# Machine Learning Workshop

Emanuela Boros
University of La Rochelle, France

4 February 2021

# Organization

1. Text pre-processing & classification
2. Machine learning basics
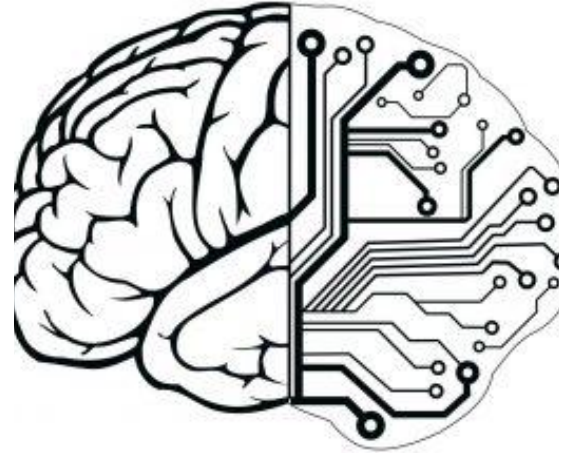3. Deep learning examples (extra)

- #ia (Discord)
- Python 3.7+
- Docker (or not, depends)
- Github (github classrooms, invitation link)
- Jupyter Notebook
- Computation libraries (numPy)
- Data libraries (pandas)
- NLP libraries (NLTK, spaCy)
- Ml libraries (scikit-learn, tensorflow, keras, etc)

# Object Classification
01

- **Human brains** are wired to **recognize patterns** and classify objects for learning and making decisions
- .. they are not able to treat every object as **unique**
- .. we don't have a **lot of memory resources** to be able to process the world around us
- → our brains develop "concepts" or mental representations of "categories of objects"

- **Classification** is **fundamental** in language, prediction, inference, decision making and all kinds of environmental interactions
- **Language**: for example, how the meaning of words in a sentence can be contextualized by earlier words or concepts

**QUICK COMMENT:**

AI can be sexist and/or racist
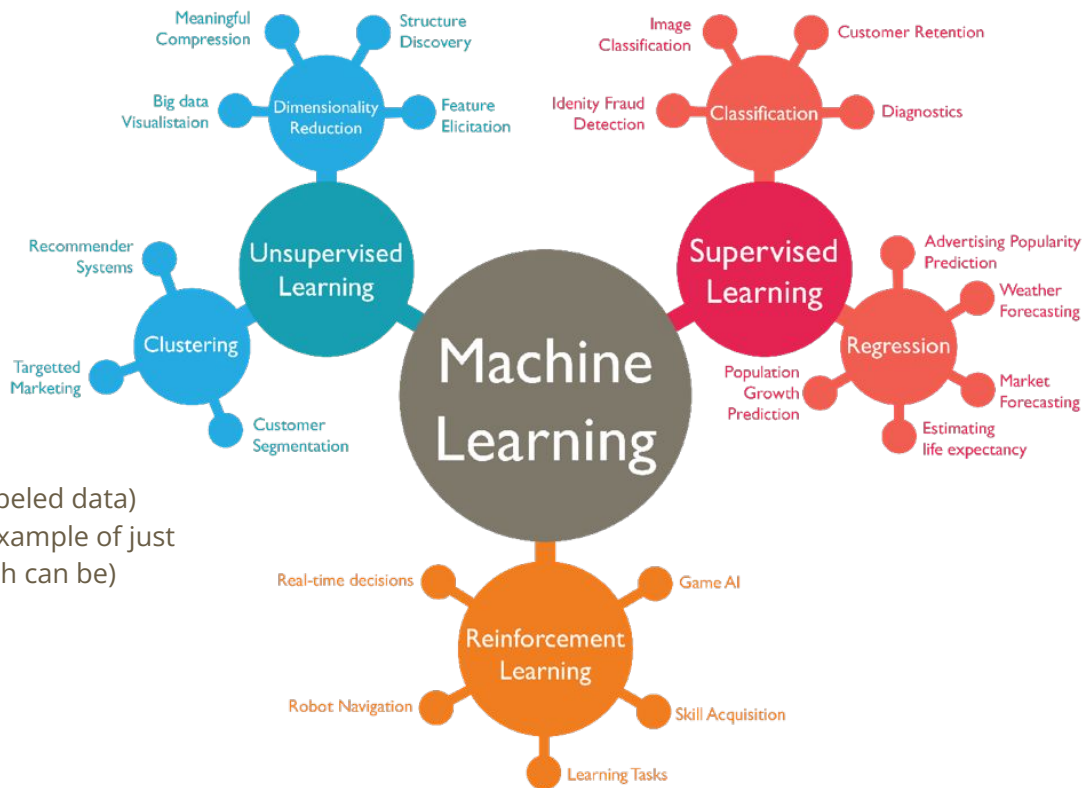Racist data? It's the human bias that is Infecting the AI development

# Object Classification

## 02

- The **classification of objects** consists of assigning a class to an object.
- These objects can be of the type:

<span style="background-color:orange;color:white">Text</span> <span style="background-color:orange;color:white">image</span> <span style="background-color:orange;color:white">audio</span> <span style="background-color:orange;color:white">video</span>

- We do classification all the time:
  - We can **recognize the way back home** from university
  - We can r**ecognize a cat that is black** even if we have only seen white and orange cats before
  - We can even distinguish between a **Chihuahua** and a **muffin**

# Concepts
## 01

- **Corpus/dataset** (corpora/datasets)
  - *Spam or not? Cat or dog?*
  - *Positive, negative or neutral?*
  - *House price prediction*
- **Machine/Deep learning types** of algorithms
  - **Supervised** learning (labeled data)
  - **Unsupervised** learning (unlabeled data)
  - **Semi-supervised** learning (labeled+unlabeled data)
  - **Reinforcement** learning (rules, a good example of just how broad the overall "learning" approach can be)
- **Machine/Deep learning evaluation**
  - **Train development test**
  - Evaluation metrics

# Dataset: Tobacco

*The US government has sued five major US tobacco companies for raising large profits by lying about the dangers of smoking. The tobacco companies agreed in 1953 to "jointly carry out a vast public relations campaign to counter the increasingly obvious evidence of a link between tobacco consumption and serious illness".*

*In this trial 6,910,192 documents were collected and digitized. In order to facilitate the use of these documents by lawyers, you are in charge of setting up an automatic classification of the types of documents:*

**Advertisement**   **Email**   **Form**   **Letter**   **Memo**

**News**   **Note**   **Report**   **Resume**   **Scientific**
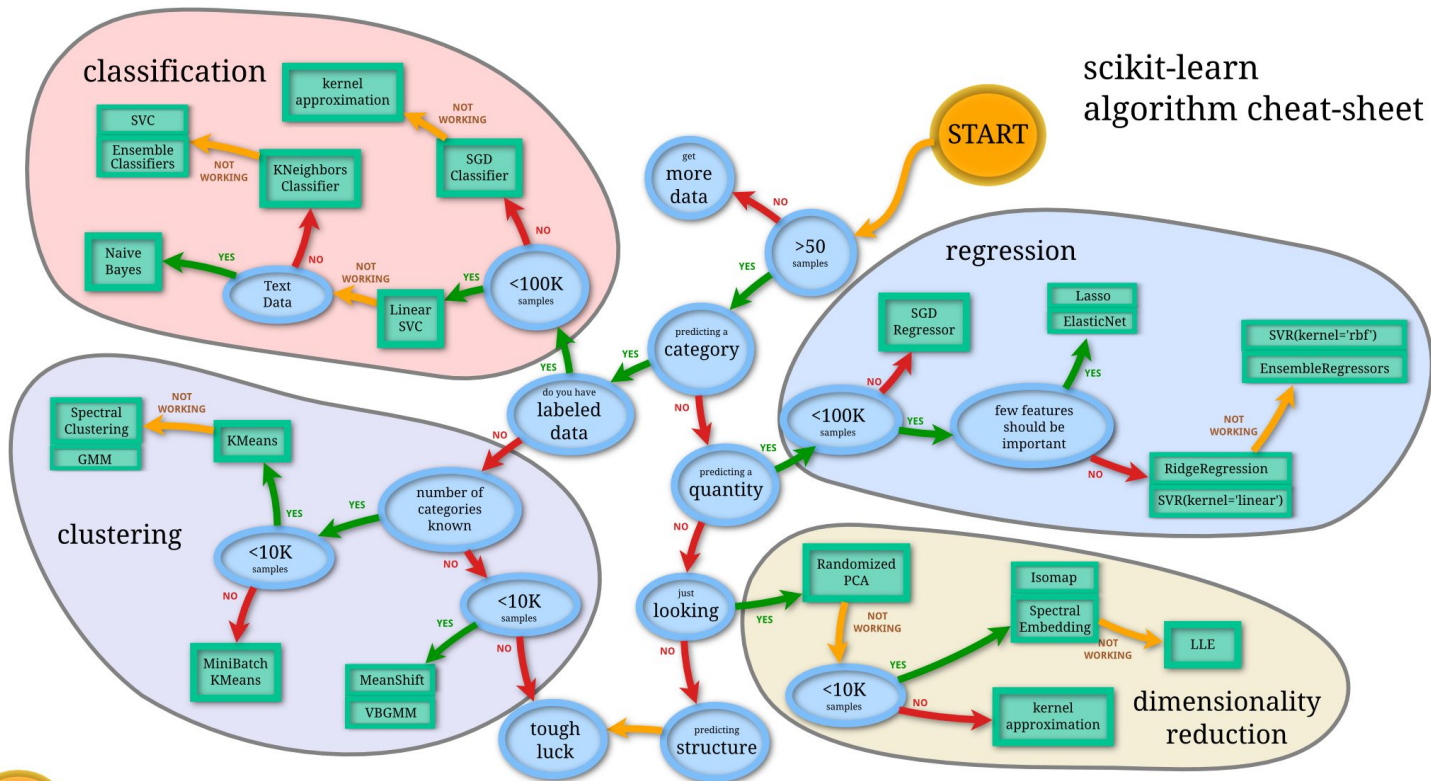


http://tc11.cvc.uab.es/datasets/Tobacco800_1
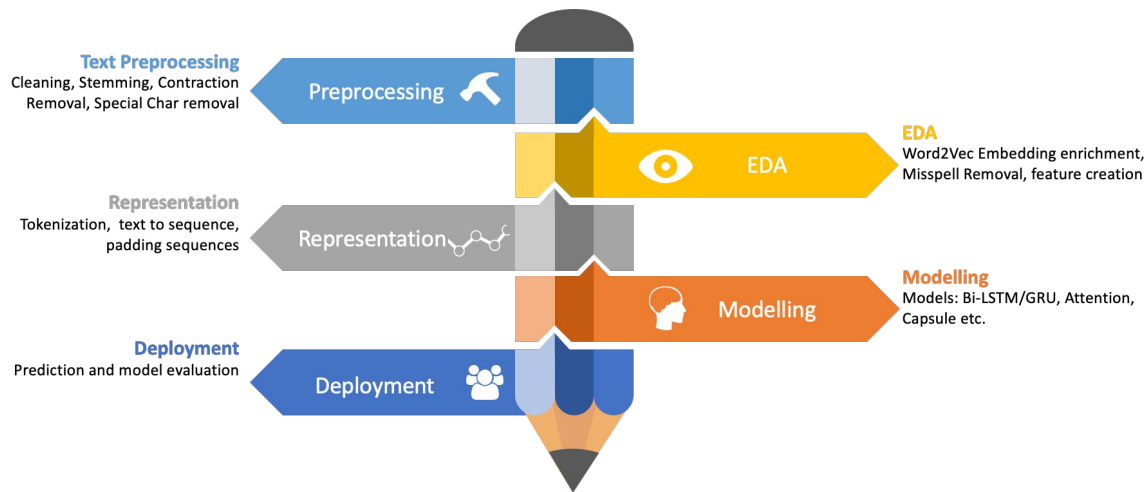David Doermann, Tobacco 800 Dataset (Tobacco800)

# Scikit-learn



scikit-learn
algorithm cheat-sheet

# Overview

1. **Recovery of the text corpus** (* .csv, * .txt, * .json, etc.)
2. **Text pre-processing**: tokenization etc
3. **Exploration** of the corpus (EDA, exploratory data analysis): frequency analysis
4. **Word representations** (bag of words, TF-IDF, word embeddings, word embeddings)
5. **Machine learning and deep learning methods**
6. **Error analysis** (evaluation, etc.)

**Text Preprocessing**
Cleaning, Stemming, Contraction Removal, Special Char removal

Preprocessing

**EDA**
Word2Vec Embedding enrichment, Misspell Removal, feature creation

EDA

**Representation**
Tokenization, text to sequence, padding sequences

Representation

**Modelling**
Models: Bi-LSTM/GRU, Attention, Capsule etc.

Modelling

**Deployment**
Prediction and model evaluation

Deployment

# Links

https://www.kaggle.com/competitions

Machine learning Coursera famous courses, Andrew Ng, https://www.coursera.org/learn/machine-learning

Machine learning Coursera (on youtube), Andrew Ng, https://www.youtube.com/watch?v=PPLop4L2eGk

The most famous book on deep learning: https://www.deeplearningbook.org/ (Ian Goodfellow, Yoshua Bengio and Aaron Courville)

# ML and DL People

**Andrew Ng**, Founder and CEO of Landing AI, Founder of deeplearning.ai.
**Fei-Fei Li**, Professor of Computer Science at Stanford University.
**Andrej Karpathy**, Senior Director of Artificial Intelligence at Tesla.
**Demis Hassabis**, Founder and CEO of DeepMind.
**Ian Goodfellow**, Director of Machine Learning at Apple.
**Yann LeCun**, Vice President and Chief AI Scientist at Facebook.
**Jeremy P. Howard**, Founding Researcher at fast.ai, Distinguished Research Scientist at the University of San Francisco.
**Ruslan Salakhutdinov**, Associate Professor at Carnegie Mellon University, Director of AI Research at Apple.
**Geoffrey Hinton**, Professor of Computer Science at the University of Toronto, VP and Engineering Fellow at Google
**Rana el Kaliouby**, CEO and Co-Founder of Affectiva.
**Daphne Koller**, Founder and CEO of insitro, Co-Founder of Coursera, Adjunct Professor of Computer Science and Pathology at Stanford.
**Alex Smola**, Director, Amazon Web Services.