

# Machine Learning for Customer Churn with Watson Studio

Darrel Pyle, Chris Tyler

Technical Evangelist Team



# Introductions



Darrel Pyle  
Technical Evangelist

 [@AnalyticsDS](https://twitter.com/AnalyticsDS)

 [AnalyticsDS](https://www.linkedin.com/company/AnalyticsDS)



 [@IBMWolfPack](https://twitter.com/IBMWolfPack)

 [team-wolfpack](https://github.com/team-wolfpack)



Chris Tyler  
Technical Evangelist

 [@chrisatylor](https://twitter.com/chrisatylor)

 [chrisatylor](https://www.linkedin.com/in/chrisatylor)

# Session Overview

## Description

The goal of this session is to **familiarize participants with the Watson Data Platform; specifically the Watson Studio and Watson Machine Learning**. This will be one within the context of analyzing customer churn.

## Audience

Intended for individuals seeking to develop a **basic understanding of data science** and machine learning.

## Pre-requisites

### Pre-requisite skills:

- Business Intelligence
- Conditioning and management of business data
- Familiarity with basic statistics

# Session Objectives

Upon completion of this session, you should be able to:

- Execute a notebook in the Watson Studio
- Deploy a Machine Learning Model
- Understand the Tools, Technology, and Processes involved in Data Science and Machine Learning

# Section 1

Introduction to Data Science

# Analytics

## A DEFINITION

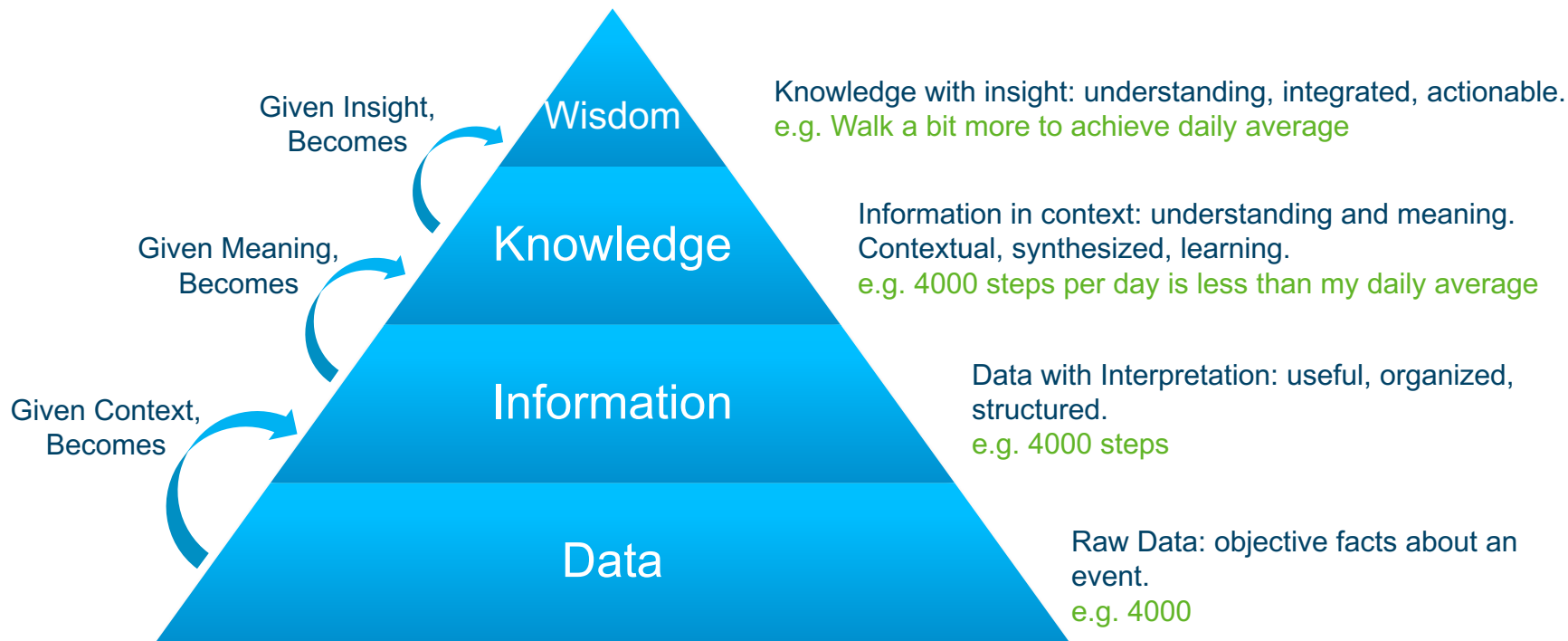
*“Analytics are the quantifiable **informational** inputs that use **past data** to identify possible trends that may provide **valuable insight** for future action.”* <sup>(2)</sup>



Pixabay. Analytics Word Cloud <sup>(1)</sup>

# Analytics

## DATA'S JOURNEY TO VALUE



Charles Sturt University. Data and Knowledge<sup>(1)</sup>

# Analytics

## DATA STRUCTURE TYPES

### Structured

Defined data type, format, and structure

Transactional Data

### Semi-Structured

Textual data with a discernable pattern, enabling parsing

Self describing XML with schema

### Quasi-Structured

Textual data with erratic data formats, can be formatted with effort, tools, and time

Clickstream data

### Unstructured

Data that has no inherent structure and is usually stored in different file types

PDF, Excel, JPG



# Data Science

## A DEFINITION

*“**Data science**, also known as data-driven science, is an interdisciplinary field about scientific methods, processes, and systems to **extract knowledge or insights from data** in various forms, either structured or unstructured.”<sup>(1)</sup>*

- \* Term first used in publication in 1974
- \* Introduced as an independent discipline in 2001

# Data Scientist

## MODERN DAY UNICORNS

- Quantitative
  - Skilled in mathematics or statistics
- Curious & Creative
  - Passionate about finding creative ways to solve problems and portray information



Tveten, J. Data Science 101 <sup>(1)</sup>

- Skeptical
  - Must be able to examine their own work critically
- Technical
  - Aptitude for software engineering, programming, and machine learning
- Communicative & Collaborative
  - Strong verbal and written skills. Must be able to articulate business value and collaborate with others.

# Data Science Team



Data Engineer

Data ingestion pipelines



Data Scientist

Wrangling, exploring, and hacking data



Quantitative Analyst

R&D advanced mathematical algorithms



Data Analyst

Test hypothesis, creates data driven reports



Front-end Developer

Develop end-user applications

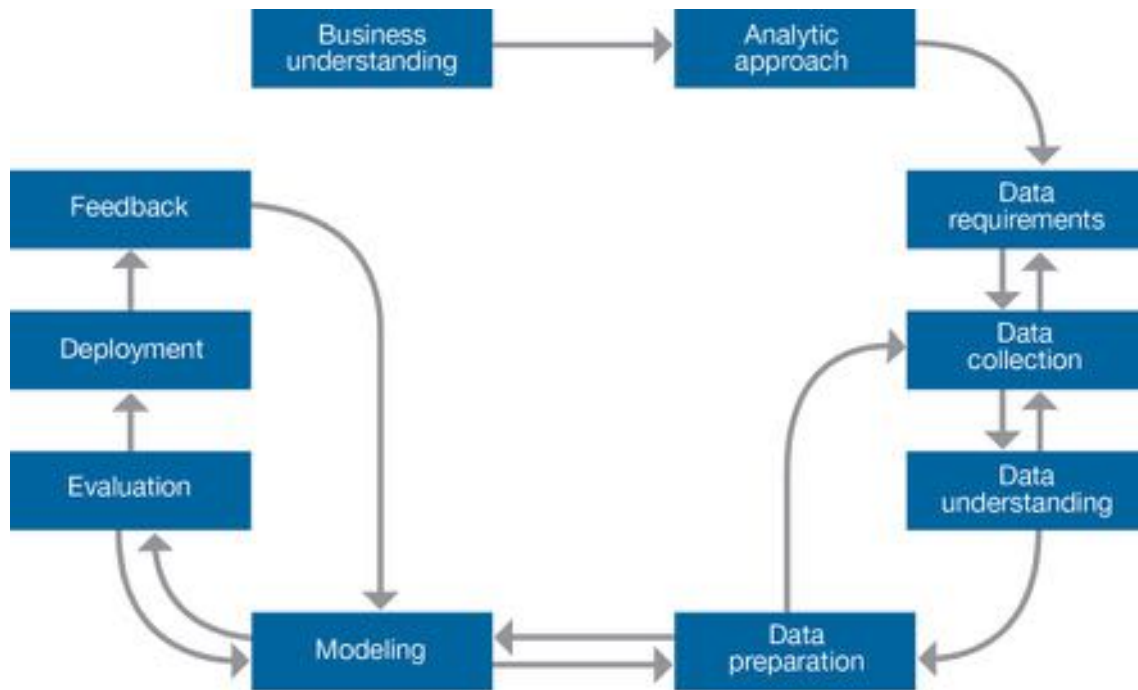


Business Expert

Identify business opportunities

# Data Science Methodology

## METHODOLOGY DIAGRAM



# Tools for Data Science

The image displays a large grid of data science tools, organized into several main categories:

- INFRASTRUCTURE**
  - HYBRID CLOUD PROVIDERS**: Includes CloudPact, Pivotal, and others.
  - HYBRID CLOUD PROVIDERS**: Includes Microsoft Azure, Google Cloud Platform, and others.
  - HYBRID CLOUD PROVIDERS**: Includes Amazon Web Services, IBM, and others.
- ANALYTICS**
  - DATA ANALYTICS PLATFORMS**: Includes Microsoft, Alteryx, and others.
  - DATA SCIENCE PLATFORMS**: Includes IBM, SAS, and others.
- APPLICATIONS - ENTERPRISE**
  - SALES**: Includes Salesforce, HubSpot, and others.
  - MARKETING**: Includes Marketo, Pardot, and others.
  - FINANCE**: Includes BlackRock, and others.
  - SECURITY**: Includes Palo Alto Networks, and others.
- APPLICATIONS - INDUSTRY**
  - MANUFACTURING**: Includes Siemens, and others.
  - TRANSPORTATION**: Includes Uber, and others.
  - HEALTHCARE**: Includes GE Healthcare, and others.
- DATA SOURCES**
  - DATA SOURCES**: Includes Twitter, and others.
  - DATA SOURCES**: Includes LinkedIn, and others.
  - DATA SOURCES**: Includes Facebook, and others.
- DATA RESOURCES**
  - DATA RESOURCES**: Includes Kaggle, and others.
  - DATA RESOURCES**: Includes GitHub, and others.
  - DATA RESOURCES**: Includes Stack Overflow, and others.

# Tools for Data Science

## Infrastructure

- Power Systems
- IBM Cloud

## Analytics

- Watson Studio
- SPSS
- Cognos

## Applications

- Watson Health
- Watson HR

## Cross-Infrastructure/Analytics

- Watson Data Platform
- IBM Cloud

## Data Sources & APIs

- Weather Company
- Watson Studio Community Data

## Open Source

- Spark
- Jupyter
- R
- Python

# Section 1: Summary

## Key points covered in this section:

- Relationship Between Data, Information, and Knowledge
- Key Characteristics of Data Science, and the Data Scientist
- Data Science Methodology
- Tools & Technology Used in Data Science by Data Scientists

# Section 2

Introduction to Machine Learning



# Overview of Machine Learning



Arthur Samuel demonstrating his Checkers program on the IBM 701 computer in 1956.

TERM FIRST COINED BY AN **IBMer** (circa 1950's)

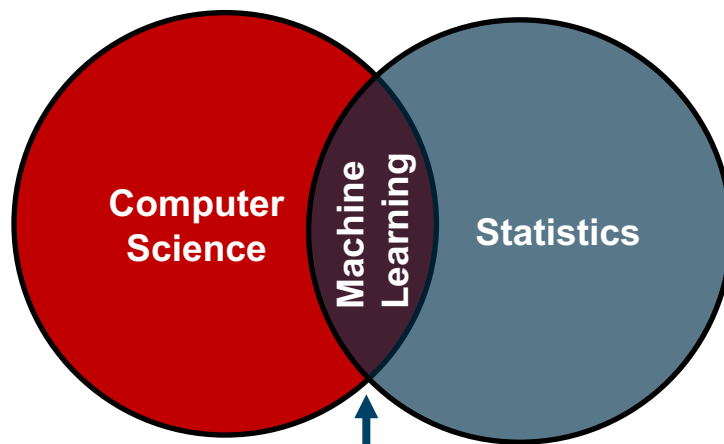
*“Machine Learning is the field of study that gives computers the ability to learn **without** being explicitly programmed.”*

*Arthur Samuel, Artificial Intelligence Pioneer  
– IBM Corporation.*

# Overview of Machine Learning

## ML IS THE NATURAL INTERSECTION OF TWO DISCIPLINES

How we build machines that solve problems.



What conclusions can be inferred from data.

How do we get computers to program themselves.

# Types of Data in Machine Learning

## Labeled



Cat



Hot  
Wings

## Unlabeled



?

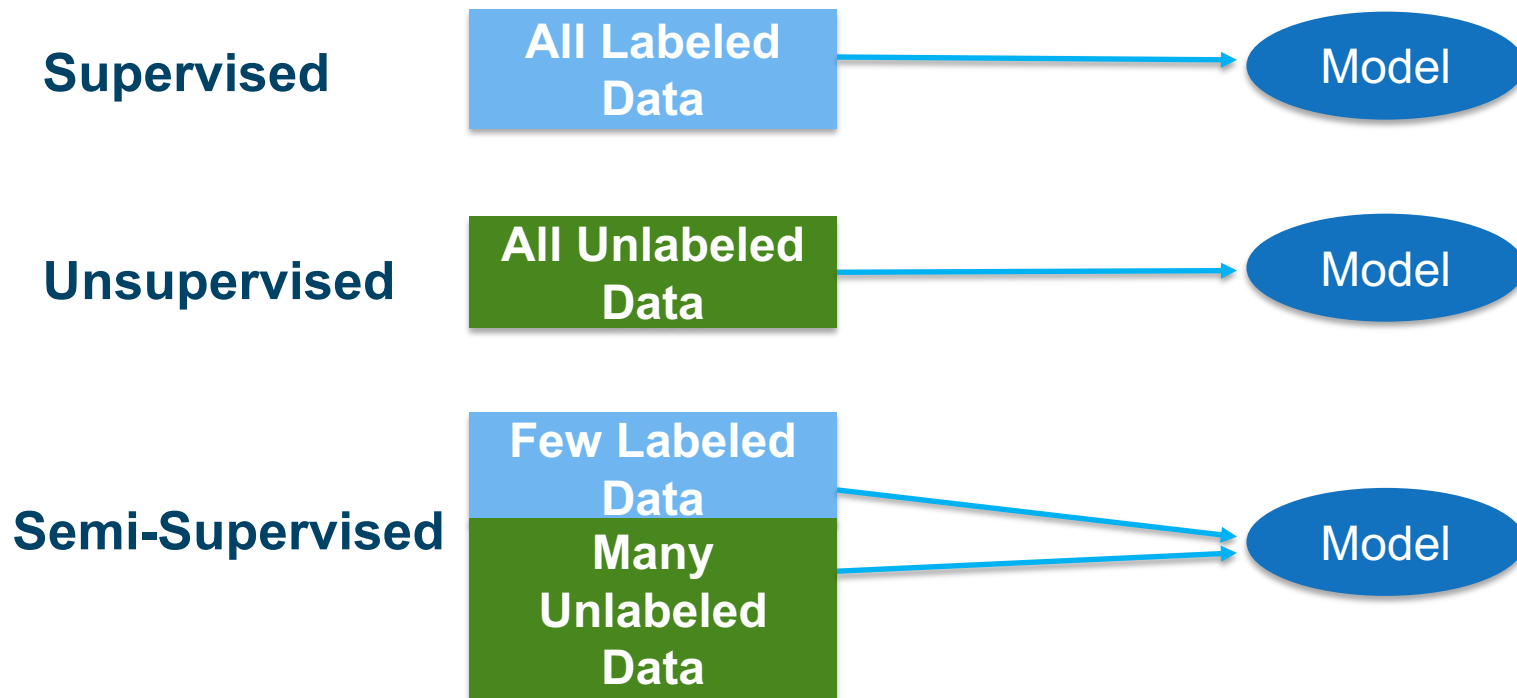


?



?

# Main Categories of ML



# Supervised Learning



Supervised Learning<sup>(2)</sup>

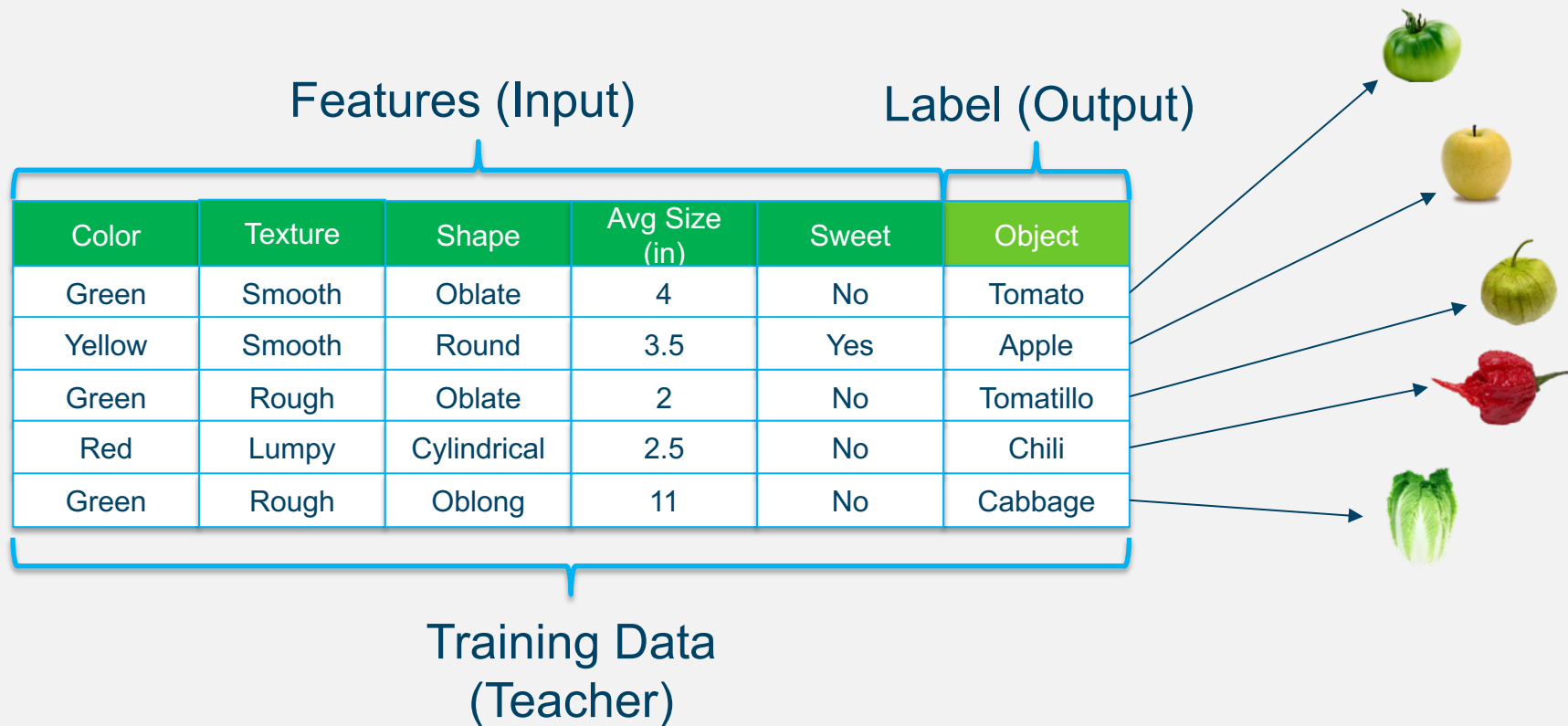
- There is a teacher
- Correct classes of training data are known
- Output is a model, or rule set

# Supervised Learning

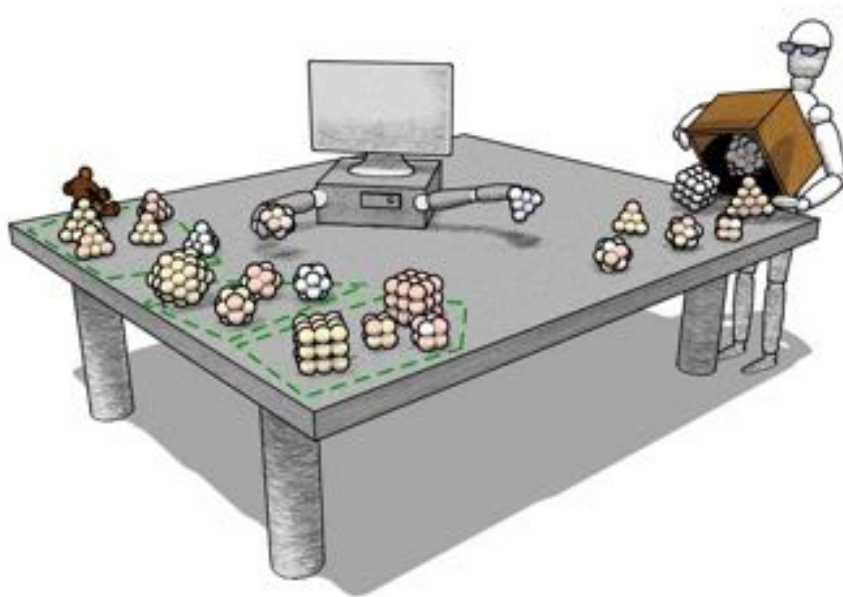
WHAT DO YOU SEE?



# Supervised Learning



# Unsupervised Learning



- There is **NO** teacher
- Correct classes of training data are **NOT** known
- Output is natural meaningful groupings



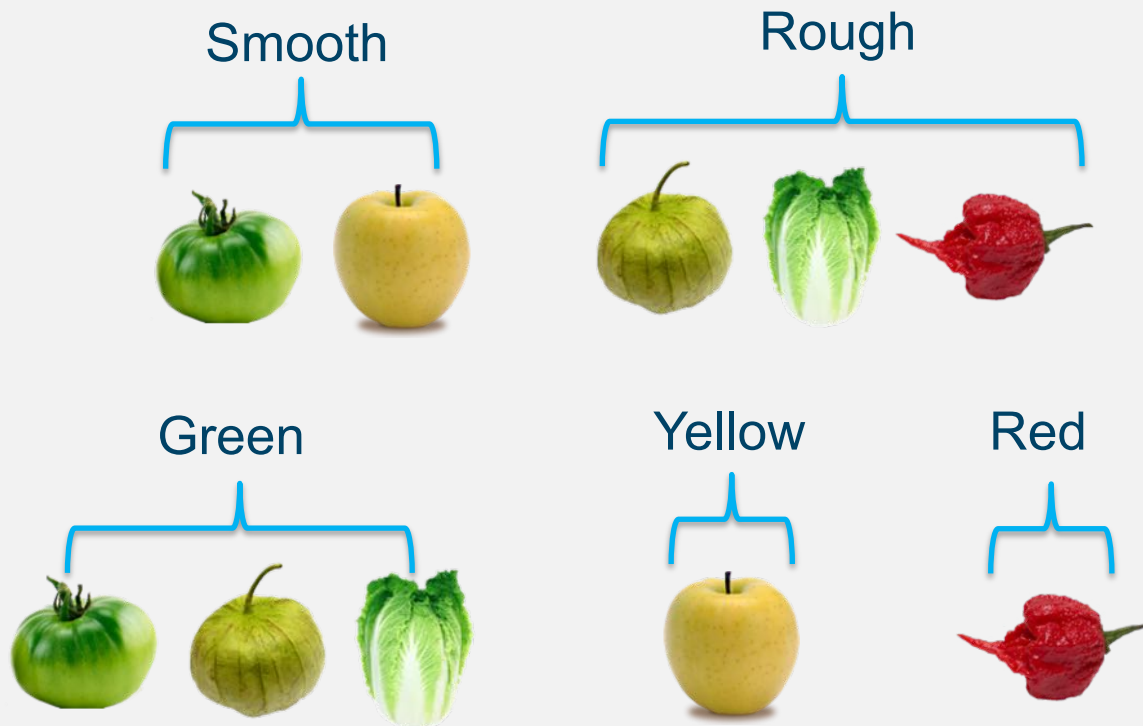
# Unsupervised Learning

**WHAT DOES A MACHINE SEE?**



# Unsupervised Learning

## GROUP AND ASSOCIATE



# Semi-supervised Learning

## SOME PRIOR KNOWLEDGE



- There is a teaching assistant
- Large amount of data
- Correct classes of training data are known for a small subset
- Output is natural “meaningful” groupings

## Section 2: Summary

Key points covered in this module:

- Machine Learning Terminology
- Labeled and Un-Labeled Data
- Main Categories of Machine Learning

# Section 3

Approaches to Machine Learning

# Main Goals of Machine Learning

## DESCRIBE

*Y = golf if rain=no and day=Saturday*

Help to understand the relationship between the inputs and the output.

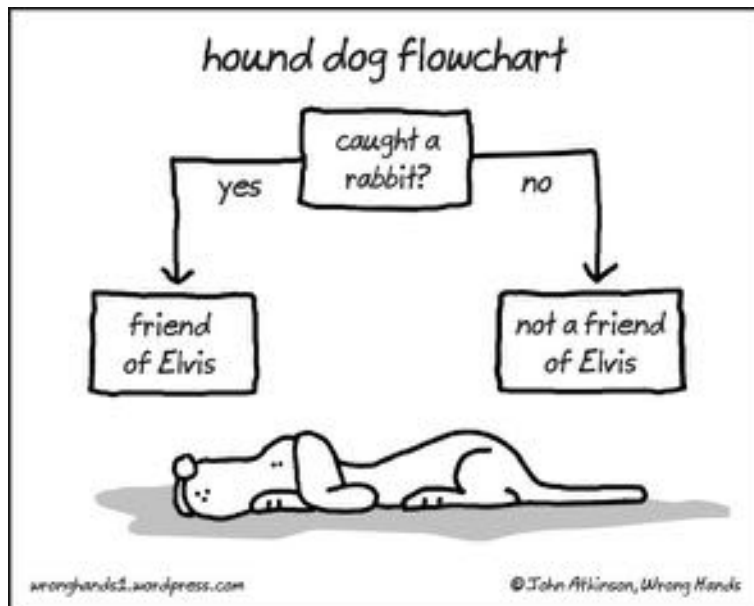
## PREDICT

Color	Texture	Shape	Avg Size (in)	Sweet	Object
Green	Smooth	Oblate	4	No	Tomato
Yellow	Smooth	Round	3.5	Yes	Apple
Green	Rough	Oblate	2	No	Tomatillo
Red	Lumpy	Cylindrical	2.5	No	Chili
Green	Rough	Oblong	11	No	Cabbage
Yellow	Smooth	Oblong	8	Yes	?

Make predictions for a new sample described by its attributes.

# Classification

## IDENTIFY GROUP MEMBERSHIP

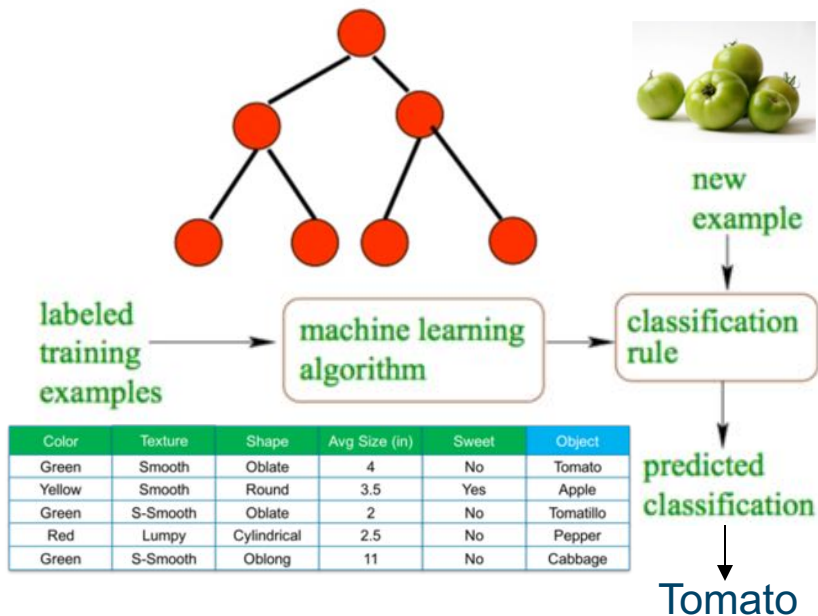


A Friend of Elvis<sup>(1)</sup>

- Supervised learning
- Predict class from observations
  - “friend of Elvis”
  - “not a friend of Elvis”
- Response variable is categorical and unordered
- Binary and nominal data

# Classification

## TECHNIQUES



### Decision Trees

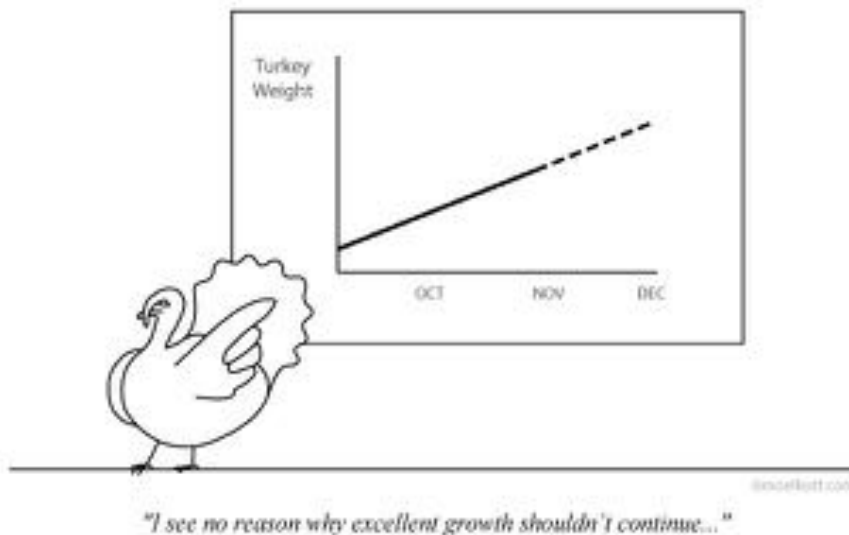
- Decision Trees
- Logistical Regression
- Naïve Bayes
- Random Forests



# Regression

## ESTIMATE OR PREDICT A RESPONSE

### THANKSGIVING PREDICTIVE ANALYTICS

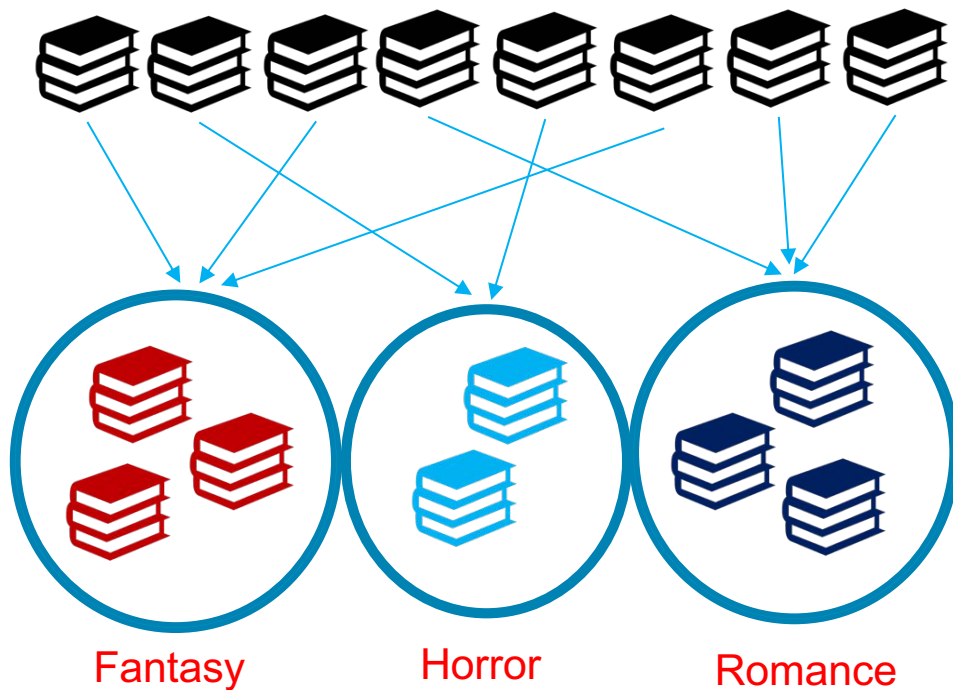


Thanksgiving Turkey. <sup>(1)</sup>

- Supervised learning
- Predict value from observations
- Response variable is numeric value, or probability of class

# Clustering

ORGANIZE DATA INTO GROUPS OF MAXIMUM COMMONALITY

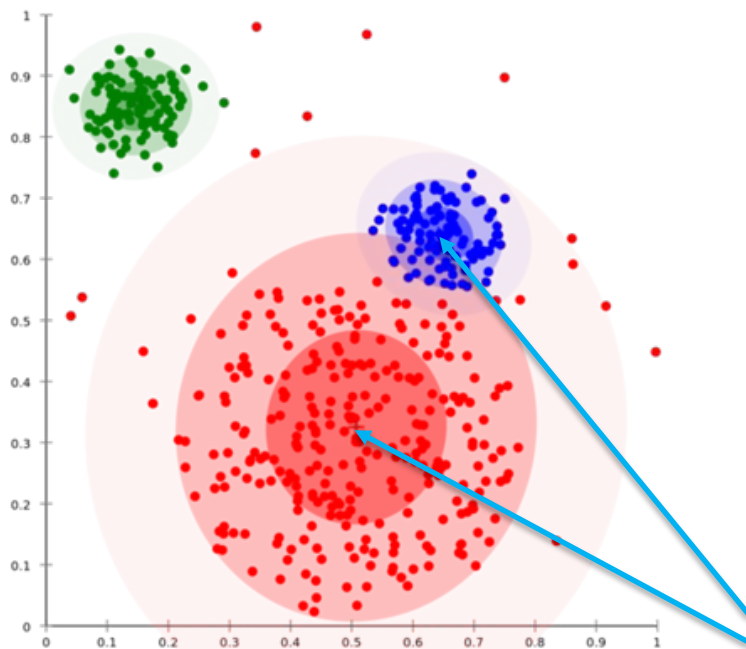


Cluster Analysis <sup>(1)</sup>

- Unsupervised learning
- Describe
- Used for exploratory mining
- Output is meaningful grouping based on similarity

# Clustering

## TECHNIQUES



Cluster Analysis <sup>(2)</sup>

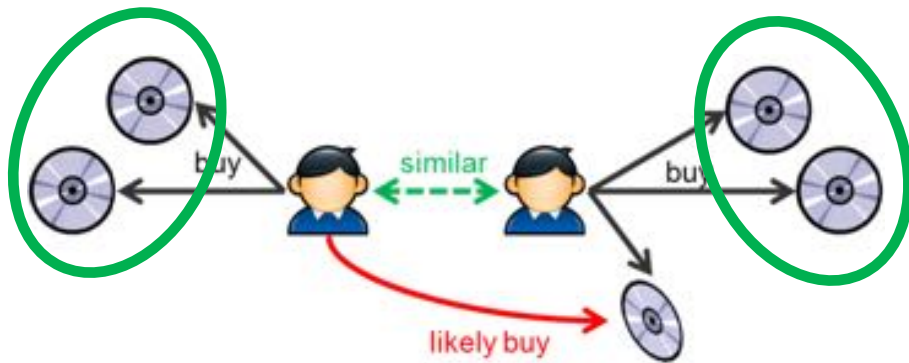
- K-Means
- K-Medians
- Expectation-Maximization
- Hierarchical Clustering

To which centroid is  
a point associated?

**K-Means**

# Associations

## DISCOVER STRONG RULES ALONG SOME MEASURE

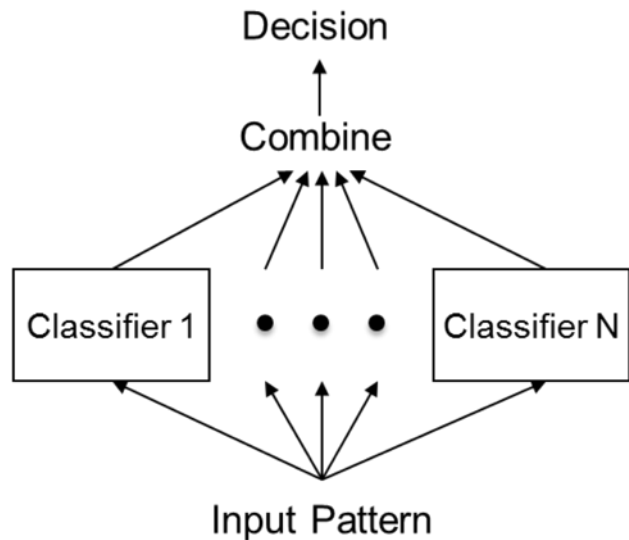


Association Rules <sup>(1)</sup>

- Supports labeled and unlabeled data
- Describes
- Data mining in large databases
- First introduced to find relationships in POS systems: Market Baskets

# Ensembles

## COMBINING MULTIPLE LEARNERS

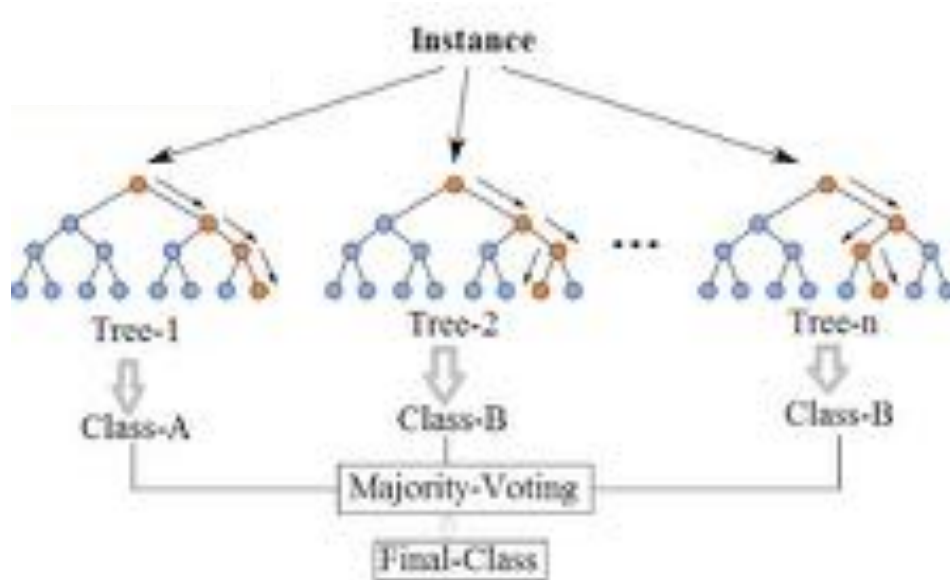


Ensemble Methods<sup>(1)</sup>

- Supervised learning
- Weak learners independently trained
- Combined predictions produce strong learner

# Ensembles

## TECHNIQUES



Random Forests<sup>(1)</sup>

## Random Forests

- Random Forests
- Bootstrap Aggregation
- Stacked Aggregation

# Method Summary

Method	Type	Goal	Input	Output	Algorithm Examples
Classification	Supervised	Predict class	Binary Nominal	Unordered categorical	Decision tree Logistic regression Naïve bayes Random forests
Regression	Supervised	Predict value Predict probability	Numeric	Numeric Probability	Linear regression Logistic regression
Clustering	Unsupervised	Describe data	Numeric	Grouping by similarity	K-Means K-Medians Expectation-Maximization Hierarchical Clustering
Association	Unsupervised	Describe data	Numeric Categorical	Associated item	Apriori
Ensemble	Supervised	Describe data Predict group Predict value Predict probability	Binary Nominal Numeric	Categorical Numerical Similarity	Random forests Bootstrap aggregation Stacked aggregation

# Section 3: Summary

Key points covered in this section:

- Main Goals of Machine Learning
- Various Approaches to Machine Learning
- Applicability of Each Approach



# Sections 4

Lab – Predicting Customer Churn

# Workshop Overview

## Key steps of the workshop:

Goal of all the models is to predict customer churn

- Import and refine data
- Import and review Jupyter Notebook
- Create SPSS Modeler flow and review output
- Create Machine Learning model, deploy, and test

# Workshop – have fun, ask questions

## Resources:

All files needed are located in GitHub

<https://github.com/team-wolfpack/Predicting-Customer-Churn-with-Watson-Data-Platform>

Watson Studio

<https://bit.ly/wplwatsonstudio>

# CALL FOR CODE®

WE ARE PROUD TO  
PUSH FOR CHANGE



@IBMWolfPack

IBM

Call for Code Founding Partner



<https://callforcode.org/>



# Appendix

Glossary

# Glossary

<b>Analytics</b>	The quantifiable informational inputs that use past data to identify possible trends that may provide valuable insight for future action.
<b>Association Rules</b>	Unsupervised learning technique to find similarities in data items based on frequent item sets. It does not predict but is used for data exploration. Most common example is market basket analysis.
<b>Classification</b>	Supervised learning technique that identifies to which of a set of categories a new observation belongs, on the basis of a training set of data containing observations whose category membership is known.
<b>Clustering</b>	Unsupervised learning technique to find similarities based on the proximity of features in a dataset. It cannot make predictions and is used for data exploration.
<b>Collinearity</b>	Refers to a linear relationship between two or more independent variables (multi-collinearity). When this exists in a data set, it can cause a model to be less accurate.
<b>Correlation</b>	Statistical relationship that involves dependence and is most often used in reference to a linear relationship. Example: Price and Demand.

# Glossary (cont.)

<b>Data</b>	Raw objective fact about an event.
<b>Data Science</b>	Interdisciplinary field about scientific methods, processes, and systems to extract knowledge or insights from data in various forms, either structured or unstructured.
<b>Dependent Variable</b>	The value of the variable is influenced by other variables. Also referred to as the outcome or target.
<b>Descriptive Statistics</b>	Methods of organizing, summarizing and presenting information about data.
<b>Ensemble</b>	Machine learning technique that combines multiple models (weak learners) to make an overall strong learner. The models may use all the same or different algorithms and the training data may be all the same or resampled. Most common example is Random Forests.
<b>Independent Variable</b>	The value of the variable is not influenced by other variables. These variables may influence dependent variables.



# Glossary (cont.)

<b>Inferential Statistics</b>	Methods to determine something about a population from a sample.
<b>Information</b>	Data with interpretation: useful, organized, and structured.
<b>Knowledge</b>	Information with context, it has understanding and meaning.
<b>Kurtosis</b>	A measure of the "tailedness" of the probability distribution of a real-valued random variable.
<b>Labeled Data</b>	Data that has been identified and assigned a label.
<b>Machine Learning</b>	The field of study that gives computers the ability to learn without being explicitly programmed.
<b>Mean</b>	The quotient of the sum of the data points and the number of data points; another name for the average.
<b>Median</b>	When the data are arranged in sorted order, the median is that data point at which 50% of the data points are either less than or greater than that data point; the data point in the middle.

# Glossary (cont.)

<b>Mode</b>	The data point that occurs most frequently. There can be more than one mode.
<b>Quasi-Structured Data</b>	Textual data with erratic data formats, can be formatted with effort, tools, and time.
<b>Semi-Structured Data</b>	Textual data with a discernable pattern, enabling parsing. E.g. Self-describing XML with schema.
<b>Semi-Supervised Learning</b>	Category of machine learning that uses few labeled data, and many unlabeled data. There is a teaching assistant, usually involves large amounts of unclassified data with few classified data. The output is a meaningful grouping of data.
<b>Skewness</b>	A measure of the asymmetry a data distribution. A measure of 0 indicates a perfect symmetry.
<b>Structured Data</b>	Defined data type, format, and structure. E.g. Transactional Data.

# Glossary (cont.)

<b>Semi-Supervised Learning</b>	Category of machine learning that uses few labeled data, and many unlabeled data. There is a teaching assistant, usually involves large amounts of unclassified data with few classified data. The output is a meaningful grouping of data.
<b>Unlabeled Data</b>	Refers to data that has no clear label or indication of what it is.
<b>Unstructured Data</b>	Data that has no inherent structure and is usually stored in different file types. E.g. PDF, Excel, and JPG.
<b>Unsupervised Learning</b>	Category of machine learning that only uses all unlabeled data. There is no teacher, the correct classes of training data are not known, and the output is a natural meaningful grouping of the data.
<b>Wisdom</b>	Knowledge with insight. Integrated understanding and actionable.