

# INTRODUCTION TO MACHINE LEARNING

Alex de Sá

[a.desa@uq.edu.au](mailto:a.desa@uq.edu.au)  
[@alexgcsa](https://twitter.com/alexgcsa)

**1st Part:** Introduction to ML

**2nd Part:** Practical ML with Python

- Synthetic example (Classification)
- Real-world example (Classification and Regression)

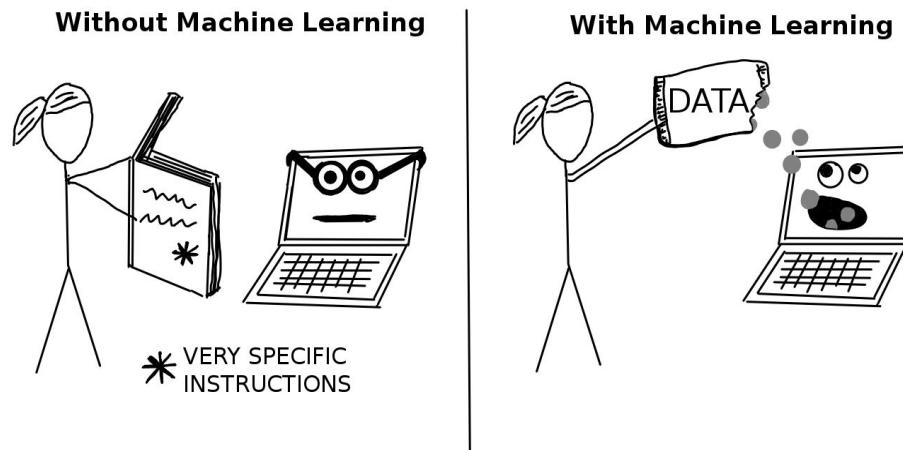
## HIGHLIGHT

You can interrupt me **at ANY time** if something is not clear

# MACHINE LEARNING

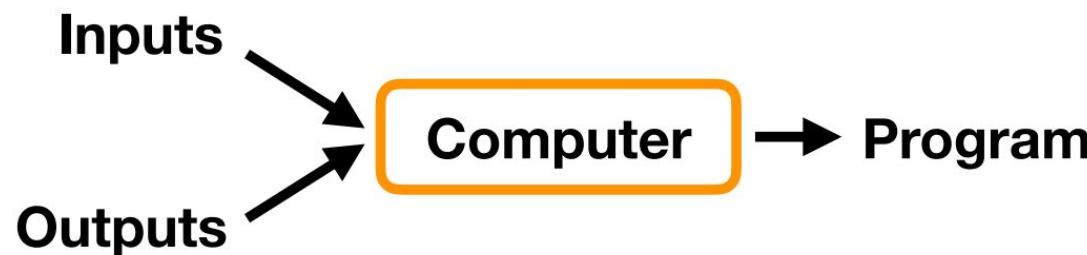
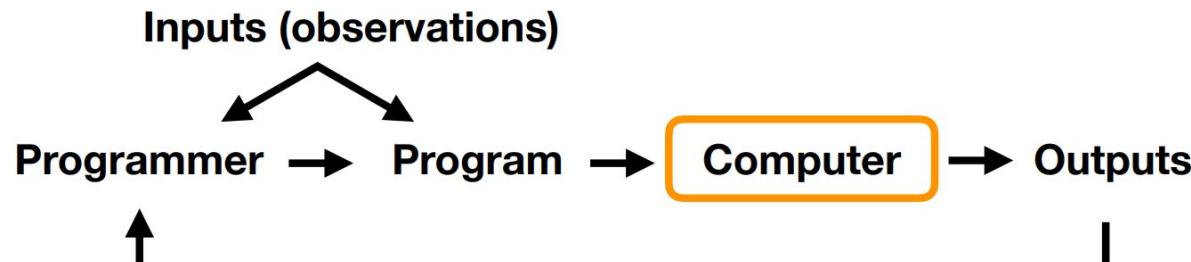
Machine learning is the field of study that gives computers the ability to learn without being explicitly programmed.

**Arthur L. Samuel, AI pioneer, 1959**



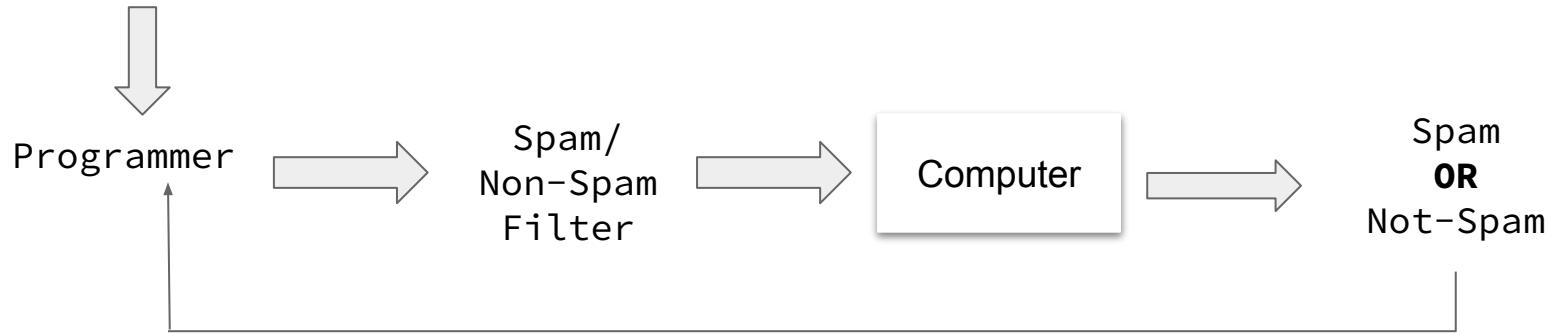
Source: Molnar, 2021

# TRADITIONAL PROGRAMMING VERSUS MACHINE LEARNING

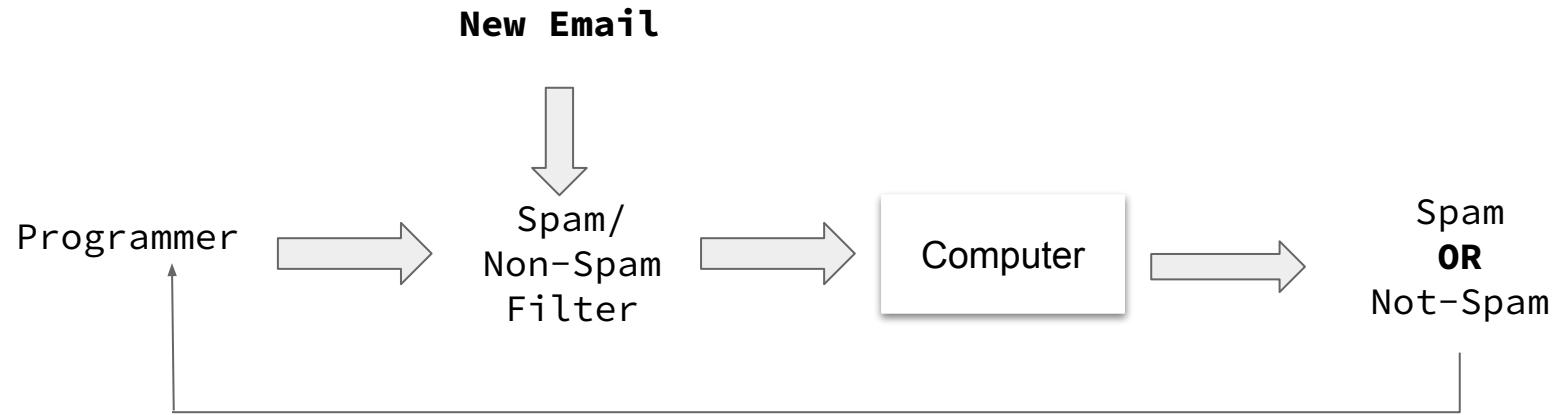


# TRADITIONAL PROGRAMMING

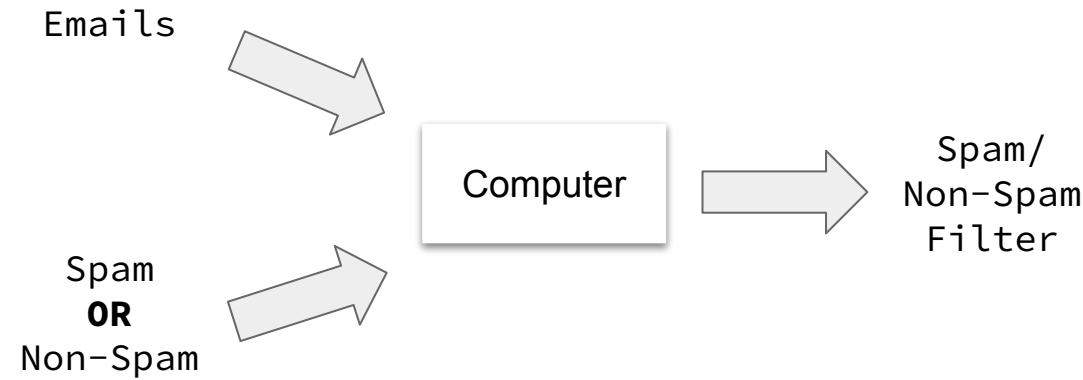
- Emails:** - Spam  
- Non-Spam



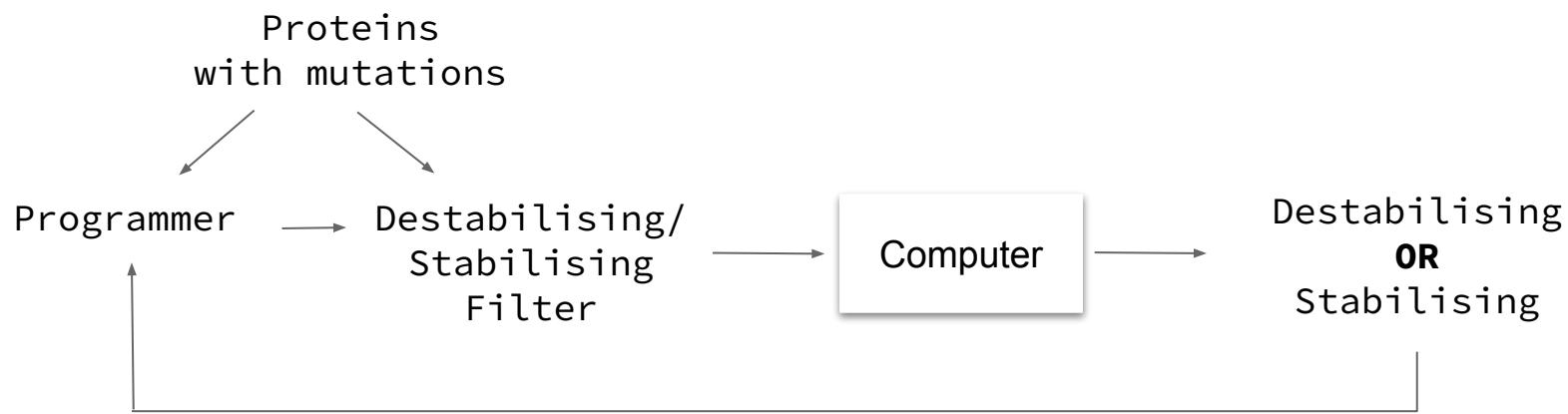
# TRADITIONAL PROGRAMMING



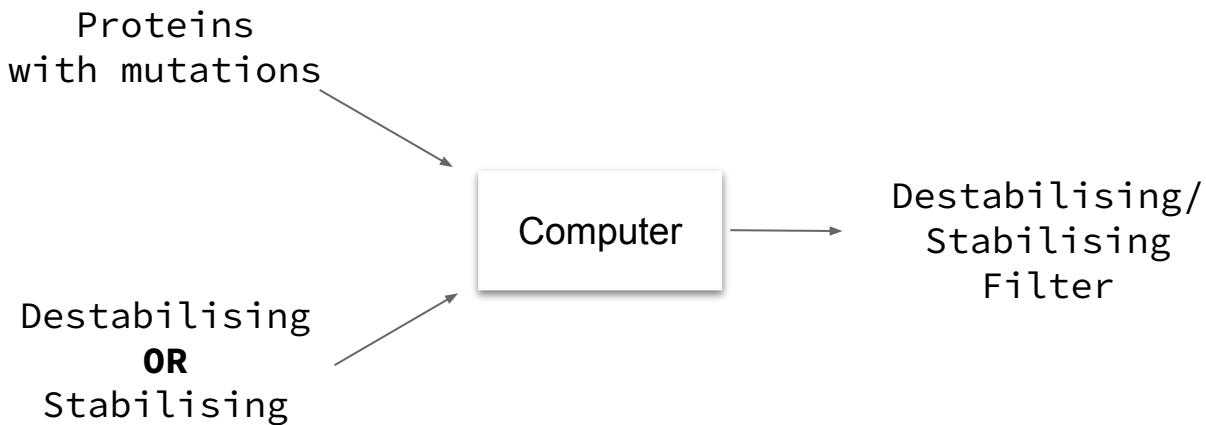
# MACHINE LEARNING



# TRADITIONAL PROGRAMMING



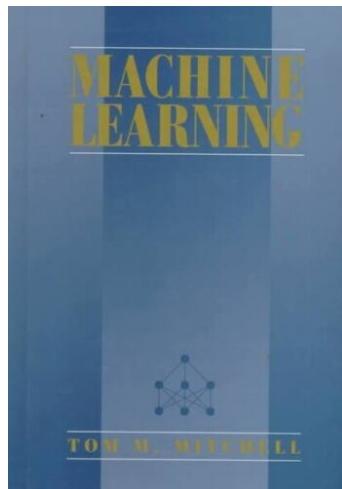
# MACHINE LEARNING



# MACHINE LEARNING DEFINITION

A computer program is said to **learn** from experience **E** with respect to some class of tasks **T** and performance measure **P**, if its performance at tasks in **T**, as measured by **P**, improves with experience **E**.”

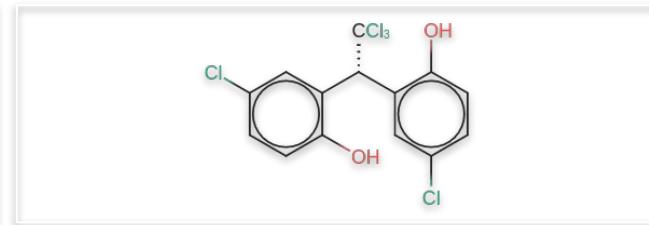
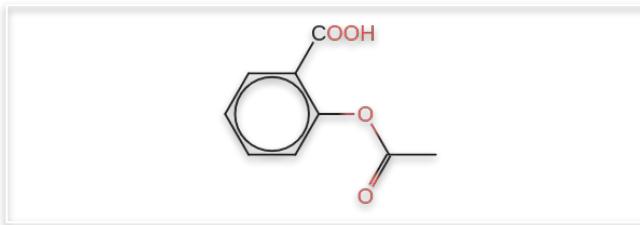
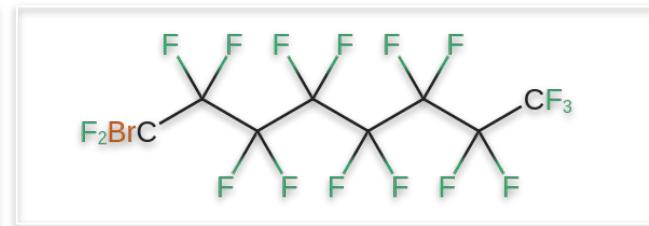
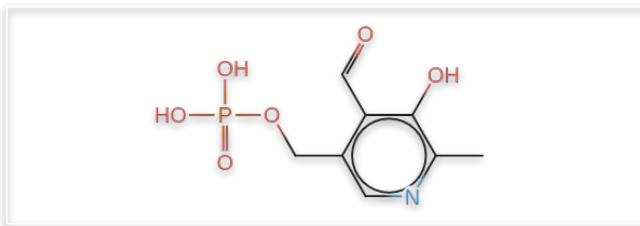
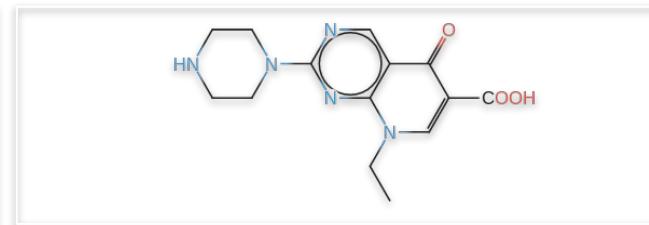
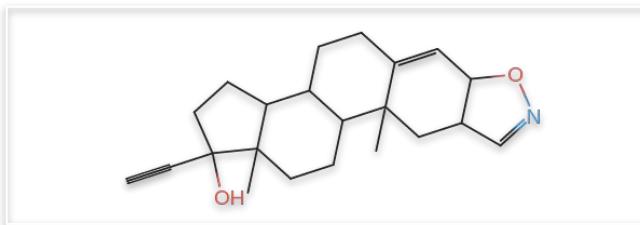
**Tom Mitchell, Professor at Carnegie Mellon University**



<https://www.cs.cmu.edu/~tom/mlbook.html>

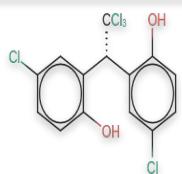
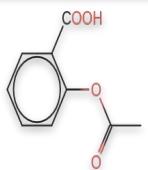
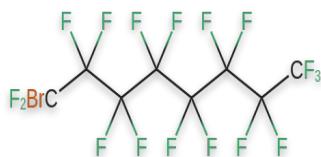
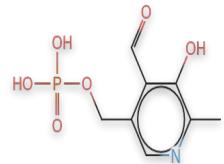
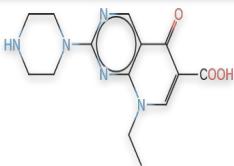
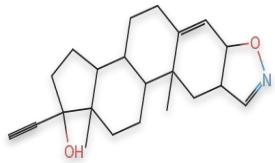
# EXAMPLE - MACHINE LEARNING DEFINITION

## Identification of Hepatotoxicity in Small Molecules



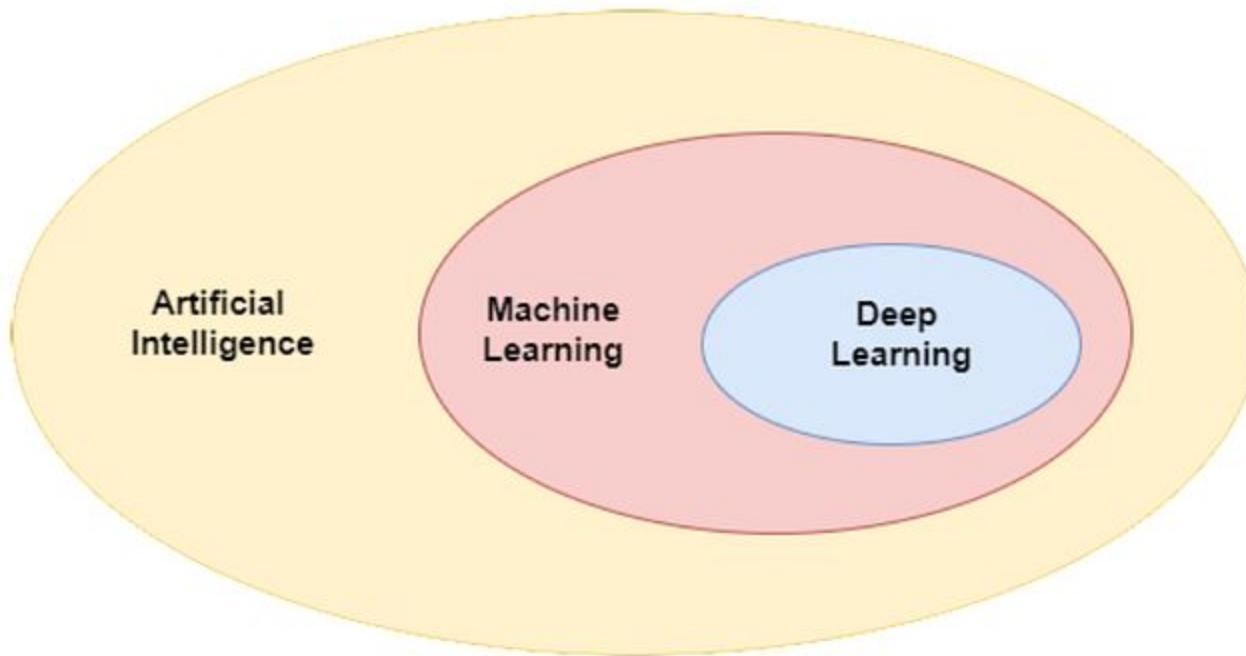
# EXAMPLE - MACHINE LEARNING DEFINITION

## Identification of Hepatotoxicity in Small Molecules

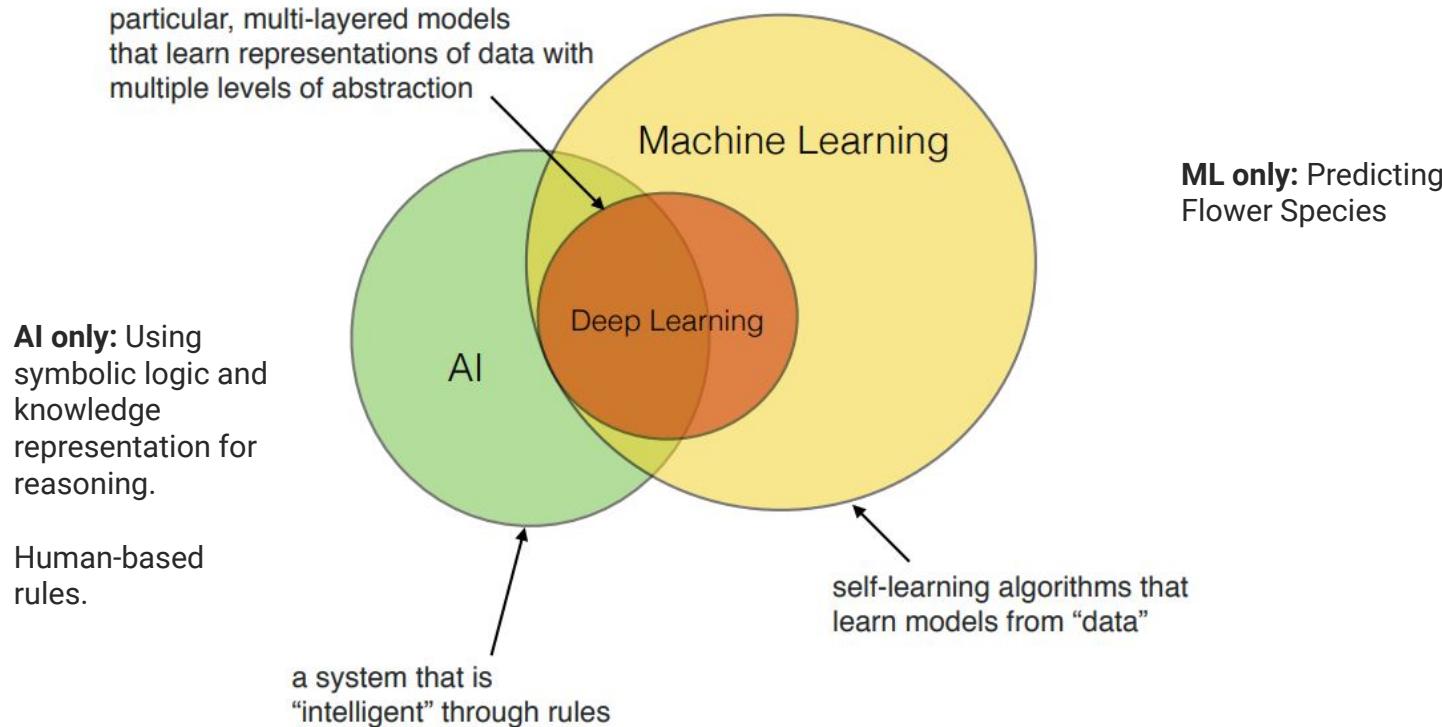


- **Task:** predicting hepatotoxicity from small molecules.
- **Performance:** percentage of molecules classified correctly as toxic.
- **Experience:** dataset of small molecules experimentally distinguishing them between toxic and non-toxic for the liver.

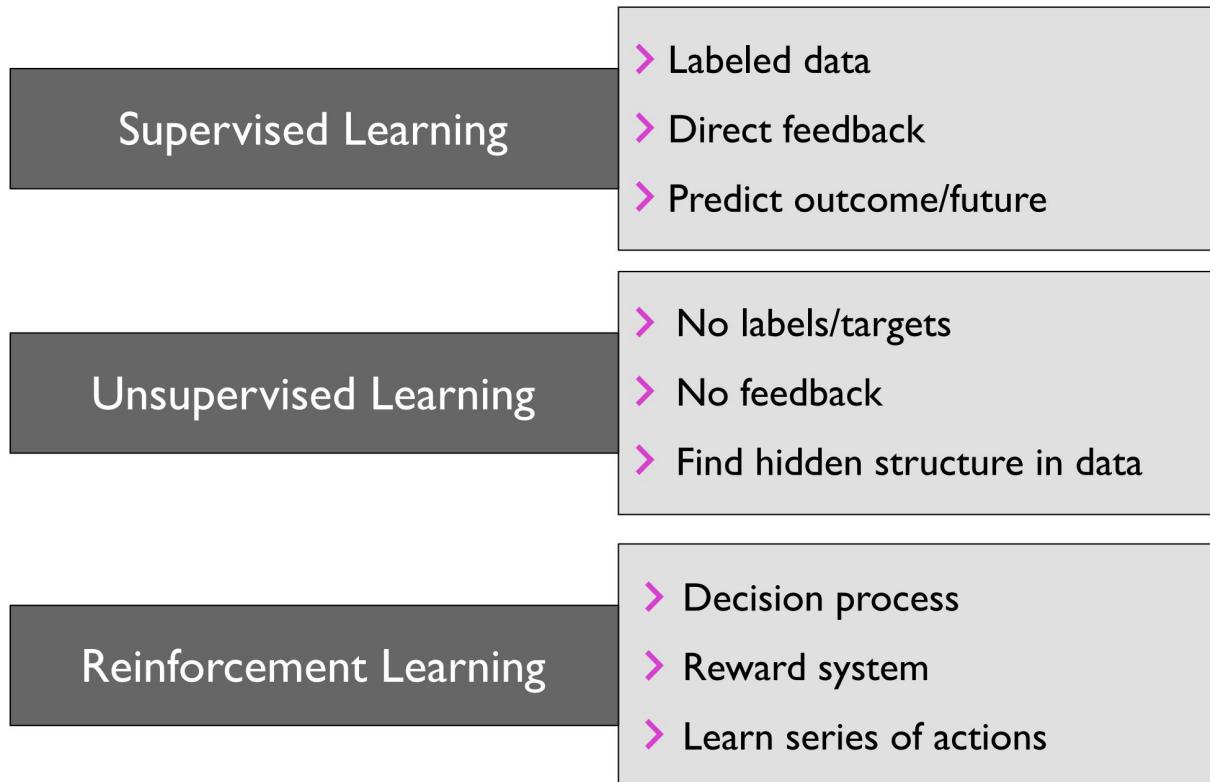
# TAXONOMY OF MACHINE LEARNING



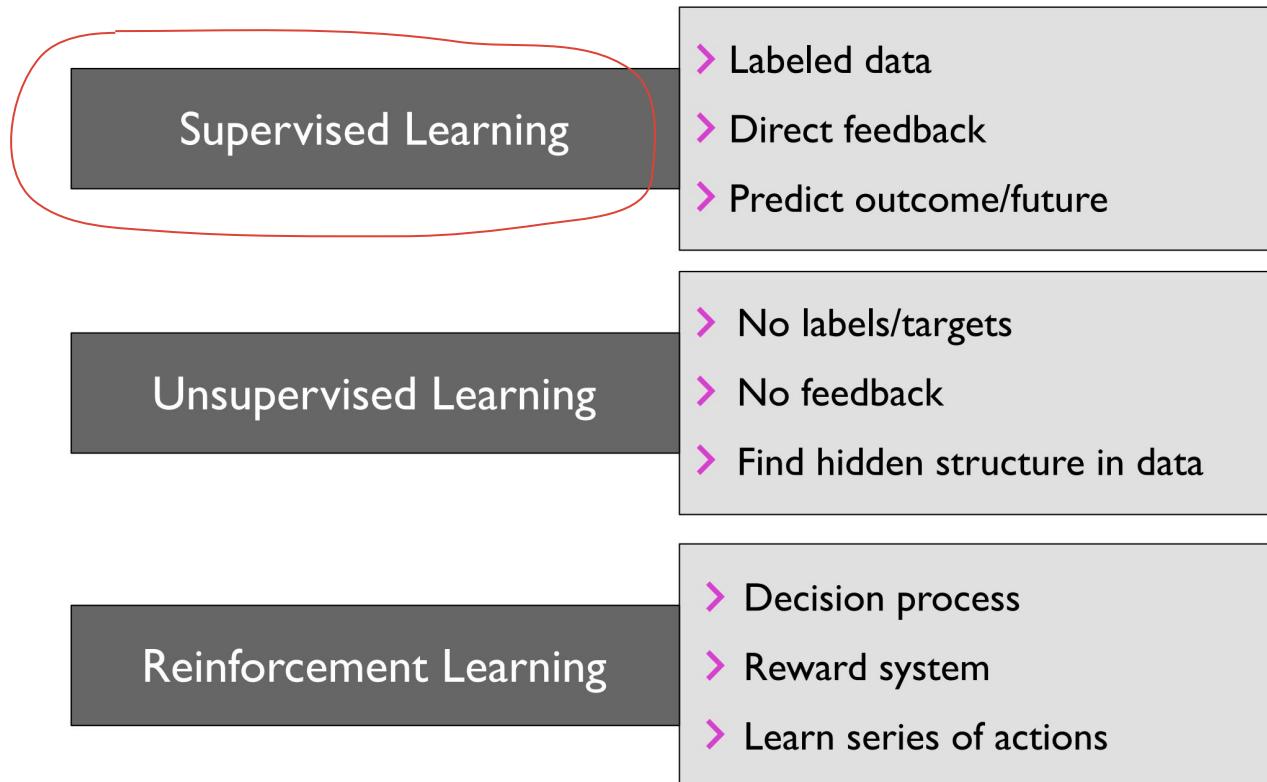
# TAXONOMY OF MACHINE LEARNING



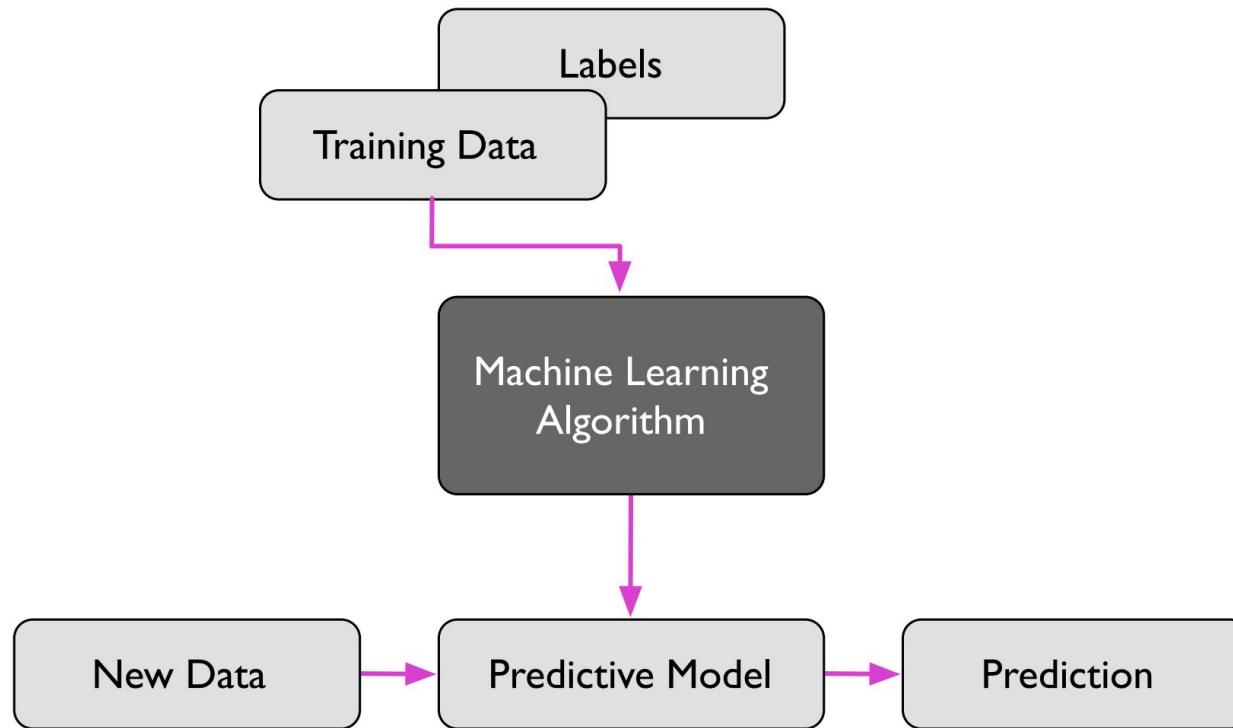
# TAXONOMY OF MACHINE LEARNING



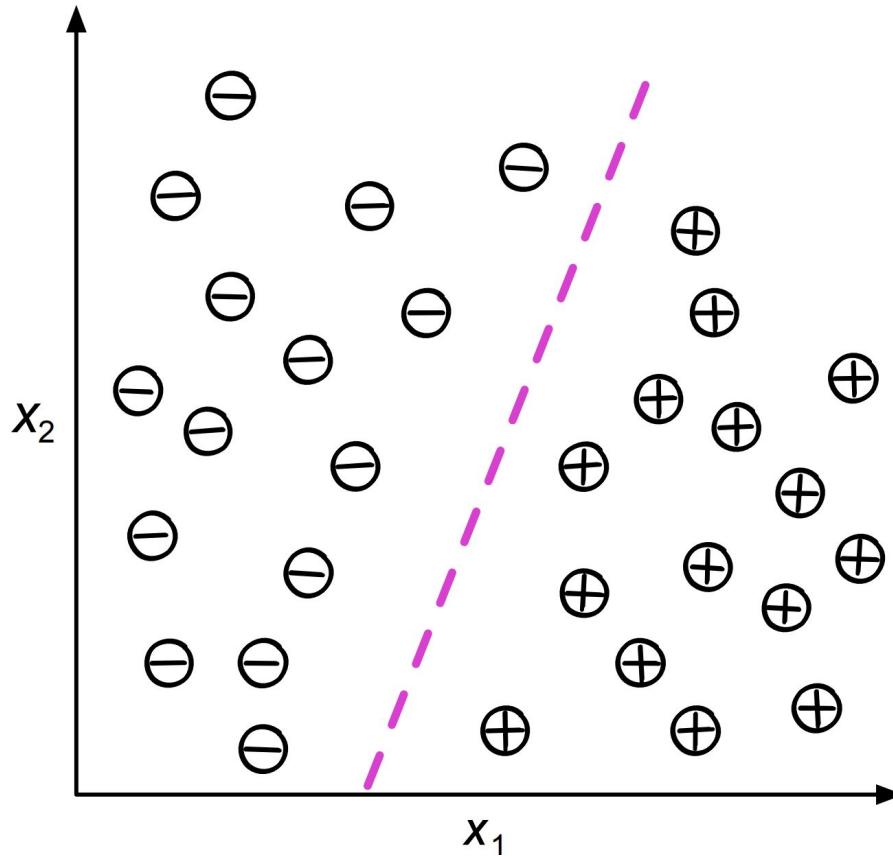
# TAXONOMY OF MACHINE LEARNING



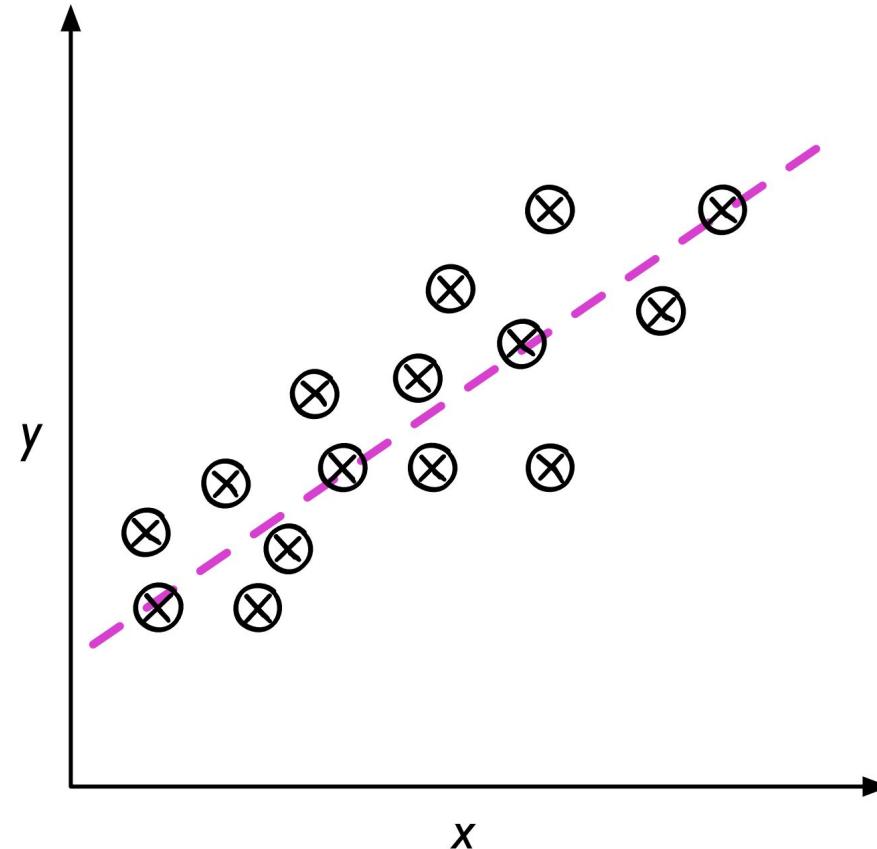
# SUPERVISED LEARNING



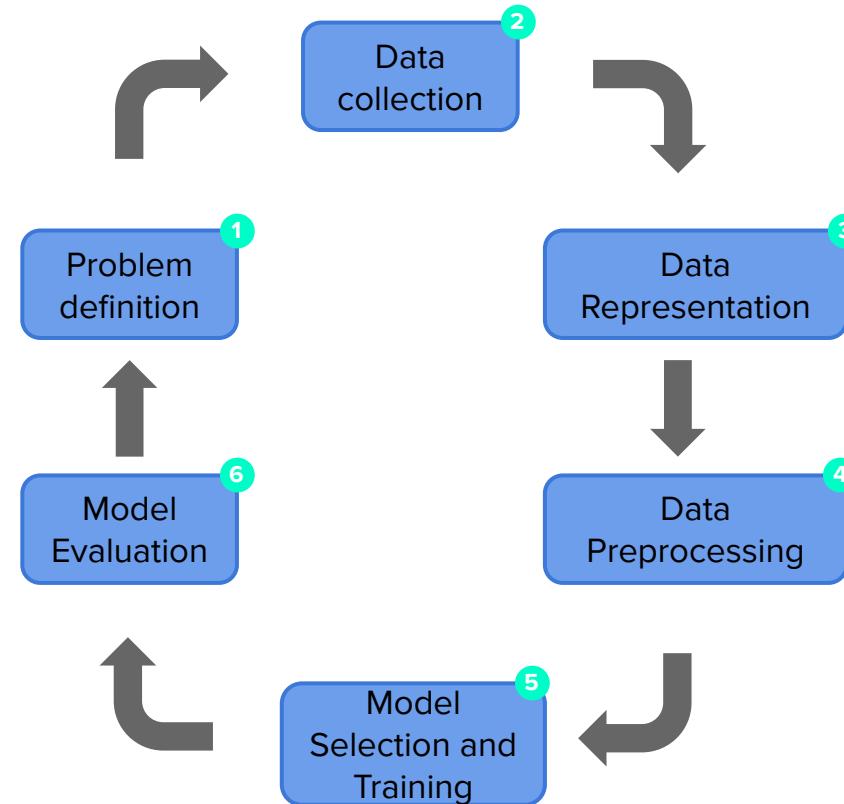
## SUPERVISED LEARNING - CLASSIFICATION



# SUPERVISED LEARNING - REGRESSION



# MACHINE LEARNING WORKFLOW

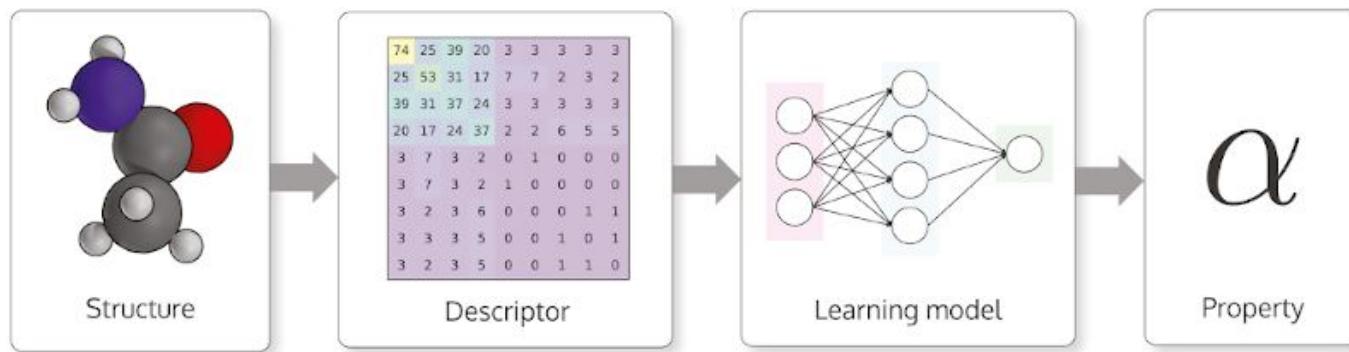


## DATA REPRESENTATION

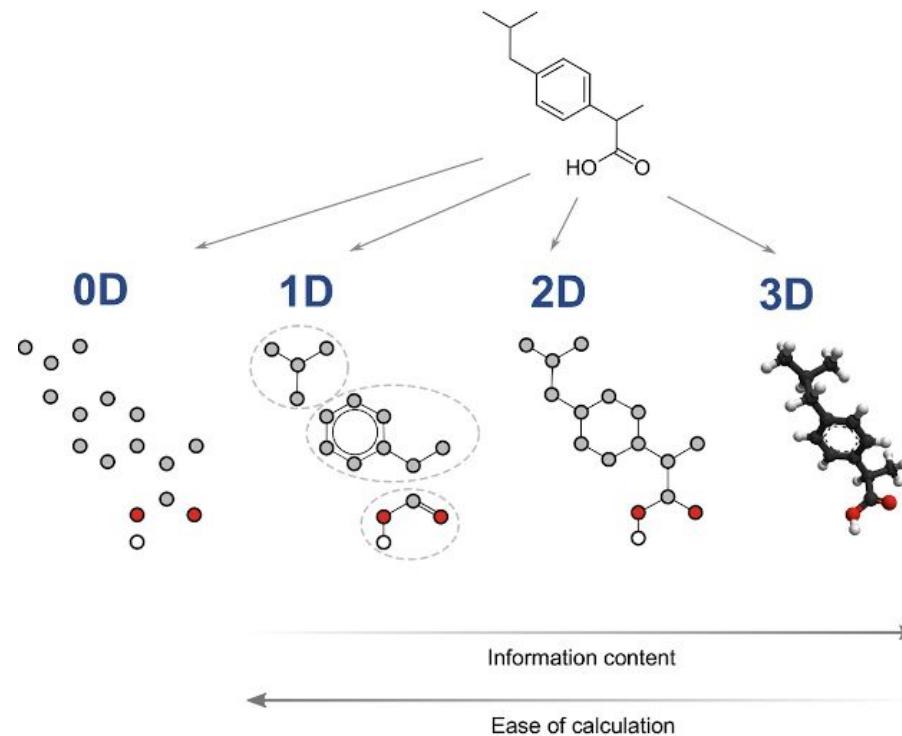
- Most of the recent Machine Learning tools (e.g., scikit-learn, etc) only accept numerical matrices (or dataframes) as inputs.
- We need to find ways to represent our biological, chemical, human data, etc in a numerical way.

## DATA REPRESENTATION - SMALL MOLECULES

Given the structure of the molecule, we are able to derive a list of descriptors aiming to characterise this input molecule to predict a given property

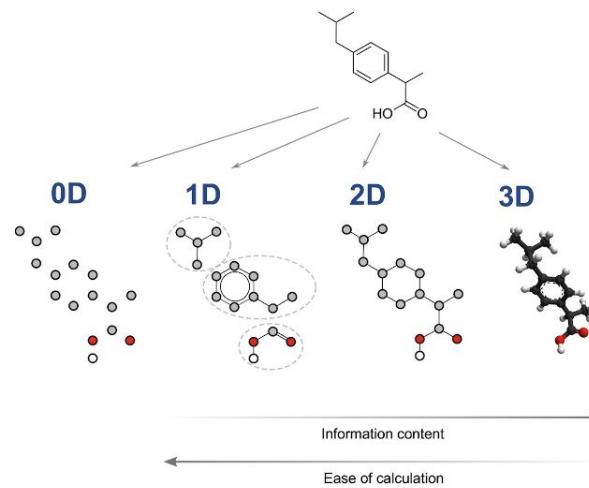


# DATA REPRESENTATION - SMALL MOLECULES

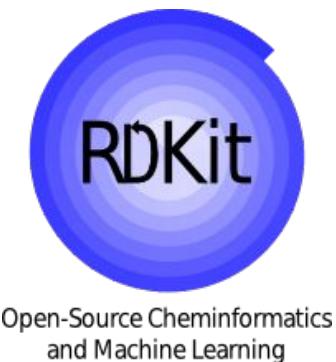


# DATA REPRESENTATION - SMALL MOLECULES

- **0D Descriptors:** No information about structure and connectivity.
  - Atom counts, or molecular weights
- **1D Descriptors:** Partial information about the structure and connectivity.
  - Fingerprints.
- **2D Descriptors:** Information on molecular topology based on the graph representation
  - Atom distance matrix.
- **3D Descriptors:** Information about the spatial coordinates of atoms of a molecule
  - 3D fingerprints.



[Chem Intelligence, 2021](#)



# DATA REPRESENTATION - PROTEINS

## Sequence-based descriptors

MAALSGGGGGAEPGQALFNGDMEPEAGAGAGAAASSAADPAIPEEVWNIKQMIKLTQEH  
IEALLDKFGGEHNPPSIYLEAYEEYTSKLDALQQREQQLLESLGNGTDFSVSSSASMDTV  
TSSSSSSLSVPSSLSVFQNPTDVARSNPKPQKPIVRVFLPNKQRTVVPARCGVTVRDS  
LKKALMMRGLIPECCAVYRIQDGKPKIGWDTDISWLTGEELHVEVLENVPLTHNFVRK



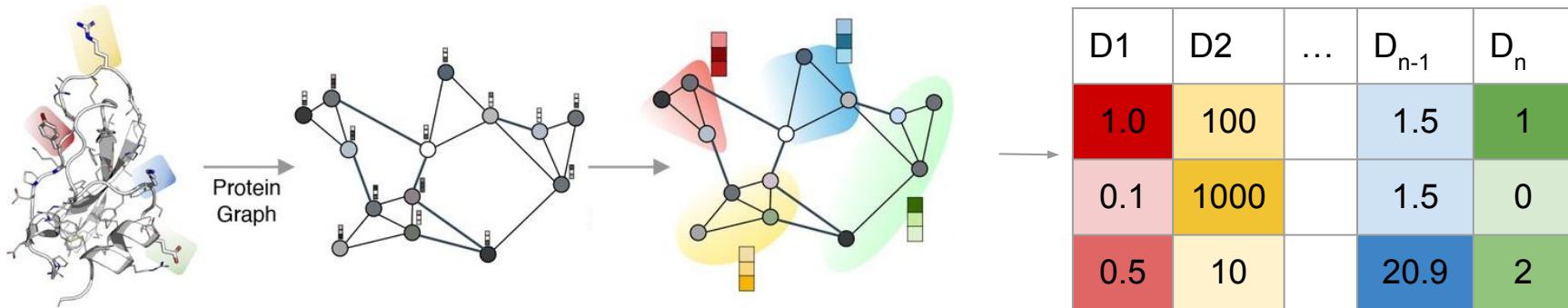
- Amino acid composition
- Physicochemical properties
- Disorder propensity scores
- ...



74	25	39	20	3	3	3	3	3	3
25	53	31	17	7	7	2	3	2	2
39	31	37	24	3	3	3	3	3	3
20	17	24	37	2	2	6	5	5	5
3	7	3	2	0	1	0	0	0	0
3	7	3	2	1	0	0	0	0	0
3	2	3	6	0	0	0	1	1	1
3	3	3	5	0	0	1	0	1	1
3	2	3	5	0	0	1	1	1	0

# DATA REPRESENTATION - PROTEINS

- Structure-based descriptors:



- residue depth
- solvent accessible surface area
- secondary structure distribution
- torsion angles
- Cumulative pair distances between pharmacophore groups

## DATA REPRESENTATION - PROTEINS

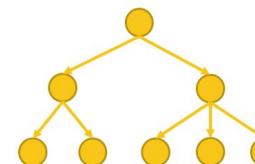
- Structure-based descriptors:

Protein ID	Feature 1	Feature 2	Feature 3	...	Feature M
Q5I0E9	0	5	0		0.150
O74717	10	3	1		0.164
P63033	5	2	1		0.457
...	7	4	1		0.973
Q8NIH1	0	12	0		0.667
D4APA9	14	3	1		0.011
O94864	2	9	1		0.199

Several **Machine Learning (ML)** models with different characteristics:

- Deep learning models
- **Tree-based Models:**
  - Decision Trees
  - Random Forest
  - Extra Trees
  - eXtreme Gradient Boosting
  - Adaptive Boosting
  - Gradient Boosting

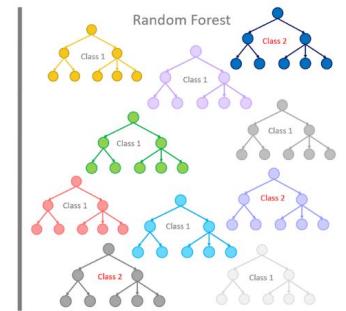
Single Decision Tree



Gradient Boosted Trees

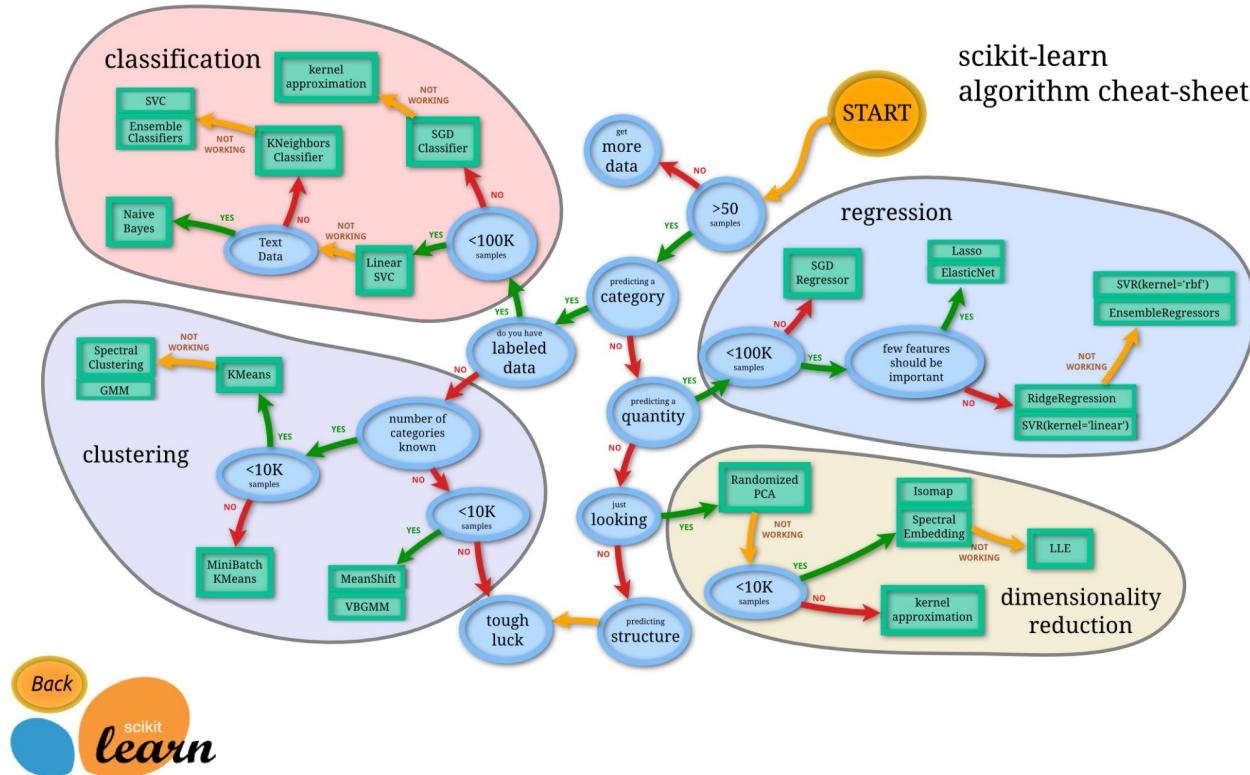


Random Forest

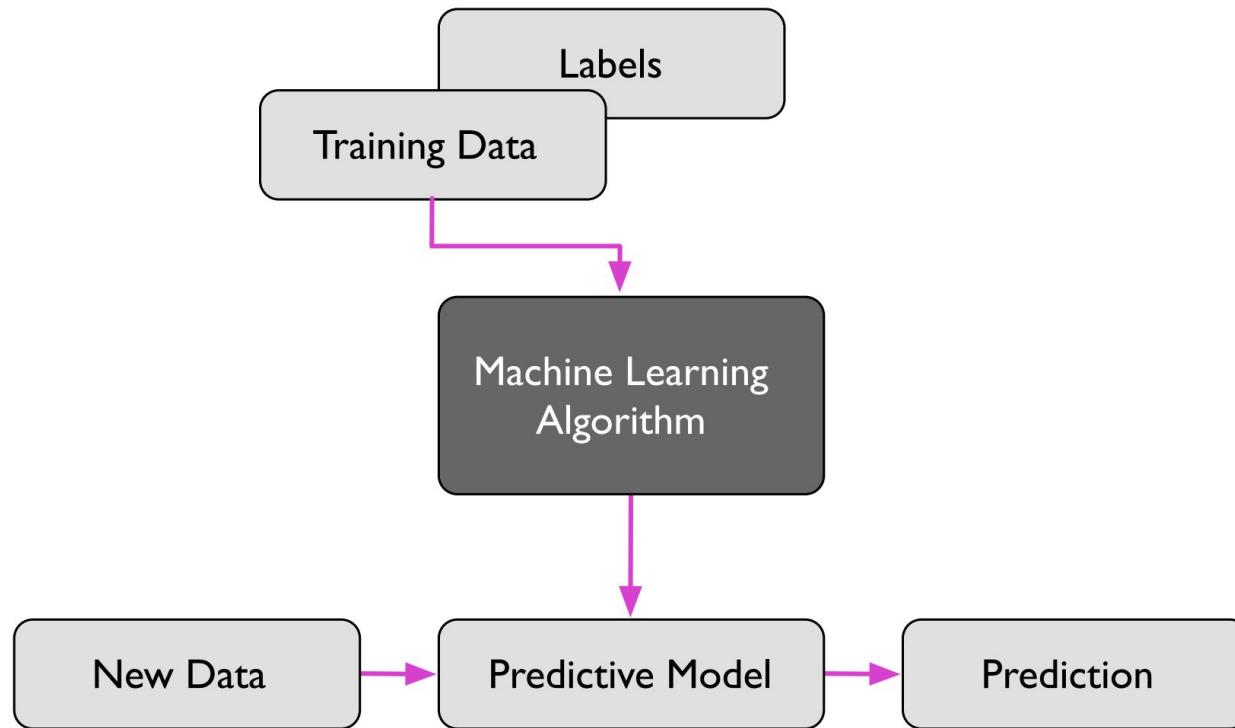


- **Support Vector Machines**
- **Naïve Bayes**
- ....

# ML ALGORITHMS



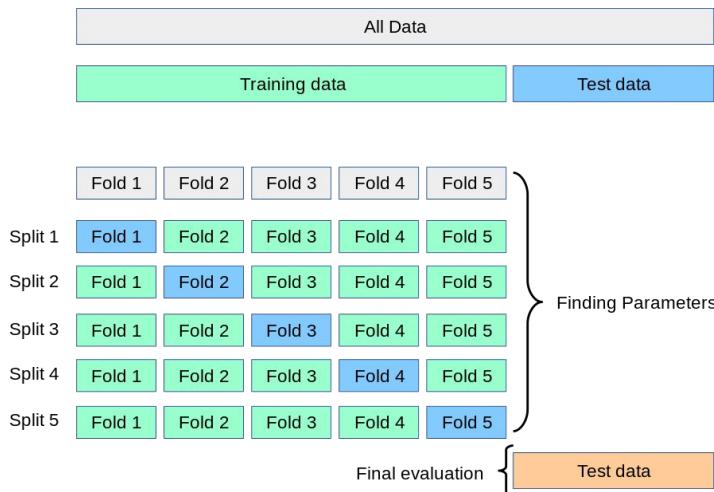
# MODEL FITTING



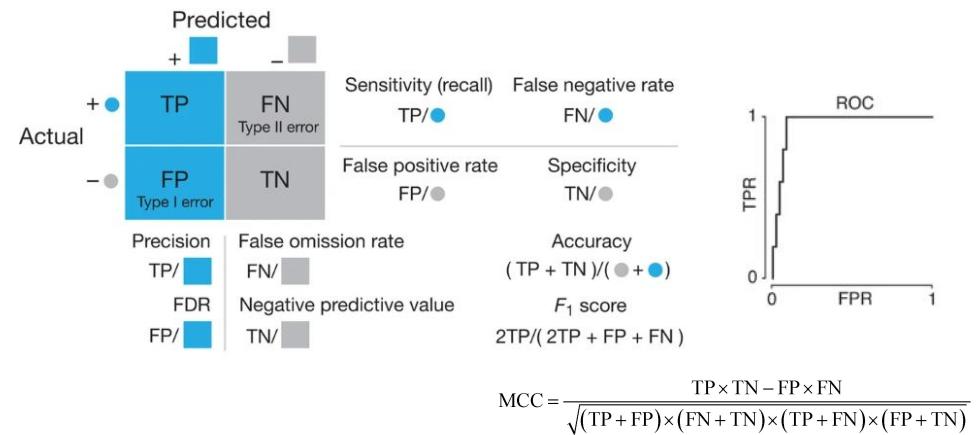
# SELECTION AND EVALUATION OF ML MODELS

## ML model's selection

### Cross-Validation and Testing



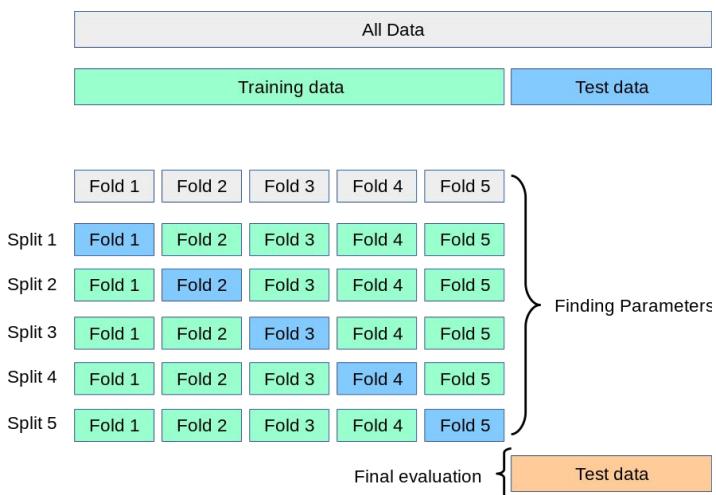
## Performance Evaluation for Classification Models



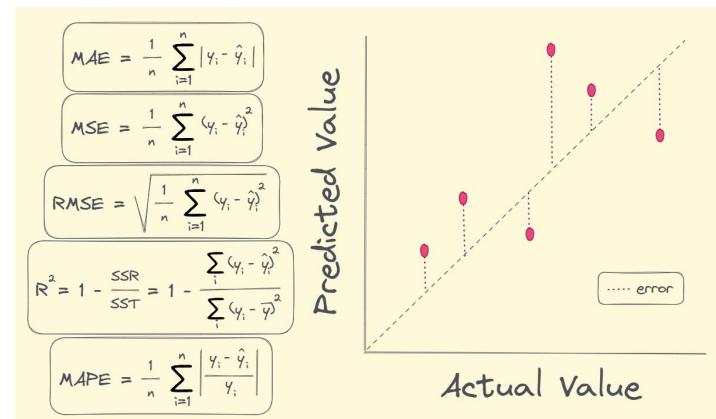
# SELECTION AND EVALUATION OF ML MODELS

## ML model's selection

### Cross-Validation and Testing



## Performance Evaluation for Regression Models

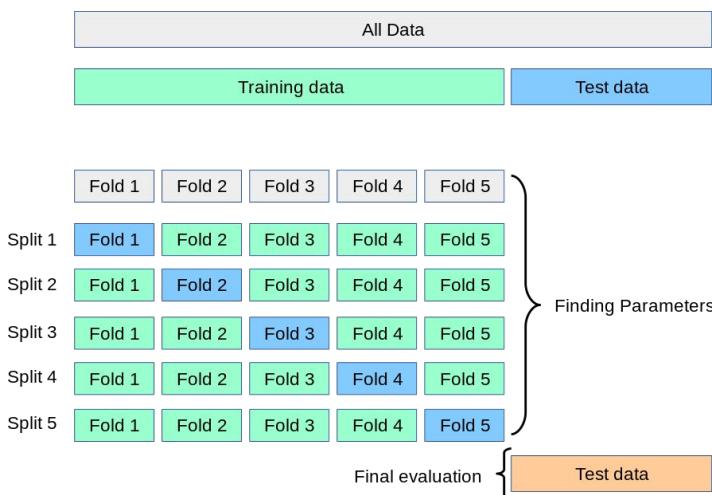


Dorfer, 2023

# SELECTION AND EVALUATION OF ML MODELS

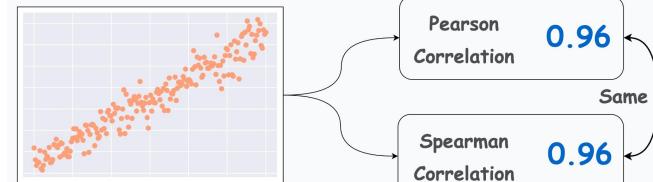
## ML model's selection

### Cross-Validation and Testing

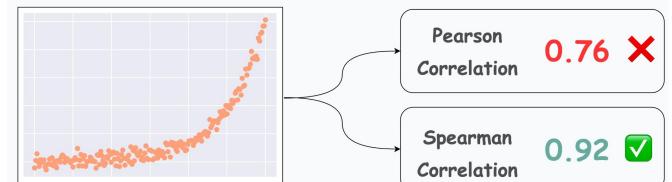


## Performance Evaluation for Regression Models

### Linear Data

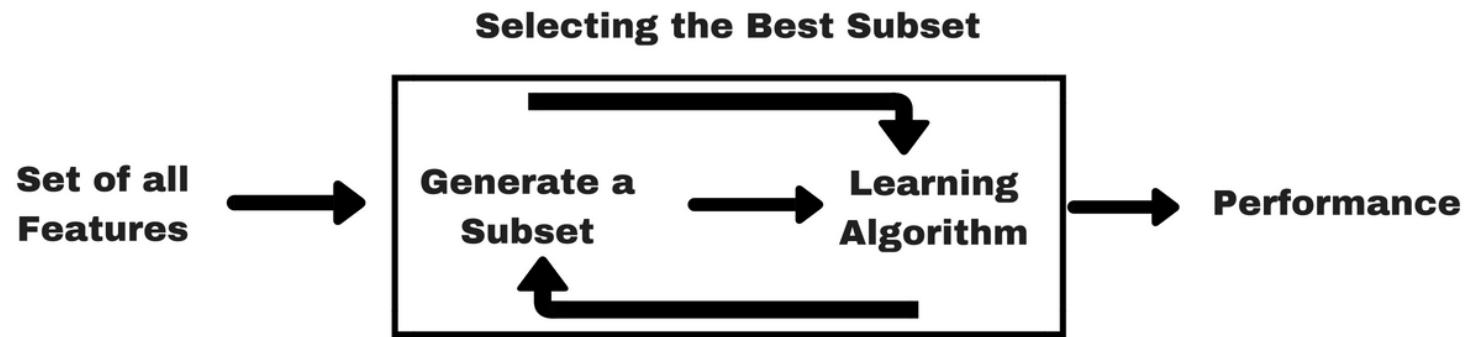


### Non-linear Data



[Chawla, 2023](#)

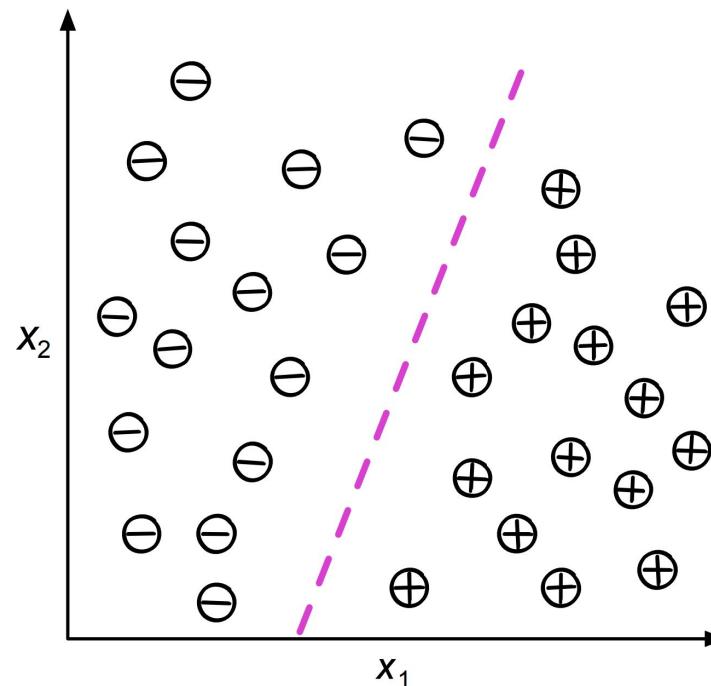
**Machine Learning (ML) coupled with feature selection:**



# Focusing on Classification

## DECISION BOUNDARY

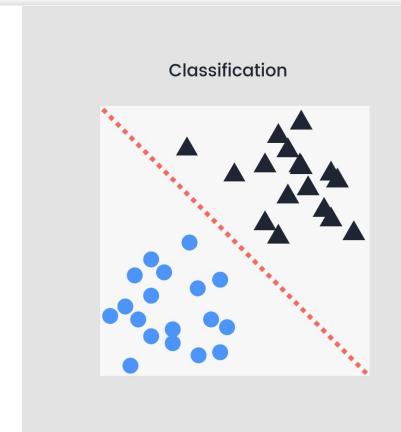
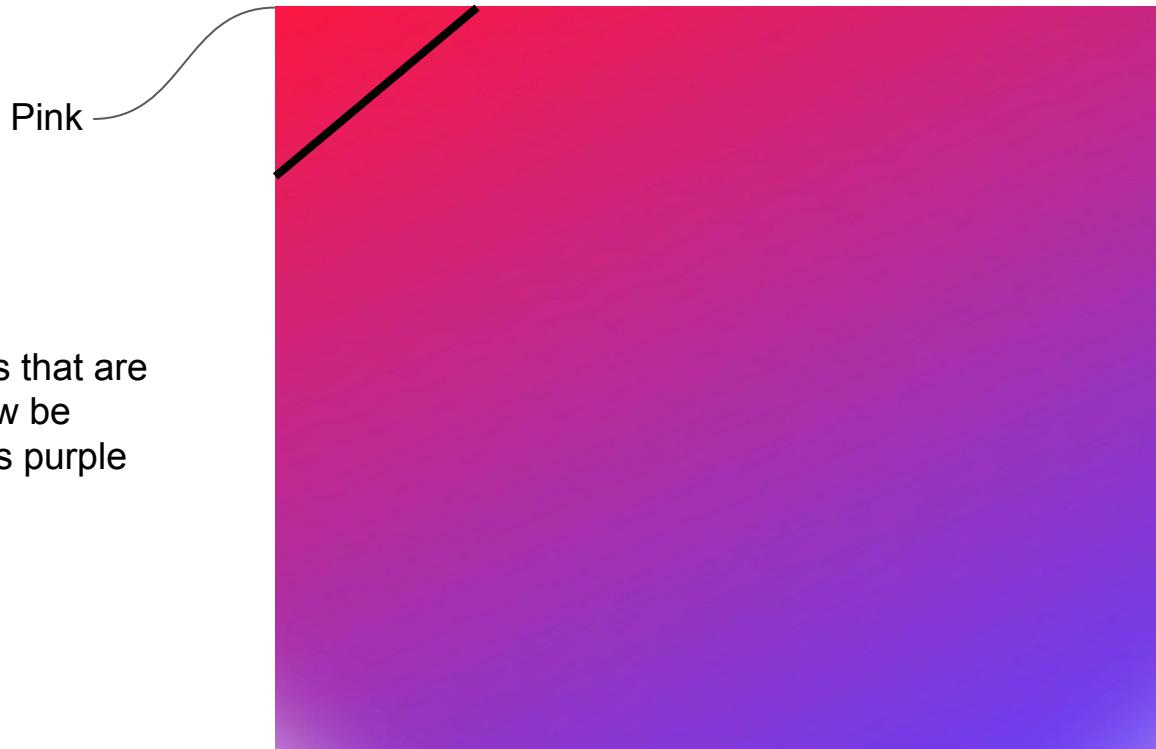
**Definition:** A decision boundary, is a surface that separates data points belonging to different class labels. ([Sahu, 2021](#))



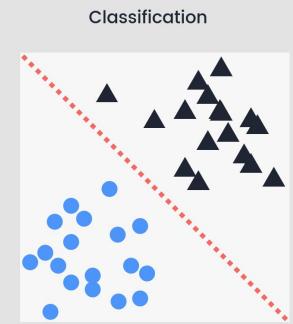
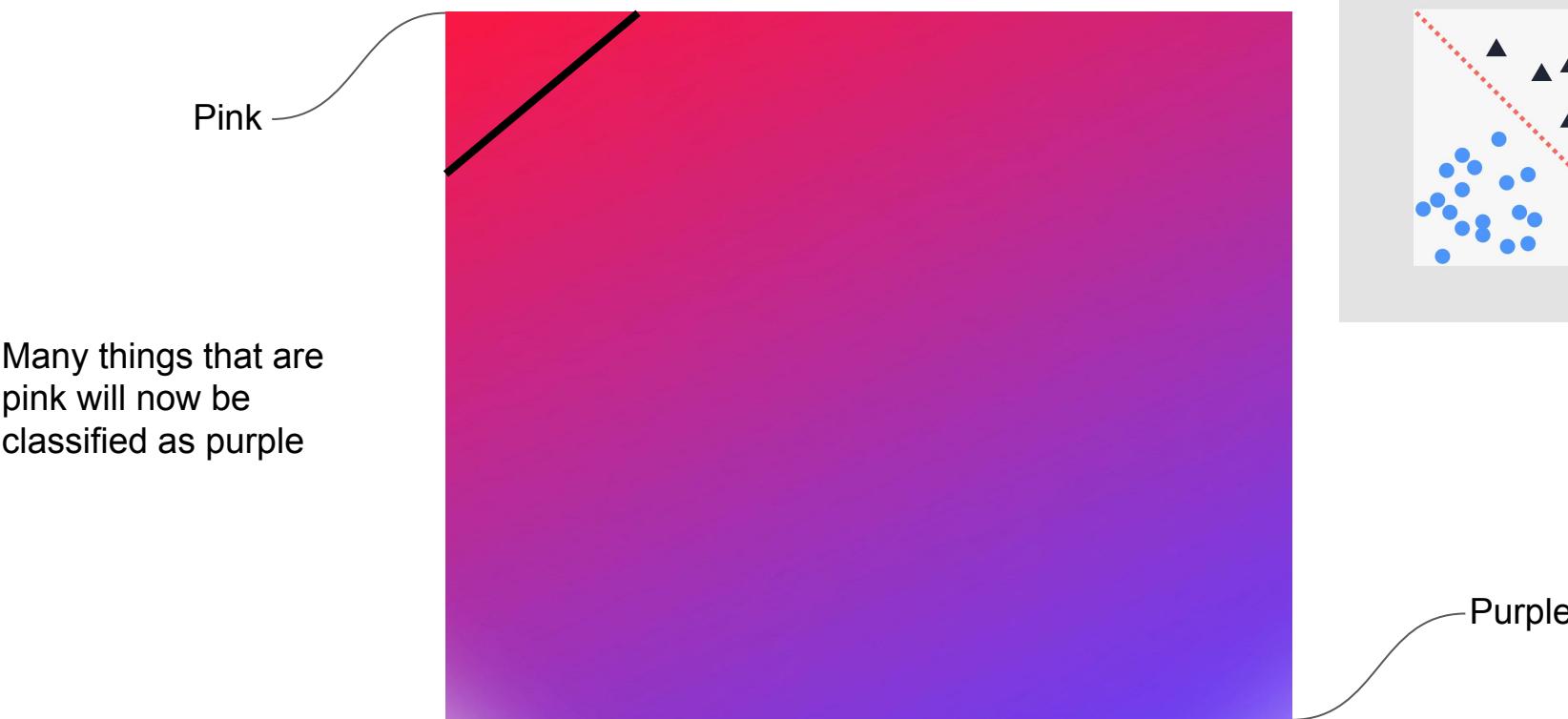
# DECISION BOUNDARY



# DECISION BOUNDARY



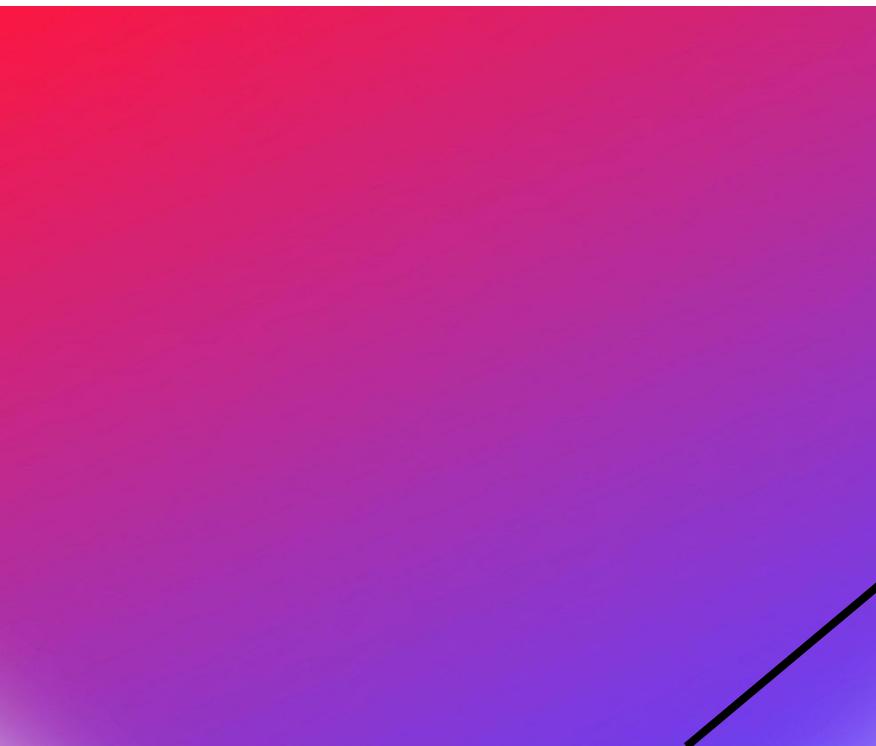
# DECISION BOUNDARY



# DECISION BOUNDARY

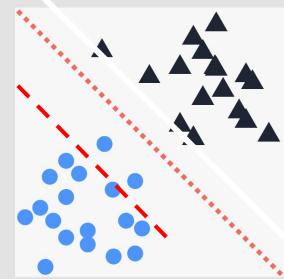
Pink

Many things that are  
purple will now be  
classified as pink

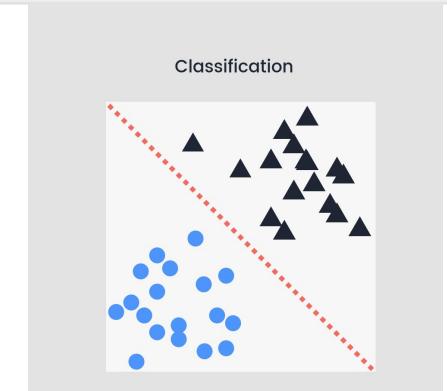
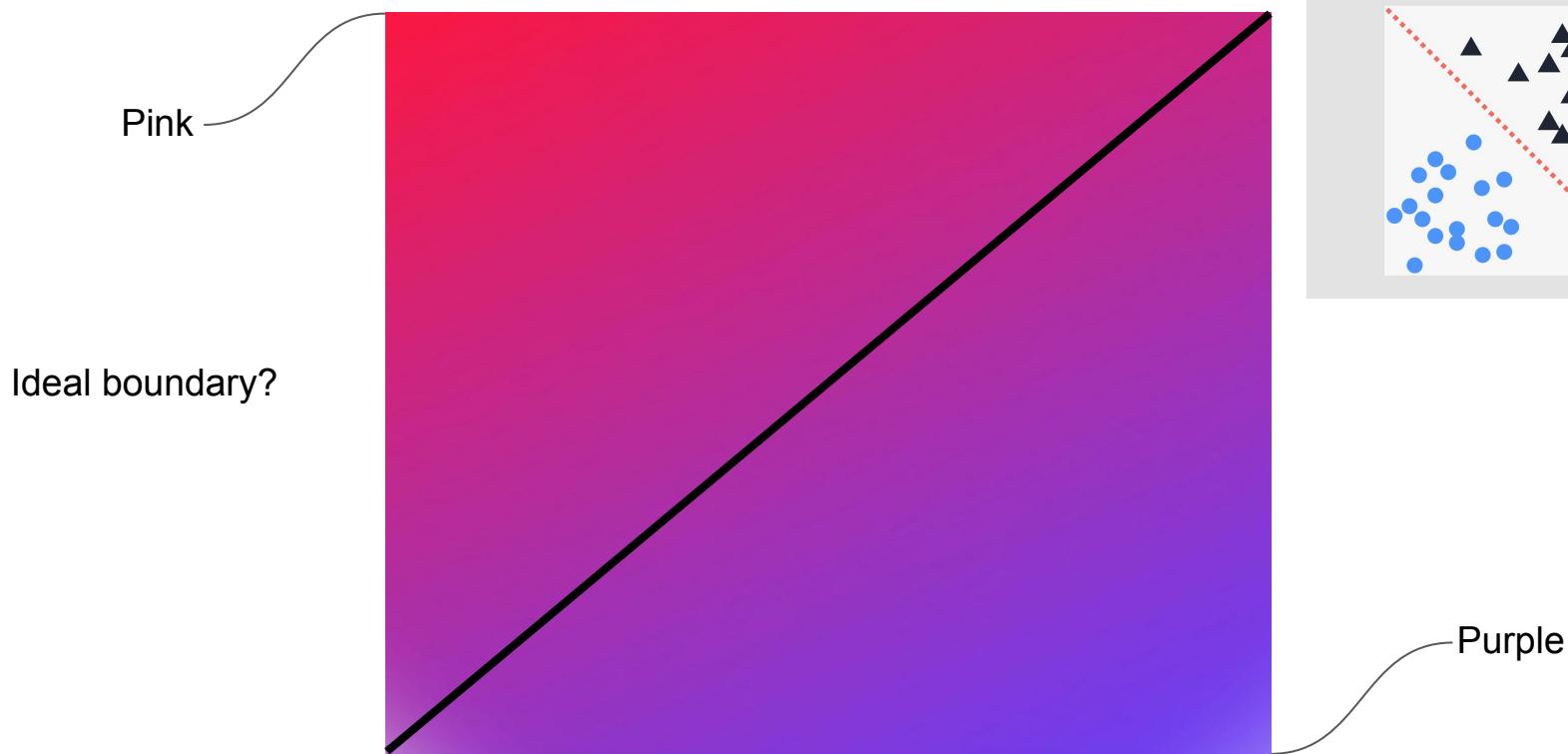


Purple

Classification

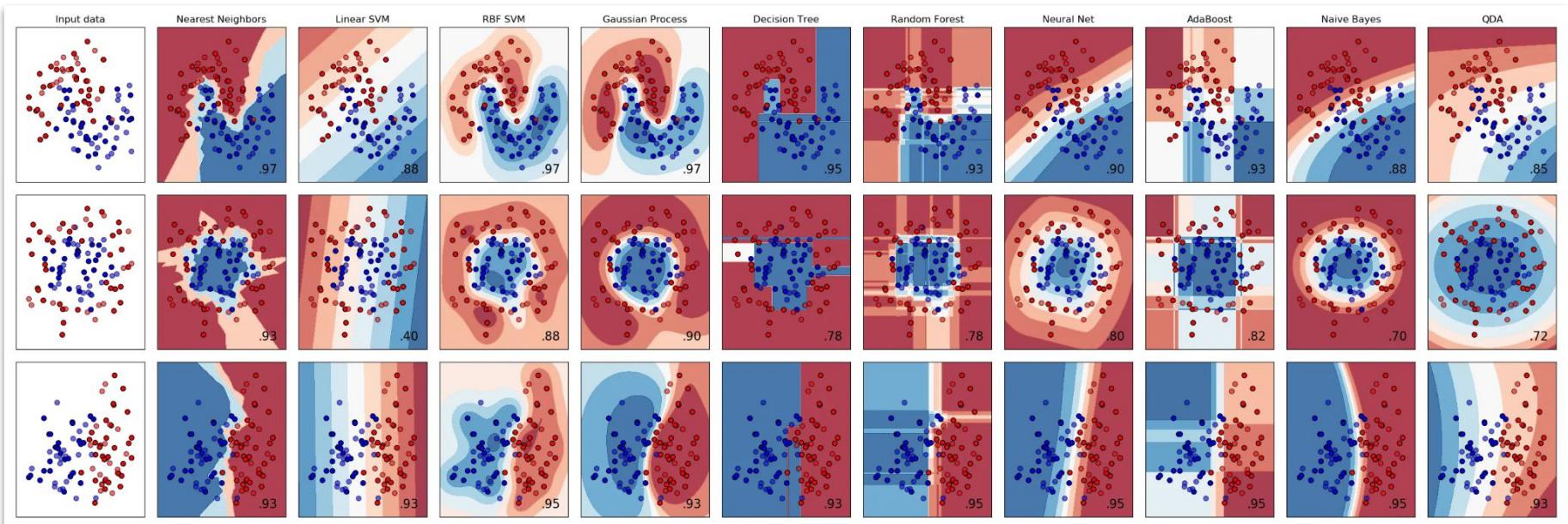


# DECISION BOUNDARY



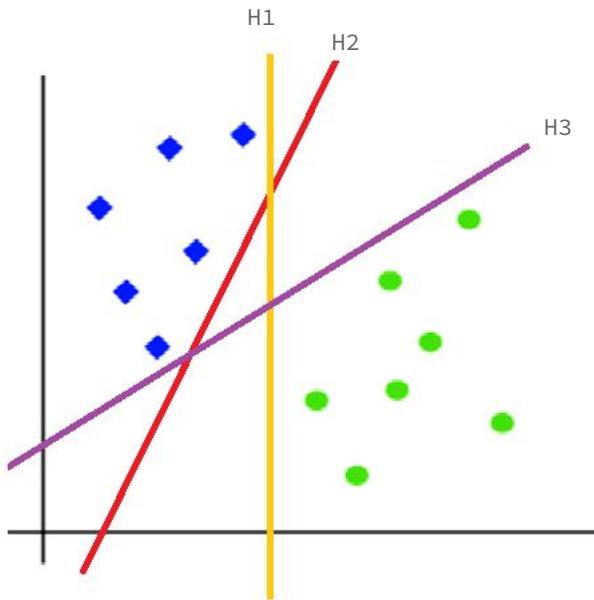
# DECISION BOUNDARY

Comparison of the decision boundaries of 10 machine learning models:

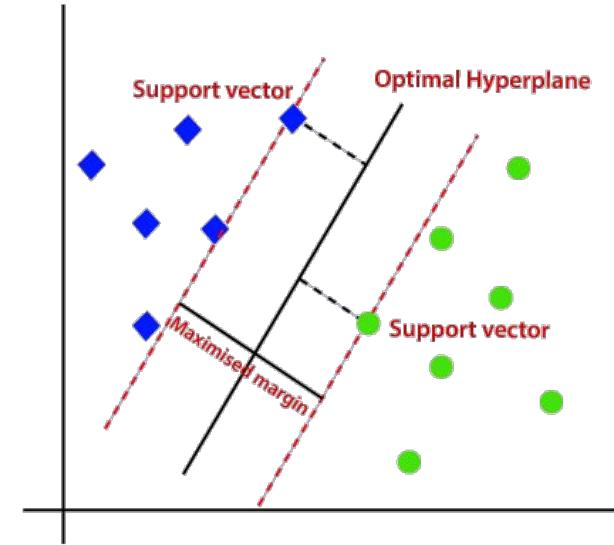
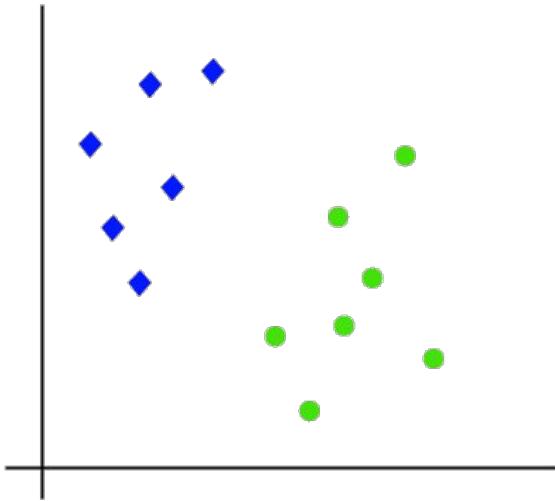


[Varoquaux and Müller](#)

# CLASSICAL SUPPORT VECTOR CLASSIFIER (SVC)

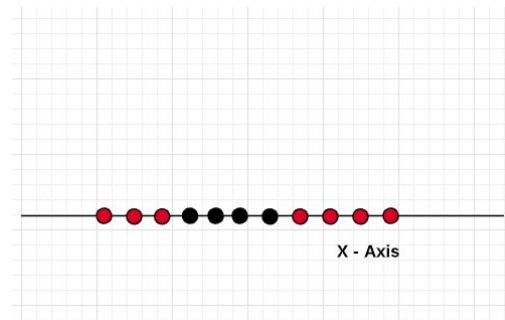


# CLASSICAL SUPPORT VECTOR CLASSIFIER (SVC)



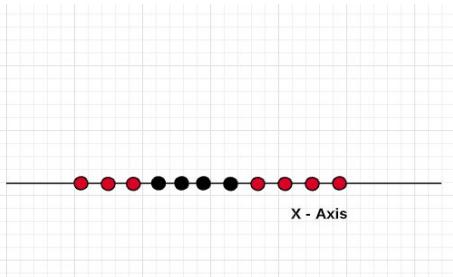
## CLASSICAL SUPPORT VECTOR CLASSIFIER (SVC)

Support Vector Machines were developed to deal with linear data.  
What happens when we take data like:



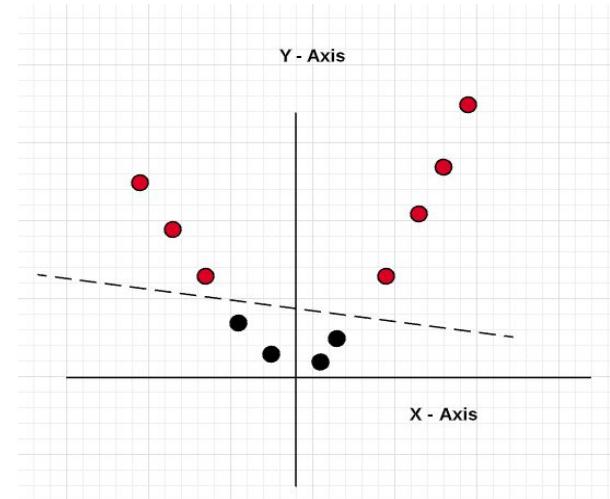
Non-linearly separable data

# CLASSICAL SUPPORT VECTOR CLASSIFIER (SVC)



$$\phi \rightarrow$$
$$y = x^2$$

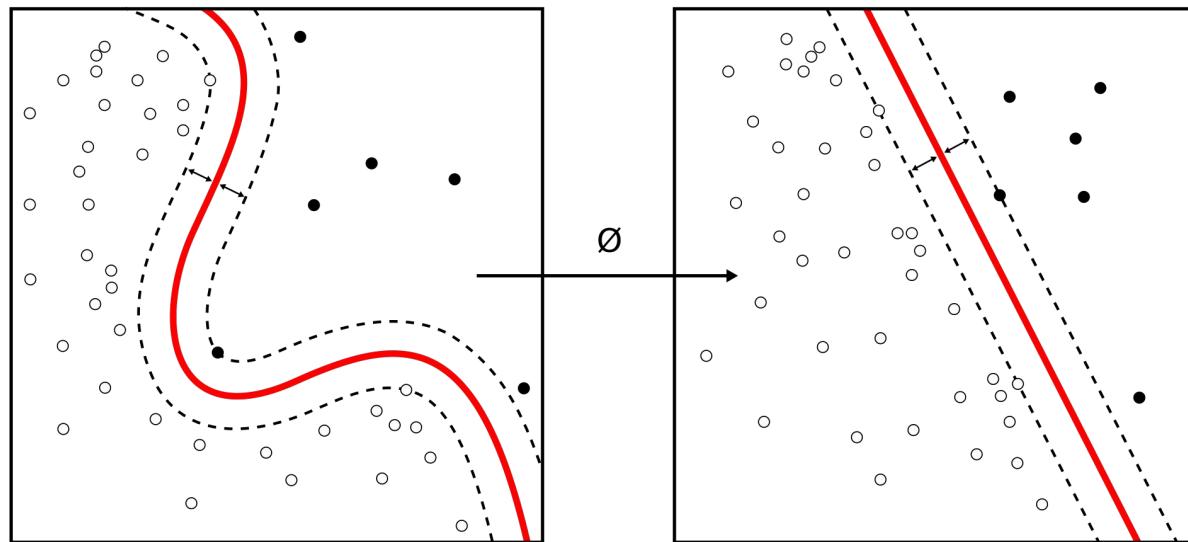
Non-linearly  
separable data



Support Vector Machines have a key component called **kernel machine**.

# CLASSICAL SUPPORT VECTOR CLASSIFIER (SVC)

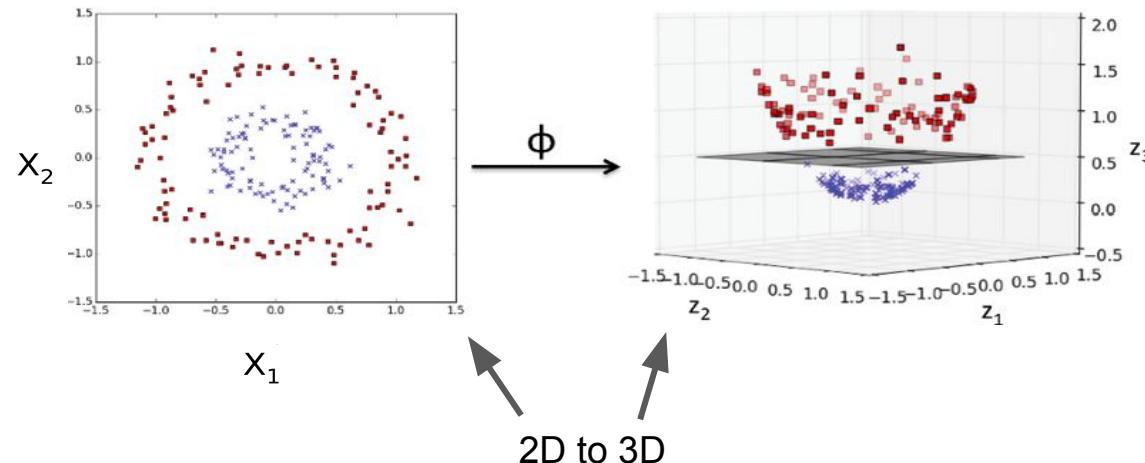
A **Kernel Function** manipulates the training data to transform a non-linear lower dimension space into a higher dimension space, which we can get a linear decision boundary



[SVM - Wikipedia](#)

# CLASSICAL SUPPORT VECTOR CLASSIFIER (SVC)

A **Kernel Function** manipulates the training data to transform a non-linear lower dimension space into a higher dimension space, which we can get a linear decision boundary



[SVM - Wikipedia](#)

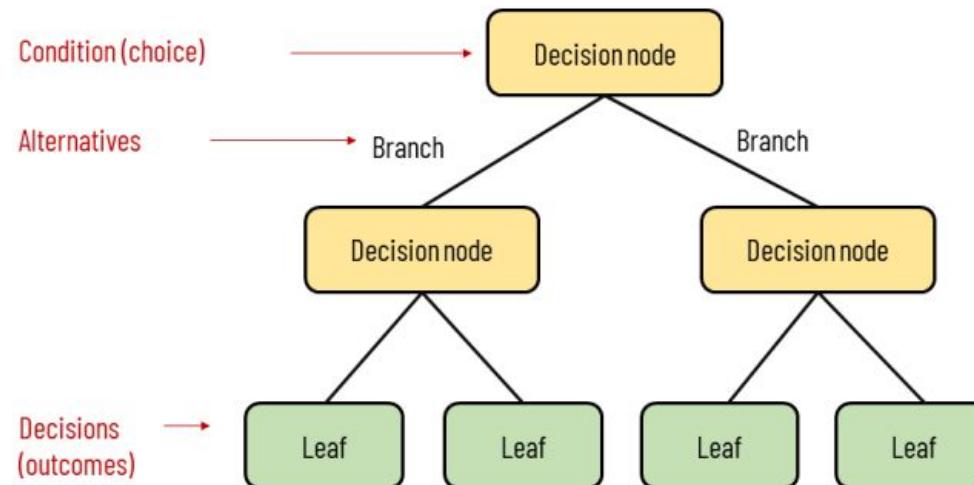
## Part 1: Introduction to Support Vector Classifier (SVC)

[https://github.com/alexgcsa/comp\\_life\\_sciences\\_2023](https://github.com/alexgcsa/comp_life_sciences_2023)

**Part 2:** Support Vector Classifier (SVC)  
for predicting small-molecule that  
inhibit protein-protein interactions

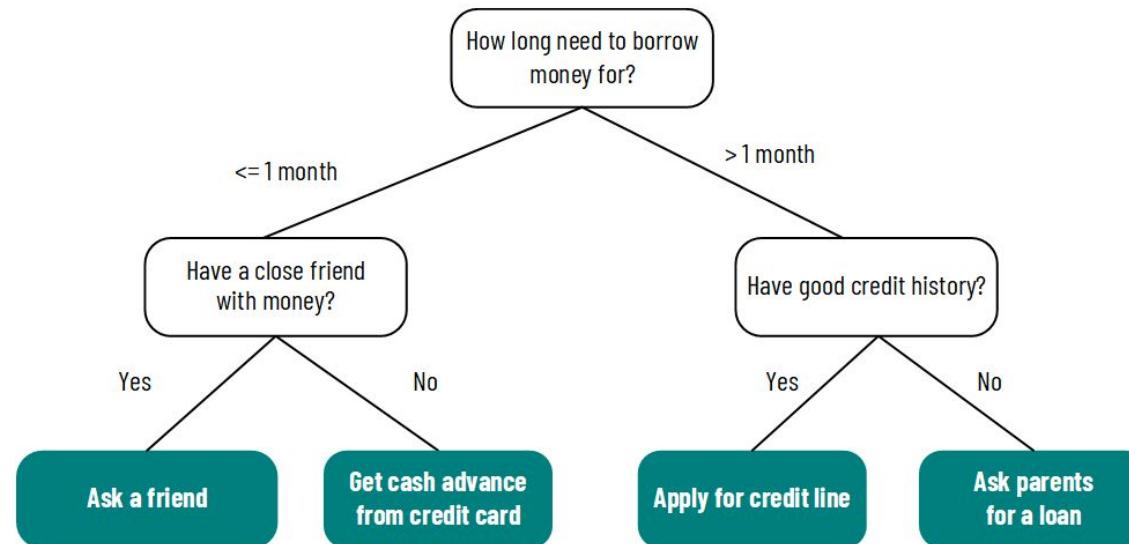
[https://github.com/alexgcsa/comp\\_life\\_sciences\\_2023](https://github.com/alexgcsa/comp_life_sciences_2023)

## Elements of a decision tree

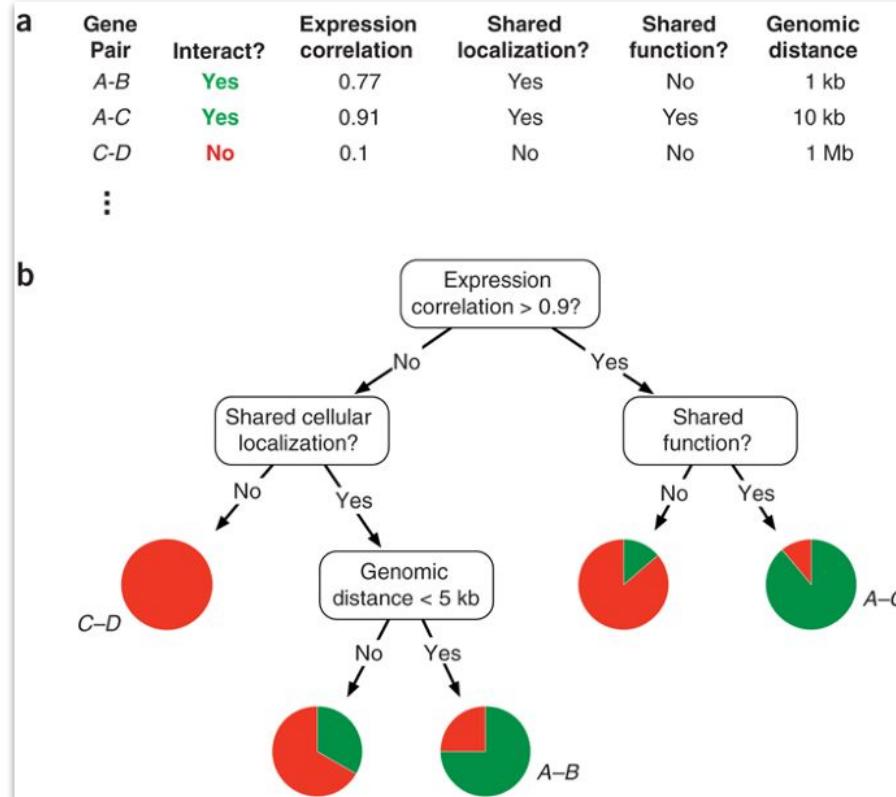


# DECISION TREES

## Decision tree: borrow \$1,000

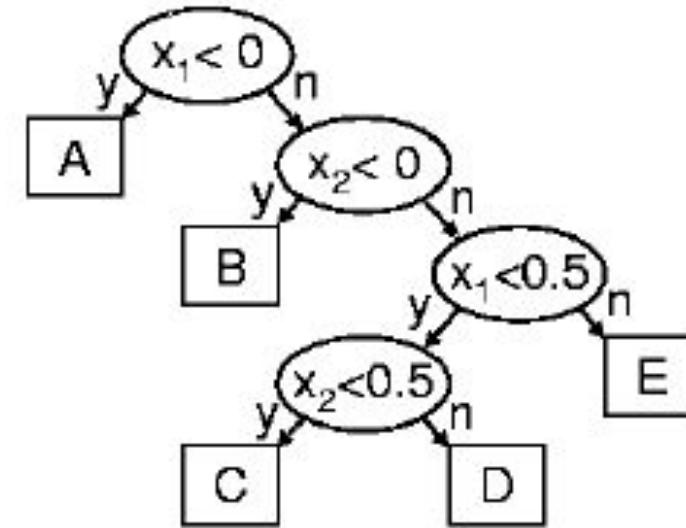
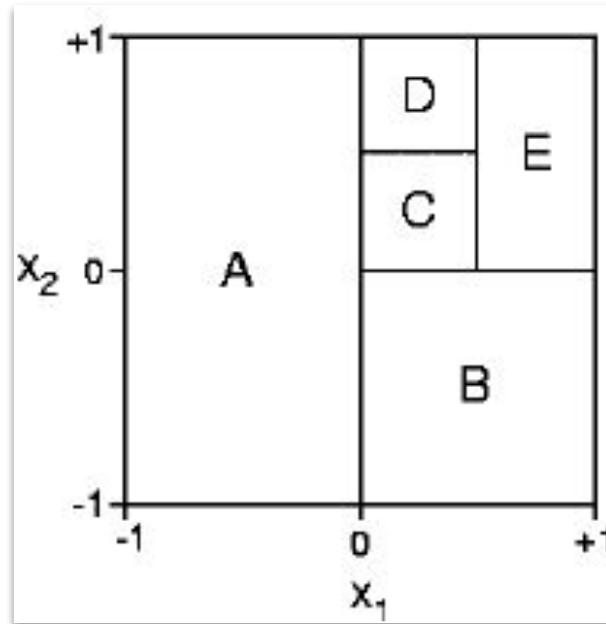


# DECISION TREES



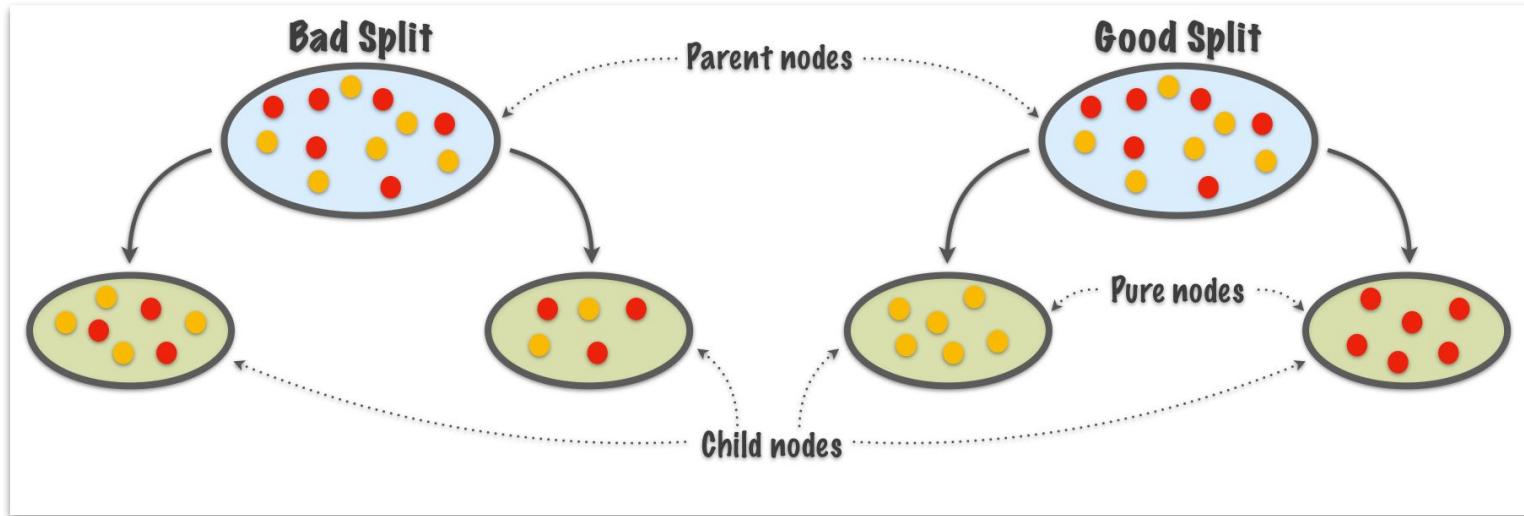
Kingsford & Salzberg, 2008

# DECISION TREES

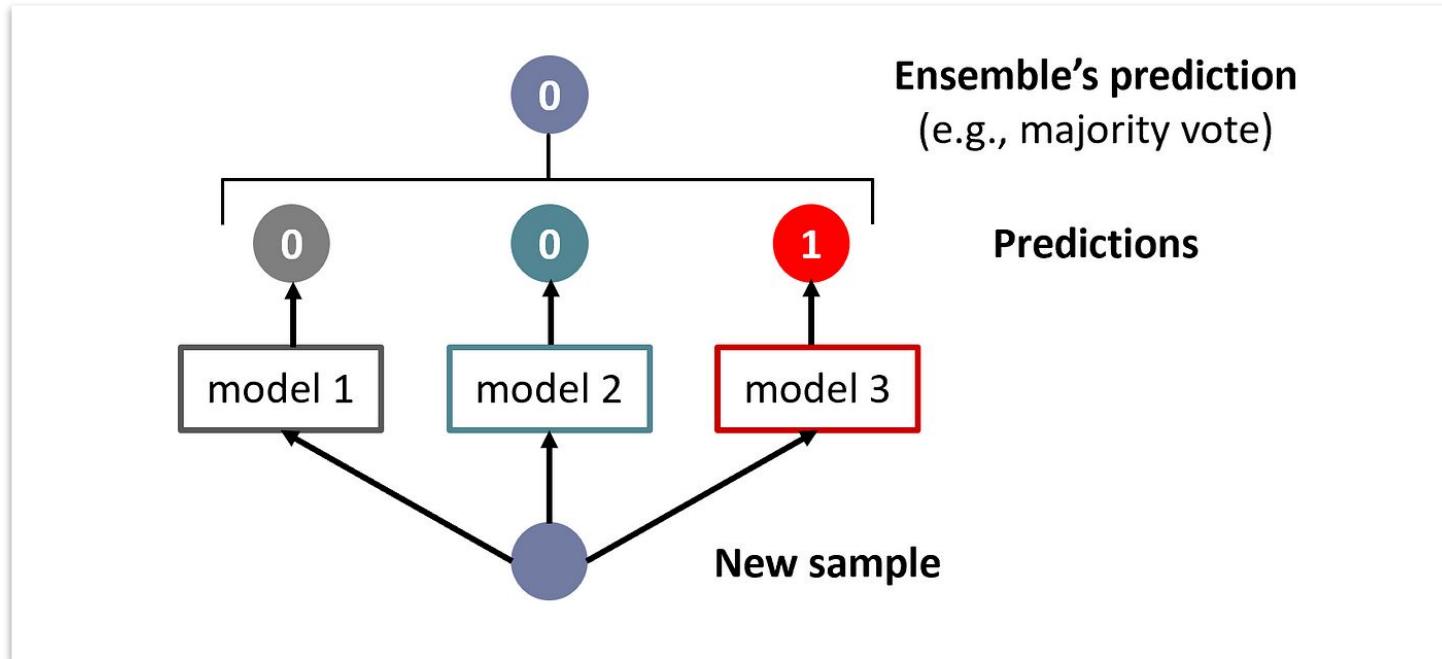


Potts and Sammut, 2000

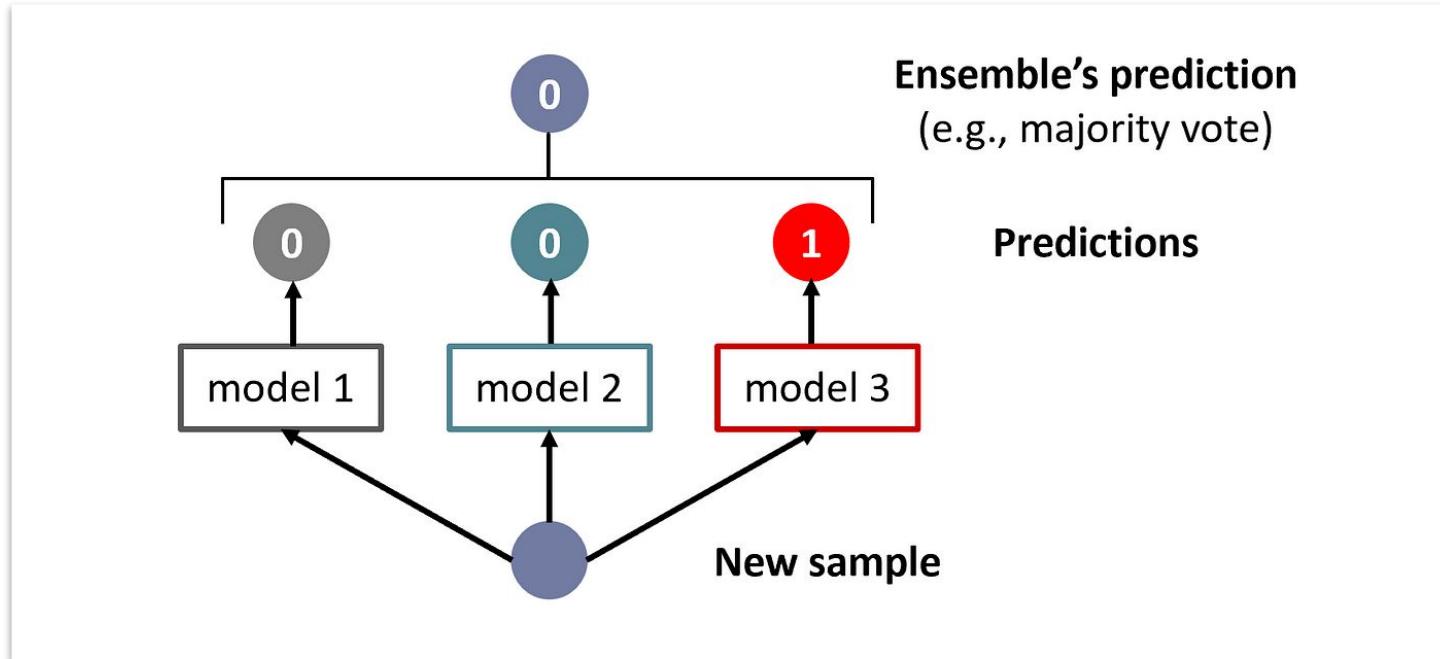
# DECISION TREES



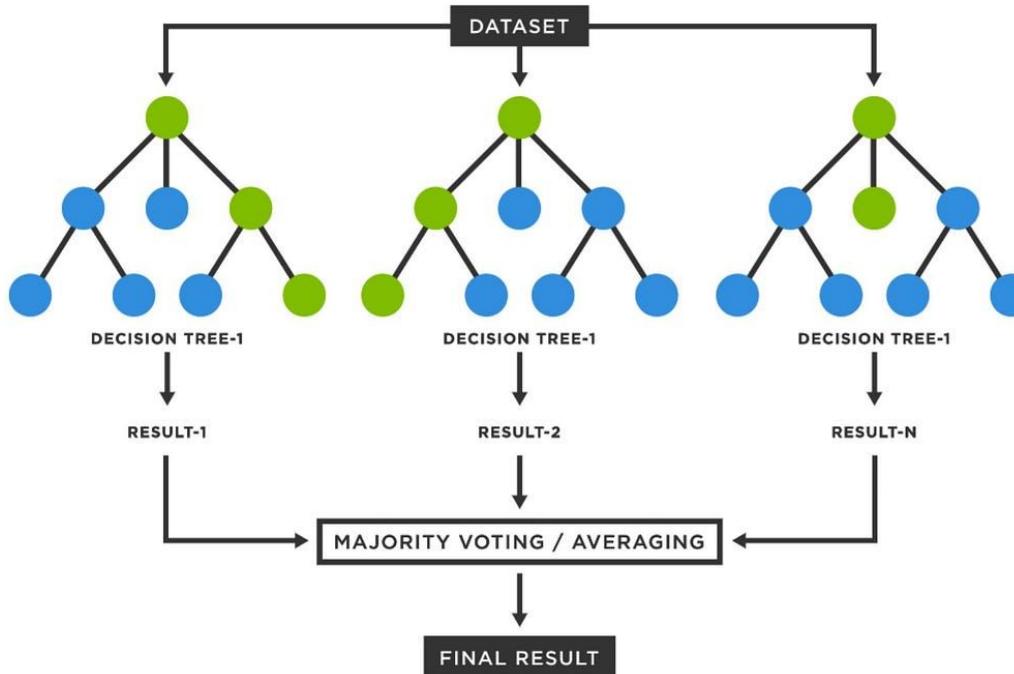
# ENSEMBLES



# RANDOM FOREST

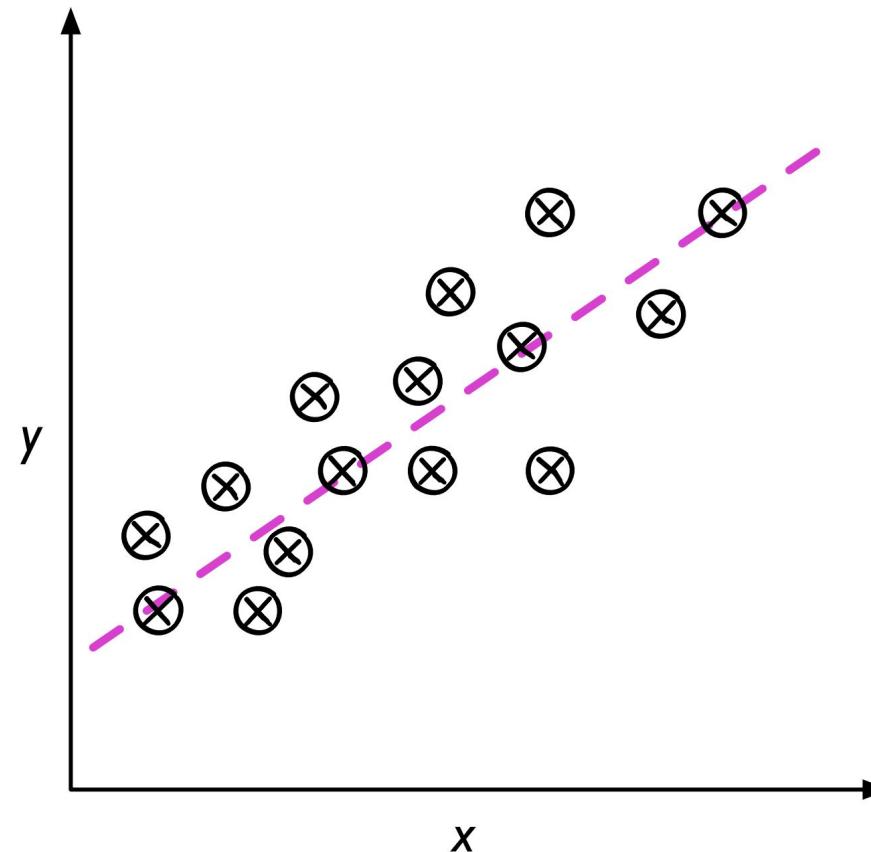


# RANDOM FOREST



# Focusing on Regression

# REGRESSION CURVE



**Part 3:** Support Vector Regressor (SVR)  
for predicting small-molecule that  
inhibit protein-protein interactions

[https://github.com/alexgcsa/comp\\_life\\_sciences\\_2023](https://github.com/alexgcsa/comp_life_sciences_2023)

# ML RESOURCES

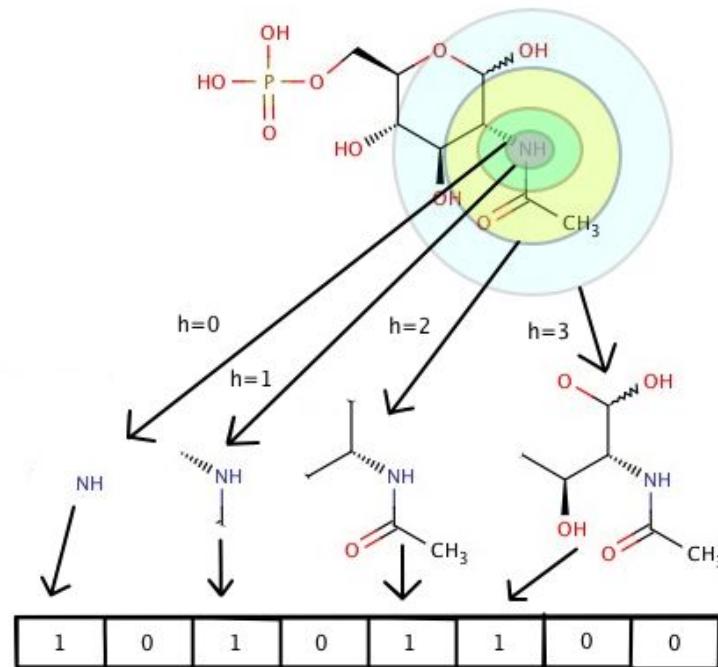
- Tom Mitchell's Book and Youtube Course:
  - <https://www.cs.cmu.edu/~tom/mlbook.html>
  - <https://www.youtube.com/watch?v=m4NlfvrRCdg&list=PLl-BBnDxtUt1hLXmIw u27P22bTi6VwMkN>
- Sebastian Raschka's Course:
  - <https://sebastianraschka.com/blog/2021/ml-course.html>
- Andrew Ng's Course:
  - <https://www.coursera.org/specializations/machine-learning-introduction>

# INTRODUCTION TO MACHINE LEARNING

Alex de Sá

[a.desa@uq.edu.au](mailto:a.desa@uq.edu.au)  
[@alexgcsa](https://twitter.com/alexgcsa)

# MORGAN OR CIRCULAR FINGERPRINT

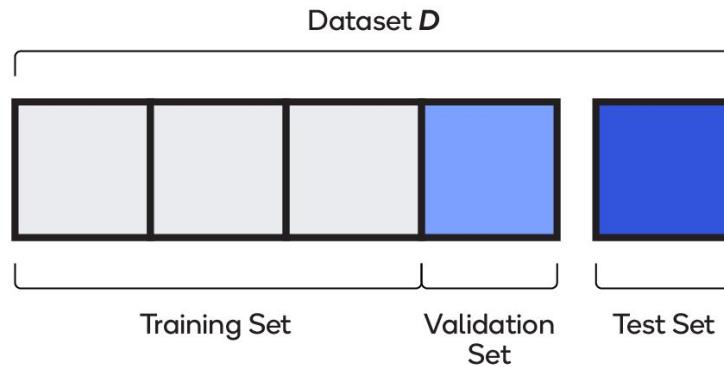


# K-FOLD CROSS-VALIDATION

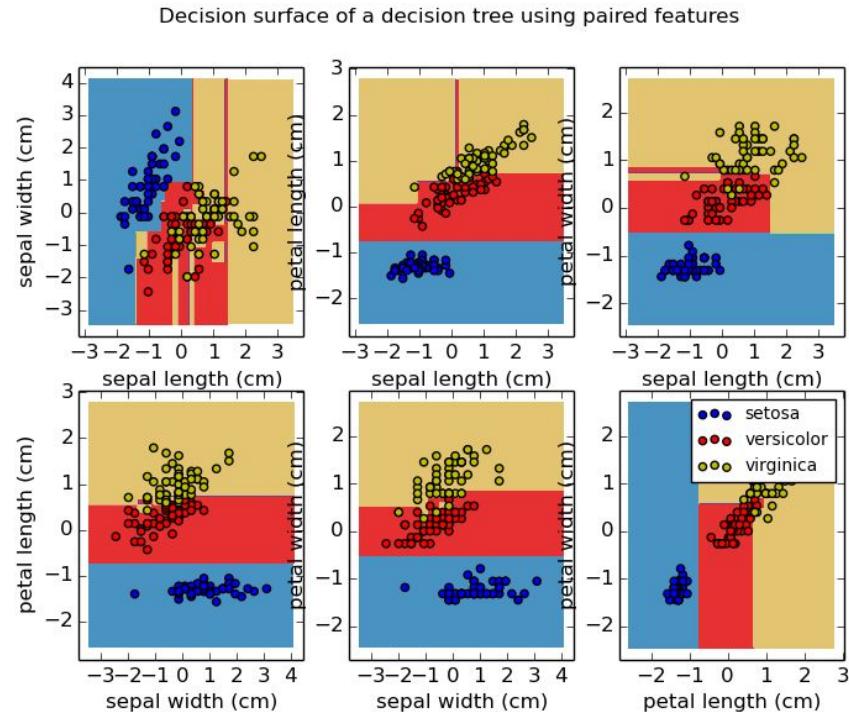
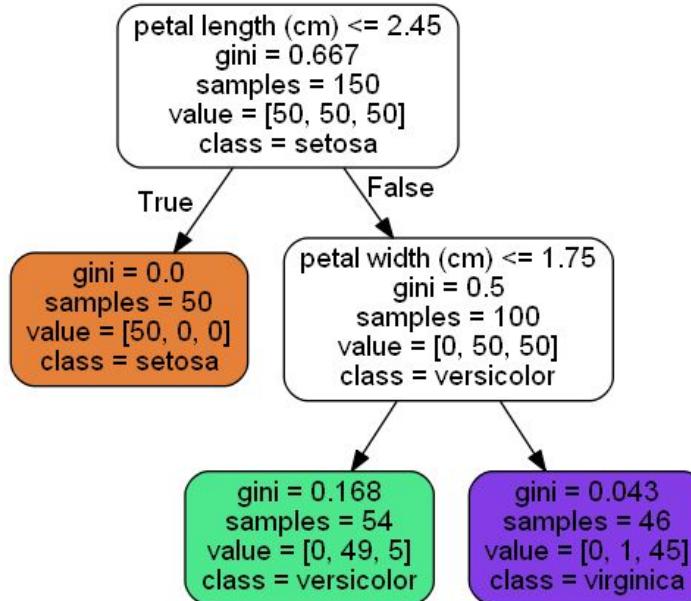


# TRAIN/TEST

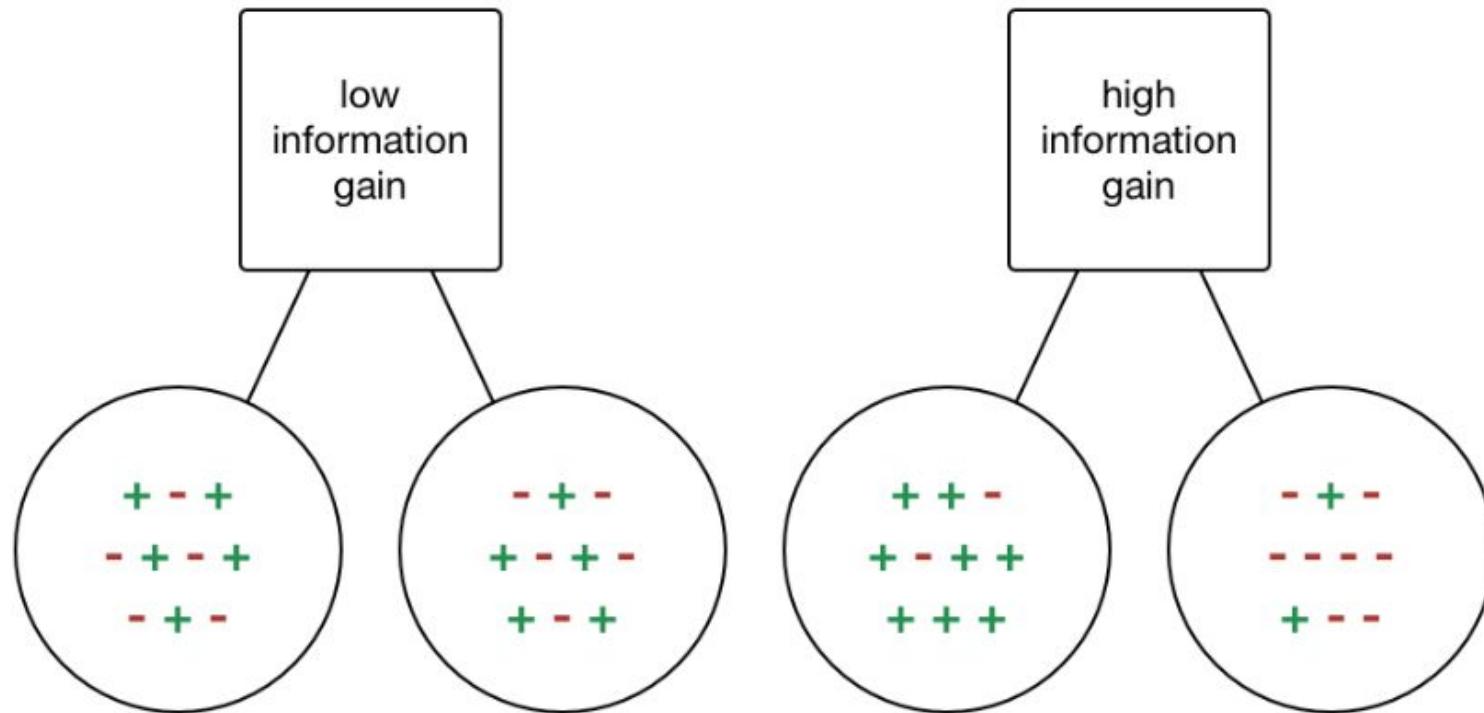
Test the model on new data, assessing its generalisation



# DECISION TREE CLASSIFIER



## DECISION TREE CLASSIFIER - SPLITTING



# CLASSIFICATION METRICS

