

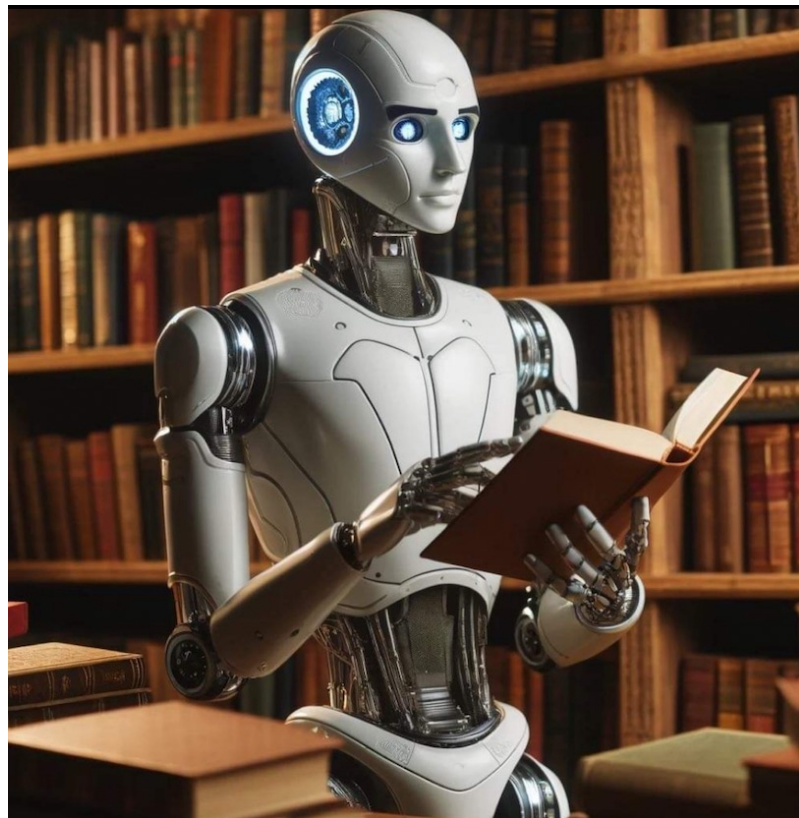


## Introduction to Machine Learning (Part III)

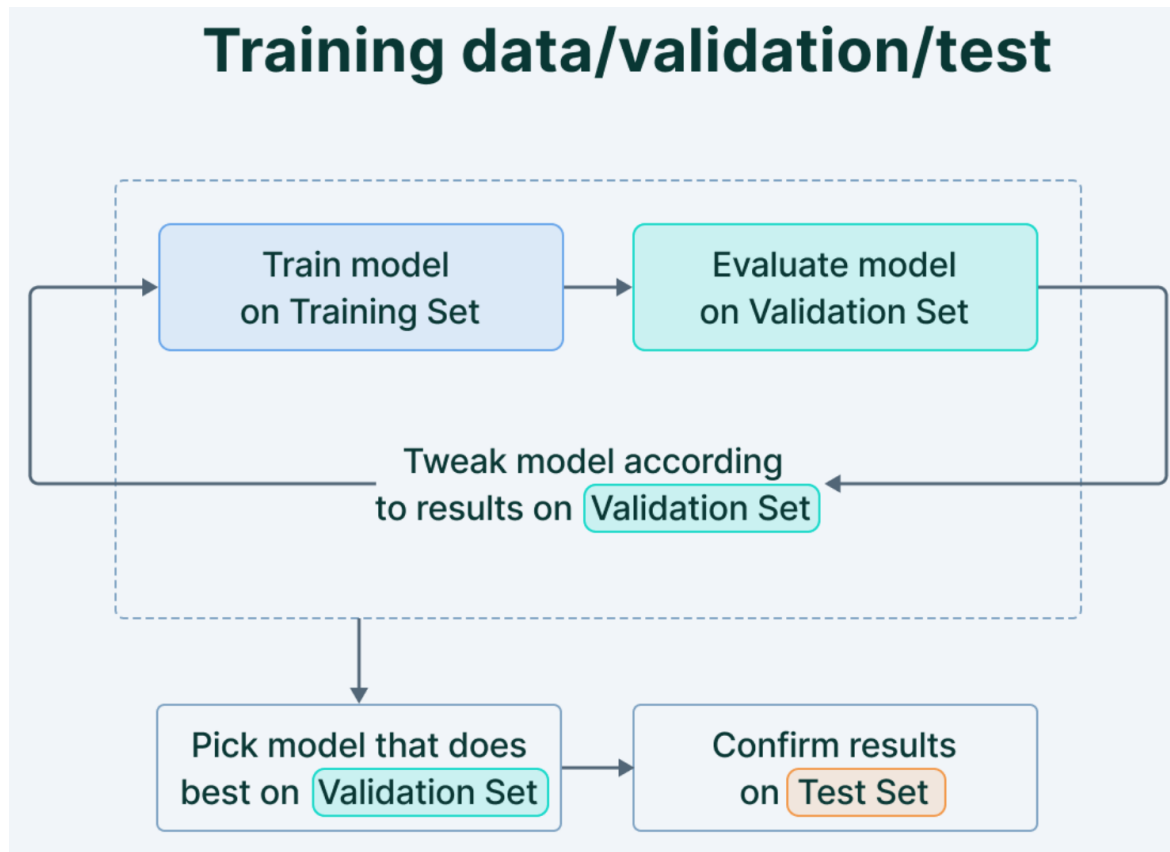
*May 23, 2024*

# Overview of tonight's lesson: ML 3

1. ML regression project!
2. Exploratory data analysis
3. Preprocessing pipeline
4. Random forest model

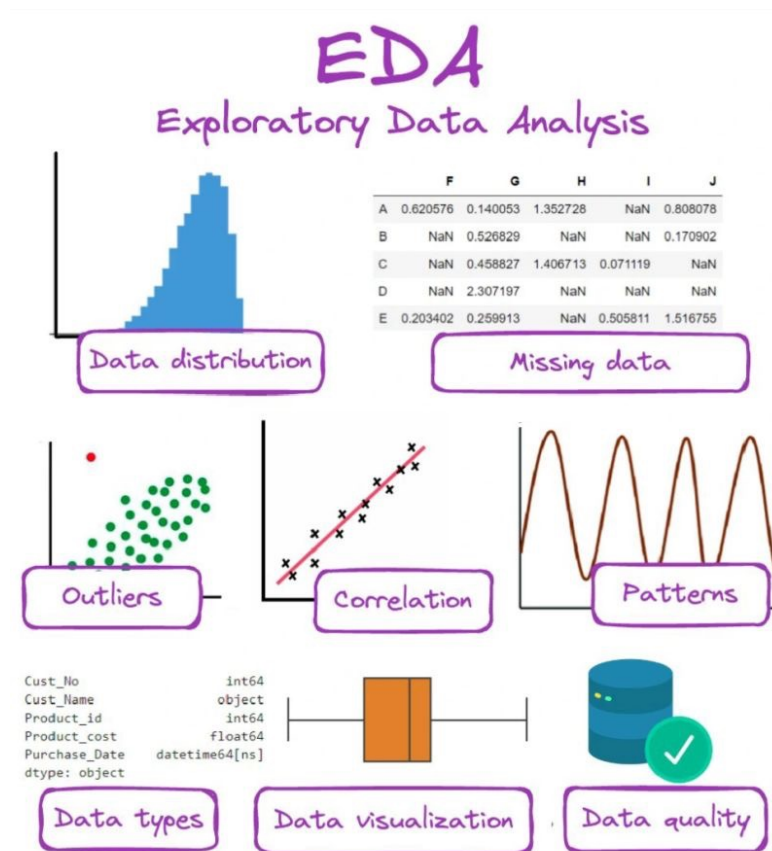


# Splitting data:



# Exploratory data analysis:

- Understand the structure and content of your data
- Discover potential outliers, missing values, inconsistencies
- Prepare data for ML analysis



# Preprocessing data:

- Transform raw data into a clean and usable format suitable for analysis or machine learning models

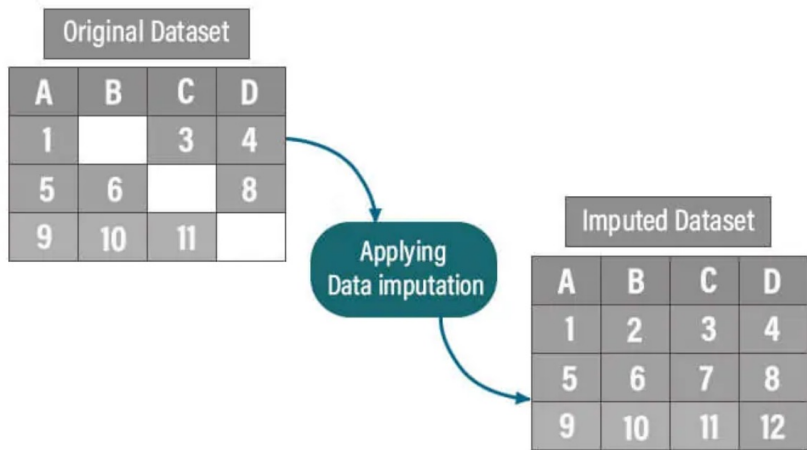
# Preprocessing data:

- Imputing the data (filling in the NaN values)
- Scaling the numerical data so that all of the features are in the same range
- Encoding the categorical variables

# Imputing the data:

- This replaces the missing values based on some method
- Methods can include using the mean/median/mode, KNN, or some other model
- Can also do this with categorical variables
- Imputation introduces some assumptions about the data!

## Data Imputation



# Scaling the numerical data :

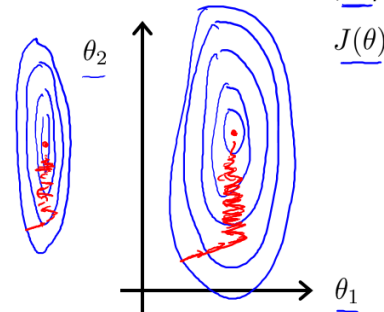
- Scale features to be in a similar range
- For example, can set the mean and standard deviation to zero

## Feature Scaling

Idea: Make sure features are on a similar scale.

E.g.  $x_1 = \text{size (0-2000 feet}^2\text{)}$  ←

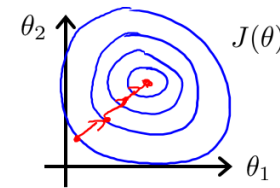
$x_2 = \text{number of bedrooms (1-5)}$  ←



$$\rightarrow x_1 = \frac{\text{size (feet}^2\text{)}}{2000} \quad \checkmark$$

$$\rightarrow x_2 = \frac{\text{number of bedrooms}}{5} \quad \checkmark$$

$$0 \leq x_1 \leq 1 \quad 0 \leq x_2 \leq 1$$

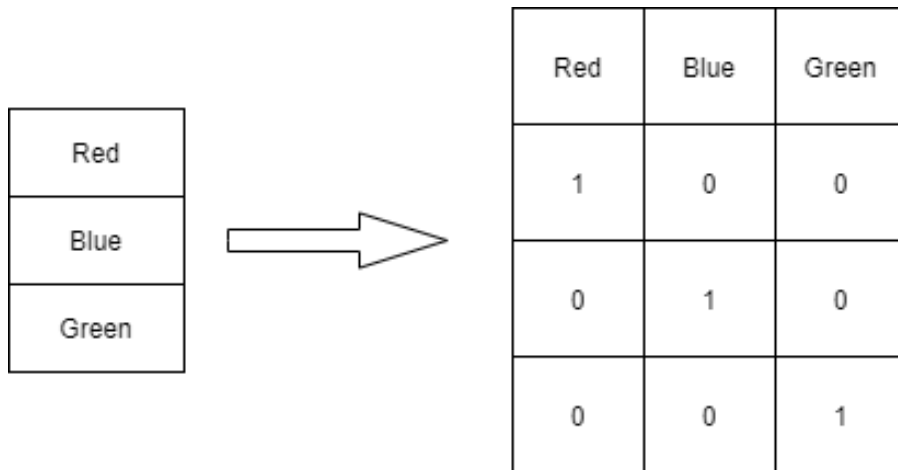


Andrew Ng



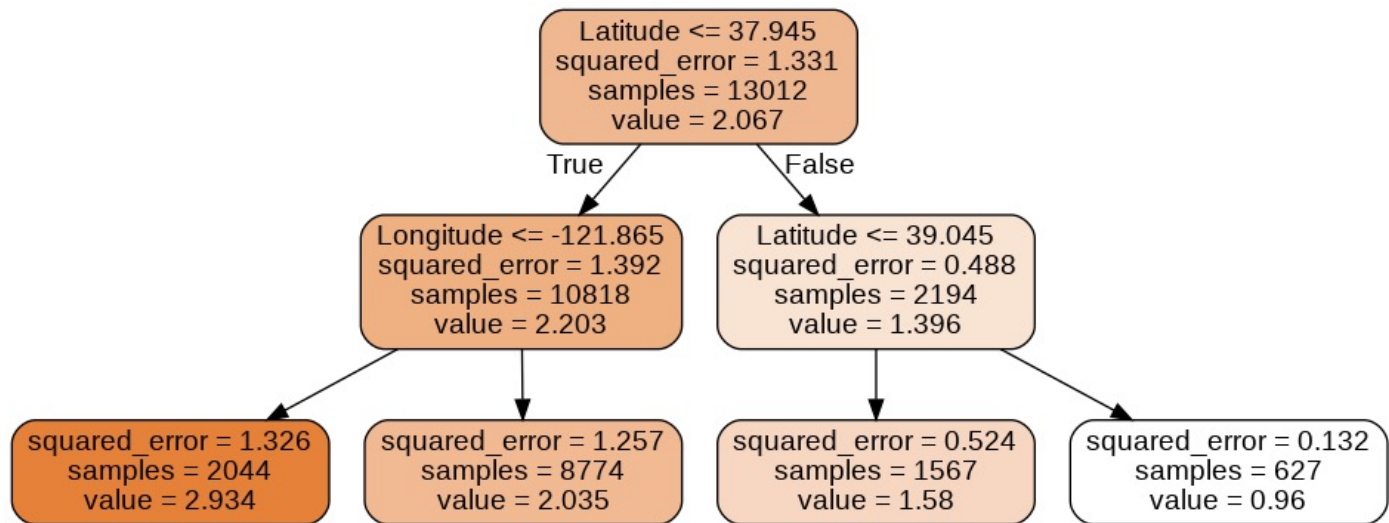
# Encoding:

- Encoding is the process of transforming categorical data (data with labels) into numerical representations.
- Need to think about: are the categories ordered or non-ordered?



# Trees and forests

- A decision tree in ML is essentially a set of if and else statements to subset data
- A random forest is a collection of decision trees where the if else statements are slightly different to get a more robust estimate



# Two important metrics for regression

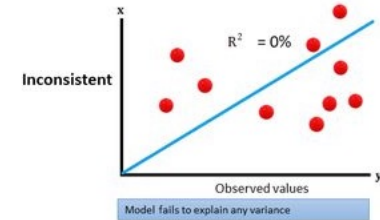
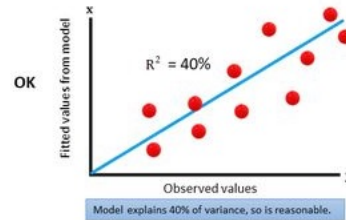
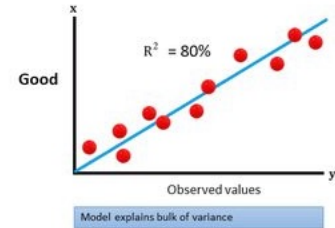
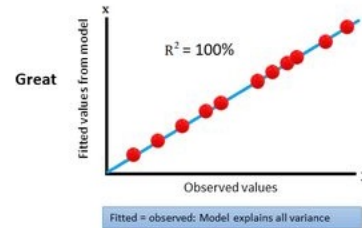
- Root mean squared error
- R square

$$RMSE = \sqrt{\sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n}}$$

$\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$  are predicted values

$y_1, y_2, \dots, y_n$  are observed values

$n$  is the number of observations



# NEXT STEPS FOR ML:

- ML&HPC!



[ace-net.ca](http://ace-net.ca)

[info@ace-net.ca](mailto:info@ace-net.ca)

[certificate@ace-net.ca](mailto:certificate@ace-net.ca)