**Artificial Intelligence : Human Intelligence exhibited by machine.**

   **Narrow AI : Computers can do a specific/one thing very well.**

   **General AI : Computers can do multiple things like humans. We are very far away from this.**

**Machine Learning : Approach to try and achieve AI through systems that can find patterns in data. Stanford Univ - Science of getting computers to act without being explicitly programmed.**

**Deep Learning :  One of the techniques to implement machine learning.**

**Data Science : Analysing Data**

# Play Ground

- [https://teachablemachine.withgoogle.com/](https://teachablemachine.withgoogle.com/)

- [https://ml-playground.com/#](https://ml-playground.com/#)

# How did we get here ?

# YouTube Recommendation Engine

- [https://ml-playground.com/#](https://ml-playground.com/#)

- X Axis - Duration of Video

- Y Axis - Likes to the Video

# Framework



Steps in a full machine learning project

What we're going to cover

Data collection

Data modelling

Deployment

What problem are we trying to solve?
1. Problem defintion

What data do we have?
2. Data

What defines success?
3. Evaluation

What features should we model?
4. Features

What kind of model should we use?
5. Modelling

What have we tried/ what else can we try?
6. Experiments

Iterative process

# 1. Problem Definition

# Types of machine learning



Machine Learning

Supervised — Unsupervised — Reinforcement

**Classification** — Is this a cat or dog ?

**Clustering** — Machine Categorizes data

**Regression** — Predicting Stock prices

**Association Rule Learning** — What customer might buy ?

Skill acquisition
Real time learning

Images: https://vas3k.com/blog/machine_learning

# When not to use machine learning ?

- Will simple hand coded instructions based system work ? If yes, then use it.

# Main types of machine learning



**Supervised Learning**

**Unsupervised Learning**

**Transfer Learning**

**Reinforcement Learning**

# Supervised learning

**Classification**

- "Is this example one thing or another?"
- Binary classification = two options
- Multi-class classification = more than two options

**Regression**

- "How much will this house sell for?"
- "How many people will buy this app?"

# Unsupervised learning



| Customer ID | Purchase 1 | Purchase 2 |
|---|---|---|
| 1 | Sunglasses | Singlet |
| 2 | Jacket | Snow boots |
| 3 | Sunscreen | Beach towel |

Cluster 1 (Summer)

Cluster 2 (Winter)

# Transfer Learning

# Reinforcement Learning



Score

Lose  -1

Win  +1

# Problem Definition

## Matching your problem


**Supervised Learning** → "I know my inputs and outputs."


**Unsupervised Learning** → "I'm not sure of the outputs but I have inputs."


**Transfer Learning** → "I think my problem may be similar to something else."

# 2.Data

"What kind of data do we have?"

# Types of Data



Rows

| ID | Weight | Sex | Blood Pressure | Chest pain | Heart disease? |
|------|--------|-----|----------------|------------|----------------|
| 4326 | 110Kg | M | 120/80 | 4 | Yes |
| 5681 | 64Kg | F | 130/90 | 1 | No |
| 7911 | 81Kg | M | 130/80 | 0 | No |

Columns

Table 1.0: Patient records

**Structured**

From: daniel@mrdbourke.com
Hey Daniel,

First of all, thank you for being so amazing.
This machine learning course is incredible.
Thank you for keeping it simple!

**Unstructured**

# Types of Data ....



**Static**

Table 1.0: Patient records
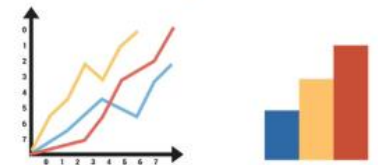


**Streaming**

# A data science workflow



Static data

Table 1.0: Patient records

pandas

$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$

Data Analysis

matplotlib

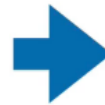Machine learning model

scikit learn

Heart disease?

# 3. Evaluation ✅

"What defines success for us?"

"For this project to be worth pursuing further, we need a machine learning model with over 99% accuracy."

Table 1.0: Patient records

Machine learning model

Heart disease?

Accuracy

97.8%

# Types of metrics

| Classification | Regression | Recommendation |
|---|---|---|
| Accuracy | Mean absolute error (MAE) | Precision at K |
| Precision | Mean squared error (MSE) | |
| Recall | Root mean squared error (RMSE) | |

# Classifying Car insurance claims



Table 2.0: Car insurance claims

(had to try a few of these)

Machine learning model

Minimum accuracy
>95%

# 4. Features

**"What do we already know about the data?"**

**Feature variables can be**
- Numerical
- Categorical

**Feature engineering**
- Deriving new features from existing one.

**Feature Coverage**
- Checking if values are correctly populated for a feature or not ? Do not use it if it is not well covered.

| | Feature variables | | | | Target variable | Derived feature | |
|---|---|---|---|---|---|---|---|
| ID | weight | Sex | Heart Rate | Chest pain | Heart disease? | visit in last year? | Most eaten food |
| 4326 | 110Kg | M | 81 | 4 | Yes | Yes | Fries |
| 5681 | 64Kg | F | 61 | 1 | No | Yes | ? . |
| 7911 | 81Kg | M | 57 | 0 | No | No | ? . |

Table 1.0: Patient records

# 5. Modelling Part 1 — 3 sets

**"Based on our problem and data, what model should we use?"**

# 3 parts to modelling

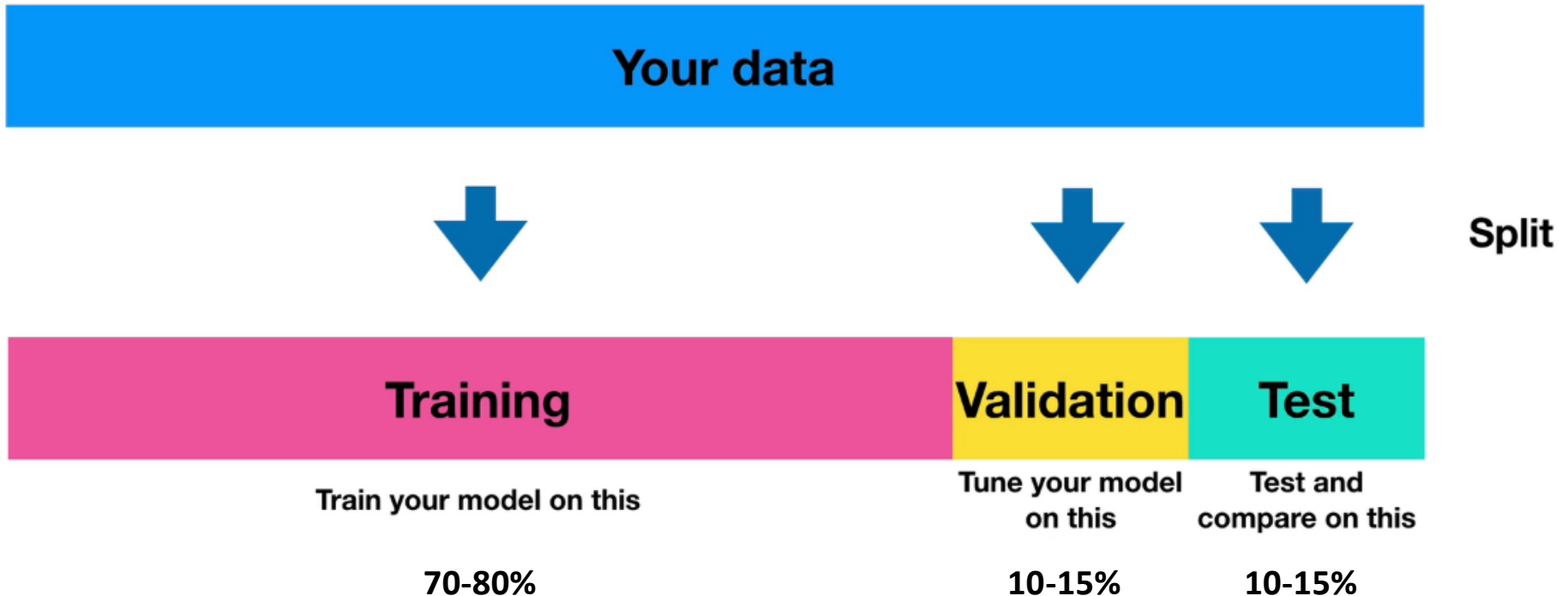1. Choosing and training a model    or

2. Tuning a model

3. Model comparison    vs.    vs.
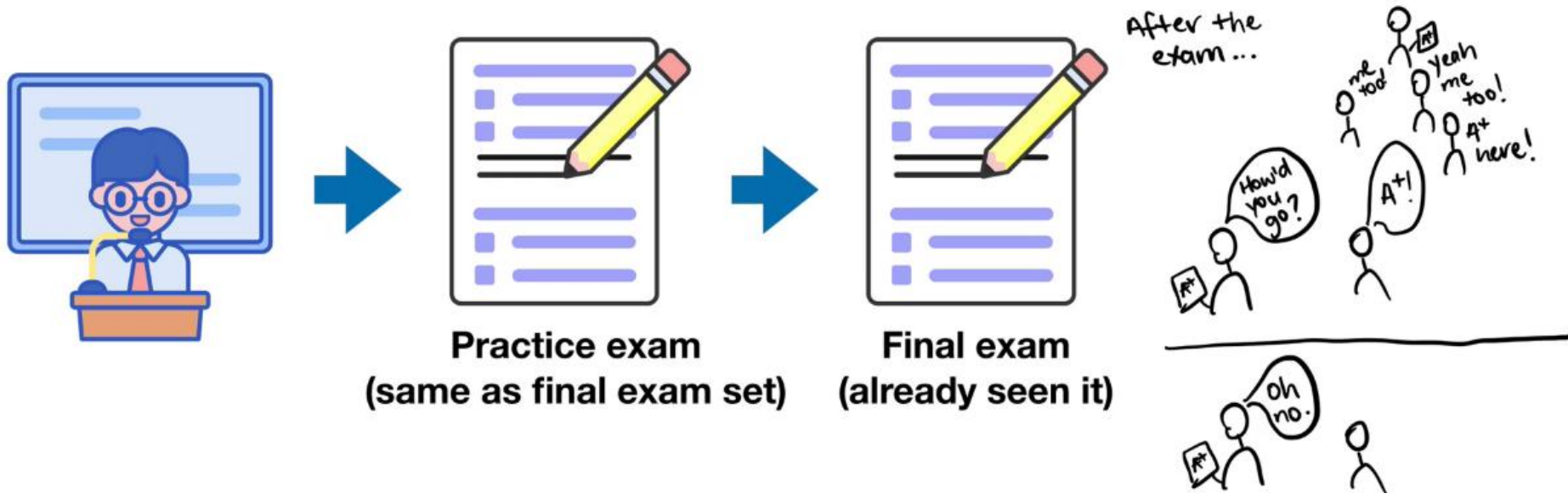
# Training, validation and test sets

# 3 sets



**Generalization** – The ability for a machine learning model to perform well on data it hasn't seen before.

# When things go wrong ?



Practice exam
(same as final exam set)

Final exam
(already seen it)

After the exam...

How'd you go?

A+!

me too

yeah me too!

A+ here!

oh no.

Machine really did not learn anything, it just memorized what it solved in training part.

# 5. Modelling Part 2 — Choosing

**"Based on our problem and data, what model should we use?"**

# Choosing a model

# Training a model



Inputs ➡ 💻 ➡ ⬤ ❌ Outputs

Inputs ➡ 💻 ➡ ⬛ ✅ Outputs

| ID | Weight | Sex | Heart Rate | Chest pain | Heart disease? |
|------|--------|-----|------------|------------|----------------|
| 4328 | 110Kg | M | 81 | 4 | Yes |
| 5681 | 64Kg | F | 61 | 1 | No |
| 7911 | 81Kg | M | 57 | 0 | No |

X (data) — y (label)

Table 1.0: Patient records

**Training Data**

# Goal - Minimize time between experiments



| Experiment | Inputs | Model | Outputs | Accuracy | Training time |
|---|---|---|---|---|---|
| 1 | | Model 1 | | 87.5% | 3 min |
| 2 | | Model 2 | | 91.3% | 92 min |
| 3 | | Model 3 | | 94.7% | 176 min |

**Sometimes for smaller %age extra of Accuracy, we end up spending lot of time. We should avoid that.**

# Remember

- Some models work better than others on different problems
- Don't be afraid to try things
- Start small and build up (add complexity) as you need

# 5. Modelling Part 3 — Tuning

"Based on our problem and data, what model should we use?"

# Tuning



Cooking time: 1 hour
Temperature: 180ºC

Cooking time: 1 hour
Temperature: 200ºC

# Tuning…

**Random Forest**

3 trees

5 trees

**Neural Networks**

2 layers

3 layers

# Remember

- Machine learning models have hyperparameters you can adjust
- A models first results aren't its last
- Tuning can take place on training or validation data sets

# 5. Modelling Part 4 — Comparison

**"How will our model perform in the real world?"**

# Model performance



| Data Set | Performance |
|----------|-------------|
| Training | 98% |
| Test | 96% |

| Underfitting (potential) | Data Set | Performance |
|---|---|---|
| | Training | 64% |
| | Test | 47% |

| Overfitting (potential) | Data Set | Performance |
|---|---|---|
| | Training | 93% |
| | Test | 99% |

# Overfitting and Underfitting



**Underfitting**

**Balanced**
(Goldilocks zone)

**Overfitting**

# Overfitting and Underfitting

# Overfitting and Underfitting

**Training Data**

**Test Data**

Data mismatch

Underfitting

# Fixes for Overfitting and Underfitting



**Underfitting**

- Try a more advanced model
- Increase model hyperparameters
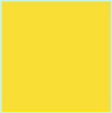- Reduce amount of features
- Train longer

**Overfitting**

- Collect more data
- Try a less advanced model

# Comparison



| Experiment | | | | Accuracy | Training time | Prediction time |
|---|---|---|---|---|---|---|
| 1 | Inputs | Model 1 | Outputs | 87.5% | 3 min | 0.5 sec |
| 2 | Inputs | Model 2 | Outputs | 91.3% | 92 min | 1 sec |
| 3 | Inputs | Model 3 | Outputs | 94.7% | 176 min | 4 sec |

# Remember

- Want to avoid overfitting and underfitting (head towards generality)
- Keep the test set separate at all costs
- Compare apples to apples
- One best performance metric does not equal best model

# 6. Experimentation

"How could we improve/what can we try next?"

# Experimentation

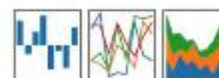- Try out a different approach for improving the machine learning model

# Tools

# Tools mapping