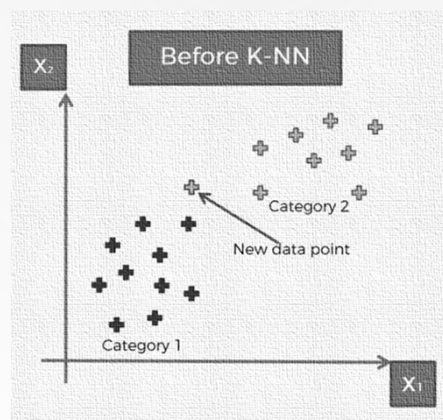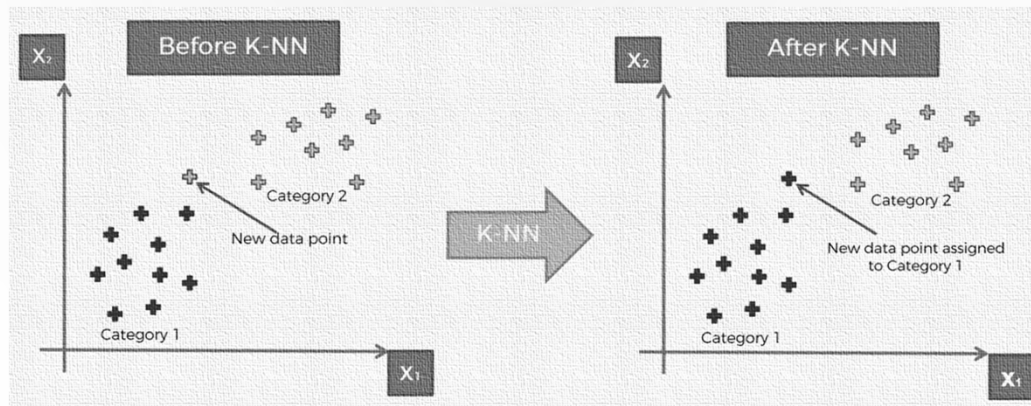# KNN – K NEAREST NEIGHBOR

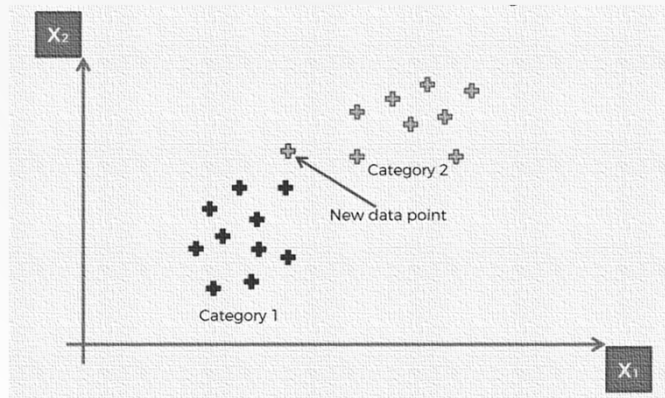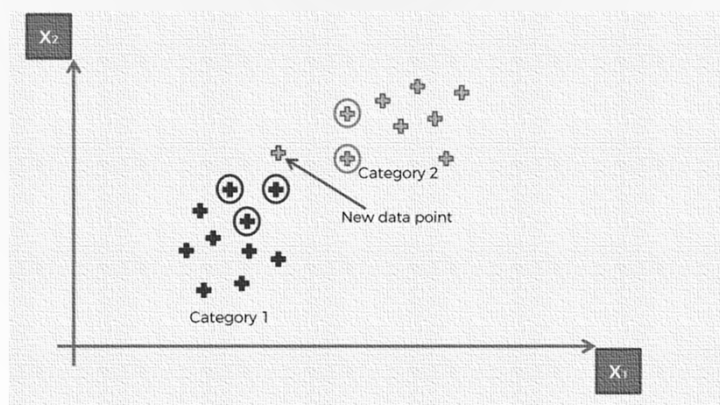MOHAN M J

---

## INTUITION

# CONTD..



# STEPS

1. Choose the number K for neighbors

2. Take the K nearest neighbors of the new data point according to the Euclidean distance

3. Among the K neighbors count the number of data points in each category

4. Assign the new data point to the category where we counted the most neighbors

# STEP 1: CHOOSE THE NUMBER OF NEIGHBORS
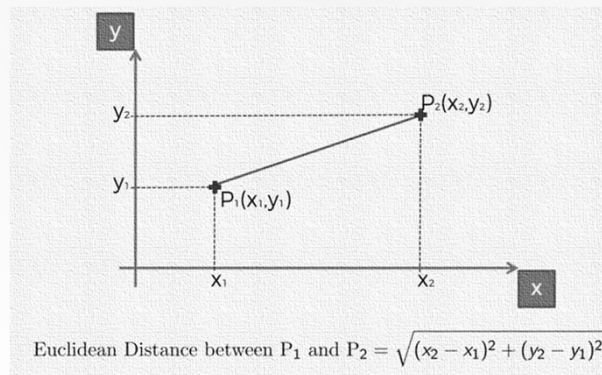
K = 5



# STEP 2: Take the K=5 nearest Neighbors of the new data point according to the Euclidian distance

# EUCLIDEAN DISTANCE

$$d(x, x') = \sqrt{(x_1 - x'_1)^2 + (x_2 - x'_2)^2 + \ldots + (x_n - x'_n)^2}$$



Euclidean Distance between $P_1$ and $P_2 = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

# STEP 3: Among the K neighbors, count the number of points in each category



Category 2

New data point

Category 1

Category 1: 3 neighbors
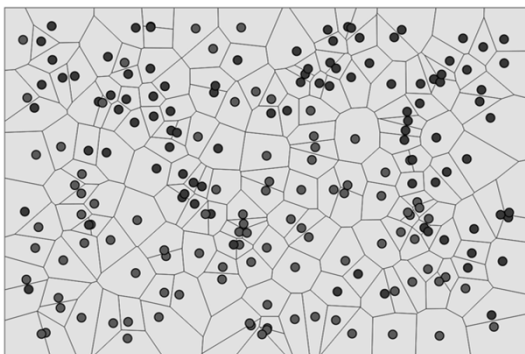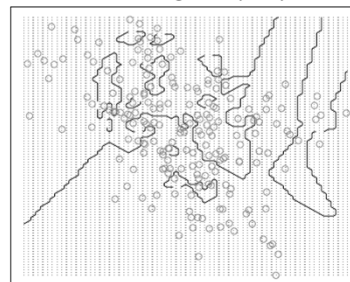
Category 2: 2 neighbors
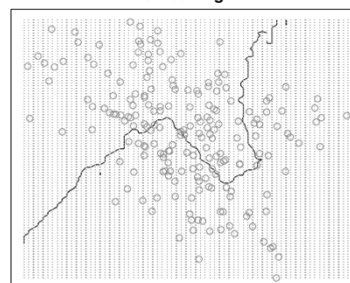
# ASSIGN TO CATEGORY



# MORE ON K

When K is small, we are restraining the region of a given prediction and forcing our classifier to be "more blind" to the overall distribution. A small value for K provides the most flexible fit, which will have low bias but high variance

# IRIS FLOWER PROBLEM

- Three species of Iris (Iris setosa, Iris virginica and Iris versicolor)
- Four features were measured from each sample: the length and the width of the sepals and petals.



# PYTHON CODE:

```
# Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
# Importing the dataset
dataset = pd.read_csv('Social_Network_Ads.csv')
dataset
X = dataset.iloc[:, [2, 3]].values
y = dataset.iloc[:, 4].values
```

## PYTHON CODE:

```python
# Splitting the dataset into the Training set and Test set
from sklearn.cross_validation import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25,
                                                    random_state = 0)
# Feature Scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
```

## PYTHON CODE:

```python
# Fitting classifier to the Training set
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors=5, metric='minkowski', p=2)
classifier.fit(X_train,y_train)
# Predicting the Test set results
y_pred = classifier.predict(X_test)
# Making the Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
cm
```
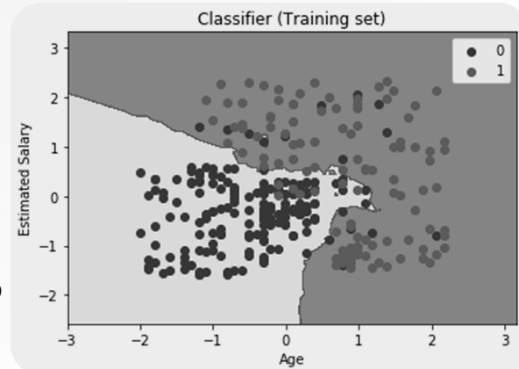
# VISUALIZING RESULTS

```
X_set, y_set = X_train, y_train

X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1,
stop = X_set[:, 0].max() + 1, step = 0.01),

                    np.arange(start = X_set[:, 1].min() - 1,
stop = X_set[:, 1].max() + 1, step = 0.01))

plt.contourf(X1, X2, classifier.predict(np.array([X1.ravel(),
X2.ravel()]).T).reshape(X1.shape),

            alpha = 0.75, cmap = ListedColormap(('pink',
'green')))

plt.xlim(X1.min(), X1.max())

plt.ylim(X2.min(), X2.max())

for i, j in enumerate(np.unique(y_set)):

    plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],

            c = ListedColormap(('r', 'green'))(i), label = j)
```

```
plt.title('Classifier (Test set)')
plt.xlabel('Age')
plt.ylabel('Estimated Salary')
plt.legend()
plt.show()
```



# VISUALISING RESULTS

```
from matplotlib.colors import ListedColormap

X_set, y_set = X_test, y_test

X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1,
stop = X_set[:, 0].max() + 1, step = 0.01),

                    np.arange(start = X_set[:, 1].min() - 1,
stop = X_set[:, 1].max() + 1, step = 0.01))

plt.contourf(X1, X2, classifier.predict(np.array([X1.ravel(),
X2.ravel()]).T).reshape(X1.shape),

            alpha = 0.75, cmap = ListedColormap(('pink',
'green')))

plt.xlim(X1.min(), X1.max())

plt.ylim(X2.min(), X2.max())

for i, j in enumerate(np.unique(y_set)):

    plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],

            c = ListedColormap(('red', 'green'))(i), label =
j)
```

```
plt.title('Classifier (Test set)')
plt.xlabel('Age')
plt.ylabel('Estimated Salary')
plt.legend()
plt.show()
```

THANK YOU