# Experiment . No. 09

**Title :** Given a bank customer, build a neural network-based classifier that can determine whether they will leave or not in the next 6 months.

**Objective :** To distinguish the feature and target set and divide the data set into training and test sets and normalize them and students should build the model on the basis of that

## Dataset Description :

The dataset contains 10,000 sample points with 14 distinct features such as ~~customfield~~ customerID, Gender, creditscore, Geography, Age, Tenure, Balance, etc
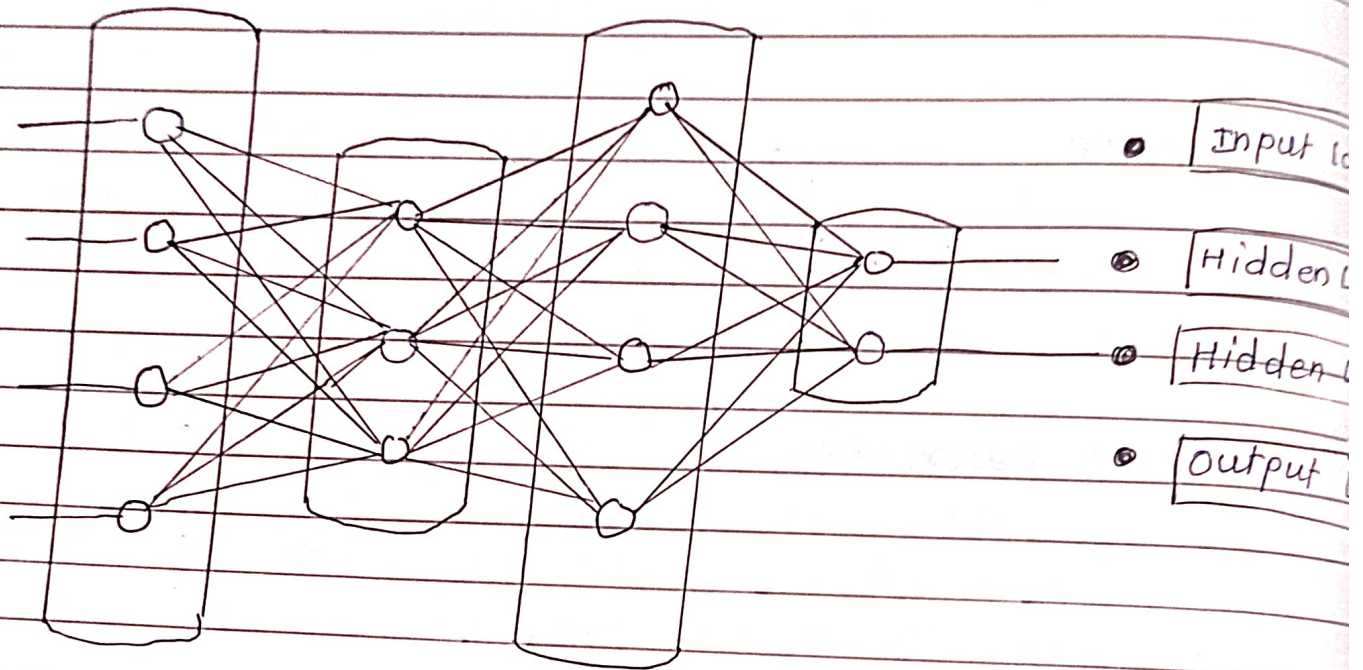
## Theory :

### Artificial Neural Network :

The term Artificial Neural network is derived from Biological neural networks that develop the structure of a human brain. similar to the human brain that has neurons interconnected to one another, artificial neural networks also have neurons that are interconnected to one another in various layers of the networks. These neurons are known as nodes

### The Architecture of an artificial neural network.

In order to define a neural network that consists of a large number of artificial neurons, which are termed units arranged in a sequence of layers.

Artificial neural network primarily consists of three layers



- ● Input la
- ◉ Hidden L
- ◉ Hidden
- ● Output

## Input Layer:

As the name suggests, it accepts inputs in several different formats provided by the programm

## Hidden Layer:

The hidden layer presents in-between input an output layers. It performs all the calculations to hidden features & patterns.

## Output Layer:

The input goes through a series of transformations using the hidden layer, which finally results in output that is conveyed using this layer.

$$\sum_{i=1}^{n} w_i * x_i + b$$

# keras :

keras is an open-source high level Neural network library, which is written in Python is capable enough to run on Theano, TensorFlow or CNTK. It was developed by one of the Google engineers, Francois chollet. It is made user-friendly, extensible and modular for facilitating faster experimentation with deep neural networks. It is not only supports Convolutional Networks and Recurrent Networks individually but also their combination.

# Tensorflow :

TensorFlow is a Google product, which is one of the most famous deep learning tools widely used in the research area of machine learning and deep learning. It came into the market on $9^{th}$ nov. 2015 under the Apache License 2.0. It is built in such a way that it can easily run on multiple CPUs and GPUs as well as on mobile operating systems. It consists of various wrappers in distinct languages such as Java, c++ or Python

# Normalization :

Normalization is a scaling technique in machine Learning applied during data preparation to change the value of numeric columns in the dataset to use a common scale. It is not necessary for all datasets in a model. It is required only when features of machine learning models have different ranges.

$$Xn = (x - x minimum) / (x maximum - x minimum)$$

where

$Xn$ = value of Normalization

$X maximum$ = maximum value of a feature

$X minimum$ = minimum value of a feature

# Normalization techniques in ML :

## ① min-max scaling -

This technique is also referred to as scal the min-max scaling method helps the dataset t shift and rescale the values of their attribut so they end up ranging between 0 & 1.

## ② standardization scaling :

Standardization scaling is also known as z-score normalization, in which values are center around the mean with a unit standard deviation we can calculate standardization by subtracting the feature value from the mean and dividing it by standard deviation

$$X' = \frac{X - \mu}{6}$$

Here, $\mu$ represents the mean of feature value, and $6$ represents the standard deviation of feature values

## confusion Matrix :

The confusion matrix is a matrix used to determine the performance of the classification models for a given set of test data. It can only be determined the performance of the classification models for the true values for test data are known. It shows the errors in the model performance in the form of a matrix, hence also known as a error metrix

|  | Actual values | |
|---|---|---|
| **predicted values** | Positive(1) | Negative(0) |
| Positive(1) | TP | FP |
| Negative(0) | FN | TN |

The above table has the following cases

* **True Negative :** model has given prediction No, and the real or actual value was also No.

* **True positive :** The model has predicted yes, and the actual value was also true

* **False Negative :** The model has predicted No, but the actual value was yes, it is also called as Type-II error

* **False Positive :** The model has predicted Yes, but the actual value was No, it is also called a Type-I error.


Calculations using confusion matrix

**Accuracy :** It defines how often the model predicts the correct output.

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

**Error rate :** It defines how often the model gives the wrong predictions.

$$Error\ rate = \frac{FP + FN}{TP + FP + FN + TN}$$

**Precision :** It can be defined as the number of corr
outputs provided by the model or out of all posit
classes that have predicted correctly by the model
how many of them were actually true.

$$precision = \frac{TP}{TP+FP}$$

**Recall :** It is defined as the out of total positive
how our model predicted correctly

$$Recall = \frac{TP}{TP+FN}$$

**F-measure —** If two models have low precision and high
recall or vice versa. so for this purpose, we
use F-score. This score helps us to evaluate the
recall & precision at the same time

$$F-Measure = \frac{2*Recall*precision}{Recall+Precision}$$

**Null-error rate :** It defines how often our model would
incorrect if it always predicted the majority class

**Roc curve :** It is a graph displaying a classifier's perfor
-ance for all possible thresholds. graph is plotted betwee
the true positive rate and false positive rate

**Conclusion :**

In this way, we build a neural n/w based classifier
that can determine whether they will leave or not in next
6 months.