

```
In [1]: # https://www.kaggle.com/code/samuelcortinhas/k-nearest-neighbours-knn-from-scratch/notebook
```

```
import numpy as np
import pandas as pd
import seaborn as sns
sns.set(style='darkgrid', font_scale=1.4)
import matplotlib.pyplot as plt
%matplotlib inline
from sklearn.model_selection import train_test_split
from scipy import stats
import time
```

```
In [2]: # Load data
data = pd.read_csv("D:/Bachelor of Technology/6th Semester/8. UCSE673 Machine Learning Lab/10th_class/knn.csv", index_c
data.drop('Unnamed: 32', axis=1, inplace=True)

# Preview data
print('Dataframe shape:', data.shape)
data.head(3)
```

Dataframe shape: (569, 31)

```
Out[2]:
```

	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean
0	M	17.99	10.38	122.8	1001.0	0.11840	0.27760	0.3001	0.14710
1	M	20.57	17.77	132.9	1326.0	0.08474	0.07864	0.0869	0.07017
2	M	19.69	21.25	130.0	1203.0	0.10960	0.15990	0.1974	0.12790

3 rows × 31 columns

```
In [3]: # Features and Labels
X = data.drop('diagnosis', axis=1)
y = data['diagnosis']

# Encode target to binary
y = (y=='M').astype('int')
```

```
In [4]: # Split dataset into training and testing data (80/20 split)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=4)
```

```
In [5]: class kNN():
    '''k-Nearest Neighbours'''
    # Initialise
    def __init__(self, k=3, metric='euclidean', p=None):
        self.k = k
        self.metric = metric
        self.p = p

    # Euclidean distance (L2 norm)
    def euclidean(self, v1, v2):
        return np.sqrt(np.sum((v1-v2)**2))

    # Manhattan distance (L1 norm)
    def manhattan(self, v1, v2):
        return np.sum(np.abs(v1-v2))

    # Minkowski distance (Lp norm)
    def minkowski(self, v1, v2, p=2):
        return np.sum(np.abs(v1-v2)**p)**(1/p)

    # Store train set
    def fit(self, X_train, y_train):
        self.X_train = X_train
        self.y_train = y_train

    # Make predictions
    def predict(self, X_test):
        preds = []
        # Loop over rows in test set
        for test_row in X_test:
            nearest_neighbours = self.get_neighbours(test_row)
            majority = stats.mode(nearest_neighbours)[0][0]
            preds.append(majority)
        return np.array(preds)
```

```

# Get nearest neighbours
def get_neighbours(self, test_row):
    distances = list()

    # Calculate distance to all points in X_train
    for (train_row, train_class) in zip(self.X_train, self.y_train):
        if self.metric=='euclidean':
            dist = self.euclidean(train_row, test_row)
        elif self.metric=='manhattan':
            dist = self.manhattan(train_row, test_row)
        elif self.metric=='minkowski':
            dist = self.minkowski(train_row, test_row, self.p)
        else:
            raise NameError('Supported metrics are euclidean, manhattan and minkowski')
        distances.append((dist, train_class))

    # Sort distances
    distances.sort(key=lambda x: x[0])

    # Identify k nearest neighbours
    neighbours = list()
    for i in range(self.k):
        neighbours.append(distances[i][1])

    return neighbours

```

```

In [6]: # Function to calculate accuracy
def accuracy(preds, y_test):
    return 100 * (preds == y_test).mean()

# Apply our kNN algorithm
for metric in ['euclidean', 'manhattan']:
    clf = kNN(k=5, metric=metric)
    clf.fit(X_train.values, y_train.values)
    preds = clf.predict(X_test.values)
    print(f'Metric: {metric}, accuracy: {accuracy(preds, y_test):.3f} %')

```

Metric: euclidean, accuracy: 87.719 %
Metric: manhattan, accuracy: 91.228 %

```

In [7]: from sklearn.neighbors import KNeighborsClassifier

clf = KNeighborsClassifier(n_neighbors=5)
clf.fit(X_train.values, y_train.values)
preds = clf.predict(X_test.values)

# Calculate test set accuracy
def accuracy(preds, y_test):
    return 100 * (preds == y_test).mean()
print(f'Sklearn accuracy: {accuracy(preds, y_test):.3f} %')

```

Sklearn accuracy: 87.719 %