

Data Science Workshop

Data Science Fundamentals

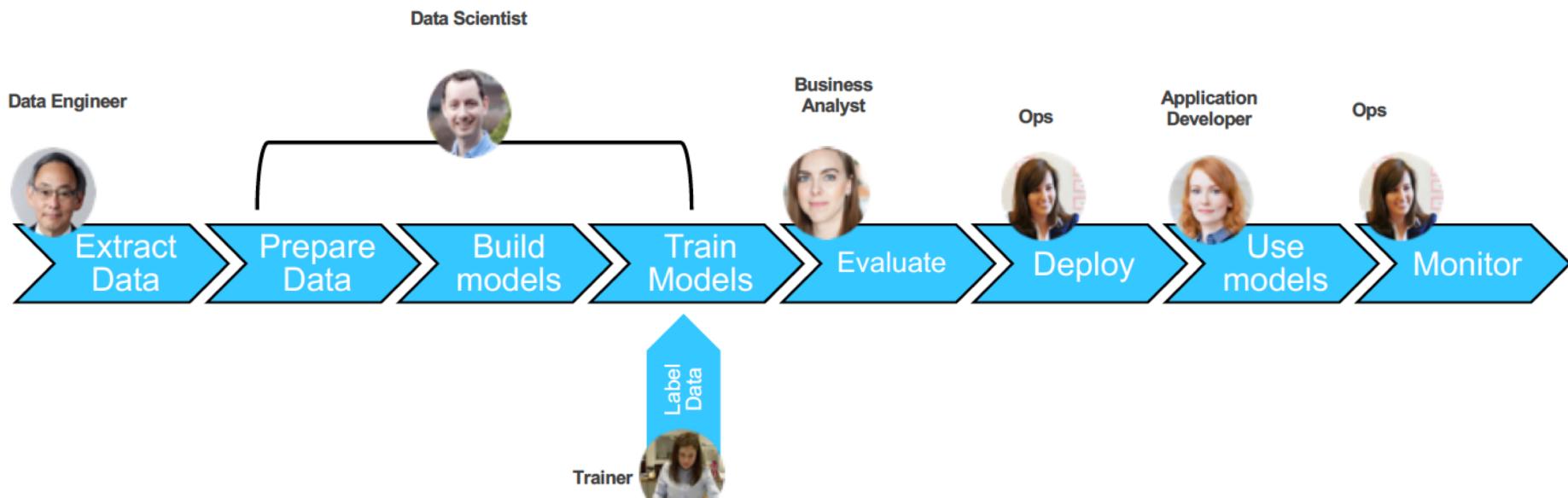
Toulon, February 2019

Emmanuel Génard – genard@fr.ibm.com
Data Scientist & Cloud Developer Advocate Europe,
IBM Business Solution Center Nice, France

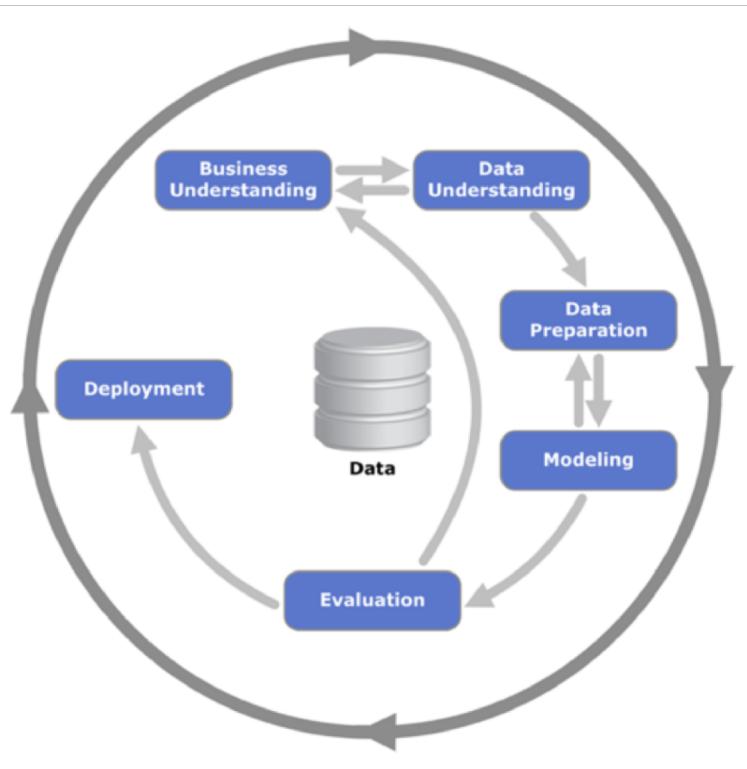
Data Science is a Team Sport

Building ML-infused apps requires multiple skillsets:

- Define an ML model
- Store, manage, update training data
- Manage lifecycle of the trained model
- Ability to do inferencing on the trained model(s)



The Data Science methodology consisted to cover the following steps



- A known Data Science methodology (Cross Industry Standard Process for Data Mining)
 - Criteria of success (Business, Technical, Change Management, ...)
 - A collaborative work between the CUSTOMER and the Data Science Expert Team, which secures the success
 - A Data Science platform

The Data Preparation step



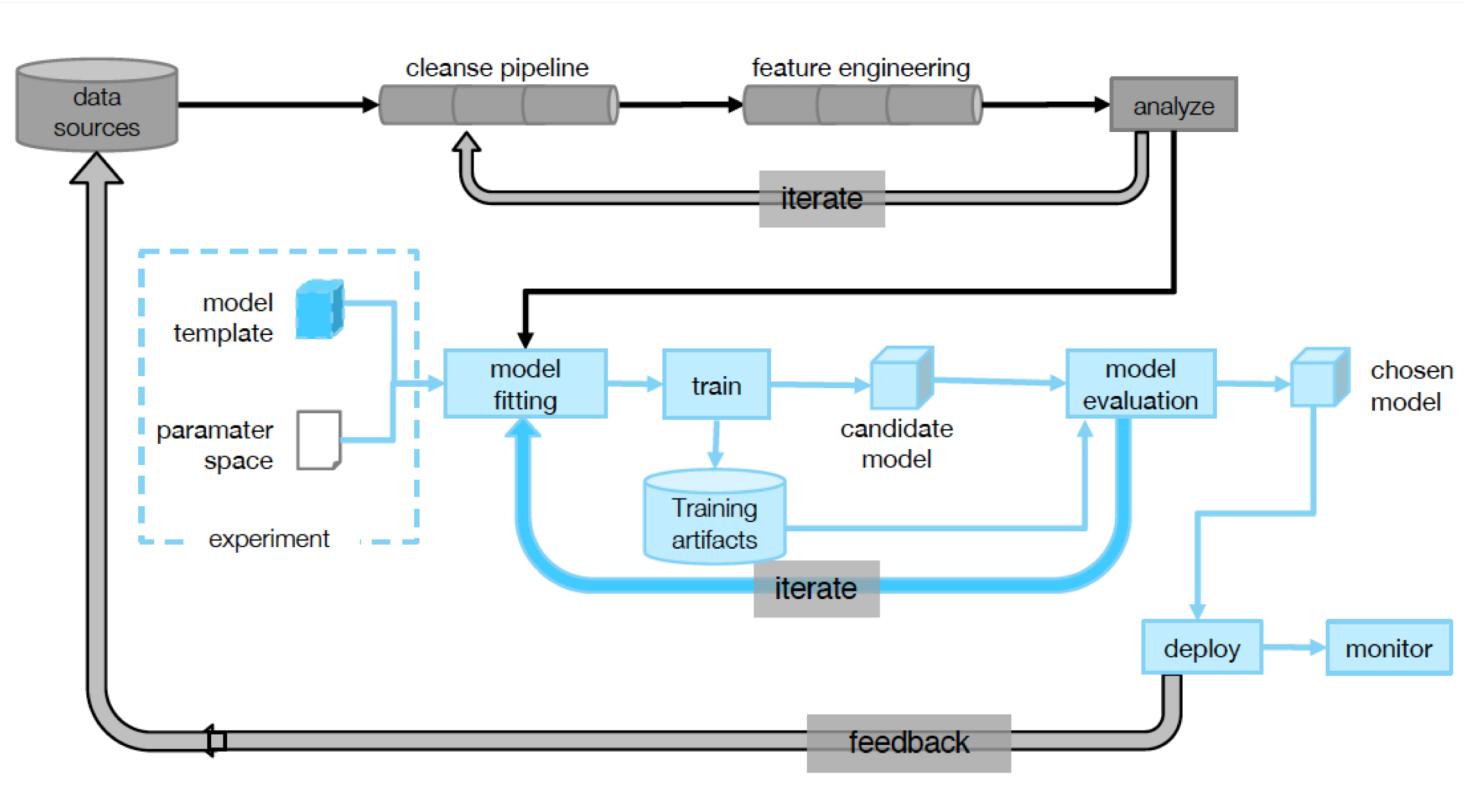
- Technical Data Cleaning (Types of variables, ...)
- Statistical Data Cleaning
- Creation of Variables/Indicators
- Selection
- Building of the Data Foundation

The Modelling step: a robust approach that helped to ensure to work in a virtuous cycle

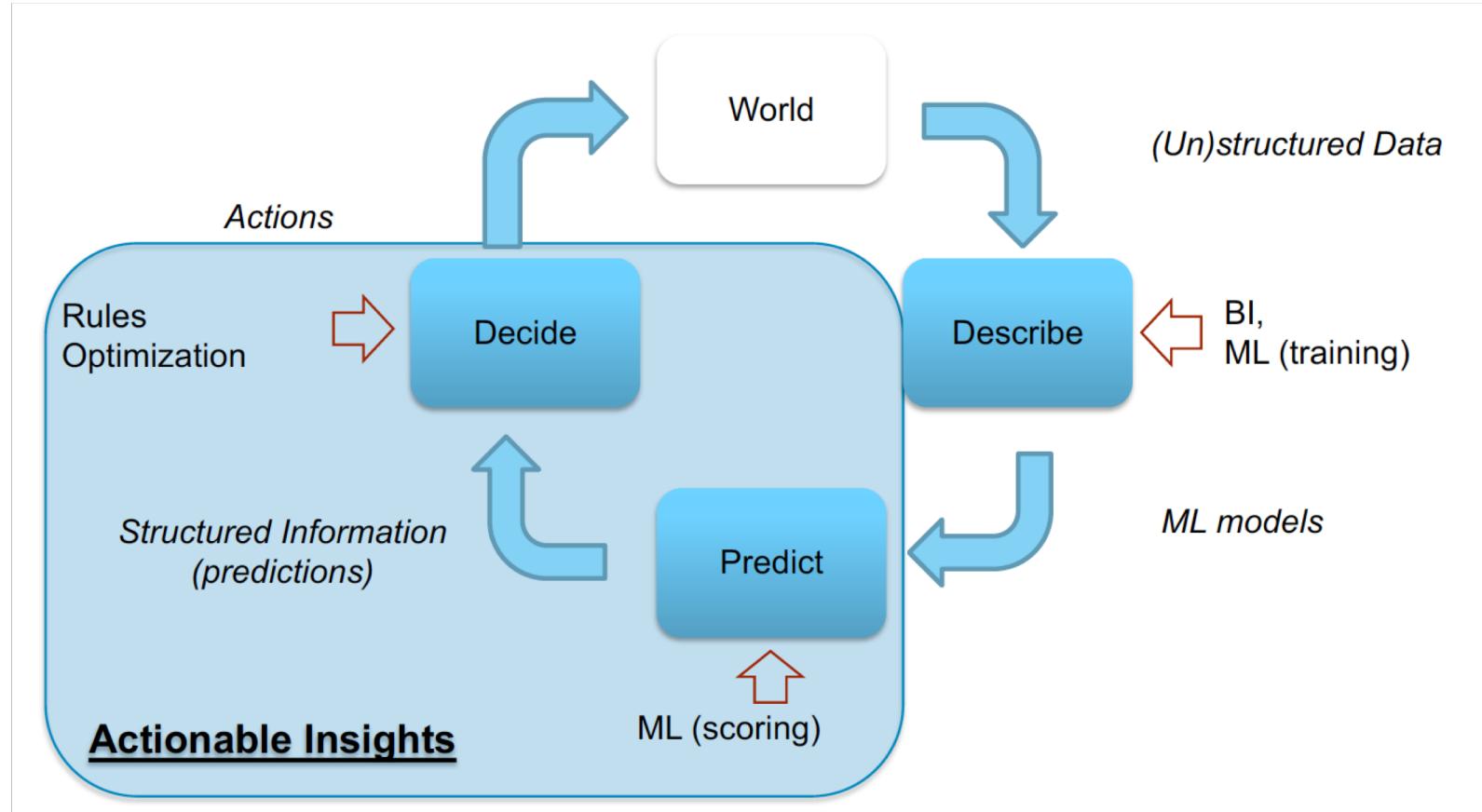


- The measurement (continuous or categorical target)
- A cross validation technique in order to validate the best algorithm and to minimize over-fitting
- An auto-modelling: that helped to know the best algorithm to use

Machine Learning architecture



Machine Learning and Decision Making



Measurement (numerous target) – RMSLE

- **RSS, for Residual Sum of Squares.** Deviations predicted from actual empirical values of data.
- **MSE, for Mean Squared Error.** The RSS is generally normalized (by the number of observations) to avoid to have a huge number (in case there is a big number of Observations).
- **RMSE, for Root Mean Squared Error.** Squared error to have the same unit than $y=f(x_i)$.
- **RMSLE for Root Mean Squared Log Error.** In case the empirical values are on a large scale, it is necessary to do a logarithmic transformation. (a error of 10 units on the value of 4 is not the same than on a value of 100!)

$$RSS = \sum_{i=1}^n (f(x_i) - y_i)^2$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2}$$

$$RMSLE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\log(f(x_i) + 1) - \log(y_i + 1))^2}$$

Measurement (numerous target) – MAE

- **Prediction Error = Actual Value - Predicted Value.** subtraction of Predicted value from Actual Value
- **Absolute Error → |Prediction Error|**
- **MAE, for Mean Absolute Error.** Mean for all recorded absolute errors (Average sum of all absolute errors). Refers to the measurement of the difference between two continuous variables.

$$mae = \frac{\sum_{i=1}^n \text{abs}(y_i - \lambda(x_i))}{n}$$

- **MAE with the logarithm**

Prediction Error = $\log(1+\text{Actual Value}) - \log(1+\text{Predicted Value})$

Measurement (continuous target) – Confusion Matrix

- In the field of machine learning and specifically the problem of statistical classification , A **confusion matrix**, also known as an error matrix, is a specific table layout that allows visualization of the performance of an algorithm. It is a table with two rows and two columns that reports the number of false positives, false negatives, true positives, and true negatives.
- The pattern built from all machines data provided the following Confusion matrix:
 - **True positives (TP)**: These are cases in which we predicted Yes (1) (the failure occurs), and the failure occurred.
 - **True negatives (TN)**: We predicted No (0), and the failure did not occur.
 - **False positives (FP)**: We predicted Yes (1), but the failure did not occur
 - **False negatives (FN)**: We predicted No (0), but the failure occurred

Measurement (continuous target) – Confusion Matrix

		True condition			
		Condition positive	Condition negative	Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Predicted condition	Predicted condition positive	True positive, Power	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$
		False negative rate (FNR), Miss rate $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$	

Performance in Machine Learning

Overfitting and Underfitting with Machine Learning Algorithms

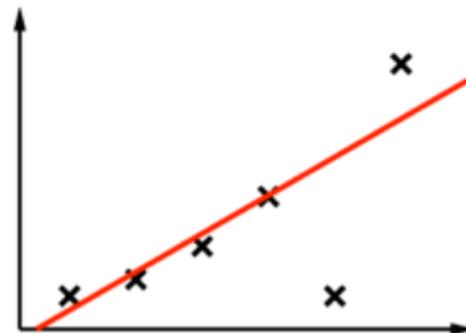
- **Supervised machine learning** is best understood as approximating a target function (f) that maps input variables (X) to an output variable (Y) ($Y = f(X)$)
- An important consideration in learning the target function from the training data is **how well the model generalizes to new data**. Generalization is important because the data we collect is only a sample, it is incomplete and noisy. Generalization refers to how well the concepts learned by a machine learning model apply to specific examples not seen by the model when it was learning.
- **The goal of a good machine learning model is to generalize well from the training data** to any data from the problem domain. This allows us to make predictions in the future on data the model has never seen.
- The cause of **poor performance** in machine learning is either **overfitting** or **underfitting** the data.
 - **Overfitting** refers to a model that models the training data too well (nonparametric and nonlinear models)
 - **Underfitting** refers to a model that can neither model the training data nor generalize to new data (obvious to detect – no discussion)

Example 1: height and weight of people

- Used for prediction: Linear regression
- Examples with only two variables: height and weight of people
 - We want to learn how to predict weight as a function of height



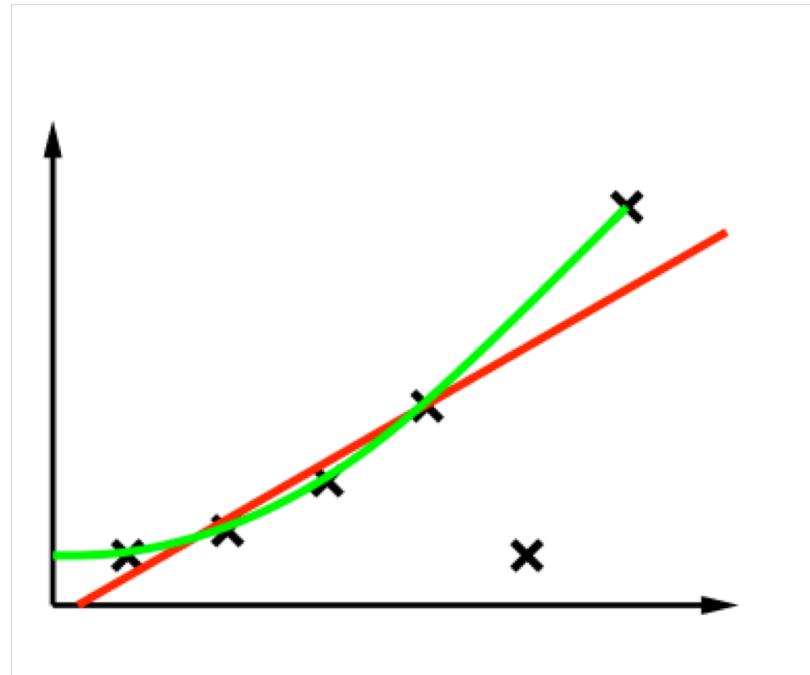
Train



Example 1: height and weight of people

Feature Engineering

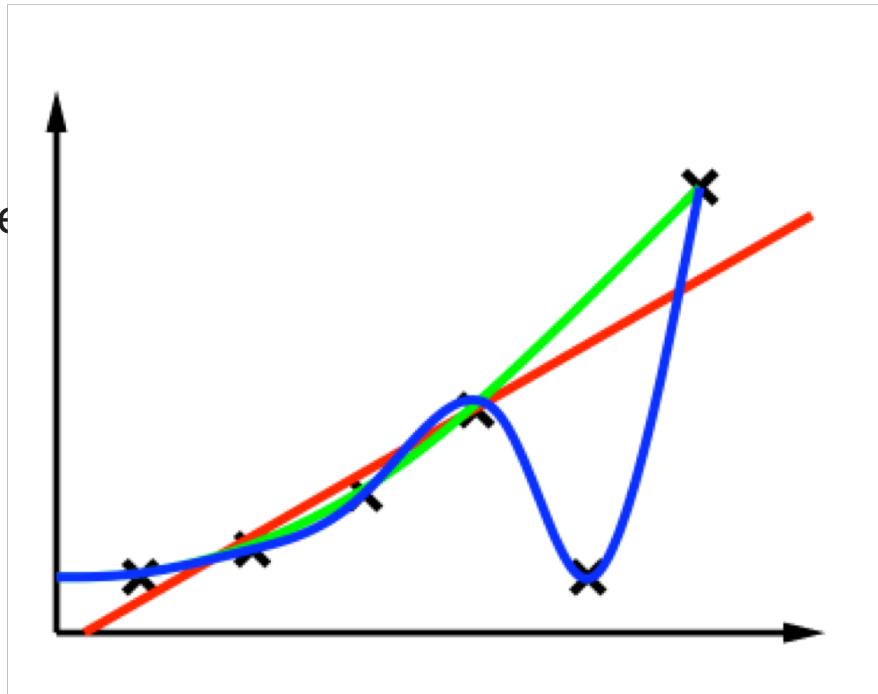
- Consider square of height
- Ignore outliers



Example 1: height and weight of people

Model Training

- Balance two goals
 - Fit train data correctly, i.e. lead to good predictions on train data
 - Have a model as simple as possible to avoid overfitting.
- Overfitting can occur with regression, see blue curve below



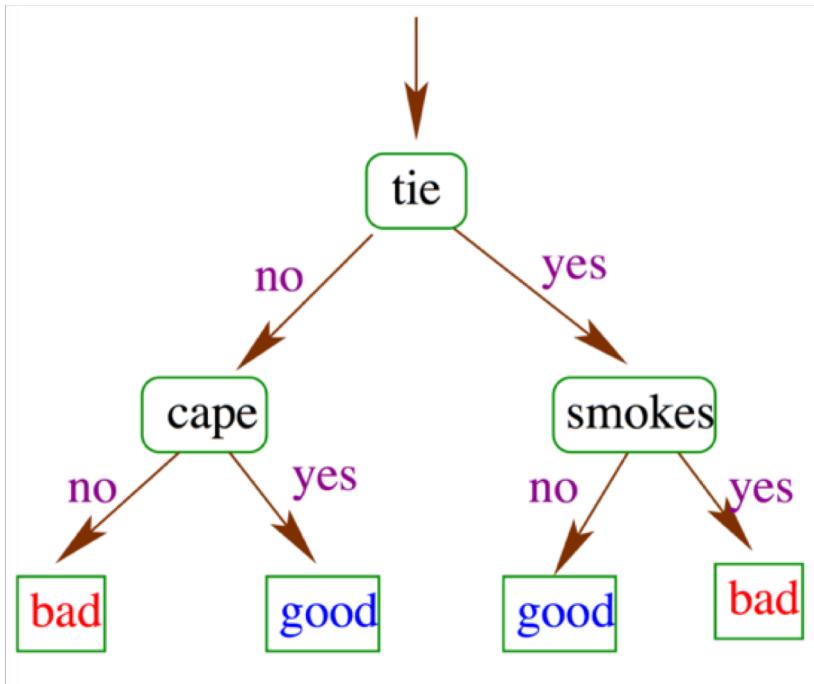
Example 2: identify people as good or bad from their appearance

- Used for prediction: Logistic regression, Decision Tree, Support Vector Machines, ...

	sex	mask	cape	tie	ears	smokes	class
<u>training data</u>							
batman	male	yes	yes	no	yes	no	Good
robin	male	yes	yes	no	no	no	Good
alfred	male	no	no	yes	no	no	Good
penguin	male	no	no	yes	no	yes	Bad
catwoman	female	yes	no	no	yes	no	Bad
joker	male	no	no	no	no	no	Bad
<u>test data</u>							
batgirl	female	yes	yes	no	yes	no	??
riddler	male	yes	no	no	no	no	??

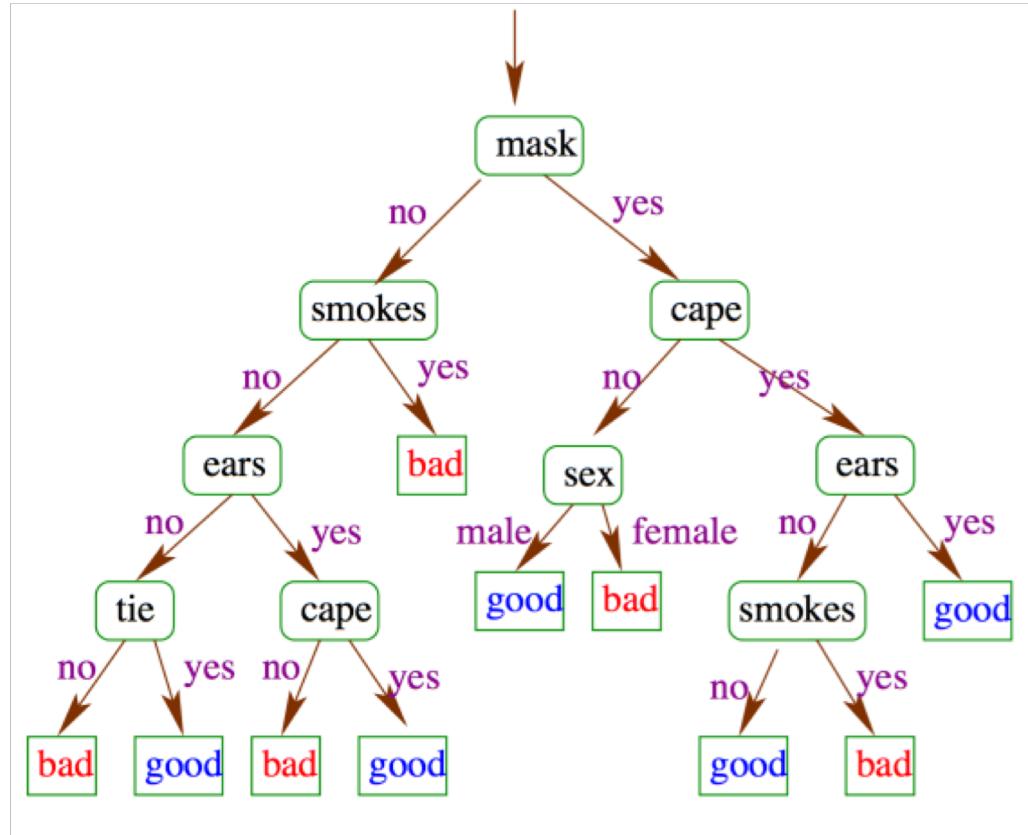
Example 2: identify people as good or bad from their appearance

- A series of tests
- Class assigned at the leaves of the tree



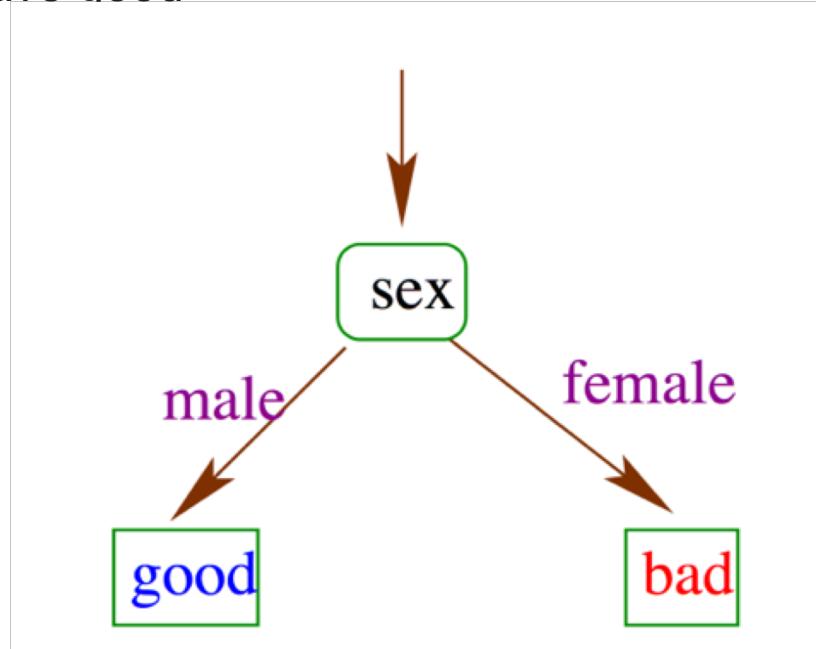
Example 2: identify people as good or bad from their appearance

- Classifies test data perfectly
- Performs poorly on new data
- Why?
 - Almost rote learning of train data
- We must favor simpler models
 - Trees with less nodes



Example 2: Underfitting

- Too simple models do not have good predictive power either



Thank You



Useful Links & Resources

External

Getting Started:

[Service Homepage](#)
[Feature Requests / Suggestions](#)

Case Studies:

[OmniEarth](#)
[Aerialtronics](#)
[BlueChasm](#)
[iTrend](#)

Tutorials & Best Practices:

[Training models with Watson Studio](#)
[Getting started with Watson + Core ML](#)
[Stacking Multiple Custom Models](#)
[Create a Calorie Counting App](#)
[Watson Visual Recognition & Twilio](#)
[Best Practices for Custom Models](#)

Code Patterns:

[Classify vehicle damage](#)
[Analyze industrial equipment for defects](#)
[Create an Android calorie-counter app](#)

External continued

Books:

[Redbook: Building Cognitive Application using IBM Watson Services vol3 – Watson Visual Recognition](#)

Blogs:

[IBM Watson on Medium](#)

Internal

[Slack Channel: #ibmvisual-recognition](#)
[Service Roadmap](#)
[IBMer key limit increase request form](#)
[ZACS portal](#)
[Digital Sales Play](#)
[Content Request & Feedback Form](#)