# KNN - Tutorial

Hisham Ihshaish
Bristol, UK
June, 2021
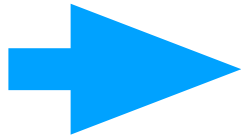
# Nah!

# Evaluate,
## in training and in testing

| | | |
|---|---|---|
| **Training Data** | | |
| **Training Labels** | **Model** | how "correct"? } **Training** |
| **Test Data** | **Prediction** | |
| **Test Labels** | **Evaluation** | how "correct"? } **Generalisation** |

# Data… splitting
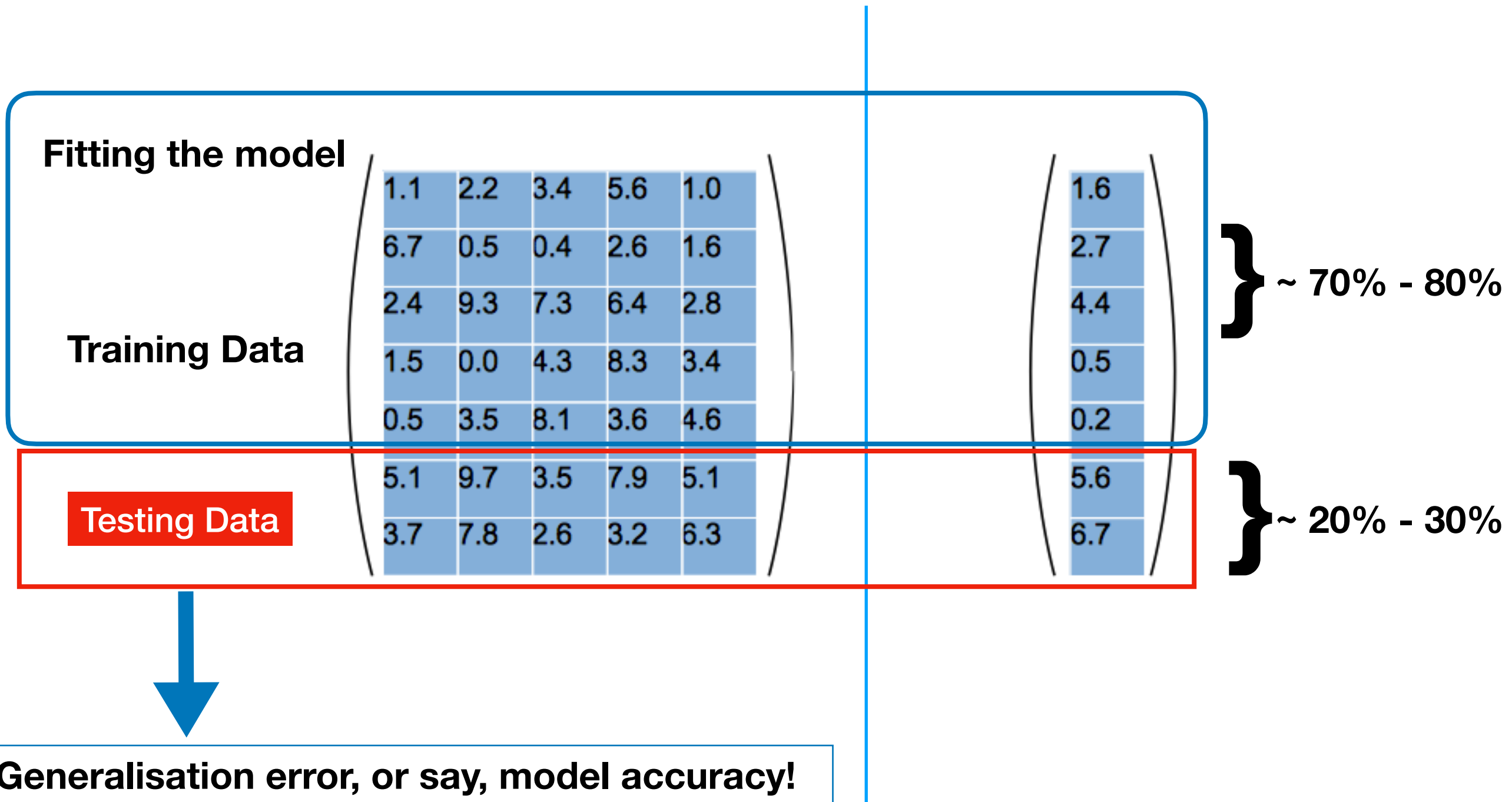
- for **training** and for **testing** (generalisation)

$$X = \begin{pmatrix} 1.1 & 2.2 & 3.4 & 5.6 & 1.0 \\ 6.7 & 0.5 & 0.4 & 2.6 & 1.6 \\ 2.4 & 9.3 & 7.3 & 6.4 & 2.8 \\ 1.5 & 0.0 & 4.3 & 8.3 & 3.4 \\ 0.5 & 3.5 & 8.1 & 3.6 & 4.6 \\ 5.1 & 9.7 & 3.5 & 7.9 & 5.1 \\ 3.7 & 7.8 & 2.6 & 3.2 & 6.3 \end{pmatrix} \qquad y = \begin{pmatrix} 1.6 \\ 2.7 \\ 4.4 \\ 0.5 \\ 0.2 \\ 5.6 \\ 6.7 \end{pmatrix}$$

# Data… splitting

**-** for **training** and for **testing** (generalisation)

**Fitting the model**

**Training Data**

| 1.1 | 2.2 | 3.4 | 5.6 | 1.0 |
|-----|-----|-----|-----|-----|
| 6.7 | 0.5 | 0.4 | 2.6 | 1.6 |
| 2.4 | 9.3 | 7.3 | 6.4 | 2.8 |
| 1.5 | 0.0 | 4.3 | 8.3 | 3.4 |
| 0.5 | 3.5 | 8.1 | 3.6 | 4.6 |

| 1.6 |
|-----|
| 2.7 |
| 4.4 |
| 0.5 |
| 0.2 |

} ~ **70% - 80%**

**Testing Data**

| 5.1 | 9.7 | 3.5 | 7.9 | 5.1 |
|-----|-----|-----|-----|-----|
| 3.7 | 7.8 | 2.6 | 3.2 | 6.3 |

| 5.6 |
|-----|
| 6.7 |

} ~ **20% - 30%**

**Generalisation error, or say, model accuracy!**

# Run example - explained

```python
import numpy as np
import pandas as pd

from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
import matplotlib.pyplot as plt
%matplotlib inline

df = pd.read_csv('/Users/test/Downloads/breast-cancer-wisconsin.data.txt')
```
*read data*

```python
df.replace('?', -99999, inplace=True)
df.drop(['id'], 1, inplace=True)
```
*preprocessing*

```python
X = np.array(df.drop(['class'], 1))
y = np.array(df['class'])
```
*identify X and y*

```python
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)


clf = KNeighborsClassifier(n_neighbors=1)
clf.fit (X_train, y_train)

print (clf.score(X_train, y_train))
print (clf.score(X_test, y_test))

data = np.array([4,3,3,2,1,2,1,1,2])

prediction= clf.predict(data.reshape(1,-1))

print(prediction)
```

$$X = \begin{pmatrix} 1.1 & 2.2 & 3.4 & 5.6 & 1.0 \\ 6.7 & 0.5 & 0.4 & 2.6 & 1.6 \\ 2.4 & 9.3 & 7.3 & 6.4 & 2.8 \\ 1.5 & 0.0 & 4.3 & 8.3 & 3.4 \\ 0.5 & 3.5 & 8.1 & 3.6 & 4.6 \\ 5.1 & 9.7 & 3.5 & 7.9 & 5.1 \\ 3.7 & 7.8 & 2.6 & 3.2 & 6.3 \end{pmatrix} \quad y = \begin{pmatrix} 1.6 \\ 2.7 \\ 4.4 \\ 0.5 \\ 0.2 \\ 5.6 \\ 6.7 \end{pmatrix}$$

# Run example - explained

```python
import numpy as np
import pandas as pd

from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
import matplotlib.pyplot as plt
%matplotlib inline

df = pd.read_csv('/Users/test/Downloads/breast-cancer-wisconsin.data.txt')
```
→ *read data*

```python
df.replace('?', -99999, inplace=True)
df.drop(['id'], 1, inplace=True)
```
→ *preprocessing*

```python
X = np.array(df.drop(['class'], 1))
y = np.array(df['class'])
```
→

$$X = \begin{pmatrix} 1.1 & 2.2 & 3.4 & 5.6 & 1.0 \\ 6.7 & 0.5 & 0.4 & 2.6 & 1.6 \\ 2.4 & 9.3 & 7.3 & 6.4 & 2.8 \\ 1.5 & 0.0 & 4.3 & 8.3 & 3.4 \\ 0.5 & 3.5 & 8.1 & 3.6 & 4.6 \\ 5.1 & 9.7 & 3.5 & 7.9 & 5.1 \\ 3.7 & 7.8 & 2.6 & 3.2 & 6.3 \end{pmatrix} \quad y = \begin{pmatrix} 1.6 \\ 2.7 \\ 4.4 \\ 0.5 \\ 0.2 \\ 5.6 \\ 6.7 \end{pmatrix}$$

```python
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
```

```python
clf = KNeighborsClassifier(n_neighbors=1)
clf.fit (X_train, y_train)

print (clf.score(X_train, y_train))
print (clf.score(X_test, y_test))

data = np.array([4,3,3,2,1,2,1,1,2])

prediction= clf.predict(data.reshape(1,-1))

print(prediction)
```

*Split into training data and testing data*

7

# Data… splitting

**-** for **training** and for **testing** (generalisation)

|     |     |     |     |     |
|-----|-----|-----|-----|-----|
| 1.1 | 2.2 | 3.4 | 5.6 | 1.0 |
| 6.7 | 0.5 | 0.4 | 2.6 | 1.6 |
| 2.4 | 9.3 | 7.3 | 6.4 | 2.8 |
| 1.5 | 0.0 | 4.3 | 8.3 | 3.4 |
| 0.5 | 3.5 | 8.1 | 3.6 | 4.6 |
| 5.1 | 9.7 | 3.5 | 7.9 | 5.1 |
| 3.7 | 7.8 | 2.6 | 3.2 | 6.3 |

**X_train**

**X_test**

| |
|-----|
| 1.6 |
| 2.7 |
| 4.4 |
| 0.5 |
| 0.2 |
| 5.6 |
| 6.7 |

**y_train**

**y_test**

**sklearn** >> **X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)**

8

# here goes what we did last week

```python
import numpy as np
import pandas as pd

from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
import matplotlib.pyplot as plt
%matplotlib inline

df = pd.read_csv('/Users/test/Downloads/breast-cancer-wisconsin.data.txt')

df.replace('?', -99999, inplace=True)
df.drop(['id'], 1, inplace=True)

X = np.array(df.drop(['class'], 1))
y = np.array(df['class'])

X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)

clf = KNeighborsClassifier(n_neighbors=1)
clf.fit (X_train, y_train)

print (clf.score(X_train, y_train))
print (clf.score(X_test, y_test))

data = np.array([4,3,3,2,1,2,1,1,2])

prediction= clf.predict(data.reshape(1,-1))

print(prediction)
```

*read data*

*preprocessing*

$$X = \begin{pmatrix} 1.1 & 2.2 & 3.4 & 5.6 & 1.0 \\ 6.7 & 0.5 & 0.4 & 2.6 & 1.6 \\ 2.4 & 9.3 & 7.3 & 6.4 & 2.8 \\ 1.5 & 0.0 & 4.3 & 8.3 & 3.4 \\ 0.5 & 3.5 & 8.1 & 3.6 & 4.6 \\ 5.1 & 9.7 & 3.5 & 7.9 & 5.1 \\ 3.7 & 7.8 & 2.6 & 3.2 & 6.3 \end{pmatrix} \qquad y = \begin{pmatrix} 1.6 \\ 2.7 \\ 4.4 \\ 0.5 \\ 0.2 \\ 5.6 \\ 6.7 \end{pmatrix}$$
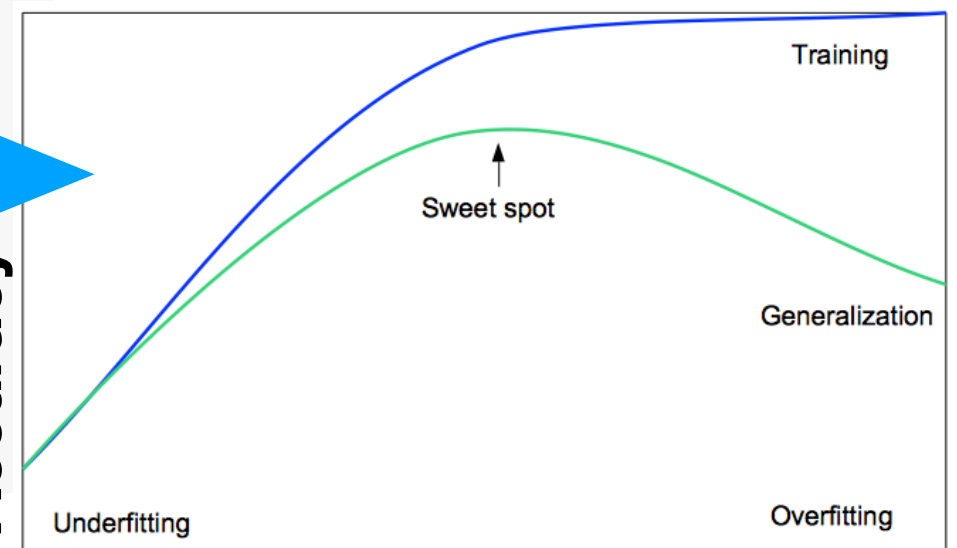
**Split into training data and testing data**

*train, only on <u>training data</u>*

*training error*

*generalisation error*



9

# note on file upload

Preparation:

- Call libraries
- Read file:

a) if on **local jupyter notebook** simply use `df = pd.read_csv('`[`/Users/test/Downloads/breast-cancer-wisconsin.data.txt`](#)`')`

b) if in colabd, import io, then upload file as shown below.

- Encapculate file in a Dataframe — (preferred) check methods of df to show first rows, shape, and statistics of your data.

```
[31]  import numpy as np
      import pandas as pd
      import io

      from sklearn.model_selection import train_test_split
      from sklearn.neighbors import KNeighborsClassifier
      import matplotlib.pyplot as plt


      #df = pd.read_csv('/Users/test/Downloads/breast-cancer-wisconsin.data.txt')

      from google.colab import files
      uploaded = files.upload()


      df = pd.read_csv(io.BytesIO(uploaded['breast-cancer-wisconsin.data.txt']))
```