



M1T2: Where Are We?

From a task perspective...

Data Science Deconstructed



We are here

SKILLS REQUIRED



FRAME THE PROBLEM

- Domain Knowledge (needs)
- Product Intuition (metrics)
- Business Strategy (priorities)
- Teamwork (people & resources)



COLLECT RAW DATA

- Database Management
 - Systems: MySQL, PostgreSQL, Oracle, MongoDB
- Querying Structured Databases
 - SQL
- Retrieving Unstructured Info
 - Informational Retrieval / Text Mining
- Distributed Storage
 - Hadoop HDFS, Spark, Flink



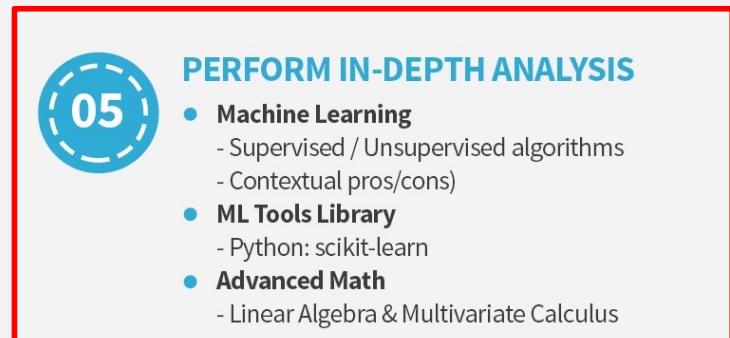
PROCESS THE DATA

- Scripting Language
 - Python or R
- Data Wrangling & Cleaning
 - Python "Pandas" library
- Distributed Processing
 - Hadoop MapReduce / Spark



EXPLORE THE DATA

- Scientific Computing
 - Python: numpy, matplotlib, scipy, pandas
- Inferential Statistics
 - hypothesis testing
 - correlation vs. causation
- Experimental Design
 - A/B tests, controlled trials



PERFORM IN-DEPTH ANALYSIS

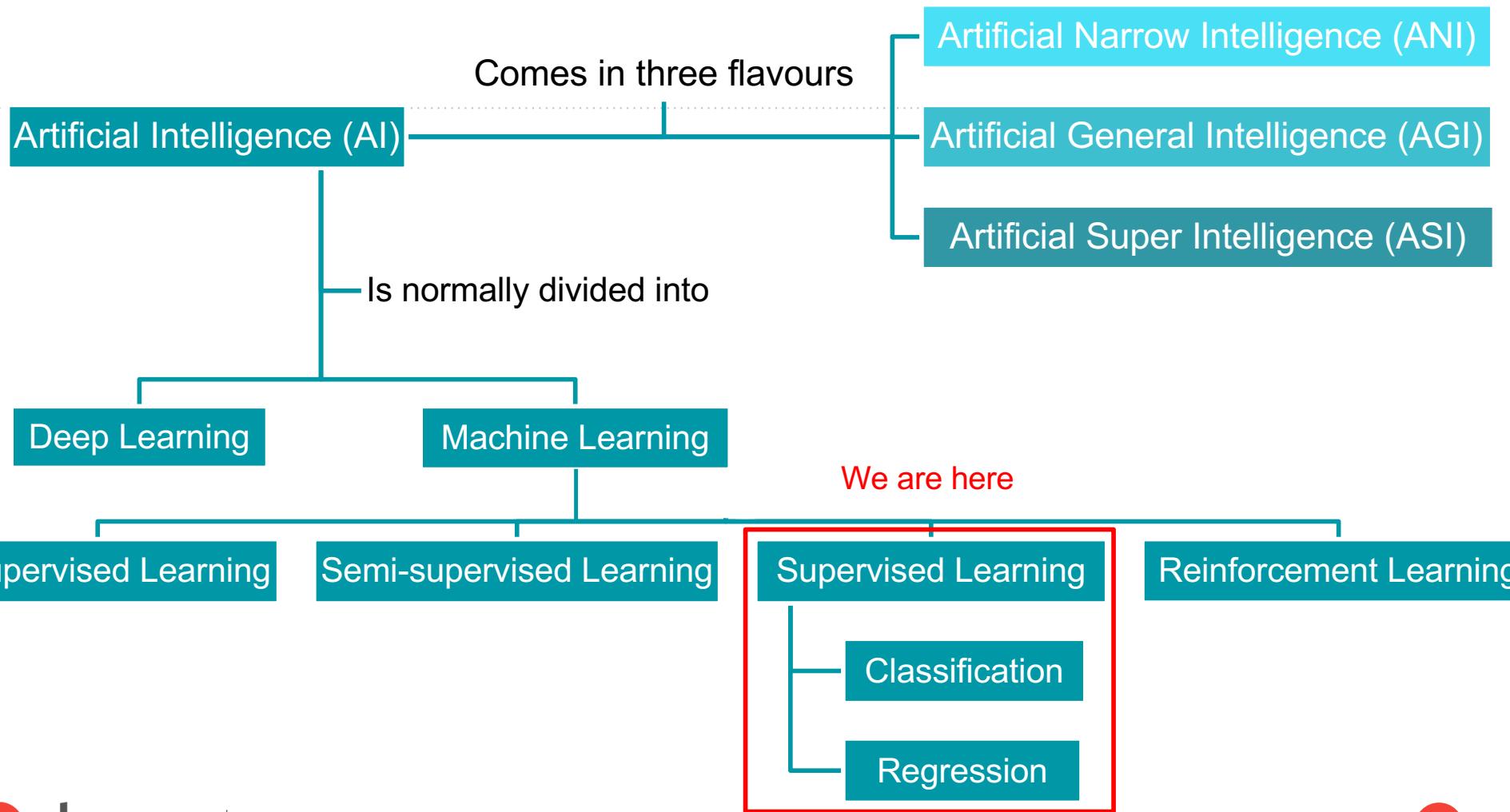
- Machine Learning
 - Supervised / Unsupervised algorithms
 - Contextual pros/cons
- ML Tools Library
 - Python: scikit-learn
- Advanced Math
 - Linear Algebra & Multivariate Calculus



COMMUNICATE RESULTS

- Business Acumen
 - Non-technical terminology
- Data Visualization Tool(s)
 - Tableau, D3.js, Google visualize, matplotlib, ggplot, seaborn
- Data Storytelling
 - presenting & speaking
 - reporting & writing

... but where are we exactly, within the field of Artificial Intelligence?



Supervised Learning

Supervised learning (SL) is the **machine learning** task of **learning** a function that maps an input to an output, based on example labelled input-output pairs.

Classification VS Regression

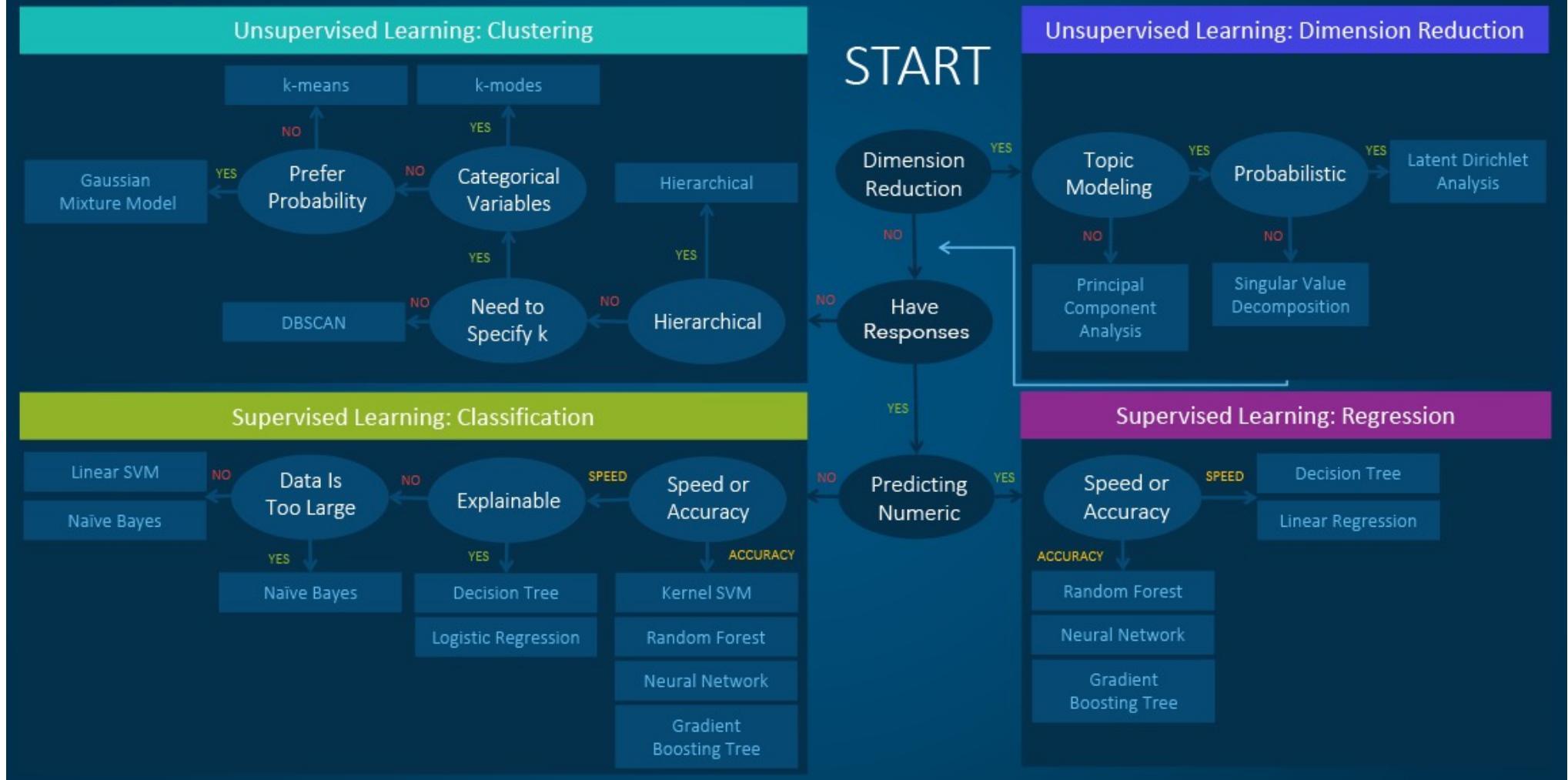
They are two major types of prediction problems Machine Learning tries to solve. Fundamentally, they differ because:

- **classification** is about predicting a label (true/false, low/medium/high, etc)
- **regression** is about predicting a quantity (age, price, height, etc)

Main Algorithms

REGRESSION	CLASSIFICATION
Decision Tree	Decision Tree
Linear Regression	Logistic Regression
Random Forest	Random Forest
Neural Network	Neural Network
Gradient Boosting Tree	Gradient Boosting Tree
Support Vector Machine	Support Vector Machine
	Kernel SVM
	Naïve Bayes
	K-Nearest Neighbours

Machine Learning Algorithms Cheat Sheet



Tips:

- As you are being bombarded with a multitude of new exciting concepts and ideas, it is very easy to get lost. Therefore, try to stay focused and explore only the concepts/tools you need for task 2.
- Run your first model ASAP and write down the first set of metrics. Once that's done, you start reiterating the process, over and over, until you get the best possible results.

Resources:

- <https://scikit-learn.org/stable/>
- <https://machinelearningmastery.com/difference-between-algorithm-and-model-in-machine-learning/>
- https://www.sas.com/en_gb/insights/articles/analytics/machine-learning-algorithms.html
- <https://blogs.sas.com/content/subconsciousmusings/2020/12/09/machine-learning-algorithm-use/>



Thanks.