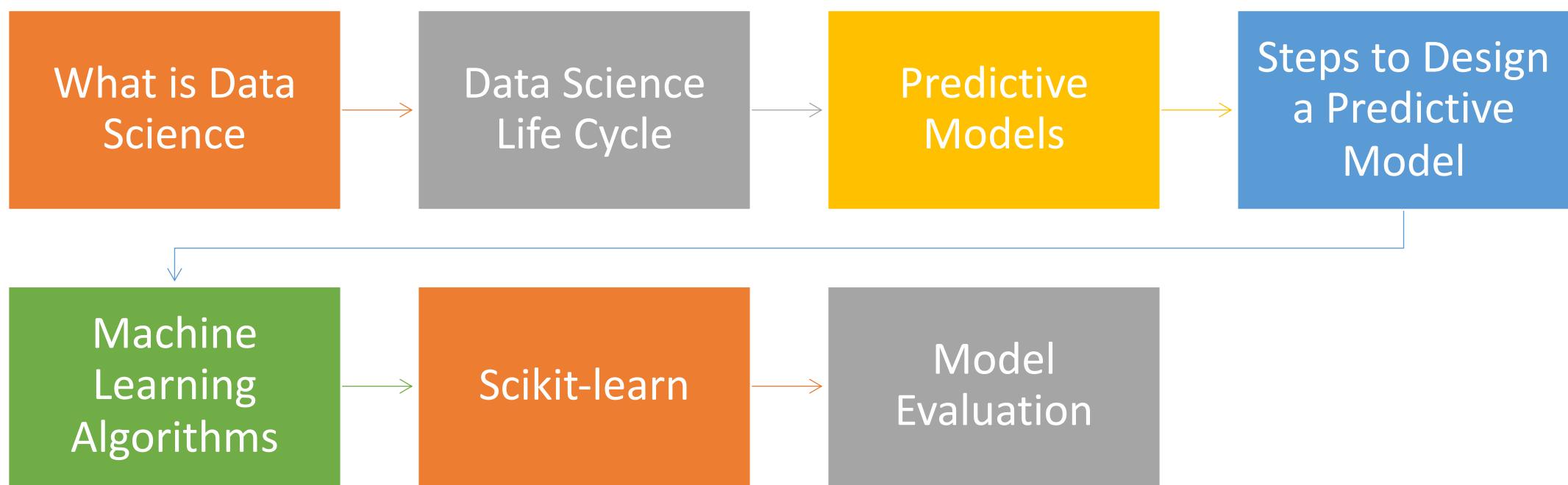


Introduction to Machine Learning



Outline



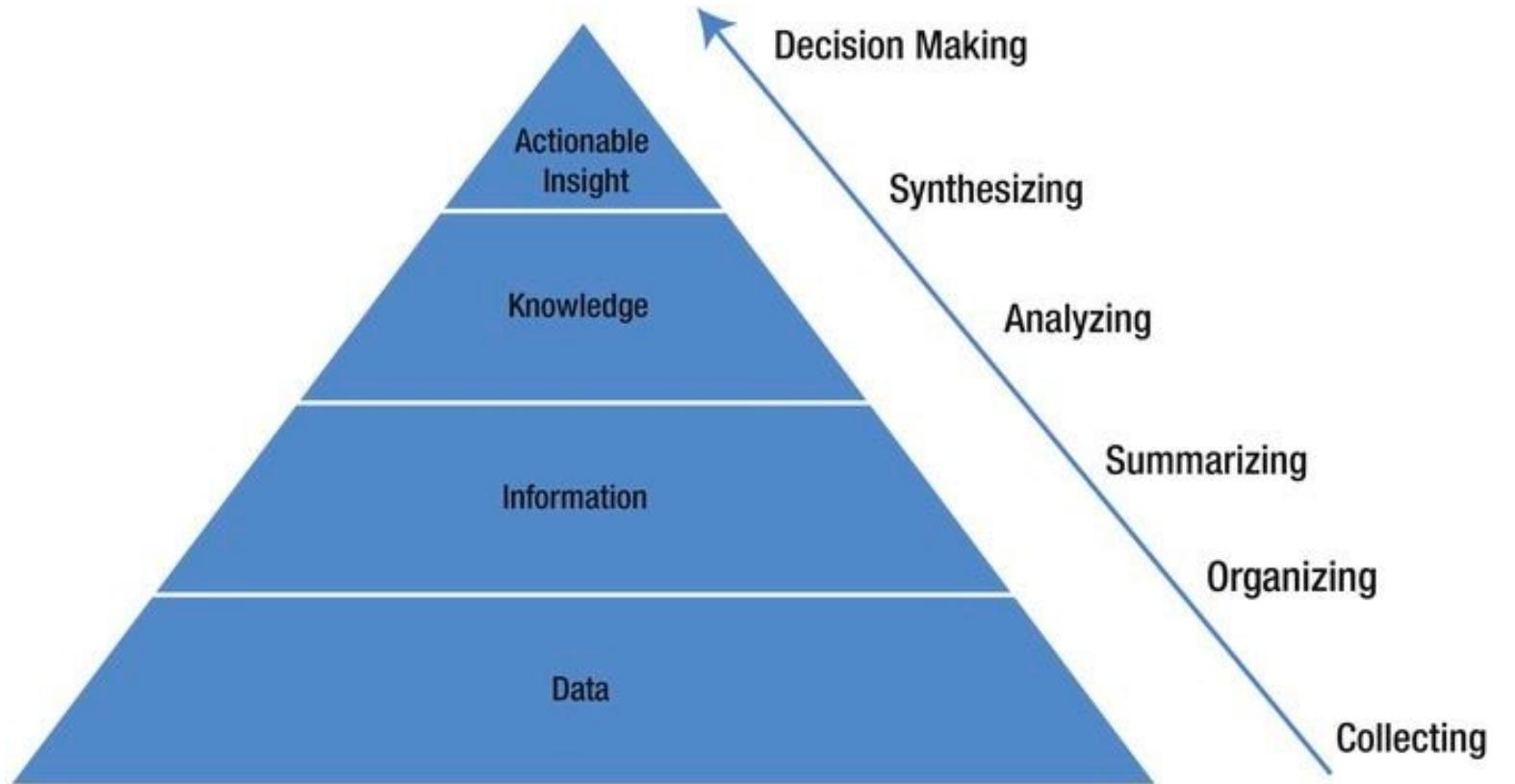
Data

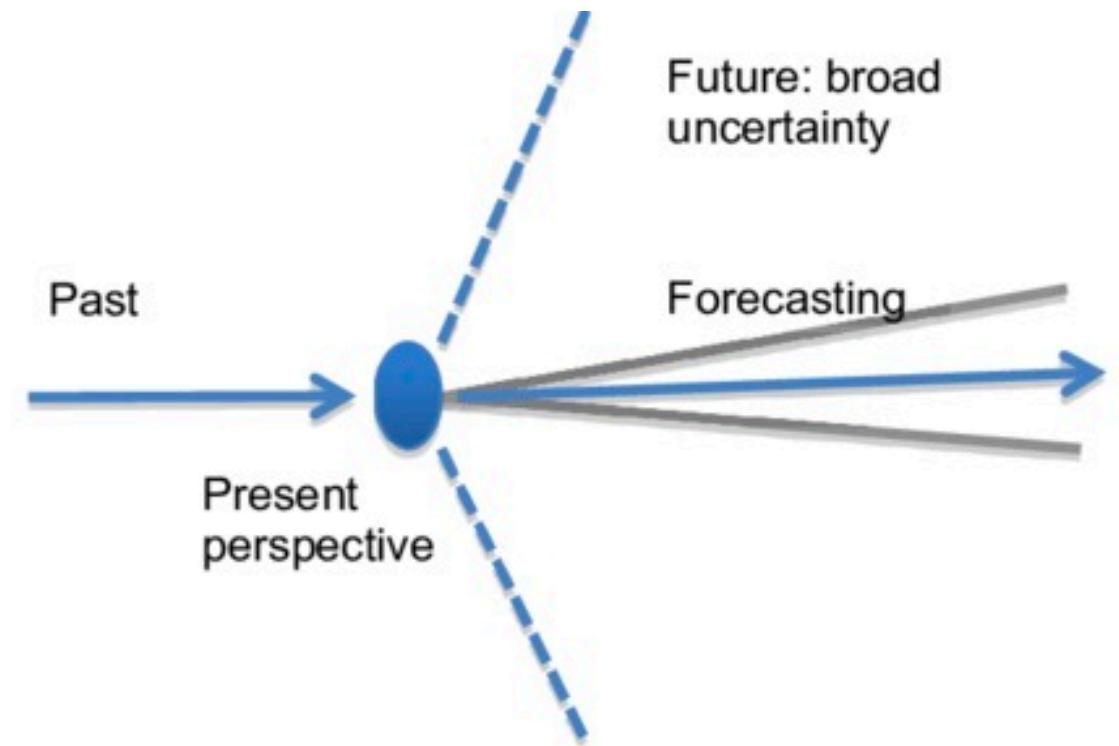
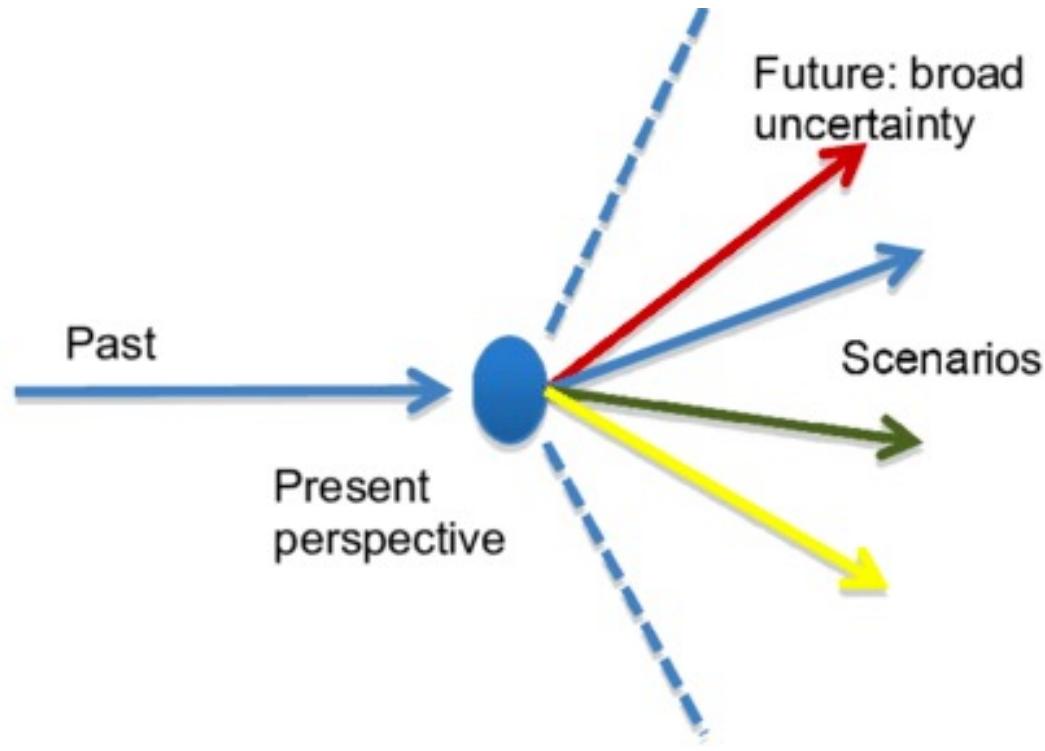
Facts and statistics collected together for reference or analysis.

Science

Systematic enterprise that builds and organizes knowledge in the form of testable explanations and predictions about the universe.

Data Science





Data Science: Uncertainty

Data Science: Questions

**What
happened?**

Descriptive Analytics

**Why did
something
happen?**

Diagnostic Analytics

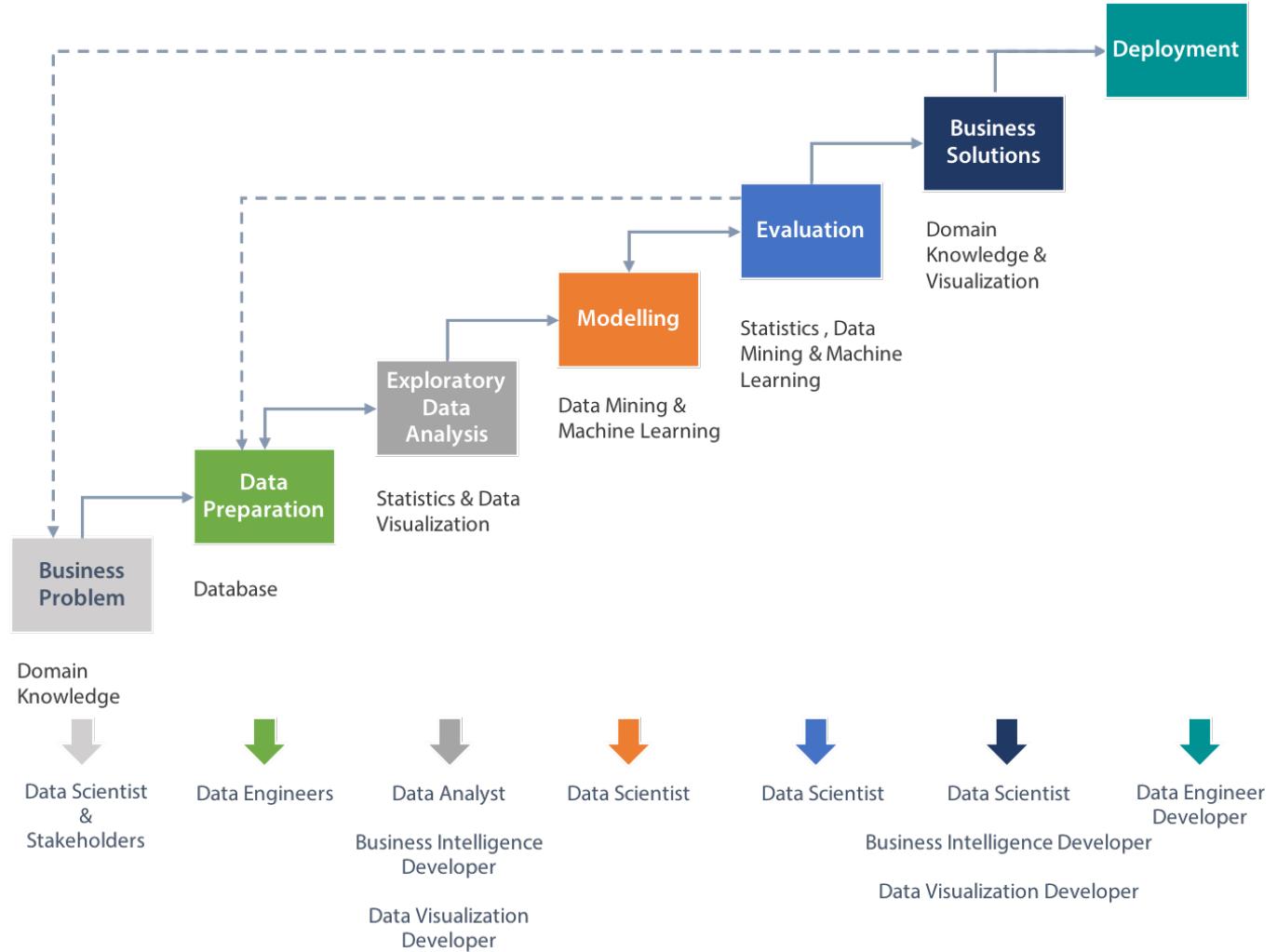
**What will
happen
next?**

Predictive Analytics

**What's the
best
decision
now?**

Prescriptive Analytics

Data Science Life Cycle



Steps to Design a Predictive Model

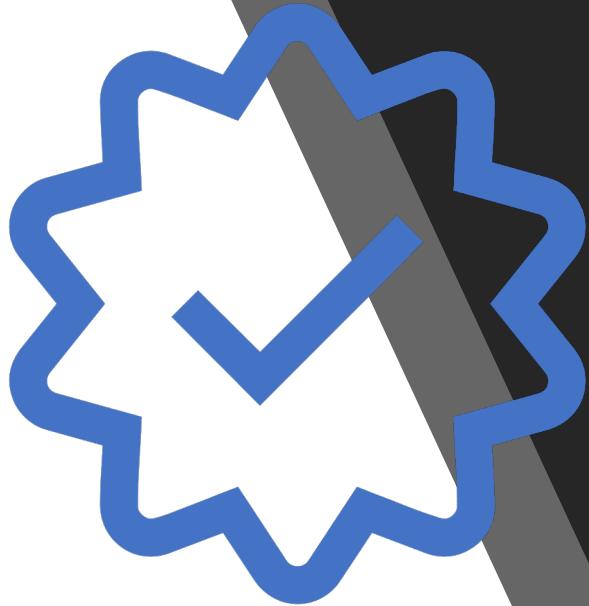


The
Center of
**Applied
Data Science**

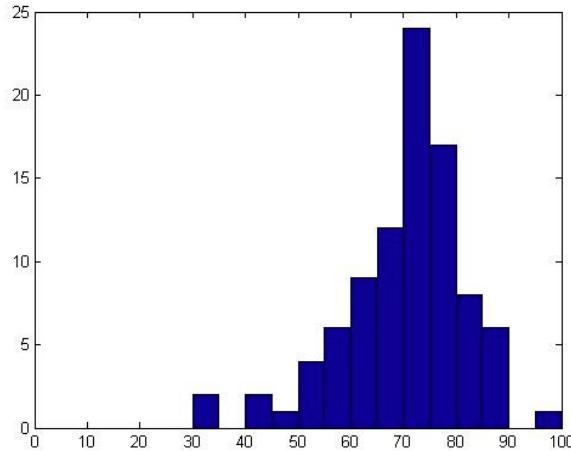


Step1: Business Problem

Who, What, When, Where, Why



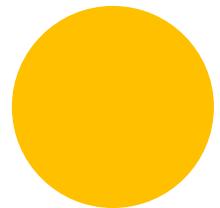
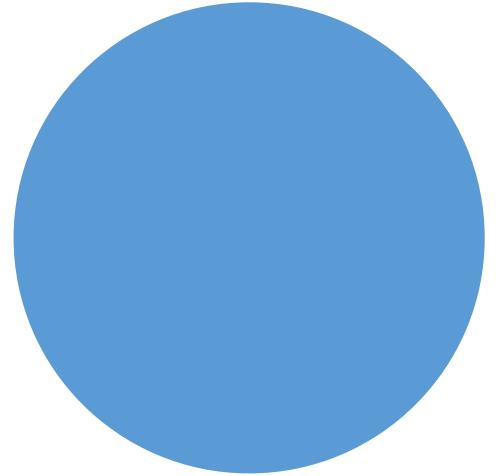
- Who has the best performance in our company?
- What products or services sell better?
- When should the company launch its next product?
- Where does the organization have big wins?
- Why are sales down despite many advertisement channels?



	A	B	C	D	E
1	StudentID	Homework	Midterm	Project	Final
2	4560	100	97	100	95
3	5540	85	90	88	90
4	6889	92	85	88	87
5	6817	65	85	87	89

Step2: Data Preparation

Import, Clean, Explore, and Visualize Data



Step3: Build and Evaluate Machine Learning Model

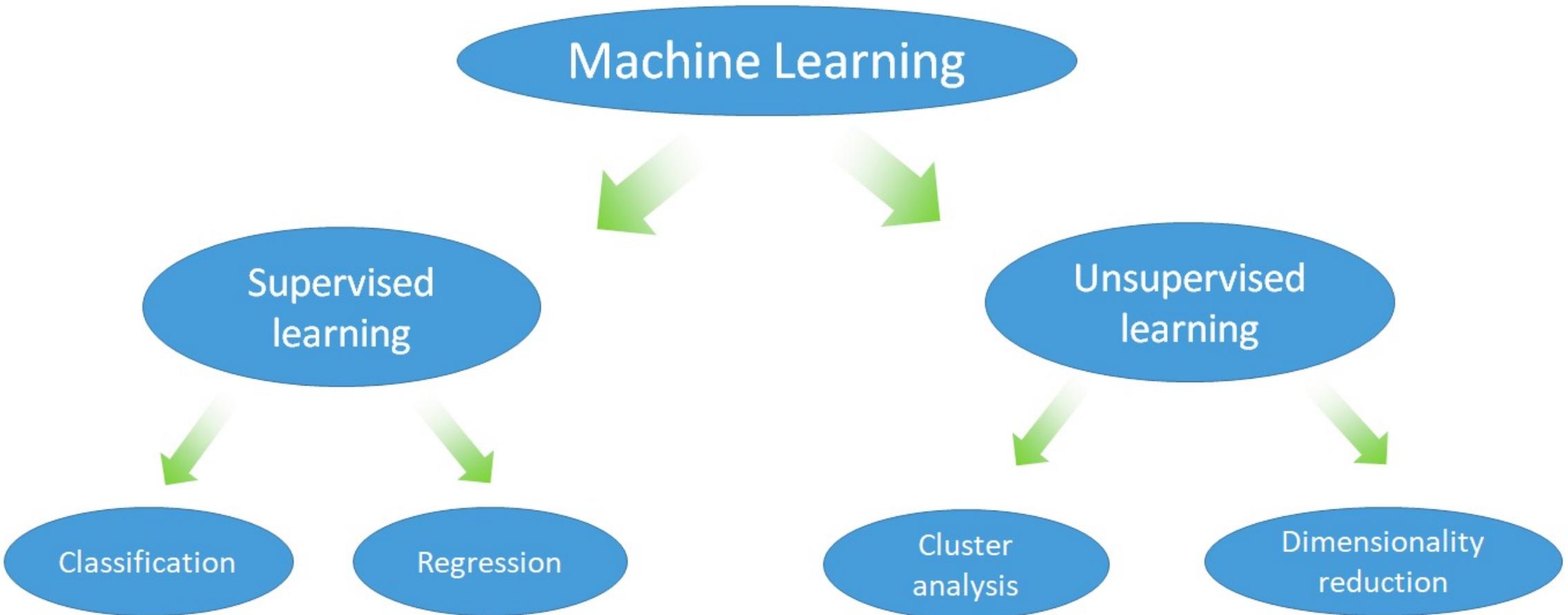
Step4: Interpret the results regarding the question

Machine Learning



Machine learning is the study of using a computer to automatically find patterns in data

Obvious patterns are already identified by domain experts, so the real advantage of machine learning is to extract non-intuitive patterns



Terminology

- The following terms are used interchangeably:
 - **Row**, observation, datapoint
 - Column, **attribute**, feature
 - **Target attribute**, target, label
- A target attribute is one we are interested in predicting
- To make this prediction we build a model
 - Models are the fundamental building blocks of data science

The diagram illustrates a data table with five columns: Name, Balance, Age, Employed, and Write-off. The 'Name' column is highlighted in blue. A bracket above the first four columns is labeled 'Attributes'. An arrow points from the text 'This is one row' to the fifth column, 'Write-off', which is labeled 'Target attribute' with a downward arrow.

Name	Balance	Age	Employed	Write-off
Mike	\$200,000	42	no	yes
Mary	\$35,000	33	yes	no
Claudio	\$115,000	40	no	no
Robert	\$29,000	23	yes	yes

This is one row

Two stages of modeling
Fit and predict

- Fit
 - Using available data (including targets), a set of rules is created
- Predict
 - For a new, unobserved datapoint (target unknown), the set of rules is applied

Returning to example

New row →

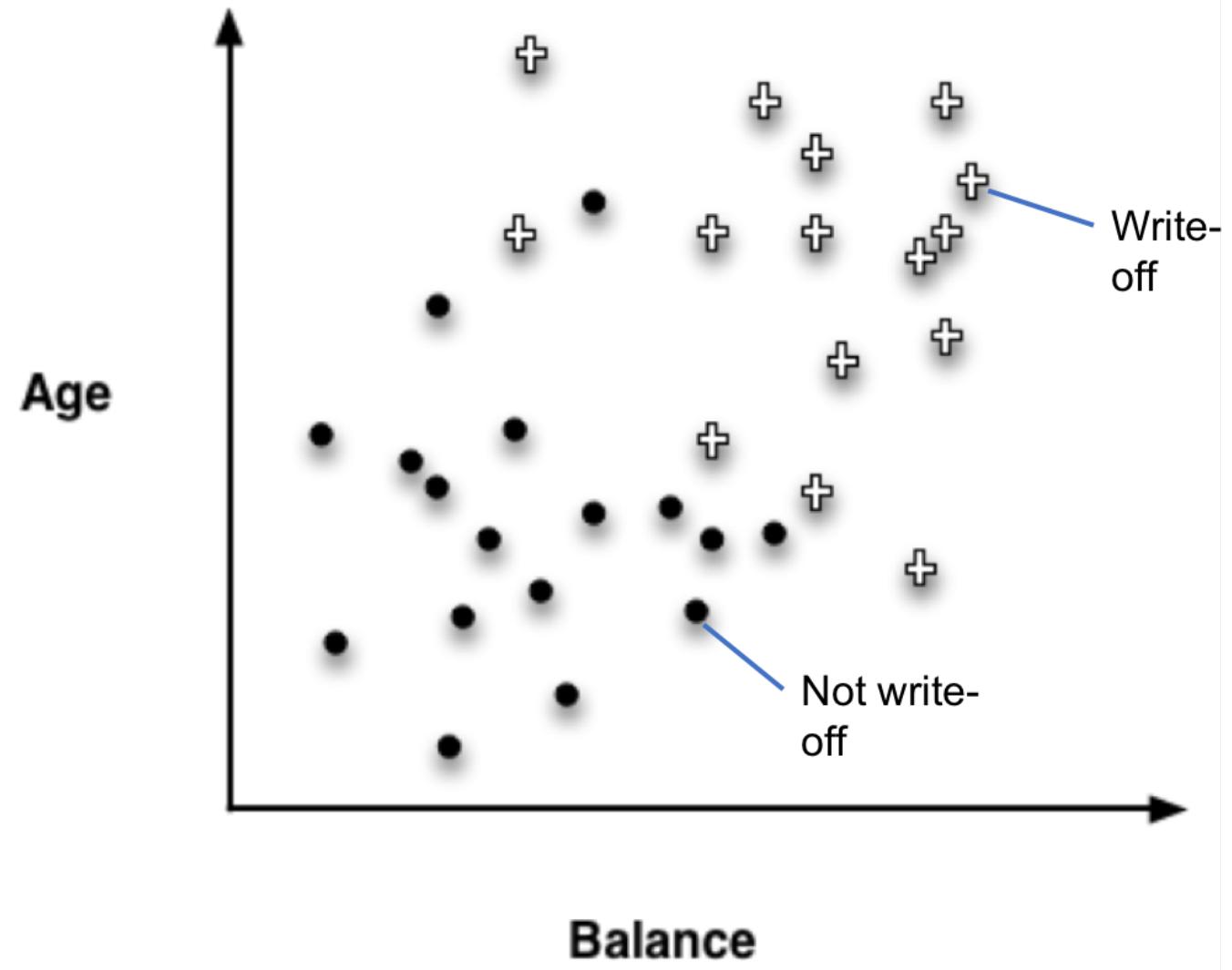
Name	Balance	Age	Employed	Write-off
Mike	\$200,000	42	no	yes
Mary	\$35,000	33	yes	no
Claudio	\$115,000	40	no	no
Robert	\$29,000	23	yes	yes
Dora	\$72,000	31	no	?

We would like for a model to predict this?

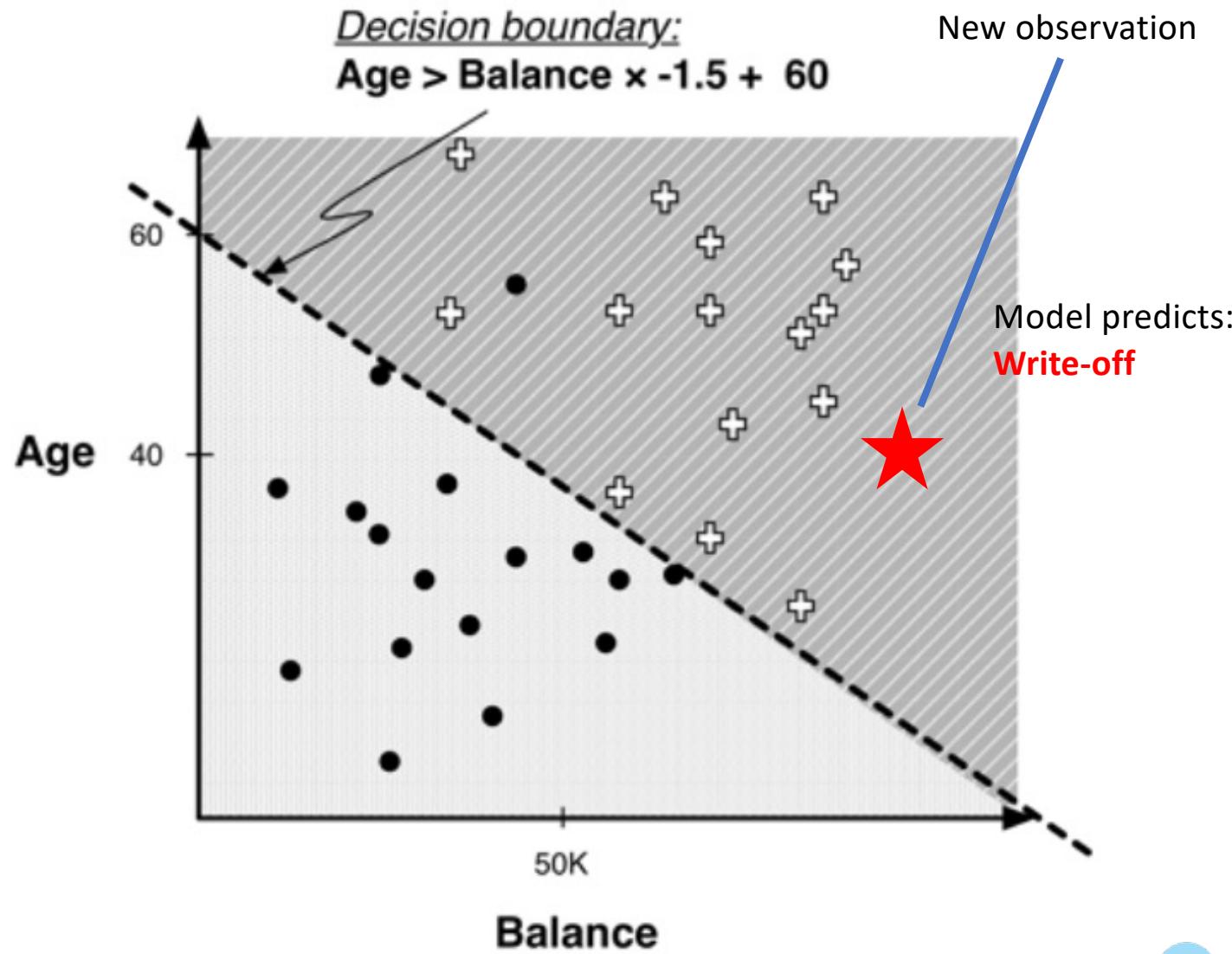
Classification

- Classification is used when the target attribute is a category
- Model is built based on target attributes of historical data
- Examples:
 - Predicting churn
 - Predicting write-off

Visualization of classification algorithms



Visualization of classification algorithms: Linear classifier

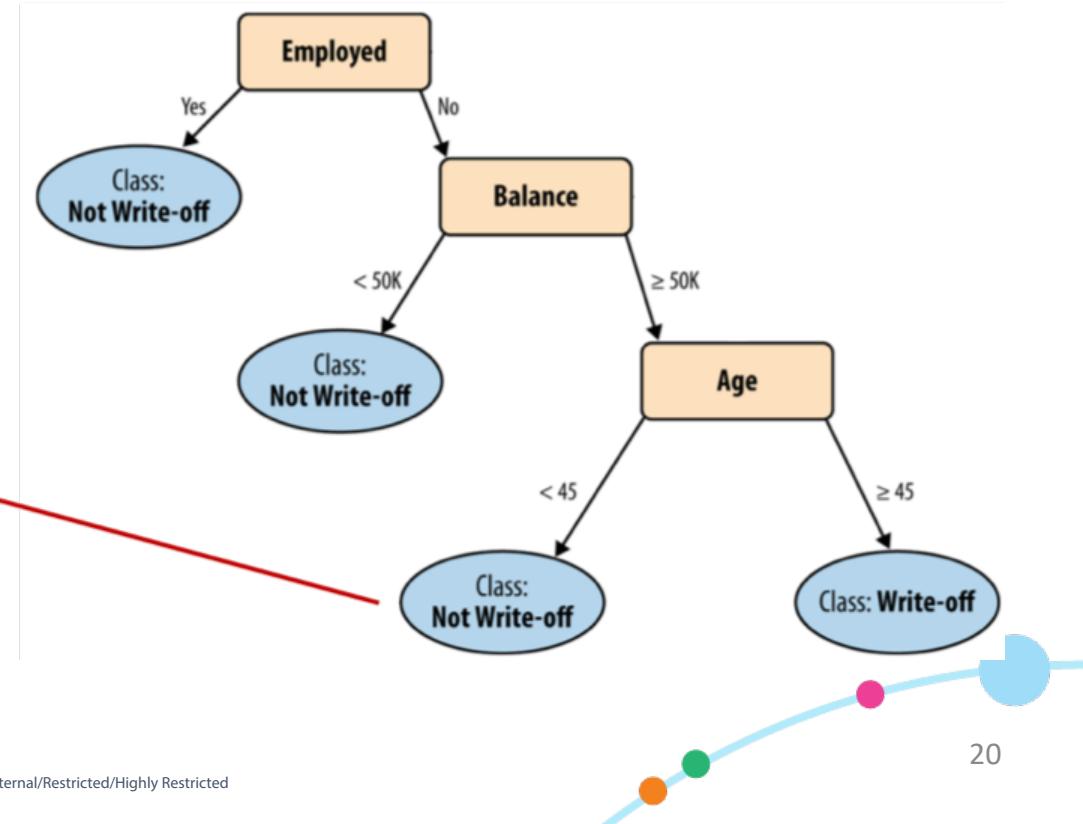


Decision Tree

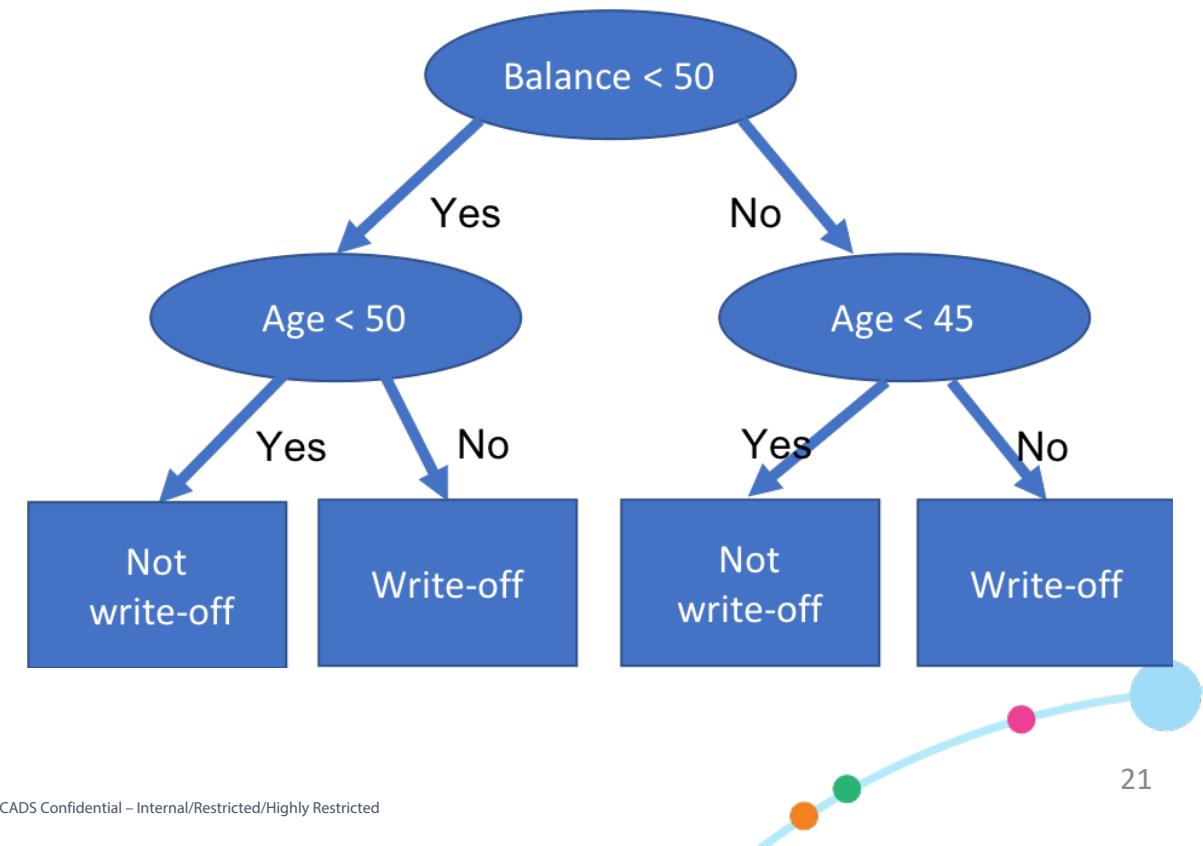
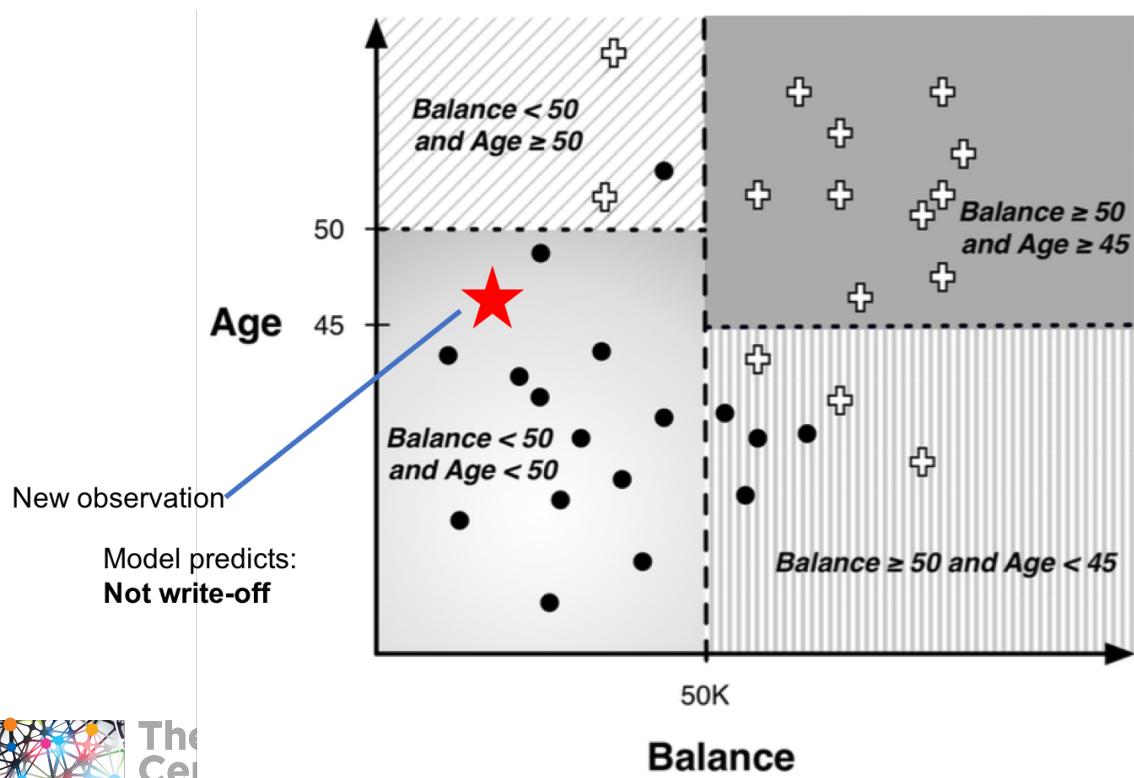
One type of model

Name	Balance	Age	Employed	Write-off
Mike	\$200,000	46	No	Yes
Mary	\$335000	33	Yes	Yes
Claudio	\$115,000	40	No	No
Robert	\$29,000	23	Yes	Yes
Dora	\$72,000	31	No	?

No

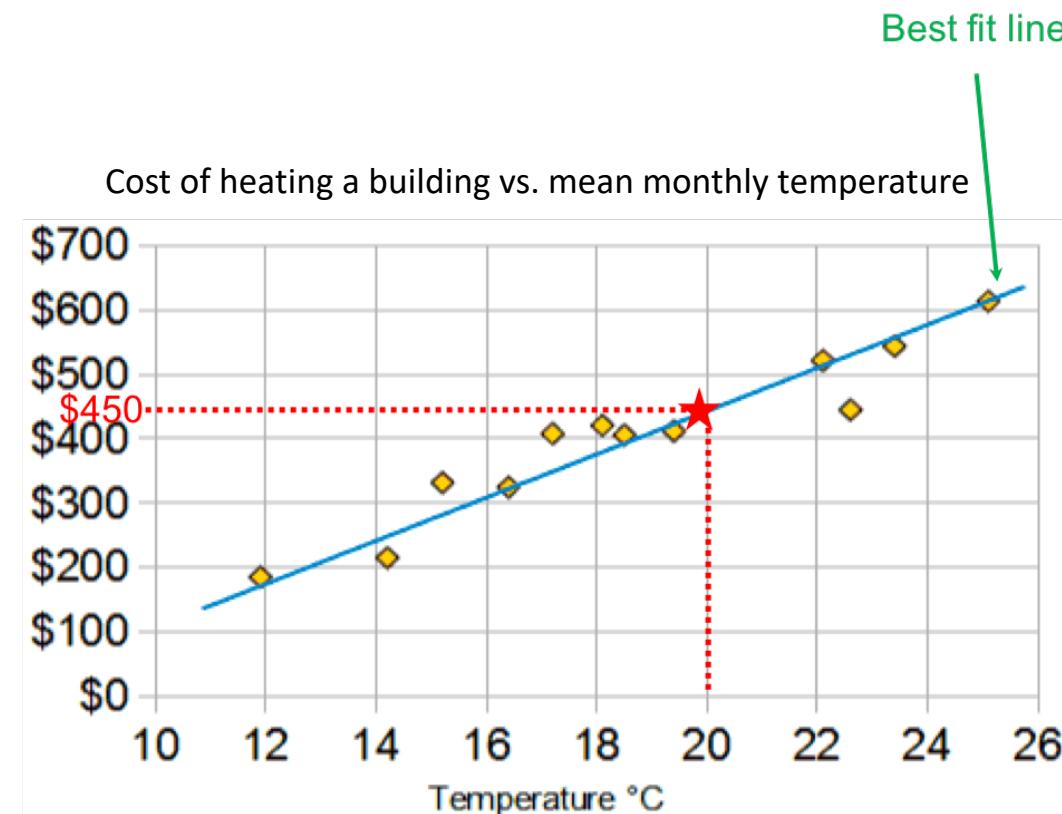


Visualization of classification algorithms: Decision tree



Regression

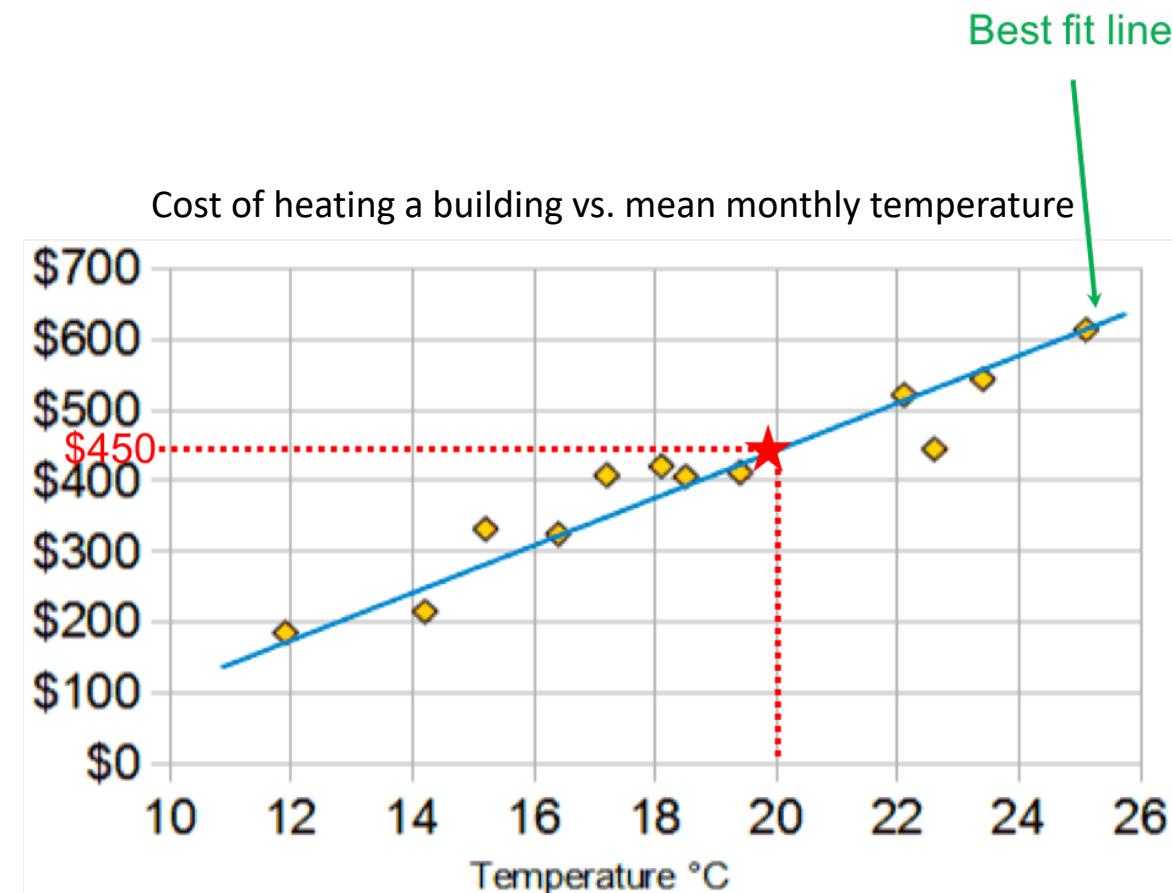
Temperature (C)	Cost
12	192
14.5	205
15.2	325
16.7	320
...	...
20	?



Question: What would be the cost for a temperature of 20°C ?

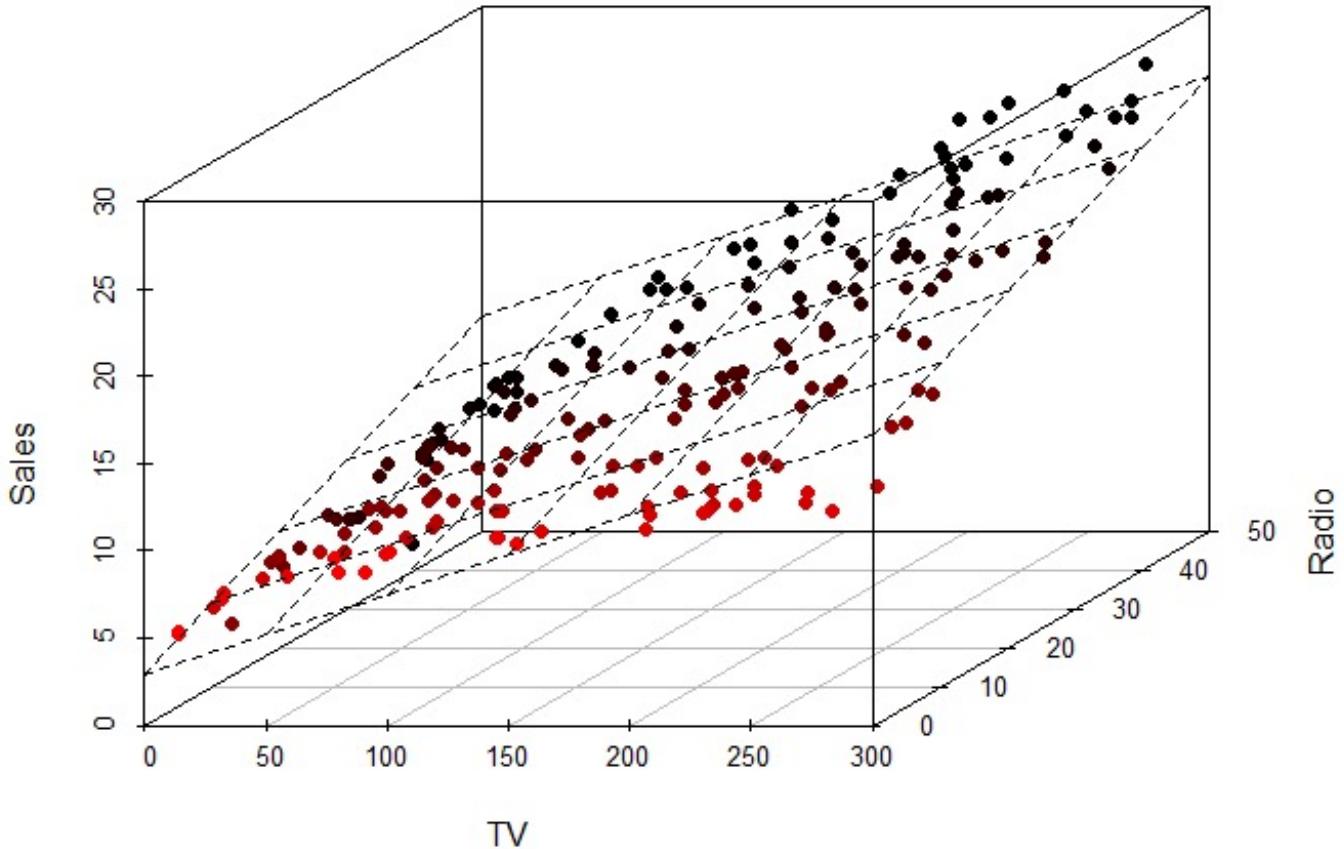
Regression

Regression predicts a target attribute which is a **real number**



Question: What would be the cost for a temperature of 20°C ?

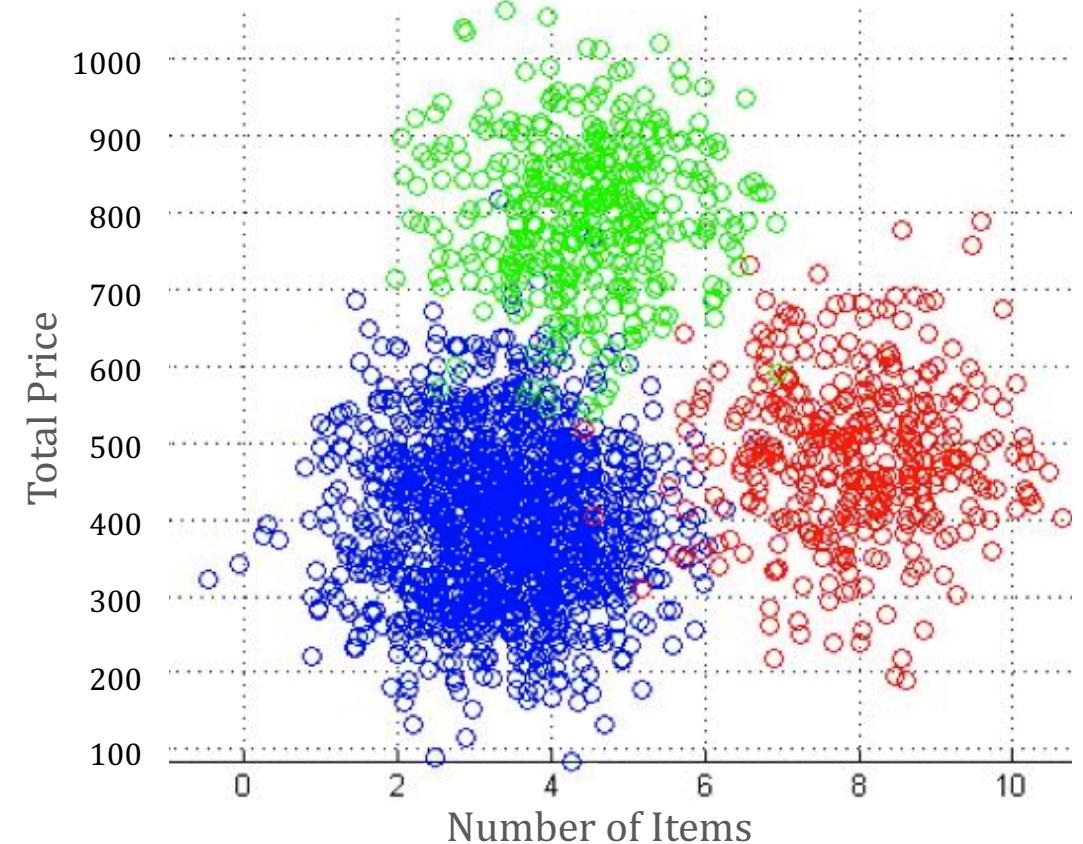
Regression



Clustering

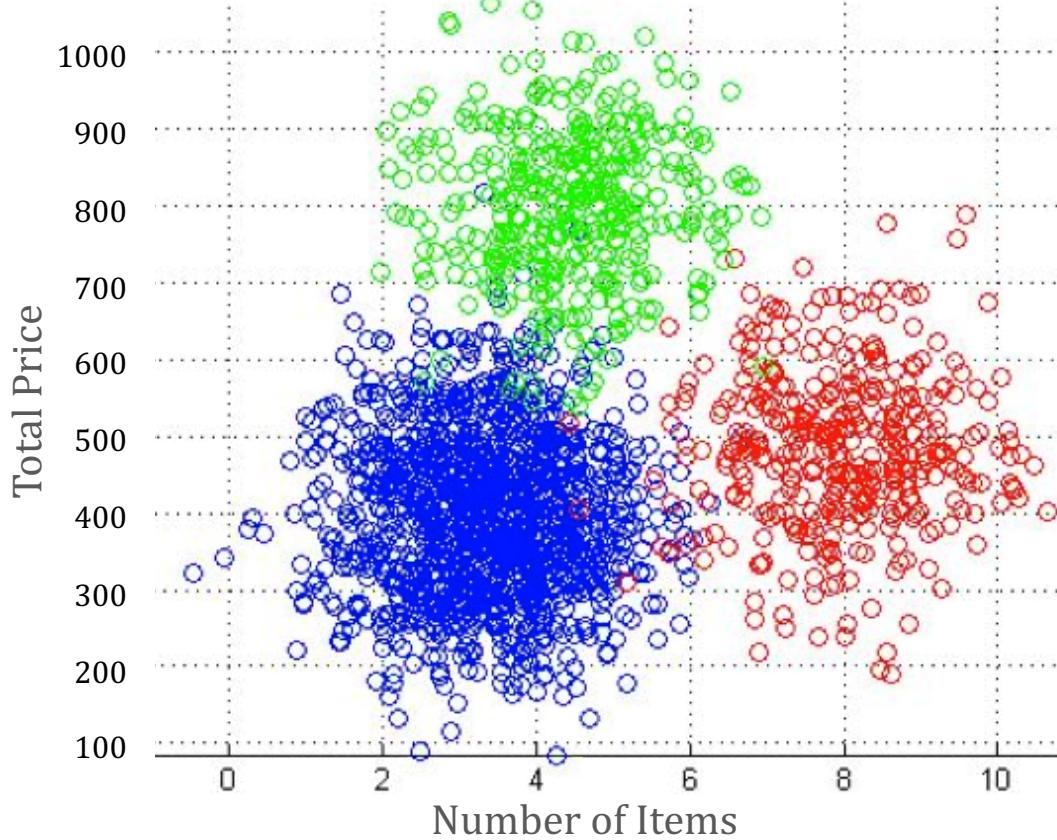
Number of Items	Total Price
2	50
4	1500
9	199
8	550
...	...

Question: How many clusters of buying behavior you observe in the data?



Clustering

- Clustering is similar to classification except that there are **no category labels in the dataset**
- Instead, we are looking for a natural grouping of the datapoints



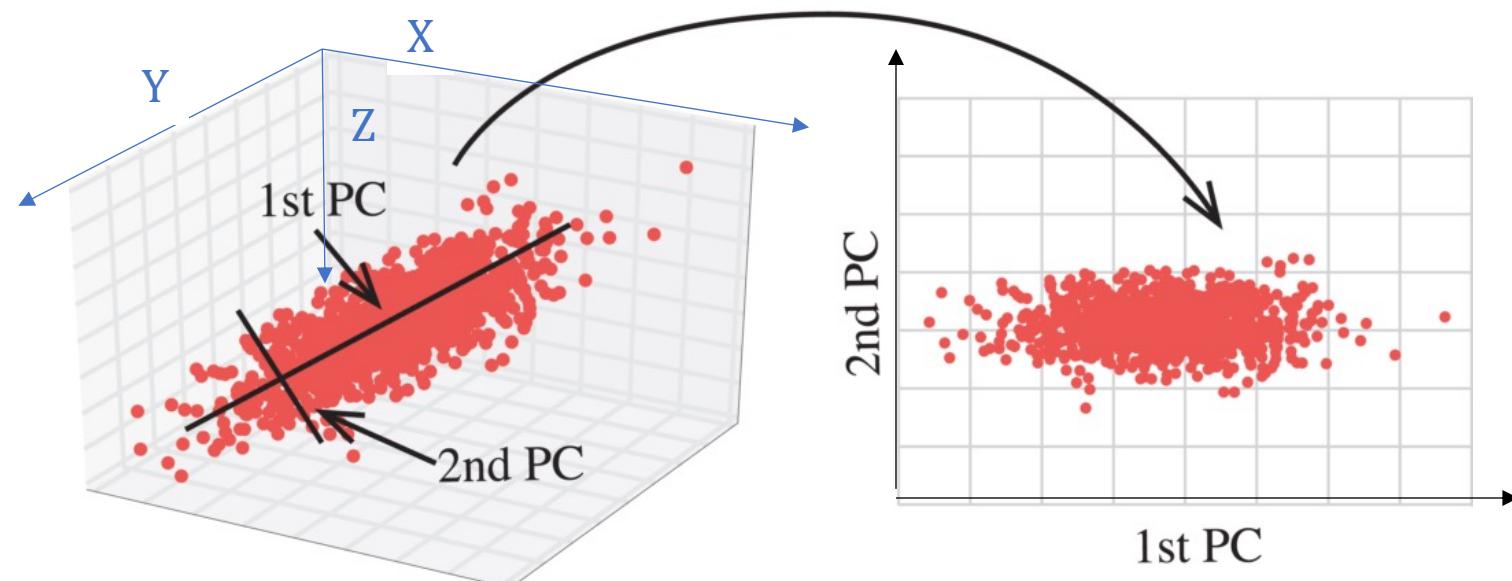
Dimensionality reduction

3 Dimension

X	Y	Z
0.5	-1	2.3
1	2.7	-3.3
0.2	1.2	3.8
...

2 Dimension

PC1	PC2
?	-?
?	?
?	?
...	...



3 Dimension

2 Dimension

Dimensionality reduction

- In some datasets, there are **hundreds or thousands of attributes** (columns)
 - **Example:** Netflix gives predictions for each user next based on their ratings of previous movies. In this case, every single movie in the database is an attribute
- A technique for **reducing the number of attributes** in a dataset
- Commonly used as a **preprocessing step**
- Collapses the attributes into new “**latent dimensions**”
- **No guarantee** that the reduced dimensions are comprehensible



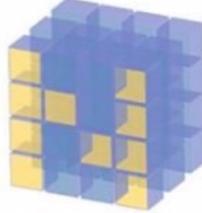
Open source project
containing a good number of
state-of-the-art ML algorithms

Good documentation, active
user community

Works well with a number of
other libraries (numpy,
pandas)



matplotlib

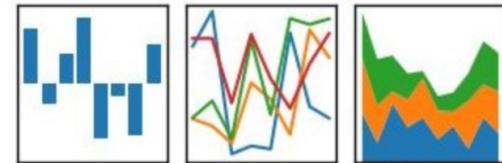


NumPy



pandas

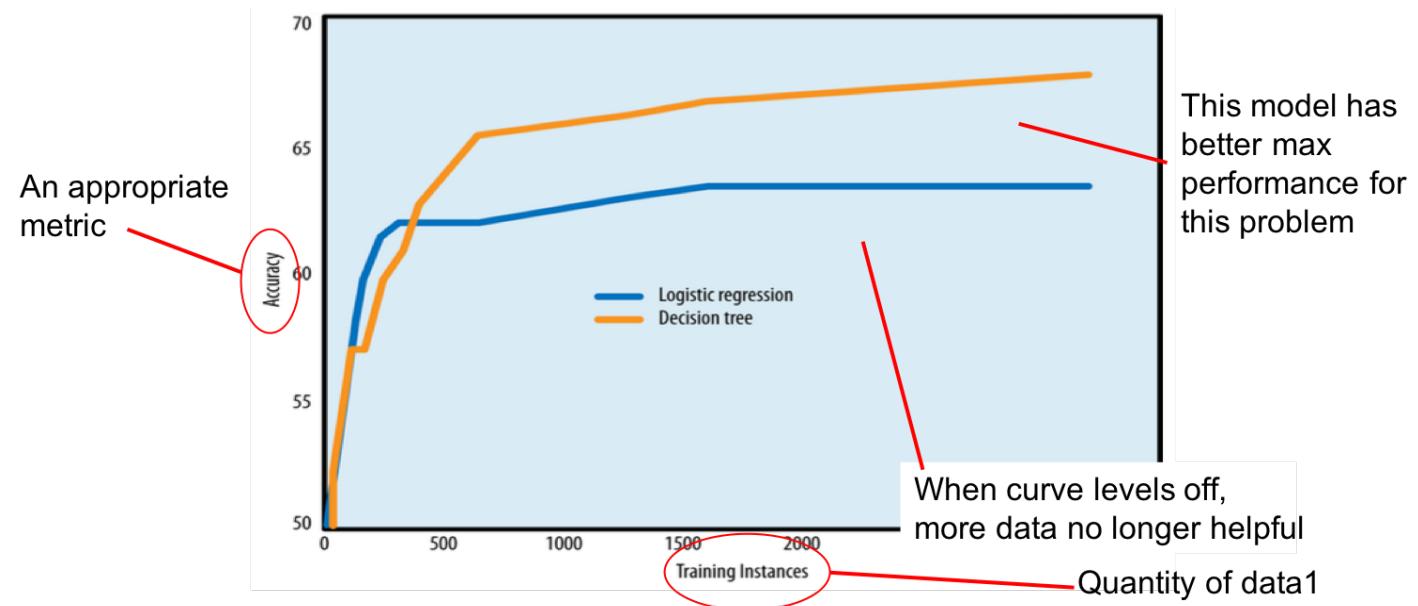
$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Other libraries and tools

Evaluating models

Example: The learning curve



The learning curve helps answer two questions

- ❖ Which model is better for this problem?
- ❖ Do we have enough data?

Thank You



The
Center of
**Applied
Data Science**

