

# PROJECT AI

672132003 Sitthiphong Prompalit

672132164 Kamchai Boonsri

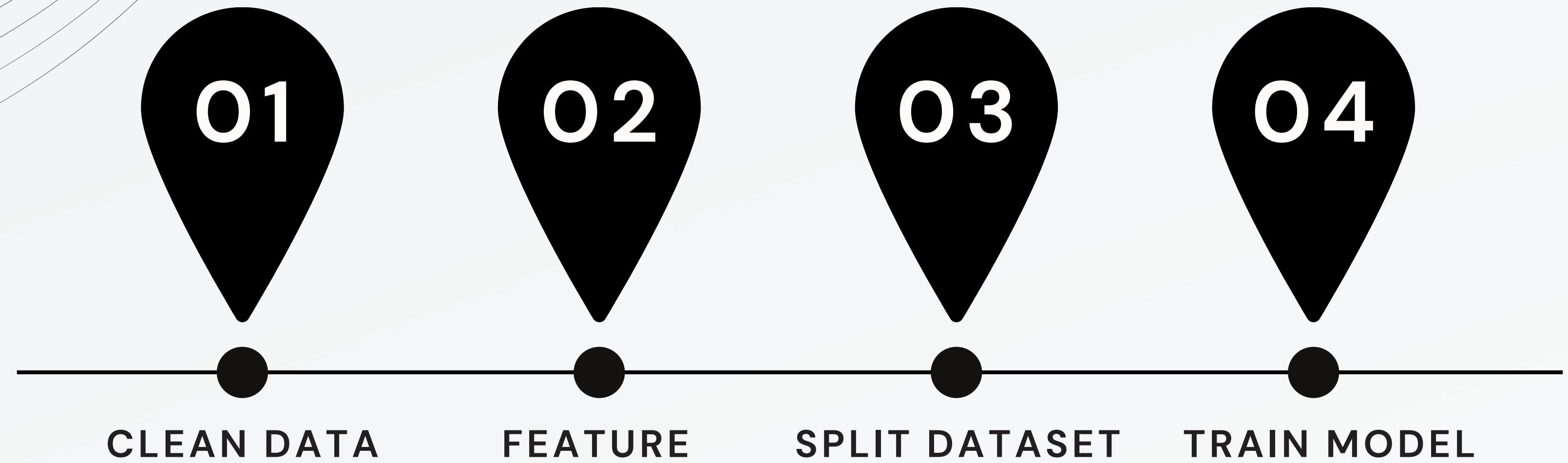
# CONTENT

- 01** SIMPLE NEURAL NETWORK MODEL
- 02** AN ALGORITHM MODEL THAT YOU LEARN FROM KAGGLE COMPETITION

# **PAHT 1**

## **SIMPLE NEURAL NETWORK MODEL**

### **TITANIC**



# IMPORT



```
1 import pandas as pd  
2 import tensorflow as tf  
3 import matplotlib.pyplot as plt  
4 import seaborn as sns
```

# CLEAN DATA



```
1 dataset = pd.read_csv('train.csv')
2 for x in dataset.index:
3     sex = str(dataset.loc[x, "Sex"])
4     embarked = str(dataset.loc[x, "Embarked"])
5     pclass = str(dataset.loc[x, "Pclass"])
6     if ('female' in sex) == True:
7         dataset.loc[x, "female"] = 1
8         dataset.loc[x, "male"] = 0
9     else:
10        dataset.loc[x, "male"] = 1
11        dataset.loc[x, "female"] = 0
12    if('C' in embarked) == True:
13        dataset.loc[x, "Cherbourg"] = 1
14        dataset.loc[x, "Queenstown"] = 0
15        dataset.loc[x, "Southampton"] = 0
16    elif('Q' in embarked) == True:
17        dataset.loc[x, "Cherbourg"] = 0
18        dataset.loc[x, "Queenstown"] = 1
19        dataset.loc[x, "Southampton"] = 0
20    elif('S' in embarked) == True:
21        dataset.loc[x, "Cherbourg"] = 0
22        dataset.loc[x, "Queenstown"] = 0
23        dataset.loc[x, "Southampton"] = 1
```



```
1 if('1' in pclass) == True:
2     dataset.loc[x, "1st"] = 1
3     dataset.loc[x, "2nd"] = 0
4     dataset.loc[x, "3rd"] = 0
5 elif('2' in pclass) == True:
6     dataset.loc[x, "1st"] = 0
7     dataset.loc[x, "2nd"] = 1
8     dataset.loc[x, "3rd"] = 0
9 elif('3' in pclass) == True:
10    dataset.loc[x, "1st"] = 0
11    dataset.loc[x, "2nd"] = 0
12    dataset.loc[x, "3rd"] = 1
13 print(dataset.head())
14
15 cleaned = dataset.dropna().replace('', None)
16
17
18 cleaned.to_csv('cleaned.csv', index=False)
19 print(cleaned.head())
```

# CLEAN DATA

	Survived	Pclass	Sex	Age	SibSp	Parch	Embarked	male	female	Cherbourg	Queenstown	Southampton	1st	2nd	3rd
1	1	1	female	38.0	1	0	C	0.0	1.0	1.0	0.0	0.0	1.0	0.0	0.0
3	1	1	female	35.0	1	0	S	0.0	1.0	0.0	0.0	1.0	1.0	0.0	0.0
6	0	1	male	54.0	0	0	S	1.0	0.0	0.0	0.0	1.0	1.0	0.0	0.0
10	1	3	female	4.0	1	1	S	0.0	1.0	0.0	0.0	1.0	0.0	0.0	1.0
11	1	1	female	58.0	0	0	S	0.0	1.0	0.0	0.0	1.0	1.0	0.0	0.0



```
1 def split_dataset(dataset, train_frac=0.3):
2     train = dataset.sample(frac=train_frac)
3     val = dataset.drop(train.index)
4     return train, val
```



```
1
2 train, val = split_dataset(cleaned)
3
4
5 xtrain = train[['female', 'male', 'Cherbourg', 'Queenstown', 'Southampton', '1st', '2nd', '3rd']]
6 ytrain = train['Survived']
7
8 xtest = val[['female', 'male', 'Cherbourg', 'Queenstown', 'Southampton', '1st', '2nd', '3rd']]
9 ytest = val['Survived']
10
11 print(xtrain.shape)
12 print(ytrain.shape)
13 print(xtest.shape)
14 print(ytest.shape)
```

(55, 8)

(55, )

(128, 8)

(128, )



```
1
2 model = tf.keras.models.Sequential()
3 model.add(tf.keras.layers.InputLayer(input_shape=(8,)))
4 model.add(tf.keras.layers.Dense(1000, activation='sigmoid'))
5 model.add(tf.keras.layers.Dense(500, activation='sigmoid'))
6 model.add(tf.keras.layers.Dense(300, activation='sigmoid'))
7 model.add(tf.keras.layers.Dense(100, activation='sigmoid'))
8 model.add(tf.keras.layers.Dense(1, activation='linear'))
9
10 print(model.summary())
11
12 loss = 'mse'
13 metric = 'mae'
14 epochs = 1000
15 model.compile(loss=loss, optimizer='sgd', metrics=[metric])
16 history = model.fit(xtrain, ytrain, epochs=epochs, batch_size=64, verbose=1, validation_data=(xtest, ytest))
17
18
```

Layer (type)	Output Shape	Param #
dense_35 (Dense)	(None, 1000)	9000
dense_36 (Dense)	(None, 500)	500500
dense_37 (Dense)	(None, 300)	150300
dense_38 (Dense)	(None, 100)	30100
dense_39 (Dense)	(None, 1)	101

=====

Total params: 690,001  
Trainable params: 690,001  
Non-trainable params: 0

---

None  
Epoch 1/1000  
1/1 [=====] - 0s 249ms/step - loss: 0.4820 - mae: 0.5059 - val\_loss: 0.2263 - val\_mae: 0.4020  
Epoch 2/1000  
1/1 [=====] - 0s 20ms/step - loss: 0.2400 - mae: 0.4157 - val\_loss: 0.2177 - val\_mae: 0.4344  
Epoch 3/1000  
...  
Epoch 999/1000  
1/1 [=====] - 0s 31ms/step - loss: 0.2258 - mae: 0.4520 - val\_loss: 0.2181 - val\_mae: 0.4442  
Epoch 1000/1000  
1/1 [=====] - 0s 58ms/step - loss: 0.2258 - mae: 0.4520 - val\_loss: 0.2181 - val\_mae: 0.4442

# ACCURACY

```
● ● ●  
1 _, accuracy = model.evaluate(xtrain, ytrain)  
2 print('Accuracy: %.2f' % (accuracy*100))
```

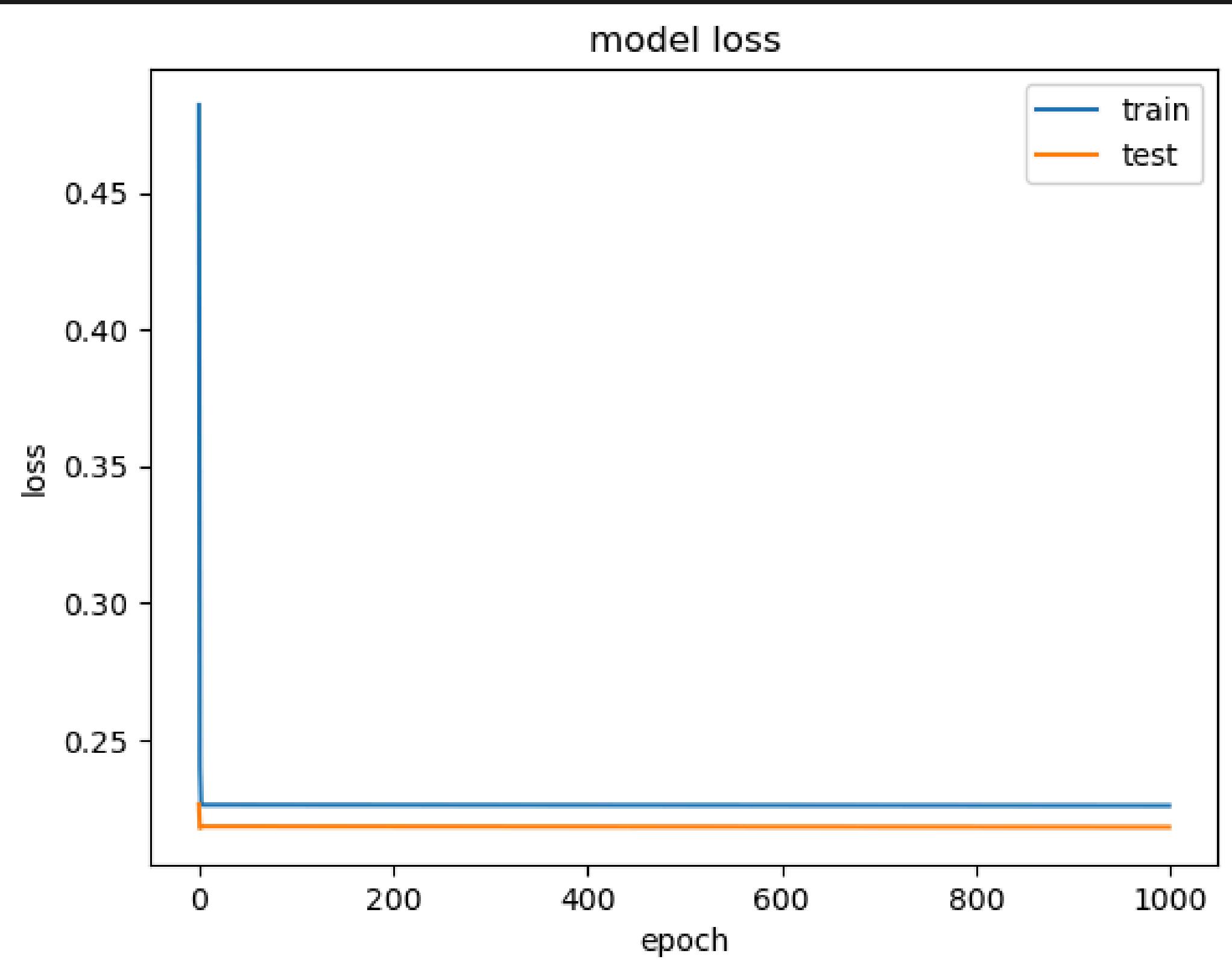
```
2/2 [=====] - 0s 2ms/step - loss: 0.2258 - mae: 0.4520  
Accuracy: 45.20
```

**45.20 %**

**Split data 0.3**

# Split Data 0.3

```
● ● ●  
1 plt.plot(history.history['loss'])  
2 plt.plot(history.history['val_loss'])  
3 plt.title('model loss')  
4 plt.ylabel('loss')  
5 plt.xlabel('epoch')  
6 plt.legend(['train', 'test'])  
7 plt.show()  
8
```



# ACCURACY



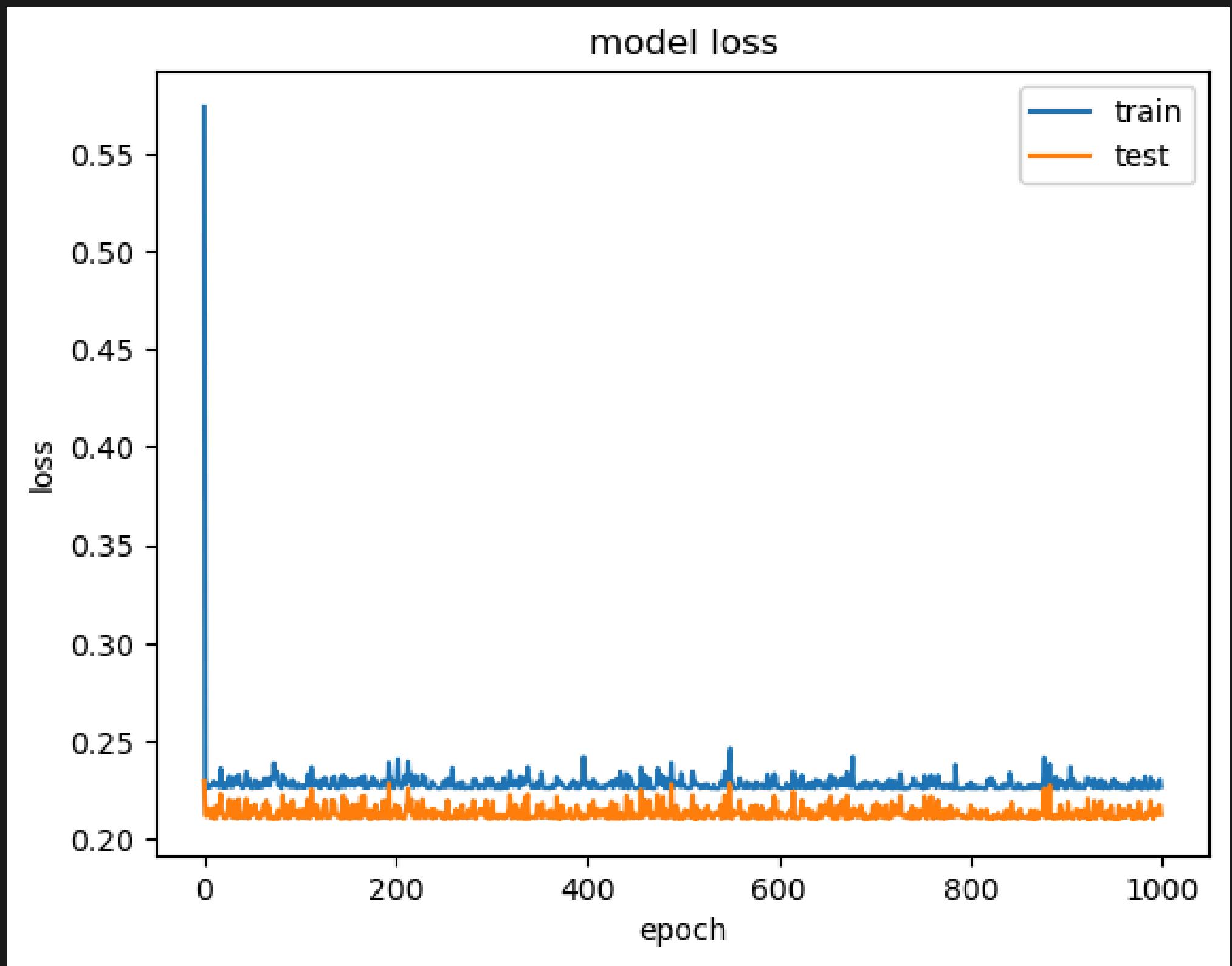
```
1 _, accuracy = model.evaluate(xtrain, ytrain)
2 print('Accuracy: %.2f' % (accuracy*100))
```

```
4/4 [=====] - 0s 2ms/step - loss: 0.2258 - mae: 0.4533
Accuracy: 45.33
```

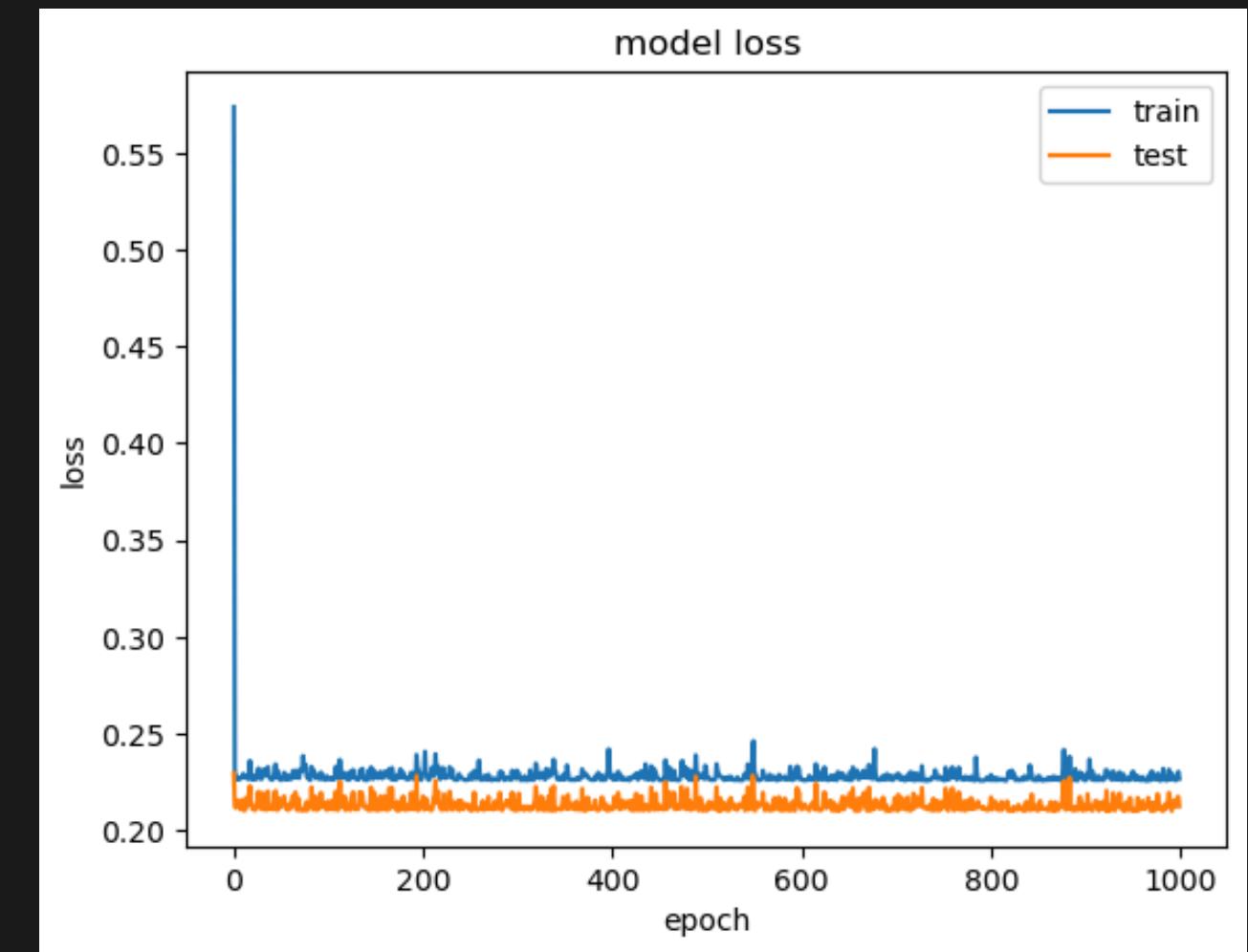
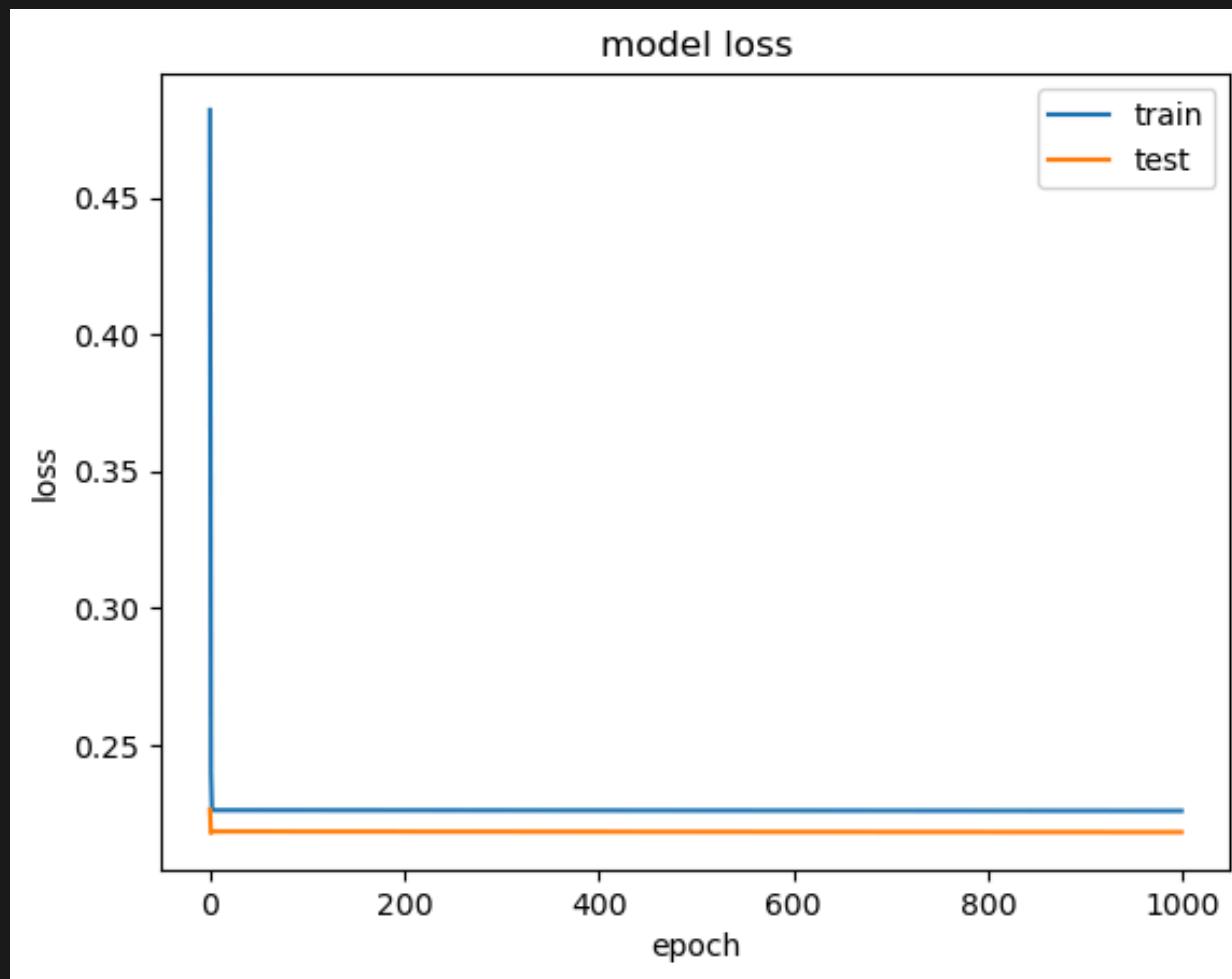
**45.33 %**

**Split data 0.6**

# Split Data 0.6



# Result



Data มีค่าความสัมพันธ์กันระหว่าง Feature น้อยเกินไป

# **PART 2**

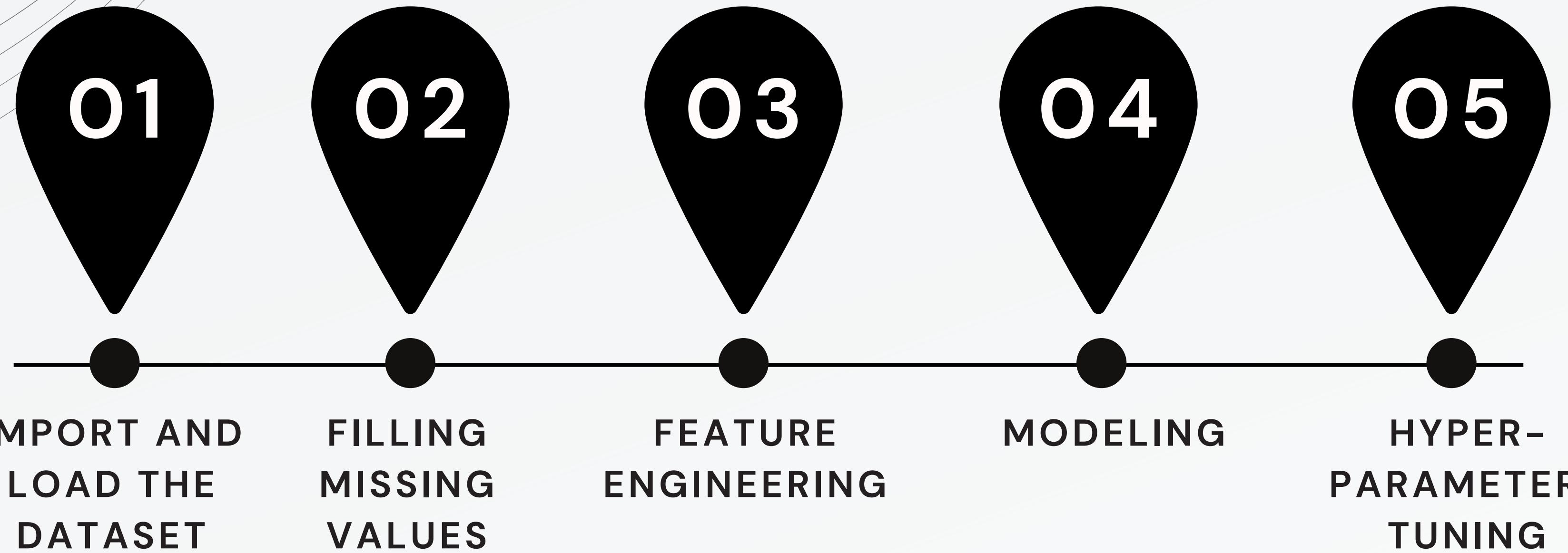
# **TITANIC FROM KAGGLE COMPETITION**

# Titanic Model with 90% accuracy by VK\_ds

The image shows a screenshot of a video thumbnail on a platform like YouTube. The thumbnail features a large orange title 'Titanic Model with 90% accuracy by VK\_ds'. Below the title, there's a profile picture of a person, the channel name 'VK\_DS', and the upload date '6Y AGO'. It also shows '62,350 VIEWS'. To the right, there are interaction metrics: '321' (likes), 'Copy & Edit' (button), '725' (dislikes), a gold play button icon, and a more options menu. The main video content area below the thumbnail shows the title 'Titanic Model with 90% accuracy' and the subtitle 'Python · Titanic - Machine Learning from Disaster'.

**Titanic Model with 90% accuracy**

Python · Titanic - Machine Learning from Disaster



# Import and load the Dataset

```
● ● ●  
1 import numpy as np  
2 import pandas as pd  
3 import matplotlib.pyplot as plt  
4 import seaborn as sns  
5 import warnings  
6 warnings.filterwarnings('ignore')  
7 %matplotlib inline
```



```
1 train_df.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S



```
1 print('__Test_DataSet__')
2 test_df.head()
```

\_\_Test\_DataSet\_\_

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	892	3	Kelly, Mr. James	male	34.5	0	0	330911	7.8292	NaN	Q
1	893	3	Wilkes, Mrs. James (Ellen Needs)	female	47.0	1	0	363272	7.0000	NaN	S
2	894	2	Myles, Mr. Thomas Francis	male	62.0	0	0	240276	9.6875	NaN	Q
3	895	3	Wirz, Mr. Albert	male	27.0	0	0	315154	8.6625	NaN	S
4	896	3	Hirvonen, Mrs. Alexander (Helga E Lindqvist)	female	22.0	1	1	3101298	12.2875	NaN	S



```
1 def missingdata(data):
2     total = data.isnull().sum().sort_values(ascending = False)
3     percent = (data.isnull().sum()/data.isnull().count()*100).sort_values(ascending = False)
4     ms= pd.concat([total, percent], axis=1, keys=['Total', 'Percent'])
5     ms= ms[ms["Percent"] > 0]
6     return ms
```



```
1 missingdata(train_df)
```



```
1 missingdata(test_df)
```

	Total	Percent
Cabin	687	77.104377
Age	177	19.865320
Embarked	2	0.224467

	Total	Percent
Cabin	327	78.229665
Age	86	20.574163
Fare	1	0.239234

# Filling missing Values



```
1 # หาค่าค่าเฉลี่ยของอายุ
2 test_df['Age'].mean()
```

= 30.272590361445783



```
1 # เติมค่าที่ขาดหายไป ในคอลัมน์ 'Embarked' ด้วยค่าที่ปรากฏบ่อยที่สุด ในคอลัมน์นั้นเอง
2 train_df['Embarked'].fillna(train_df['Embarked'].mode()[0], inplace = True)
```



```
1 # เติมค่าที่ขาดหายไป ในคอลัมน์ 'Fare' ด้วยค่าที่ปรากฏบ่อยที่สุด ในคอลัมน์นั้นเอง
2 test_df['Fare'].fillna(test_df['Fare'].median(), inplace = True)
```



```
1 # "Cabin" มีข้อมูลขาดหายมากกว่า 75% ทั้งใน Train และ Test เขาจึงตัด feature นี้ออกไป
2 drop_column = ['Cabin']
3 train_df.drop(drop_column, axis=1, inplace = True)
4 test_df.drop(drop_column, axis=1, inplace=True)
```



```
1 # "Age" ในชุดข้อมูล Train และ Test มีข้อมูลขาดหายมากกว่า 15% เขายังเติมข้อมูลที่ขาดหายด้วยค่ามัธยฐาน (median)
2 test_df['Age'].fillna(test_df['Age'].median(), inplace = True)
3 train_df['Age'].fillna(train_df['Age'].median(), inplace = True)
```



```
1 # ເຊື້ອ missing value ອຶກຄ້າ  
2 print('check the nan value in train data')  
3 print(train_df.isnull().sum())  
4 print('__'*30)  
5 print('check the nan value in test data')  
6 print(test_df.isnull().sum())  
7
```

```
check the nan value in train data  
PassengerId      0  
Survived         0  
Pclass           0  
Name             0  
Sex              0  
Age              0  
SibSp            0  
Parch            0  
Ticket           0  
Fare             0  
Embarked         0  
dtype: int64
```

---

```
check the nan value in test data  
PassengerId      0  
Pclass           0  
Name             0  
Sex              0  
Age              0  
SibSp            0  
Parch            0  
Ticket           0  
Fare             0  
Embarked         0  
dtype: int64
```

# Feature Engineering



```
1 # รวมข้อมูล train และ test เข้าด้วยกัน  
2 all_data=[train_df,test_df]
```



```
1 # สร้าง feature ใหม่ ด้วยการรวม SibSp และ Parch  
2 # SibSp จำนวนพี่น้อง/คู่สมรสบนเรือไททานิค  
3 # Parch พ่อแม่/ลูกบนเรือไททานิค  
4 # บวกด้วย 1 (ตัวเอง)  
5 # FamilySize เป็น feature รวมจำนวนหัวหนี้ในครอบครัวของคนหนึ่งคน  
6  
7 for dataset in all_data:  
8     dataset['FamilySize'] = dataset['SibSp'] + dataset['Parch'] + 1
```

```
1 import re
2 # function หาคำนำหน้าชื่อ
3 def get_title(name):
4     title_search = re.search(' ([A-Za-z]+)\.', name)
5     if title_search:
6         return title_search.group(1)
7     return ""
8
9 # สร้างชื่อ "Title" ซึ่งเก็บคำนำหน้าชื่อ
10 for dataset in all_data:
11     dataset['Title'] = dataset['Name'].apply(get_title)
12
13 # รวมคำนำหน้าชื่อที่ไม่พบบ่อยทั้งหมดเป็นกลุ่มเดียวกันชื่อ "Rare"
14 for dataset in all_data:
15     dataset['Title'] = dataset['Title'].replace(['Lady', 'Countess','Capt', 'Col','Don',
16                                                 'Dr', 'Major', 'Rev', 'Sir', 'Jonkheer', 'Dona'], 'Rare')
17     dataset['Title'] = dataset['Title'].replace('Mlle', 'Miss')
18     dataset['Title'] = dataset['Title'].replace('Ms', 'Miss')
19     dataset['Title'] = dataset['Title'].replace('Mme', 'Mrs')
```

```
1 # สร้าง "Age_bin" แบ่งค่าในคอลัมน์ "Age" ออกเป็นกลุ่มต่างๆ ตามช่วงอายุที่กำหนด
2 for dataset in all_data:
3     dataset['Age_bin'] = pd.cut(dataset['Age'], bins=[0,12,20,40,120], labels=['Children','Teenage','Adult','Elder'])
```



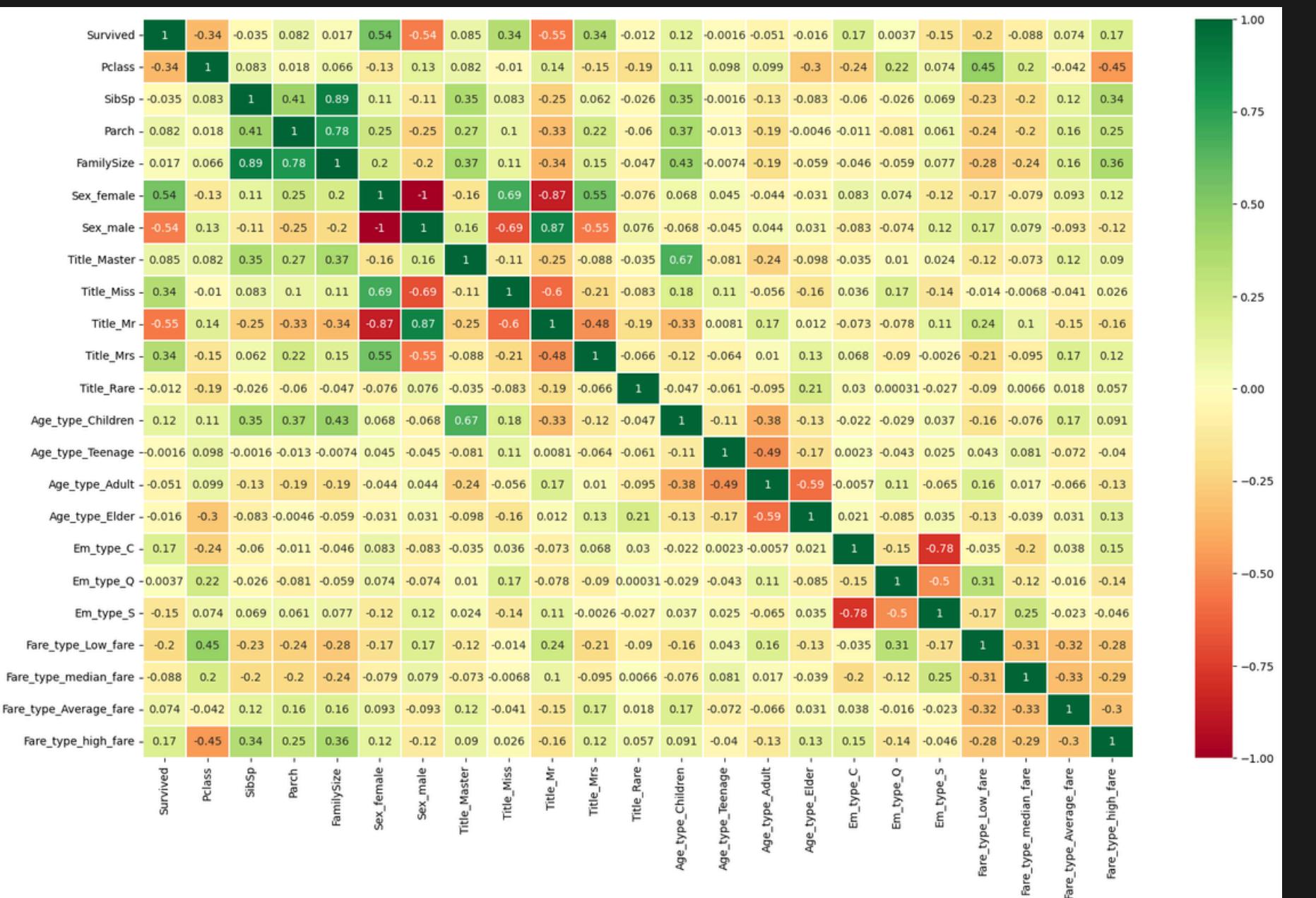
```
1 # get_dummies เป็น function สำหรับแปลงข้อมูลเชิงหมวดหมู่ (categorical data) ให้เป็นข้อมูลเชิงตัวเลข (numerical data)
2 traindf = pd.get_dummies(traindf, columns = ["Sex","Title","Age_bin","Embarked","Fare_bin"],
                           prefix=["Sex","Title","Age_type","Em_type","Fare_type"])
3
```



```
1 testdf = pd.get_dummies(testdf, columns = ["Sex","Title","Age_bin","Embarked","Fare_bin"],
                           prefix=["Sex","Title","Age_type","Em_type","Fare_type"])
2
```

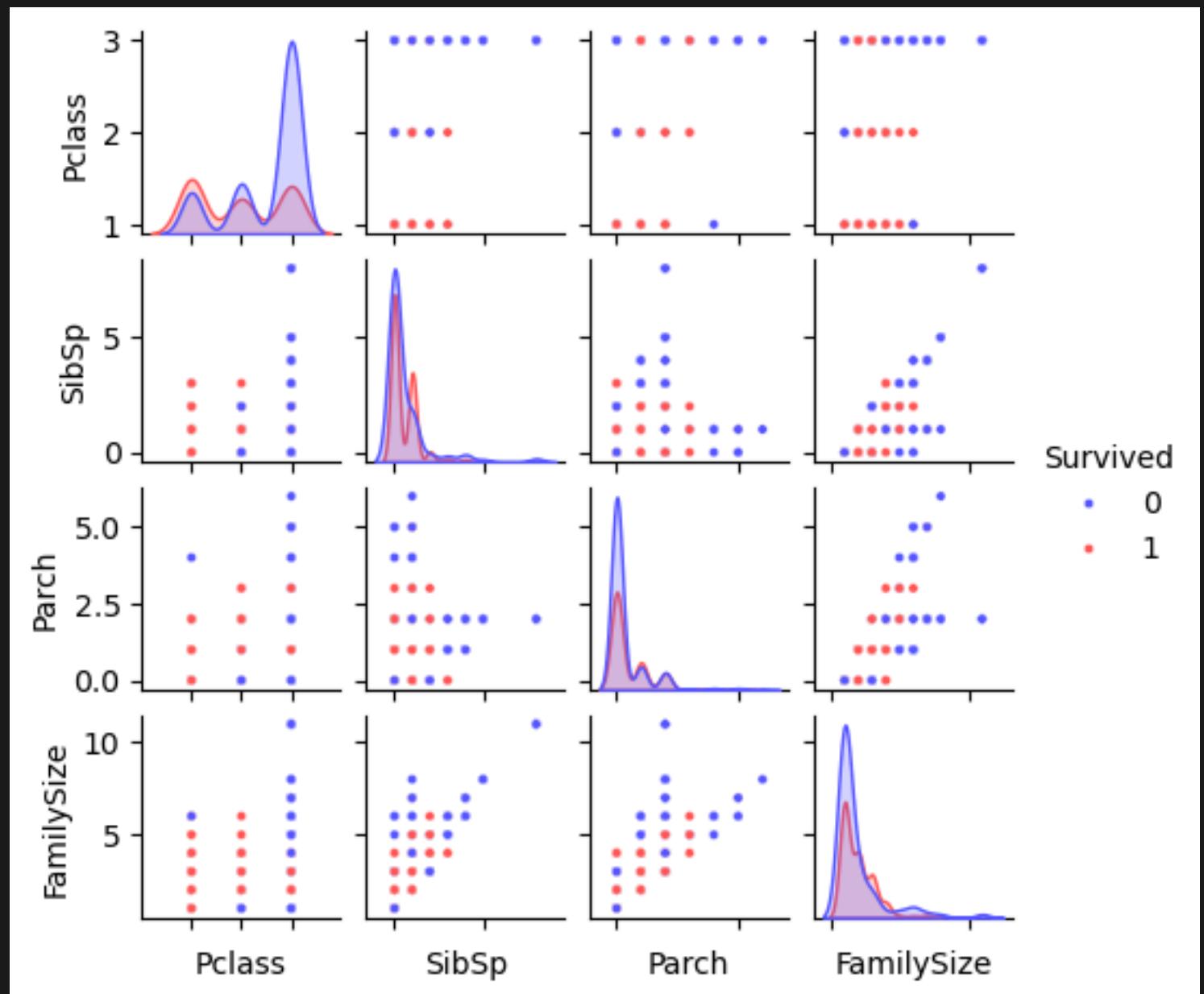


```
1 # แสดงความสัมพันธ์ระหว่าง feature ทั้งหมด
2 sns.heatmap(traindf.corr(), annot=True, cmap='RdYlGn', linewidths=0.2) #data.corr()-->correlation matrix
3 fig=plt.gcf()
4 fig.set_size_inches(20,12)
5 plt.show()
```





```
1 # ใช้ pairplot สร้างกราฟเพื่อสังเกตการกระจายของข้อมูลระหว่าง feature และ lable (Survived) ต่างๆ โดยใช้ Seaborn ช่วย
2 g = sns.pairplot(data=train_df, hue='Survived', palette = 'seismic', size=1.2,diag_kind = 'kde',diag_kws=dict(shade=True),plot_kws=dict(s=10) )
3 g.set(xticklabels=[])
```



# Modelling

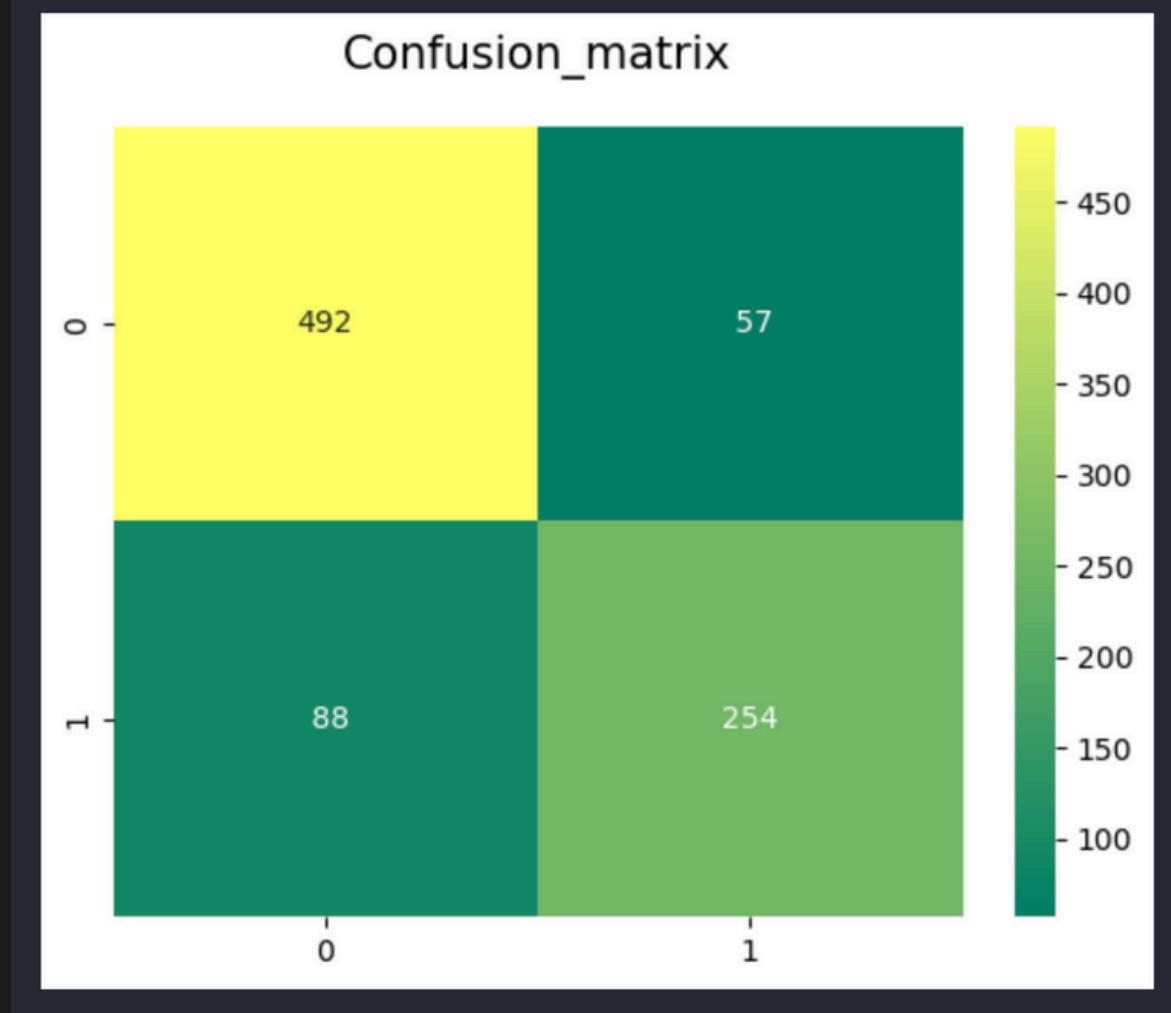


```
1 from sklearn.model_selection import train_test_split #for split the data
2 from sklearn.metrics import accuracy_score #for accuracy_score
3 from sklearn.model_selection import KFold #for K-fold cross validation
4 from sklearn.model_selection import cross_val_score #score evaluation
5 from sklearn.model_selection import cross_val_predict #prediction
6 from sklearn.metrics import confusion_matrix #for confusion matrix
7 all_features = traindf.drop("Survived",axis=1)
8 Targeted_feature = traindf["Survived"]
9 X_train,X_test,y_train,y_test = train_test_split(all_features,Targeted_feature,test_size=0.3,random_state=42)
10 X_train.shape,X_test.shape,y_train.shape,y_test.shape
```

# Random Forest Classifier

```
● ● ●  
1 # Random Forests  
2 from sklearn.ensemble import RandomForestClassifier  
3 model = RandomForestClassifier(criterion='gini', n_estimators=700,  
4                                 min_samples_split=10,min_samples_leaf=1,oob_score=True,  
5                                 random_state=1,n_jobs=-1)  
6 model.fit(X_train,y_train)  
7 prediction_rm=model.predict(X_test)  
8 print('-----The Accuracy of the model-----')  
9 print('The accuracy of the Random Forest Classifier is',round(accuracy_score(prediction_rm,y_test)*100,2))  
10 kfold = KFold()  
11 result_rm=cross_val_score(model,all_features,Targeted_feature,cv=10,scoring='accuracy')  
12 print('The cross validated score for Random Forest Classifier is:',round(result_rm.mean()*100,2))  
13 y_pred = cross_val_predict(model,all_features,Targeted_feature,cv=10)  
14 sns.heatmap(confusion_matrix(Targeted_feature,y_pred),annot=True,fmt='3.0f',cmap="summer")  
15 plt.title('Confusion_matrix', y=1.05, size=15)
```

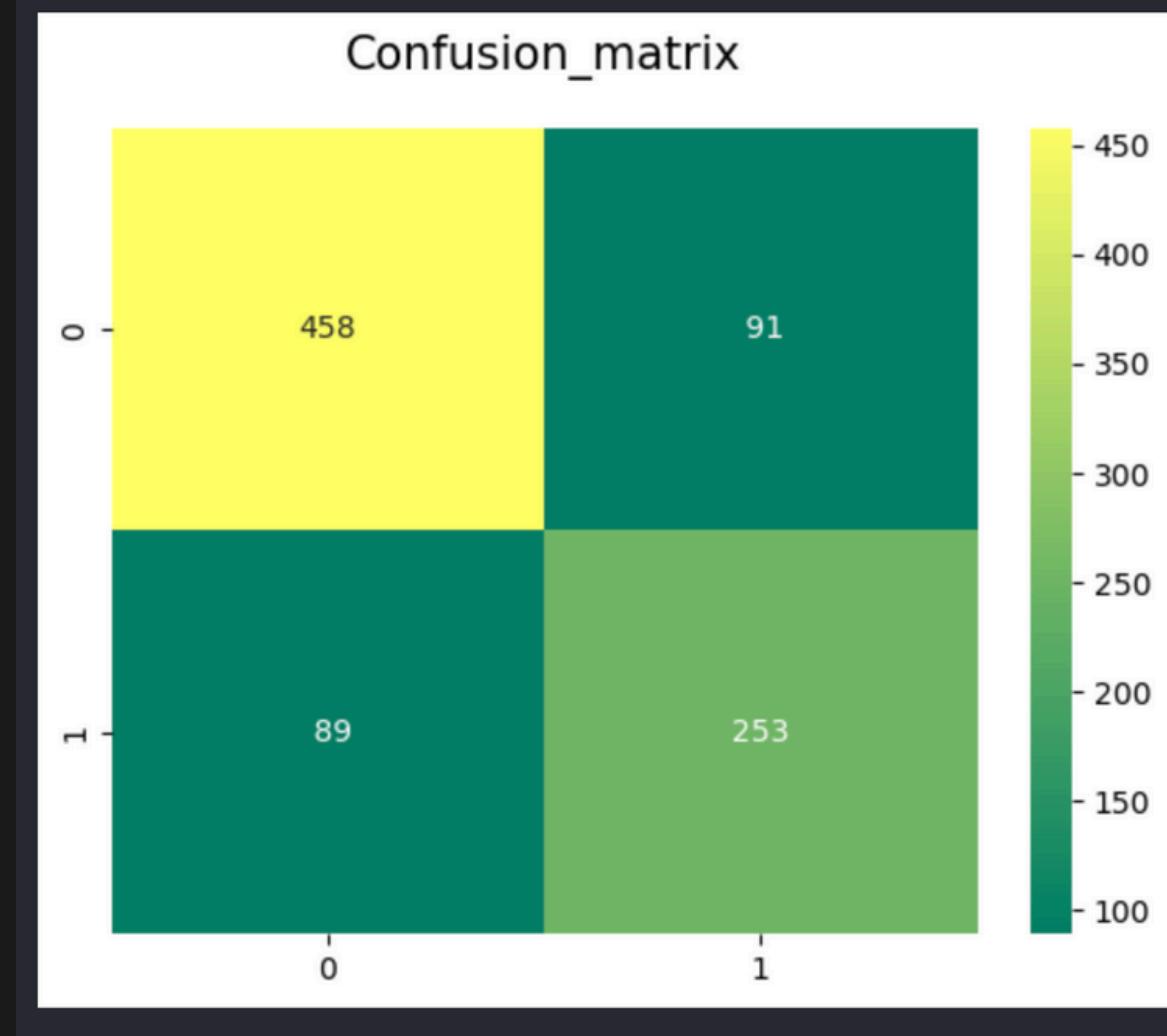
```
-----The Accuracy of the model-----  
The accuracy of the Random Forest Classifier is 82.46  
The cross validated score for Random Forest Classifier is: 83.73  
Text(0.5, 1.05, 'Confusion_matrix')
```



# Gaussian Naive Bayes

```
● ● ●  
1 # Gaussian Naive Bayes  
2 from sklearn.naive_bayes import GaussianNB  
3 model= GaussianNB()  
4 model.fit(X_train,y_train)  
5 prediction_gnb=model.predict(X_test)  
6 print('-----The Accuracy of the model-----')  
7 print('The accuracy of the Gaussian Naive Bayes Classifier is',round(accuracy_score(prediction_gnb,y_test)*100,2))  
8 kfold = KFold() # k=10, split the data into 10 equal parts  
9 result_gnb=cross_val_score(model,all_features,Targeted_feature,cv=10,scoring='accuracy')  
10 print('The cross validated score for Gaussian Naive Bayes classifier is:',round(result_gnb.mean()*100,2))  
11 y_pred = cross_val_predict(model,all_features,Targeted_feature,cv=10)  
12 sns.heatmap(confusion_matrix(Targeted_feature,y_pred),annot=True,fmt='3.0f',cmap="summer")  
13 plt.title('Confusion_matrix', y=1.05, size=15)
```

```
-----The Accuracy of the model-----  
The accuracy of the Gaussian Naive Bayes Classifier is 79.48  
The cross validated score for Gaussian Naive Bayes classifier is: 79.8  
Text(0.5, 1.05, 'Confusion_matrix')
```



# Hyper-Parameters Tuning

```
● ● ●  
1 # Random Forests  
2 from sklearn.ensemble import RandomForestClassifier  
3 random_forest = RandomForestClassifier(  
4     bootstrap=True, class_weight=None, criterion='gini',  
5     max_depth=None, max_leaf_nodes=None,  
6     min_impurity_decrease=0.0,  
7     min_samples_leaf=1, min_samples_split=2,  
8     min_weight_fraction_leaf=0.0, n_estimators=400, n_jobs=1,  
9     oob_score=False, random_state=None, verbose=0,  
10    warm_start=False  
11 )  
12 random_forest.fit(train_X, train_Y)  
13 Y_pred_rf = random_forest.predict(test_X)  
14 random_forest.score(train_X,train_Y)  
15 acc_random_forest = round(random_forest.score(train_X, train_Y) * 100, 2)  
16  
17 print("Important features")  
18 pd.Series(random_forest.feature_importances_,train_X.columns).sort_values(ascending=True).plot.barh(width=0.8)  
19 print('__'*30)  
20 print(acc_random_forest)  
21
```

# Result After Tuning

