

INTRO TO ML

The Hello World of Machine Learning

WHO AM I?

WHO AM I?

Backend and Feature Developer for

“Google Analytics for healthcare data”

Machine Learning and Touch Rugby Enthusiast



christine.winter.42@gmail.com



github.com/christinewinter



LIGHTNING TALKS & STUDY GROUP

Wednesday, June 19, 2019

5:30 PM to 8:30 PM

Foo Café (Hammarby kaj 10D • Stockholm)

www.meetup.com/PyLadiesStockholm/events



pyladies

women who program in Python

WHO IS LISTENING?

WHO KNOWS WHAT MACHINE LEARNING IS?

WHO PROGRAMMED IN PYTHON BEFORE?

WHAT IS LEARNING ?

WHAT IS LEARNING ?

The acquisition of knowledge or skills through study, experience, or being taught.

Oxford dictionary

WHAT IS LEARNING ?

The acquisition of knowledge or skills through study, experience, or being taught.

Oxford dictionary

Understand what's going on with different methods and use this information the next time.

My interpretation

HOW DO WE LEARN ?

HOW DO WE LEARN ?

Get instructions from others

HOW DO WE LEARN ?

Get instructions from others

Interpret a situation and try to define rules or draw conclusions alone

HOW DO WE LEARN ?

Get instructions from others

Interpret a situation and try to define rules or draw conclusions alone

Learn from previous mistakes or success

HOW DO MACHINES LEARN ?

HOW DO MACHINES LEARN ?

Get instructions from others

Interpret a situation and try to define rules or draw conclusions alone

Learn from previous mistakes or success

HOW DO MACHINES LEARN ?

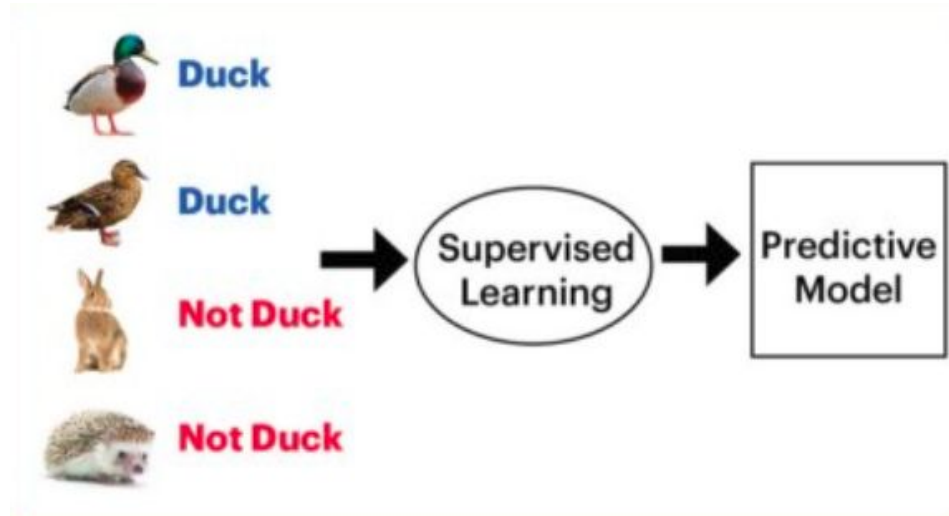
Get instructions from others

→ **supervised learning**

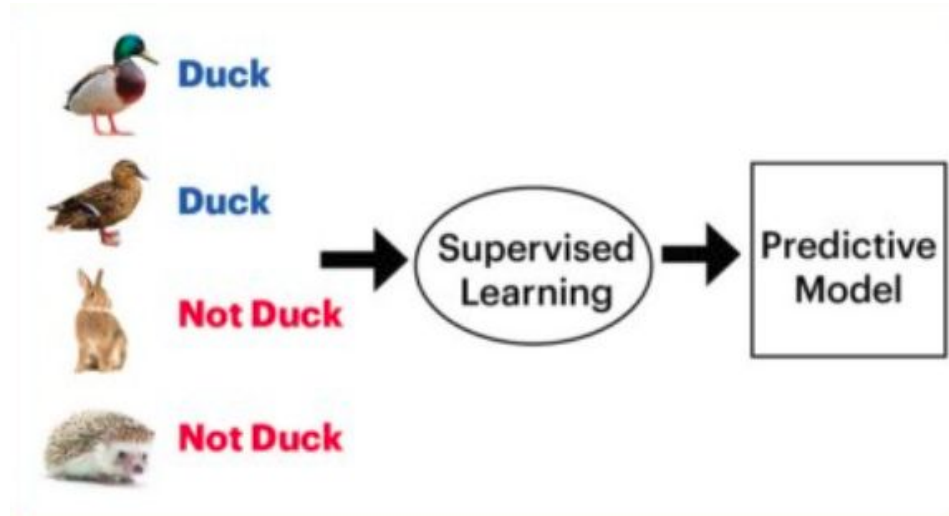
Interpret a situation and try to define rules or draw conclusions alone

Learn from previous mistakes or success

SUPERVISED



SUPERVISED



HOW DO MACHINES LEARN ?

Get instructions from others

→ **supervised learning**

Interpret a situation and try to define rules or draw conclusions alone

Learn from previous mistakes or success

HOW DO MACHINES LEARN ?

Get instructions from others

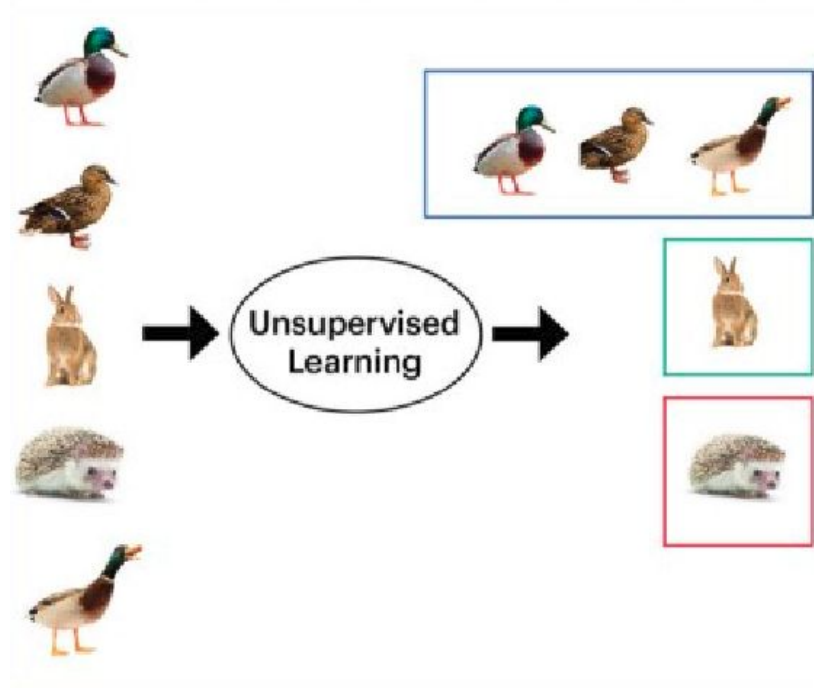
→ supervised learning

Interpret a situation and try to define rules or draw conclusions alone

→ **unsupervised learning**

Learn from previous mistakes or success

UNSUPERVISED



HOW DO MACHINES LEARN ?

Get instructions from others

→ supervised learning

Interpret a situation and try to define rules or draw conclusions alone

→ **unsupervised learning**

Learn from previous mistakes or success

HOW DO MACHINES LEARN ?

Get instructions from others

→ supervised learning

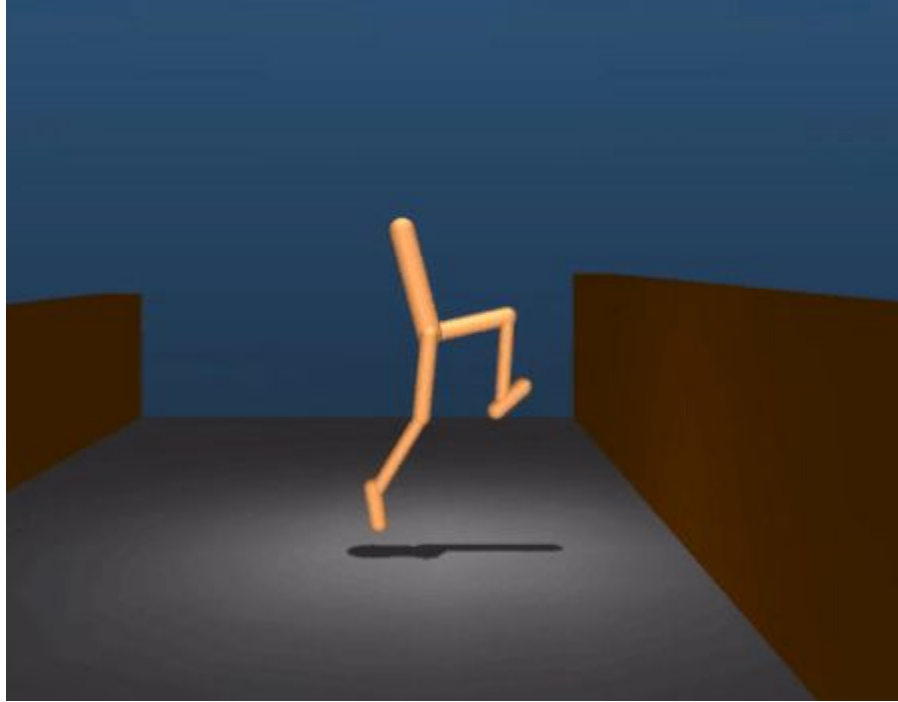
Interpret a situation and try to define rules or draw conclusions alone

→ unsupervised learning

Learn from previous mistakes or success

→ **reinforcement learning**

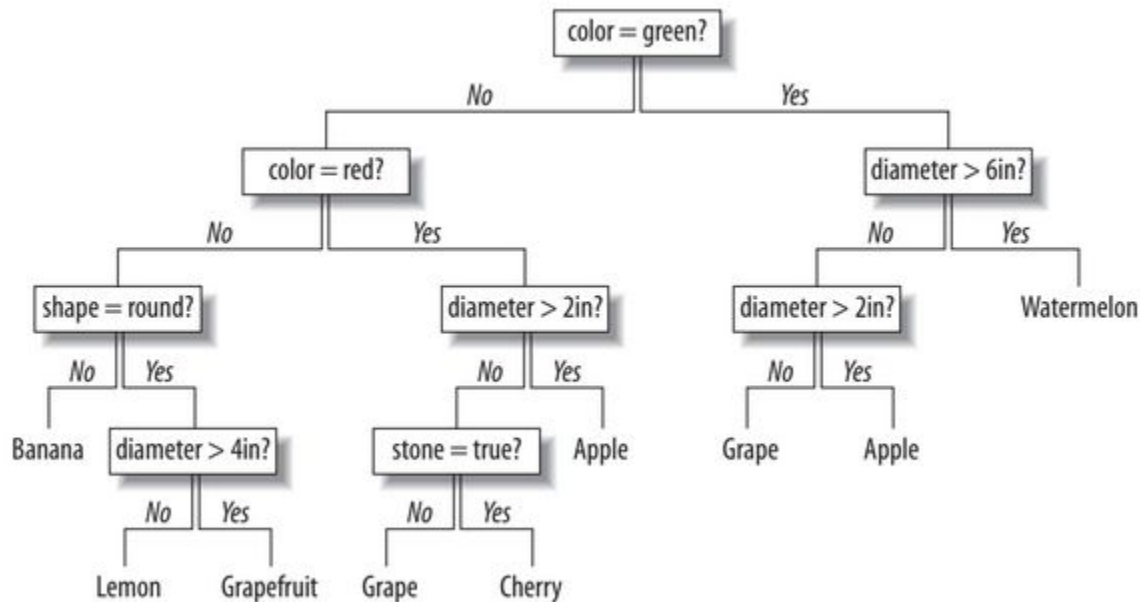
REINFORCEMENT



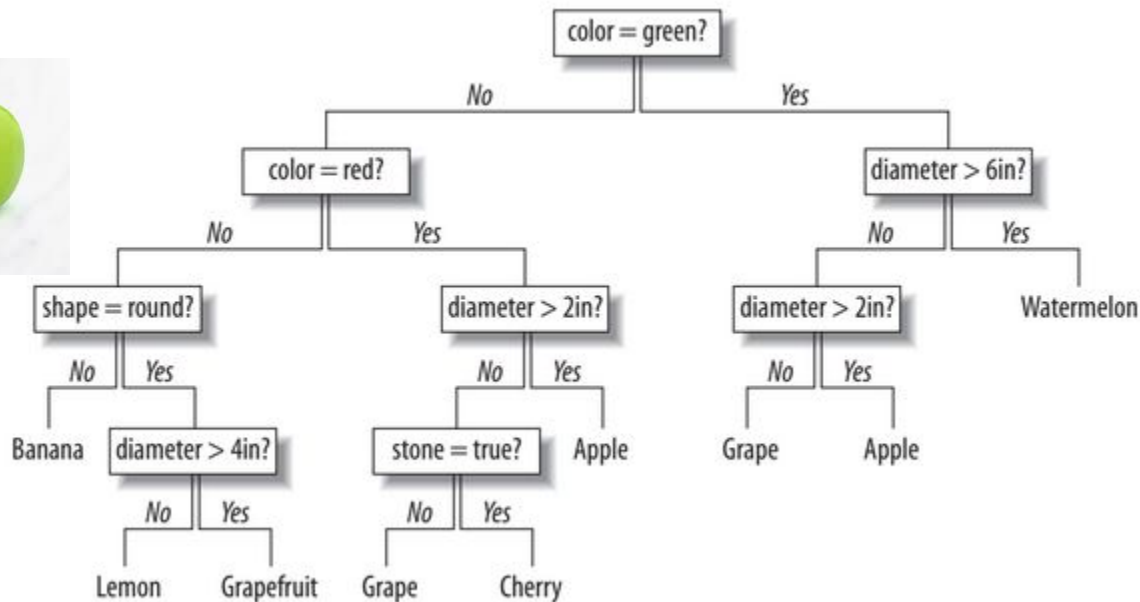
<https://www.theverge.com/tldr/2017/7/10/15946542/deepmind-parkour-agent-reinforcement-learning>

SO, WHAT WILL WE DO?

SO, WHAT WILL WE DO? - LEARN WITH DECISION TREES



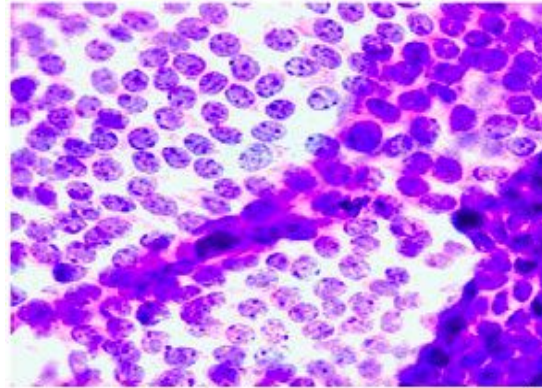
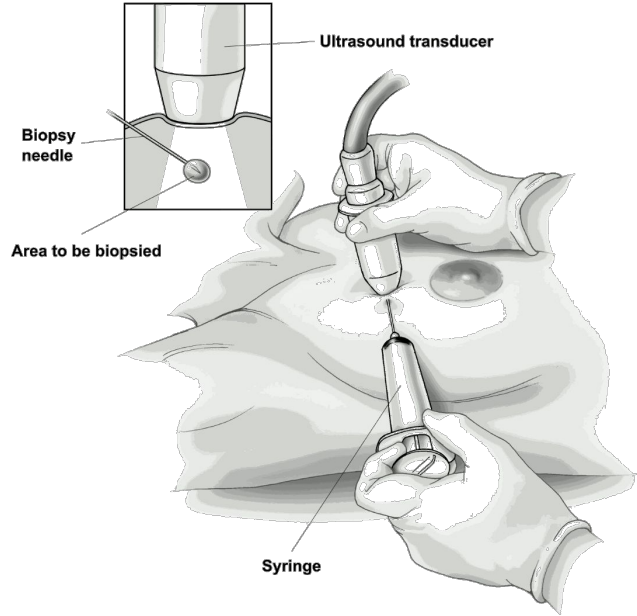
SO, WHAT WILL WE DO? - LEARN WITH DECISION TREES



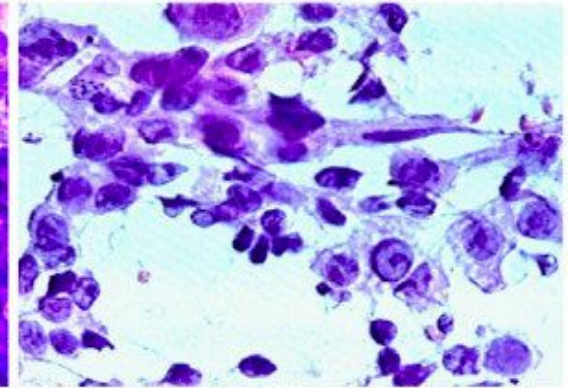
THE DATA

- ❖ Breast Cancer Wisconsin (Diagnostic) Data Set acquired by researchers University of Wisconsin (Dr. W.H. Wolberg, N. Street, O.L. Mangasarian)
- ❖ Features computed from a digitized image of a fine needle aspirate (FNA) of a breast mass
- ❖ Describe characteristics of cell nuclei present in the image

THE DATA



Smear with BENIGN diagnosis – uniform nucleus of cells, symmetrical, homogeneous, with areas within normal size



Smear with MALIGNANT diagnosis – nucleus of cells without uniformity, asymmetrical, not homogeneous (multiple sizes) and with areas above normal size

Fine needle aspiration using ultrasound

© Sam and Amy Collins

LETS GO !

SETTING UP

Local: clone repo, install requirements.txt

https://github.com/christinewinter/intro_ml_bc

Remote:

Go to: <https://mybinder.org/>

Insert:

https://github.com/christinewinter/intro_ml_bc

Will take ~ 5 minutes



Turn a Git repo into a collection of interactive notebooks

Have a repository full of Jupyter notebooks? With Binder, open those notebooks in an executable environment, making your code immediately reproducible by anyone, anywhere.

Build and launch a repository


GitHub repository name or URL
 GitHub ▾

Git branch, tag, or commit

Path to a notebook file (optional)
 File ▾

Copy the URL below and share your Binder with others:

📄

Copy the text below, then paste into your README to show a binder badge:  [launch](#) [binder](#) ▶

Waiting Already built!

Build logs show

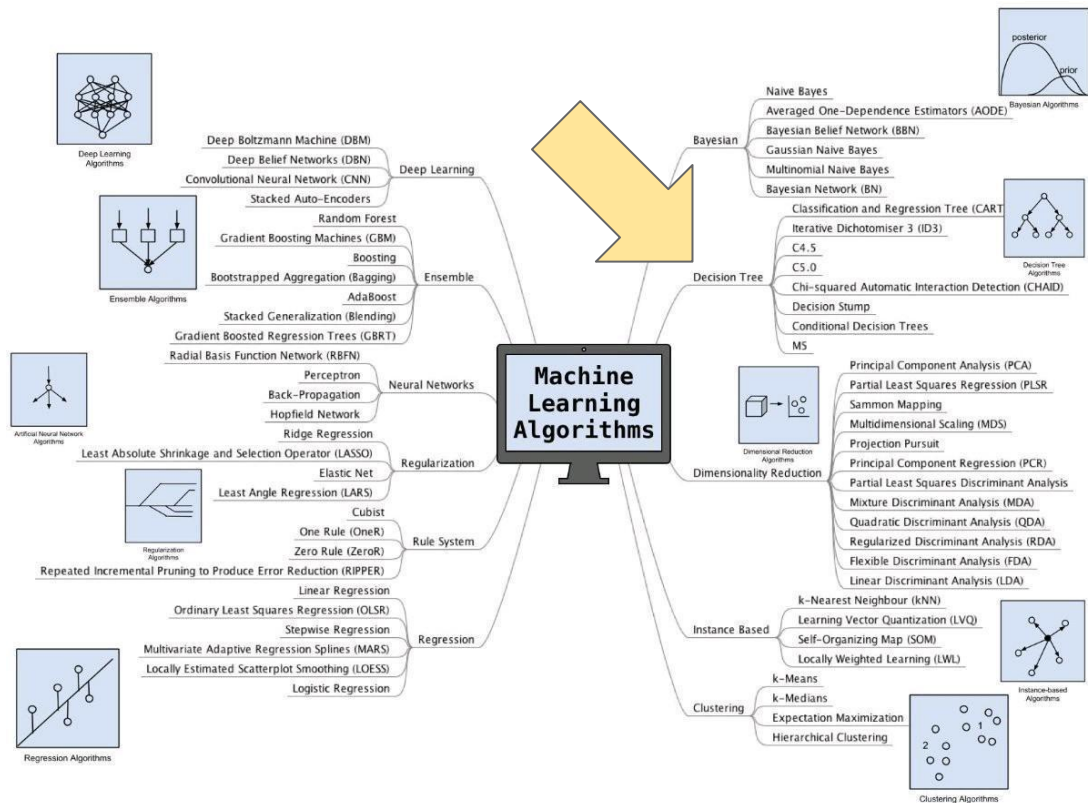
THE ACTUAL MACHINE LEARNING PART

```
model = DecisionTreeClassifier()  
model.fit(X_train, y_train)  
y_pred = model.predict(X_test)  
metrics.accuracy_score(y_test, y_pred)  
y_pred = model.predict(X_unknown)
```



scikit-learn.org

DECISION TREES ARE JUST ONE WAY TO LEARN



WHAT IS IMPORTANT FOR ML ?

WHAT IS IMPORTANT FOR ML ?

- ★ Data driven
 - Density, diversity, structure

WHAT IS IMPORTANT FOR ML ?

- ★ Data driven
 - Density, diversity, structure
- ★ Problem definition
 - Supervised, unsupervised, reinforcement learning

WHAT IS IMPORTANT FOR ML ?

- ★ Data driven
 - Density, diversity, structure
- ★ Problem definition
 - Supervised, unsupervised, reinforcement learning
- ★ Model choice
 - Decision tree, deep neural networks, support vector machines...

WHAT IS IMPORTANT FOR ML ?

- ★ Data driven
 - Density, diversity, structure
- ★ Problem definition
 - Supervised, unsupervised, reinforcement learning
- ★ Model choice
 - Decision tree, deep neural networks, support vector machines...
- ★ Metrics to evaluate the model
 - Prediction score, prevention of overfitting

WHAT IS IMPORTANT FOR ML ?

- ★ Data driven
 - Density, diversity, structure
- ★ Problem definition
 - Supervised, unsupervised, reinforcement learning
- ★ Model choice
 - Decision tree, deep neural networks, support vector machines...
- ★ Metrics to evaluate the model
 - Prediction score, prevention of overfitting
- ★ Feature selection
 - Relevance, correlations

WHAT IS IMPORTANT FOR ML ?

- ★ Data driven
 - Density, diversity, structure
- ★ Problem definition
 - Supervised, unsupervised, reinforcement learning
- ★ Model choice
 - Decision tree, deep neural networks, support vector machines...
- ★ Metrics to evaluate the model
 - Prediction score, prevention of overfitting
- ★ Feature selection
 - Relevance, correlations
- ★ Iteration & architecture
 - Implementation, backup

WHAT YOU LEARNED TODAY:

WHAT YOU LEARNED TODAY:

★ Types of learning:

- supervised, unsupervised, reinforcement

WHAT YOU LEARNED TODAY:

- ★ Types of learning:

- supervised, unsupervised, reinforcement

- ★ One machine learning model:

- Decision trees

WHAT YOU LEARNED TODAY:

- ★ Types of learning:
 - supervised, unsupervised, reinforcement
- ★ One machine learning model:
 - Decision trees
- ★ Data is the key component
 - Clean up is 80 % of the work

WHAT YOU LEARNED TODAY:

- ★ Types of learning:
 - supervised, unsupervised, reinforcement
- ★ One machine learning model:
 - Decision trees
- ★ Data is the key component
 - Clean up is 80 % of the work
- ★ Implementation
 - Python, jupyter, Scikit learn

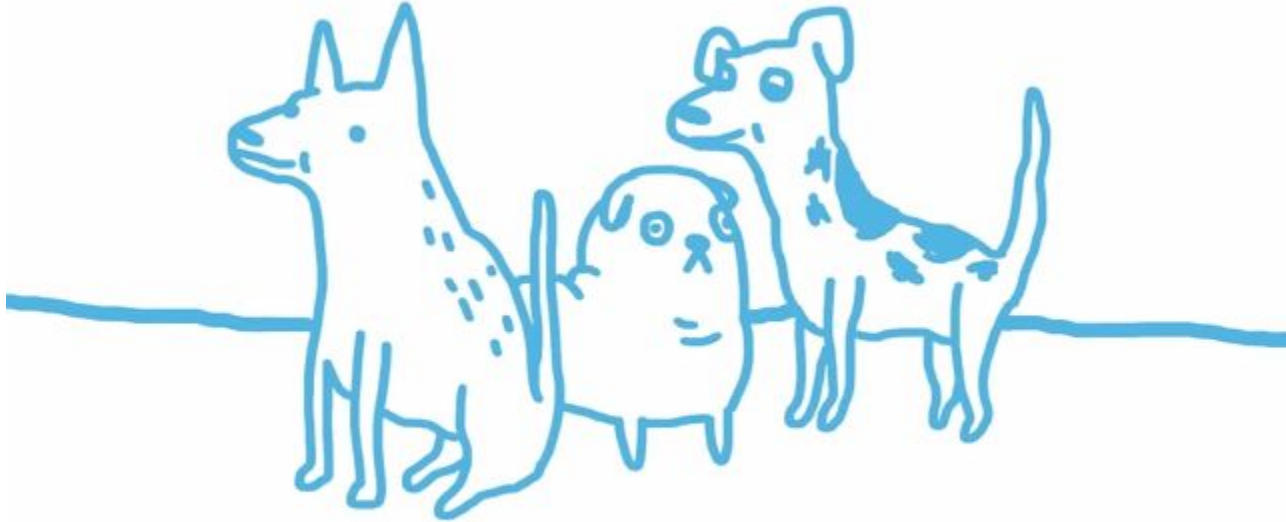
BUZZWORD BINGO :

A word cloud of AI and ML buzzwords arranged in a bingo card grid. The words are of varying sizes and orientations, scattered across the grid. The words include:

- Artificial Neural Network
- Machine Learning - ML
- Artificial Intelligence - AI
- Cloud AI
- Data Mining
- Deep Learning - DL
- AI in IoT
- Artificial Neocortex
- DevOps
- AIOps - automated
- BIG data
- Data driven
- Tensorflow
- Keras
- CNN
- DNN
- ANN

Artificial Neural Network	Machine Learning - ML	Artificial Intelligence - AI		
DNN	Deep Learning - DL	Cloud AI		Data Mining
ANN	AI in IoT	AIOps - automated	BIG data	
CNN	Tensorflow	DevOps	Artificial Neocortex	Data driven
Keras				

QUESTIONS ?



GINI IMPURITY

Gini impurity is the expected error rate if one of the results from a set is randomly applied to one of the items in the set.

If every item in the set is in the same category, the guess will always be correct, so the error rate is 0. If there are four possible results evenly divided in the group, there's a 75 percent chance that the guess would be incorrect, so the error rate is 0.75.

This function calculates the probability of each possible outcome by dividing the number of times that outcome occurs by the total number of rows in the set. It then adds up the products of all these probabilities. This gives the overall chance that a row would be randomly assigned to the wrong outcome.

The higher this probability, the worse the split. A probability of zero is great because it tells you that everything is already in the right set.