

Aprendizado de Máquina e Reconhecimento de Padrões 2021.2



The Machine Learning Landscape

Prof. Samuel Martins (Samuka)
samuel.martins@ifsp.edu.br



What is Machine Learning?

What is Machine Learning?



“The field of study that gives computers the ability to learn without being explicitly programmed.”

Arthur Samuel

What is Machine Learning?



Arthur Samuel

“The field of study that gives computers the ability to learn without being explicitly programmed.”



Tom Mitchell

*“A computer program is said to **learn** from **experience E** with respect to some **class of tasks T** and **performance measure P**, if its performance at tasks in **T**, as measured by **P**, improves with **experience E**.”*

What is Machine Learning?



Arthur Samuel

“The field of study that gives computers the ability to learn without being explicitly programmed.”



Tom Mitchell

Ex: Spam detector.

T = Predict whether an email is spam or not.

E = Collection of emails labeled as spam or not.

P = Spam classifier accuracy.

What is Machine Learning?



Arthur Samuel

“The field of study that gives computers the ability to learn without being explicitly programmed.”



(Train)
Data

“A computer program is said to *learn* from **experience E** with respect to some **class of tasks T** and **performance measure P**, if its performance at tasks in **T**, as measured by **P**, improves with **experience E**.”



Tom Mitchell

Ex: Spam detector.

T = Predict whether an email is spam or not.

E = Collection of emails labeled as spam or not.

P = Spam classifier accuracy.



facebook

IBM

NETFLIX

UBER

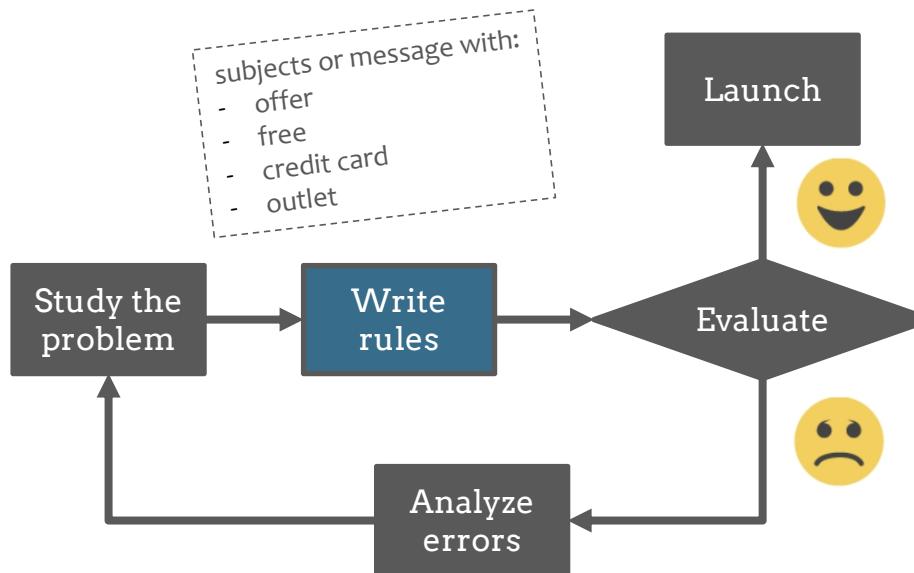
Google

amazon



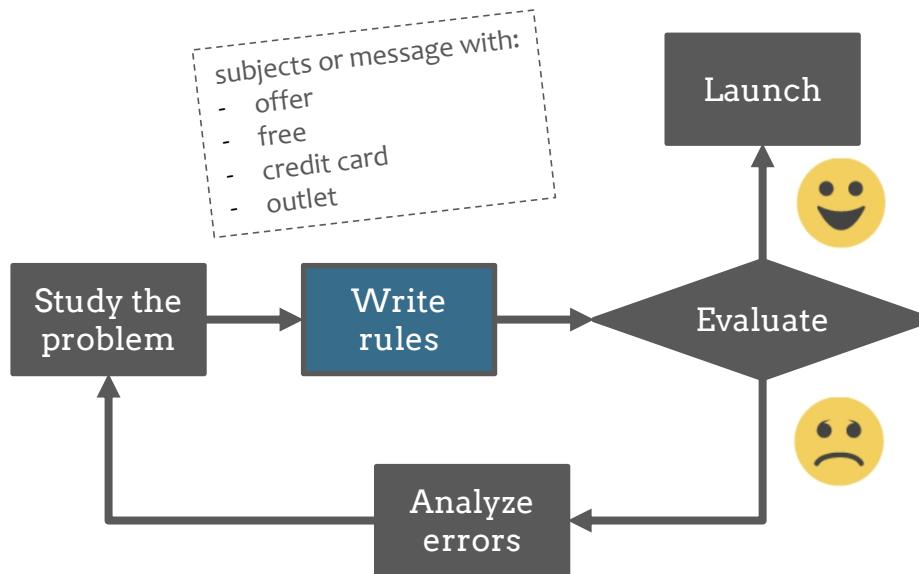
Why Use Machine Learning?

Why Use Machine Learning?

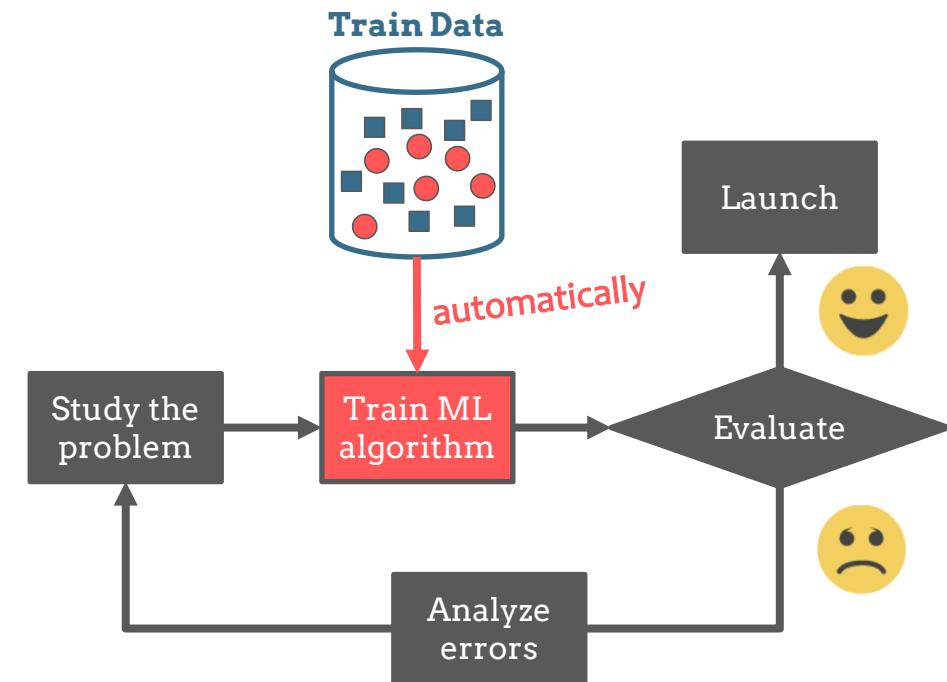


The traditional approach

Why Use Machine Learning?



The traditional approach



The machine learning approach

Why Use Machine Learning?

The traditional approach

Study the problem

sub

To Joan Perez
From L Spencer <lspencer293@gmail.com>

Response Required: Project ABC

Hi Joan,

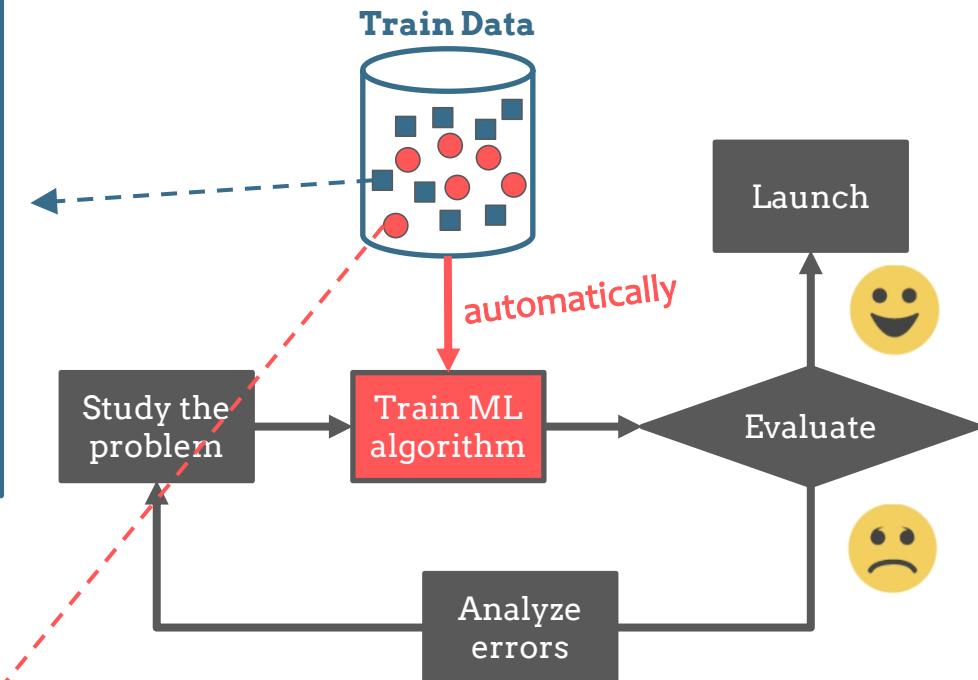
Your risk analysis for Project ABC was very helpful. Thanks for being so thorough.

However, you are behind on Phase 1. It was due yesterday, and we need it as soon as possible so we can move on to the next project phase.

Please let me know the revised completion date for Phase 1 by the end of the day. If you are having trouble with this phase, I would be happy to help answer any questions you may have. You can reach me at: 010-555-0100

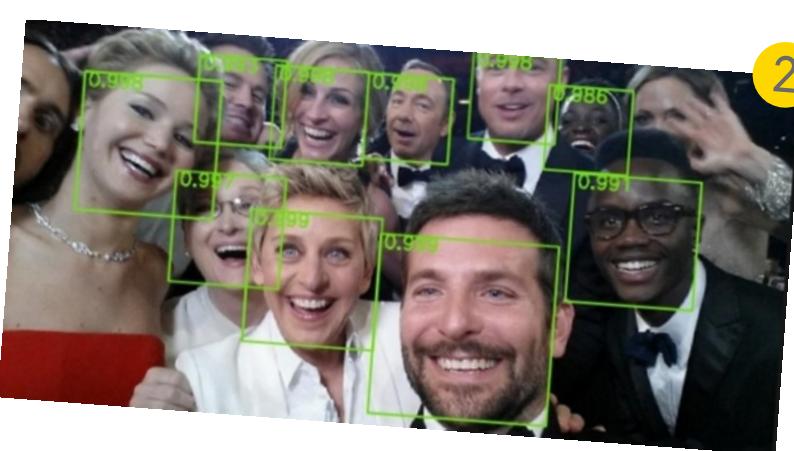
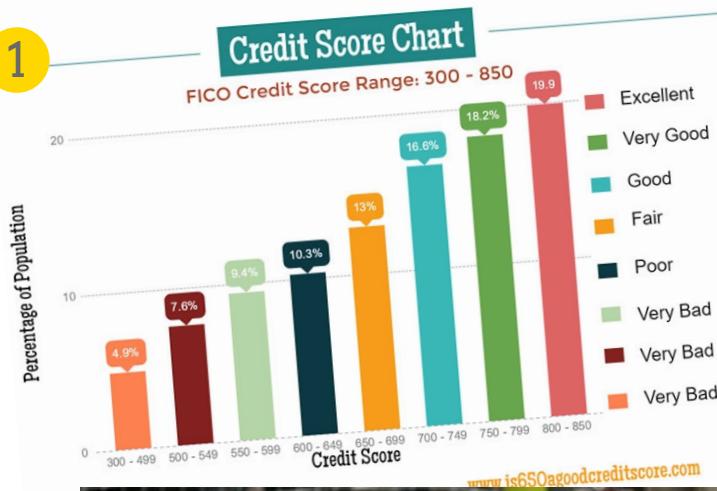
eai Samuel, tudo?
venho trazer UMA SURPRESA, ACONTECIMENTO RARÍSSIMO, MARAVILHOSO, que é o [OUTLET da dobra](#) online.

OUTLET SURPRESA DA DOBRA



The machine learning approach

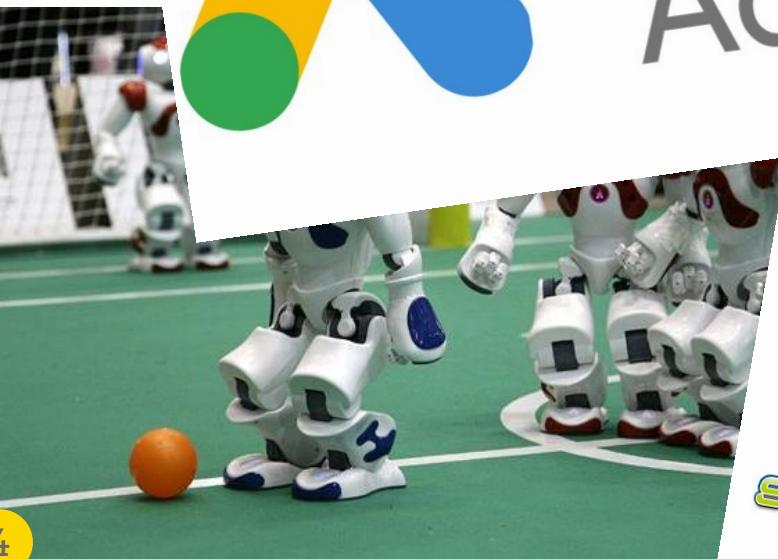
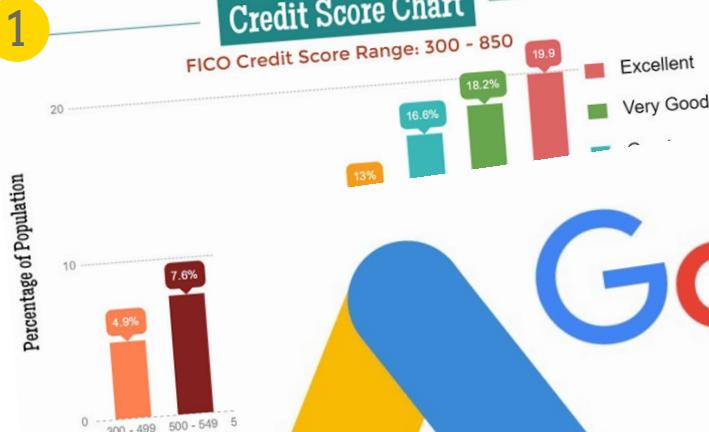
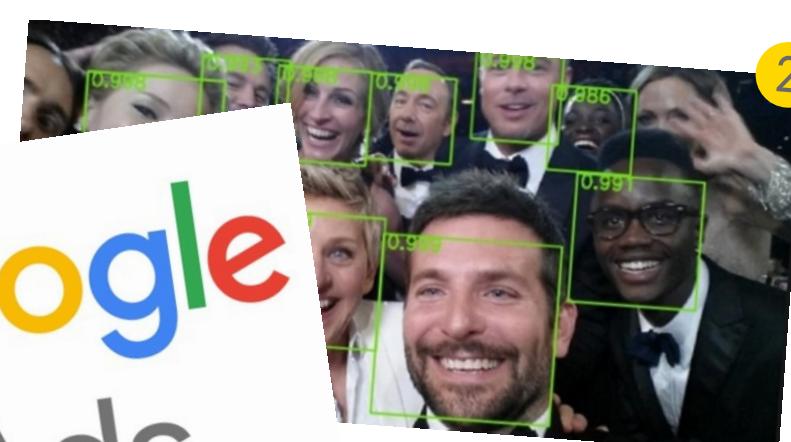
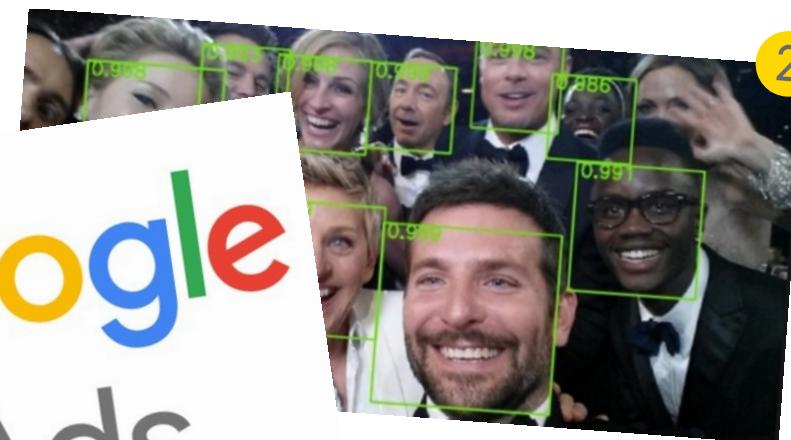
1



3



1

4
KINECT
SPORTS

3

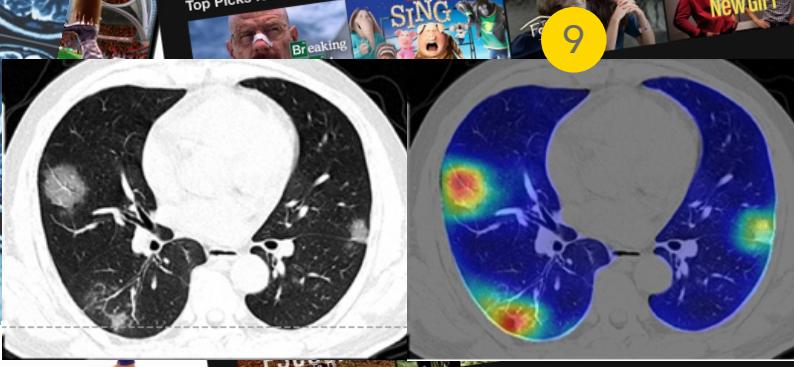


5



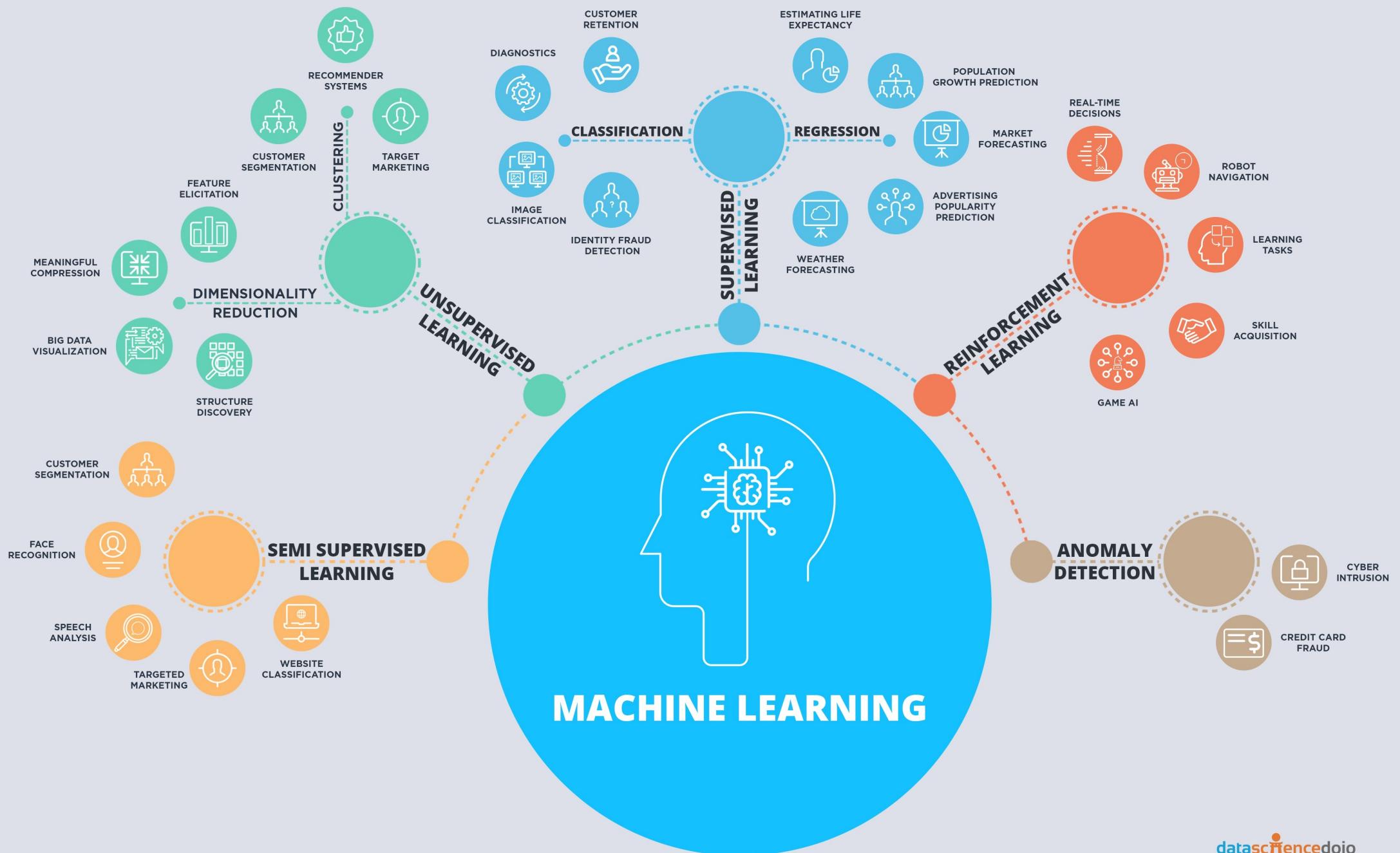
7

6



<https://internetofbusiness.com/nvidia-ai-synthetic-mri-diagnostics/>

Types of Machine Learning Systems



Supervised Learning

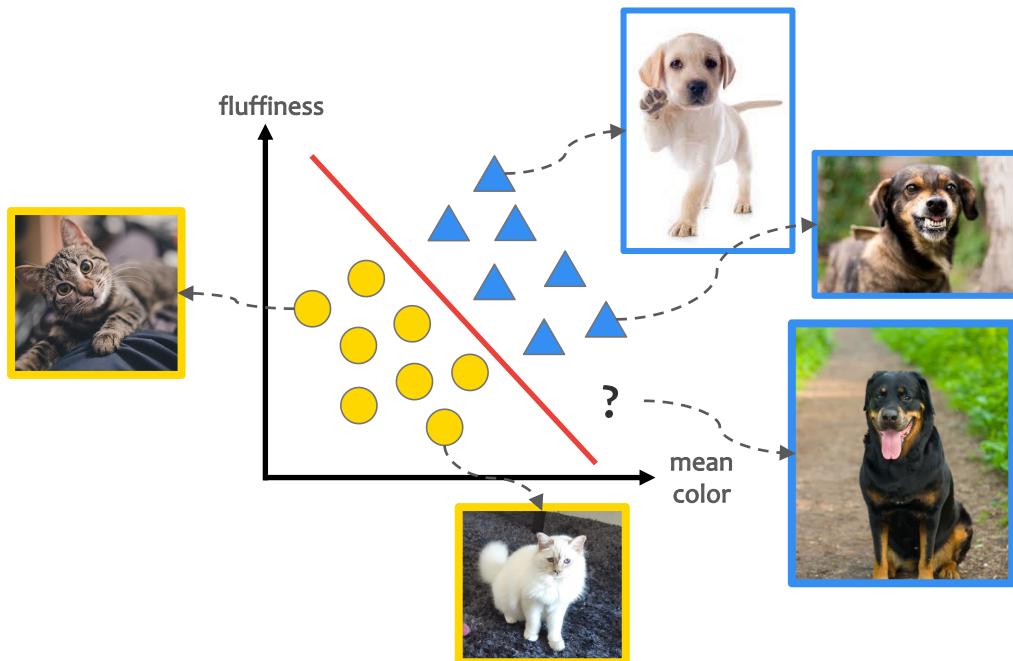
Supervised Learning

Classification

The training set includes **annotations / labels / classes**

(usually) annotated by humans:
time-consuming and costly

Prediction of **categorical variables (labels)**



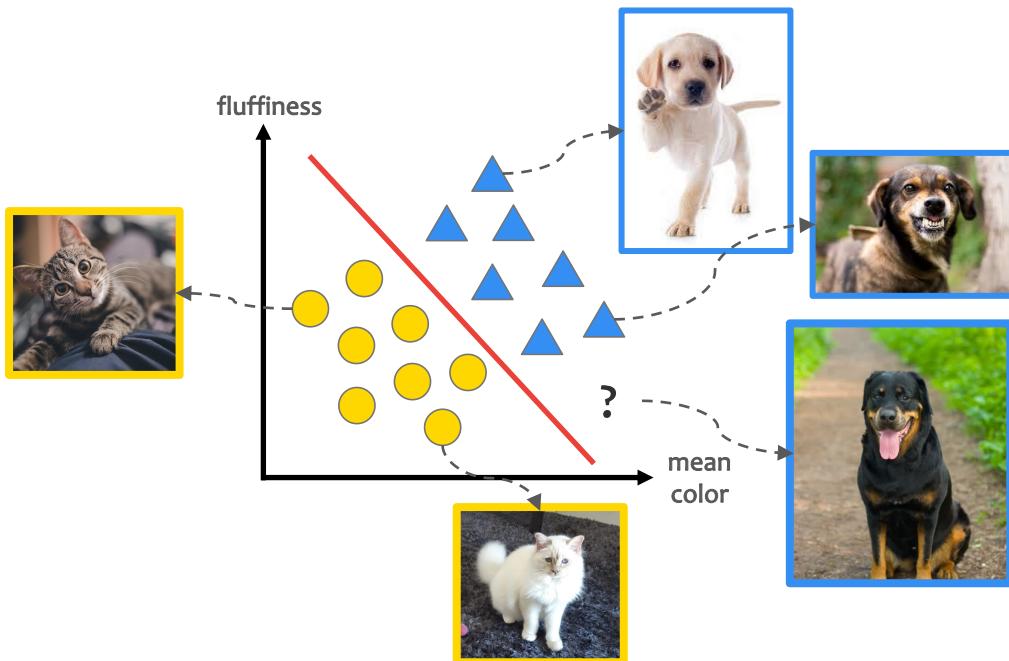
- k-Nearest Neighbors (KNN)
- Logistic Regression
- Support Vector Machines (SVM)
- Decision Trees and Random Forest
- Neural Networks
- ...

Supervised Learning

Classification

The training set includes **annotations / labels / classes**

Prediction of **categorical variables (labels)**

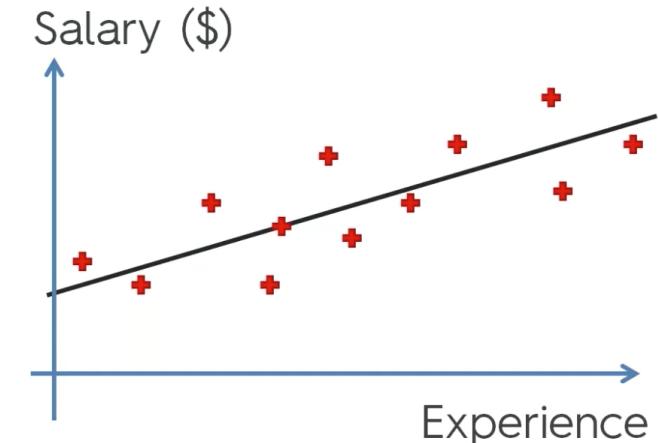


- k-Nearest Neighbors (KNN)
- Logistic Regression
- Support Vector Machines (SVM)
- Decision Trees and Random Forest
- Neural Networks
- ...

Regression

Prediction of **numeric variables.**

Experience	Salary
60	10,000.00
75	12,000.00
100	15,000.00
...	...



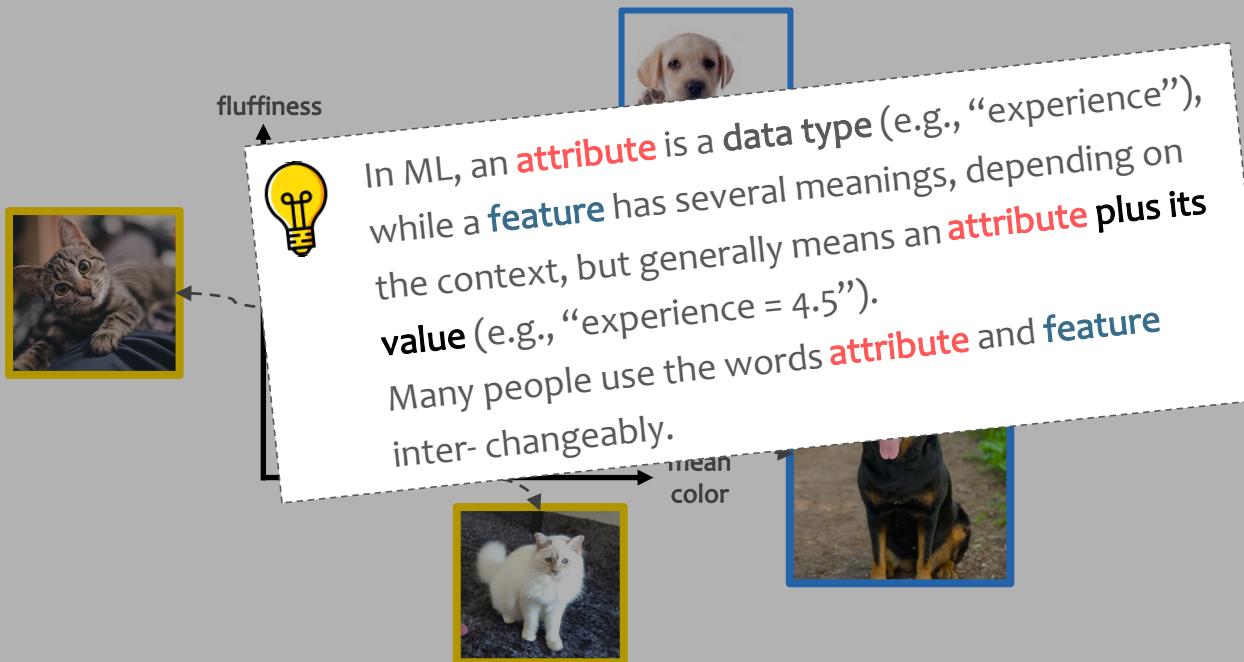
- Linear Regression
- Polynomial Regression
- KNN Regression
- SVM Regression
- Decision Tree Regression
- ...

Supervised Learning

Classification

The training set includes **annotations / labels / classes**

Prediction of **categorical variables (labels)**

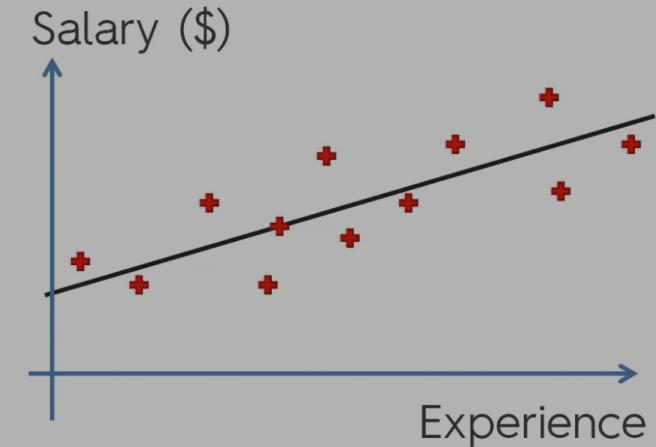


- k-Nearest Neighbors (KNN)
- Logistic Regression
- Support Vector Machines (SVM)
- Decision Trees and Random Forest
- Neural Networks
- ...

Regression

Prediction of **numeric variables**.

Experience	Salary
60	10,000.00
75	12,000.00
100	15,000.00
...	...



- Linear Regression
- Polynomial Regression
- KNN Regression
- SVM Regression
- Decision Tree Regression
- ...

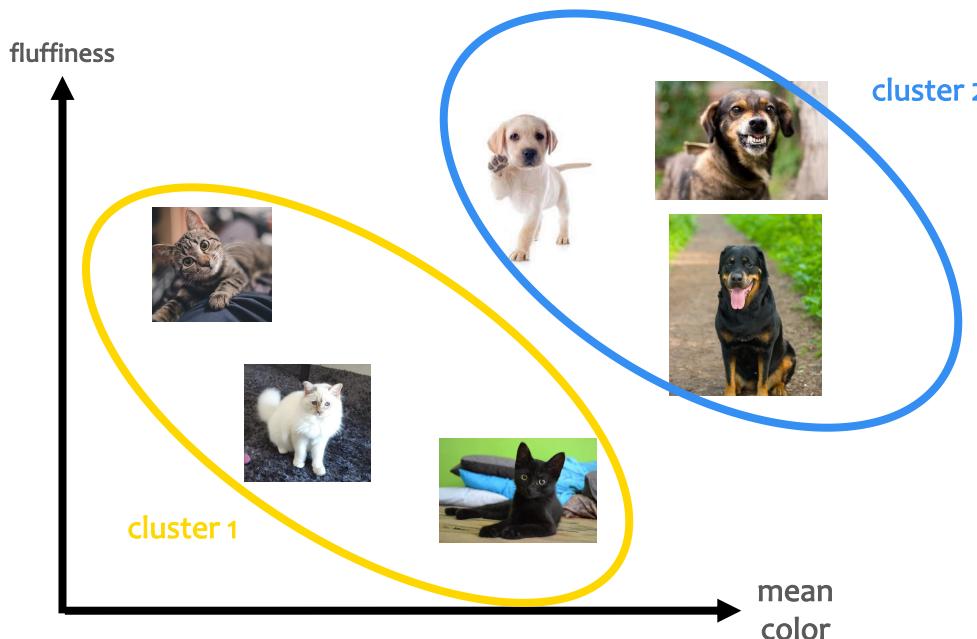
Unsupervised Learning

no supervision → no **labels / annotations**

Unsupervised Learning

no supervision → no labels / annotations

Clustering

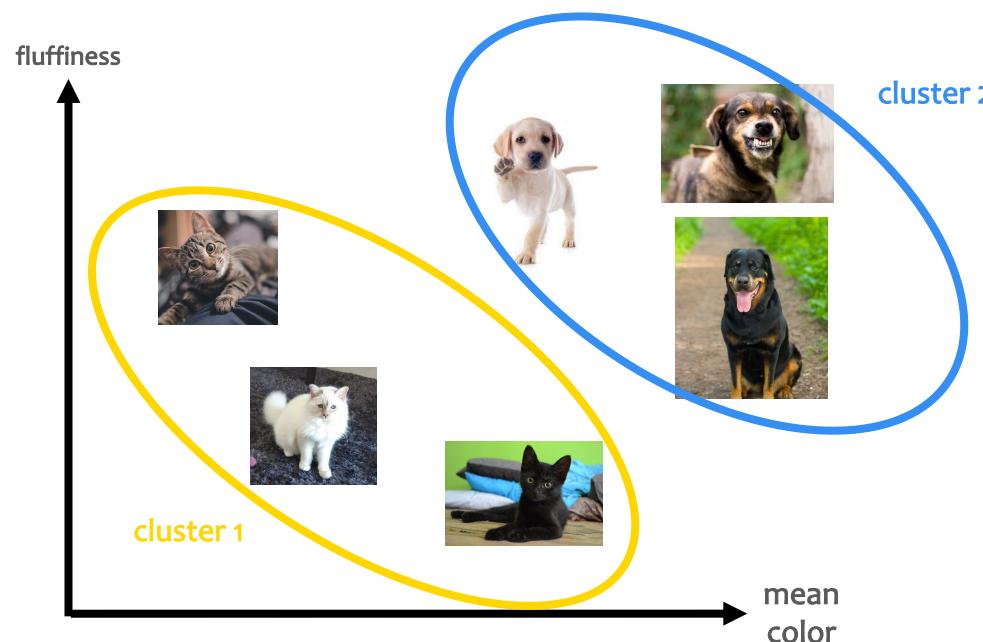


- k-Means
- DBSCAN
- Mean-shift
- OPF Clustering
- Hierarchical Clustering
- ...

Unsupervised Learning

no supervision → no labels / annotations

Clustering

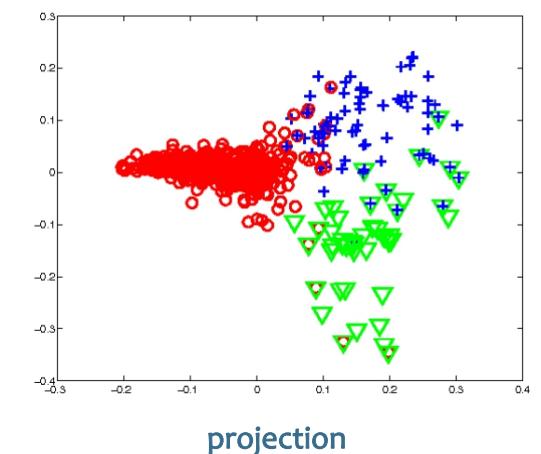
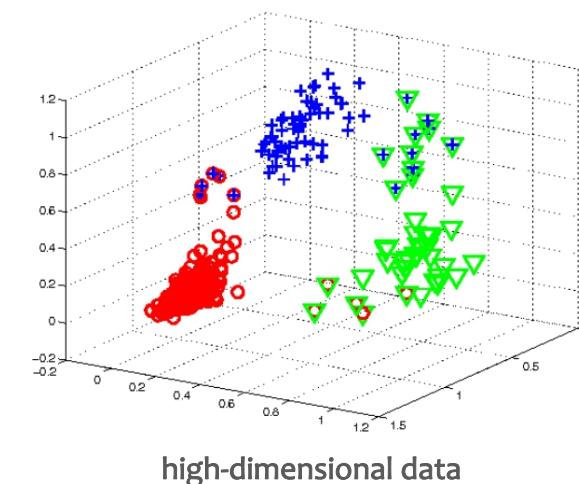


- k-Means
- DBSCAN
- Mean-shift
- OPF Clustering
- Hierarchical Clustering
- ...

Dimensionality Reduction / Projection

Project a high-dimensional data on to a lower-dimensional feature space without losing much information.

Projection on to a 2D or 3D space is used to visualize and understand the data.

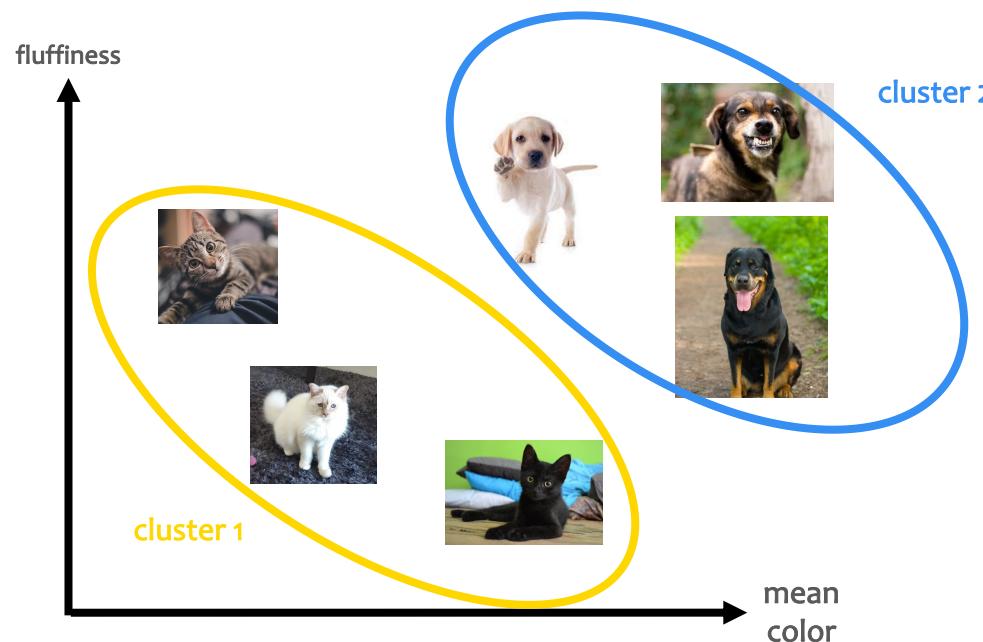


- Principal Component Analysis (PCA)
- Locally Linear Embedding (LLE)
- t-SNE
- UMAP
- ...

Unsupervised Learning

no supervision → no labels / annotations

Clustering

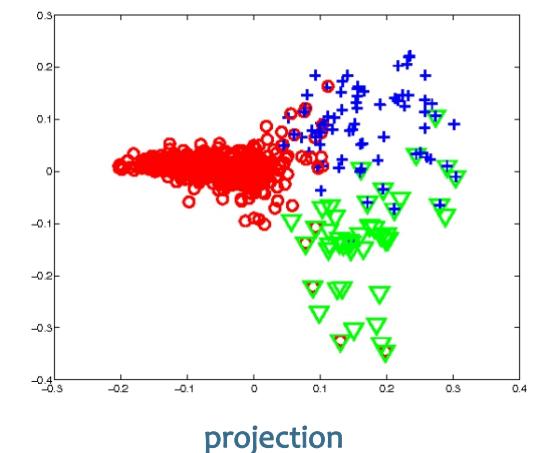
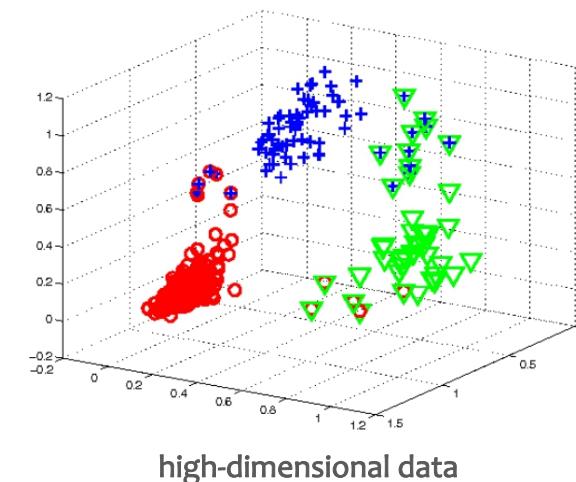


- k-Means
- DBSCAN
- Mean-shift
- OPF Clustering
- Hierarchical Clustering
- ...

Dimensionality Reduction / Projection

Project a high-dimensional data on to a lower-dimensional feature space without losing much information.

Projection on to a 2D or 3D space is used to visualize and understand the data.



There are some algorithms that use **label information** during projection.
E.g.: Linear Discriminant Analysis (LDA)

Principal Component Analysis (PCA)
t-SNE (tSNE)
Locally Linear Embedding (LLE)

• UMAP
• ...

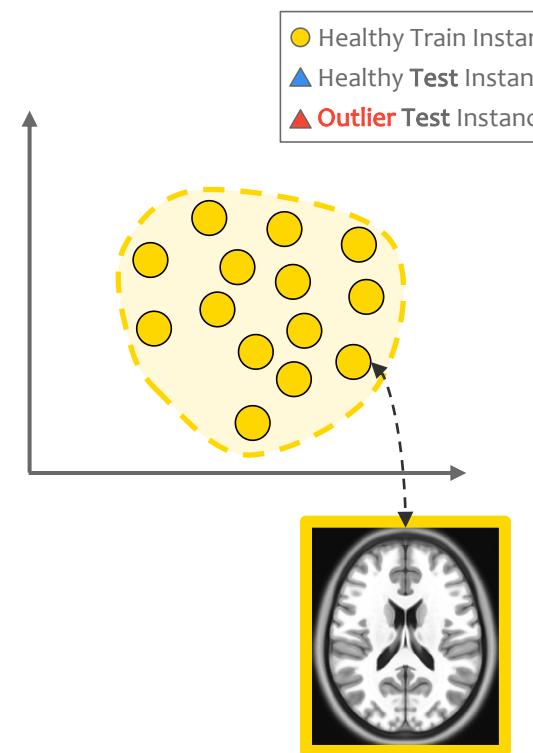
Unsupervised Learning

Anomaly/Novelty/Outlier Detection

The algorithm is fed with training data containing **only** a **single class** (normal, control, healthy).

New unseen instances (test) that **DOES NOT** look a **normal one** is an **outlier**, being considered as an **anomaly**.

The **training set** must be very **clean** and **representative**.



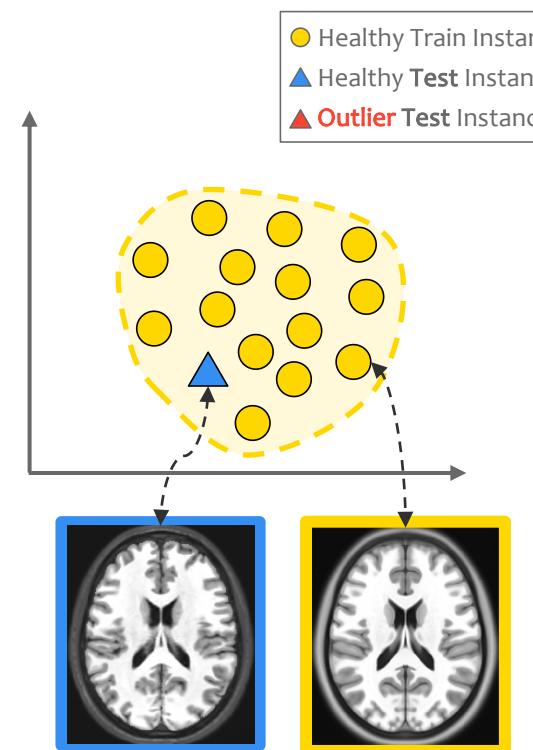
Unsupervised Learning

Anomaly/Novelty/Outlier Detection

The algorithm is fed with training data containing **only** a **single class** (normal, control, healthy).

New unseen instances (test) that **DOES NOT** look a **normal one** is an **outlier**, being considered as an **anomaly**.

The training set must be very **clean** and **representative**.



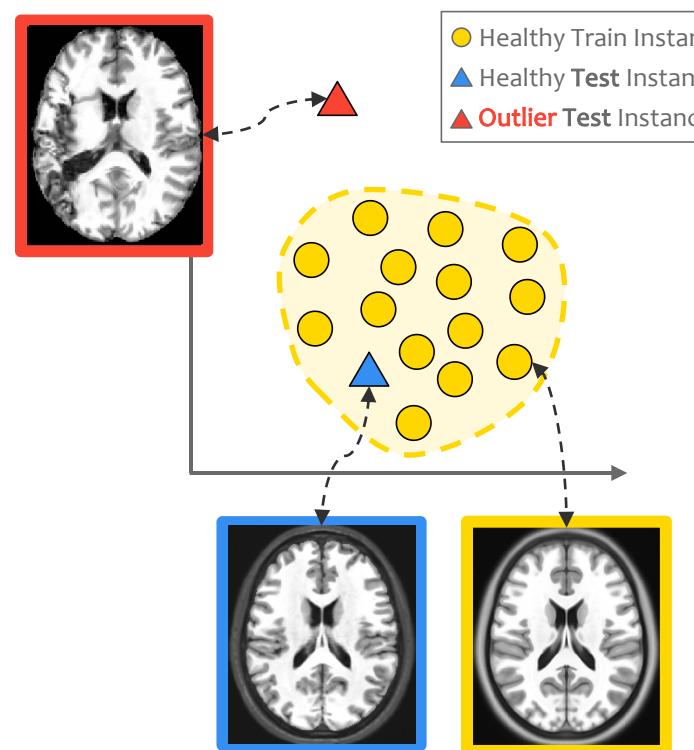
Unsupervised Learning

Anomaly/Novelty/Outlier Detection

The algorithm is fed with **training data** containing **only** a **single class** (normal, control, healthy).

New unseen instances (test) that **DOES NOT** look a **normal one** is an **outlier**, being considered as an **anomaly**.

The **training set** must be very **clean** and **representative**.



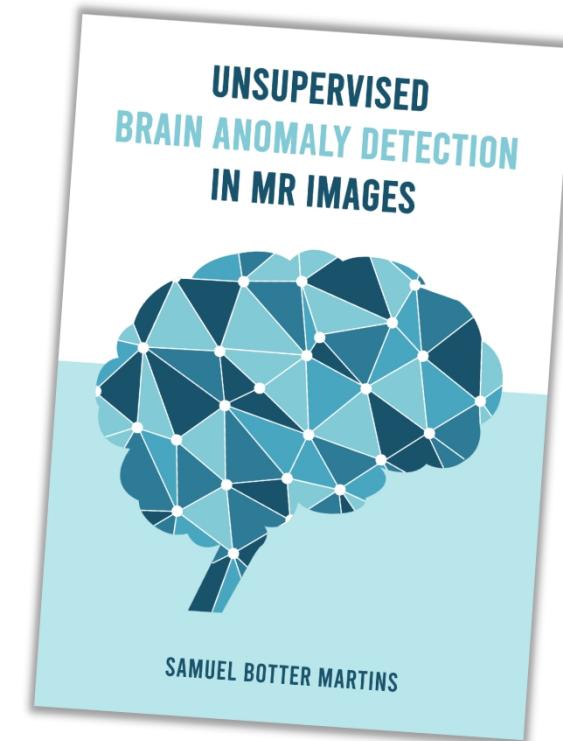
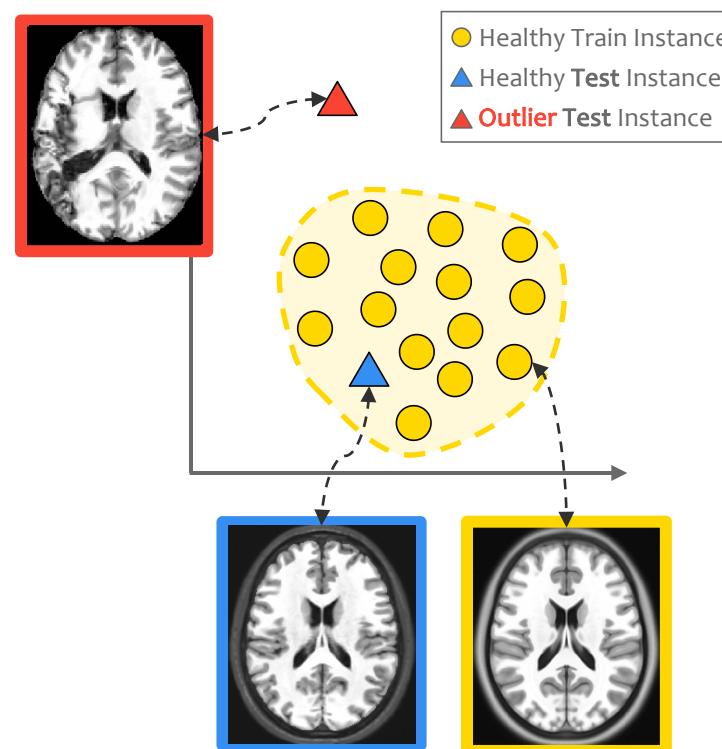
Unsupervised Learning

Anomaly/Novelty/Outlier Detection

The algorithm is fed with training data containing **only** a **single class** (normal, control, healthy).

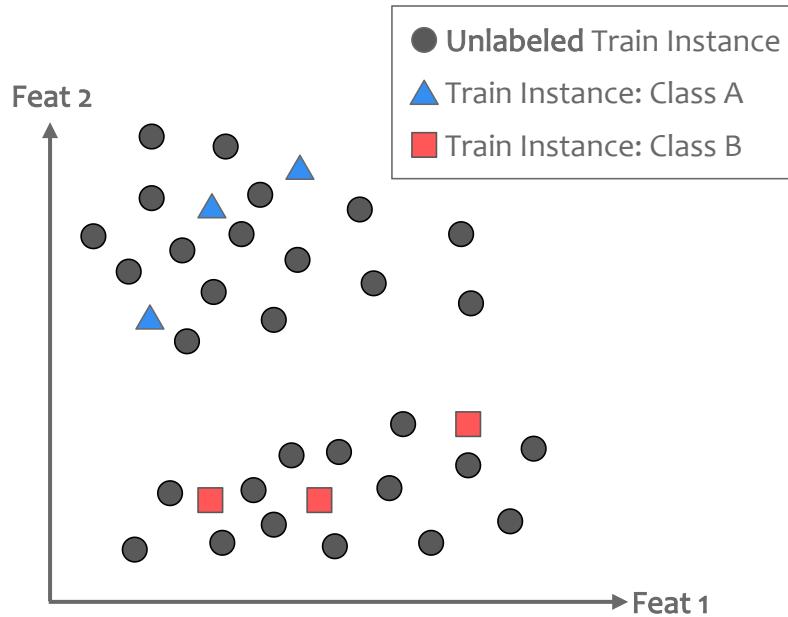
New unseen instances (test) that **DOES NOT** look a **normal one** is an **outlier**, being considered as an **anomaly**.

The training set must be very clean and representative.

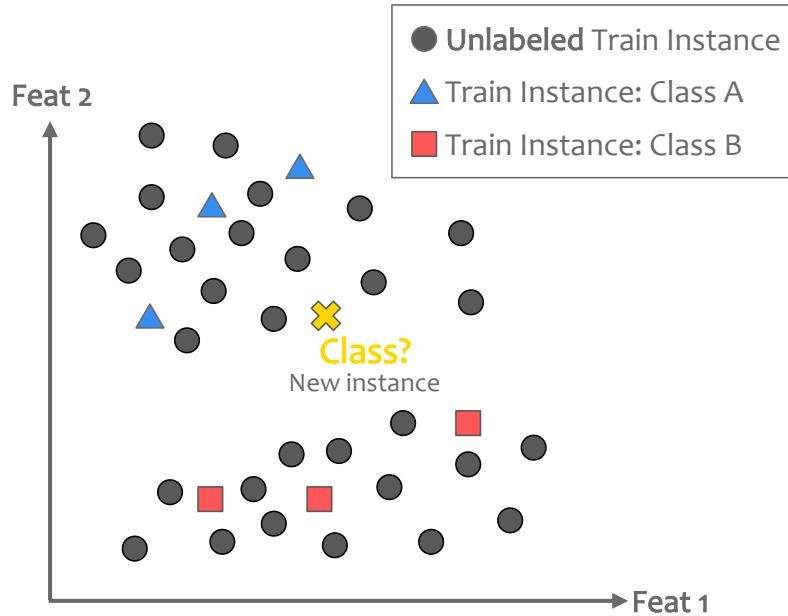


<https://research.rug.nl/en/publications/unsupervised-brain-anomaly-detection-in-mr-images>

Semisupervised Learning

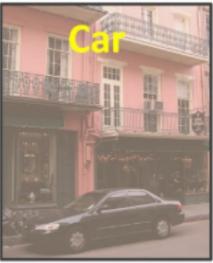
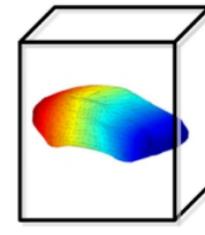


Semisupervised Learning



Weakly Supervised Learning

Supervision with **noisy labels**.

	Image level	Bounding-box level	Pixel level	Voxel/mesh level	annotation
Label cost	1 sec per class	10 sec per instance	78 sec per instance	Not acquirable without expertise	
Label example					

Zhang, Dingwen, et al. "Weakly Supervised Object Localization and Detection: A Survey." IEEE transactions on pattern analysis and machine intelligence (2021).

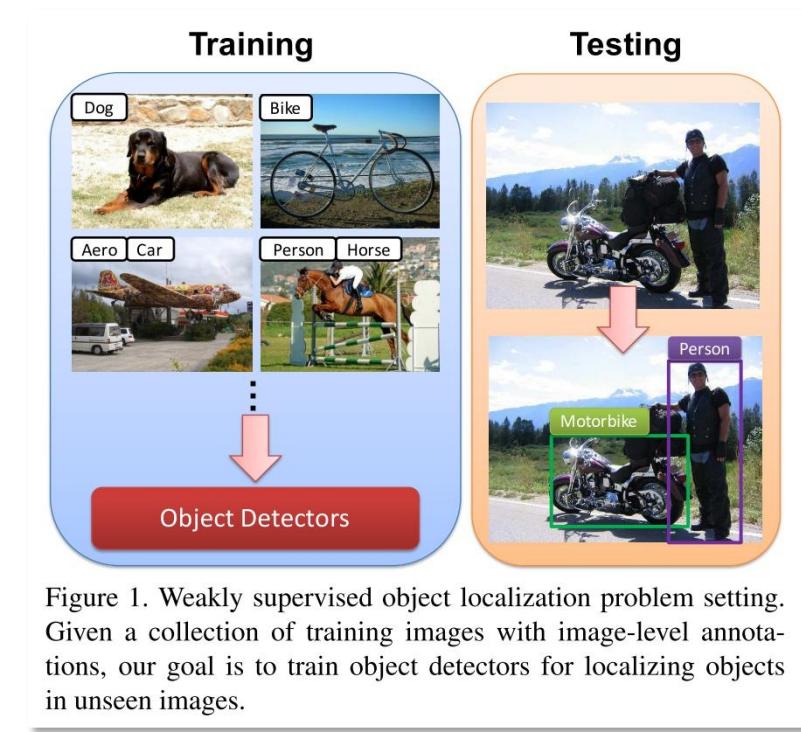
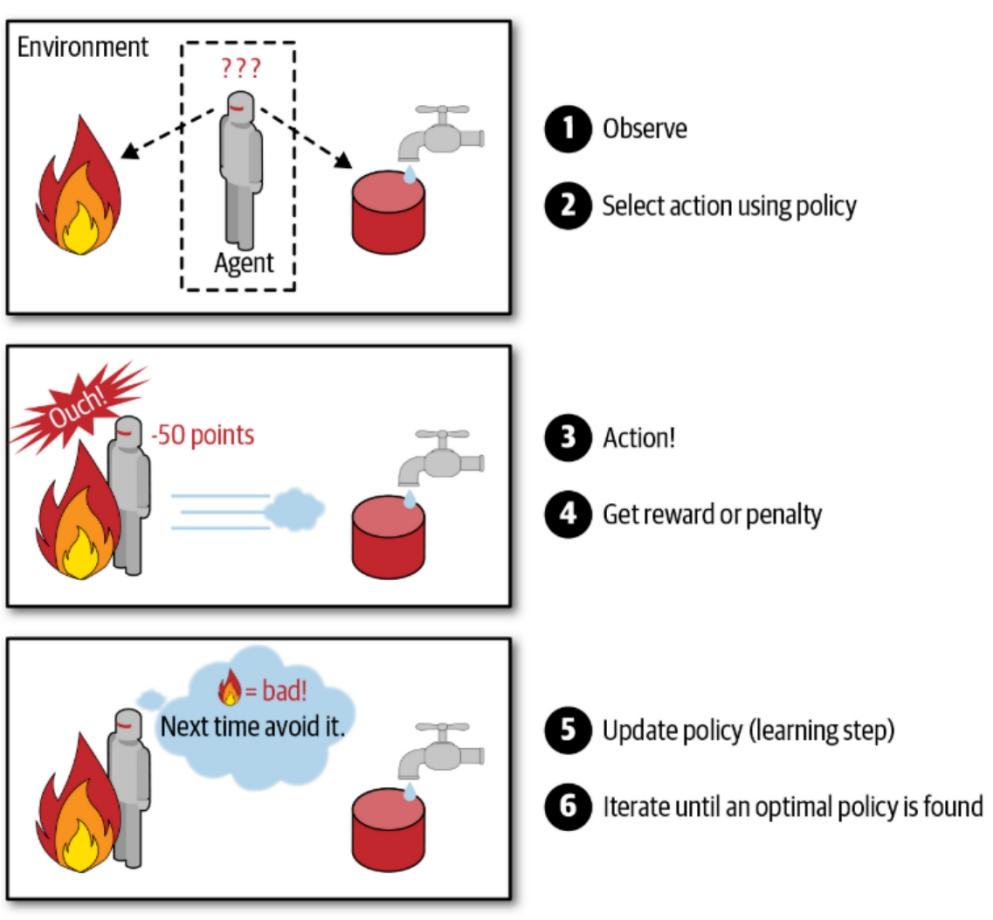


Figure 1. Weakly supervised object localization problem setting. Given a collection of training images with image-level annotations, our goal is to train object detectors for localizing objects in unseen images.

Li, Dong, et al. "Weakly supervised object localization with progressive domain adaptation." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016.

Reinforcement Learning

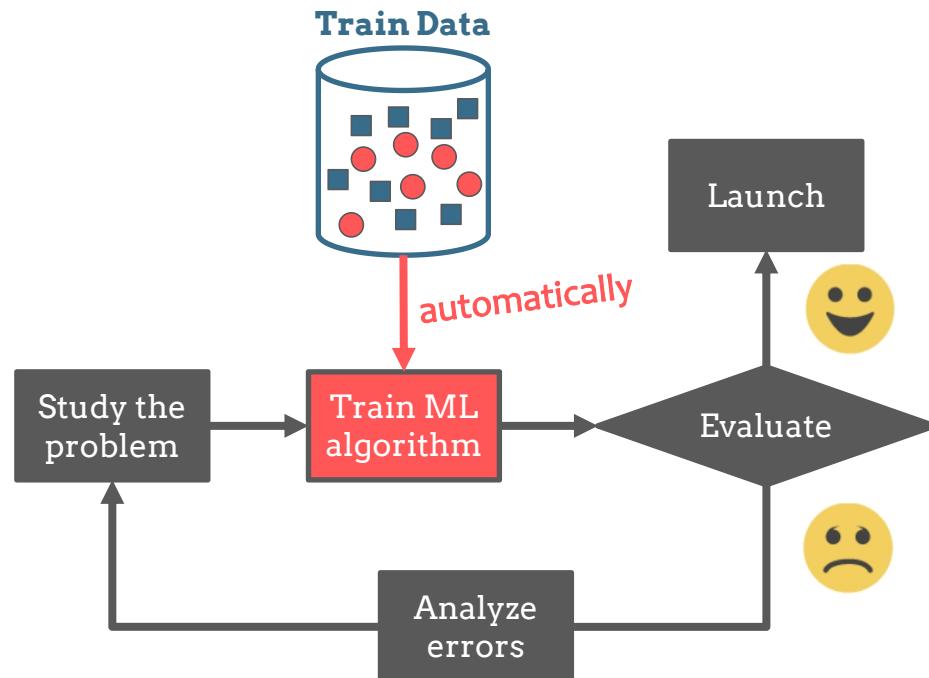


Géron, Aurélien. Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems. O'Reilly Media, 2019.

Batch Learning

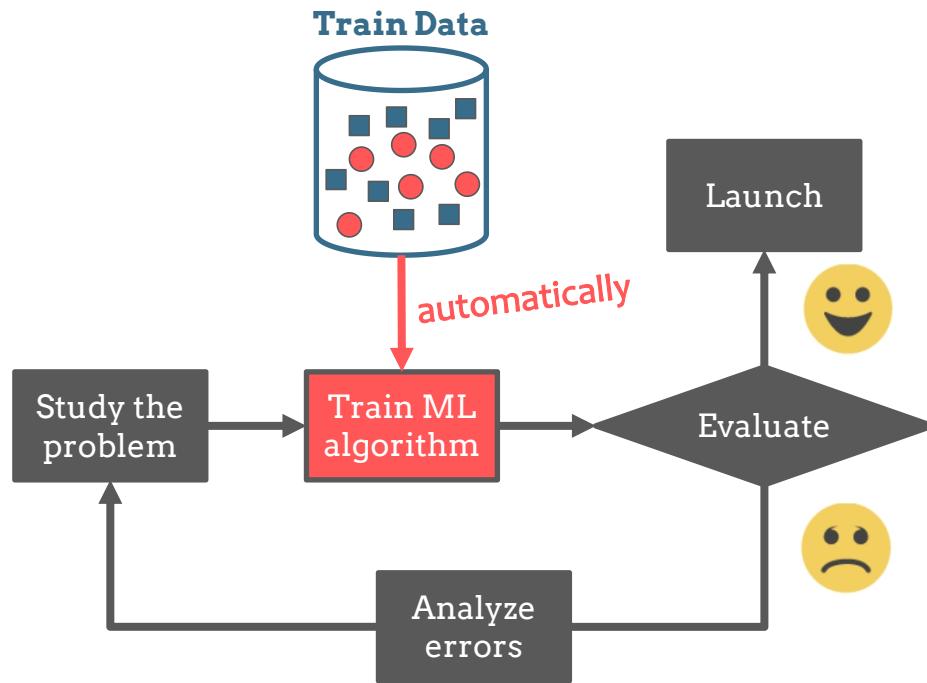
- The system use all available data for *offline training* once and launched into production **without learning anymore** → **no incremental updates**.
- Time-consuming and computational expensive.
- Knowing **new data** requires to **train a new version** of the system from scratch.

Online Learning



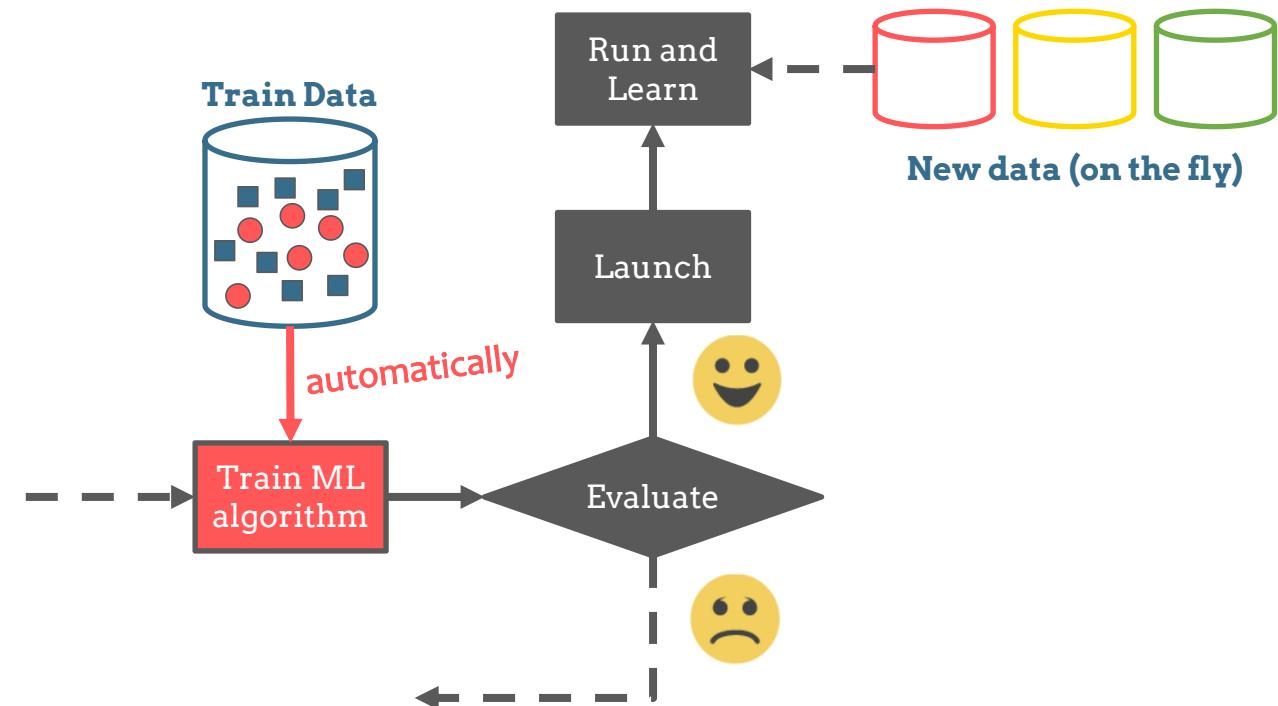
Batch Learning

- The system use all available data for **offline training** once and launched into production **without learning anymore** → **no incremental updates**.
- Time-consuming and computational expensive.
- Knowing **new data** requires to **train a new version** of the system from scratch.



Online Learning

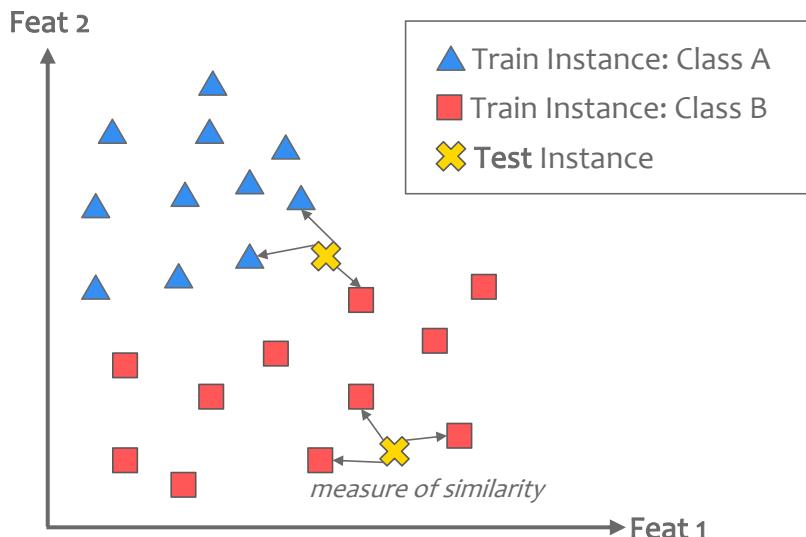
- Train the system **incrementally** by feeding it data instances sequentially, either individually or in small groups called mini-batches;
- Each learning step is **fast and cheap**, so the system can learn about new data **on the fly**, as it arrives.



Instance-based learning

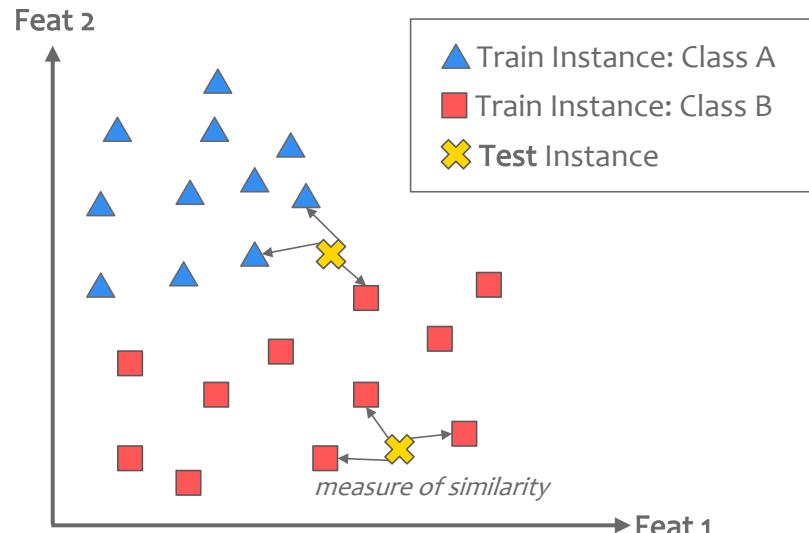
- The system learns the examples **by heart**, then generalizes to **new cases** by using a **similarity measure** to compare them to the learned examples (or a subset of them) → “**there is no training**”.

Model-based learning



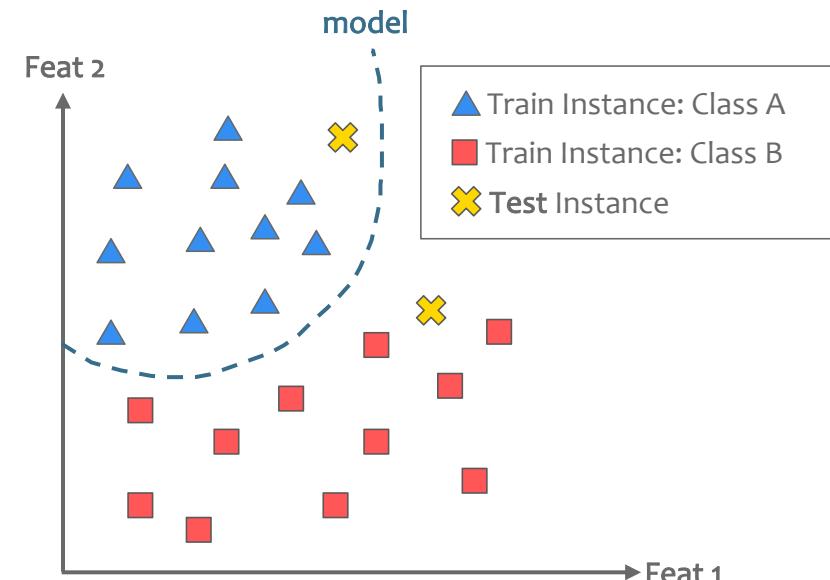
Instance-based learning

- The system learns the examples **by heart**, then generalizes to **new cases** by using a **similarity measure** to compare them to the learned examples (or a subset of them) → “**there is no training**”.



Model-based learning

- The system **trains a model** (function, neural network, ...) from a training data to make ***predictions***.

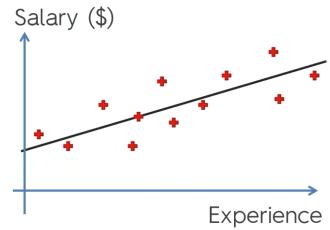


Main Challenges of Machine Learning

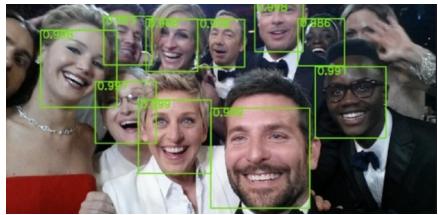
Insufficient Quantity of Train Data

Insufficient Quantity of Train Data

ML algorithms requires a considerable quantity of train data to work properly.



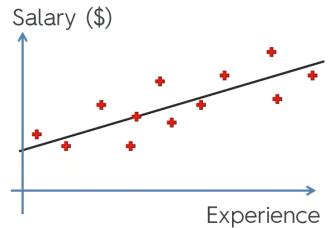
simple problems
(thousands of examples)



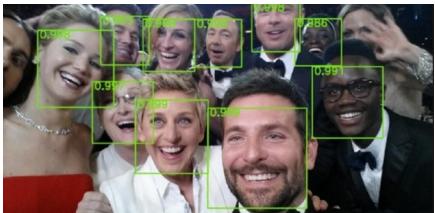
complex problems
(millions of examples)

Insufficient Quantity of Train Data

ML algorithms requires a considerable quantity of train data to work properly.



simple problems
(thousands of examples)



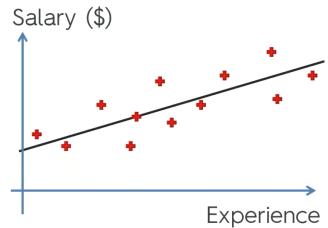
complex problems
(millions of examples)



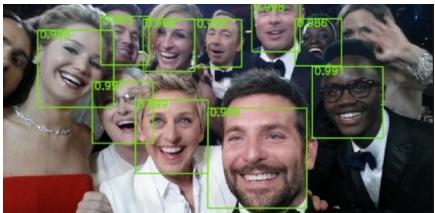
Some authors consider that “data matters more than algorithms for complex problems”.

Insufficient Quantity of Train Data

ML algorithms requires a considerable quantity of train data to work properly.



simple problems
(thousands of examples)



complex problems
(millions of examples)

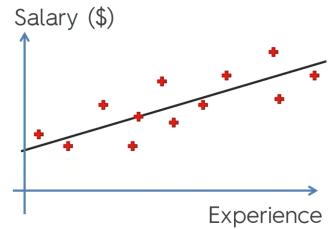


Some authors consider that “data matters more than algorithms for complex problems”.

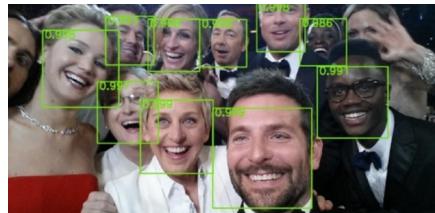
Poor-Quality Data

Insufficient Quantity of Train Data

ML algorithms require a considerable quantity of train data to work properly.



simple problems
(thousands of examples)

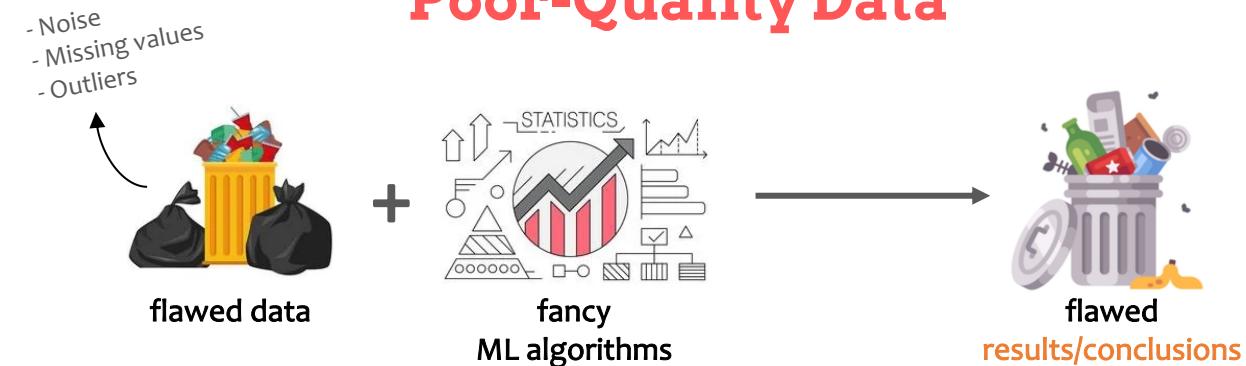


complex problems
(millions of examples)



Some authors consider that “data matters more than algorithms for complex problems”.

Poor-Quality Data



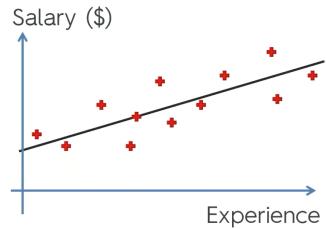
Garbage In



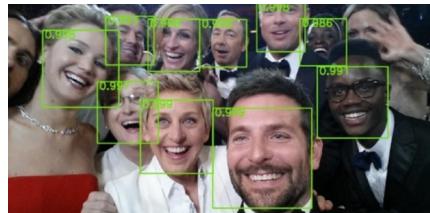
Garbage Out

Insufficient Quantity of Train Data

ML algorithms require a considerable quantity of train data to work properly.



simple problems
(thousands of examples)



complex problems
(millions of examples)



Some authors consider that “data matters more than algorithms for complex problems”.

Nonrepresentative Train Data

Nonrepresentative train data refers to data that does not accurately reflect the real-world distribution of the population it is intended to represent.

- Noise
- Missing values
- Outliers



flawed data



fancy
ML algorithms



flawed
results/conclusions

Poor-Quality Data

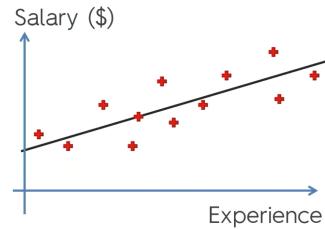
Garbage In



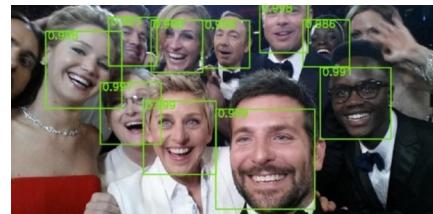
Garbage Out

Insufficient Quantity of Train Data

ML algorithms require a considerable quantity of train data to work properly.



simple problems
(thousands of examples)



complex problems
(millions of examples)



Some authors consider that “data matters more than algorithms for complex problems”.

Nonrepresentative Train Data

It is crucial to use a **training set** that is **representative** of the cases you want to generalize to.

E.g.: Predict the risk of infant mortality of a **city population**.



fancy neighborhood



poor neighborhood

Poor-Quality Data



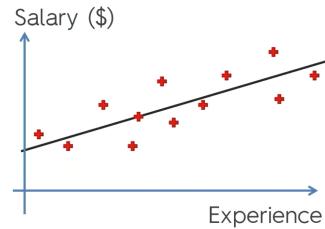
Garbage In



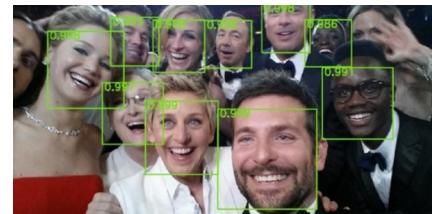
Garbage Out

Insufficient Quantity of Train Data

ML algorithms require a considerable quantity of train data to work properly.



simple problems
(thousands of examples)



complex problems
(millions of examples)



Some authors consider that “data matters more than algorithms for complex problems”.

Nonrepresentative Train Data

It is crucial to use a **training set** that is **representative** of the cases you want to generalize to.

E.g.: Predict the risk of infant mortality of a **city population**.



fancy neighborhood



poor neighborhood

- Noise
- Missing values
- Outliers



flawed data



fancy
ML algorithms



flawed
results/conclusions

Poor-Quality Data

Irrelevant Features

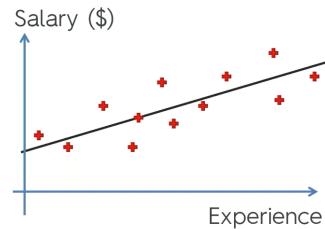
Garbage In



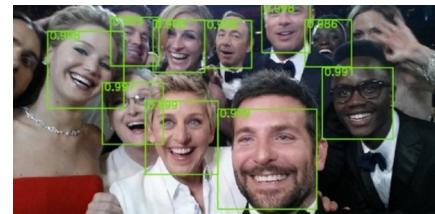
Garbage Out

Insufficient Quantity of Train Data

ML algorithms require a considerable quantity of train data to work properly.



simple problems
(thousands of examples)



complex problems
(millions of examples)



Some authors consider that “data matters more than algorithms for complex problems”.

Nonrepresentative Train Data

It is crucial to use a **training set** that is **representative** of the cases you want to generalize to.

E.g.: Predict the risk of infant mortality of a **city population**.



fancy neighborhood



poor neighborhood

- Noise
- Missing values
- Outliers



flawed data



fancy
ML algorithms



flawed
results/conclusions

Poor-Quality Data

Irrelevant Features

E.g.: Predict employees' salaries.

Experience	Education	Soccer Team	Favorite Color	Salary
60	Undergraduation	Barcelona	Blue	10,000.00
75	MSc	Real Madrid	Red	12,000.00
100	PhD	Botafogo	Black	15,000.00
...

Feature engineering = feature selection + feature extraction

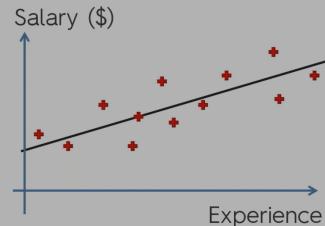
Garbage In



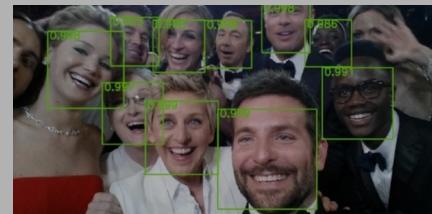
Garbage Out

Insufficient Quantity of Train Data

ML algorithms require a considerable quantity of train data to work properly.



simple problems
(thousands of examples)



complex problems
(millions of examples)



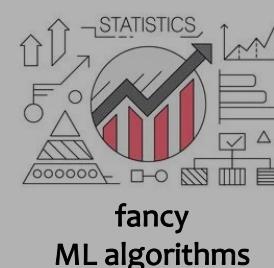
Some authors consider that “data matters more than algorithms for complex problems”.



Other challenges of ML, such as **overfitting** and **underfitting**, will be discussed later in this course.

Poor-Quality Data

- Noise
- Missing values
- Outliers



flawed data

fancy
ML algorithms

Garbage In → Garbage Out

Nonrepresentative Train Data

It is crucial to use a **training set** that is **representative** of the cases you want to generalize to.

E.g.: Predict the risk of infant mortality of a **city population**.



neighborhood



poor neighborhood

Irrelevant Features

E.g.: Predict employees' salaries.

Experience	Education	Soccer Team	Favorite Color	Salary
60	Undergraduation	Barcelona	Blue	10,000.00
75	MSc	Real Madrid	Red	12,000.00
100	PhD	Botafogo	Black	15,000.00
...

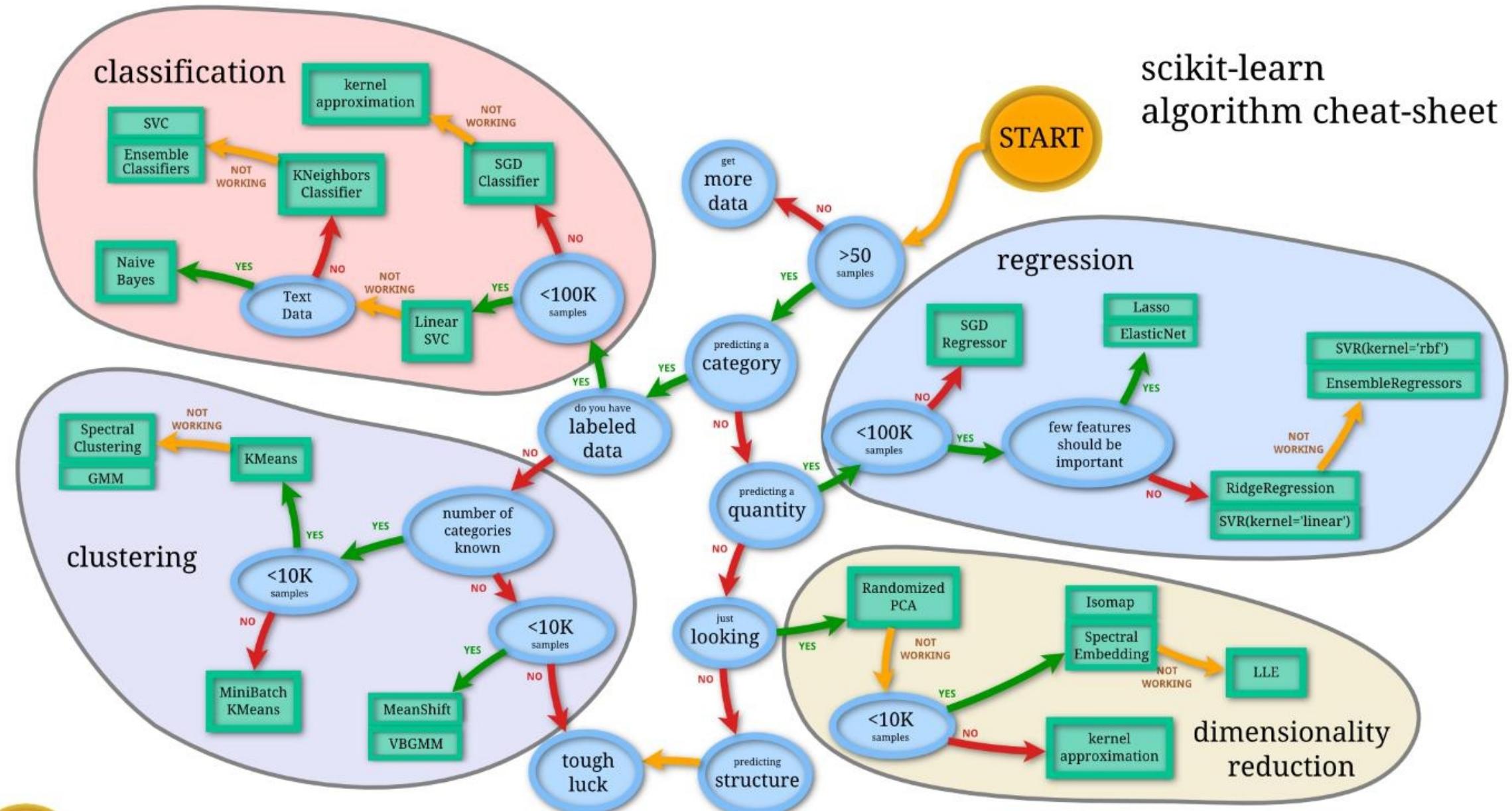
Feature engineering = feature selection + feature extraction

No Free Lunch Theorem

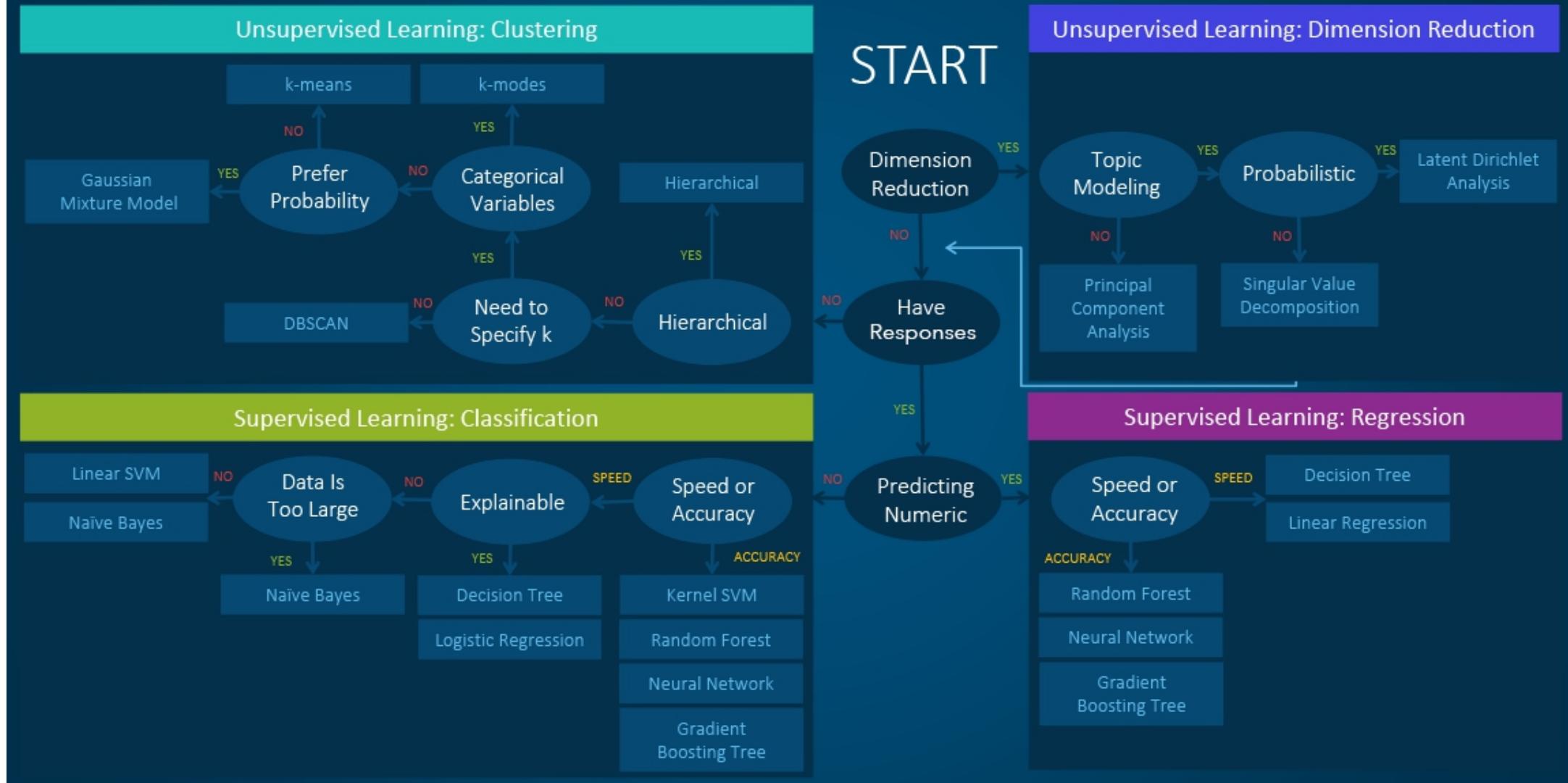
- There is no model that is *a priori* guaranteed to work better for some problems/datasets;
- The only way to know for sure which model is best is to evaluate them all.
- Since this is not possible, in practice you make some reasonable assumptions about the data and evaluate only a few reasonable models.

Wolpert, David H. "The lack of a priori distinctions between learning algorithms." Neural computation 8.7 (1996): 1341-1390.

scikit-learn algorithm cheat-sheet

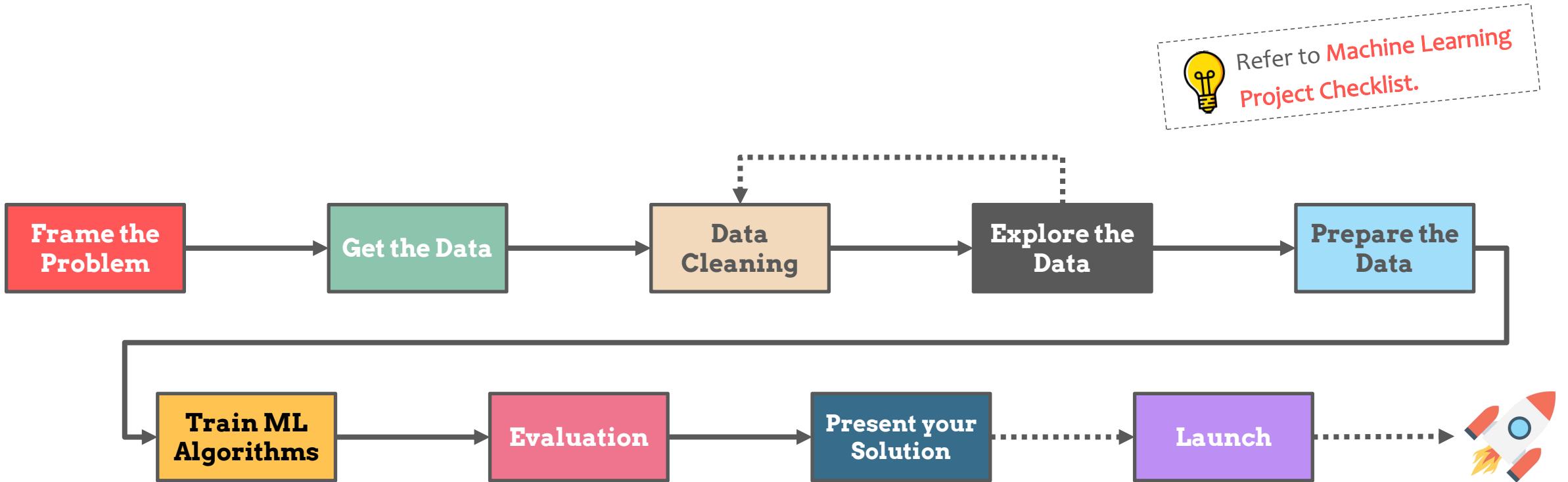


Machine Learning Algorithms Cheat Sheet



The Machine Learning Pipeline

An adapted version from: Géron, Aurélien. *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems.* O'Reilly Media, 2019.



Aprendizado de Máquina e Reconhecimento de Padrões 2021.2



The Machine Learning Landscape

Prof. Samuel Martins (Samuka)
samuel.martins@ifsp.edu.br

