# Dados e Aprendizagem Automática

## Intro to Data Science & Python/Scikit-learn

DAA @ MEI-1º/MiEI-4º – 1º Semestre

César Analide, Bruno Fernandes, Filipa Ferraz, Filipe Gonçalves, Victor Alves

Part I

# Contents

- Data Types

- Mean, Median & Mode

- Standard Deviation & Variance

- Probability Density Functions

- Percentiles

- Covariance & Correlation

- Virtual Environment

- Environment Setup

- Hands On

# Data Types

# Data Types

- Major Types of Data:
  - Numerical
  - Categorical
  - Ordinal

# Data Types

**Numerical**

- Represents some sort of quantitative measurement
  - Heights of people, page load times, stock prices, etc.

- Discrete Data
  - Integer based; often counts of some event.
    - How many purchases did a customer make in a year?
    - How many times did I flip "heads"?

- Continuous Data
  - Has an infinite number of possible values
    - How much time did it take for a user to check out?
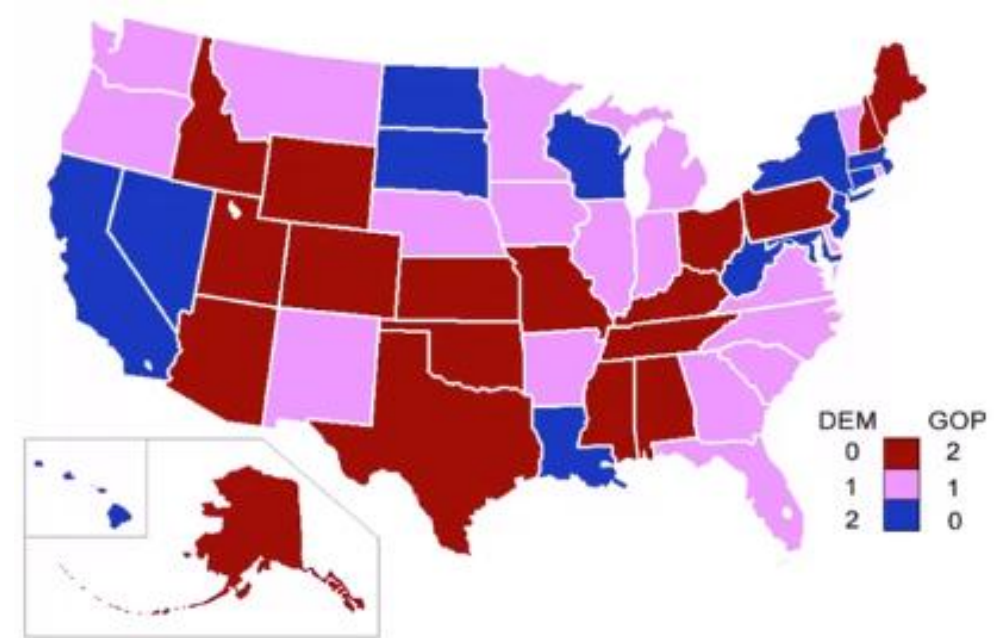    - How much rain fell on a given day?

# Data Types

## Categorical

☐ Qualitative data that has no inherent mathematical meaning

    ○ Gender, Yes/No (Binary Data), Race, State of Residence, Product Category, Political Party, etc.

☐ You can assign numbers to categories in order to represent them more compactly, but the numbers don't have mathematical meaning

# Data Types

**Ordinal**

- A mixture of numerical and categorical
- Categorical data that has mathematical meaning
- Example: movie ratings on a 1-5 scale.
  - Ratings must be 1,2,3,4 or 5
  - These values have mathematical meaning; 1 means it's a worse movie than a 2.

# Data Types

**Quick Quiz:**

- Are the following types of data numerical, categorical, or ordinal?
  - How much gas is in your gas tank?
  - A rating of your overall health where the choices are 1,2,3 or 4, corresponding to "poor", "moderate", "good" and "excellent"
  - The nationalities of your classmates
  - Ages in years
  - Money spent in a store

# Mean, Median & Mode

# Mean, Median & Mode

**Mean**

- AKA Average
- Sum / number of samples
- Example:
  - Number of children in each house on my street:

0, 2, 3, 2, 1, 0, 0, 2, 0

The MEAN is (0+2+3+2+1+0+0+2+0) / 9 = **1.11**

# Mean, Median & Mode

## Median

☐ Sort the values, and take the value at the midpoint.

☐ Example:

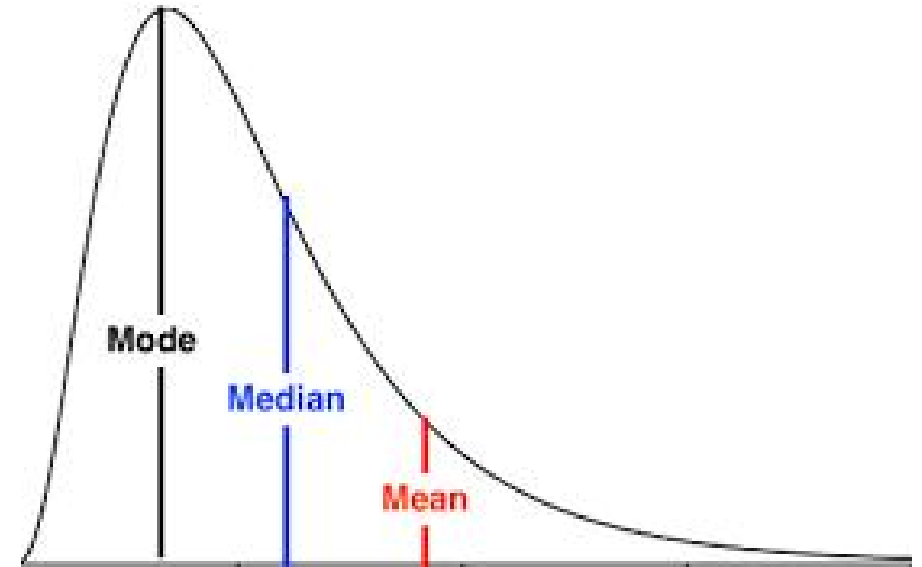0, 2, 3, 2, 1, 0, 0, 2, 0
Sort it:
0, 0, 0, 0, 1, 2, 2, 2, 3
⬆

☐ If you have an even number of samples, take the average of the two in the middle.

# Mean, Median & Mode

## Median

□ Median is less susceptible to outliers than the mean

- o Example: mean household income in the USA is $72,641, but the median is only $51,939 – because the mean is skewed by a handful of billionaires.

- o Median represents better the "typical" American in this example.

# Mean, Median & Mode

## Mode

☐ The most common value in a dataset

- o Not relevant to continuous numerical data

☐ Number of kids in each house example:

0, 2, 3, 2, 1, 0, 0, 2, 0
How many of each value are there?
0: 4, 1: 1, 2: 3, 3: 1
The MODE is 0

# Standard Deviation & Variance

# Standard Deviation & Variance

An example of a histogram…

# Standard Deviation & Variance

**Variance measures how "spread-out" the data is.**

☐ Variance ($\delta^2$) is simply the average of the squared differences from the mean

☐ Example: What is the variance of the data set (1, 4, 5, 4, 8)?

- First find the mean: (1+4+5+4+8) / 5 = 4.4
- Now find the difference from the mean: (-3.4, -0.4, 0.6, -0.4, 3.6)
- Find the squared differences: (11.56, 0.16, 0.36, 0.16, 12.96)
- Find the average of the squared differences:
- $\delta^2$ = (11.56+0.16+0.36+0.16+12.96) / 5 = 5.04

# Standard Deviation & Variance

Standard Deviation δ is the square root of the variance.

- Case Study = (1,4,5,4,8)
- Mean = 4.4
- $δ^2$ = 5.04
- δ = 2.24

- Stand. Dev. Is usually used as a way to identify outliers.
- Data points that lie more than one standard deviation from the mean can be considered unusual.

- You can talk about how extreme a data point is by talking about "how many sigmas" away from the mean it is.

# Probability Density Functions

# Probability Density Functions

## Example: a "normal distribution"

- Gives you the probability of a data point falling within some given range of a given value

- Based on histogram values, a normal probability density function can be calculated

# Probability Density Functions

## Example: Probability Mass Function

- Used for discrete data
- Based on histogram values, a normal probability density function can be calculated

# Percentiles

# Percentiles

## Percentiles

- In a dataset, what's the point at which X% of the values are less than that value?

- Example: income distribution
  - Take all incomes from a country's population and sort them
  - 99th percentile represents the income amount in which 99% of the population gains less then that value (i.e., $506,553)

# Percentiles

## Percentiles in a normal distribution

□ Between Quartil 1 & Quartil 3 represents 50% of the data distribution

□ IQR (Inter-Quartil Range) represents the area in the middle of the distribution (where data is more focused)

# Covariance & Correlation

# Covariance & Correlation

## Covariance

□ Measures how two variables vary in tandem from their means.

□ i.e. how 2 attributes depend on each other (left plot – low covariance / right plot – high covariance)

# Covariance & Correlation

## Measuring covariance

- Think of the datasets for the two variables as high-dimensional vectors

- Convert these to vectors of variances from the mean

- Take the dot product (cosine of the angle between them) of the two vectors

- Divide by the population size

Population Covariance Formula

$$Cov(x,y) = \frac{\Sigma(x_i - \bar{x})(y_i - \bar{y})}{N}$$

Sample Covariance

$$Cov(x,y) = \frac{\Sigma(x_i - \bar{x})(y_i - y)}{N-1}$$

# Covariance & Correlation

**Interpreting <span style="color:red">covariance</span> is hard**

- Small covariance (close to 0) means there isn't much correlation between the two variables
- Large covariance (far from 0 – can be negative for inverse relationships) means that there is a correlation

**Interpreting <span style="color:red">correlation</span> is easier**

- Normalization value of covariance divided by the standard deviations of both variables
  - Correlation of -1: perfect inverse correlation
  - Correlation of 0: no correlation
  - Correlation of 1: perfect correlation

# Covariance & Correlation

**Correlation does not imply causation!**

- Only a controlled, randomized experiment can give you insights on causation.

- Use correlation to decide what experiments to conduct!

# Covariance & Correlation

STRONG POSITIVE CORRELATION

WEAK POSITIVE CORRELATION

STRONG NEGATIVE CORRELATION

WEAK NEGATIVE CORRELATION

MODERATE NEGATIVE CORRELATION

NO CORRELATION

# Virtual Environments

# Virtual Environments

- Virtual Environments allow you to set up virtual installations of Python and libraries on your computer

- You can have multiple versions of Python or libraries and easily activate or deactivate these environments

- Let's see some examples of why you may want to do this

# Virtual Environments

- Sometimes you'll want to program in different versions of a library
- For example:
  - You develop a program with SciKit-Learn 0.17
  - SciKit-Learn 0.18 is released
  - You want to explore 0.18 but don't want your old code to break

- Sometimes you'll want to make sure your library installations are in the correct location

- For example:
  - You want multiple versions of Python on your computer
  - You want one environment with Python 2.7 and another with Python 3.6

# Virtual Environments

- Anaconda has a built-in virtual environment manager that makes the whole process really easy
- Check out the resource link for the official documentation:
    - http://conda.pydata.org/docs/using/envs.html

- Command Prompt Example (create env. and activate it):

```
conda create --name mypython3version python=3.6 numpy
conda info --evns
activate mypython3version
python
import numpy as np
import pandas as pd      -> Error
quit()
conda install pandas
deactivate
```

# Anaconda Distribution

The World's Most Popular Python/R Data Science Platform

- FOSS

- Share, collaborate on, and reproduce projects

- Highly supported by the community

- Conda, a package, dependency and environment manager
  - Easily create, save, load and switch between environments
  - Easily install, update and run any package (and its dependencies… automatically!)


- Anaconda provides two user clients
  - Anaconda Navigator
  - Anaconda Prompt (or the terminal on Linux and macOS)

# ANACONDA

### Conda basics

| | |
|---|---|
| Verify conda is installed, check version number | `conda info` |
| Update conda to the current version | `conda update conda` |
| Install a package included in Anaconda | `conda install PACKAGENAME` |
| Run a package after install, example Spyder* | `spyder` |
| Update any installed program | `conda update PACKAGENAME` |
| Command line help | `COMMANDNAME --help`<br>`conda install --help` |

### Using environments

| | |
|---|---|
| Create a new environment named py35, install Python 3.5 | `conda create --name py35 python=3.5` |
| Activate the new environment to use it | `WINDOWS:    activate py35`<br>`LINUX, macOS: source activate py35` |
| Get a list of all my environments, active environment is shown with * | `conda env list` |
| Make exact copy of an environment | `conda create --clone py35 --name py35-2` |
| List all packages and versions installed in active environment | `conda list` |
| List the history of each change to the current environment | `conda list --revisions` |
| Restore environment to a previous revision | `conda install --revision 2` |
| Save environment to a text file | `conda list --explicit > bio-env.txt` |
| Delete an environment and everything in it | `conda env remove --name bio-env` |
| Deactivate the current environment | `WINDOWS: deactivate`<br>`macOS, LINUX: source deactivate` |
| Create environment from a text file | `conda env create --file bio-env.txt` |
| Stack commands: create a new environment, name it bio-env and install the biopython package | `conda create --name bio-env biopython` |

### Installing and updating packages

| | |
|---|---|
| Install a new package (Jupyter Notebook) in the active environment | `conda install jupyter` |
| Run an installed package (Jupyter Notebook) | `jupyter-notebook` |
| Install a new package (toolz) in a different environment (bio-env) | `conda install --name bio-env toolz` |
| Update a package in the current environment | `conda update scikit-learn` |
| Install a package (boltons) from a specific channel (conda-forge) | `conda install --channel conda-forge boltons` |
| Install a package directly from PyPI into the current active environment using pip | `pip install boltons` |
| Remove one or more packages (toolz, boltons) from a specific environment (bio-env) | `conda remove --name bio-env toolz boltons` |

### Managing multiple versions of Python

| | |
|---|---|
| Install different version of Python in a new environment named py34 | `conda create --name py34 python=3.4` |
| Switch to the new environment that has a different version of Python | `Windows:    activate py34`<br>`Linux, macOS:    source activate py34` |
| Show the locations of all versions of Python that are currently in the path<br>NOTE: The first version of Python in the list will be executed. | `Windows:    where python`<br>`Linux, macOS: which -a python` |
| Show version information for the current active Python | `python --version` |

# Environment Setup

# Environment Setup

- This course will use Jupyter Notebooks/spyder for teaching and to provide notes
  - **Note**: you are free to use **whatever development environment you prefer (e.g., Spyder, PyCharm, ..)**

- We will be using the Python 3.6 for this course through the Anaconda Distribution
- Now let's go over your installation options for Jupyter Notebook!

# Environment Setup

- ☐ For experienced users who already have Python
  - o As an existing Python user, you may wish to install Jupyter and required APIs using Python package manager pip, instead of Anaconda
  - o Just go to your command prompt or terminal and use:

**pip install jupyter**

- ☐ For new users, we highly recommend installing Anaconda
  - o Anaconda conveniently installs Python, the Jupyter Notebook, and other commonly used packages for scientific computing and data science
  - o Let's go to [www.jupyter.org](www.jupyter.org) to walkthrough the installation steps!

# Environment Setup

# Environment Setup
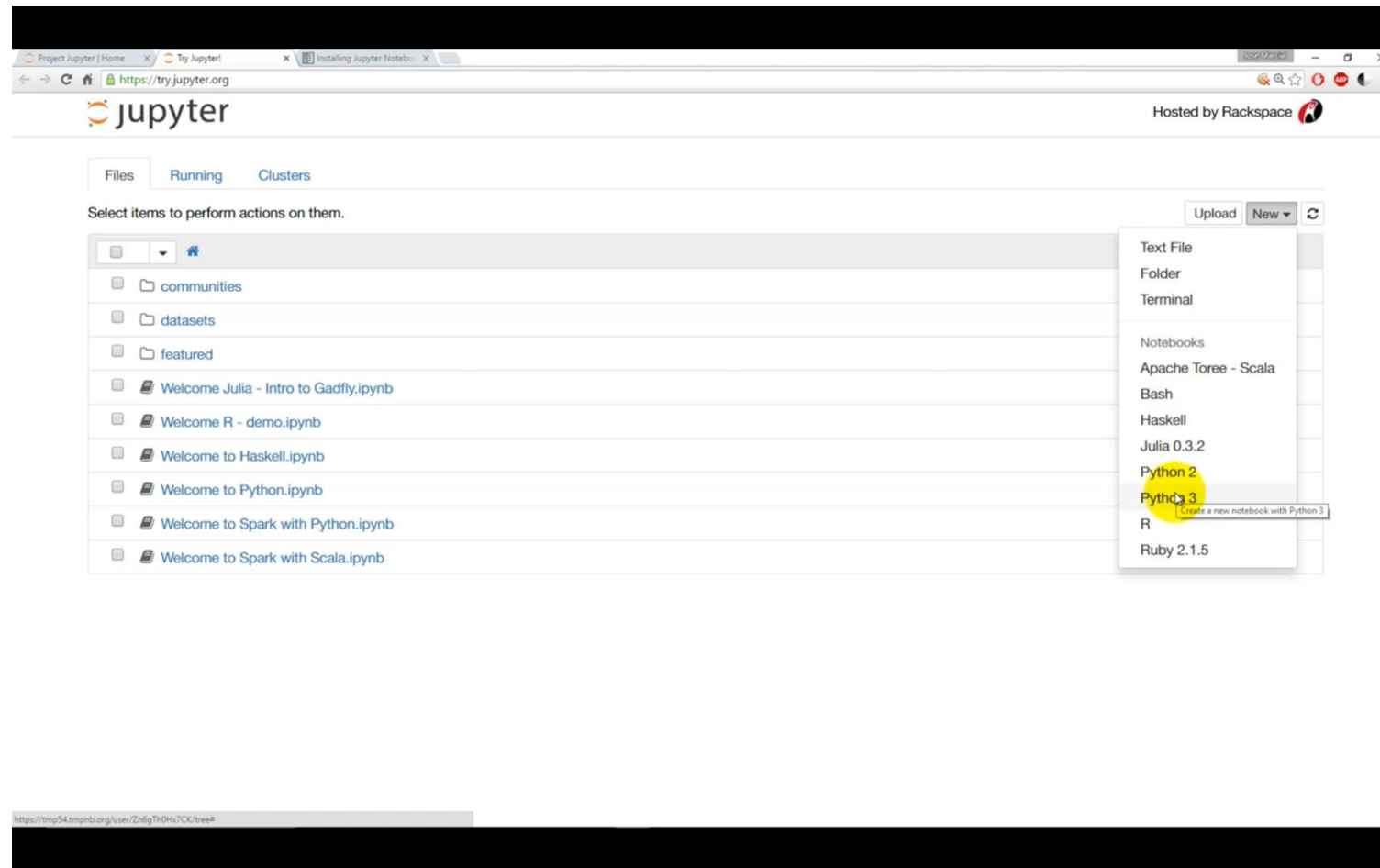
# Environment Setup

# Environment Setup

# Environment Setup

# Hands On

# Hands On

**T1**

- We will use scikit-learn/sklearn (Anaconda – Python Management Environment). Download and install the Anaconda Python package for your respective platform (Windows, Mac OS, Linux). The platform is available at https://www.anaconda.com/

  - Anaconda – Python 3.6

  - Deep Learning Libraries **not** required (Theano, Tensorflow, Keras)

  - Required to install Python IDE

  - https://machinelearningmastery.com/setup-python-environment-machine-learning-deep-learning-anaconda/

# Hands On

**T2**

- Start Anaconda and create a virtual Python3.6 environment:
  - Open Terminal & Execute:
    - *conda create --name ==env== python==3.6== numpy pandas xlrd xlwt matplotlib seaborn scikit-learn ==jupyterlab==*
  - To install packages, enter the env. and execute: *conda install PACKAGENAME*
  - To work inside the python environment, execute: *conda activate ==env==*
  - To exit python environment, execute: *conda deactivate*

# Hands On

**T2**

- In this environment, the following libraries must be installed:

  a. Numpy

  b. Pandas

  c. Xlrd

  d. Xlwt

  e. Matplotlib

  f. Seaborn

  g. Scikit-learn

  h. Jupyterlab

# Hands On

**T3**

☐ Activate the created virtual environment and check the installed libraries; validate the installation of the set of libraries presented in T2

**T4**

☐ Briefly check the documentation for each library mentioned in question T2. Identify its relevance in the context of Machine Learning algorithm development.