# How to start your journey as a Data Scientist?

PARVANEH SHAFIEI

WOMAN IN DATA SCIENCE – TURIN 26TH FEBRUARY

# Who I am

- Past Web & Software developer

- Master: computer science in Polimi

# Where are Rladies now?

**+70 CHAPTERS**

**+20 COUNTRIES**

Ladies
Milan

- **Founded July 2017**
- **Near 250 Rladiers**

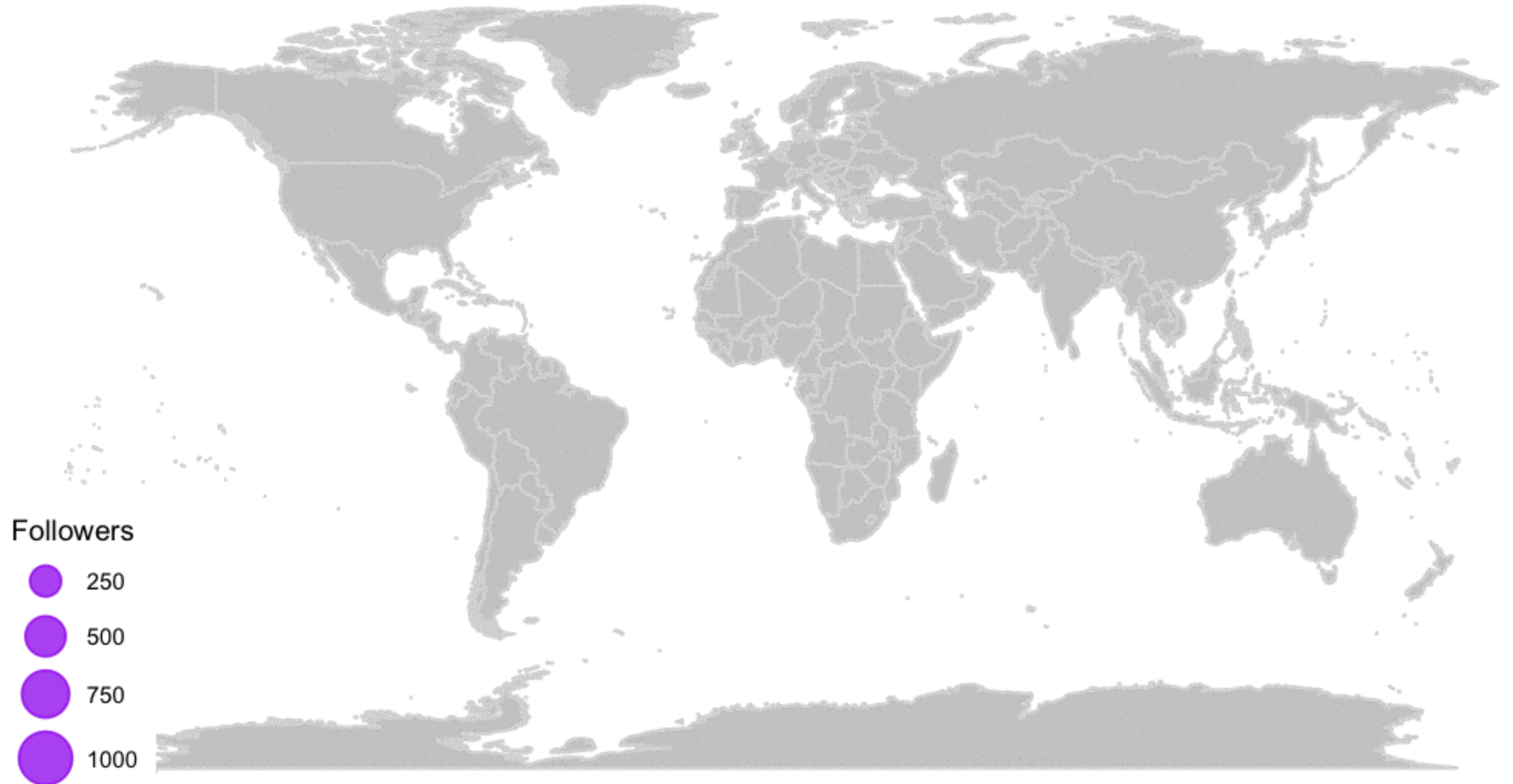**14%**    R package developers are female

**26%**    Only make up data professionals

**20%**    Of all tech startups across the world are founded by women

# 2011-09-01

## Twitter followers by each chapter

**Followers**

- 250
- 500
- 750
- 1000

Made with 🧡 by Daniela Vázquez

# MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21th century, requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

## MATH & STATISTICS

☆ Machine learning
☆ Statistical modeling
☆ Experiment design
☆ Bayesian inference
☆ Supervised learning: decision trees, random forests, logistic regression
☆ Unsupervised learning: clustering, dimensionality reduction
☆ Optimization: gradient descent and variants

## PROGRAMMING & DATABASE

☆ Computer science fundamentals
☆ Scripting language e.g. Python
☆ Statistical computing packages, e.g., R
☆ Databases: SQL and NoSQL
☆ Relational algebra
☆ Parallel databases and parallel query processing
☆ MapReduce concepts
☆ Hadoop and Hive/Pig
☆ Custom reducers
☆ Experience with xaaS like AWS

## DOMAIN KNOWLEDGE & SOFT SKILLS

☆ Passionate about the business
☆ Curious about data
☆ Influence without authority
☆ Hacker mindset
☆ Problem solver
☆ Strategic, proactive, creative, innovative and collaborative

## COMMUNICATION & VISUALIZATION

☆ Able to engage with senior management
☆ Story telling skills
☆ Translate data-driven insights into decisions and actions
☆ Visual art design
☆ R packages like ggplot or lattice
☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

''A data scientist is better at statistic than a software engineer and is better at software engineering than a statistician''

# Who ia Data scientist?

▶ **Passionate** about data

▶ **Thinker**

▶ A **listener** to understand business problem

▶ A good communicator and **storyteller**

▶ Have **technical skills** for coding and analyzing data

▶ **Obsess** with solving problems

## Data scientists

- Use analytics & technical skills to extract, analyze and model data

## Statistician

- Understand statistic and apply it on real problems

## Data engineer

- Responsible for architecture of data
- Ensure the flow of data within servers & applications

Be ready for the journey!

# Learning in data science is not a linear path!

▶ Data science is **not about just one specific skills**

▶ You must know all things and **be expert in one** of them

▶ It is a **fast evolving field**

▶ It does **not need** to be **expert in the domain**

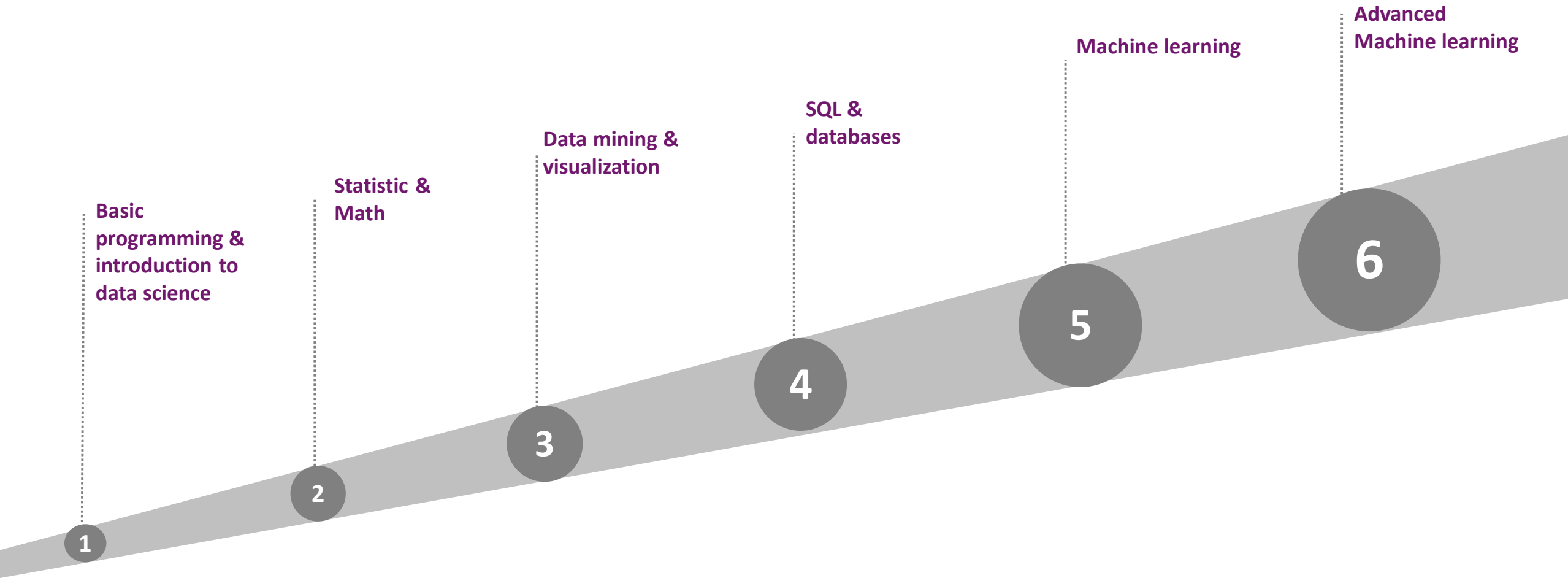# Where to start learning or add more skill sets to your toolbox?

## Online courses

**DataCamp**

**coursera**

**edX**

**UDACITY**

**MITOPENCOURSEWARE**
MASSACHUSETTS INSTITUTE OF TECHNOLOGY

**fast.ai**

## Datasets for practices

**kaggle**

**UCI** Machine Learning Repository

**MLCOMP**

**UCI KDD Archive**

**Google Trends Datastore**

**CrowdFlower**

**DATA.GOV.UK** Beta
Opening up Government

# There is a starting point but no end road!

**Advanced Machine learning**

**Machine learning**

**SQL & databases**

**Data mining & visualization**

**Statistic & Math**

**Basic programming & introduction to data science**

1

2

3

4

5

6

# Beginners steps

**Basic programming & introduction to data science**

| | |
|---|---|
| • R<br>• Python | • Benefits & potential of the languages<br>• Answer possible questions in the field<br>• How to code and use the basics such as libraries, functions,..<br>• How to load & read & manipulate data |

**DataCamp**

**Introduction to R**

**Intro to Python for Data Science**

**coursera**

R Programming

Introduction to Data Science in Python

**edX**

**Introduction to Python: Absolute Beginner**

**Introduction to R for Data Science**

# Beginners steps

**Statistic & Math**

- Learn fundamental concepts of statistics such as p-values, variance, correlation, statistical hypothesis..
- Evaluate various types of data and how to interpret their structure
- How to apply various statistical methods on the data

**DataCamp**

**A Hands-on Introduction to Statistics with R**

**Intro to Statistics with R: Introduction**

**coursera**

Statistical Inference

Introduction to Probability and Data

**edX**

**Statistics and R**

**Introductory Statistics : Basic Ideas and Instruments for Statistical Inference**

# Beginners steps

**Data mining & visualization**

- Handle anomalies in data such as missing values, outliers,..
- Explore correlations among variables
- Apply feature engineering on the data
- Create graph & visualization to demonstrate findings in the data

**DataCamp**

**Data Visualization with ggplot2**

**Introduction to Data Visualization with Python**

**coursera**

Exploratory Data Analysis

Data Management and Visualization

**edX**

**Analyzing and Visualizing Data with Power BI**

**Data Analysis: Visualization and Dashboard Design**

# Beginners steps

## SQL & databases

- Understand how relational databases are working
- How to interact with databases for fetching, saving and manipulation of data

**DataCamp**

**coursera**

**edX**

**Intro to SQL for Data Science**

Using Databases with Python

**Querying Data with Transact-SQL**

**Introduction to Databases in Python**

Managing Big Data with MySQL

# Beginners steps

**Machine learning**

- Learn various algorithms such as regression, random forest, classification tree, etc. and their concepts and understand where to use them
- How to apply predictive modeling on set of data
- Learn how to train various model and what are the metrics of trained models and how to compare them

**DataCamp**

**Machine Learning Toolbox**

**Introduction to Machine Learning**

**coursera**

Machine Learning

Practical Machine Learning

**edX**

**Applied Machine Learning**

**Principles of Machine Learning**

# Intermediate & advanced steps

## Advanced machine learning

- Learn how to manipulate unstructured data
- Learn and understand advanced topics such as deep learning, social network analysis, text mining, time series processing,..

**DataCamp**

**Text Mining: Bag of Words**

**Manipulating Time Series Data in R with xts & zoo**

**coursera**

Deep Learning Specialization

Practical Time Series Analysis

**edX**

**Deep Learning Explained**

**Graph Algorithms**

# Being specialist in one field?

**Other topics**

- Depends on the field, business and type of problems

**DataCamp**

**coursera**

**Building Web Applications in R with Shiny**

**Network Analysis in R**

**Building Chatbots in Python**

**Credit Risk Modeling in R**

Bayesian Methods for Machine Learning

Practical Reinforcement Learning

# Kaggle learning

## Hands-On Data Science Education

Learn the basics to confidently start a new career or upgrade your skills.



Machine Learning

R

Data Visualisation

Deep Learning

SQL

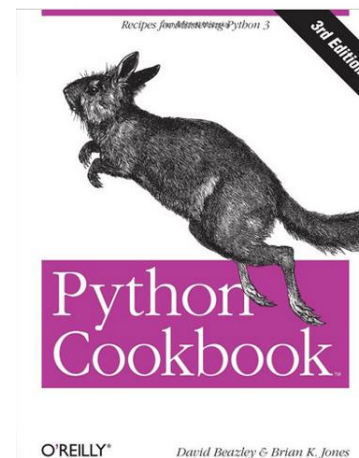If you have to read just one single book / or watch just 15 hours videos
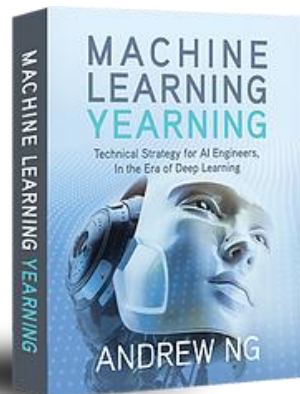
# Free books

# Websites
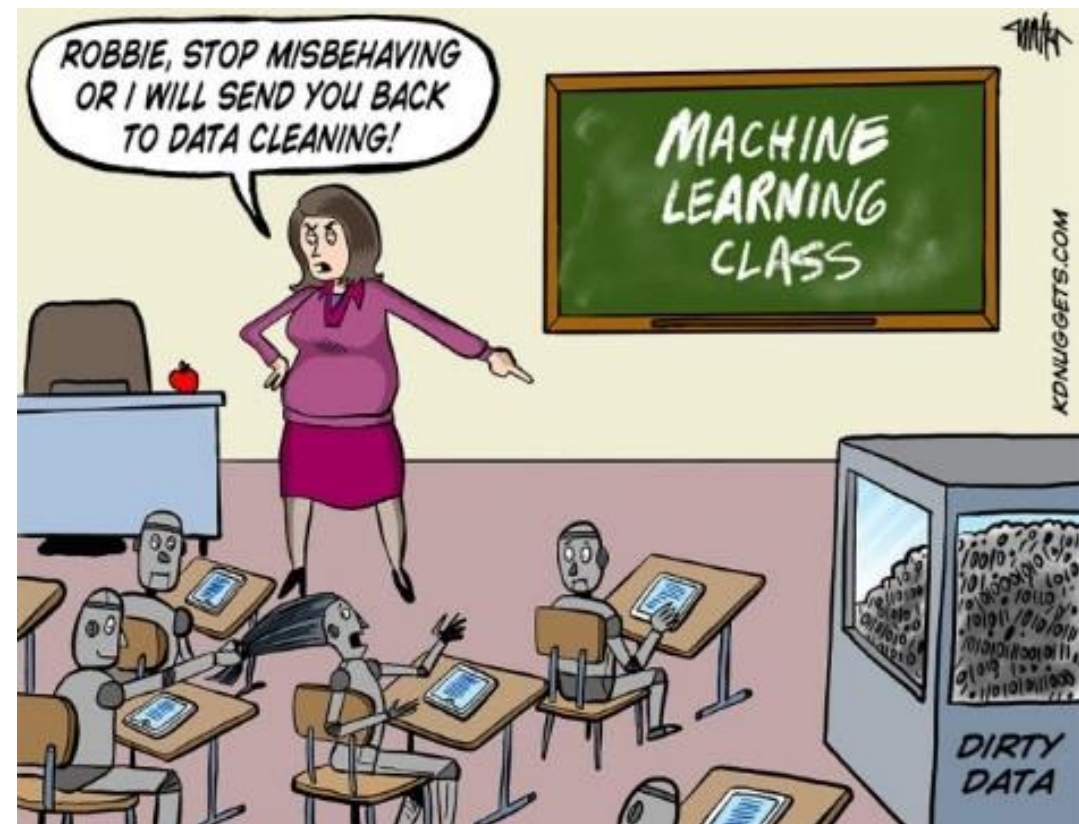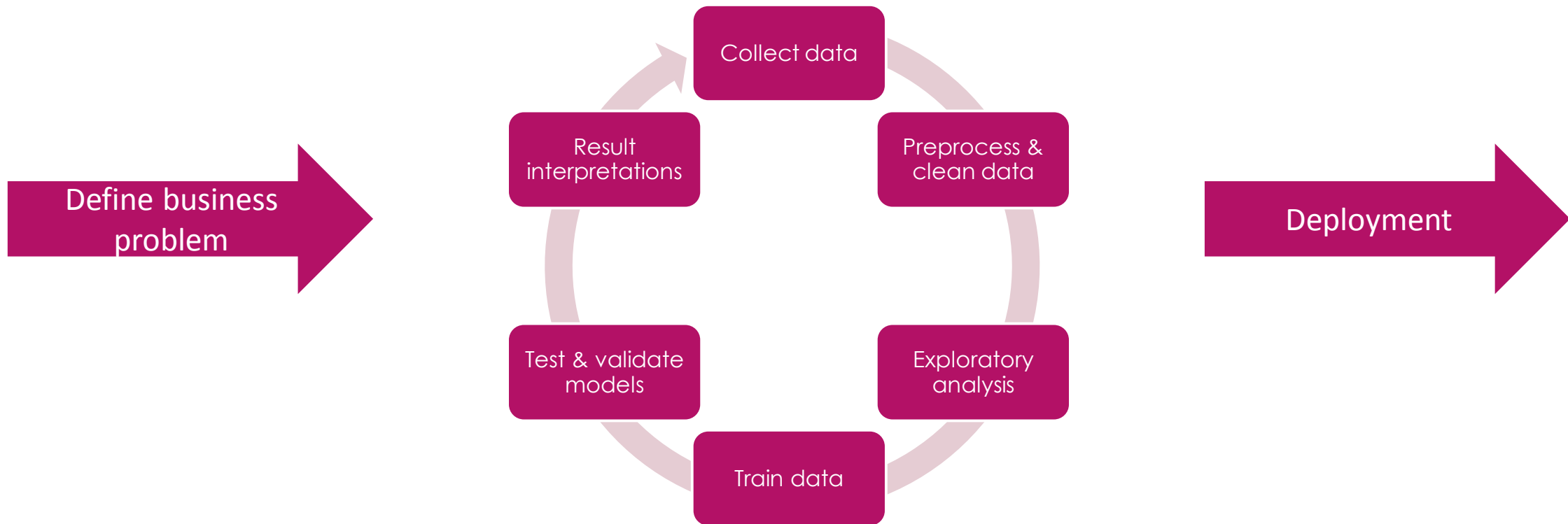
# You don't have any experience?

▶ Do the action!

▶ Apply what you learn on some use cases from Kaggle or any other data set that you can

▶ Do the competitions in Kaggle, share the code in blogs, GitHub account, LinkedIn,...

▶ You can find internship, asking companies directly
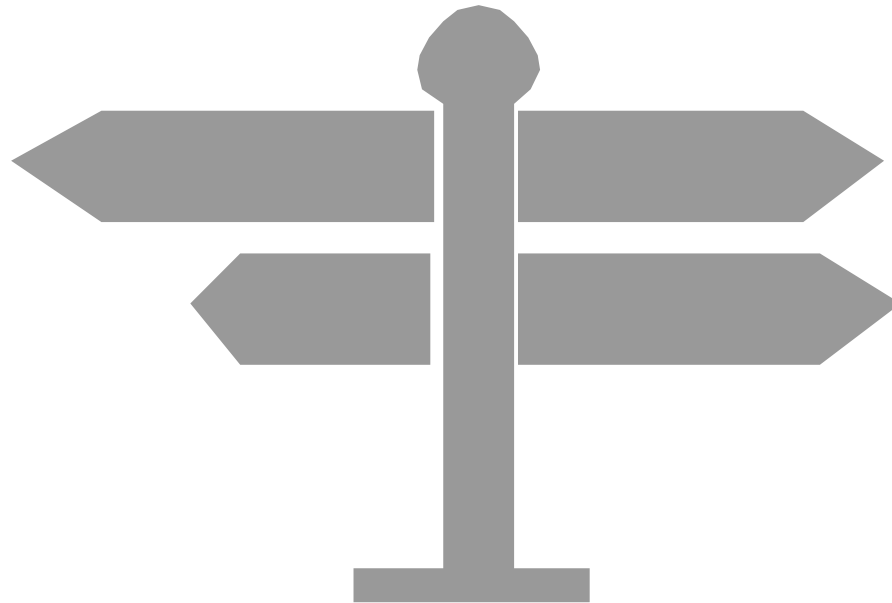
▶ Ask mentors to follow you in the way

# Ok, I'm all set, but how can I tackle a real problem?

Define business problem →

- Collect data
- Preprocess & clean data
- Exploratory analysis
- Train data
- Test & validate models
- Result interpretations

→ Deployment

# Define the business problem

- Listen to business's need
- You must be a good listener

- The business you are facing do not know their problem yet? That is even fine too!
- A lot of time, they don't have any idea where to start

# Collect data

Data are in a structured format such as database, csv file,...

**Structured data**

**+**

Data are in a unstructured format such as text, images, audio and video

**Unstructured data**

**Near 80% of data in organization are unstructured**

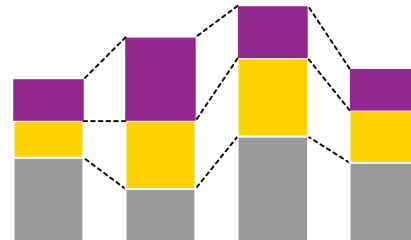# Preprocessing & cleaning data

▶ **Never trust** your data

▶ **Check consistency & structure** of the data to be sure about its quality

▶ Fill **missing values**

▶ **Be aware of noisy data**  and treat them well!

▶ Check **outliers** & remove them

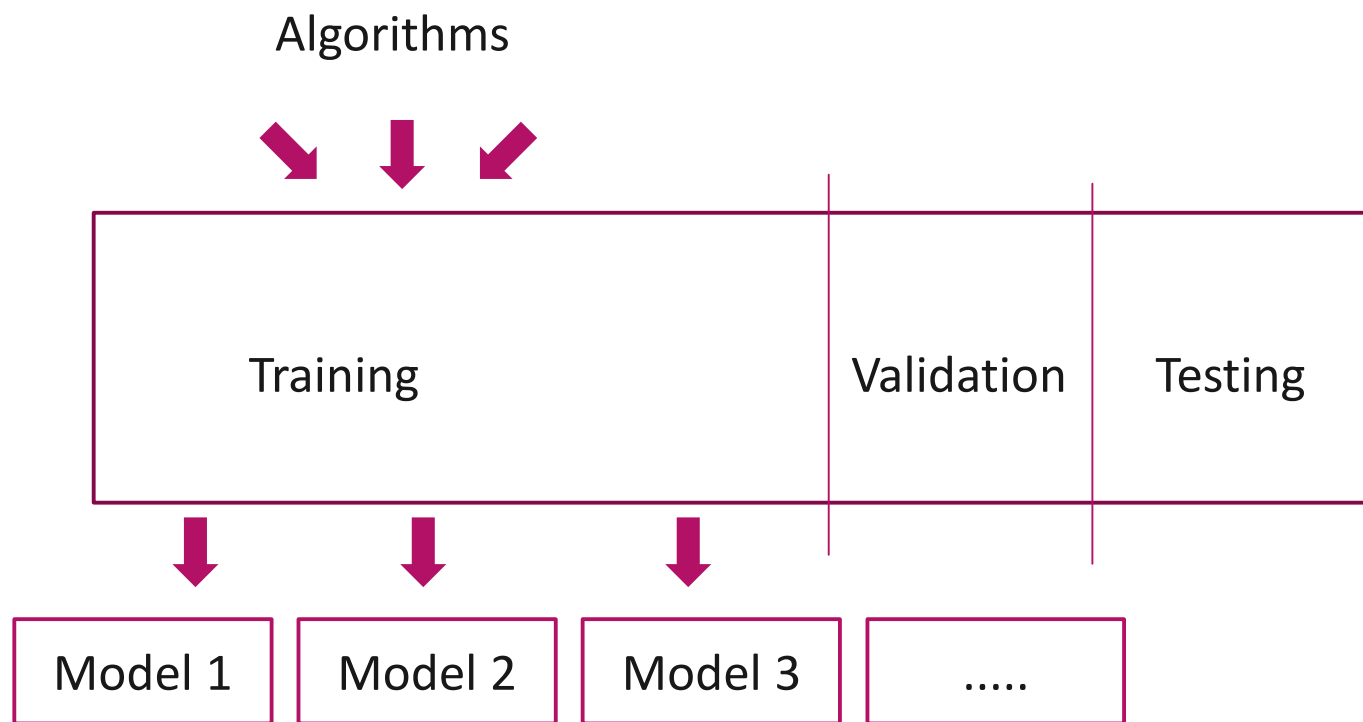▶ Apply **data transformation** such as normalization or aggregation

# Exploratory analysis

▶ **Understand and summarize the** data

▶ **Explore** for insights & hidden patterns

▶ Investigate the correlations among various variables

▶ Create some **hypothesis** based on the findings

▶ **Understand which model & technique to** use for analyzing the data

▶ **Feature engineering** can happen in this stage too

**Think about the questions that you are going to answer**

# Model training, testing & validation

Algorithms

Training | Validation | Testing

Model 1 | Model 2 | Model 3 | .....

**60%, 20%, 20%**

# Algorithms?

**Classification / supervised learning**

- Linear regression
- Logistic regression
- Decision tree
- Random forest
- PCA

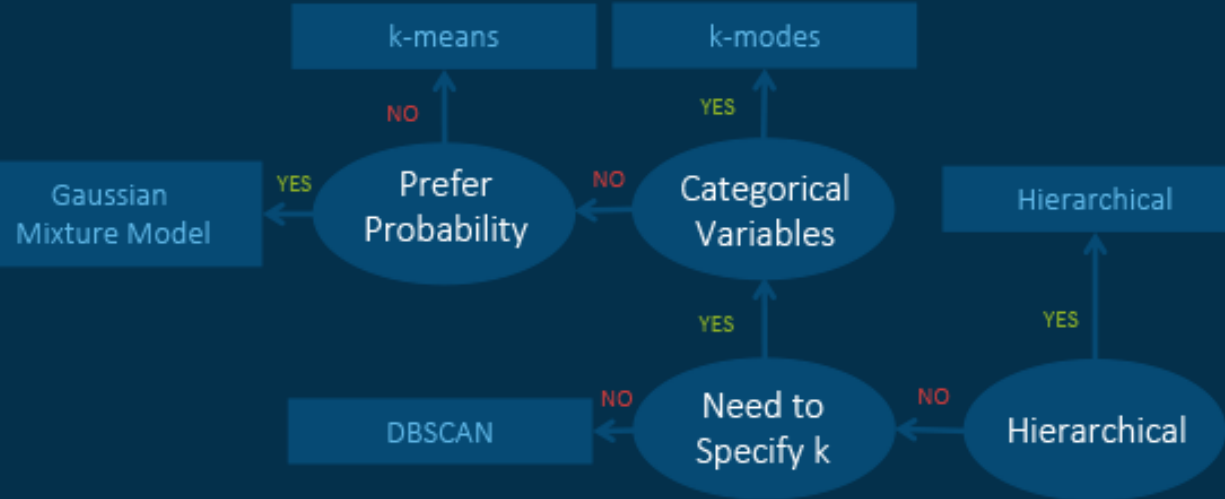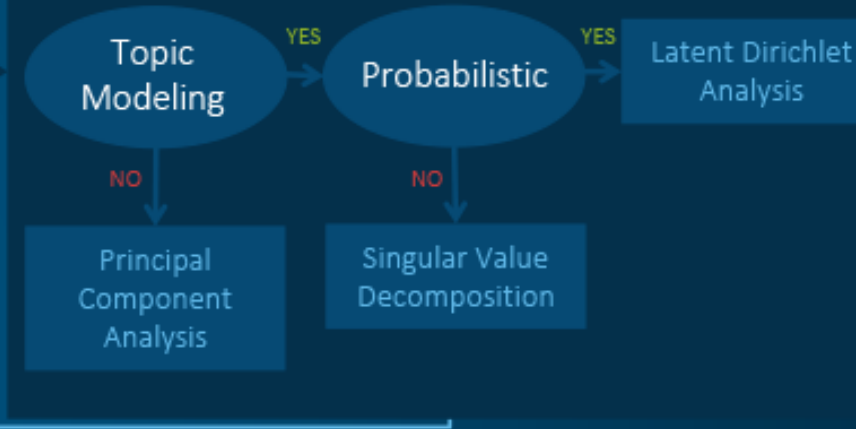**Clustering/ unsupervised learning**

- K-means
- PAM
- Hierarchy
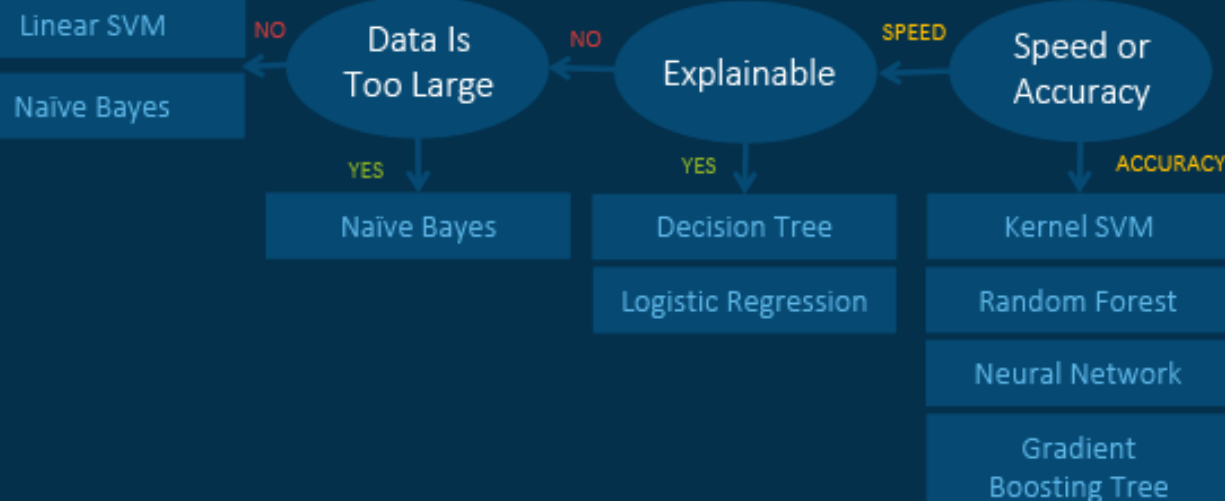- Density based

# Machine Learning Algorithms Cheat Sheet

## Unsupervised Learning: Clustering

k-means

k-modes

NO → Prefer Probability

YES → Categorical Variables

YES → Gaussian Mixture Model

NO

Hierarchical

YES → Need to Specify k

NO → DBSCAN

YES → Hierarchical

NO

## START

Dimension Reduction

YES

NO → Have Reponses

## Unsupervised Learning: Dimension Reduction

Topic Modeling

YES → Probabilistic

YES → Latent Dirichlet Analysis

NO → Principal Component Analysis

NO → Singular Value Decomposition

## Supervised Learning: Classification

Linear SVM

Naïve Bayes

NO → Data Is Too Large

NO → Explainable

SPEED → Speed or Accuracy

YES → Naïve Bayes

YES → Decision Tree

Logistic Regression

ACCURACY → Kernel SVM

Random Forest

Neural Network

Gradient Boosting Tree

NO → Predicting Numeric

YES

## Supervised Learning: Regression

Speed or Accuracy

SPEED → Decision Tree

Linear Regression

ACCURACY → Random Forest

Neural Network

Gradient Boosting Tree

# Deployment

"More data beats better models. Better data beats more data." — Riley Newman