

KNN (K nearest neighbours)

MODULES

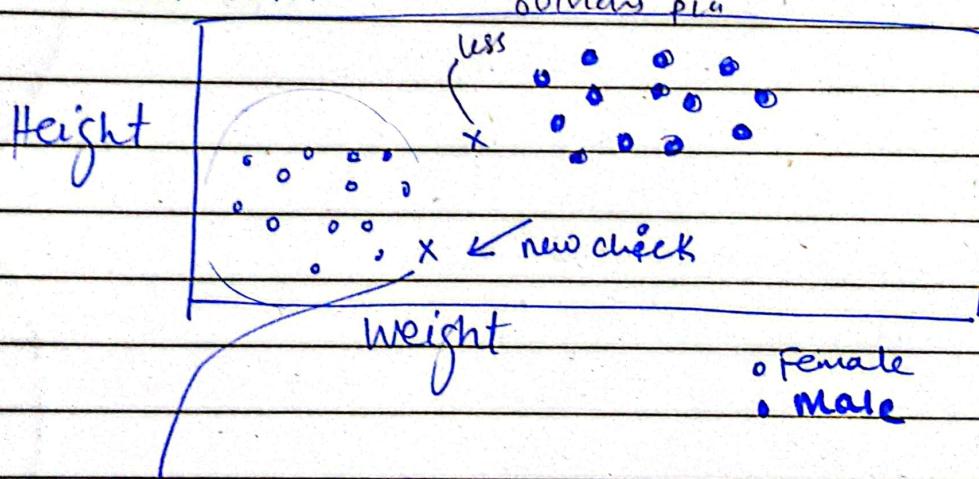
DATE: 24/12/2023

- KNN (K nearest neighbours) is one of the simplest algorithms we will learn about!
 - Section Overview:
 - KNN Theory and Intuition
 - KNN Classification Coding Example
 - KNN exercise overview
 - KNN Exercise Solution:-
 - While KNN can be used for regression tasks, its performance can be quite poor and less efficient than other algorithm, so we've decided not to exhibit its use for regression.
 - However if you do want to use it for regression it is very easy to swap in the KNN Regressor model with scikit-learn
 - K mean algorithm is different than KNN
- ISLR Relevant Reading:-
- Chapter 2

2- KNN Classification Theory and Intuition

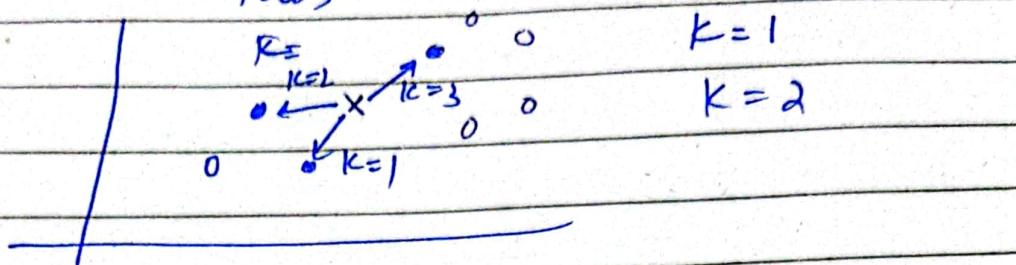
- K nearest neighbours is one of the simplest machine learning algorithms.
- It's simply assign a label to new data based on the distance between the old data and new data.
- Let's go through the intuition with an example use case.
- Sexing chicks is still a very manual process.

- Imagine a height and weight datasets.



- We intuitively "know" this is likely female - Intuition comes from distance to points.
- What about less obvious points, how many points should be consider

MATERIALS DATE: _____
Let's imagine a situation like this



when $K=4,2$ there could be a possibility tie how to chose when this condition occurs.

Tie consideration and options:

- Always chose an odd tie-
- In case of tie, simply reduce K by 1 until tie is broken-
- Randomly break tie
- choose the nearest class

What does scikit-learn do in case of tie?

- Warning: Regarding the Nearest Neighbours algorithms, if it is found that two neighbour, neighbours $K+1$ and K , ~~are~~ ^{have} identical distances but different labels, the result will depend on the order of the training data -

- In case of tie, the answer will be the class that happens to appear first in the sets of neighbour.
- Results are ordered by distance, so it choose the class of the closest point.

M T W T F S

DATE: _____

- We Want K value that minimizes error: $\text{Error} = 1 - \text{Accuracy}$

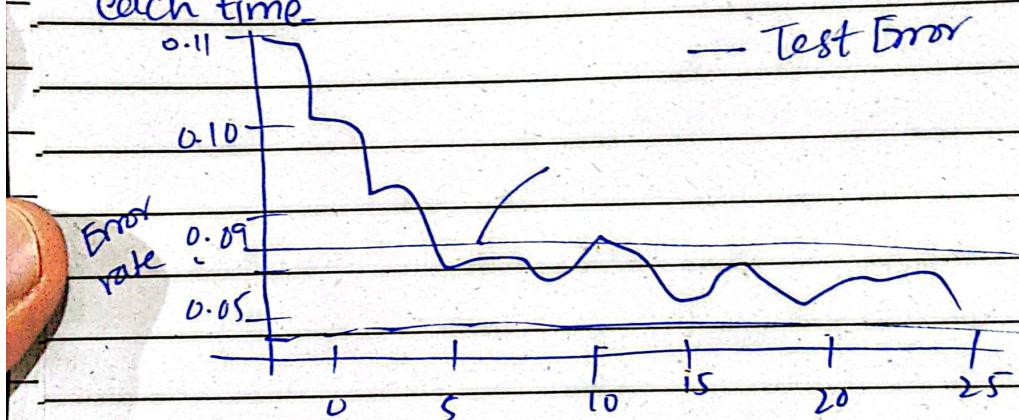
- Two methods:

- Elbow methods.

- Cross validation a grid search of multiple K value and chose K that result in lowest error or highest accuracy.

- Elbow Method:-

Train your model on multiple K-values and plot the error ~~on test~~ after testing each time.



- Cross validation only takes into account the K value with the lowest error rate across multiple folds.
- This could result in a more complex model (higher value of K).
- Consider the context of problem to decide if K values are an issue.

- KNN Algorithm
 - choose K values
 - Sort feature vectors (N-dimensional space) by a distance metric.
 - Choose class based on K nearest feature vectors.

- KNN considerations:-

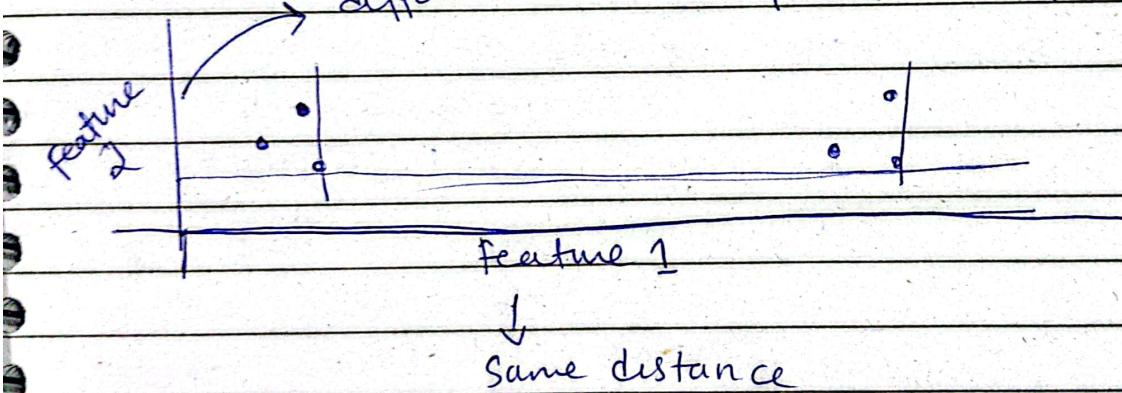
- Distance Metrics

- Many ways to measure distance:

- Minkowski
- Euclidean
- Manhattan
- Chebychev

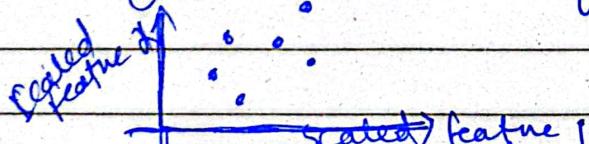
- Scaling for Distance

different distance for feature 2



- Scaling for DB feature could have vastly different ranges!

- Scaling is a ^{value} necessary for KNN.



- While KNN algorithm is relatively simple, keep in mind the following considerations:
 - choosing the optimal K value.
 - Scaling features.
- Let's continue to explore how to explore KNN for classification-

03. KNN Classification

Coding part 1: Data and Model

Coding

Import numpy as np
import pandas as pd

df = pd.read_csv("gene_expression")

df.head()

sns.scatterplot(data=df, x="Gene One",
y="Gene Two",
hue="Type",
cancer points
alpha=0.8)

plt.xlim(2, 4)

style="cancer-pew"

plt.ylim(4, 6)

sns.pairplot(data=df, hue="Patient
Group Cancer")

```
from sklearn.model_selection import train_
from sklearn.preprocessing import StandardScaler
x = df.drop("patient_cancer", axis=1)
y = df["patient_cancer"]
scaler = train_test_split(x, y, test_size=0.2, random_state=42)
Scalor = StandardScaler()
Scaled_X_train = Scalor.fit_transform(X_train)
Scaled_X_test = Scalor.transform(X_test)
from sklearn.neighbors import KNeighborsClassifier
help(KNeighborsClassifier)
knn_model = KNeighborsClassifier(n_neighbors=1)
knn_model.fit(Scaled_X_train, y_train)
y_preds = knn_model.predict(Scaled_X_test)
from sklearn.metrics import confusion_
matrix,
classification_report
confusion_matrix(y_test, y_preds)
print(classification_report(y_test, y_preds))
df['cancer present'].value_counts()
```

04

KNN Classification:

Coding part Two: choosing K

- A Pipeline object in scikit-learn can set up a sequence of repeated operations, such as a scaler and a model.
- This way only the pipeline needs to be called, instead of having to repeatedly call a scaler and a model.

1

test_error_rates = [];

for k in range(1, 30):

knn_model = KNeighborsClassifier(n_neighbors=k)

knn_model.fit(scaled_X_train, y_train)

y_preds_test = knn_model.predict(scaled_X_test)

test_error = 1 - accuracy_score(y_test, y_preds_test)

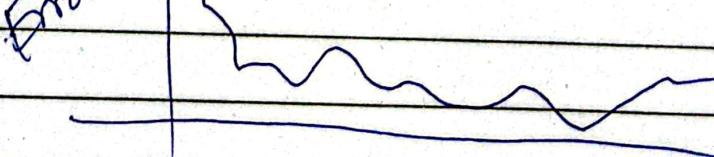
test_errors.append(test_error)

plt.plot(range(1, 30), test_error_rates)

plt.xlabel("K")

plt.ylabel("Error")

Error



Setting up pipeline

```
from sklearn.preprocessing import StandardScaler
```

```
knn = KNeighborsClassifier()
```

```
knn.get_params().keys()
```

```
operations = [ ('scaler', StandardScaler()), ('knn', KNeighborsClassifier()) ]
```

```
from sklearn.pipeline import Pipeline
```

```
pipe = Pipeline(operations)
```

```
from sklearn.model_selection import GridSearchCV
```

```
k_value = list(range(1, 20))
```

```
param_grid = { 'knn__n_neighbors': k_values,
```

```
full_cv_classifier = GridSearchCV(pipe, param_grid,  
cv=5,  
scoring='accuracy')
```

```
full_cv_classifier.fit(X_train, y_train)
```

```
full_cv_classifier.best_estimators_.get_params()
```

```
full_cv_classifier.predict(X_test)
```

```
full_preds = full_cv_classifier.predict(X_test)
```