# Machine Learning per la Finanza

Docenti: Zelda Marino e Paolo Zanetti

Email: {zelda.marino;paolo.zanetti}@uniparthenope.it

Data: 16/04/2021
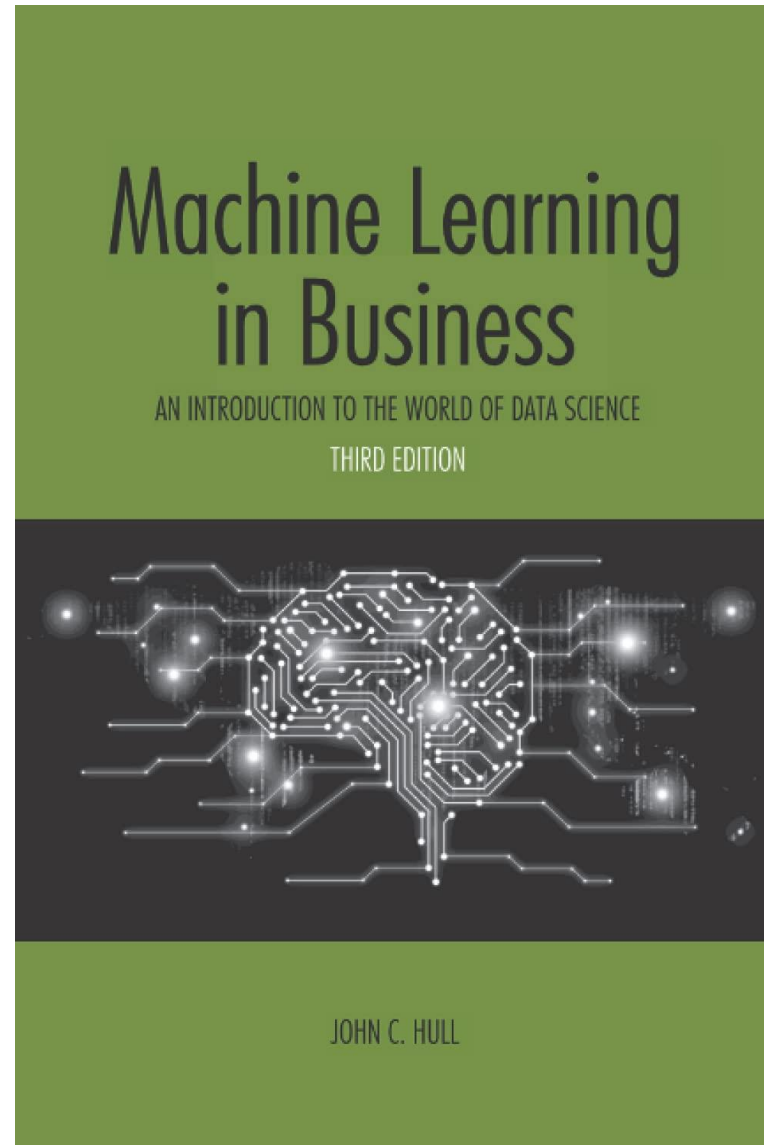
# Marzo

| L | M | M | G | V | S | D | Tot Ore |
|---|---|---|---|---|---|---|---------|
| 27 | 28 | 1 | 2 | 3 | 4 | 5 | 7 |
| 6 | 7 | 8 | 9 | 10 | 11 | 12 | 14 |
| 13 | 14 | 15 | 16 | 17 | 18 | 19 | 21 |
| 20 | 21 | 22 | 23 | 24 | 25 | 26 | 28 |
| 27 | 28 | 29 | 30 | 31 | | | 35 |

# Aprile

| L | M | M | G | V | S | D | Tot Ore |
|---|---|---|---|---|---|---|---------|
| | | | | | 1 | 2 | |
| 3 | 4 | 5 | 6 | 7 | 8 | 9 | |
| 10 | 11 | 12 | 13 | 14 | 15 | 16 | |
| 17 | 18 | 19 | 20 | 21 | 22 | 23 | 42 |
| 24 | 25 | 26 | 27 | 28 | 29 | 30 | 49 |

# Maggio

| L | M | M | G | V | S | D | Tot Ore |
|---|---|---|---|---|---|---|---------|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 54 |
| 8 | 9 | 10 | 11 | 12 | 13 | 14 | 61 |
| 15 | 16 | 17 | 18 | 19 | 20 | 21 | 68 |
| 22 | 23 | 24 | 25 | 26 | 27 | 28 | 72 |
| 29 | 30 | 31 | | | | | |

| # | Lezione Marino | # | prove intercorso |
|---|----------------|---|------------------|
| # | Lezione Zanetti | # | recuperi |
| # | Festa accademica | | |

| DATE ESAME | | LUN | 15.00 | 17.00 |
|------------|---|-----|-------|-------|
| 5 | giugno | GIOV | 8.30 | 10.30 |
| 19 | giugno | VEN | 11.30 | 14.30 |
| 10 | luglio | | | |
| 11 | settembre | | | |

Machine Learning in Business

AN INTRODUCTION TO THE WORLD OF DATA SCIENCE

THIRD EDITION

JOHN C. HULL

# What is Data Science? Ask google…

# The Virtuous Circle of Machine Learning and AI

# Making data for you



Use data to better describe the present or better predict the future

A DAY IN DATA

The exponential growth of data is undisputed, but the numbers behind this explosion – fuelled by internet of things and the use of connected devcies – are hard to comprehend, particularly when looked at in the context of one day

**500m** tweets are sent every day — Twitter

**4PB** of data created by Facebook, including
350m photos
100m hours of video watch time
Facebook Research

**320bn** emails to be sent each day by 2021

**306bn** emails to be sent each day by 2020

**294bn** billion emails are sent — Radicati Group

**3.9bn** people use emails

**4TB** of data produced by a connected car — Intel

**65bn** messages sent over WhatsApp and two billion minutes of voice and video calls made — Facebook

**95m** photos and videos are shared on Instagram — Instagram Business

**463EB** of data will be created every day by 2025 — IDC

**28PB** to be generated from wearable devices by 2020 — Statista

Searches made a day — 5bn
Searches made a day from Google — 3.5bn
Smart Insights

ACCUMULATED DIGITAL UNIVERSE OF DATA
4.4ZB — 2013
44ZB — 2020
PwC

**DEMYSTIFIYING DATA UNITS**
From the more familiar 'bit' or 'megabyte', larger units of measurement are more frequently being used to explain the masses of data

| Unit | | Value | Size |
|------|------|-------|------|
| b | bit | 0 or 1 | 1/8 of a byte |
| B | byte | 8 bits | 1 byte |
| KB | kilobyte | 1,000 bytes | 1,000 bytes |
| MB | megabyte | 1,000$^2$ bytes | 1,000,000 bytes |
| GB | gigabyte | 1,000$^3$ bytes | 1,000,000,000 bytes |
| TB | terabyte | 1,000$^4$ bytes | 1,000,000,000,000 bytes |
| PB | petabyte | 1,000$^5$ bytes | 1,000,000,000,000,000 bytes |
| EB | exabyte | 1,000$^6$ bytes | 1,000,000,000,000,000,000 bytes |
| ZB | zettabyte | 1,000$^7$ bytes | 1,000,000,000,000,000,000,000 bytes |
| YB | yottabyte | 1,000$^8$ bytes | 1,000,000,000,000,000,000,000,000 bytes |

*A lowercase "b" is used as an abbreviation for bits, while an uppercase "B" represents bytes.

https://www.weforum.org/agenda/2019/04/how-much-data-is-generated-each-day-cf4bddf29f/

RACONTEUR

Speed up calculations with 1000s of processors

NVIDIA TITAN V

NVIDIA's SUPERCOMPUTING GPU ARCHITECTURE, NOW FOR YOUR PC

NVIDIA TITAN V is the most powerful Volta-based graphics card ever created for the PC. NVIDIA's supercomputing GPU architecture is now here for your PC, and fueling breakthroughs in every industry.

Scale computations with infinite compute power

amazon webservices™

Microsoft Azure

Google Cloud Platform

# What is Machine Learning

- Machine learning is a branch of AI

- The idea underlying machine learning is that we give a computer program access to lots of data and let it learn about relationships between variables and make predictions

- Some of the techniques of machine learning date back to the 1950s but improvements in computer speeds and data storage costs have now made machine learning a practical tool

# Artificial Intelligence

a branch of computer science dealing with the simulation of intelligent behavior in computers; (2) the capability of a machine to imitate intelligent human behavior.

Artificial Intelligence

Machine Learning

Deep Learning

the study of algorithms related to artificial neural networks that contain many blocks stacked on each other.

# Data science skills

- **Technology**: manage data (structured and unstructured, big data)

- **Machine Learning**: Mathematics and Statistics

- **Programming**: open source Data Science programming suite (Python, R, ….)

- **Business knowledge:** transfer information got from data to Analysts

- **Communication**: data visualization tools

# Machine Learning vs. Automation

- Computers have been used to automate many business decisions (payroll, sending out invoices, summarizing sales by region, etc)

- This is digitization: the third industrial revolution

- Machine learning is central to the fourth industrial revolution where computers are used to create intelligence

# The 4th Industrial revolution is Here!



Source: *Christoph Roser at AllAboutLean.com*

As per Wikipedia*, "The 4th Industrial Revolution ….. marked by emerging technology breakthroughs in a number of fields, including robotics, **artificial intelligence**, nanotechnology, quantum computing, biotechnology, the Internet of Things, the Industrial Internet of Things (IIoT), decentralized consensus, fifth-generation wireless technologies (5G), additive manufacturing/3D printing and fully autonomous vehicles."

* https://en.wikipedia.org/wiki/Fourth_Industrial_Revolution

# Example: Loan Applications (digitization vs. ML)

- If loan officers applied certain known rules we could digitize their activities

- If we did not know the rules used, we could use ML to determine them

- But we could go one step further and use ML to improve upon the rules for accepting or rejecting loans

# The new world of statistics

- Huge data sets

- Fantastic improvements in computer processing speeds and data storage costs

- Machine learning tools are now feasible

- Can now develop non-linear prediction models, find patterns in data in ways that were not possible before, and develop multi-stage decision strategies

- New terminology: features, labels, activation functions, target, bias, supervised/unsupervised learning……

# Traditional Statistics vs Machine Learning (Figure 1.1)



**Statistics**

Develop Hypothesis

↓

Collect Data

↓

Test Hypothesis

**Machine Learning**

Collect Data

↓

Try different models

↓

Find Patterns or Develop a Predictive Tool

- **Describe the current state of an organization or process**

- **Detect anomalous event**

- **Diagnose the causes of events and behaviors**

- **Predict future events**

Dynamic pricing

Predicting flight delay

01 TRAVEL

Disease prediction

Medication effectiveness

Cross selling

Upselling

Predicting lifetime valueof customer

02 MARKETING

Sentiment analysis

Digital marketing

Churn

03 HEALTHCARE

Discount offering

Demand forecasting

04 SOCIAL MEDIA

Self driving cars

06 AUTOMATION

Pilotless aircrafts, drones

05 SALES

Claims prediction

07 CREDIT & INSURANCE

Fraud & risk detection

**Average Years of Experience**

| | |
|---|---|
| 0 - 1 | 8% |
| 2 - 4 | 80% |
| 5 - 7 | 6% |
| 8+ | 6% |

**Common Skill Sets**

- ✓ Machine Learning
- ✓ Python
- ✓ Hadoop SPARK
- ✓ SQL
- ✓ Statistics
- ✓ Natural Language Processing
- ✓ Algorithms
- ✓ Programming Languages

## Data Scientist Salaries

**Average Base Pay**

# $114,673 /yr

Same as national average

Not including cash compensation

**See More Insights**

$81K          Median: $115K          $163K

# The Data scientist

# Data science in sports

# Data Science in Sports



The market for baseball players was so inefficient
...
that superior management  could run circles
around taller piles of cash
- Michael Lewis

Legendary 2002 season for Oakland Athletics.

Manager Billy Beane put together an unexpected
team using data science.

# Data Science in Drug Discovery

# Data Science in Drug Discovery



https://www.nature.com/articles/d41586-018-05267-x

https://blog.benchsci.com/startups-using-artificial-intelligence-in-drug-discovery

# Data science for politics

**Opinions**    Editorial Board    The Opinions Essay    Global Opinions    Post Opinión    First 100 Days    Reimag

**Opinions**

# Obama, the 'big data' president

By **Nancy Scola**
June 14, 2013

*Nancy Scola is a journalist covering technology and politics. From 2001 to 2005, she served on the staff of the House government oversight committee.*

n the political world, the promise of data — whether it's Nate Silver's spot-on election predictions or President Obama's clearinghouse of government information, Data.gov — is that we no longer have to take so much on faith. "What do the data show?" is the new "What do you think?," the new "Is this a good idea?"

FiveThirtyEight

# Who will win the presidency?

**Chance of winning**

Hillary Clinton
**71.4%**

Donald Trump
**28.6%**

https://projects.fivethirtyeight.com/2016-election-forecast/

# Data Science in Politics



https://fivethirtyeight.com/tag/2018-election/

# Data Science in Commerce

# Data Science in Commerce

# Netfix challenge

# ML and AI is revolutionizing finance

# Market impact at the speed of light!

# Machine Learning & AI in finance: A paradigm shift

Quant

| |
|---|
| Stochastic Models |
| Factor Models |
| Optimization |
| Risk Factors |
| P/Q Quants |
| Derivative pricing |
| Trading Strategies |
| Simulations |
| Distribution fitting |

Data Scientist

| |
|---|
| Real-time analytics |
| Predictive analytics |
| Machine Learning |
| RPA |
| NLP |
| Deep Learning |
| Computer Vision |
| Graph Analytics |
| Chatbots |
| Sentiment Analysis |
| Alternative Data |

«Financial technologies of «fintech» is used to describe a variety of

**innovative business models**

and

**emerging technologies**

that have the potential to transform the financial service industry»

# Fintech funding more than doubled QoQ
## Global VC-backed fintech funding trends, Q1'18 – Q1'21



Fintech investment trends

# Europe saw the largest QoQ increase in funding

## FINTECH INVESTMENT TRENDS

Quarterly funding ($M) by continent, Q1'20 – Q1'21

**Q1 2020:** $0.3, $43, $223, $2,261, $4,763, $4,198

**Q2 2020:** $15, $313, $405, $2,048, $1,842, $5,731

**Q3 2020:** $15, $42, $714, $2,234, $1,765, $5,666

**Q4 2020:** $45, $334, $534, $1,802, $1,908, $6,211

**Q1 2021:** $45, $193, $999, $5,049, $3,668, $12,822

+180% QoQ in Europe

**Legend:** North America ■ Asia ■ Europe ■ South America ■ Australia ■ Africa

# What the state of fintech covers

**PAYMENTS**

Payments processing, card developers, money transfer platforms, and tracking software

**BANKING**

Digital-first banks or companies digitizing banking services for credit and debit

**DIGITAL LENDING**

Companies creating new solutions for personal or commercial lending

**WEALTH MANAGEMENT**

Personal finance tools, investment and wealth management platforms, and analytics tools

**INSURANCE**

Companies selling or distributing insurance digitally or providing data analytics and software for (re)insurers

**CAPITAL MARKETS**

Sales and trading, analysis, and infrastructure tools for financial institutions

**SMB**

Companies focused on providing solutions to small- and medium-sized businesses

**REAL ESTATE**

Mortgage lending, transaction digitization, and financing platforms

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

Algorithmic trading (or simply algo trading) is the use of algorithms to conduct trades autonomously.

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

Robo-advisors, algorithms built to calibrate a financial
portfolio to the goals and risk tolerance of the user. Additionally, they provide
automated financial guidance and service to end investors and clients.

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

Fraud is a massive problem for financial institutions and one of the foremost reasons to leverage machine learning in finance.

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

Underwriting could be described as a perfect job for machine learning in finance, and indeed there is a great deal of worry in the industry that machines will replace a large swath of underwriting positions that exist today.

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

Automation is patently well suited to finance. It reduces the strain that repetitive, low-value tasks put on human employees. It tackles the routine, everyday processes, freeing up teams to finish their high-value work.

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

All aspects of understanding and controlling

risk are being revolutionized through

the growth of solutions driven by machine

learning

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

Asset price prediction is considered the most frequently discussed and most sophisticated area in finance.

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

The classic derivative pricing models are built on several impractical assumptions to reproduce the empirical relationship between the underlying input data (strike price, time to maturity, option type) and the price of the derivatives observed in the market.

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

Sentiment analysis involves the perusal of enormous volumes of unstructured data, such as videos, transcriptions, photos, audio files, social media posts, articles, and business documents, to determine market sentiment.

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

Trade settlement is the process of transferring securities into the account of a buyer and cash into the seller's account following a transaction of a financial asset.

# Current and Future Applications in Finance

- Algorithmic Trading
- Portfolio Management and Robo-Advisors
- Fraud Detection
- Loans/Credit Card/Insurance Underwriting
- Automation and Chatbots
- Risk Management
- Asset Price Prediction
- Derivative Pricing
- Sentiment Analysis
- Trade Settlement
- Money Laundering

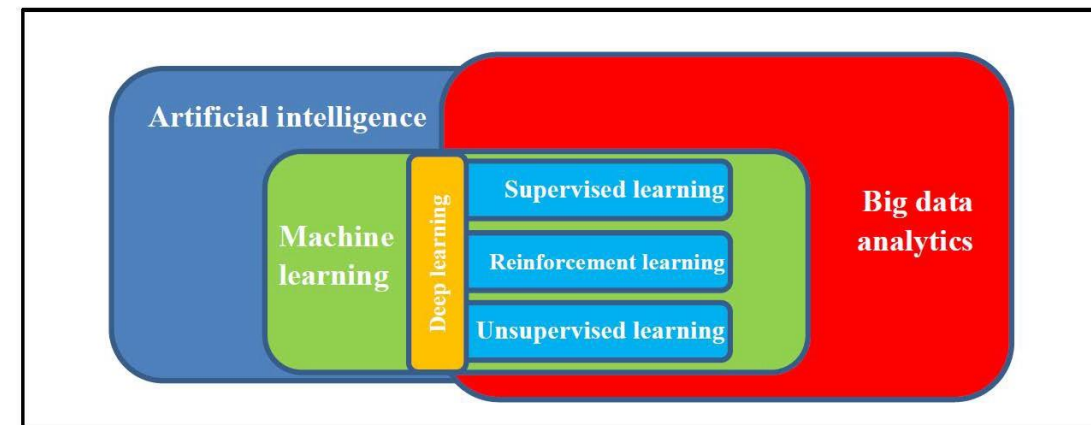A United Nations report estimates that the amount of money laundered worldwide per year is 2%–5% of global GDP

# An intuitive introduction to AI and ML
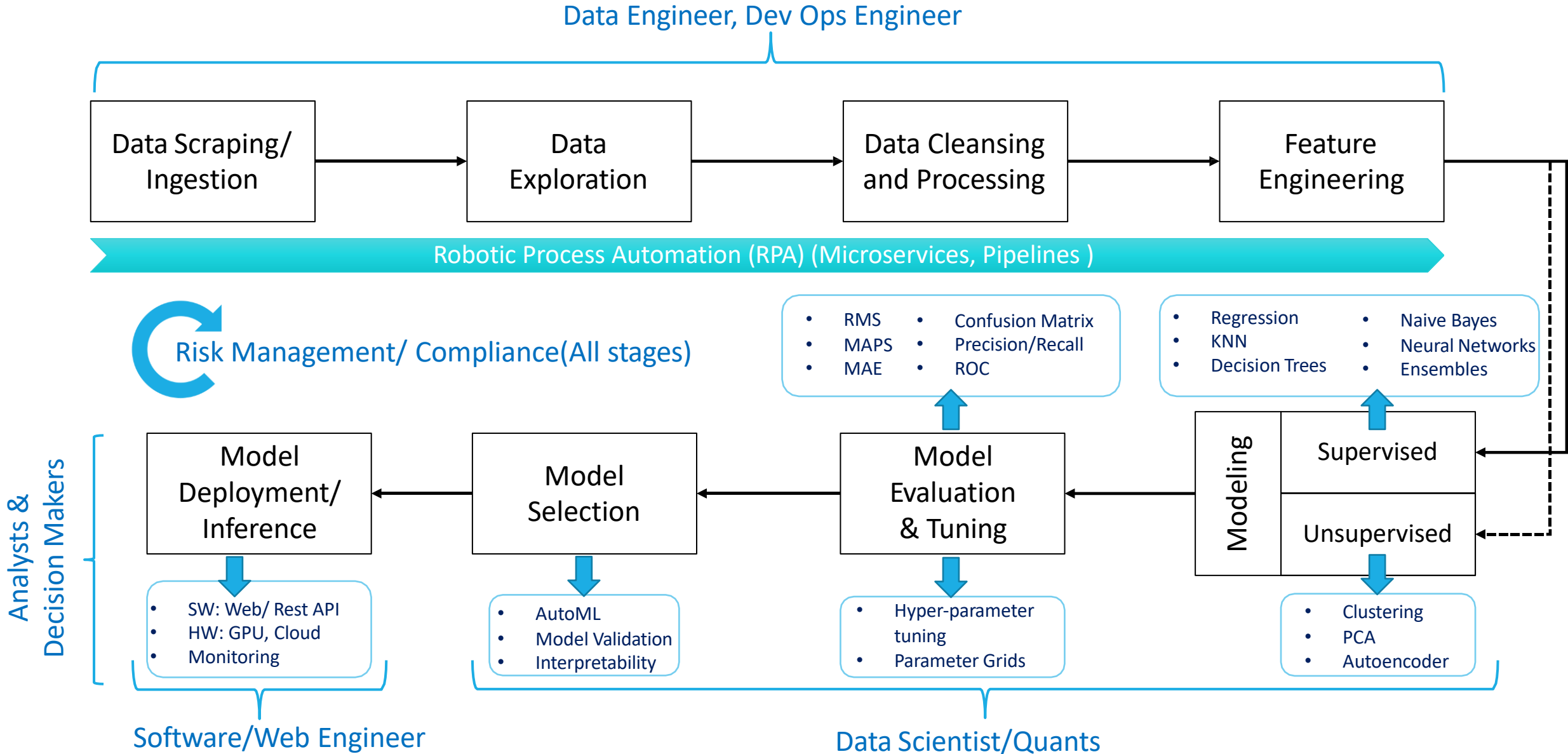
# Definitions: Machine Learning and AI

- Machine learning is the scientific study of algorithms and statistical models that computer systems use to effectively perform a specific task without using explicit instructions, relying on patterns and inference instead[1]

- Artificial intelligence is intelligence demonstrated by machines, in contrast to the natural intelligence displayed by humans and animals[1]

1. https://en.wikipedia.org/wiki/Machine_learning
2. Figure Source: http://www.fsb.org/wp-content/uploads/P011117.pdf



Figure 1: A schematic view of AI, machine learning and big data analytics

# Machine Learning Workflow

Data Engineer, Dev Ops Engineer

```
Data Scraping/        →  Data           →  Data Cleansing    →  Feature
Ingestion                Exploration        and Processing       Engineering
```

Robotic Process Automation (RPA) (Microservices, Pipelines )

Risk Management/ Compliance(All stages)

- RMS
- MAPS
- MAE
- Confusion Matrix
- Precision/Recall
- ROC

- Regression
- KNN
- Decision Trees
- Naive Bayes
- Neural Networks
- Ensembles

Analysts & Decision Makers

```
Model              ←  Model        ←  Model            ←  Modeling
Deployment/           Selection        Evaluation          Supervised
Inference                              & Tuning            Unsupervised
```

- SW: Web/ Rest API
- HW: GPU, Cloud
- Monitoring

- AutoML
- Model Validation
- Interpretability

- Hyper-parameter tuning
- Parameter Grids

- Clustering
- PCA
- Autoencoder

Software/Web Engineer

Data Scientist/Quants

# Key steps involved

1. Data

2. Goals

3. Machine learning algorithms

4. Process

5. Performance evaluation

# Data

## Quantitative data

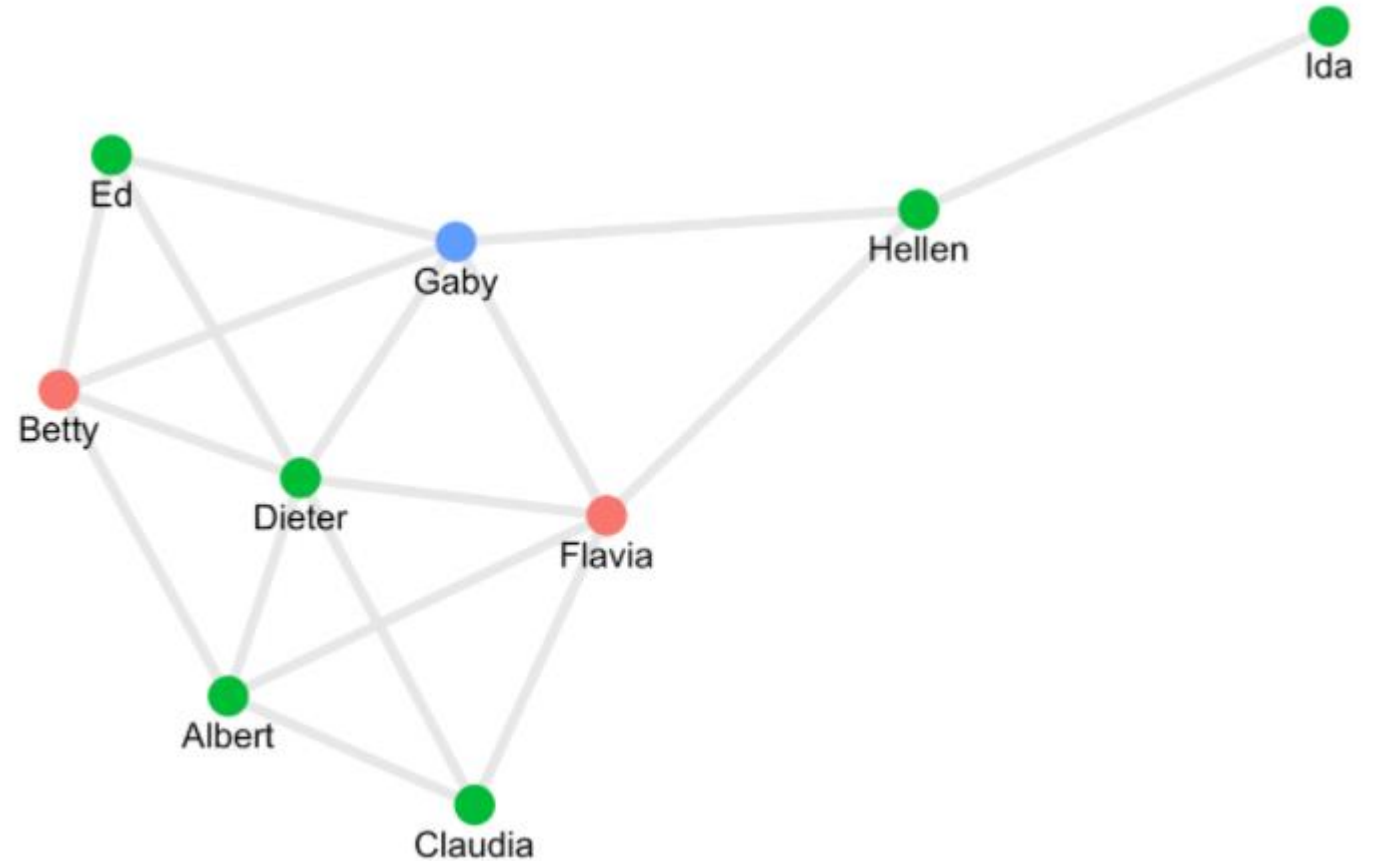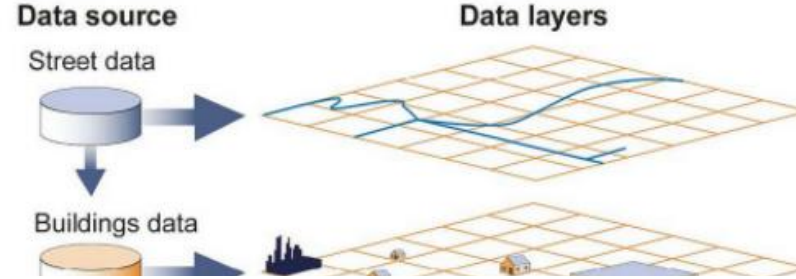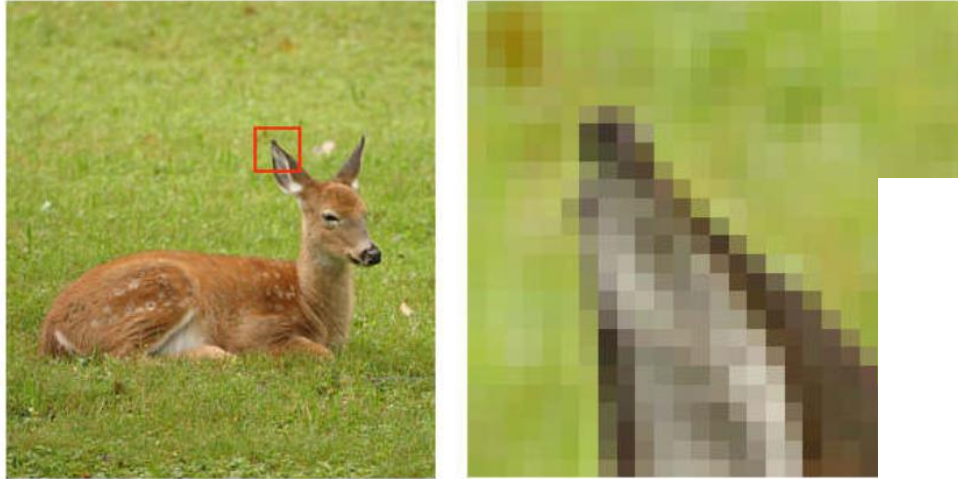- Deals with numbers

- Data can be measured

## Qualitative data

- Deals with descriptions

- Data can be observed but not measured

# Quantitative vs qualitative data



"Great evening, extremely good value"

⊙⊙⊙⊙⊙ Review of L'Ange 20 Restaurant

I went to this place with my boyfriend for a special occa
were greeted warmly by Christopher who guided us thro
delicious and I only wish that we could have had room f
excellent compared to other prices we had seen and we
hard to match during the rest of our stay.

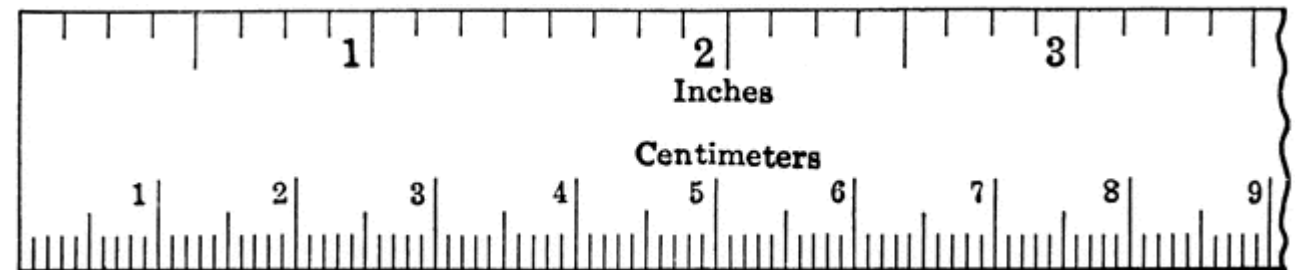I had the lamb which I can highly recommend. When we

A variable could be:

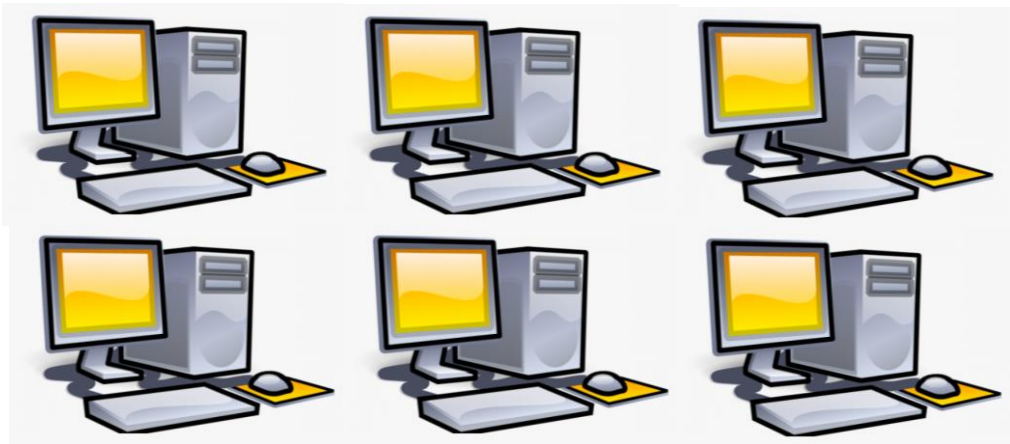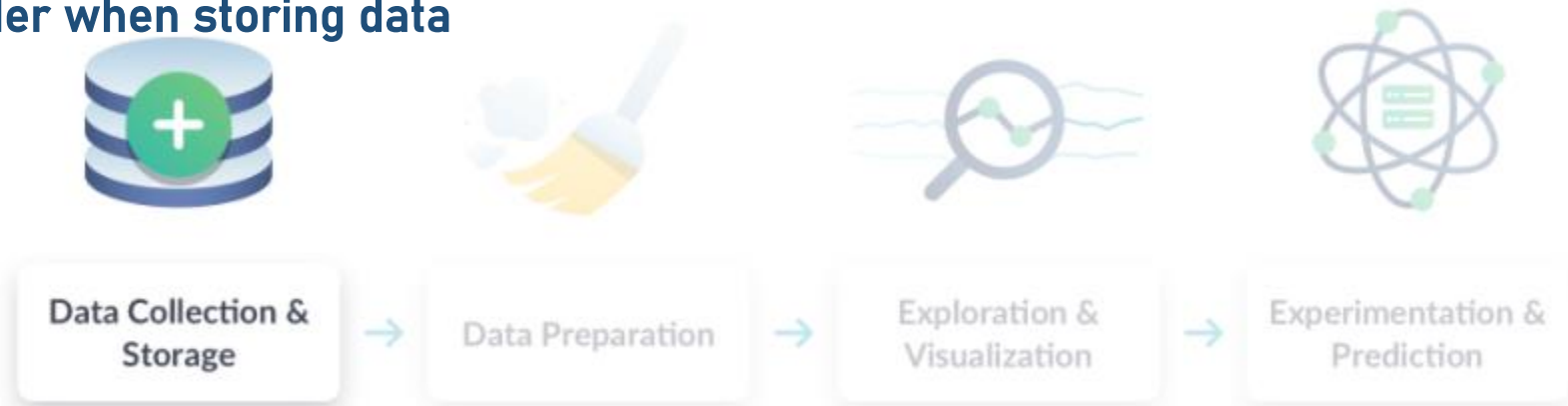- **Categorical**
  - Yes/No flags
  - AAA,BB ratings for bonds

- **Numerical**
  - 35 mpg
  - $170K salary

# Data storage and retrieval

## Things to consider when storing data

- Location
- Data type
- Retrieval

Data Collection & Storage → Data Preparation → Exploration & Visualization → Experimentation & Prediction

# Types of data storage

## Unstructured

- Email

- Text

- Video and audio files

- Web pages

- Social media

Tabular

| Customer Name | Customer Address | ... |
|---|---|---|
| Jane Doe | 123 Maple St. | ... |

## Relational Database

| Data Type | Query Language |
|---|---|
| Document Database | NoSQL |
| Relational Database | SQL |

## Document Database

# Longitudinal

- Observations are dependent
- Temporal-continuity is required

|  | 2009 |
|---|---|
| January | 339 778 |
| February | 343 438 |
| March | 339 228 |
| April | 338 344 |
| May | 339 873 |
| June | 342 912 |
| July | 342 489 |
| August | 350 800 |
| September | 343 687 |
| October | 347 641 |
| November | 354 467 |
| December | 354 085 |

# Cross-sectional

- Observations are independent

| Ten Highest-Yielding Dow Stocks | | | |
|---|---|---|---|
| December 31, 2007 | 2007 Dividends | Price | Yield |
| 1 Citigroup | 2.1600 | 29.44 | 7.34% |
| 2 Pfizer | 1.1600 | 22.73 | 5.10% |
| 3 Altria Group | 3.0500 | 75.58 | 4.04% |
| 4 General Motors | 1.0000 | 24.89 | 4.02% |
| 5 Verizon | 1.6450 | 43.69 | 3.77% |
| 6 du Pont | 1.5200 | 44.09 | 3.45% |
| 7 AT&T | 1.4200 | 41.56 | 3.42% |
| 8 Home Depot | 0.9000 | 26.94 | 3.34% |
| 9 JP Morgan Chase | 1.4400 | 43.65 | 3.30% |
| 10 General Electric | 1.1500 | 36.96 | 3.11% |

# Goals

- **Descriptive Statistics**
  - Goal is to describe the data at hand
  - Backward-looking
  - Statistical techniques employed here
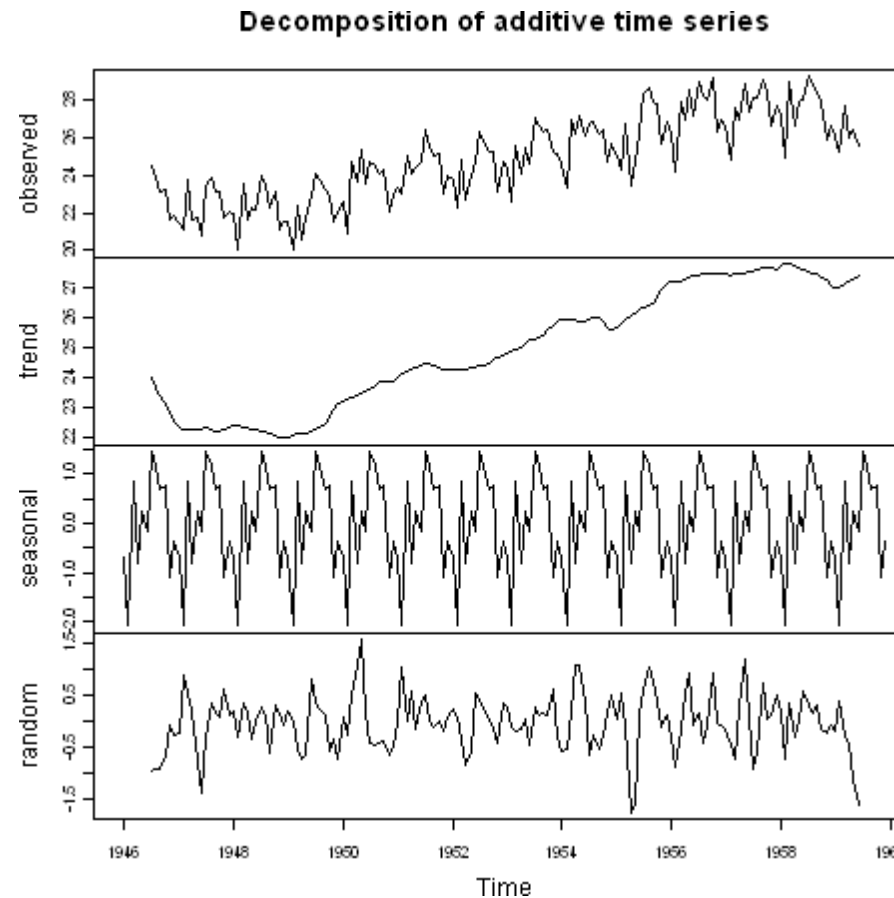
- **Predictive Analytics**
  - Goal is to use historical data to build a model for prediction
  - Forward-looking
  - Machine learning & AI techniques employed here

- How do you summarize numerical variables ?
- How do you summarize categorical variables ?
- How do you describe variability in numerical variables ?
- How do you summarize relationships between categorical and numerical variables ?
- How do you summarize relationships between 2 numerical variables?

# Goal is to extract the various components
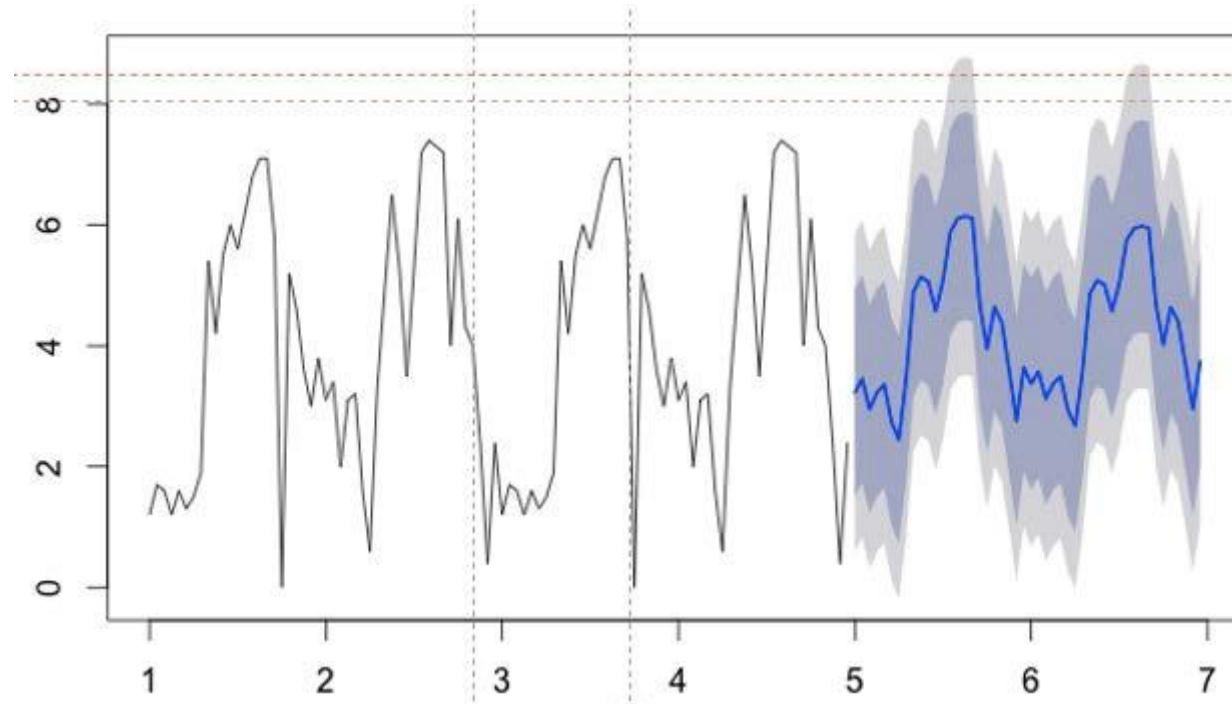

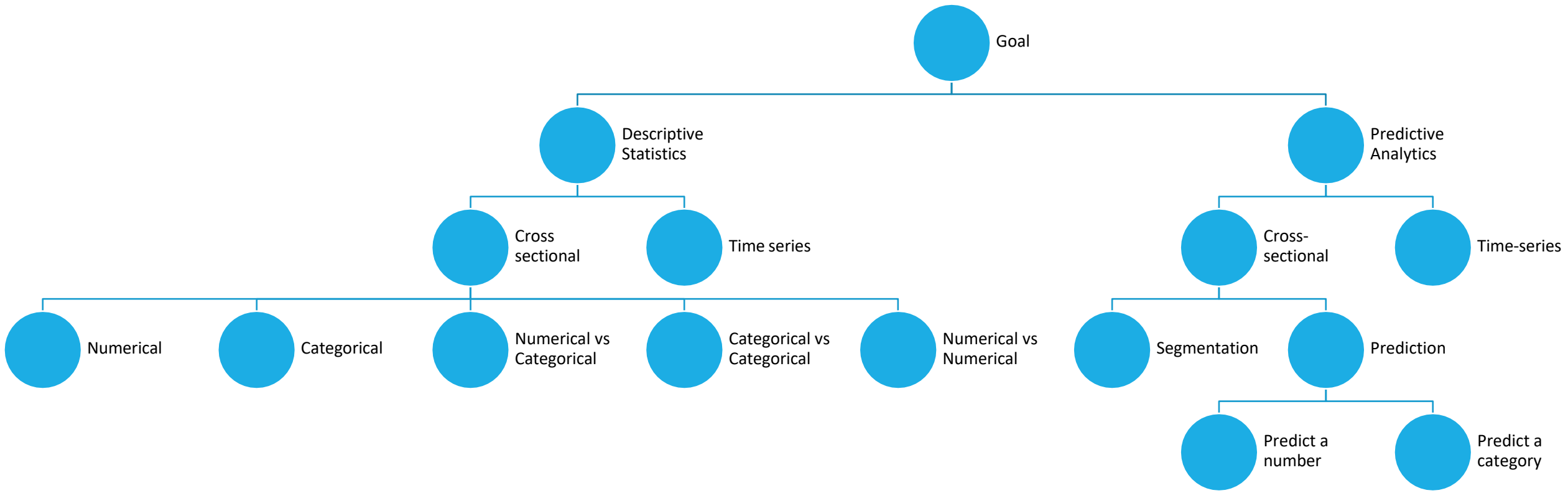Decomposition of additive time series

- Given a dataset, build a model that captures the similarities in different observations and assigns them to different buckets.

- Given a set of variables, predict the value of another variable in a given data set
  - Predict salaries given work experience, education etc.

  - Predict whether a loan would be approved given fico score, current loans, employment status etc.

- Given a time series dataset, build a model that can be used to forecast values in the future
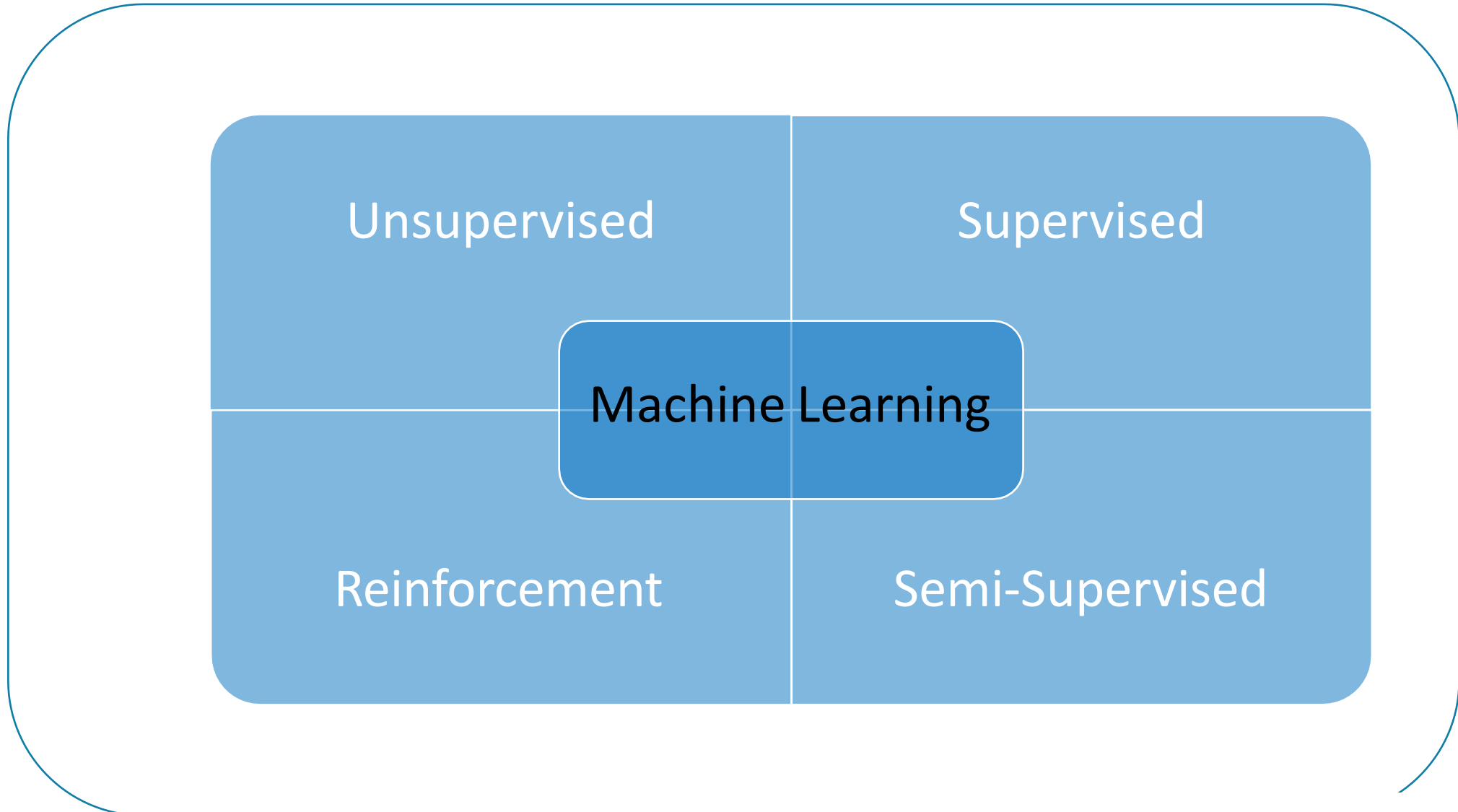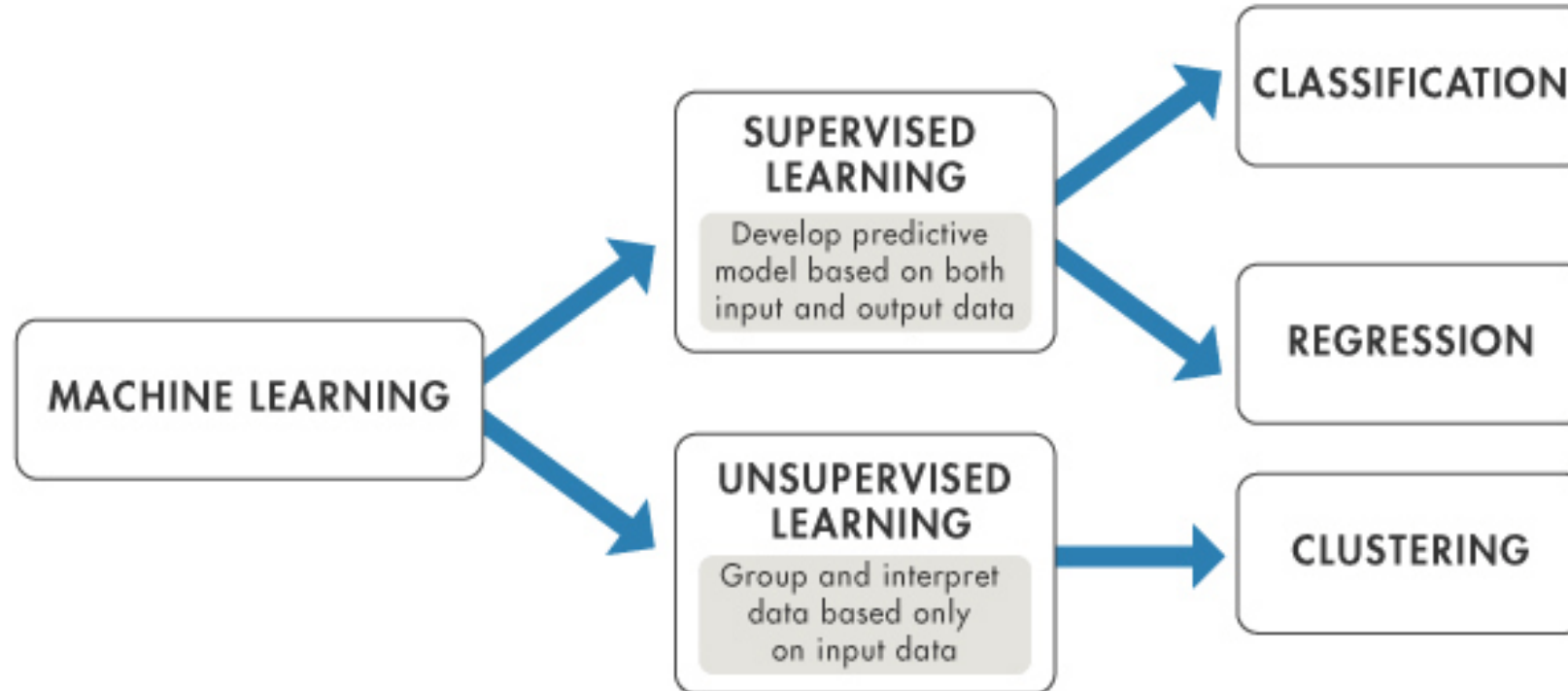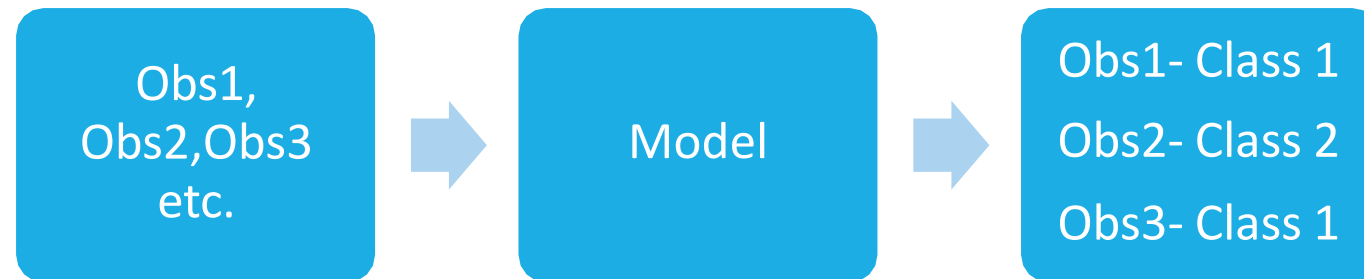
# Machine Learning algorithms

## Supervised Algorithms

▫ Given a set of variables $\underline{x}$, predict the value of another variable $y$ in a given data set such that

$$x1,x2,x3... \Rightarrow \text{Model } F(X) \Rightarrow y$$

▫ If y is numeric => **Prediction**

▫ If y is categorical => **Classification**

▫ **Example:** Given that a customer's Debt-to-Income ratio increased 20%, what are the chances he/she would default in 3 months?
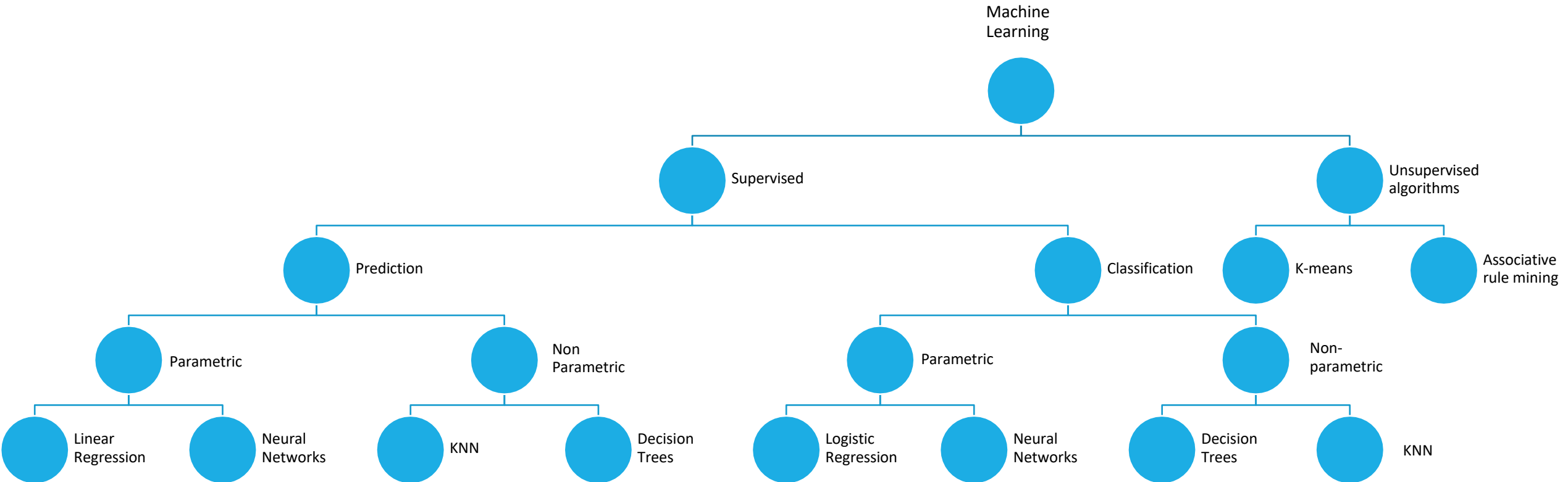
# Unsupervised Algorithms

▫ Given a dataset with variables $\underline{x}$, build a model that captures the similarities in different observations and assigns them to different buckets => **Clustering**



| Obs1, Obs2,Obs3 etc. | → | Model | → | Obs1- Class 1 Obs2- Class 2 Obs3- Class 1 |

▫ Example: Given a list of emerging market stocks, can we segment them into three buckets?

A mostly complete chart of
# Neural Networks

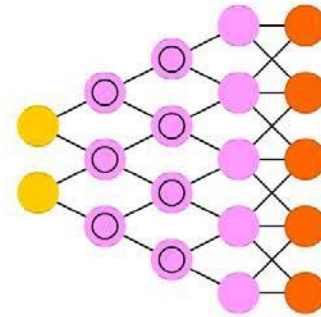http://www.asimovinstitute.org/neural-network-zoo/

```
┌──────────┐    ┌──────────┐    ┌──────────┐    ┌──────────┐    ┌──────────┐    ┌──────────┐
│   Data   │ →  │   Data   │ →  │ Feature  │ →  │ Training │ →  │  Model   │ →  │  Model   │
│ ingestion│    │ cleaning │    │engineering│   │and testing│   │ building │    │selection │
└──────────┘    └──────────┘    └──────────┘    └──────────┘    └──────────┘    └──────────┘
```

Split historical data into training and testing sets

# Performance evaluation

# Model performance

- Over fitting

- Cross validation

- Evaluation metrics

# Overfitting

Assume that the "true" *f* is given by the black curve. The others are three possible estimates for *f*. The orange line is the linear regression fit. The blue and green curves were produced using flexible methods. The green curve is the most flexible and it is the best one in matching the data; however, we observe that it fits the true *f* poorly.

# Overfitting

This situation is referred to as *overfitting*. This happens because the procedure is working too hard to find patterns in the training data, and may be picking up some patterns that are caused by random chance rather than by true properties of the unknown function. The picture shows overfitting when the true model is linear.

# Evaluation Metrics

## Regression

- Mean absolute error (MAE)
- Mean squared error (MSE)
- R squared ($R^2$)
- Adjusted R squared (Adj-$R^2$)

## Classification

- Accuracy
- Precision
- Recall
- Area under curve (AUC)
- Confusion matrix

# Model accuracy

In order to evaluate the performance of a statistical learning method on a given data set, we have to measure how accurately its predictions match the observed data. This is usually done by estimating the error on the training set.

- Linear regression – *Mean Squared Error* :

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

- Classification – Indicator function:

$$MSE = \frac{1}{n}\sum_{i=1}^{n} I(y_i \neq \hat{y}_i)$$

One is generally interested in the accuracy of the predictions obtained applying the method to previously unseen test data.

This theoretically corresponds to minimizing the average prediction error for large number of test observations (*average test MSE*), that could be not available.

It can be shown that the *expected test MSE*, that is the average test MSE that we would obtain if we repeatedly estimated $f$ using a large number of training sets, and tested each at a fixed point can be decomposed into the sum of three fundamental quantities:

1. variance, which refers to the amount by which the estimated function would change if we used a different training data set;

2. bias, which refers to the error that is introduced by approximating the "true" model by a simpler one;

3. irreducible error.

As a general rule, <u>as we use more flexible methods, the variance will increase and the bias will decrease</u>.

# Evaluation Metrics

For simplicity, we will mostly discuss things in terms of a binary classification problem; some common terms are:

- True positives (TP) $\rightarrow$ Predicted positive and are actually positive.

- False positives (FP) $\rightarrow$ Predicted positive and are actually negative.

- True negatives (TN) $\rightarrow$ Predicted negative and are actually negative.

- False negatives (FN) $\rightarrow$ Predicted negative and are actually positive.

$$\text{Precision} = \frac{\text{True positive}}{\text{Actual results}} \quad \text{or} \quad \frac{\text{True positive}}{\text{True positive + False positive}}$$

$$\text{Recall} = \frac{\text{True positive}}{\text{Predictive results}} \quad \text{or} \quad \frac{\text{True positive}}{\text{True positive + False negative}}$$

$$\text{Accuracy} = \frac{\text{True positive + True negative}}{\text{Total}}$$

# Confusion matrix

# Model selection

| | Linear regression | Logistic regression | SVM | CART | Gradient boosting | Random forest | Artificial neural network | KNN | LDA |
|---|---|---|---|---|---|---|---|---|---|
| Simplicity | ✔ | ✔ | ✔ | ✔ | ✘ | ✘ | ✘ | ✔ | ✔ |
| Training Time | ✔ | ✔ | ✘ | ✔ | ✘ | ✘ | ✘ | ✔ | ✔ |
| Handle non-linearity | ✘ | ✘ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |
| Robust to overfitting | ✘ | ✘ | ✔ | ✘ | ✘ | ✔ | ✘ | ✔ | ✘ |
| Large datasets | ✘ | ✘ | ✘ | ✔ | ✔ | ✔ | ✔ | ✘ | ✔ |
| Many features | ✘ | ✘ | ✔ | ✔ | ✔ | ✔ | ✔ | ✘ | ✔ |
| Model interpretation | ✔ | ✔ | ✘ | ✔ | ✔ | ✔ | ✘ | ✔ | ✔ |
| Feature scaling needed | ✘ | ✘ | ✔ | ✘ | ✘ | ✘ | ✘ | ✘ | ✘ |

# Main Python libraries for Data Science

- **NumPy**: storage and manipulation of dense data arrays.

- **Pandas**: DataFrame object for storage and manipulation of labeled/columnar data.

- **Matplotlib**: capabilities for a flexible range of data visualizations.

- **Scikit-learn**: machine learning.

# Data cleaning

# Data preparation



Data Collection & Storage → **Data Preparation** → Exploration & Visualization → Experimentation & Prediction

## Why prepare data?

- Preparation is done to prevent:

- Errors

- Incorrect results

- Biasing algorithms

# Tidy Data

| | Sara | Lis | Hadrien | Lis |
|---|---|---|---|---|
| Age | "27" | "30" | | "30" |
| Size | 1.77 | 5.58 | 1.80 | 5.58 |
| Country | "Belgium" | "USA" | "FR" | "USA" |

| Name | Age | Size | Country |
|---|---|---|---|
| Sara | "26" | 1.78 | "Belgium" |
| Lis | "30" | 5.58 | "USA" |
| Hadrien | | 1.80 | "FR" |
| Lis | "30" | 5.58 | "USA" |

# Tidy Data

| Name | Age | Size | Country |
|---|---|---|---|
| Sara | "26" | 1.78 | "Belgium" |
| Lis | "30" | 5.58 | "USA" |
| Hadrien | | 1.80 | "FR" |
| Lis | "30" | 5.58 | "USA" |

| Name | Age | Size | Country |
|---|---|---|---|
| Sara | "27" | 1.77 | "Belgium" |
| Lis | "30" | 5.58 | "USA" |
| Hadrien | | 1.80 | "FR" |

# Tidy Data

| Name | Age | Size | Country |
|---|---|---|---|
| Sara | "27" | 1.77 | "Belgium" |
| Lis | "30" | 5.58 | "USA" |
| Hadrien | | 1.80 | "FR" |

| ID | Name | Age | Size | Country |
|---|---|---|---|---|
| 0 | Sara | "27" | 1.77 | "Belgium" |
| 1 | Lis | "30" | 5.58 | "USA" |
| 2 | Hadrien | | 1.80 | "FR" |

# Homogeneity

| ID | Name | Age | Size | Country |
|----|------|-----|------|---------|
| 0 | Sara | "27" | 1.77 | "Belgium" |
| 1 | Lis | "30" | 5.58 | "USA" |
| 2 | Hadrien | | 1.80 | "FR" |

| ID | Name | Age | Size | Country |
|----|------|-----|------|---------|
| 0 | Sara | "27" | 1.77 | "Belgium" |
| 1 | Lis | "30" | 1.70 | "USA" |
| 2 | Hadrien | | 1.80 | "FR" |

# Homogeneity

| ID | Name | Age | Size | Country |
|----|------|-----|------|---------|
| 0 | Sara | "27" | 1.77 | "Belgium" |
| 1 | Lis | "30" | 1.70 | "USA" |
| 2 | Hadrien | | 1.80 | "FR" |

| ID | Name | Age | Size | Country |
|----|------|-----|------|---------|
| 0 | Sara | "27" | 1.77 | "BE" |
| 1 | Lis | "30" | 1.70 | "US" |
| 2 | Hadrien | | 1.80 | "FR" |

# Data types

| ID | Name | Age | Size | Country |
|----|------|------|------|---------|
| 0 | Sara | "27" | 1.77 | "BE" |
| 1 | Lis | "30" | 1.70 | "US" |
| 2 | Hadrien | | 1.80 | "FR" |

| ID | Name | Age | Size | Country |
|----|------|------|------|---------|
| 0 | Sara | 27 | 1.77 | "BE" |
| 1 | Lis | 30 | 1.70 | "US" |
| 2 | Hadrien | | 1.80 | "FR" |

# Data types

| ID | Name | Age | Size | Country |
|----|------|-----|------|---------|
| 0 | Sara | 27 | 1.77 | "BE" |
| 1 | Lis | 30 | 1.70 | "US" |
| 2 | Hadrien | | 1.80 | "FR" |

| ID | Name | Age | Size | Country |
|----|------|-----|------|---------|
| 0 | Sara | 27 | 1.77 | "BE" |
| 1 | Lis | 30 | 1.70 | "US" |
| 2 | Hadrien | 28 | 1.80 | "FR" |

## Missing values

**Reasons:**
- Data entry
- Error
- Valid missing value

**Solutions:**
- impute
- drop
- keep