

# Project 5: Making an AI

---

## 1. Table of Display

	Topic	Method
	What is ML	Summary
	Intro to Scikit-Learn	Summary w/Code
	Hyperparameters and Model Validation	?
	Feature Engineering	?
	DTrees and Random Forests	?
	Face Detection Pipeline	Final Project

## 2. Table Of Contents

1. [1. Table of Display](#)
2. [2. Table Of Contents](#)
3. [3. Background](#)
  1. [3.1. What is Machine Learning?](#)
    1. [3.1.1. Supervised learning](#)
      1. [3.1.1.1. Classification](#)
      2. [3.1.1.2. Regression](#)
    2. [3.1.2. Unsupervised Learning](#)
      1. [3.1.2.1. Clustering](#)
      2. [3.1.2.2. Dimensionality Reduction](#)
  2. [3.2. Scikit-Learn](#)
    1. [3.2.1. Data Representation](#)
    2. [3.2.2. Essentials of Scikit](#)
      1. [3.2.2.1. Using Scikit-Learn](#)
  3. [3.3. Model Validation](#)
  4. [Model Selection](#)
  5. [Hyperparameters](#)
    1. [Holdout Sets](#)
  6. [3.5. Feature Engineering](#)
  7. [3.6. Special Topic](#)
4. [4. Project: Face Detection Pipeline](#)
  1. [4.1. Output](#)
  2. [4.2. Code](#)

## 3. Background

### 3.1. What is Machine Learning?

Machine Learning (ML) is the art of designing mathematical models of data which can then be taught via tunable parameters. Because these algorithms are designed to assist in understanding data, there can be some debate of whether ML could be considered a branch of Artificial Intelligence (AI) anymore.

Because ML works with big data which vary greatly in both size, complexity there must also be various methods of analyzing this data. The two general categories of ML are **Supervised** and **Unsupervised** learning. We will explore these methods further in this paper.

### 3.1.1. Supervised learning

Supervised learning takes data and labels associated with the data to model their relationship. This ML model is used to apply labels to novel data. Supervised learning is commonly subdivided into **classification** and **regression**.

#### 3.1.1.1. Classification

Classification models use discrete categories such as a status. This type of model may be used to identify objects in an image or seek for when a part might be damaged within a piece of machinery.

This model will require the developer to first make a labeled dataset. Then you must design model inputs and the general assumptions you can provide. After this you can provide a set of model parameters which can be adjusted by the model during the training stage. The end result is that when introduced to novel data the classical model will provide a predictive label.

#### 3.1.1.2. Regression

Regression models place their categories in a continuous spectrum. This can be used to seek the probability of damage to a piece of machinery or seek more complex relationships such as genetics to regions.

Regression models could be treated as the opposite of [Dimensionality Reduction](#) models. They will extract a new unknown relationship and create a new dimension of labels.

### 3.1.2. Unsupervised Learning

Unsupervised learning results from not having a labeled dataset. This type of learning is used to find relationships within a data set. Unsupervised Learning can be further subdivided into **Clustering** and **Dimensionality Reduction** models.

#### 3.1.2.1. Clustering

Clustering models act similarly to classification models except they are seeking the distinct groups within the dataset. This can be used to seek out groups within massive datasets that humans may never see.

One method to do clustering is through the  $k$ -means model. This method finds the center of data clusters, it must be provided with a value  $k$  or this value may be tunable. This model seeks the position of centers that has the minimum distance between all points in the dataset.

#### 3.1.2.2. Dimensionality Reduction

Dimensionality Reduction as its name suggests is designed to simplify a data set into the smallest structure possible. This type of model could be used to seek important parameters to watch in larger systems which may take years for a human to analyze. This method allows us to infer new structures that may not be labels (or may not exist).

A dimensionality reduction model typically will remove one or more of the layers of data. This could be used to review a large dataset and place the data neatly into graphs that humans can easily understand, interpret and apply.

## 3.2. Scikit-Learn

While developing a ML algorithm from scratch is possible it could be cumbersome and likely has already been done. Thus to prevent reinventing the wheel, you can use the module **Scikit-Learn**. Scikit-Learn contains many of the popular algorithms in one consistent API.

### 3.2.1. Data Representation

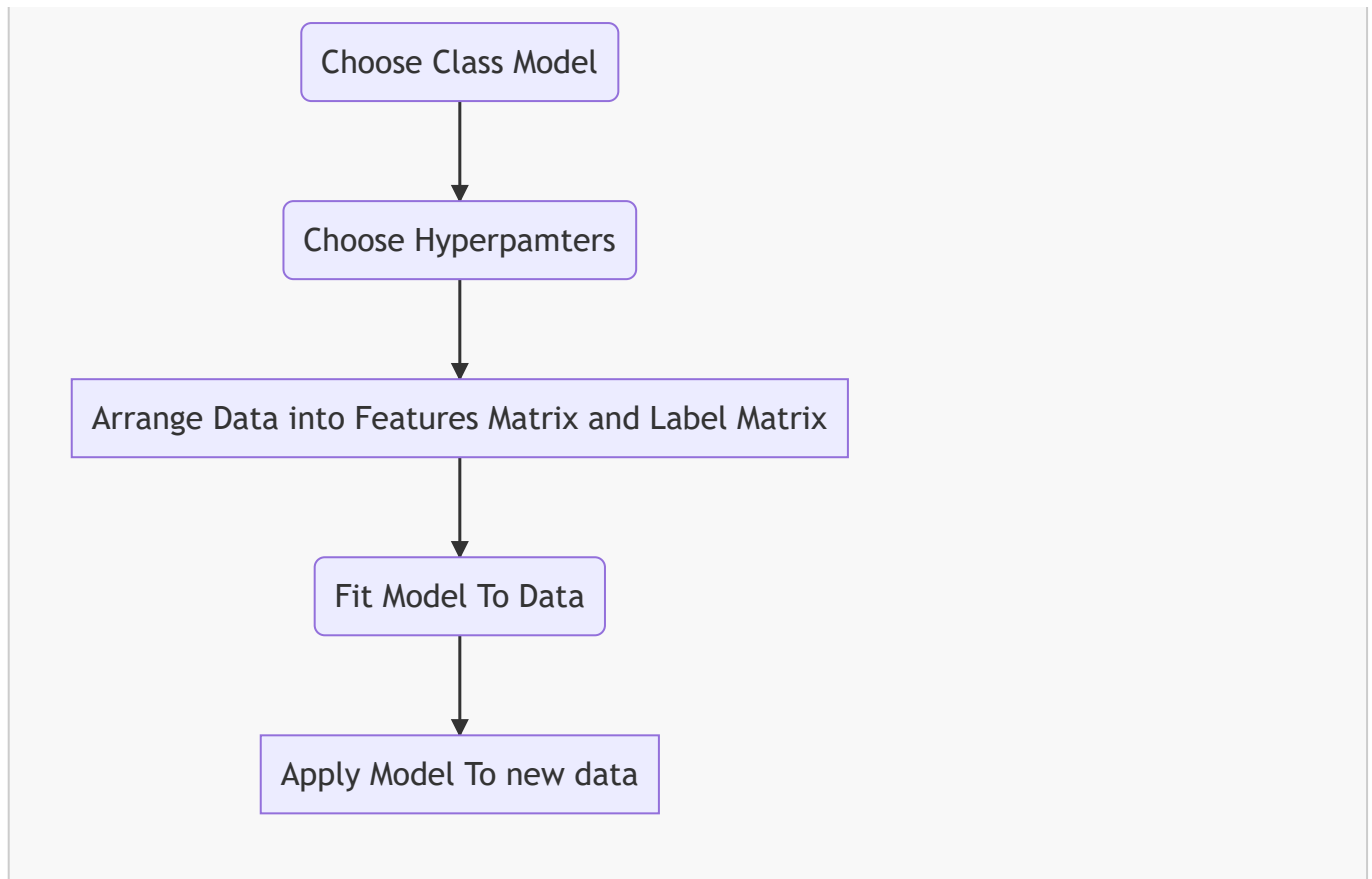
Scikit-Learn describes data with a *features matrix* this matrix can be stored in a NumPy array or Pandas DataFrame. The rows of the features matrix are often called samples and the columns are considered features. Thus the features matrix is considered to have the shape `[n_samples, n_features]`.

The features matrix is then coupled with a *label array*. This array is typically a Numpy array or Pandas Series of `n_samples`, but it can be two-dimensional with the shape `[n_samples, n_targets]`. The label array could be considered the dependent variable the model is attempting to predict.

### 3.2.2. Essentials of Scikit

Scikit is designed to a *consistent* API used for *inspection* of data. It also provides a simplistic *composition* through the simplification of complex ML algorithms into their foundational parts. It also uses common data structures and provides excellent user-oriented defaults.

This results in the following workflow for algorithm design.



#### 3.2.2.1. Using Scikit-Learn

Each model class is stored as a class within Scikit-Learn and are typically imported as a single object. For instance to import the *LinearRegression* module you would use

```
from sklearn.linear_model import LinearRegression
```

### 3.3. Model Validation

After selecting a model and selecting its hyperparameters you must then *validate* the model. The two methods of model validation we will discuss are **holdout sets** and **cross validation**. Once you have validated a model you must then select the best model by balancing the bias and variance of that model for the data you are using.

First let us setup a model to validate:

Model Selection

Hyperparameters

**Holdout Sets**

### 3.5. Feature Engineering

### 3.6. Special Topic

## 4. Project: Face Detection Pipeline

### 4.1. Output

### 4.2. Code