



Module 1: Steps of developing a Machine Learning Application

These steps form the backbone to any machine learning process and knowing them will make your life much easier when trying to build ML models

1. Data Collection

Machine learning requires training data, a lot of it. This data can either be labelled meaning Supervised Learning or not labelled meaning Unsupervised Learning.

Accuracy of the model depends on the quality and quantity of the data. The outcome of this step is generally a representation of data which will be used for training.

Using pre-collected data, by way of datasets from sites like Kaggle, UCI, etc. forms the basis of your Machine learning project. You may also collect data through user-surveys, analysis reports, trends, usage metrics, etc.

2. Data Preparation

We cannot work on raw data. **Data needs to be processed by normalization, removing duplicates, errors and biases.**

Visualising data can be helpful in **searching for patterns and outliers** to check if the data collected is right or if it contains missing values. This can be done using libraries like seaborn, matplotlib, etc. Visualize data to help detect relevant relationships between **variables or class imbalances, or perform other exploratory analysis.**

After performing **data wrangling**, we need to prepare the data for training. Cleaning of data is done that involves steps like removing duplicates, dealing with missing values, type conversions, correcting errors, normalizing the data, etc.

Not all the above steps are needed to be performed as it depends entirely on the data collected. **Some datasets may not require data preparation at all while for some data preparation step takes majority of their ML model build time.**



We can also Randomize data, which **erases the effects of the particular order in which we collected and/or otherwise prepared our data**. Later we can split the data into training, testing and evaluation sets

3. Choose a Model / Algorithm

The third step consists of **selecting the right model**. There are many models which can be used for many different purposes. Once the model is selected, it needs to meet the business goal.

We need to have an idea about the preparation the model requires along with its accuracy and scalability. **Having a complex model does not mean a better model**.

Common machine learning algorithms include Decision Trees, Random Forest, Linear Regression, Support Vector Machines (SVM), Logistic Regression, K-means, Principal Component Analysis (PCA), Naïve Bayes, and Neural Networks. Different algorithms need to be applied to different tasks, you need to choose the correct one for your use case.

4. Training the Model

Training a model forms the basis of machine learning. The goal is to use our training data and improve the predictions of our model.

Every cycle in training a model involves updating the **weights** and **biases** in each training step. We can use labelled sample data in case supervised machine learning and unlabelled sample data for unsupervised learning.

The goal of training is to evaluate and further improve our model accuracy and performance. Training happens in the form of iterations which is called a **training step**.

5. Evaluate the Model

After training the model comes evaluating the model. The larger the number of variables in the real world, the bigger the training and test data should be.

Performance metrics are used to measure the performance of the model. These include precision, recall, accuracy, specificity, etc.



The model is then tested against previously unseen data. The unseen data is meant to act as representative of model performance in the real world, but still helps tune the model (as opposed to test data, which does not).

A 70/30 split, or similar, is considered a good train/eval split, which depends on things like **data availability, dataset features, domain, etc.**

6. Parameter Tuning

The original model parameters need to be tested after evaluating your model. By increasing the training, it can lead to better results.

Parameter tuning is an experimental process and hence we need to define when to stop parameter tuning otherwise it will continue to tweak the model.

Hyperparameter tuning is an art and one that requires patience & experience. Once the model parameters are tuned it can give us better results. Some common hyperparameters include: **number of training steps, learning rate, initialization values and distribution, etc.**

7. Make Predictions

After the processes of collecting data, preparing the data, selecting a machine learning algorithm, training the model and evaluating the model & tuning the parameters, **we need to make predictions.**



Photo by Carlos Muza on Unsplash

Our machine learning model can make predictions ranging from image recognition to predictive analytics to natural language processing.

After building the model needs to be tested on a testing set to check how the model performs on unseen data. It helps to further evaluate the model and provides better approximation.