

Machine Learning 3(2-2-5), 2/2565

# Unit 1

## What is “Data Science” and “Machine Learning”

อ.อานันท์ ไม้ประดิษฐ์

**What is  
DATA SCIENCE?**

# Data + Science

What is Data?

Data is new electricity.

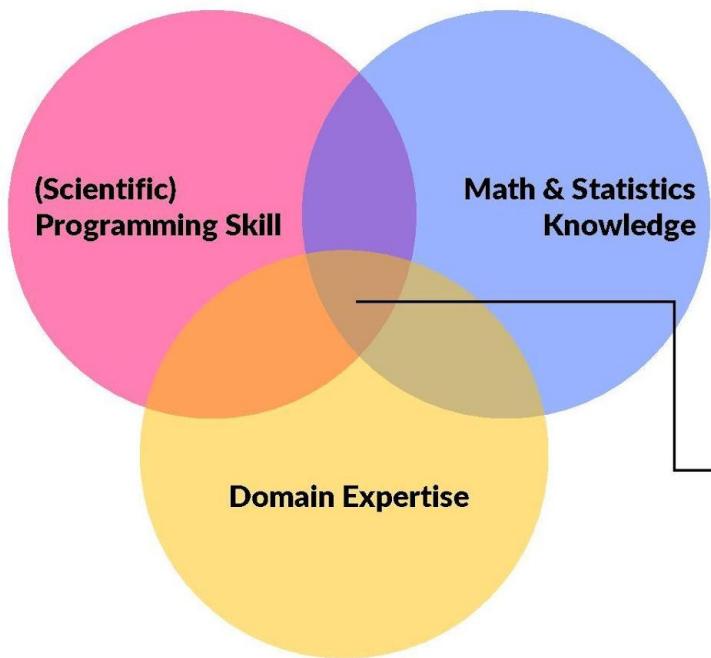
Science as we know it.

Scientists use data to test hypothesis, theory formation.



Andrew Ng

# DS Venn Diagram



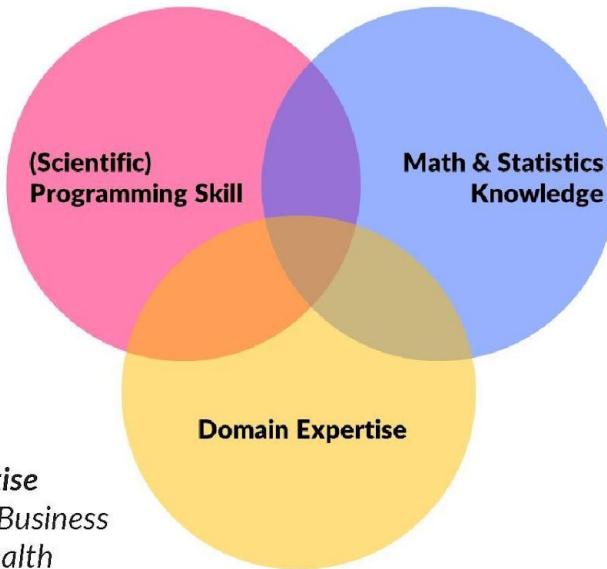
**Data Science**  
*is the intersection of three skills*

<http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>

# Essential Knowledge

## Programming skill

- R
- Python
- SQL
- Spark
- Hadoop



## Domain Expertise

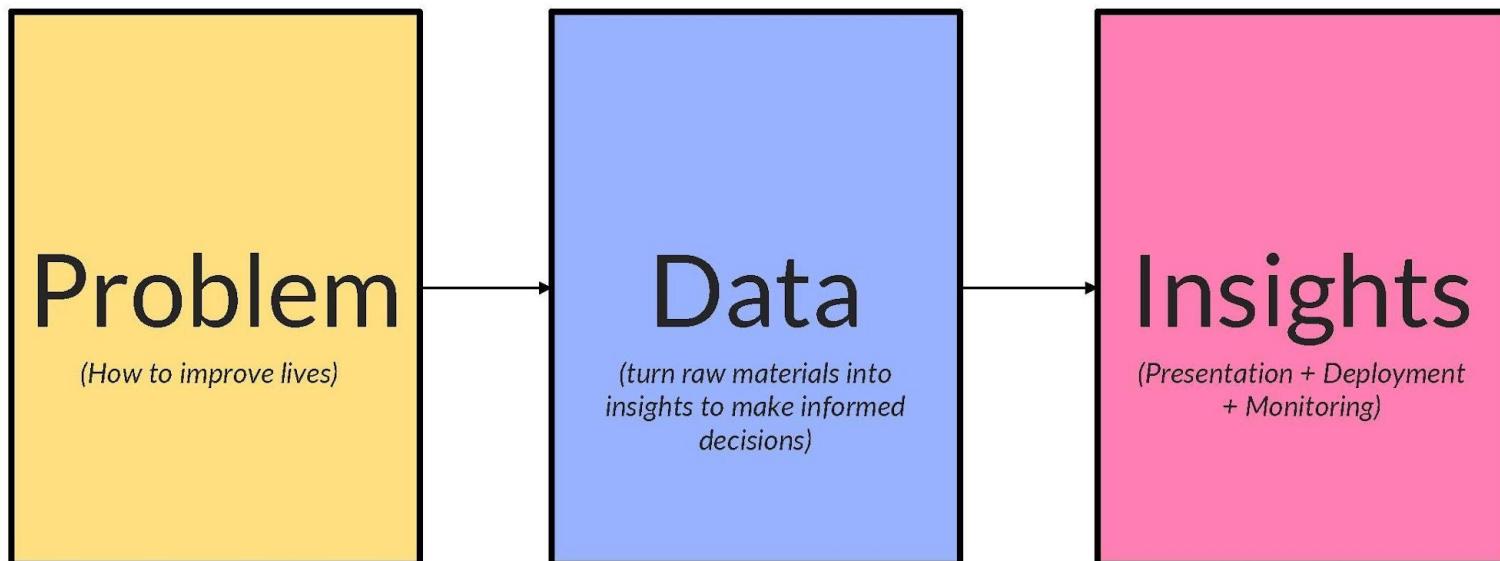
- Marketing | Business
- Medical | Health
- Politic
- Engineer
- Could be anything

## Math & Statistics

- Linear algebra
- Calculus
- Descriptive statistics
- Inferential statistics
- Bayesian statistics
- Probability

Always

Good data science starts with good question

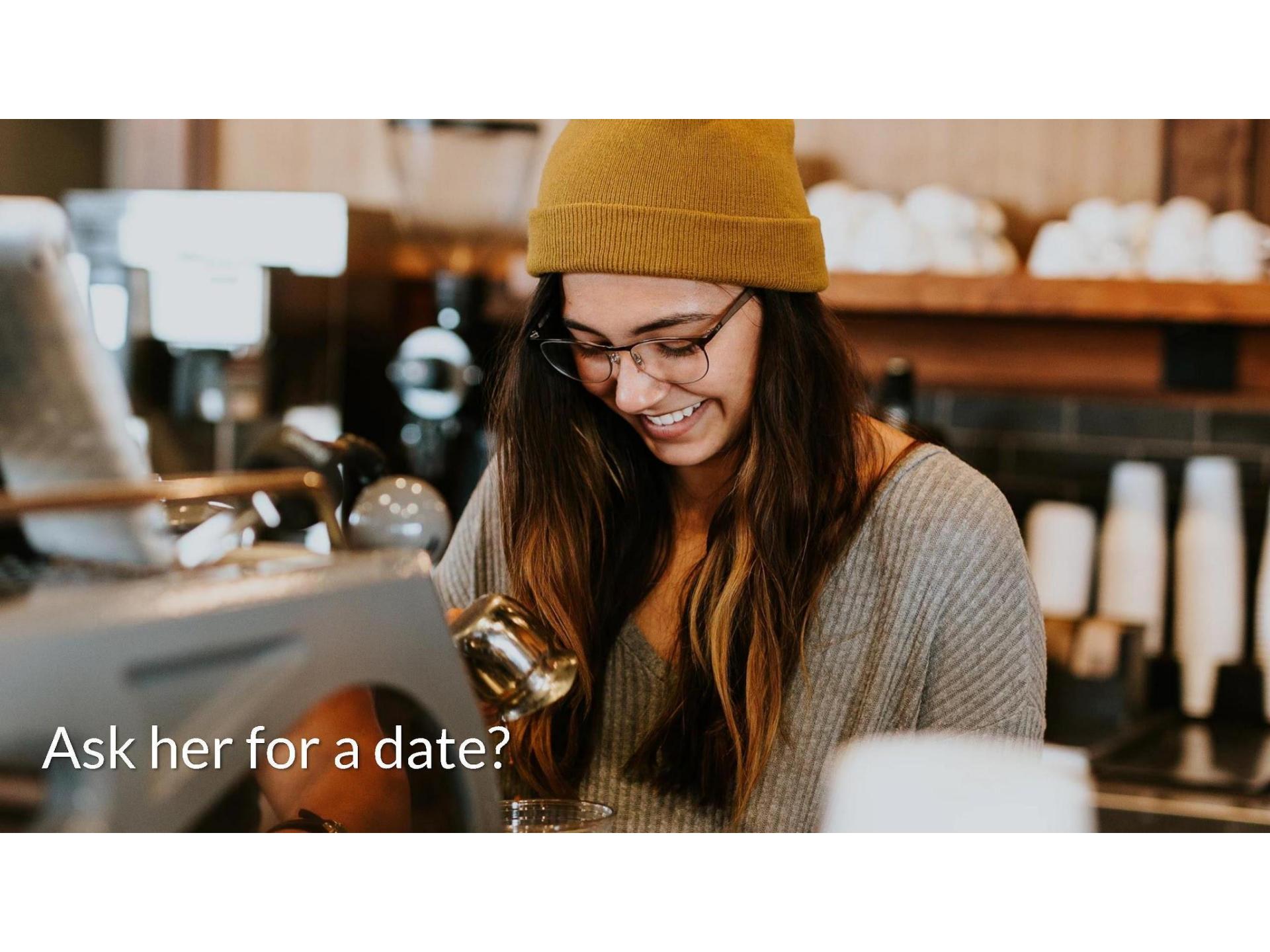


**What is  
MACHINE LEARNING?**

# Can you guess?

x (input)	1	2	3	4	5	6	7	8	9	10
y (output)	1	4	9	16	25	36	49	64	81	<input type="text"/>

You guess the function is  $x^2$

A young woman with long brown hair, wearing a yellow beanie and glasses, is smiling while pouring coffee from a gold-colored pour-over device into a glass cup. She is wearing a grey ribbed sweater. The background shows a blurred coffee shop interior with shelves and equipment.

Ask her for a date?

# Right! Can you guess this one?

	X1	X2	X3	X4	Y
Occasion 1	Restaurant	Many	Friday	Good	Yes
Occasion 2	Pub   Bar	Many	Saturday	OK	Yes
Occasion 3	Pub   Bar	Few	Monday	OK	No
Occasion 4	Restaurant	Few	Friday	Bad	Yes
Occasion 5	Restaurant	Many	Friday	Good	No
Occasion 6	Pub   Bar	Few	Saturday	OK	Yes
Occasion 7	Pub   Bar	Many	Monday	Good	No
Occasion 8	Restaurant	Few	Saturday	OK	No
Occasion 9	Pub   Bar	Many	Friday	Bad	No
Occasion 10	Restaurant	Few	Friday	Good	<input type="text"/>

Yes or No?

**Machine Learning**  
is *function approximation*

# Machine Learning is *inductive reasoning*

Learning from  
experience



## ML algorithms trying to MAP inputs (x's) to output (y)

	X1	X2	X3	X4	Y
Occasion 1	Restaurant	Many	Friday	Good	Yes
Occasion 2	Pub   Bar	Many	Saturday	OK	Yes
Occasion 3	Pub   Bar	Few	Monday	OK	
Occasion 4	Restaurant	Few	Friday	Bad	
Occasion 5	Restaurant	Many	Friday	Good	
Occasion 6	Pub   Bar	Few	Saturday	OK	
Occasion 7	Pub   Bar	Many	Monday	Good	
Occasion 8	Restaurant	Few	Saturday	OK	
Occasion 9	Pub   Bar	Many	Friday	Bad	
Occasion 10	Restaurant	Few	Friday	Good	



**It's good to have lots data but ...**

**REPRESENTATION**

**is the quality we desire**

# ML Example

- Predict breast cancer
- Weather forecast
- Watson analytics
- Google flu trend
- Obama winning the US election



*Predicting crime before it happens*

# A proper definition of ML

“

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



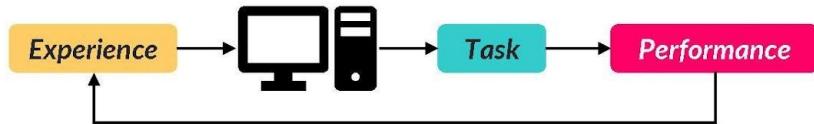
Arthur Samuel (1959)

# A more modern def in 1998

Well posed Learning Problem: A computer program is said to learn from experience **E** with respect to some **task T** and some **performance measure P**, if its performance on **T**, as measured by **P**, improves with experience **E**.



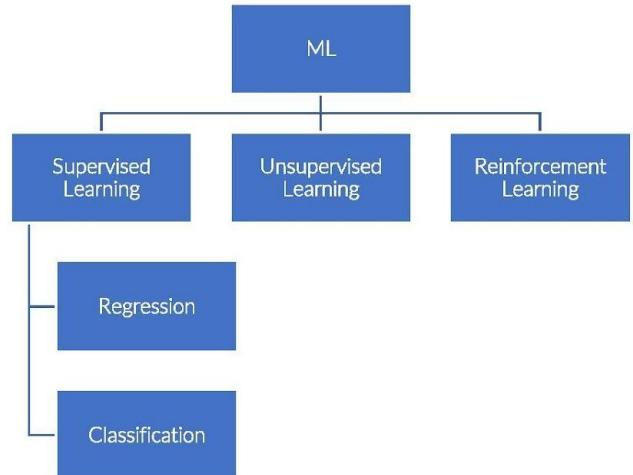
Tom Mitchell (1998)



Can we improve computer's ability to perform task **T** over time (without being explicitly programmed)?

# Types of ML algorithms

1. Supervised Learning
  - i. Regression
  - ii. Classification
2. Unsupervised Learning
3. Reinforcement Learning



# What's Artificial Intelligence



AI

ML

DEEP LEARNING

Function approximation

$$Y = f(X)$$

**Model = Algorithm (Data)**

Predicting new unseen data

In the past

We build models to represent the world (problem we're trying to solve)



# **And the goal of Machine Learning is GENERALIZATION**

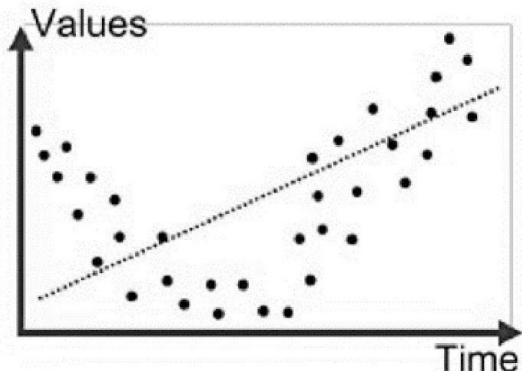
Beyond (past) data you used to build the model

**Generalization means we make sure  
our models **don't overfit** the data.**

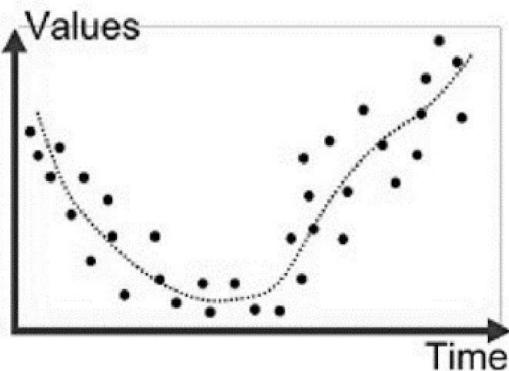
Also don't underfit

We mean  
data in the  
past

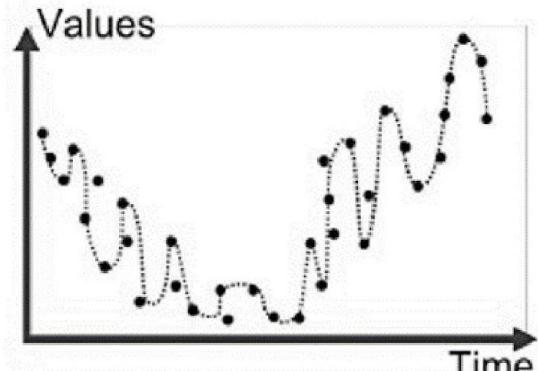
If your model ***can't generalize to new unseen data***, then you're in trouble.



*Underfitting*



*Good fit*



*Overfitting*

[https://cdn-images-1.medium.com/max/1600/1\\*6vPGzBNppqMHllg1o\\_se8Q.png](https://cdn-images-1.medium.com/max/1600/1*6vPGzBNppqMHllg1o_se8Q.png)

This model looks okay as it fit underlying patterns in the data well. Though, not 100% correct ?

# SUMMARY

- Data Science in Thailand is still very young
- ML is data scientist's arsenal tool
- ML is **Function Approximation**, learning from experience
- **Representation & Generalization** (avoid overfitting | underfitting)
- 3 types of ML algorithms including supervised, unsupervised and reinforcement learning
- Regression predicts numbers | Classification predicts categories

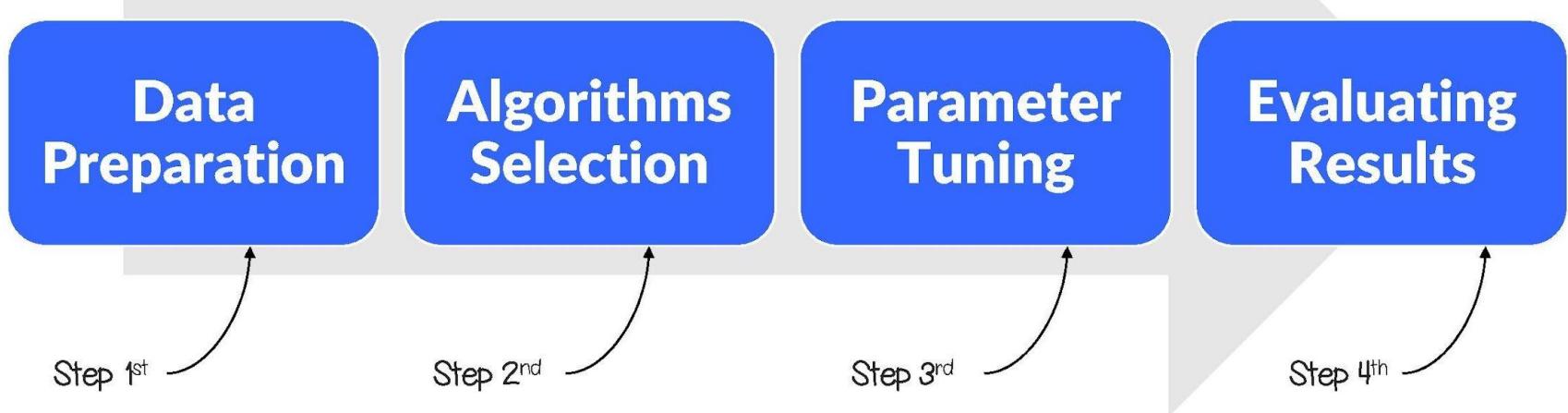
Machine Learning 3(2-2-5), 2/2565

# Unit 2

## Unit 2: Four key steps in Machine Learning

อ.อานันท์ ไม้ประดิษฐ์

# Four key steps in ML



# DATA PREPARATION

# Data Format

The diagram illustrates a data format table with various annotations:

- Variables | Features:** A horizontal double-headed arrow pointing across the columns X1, X2, X3, X4, and Y.
- Target Variable:** A horizontal double-headed arrow pointing specifically at column Y.
- Data Points | Observations:** A vertical double-headed arrow pointing down through all rows of the table.

	X1	X2	X3	X4	Y
Occasion 1	Restaurant	Many	Friday	Good	Yes
Occasion 2	Pub   Bar	Many	Saturday	OK	Yes
Occasion 3	Pub   Bar	Few	Monday	OK	No
Occasion 4	Restaurant	Few	Friday	Bad	Yes
Occasion 5	Restaurant	Many	Friday	Good	No
Occasion 6	Pub   Bar	Few	Saturday	OK	Yes
Occasion 7	Pub   Bar	Many	Monday	Good	No
Occasion 8	Restaurant	Few	Saturday	OK	No
Occasion 9	Pub   Bar	Many	Friday	Bad	No
Occasion 10	Restaurant	Few	Friday	Good	No

# **Variable Types**

1. Binary (1/0)
2. Categorical
3. Integer
4. Continuous

# Feature Selection

- Choose relevant set of features that represent the problem
- Curse of dimensionality:  $2^n$

# Feature Engineering

- Recode
- Combine multiple variables (PCA)
- Extract new information from existing feature

# Gentle intro to Feature Engineering

ID	Name	DOB	Math	Stats	Art	Music	Design
1	Mr. David Beckham	13 / 05 / 1975	85	90	50	42	49
2	Ms. Angelina Jolie	20 / 09 / 1988	50	60	90	95	88
3	Mr. Khalid Khan	05 / 01 / 1950	32	65	70	72	65
4	Mr. Hibino Takamoto	31 / 12 / 1995	75	76	80	50	30

**Question:**

**What new variable can we create from this dataset?**

**Think at least 5 new variables**

We can use  
`dplyr::mutate( )` to  
create new variables

# Missing Data

1. Approximated
2. Predicted (Supervised Learning)
3. Removed

Mean | Median  
imputation

If we have large data points and NA is not much, we may consider removing them

# Approximated or Removed?

```
Console ~/Desktop/R PROGRAMMING/
> head(complete.cases(biopsy), 200)
 [1] TRUE  TRUE 
[14] TRUE  FALSE TRUE 
[27] TRUE  TRUE 
[40] TRUE  FALSE TRUE  TRUE 
[53] TRUE  TRUE 
[66] TRUE  TRUE 
[79] TRUE  TRUE 
[92] TRUE  TRUE 
[105] TRUE  TRUE 
[118] TRUE  TRUE 
[131] TRUE  TRUE  TRUE  TRUE  TRUE  TRUE  TRUE  TRUE  TRUE  FALSE TRUE  TRUE  TRUE 
[144] TRUE  TRUE  FALSE TRUE  TRUE  TRUE  TRUE  TRUE  TRUE  TRUE  TRUE  TRUE  TRUE 
[157] TRUE  TRUE  FALSE TRUE  TRUE  TRUE  TRUE  FALSE TRUE  TRUE  TRUE  TRUE  TRUE 
[170] TRUE  TRUE 
[183] TRUE  TRUE 
[196] TRUE  TRUE  TRUE  TRUE  TRUE
```

R functions:

`complete.cases(df)` สำหรับตรวจสอบว่าข้อมูลมี NA หรือเปล่า?

`tidyverse::drop_na(df)` สำหรับ drop case ที่มีข้อมูลไม่ครบถ้วน

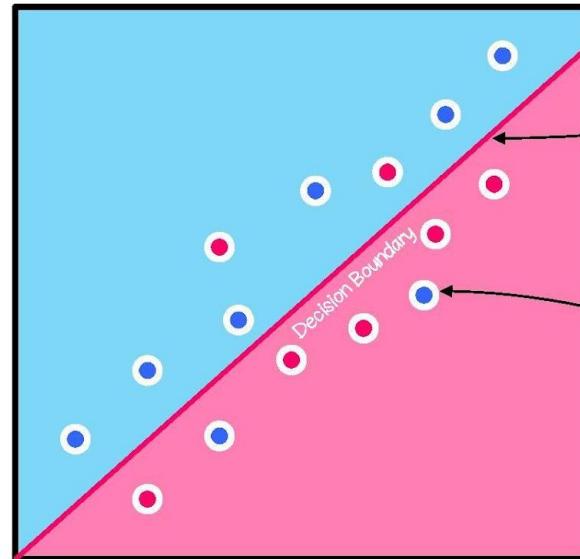
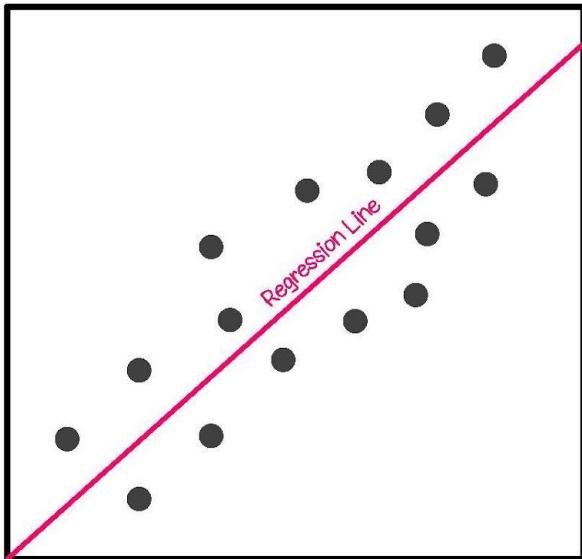
# ALGORITHM SELECTION

# There are 3 types of algorithms

	Algorithms	Main Tasks
<b>Unsupervised learning</b>	K-Means Principal Component Analysis (PCA) Association Rules Social Network Analysis (SNA)	Tell me what patterns exist in my data - <b>descriptive</b>
<b>Supervised learning</b>	Linear Regression Logistic Regression Decision Tree Random Forests KNN Support Vector Machine Neural Networks	Use the pattern in my data to make predictions - <b>predictive</b>
<b>Reinforcement learning</b>	Multi-Armed Bandits (UBC, Thompson Sampling)	Use the pattern in my data to make predictions and - <b>improve</b> - these predictions as more results come in

Source: Data Science For The Layman (2017)

# Regression vs. Classification



**Decision Boundary:**  
The line that best split  
red vs. blue data points

100% in reality may not  
be feasible, some blue  
points are mistakenly  
classified as red

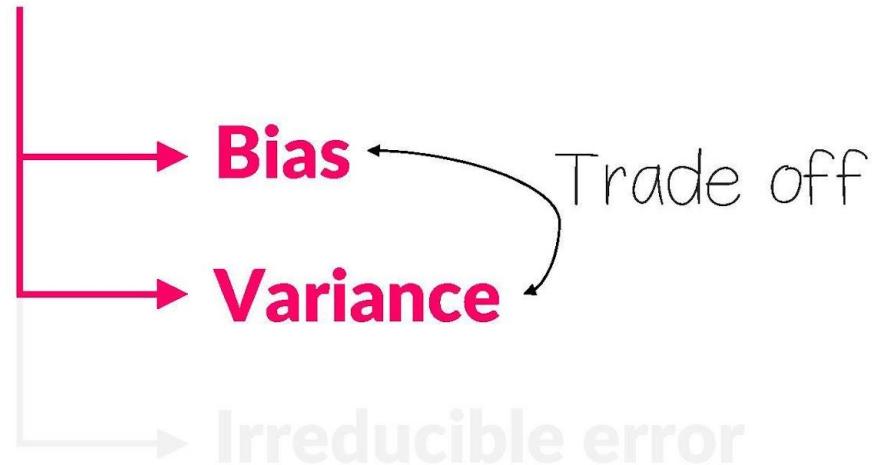
# NO FREE LUNCH

- Speed
- Accuracy
- Bias & Variance Tradeoff



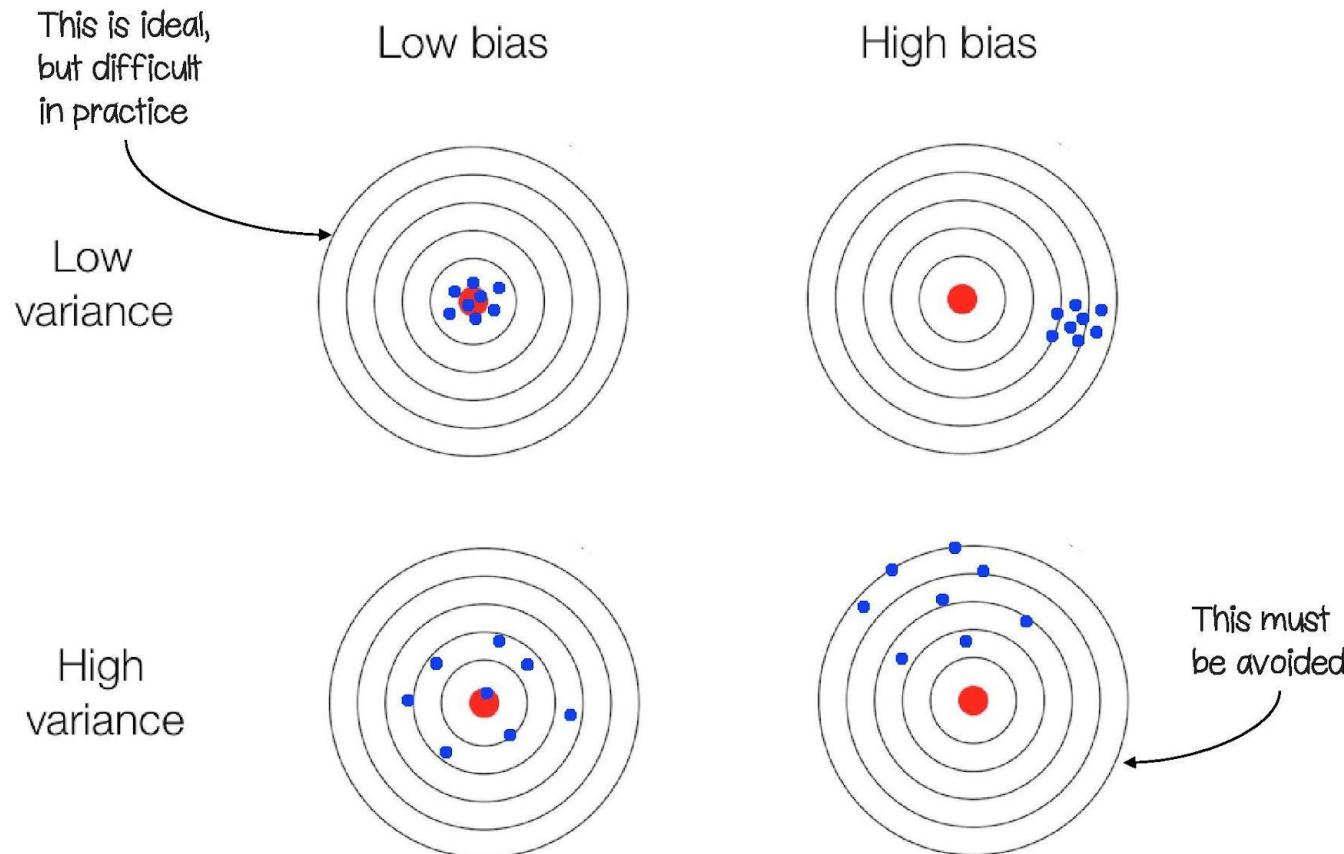
Our model consists of 3 types of error

$$Y = f(X) + \text{Error}$$



Low variance algorithm(high bias): → Less complex, consistent but inaccurate on average → underfit

Low bias algorithm(high variance): → More complex, accurate on average but inconsistent → overfit



# Example Models

<b>High Bias</b> <i>Having lots of assumptions to do function approximation</i>	<b>Linear Regression</b> <b>Logistic Regression</b> <b>Linear Discriminant Analysis</b>	}	<b>Low Variance</b> <i>New data does not cause model to change much</i>
<b>Low Bias</b> <i>Having little or few assumptions about the function approximation</i>	<b>Decision Tree</b> <b>KNN</b> <b>SVM</b>		<b>High Variance</b> <i>New data causes model to change much i.e. prediction results vary significantly compared to low variance algorithms</i>

# PARAMETER TUNING

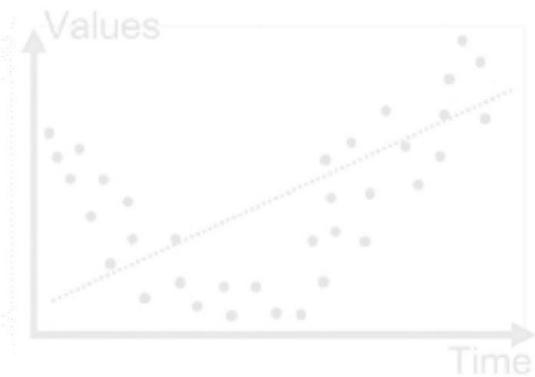
**Music Composer is  
a very good example**



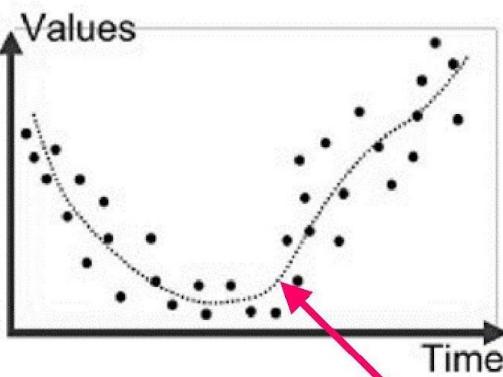
**There are so many learners.**

**Even the same learner can have  
different performances.**

# We tune parameter to get the best results (**GOOD FIT**) from our learner.



*Underfitting*



*Good fit*



*Overfitting*

The parameters that can be tuned called "**HYPERPARAMETERS**"

[https://cdn-images-1.medium.com/max/1600/1\\*6vPGzBNppqMHlg1o\\_se8Q.png](https://cdn-images-1.medium.com/max/1600/1*6vPGzBNppqMHlg1o_se8Q.png)

# Some models have **hyperparameters**, some don't.

## Stochastic Gradient Boosting

```
method = 'gbm'
```

Type: Regression, Classification

Tuning parameters:

- n.trees (# Boosting Iterations)
- interaction.depth (Max Tree Depth)
- shrinkage (Shrinkage)
- n.minobsinnode (Min. Terminal Node Size)



*Hyperparameters of GBM algorithms*

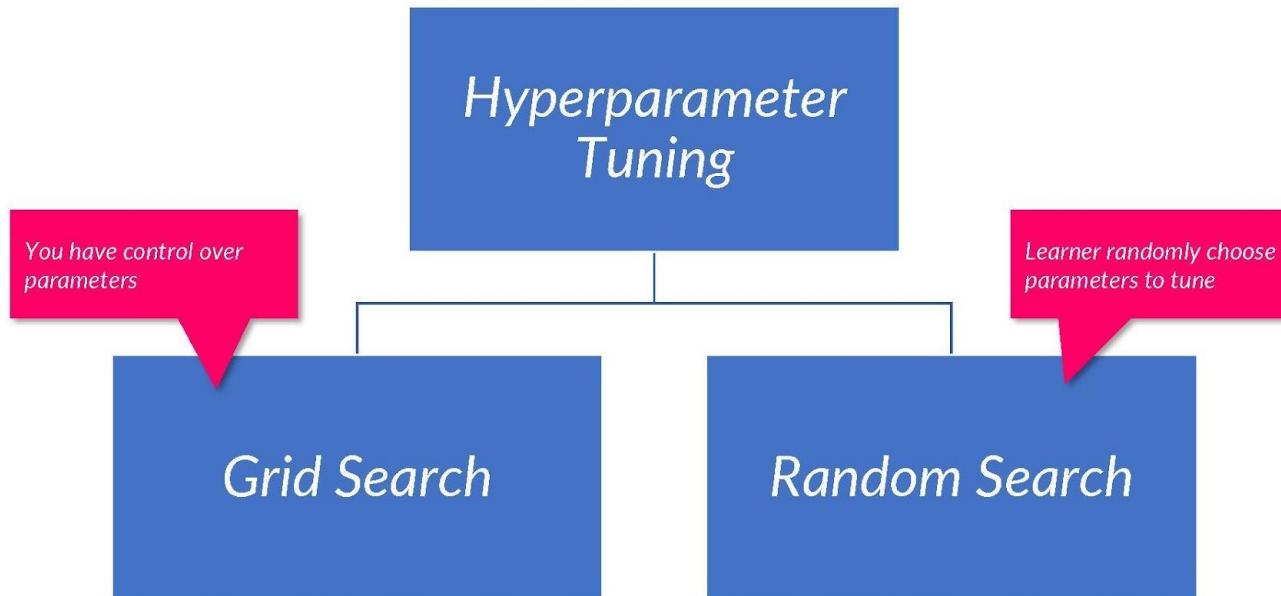
Required packages: gbm , plyr

A model-specific variable importance metric is available.

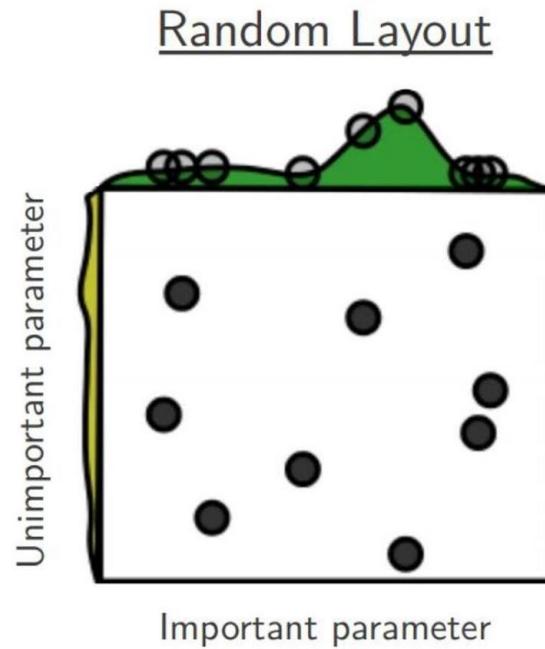
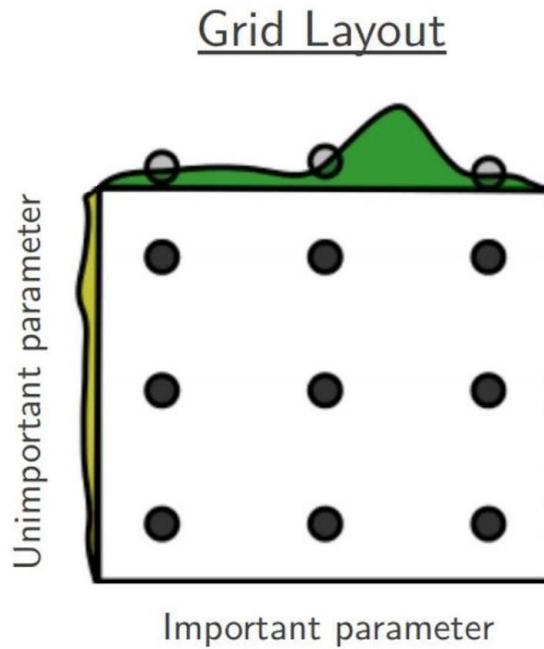
<http://topepo.github.io/caret/train-models-by-tag.html>

# **SEARCH** for the best tune

ML is optimization problem



**It really depends on the learners. Sometimes it's more efficient to use random search, sometimes it's not.**



<http://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf>

# EVALUATING RESULTS

# Performance Metrics

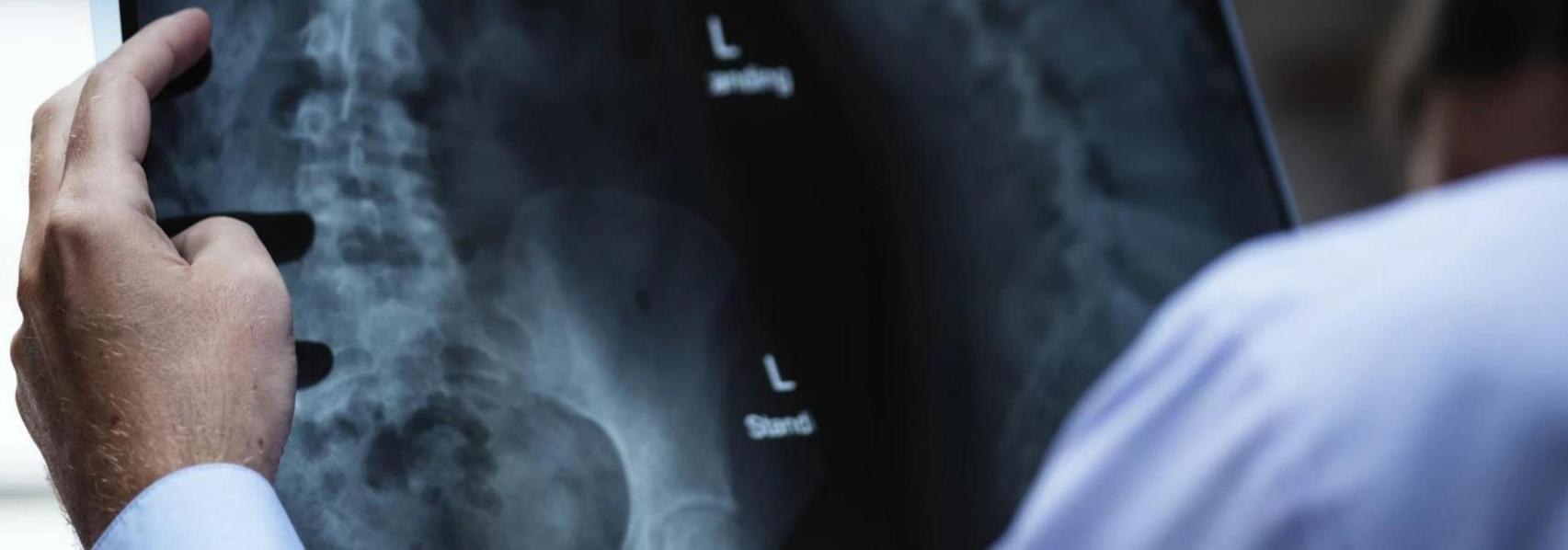
## Classification problems

- *Confusion Matrix*
- *% Accuracy*
- *Area Under Curve (AUC)*
- *Logarithmic Loss (Log Loss)*
- *Cohen's KAPPA #for imbalanced problems*

## Regression problems

- *Mean Absolute Error (MAE)*
- *Root Mean Squared Error (RMSE)*

# Cancer Detection



# Confusion Matrix

Total n = 100	Reference	
Predicted	Cancer	No Cancer
Cancer	60	3
No Cancer	2	35

Accuracy is the sum of diagonal divided by total sample size

$$60+35 / 100 = 95\%$$

# Confusion Matrix

Total n = 100		Reference	
Predicted	Cancer	No Cancer	
Cancer	60 ( <b>A</b> )	3 ( <b>B</b> )	
No Cancer	2 ( <b>C</b> )	35 ( <b>D</b> )	

ถ้าผู้ป่วยเป็นมะเร็ง เราหายถูกเท่าไร? (%)

$$\text{Sensitivity} = \mathbf{A} / \mathbf{A+C}$$

ถ้าผู้ป่วยไม่ได้เป็นมะเร็ง เราหายถูกเท่าไร? (%)

$$\text{Specificity} = \mathbf{D} / \mathbf{B+D}$$

ผล prediction "cancer" ของเราถูกกี่ %

$$\text{Precision} = \mathbf{A} / \mathbf{A+B}$$

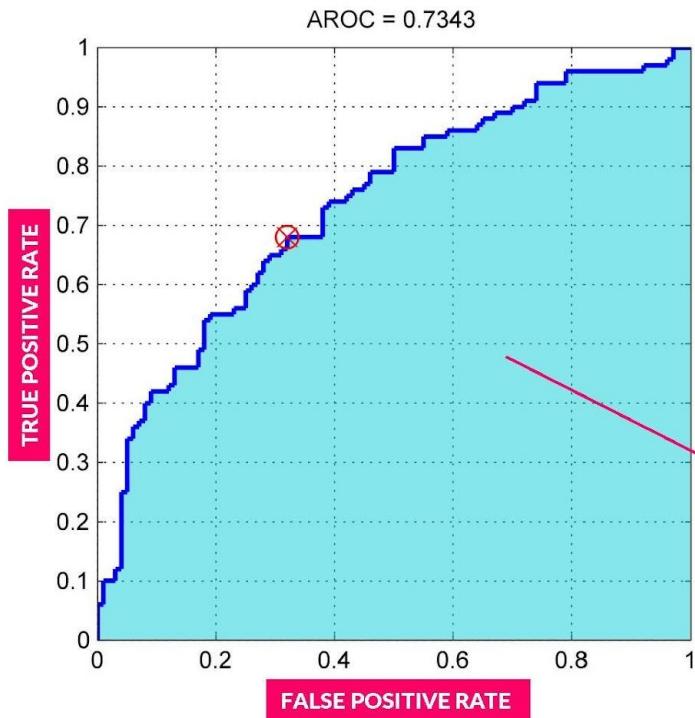
ในผู้ป่วยโรคมะเร็งทั้งหมด เรายำนาຍถูกเท่าไร? (%)

$$\text{Recall} = \mathbf{A} / \mathbf{A+C}$$

Measure of overall test's accuracy

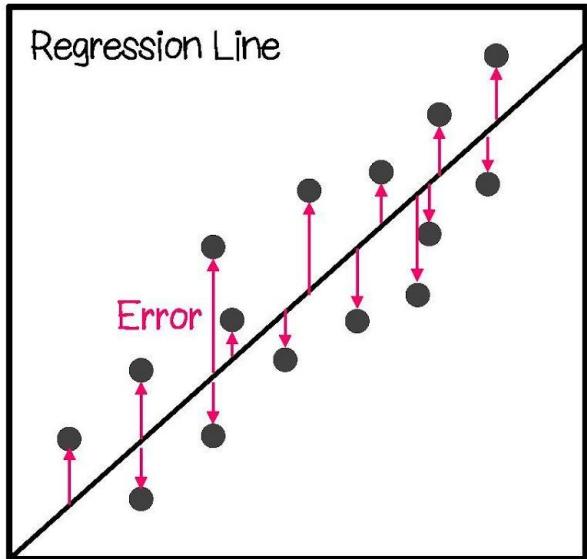
$$\text{F1} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

# Area Under Curve (AUC)



*For binary classification, sometimes we want to maximize AUC*

# Root Mean Squared Error (RMSE)



Mathematical Expression always look scary

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (Actual - Predicted y_i)^2}{N}}$$

$$RMSE = \sqrt{\frac{\text{sum(error}^2)}{N}}$$

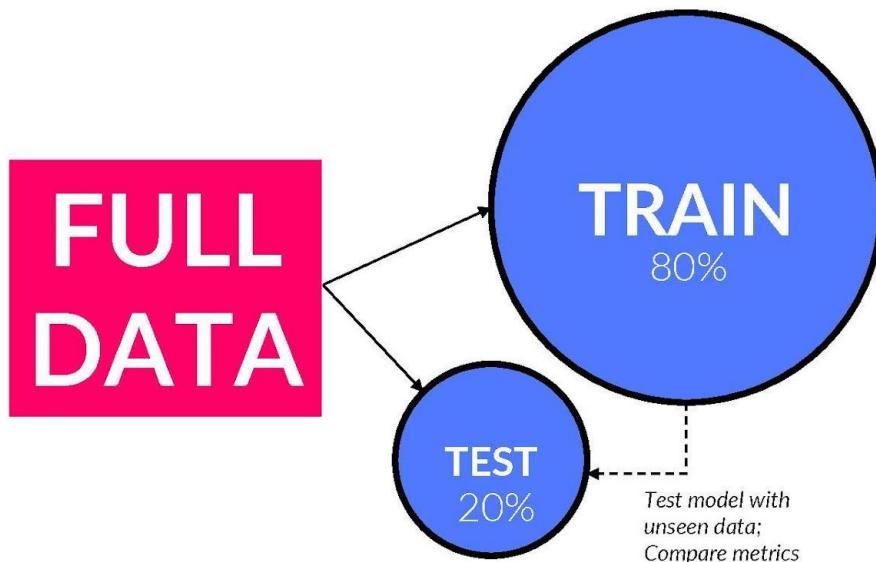
**Remember this one**

$$RMSE = \sqrt{\text{mean } (error^2)}$$

**Recall the goal of ML is  
“GENERALIZATION”**

**We must test our model on new  
unseen dataset**

# Is our model good to go?



*Training set (build model)*

*Testing set (evaluate model)*

## Decision Criterion:

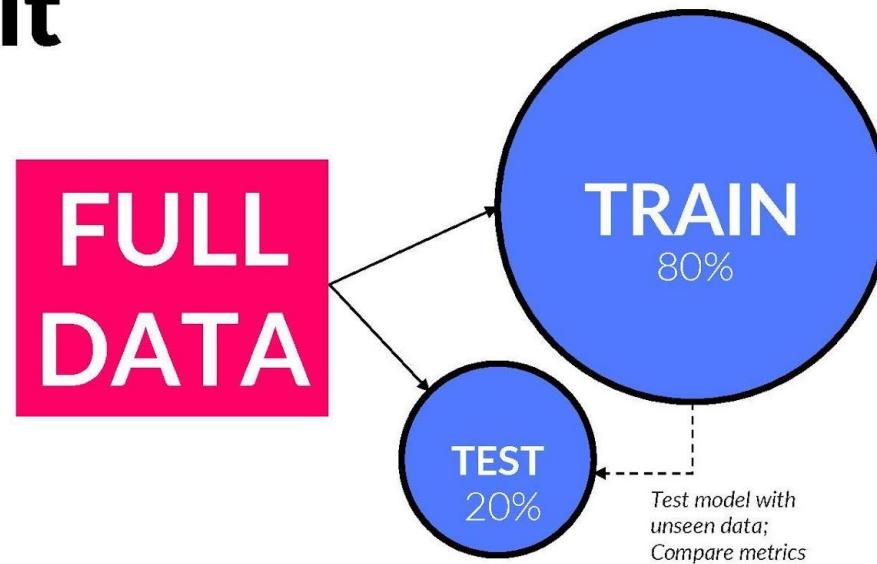
*The model **passes the test** if  
tested performance is similar  
trained performance [metric]*

# Resampling techniques

1. Train-test split
2. *k*-Fold Cross-Validation
3. Repeated *k*-Fold Cross-Validation
4. Bootstrap
5. Leave One Out Cross-Validation

*This method is preferred  
but it's very expensive for  
large dataset*

# The most efficient method is train-test split



# K-Fold Cross Validation

FULL DATASET

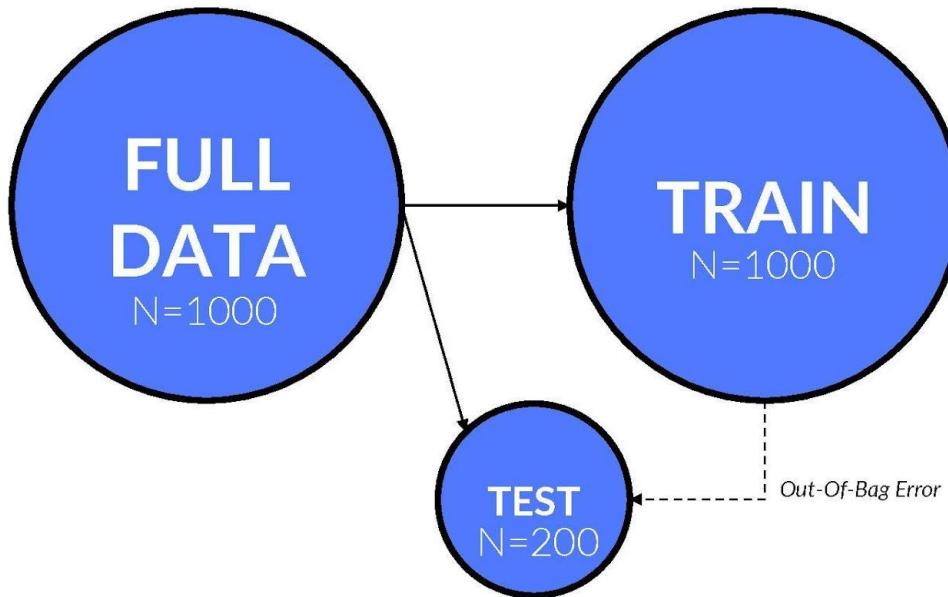


[1]	Train	Train	Train	Train	Test
[2]	Train	Train	Train	Test	Train
[3]	Train	Train	Test	Train	Train
[4]	Train	Test	Train	Train	Train
[5]	Test	Train	Train	Train	Train

People normally  
do 5, 10 K-Fold

# Bootstrap Sampling

(aka. Random sampling with replacement)

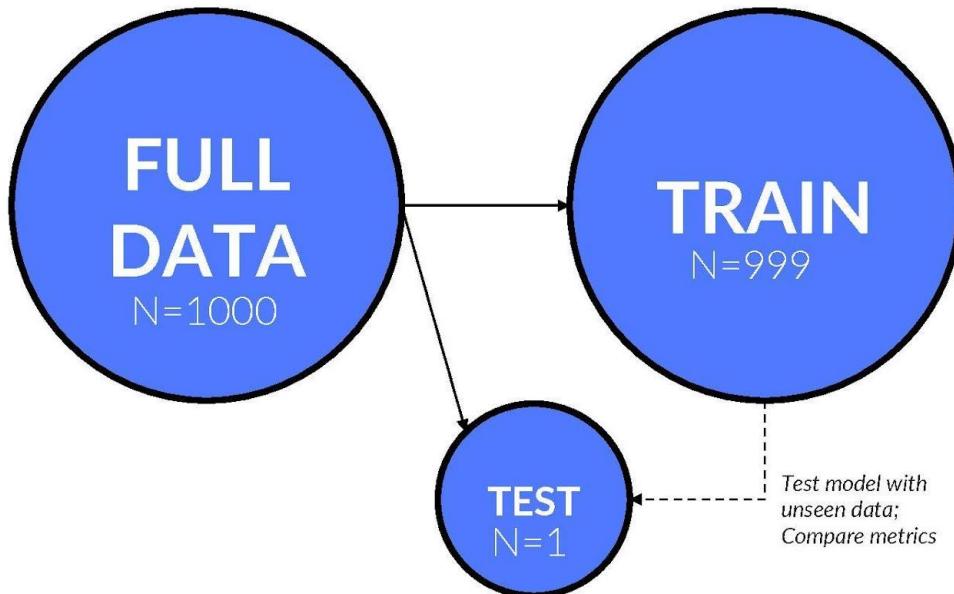


## Explanation:

*1000 may come from 800 unique data points + 200 repeated data points*

# Leave One Out CV

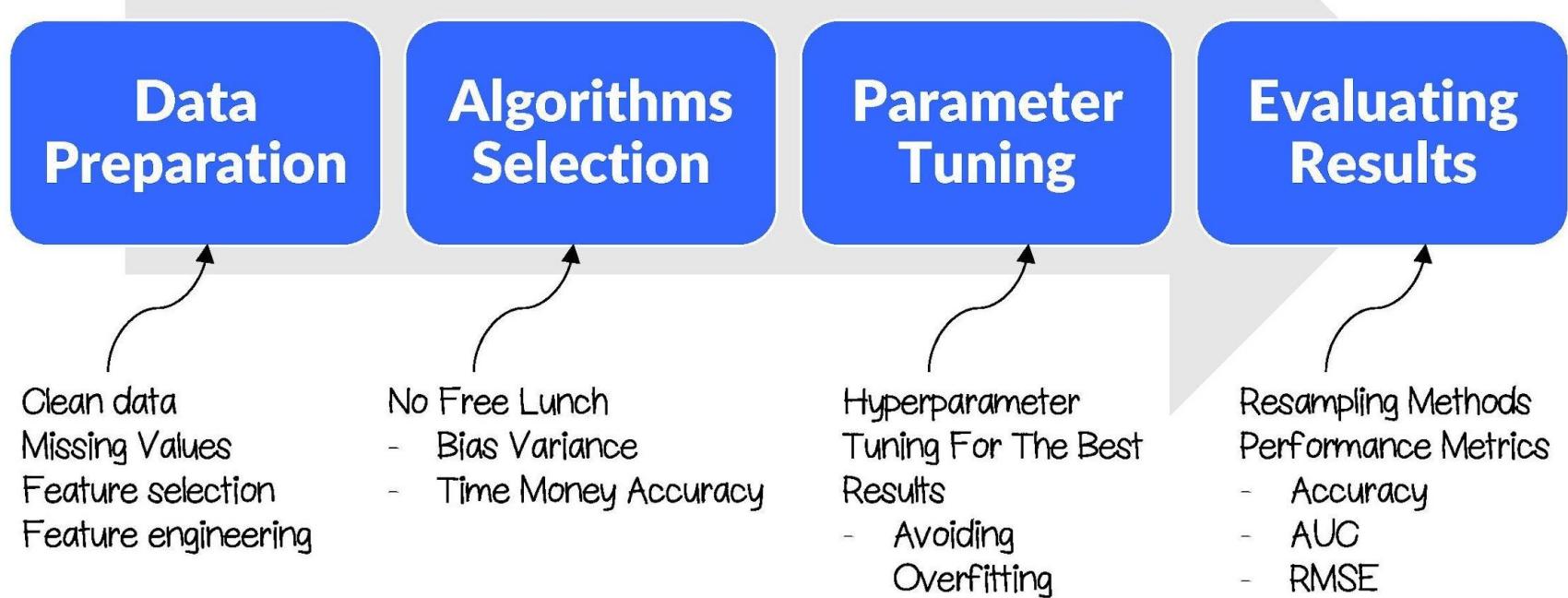
This is the most expensive method – time consuming



## Explanation:

Leave 1 data point to test the model. For  $n=1000$ , this means we need to run model 1000 times

# #SIMPLE :)



# SUMMARY

- Data preparation takes 70-80% of our time.
  - “Garbage in, Garbage out”
  - Feature selection | Feature engineering | Missing Value
- No free lunch theorem
  - Bias vs. Variance Tradeoff
  - Speed vs. Accuracy Tradeoff
- [Hyper]Parameter tuning to get the best out of our learning algorithms
- Resampling techniques used in ML research incl. train-test, k-fold CV, repeated k-fold CV, bootstrap, and LOOCV

# OKAY !! Let's recap a bit

- ML is **function approximation**
- Learn from past data to **predict future (unseen) data point**
- The goal is **representation & generalization** (avoid overfitting)
- Regression or Classification problems depend **target variable ( $y$ )**
- **No Free Lunch**: bias vs. variance tradeoff
- Algorithms can be tuned via **hyperparameter**
- **K-Fold cross validation** is popular resampling technique
- Our model passes the test if **training vs. testing performances are similar** (not overfit)

# **ALGORITHMS**

## **YOU SHOULD KNOW**

# We cover 8 algorithms

Linear Regression	Regression
Logistic Regression	Classification
Decision Tree	Regression & Classification
Random Forest	(mainly) Classification
KNN	(mainly) Classification
Support Vector Machine	(mainly) Classification
Neural Networks	Regression & Classification
K-Means	Unsupervised Learning

Supervised Learning

NN is special algorithm, it works like reinforcement learning

Machine Learning 3(2-2-5), 2/2565

# Unit 3

## Supervised Learning

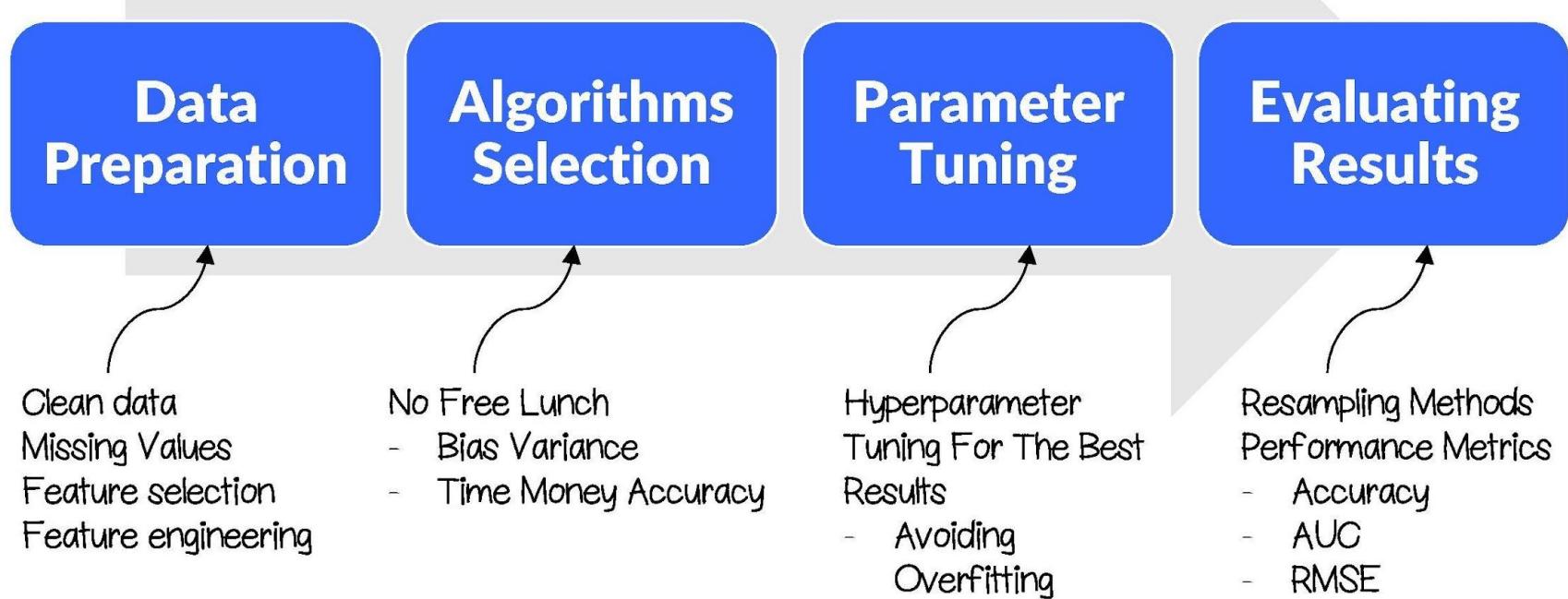
อ.อานันท์ ไม้ประดิษฐ์

# LINEAR REGRESSION

Hello World in  
Machine Learning

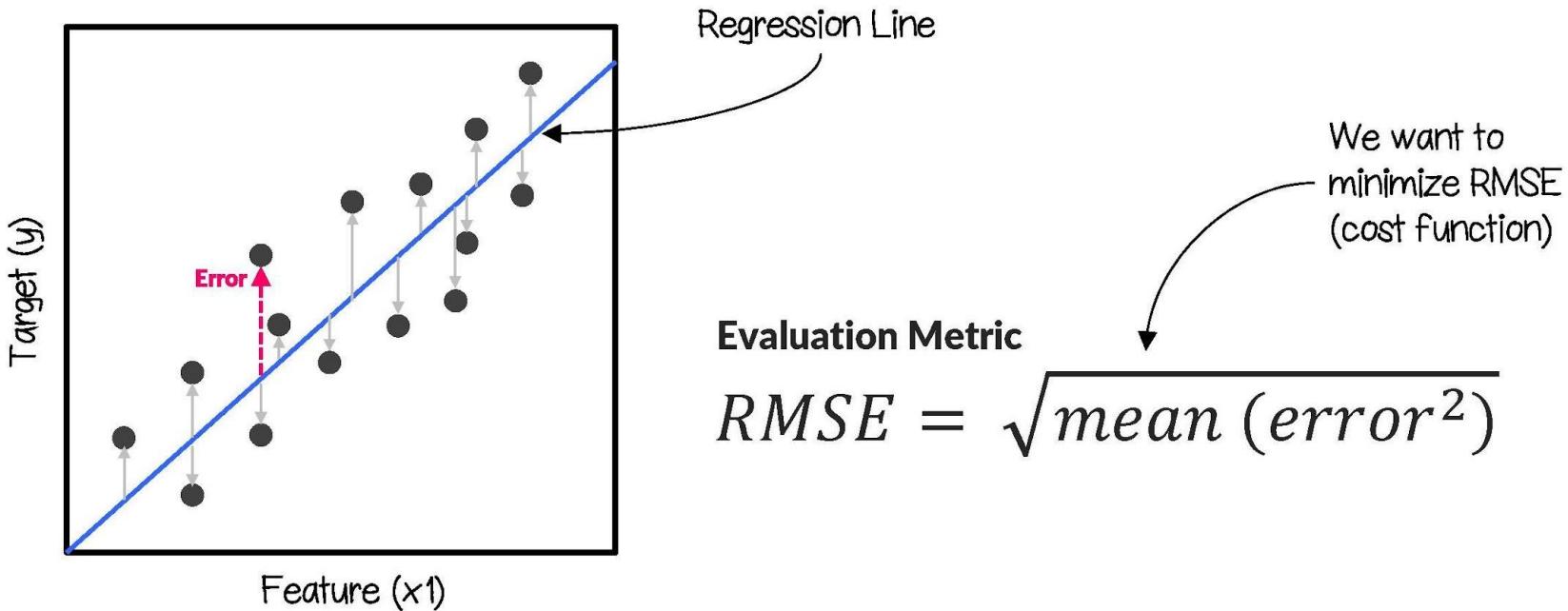


# Recall these steps

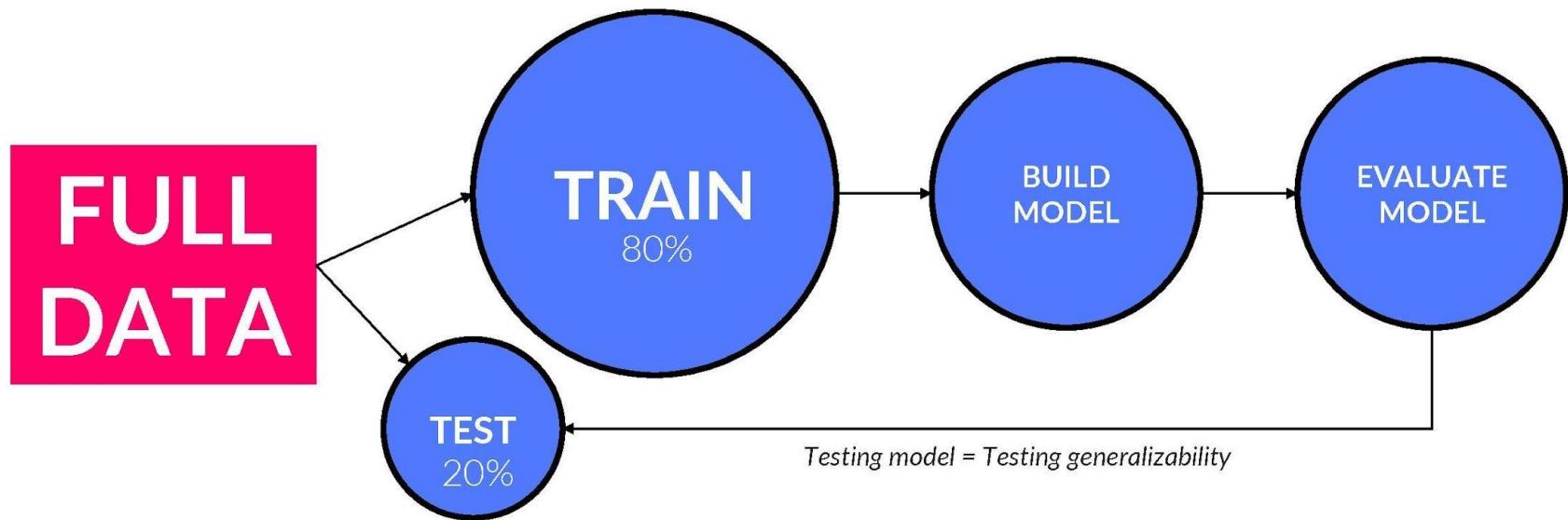


# Algorithm explained:

How to draw a single straight line that fit best with given data points



Let's do a **train-test split** for linear regression problem



## **Advantages:**

- *Easy to train*
- *Always find the optimum solutions (least squared)*
- *Foundation of advanced algorithms like neural nets*

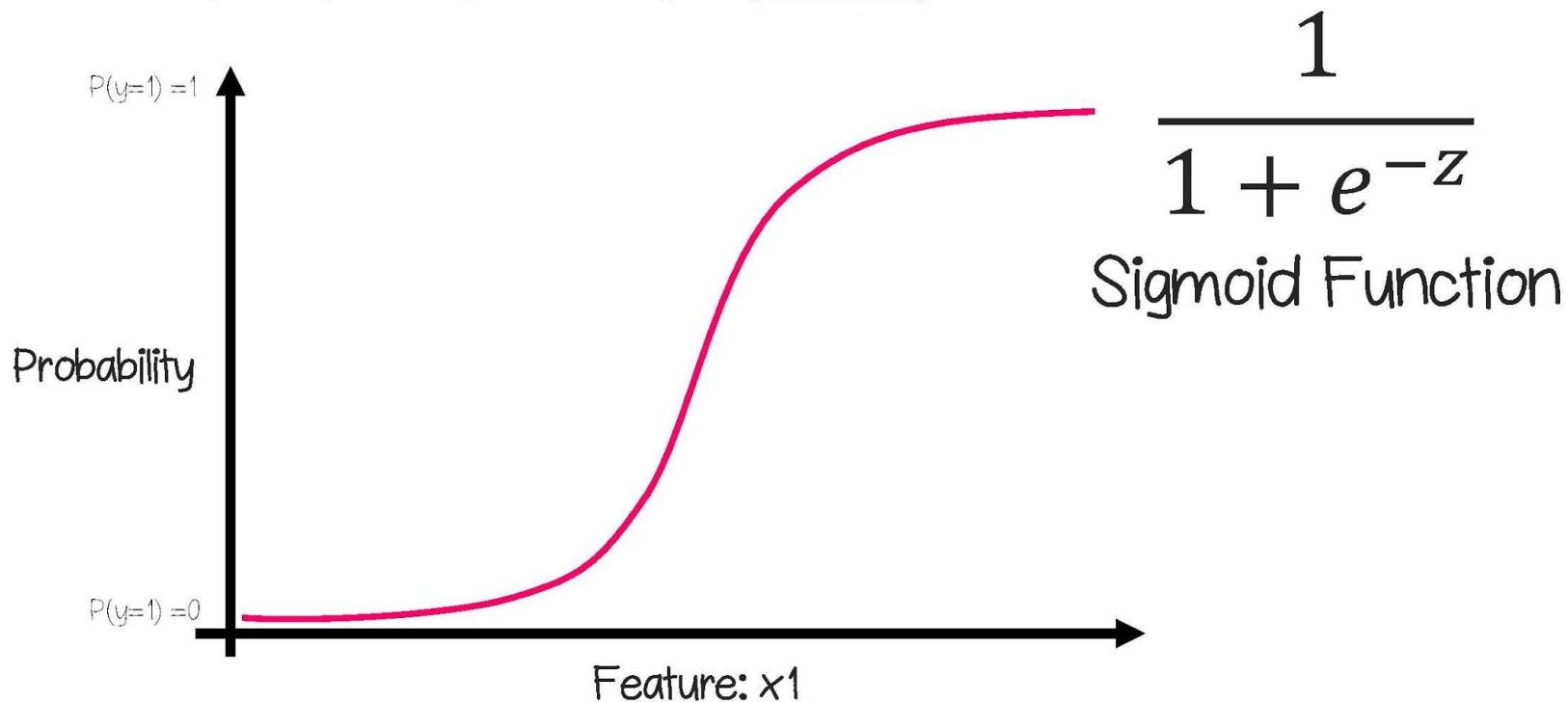
## **Disadvantages:**

- *Only work for regression problems (continuous y)*
- *Full of assumptions (high bias algorithms)*

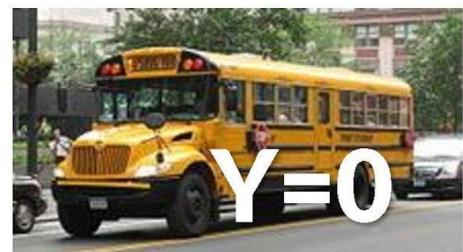
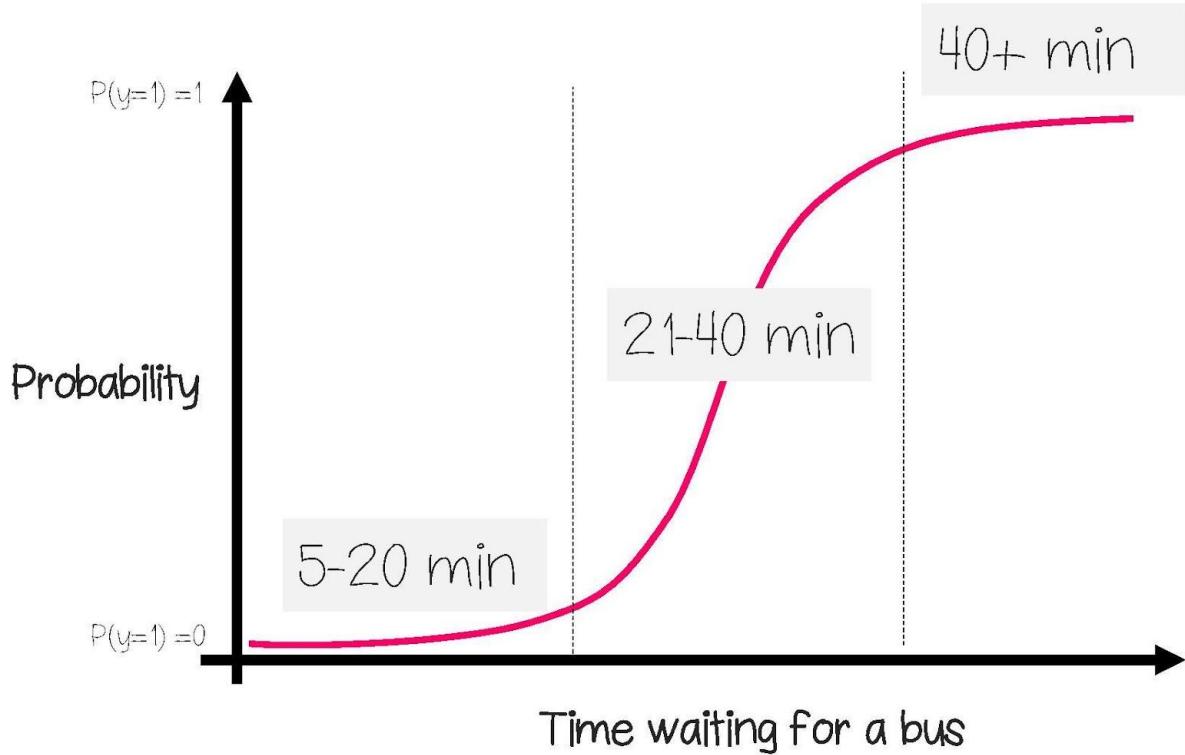
# LOGISTIC REGRESSION

# Algorithm explained:

Classify data points by calculating the **probability** that Y=1



# Should I get taxi or a bus?



# The model is very similar to linear regression

(assuming linear relationship between x and y)

$z = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + b_4x_4 \rightarrow$  then  $\frac{1}{1+e^{-z}}$  to get probability scores

ID	X1 (minute)	X2 (crowd)	X3 (time)	X4 (weather)	P(Y=1) (get Taxi)	Prediction
1	50	Many	Friday	Good	.95	Yes
2	20	Many	Saturday	OK	.32	No
3	20	Few	Monday	OK	.31	No
4	25	Few	Friday	Bad	.85	Yes
5	30	Many	Friday	Good	.68	Yes
6	35	Few	Saturday	OK	.25	No
7	12	Many	Monday	Good	.52	Yes
8	5	Few	Saturday	OK	.12	No
9	10	Many	Friday	Bad	.20	No
10	65	Few	Friday	Good	.98	Yes



We set threshold = 0.5  
If  $p(y=1) \geq 0.5$ , then we predicted YES get taxi!



## Advantages:

- One of the most powerful algorithms to date
- Easy to train
- Can be developed to create complex decision boundaries (e.g.  $x^2$   $y^2$ )

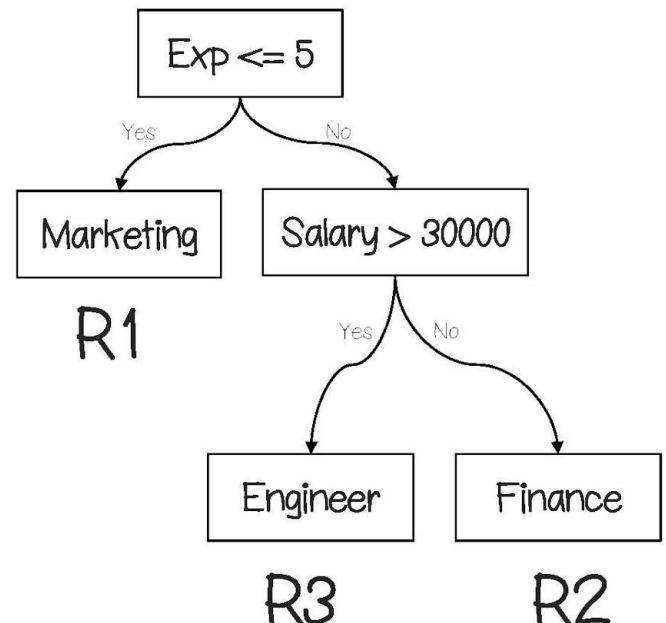
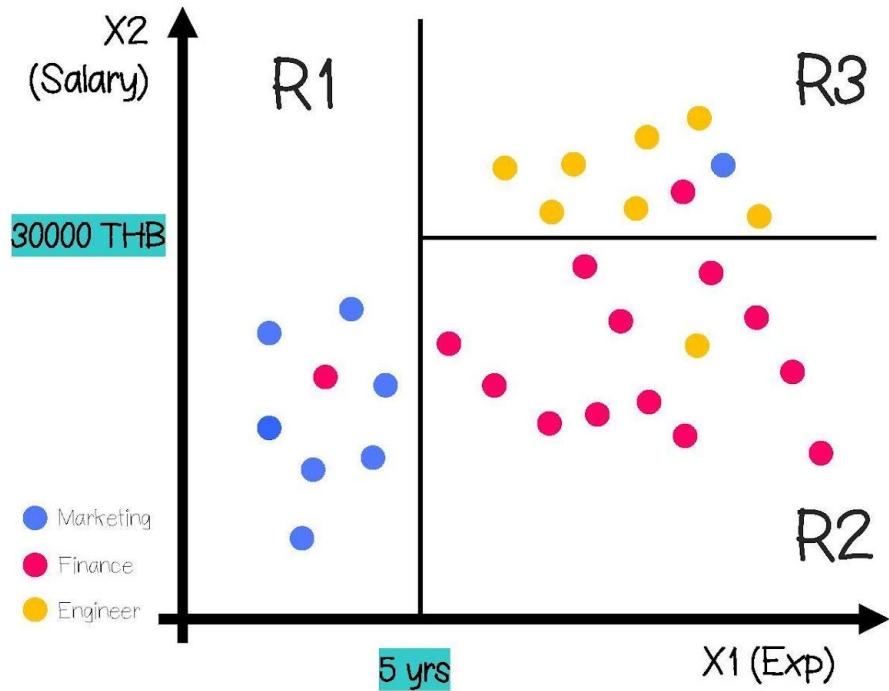
## Disadvantages:

- Also a lot of assumption e.g. no outlier, no noise
- Logistic regression is also a linear algorithm
- Prone to multicollinearity (check correlation first)  
(sigmoid function transform an output into non-linear)

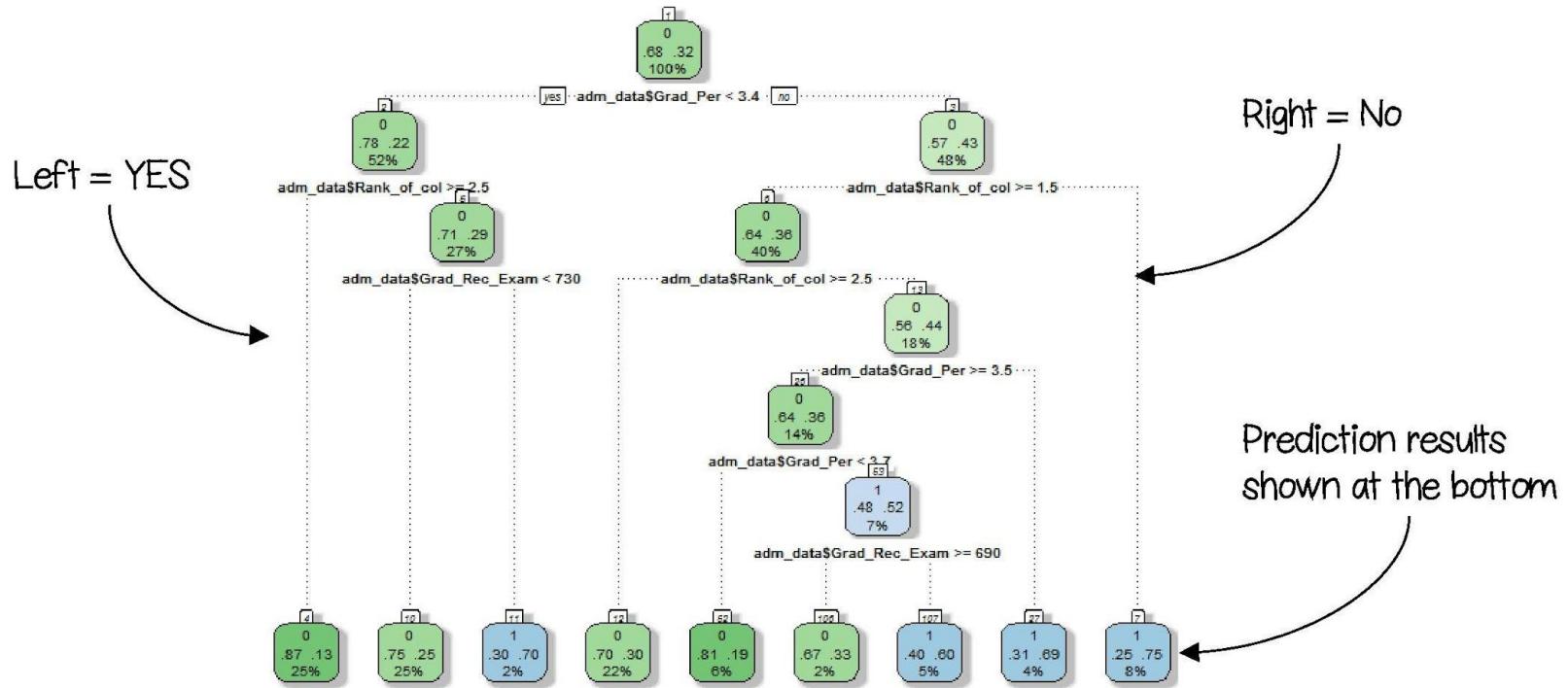
# DECISION TREE

# Algorithm explained:

Divide space into regions that best predict data points (classes)



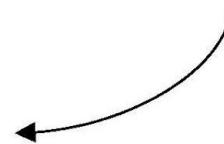
# Decision Tree is easily presented using diagram



## Advantages:

- *Easy to train*
- *Easy to interpret results*
- *Widely used in many disciplines*

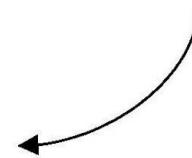
The main reason why people use decision tree



## Disadvantages:

- *Prone to overfitting: Greedy algorithms*

Locally optimal decision are made at each node, but can't guarantee globally optimal decision trees



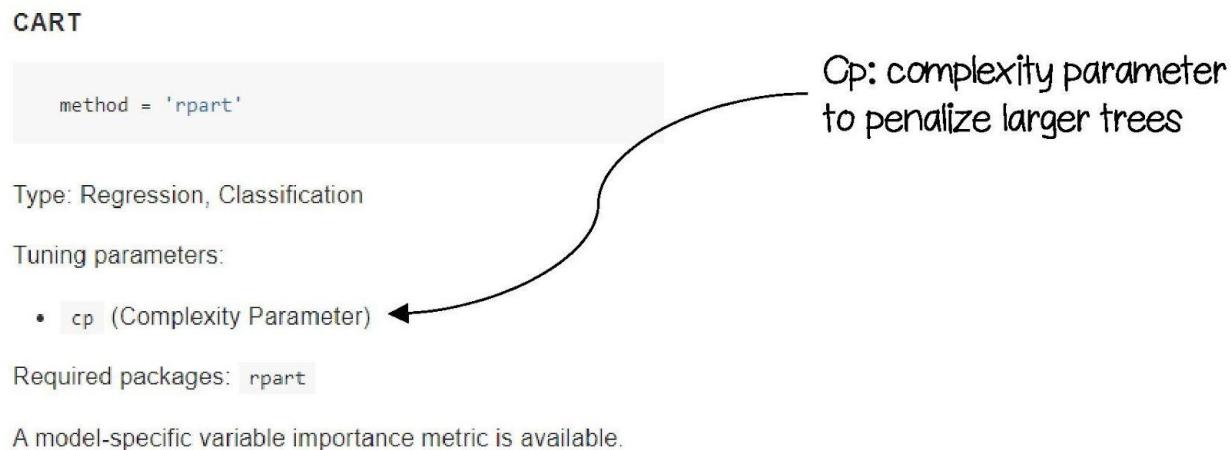
[https://en.wikipedia.org/wiki/Decision\\_tree\\_learning](https://en.wikipedia.org/wiki/Decision_tree_learning)

# Required Packages:

`library(caret)` – train & test models

`library(caTools)` – split data into training & testing sets

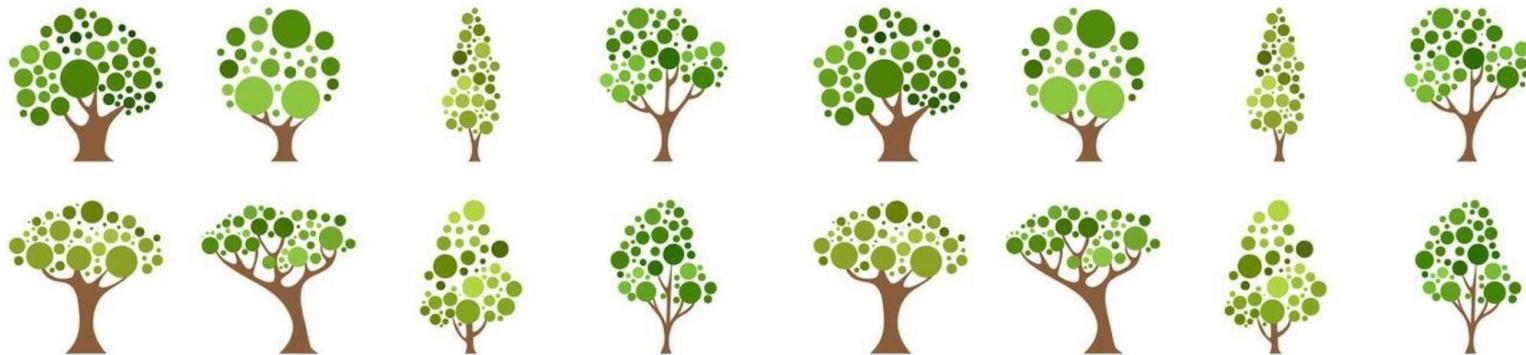
`library(rpart)` – train decision tree (CART): recursive partitioning and regression tree



# RANDOM FOREST

# Algorithm explained:

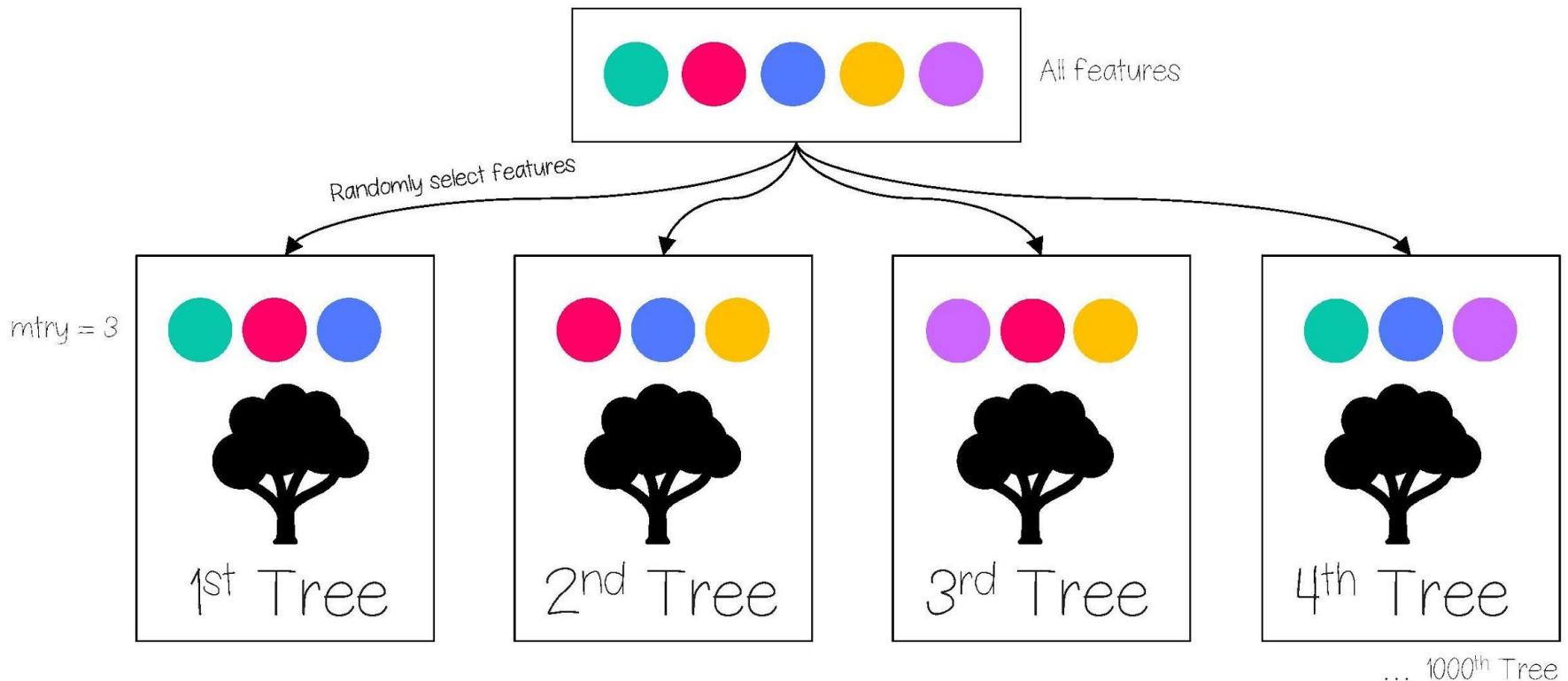
Grow thousand trees and average their prediction results



Grow **uncorrelated** 1000 trees

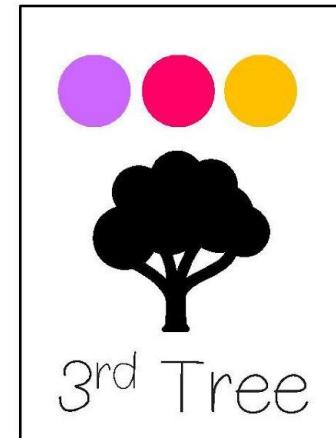
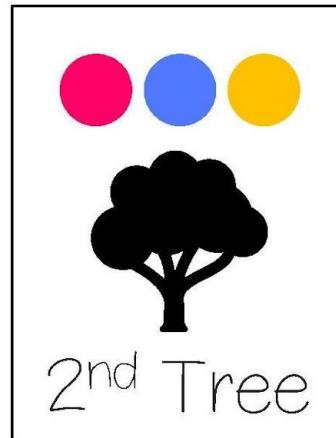
# How do we grow 1000 trees?

Randomness is your friend, hyperparameter is **mtry**



# Wisdom of the crowds

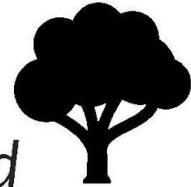
Find the average prediction results of all grown trees



... 1000<sup>th</sup> Tree

## **Advantages:**

- *Amongst the best learners in the world*
- *Easy to train (Improvement from Decision Tree)*



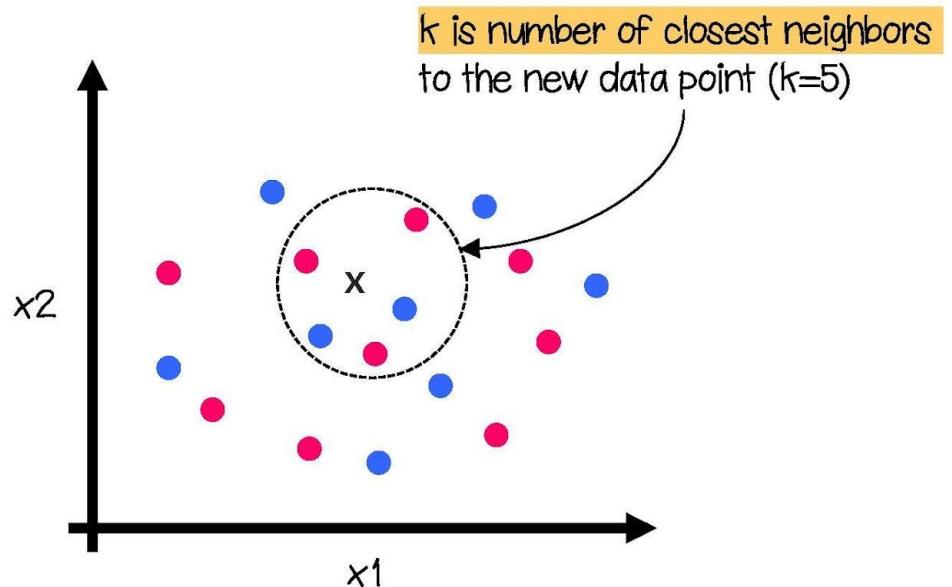
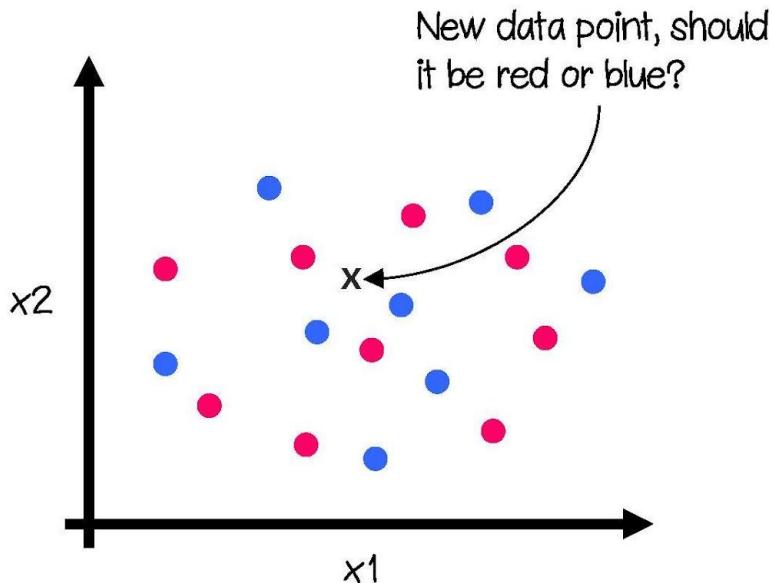
## **Disadvantages:**

- *Can be slow to train for large dataset*
- *Black Box Algorithm (not easy to explain)*

# KNN (K-NEAREST NEIGHBOR)

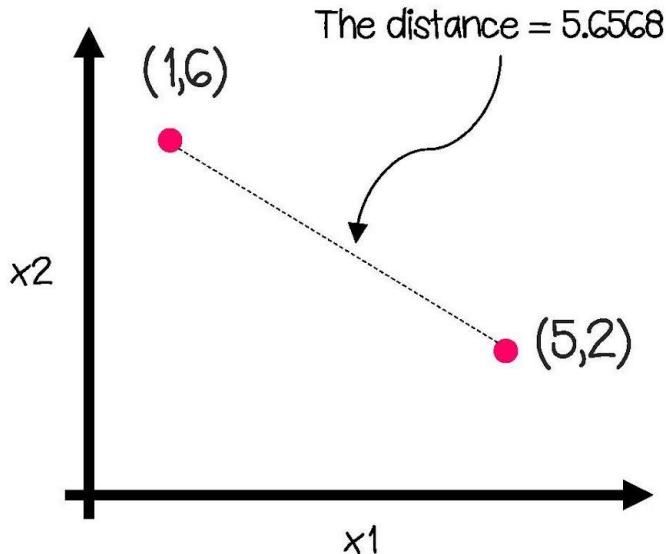
# Algorithm explained:

Predict new data according to the distance from its neighbors



# Euclidean Distance:

The most popular distance calculation method



$$d = \sqrt{\sum (x_i - y_i)^2}$$

$$d = \sqrt{(1 - 5)^2 + (6 - 2)^2}$$

$$d = \sqrt{16 + 16}$$

$$\boxed{d = 5.656854}$$

Very simple  
to compute

# Another Example

Consider the nutrient dataset provided with the `flexclust` package. The dataset contains measurements on the nutrients of 27 types of meat, fish, and fowl. The first few observations are given by

```
> data(nutrient, package="flexclust")
> head(nutrient, 4)
```

	energy	protein	fat	calcium	iron
BEEF BRAISED	340	20	28	9	2.6
HAMBURGER	245	21	17	9	2.7
BEEF ROAST	420	15	39	7	2.0
BEEF STEAK	375	19	32	9	2.6

We can calculate many dimensions (more than two)

and the Euclidean distance between the first two (beef braised and hamburger) is

$$d = \sqrt{(340 - 245)^2 + (20 - 21)^2 + (28 - 17)^2 + (9 - 9)^2 + (2.6 - 2.7)^2} = 95.64$$

## **Advantages:**

- *Very simple*
- *Easy to train (remember all data)*

## **Disadvantages:**

- *Worst for large dataset*
- *Not good with high dimensional data (many features)*
- *Don't work well with categorical data*

# Example Code:

```
5 ## load library
6 library(caret)
7
8 ## train model
9 ## define trControl
10 my_trcontrol = trainControl(method = "cv",
11                             number = 3,
12                             verboseIter = TRUE,
13                             search = "random")
14
15 knn_model <- train(Species ~ . , data = m,
16                      method = "knn",
17                      metric = "Accuracy",
18                      tuneLength = 10,
19                      trControl = my_trcontrol)
```

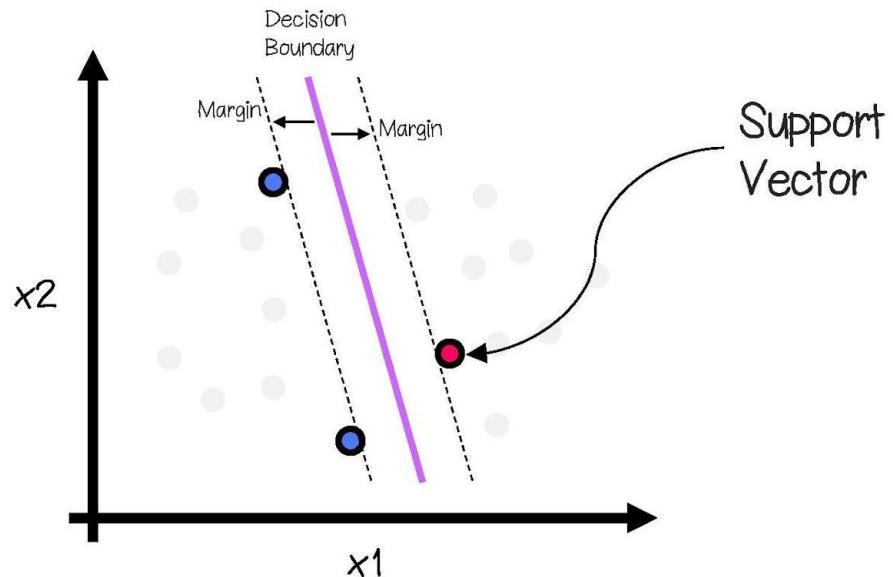
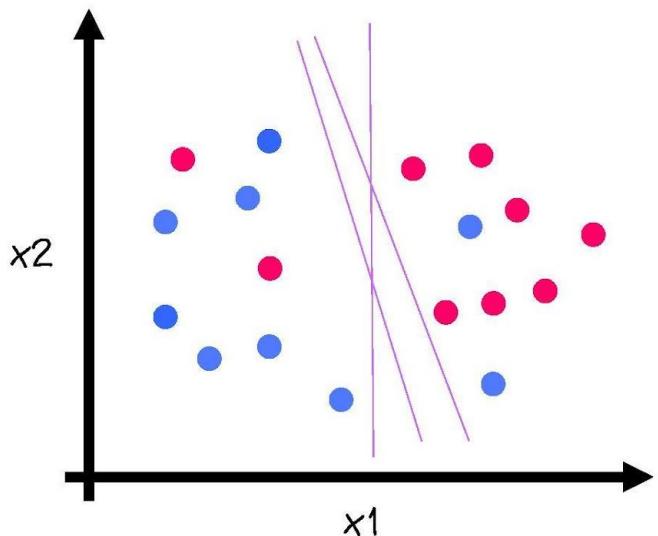
Define our technique  
(to train the model)

Train the model using  
random search +  
tuneLength = 10

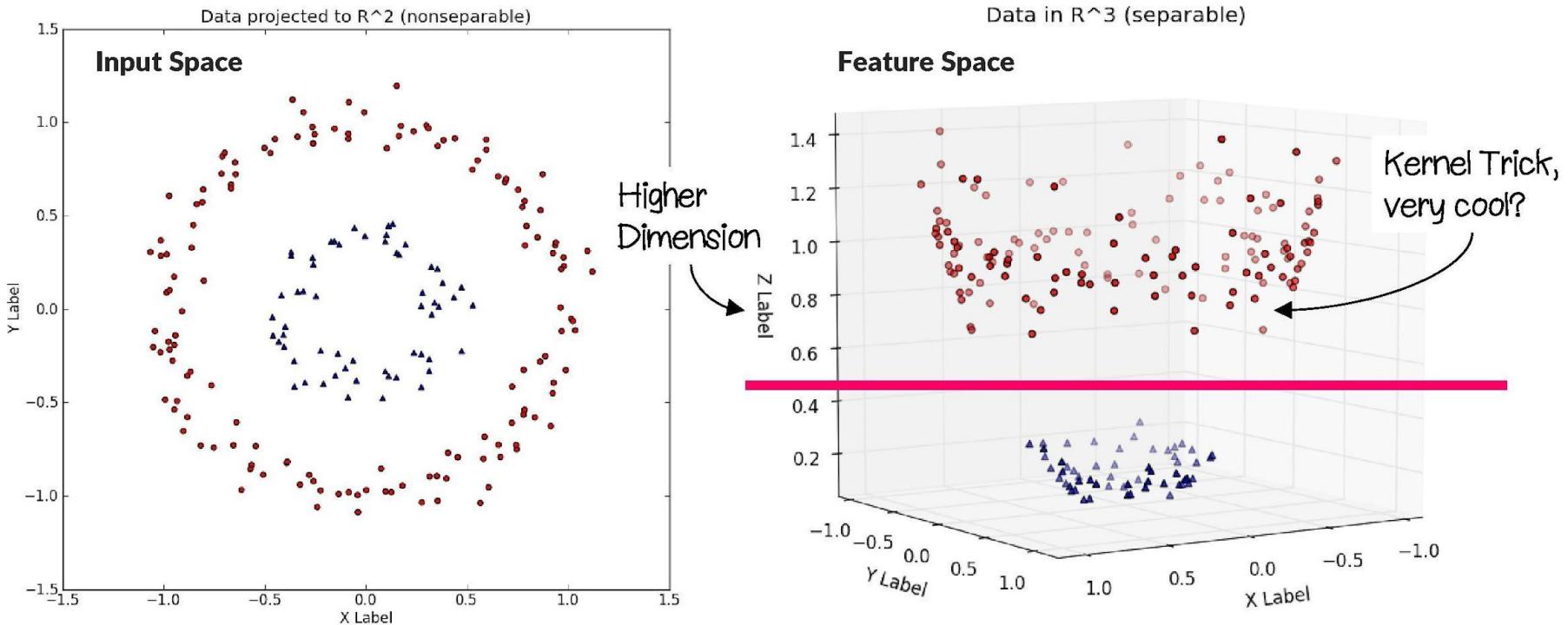
# SUPPORT VECTOR MACHINE

# Algorithm explained:

We want to maximize the **margin** from decision boundary

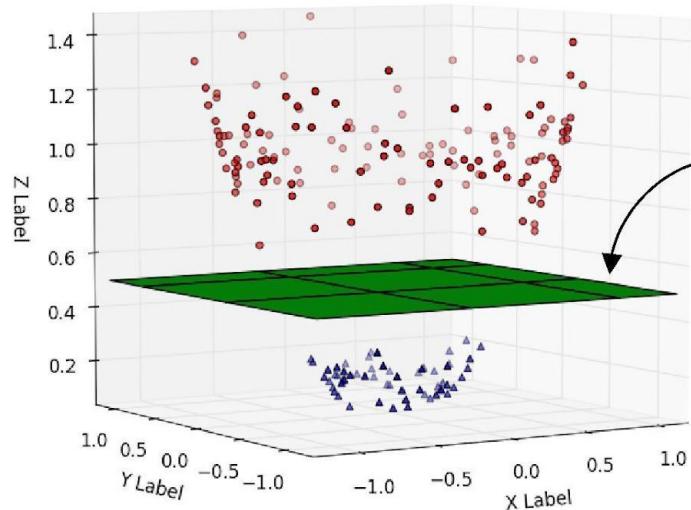


# But can you divide this with a single straight line?

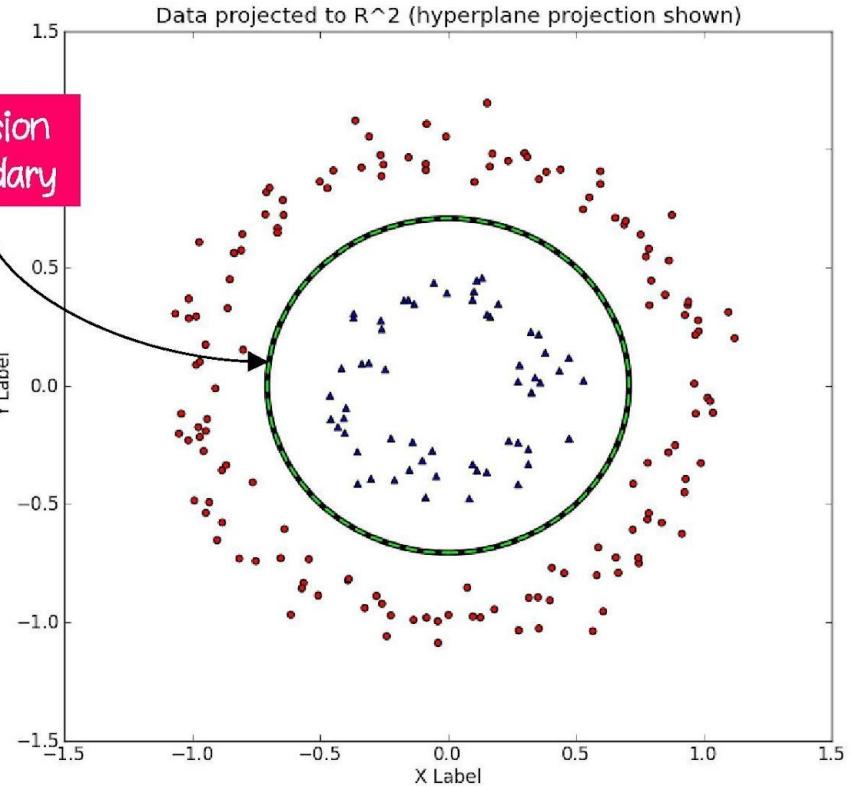


# Impressive Result!!

Data in  $\mathbb{R}^3$  (separable w/ hyperplane)

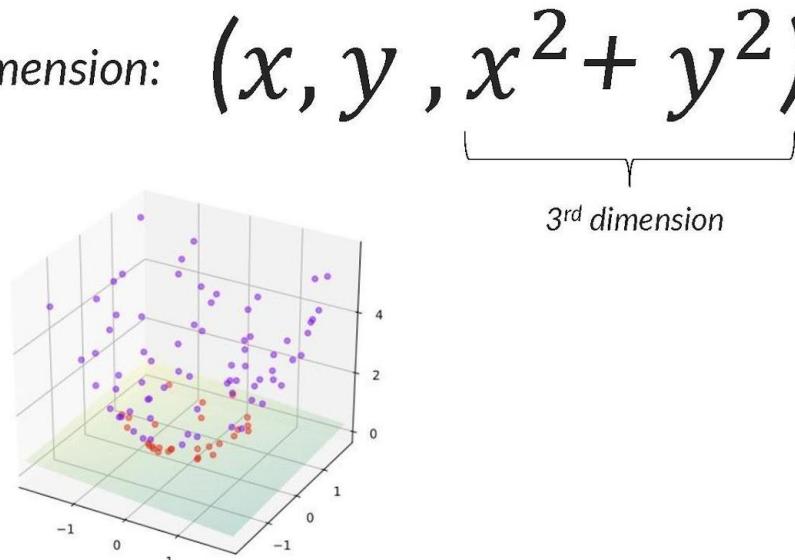
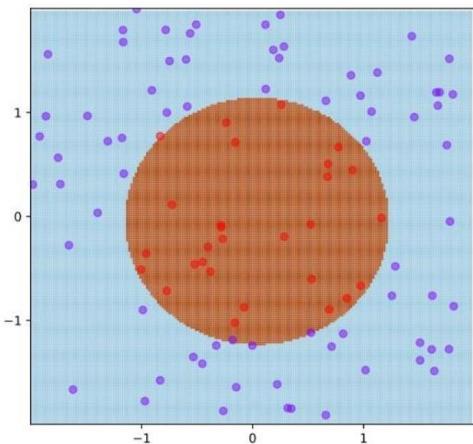


Data projected to  $\mathbb{R}^2$  (hyperplane projection shown)



# Let's look closely at kernel trick

- we have  $(x,y)$  coordinates
- we project a new higher dimension:  $(x, y, x^2 + y^2)$



[https://en.wikipedia.org/wiki/Kernel\\_method](https://en.wikipedia.org/wiki/Kernel_method)

## **Advantages:**

- *Work with Regression | Classification problems*
- *Superior Performance*
- *Having lots of available kernel tricks*

<https://topepo.github.io/caret/train-models-by-tag.html#support-vector-machines>

## **Disadvantages:**

- *Expensive*
- *It's lazy algorithm (remember data like KNN)*

# Required Packages:

`library(caret)` – train & test models

`library(caTools)` – split data into training & testing sets

`library(e1071)` – load some algorithms e.g. SVM

## Support Vector Machines with Linear Kernel

```
method = 'svmLinear2'
```

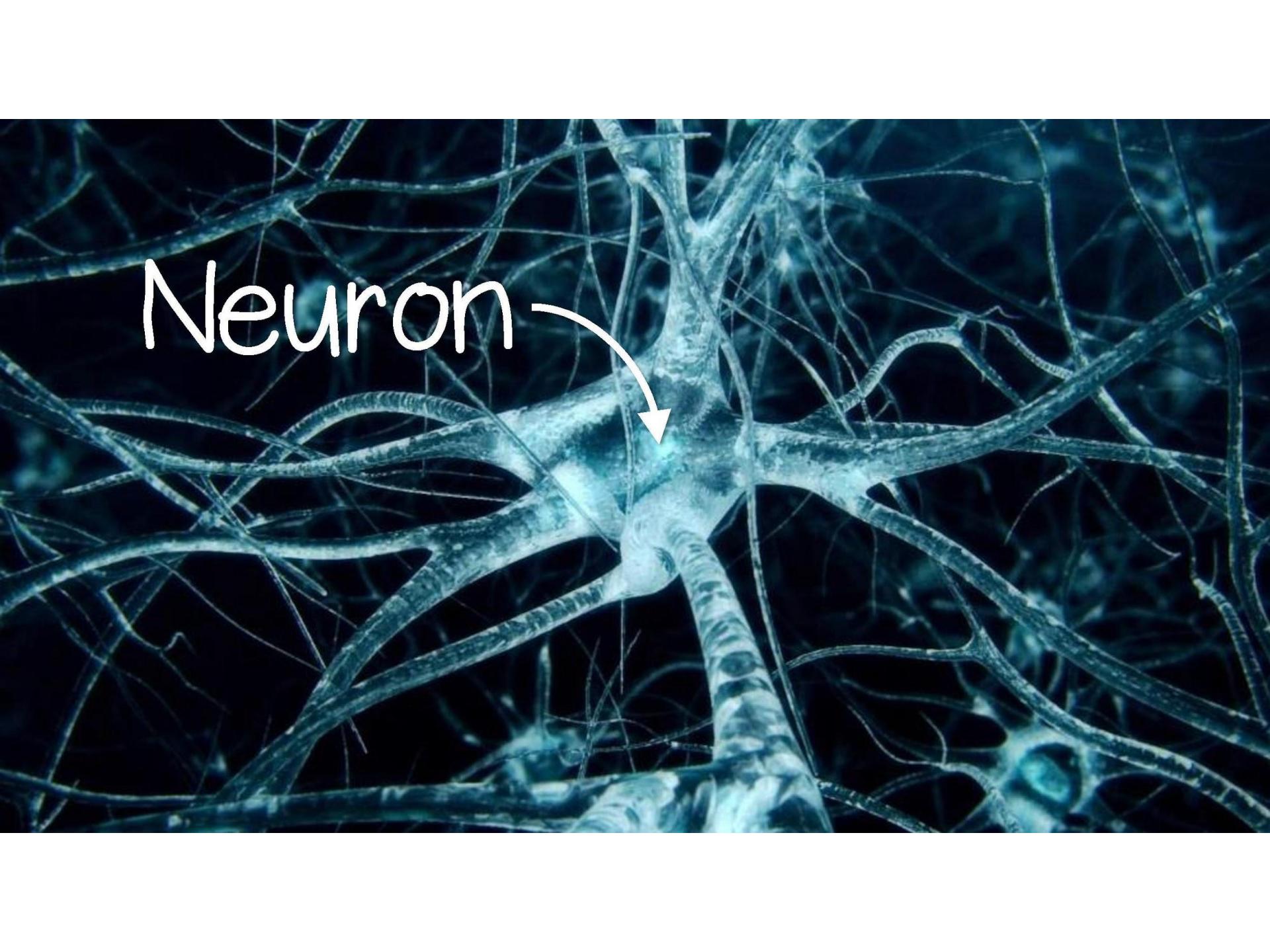
Type: Regression, Classification

Tuning parameters:

- `cost` (Cost)

Required packages: `e1071`

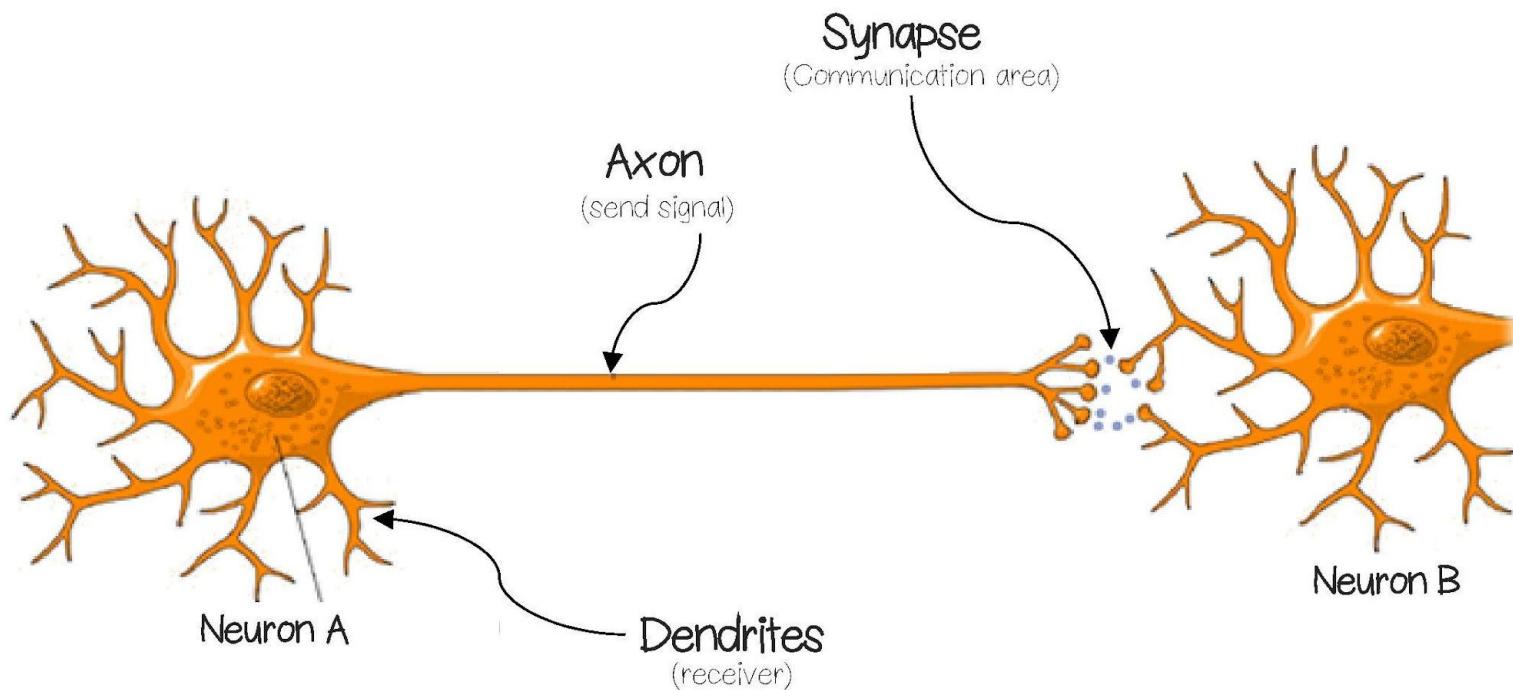
# NEURAL NETWORK

A dense network of blue glowing neurons against a black background.

Neuron



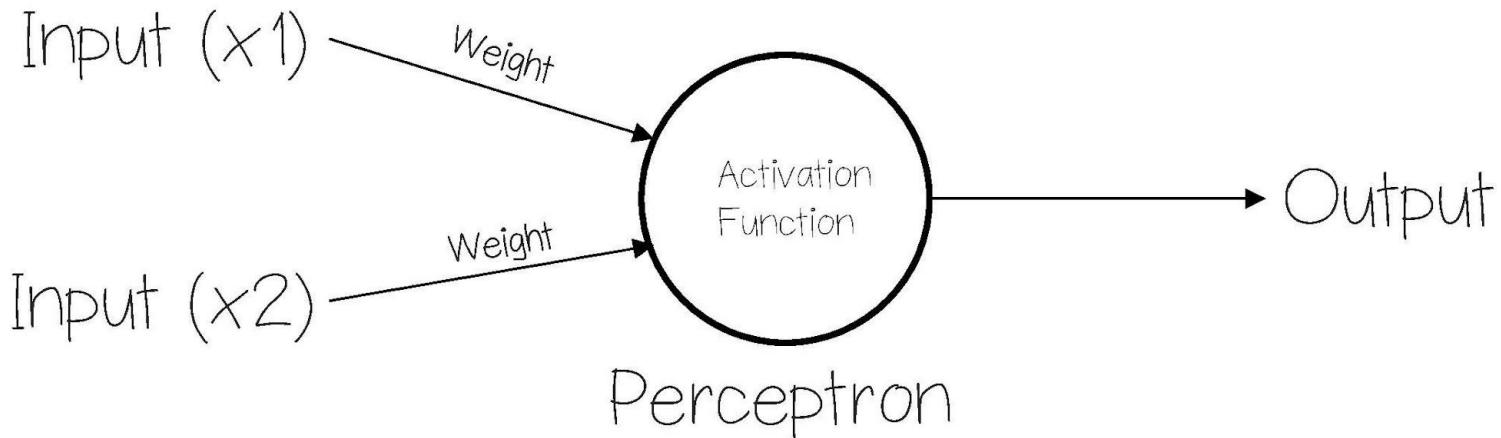
# How your brain cells communicate





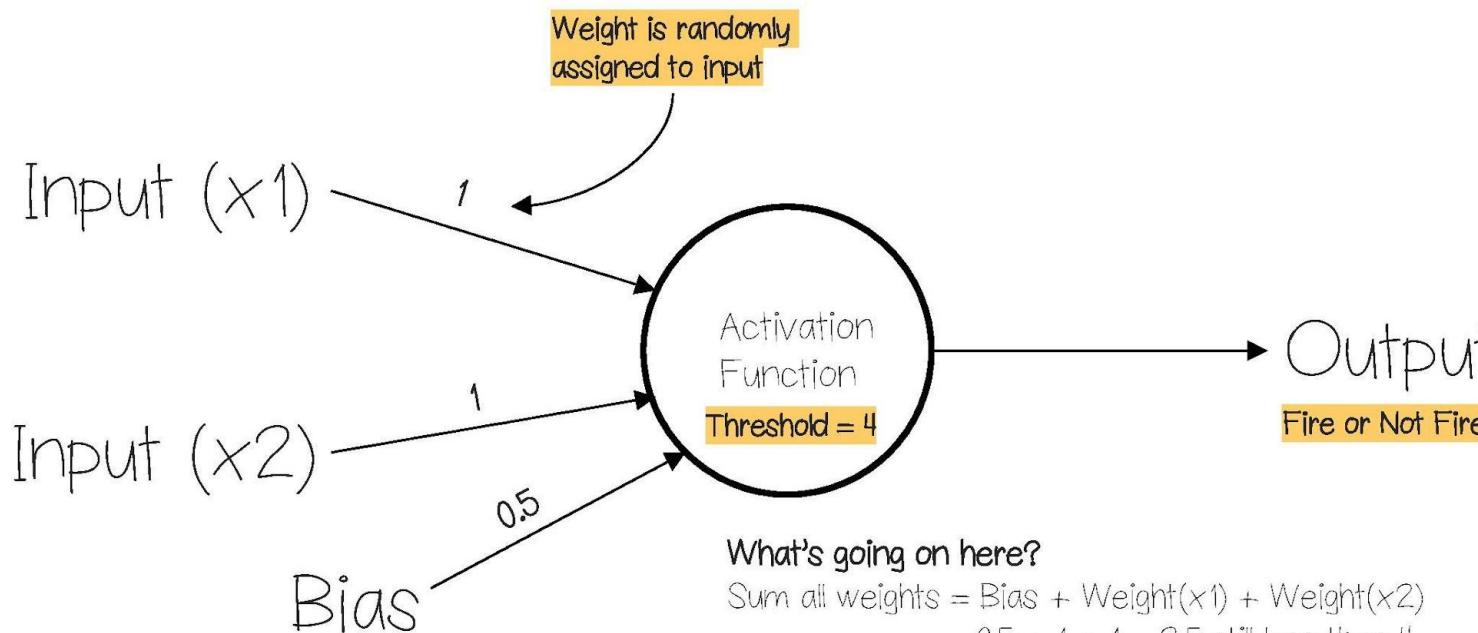
# Algorithm Explained:

Neurons fire only when it's the signal strong enough to activate them



# Algorithm Explained:

Neurons fire only when it's the signal strong enough to activate them

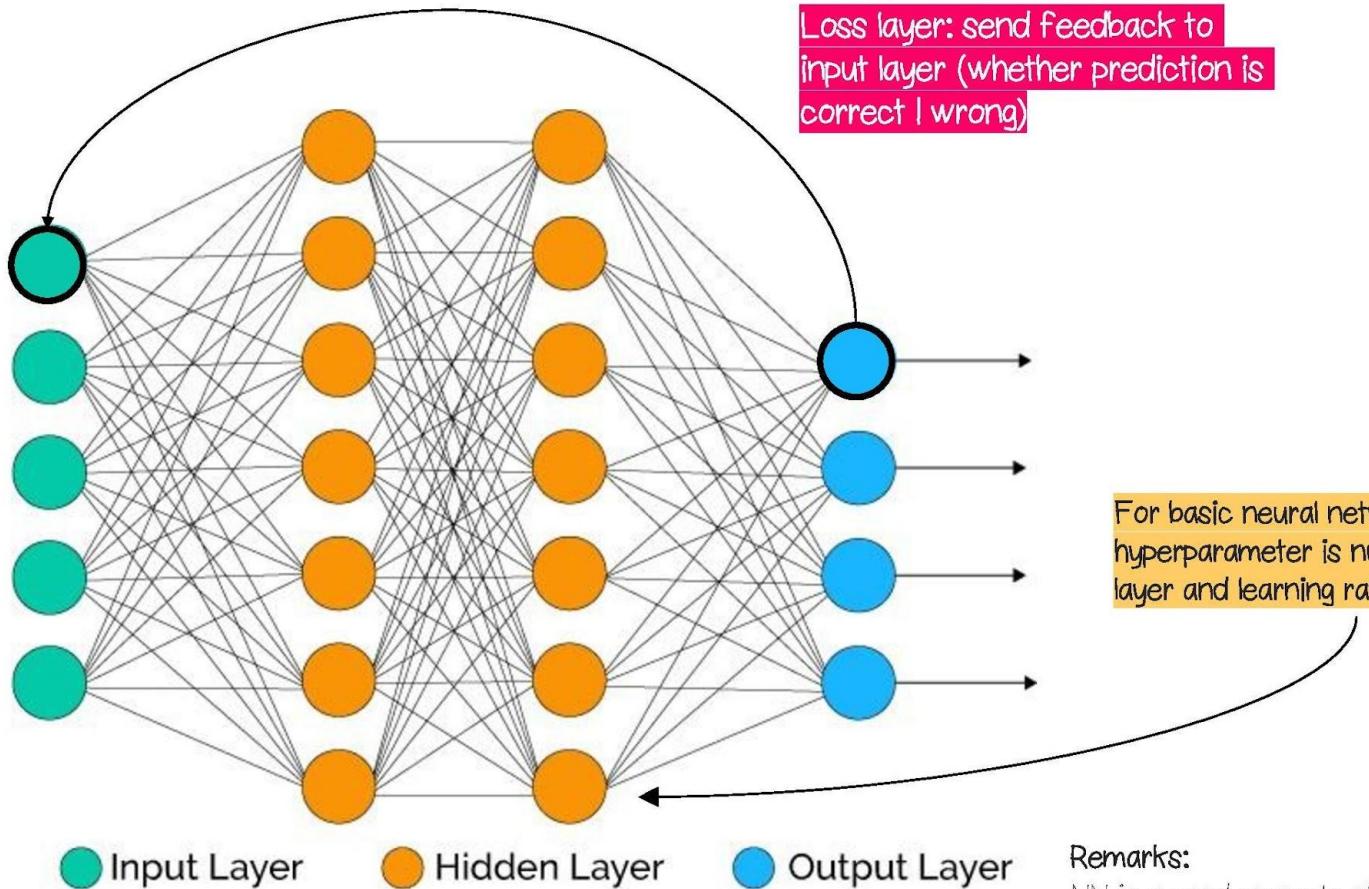


What's going on here?

$$\begin{aligned} \text{Sum all weights} &= \text{Bias} + \text{Weight}(x_1) + \text{Weight}(x_2) \\ &= 0.5 + 1 + 1 = 2.5 \text{ still less than } 4 \end{aligned}$$

: neuron not activate

: perceptron learn from this experience and adjust the weight

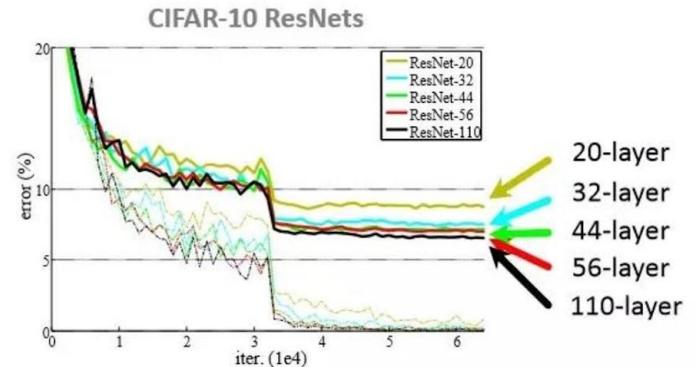
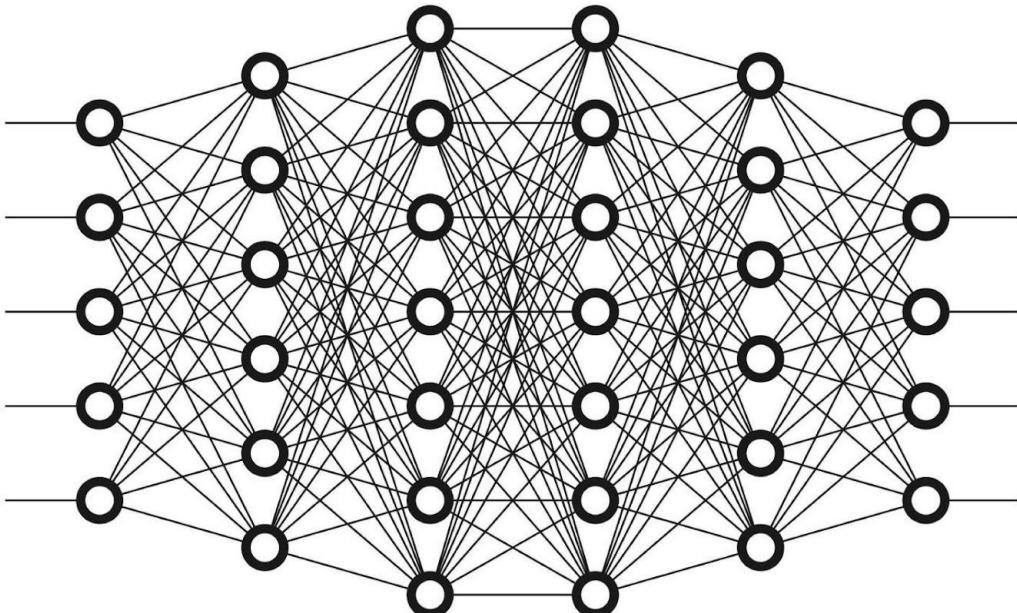


Remarks:

NN is a good example of supervised learning that also a reinforcement learning (learn from experience)

# Deep Learning

Expands the number of hidden layers



Hieu Pham, Has done some machine learning  
Updated Apr 29, 2017

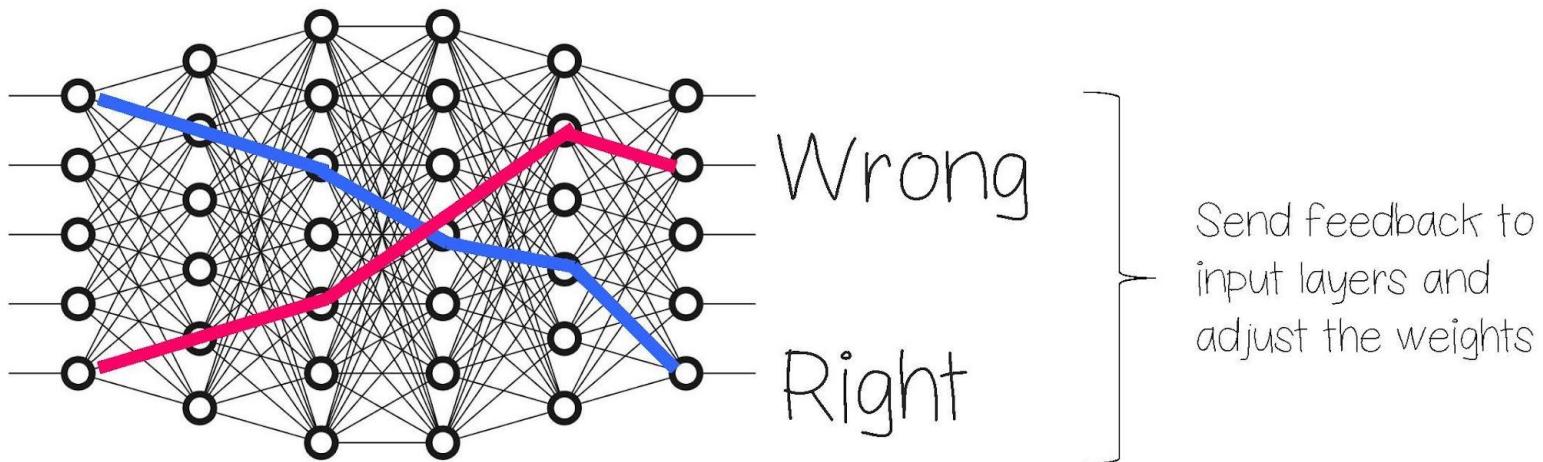


In this paper [\(Huang et al, 2016\)](#), the authors report that they “can increase the depth of residual networks **even beyond 1200 layers** and still yield meaningful improvements in test error (4.91% on CIFAR-10)”.

I have never heard of anything deeper than that.

# Backpropagation

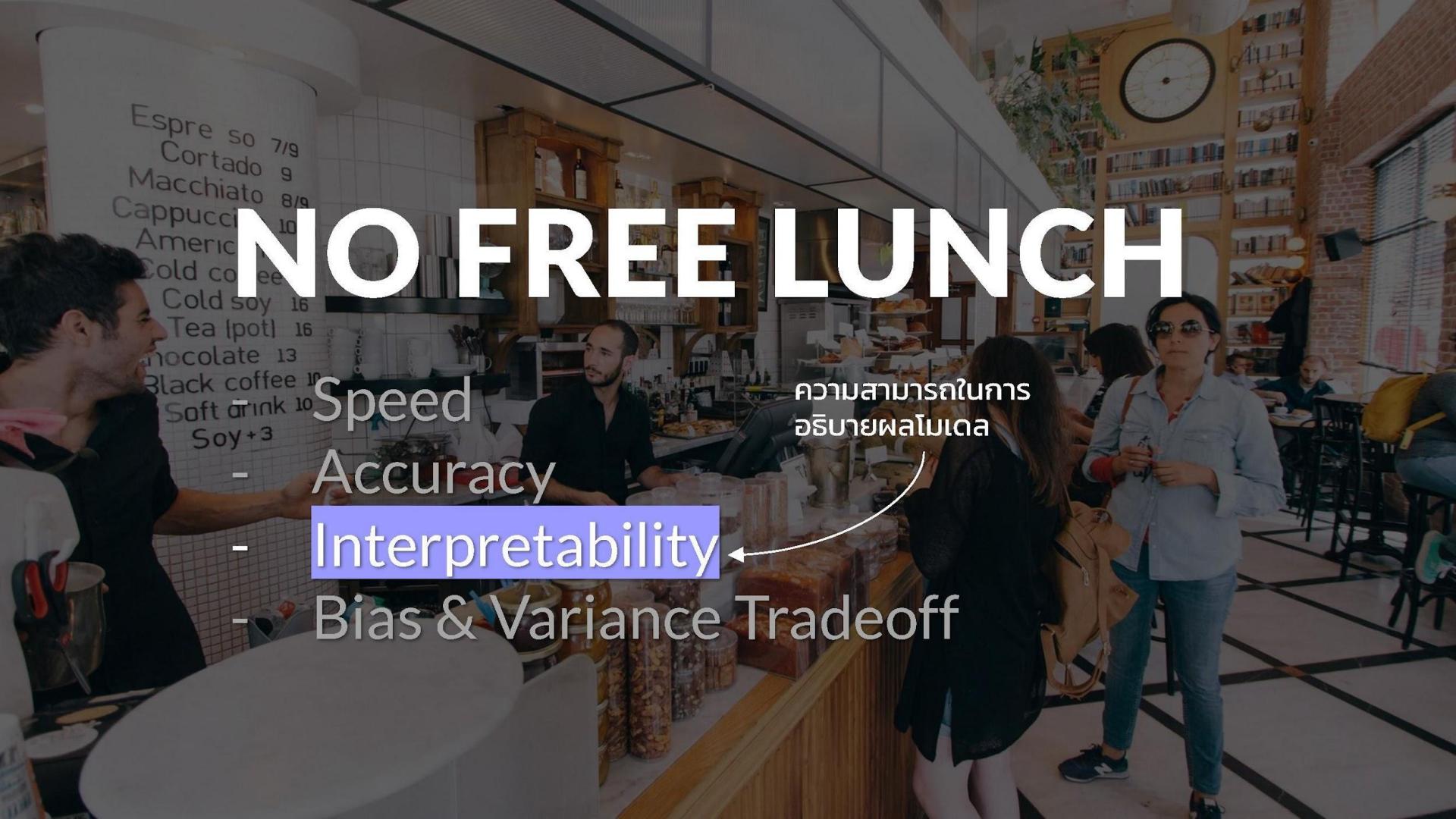
The backbone of artificial neural networks



# NO FREE LUNCH

- Speed
- Accuracy
- Interpretability
- Bias & Variance Tradeoff

ความสามารถในการ  
อธิบายผลโดยเดล



## Advantages:

- Amongst the very best performers | learners in the world
- Can be developed into deep learning algorithms
- Work with both regression and classification problems

## Disadvantages:

- Requires large data set
- Black Box
- Prone to overfitting problem
- Need more power if more hidden layers

Time to buy a  
graphic card :D



# Required Packages:

`library(caret)` – train & test models

`library(caTools)` – split data into training & testing sets

`library(nnet)` – used to train neural network in caret

## Neural Network

```
method = 'nnet'
```

Type: Classification, Regression

Tuning parameters:

- `size` (#Hidden Units)
- `decay` (Weight Decay)

Required packages: `nnet`

A model-specific variable importance metric is available.

Machine Learning 3(2-2-5), 2/2565

# Unit 4

## Unsupervised Learning

อ.อานันท์ ไม้ประดิษฐ์

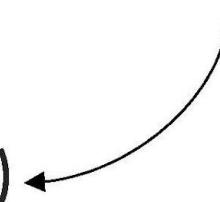
# UNSUPERVISED: K-MEANS

# **Recall the main task of unsupervised algorithm:**

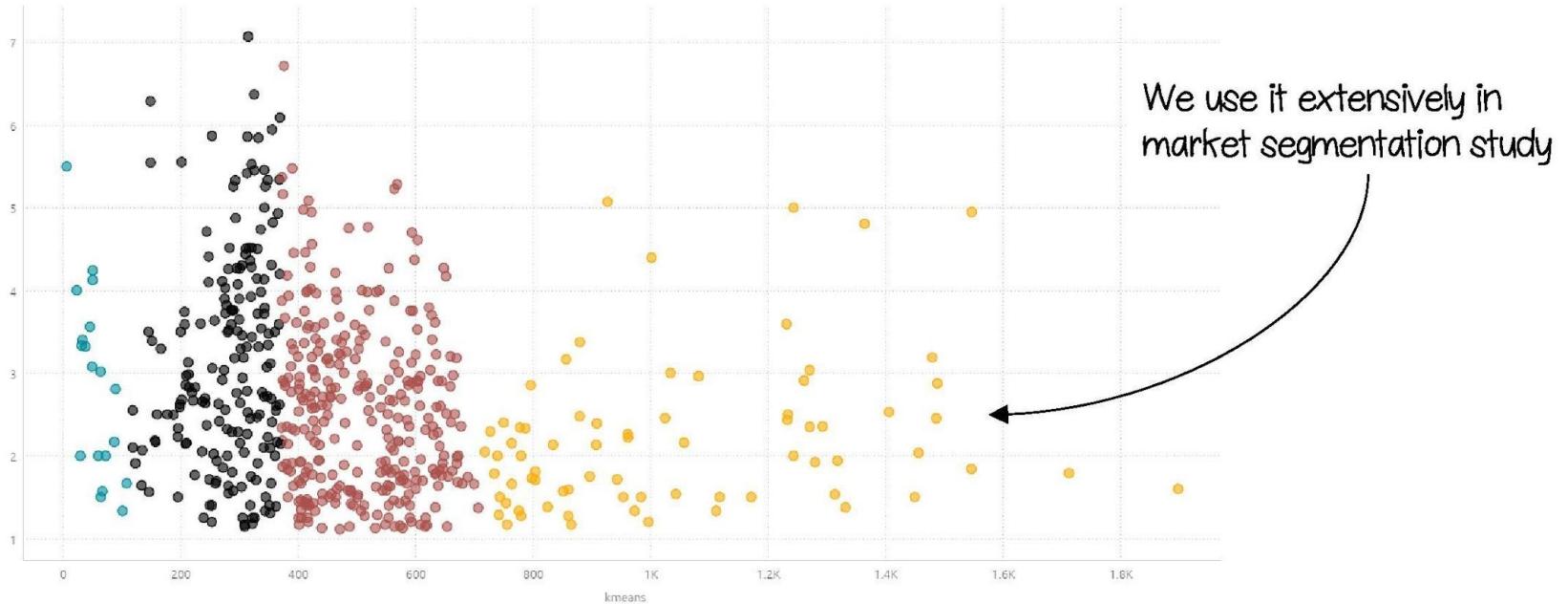
*To find patterns in our data set*

*(Descriptive – no right | wrong answer)*

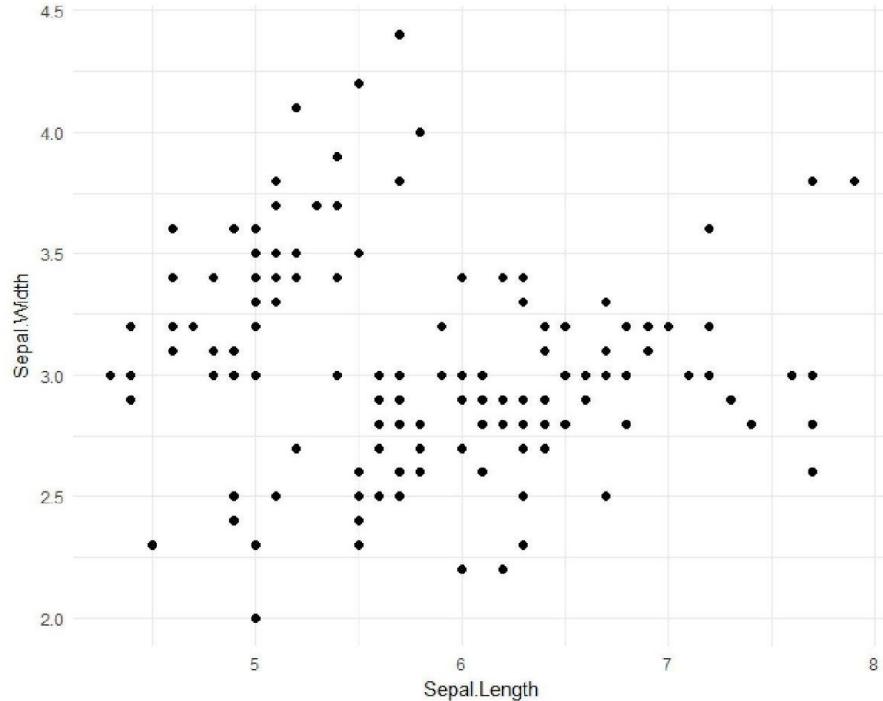
It's very subjective. But there are some advanced techniques to guide our decisions (not discussed here).



# K-Means is very popular unsupervised learning algorithms



# K-Means is **very efficient**, easily done in one line

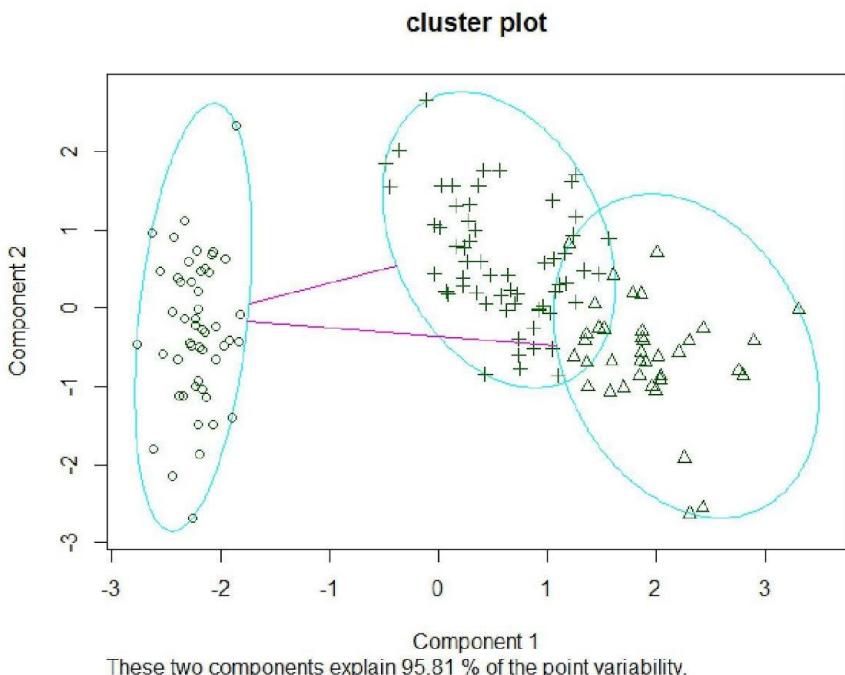


```
> library(ggplot2)
> ggplot(iris, aes(x=Sepal.Length, y=Sepal.Width)) +
  geom_point(size = 2) + theme_minimal()

> set.seed(2)
> km_fit <- kmeans(iris[1:4], centers = 3)
```

Specify (guess)  
number of clusters

# In reality, we don't know how accurate our model is.



```
> library(clusplot)
> clustplot(iris[1:4], km_fit$cluster, main = "cluster plot")
> # cross-check with actual class in iris dataframe
> table(km_fit$cluster, iris$Species)
```

K-means Cluster

Actual Class

## **Advantages:**

- **Very easy to train**
- **Often produce reasonable classification results**

## **Disadvantages:**

- **Number of clusters sometimes difficult to determine**
- **Not all features are numeric (can't find mean)**

# Reinforcement Learning: Q-Learning

# So far... Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification,  
regression, object detection,  
semantic segmentation, image  
captioning, etc.



→ Cat

Classification

# So far... Unsupervised Learning

**Data:**  $x$

Just data, no labels!

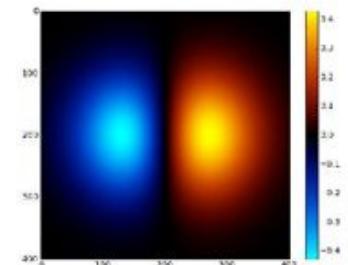
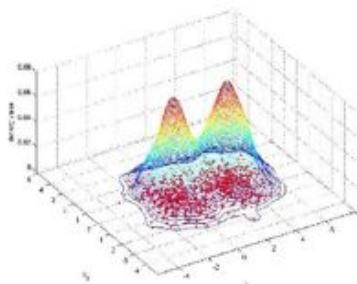
**Goal:** Learn some underlying hidden *structure* of the data

**Examples:** Clustering, dimensionality reduction, feature learning, density estimation, etc.



Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation

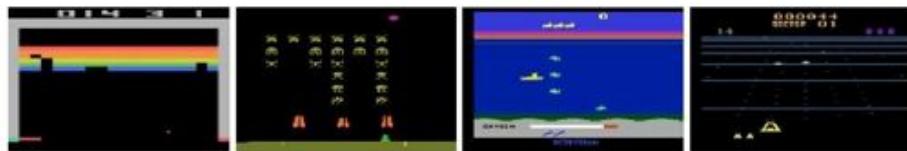
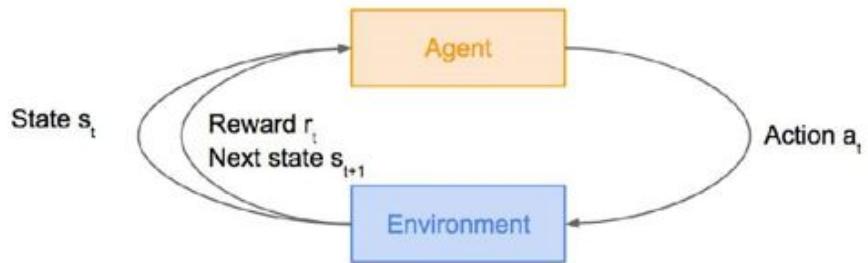


2-d density estimation

# Today: Reinforcement Learning

Problems involving an **agent** interacting with an **environment**, which provides numeric **reward** signals

**Goal:** Learn how to take actions in order to maximize reward



# Overview

- What is Reinforcement Learning?
- Markov Decision Processes
- Q-Learning
- Policy Gradients

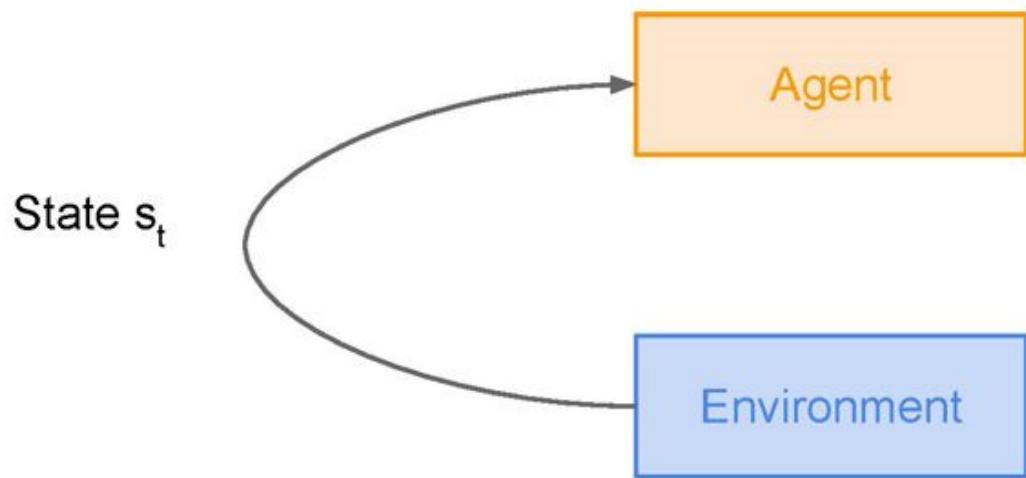
# Reinforcement Learning



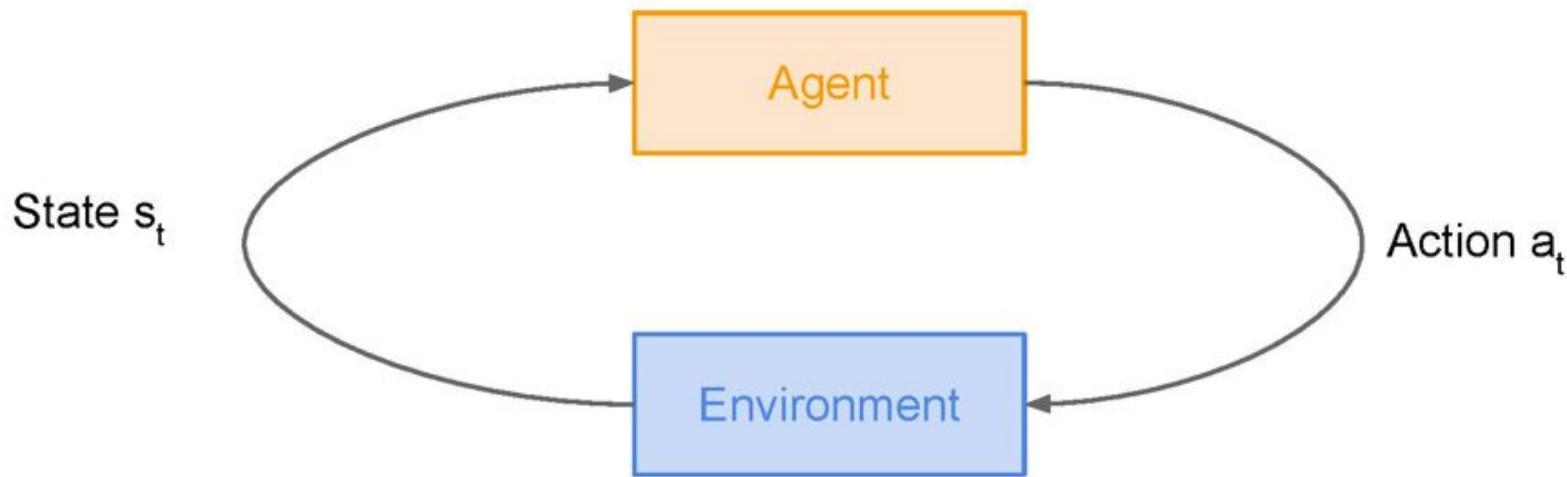
Agent

Environment

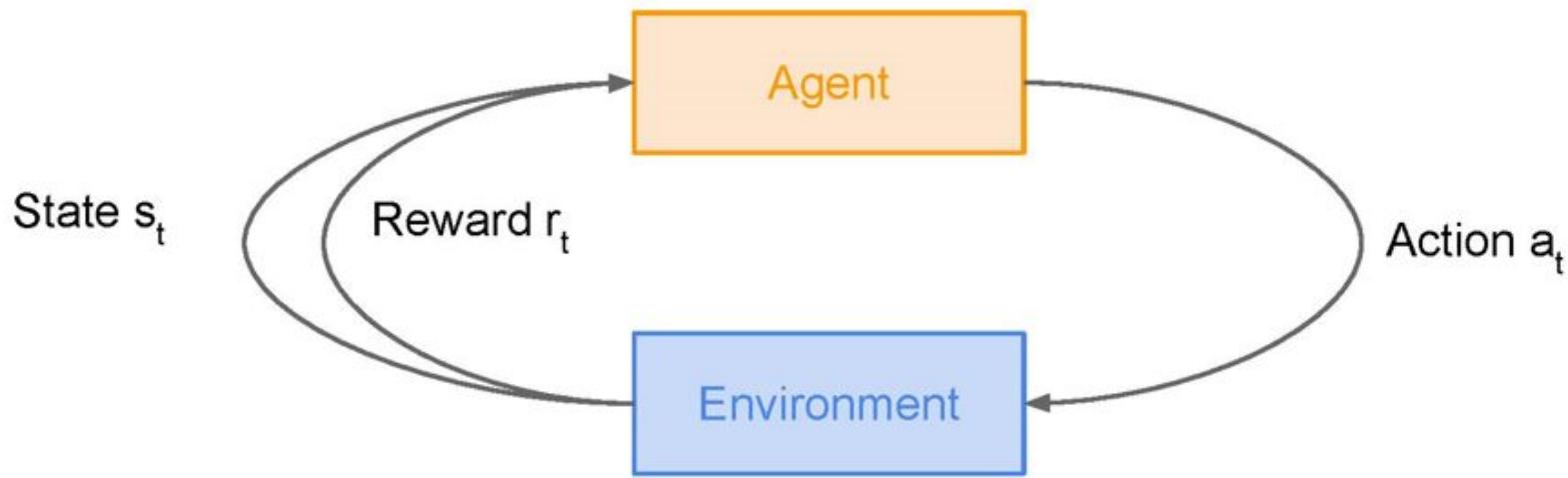
# Reinforcement Learning



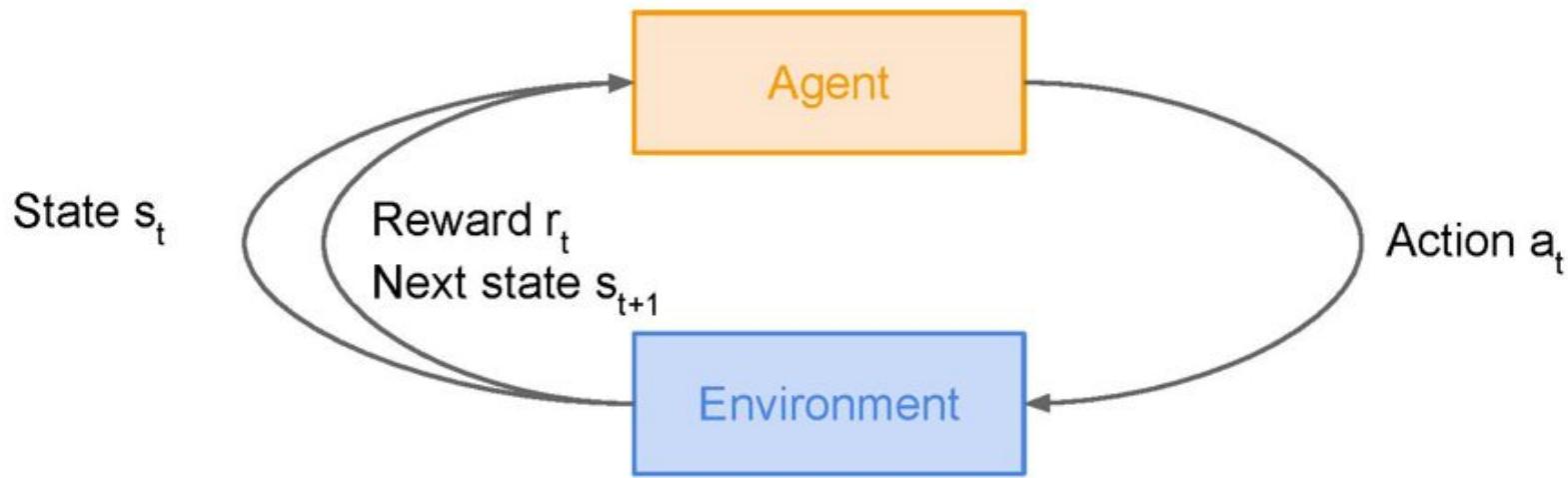
# Reinforcement Learning



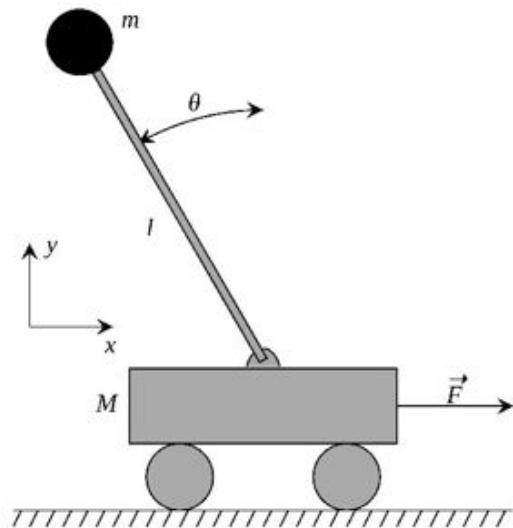
# Reinforcement Learning



# Reinforcement Learning



# Cart-Pole Problem



**Objective:** Balance a pole on top of a movable cart

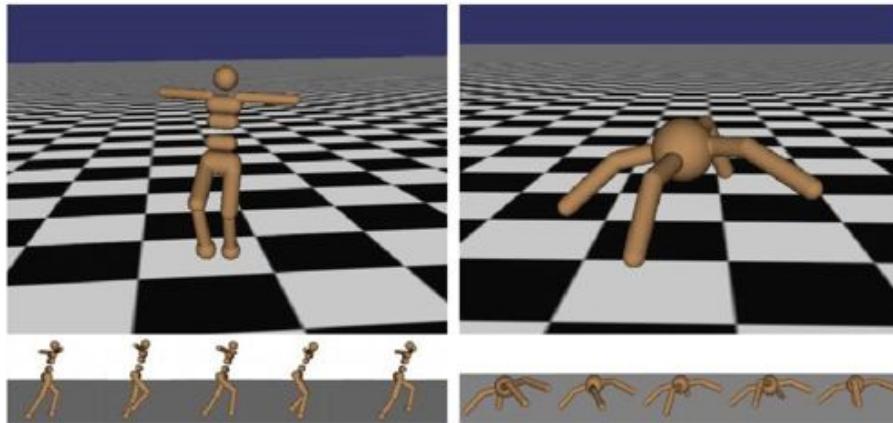
**State:** angle, angular speed, position, horizontal velocity

**Action:** horizontal force applied on the cart

**Reward:** 1 at each time step if the pole is upright

[Cart-Pole Video](#)

# Robot Locomotion



**Objective:** Make the robot move forward

**State:** Angle and position of the joints

**Action:** Torques applied on joints

**Reward:** 1 at each time step upright + forward movement

Learning to Walk via Deep Reinforcement Learning

# Atari Games



**Objective:** Complete the game with the highest score

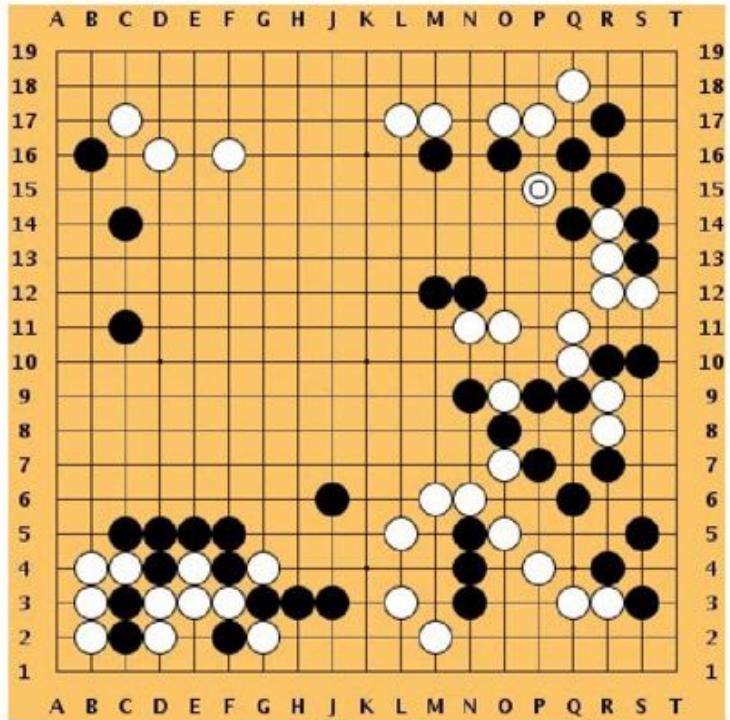
**State:** Raw pixel inputs of the game state

**Action:** Game controls e.g. Left, Right, Up, Down

**Reward:** Score increase/decrease at each time step

Google DeepMind's Deep Q-learning playing Atari Breakout!

# Go



**Objective:** Win the game!

**State:** Position of all pieces

**Action:** Where to put the next piece down

**Reward:** 1 if win at the end of the game, 0 otherwise



# DOTA 2

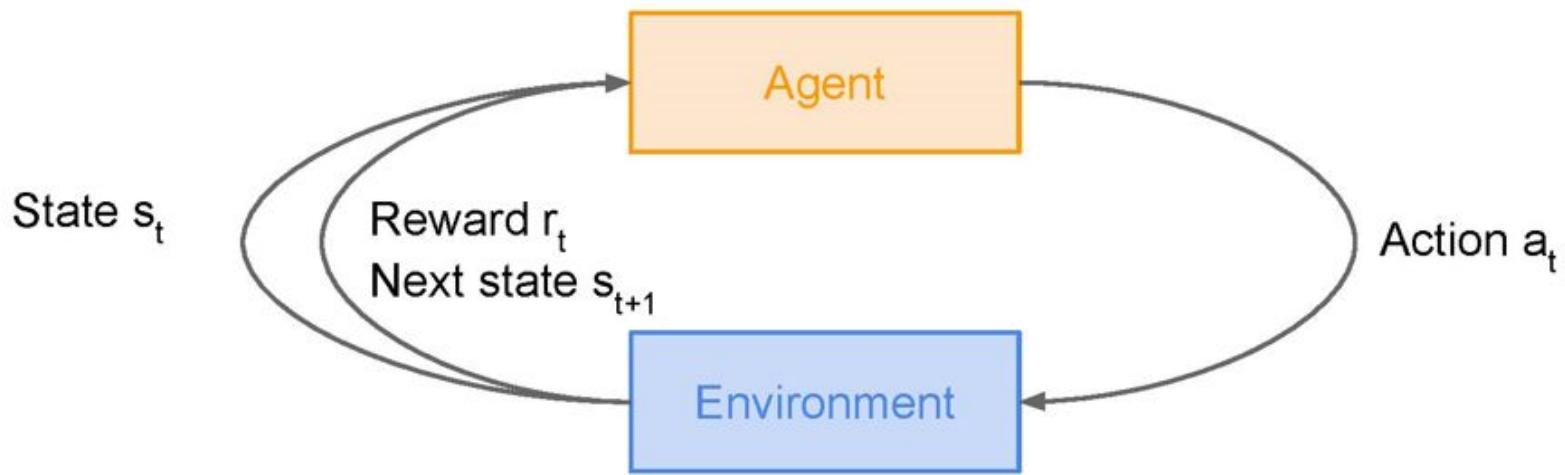
VALVE

[OpenAI Five: Dota Gameplay](#)



I tried to make a Valorant AI using computer vision

## How can we mathematically formalize the RL problem?



# Markov Decision Process

- Mathematical formulation of the RL problem
- **Markov property:** Current state completely characterises the state of the world

Defined by:  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

$\mathcal{S}$  : set of possible states

$\mathcal{A}$  : set of possible actions

$\mathcal{R}$  : distribution of reward given (state, action) pair

$\mathbb{P}$  : transition probability i.e. distribution over next state given (state, action) pair

$\gamma$  : discount factor

# Markov Decision Process

- At time step  $t=0$ , environment samples initial state  $s_0 \sim p(s_0)$
- Then, for  $t=0$  until done:
  - Agent selects action  $a_t$
  - Environment samples reward  $r_t \sim R(\cdot | s_t, a_t)$
  - Environment samples next state  $s_{t+1} \sim P(\cdot | s_t, a_t)$
  - Agent receives reward  $r_t$  and next state  $s_{t+1}$
- A policy  $\pi$  is a function from  $S$  to  $A$  that specifies what action to take in each state
- **Objective:** find policy  $\pi^*$  that maximizes cumulative discounted reward:  $\sum_{t \geq 0} \gamma^t r_t$

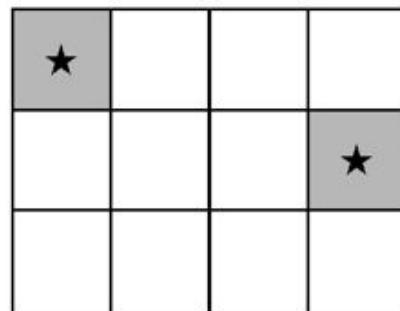
# A simple MDP: Grid World

actions = {

1. right →
2. left ←
3. up ↑
4. down ↓

}

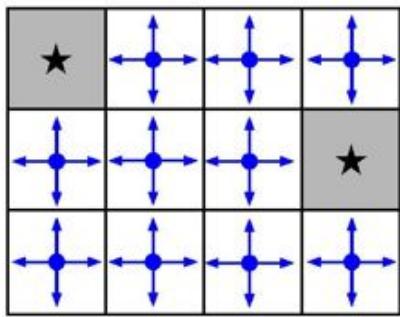
states



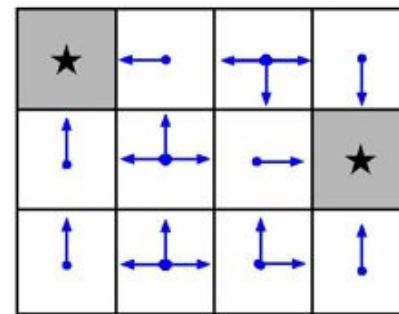
Set a negative “reward”  
for each transition  
(e.g.  $r = -1$ )

**Objective:** reach one of terminal states (greyed out) in  
least number of actions

# A simple MDP: Grid World



Random Policy



Optimal Policy

## The optimal policy $\pi^*$

We want to find optimal policy  $\pi^*$  that maximizes the sum of rewards.

How do we handle the randomness (initial state, transition probability...)?

## The optimal policy $\pi^*$

We want to find optimal policy  $\pi^*$  that maximizes the sum of rewards.

How do we handle the randomness (initial state, transition probability...)?  
Maximize the **expected sum of rewards!**

Formally:  $\pi^* = \arg \max_{\pi} \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | \pi \right]$  with  $s_0 \sim p(s_0), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$

## Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths)  $s_0, a_0, r_0, s_1, a_1, r_1, \dots$

How good is a state?

The **value function** at state  $s$ , is the expected cumulative reward from following the policy from state  $s$ :

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | s_0 = s, \pi \right]$$

## Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths)  $s_0, a_0, r_0, s_1, a_1, r_1, \dots$

### How good is a state?

The **value function** at state  $s$ , is the expected cumulative reward from following the policy from state  $s$ :

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | s_0 = s, \pi \right]$$

### How good is a state-action pair?

The **Q-value function** at state  $s$  and action  $a$ , is the expected cumulative reward from taking action  $a$  in state  $s$  and then following the policy:

$$Q^\pi(s, a) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

# Bellman equation

The optimal Q-value function  $Q^*$  is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

# Bellman equation

The optimal Q-value function  $Q^*$  is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

$Q^*$  satisfies the following **Bellman equation**:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q^*(s', a') | s, a \right]$$

**Intuition:** if the optimal state-action values for the next time-step  $Q^*(s', a')$  are known, then the optimal strategy is to take the action that maximizes the expected value of  $r + \gamma Q^*(s', a')$

# Bellman equation

The optimal Q-value function  $Q^*$  is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

$Q^*$  satisfies the following **Bellman equation**:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q^*(s', a') | s, a \right]$$

**Intuition:** if the optimal state-action values for the next time-step  $Q^*(s', a')$  are known, then the optimal strategy is to take the action that maximizes the expected value of  $r + \gamma Q^*(s', a')$

The optimal policy  $\pi^*$  corresponds to taking the best action in any state as specified by  $Q^*$

# Solving for the optimal policy

**Value iteration** algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s, a) = \mathbb{E} \left[ r + \gamma \max_{a'} Q_i(s', a') | s, a \right]$$

$Q_i$  will converge to  $Q^*$  as  $i \rightarrow \infty$

**What's the problem with this?**

Not scalable. Must compute  $Q(s, a)$  for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!

**Solution:** use a function approximator to estimate  $Q(s, a)$ . E.g. a neural network!

# Solving for the optimal policy: Q-learning

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s, a; \theta) \approx Q^*(s, a)$$

function parameters (weights)

If the function approximator is a deep neural network => **deep q-learning!**

# Solving for the optimal policy: Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q^*(s', a') | s, a \right]$$

Forward Pass

Loss function:  $L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} [(y_i - Q(s, a; \theta_i))^2]$

where  $y_i = \mathbb{E}_{s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a \right]$

Iteratively try to make the Q-value close to the target value ( $y_i$ ) it should have, if Q-function corresponds to optimal  $Q^*$  (and optimal policy  $\pi^*$ )

Backward Pass

Gradient update (with respect to Q-function parameters  $\theta$ ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[ r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right] \nabla_{\theta_i} Q(s, a; \theta_i)$$