

LAB-06:-

A Random Forest Algorithm:-

- For each tree $t=1$ to T ,
 - Draw a boot strap sample D_t of size n from training data D (sam with replacement)
 - Train a decision tree h_t on D_t at each node.
 - Randomly select m features from the full feature set.
 - Choose the best feature & split point among m features using Impurity, Entropy or MSE (for regression)
 - Split the node into child nodes.

At Aggregate the predictions:

For Classification: Use majority voting over all trees:

$$\hat{y} = \text{mode}(h_1(x), h_2(x), \dots, h_T(x))$$

For regression: Use mean prediction:

$$\hat{y} = \frac{1}{T} \sum_{t=1}^T h_t(x)$$

#Code

```
class import RandomForest:
    def __init__(self, n-trees=10, max-feature
                  sqrt'):
        self.n-trees = n-trees
        self.max-features = max-features
        self.trees = []
        def boot strap sample (self, x, y):
            n-sample = x.shape[0]
```



```

size = n-sample . replace = True
return x[indices], y[indices]
def get_max_features(self, n_features):
    if self.max_features == 'sqrt':
        return int(np.sqrt(n_features))
    elif isinstance(self.max_features, int):
        return self.max_features
def fit(self, X, y):
    self.tree = []
    n_features = X.shape[1]
    max_feature = self.get_max_features

```

★ => K-Means Algorithm:-

01. Initialize centroids:

Randomly choose K data points from x as initial cluster centroids:

$$u_1, u_2, \dots, u_K$$

02. Repeat until convergence.

@ Repeat each data point to the nearest centroid. For each point x_i , find the closest centroid u_j based on distance

$$u_j = \frac{1}{|C_i|} \sum_{x_i \in C_i} x_i$$

03. check for convergence

- * if cluster assignments don't change
- * Centroids don't move significantly, then stop

Code:-

```
import matplotlib.pyplot as plt
x, _ = make_blobs (n-sample=300,
center=3, random_state=42)
model = KMeans (n-cluster=3)
model.fit(x)
plt.scatter(x[:,0], x[:,1])
(= model.labels_, cmap= viridis)
plt.show()
```

Principle component Model Algorithm

- 1) Standardize the data
PCA is affected by scale of feature.
∴ data should be standardized

$$\hat{x}_j = \frac{x_j - \mu_j}{\sigma_j}$$

- 2) Compute the Covariance matrix

$$\Sigma = \frac{1}{n-1} x^T x$$

- 3) Compute the Eigen value and Eigen value of Covariance matrix.

$$\Sigma v = \lambda v$$

- 4) Sort the eigen values & vectors in descending order.

- ⑤ select the top i.e eigen vectors
- ⑥ project the data onto new space
- ⑦ Output the transformed data

code

```
data = load - iris()
x = data.data
pca = PCA(n-components = 2)
x_pca = pca.fit_transform(x)

x_pca = pca.fit_transform(x)
plt.scatter(x_pca[:, 0], x_pca[:, 1],
            c = data.target)
print('Explained variance by each
      component: {3}')
```