# Introduction to Data Mining

Data Mining and Data integration in Biomedicine
Master in Bioinformatics

Janet Piñero
Medbioinformatics Solutions SL
2025-2026

# Outline

Brief description of the course

What is Data Mining?

Data

Common Data Mining Tasks

Best Practices in Data science projects

Hands-on sessions

Bibliography

# Lectures / Hands-on sessions

**DMI**

1. Introduction to data mining and logistics

2. Exploratory Data Analysis (EDA)

3. Data Visualization

4. Unsupervised Learning (I)

5. Unsupervised Learning (II)

6. Supervised Learning (I)

7. Supervised Learning (II)

8. **Introduction to Agentic AI, Francesco Ronzano**

9. **Core capabilities of AI agents, Francesco Ronzano**

10. **Python-lab: building an AI agent, Francesco Ronzano**

11. Students Presentations (1) (16/02)

12. Students Presentations (2)  (18/02)

13. **General Discussion**

14. **Invited speaker: Ferran Muiños, Nonparametric regression models and linear & mixed-effects models**

3

# Evaluation

**Work in pairs!**

Exercises performed during the course -> 80 % (GitHub Classroom)

- Hands-on session 1: Lectures 1, 2 and 3 (Intro, EDA, Visualization) -> 25%

- Hands-on session 2: Lectures 4 and 5 (Unsupervised Learning) -> 25%

- Hands-on session 3: Lectures 6 and 7 (Supervised Learning) -> 25%

10 minute presentation of a scientific article -> 25 %
   +5 minutes for questions

# About the Hands-on sessions

# Evaluation rubric for the hands-on sessions

**DMI**

✓ Does the GitHub repository contain all the necessary files (data and code)?

✓ Is the GitHub repository well-organized?

✓ **Does the R Markdown file (example, index.Rmd) and any additional code run the analysis without errors and generate the expected output (e.g., index.html)?**

✓ Does the analysis described in the resulting index.html file conform to the requested sectioning?

✓ Does the introduction clearly explain the question being addressed, the data used, and the number of observations and variables involved?

✓ Do the plots provide meaningful summaries of the data? Are the axes in the plots labeled in plain language and sized appropriately for readability?

✓ Does the R Markdown file (index.Rmd) and any other additional R Markdown files include R code interspersed with explanatory Markdown text? Is the code easy to read and understand?

✓ **Is the report clean and organized?**

# Evaluation rubric for the hands-on sessions

# Presentation of a scientific article

**Select wisely!**

**DMI**

10 minute presentation of a scientific article -> 20 %

    +5 minutes for questions

/ DMI: Data Mining and Data Integration in Biomedicine 2024-32548-T1
/ Publications for presentations

📁 **Publications for presentations**

Carpeta    Paràmetres    Més ⌄

Edita

- 2000 Machine learning for survival analysis a case study on recurrence of prostate cancer.pdf
- 2000 Prospective evaluation of logistic regression models for the diagnosis of ovarian cancer.pdf
- 2001 Validating clustering for gene expression data.pdf
- 2002 Robust cluster analysis of microarray gene expression data with the number of clusters determined biologically.pdf
- 2003 Ensemble machine learning on gene expression data for cancer classification.pdf
- 2005 Novel Hybrid Hierarchical-K-means Clustering Method H-K-means for microarray data.pdf
- 2007 Blood gene expression signatures predict exposure levels.pdf
- 2008 A comprehensive comparison of random forests and support vector machines for microarray-based cancer classification.pdf
- 2009 Cancer-Specific High-Throughput Annotation of Somatic Mutations_ Computational Prediction of Driver Missense Mutations.pdf
- 2010 Enhanced_Prediction_of_Heart_Disease_with_Feature_Subset_Selection_using_Genetic_Algorithm.pdf
- 2010 Missing data imputation using statistical and machine learning methods in a real breast cancer problem.pdf
- 2011 Microarray-based cancer prediction using single genes.pdf
- 2011 Predicting disease risks from highly imbalanced data using random forest.pdf

# Presentation of a scientific article

**Select wisely!**

**DMI**

## Guide for the presentations

- What is the objective of the article?

- What programming languages are used?

- What is the quality and size of the data? Is it publicly available?

- What preprocessing steps are applied to the data? How are they reported in the article?

- Regarding the algorithms employed: How are they described? Are the parameters clearly reported?

- For machine learning models, is the training process clearly described? How is the data split?

- How is the performance of the algorithm evaluated?

- Provide a critical assessment of the figures and tables. Are they clear, accurate, and well-labeled?

- How is the supplementary material used? What information is included in the supplementary material?

- Is the code available? Is the analysis reproducible based on the provided information?

- Are there any limitations or biases in the study? How are they addressed?

# Presentation of a scientific article

**DMI**

The chosen scientific article should be sent via email (janet.pinero@upf.edu) no later than **February 01**, along with the names of the students. The list of the publications that have been already selected will be updated here <u>here</u>.

**The presentation should be sent by email no later than**
**February 15** (janet.pinero@upf.edu)

# What is Data Mining?

# Why do we need to analyze data?

- To understand what has happened or what is happening;

- To predict what is likely to happen, either in the future or in other circumstances we haven't seen yet;

- To guide us in making decisions.

# Drivers of Data Mining

- The Explosive Growth of Data: from terabytes to petabytes

- Easy data collection and increase of data availability

- Automated data collection tools, database systems, Web, computerized society

- Advances in computing power

Abundant data, limited knowledge.

# Data, information, and knowledge

**Data** – observations, signals, measurements:  *Body temperature 103*
**Information** –  a set of data in context, data with meaning: *The patient is having a fever*
**Knowledge** – justifiable beliefs based on data and information: i*f a patient temperature is > 100 F, he might a fever (or hyperthermia).*

**Data**  –   a patient's blood pressure is 100/50.
**Information**  –   that patient has a ten year history of blood pressures of 150/100.
**Knowledge**  –   the patient has a known history of coronary artery disease and is now experiencing chest pain. The sudden drop in blood pressure could indicate a serious myocardial infarction in progress.

 **Wisdom** is acting on knowledge in an appropriate way



The DIKW pyramid

# What is Data Mining?

**DMI**

❖ "**Data** mining is the science of extracting **useful knowledge** from huge **data** repositories"

❖ "**Data** Mining refers to a **set of methods** applicable to large and complex **databases** to eliminate the randomness and discover the hidden pattern. Data mining is about tools, methodologies, and theories for r**evealing patterns in data** — which is a critical step in **knowledge** discovery."

❖ "Data mining (aka knowledge discovery from data, or **data science**) consists on the extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) **patterns or knowledge** from huge amount of **data**"

❖ Data mining is the **process** of deciphering **meaningful insights** from existing **datasets**

❖ "Finding hidden information in a **database**"

# What is Data Mining?

- Many Definitions

- Non-trivial extraction of implicit, previously unknown and potentially useful information from data

- Exploration & analysis, by automatic or **semi**-automatic means, of large quantities of data in order to discover meaningful patterns

- It's a process, which means that it *can* be iterative

# Data Mining: Confluence of Multiple Disciplines

**DMI**

# Why Confluence of Multiple Disciplines?

**DMI**

**Tremendous amount of data**

- Algorithms must be *scalable* to handle big data

**High-dimensionality of data**

- Microarray data may have tens of thousands of *dimensions*

**High complexity and heterogeneity of data**

- Data streams and sensor data
- Time-series data, temporal data, sequence data
- Structure data, graphs, social and information networks
- Spatial, spatiotemporal, multimedia, text and Web data
- Software programs, scientific simulations

**Need for new and sophisticated applications**

# Myths about Data Science

**DMI**

- Data mining is an automatic process

- Every data mining process needs big data and deep learning

- Modern data science tools are easy to use

- There is always a hidden gem in the data

# General steps of the data mining process  DMI

Most data science projects follow a common set of steps:

1. Define the question/goal 🤔
2. Data: Obtain the data 📊
3. **Exploratory data analysis** 😅
4. **Model**: Statistical prediction/modeling 😁
5. Interpret results 💡
6. Synthesize/write up results ✍️
7. Create reproducible code 😒

These steps might vary somewhat depending on the context of the project, but are remarkably consistent across projects in both industry and academia.

# The Data Mining process

# Data

"A map is not the territory it represents, but, if correct, it has a similar structure to the territory which accounts for its usefulness" Alfred Korzybski

# Raw data

**DMI**

**It is critical that you include the rawest form of the data that you have access to.**

Some examples of the raw form of data are as follows.

1. The strange binary file your measurement machine spits out

2. The unformatted Excel file with 10 worksheets your colaborator sent you

3. The complicated JSON data you got from scraping the Twitter API

4. The hand-entered numbers someone collected looking through a microscope

5. The file that you downloaded from a particular resource

6. The supplementary material of a publication (if you are lucky!)

# Data: Structured vs unstructured

**DMI**

**Unstructured data sources**

**Structured data sources**

# Data: Structured vs unstructured

**DMI**

**Unstructured data sources**
○ No logical data model, no particular logic and structure is defined to collect, store, and expose data

**Structured data sources**
○ High degree of organization. There is a specific data model

## Data: Structured vs unstructured

**Unstructured data sources**

○ No logical data model, no particular logic and structure is defined to collect, store, and present data:

    ○ Medical images: X-rays, MRIs, CT scans, pathology slides

    ○ Clinical notes: doctors' narratives, discharge summaries

    ○ Biomedical signals: audio recordings of heart or lung sounds

    ○ Research articles: full-text biomedical papers

# Data: Structured vs unstructured

**DMI**

**Structured data sources**

○ High degree of organization. There is a specific data model: **an R data frame is a typical example of structured data, where you can find columns and rows, and every column has a specific type of data among all records**

# Data: Structured vs unstructured

**DMI**

## Examples of structured and unstructured data in electronic health records (EHRs)

| Time | dd/mm/yyyy<br>14 September 2017<br>Prescription start date: dd/mm/yy | 'Earlier today Mr X experienced chest pain'<br>'Operation scheduled on Tuesday'. |
| --- | --- | --- |

# Data: Structured vs unstructured

**DMI**

## Examples of structured and unstructured data in electronic health records (EHRs)

| Diagnosis | The patient has diabetes without complications | C0271635 |
|-----------|-----------------------------------------------|----------|

# Data: Types of data sources

**DMI**

**<u>Examples of structured and unstructured data in electronic health records (EHRs)</u>**

| Data item | Structured | Unstructured |
|---|---|---|
| Time | dd/mm/yyyy<br>14 September 2017<br>Prescription start date: dd/mm/yy | 'Earlier today Mr X experienced chest pain'<br>'Operation scheduled on Tuesday'. |
| Symptoms | N242300 Neuropathic pain<br>1B1B.00 Cannot sleep—insomnia | '…c/o shooting pain in upper right leg during the night, disturbing her sleep' |
| Diagnosis | C0271635 Type 2 diabetes without complication | 'The patient has diabetes without complications' |
| Prescription | 01040200 (BNF code for codeine phosphate 60 mg tablets) | 'Px codeine 60 mg PO qid×7 days' |
| Referral | 8H4D.00 Referral to psychogeriatrician | Rev 4 w ?refer pyscho ger |
| Test | 43F1.00 Rheumatoid factor positive | Rheumatoid factor was 42 IU/mL which is a positive result |

# Types of Data Sets

# Types of Data Sets: Record Data

**DMI**

❑ Relational records
  ❑ Relational tables, highly structured
❑ Data matrix, e.g., numerical matrix

|  | China | England | France | Japan | USA | Total |
|---|---|---|---|---|---|---|
| Active Outdoors Crochet Glove | | 12.00 | 4.00 | 1.00 | 240.00 | 257.00 |
| Active Outdoors Lycra Glove | | 10.00 | 6.00 | | 323.00 | 339.00 |
| InFlux Crochet Glove | 3.00 | 6.00 | 8.00 | | 132.00 | 149.00 |
| InFlux Lycra Glove | | 2.00 | | | 143.00 | 145.00 |
| Triumph Pro Helmet | 3.00 | 1.00 | 7.00 | | 333.00 | 344.00 |
| Triumph Vertigo Helmet | | 3.00 | 22.00 | | 474.00 | 499.00 |
| Xtreme Adult Helmet | 8.00 | 8.00 | 7.00 | 2.00 | 251.00 | 276.00 |
| Xtreme Youth Helmet | | 1.00 | | | 76.00 | 77.00 |
| Total | 14.00 | 43.00 | 54.00 | 3.00 | 1,972.00 | 2,086.00 |

❑ Document data: Term-frequency vector (matrix) of text documents

Person:

| Pers_ID | Surname | First_Name | City |
|---|---|---|---|
| 0 | Miller | Paul | London |
| 1 | Ortega | Alvaro | Valencia |
| 2 | Huber | Urs | Zurich |
| 3 | Blanc | Gaston | Paris |
| 4 | Bertolini | Fabrizio | Rom |

no relation

Car:

| Car_ID | Model | Year | Value | Pers_ID |
|---|---|---|---|---|
| 101 | Bentley | 1973 | 100000 | 0 |
| 102 | Rolls Royce | 1965 | 330000 | 0 |
| 103 | Peugeot | 1993 | 500 | 3 |
| 104 | Ferrari | 2005 | 150000 | 4 |
| 105 | Renault | 1998 | 2000 | 3 |
| 106 | Renault | 2001 | 7000 | 3 |
| 107 | Smart | 1999 | 2000 | 2 |

|  | team | coach | play | ball | score | game | win | lost | timeout | season |
|---|---|---|---|---|---|---|---|---|---|---|
| Document 1 | 3 | 0 | 5 | 0 | 2 | 6 | 0 | 2 | 0 | 2 |
| Document 2 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| Document 3 | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 |

32

# Types of Data Sets: Networks

- ❑ Protein interaction networks

- ❑ Comorbidity networks

- ❑ Molecular Structures

- ❑ Social or information networks

# Types of Data Sets: Ordered data

❑ Video data: sequence of images

❑ Temporal data: time-series

❑ Genetic sequence data

# Characteristics of Structured Data Sets

**DMI**

❑  Dimensionality: Number of attributes

❑  Sparsity: The number of empty cells

❑  Resolution: Patterns depend on the scale

❑  Distribution: Centrality and dispersion

# Data Objects

- Data sets are made up of data objects
- A data object represents an entity
- Examples:

    sales database:  customers, store items, sales

    medical database: patients, diseases, treatments

    university database: students, professors, courses
- Also called **samples , examples, instances, data points, objects**

# Data Attributes

**DMI**

- Data objects are described by attributes (also called features, dimensions)
- A data field, representing a characteristic or feature of a data object.
    - E.g., customer_ID, name, address, height, weight
- Correct representation of data:
    - Database rows → data objects
    - columns → attributes

# Attribute types:  scales of measurement

Qualitative (Categorical) Data

- Nominal:
    - Unordered set of categories identified only by name, no ranking, no magnitude
    - Hair color = {auburn, black, blond, brown, grey, red, white}, marital status, occupation, ID numbers, zip codes, blood type = {A, B, AB, O}
- Binary
    - A special case of nominal data with only two states (0 and 1)
    - Symmetric binary: both outcomes equally important
        - e.g., gender
    - Asymmetric binary: outcomes not equally important.
        - e.g., medical test (positive vs. negative)
        - Convention: assign 1 to most important outcome (e.g., HIV positive)
- Ordinal
    - Values have a meaningful order (ranking) but magnitude between successive values is not known
    - Size = {small, medium, large}, Stage of cancer  (stage I, II, III, IV), hypertension categories = {mild, moderate, severe}

# Attribute types: scales of measurement

**DMI**

Quantitative (Numeric) Data

**Interval**

- ■ Measured on a scale of **equal-sized units**

- ■ Values have order

    - ● E.g., *temperature in C˚or F˚, calendar dates, IQ scores, z-scores*

- ■ No true zero-point (zero doesn't mean the absence of value)

**Ratio**

- ■ Inherent **zero-point** ( represents the complete absence of any amount of the variable.)

- ■ We can speak of values as being an order of magnitude larger than the unit of measurement

    - ● E.g., *length, counts, heart rate, age, weight, blood concentration measurements, dosage of drugs*

# Summary

```
                          Variables
                              |
              +---------------+---------------+
              |                               |
         Qualitative                    Quantitative
              |                               |
        +-----+-----+                   +------+------+
        |           |                   |             |
```

| Categorical Data | Nominal | Ordinal | Interval | Ratio | Discrete or continuous |
|---|---|---|---|---|---|

| | Nominal | Ordinal | Interval | Ratio |
|---|:---:|:---:|:---:|:---:|
| Frequency distribution | ✔ | ✔ | ✔ | ✔ |
| Median and percentiles | ✘ | ✔ | ✔ | ✔ |
| Add or subtract | ✘ | ✘ | ✔ | ✔ |
| Mean, standard deviation | ✘ | ✘ | ✔ | ✔ |
| Ratios, coefficient of variation | ✘ | ✘ | ✘ | ✔ |

# Tidy data

```
table1
#> # A tibble: 6 × 4
#>   country      year  cases population
#>   <chr>       <dbl>  <dbl>      <dbl>
#> 1 Afghanistan  1999    745   19987071
#> 2 Afghanistan  2000   2666   20595360
#> 3 Brazil       1999  37737  172006362
#> 4 Brazil       2000  80488  174504898
#> 5 China        1999 212258 1272915272
#> 6 China        2000 213766 1280428583
```

```
table2
#> # A tibble: 12 × 4
#>   country      year type          count
#>   <chr>       <dbl> <chr>         <dbl>
#> 1 Afghanistan  1999 cases           745
#> 2 Afghanistan  1999 population  19987071
#> 3 Afghanistan  2000 cases          2666
#> 4 Afghanistan  2000 population  20595360
#> 5 Brazil       1999 cases         37737
#> 6 Brazil       1999 population 172006362
```

```
table3
#> # A tibble: 6 × 3
#>   country      year rate
#>   <chr>       <dbl> <chr>
#> 1 Afghanistan  1999 745/19987071
#> 2 Afghanistan  2000 2666/20595360
#> 3 Brazil       1999 37737/172006362
#> 4 Brazil       2000 80488/174504898
#> 5 China        1999 212258/1272915272
#> 6 China        2000 213766/1280428583
```

```
# Spread across two tibbles
table4a  # cases
#> # A tibble: 3 × 3
#>   country     `1999` `2000`
#>   <chr>        <dbl>  <dbl>
#> 1 Afghanistan    745   2666
#> 2 Brazil       37737  80488
#> 3 China       212258 213766
table4b  # population
#> # A tibble: 3 × 3
#>   country        `1999`      `2000`
#>   <chr>           <dbl>       <dbl>
#> 1 Afghanistan  19987071    20595360
#> 2 Brazil      172006362   174504898
#> 3 China      1272915272  1280428583
```

Each dataset shows the same values of four variables: country, year, population, and number of documented cases of TB (tuberculosis)

https://vita.had.co.nz/papers/tidy-data.pdf     41

# Tidy data

**DMI**

```
table1
#> # A tibble: 6 × 4
#>    country      year  cases population
#>    <chr>       <dbl>  <dbl>      <dbl>
#> 1 Afghanistan  1999    745   19987071
#> 2 Afghanistan  2000   2666   20595360
#> 3 Brazil       1999  37737  172006362
#> 4 Brazil       2000  80488  174504898
#> 5 China        1999 212258 1272915272
#> 6 China        2000 213766 1280428583
```

```
table2
#> # A tibble: 12 × 4
#>    country       year type           count
#>    <chr>        <dbl> <chr>          <dbl>
#> 1 Afghanistan  1999 cases            745
#> 2 Afghanistan  1999 population  19987071
#> 3 Afghanistan  2000 cases           2666
#> 4 Afghanistan  2000 population  20595360
#> 5 Brazil       1999 cases          37737
#> 6 Brazil       1999 population 172006362
```

Multiple variables are stored in one column

"If you find yourself using a column to describe what another column means, your data probably isn't tidy."

```
table3
#> # A tibble: 6 × 3
#>    country      year rate
#>    <chr>       <dbl> <chr>
#> 1 Afghanistan  1999 745/19987071
#> 2 Afghanistan  2000 2666/20595360
#> 3 Brazil       1999 37737/172006362
#> 4 Brazil       2000 80488/174504898
#> 5 China        1999 212258/1272915272
#> 6 China        2000 213766/1280428583
```

Multiple variables are stored in one column

```
# Spread across two tibbles
table4a  # cases
#> # A tibble: 3 × 3
#>    country     `1999` `2000`
#>    <chr>        <dbl>  <dbl>
#> 1 Afghanistan    745   2666
#> 2 Brazil       37737  80488
#> 3 China       212258 213766
table4b  # population
#> # A tibble: 3 × 3
#>    country         `1999`      `2000`
#>    <chr>            <dbl>       <dbl>
#> 1 Afghanistan   19987071    20595360
#> 2 Brazil       172006362   174504898
#> 3 China       1272915272 1280428583
```

Column headers are values, not variable names

A single observational unit is stored in multiple tables

https://vita.had.co.nz/papers/tidy-data.pdf   42

## Tidy data

```
table1
#> # A tibble: 6 × 4
#>   country      year  cases population
#>   <chr>       <dbl>  <dbl>      <dbl>
#> 1 Afghanistan  1999    745   19987071
#> 2 Afghanistan  2000   2666   20595360
#> 3 Brazil       1999  37737  172006362
#> 4 Brazil       2000  80488  174504898
#> 5 China        1999 212258 1272915272
#> 6 China        2000 213766 1280428583
```

There are three interrelated rules which make a dataset tidy:

1. Each variable must have its own column.
2. Each observation must have its own row.
3. Each value must have its own cell.

"Happy families are all alike; every unhappy family is unhappy in its own way." –– Leo Tolstoy
"Tidy datasets are all alike, but every messy dataset is messy in its own way." –– Hadley Wickham

https://vita.had.co.nz/papers/tidy-data.pdf

# The Data Mining process

Raw data not always the best data set for data mining

https://doi.org/10.1609/aimag.v17i3.1230

# Data preprocessing and transformation

**DMI**

- Dealing with missing data

  ○ replace all missing values with a single global constant (highly application dependent)

  ○ replace a missing value with its feature mean

  ○ replace a missing value with its feature mean for the given class (this approach is possible only for classification problems)

  ○ generate a predictive model to predict each of the missing values

- Outlier detection

  ○ Detect and eventually remove outliers as a part of the preprocessing phase

  ○ Develop robust modeling methods that are insensitive to outliers.

- Data transformation

  ○ Scaling, encoding, normalization, and selecting features

# The Data Mining process

# Types of analysis

patterns, trends, or relationships between variables.

Did they summarize the data? — *Yes* → Did they report the summaries without interpretation? — *No* → Did they quantify whether the discoveries are likely to hold in a new sample? — *Yes* → Are they trying to figure out how changing the average of one measurement affects another?

*No* ↓ Not a data analysis

*Yes* ↓ Descriptive

*No* ↓ Exploratory

*No* / *Yes* →

Is the effect they are looking for an average effect or a deterministic effect?

Are they trying to predict measurement(s) for individuals?

*No* / *Yes*

*Average* / *Deterministic*

Inferential    Predictive    Causal    Mechanistic

http://science.sciencemag.org/content/347/6228/1314

make a statement about something outside the dataset

Use a set of features to predict another feature on a given sample

what happens to one measurement if you make another change

not only understand that there is an effect, but how that effect operates

47

# Data Mining Tasks

**DMI**

<u>**The two high-level primary goals of data mining tend to be description and prediction.**</u>

**Descriptive Methods**

- Find human-interpretable patterns that describe the data (correlations, trends, clusters, trajectories, and anomalies).

**Predictive Methods**

- Use some features (variables) to predict and unknown or future value of other variable (target variable).

> Many data analyses answer multiple types of questions.
> The type of question you ask is determined in part by the data available to you

# Basic data mining tasks

Regression

Classification

Cluster Analysis

Predictive Modeling

Data

| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|----|-----------|----------------|---------------|--------------------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 80K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

Association Analysis

Anomaly Detection

Milk

DIAPER

Introduction to Data Mining by Pang-Ning Tan, Michael Steinbach and Vipin Kumar, Addison Wesley,

# Basic data mining tasks

**DMI**



**Predictive modeling:** building a model for the target variable as a function of the explanatory variables.

# Classification

Breast Cancer (Malignant / Benign)

Classification is the process of predicting the class of given data points.
Best used when the output has finite and discrete values
Classes are sometimes called as targets/ labels or categories.

**Types of Classification Algorithms**

- Linear Models
  - **Logistic Regression**
  - Support Vector Machines

- Nonlinear models
  - **K-nearest Neighbors (KNN)**
  - Kernel Support Vector Machines (SVM)
  - Naïve Bayes
  - **Decision Tree Classification**
  - **Random Forest Classification**

# Basic data mining tasks

**DMI**



Regression

**Predictive modeling: building a model for the target variable as a function of the explanatory variables.**

Classification

| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|----|------------|----------------|---------------|--------------------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

Data

Cluster Analysis

Predictive Modeling

Association Analysis

Anomaly Detection

Milk

DIAPER

# Regression

**Set of machine learning methods that allow us to predict a continuous outcome variable (y) based on the value of one or multiple predictor variables (x).**

**Examples of regression algorithms**

- Simple Linear Regression
- Polynomial Regression
- Support Vector Regression
- Decision Tree Regression
- Random Forest Regression

Some algorithms can be used for both classification and regression with small modifications, such as decision trees and artificial neural networks.

53

# Basic data mining tasks

Descriptive modeling:
the objective is to derive
patterns (correlations,
trends, clusters, trajectories,
and anomalies)

# Basic data mining tasks

**DMI**



Descriptive modeling:
the objective is to derive
patterns (correlations,
trends, clusters, trajectories,
and anomalies)

# Clustering

the goal of clustering is to create clusters (groups with the same or similar features)



Intracluster distances are minimized

Intercluster distances are maximized

# Basic data mining tasks

**DMI**



Association analysis is used to discover patterns that describe strongly associated features in the data.

# Basic data mining tasks

**DMI**



Anomaly detection focuses on finding deviations in data according to similar data collected in the past or by typical values of these data.

# Applications of Data Mining in Biomedicine

59

**DMI**

# Best Practices for Data Science Projects

# Best Practices for Data Science Projects

**DMI**

- Develop your own system
- Put everything in a common directory
- Be consistent
  - directory structure; names
- Separate raw from processed data
- Separate code from data
- It should be obvious what code created what files, and what the dependencies are.
- **No hand-editing of data files !!!!!**
- All aspects of data cleaning should be in scripts
- No manual steps!  Each aspect of the analysis should be in scripts.
- Save your seeds for random number generation.
- Make regular backups of your data. In multiple locations.
- Use relative paths (../blah), not absolute paths ( ~/blah [or] /users/blah)

> "Your closest collaborator is you six months ago, but you don't reply to emails."

# Best Practices for Data Science Projects

## Organizing a data mining project

- README.md
- .gitignore
- project.Rproj
- data/
  - raw_data/
  - processed_data/
  - codebook.md
- code/
  - raw_code/
  - final_code/
- figures/
  - exploratory_figures/
  - explanatory_figures/
- products/
  - writing/

```
Q  Go to file
>  📁 .github
∨  📁 assignment-1-data
   ∨  📁 covid_data
         📄 covid_data.XLSX
         📄 covid_metadata.XLSX
   ∨  📁 heart_disease_data
         📄 heart_disease_dataset.csv
         📄 heart_disease_description.txt
   📄 Hands_on_I.Rmd
   📄 Hands_on_I.html
   📄 README.md
```

"File organization and naming are powerful weapons against chaos."
— @JennyBryan

Dear past-Hadley: PLEASE COMMENT YOUR CODE BETTER. Love present-Hadley
— @hadleywickham

# Best Practices for Data Science Projects

Naming files



https://xkcd.com/1459/

# Best Practices for Data Science Projects

## Naming files

- **Machine readable:** File names **shouldn't contain spaces, special characters**, and should include important pieces of information about the file contents (sometimes called slugs) separated by underscores or dashes. It is often easier if file names are entirely lowercase letters.

- **Human readable:** files should be labeled with names that make it easy for a person to follow along. So err on the side of **long file names** rather than abbreviations, numeric only file names, or obscure file names.

- **Be nicely ordered:** it is useful to see the files in the order you want them to be run. One way to do this is to name files alpha-numerically so that they will appear in the right order on your computer. One way to do this is to order the files and then append a number to the beginning of each file name (like 01_data_cleaning.R, 02_exploratory_analysis.R etc. so that the files will be ordered correctly)

# Best Practices for Data Science Projects

**Naming files**

| Bad Naming |
|---|
| 2013 my report.md |
| malik's_report.md |
| 01_zoë_report.md |
| AdamHooverReport.md |
| executivereportpepsiv1.md |

# Best Practices for Data Science Projects

## Naming files

| Bad Naming | Good Naming |
|---|---|
| 2013 my report.md | 2013_my_report.md |
| malik's_report.md | maliks_report.md |
| 01_zoë_report.md | 01_zoe_report.md |
| AdamHooverReport.md | adam-hoover-report.md |
| executivereportpepsiv1.md | executive_report_pepsi_v1.md |

# Best Practices for Data Science Projects

## Naming variable and functions

Variable and function names should be lowercase. Use an underscore (_) or hyphen to separate words within a name. **Be consistent!** Generally, variable names should be nouns and function names should be verbs. Strive for names that are concise and meaningful (this is not easy!).

```
# Good ✔            # Bad ✘
day_one             first_day_of_the_month
day_1               DayOne
                    dayone
                    djm1
```

Where possible, avoid using names of existing functions and variables. Doing so will cause confusion for the readers of your code.

```
c <- 10 ✘
mean <- function(x) sum(x) ✘
```

67

# Best Practices for Data Science Projects

## Automating analysis



"I SPEND A LOT OF TIME ON THIS TASK.
I SHOULD WRITE A PROGRAM AUTOMATING IT!"

https://xkcd.com/1319/

**Automate a pipeline**

- to reproduce previous results.
- to recreate results deleted by mistake
- to rerun the pipeline with updated software.
- to run the same pipeline on a new data set.

# Best Practices for Data Science Projects

## Version Control



"Version control is a truly vital concept that has unfortunately been implemented by madmen."

# Best Practices for Data Science Projects

## Version Control: Git

- A version control system: manages the evolution of a set of files – called a repository – in a sane, highly structured way ("the *Track Changes* features from Microsoft Word on steroids." Bryan, 2017 )
- Originated to help groups of developers work collaboratively on big software projects.
- Re-purposed by the data science community to manage source code, and other files such as of data, figures, and reports.
- Help communicating and collaborating with other people.

# Best Practices for Data Science Projects

## Version Control: Github

- GitHub complements Git by providing a slick user interface and distribution mechanism for Git repositories.
- Git is the software you will use locally to record changes to a set of files. GitHub is a hosting service that provides a Git-aware home for such projects on the internet.
- GitHub is like DropBox or Google Drive, but more structured, powerful, and programmatic.
- Available at  https://github.com
- Create your account if you don't have one yet!

# Best Practices for Data Science Projects

## Version Control: Git

In the root directory, create a file called  .gitignore

And write in the file:

.Rproj.user

.Rhistory

.RData

.Ruserdata

<span style="color:red">data/ -> for other cases, not the practical sessions of DMI!!</span>

# Best Practices for Data Science Projects

**DMI**

## Version Control: Git

*git add Hands_on_I.Rmd*

*git add Hands_on_I.html*

*git commit –m "first commit"*

*git status*

*git push*

To check repository status and history:
**git status** ->shows if there are  changes to any files that have not been committed
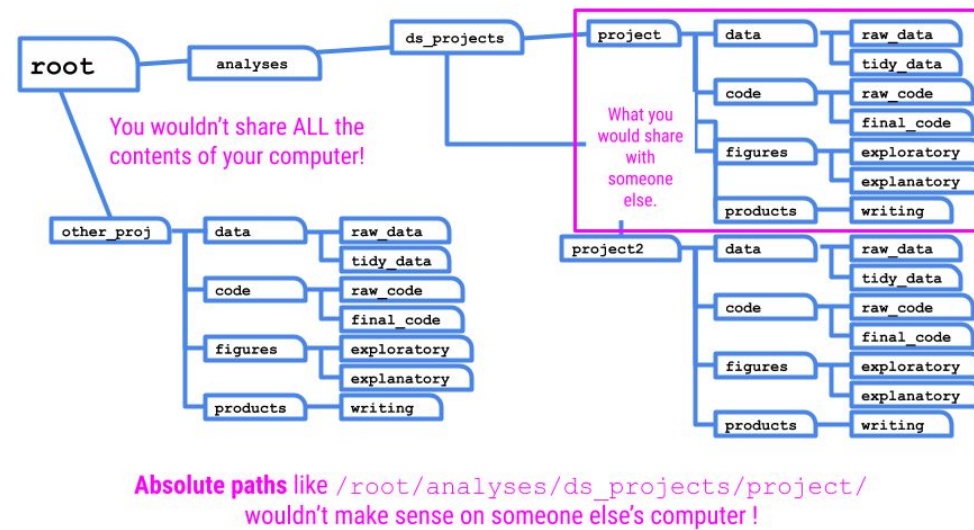**git diff file** -> shows differences between last snapshot and the current file
**git restore file** -> discards local changes

# Best Practices for Data Science Projects

### Dealing with paths

# Handling expectations

## Observing many researchers using the same data and hypothesis reveals a hidden universe of uncertainty

Results from our controlled research design in a large-scale crowdsourced research effort involving 73 teams demonstrate that **analyzing the same hypothesis with the same data can lead to substantial differences in statistical estimates and substantive conclusions.** In fact, <u>no two teams arrived at the same set of numerical results or took the same major decisions during data analysis</u>. Our finding of outcome variability echoes those of recent studies involving many analysts undertaken across scientific disciplines. The study reported here differs from these previous efforts because it attempted to catalog every decision in the research process within each team and use those decisions and predictive modeling to explain why there is so much outcome variability. Despite this highly granular decomposition of the analytical process, we could only explain less than 2.6% of the total variance in numerical outcomes. We also tested if expertise, beliefs, and attitudes observed among the teams biased results, but they explained little. Even highly skilled scientists motivated to come to accurate results varied tremendously in what they found when provided with the same data and hypothesis to test.

https://doi.org/10.1073/pnas.2203150119

# What distinguishes data science from statistics?

# What distinguishes data science from statistics?

**DMI**

| Data Mining | Statistics |
|---|---|
| Explore and gather data first, builds model to detect patterns and make theories. | Usually starts with a hypothesis |
| Data used is Numeric or Non numeric. | Data used is Numeric. |
| Inductive Process (Generation of new theory from data) | Deductive Process (Does not involve making any predictions) |
| Data collection is less important | Data collection is more important. |
| Data cleaning is important | Clean data is used to apply statistical method. |
| Suitable for large data sets | Suitable for smaller data sets |
| Algorithm "learns" from data, no explicit programming rules | Formalization of the relationships in data as a mathematical equation |
| Examples: Classification, Clustering, Neural network, Association, Estimation, Sequence based analysis, Visualization | Examples: Descriptive Statistical, Inferential Statistical |

# Bibliography

DMI

**Introduction to Data Mining**, book by Michael Steinbach, Pang-Ning Tan, and Vipin Kumar

**The Elements of Statistical Learning**, book by Jerome H. Friedman, Robert Tibshirani, and Trevor Hastie

**Excuse me, do you have a moment to talk about version control?** (2017) preprint by Jennifer Bryan

Tidy Data, paper by Hadley Wickham

## Github assignments

**DMI**

Go to this page and accept the assignment

Go to this sheet and enter your GitHub handle

Choose your scientific article with your partner and update it in the sheet

# Hands-on sessions

We will work with the publication **Proteomic and Metabolomic Characterization of COVID-19 Patient Sera**