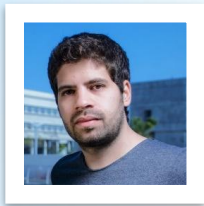


# Introduction to Machine Learning

Lior Sidi & Noa Lubin

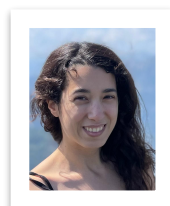


# Meet us



## Lior Sidi

- Data Science Tech Lead at Wix
- **Ex:** AutoML startup, Consultant, Deutsche Telekom Labs.
- **Domains:** NLP, RecSys, Design, User Centric AI (B2C products).
- BSc and MSc from BGU SISE.
- [Scholar](#) / [Linkedin](#)



## Noa Lubin

- Machine Learning Team Lead at Diagnostic Robotics
- **Ex:** NASA, Amazon, Elbit, IAI
- **Domains:** NLP, health, space
- BSc EE Technion
- MSc CS NLP Bar Ilan
- [Linkedin](#)



*“If you invent a breakthrough in AI, so machines can learn, that is worth 10 Microsofts”*

— Bill Gates, Former Chairman, Microsoft. 2004

# General Outline

1. [What is ML?](#)
2. [What are the types of ML Problems?](#)
3. [How is ML in Practice?](#)
4. [How to estimate Model Performance?](#)
5. [How to prepare data for ML?](#)
6. [What type of ML algorithms are there?](#)
7. [How to improve ML models?](#)

# 1. What is ML?

# How is Machine Learning Relates to Programming

## Traditional Programming



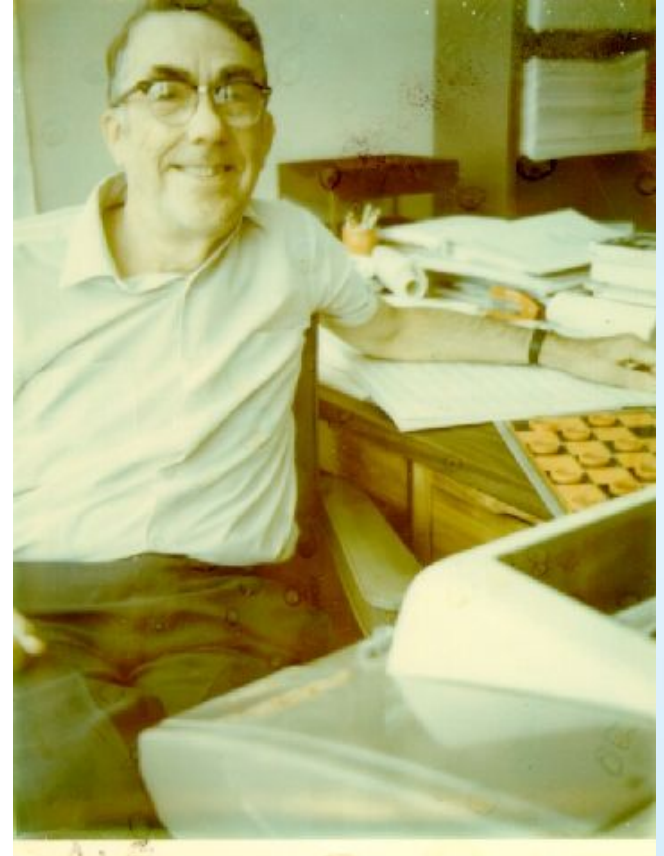
- Machines Follow Instruction
- Humans Learn From Experience

# The Arthur Samuel's Checkers

“Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.”

— Arthur Samuel (1959)

***Machine Learning is the ability to generalize from experience onto unseen example***



# How is Machine Learning Relates to Programming

## Traditional Programming



## Machine Learning





# Agriculture Metaphor of ML

“Machine learning is like farming or gardening. Seeds is the algorithms, nutrients is the data, the gardener is you and plants is the programs.”

[A Study on Machine Learning: Overviews and Applications International conference on Recent Trends in Artificial Intelligence, IOT, Smart Cities & Applications \(ICAISC-2020\)](#)

**Seeds** = Algorithms

**Nutrients** = Data

**Gardener** = You

**Plants** = Programs



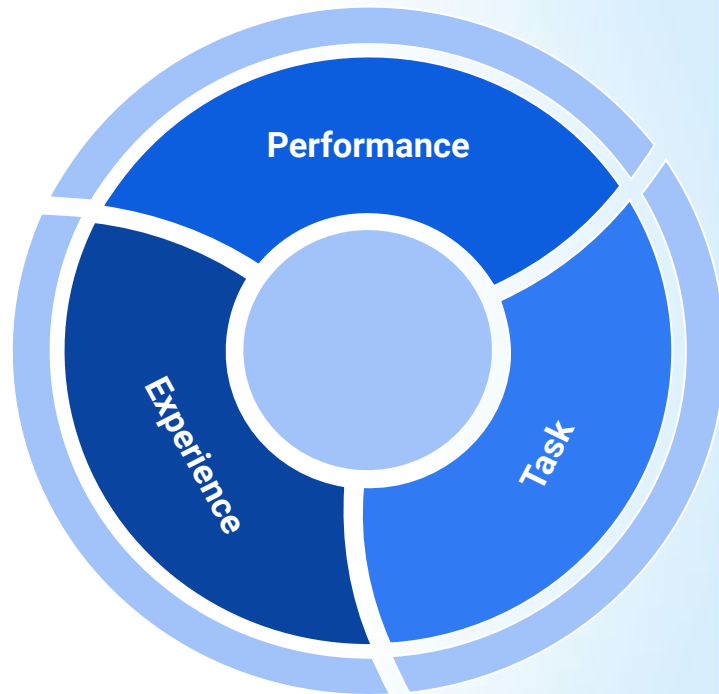
# The PTE View on ML

Machine Learning is the study of algorithms that

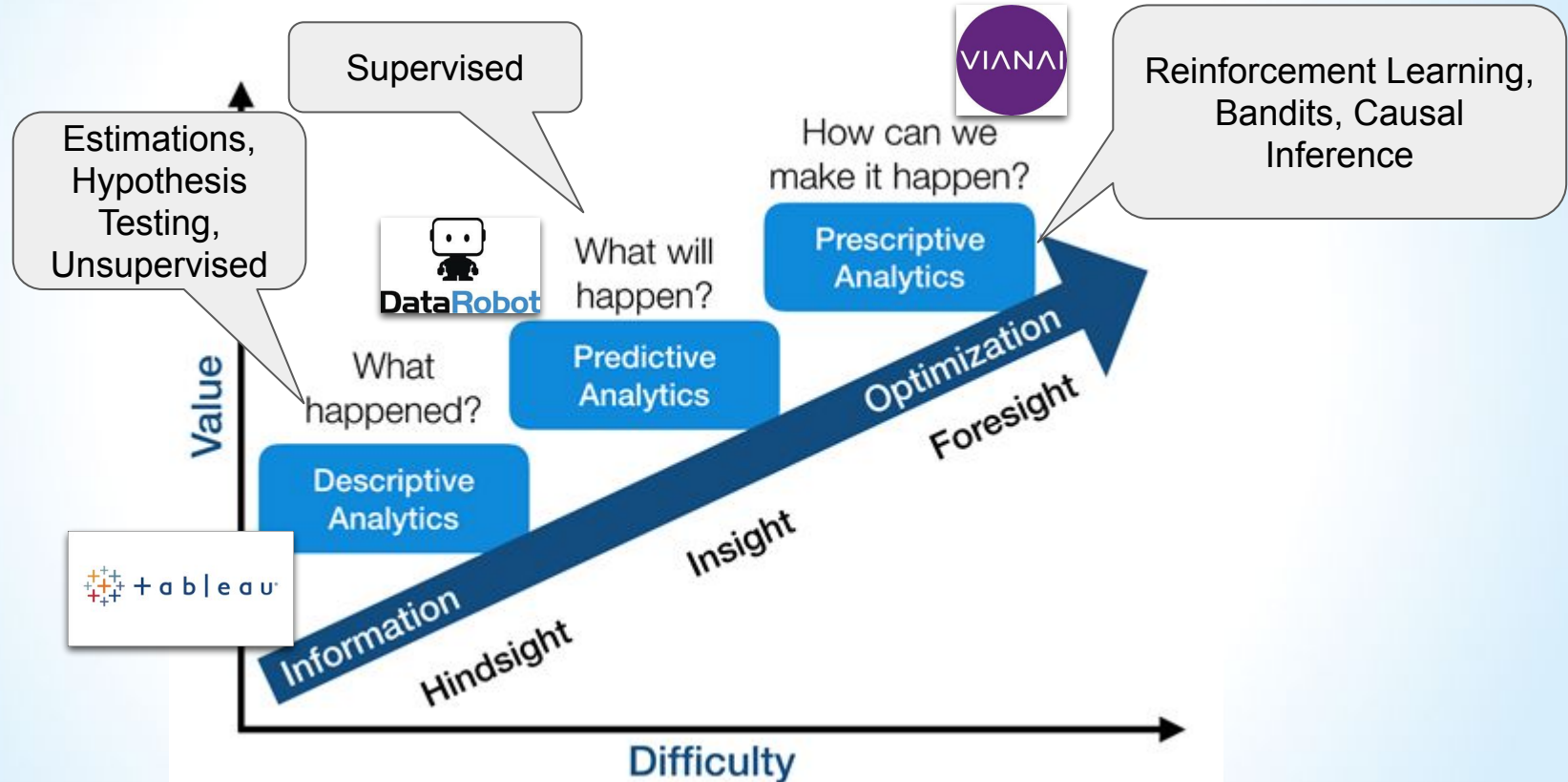
- improve their **performance** (P)
- at some **task** (T)
- with **experience** (E).

A well-defined learning task is given by  $\langle P, T, E \rangle$ .

e.g., predict house pricing (T), according to last year's data (E), minimizing the square error (P).

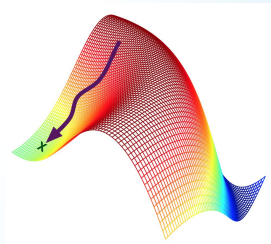
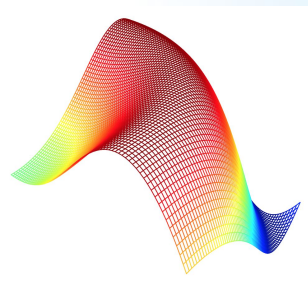
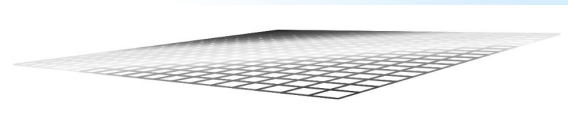


# Descriptive vs. Predictive vs. Prescriptive



# Key Elements of Machine Learning

1. **Representation**: how to represent knowledge.  
**The Model.**
2. **Evaluation**: the way to evaluate candidate representations (programs/hypotheses).
3. **Optimization**: how to find the optimal representation in the evaluation terms.



# Application of Machine Learning

Sample applications of machine learning:

- The
- List
- Is
- Just
- Too
- Long

## Artificial Intelligence is the New Electricity — Andrew Ng



Synced **Following**  
Apr 28, 2017 · 6 min read



### Abstract

On Wednesday, January 25, Andrew Ng — former Baidu Chief Scientist, Coursera co-founder, and Stanford Adjunct Professor — gave a talk at the Stanford MSx Future Forum. During the talk, Professor Ng shared his opinion on AI. He mainly discussed how artificial intelligence (AI) is transforming industry and business.

[Machine Learning is the new Electricity](#)

## 2. What are the types of ML Problems?

# Types of ML

|                        |   |
|------------------------|---|
| Supervised Learning    | <ul style="list-style-type: none"><li>&gt; Labeled data</li><li>&gt; Direct feedback</li><li>&gt; Predict outcome/future</li></ul>    |
| Unsupervised Learning  | <ul style="list-style-type: none"><li>&gt; No labels</li><li>&gt; No feedback</li><li>&gt; Find hidden structure in data</li></ul>    |
| Reinforcement Learning | <ul style="list-style-type: none"><li>&gt; Decision process</li><li>&gt; Reward system</li><li>&gt; Learn series of actions</li></ul> |

# Additional Types of Learning

- **Semi-Supervised learning**

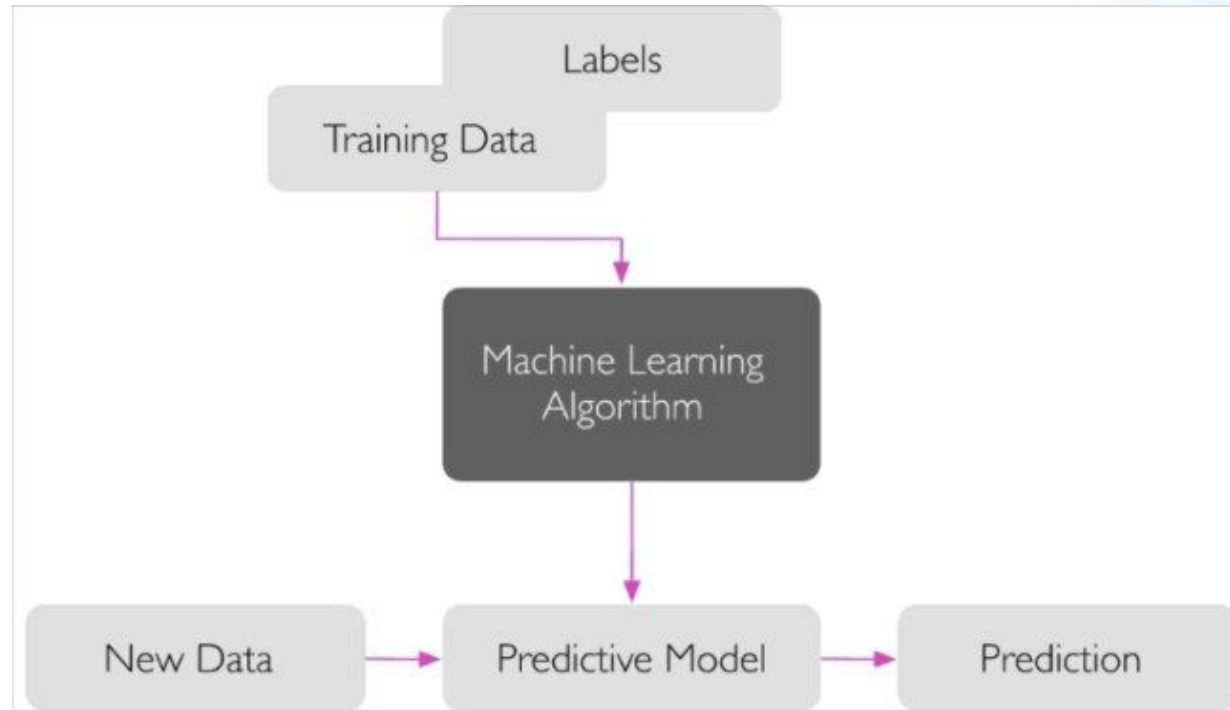
Training data includes a few desired outputs.

- **Causal Inference Learning**

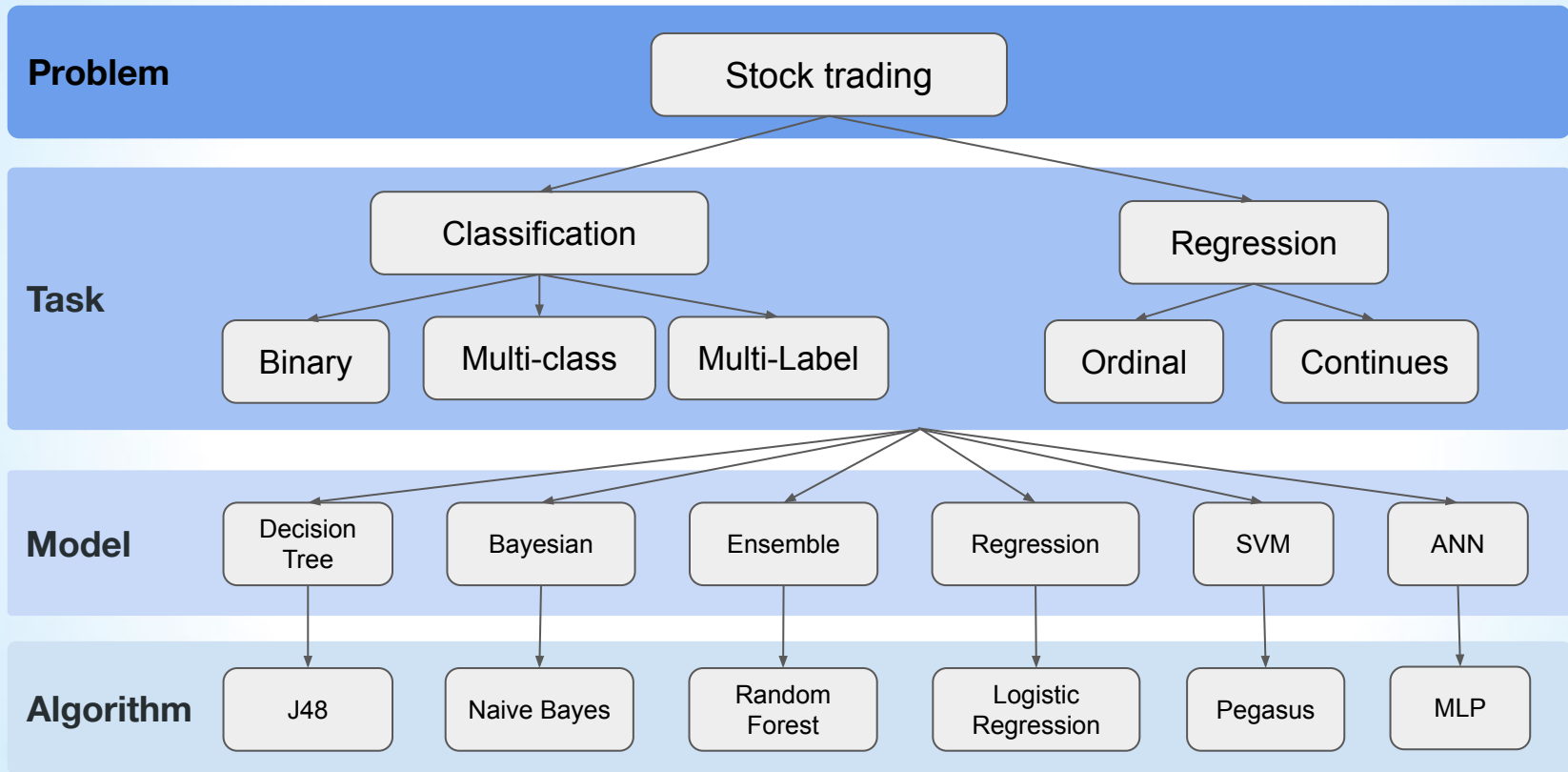
Discover and estimate the causal relationship between variables



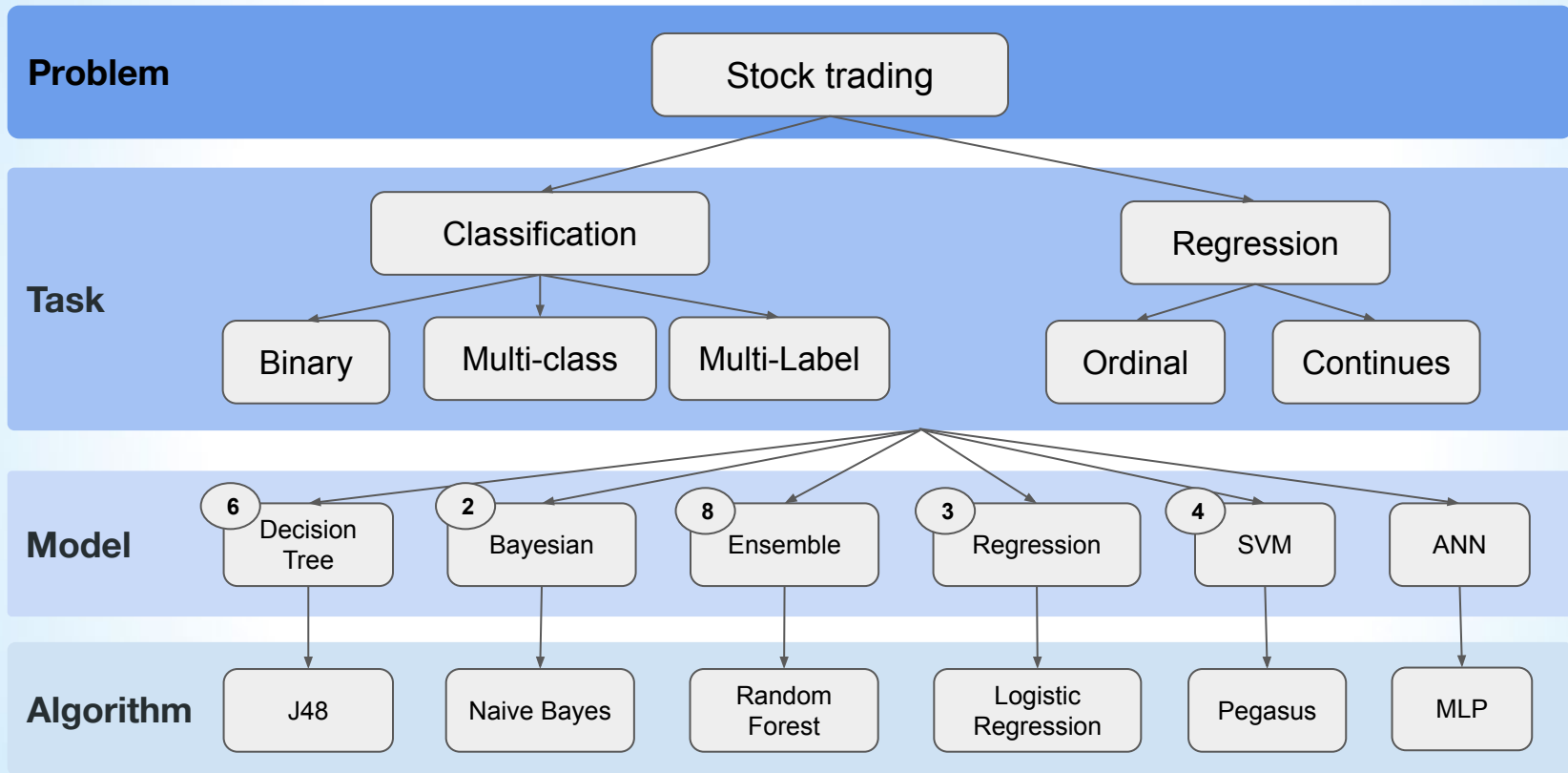
# Supervised



# Supervised tasks



# Supervised tasks



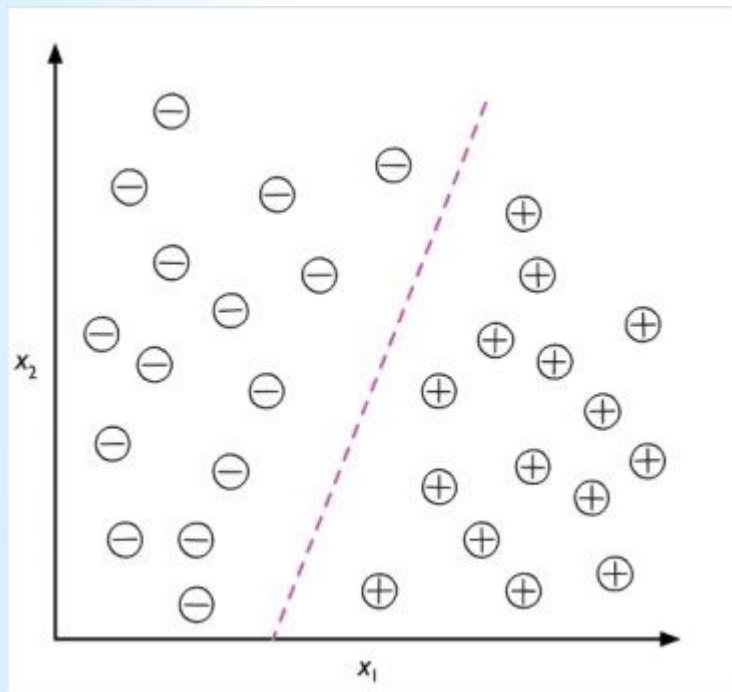
# Syllabus

| week                                   | Topics   |
|--|--|
| <b>1</b><br><b>Intro</b>               | Introduction, definition of ML, types of ML, KNN   |
| <b>2</b><br><b>Naive Bayes</b>         | <ul style="list-style-type: none"> <li>-Refresher on the Bayes theorem</li> <li>-Naive Bayes theory</li> <li>-Spam detection use case including preprocessing phase</li> </ul>   |
| <b>3</b><br><b>Linear Regression</b>   | <ul style="list-style-type: none"> <li>- Ordinary Linear Regression - Analytical solution and Gradient Descent solution</li> <li>- Common Pitfalls - weight is not importance</li> <li>- L2 regularization on OLS</li> <li>- L1 regularization on OLS</li> </ul> |
| <b>4</b><br><b>Logistic Regression</b> | <ul style="list-style-type: none"> <li>- Regression vs Classification</li> <li>- Classifier performance measurement</li> <li>- Binary Logistic Regression</li> <li>- Multi-Class classification</li> </ul>   |

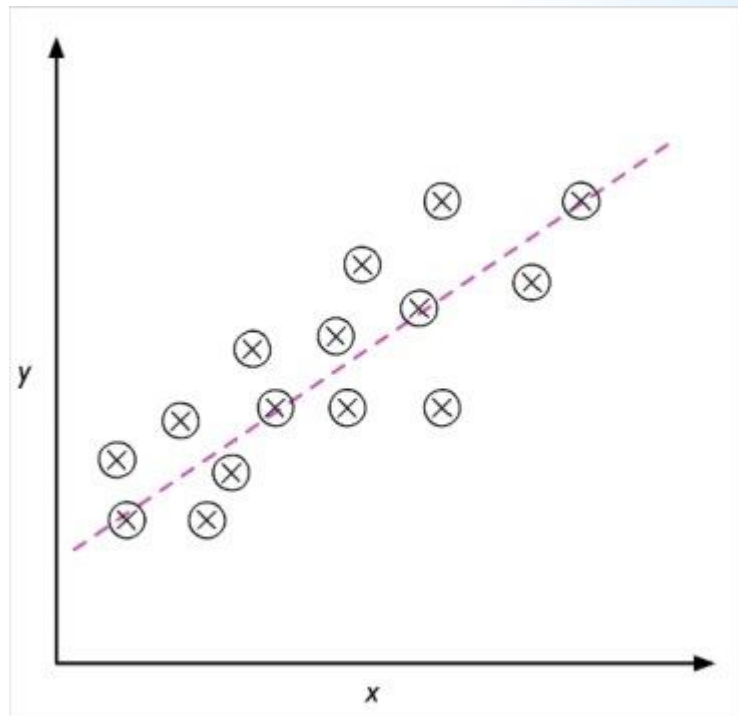
| week                              | Topics   |
|-----------------------------------|--|
| <b>5</b><br><b>SVM</b>            | - Support vector machines  |
| <b>6</b><br><b>Decision trees</b> | <ul style="list-style-type: none"> <li>- Decision Tree as a Greedy Method</li> <li>- Optimization Criteria: Gini &amp; Entropy+ L2</li> <li>- Depth, Leaves and other Hyper Parameters</li> </ul>                                      |
| <b>7</b><br><b>End2End ML</b>     | <ul style="list-style-type: none"> <li>- Bias-variance tradeoff, validation set, cross-validation</li> <li>- Overfitting and underfitting</li> <li>- Regularization and hyperparameter tuning</li> <li>- Features selection</li> </ul> |
| <b>8-9</b><br><b>Ensemble</b>     | Intro to Ensemble Methods <ul style="list-style-type: none"> <li>- Aggregation</li> <li>- Bagging</li> <li>- Stacking</li> <li>- Boosting</li> <li>- Gradient Boosting</li> </ul>  |

# Regression Vs Classification

Classification

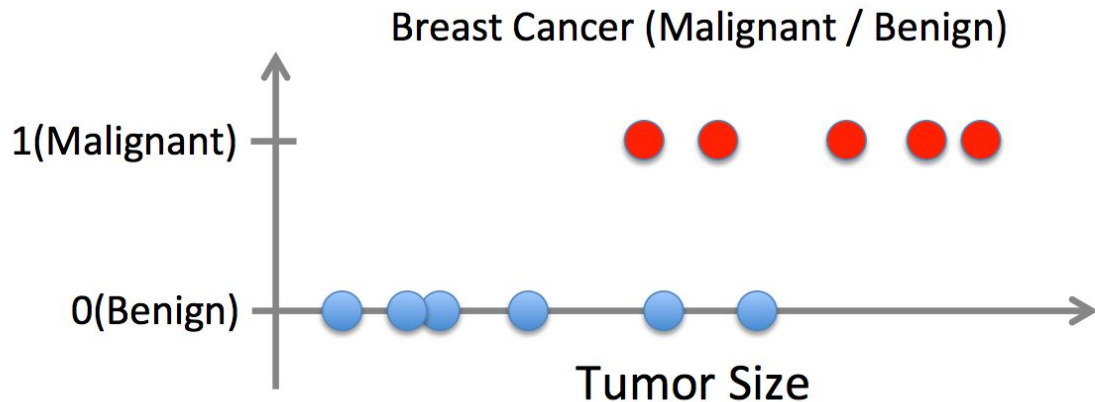


Regression



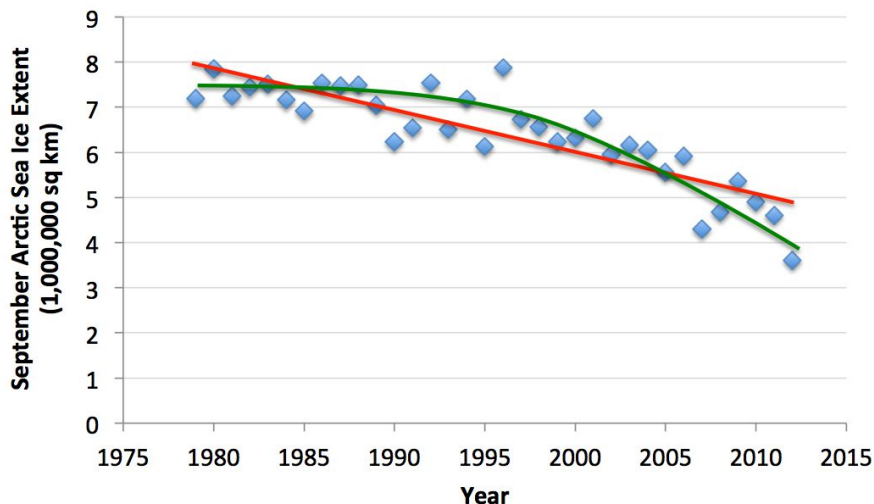
# Supervised Classification

- Given  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Learn a function  $f(x)$  to predict  $y$  given  $x$ 
  - $y$  is categorical == classification



# Supervised Regression

- Given  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Learn a function  $f(x)$  to predict  $y$  given  $x$ 
  - $y$  is real-valued == regression



# Type of Supervised Learning

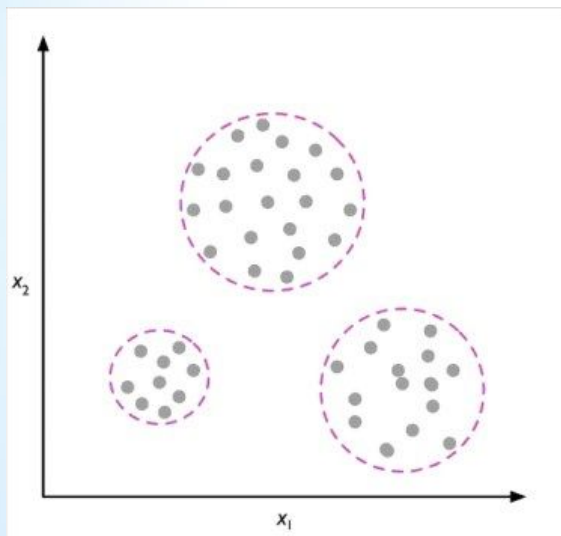
- Output Data Type
  - Discrete - Classification
  - Continuous - Regression
  - Structured - Structure Prediction (e.g. sentence tagging)
- Classification Types:
  - Multi-class
  - Multi-label
- Regression Types:
  - Dimensions - Uni/Multivariate
  - Monotonic - Isotonic Regression
  - Loss Function - MSE, MAE
  - Regularization Type - Lasso, Ridge, etc
  - Percentile - Quantile Regression

|          |     |
|----------|-----|
| This     | DT  |
| is       | VBZ |
| a        | DT  |
| tagged   | JJ  |
| sentence | NN  |

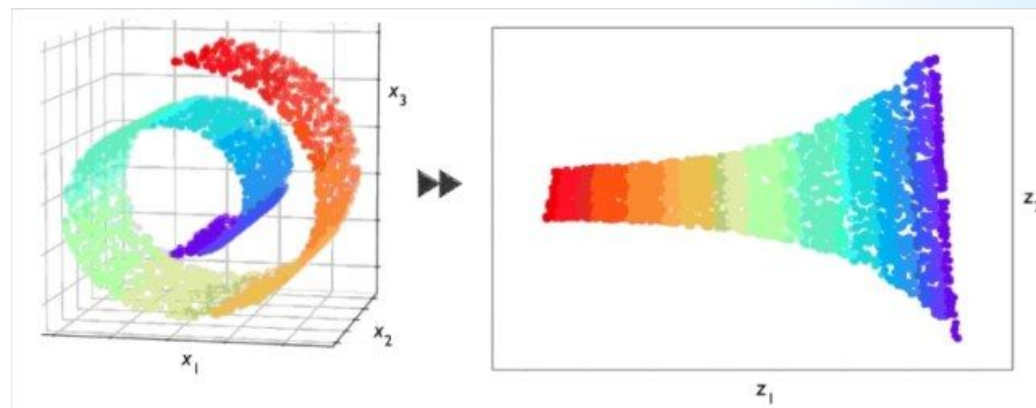


# Unsupervised Learning

Clustering

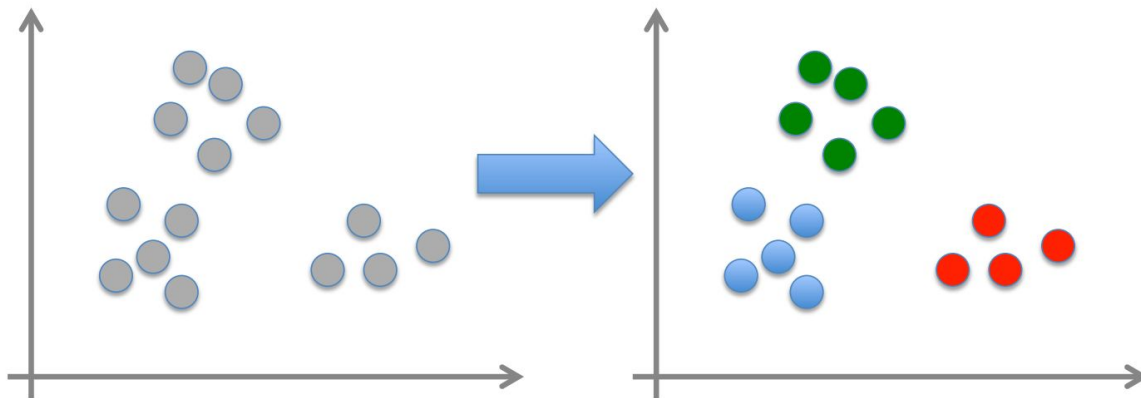


Dimension Reduction

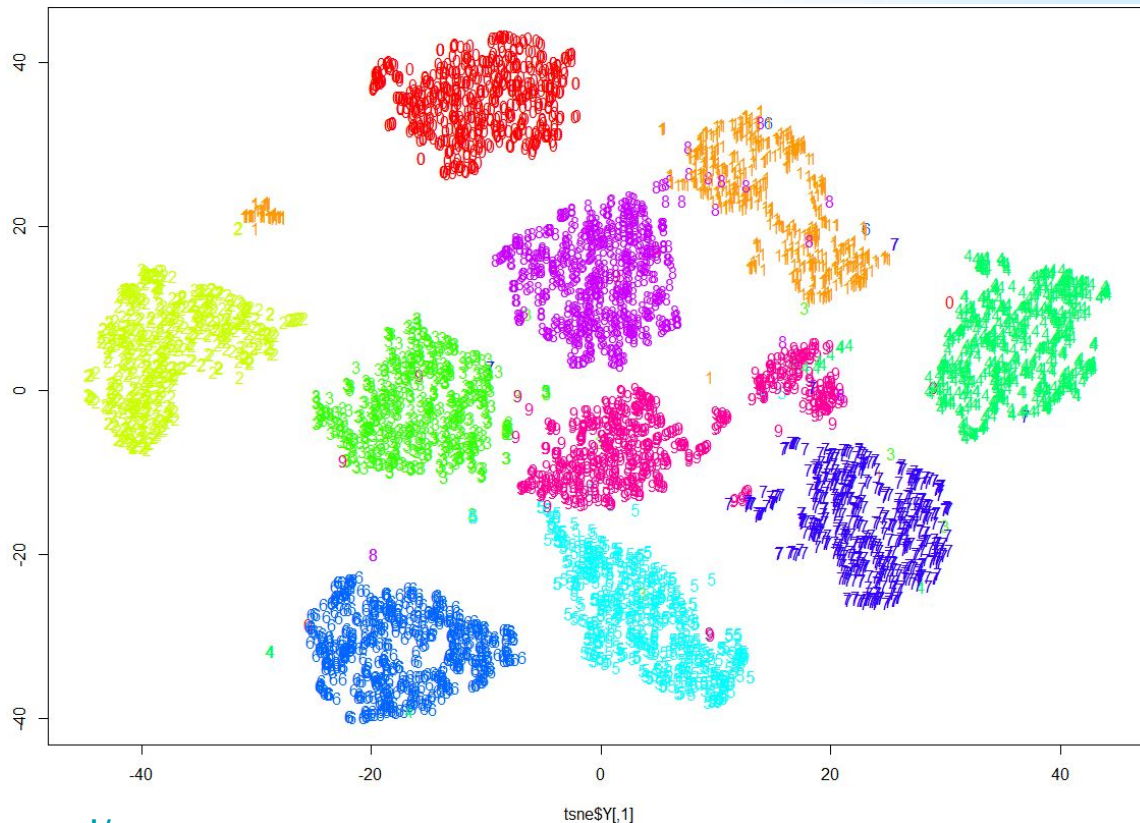


# Unsupervised

- Given  $x_1, x_2, \dots, x_n$  (without labels)
- Output hidden structure behind the  $x$ 's
  - E.g., clustering



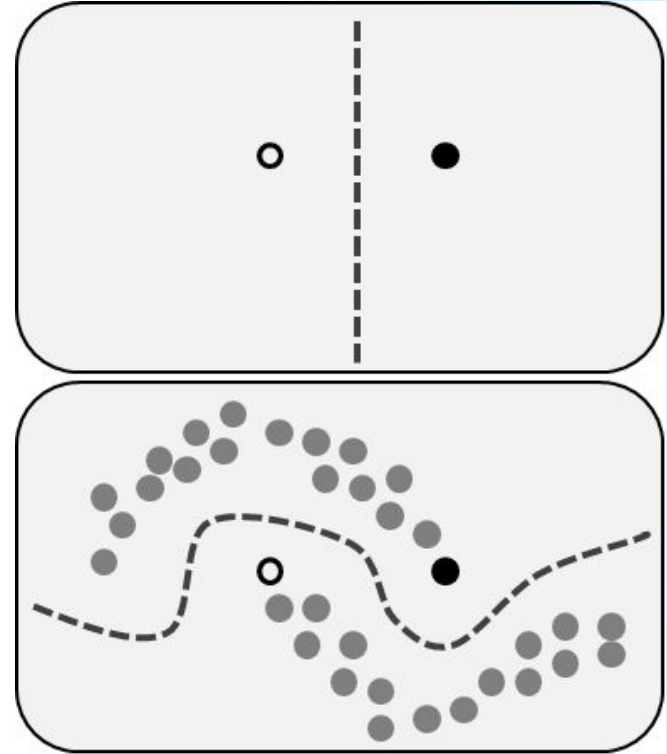
# Unsupervised Example



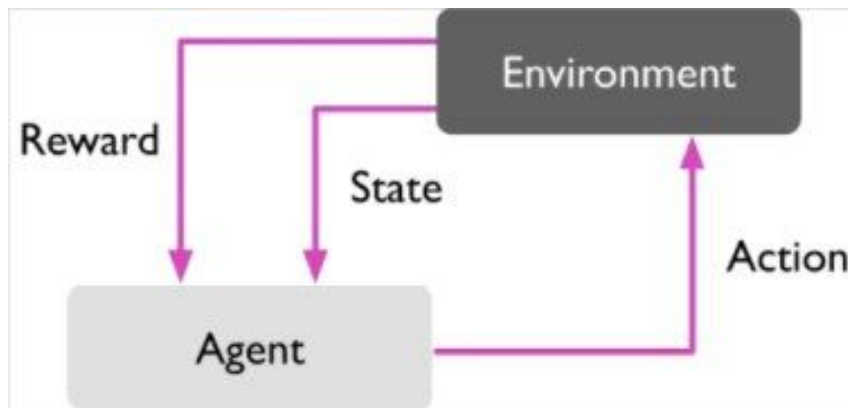
# Semi-supervised Learning Example

Antivirus:

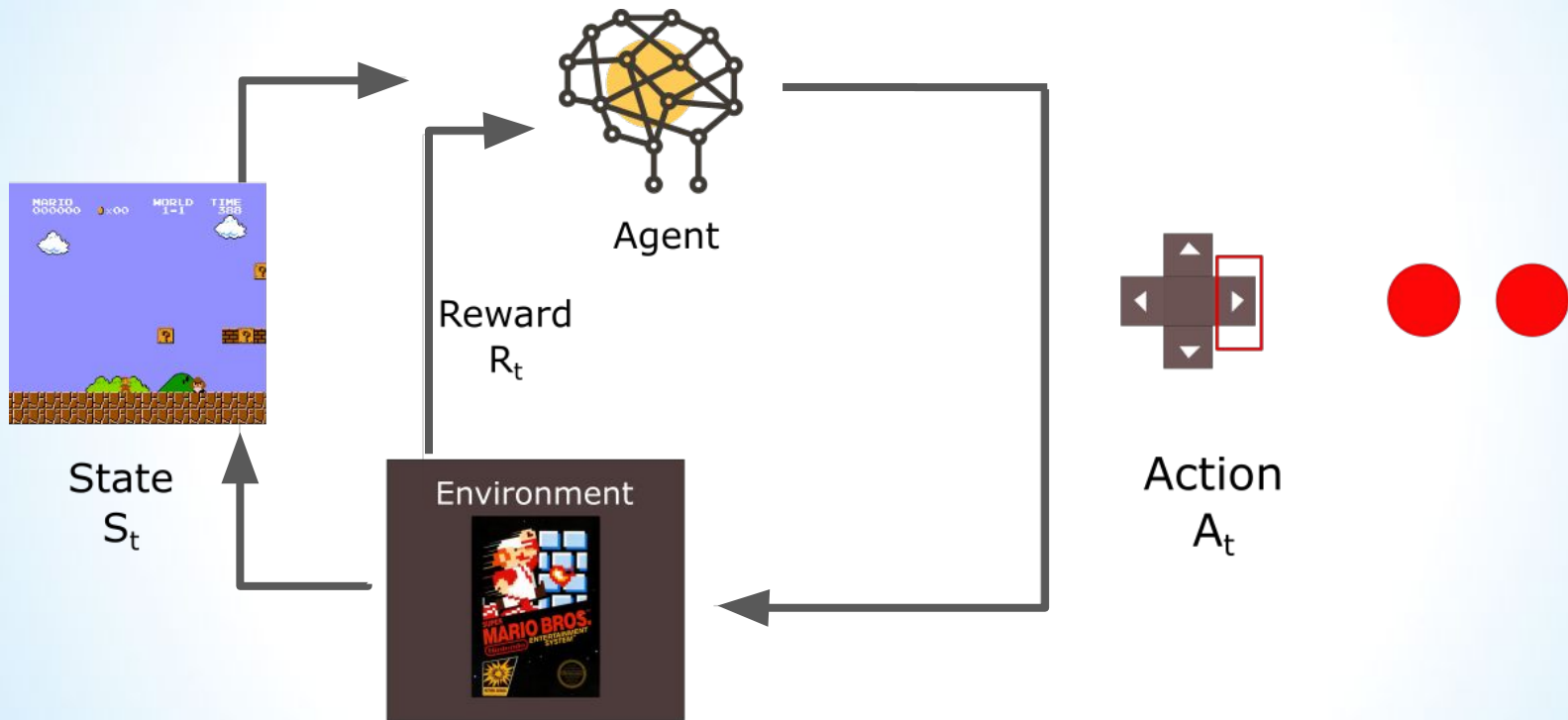
- Some files are known to be safe
- Some are known to be malicious
- Most are unknown



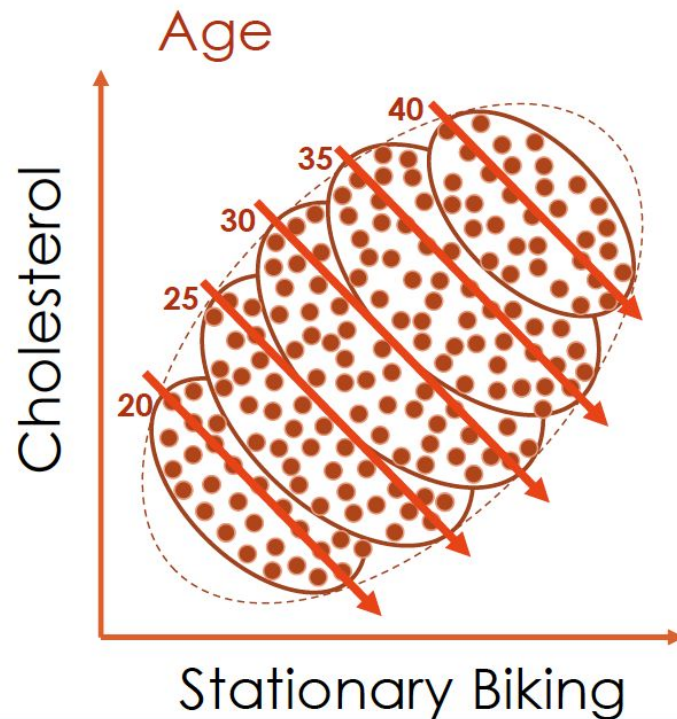
# Reinforcement Learning



# RL Example



# Causal Inference



### 3. How is ML in Practice?

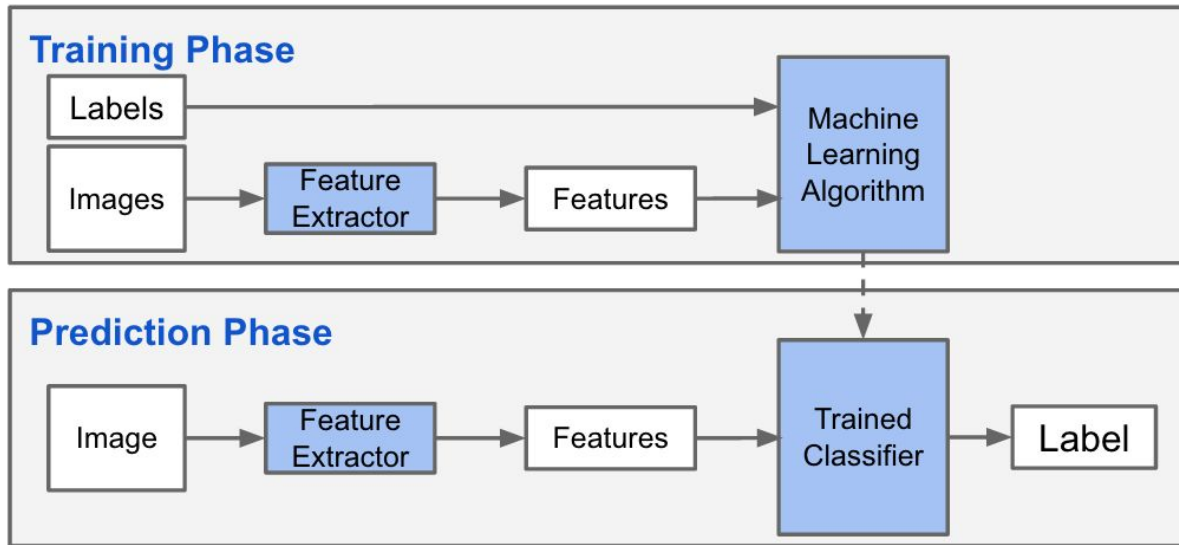


THE #1 DEEP LEARNING EXCUSE  
FOR LEGITIMATELY SLACKING OFF:

"MY MODEL IS TRAINING."



# Train vs Test Pipelines



Machine Learning Phases

# ML $\in$ DS

1. Understand the domain and explore the data
2. Data integration, selection, cleaning and pre-processing.
3. Learning models  $\leftarrow$  ML part
4. Interpreting results.
5. Consolidating and deploying discovered knowledge.



It is not a one-shot process, it is a cycle. You need to run the loop until you get a result that you can use in practice. Also, the data can change, requiring a new loop.

# Hidden Technical Debt in Machine Learning Systems

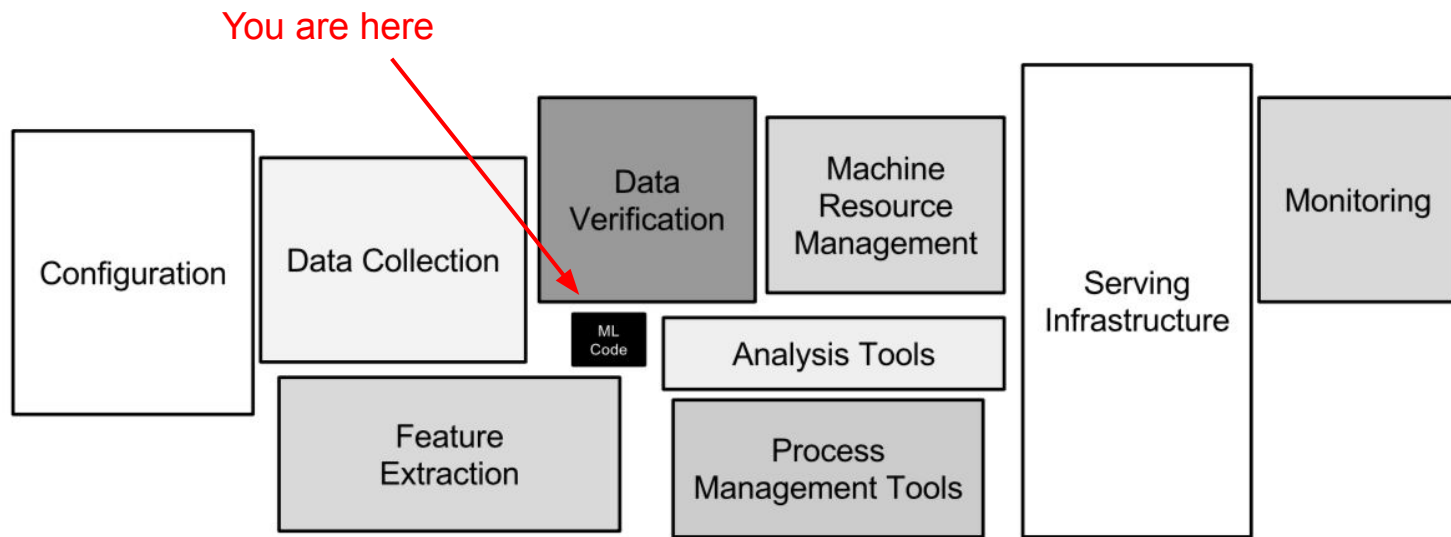
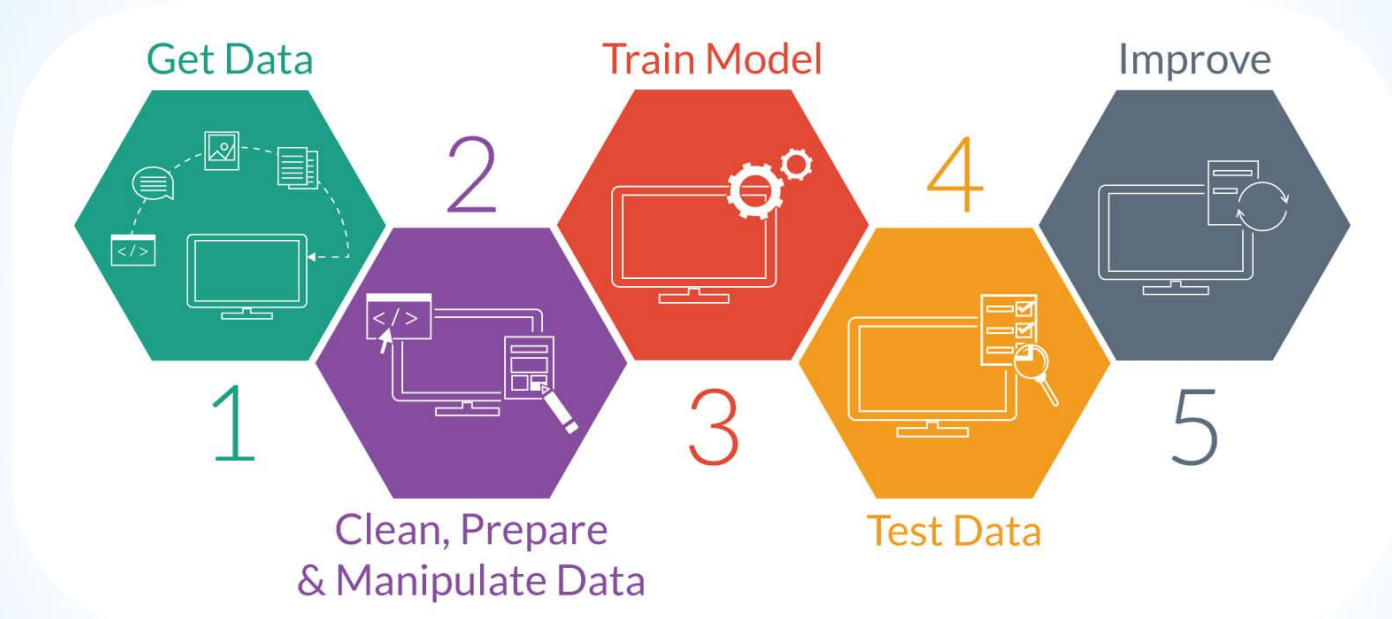


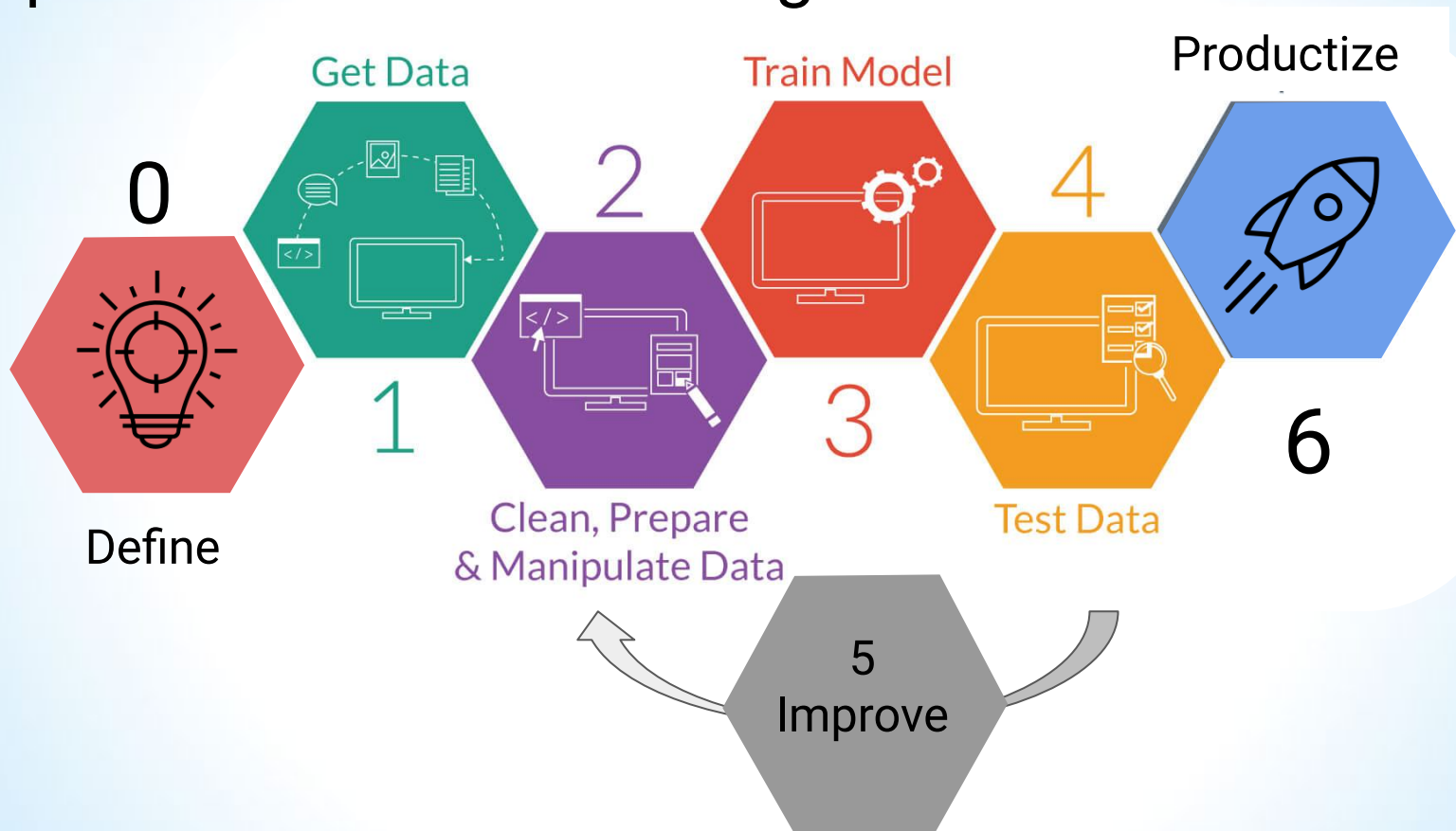
Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

## IV. DS project management

# Steps to Predictive Modeling



# Steps to Predictive Modeling



# “Tachles” - In Reality

1. Build a Baseline model as soon as possible!
  - a. Understand the problem.
  - b. Have a SIMPLE data validation pipeline.
  - c. Train the SIMPLEST model.
2. Iteratively improve the model
  - a. Incremental - add features, tune models, clean data etc.
  - b. Go Wild - re-model the data, add sophisticated features, use SOTA approaches.



# Task

Given a **KPI** define a feature in the product.

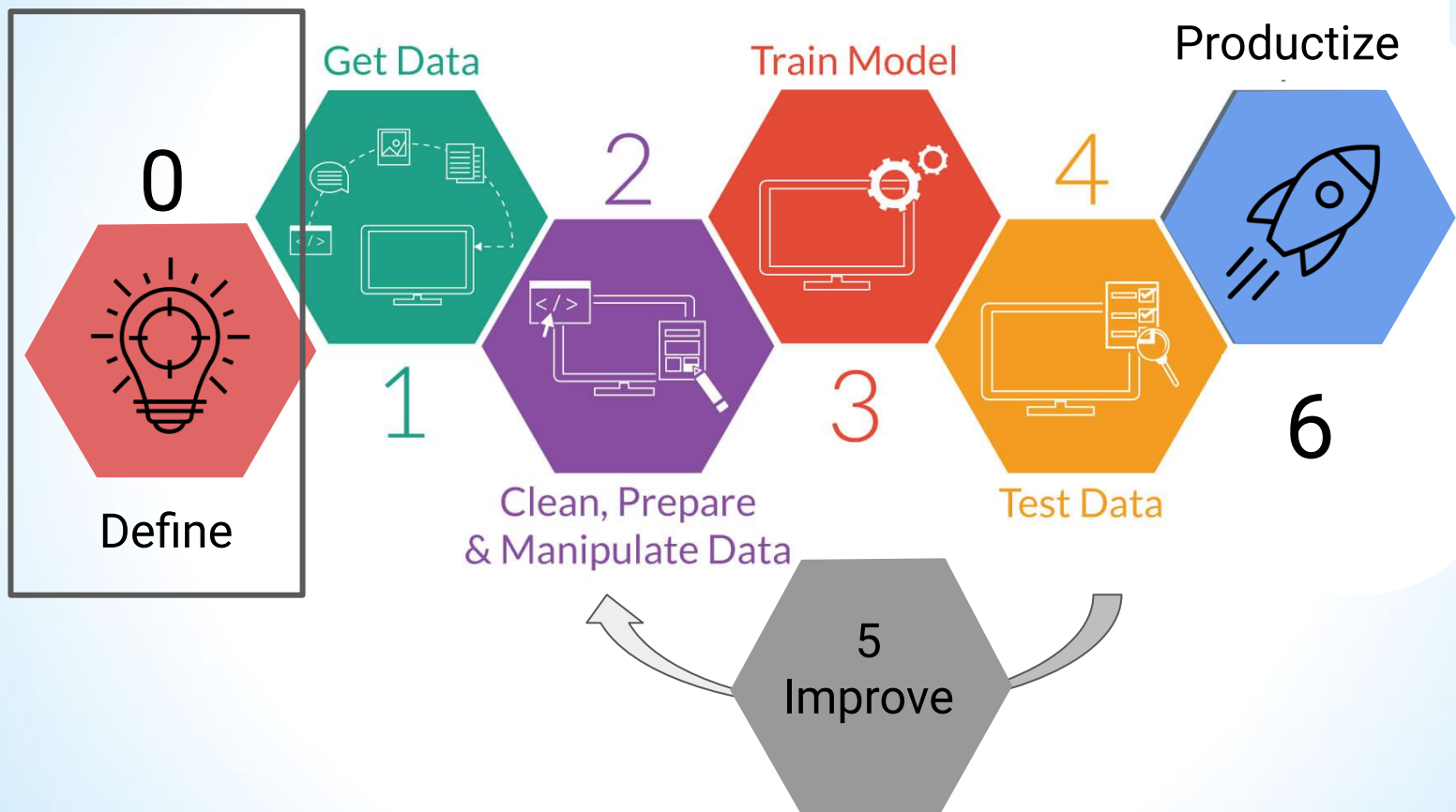
## **Transaction made weekly by user**

1. Stock recommendation - Which Stock should I buy today?
2. Decision support - Should I buy Apple stock?

## **Hours spend daily by customer**

1. Decision support - What will be the Microsoft stock price in the next 20 days?
2. Exploration - Which stock is similar to Amazon stock?

# Steps to Predictive Modeling



# Define Task

## Which Stock should I buy today?

- Story: Beginning of day a customer will get a list of stocks

## Steps

1. Literature overview
2. Data source
3. Data modeling: Entities
4. Labeling function: Rules, manual, data source

# Define Data

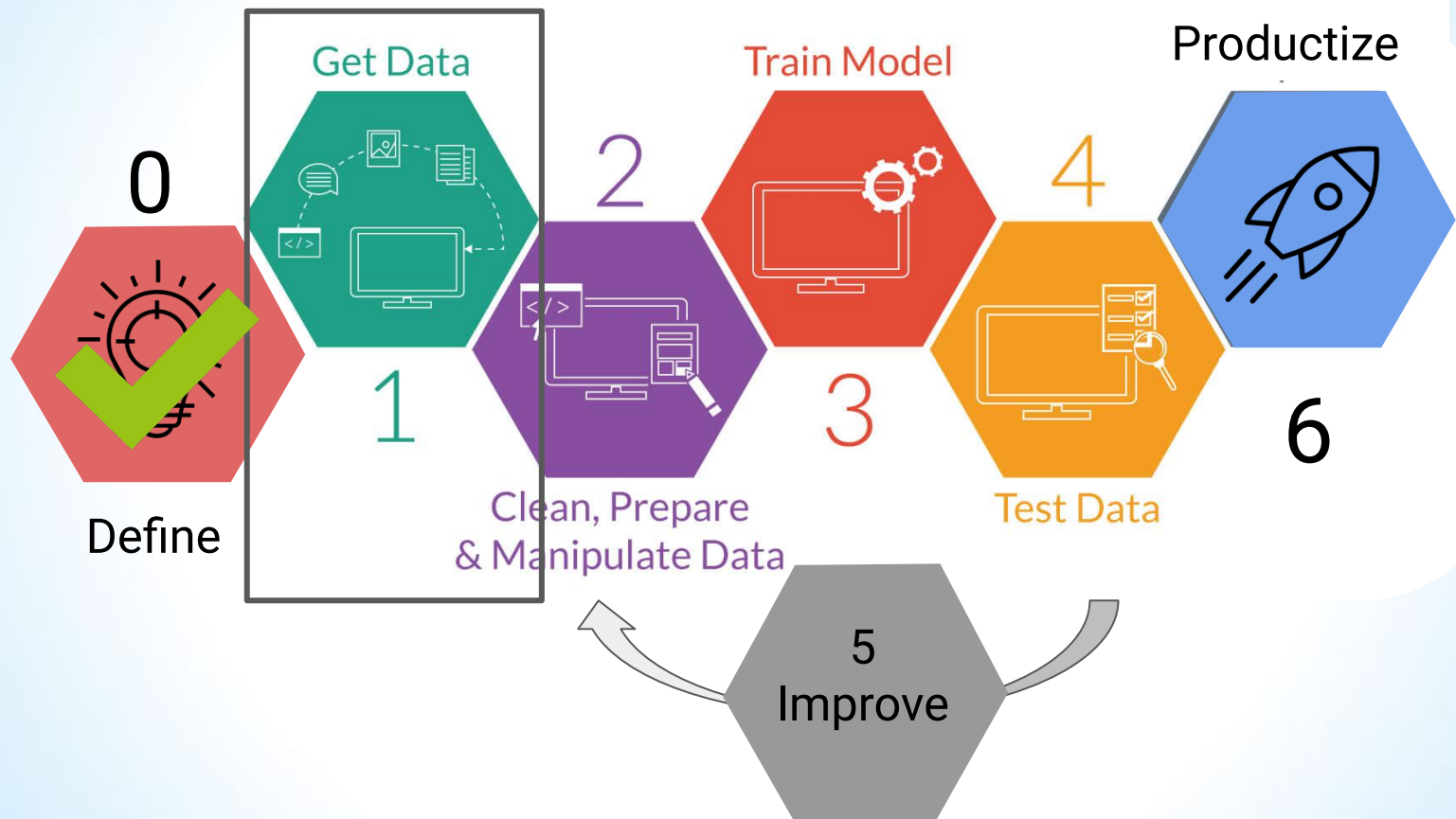
## Which Stock should I buy today?

- Story: Beginning of day a customer will get a list of stocks

## Data sources - based on “what” we make a decision?

- **Data** - All stock of S&P 500 - after 2009 and before 2020.
- **Entity** - stock
- **Horizon** - next Day

# Steps to Predictive Modeling



# Ground Truth

## Which Stock should I buy today?

- Story: Beginning of day a customer will get a list of stocks

## Data sources - based on “what” we make a decision?

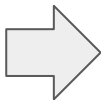
- **Data** - All stock of S&P 500 - after 2009 and before 2020.
- **Entity** - stock
- **Horizon** - next Day

## Labeling - What is a good stock?

- A stock with low volatility in the next day -> low Sharp ratio
- A stock with high revenue potential -> 4% price increment

# Labeling function

| Date     | Stock | Price |
|----------|-------|-------|
| 1/1/2020 | TSLA  | 50    |
| 1/1/2020 | AMZ   | 63    |
| 1/1/2020 | APPL  | 42    |
| 2/1/2020 | TSLA  | 55    |
| 2/1/2020 | AMZ   | 60    |
| 2/1/2020 | APPL  | 39    |
| 3/1/2020 | TSLA  | 60    |

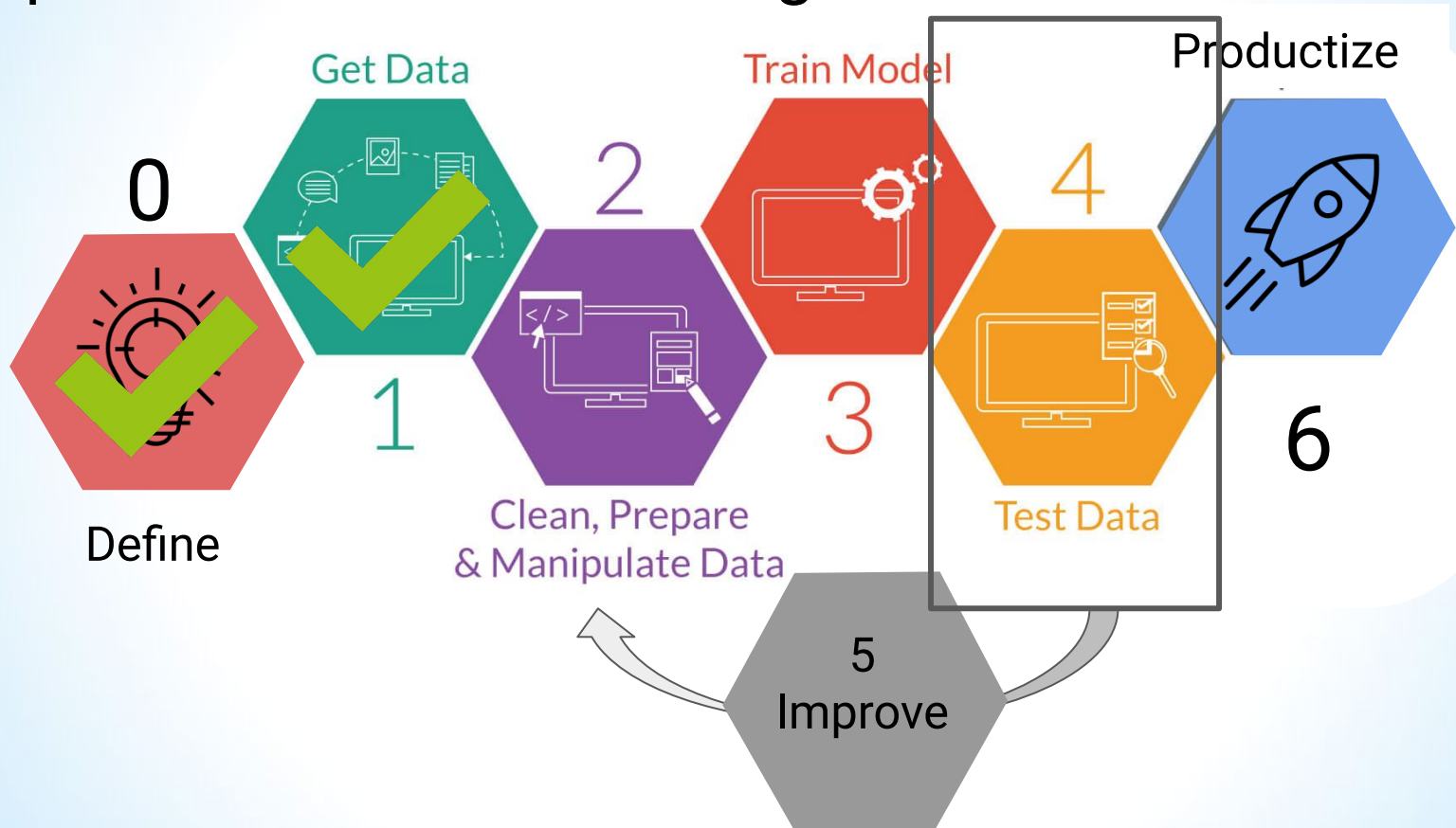


| Date     | Stock | Price | Volatile | Profit | Label |
|----------|-------|-------|----------|--------|-------|
| 1/1/2020 | TSLA  | 50    | Low      | High   | 1     |
| 1/1/2020 | AMZ   | 63    | Low      | Low    | 0     |
| 1/1/2020 | APPL  | 42    | High     | Low    | 0     |
| 2/1/2020 | TSLA  | 55    | High     | High   | 0     |
| 2/1/2020 | AMZ   | 60    | High     | Low    | 0     |
| 2/1/2020 | APPL  | 39    | Low      | Low    | 0     |
| 3/1/2020 | TSLA  | 60    | Low      | High   | 1     |

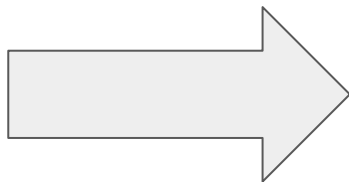
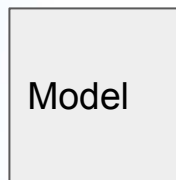
## 4. How to Estimate Model's Performance



# Steps to Predictive Modeling



# How To Evaluate?



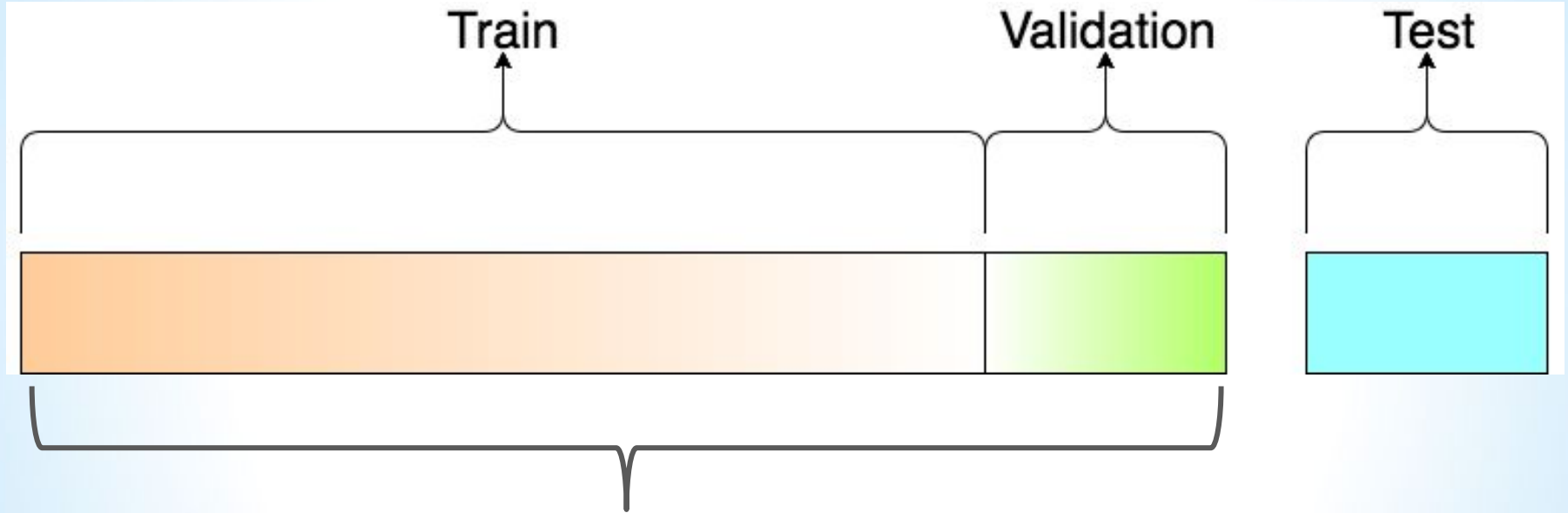
## **Tesla**

- Buy Confidence 80%
- Not Buy Confidence 20%

## **Amazon**

- Buy Confidence 49%
- Not Buy Confidence 51%

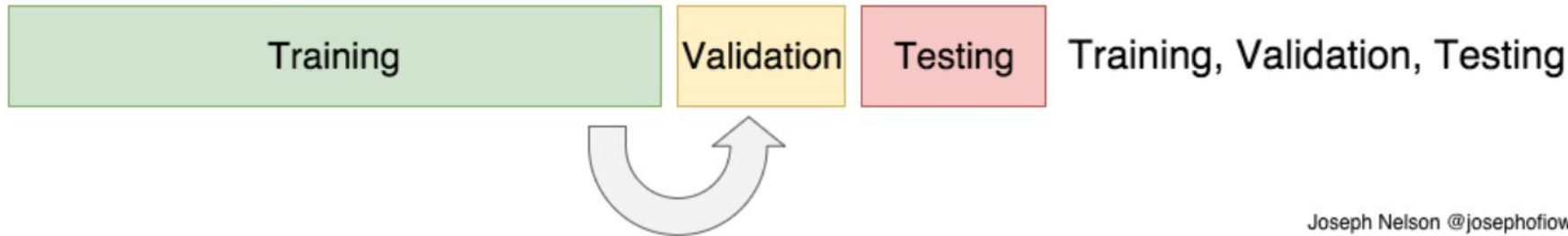
# Estimating Performance



E.g. Choose best model  
Hyper parameter optimization

# Estimating Performance - Data is Abundant

Data Permitting:



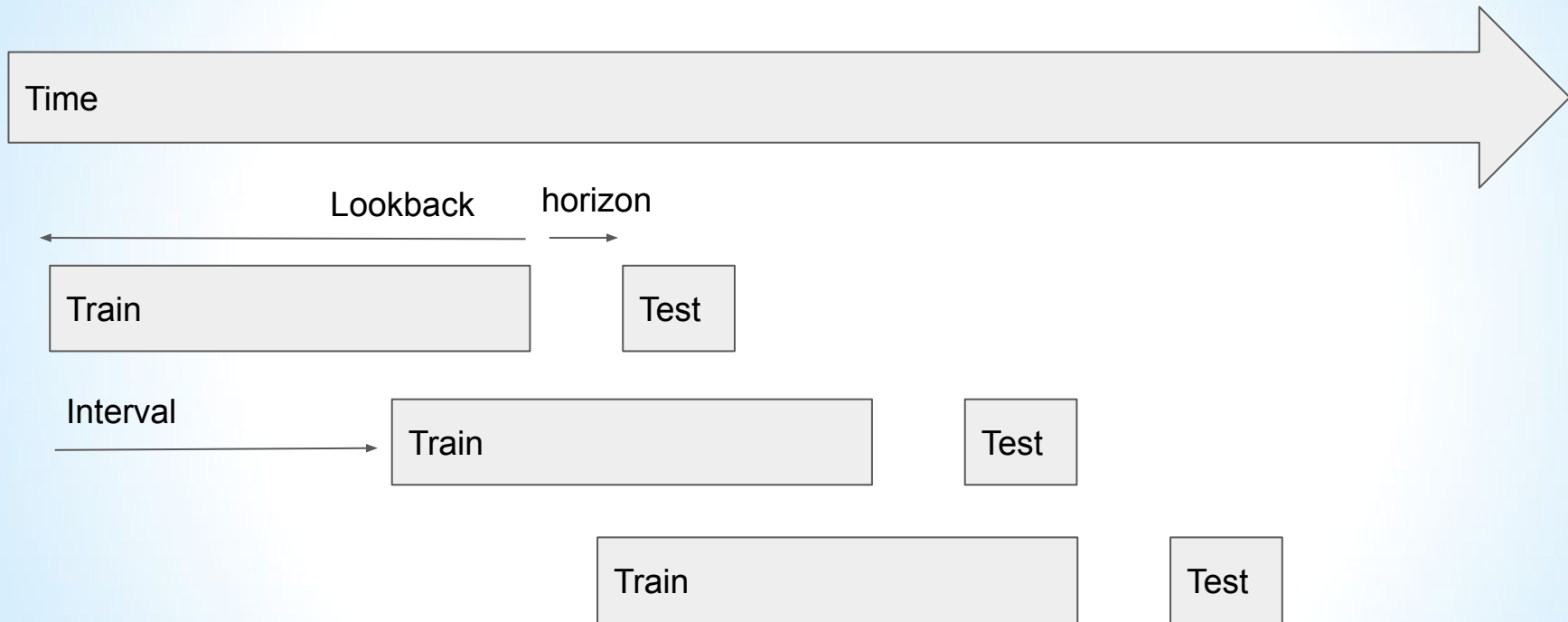
Joseph Nelson @josephofiowa

Datasets distribution: Training  $\leftrightarrow$  Validation  $\Rightarrow$  Test  $\sim$  Real world = Random

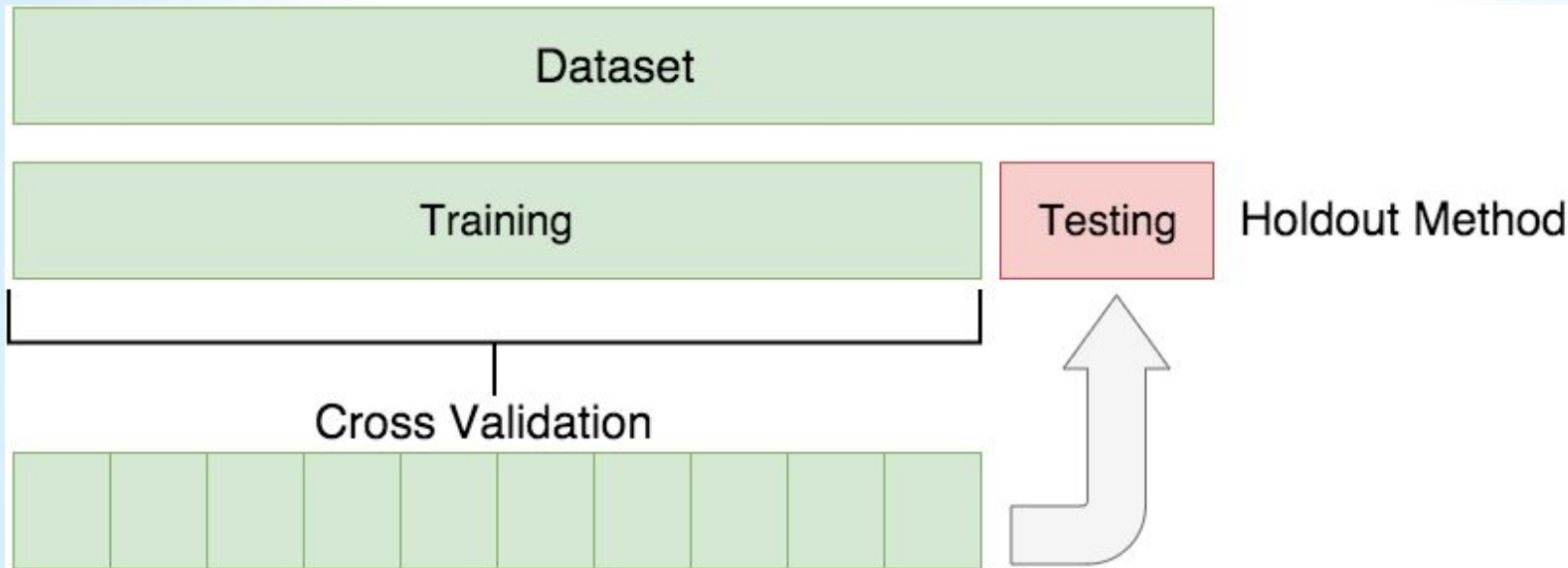
Validation used for Hypertuning and model Calibration

Testing used for final evaluation

# Rolling window cross validation



# Cross Validation



# Classification Metrics

|            |                                     | Label                       |                                   |  |
|------------|-------------------------------------|-----------------------------|-----------------------------------|--|
|            |                                     | Condition Positive<br>(Buy) | Condition Negative<br>(Don't Buy) |  |
| Classifier | Predict Positive<br>(should buy)    |                             |                                   |  |
|            | Predict Negative<br>(shouldn't buy) |                             |                                   |  |
|            |                                     |                             |                                   |  |

# Classification Metrics

|            |                                     | Label                       |                                   |  |
|------------|-------------------------------------|-----------------------------|-----------------------------------|--|
|            |                                     | Condition Positive<br>(Buy) | Condition Negative<br>(Don't Buy) |  |
| Classifier | Predict Positive<br>(should buy)    | True Positive<br>(TP) = 20  |                                   |  |
|            | Predict Negative<br>(shouldn't buy) |                             | True Negative<br>(TN) = 1820      |  |
|            |                                     |                             |                                   |  |



# Classification Metrics

|            |                                     | Label                       |                                   |  |
|------------|-------------------------------------|-----------------------------|-----------------------------------|--|
|            |                                     | Condition Positive<br>(Buy) | Condition Negative<br>(Don't Buy) |  |
| Classifier | Predict Positive<br>(should buy)    | True Positive<br>(TP) = 20  | False Positive<br>(FP) = 180      |  |
|            | Predict Negative<br>(shouldn't buy) | False Negative<br>(FN) = 10 | True Negative<br>(TN) = 1820      |  |
|            |                                     |                             |                                   |  |

# Classification Metrics

|            |                                     | Label                       |                                   |  |
|------------|-------------------------------------|-----------------------------|-----------------------------------|--|
|            |                                     | Condition Positive<br>(Buy) | Condition Negative<br>(Don't Buy) |  |
| Classifier | Predict Positive<br>(should buy)    | True Positive<br>(TP) = 20  | False Positive<br>(FP) = 180      | Positive predictive value   Precision<br>$TP / (TP + FP)$<br>$= 20 / (20 + 180)$<br>$= 10\%$ |
|            | Predict Negative<br>(shouldn't buy) | False Negative<br>(FN) = 10 | True Negative<br>(TN) = 1820      | Negative predictive value<br>$TN / (FN + TN)$<br>$= 1820 / (10 + 1820)$<br>$\approx 99.5\%$  |
|            |                                     |                             |                                   |  |

# Classification Metrics

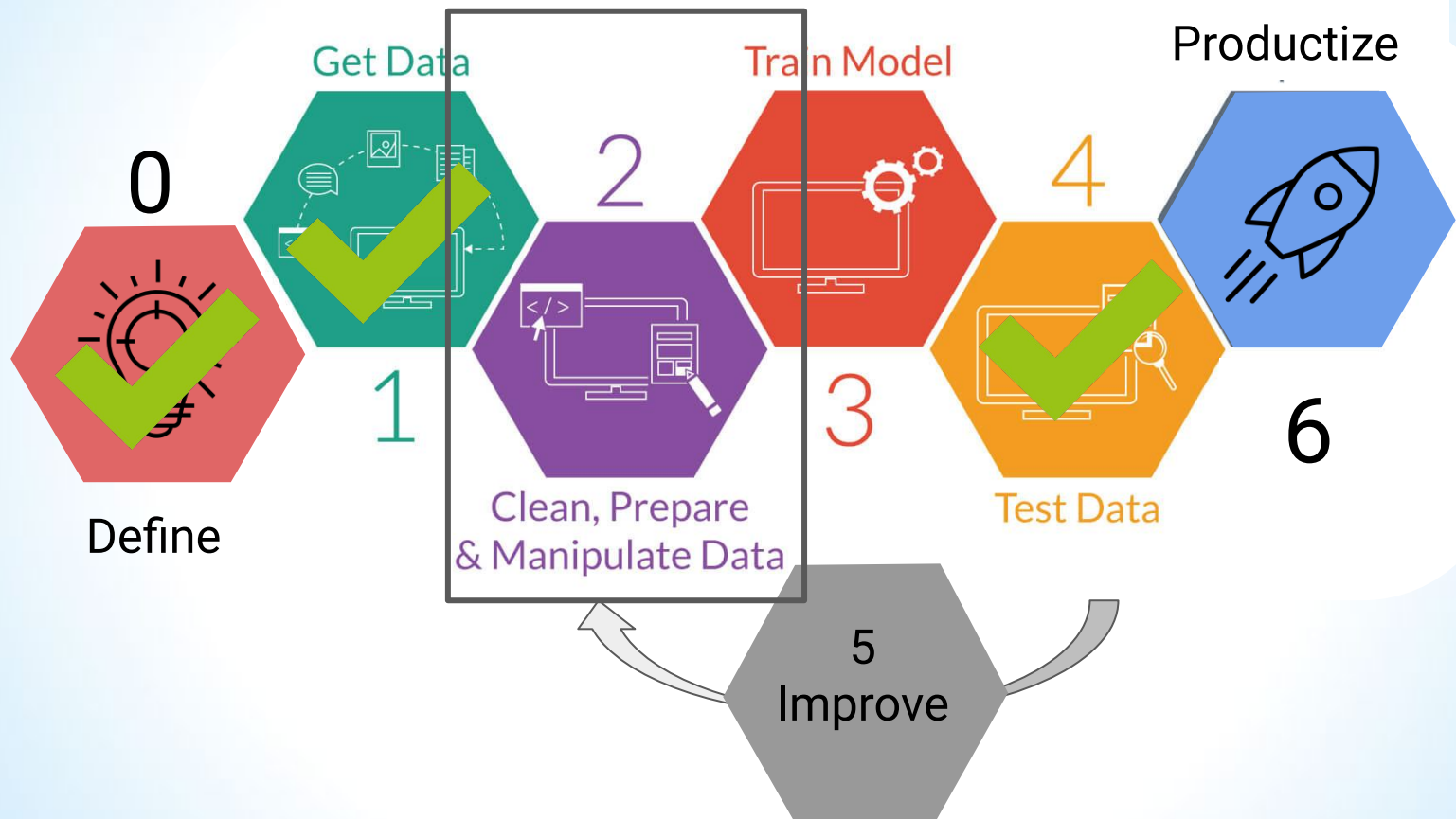
|            |                                     | Label   |                                   |  |
|------------|-------------------------------------|---|-----------------------------------|--|
|            |                                     | Condition Positive<br>(Buy)   | Condition Negative<br>(Don't Buy) |  |
| Classifier | Predict Positive<br>(should buy)    | True Positive<br>(TP) = 20  | False Positive<br>(FP) = 180      | Positive predictive value   Precision<br>$TP / (TP + FP)$<br>$= 20 / (20 + 180)$<br>$= 10\%$ |
|            | Predict Negative<br>(shouldn't buy) | False Negative<br>(FN) = 10   | True Negative<br>(TN) = 1820      | Negative predictive value<br>$TN / (FN + TN)$<br>$= 1820 / (10 + 1820)$<br>$\approx 99.5\%$  |
|            |                                     | True Positive Rate   Recall   Sensitivity<br>$TP / (TP + FN)$<br>$= 20 / (20 + 10)$<br>$\approx 67\%$ |                                   | Specificity<br>$TN / (FP + TN)$<br>$= 1820 / (180 + 1820)$<br>$= 91\%$                       |

# Classification Metrics

|            |                                     | Label   |  |   |
|------------|-------------------------------------|---|--|---|
|            |                                     | Condition Positive<br>(Buy)   | Condition Negative<br>(Don't Buy)                                      |   |
| Classifier | Predict Positive<br>(should buy)    | True Positive<br>(TP) = 20  | False Positive<br>(FP) = 180   | Positive predictive value   Precision<br>$TP / (TP + FP)$<br>$= 20 / (20 + 180)$<br>$= 10\%$                          |
|            | Predict Negative<br>(shouldn't buy) | False Negative<br>(FN) = 10   | True Negative<br>(TN) = 1820   | Negative predictive value<br>$TN / (FN + TN)$<br>$= 1820 / (10 + 1820)$<br>$\approx 99.5\%$                           |
|            |                                     | True Positive Rate   Recall   Sensitivity<br>$TP / (TP + FN)$<br>$= 20 / (20 + 10)$<br>$\approx 67\%$ | Specificity<br>$TN / (FP + TN)$<br>$= 1820 / (180 + 1820)$<br>$= 91\%$ | Accuracy<br>$(TP + TN) / (TP + TN + FP + FN)$<br><b>F1 score</b><br>$(2 * Precision * Recall) / (Precision + Recall)$ |

## 5. How to prepare data for ML

# Steps to Predictive Modeling



# Classical vs Deep Learning Framework



Traditional Machine Learning Flow

# Modeling (data & learning)

## Baseline

Only last price value:

- Last day increment
- Moving average

## Improvements

Based on last year data predict next day performance

- Aggregate last week data
- Extract technical indicator features



# Feature extraction - Data modeling

| Date     | Stock | Price | Label |
|----------|-------|-------|-------|
| 1/1/2020 | TSLA  | 50    | 1     |
| 1/1/2020 | AMZ   | 63    | 0     |
| 1/1/2020 | APPL  | 42    | 0     |
| 2/1/2020 | TSLA  | 55    | 0     |
| 2/1/2020 | AMZ   | 60    | 0     |
| 2/1/2020 | APPL  | 39    | 0     |
| 3/1/2020 | TSLA  | 60    | 1     |

Group by

Stock, week.

| Days price increase | Price Change | RSI | Sector | Label |
|---------------------|--------------|-----|--------|-------|
| 4                   | 0.12         | 0.4 | Auto   | 1     |

# Feature extraction

| Days price increase | Price Change | RSI | Sector    | Label |
|---------------------|--------------|-----|-----------|-------|
| 4                   | 0.12         | 0.4 | Auto      | 1     |
| 2                   | 0.35         | 0.7 | Software  | 0     |
| 4                   | 0.8          | 0.5 | Energy    | 0     |
| 3                   | 0.22         | 0.3 | Materials | 1     |
| 5                   | 0.3          | 0.6 | Health    | 0     |
| 1                   | 0.1          | 0.3 | Telco     | 1     |

Extract features

| Days price increase norm | Price Change <0.3 | RSI | Sector_Auto |
|--------------------------|-------------------|-----|-------------|
| 4 / 7                    | 0                 | 0.4 | 1           |
| 2 / 7                    | 1                 | 0.7 | 0           |
| 4 / 7                    | 1                 | 0.5 | 0           |
| 3 / 7                    | 0                 | 0.3 | 0           |
| 7 / 7                    | 1                 | 0.6 | 0           |
| 1 / 7                    | 0                 | 0.3 | 0           |

# Feature Selection

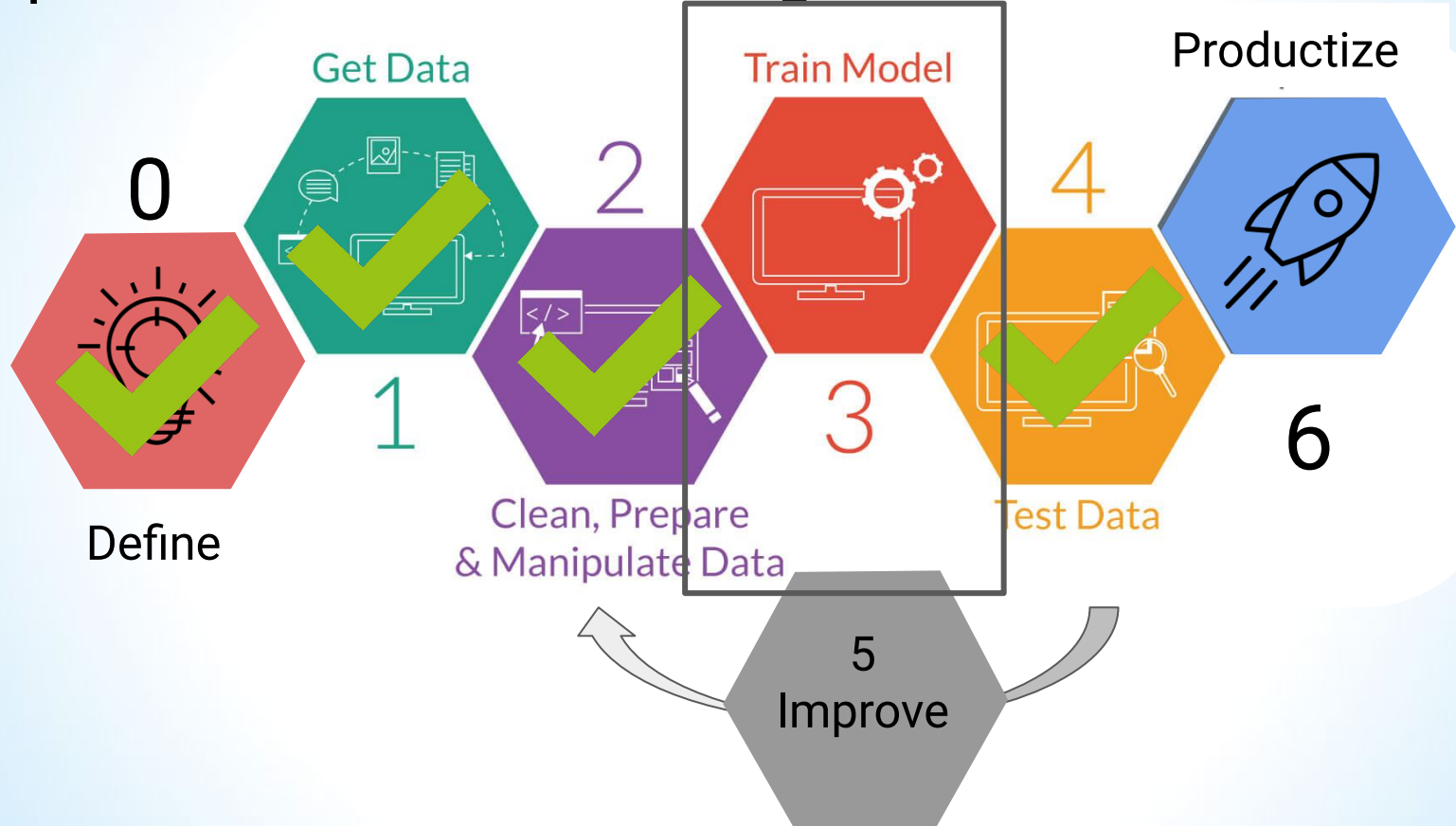
| Days price increase | Price Change | RSI | Sector    | Label |
|---------------------|--------------|-----|-----------|-------|
| 4                   | 0.12         | 0.4 | Auto      | 1     |
| 2                   | 0.35         | 0.7 | Software  | 0     |
| 4                   | 0.8          | 0.5 | Energy    | 0     |
| 3                   | 0.22         | 0.3 | Materials | 1     |
| 5                   | 0.3          | 0.6 | Health    | 0     |
| 1                   | 0.1          | 0.3 | Telco     | 1     |

Extract features

| Days price increase norm | Price Change <0.3 | RSI | Sector_Auto |
|--------------------------|-------------------|-----|-------------|
| 4 / 7                    | 0                 | 0.4 | 1           |
| 2 / 7                    | 1                 | 0.7 | 0           |
| 4 / 7                    | 1                 | 0.5 | 0           |
| 3 / 7                    | 0                 | 0.3 | 0           |
| 7 / 7                    | 1                 | 0.6 | 0           |
| 1 / 7                    | 0                 | 0.3 | 0           |

6. What types of ML Algorithms are there?

# Steps to Predictive Modeling



# Parametric vs. Non-parametric Models

Almost all models for machine learning have “parameters” or “weights” that need to be learned.

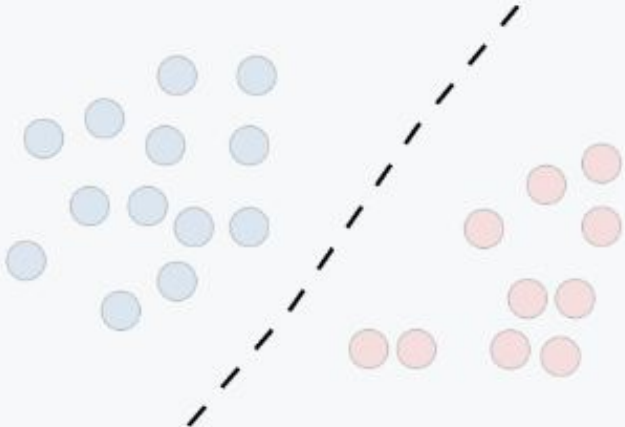
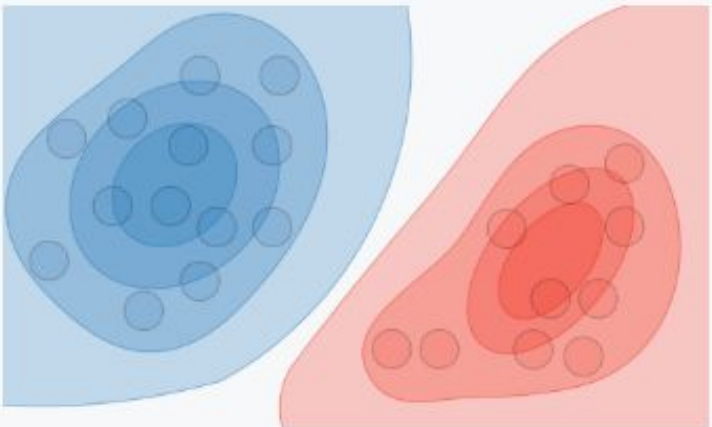
| Parametric Models  | Nonparametric models   |
|--|--|
| The number of parameters is constant, or independent of the number of training examples. | The number of parameters grows with the number of training examples. |

# Can you think of an example?

Can you think of an example for parametric and non-parametric method?

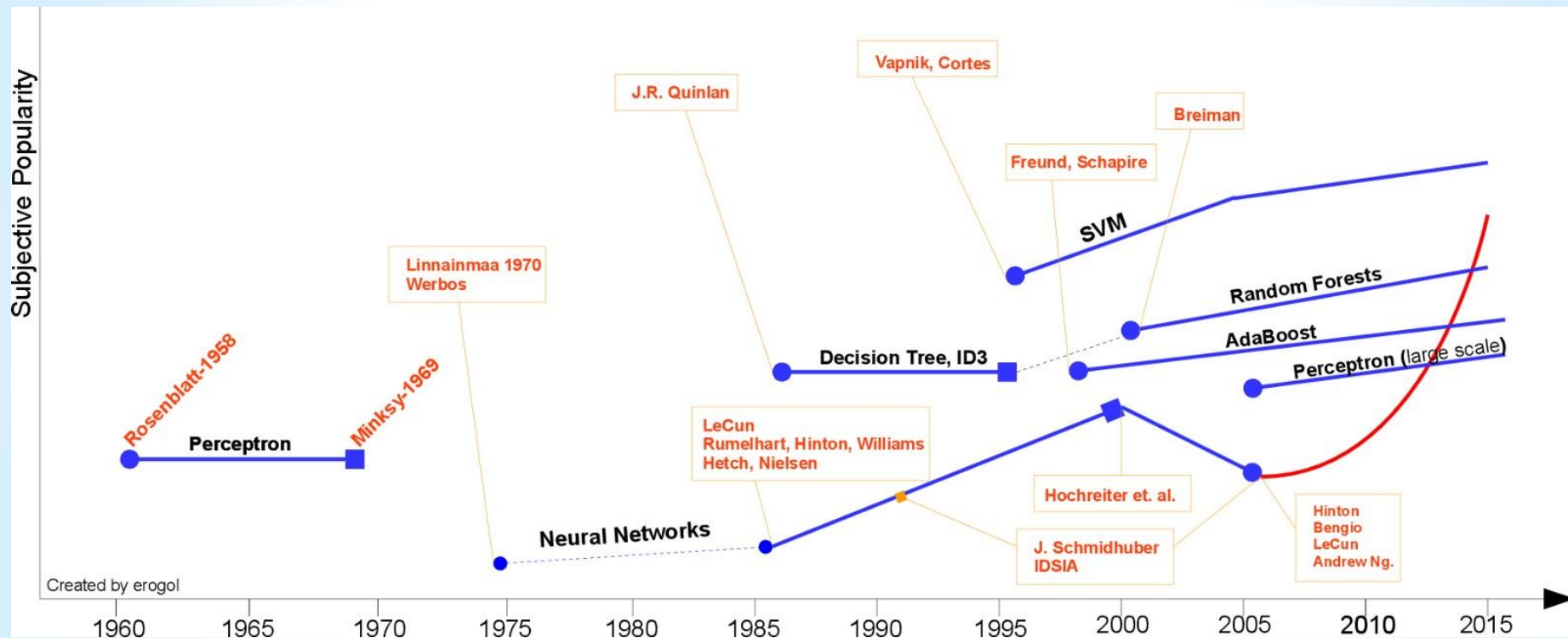
**EXAMPLE**

# Generative vs. Discriminative

|                | Discriminative model   | Generative model  |
|----------------|--|---|
| Goal           | Directly estimate $P(y x)$   | Estimate $P(x y)$ to then deduce $P(y x)$   |
| What's learned | Decision boundary  | Probability distributions of the data   |
| Illustration   |  |  |
| Examples       | Regressions, SVMs  | GDA, Naive Bayes  |



# The Brief History of Machine Learning



# Classical vs Deep Learning Framework



Traditional Machine Learning Flow

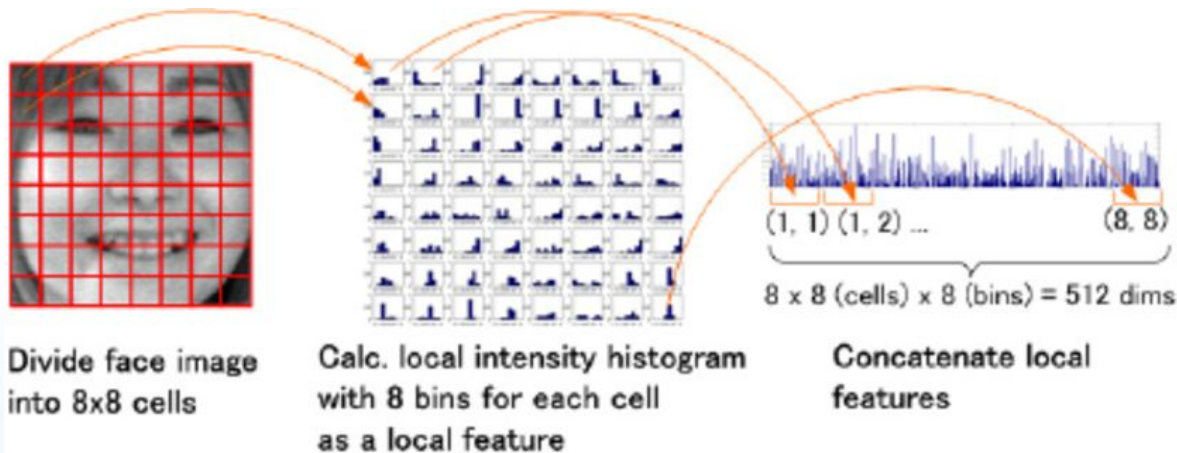


Deep Learning Flow

# Feature Extraction

“Algorithm which transforms raw data into numeric values which can be used as input to a learning algorithm. Usually helps with **reducing** and **fixing** dimensionality.”

e.g.



# Some Realities on DL

Don't be fool by the hype

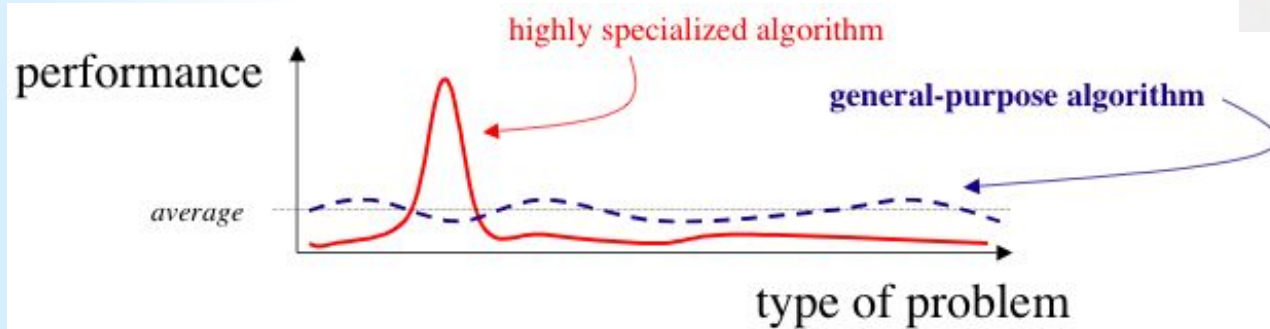
1. Can be beaten by GBT for tabular data (CatBoost, XGBoost, LightGBM).
2. No Feature Engineering - Yuppie.  
Yet... Network Architecture Search (NAS), Annoying GPU issues, Loss design, hours of training
3. Overkill sometimes and infeasible  
Example: try to train a DL to predict if a number is even or odd...



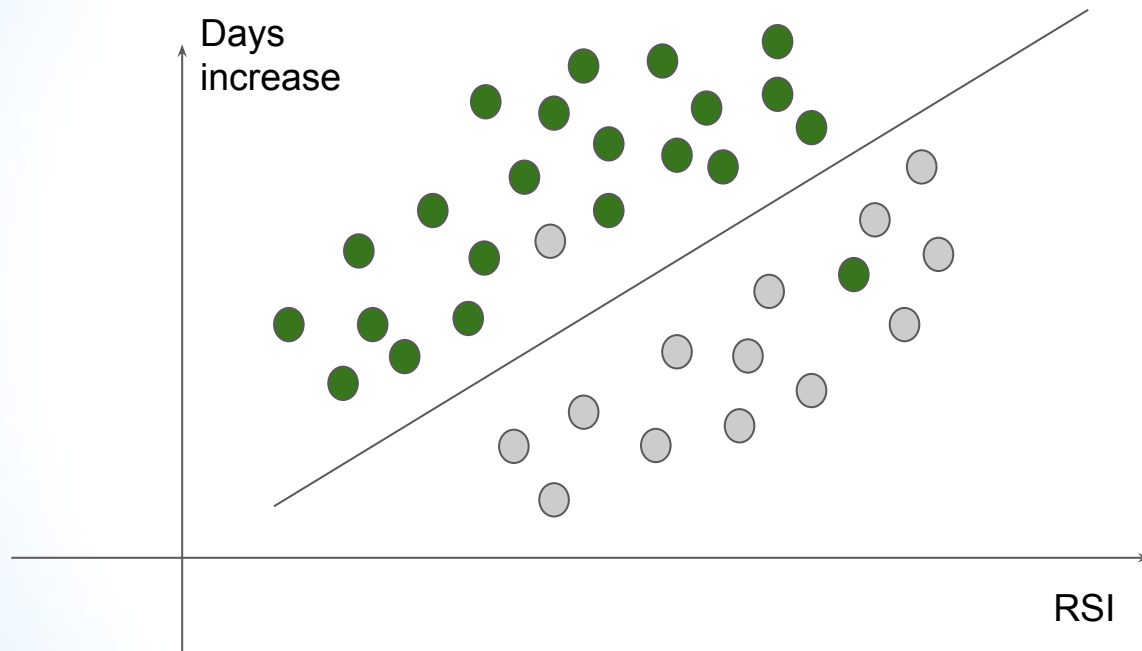
# No Free Lunch Theorem - Best Model Does Not Exist

A superior black-box optimisation strategy, which is better than anything else for any kind of problem, is impossible.

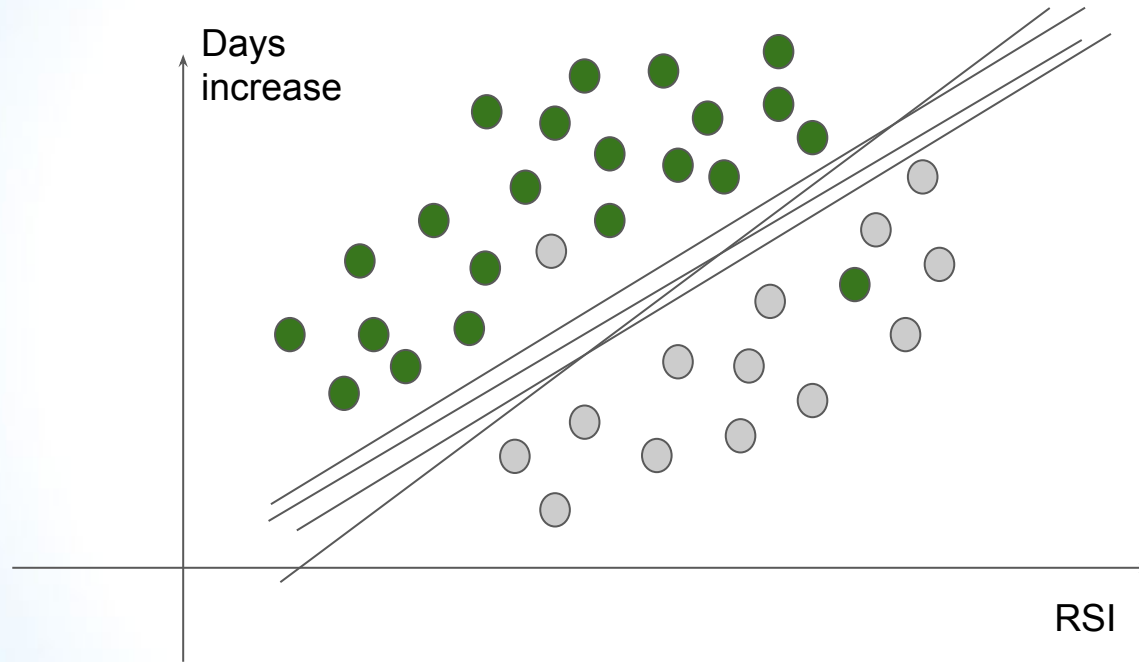
**Deep cannot be always better.**



# Modeling - Linear

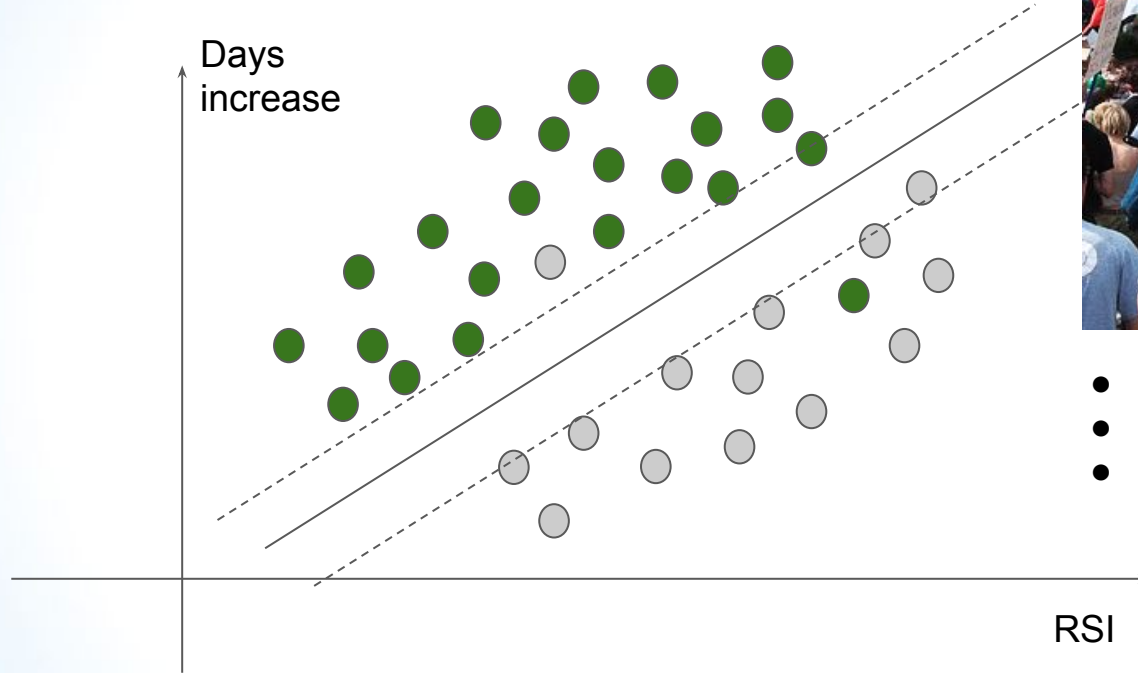


# Modeling - Linear





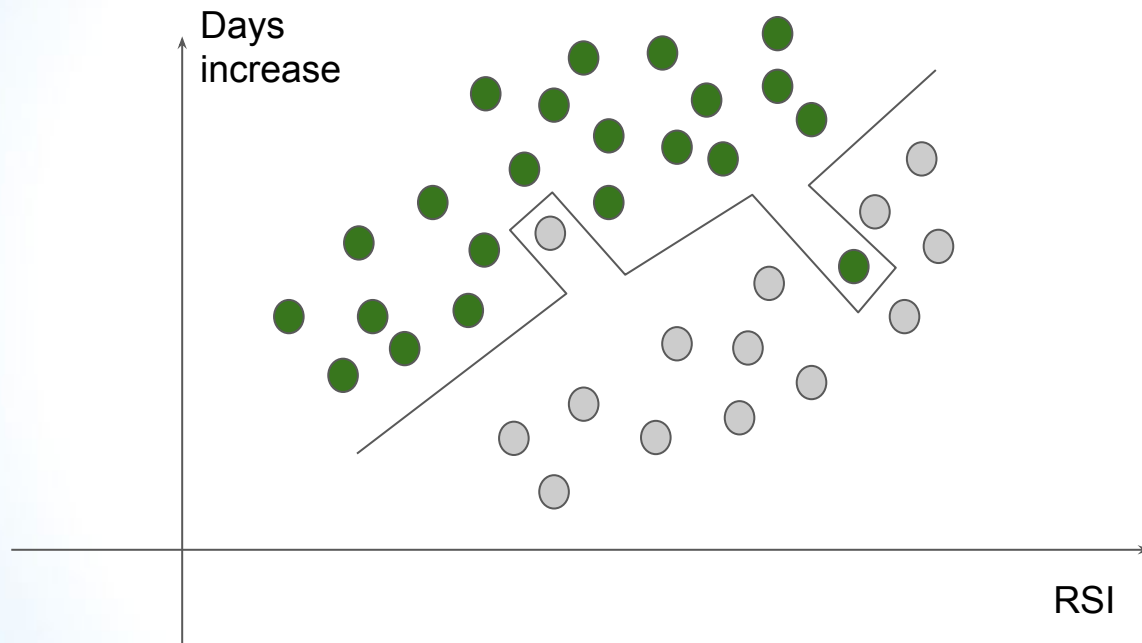
# Modeling - Maximum Margin



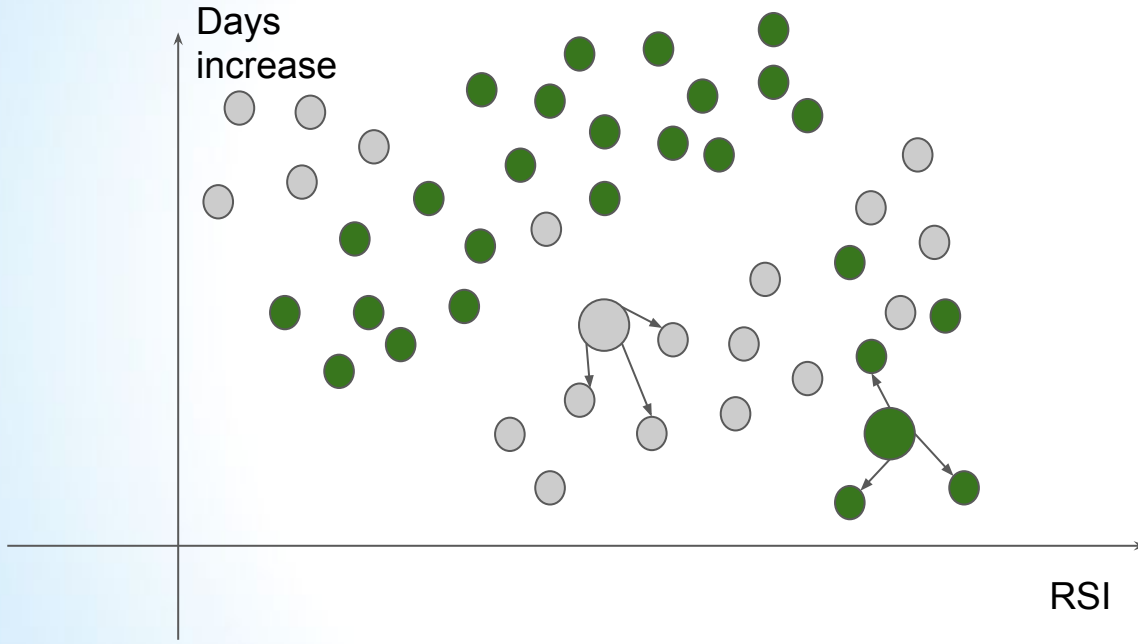
- A Linear classifier
- Maximize the margin
- Simple Linear SVM - LSVM



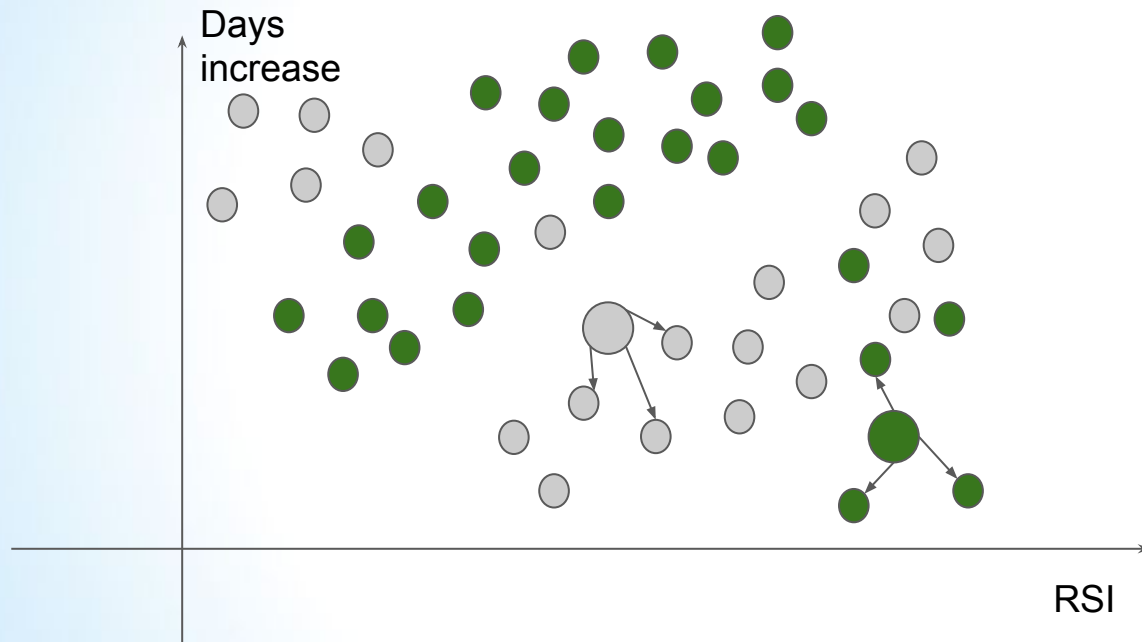
# Modeling



# Modeling - K nearest neighbors

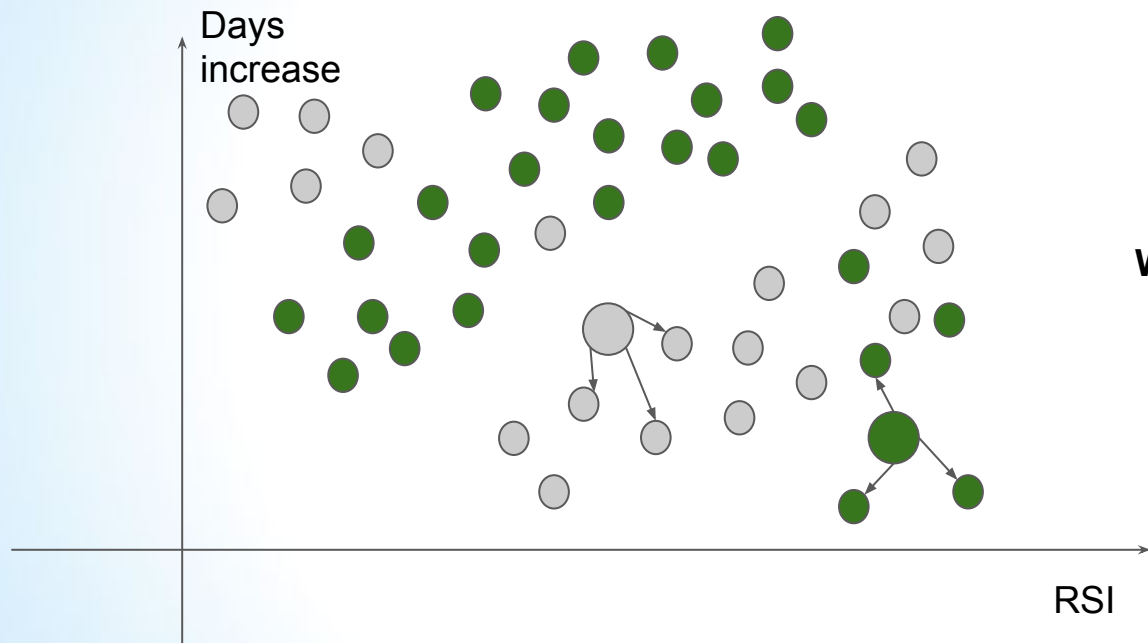


# Modeling - K nearest neighbors



1. To classify a new input vector  $x$
2. Examine the  $k$  closest training data points to  $x$
3. Assign the object to the most frequently occurring class

# Modeling - K nearest neighbors



1. To classify a new input vector  $x$
2. Examine the  $k$  closest training data points to  $x$
3. Assign the object to the most frequently occurring class

## What about?

- $K$  is Odd vs Even  $K$ ?
- How can we apply Voting?

# KNN Best Practices

## When to Consider

- Less than 20 attributes per instance
- Lots of training data

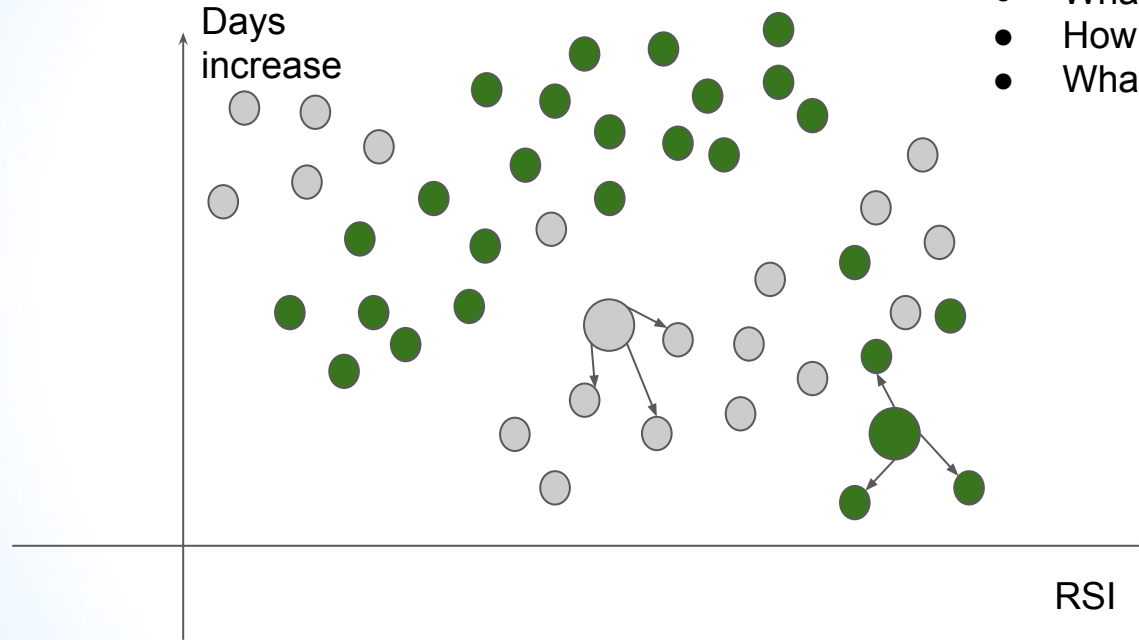
## Advantages

- Training is very fast
- Learn complex target functions
- Do not lose information

## Disadvantages

- Slow at query time
- Easily fooled by irrelevant attributes

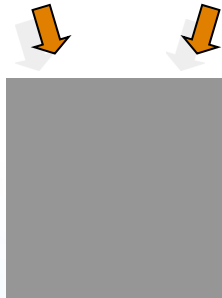
# Modeling - K nearest neighbors



- What is the model?
- How to measure distance?
- What is the training error of 1nn?

# Defining Distance Measures

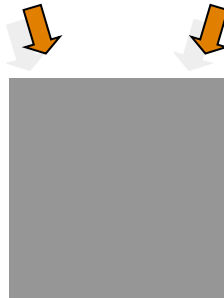
**Definition:** Let  $O_1$  and  $O_2$  be two objects from the universe of possible objects. The distance (dissimilarity) between  $O_1$  and  $O_2$  is a real number denoted by  $D(O_1, O_2)$



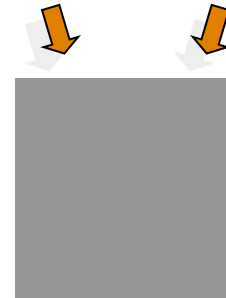
0.23

Peter

Piotr



3



342.7

# Distance function behavior

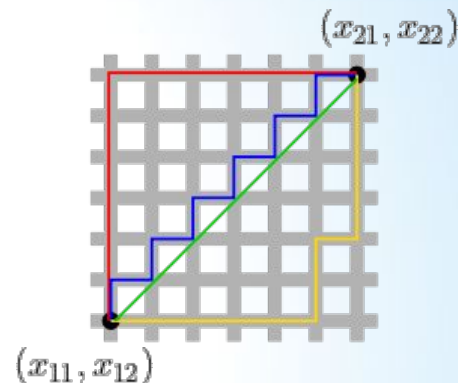
- $\text{dis}(x,y) \geq 0$
- $\text{dis}(x,y) = 0$  iff  $x=y$
- $\text{dis}(x,y) = \text{dis}(y,x)$
- $\text{dis}(x, z) \leq \text{dis}(x, y) + \text{dis}(y, z)$



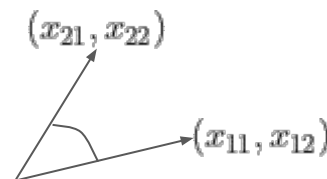
# Distance Function

$$L1(X_1, X_2) = \text{ManhattanDistance}\left(\begin{bmatrix} x_{11} \\ x_{1i} \\ x_{1n} \end{bmatrix}, \begin{bmatrix} x_{21} \\ x_{2j} \\ x_{2n} \end{bmatrix}\right) = \sum_{i=1}^n |x_{1i} - x_{2i}|$$

$$L2(X_1, X_2) = \text{EuclideanDistance}\left(\begin{bmatrix} x_{11} \\ x_{1i} \\ x_{1n} \end{bmatrix}, \begin{bmatrix} x_{21} \\ x_{2i} \\ x_{2n} \end{bmatrix}\right) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}$$



$$\text{CosineSimilarity}(X_1, X_2) = \frac{\sum_{i=1}^n (x_{1i} * x_{2i})}{\sqrt{\sum_{i=1}^n x_{1i}^2} * \sqrt{\sum_{i=1}^n x_{2i}^2}}$$



# Using euclidean distance

| Price Change <0.3 | RSI | Sector_Auto | Label |
|-------------------|-----|-------------|-------|
| 0                 | 0.4 | 1           | 1     |
| 1                 | 0.7 | 0           | 0     |
| 1                 | 0.5 | 0           | 0     |
| 0                 | 0.3 | 0           | 1     |
| 1                 | 0.6 | 0           | 0     |
| 0                 | 0.3 | 0           | 1     |

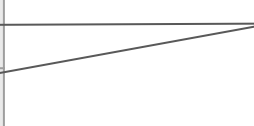
$$\sqrt{(0 - 1)^2 + (0.4 - 0.7)^2 + (1 - 0)^2} = \sqrt{1 + 0.09 + 1} = 1.44$$


$$\sqrt{(0 - 0)^2 + (0.4 - 0.3)^2 + (1 - 0)^2} = \sqrt{0 + 0.01 + 1} = 1.004$$

What do you think about the distance values?

# Using Manhattan distance

| Price Change <0.3 | RSI | Sector_Auto | Label |
|-------------------|-----|-------------|-------|
| 0                 | 0.4 | 1           | 1     |
| 1                 | 0.7 | 0           | 0     |
| 1                 | 0.5 | 0           | 0     |
| 0                 | 0.3 | 0           | 1     |
| 1                 | 0.6 | 0           | 0     |
| 0                 | 0.3 | 0           | 1     |


$$|0 - 1| + |0.4 - 0.7| + |1 - 0|$$
$$= 1 + 0.3 + 1 = 2.3$$


$$|0 - 0| + |0.4 - 0.3| + |1 - 0|$$
$$= 0 + 0.1 + 1 = 1.01$$

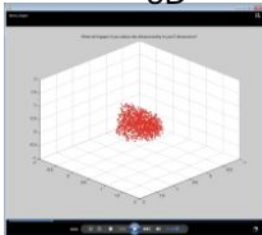
Any ideas about issues with using absolute ?

# Curse of Dimensionality



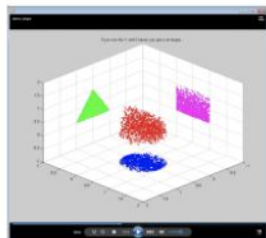


A cloud of points in  
3D

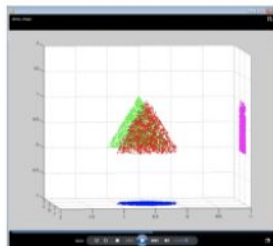


Can be projected into  
2D

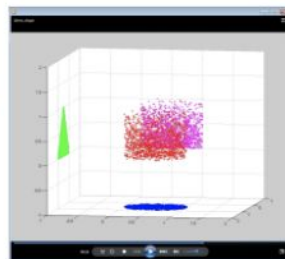
XY or XZ or YZ



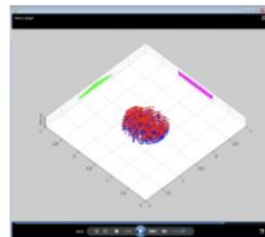
In 2D XZ we see  
a triangle



In 2D YZ we see  
a square

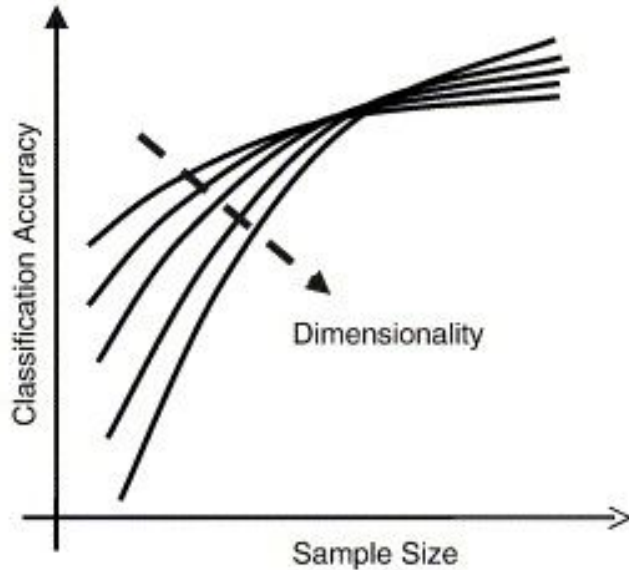


In 2D XY we see  
a circle

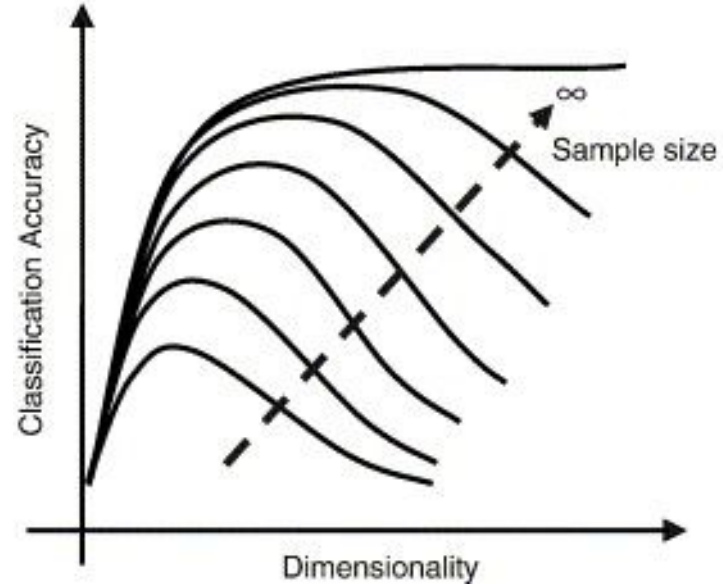


# Hughes phenomenon (1968) (Peaking Paradox)

a. Curse of Dimensionality

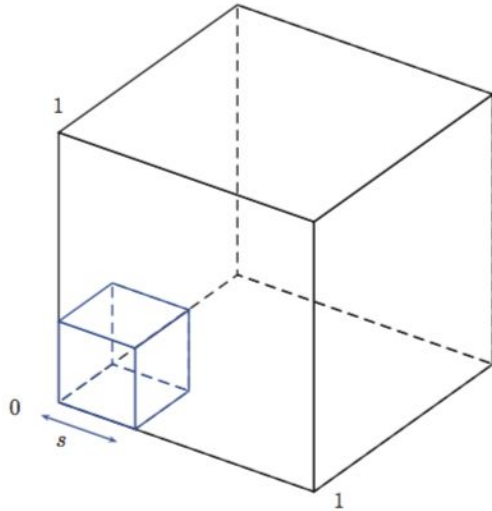


b. Hughes phenomenon

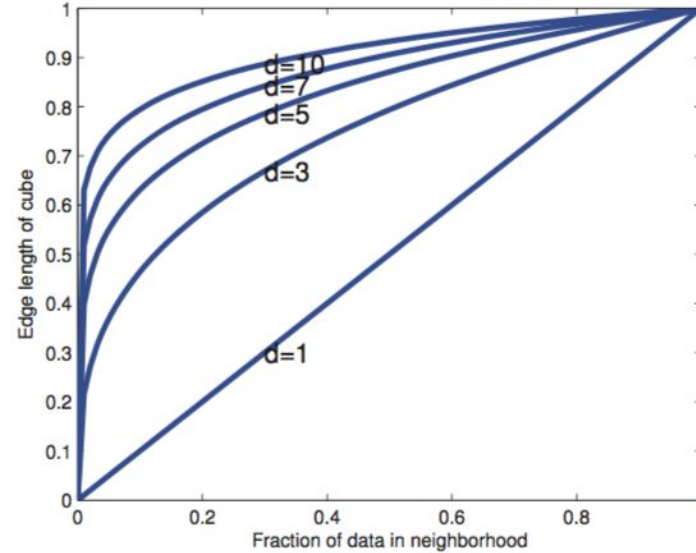




# Why Nearest Neighbours Fails in High Dimensions?



(a)

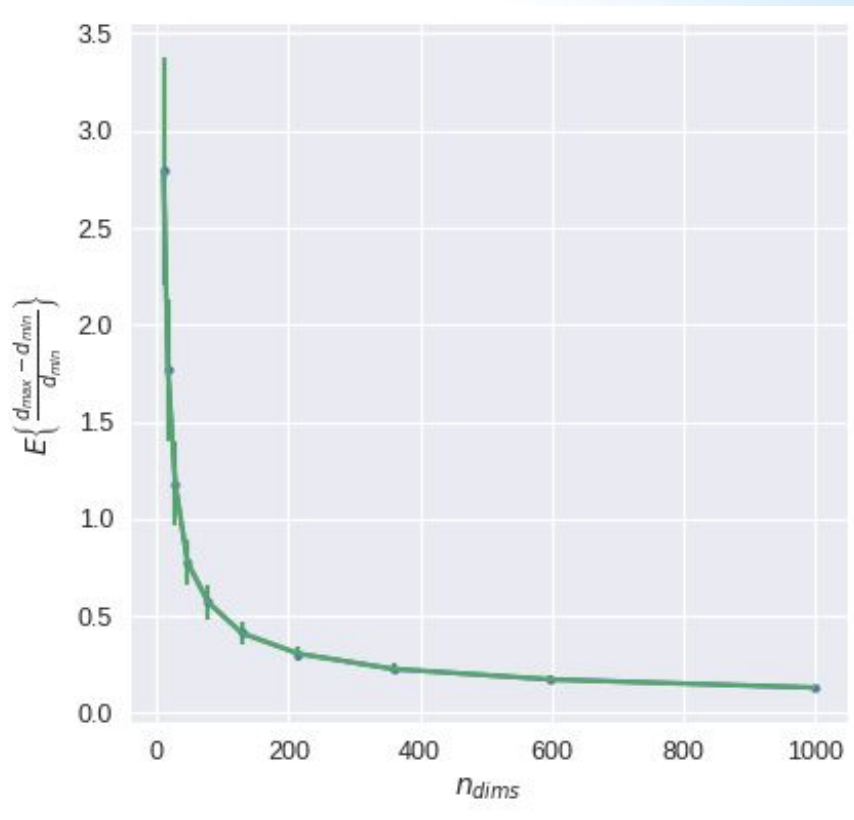


(b)

# Beyer et. al. Theorem

The difference between the maximum and minimum distances to a given query point does not increase as fast as the nearest distance to any point in high dimensional space.

This makes a proximity query meaningless and unstable because there is poor discrimination between the nearest and furthest neighbor.



# Example: Detecting Suspicious URL Names

- DDOS
- Botnets
- Derive by download
- Phishing - How phishing is different from other attacks?

## Task:

- Build an analytics tools to explore and detect NEW types of malicious URLs.
- What we wish to gain? Precision / Recall
- How you can represent domain Names?

[Appendix - helix](#)

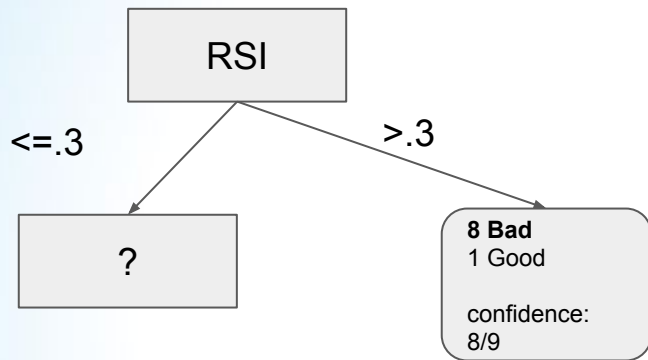
# Review Homework

Part 1 - Implement k-Nearest Neighbours (KNN)

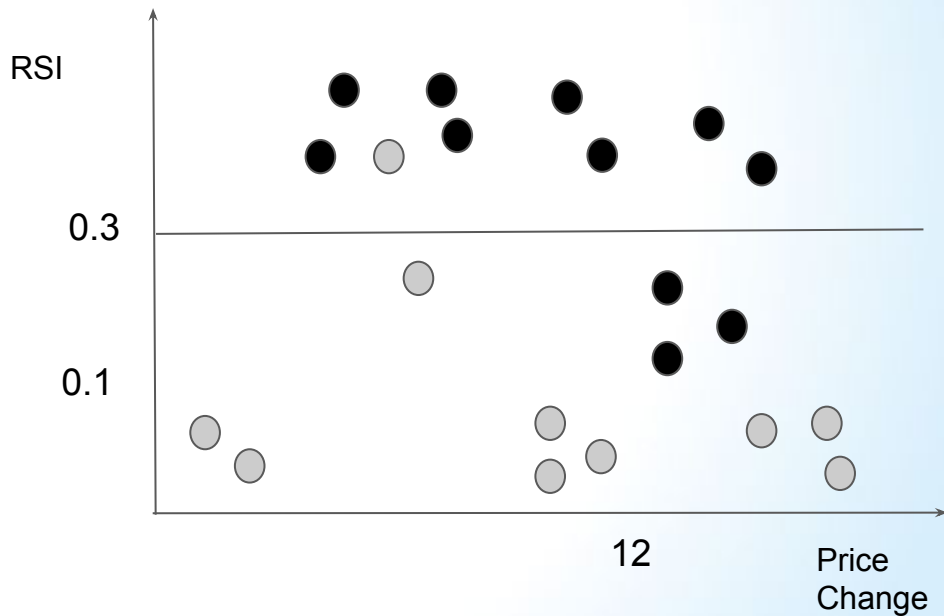
Part 2.1 - Learn and evaluate kNN algorithm on artificial data + Analyse the properties of KNN

Part 3 - bonus

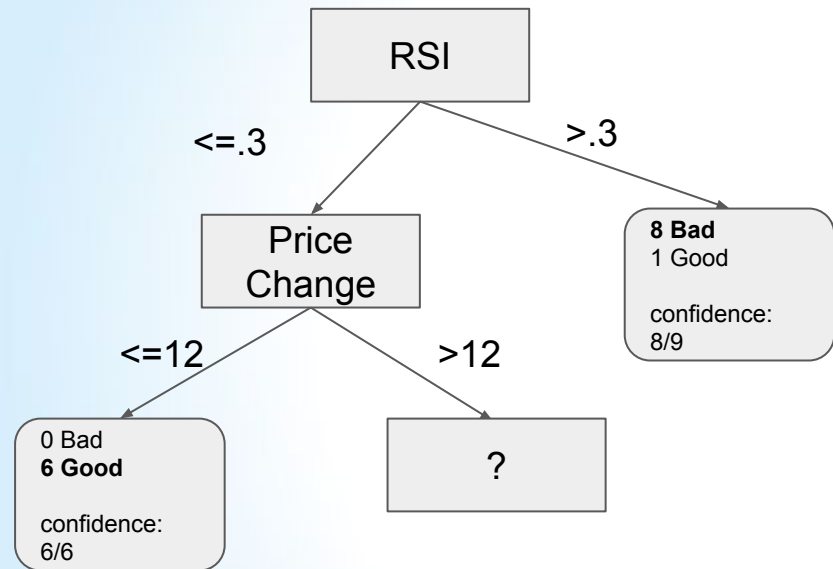
# Decision Tree - Adding splits



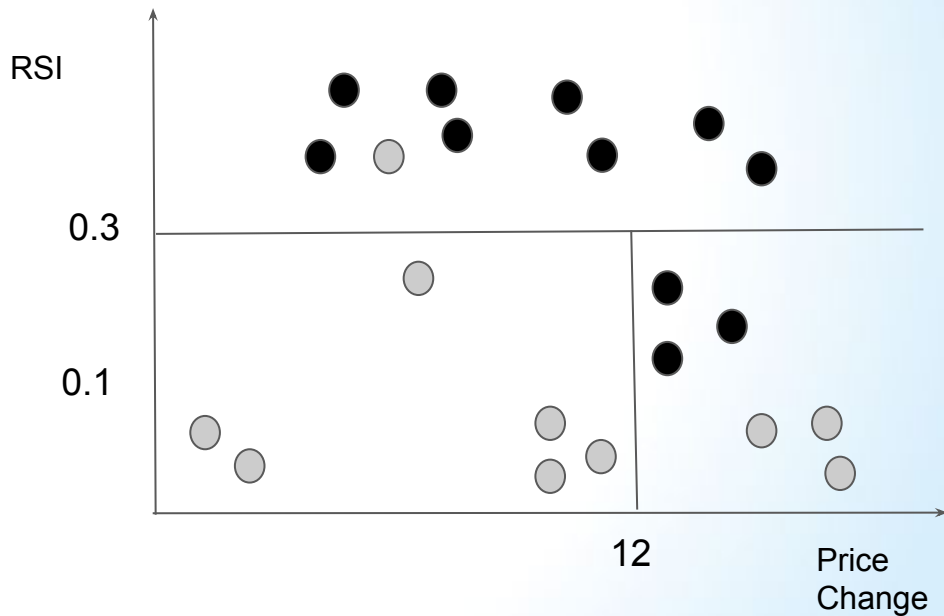
- Nodes = attribute decision
- Leaf nodes = classification decision



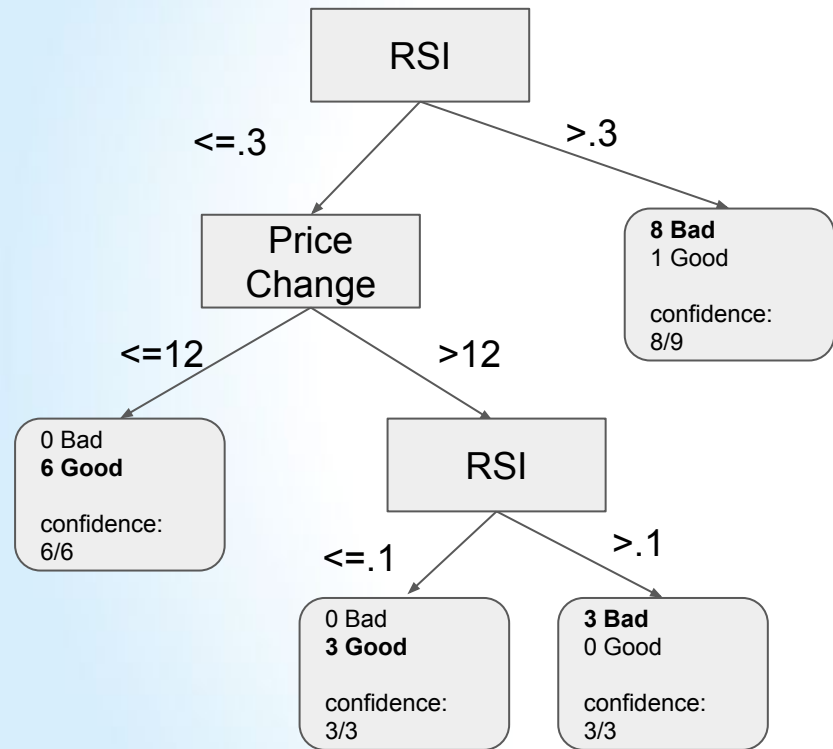
# Decision Tree - Adding splits



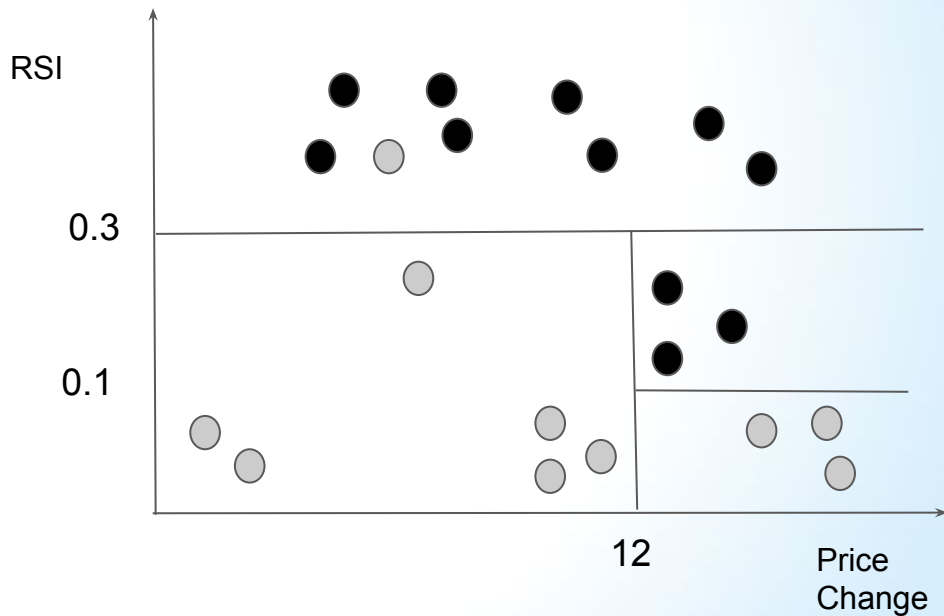
- Nodes = attribute decision
- Leaf nodes = classification decision



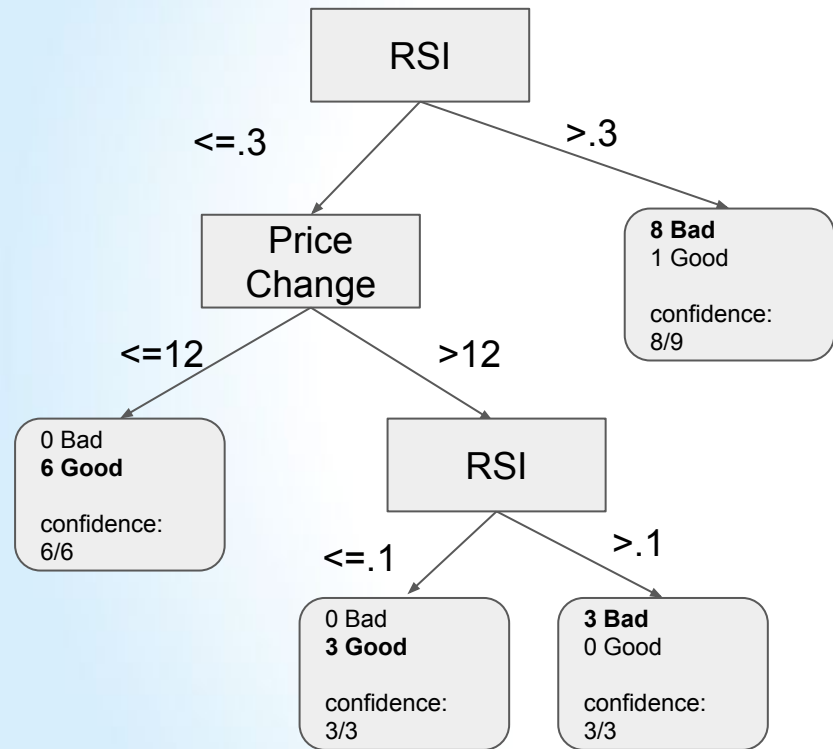
# Decision Tree - Adding splits



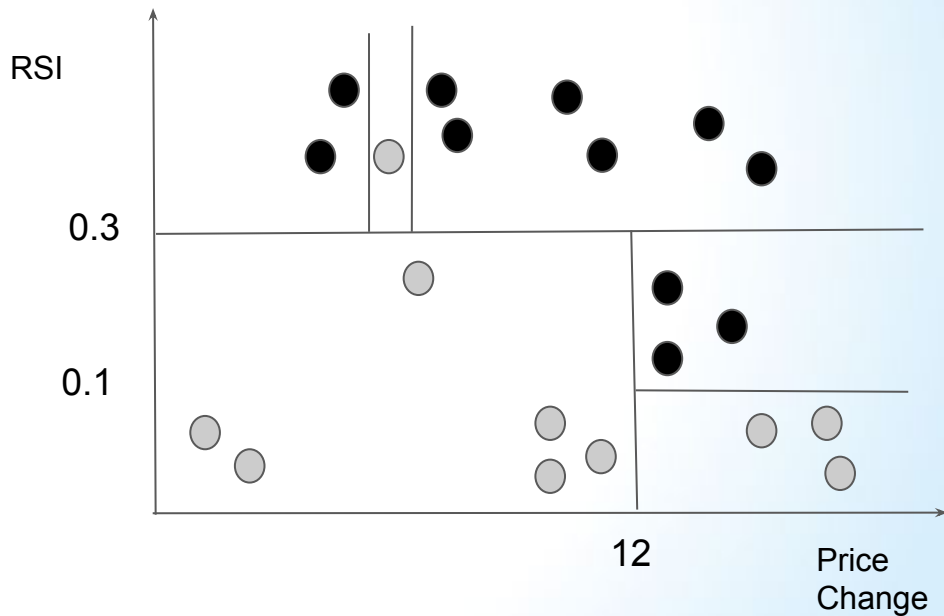
- Nodes = attribute decision
- Leaf nodes = classification decision



# Decision Tree - Adding splits



- Nodes = attribute decision
- Leaf nodes = classification decision



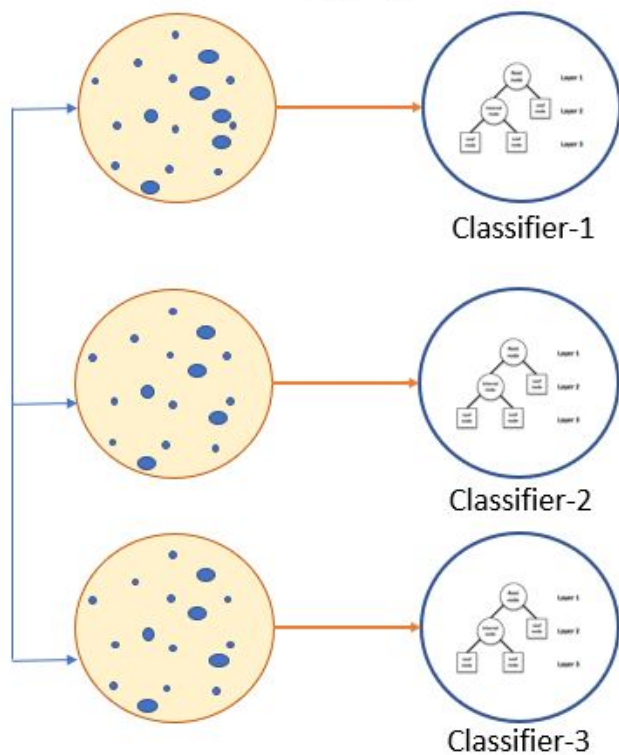


# Ensemble Learning

Use multiple learning algorithms to obtain better predictive performance than could be obtained from any of the constituent learning algorithms alone.

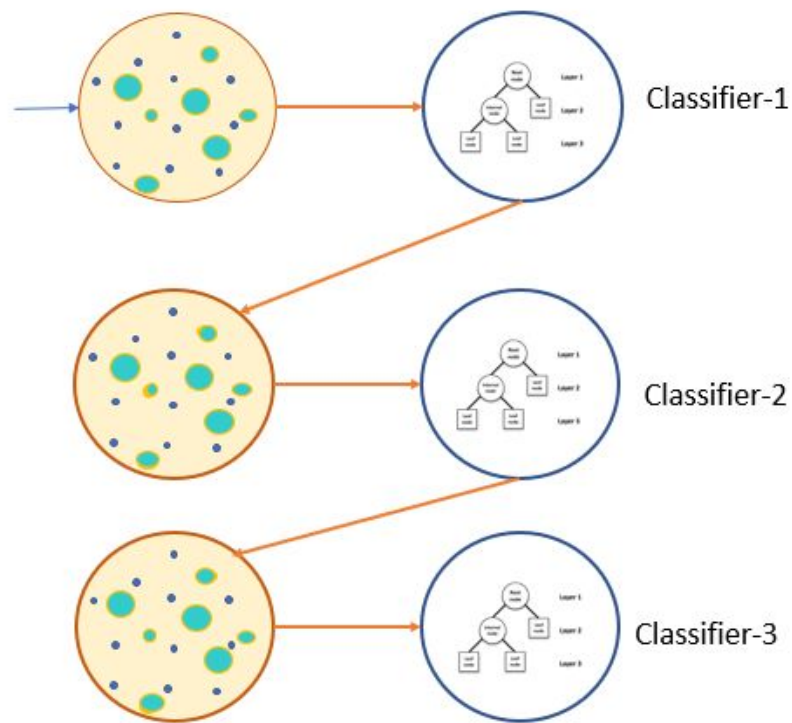
Empirically, ensembles tend to yield better results when there is a significant diversity among the models.

## Bagging



**Parallel**

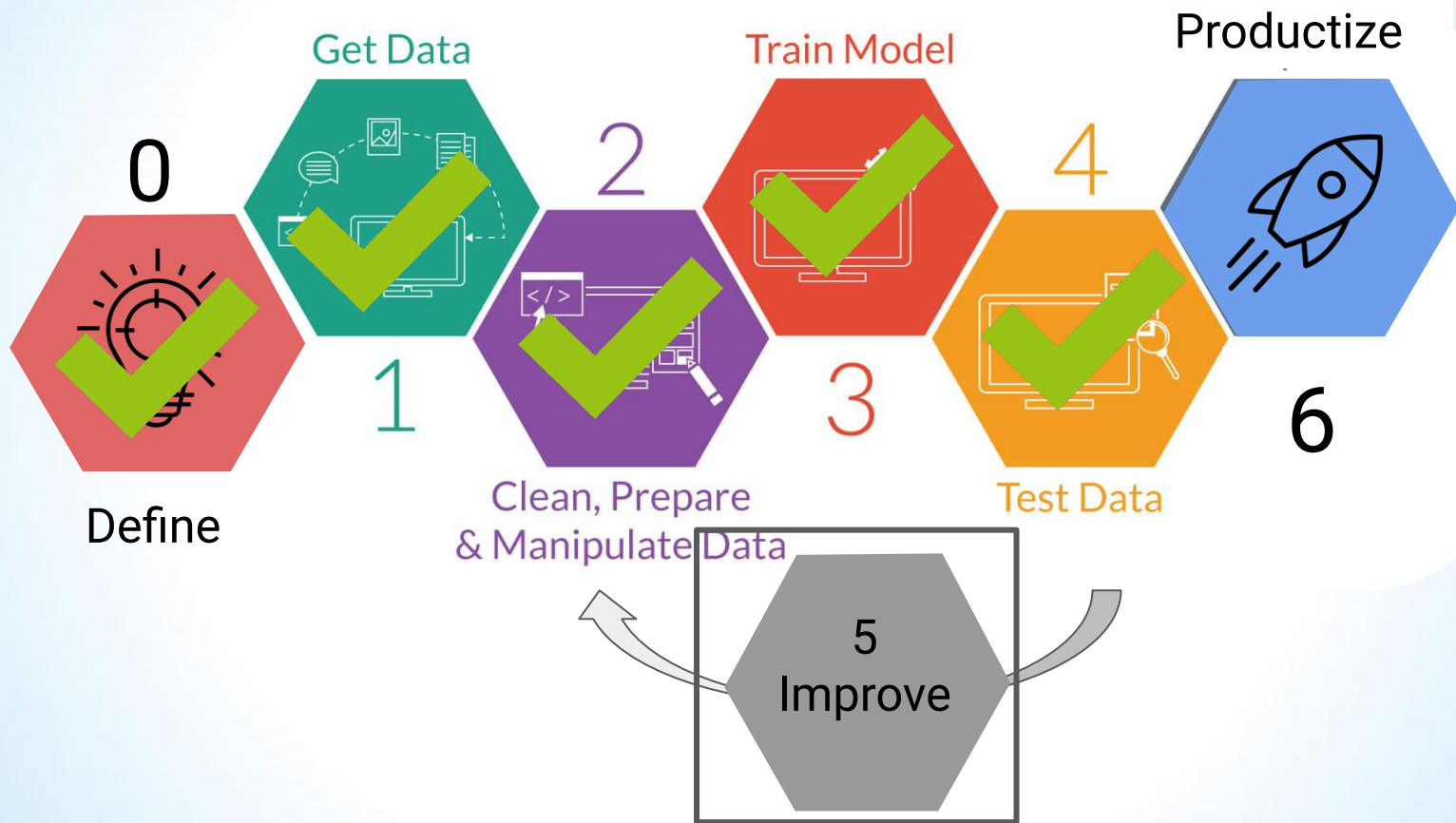
## Boosting



**Sequential**

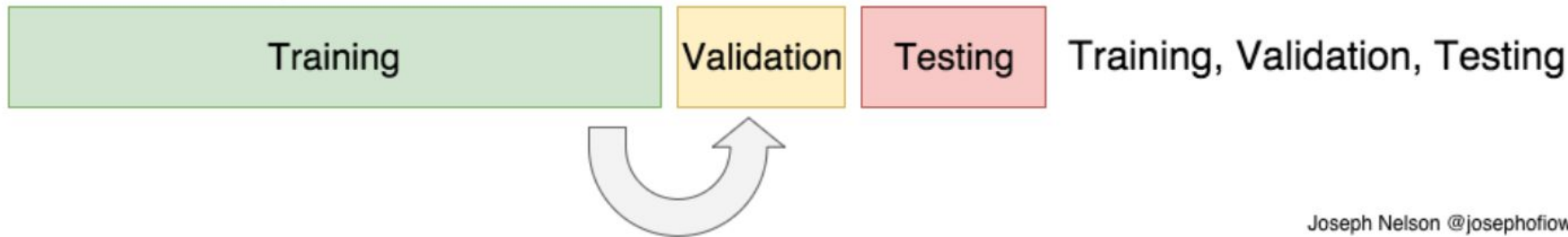
## 7. How to improve ML models?

# Steps to Predictive Modeling



# Estimating Performance - Data is Abundant

Data Permitting:



Joseph Nelson @josephofiowa

Datasets distribution: Training  $\leftrightarrow$  Validation  $\Rightarrow$  Test  $\sim$  Real world = Random

Hypertuning, Calibration: Validation

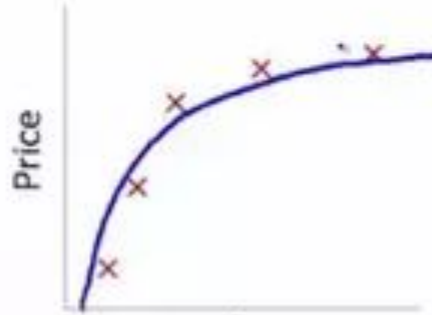
Evaluation: Testing

# Bias Variance Tradeoff - Regression



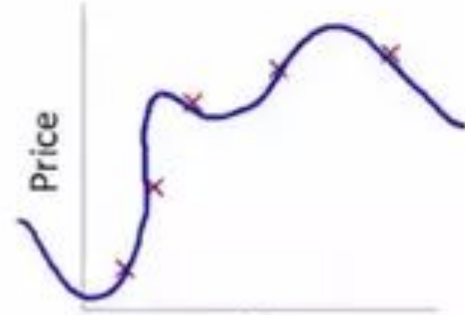
$$\theta_0 + \theta_1 x$$

High bias  
(underfit)



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

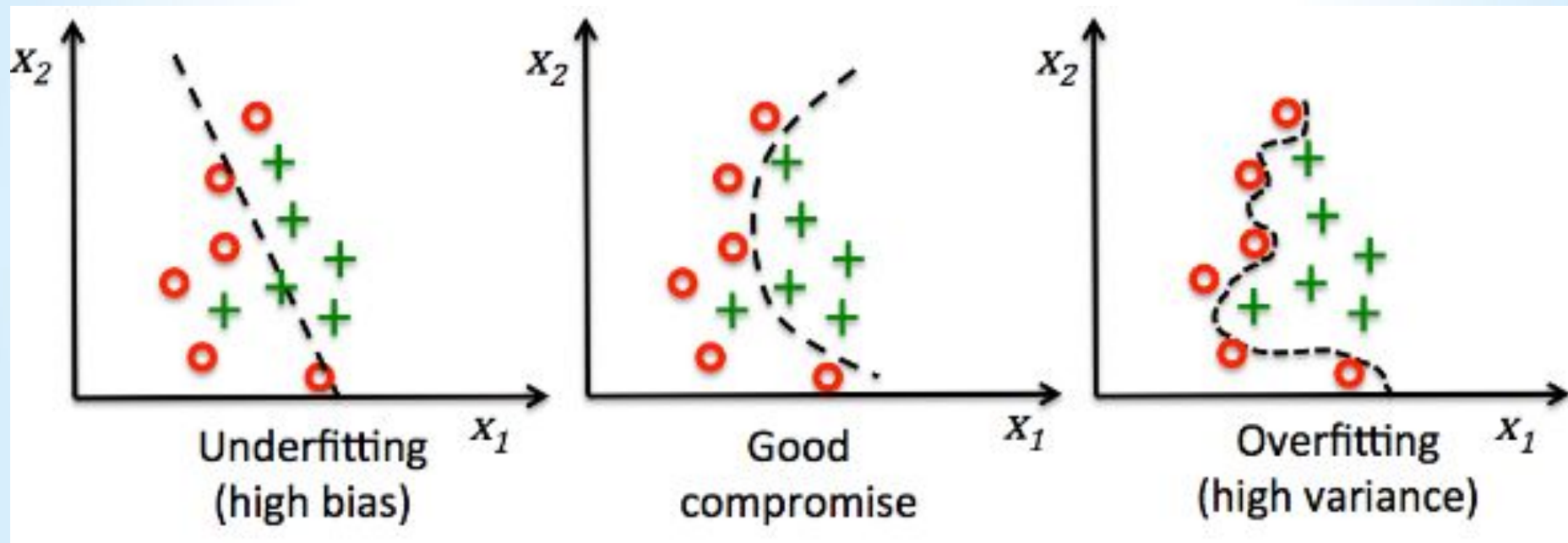
"Just right"



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

High variance  
(overfit)

# Bias Variance Tradeoff - Classification



# Total Error

Assume a simple model  $y = f(x) + \epsilon$ ,  $E(\epsilon) = 0$ ,  $\text{Var}(\epsilon) = \sigma_\epsilon^2$ ,

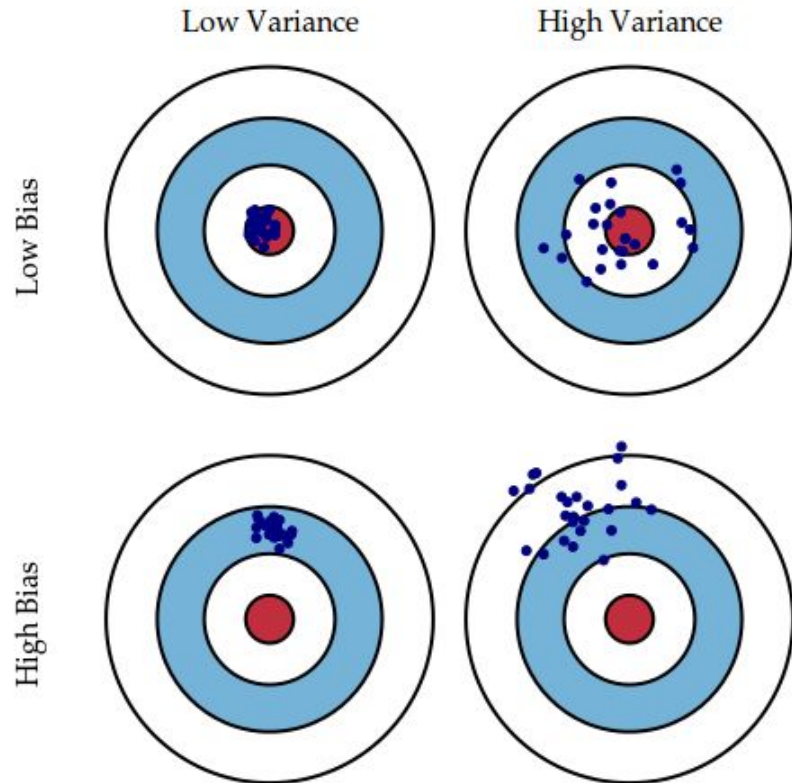
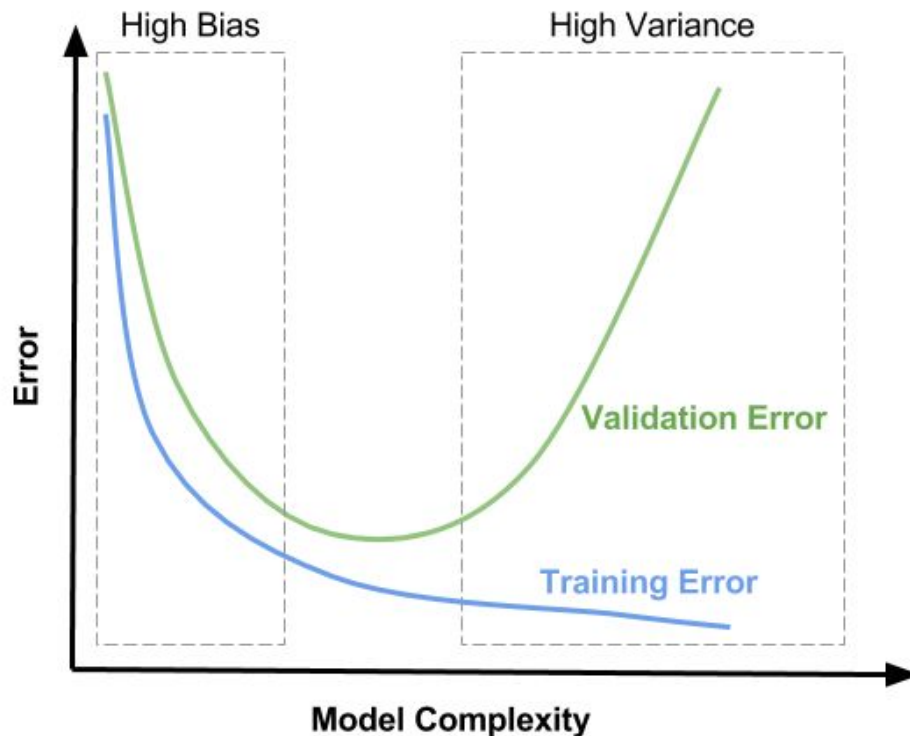
$$\begin{aligned}\text{Err}(x_0) &= E[(y - h(x_0))^2 | X = x_0] \\ &= \sigma_\epsilon^2 + [Eh(x_0) - f(x_0)]^2 + E[h(x_0) - Eh(x_0)]^2 \\ &= \sigma_\epsilon^2 + \text{Bias}^2(h(x_0)) + \text{Var}(h(x_0)) \\ &= \text{Irreducible Error} + \text{Bias}^2 + \text{Variance}\end{aligned}$$

Optional pencil and paper exercise: prove it in details

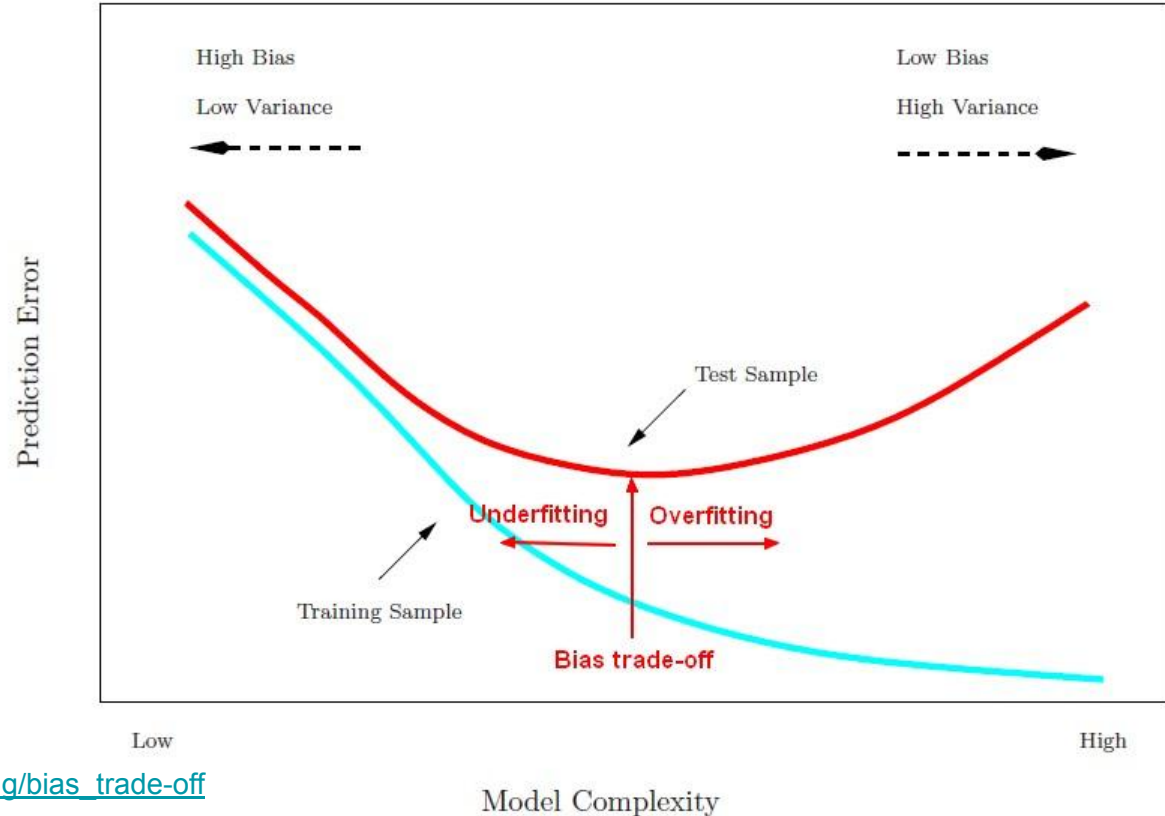


# Bias Variance Tradeoff

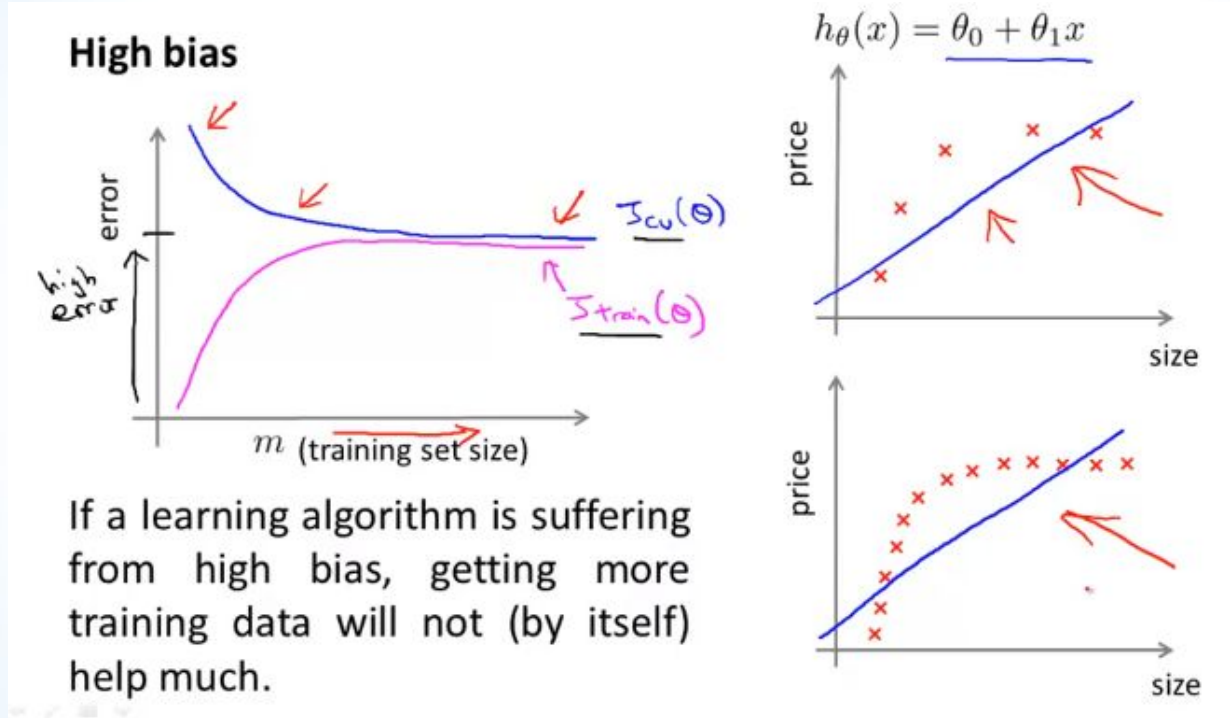
**IMPORTANT  
NOTICE**



# Over and Underfitting

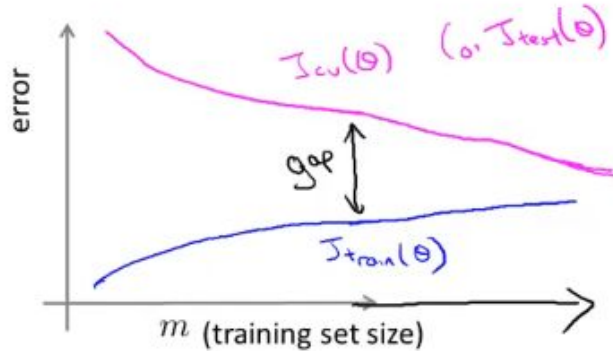


# Bias Variance Analysis - Learning Curve



# Bias Variance Analysis - Learning Curve

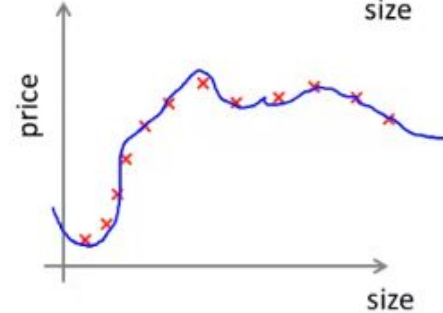
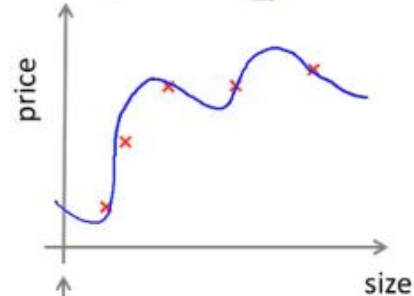
## High variance



If a learning algorithm is suffering from high variance, getting more training data is likely to help. ←

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$$

(and small  $\lambda$ )



# Hypertuning

Choosing the set of “optimal” hyperparameters for a training algorithm

**KNN:**

**Decision Tree:**

# Hypertunning

Choosing the set of “optimal” hyperparameters for a training algorithm

## KNN:

- Number of K
- Distance function
- Distance function params
- Post processing

## Decision Tree:

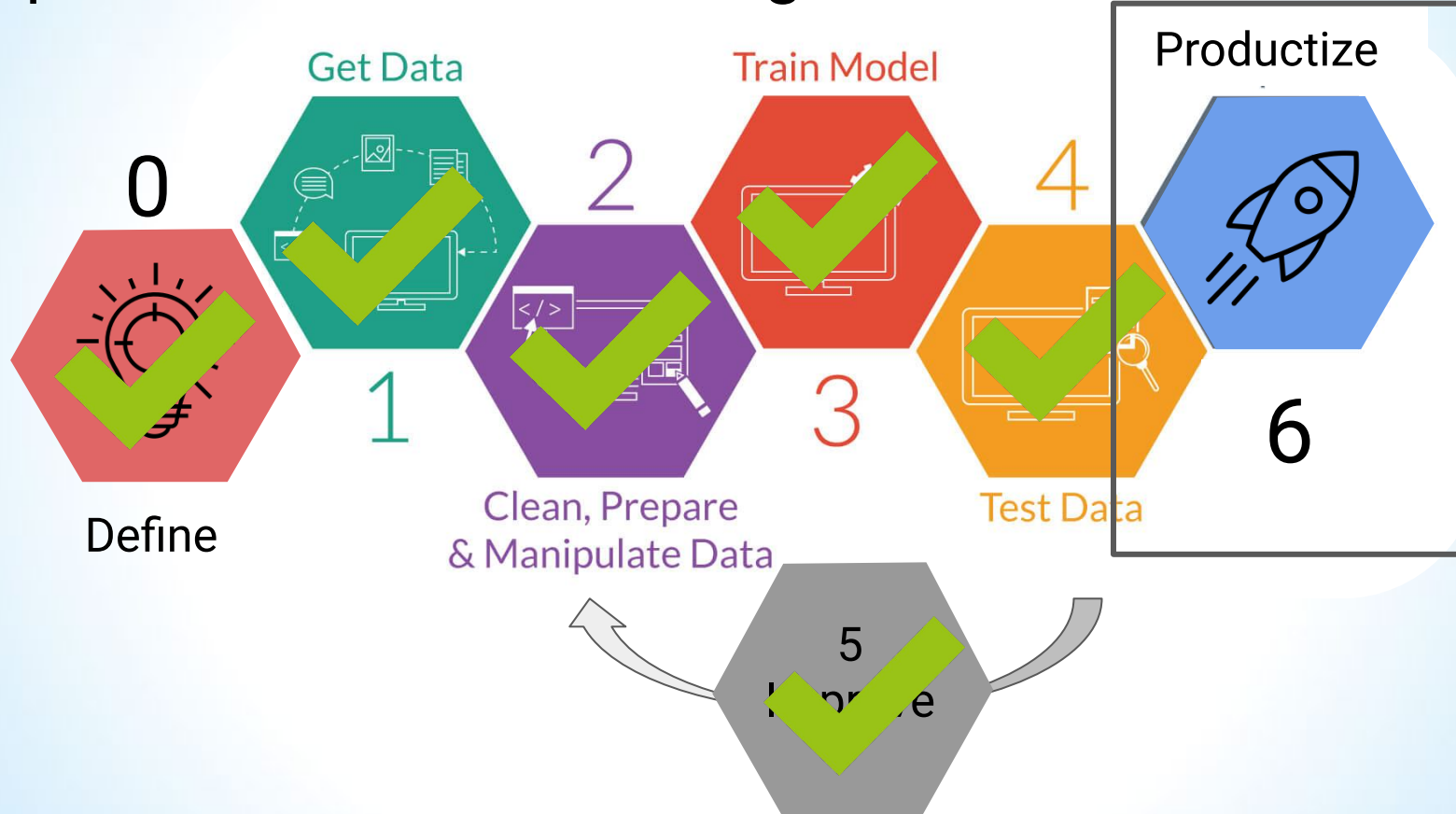
- Split criteria
- Tree depth
- Minimum instances per split
- Data discretization

# Review Homework

Part 2.2 - Finding the optimal  $k$

Part 2.3 - Using cross validation

# Steps to Predictive Modeling



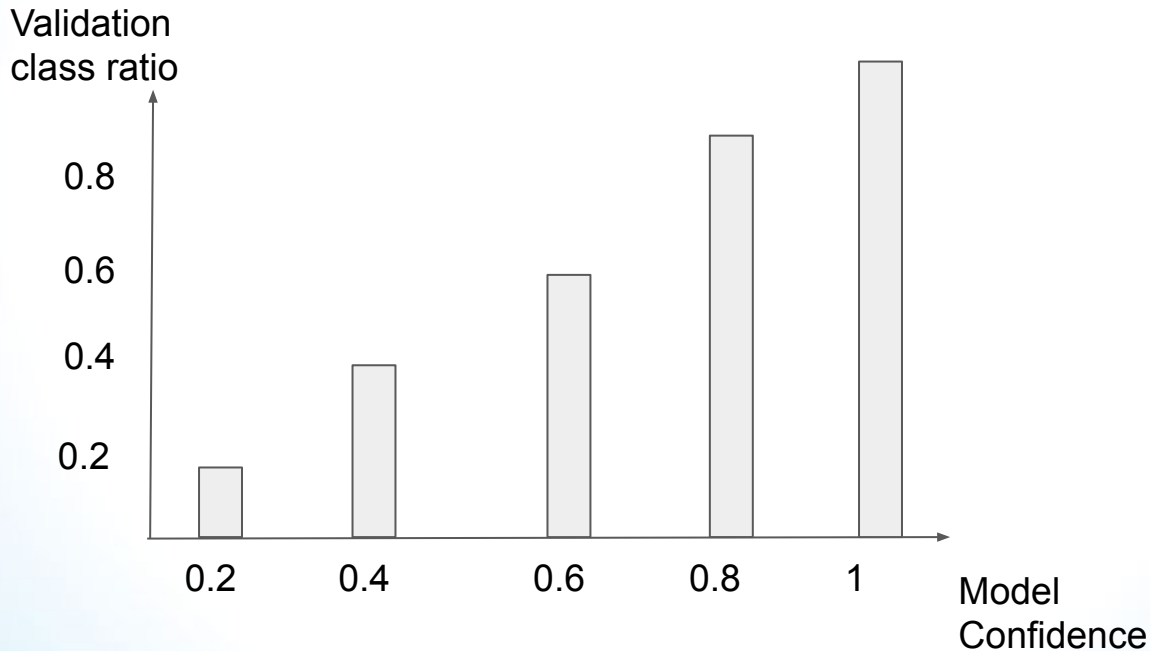


# Productization

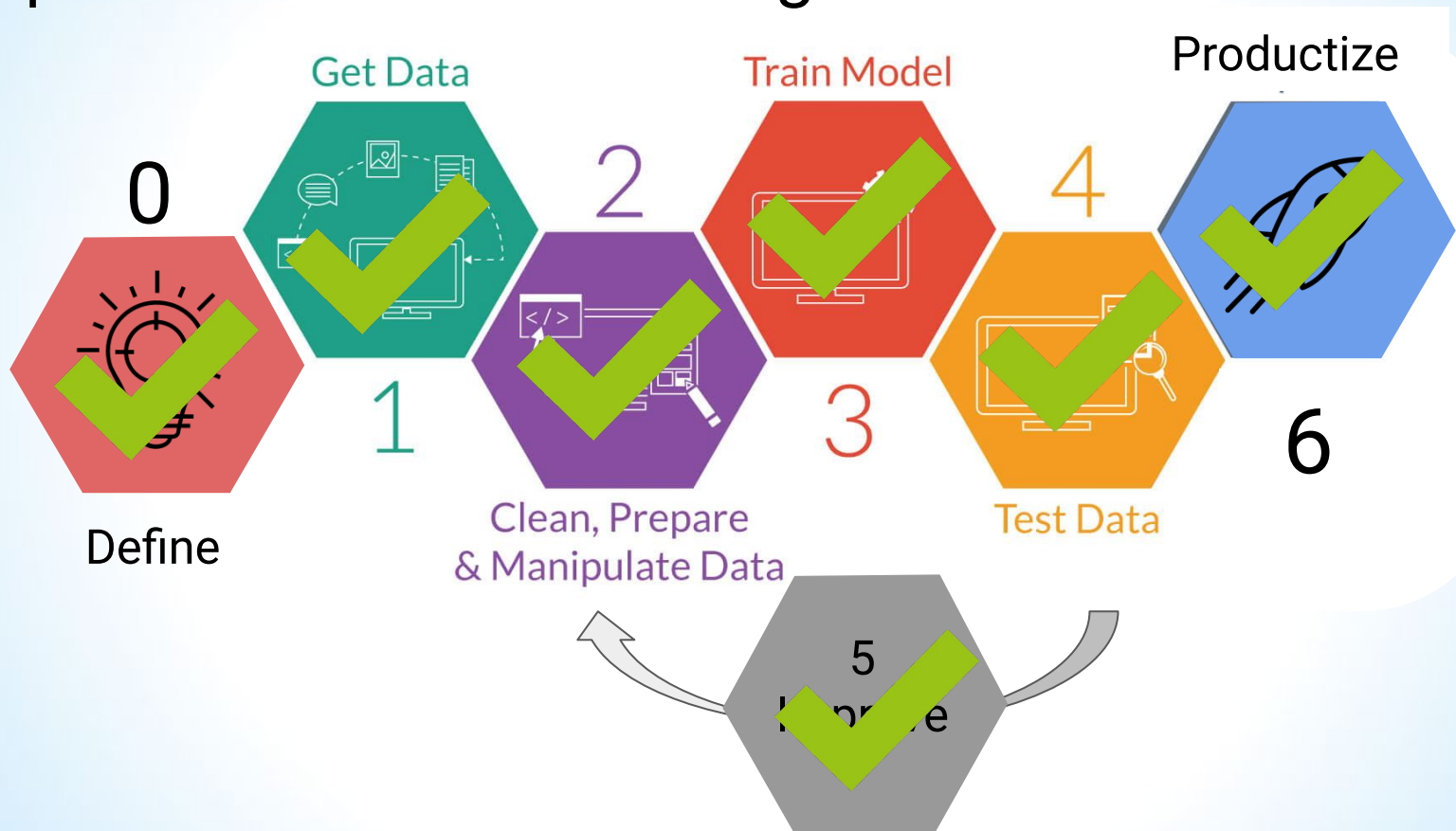
- Make the Model (and the Data) accessible for prediction
- Monitor model's performance - what type of performance?
- Concept drift - changes in data, labels and their relation
- Feedback loop - label more data and retrain the model

# Calibration

Calibration in classification means turning transform classifier scores into class membership probabilities.



# Steps to Predictive Modeling



# TL,DR

- Data Science is a practical profession
- There are many topics -> understand over memorize
- Understand the product task + Have a clear end2end intuition of the project
- Build a Baseline model as soon as possible!
- Keep Asking questions
- There is no silver bullet / free lunch theory

# Reading Materials

1. [A few useful things to know about machine learning.pdf](#)
2. [CIS 419:519 Introduction to Machine Learning.pdf](#)
3. [Empirical Risk Minimization.pdf](#)
4. [Introduction to Statistical Learning Theory.pdf](#)
5. [On the Surprising Behavior of Distance Metrics in High Dimensional Space.pdf](#)
6. [Statistical learning theory - a primer.pdf](#)
7. [Statistical Machine Learning- Introduction.pdf](#)
8. [2012b A Geometrical Explanation of Stein Shrinkage.pdf](#)
9. [INADMISSIBILITY OF THE USUAL ESTIMATOR FOR THE MEAN OF MULTIVARIATE NORMAL DISTRIBUTION - STEIN.pdf](#)
10. [THE USE OF MULTIPLE MEASUREMENTS IN TAXONOMIC PROBLEMS - FISHER - 1936 - Annals of Eugenics - Wiley Online Library.pdf](#)

QA