

DSE8-Prediction



DSE - Data Science Course by Data Scientists

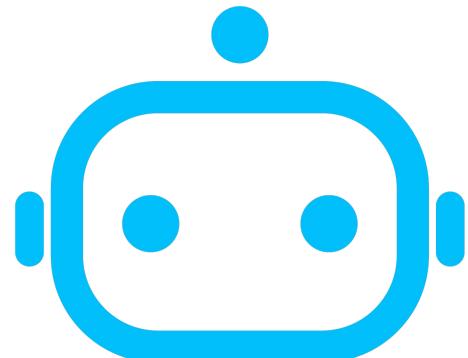
About myself



Head of Digital Human Witthawin Sripheanpol

- B.Eng. Institute of Field Robotics, King Mongkut's University of Technology Thonburi
- Project Manager & Lead Lecturer iGenius Robot Education, Thailand
- Senior Data Scientist & Machine Learning Engineer Botnoi Group, Thailand
- Head of Digital Human Botnoi Group, Thailand
- Co-Founder FB: HowKnow I Know





BOTNOI

We localise AI

Agenda

1. What is Predictive Model?
2. How it work?
3. How to do Predictive Model on Production

[Workshop] Streamlit & Production

1

What Is Predictive Model?

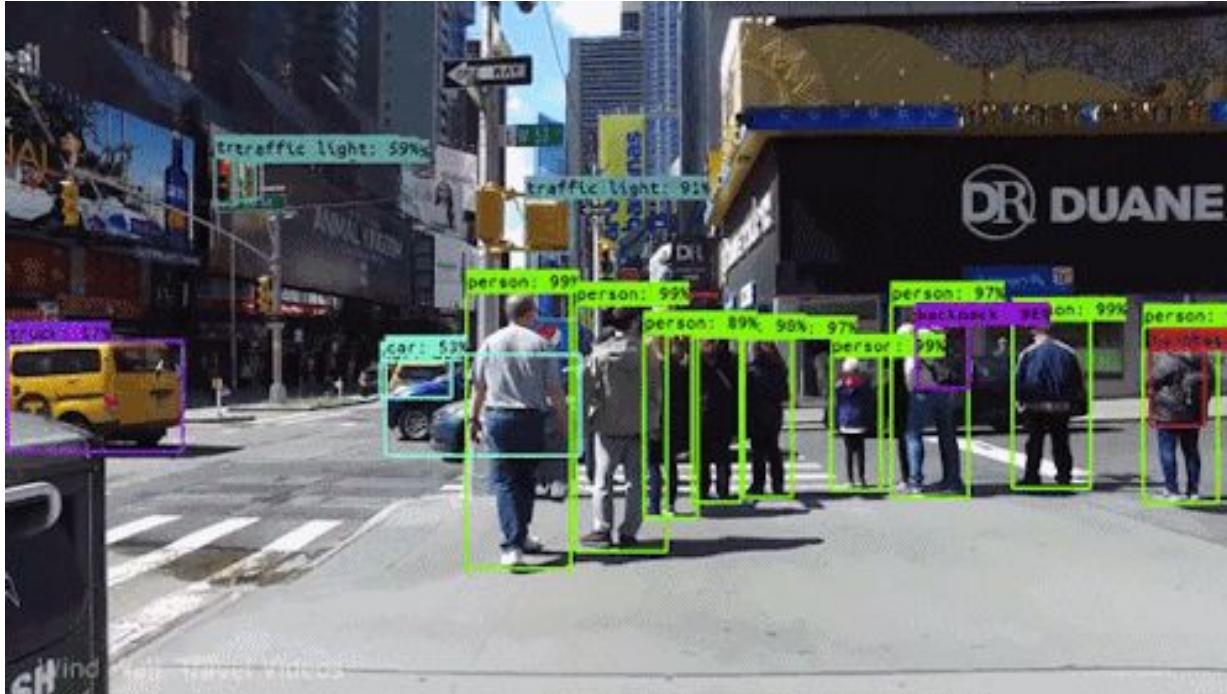




How to ขุดตันไม้
ขอวยอย่างไร ให้ได้ผล

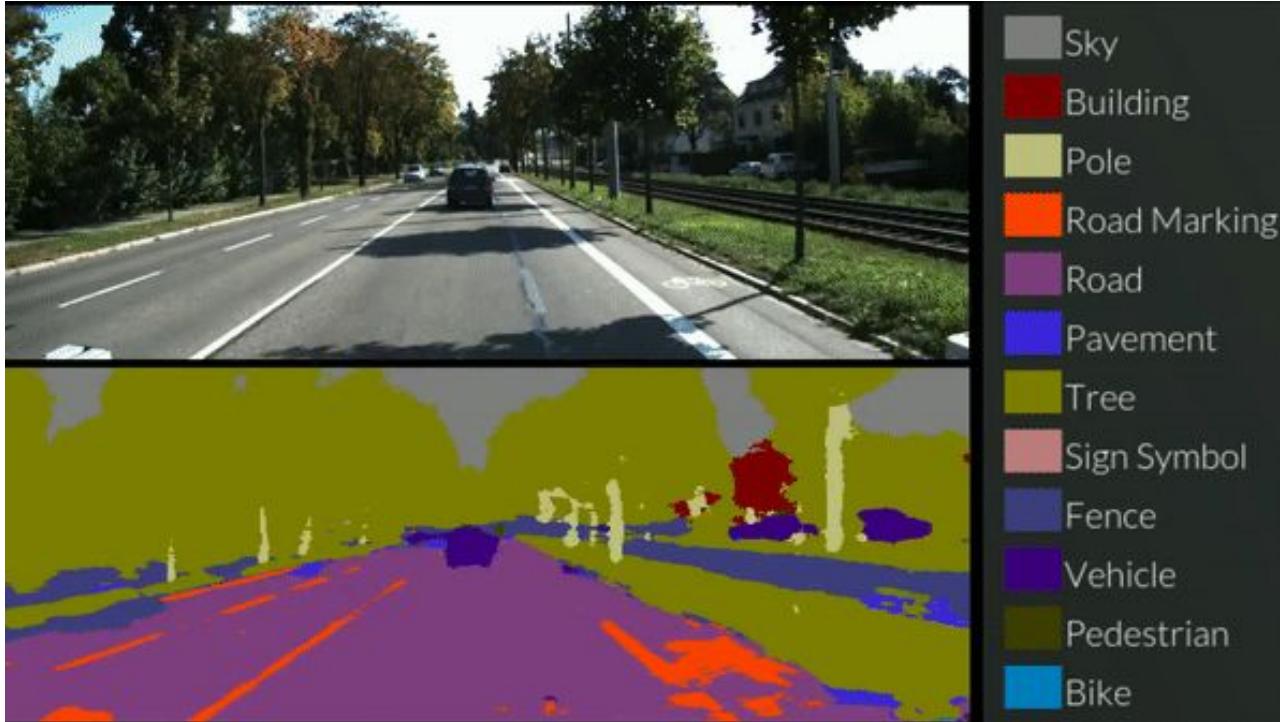
Computer Vision

(Object Detection)



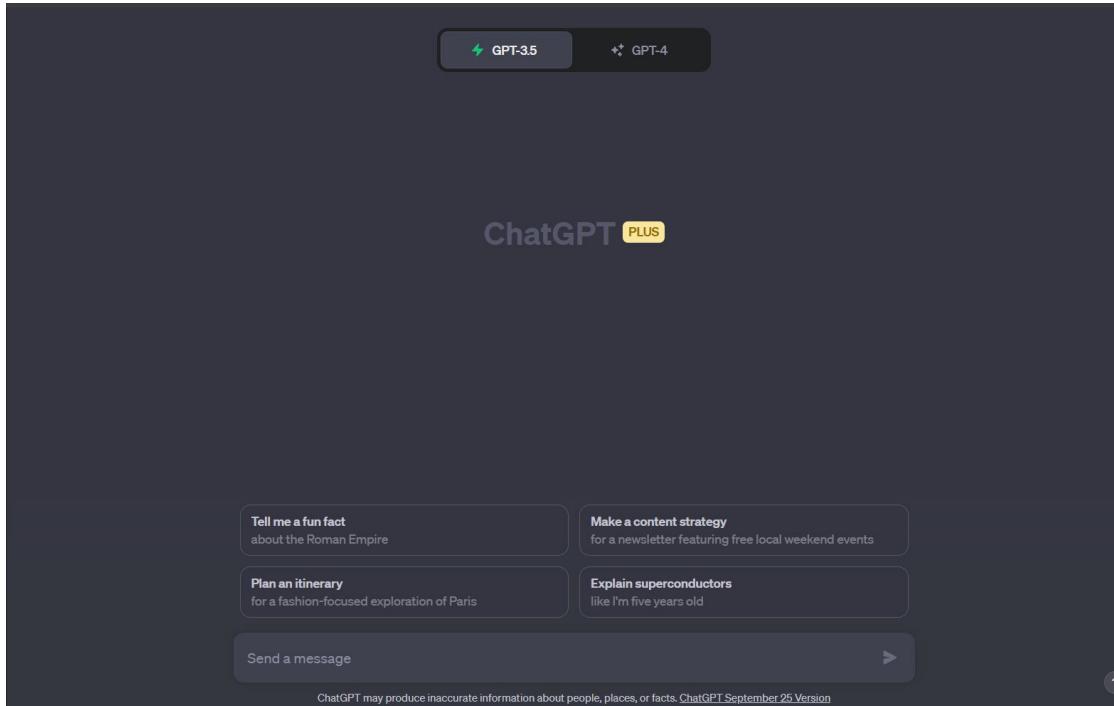
Computer Vision

(Image Segmentation)



Natural Language Processing

(AI understand text, message)



Game & Simulation

(AI understand Mechanics and Motion)



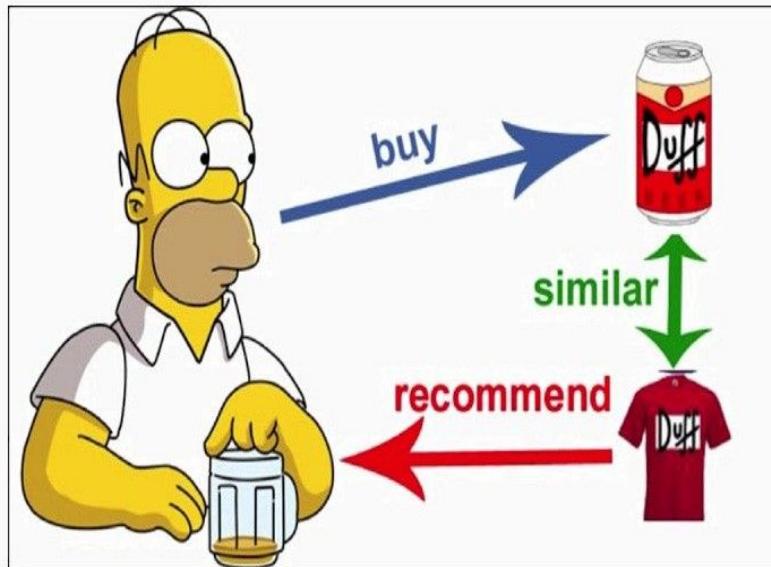
Robotics & Automation

(AI understand Mechanics and Motion)



Recommendation

(AI understand behavior)

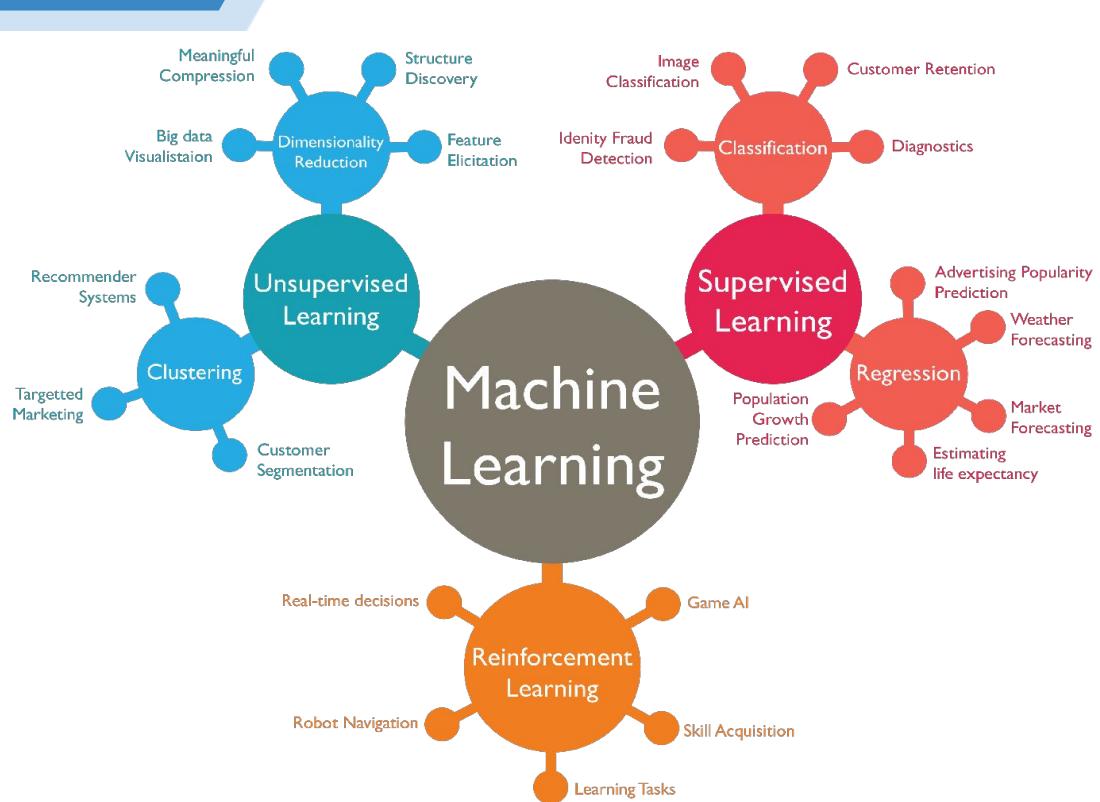
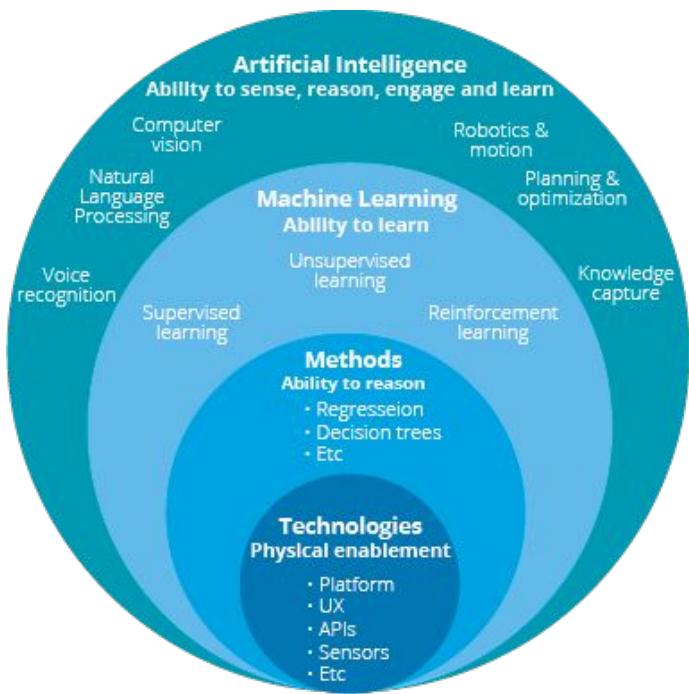


2

How it work?



Recap : Machine Learning

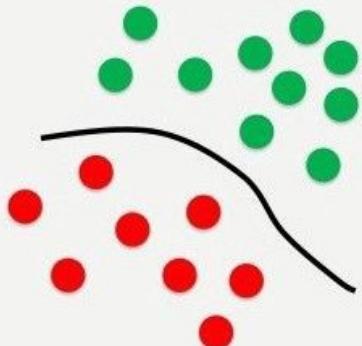


Supervised Learning

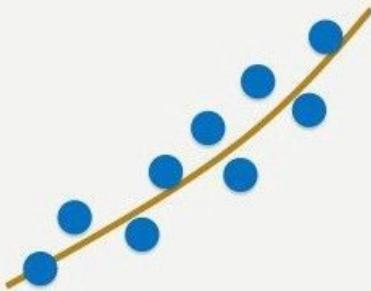
Classification vs Regression

CLASSIFICATION VS REGRESSION

Classification

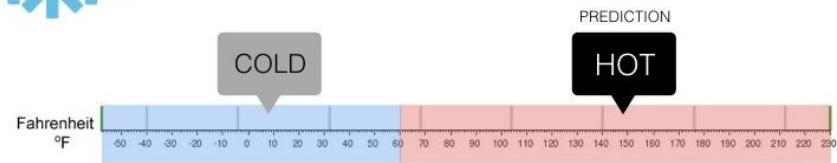


Regression



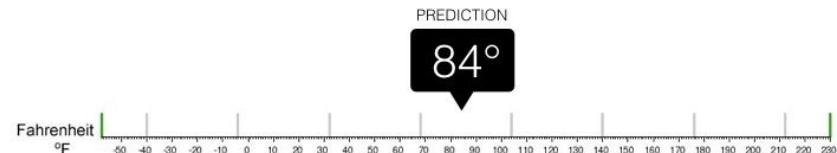
Classification

Will it be Cold or Hot tomorrow?



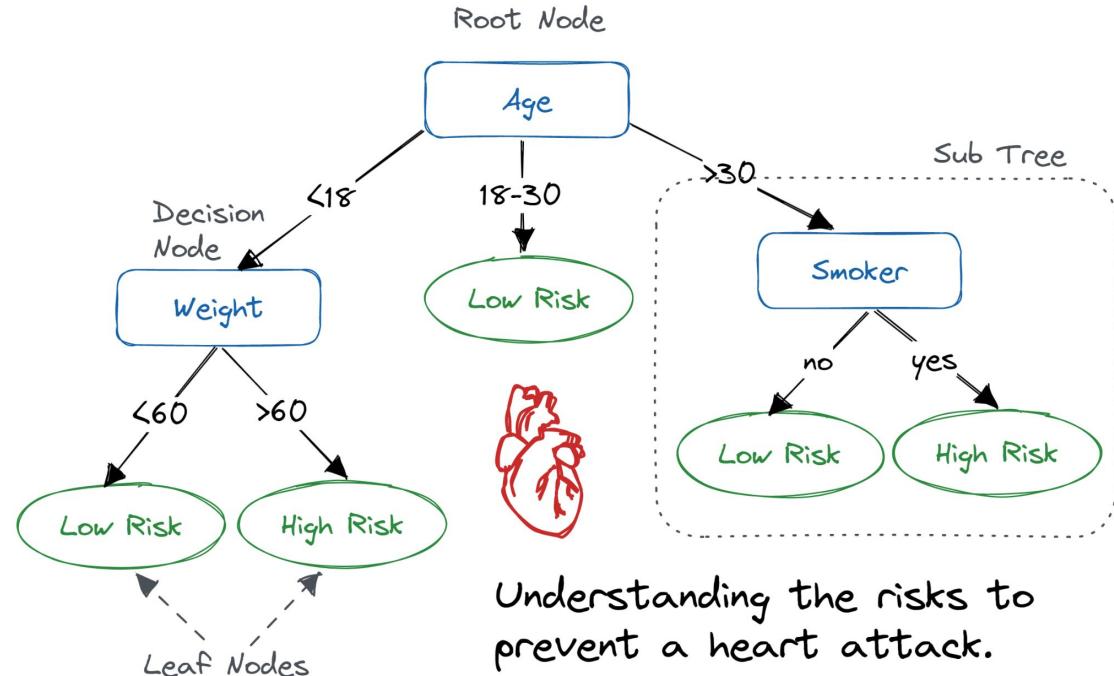
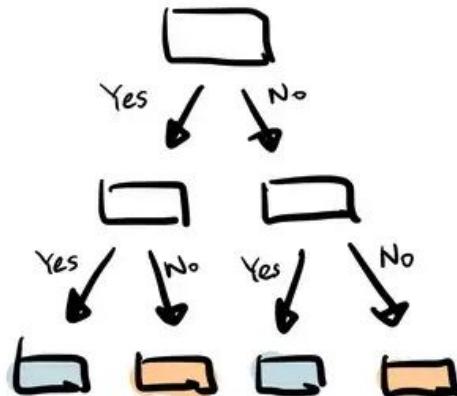
Regression

What is the temperature going to be tomorrow?



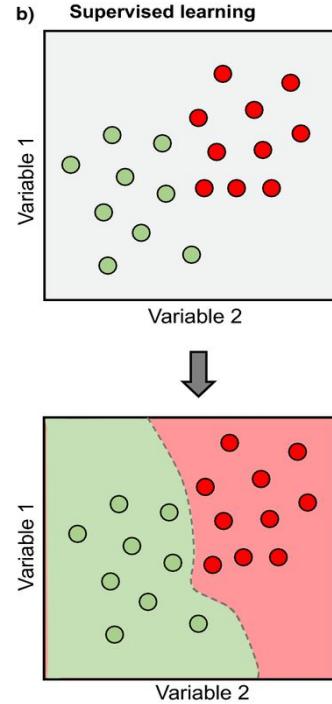
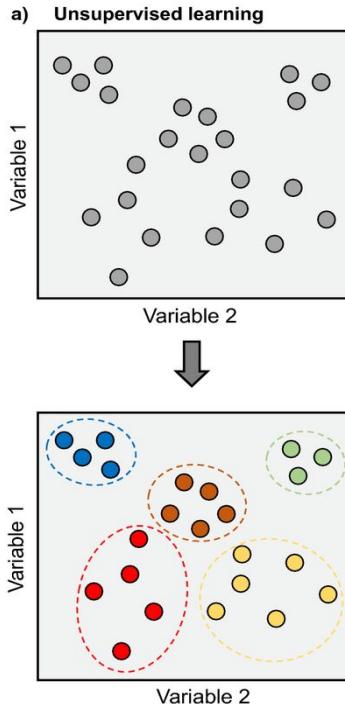
Classification Algorithm

(Basic Algorithm : Decision Tree)



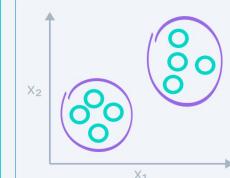
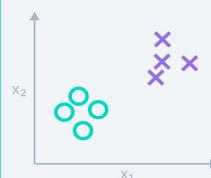
Unsupervised Learning

(unknown label)



Supervised learning
Input data is labeled
Has a feedback mechanism
Data is classified based on the training dataset
Divided into Regression & Classification
Used for prediction
Algorithms include: decision trees, logistic regressions, support vector machine
A known number of classes

Unsupervised learning
Input data is unlabeled
Has no feedback mechanism
Assigns properties of given data to classify it
Divided into Clustering & Association
Used for analysis
Algorithms include: k-means clustering, hierarchical clustering, apriori algorithm
A unknown number of classes

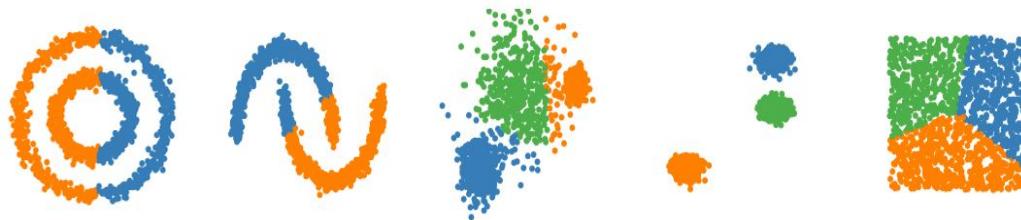


V7 Labs

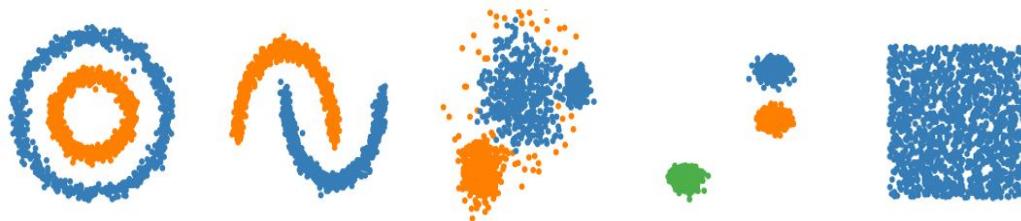
Clustering Algorithm

(Basic Algorithm for grouping data)

k-means

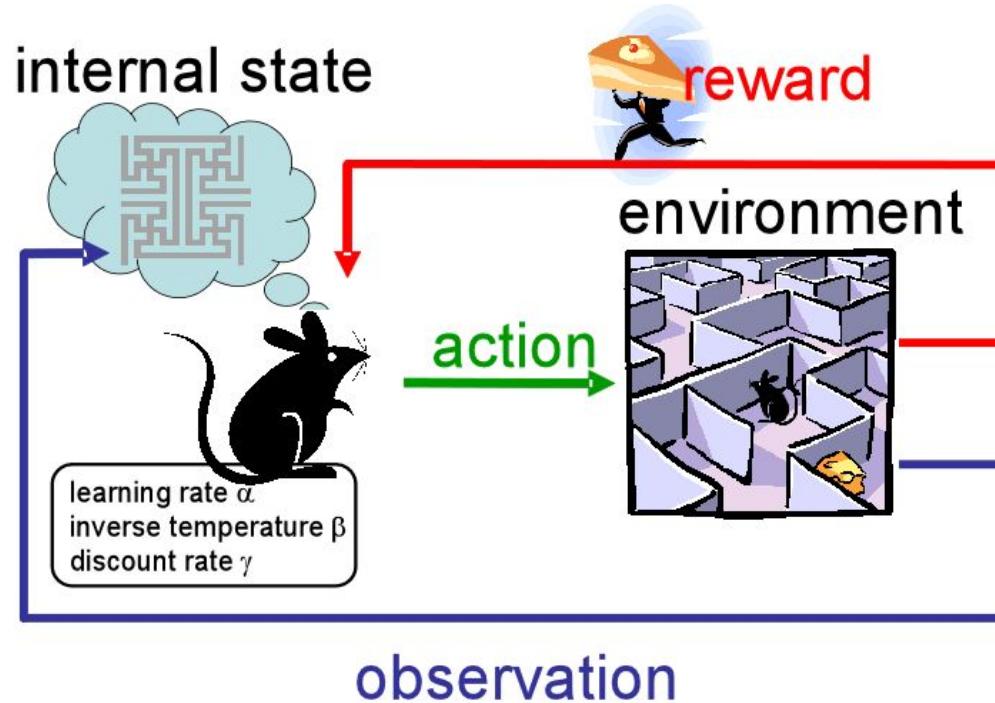


DBSCAN



Reinforcement Learning

(I don't know but I'll maximize reward)



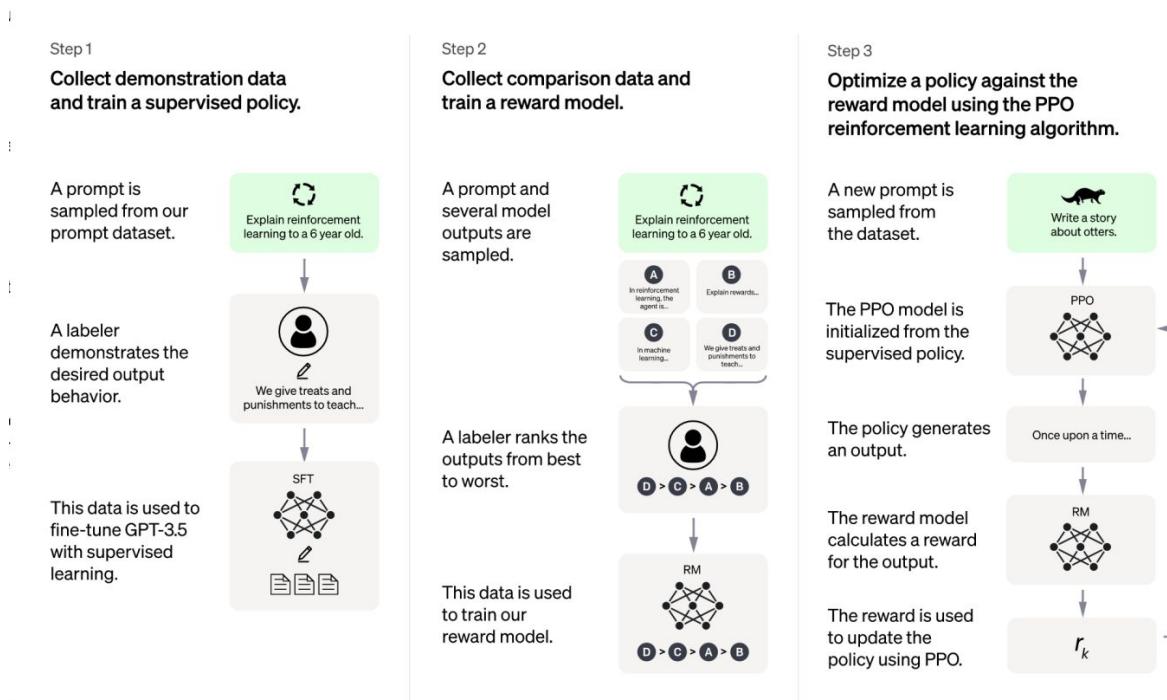
Reinforcement Learning

(Sample Computing)



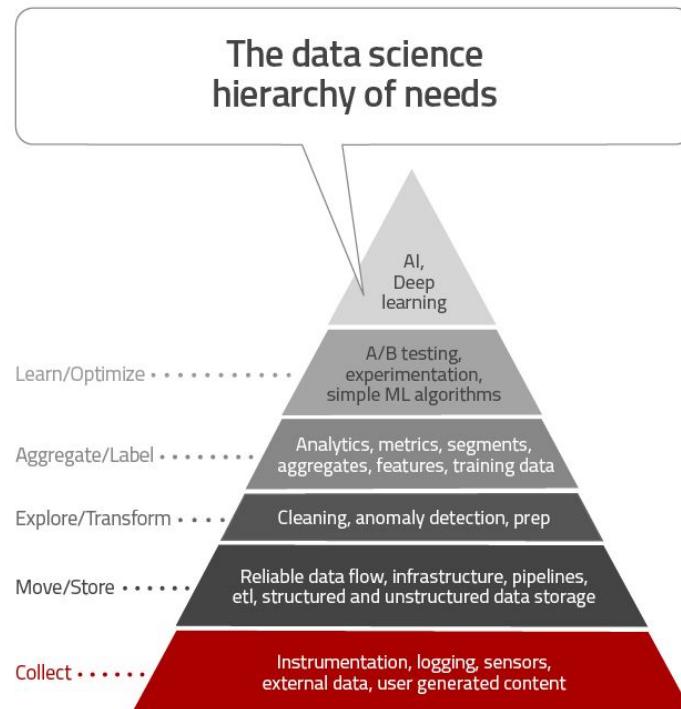
Reinforcement Learning

(Deep RL in GPT & Other Large Model)



Data to AI Application

(AI driven by Data)



SOURCE: Monica Rogati © August 2017 The Financial Brand

3

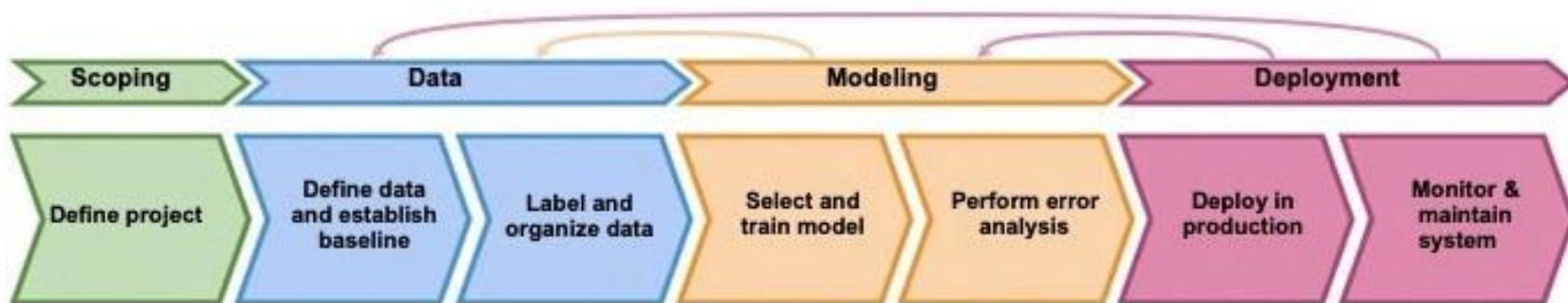
How to do Prediction



Machine Learning Pipeline

(ML & Prediction Lifecycle)

The ML Project Lifecycle



Problem

(Define Project for solve problem)



Data Collection

(collect many data)

Structured data



Characteristics

- Predefined data models
- Easy to search
- Text-based
- Shows what's happening

Resides in

- Relational databases
- Data warehouses

Stored in

- Rows and columns

Examples

- Dates, phone numbers, social security numbers, customer names, transaction info

Unstructured data



Characteristics

- No predefined data models
- Difficult to search
- Text, pdf, images, video
- Shows the why

Resides in

- Applications
- Data warehouses and lakes

Stored in

- Various forms

Examples

- Documents, emails and messages, conversation transcripts, image files, open-ended survey answers

Semi-structured data



Characteristics

- Loosely organized
- Meta-level structure that can contain unstructured data
- HTML, XML, JSON

Resides in

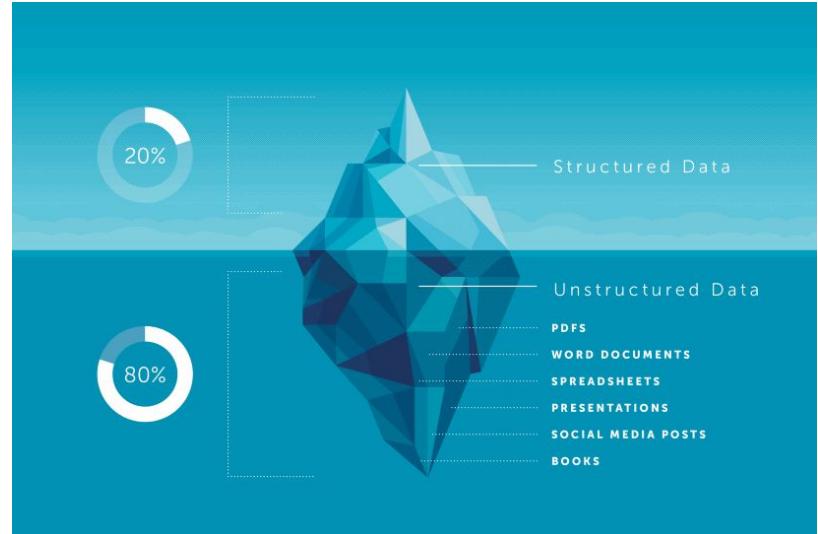
- Relational databases
- Tagged-text format

Stored in

- Abstracts & figures

Examples

- Server logs, tweets organized by hashtags, emails sorting by folders (inbox; sent; draft)



LEVITY

Data Collection

(Public Dataset)



Hugging Face

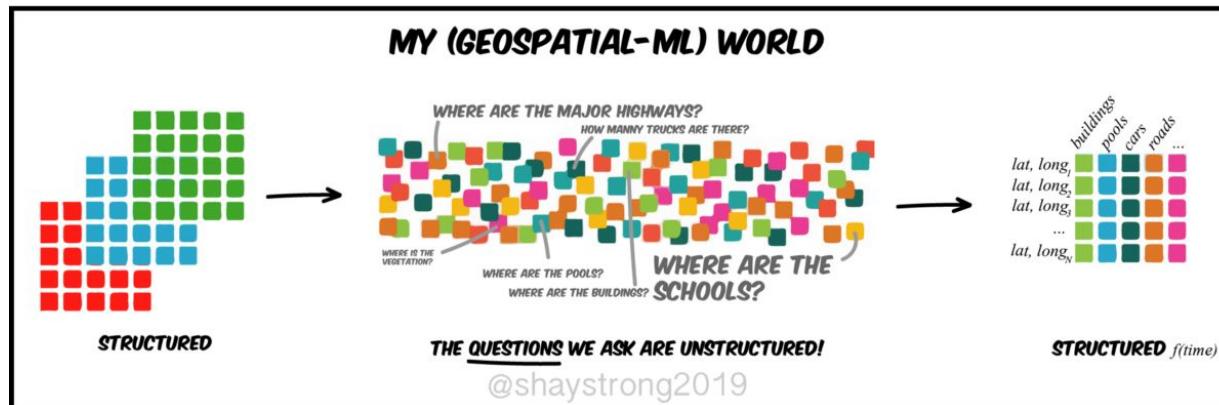
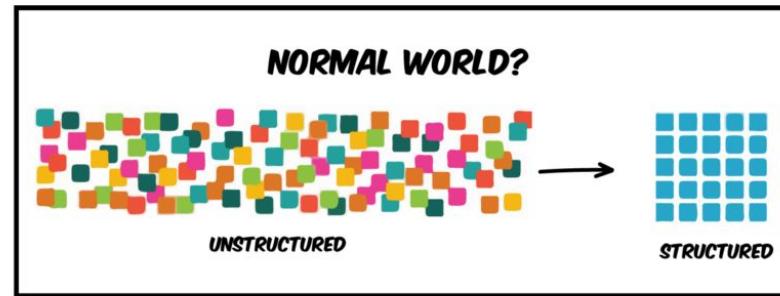
Google

kaggle

<https://pub.towardsai.net/best-datasets-for-machine-learning-data-science-computer-vision-nlp-ai-c9541058cf4f>

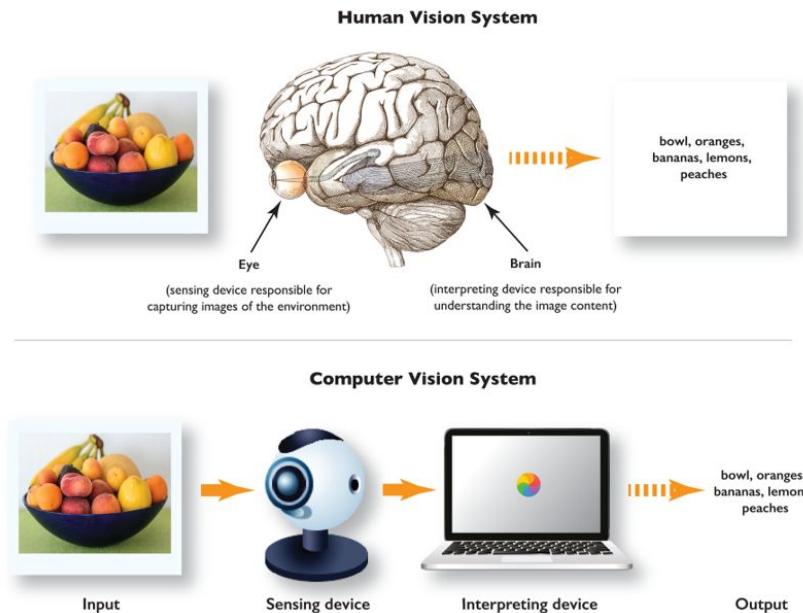
Data Preparation (ETL)

(Extract + Transform + Load)



Computer Vision

(AI understand Media)

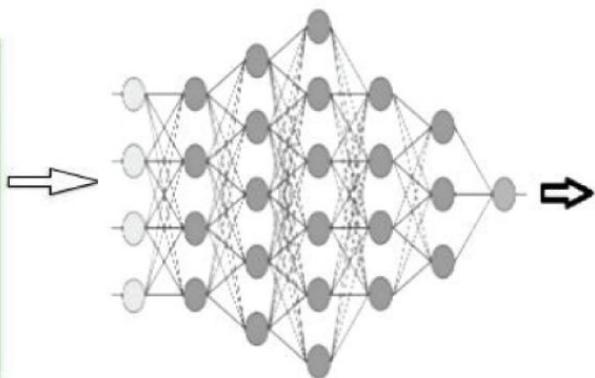


0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6
4	5	3	0	1	2	3	4	5
7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0

Natural Language Processing

(AI understand text, message)

word
The
Cardinals
will
win
the
world
series

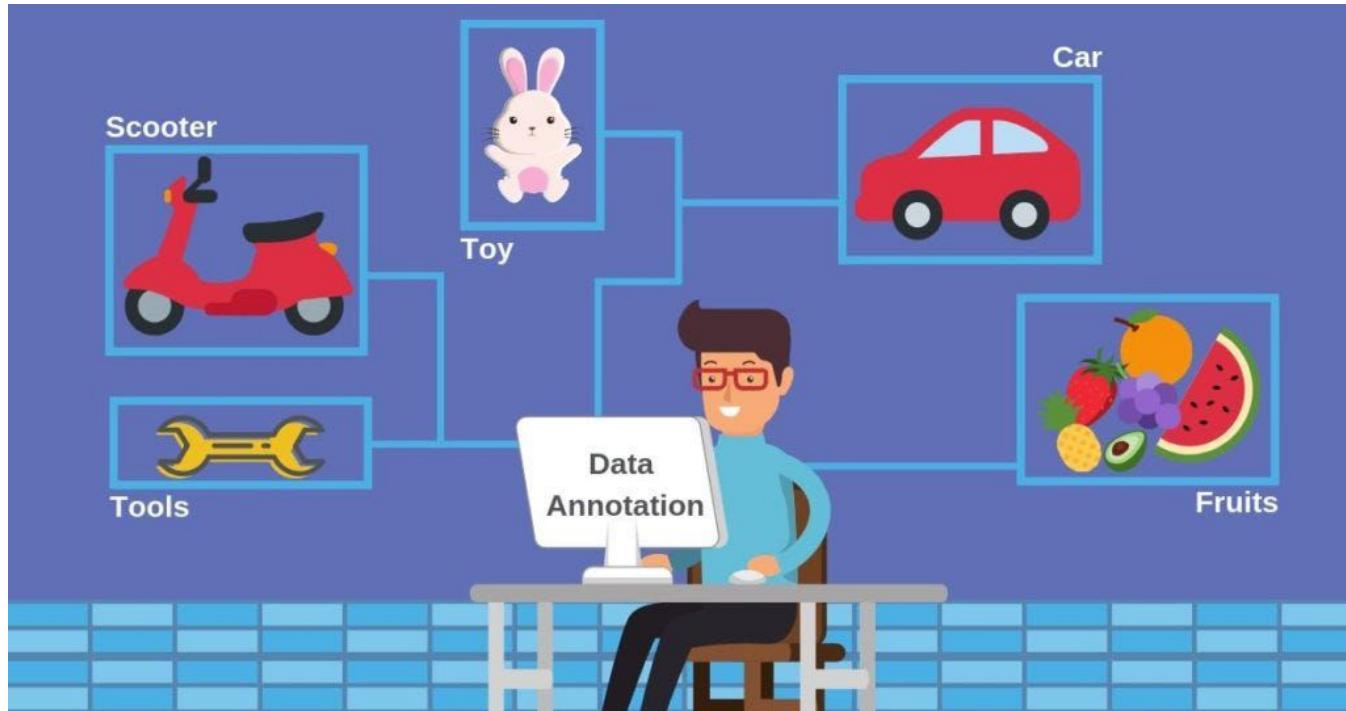


word	vector
The	(0.12, 0.23, 0.56)
Cardinals	(0.24, 0.65, 0.72)
will	(0.38, 0.42, 0.12)
win	(0.57, 0.01, 0.02)
the	(0.53, 0.68, 0.91)
world	(0.11, 0.27, 0.45)
series	(0.01, 0.05, 0.62)

sentence vector
(0.28, 0.33, 0.49)

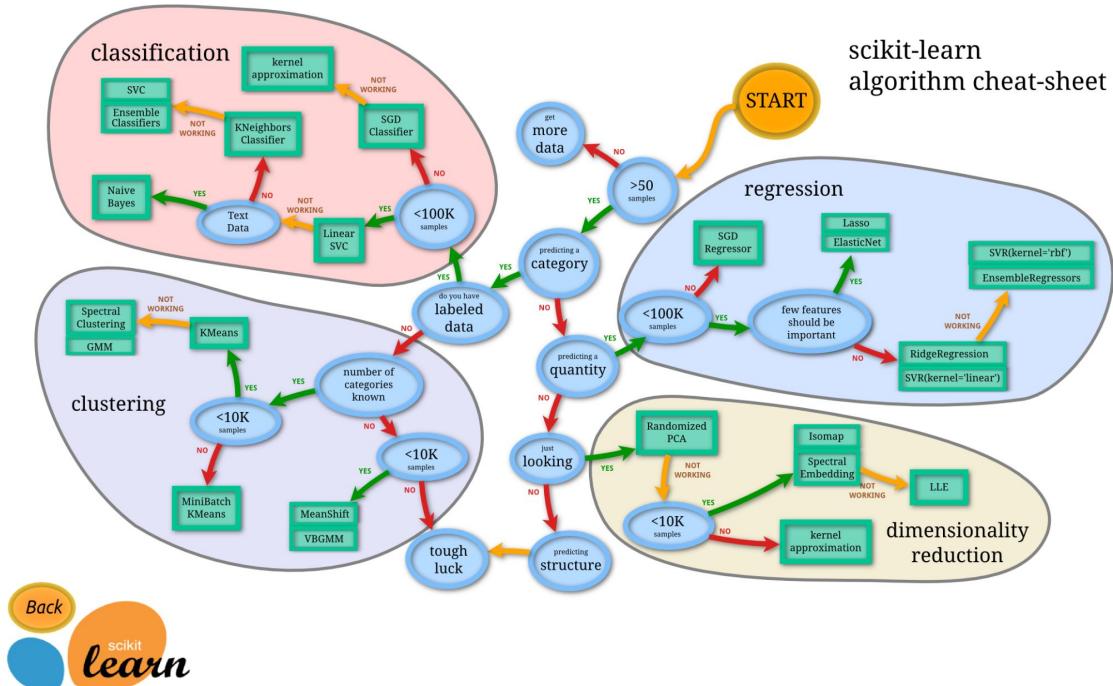
Data Labeling

(make choice or answer)



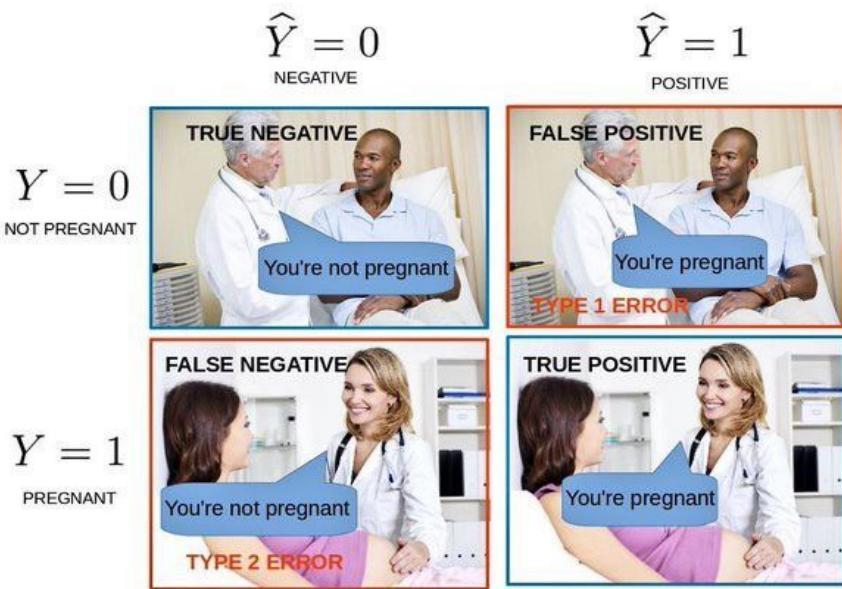
Machine Learning Algorithm

(Select Algorithm)



Performance Classification

(Confusion Matrix)



- True Positive (TP) : Actual=0 & Predict=0
- True Negative (TN) : Actual=0 & Predict=1
- False Positive (FP) : Actual=1 & Predict=0
- False Negative (FN) : Actual=1 & Predict=1

accuracy (ACC)

$$ACC = \frac{TP + TN}{P + N} = \frac{TP + TN}{TP + TN + FP + FN}$$

precision or positive predictive value (PPV)

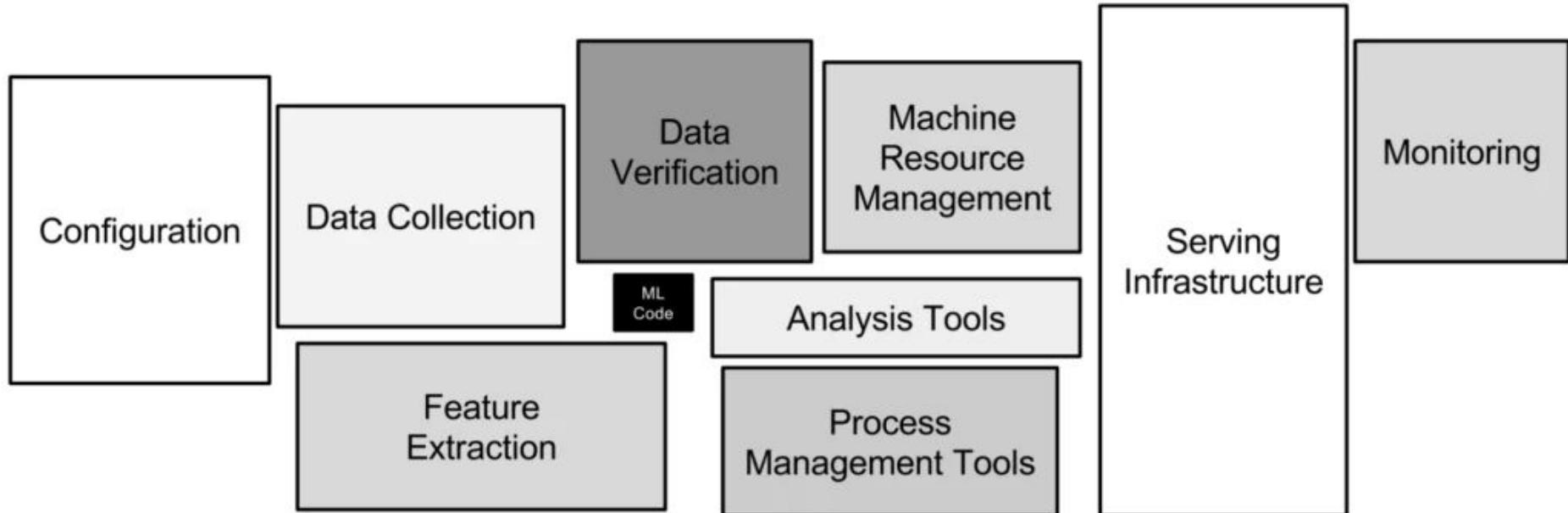
$$PPV = \frac{TP}{TP + FP} = 1 - FDR$$

sensitivity, recall, hit rate, or true positive rate (TPR)

$$TPR = \frac{TP}{P} = \frac{TP}{TP + FN} = 1 - FNR$$

Machine Learning Pipeline

(ML is 10% of Project)

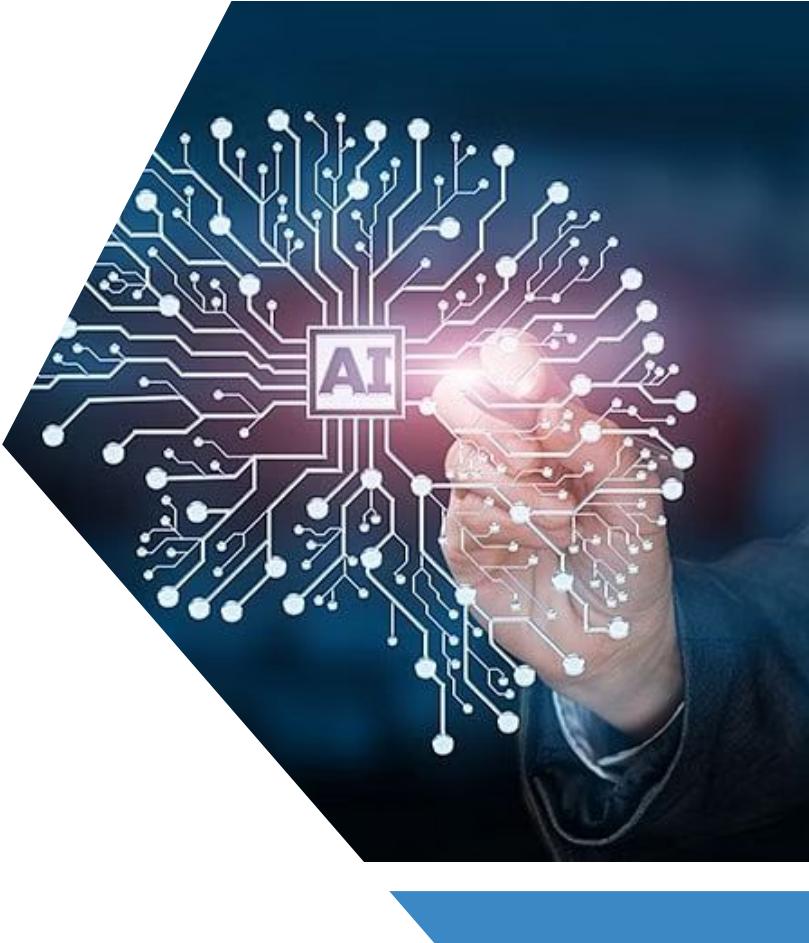


Predictive Production



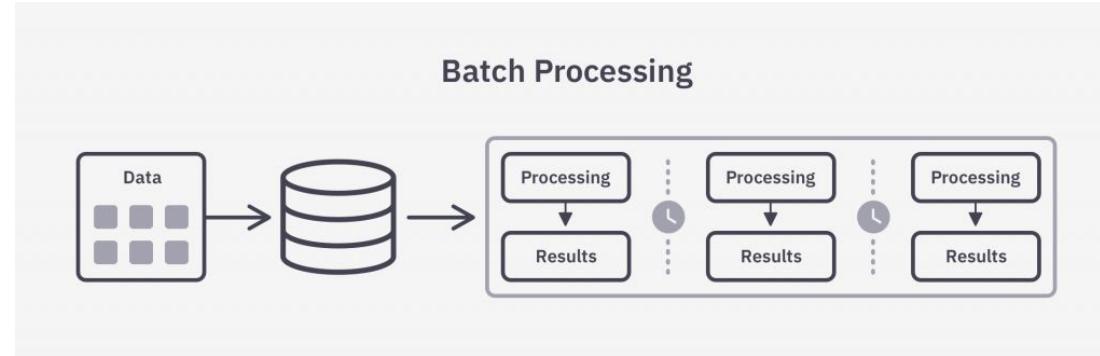
“The study needs to match application”

Bonus



Batch vs Stream Processing

(life cycle production)



Data Stream Processing



Test Production

(POC on Streamlit Cloud)

Latest news: In-browser editing, st.status to see your app's processes, custom chatbots with Llamaindex 🎉

Cloud Gallery Components Generative AI Community Docs Blog Sign in Sign up

Community Cloud

Deploy, manage, and share your apps with the world, directly from Streamlit — all for free.

Get Started

NYC Ridesharing Data

Select hour of pickup
0 0 23

All New York City from 0:00 and 1:00

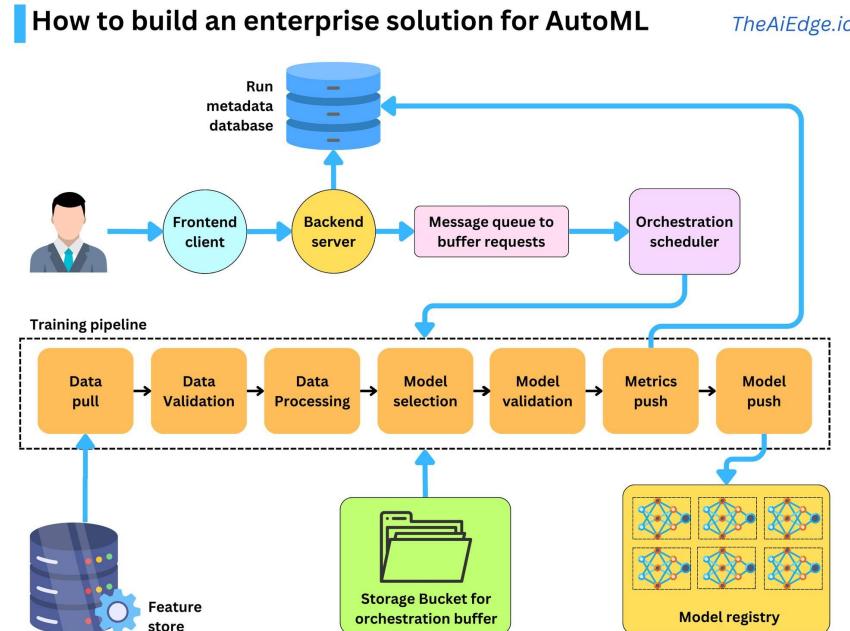
Examining how Uber pickups vary over time in New York City's and at its major regional airports. By sliding the slider on the left you can view different slices of time and explore different transportation trends.

La Guardia Airport JFK Airport Newark Airport

© Botnoi Group

AutoML

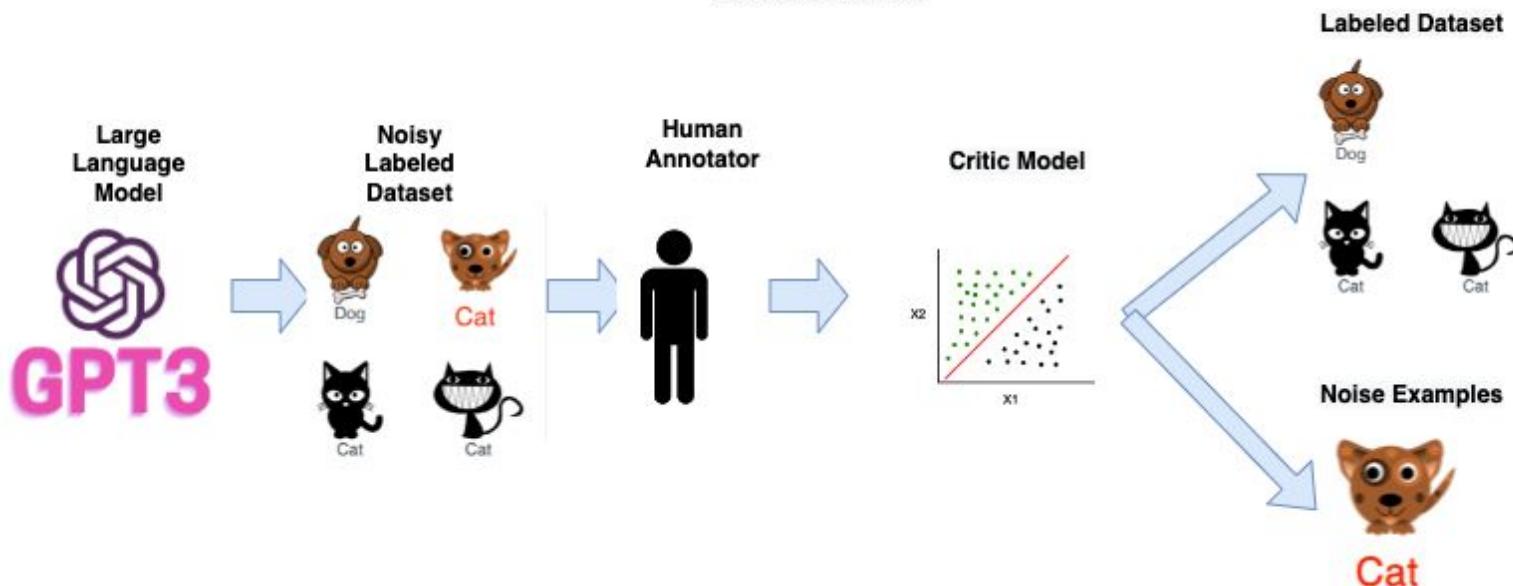
(Advance ML Pipeline)



How LLM Labeling Data

(self-supervised learning)

LM Augmented Data Annotation



Trend of AI Generative

(AI Trend 2023)

