

# Introduction to machine learning and **scikit-learn**

## *Part I: Supervised learning*

QLSC 612 | 29 May 2025

*By*

*Michelle Wang & Mohammad Torabi*

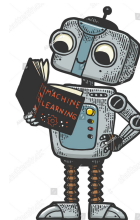
*(reusing some of Nikhil Bhagwat's slides)*



**McGill**  
UNIVERSITY



**neuro**  
Montreal Neurological  
Institute-Hospital

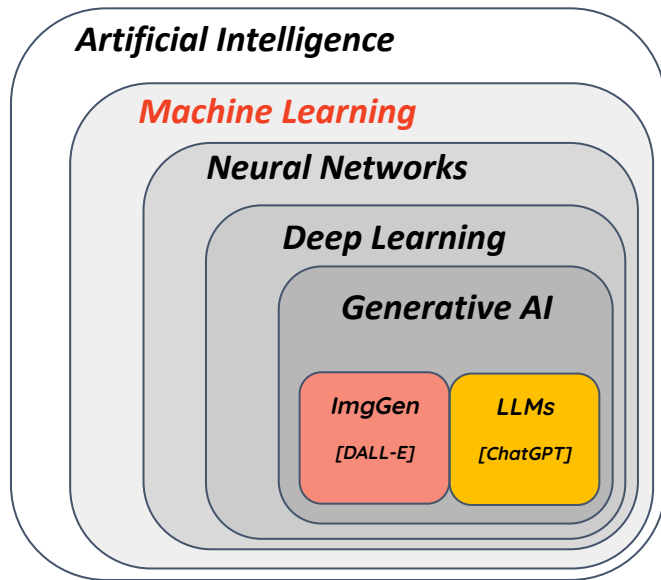


# Outline

- Machine learning overview
- Supervised learning
  - Goal
  - Example models
  - Supervised learning with scikit-learn interface
  - Model evaluation
- Deep learning (brief)

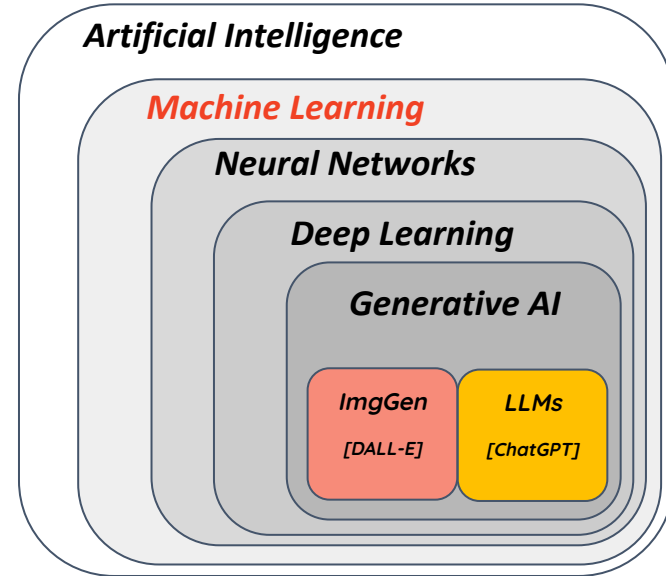
# Machine-learning - what, why, and when?

- What is Machine learning (ML)?
  - ML is the study of computer algorithms that improve automatically through “experience” and by the use of data.



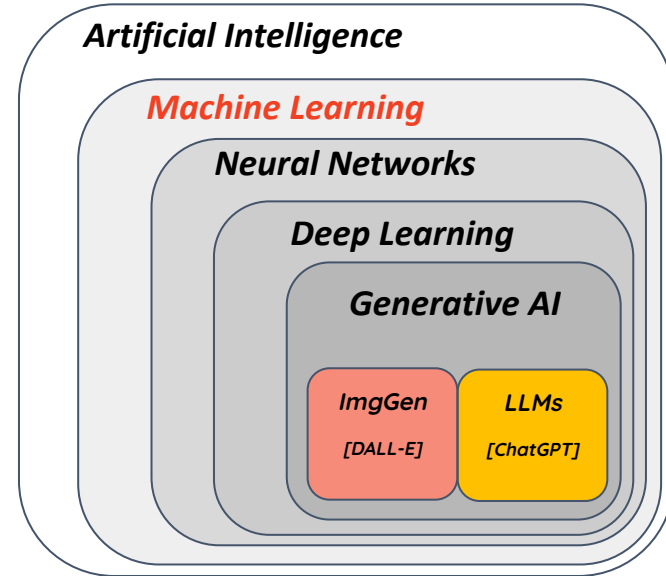
# Machine-learning - what, why, and when?

- What is Machine learning (ML)?
  - ML is the study of computer algorithms that improve automatically through “experience” and by the use of data.
- Why is it useful - especially in life sciences?
  - Biology, medicine, environmental sciences comprise phenomena (e.g. a disease) with large number of variables.
  - We want to model complex relationships within these variables and **make accurate predictions**.



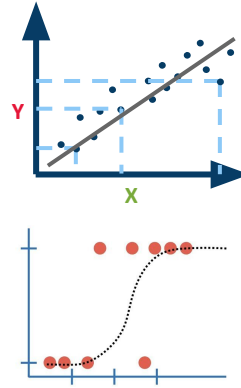
# Machine-learning - what, why, and when?

- What is Machine learning (ML)?
  - ML is the study of computer algorithms that improve automatically through “experience” and by the use of data.
- Why is it useful - especially in life sciences?
  - Biology, medicine, environmental sciences comprise phenomena (e.g. a disease) with large number of variables.
  - We want to model complex relationships within these variables and **make accurate predictions**.
- When do I use it?
  - You are interested in 1) prediction tasks or 2) low-dimensional representation.
  - **You have sufficient data.**



# Types of ML Algorithms

- **Supervised** → labels are known
  - Regression → labels are continuous
  - Classification → labels are discrete

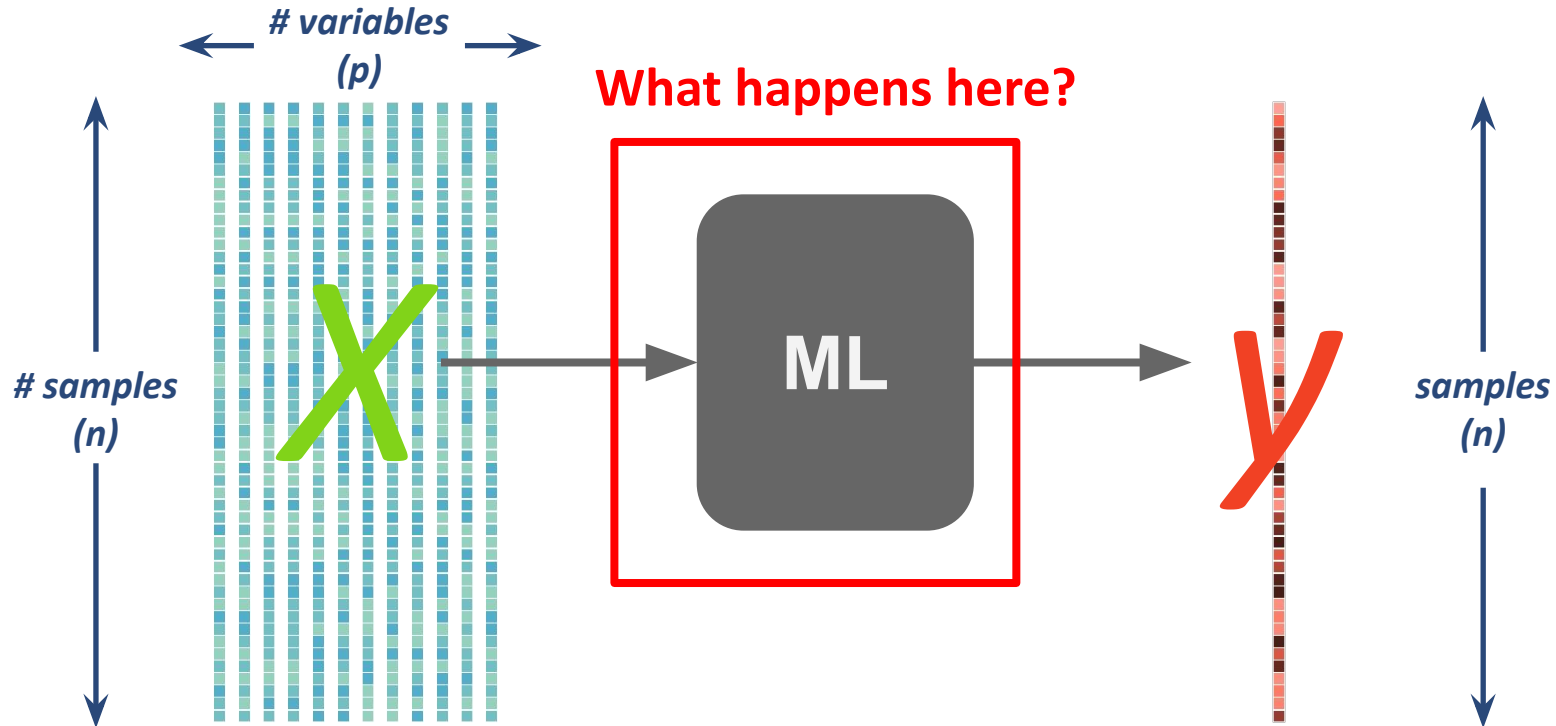


- **Unsupervised** → labels are unknown
  - Associations, dimensionality reduction, clustering
  - Covered in Part 2

# Machine learning terminology

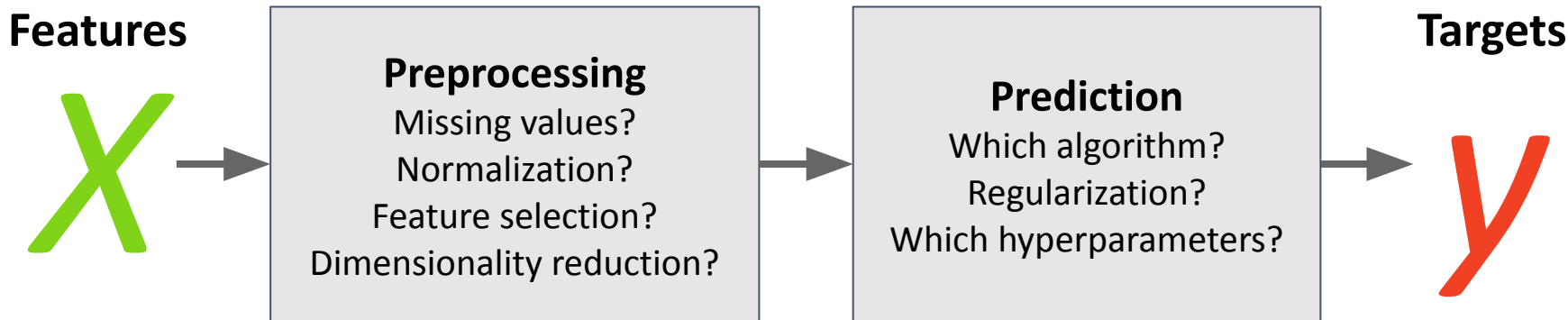
Input (features, etc.)

Output (labels, targets, etc.)



# A typical supervised learning workflow

*Decision points when developing a model*



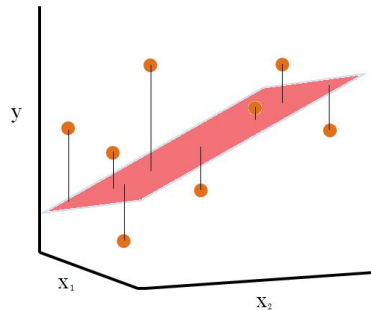


# Supervised Learning: Models

- **Goal:** Learn parameters (or weights) of a model that maps  $x$  to  $y$

# Supervised Learning: Models

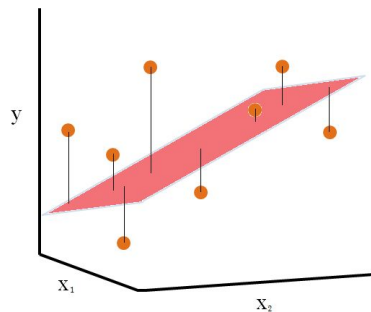
- **Goal:** Learn parameters (or weights) of a model that maps  $\mathbf{x}$  to  $\mathbf{y}$
- Example models (see also [scikit-learn documentation](#)):
  - Linear / Logistic regression



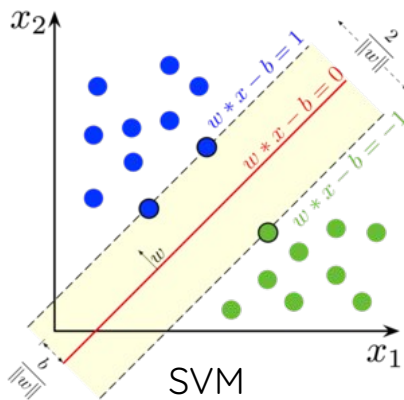
Linear Regression

# Supervised Learning: Models

- **Goal:** Learn parameters (or weights) of a model that maps  $x$  to  $y$
- Example models (see also [scikit-learn documentation](#)):
  - Linear / Logistic regression
  - Support vector machines



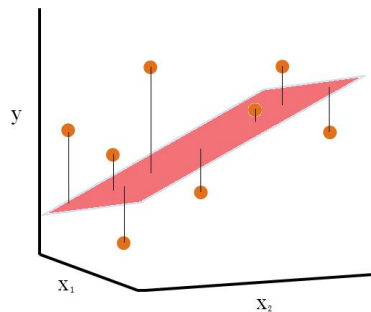
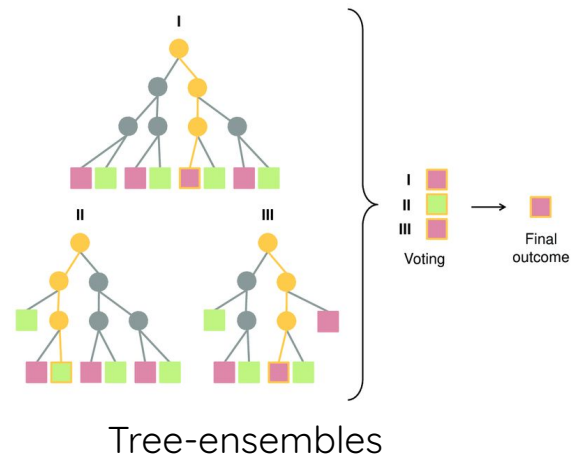
Linear Regression



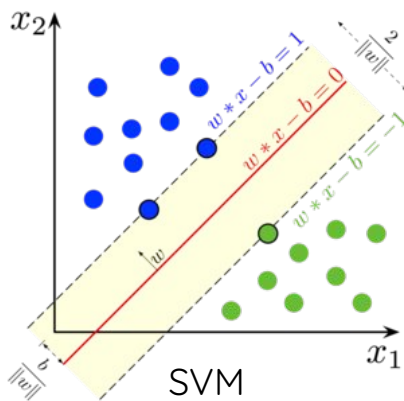
SVM

# Supervised Learning: Models

- **Goal:** Learn parameters (or weights) of a model that maps  $x$  to  $y$
- Example models (see also [scikit-learn documentation](#)):
  - Linear / Logistic regression
  - Support vector machines
  - Tree-ensembles: random forests, gradient boosting



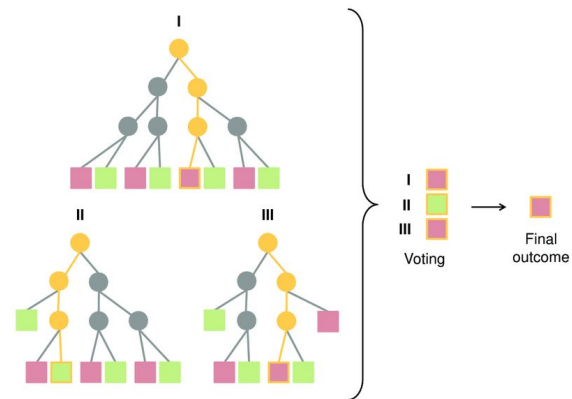
Linear Regression



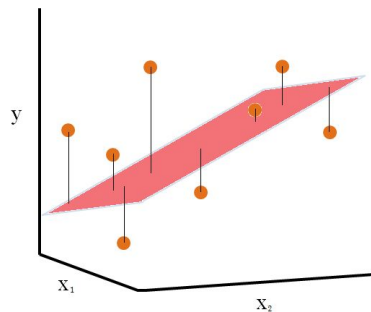
SVM

# Supervised Learning: Models

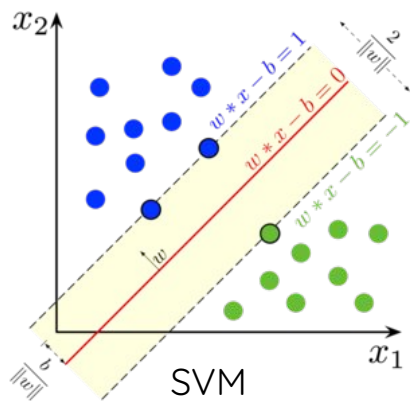
- **Goal:** Learn parameters (or weights) of a model that maps  $x$  to  $y$
- Example models (see also [scikit-learn documentation](#)):
  - Linear / Logistic regression
  - Support vector machines
  - Tree-ensembles: random forests, gradient boosting
  - Artificial Neural networks



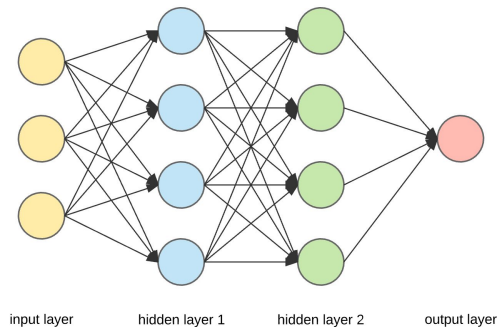
Tree-ensembles



Linear Regression



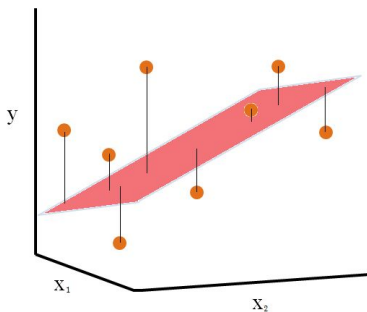
SVM



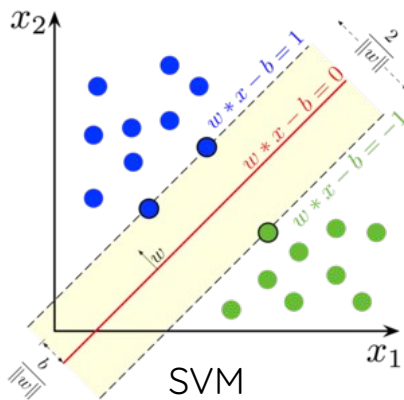
ANN

# Supervised Learning: Models

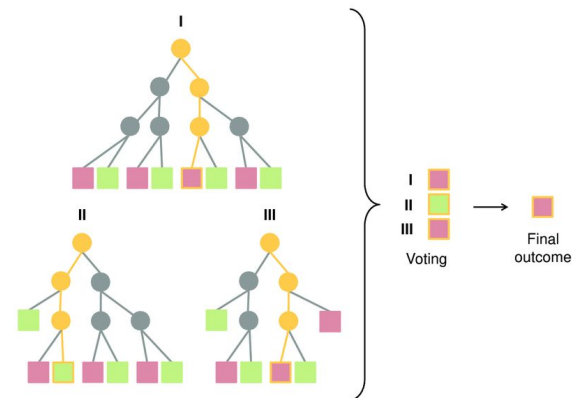
- **Goal:** Learn parameters (or weights) of a model that maps  $x$  to  $y$
- Example models (see also [scikit-learn documentation](#)):
  - Linear / Logistic regression
  - Support vector machines
  - Tree-ensembles: random forests, gradient boosting
  - Artificial Neural networks
- **How are these models different from one another?**



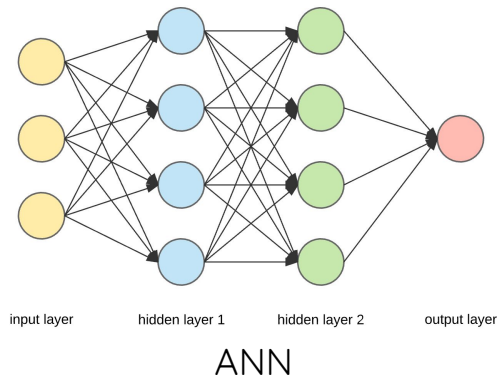
Linear Regression



SVM



Tree-ensembles

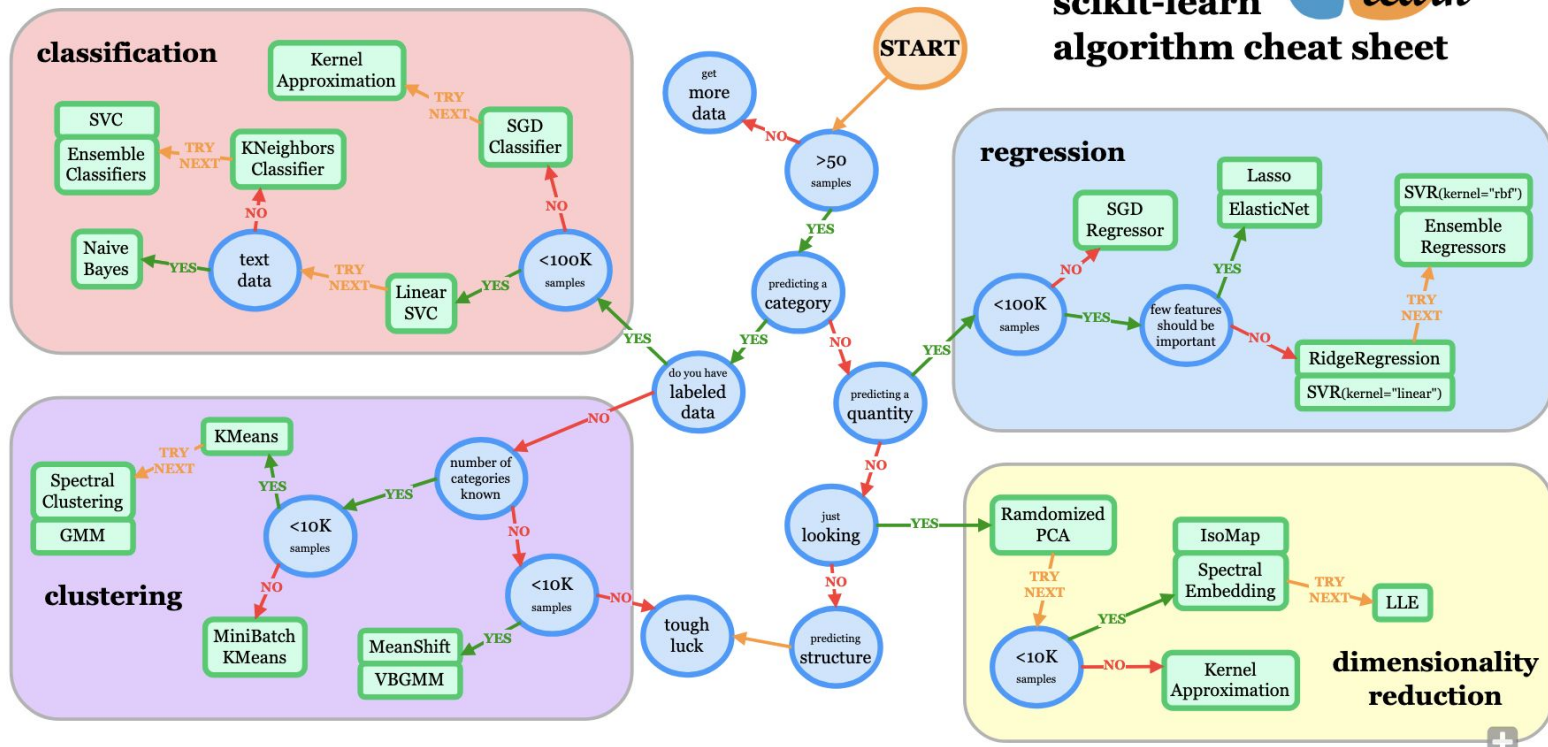


ANN

# Some models make more sense in some situations

[https://scikit-learn.org/stable/machine\\_learning\\_map.html](https://scikit-learn.org/stable/machine_learning_map.html)

## scikit-learn algorithm cheat sheet



# Model fitting is easy with scikit-learn

Example with **linear regression**

```
# import
from sklearn.linear_model import Lasso
```

```
# data
X = [[0, 0], [1, 1]]
y = [0, 1]
```

```
# instantiate the model
model = Lasso()
```

Change this to use different  
models/hyperparameters

```
# fit the model with data
model.fit(X, y)
```

```
# predict on new data
y_pred = model.predict([[1, 0]])
```



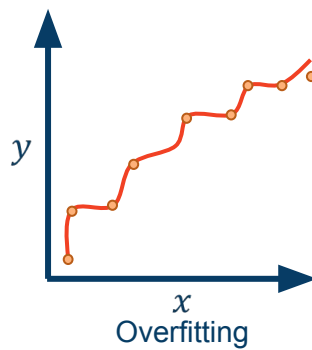
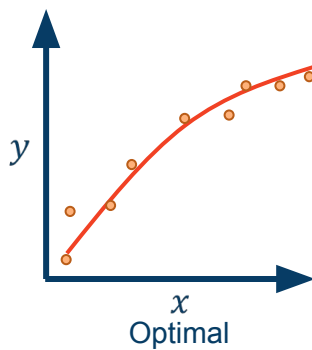
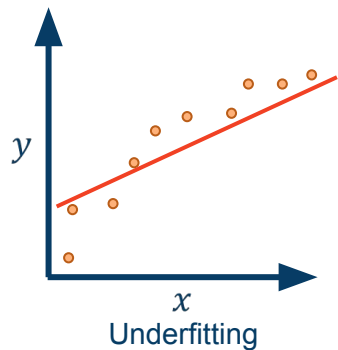
# I fitted my model, now what?

- Model evaluation metrics
  - **Regression:**  $R^2$ , mean squared error, mean absolute error, etc.
  - **Classification:** balanced accuracy, [AUROC](#), confusion matrix, etc.
  - See [https://scikit-learn.org/stable/modules/model\\_evaluation.html](https://scikit-learn.org/stable/modules/model_evaluation.html) for more
- How does the model perform
  - On the data it was trained on?
  - On previously unseen data?

**We want good generalizability on new data**

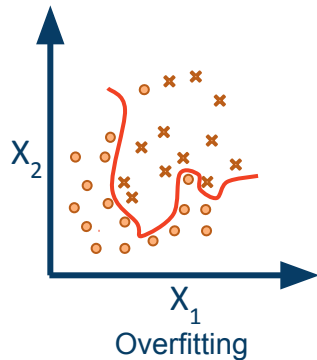
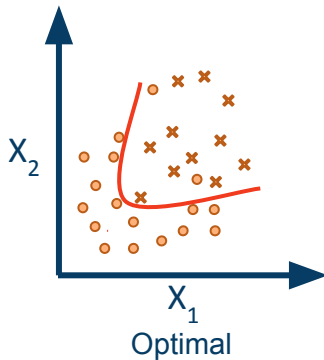
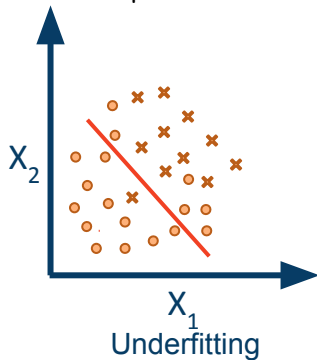
# Models can overfit (or underfit)

Example: **regression**



● Training data

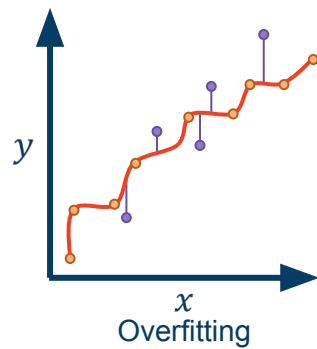
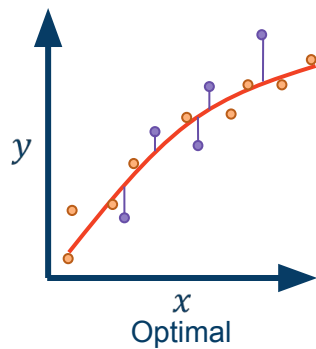
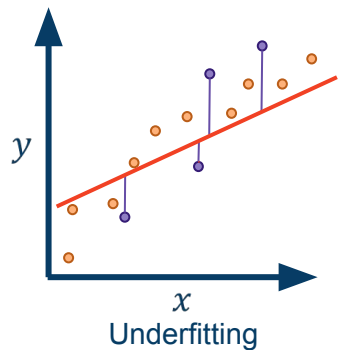
Example: **classification**



○● Training data

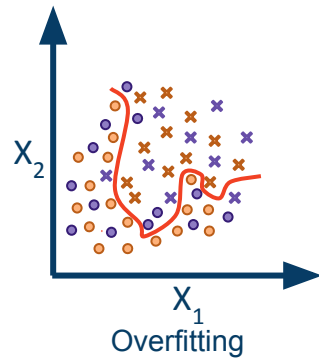
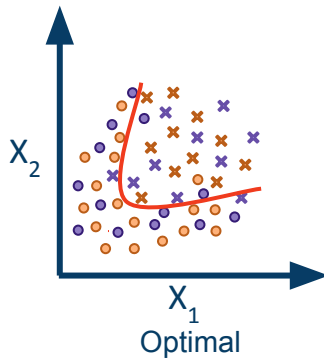
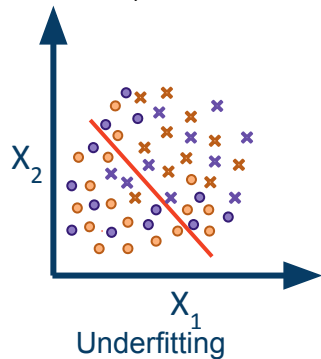
# Models can overfit (or underfit)

Example: **regression**



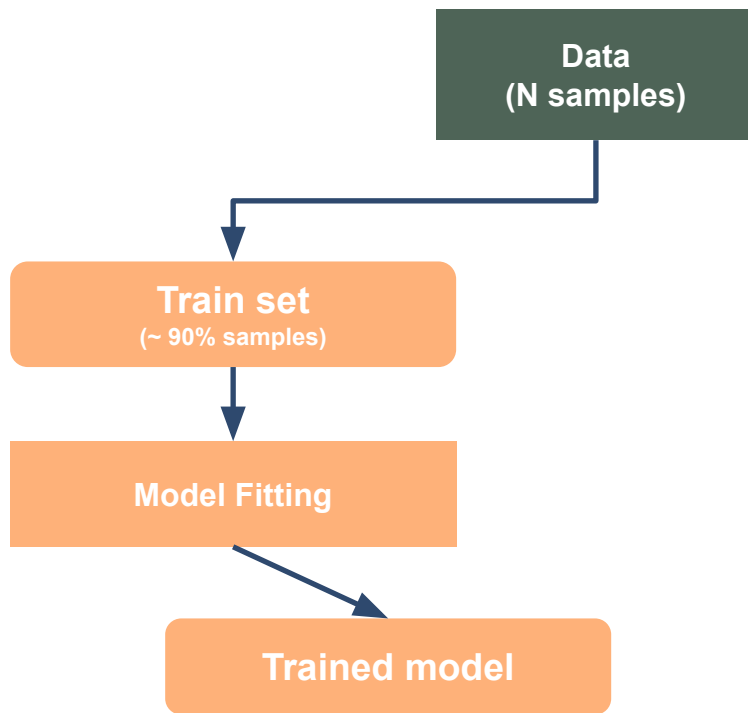
● Training data  
● New data

Example: **classification**

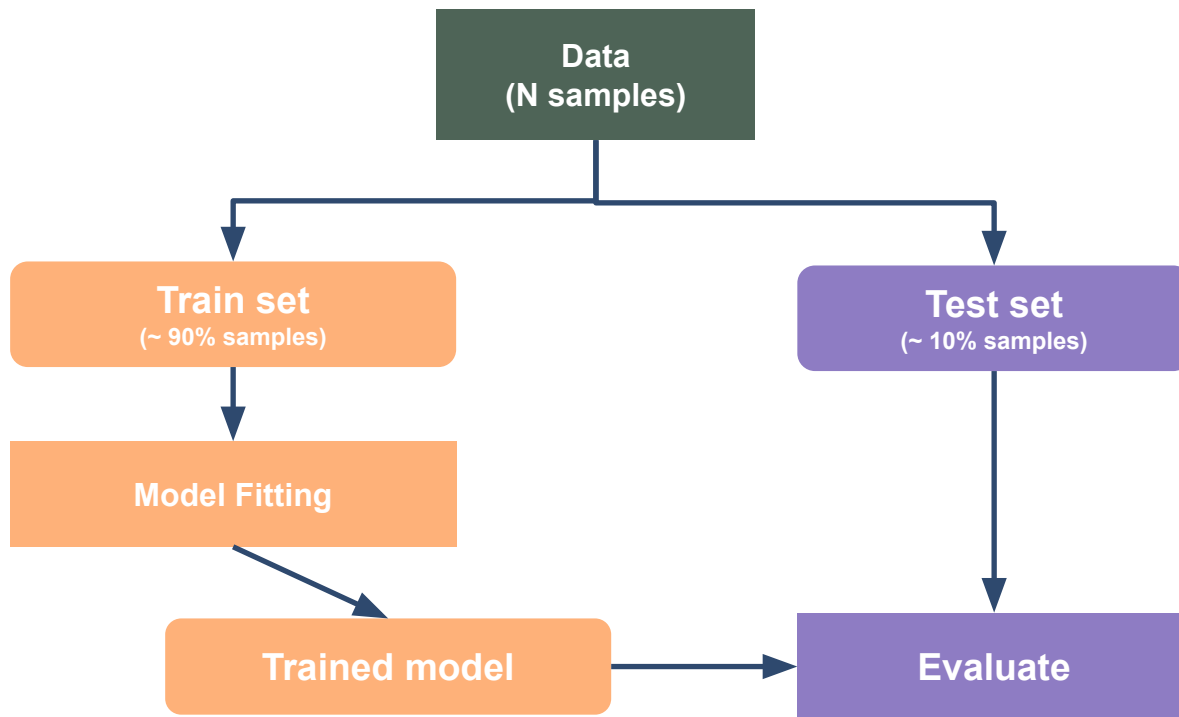


✕ ● Training data  
✕ ● New data

# Split data into train and test sets

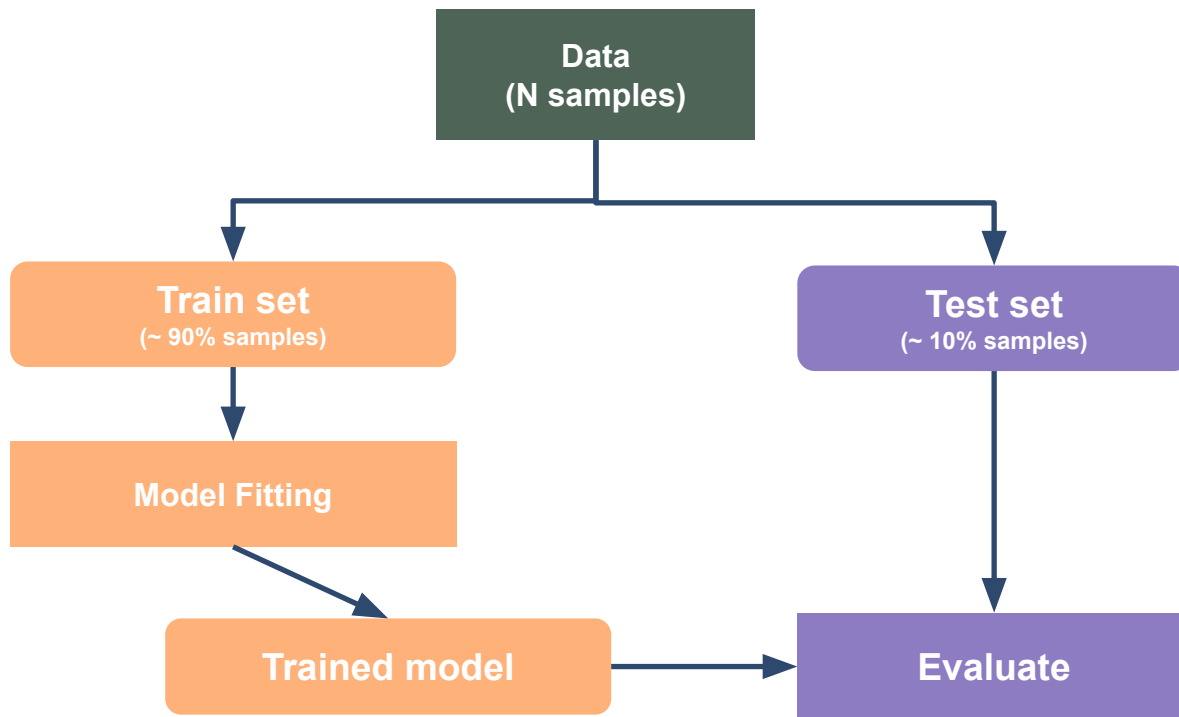


# Split data into train and test sets



# Exercise 1

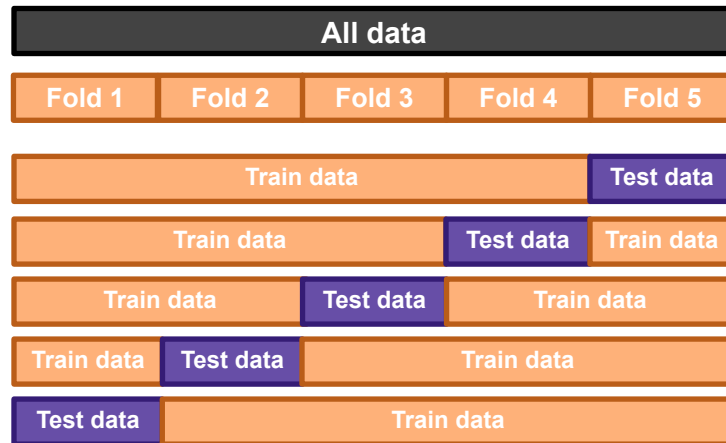
# Split data into train and test sets



**How to sample the train and test sets?**

# K-fold cross-validation

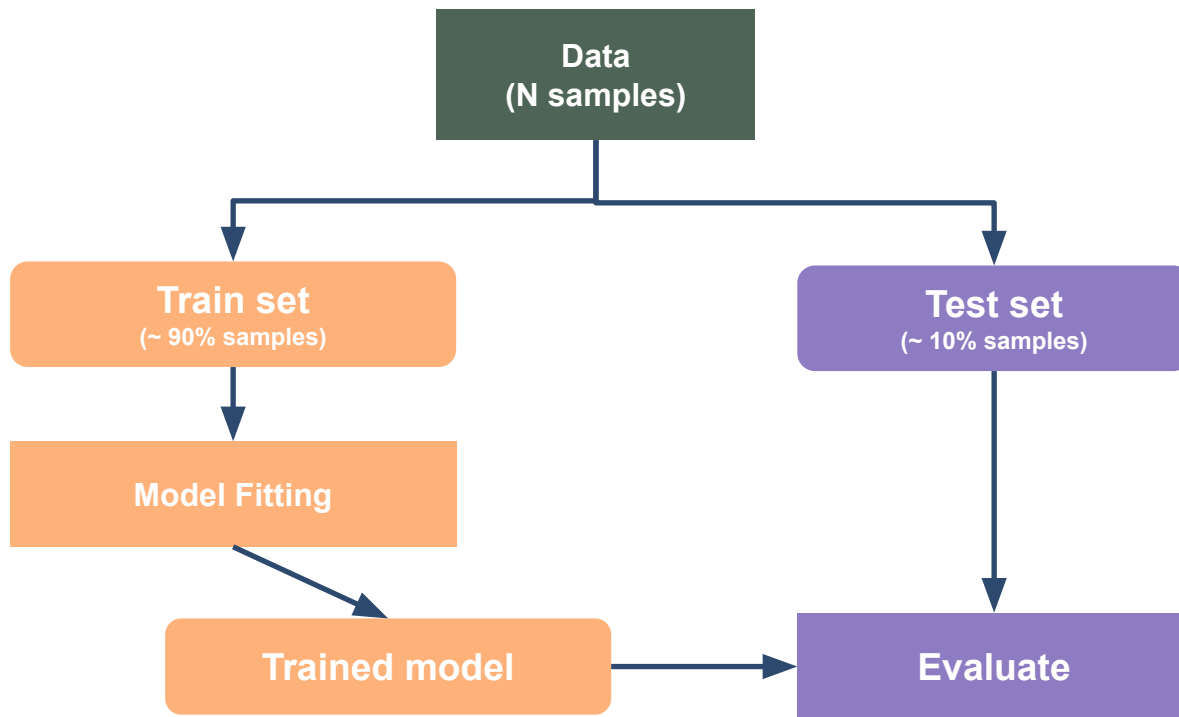
- How do we sample train and test sets?
  - Train set: learn model parameters
  - Test set (a.k.a held-out sample): Evaluate model performance
  - Repeat for different Train-Test splits
  - Report performance statistics over all test folds



Alternative method: shuffle-split ([https://scikit-learn.org/stable/modules/cross\\_validation.html#shufflesplit](https://scikit-learn.org/stable/modules/cross_validation.html#shufflesplit))



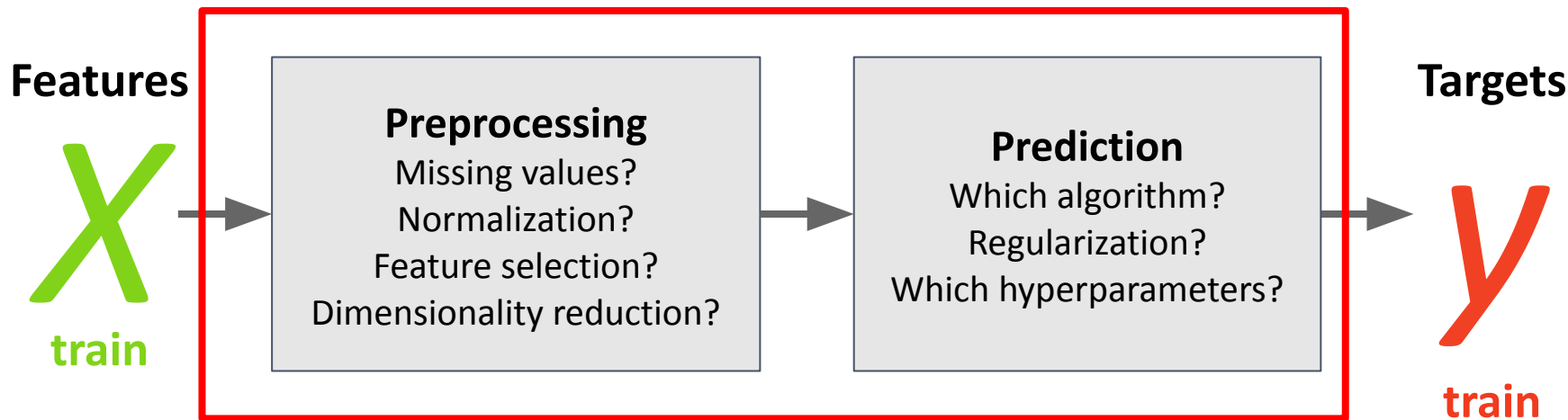
# Split data into train and test sets



**Be careful about data leakage/double-dipping!**

# A typical supervised learning workflow

*Decision points when developing a model*

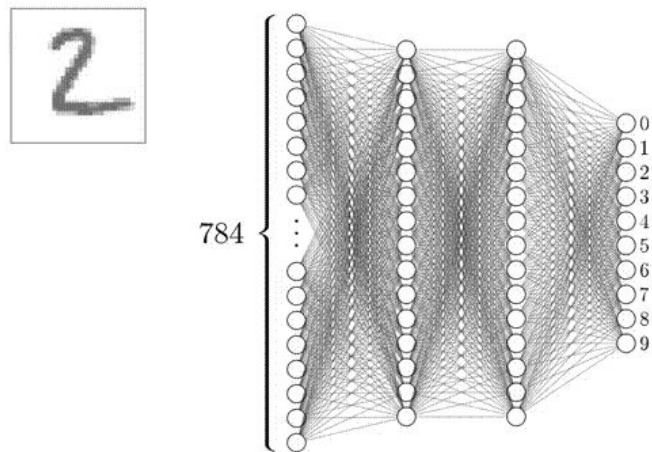


**Do not use test data to make these decisions!**  
**(z-score mean/std., hyperparameter tuning, etc.)**

Questions?

# Deep-learning

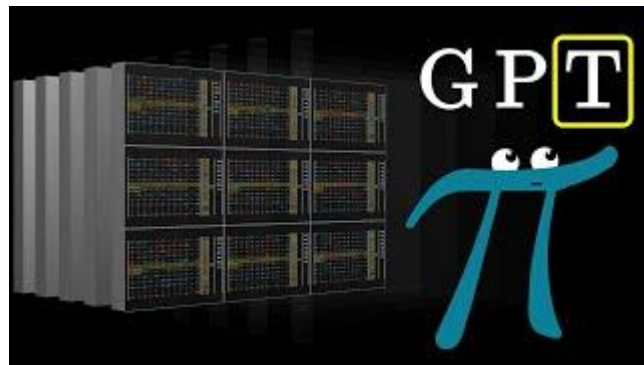
- Why the buzz?
  - Works amazing on spatio-temporal input
  - Highly flexible → universal function approximator



ANN for handwritten-digit images  
(gif source: [3b1b](#))

# Deep-learning

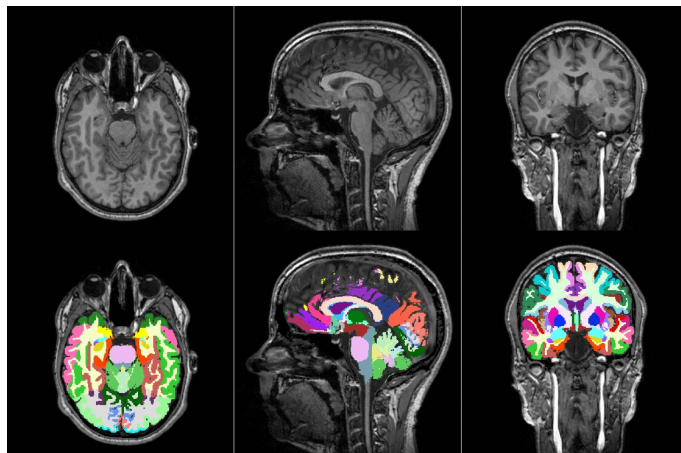
- Why the buzz?
  - Works amazing on spatio-temporal input
  - Highly flexible → universal function approximator
- What are the challenges?
  - Large number of parameters (175B!) → data hungry
  - Large number of hyper-parameters → difficult to train



LLM Transformers  
(gif source: [3b1b](#))

# Deep-learning

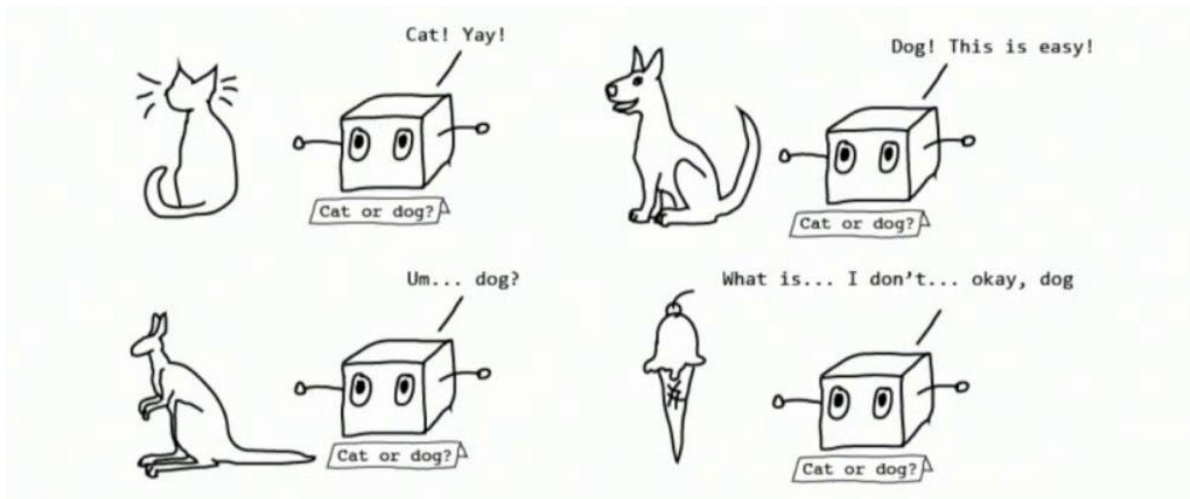
- Why the buzz?
  - Works amazing on spatio-temporal input
  - Highly flexible → universal function approximator
- What are the challenges?
  - Large number of parameters (175B!) → data hungry
  - Large number of hyper-parameters → difficult to train
- When do I use it?
  - If you have highly-structured input, eg. medical images.
  - You have a lot of data and computational resources.



Source:  
<https://github.com/fepegar/torchio>

# Pitfalls and Challenges

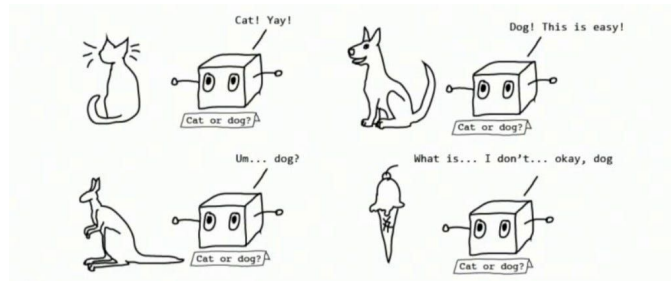
- Models do not generalize even after good CV performance
  - Implicit double-dipping
  - Dataset biases (eg. North-American demographics)
  - Noisy labels (eg. diagnosis definitions)
  - Data distribution shifts (eg. assay, scanner upgrades)



# Pitfalls and Challenges

- Models do not generalize even after good CV performance

- Implicit double-dipping
- Dataset biases (eg. North-American demographics)
- Noisy labels (eg. diagnosis definitions)
- Data distribution shifts (eg. assay, scanner upgrades)



- Unnecessary complexity

- Do I really need a giant deep-net or a simple linear model would do?





# ML Novice Checklist

- Data

- What is my `n_features` and `n_samples`?
- Am I [encoding](#) categorical data correctly?
- Am I using information (e.g. mean) from test set to preprocess (eg. z-score) the data?

# ML Novice Checklist

## ○ Data

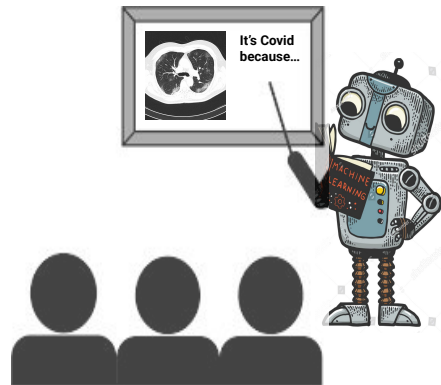
- What is my n\_features and n\_samples?
- Am I [encoding](#) categorical data correctly?
- Am I using information (e.g. mean) from test set to preprocess (eg. z-score) the data?

## ○ Model

- Do my performance metrics capture the practical use-case of interest?
- What is the null / dummy model performance?
  - Classification: Predict majority class all the time
  - Regression: Predict the median value all the time
- Am I interpreting model parameters (i.e. weights) correctly?

# Takeaways

- Supervised ML is useful for **predictions** but **not really for explanations**
  - eg. image segmentation, prognosis, drug development
- Our job is to ensure **generalizability** of these models
  - Multitude of validations
  - Understanding model biases and limitations
- **Engineering tools** vs *Scientific discovery*
  - Interpretability and explainability



Explainable AI

# Useful resources

- [https://scikit-learn.org/stable/user\\_guide.html](https://scikit-learn.org/stable/user_guide.html)
- **nilearn**, Python package for machine learning for brain images:  
<https://nilearn.github.io/stable/index.html>
- **skrub**, Python package machine learning for tabular data:  
<https://skrub-data.org/stable/>
- [https://inria.github.io/scikit-learn-mooc/ml\\_concepts/slides.html](https://inria.github.io/scikit-learn-mooc/ml_concepts/slides.html)
- <https://www.3blue1brown.com/topics/linear-algebra>
- 3Blue1Brown Gradient Descent:  
<https://www.youtube.com/watch?v=IHZwWFHWa-w>