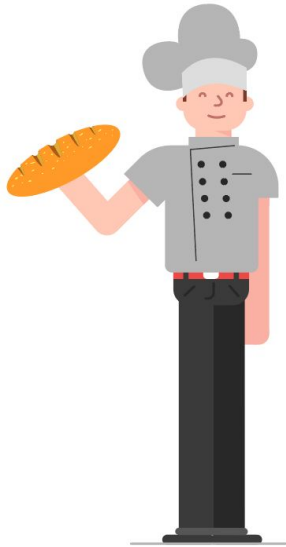


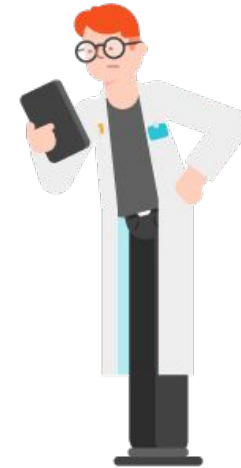
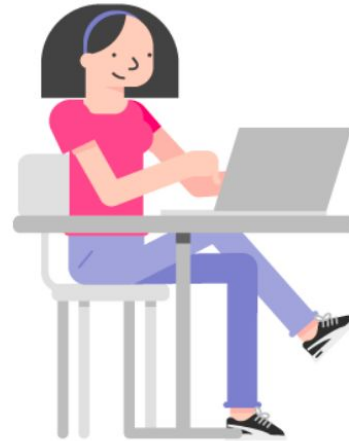
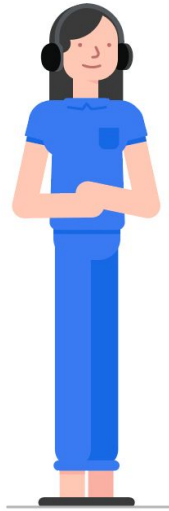
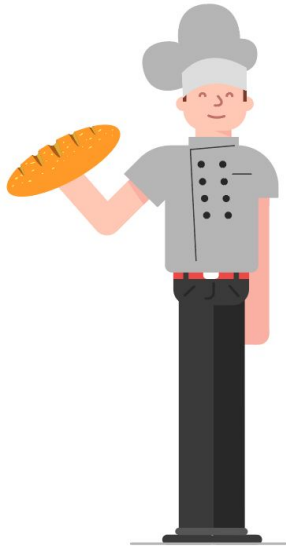
Modelling Methodologies

Sebastien Perez Vasseur

This story starts with people ...

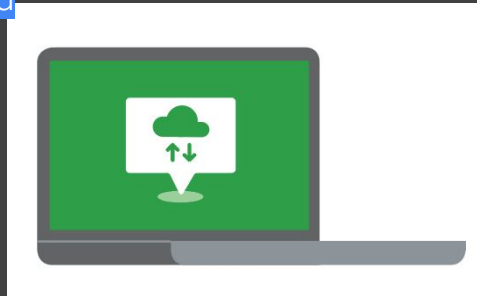


They already have solutions for their business

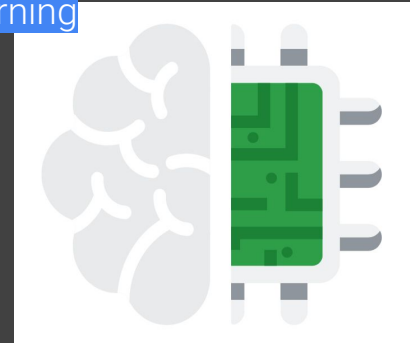


But the world has changed !

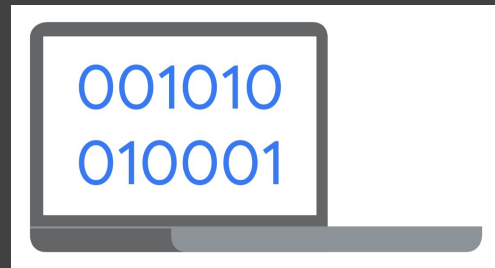
Cloud



Machine Learning



Data



IoT y sensors



In summary



What is Machine Learning ?

"Machine learning is a field of computer science that uses statistical techniques to give computer systems the ability to **"learn" with data**, without being explicitly programmed."

Wikipedia

Cognitive Tasks

In the same way, physical tasks were automated in the 70s, cognitive tasks are automated with Machine Learning.

The tasks are simple and need to be assembled in order to produce results.



There can be too much automation ...

Too much automation can lead to a bad quality of output



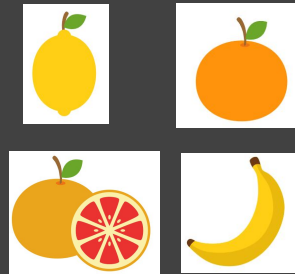
<https://techcrunch.com/2019/03/05/elon-musk-wasnt-wrong-about-automating-the-model-3-assembly-line-he-was-just-ahead-of-his-time/>

Types of learning

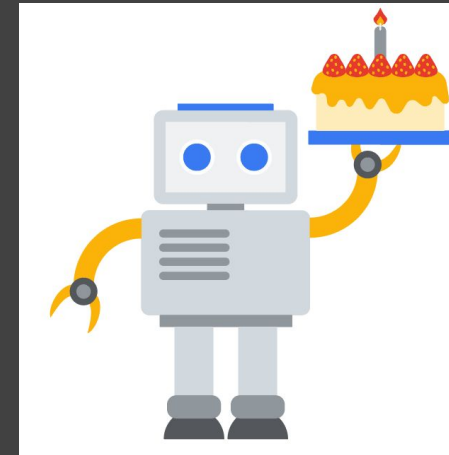
Generalization
by relation



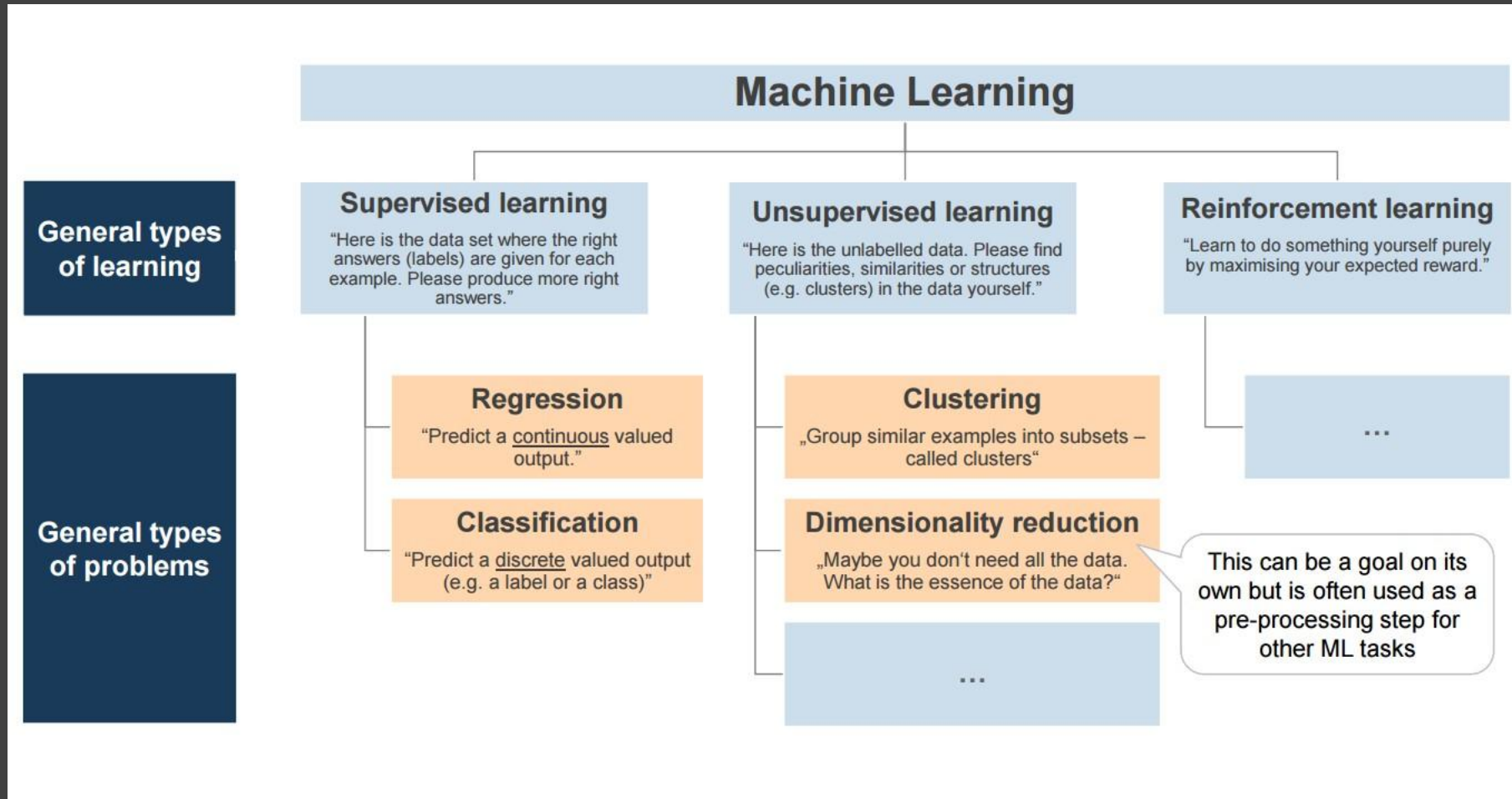
Comparison



Reinforcement



Which types of cognitive tasks ?



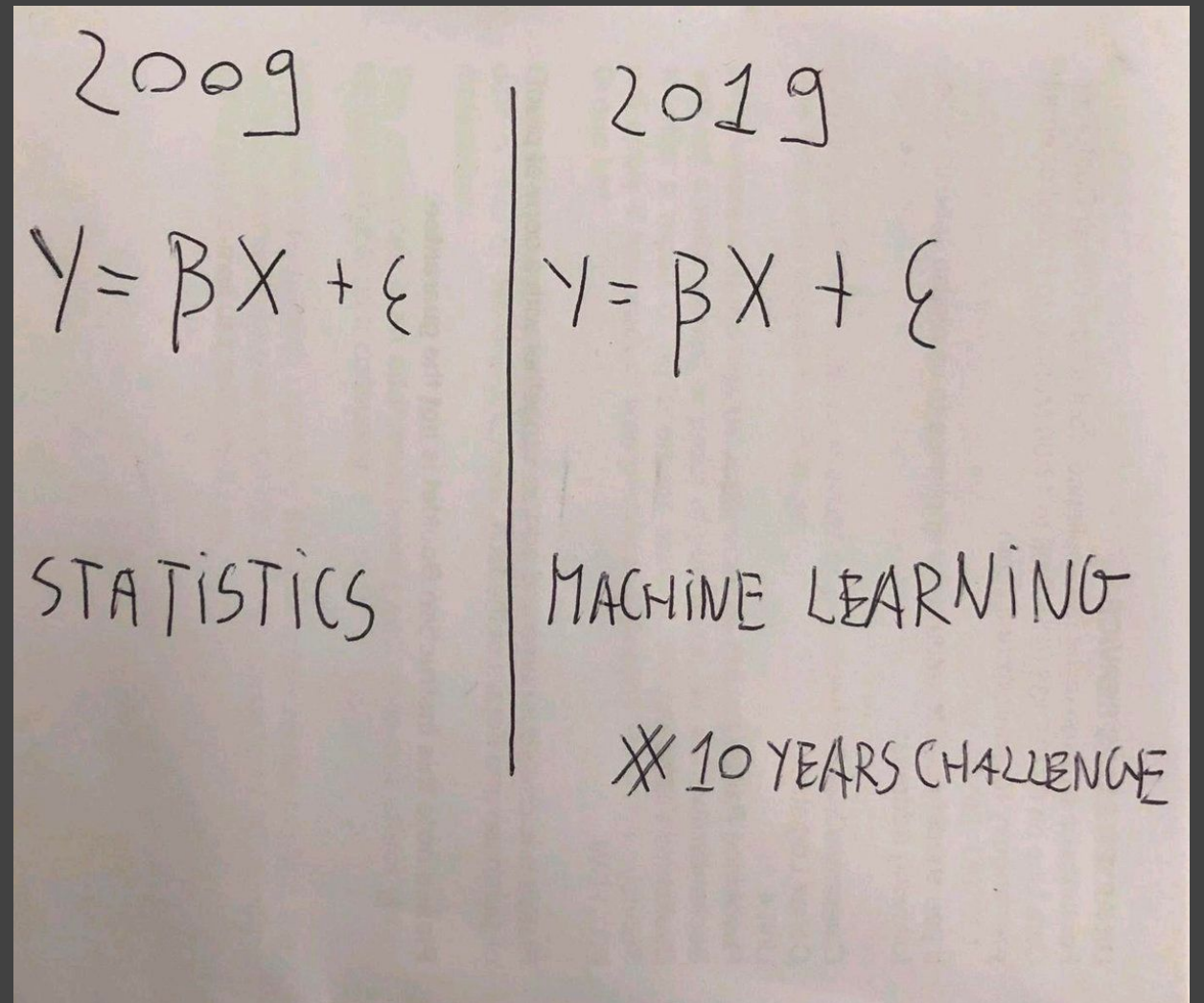
Which types of cognitive tasks ?

- Supervised Learning (Learning from **past** elements):
 - Regression: Predicting a number
 - Classification: Predicting a label
- Unsupervised (Learning by **comparison**):
 - Clustering: Finding elements alike
 - Dimensionality Reduction: Explaining elements with less attributes

Machine Learning

- Vs Artificial Intelligence:
 - Artificial Intelligence may use machine learning to explain how humans think.
- Vs Deep Learning
 - Deep Learning is a type of machine learning models
- Vs Statistics
 - Statistics are a field of Mathematics. Some aspects of Statistics have been renamed as Machine Learning.

#10yearchallenge



What is ML about ?

ML Example: Predicting House Prices

Problem Statement:

We would like to predict the price of a house according to its characteristics.

Available Data:

- Catastro
- House Websites: Idealista, CasaNueva, ...
- ...



Machine Learning is about features

- Machine Learning requires explaining reality as a table of features

Data

```
100100011101000000101000110111010110
100100111101110000001111100110100100
100001101101111101010011100001101001
111111010000110111001010111100001011
11001111101111111100100001110110110
010000110100110110000110000100010000
010101110011001111011001110100010111
001000010101100101000001000010011110
0111010011111100101110101010111100
100010000101100010101101010111000101
010010000100101011110011100001010000
010110000010011101010010101110110001
011011111010111100010100010100010000
011010011011011010001000101111001101
000101000001100110001100100010010110
100101010100010011100101010101111101
```



Id	LotArea	GrLivArea	GarageArea	PoolArea	Sale Price
0	8450	196.0	1710	548	208500
1	9600	0.0	1262	460	181500
2	11250	162.0	1786	608	223500
3	9550	0.0	1717	642	140000
4	14260	350.0	2198	836	250000

Features Description

Id	LotArea	GrLivArea	GarageArea	PoolArea	Sale Price
0	8450	196.0	1710	548	208500
1	9600	0.0	1262	460	181500
2	11250	162.0	1786	608	223500
3	9550	0.0	1717	642	140000
4	14260	350.0	2198	836	250000



Each line represents a house

Target

Machine Learning is about models

- Machine Learning use **models** to explain relationships between features.
- In this case, we look for a function that explains the saleprice:

$$\text{SalePrice} = f(\text{Features})$$

Example 1 of Model: Linear Regression

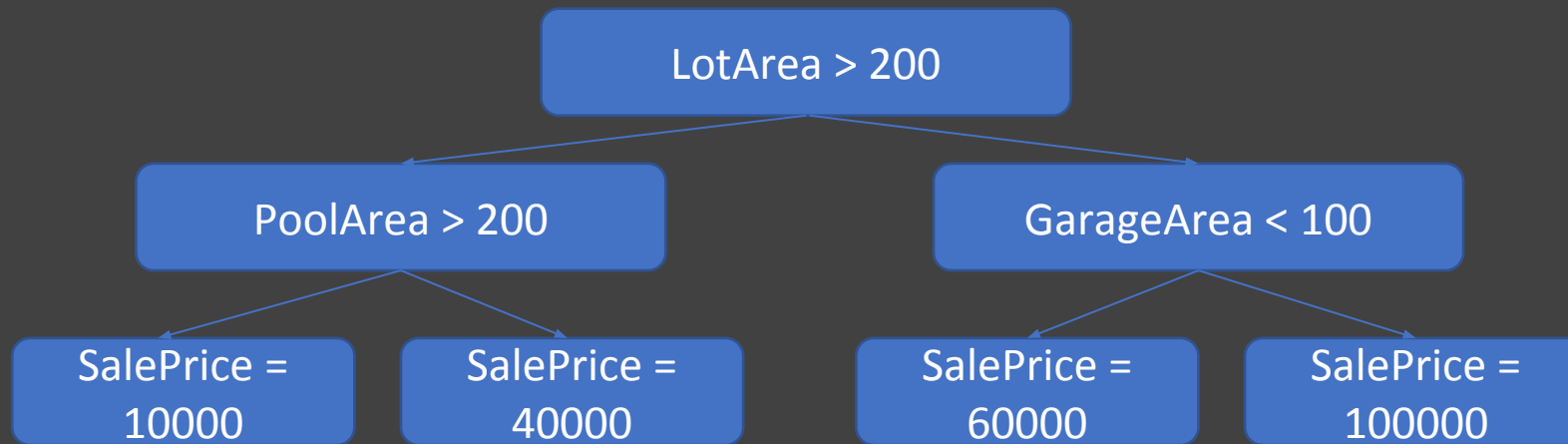
- A linear regression tries to find a linear function, in this case:

$$\text{SalePrice} = a \times \text{LotArea} + b \times \text{GrLivArea} + c \times \text{GarageArea} + d \times \text{PoolArea}$$

- We need to find the best a , b , c and d to have the best relationship.

Example 2 of Model: Decision Tree

- A decision tree tries to find a path to explain the target:



- We need to find the best splits for our predictions.

And there are many more ...

- Every Data Scientist has different types of models for a given cognitive task:
 - Linear Regression
 - Decision Tree
 - Random Forest
 - K-neighbors
 - SVM
 - Neural Networks, Deep Learning, ...
 - And the list is growing thanks to the active research.

Parameters of the model

The Data Scientist works to find the best parameters for the model.

Linear Regression: a, b, c and d .

Decision Tree: Number of Splits, ...

How do you find the best parameters ?

Training

- Each model has a particular way to find the best parameters.
- We say the model is trained when it has found the best parameters with the available data
- We called this Data the Training Data

Machine Learning is about Metrics

There are several metrics that can be chosen for this task. For example:

- Bias: Average of the errors
- MAE: Average of the absolute value of errors
- RMSE: Square root of Average of the square of errors ...

We chose the model parameters that provide the best METRIC. In this case, let's use MAE:

$$\frac{\text{RealPrice1-PredictedPrice1) + (Real Price2-PredictedPrice2) + ...}{\text{Number Of Houses in Test Set}}$$

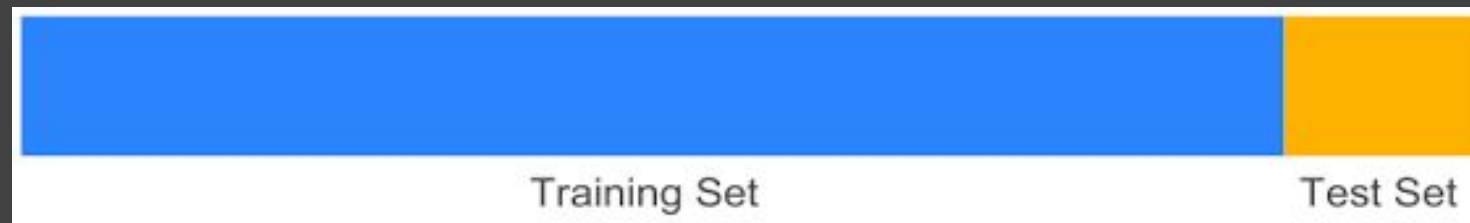
Evaluate Models: Training – Test Split

- We take the available features table and we split in 2 sets:

1 set to calculate parameters : Training Set

1 set to calculate the metrics : Test Set

In the Test Set, we compare the real price of the house with the predicted price.



Final Model

- Once we have the best parameters for a type of model, we can rank the models with the best metric they had:
- Linear Regression - MAE: 15
- Decision Tree – MAE: 10
- And we ultimately chose the model with the best metric.

Model Usage

- We can now use this model to “predict” the SalePrice of a new house from its characteristics:

$$\text{SalePrice} = F(\text{House features})$$

Summary: Machine Learning is about features, models and metrics

- Machine Learning requires explaining reality as a table of features
- We then use models to explain relationships between features for the given ML task
- Each model is evaluated according to a metric

ML Example: Detecting diabetes in patients

Classification

- Problem: Detecting diabetes in patients
- Features:

Pregnancies	Glucose	Blood Pressure	Skin Thickness	Insulin	BMI	Diabetes Pedigree Function	Age	Target Outcome
6	148	72	35	0	33.6	0.627	50	1
1	85	66	29	0	26.6	0.351	31	0
8	183	64	0	0	23.3	0.672	32	1
1	89	66	23	94	28.1	0.167	21	0
0	137	40	35	168	43.1	2.288	33	1
5	116	74	0	0	25.6	0.201	30	0
3	78	50	32	88	31	0.248	26	1
10	115	0	0	0	35.3	0.134	29	0

Each line represents a patient

- Metric: Accuracy (Percentage of correctly guessed patient's disease)

ML Example: Detecting groups of similar customers

Segmentation

- Problem: Detecting groups of similar customers

- Features:

CustomerID	Departure Date	Age	Distance	Country of POS
1	22/12/18	20	1000	Spain
2	21/04/18	25	100	France
3	24/05/18	30	5000	Germany
4	30/03/19	55	500	Spain
5	14/03/20	32	200	Germany
6	22/12/18	45	400	Bulgaria
7	20/01/18	43	10000	France
8	12/06/19	12	5000	United Kingdom
9	25/03/18	39	1300	Spain
10	31/07/18	67	700	France

Each line represents a customer

- Method: We find groups of similar customers based on their proximity
- Metric: Silhouette Index (how far are customers from other clusters than theirs)

ML Canvas

Translating a Business Problem to Machine Learning

Any Business Problem can be analyzed through the following axes:

- Business Value and Measure of Success (1)
- Data Sources (2)
- Machine Learning task + Metric (3)
- Features Extraction + Model Creation (4)
- Industrialization and Maintenance (5)

Opportunity:

Estimated Value:

Estimated Cost:

Value Proposition

Business Description
Resulting Action
Measure of Success (KPI)
ML Initiative Cost

Machine Learning

ML task
Methods of evaluation

Ranking of Models

List of the models and their metric

Industrialization

Description
Implementation Cost
Maintenance Cost
Model Specific Cost

Data

Sources
Frequency
Legal
History

Acceptable Quality
Main Features

By:

Iteration:

Date:

Collaboration with Data Scientist

Interaction with Data Scientists involves then discussing about data:

- Data Sources
- History
- Features
- Frequency
- Quality
- Volume

And about measuring success:

- Evaluation of success
- How does it translate to a metric

Data Science Process

Once the first version of the canvas is created, DS perform the following actions:

- Extraction of features (2)
- Selection of ML task and evaluation (3)
- Model Creation (4)
- Quantity of Data > Parameters
- Loop until satisfied

The whole process is iterative

- Product Managers and Data Scientists can revise their initial plans based on discoveries made during the different steps:
 - Insufficient Data
 - Metric not good enough for the task
 - ...

Opportunity:

Estimated Value:

Estimated Cost:

Value Proposition

Machine Learning

Ranking of Models

Industrialization

Business Description
Resulting Action
Measure of Success (KPI)
ML Initiative Cost

ML task
Methods of evaluation

List of the models and their metric

Description
Implementation Cost
Maintenance Cost
Model Specific Cost

Data

Sources
Frequency
Legal
History

Acceptabl
Main Feat

DS
+
PM

By:

Iteration:

Date:

Deliverables

- Model(s) with satisfactory metric:
 - A suitable model (or type of model) is chosen for the task at hand.
 - This model can now do the necessary predictions, but needs to be assembled ...
- Industrial Implementation Discussion:
 - Product Managers, Data Scientists and Development can start discussing how the model will be used, trained, ...

Opportunity:

Estimated Value:

Estimated Cost:

Value Proposition

Business Description
Resulting Action
Measure of Success (KPI)
ML Initiative Cost

Machine Learning

ML task
Methods of evaluation

Data

Sources
Frequency
Legal
History

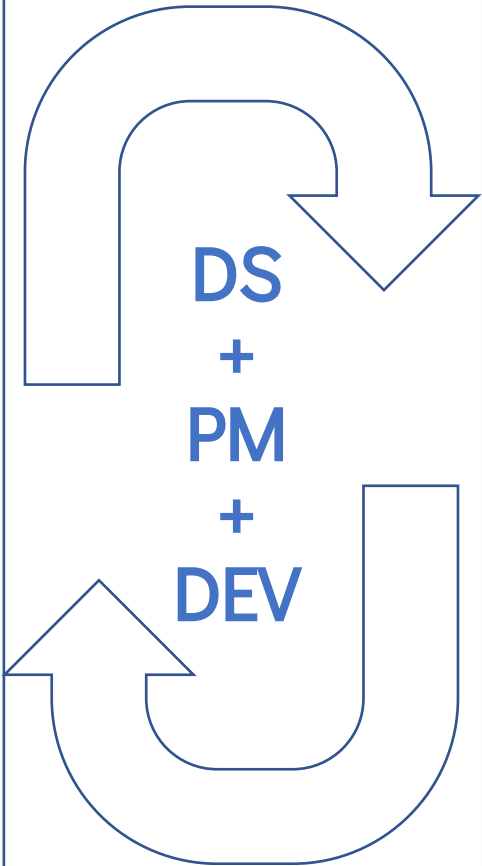
Acceptable Quality
Main Features

Ranking of Models

List of the models and their metric

Industrialization

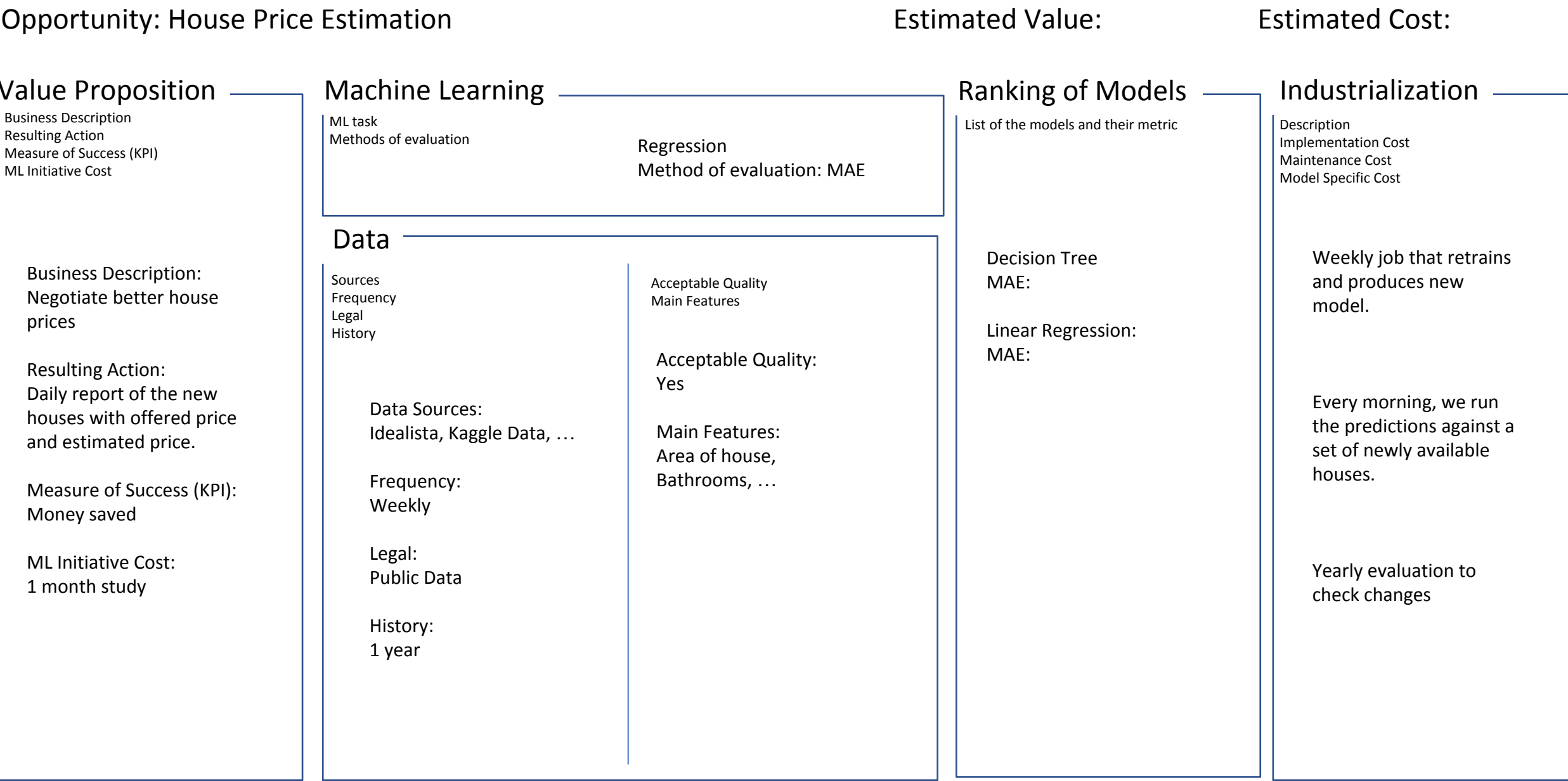
Description
Implementation Cost
Maintenance Cost
Model Specific Cost



By:

Iteration:

Date:



Hands on Exercises

Hands on Exercises

Let's try to do that exercise for those 4 use cases:

- Weather Forecast
- Movie Recommendation
- Email Spam Filters
- Pollution Prediction
- Customer Segmentation

Activity 1: Define your Business (10mn)

Steps:

1. List businesses that are related to the defined problem
2. Choose 1 business

Example: Predict people income

- Bank
- Marketing Company
- Dating App
- ...

Activity 2: Define the Business problem (10mn)

You need to define the business problem

Example: Bike Rental Company

- Rain spoils bike trips
- Warm days register higher occupancy

Activity 3: Business Success (10mn)

We need to define the KPI or milestones the business is seeking

Examples:

- Higher revenue
- Increase conversion rate by x%
- ...

A word on metrics

You still lack knowledge on metrics. However, try to think on creative ways to measure the predictions of models:

- Regression:
 - Percentage of error, Above reference, ...
 - Are negative or positive errors as bad ?
- Classification:
 - Are mistakes as costly for each label ?
 - Is recall more important than precision (Detecting of one label) ?

Activity 4: Define the ML problem and the metric (5mn)

Once we know the business problem, it's time to define the Machine Learning task and the metric associated to it.

Example:

- House Prices: Regression. Metric: MAPE
- Credit: Classification. Metric: Unpaid credit precision
- ...

A word on Data Sources

Data Sources need to be completed in a small table with:

Name	Frequency	Legal	History

Activity 5: Find Data Sources (15mn)

Let's try to find:

- Public Data
 - Google Datasets
 - Government Websites
 - kaggle
 - ...
- Private Data (Purchased or created)

Activity 6: Features Recognition (15mn)

Steps:

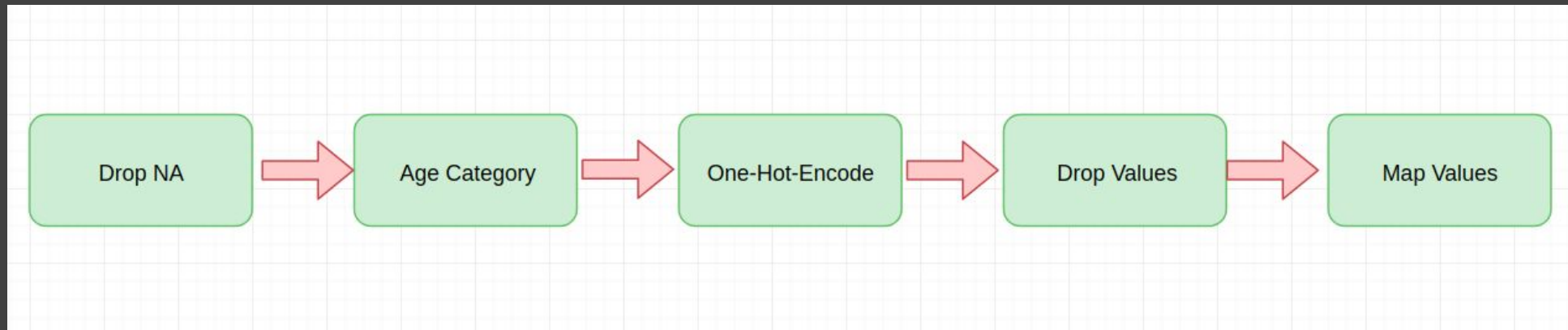
- Analyze the files in the datasources
- Get all the columns
- List the features for the ML table

Activity 7: Draw the data Flow (15mn)

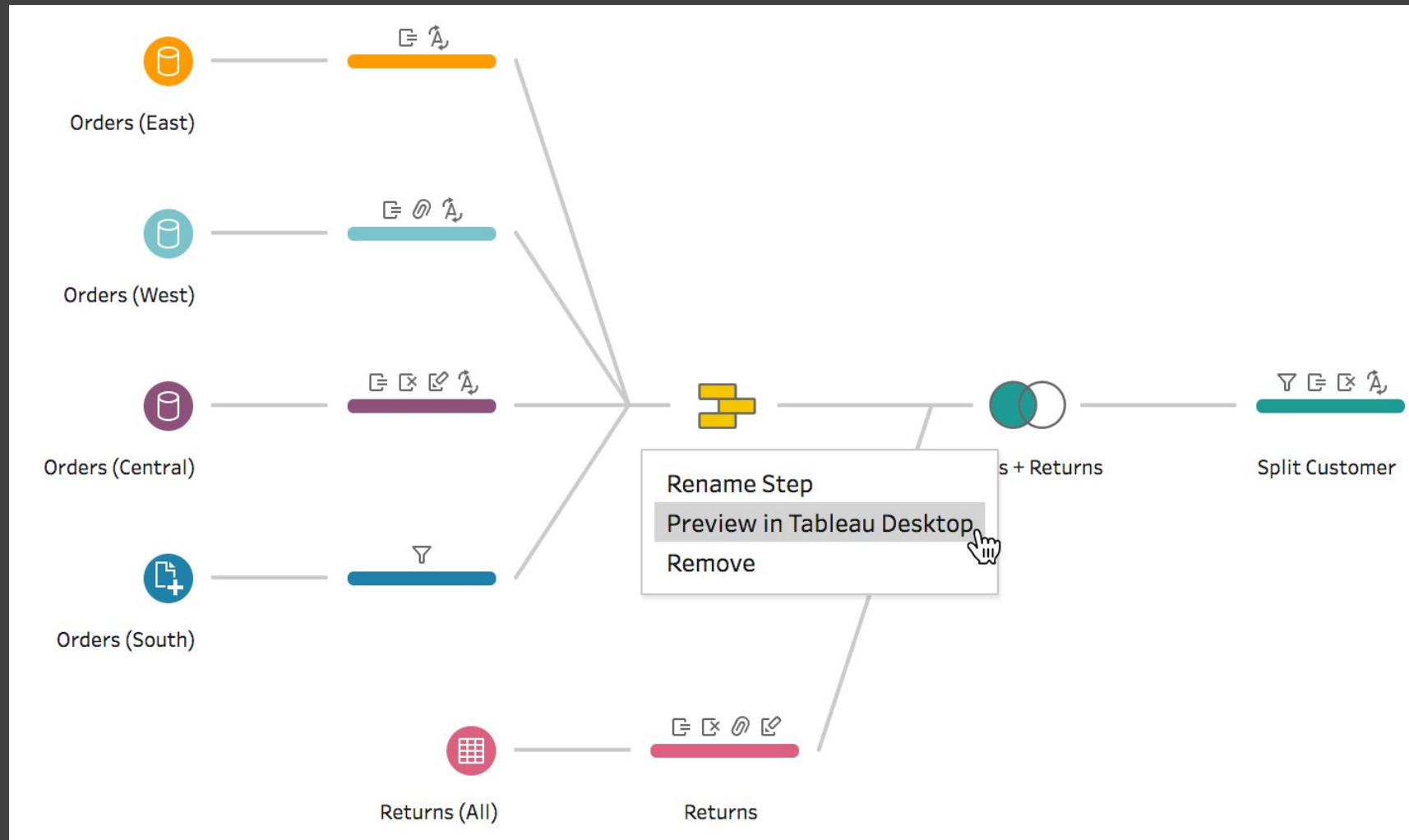
The ML table is created through a data flow with basic steps. (List of steps is sent on Basecamp Forum).

Steps:

- Define initial state (tables) and final state (ML table).
- List the different operations that are performed
- Draw the operations in a piece of paper



Another data flow



Activity 8: Industrialization (20mn)

Steps:

- When is the model used ?
- Which decision is taken with the prediction ?
- Draw a system diagram with the use of the model

Final Presentation (5mn per team)

It's time to present:

- your ML Canvas
- Data Flow diagram
- Industrialization Diagram