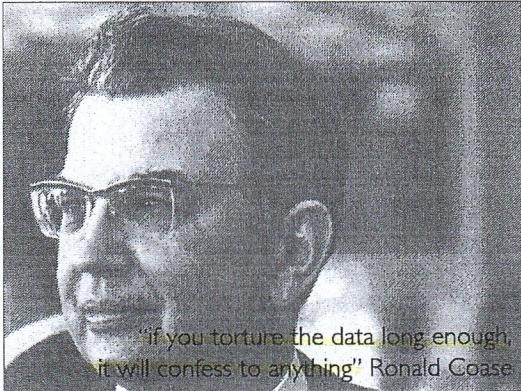
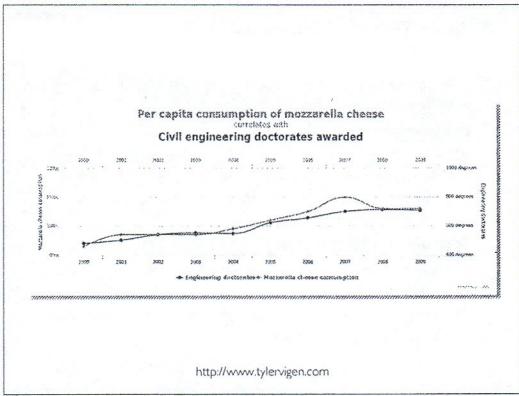
- **Completeness**
 - Find all the interesting patterns
 - Can a data mining system find all the interesting patterns?
 - Association vs. classification vs. clustering
- **Optimization**
 - Search for only interesting patterns:
 - Can a data mining system find only the interesting patterns?
 - Two approaches: (1) first general all the patterns and then filter out the uninteresting ones; (2) generate only the interesting patterns—mining query optimization

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

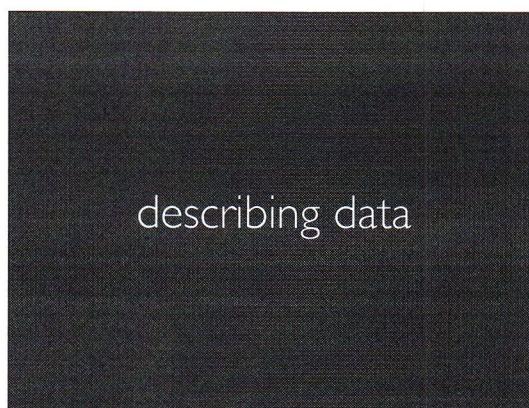
Pitfalls





Data Representation
Data Mining and Text Mining

Prof. Perluca Lanzi



The Weather Dataset

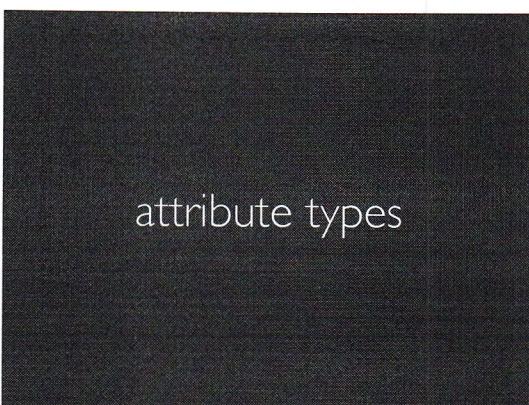
3

label

Outlook	Temperature	Humidity	Windy	Play
Sunny	85	85	False	no
Sunny	80	90	True	no
Overcast	83	86	False	yes
Rainy	70	96	False	yes
Rainy	68	80	False	yes
Rainy	65	70	True	no
Overcast	64	65	True	yes
Sunny	72	95	False	no
Sunny	69	70	False	yes
Rainy	75	80	False	yes
Sunny	75	70	True	yes
Overcast	72	90	True	yes
Overcast	81	75	False	yes
Rainy	71	91	True	no

Prof. Perluca Lanzi POLITECNICO DI MILANO

- Instances, Attributes, Concepts
- 4
- **Instances (aka observations, cases, records, items, examples)**
 - The atomic elements of information from a dataset
 - Each row in previous table corresponds to an instance
 - **Attributes (aka variables, features, independent variables)**
 - Measures aspects of an instance
 - Each instance is composed of a certain number of attributes
 - Each column in previous table contains values of an attribute
 - **Concept (aka class, target variable, dependent variable)**
 - Special content inside the data
 - Kind of things that can be learned
 - Intelligible and operational concept description
 - Last column of previous table was the class
- Prof. Perluca Lanzi POLITECNICO DI MILANO



Attributes

• Numeric Attributes

- Real-valued or integer-valued domain
- Interval-scaled when only differences are meaningful
(e.g., temperature)
- Ratio-scaled when differences and ratios are meaningful
(e.g., Age)

• Categorical Attributes

- Set-valued domain composed of a set of symbols
- Nominal when only equality is meaningful
(e.g., domain(Sex) = { M, F })
- Ordinal when both equality (are two values the same?) and inequality (is one value less than another?) are meaningful
(e.g., domain(Education) = { High School, BS, MS, PhD })

Prof Pierluca Lanzi

POLITECNICO DI MILANO

Numerical Attributes

- Not only ordered but measured in fixed and equal units
- Examples
 - Attribute "temperature" expressed in degrees
 - Attribute "year"
- Characteristics
 - Difference of two values makes sense
 - Sum or product doesn't make sense
 - Zero point is not defined
- Sometimes they are divided into "discrete" and "continuous"

Prof Pierluca Lanzi

POLITECNICO DI MILANO

Nominal Attributes (or Categorical)

- Values are distinct symbols that serve only as labels or names
- Example
 - Attribute "outlook" from weather data
 - Values: "sunny", "overcast", and "rainy"
- Characteristics
 - No relation is implied among nominal values
 - No ordering
 - No distance measure
 - Only equality tests can be performed

Prof Pierluca Lanzi

POLITECNICO DI MILANO

Other Types of Attributes

- **Ordinal Attributes** → basically : categorical + ordering
 - Categorical attributes with an imposed order on values
 - No distance between values defined
 - For instance, temperature encoded as "hot", "mild", and "cool"
 - Size encoded as "small", "medium", "large", and "jumbo"
- **Ratio Attributes**
 - Numerical attributes for which the measurement scheme defines a zero point (e.g., an attribute representing distance)
- **Binary Attributes**
 - Represented by just two values 0/1

Prof Pierluca Lanzi

POLITECNICO DI MILANO

example

Contraceptive Method Choice Data Set
https://archive.ics.uci.edu/ml/datasets/Contraceptive+Method+Choice

	Age	WifeEducation	HusbandEducation	NumberOfChildren	WifeWorking	HusbandWorking	HusbandEverWorked	StandardOfLiving	MediaExposure	MethodUsed
0	26	1	3	3	1	1	2	3	0	1
1	27	1	3	3	1	1	2	4	0	1
2	43	2	3	2	1	1	2	3	0	1
3	42	3	3	1	1	1	3	3	0	1
4	36	2	3	6	1	1	3	7	0	1
5	35	2	3	6	1	1	3	3	0	1
6	34	2	3	6	1	1	3	3	0	1
7	31	3	3	5	1	0	3	2	0	3
8	27	2	3	3	1	1	2	4	0	1
9	45	1	3	6	1	1	2	3	1	1
10	35	2	3	2	1	1	2	4	0	1
11	22	1	3	2	1	1	1	4	3	5
12	44	4	1	1	0	1	1	4	0	1
13	43	2	4	5	1	0	3	3	0	1
14	38	5	2	5	1	1	2	3	0	1
15	35	4	2	5	1	1	2	3	0	1
16	45	1	2	5	1	1	2	4	0	1
17	25	2	2	6	0	1	2	4	0	1
18	37	2	2	5	1	1	3	3	0	1
19	29	2	1	5	1	1	2	1	1	5

What are the attribute types of these variables?

Prof Perluca Lanzi

POLITECNICO DI MILANO

Contraceptive Method Choice Data Set (Data Description)

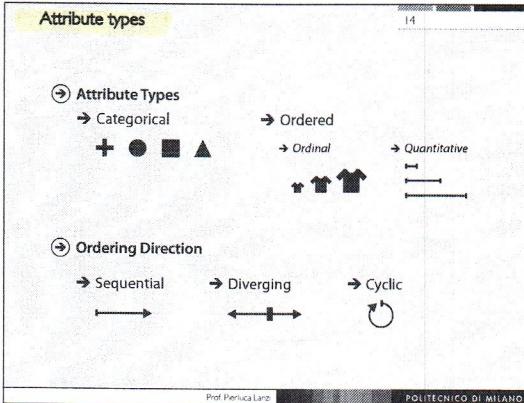
- Wife's age (numerical)
- Wife's education (categorical) 1=low, 2, 3, 4=high → it's actually a sorting (and to ordinal)
- Husband's education (categorical) 1=low, 2, 3, 4=high → same (↑)
- Number of children ever born (numerical)
- Wife's religion (binary) 0=Non-Islam, 1=Islam → it's actually nominal (categorical)
- Wife's now working? (binary) 0=Yes, 1=No → same (↑)
- Husband's occupation (categorical) 1, 2, 3, 4 → this is categorical (nor numerical, nor ordinal)
- Standard-of-living index (categorical) 1=low, 2, 3, 4=high → categorical
- Media exposure (binary) 0=Good, 1=Not good
- 10. Contraceptive method used (class attribute)
1=No-use, 2=Long-term, 3=Short-term

What are the attribute types of these variables (now)?

Prof Perluca Lanzi

POLITECNICO DI MILANO

another perspective



Prof Perluca Lanzi

POLITECNICO DI MILANO

Attribute types

④ Attribute Types

- Categorical
- + ● ■ ▲
- e.g., gender, race, eye color

→ Ordered

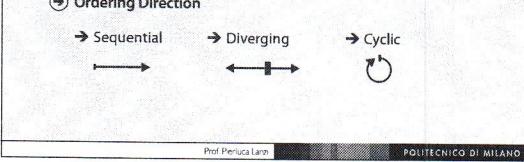
→ Ordinal

→ Quantitative

④ Ordering Direction

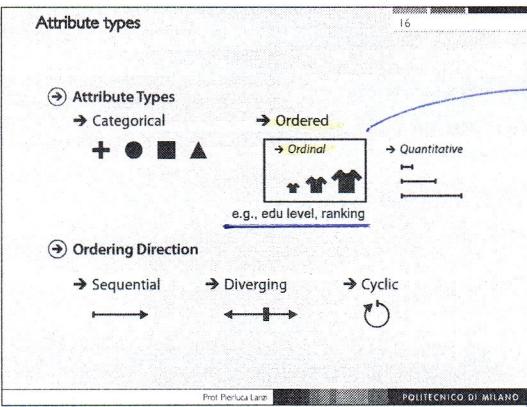
- Sequential
- Diverging
- Cyclic

having a categorical attribute is like having a shape: we just can show them

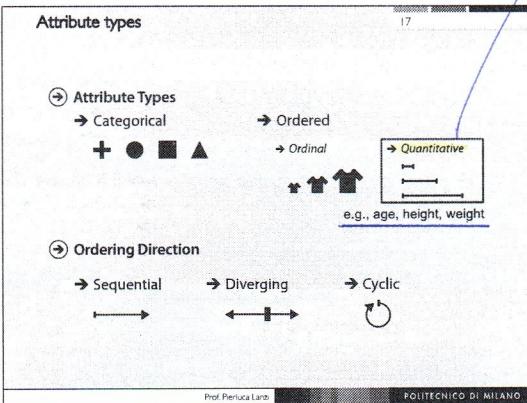


Prof Perluca Lanzi

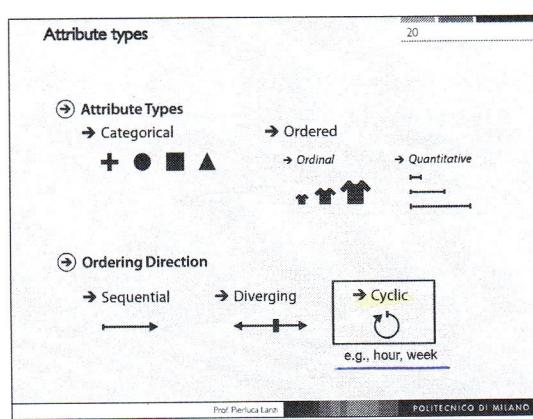
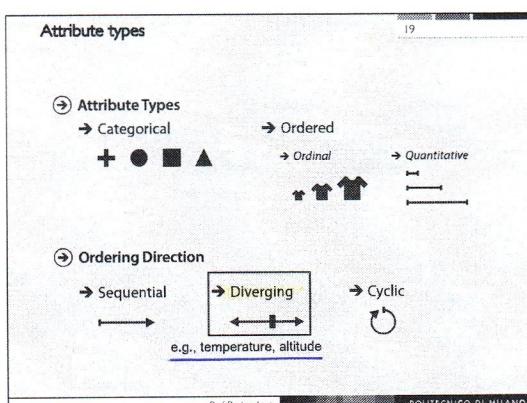
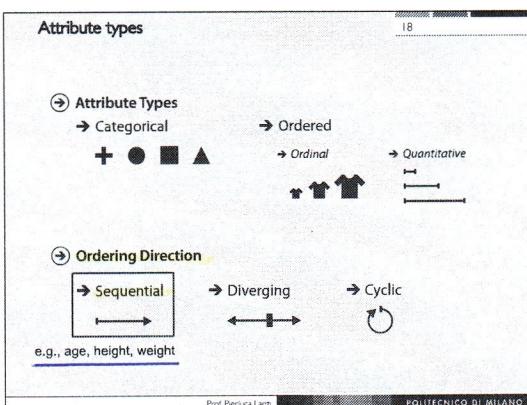
POLITECNICO DI MILANO



When we have ordinals (ordered) then we can sort them or create a representation that compares them (quantitative).



The fact that we can have our order and a quantitative variables, we can have different ordering directions:



hierarchies

day/month/year

Some attributes may have an internal hierarchical structure

For example, dates, mail addresses, spatial regions, taxonomies, etc.

who is writing
is a student

home.cognome @ mail.polimi.it
↓
more specific
from politecnico the mail
is from Italy

missing values

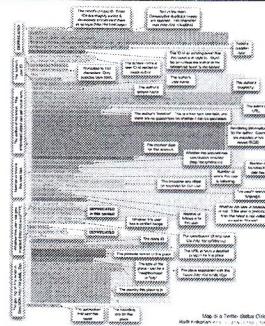
Why Missing Values Exist?

- Faulty equipment, incorrect measurements, missing cells in manual data entry, censored/anonymous data
 - Review scores for movies, books, etc.
 - Very frequent in questionnaires for medical scenarios
 - Censored or anonymous data
 - Data from rich representations for which many of the fields might not be used (e.g., Twitter data)

2

Prof. Pierluca Lanzì

POLITECNICO DI MILANO



- Missing value may have significance in itself
 - E.g. missing test in a medical examination or an empty field in a questionnaire
- They are frequently indicated by out-of-range entries (e.g. max/min float), NaN or special values* (e.g., zero)
- Most schemes assume that is not the case and "missing" may need to be coded as additional value
- Does absence of value have some significance?
 - If it does, "missing" is a separate value
 - If it does not, "missing" must be treated in a special way

What Types of Missing Values?

- Missing not at random (MNAR)
 - Distribution of missing values depends on missing value
 - E.g. respondents with high income less likely to report it
- Missing at random (MAR)
 - Distribution of missing values depends on observed attributes, but not missing value
 - E.g. men less likely than women to respond to question about mental health

Here it depends on the value
(if someone earns more then
it's more likely that this someone
is not going to tell).

Here it depends on the whole
column of the attribute: if
someone has a certain feature
(e.g. being male) then it's more
likely to not respond. So, basically
it depends on the type of the target,
not on some specific values of
the target.

What Types of Missing Values?

- Missing completely at random (MCAR)
 - Distribution of missing values does not depend on observed attributes or missing value
 - E.g. survey questions randomly sampled from larger set of possible questions
- Identifying MNAR and MAR can be difficult often requires domain knowledge

we consider all the instances with missing values
and we check if there are some dependencies.
of mental health questions we
notice that most of missing values
are men → it's an information)

Dealing with Missing Values

- Use what you know
 - Why data is missing
 - Distribution of missing data
- Decide on the best strategy to yield the least biased estimates
 - Deletion Methods (listwise deletion, pairwise deletion)
 - Single Imputation Methods (mean/mode substitution, dummy variable method, single regression)
 - Model-Based Methods (maximum Likelihood, multiple imputation)

(Possible issue:)

If the missing values are MNAR
we're losing informations. If we
eliminate the missing values and
it turns out that all people earning
a lot did not say it then we're
erasing a part of the population
from our analysis.

Strategies for missing values handling

- The handling of missing data depends on the type
- Discarding all the examples with a missing values
 - Simplest approach
 - Allows the use of unmodified data mining methods
 - Only practical if there are few examples with missing values. Otherwise, it can introduce bias
- Fill in the missing value manually ☺
 - Use a special value for it
 - Add an attribute that indicates if value is missing or not
 - Greatly increases the difficulty of the data mining process
- Convert the missing values into a new value
 - Use a special value for it
 - Assign a value to the missing one, based on the rest of the dataset. Use the unmodified data mining methods.
- Imputation methods

1.

Do Not Impute (DNI)

31

- Simply use the default policy of the data mining method
- Works only if the policy exists
- Some methods can work around missing data

Prof. Perluca Lanza

POLITECNICO DI MILANO

1.1.

List-wise Deletion (Complete Case Analysis)

32

- Only analyze cases with available data on each variable
- Simple, but reduces the data
- Comparability across analyses
- Does not use all the information
- Estimates may be biased if data not MCAR

we need massive amount of data

Gender	8th grade math test score	12th grade math score
F	45	—
M	—	99
F	55	86
F	85	88
F	80	75
—	81	82
F	75	80
M	95	—
M	86	90
F	70	75
F	—	85

delete all the instances with missing values

1.2.

Pairwise deletion (Available Case Analysis)

33

- Analysis with all cases in which the variables of interest are present
- Example
 - When using only the first two variables, the missing values of the third variable are not considered
- Advantage
 - Keeps as many cases as possible for each analysis
 - Uses all information possible with each analysis
- Disadvantage
 - Can't compare analyses because sample different each time

If we do 2 different analysis and we consider 2 different set of variables then the eliminated instances will be different in the 2 sets → non comparable analysis

Gender	8th grade math test score	12th grade math score
F	45	—
M	—	99
F	55	86
F	85	88
F	80	75
—	81	82
F	75	80
M	95	—
M	86	90
F	70	75
F	—	85

If we use only 2 attributes out of 3 then we consider missing values only in those 2 that we use.
(consider and delete the instances)

2.

Single Imputation Methods

34

- Mean/mode substitution (most common value)**
 - Replace missing value with sample mean or mode
 - Run analyses as if all complete cases
 - Advantages: Can use complete case analysis methods
 - Disadvantages: Reduces variability
- Dummy variable control**
 - Create an indicator for missing value (1=value is missing for observation; 0=value is observed for observation)
 - Impute missing values to a constant (such as the mean)
 - Include missing indicator in the algorithm
 - Advantage: uses all available information about missing observation
 - Disadvantage: results in biased estimates, not theoretically driven
- Regression Imputation**
 - Replaces missing values with predicted score from a regression equation.

if we don't want to lose the information that the missing data has been replaced
(so if we pass the dataset to someone else, also this someone else will know that there were missing data)

Imputation methods

35

- Extract a model from the dataset to perform the imputation
 - Suitable for MCAR and, to a lesser extent, for MAR
 - Not suitable for NMAR type of missing data
- For NMAR we need to go back to the source of the data to obtain more information
- Survey of imputation methods available at
 - <http://sci2s.ugr.es/MVDM/index.php>
 - <http://sci2s.ugr.es/MVDM/biblio.php>

inaccurate values

Inaccurate Values

37

- Data has not been collected for mining it
- Errors and omissions that don't affect original purpose of data (e.g. age of customer)
- Typographical errors in nominal attributes, thus values need to be checked for consistency
- Typographical and measurement errors in numeric attributes, thus outliers need to be identified
- Errors may be deliberate (e.g. wrong zip codes)

Prof Pierluca Lanzi

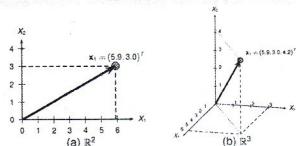
POLITECNICO DI MILANO

the geometric view

The Geometrical View of the Data

39

- When the data contains only numerical values
 - Every row can be viewed as a point in a d-dimension space
 - Every column as a point in a n-dimensional space



Prof Pierluca Lanzi

POLITECNICO DI MILANO

From categorical attributes
to numerical ones and Vice-versa

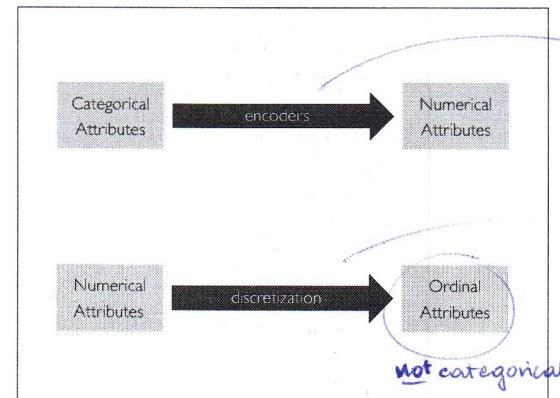
Why do we care about data types?

They influence the type of statistical analyses
and visualization we can perform

Some algorithms and functions
fit some specific data types best

Check for valid values

Deal with missing values, etc.



we have to map labels
into numbers, there are 3
ways to do it : *

It basically means to create
intervals and use them as
ordinal attributes

Scikit-Learn Encoders 43

- **LabelEncoder**
 - Encodes target labels with values between 0 and $n_labels - 1$
- **OneHotEncoder**
 - Performs a one-hot encoding of categorical features.
- **OrdinalEncoder**
 - Performs an ordinal (integer) encoding of the categorical features
- ...

Prof. Pierluigi Lanzi POLITECNICO DI MILANO

1. Label Encoder

Map a categorical variable described by n values
into a numerical variables with values from 0 to n-1

For example, attribute Outlook would be replaced by a
numerical variables with values 0, 1, and 2

→ Remember to keep
tracking of the encoding

The diagram shows a transformation of the 'Outlook' attribute. On the left, a vertical list of categories includes 'Outlook', 'Sunny', 'Sunny', 'Overcast', 'Rainy', 'Rainy', 'Overcast', 'Sunny', 'Sunny', 'Rainy', 'Sunny', 'Overcast', 'Overcast', and 'Rainy'. A large arrow points to the right, where the same data is shown as a vertical list of numerical values: 2, 2, 0, 1, 1, 1, 0, 2, 2, 1, 2, 0, 0, 1. Below this, the caption reads 'Label Encoder for the Outlook attribute.'

Label Encoder for the Outlook attribute.

Warning

By replacing a label with a number might influence the process in unexpected ways

In the example, by assigning 0 to overcast and 2 to sunny we give a higher weight to the latter

What happens if we then apply a regression model?

Would the result change with different assigned values?

If we apply label encoding, we should store the mapping used for each attribute to be able to map the encoded data into the original ones

we should always check that different labels encoding do not create too much different results (if so, the labels encoding is influencing and this is not good)

2. One Hot Encoding

Map each categorical attribute with n values into n binary 0/1 variables

Each one describing one specific attribute values

For example, attribute Outlook is replaced by three binary variables Sunny, Overcast, and Rainy

Outlook	Outlook_Outcast	Outlook_Rainy	Outlook_Sunny
Sunny	0	0	1
Overcast	1	0	0
Rainy	0	1	0
Rainy	0	1	0
Overcast	1	0	0
Sunny	0	0	1
Rainy	0	1	0
Sunny	0	0	1
Overcast	1	0	0
Overcast	1	0	0
Rainy	0	1	0

One Hot Encoding for the Outlook attribute.

Warning

One hot encoding assign the same numerical value (1) to all the labels

But it can generate a massive amount of variables when applied to categorical variables with many values

We will discuss discretization later ...

(3.)

data format

52

- Most commercial tools have their own proprietary format
- Most tools import excel files and comma-separated value files

Year;Make;Model;Length 1997;Ford;E350;2,34 2000;Mercury;Cougar;2,38	Year;Make;Model;Length 1997;Ford;E350;2,34 2000;Mercury;Cougar;2,38
---	---

Prof Perucca Lanz POLITECNICO DI MILANO

53

Attribute-Relation File Format (ARFF)

```
% ARFF file for weather data with some numeric features
%
@relation weather

@attribute outlook {sunny, overcast, rainy}
@attribute temperature numeric
@attribute humidity numeric
@attribute windy {true, false}
@attribute play? {yes, no}

@data
sunny, 85, 85, false, no
sunny, 80, 90, true, no
overcast, 83, 86, false, yes
...
```

<http://www.cs.waikato.ac.nz/~ml/weka/arff.html>

Prof Perucca Lanz POLITECNICO DI MILANO

54

Missing Values in ARFF

```
@relation labor
@attribute duration real
@attribute wage-average-first-per-real
@attribute wage-increase-second-year-real
@attribute wage-increase-third-year-real
@attribute 'cost of living-adjustment' {"none","1cf","1cf"}
@attribute working-hours real
@attribute pension {"none","ret_lake","empl_contr"}
@attribute standby-pay real
@attribute shift-differential real
@attribute education-allevance {"yes","no"}
@attribute statutory-holidays real
@attribute 'long-term disability-insurance' {"none","generous"}
@attribute "long-term disability-insurance" {yes,no}
@attribute "contributes-to-health-plan" {"none","half","full"}
@attribute "contributes-to-health-plan" {yes,no}
@attribute "contributes-to-health-plan" {"none","half","full"}
@attribute "diss" {bad,good}
@data
1.51 1.442 1.22 1.11 average:3 Eyes:2 good
2.4 2.3 2.25 ret_lake:1 1 cf:1 below-average:1 half:full good
4.4 4.2 4.25 empl_contr:3 3.5 3.11 yes,no:yes, yes,half,yes,half,good
3.37453c:1 1 1 1 yes:1 1 1 Eyes:2 good
```

Prof Perucca Lanz POLITECNICO DI MILANO

55

DSPL: Dataset Publishing Language

- Open format by Google available at <http://code.google.com/apis/publicdata/> <https://code.google.com/archive/p/dspl/downloads>
- Use existing data: add an XML metadata file existing CSV
- Read by the Google Public Data Explorer, which includes animated bar chart, motion chart, and map visualization
- Allow linking to concepts in other datasets
- Geo-enabled: allows adding latitude and longitude data to your concept definitions

Prof Perucca Lanz POLITECNICO DI MILANO

model representation

Predictive Model Markup Language

- XML-based markup language developed by the Data Mining Group (DMG) to provide a way for applications to define models related to predictive analytics and data mining
- The goal is to share models between applications
- Vendor-independent method of defining models
- Allow to exchange of models between applications.
- PMML Components: data dictionary, data transformations, model, mining schema, targets, output

57

Prof. Pierluca Lanzi POLITECNICO DI MILANO

data repositories

Publicly Available Datasets

- **UCI repository**
 - <http://archive.ics.uci.edu/ml/>
 - Probably the most famous collection of datasets
- **Kaggle**
 - <http://www.kaggle.com/>
 - It is not a static repository of datasets, but a site that manages Data Mining competitions
 - Example of the modern concept of crowdsourcing
- **KDNuggets**
 - <http://www.kdnuggets.com/datasets/>
- **Google Data Search**
 - <https://datasetsearch.research.google.com>

59

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Readings

- "Data Mining and Analysis" – Chapter 1
- "Mining of Massive Datasets" – Chapter 1
- Survey of imputation methods
<http://sci2s.ugr.es/MVDM/index.php>

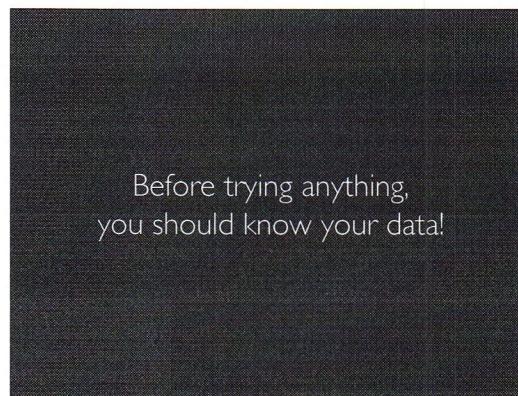
60

Prof. Pierluca Lanzi POLITECNICO DI MILANO

POLITECNICO DI MILANO

Data Exploration and Visualization
Data Mining and Text Mining

Prof. Pierluca Lanzi



What is Data Exploration?

• Preliminary exploration of the data aimed at identifying their most relevant characteristics

• What the key motivations?

- Help to select the right tool for preprocessing and data mining
- Exploit humans' abilities to recognize patterns not captured by automatic tools

• Related to Exploratory Data Analysis (EDA)

- Invented/advocated by statistician John Tukey

3

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Exploratory Data Analysis

• "An approach of analyzing data to summarize their main characteristics without using a statistical model or having formulated a prior hypothesis."

• "Exploratory data analysis was promoted by John Tukey to encourage statisticians to explore the data, and possibly formulate hypotheses that could lead to new data collection and experiments."

• Resources

- Seminal book: "Exploratory Data Analysis" by Tukey
https://books.google.it/books/about/Exploratory_Data_Analysis.html?id=U19dAAAAIAAJ
- Online introduction in Chapter 1 of NIST Engineering Statistics Handbook
<http://www.itl.nist.gov/dv898/handbook/index.htm>
- Wikipedia article
http://en.wikipedia.org/wiki/Exploratory_data_analysis

4



Prof. Pierluca Lanzi POLITECNICO DI MILANO

What Data Exploration Techniques?

• Exploratory Data Analysis (as originally defined by Tukey), was mainly focused on

- Visualization
- Clustering and anomaly detection (viewed as exploratory techniques)
- In data mining, clustering and anomaly detection are major areas of interest, and not thought of as just exploratory

• In this section, we focus on data exploration using

- Summary statistics
- Visualization

• We will return to data exploration when discussing clustering

5

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Summary Statistics

Summary Statistics

- What are they?
 - Numbers that summarize properties of the data
- Summarized properties include
 - Location, mean, spread, skewness, standard deviation, mode, percentiles, etc.
- Most summary statistics can be calculated in a single pass

7

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Frequency and Mode

- The frequency of an attribute value
 - The percentage of time the value occurs in the data set
 - For example, given the attribute 'gender' and a representative population of people, the gender 'female' occurs about 50% of the time.
- The mode of an attribute is the most frequent attribute value
- The notions of frequency and mode are typically used with categorical data

8

Prof. Pierluca Lanzi POLITECNICO DI MILANO

Measures of Location: Mean and Median

9

- The mean is the most common measure of the location of a set of points
- However, the mean is very sensitive to outliers
- Thus, the median or a trimmed mean is also commonly used

$$\text{mean}(x) = \bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

$$\text{median}(x) = \begin{cases} x_{(r+1)} & \text{if } m \text{ is odd, i.e., } m = 2r + 1 \\ \frac{1}{2}(x_{(r)} + x_{(r+1)}) & \text{if } m \text{ is even, i.e., } m = 2r \end{cases}$$

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Percentiles

10

- For continuous data, the notion of a percentile is very useful
- p-th percentile
 - Given an ordinal or continuous attribute x and a number p
 - p-th percentile is a value x_p of x such that p% of the observed values of x are less than x_p
- For instance, the 50th percentile is the value $x_{50\%}$ such that 50% of all values of x are less than $x_{50\%}$

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Trimean & Truncated Mean

- **Trimean**

- It is the weighted mean of the first, second and third quartile

$$TM = \frac{x_{25} + 2x_{50} + x_{75}}{4}$$

- **Truncated Mean**

- Discards data above and below a certain percentile
- For example, below the 5th percentile and above the 95th percentile

- **Interquartile Mean**

- Truncate data at 25th and 75th percentile
- If the data (x_1, \dots, x_n) is sorted by value we have:

$$X_{IQM} = \frac{2}{n} \sum_{i=0.25n+1}^{0.75n} x_i$$

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Measures of Spread: Range and Variance

12

- Range is the difference between the max and min

- The variance or standard deviation is the most common measure of the spread of a set of points

$$\text{variance}(x) = s_x^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$$

- However, this is also sensitive to outliers, so that other measures are often used

$$AAD(x) = \frac{1}{m} \sum_{i=1}^m |x_i - \bar{x}|$$

$$MAD(x) = \text{median}(|x_1 - \bar{x}|, \dots, |x_m - \bar{x}|)$$

$$\text{interquartile range}(x) = x_{75\%} - x_{25\%}$$

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

useful for variable selection: if a variable has low variance we can consider to get rid of it since it's not informative! Extreme case: constant values, null variance.

Correlation Analysis

13



Karl Pearson

- Given two attributes, measure how strongly one attribute implies the other, based on the available data
- Use correlation measures to estimate how predictive one attribute is of another
- Linear correlation
 - We often look for linear relationship between variables not all variables are linearly related, of course!
 - Linear correlation measures are symmetric
 - They can be positive (high values of one attribute are likely given high values of the other)
 - Or negative (high values are predictive of low values of the other variable)
 - Latter usually referred to as anti-correlation

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Correlation Analysis

14

- Given two attributes it measure how strongly one attribute implies the other, based on the available data
- **Numerical Variables**
 - For two numerical variables, we can compute the correlation coefficient, Pearson's product moment coefficient
- **Ordinal Variables**
 - We can compute Spearman rank correlation coefficient
- **Categorical Variables**
 - For two categorical variables, A and B, we can compute χ^2 (chi-square) statistic test which tests the hypothesis that A and B are independent
- **Binary Variables**
 - Compute Point-biserial correlation

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Correlation vs Causation

15

- **Correlation does not imply causation**

- Just because value of one attribute is highly predictive of value of other doesn't mean that forcing the first variable to take on a particular value will cause the second to change

- **Causality has a direction, while correlation typically doesn't**

- Correlation between high income and owning a Ferrari
- Giving a person a Ferrari doesn't affect their income
- But increasing their income may make them more likely to buy a Ferrari

- **Confounding variables can cause attributes to be correlated**

- High heart rate and sweating are correlated with each other since they tend to both happen during exercise (confounder)
- Causing somebody to sweat by putting them in sauna won't necessarily raise their heart rate (it does a little, but not as much as exercise)
- And giving them beta-blockers to lower their heart rate might not prevent sweating (it might a little, but again not like stopping exercising)

correlations might be misleading, we shouldn't assign meanings

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Outliers 16

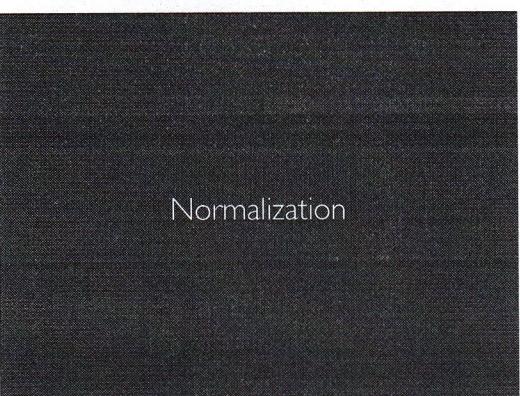
- **What are outliers?**
 - Data objects that do not comply with the general behavior or model of the data, that is, values that appear as anomalous
 - Most data mining methods consider outliers noise or exceptions.
- **Outliers may be detected using**
 - Manual inspection and knowledge of reasonable values.
 - Statistical tests that assume a distribution or probability model for the data
 - Distance measures where objects that are a substantial distance from any other cluster are considered outliers
 - Deviation-based methods identify outliers by examining differences in the main characteristics of objects in a group

Prof. Pierluca Lanzi POLITECNICO DI MILANO

How Do We Manage Outliers? 17

- Outliers are typically filtered out by eliminating the data points containing them
- **Trimming**
 - Eliminate the outlier data values
- **Winsorizing**
 - A 10% Winsorizing, consider the 5th and 95th percentiles
 - Set the values below the 5th percentile to the 5th percentile itself
 - Set the values above the 95th percentile to the 95th percentile itself
- Note that, in some applications, outliers are the focus on the analysis. Fraud detection is a typical example of this.

Prof. Pierluca Lanzi POLITECNICO DI MILANO

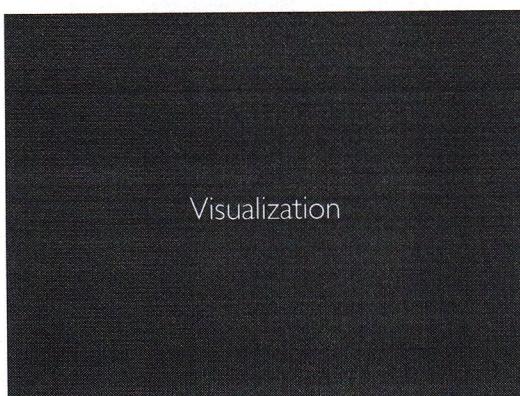


Normalization 19

- When attributes have vastly different scales (e.g., age vs income), it is necessary to normalize them
- Range normalization
 - Converts all values to the range [0,1]

$$x'_i = \frac{x_i - \min_i x_i}{\max_i x_i - \min_i x_i}$$
- Standard Score Normalization
 - Forces variables to have mean of 0 and standard deviation of 1
 - So if data was normally distributed, most of it (68%) will lie in range [-1,+1]
$$x'_i = \frac{x_i - \mu}{\sigma}$$

Prof. Pierluca Lanzi POLITECNICO DI MILANO



Why visualize data? Anscombe's Quartet

21

	x^1	y	x^2	y	x^3	y	x^4	y
1.0	8.04	10.0	9.14	10.0	2.46	9.0	6.58	
1.0	9.08	8.0	8.14	6.0	2.46	8.0	6.98	
1.0	7.83	10.0	8.74	10.0	1.97	8.0	7.71	
1.0	8.81	8.0	8.84	6.0	1.97	8.0	8.77	
1.0	8.33	11.0	9.26	11.0	7.81	8.0	8.42	
1.0	9.96	14.0	9.16	14.0	8.84	8.0	7.04	
1.0	7.27	10.0	7.58	10.0	1.87	8.0	7.53	
1.0	4.25	4.0	3.10	4.0	5.38	10.0	12.50	
1.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	
1.0	7.0	4.0	7.0	4.0	7.0	4.0	5.54	
1.0	5.0	3.0	5.0	3.0	5.0	3.0	4.0	
1.0	5.0	5.0	4.74	5.0	5.73	4.0	6.89	
Mean	9.0	7.3	9.0	7.5	9.0	7.3	9.0	7.5
Variance	0.75	0.75	0.75	0.75	0.75	0.75	0.75	0.75
Correlation	0.818	0.818	0.818	0.818	0.818	0.818	0.818	0.818

https://en.wikipedia.org/wiki/Anscombe%27s_quartet

Prof. Pierluca Lanzo POLITECNICO DI MILANO

Why visualize data? Anscombe's Quartet

22

	x^1	y	x^2	y	x^3	y	x^4	y
1.0	8.04	10.0	9.14	10.0	2.46	9.0	6.58	
1.0	9.08	8.0	8.14	6.0	2.46	8.0	6.98	
1.0	7.83	10.0	8.74	10.0	1.97	8.0	7.71	
1.0	8.81	8.0	8.84	6.0	1.97	8.0	8.77	
1.0	8.33	11.0	9.26	11.0	7.81	8.0	8.42	
1.0	9.96	14.0	9.16	14.0	8.84	8.0	7.04	
1.0	7.27	10.0	7.58	10.0	1.87	8.0	7.53	
1.0	4.25	4.0	3.10	4.0	5.38	10.0	12.50	
1.0	10.0	10.0	10.0	10.0	10.0	10.0	10.0	
1.0	7.0	4.0	7.0	4.0	7.0	4.0	5.54	
1.0	5.0	3.0	5.0	3.0	5.0	3.0	4.0	
1.0	5.0	5.0	4.74	5.0	5.73	4.0	6.89	
Mean	9.0	7.3	9.0	7.5	9.0	7.3	9.0	7.5
Variance	0.75	0.75	0.75	0.75	0.75	0.75	0.75	0.75
Correlation	0.818	0.818	0.818	0.818	0.818	0.818	0.818	0.818

https://en.wikipedia.org/wiki/Anscombe%27s_quartet

mean, variance and correlation are the same in all the 4 cases

All the summary statistics are the same but the data points are very different

Prof. Pierluca Lanzo POLITECNICO DI MILANO

Visualization

23

- Visualization is the conversion of data into a visual or tabular format so that the characteristics of the data and the relationships among data items or attributes can be analyzed or reported
- Data visualization is one of the most powerful and appealing techniques for data exploration
 - Humans have a well-developed ability to analyze large amounts of information that is presented visually
 - Can detect general patterns and trends
 - Can detect outliers and unusual patterns

Prof. Pierluca Lanzo

POLITECNICO DI MILANO

Bar Plots

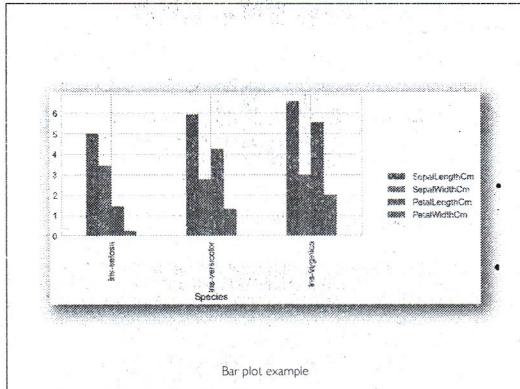
Bar Plots

25

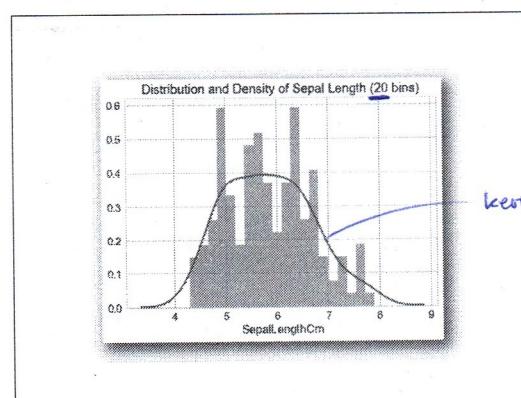
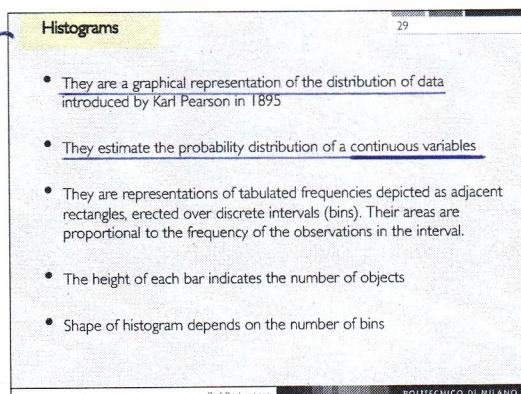
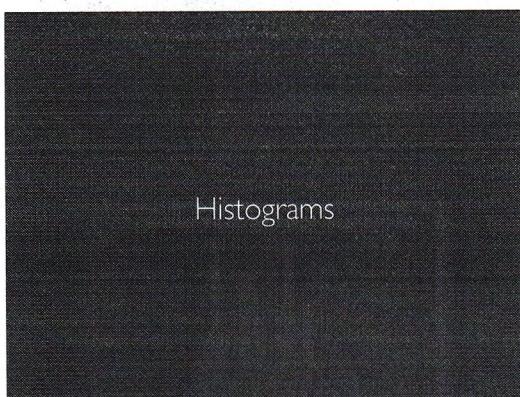
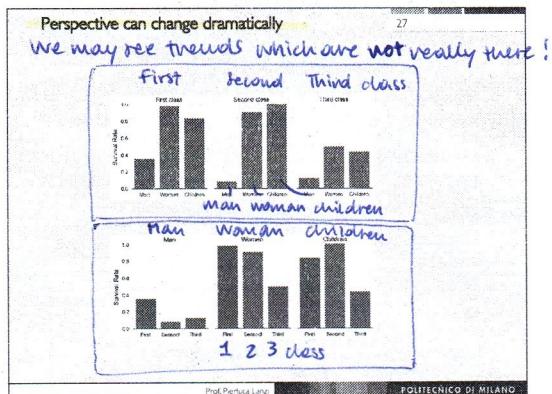
- They use horizontal or vertical bars to compare categories.
- One axis shows the compared categories, the other axis represents a discrete value
- Some bar graphs present bars clustered in groups of more than one (grouped bar graphs), and others show the bars divided into subparts to show cumulative effect (stacked bar graphs)

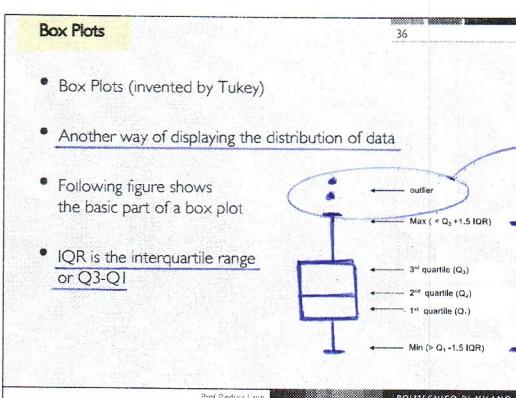
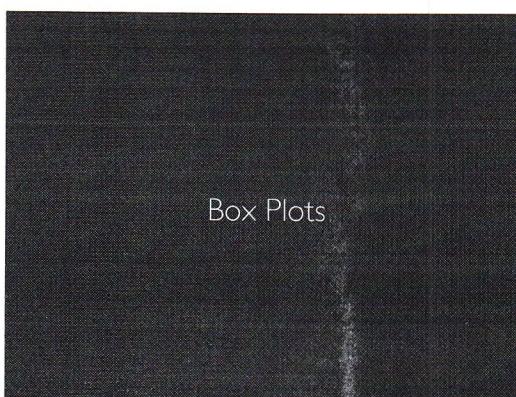
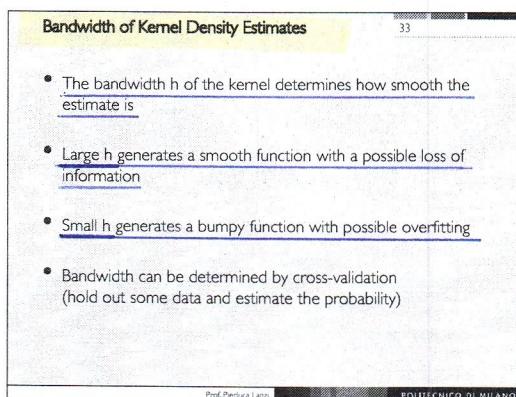
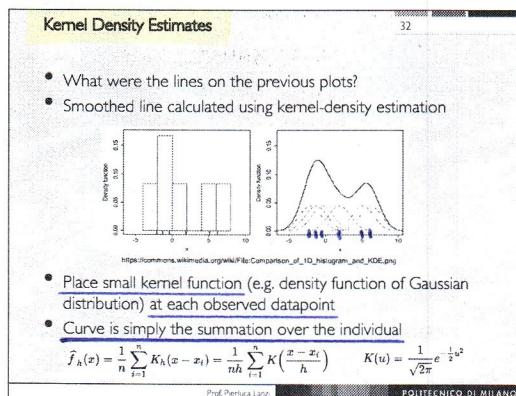
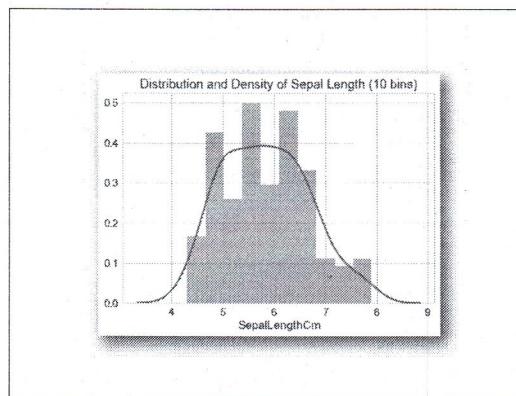
Prof. Pierluca Lanzo

POLITECNICO DI MILANO



Bar plot example



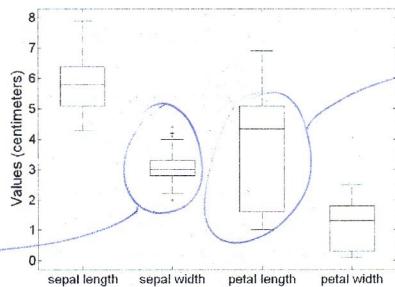


The fact that boxplots name them outliers doesn't mean that they really are. This has nothing to do with outliers detection.

$$Q_1 - \frac{3}{2} \text{ IQR}$$

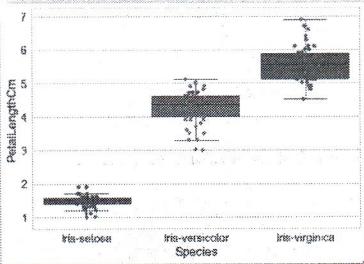
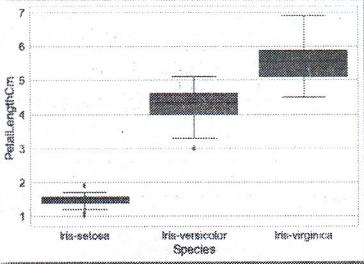
Box Plot Example

37



Here we can see
that the median is
almost at the center
→ symmetric data

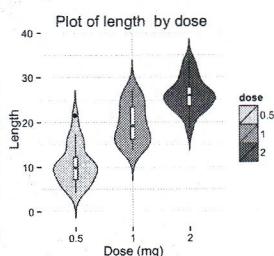
Here instead the median is far
from being in the center, so data
are not symmetric



Violin Plots

40

- Uses the kernel density estimate (smoothed histogram) instead of just a box at quartiles



Source: <http://www.sibla.unimib.it/~vito/pd2/violin.pdf> (quick start guide to R, violin plots)

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

Scatter Plots

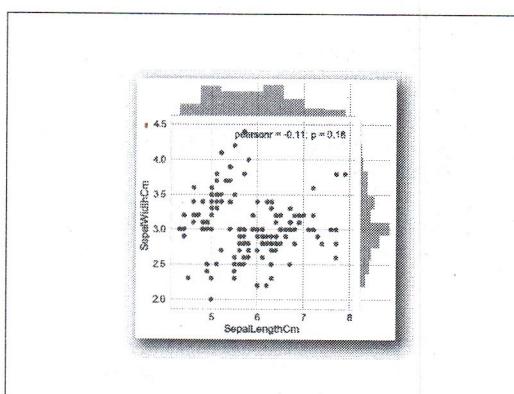
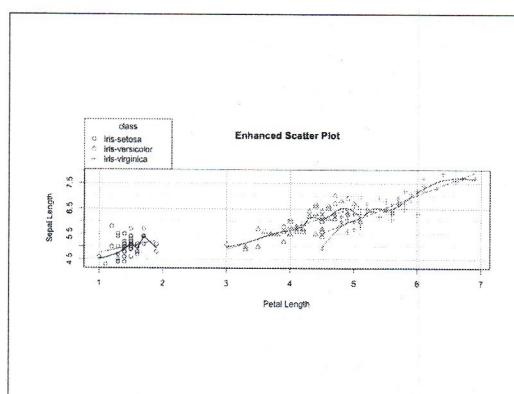
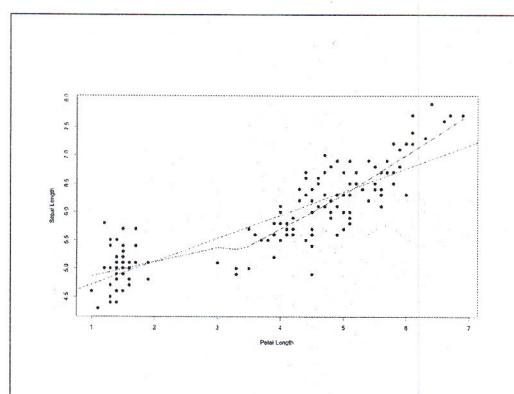
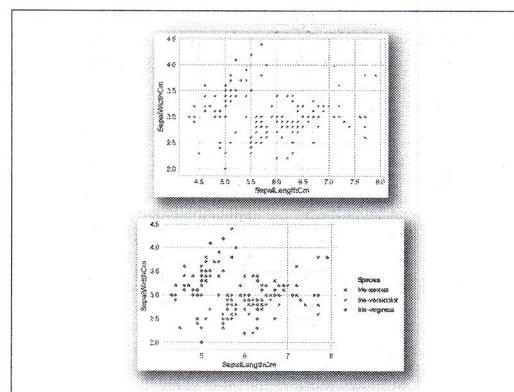
Scatter Plots

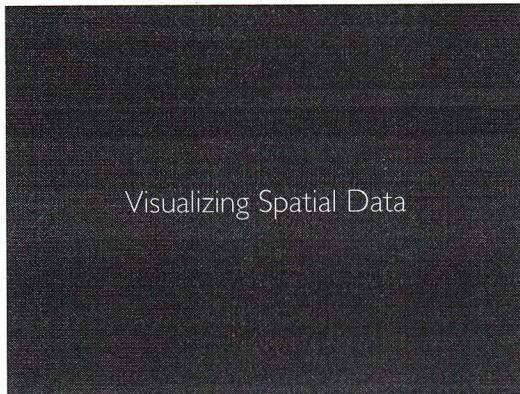
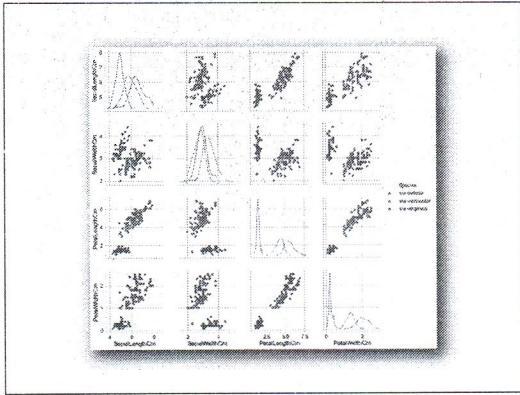
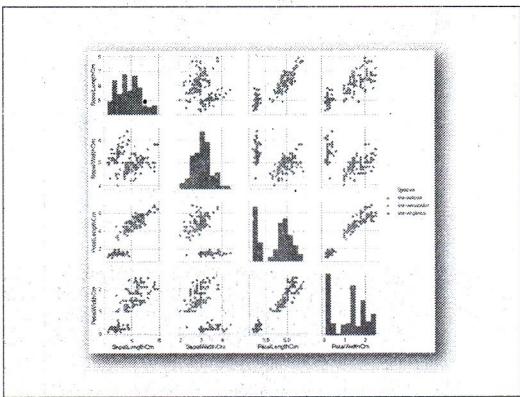
42

- Used to compare two (or more) attributes
 - Attributes values used to determine the position of the point
 - Two-dimensional scatter plots most common, but three-dimensional plots also used
- Often additional attributes can be displayed by using the size, shape, and color of the markers that represent the objects
- It is useful to have arrays of scatter plots can compactly summarize the relationships of several pairs of attributes
- Examples: <http://www.statmethods.net/graphs/scatterplot.html>

Prof. Pierluca Lanzi

POLITECNICO DI MILANO





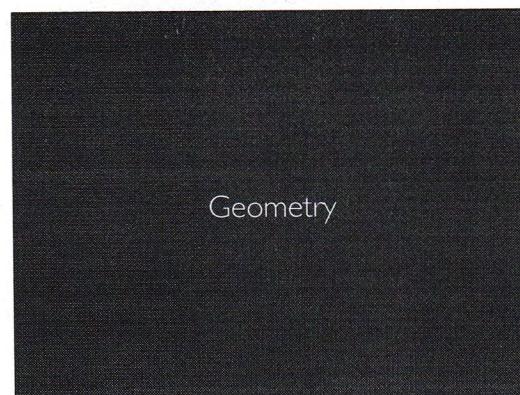
Visualizing Spatial Data

- **Geometry**
 - Geographic
 - Others
- **Scalar fields (one value per cell)**
 - Isocontours
 - Direct volume rendering
- **Vector and tensor fields (many values per cell)**
 - Flow glyphs
 - Geometric (sparse seeds)
 - Textures (dense seeds)
 - Features (globally derived)

50

Prof. Pierluigi Lanzi

POLITECNICO DI MILANO



Dot Maps 52

- **What?**
 - Geometry (position)
- **Why?**
 - Locate data in space
- **Remarks**
 - Scale up to hundreds of items
 - Color/shape can encode an additional categorical attribute (reduce scalability)

Prof. Pierluca Lanz | POLITECNICO DI MILANO

Bubble Maps 53

- **What?**
 - Geometry (position)
 - One quantitative attribute
- **Why?**
 - Locate data in space
 - Lookup and compare
- **Remarks**
 - Scale up to hundreds of items
 - Color can encode an additional cat (interaction with size).

Prof. Pierluca Lanz | POLITECNICO DI MILANO

Choropleth Map 54

- **What?**
 - Geometry (position)
 - One quantitative attribute
- **Why?**
 - Locate data in space
 - Lookup and compare
- **Remarks**
 - Scale up to ~1000 items
 - Hue can encode an additional categorical attribute (better if binary)

Prof. Pierluca Lanz | POLITECNICO DI MILANO

Bias of Population-related Maps 55

- Using absolute values is dangerous!
- Any map would show population distributionS
- How to deal with this?
 - Visualize per capita (relative)
 - Use statistical models
- See the example at <https://xkcd.com/1138>

REFRESH #208
GEOGRAPHIC PROFILE MAPS WHICH ARE
BASICALLY JUST POPULATION MAPS

Prof. Pierluca Lanz | POLITECNICO DI MILANO

Surprise Maps 56

- **What?**
 - Geometry (position)
 - One quantitative attribute
 - One derivative attribute
- **Why?**
 - Locate data in space
 - Lookup and compare
- **Remarks**
 - The surprise is computed as a function prior and posterior probability of data distribution.
 - Prior probability is generated with a family of standard models.
- Illustration from <https://medium.com/@uwdata/surprise-maps-showing-the-unexpected-e92b67398865>

Signed Surprise

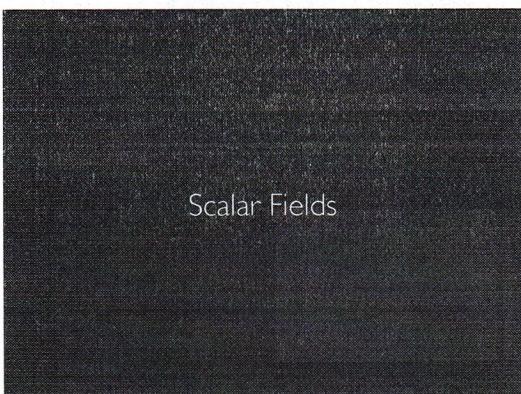
Prof. Pierluca Lanz | POLITECNICO DI MILANO

Connection Maps

57

- **What?**
 - Network and positions
- **Why?**
 - Lookup path
 - Identify patterns
- **Remarks**
 - Size of links can encode an additional ordered attribute (3-4 bins at max)

Prof. Pierluca Lenzi POLITECNICO DI MILANO



Isocontour Maps

59

- **What?**
 - Geographic data
 - One quantitative attribute
 - Derived positions
- **Why?**
 - Shape
- **Remarks**
 - The lines are computed from the values of scalar field
 - Area can be filled and color encoded

Prof. Pierluca Lenzi POLITECNICO DI MILANO

Isocontour Plots

60

- **What?**
 - 2D spatial field
 - One quantitative attribute
 - Derived geometry
- **Why?**
 - Shape and patterns
- **Remarks**
 - The lines are computed from the values of scalar field
 - Area can be empty or filled and color encoded

Prof. Pierluca Lenzi POLITECNICO DI MILANO

Isosurface Plots

61

- **What?**
 - 3D spatial scalar field
 - 1 quantitative attribute
 - derived geometry
- **Why?**
 - Shape
- **Remarks**
 - Tree of isosurfaces: positions computed for specific values of the scalar field

Prof. Pierluca Lenzi POLITECNICO DI MILANO



Glyph Flow 63

- **What?**
 - 2D vector fields
- **Why?**
 - Shape and patterns
 - Identify critical points
- **Remarks**
 - Different glyphs can be used to represent vectors
 - Density of grid and jittering

Prof. Pierluca Lanza POLITECNICO DI MILANO

Geometric Flow 64

- **What?**
 - 2D/3D vector field
 - Derived geometry
- **Why?**
 - Shape and patterns
 - Identify critical points
- **Remarks**
 - Seeding strategy affects the outcome
 - Usage of clustering and color coding improves readability

Prof. Pierluca Lanza POLITECNICO DI MILANO

Texture Flow 65

- **What?**
 - 2D vector field
- **Why?**
 - Shape and patterns
 - Identify critical points
- **Remarks**
 - Similar to glyphs flow, but computes the flow of a continuous distribution of particles

Prof. Pierluca Lanza POLITECNICO DI MILANO

Feature flow: anatomy 66

- **What?**
 - 2D vector fields
- **Why?**
 - Shape and patterns
 - Identify critical points
- **Remarks**
 - Similar to glyphs flow
 - But seeding is based on global computing strategy to identify areas with similar behaviors

Prof. Pierluca Lanza POLITECNICO DI MILANO

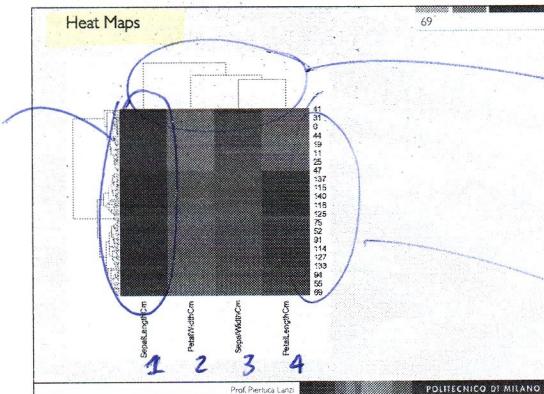
More than Two Dimensions at Once

Two main approaches

Visualize all the dimensions at once
(e.g., heatmaps, spider plots, and Chernoff)

Project the data into a smaller space
and visualize the projected data

We can see that these values are very different from the others



From here we can obtain a clustering. From this clustering we can see that 3 and 4 are the most similar. Then 2 is the most similar to them. And the least similar is 1.

we see that there is a sorting

Spider plots, Radar Plots, and Star Plots

• Information radiates outward from central point

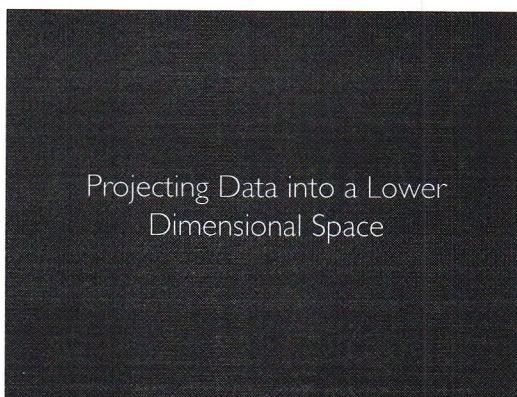
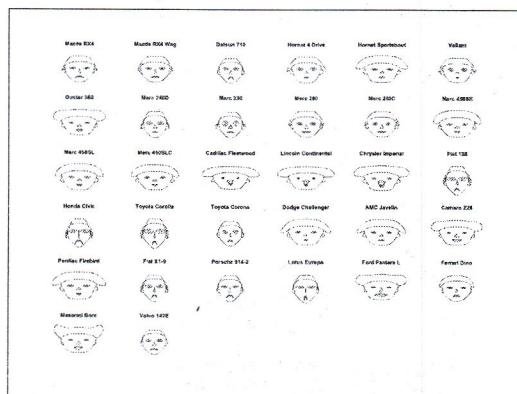
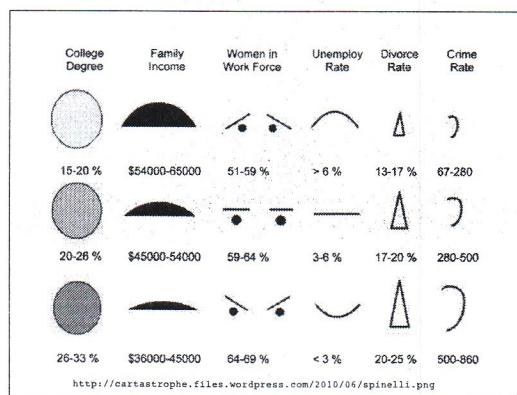
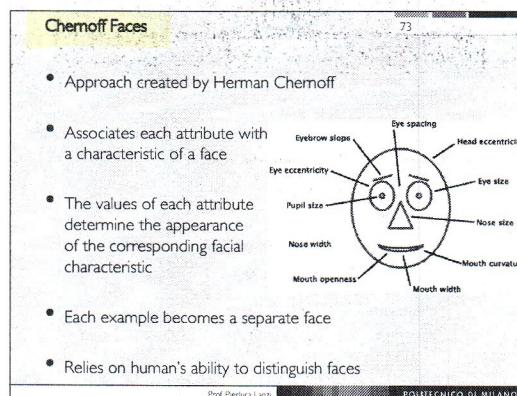
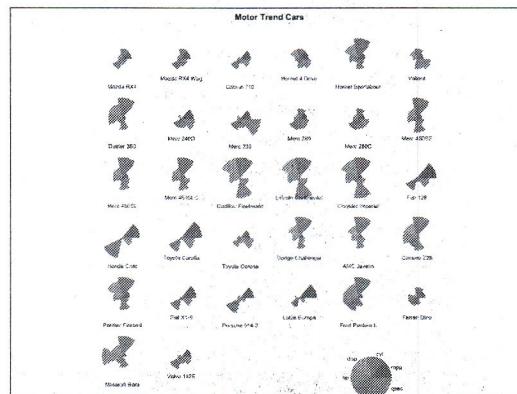
• The line connecting the values of an object is a polygon



Example Sector Chart by David Clement

We have a vertex for every variable (feature). If we plot all the datapoints we should be able to detect some patterns.





When projecting high-dimensional data into fewer dimensions we can either

1. Find a linear projection

e.g. use Principle Component Analysis

2. Find a non-linear projection

e.g. use t-distributed Stochastic Neighbor Embeddings (t-SNE)

1.

Principal Component Analysis

Principal Component Analysis (PCA)

80

- Typically applied to reduce the number of dimensions of data (feature selection)
- The goal of PCA is to find a projection that captures the largest amount of variation in data
- Given N data vectors from n-dimensions, find $k < n$ orthogonal vectors (the principal components) that can be used to represent data
- Works for numeric data only and it is affected by scale so data usually need to be rescaled before applying PCA

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

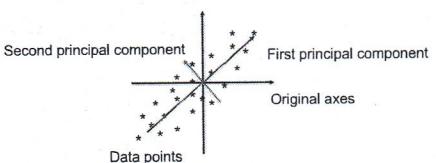
Principal Component Analysis (PCA)

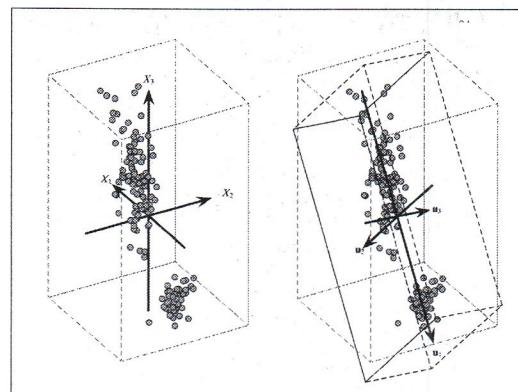
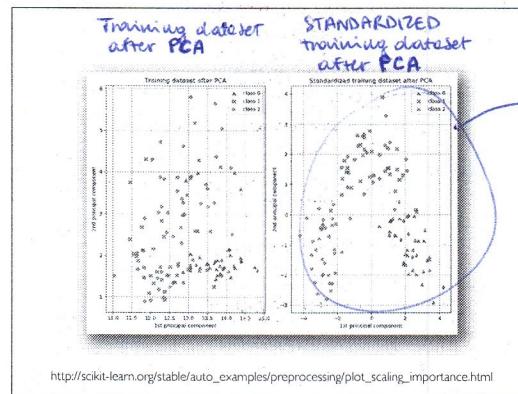
81

- Steps to apply PCA
 - Normalize input data
 - Compute k orthonormal (unit) vectors, i.e., principal components
 - Each input data point can be written as a linear combination of the k principal component vectors
- The principal components are sorted in order of decreasing "significance" or strength
- Data size can be reduced by eliminating the weak components, i.e., those with low variance.
- Using the strongest principal components, it is possible to reconstruct a good approximation of the original data)

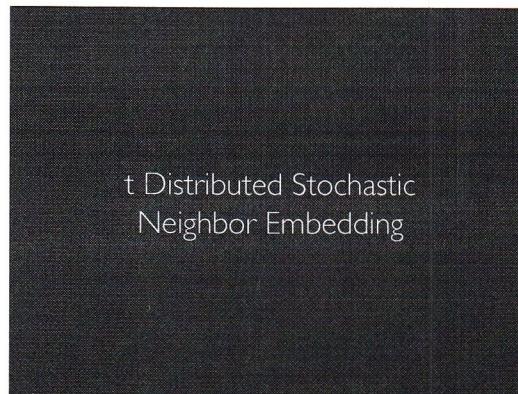
Prof. Pierluca Lanzi

POLITECNICO DI MILANO

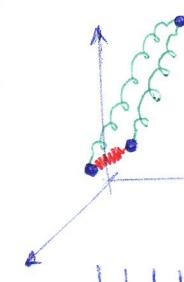




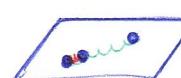
2.



Consider a set of points in a multiple-dim space. We put a strong spring ("molla") between two close points and a weak spring between two far points. We do this for all the points combinations.



Once we're done with the springs we let the points fall down on a 2 dimensional plane.



Once the spring has calmed down we see that the closest points in higher dimension are close also in this projection (because of the strong spring)

Note: there is a stochastic component, if we repeat the procedure K times we obtain K different results

Non-linear Dimensionality Reduction

- Data in high dimensions never fills the entire space and always lives within some lower-dimensional manifold
- t-SNE is a non-linear dimensionality reduction technique used to map high-dimensional data into 2 or 3 dimensions
- Points from original space mapped onto "map points" in 2D/3D
- Unlike PCA, the mapped points are not linear combination of original attribute values and the axes of mapped space are not linear combination (rotation) of original axes

86

original axes

PCA's x-axis (first principal component)

t-SNE's x-axis (non-linear mapping)

original axes

86

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

t Distributed Stochastic Neighbor Embedding

- t-SNE tries hard to preserve local distances to nearby points
- Unlike PCA which tries to preserve global (long range) distances between points as much as possible
- t-SNE converts distances between data points to joint probabilities then models original points by mapping them to low dimensional map points such that position of map points conserves the structure of the data
- i.e. similar data points are modeled by nearby map points while dissimilar data points are modeled by distant map points

87

Prof. Pierluca Lanzi

POLITECNICO DI MILANO

- The t-SNE algorithm has two main steps

- Define a probability distribution over pairs of high-dimensional data points so that:
 - Similar data points have a high probability of being picked
 - Dissimilar points have an extremely small probability of being picked
- Define a similar distribution over the points in the map space
 - Minimize the Kullback–Leibler divergence between the two distributions with respect to the locations of the map points
 - To minimize the score, it applies gradient descent

- Assume that map points are all connected with springs.
- The stiffness of a spring connecting two points depends on the mismatch between the similarity of the two data points and the similarity of the two map points
- Let the system evolve according to the laws of physics
 - If two map points are far apart while the data points are close, they are attracted together
 - If they are nearby while the data points are dissimilar, they are repelled
- The final mapping is obtained when the equilibrium is reached.

- Optical Recognition of Handwritten Digits Data Set from the UCI machine learning repository
<http://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits>
- Contains 1797 images with 8x8 pixels each



Iterations of the t-SNE algorithm over the Optical Recognition of Handwritten Digits Data Set
<https://github.com/oreillymedia/t-SNE-tutorial>

- The parameter perplexity says (loosely) how to balance attention between local and global aspects of the data
- Different initializations will lead to different results
- Should be applied to data with a "reasonable" number of dimensions (e.g. 30-50)
- If the data have more dimensions, another dimensionality reduction algorithm should be applied

Force-directed Layout

93

- Idea of mapping complicated data into 2D is not limited to high dimensional data
- We can map any graph of data points into 2D provided we have some (dis)similarity value between pairs of nodes
 - Such as the Euclidean distance between them in higher dimensional space
 - Or their joint probability under a Gaussian kernel (in case of t-SNE)
 - Or Pearson's correlation, Spearman's Rank correlation, chi-squared, etc.
- It works by moving points around in the mapped 2D space until convergence
- Technique is called force-directed layout. There are many algorithms implemented in Gephi (<https://gephi.org>). On next slide we have an example of a graph generated using it.

Prof Pierluca Lanzi

POLITECNICO DI MILANO

Other Visualizations

Visualization of retweets during the
2013 Boston Marathon on April 15, 2013

obtained using Gephi and the Twitter API

Some Interesting Pages

96

- <https://python-graph-gallery.com>
- <https://python-graph-gallery.com/category/seaborn/>
- <https://github.com/matplotlib/AnatomyOfMatplotlib>
- <https://github.com/rasbt/matplotlib-gallery>
- <https://seaborn.pydata.org/examples/index.html>
- <http://bokeh.pydata.org/en/latest/>

Prof Pierluca Lanzi

POLITECNICO DI MILANO

Study Material

97

- "Data Mining and Analysis" – Chapter 2 & 3
- **t-Distributed Stochastic Neighbor Embedding (t-SNE)**
 - <http://vdmaaten.github.io/tsne/>
 - <https://www.youtube.com/watch?v=RjVl80Gg3IA>
 - <http://alexanderfabisch.github.io/t-sne-n-scikit-learn.html>
 - <http://mrt.csail.mit.edu/papers/volume9/vandermaaten08a/vandermaaten08a.pdf>
 - <https://distill.pub/2016/misread-tsne/>
- **Principal Component Analysis**
 - "Data Mining and Analysis" – Chapter 7
 - http://sebastianraschka.com/Articles/2015_pca_in_3_steps.html
- **Data Visualization**
 - PhD Course by Daniele Loiacono

Prof Pierluca Lanzi

POLITECNICO DI MILANO