# Week - 4 Reading assignment

## Train-Test Split for Evaluating Machine Learning Algorithms

The train-test split is a technique for evaluating the performance of a machine learning algorithm.

It can be used for classification or regression problems and can be used for any supervised learning algorithm.

The procedure involves taking a dataset and dividing it into two subsets. The first subset is used to fit the model and is referred to as the training dataset. The second subset is not used to train the model; instead, the input element of the dataset is provided to the model, then predictions are made and compared to the expected values. This second dataset is referred to as the test dataset.

- **Train Dataset**: Used to fit the machine learning model.
- **Test Dataset**: Used to evaluate the fit machine learning model.

The objective is to estimate the performance of the machine learning model on new data: data not used to train the model.

The train-test procedure is appropriate when there is a sufficiently large dataset available.

If you have insufficient data, then a suitable alternate model evaluation procedure would be the k-fold cross-validation procedure.

You must choose a split percentage that meets your project's objectives with considerations that include:

- Computational cost in training the model.
- Computational cost in evaluating the model.
- Training sets representativeness.
- Test set representativeness.

**Latent Feature and Manifold**

It's not always possible to have examples for each variable that we are trying to model. Sometimes the unobserved variables make the learning more difficult. These models tend to explain the complex relationships between several variables by drawing simple relations between the unobserved and the variables. A latent variable model p is a probability distribution over 2 sets of variables x and z

$$p(x, z; \theta)$$

Where x -> these variables are observed during learning time in Dataset z -> they are never observed. These models can be directed and undirected. Ex: Gaussian Mixture Models These are the most widely used models in machine learning.

**Latent Variable Models are useful for these two reasons:**

1. One reason is that some variables could be naturally unobserved. Such as if we are doing some clinical research it is very much possible that we don't have samples for few of the cases. So this model helps to work with those missing data.

2. It is also very helpful to leverage the prior knowledge we have when defining the model. This way if we already have the domain knowledge we will be able to design the model to capture those values.
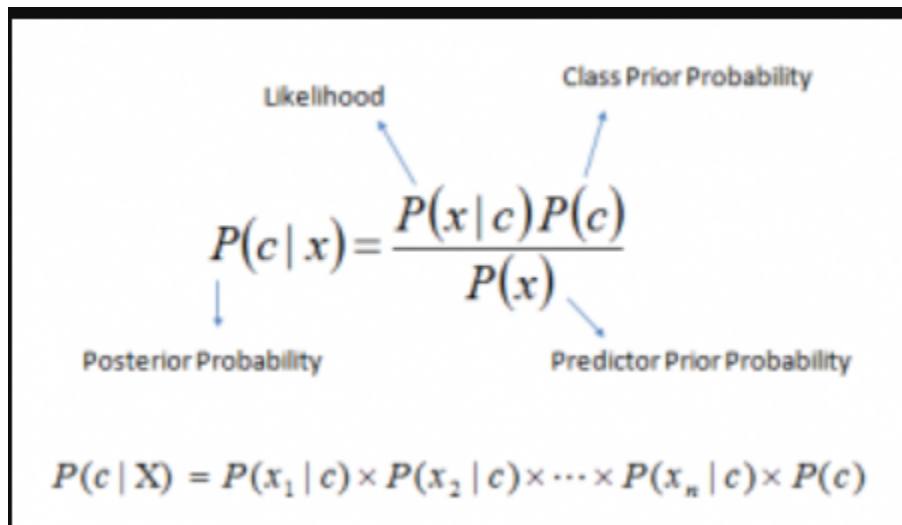
**Linear Factor Analysis:** Since the idea behind Linear LNF is to draw the relationship between variables in X and the unobserved in latent Y plus noise, here the number of factors q is less than p. The aim is to determine the smallest value of q for which the model is adequate.

*Naive Bayesian algorithm:*

It is a classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

The Naive Bayes model is easy to build and particularly useful for very large data sets. Along with simplicity, Naive Bayes is known to outperform even highly sophisticated classification methods.

Bayes theorem provides a way of calculating posterior probability $P(c|x)$ from $P(c)$, $P(x)$ and $P(x|c)$. Look at the equation below:



$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \cdots \times P(x_n|c) \times P(c)$$

**Linear Regression:**

Linear Regression is **a machine learning algorithm based on supervised learning**. It performs a regression task. Regression models a target prediction value based on independent variables. It is mostly used for finding out the relationship between variables and forecasting.

Linear regression performs the task to predict a dependent variable value (y) based on a given independent variable (x). So, this regression technique finds out a linear relationship between x (input) and y(output). Hence, the name is Linear Regression.

**PCA:**

Principal Component Analysis, or PCA, is a dimensionality-reduction method that is often used to reduce the dimensionality of large data sets, by transforming a large set of variables into a smaller one that still contains most of the information in the large set.

Reducing the number of variables of a data set naturally comes at the expense of accuracy, but the trick in dimensionality reduction is to trade a little accuracy for simplicity. Because smaller data sets are easier to explore and visualize and make analyzing data much easier and faster for machine learning algorithms without extraneous variables to process.

So to sum up, the idea of PCA is simple — reduce the number of variables of a data set, while preserving as much information as possible.

**Logistic regression:**

- Logistic regression is one of the most popular Machine Learning algorithms, which comes under the Supervised Learning technique. It is used for predicting the categorical dependent variable using a given set of independent variables.

- Logistic regression predicts the output of a categorical dependent variable. Therefore the outcome must be a categorical or discrete value. It can be either Yes or No, 0 or 1, true or False, etc. but instead of giving the exact value as 0 and 1, it gives the probabilistic values which lie between 0 and 1.

- Logistic Regression is much similar to the Linear Regression except that how they are used. Linear Regression is used for solving Regression problems, whereas Logistic regression is used for solving the classification problems.

- In Logistic regression, instead of fitting a regression line, we fit an "S" shaped logistic function, which predicts two maximum values (0 or 1).

- The curve from the logistic function indicates the likelihood of something such as whether the cells are cancerous or not, a mouse is obese or not based on its weight, etc.

- Logistic Regression is a significant machine learning algorithm because it has the ability to provide probabilities and classify new data using continuous and discrete datasets.

- Logistic Regression can be used to classify the observations using different types of data and can easily determine the most effective variables used for the classification.