

# Download Data from Kaggle

```
In [29]: !pip install --upgrade --force-reinstall --no-deps kaggle

Processing /root/.cache/pip/wheels/af/6a/72/630b7499ff85a4a593377a87ecf55f7d08af42f0d
e9b6303/kaggle-1.5.12-cp37-none-any.whl
Installing collected packages: kaggle
  Found existing installation: kaggle 1.5.12
    Uninstalling kaggle-1.5.12:
      Successfully uninstalled kaggle-1.5.12
  Successfully installed kaggle-1.5.12

In [30]: from google.colab import files
files.upload()
```

Choose Files No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving kaggle.json to kaggle (/kaggle) (kaggle.json: b'{"username": "renata Lucia", "key": "7f35d1eff823543d4d4449562279d1ad b"}')

```
In [32]: #Make directory named kaggle and copy kaggle.json file there.
!mkdir ~/.kaggle

! cp kaggle.json ~/.kaggle/

#Change the permissions of the file.
! chmod 600 ~/.kaggle/kaggle.json

mkdir ~ cannot create directory '/root/.kaggle': File exists
```

```
In [33]: !kaggle competitions download -c tabular-playground-series-apr-2021

tabular-playground-series-apr-2021.zip: Skipping, found more recently modified local copy (use --force to force download)
```

```
In [34]: !ls /content

'kaggle (1).json'  playground  tabular-playground-series-apr-2021.zip
kaggle.json       sample_data
```

```
In [35]: !unzip /content/tabular-playground-series-apr-2021.zip -d playground

Archive: /content/tabular-playground-series-apr-2021.zip
replace playground/sample_submission.csv? [y]es, [n]o, [A]ll, [N]one, [r]ename: y
inflating: playground/sample_submission.csv
replace playground/test.csv? [y]es, [n]o, [A]ll, [N]one, [r]ename: y
inflating: playground/test.csv
replace playground/train.csv? [y]es, [n]o, [A]ll, [N]one, [r]ename: y
inflating: playground/train.csv
```

## Exploratory Data Analysis

```
In [36]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

```
In [45]: df_train = pd.read_csv('/content/playground/train.csv')
print(df_train.columns)
print(len(df_train.columns))
print(df_train.head(100))

Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp', 'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'], dtype='object')
12
PassengerId  Survived  Pclass      Name      Sex      Age  SibSp  Parch      Ticket  Fare  Cabin  Embarked
0            0         1         1  O'Connor, Frankie  male   NaN      2      0      21745  53.00  NaN      S
09245        27.14  C12239      0      S      Bryan, Drew  male   NaN      0      0      27323  13.35  NaN      S
1            1         3         3  Owens, Kenneth  male  0.33      1      2      7129  21.01  NaN      S
2            2         3         3  Kramer, James  male  19.00      0      0      4      A.  57703  71.29  NaN      0
3            3         3         3  Bond, Michael  male  25.00      0      0      4      10866  13.04  NaN      S
10866        13.04  NaN      1      S      Bond, Michael  male  25.00      0      0      4      27635  7.76  NaN      S
...         ...         ...         ...         ...         ...         ...         ...         ...         ...
95         ...         ...         ...         ...         ...         ...         ...         ...         ...         ...
25538        32.24  NaN      1      S      Jordan, Emma  female  28.00      0      0      1      67707  26.46  NaN      1
97         ...         ...         ...         ...         ...         ...         ...         ...         ...         ...
11617        131.28  D7415  1      Q      Elliott, Monica  female  56.00      0      0      0      PC  98
NaN         27.88  NaN      1      S      Bohn, Teresa  female  8.00      0      0      0      99
NaN         27.88  NaN      1      S      Kanter, Betty  female  59.00      0      0      0      48662  278.56  C16097  C
[100 rows x 12 columns]
```

```
In [38]: df_test = pd.read_csv('/content/playground/test.csv')
print(df_test.columns)
print(len(df_test.columns))
print(df_test.head(100))

Index(['PassengerId', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp', 'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'], dtype='object')
11
PassengerId  Pclass      Name      Sex      Age  SibSp  Parch  Ticket  Fare  Cabin  Embarked
0            NaN      3  Holliday, Daniel  female  19.0      0      0      24745  63.01  NaN      S
1            NaN      3  Nguyen, Lorraine  female  53.0      0      0      13264  5.81  NaN      S
2            NaN      1  Harris, Heather  female  19.0      0      0      25990  38.91  NaN      S
BL5315        100003      2  Larsen, Eric  male  25.0      0      0      314011  12.93  NaN      S
4            100004      1  Cleary, Sarah  female  17.0      0      2      26203  26.89  NaN      S
B2515        100003      2  Guzman, Donald  male  4.0      1      0      33068  8.48  NaN      S
95         ...         ...         ...         ...         ...         ...         ...         ...         ...         ...
100009        100096      1  Moleen, Robin  male  52.0      0      0      292128  12.76  NaN      C
96         ...         ...         ...         ...         ...         ...         ...         ...         ...         ...
NaN         100000      2  Lopez, Darlene  female  20.0      0      0      447091  6.17  NaN      S
NaN         100008      3  Tobe, Enrique  male  34.0      0      0      25264  11.64  NaN      S
Al0911        100099      3  Ricker, Matthew  male  3.0      1      1      48876  31.12  NaN      S
99         ...         ...         ...         ...         ...         ...         ...         ...         ...         ...
NaN         100099      3  Ricker, Matthew  male  3.0      1      1      48876  31.12  NaN      S
[100 rows x 11 columns]
```

```
In [46]: print(df_train.shape)
print(df_test.shape)

(100000, 12)
(100000, 11)
```

```
In [47]: pd.set_option('display.max_columns', 12)
pd.set_option('display.width', 1000)
df_train.head()
```

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	0	1	1 O'Connor, Frankie	male	NaN	2	0	209245	2714	C12239	
1	1	0	3 Bryan, Drew	male	NaN	0	0	27323	13.35	NaN	
2	2	0	3 Owens, Kenneth	male	0.33	1	2	7129	21.01	NaN	
3	3	0	3 Kramer, James	male	19.00	0	0	4	10866	13.04	NaN
4	4	1	3 Bond, Michael	male	25.00	0	0	427635	7.76	NaN	

## Feature distributions for train and test datasets

```
In [66]: # count occurrences of each feature value (train dataset)
df_counts_train = {}
for col in df_train.columns:
    if col=="PassengerId":
        continue
    if col=="Survived":
        continue
    if col=="Name":
        continue
    s = df_train[col].value_counts()
    df_feat_dist = pd.DataFrame({'col':s.index, "sum":s.values})
    df_counts_train[col] = df_feat_dist
```

```
In [69]: # count occurrences of each feature value (test dataset)
df_counts_test = {}
for col in df_test.columns:
    if col=="PassengerId":
        continue
    if col=="Survived":
        continue
    if col=="Name":
        continue
    s = df_test[col].value_counts()
    df_feat_dist = pd.DataFrame({'col':s.index, "sum":s.values})
    df_counts_test[col] = df_feat_dist
```

## Plot distributions of features with categorical values

```
In [70]: # importing pandas library
import pandas as pd
# import matplotlib library
import matplotlib.pyplot as plt

for col, df in df_counts_train.items():
    print(col)
    if col not in ["Sex", "Pclass", "SibSp", "Parch", "Embarked"]:
        continue

    xlabels = []
    y1values = []
    y2values = []
    for item in df_counts_train[col]:
        xlabels.append(item)

        y1values.append(df_counts_train[col].loc[df_counts_train[col][col]==item]["sum"].values)
        y2values.append(df_counts_test[col].loc[df_counts_test[col][col]==item]["sum"].values)

# creating dataframe
df = pd.DataFrame({
    col: xlabels,
    'Train': y1values,
    'Test': y2values
})

# plotting graph
df.plot(x=col, y=["Train", "Test"], kind="bar")
```

Pclass

Sex

Age

SibSp

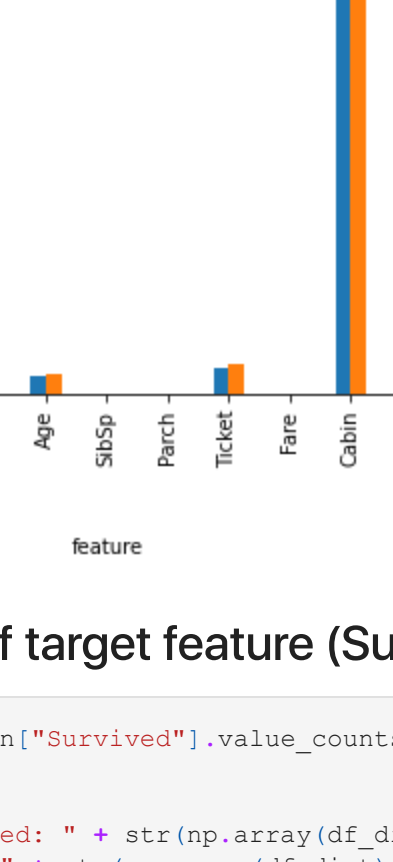
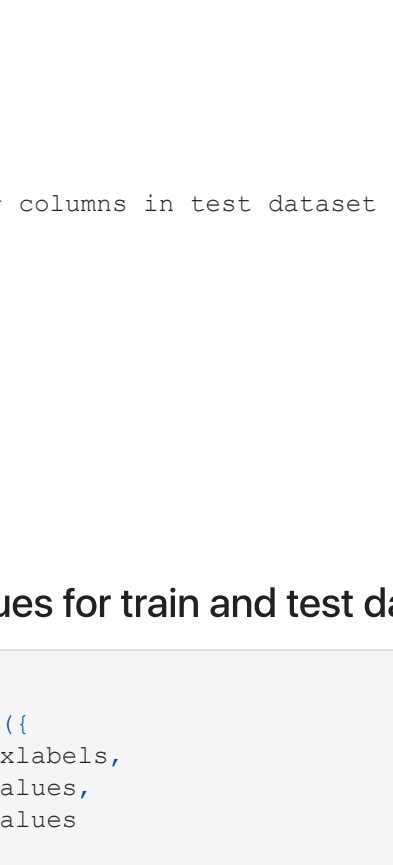
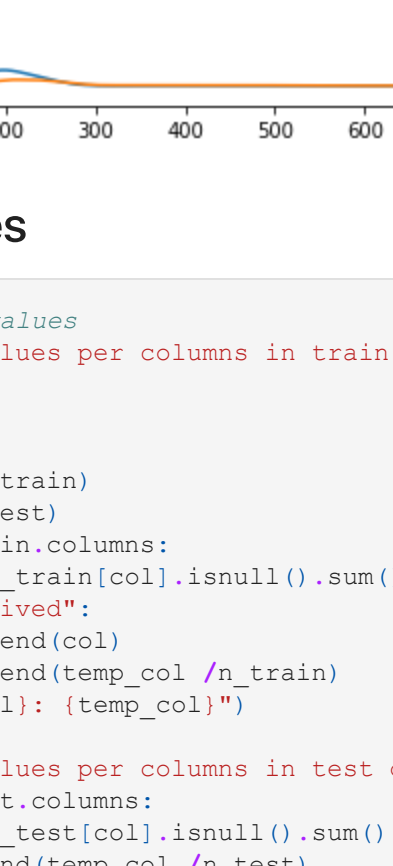
Parch

Ticket

Fare

Cabin

Embarked



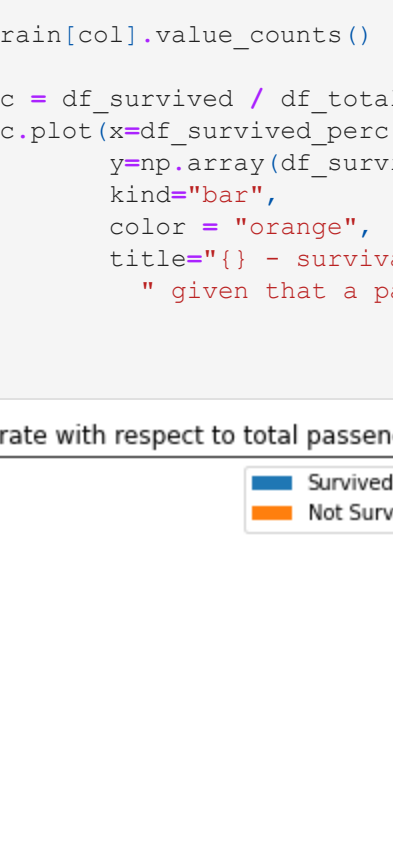
## Plot distribution of features with real values

The probability density function of each feature is estimated using KDE (Kernel Density Estimation)

```
In [ ]: df_age = pd.concat([df_train["Age"], df_test["Age"]], keys=["train", "test"], names=[
df_fare = pd.concat([df_train["Fare"], df_test["Fare"]], keys=["train", "test"], names=

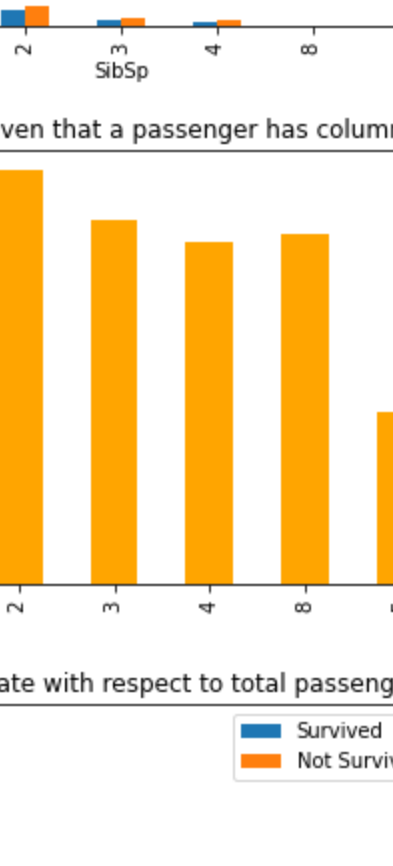
In [ ]: df_age.groupby('type').Age.plot(kind='kde', legend=True, xlim= [df_age["Age"].min(), c

Out [ ]: type
test AxesSubplot(0.125,0.125,0.775x0.755)
train AxesSubplot(0.125,0.125,0.775x0.755)
Name: Age, dtype: object
```



```
In [ ]: df_fare.groupby('type').Fare.plot(kind='kde', legend=True, xlim= [df_fare["Fare"].min

Out [ ]: type
test AxesSubplot(0.125,0.125,0.775x0.755)
train AxesSubplot(0.125,0.125,0.775x0.755)
Name: Fare, dtype: object
```



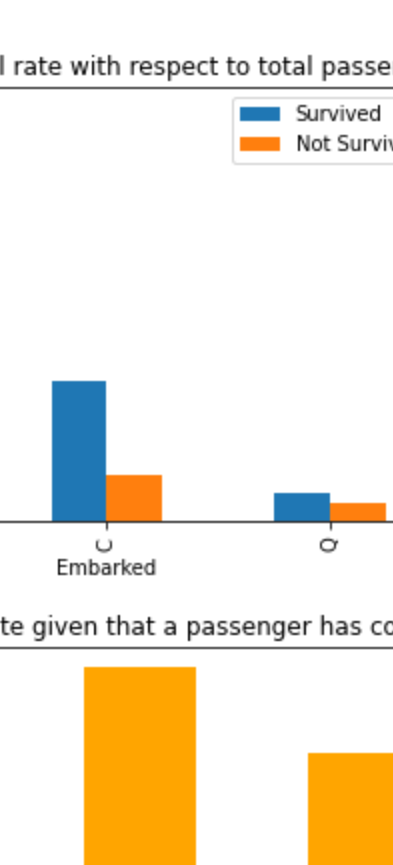
## Missing values

```
In [41]: # Check missing values
print('Missing values per columns in train dataset')
xlabels = []
yvalues = []
n_train = len(df_train)
n_test = len(df_test)
for col in df_train.columns:
    temp_col = df_train[col].isnull().sum()
    if col=="Survived":
        xlabels.append(temp_col/n_train)
        yvalues.append(temp_col/n_train)
    print(f'{col}: {temp_col}')
print()
print('Missing values per columns in test dataset')
for col in df_test.columns:
    temp_col = df_test[col].isnull().sum()
    y2values.append(temp_col/n_test)
    print(f'{col}: {temp_col}')
```

## Plot missing values for train and test datasets

```
In [42]: # Plot chart
df = pd.DataFrame({
    'feature': xlabels,
    'train': y1values,
    'test': y2values
})

df.plot(x='feature',
        kind='bar',
        title = "Missing values %")
plt.show()
```

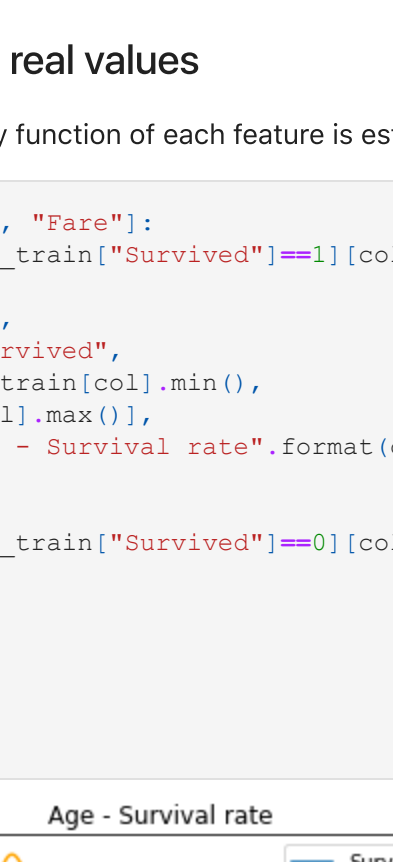


## Distribution of target feature (Survived)

```
In [48]: df_dist = df_train["Survived"].value_counts()
print(df_dist)
print('Not Survived: ' + str(np.array(df_dist)[0]/np.array(df_dist).sum()) + "%")
print('Survived: ' + str(np.array(df_dist)[1]/np.array(df_dist).sum()) + "%")

0 57226
1 42774
Name: Survived, dtype: int64
Not Survived: 0.57226%
Survived: 0.42774%
```

```
In [50]: plt.bar(["Not Survived", "Survived"], np.array(df_dist))
plt.show()
```



## Survival rates according to parameters

### Parameters with categorical values

Two charts are plotted for each feature.

The first chart compares survival and non-survival rates.

The second chart shows survival rates given that a passenger belongs to a certain class.

Take the "Sex" charts as example. The first shows that around 30% of the passengers that survived are female. The second chart shows that around 70% of the females survived.

```
In [65]: total_passengers = len(df_train)
for col in df_train.columns:
    if col not in ["Sex", "Pclass", "SibSp", "Parch", "Embarked", "Cabin_code"]:
        continue

    df_survived = df_train.loc[df_train["Survived"] == 1][col].value_counts()
    df_survived_perc = df_survived / total_passengers
    df_not_survived = df_train.loc[df_train["Survived"] == 0][col].value_counts()
    df_not_survived_perc = df_not_survived / total_passengers

    xlabels = []
    y1values = []
    y2values = []
    for k, v in df_survived_perc.items():
        xlabels.append(k)
        y1values.append(df_survived_perc[k])
        y2values.append(df_not_survived_perc[k])

    df = pd.DataFrame({
        col: xlabels,
        'Survived': y1values,
        'Not Survived': y2values
    })

    df.plot(x=col,
            y=["Survived", "Not Survived"],
            kind="bar",
            title = ("Survival rate" + df_survived_perc[0].format(col) +
                    " given that a passenger has column value x")
    plt.show()

    df_total = df_train[col].value_counts()
    df_survived_perc = df_survived / df_total
    df_survived_perc.plot(x=df_survived_perc.index,
                          y=df_survived_perc.values,
                          kind="bar",
                          color = "orange",
                          title = ("Survival rate" + df_survived_perc[0].format(col) +
                                  " given that a passenger has column value x")
    plt.show()
    print()
```

SibSp - Survival rate with respect to total passengers

SibSp - survival rate given that a passenger has column value x

Parch - Survival rate with respect to total passengers

Parch - survival rate given that a passenger has column value x

Embarked - Survival rate with respect to total passengers

Embarked - survival rate given that a passenger has column value x

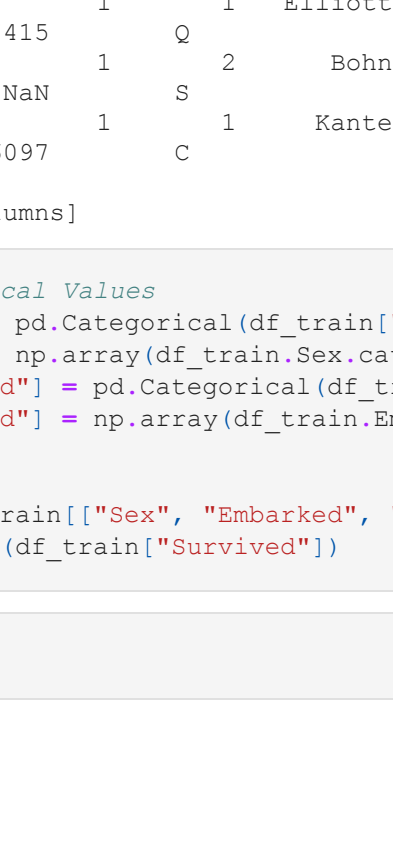
### Parameters with real values

The probability density function of each feature is estimated using KDE (Kernel Density Estimation)

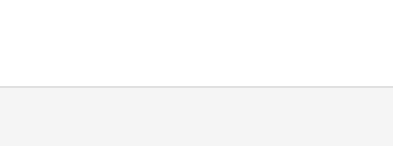
```
In [63]: for col in ["Age", "Fare"]:
    df_train.loc[df_train["Survived"]==1][col].plot(
        kind='kde',
        legend=True,
        label = "Survived",
        color = "orange",
        xlim = [df_train[col].min(),
                df_train[col].max()],
        title = ("Survival rate" + df_train[col].max() +
                " given that a passenger has column value x")
    )

    df_train.loc[df_train["Survived"]==0][col].plot(kind='kde',
                                                    legend=True,
                                                    label = "Not Survived",
                                                    color = "orange",
                                                    xlim = [df_train[col].min(),
                                                            df_train[col].max()])
    plt.show()
```

Age - Survival rate



Fare - Survival rate





```
# Imputation of nan values
df_train['Embarked'] = df_train['Embarked'].replace(-1, np.nan)
from sklearn.impute import KNNImputer
imputer = KNNImputer(n_neighbors=2, weights='uniform')
X = imputer.fit_transform(X)
```

```
# Split into train and validation datasets
X_train, X_valid, y_train, y_valid = train_test_split(X, labels, test_size=0.33)
```

## Define and compile Classification model

```
# Define the model
model = Sequential()
model.add(Dense(20,
                kernel_initializer='random_normal',
                bias_initializer='zeros',
                input_dim=3,
                kernel_regularizer=tf.keras.regularizers.l2(0.01),
                activation='relu'))
model.add(Dense(1, activation='sigmoid'))
```

```
model.summary()
```

Model: "sequential_2"		
Layer (type)	Output Shape	Param #
=====		
dense_4 (Dense)	(None, 20)	80
=====		
dense_5 (Dense)	(None, 1)	21
=====		
Total params:	101	
Trainable params:	101	
Non-trainable params:	0	

```
model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
```

```
opt = tf.optimizers.Adam(learning_rate=0.0001)
model.compile(optimizer=opt, loss='binary_crossentropy', metrics=['accuracy'])
```

## Train the model

```
history = model.fit(X_train,
                    y_train,
                    epochs=10,
                    batch_size=32,
                    verbose=1,
                    validation_split=0.2
                    # validation_data=(X_valid, y_valid)
                    )
```

Epoch 1/10 [=====] - 2s lms/step - loss: 0.6675 - accuracy: 0.5692 - val\_loss: 0.6551 - val\_accuracy: 0.5695  
Epoch 2/10 [=====] - 2s lms/step - loss: 0.6503 - accuracy: 0.5707 - val\_loss: 0.6461 - val\_accuracy: 0.5695  
Epoch 3/10 [=====] - 2s lms/step - loss: 0.6424 - accuracy: 0.5706 - val\_loss: 0.6389 - val\_accuracy: 0.5700  
Epoch 4/10 [=====] - 2s lms/step - loss: 0.6363 - accuracy: 0.5693 - val\_loss: 0.6319 - val\_accuracy: 0.6263  
Epoch 5/10 [=====] - 2s lms/step - loss: 0.6286 - accuracy: 0.6288 - val\_loss: 0.6251 - val\_accuracy: 0.6261  
Epoch 6/10 [=====] - 2s lms/step - loss: 0.6196 - accuracy: 0.6336 - val\_loss: 0.6184 - val\_accuracy: 0.6614  
Epoch 7/10 [=====] - 2s lms/step - loss: 0.6134 - accuracy: 0.6710 - val\_loss: 0.6117 - val\_accuracy: 0.6614  
Epoch 8/10 [=====] - 2s lms/step - loss: 0.6076 - accuracy: 0.6701 - val\_loss: 0.6051 - val\_accuracy: 0.6764  
Epoch 9/10 [=====] - 2s lms/step - loss: 0.5978 - accuracy: 0.6949 - val\_loss: 0.5988 - val\_accuracy: 0.6865  
Epoch 10/10 [=====] - 2s lms/step - loss: 0.5934 - accuracy: 0.6971 - val\_loss: 0.5929 - val\_accuracy: 0.6865

## Evaluate the model on validation data

```
loss = model.evaluate(X_valid, y_valid, verbose=0)
print(loss)
```

[0.588466465473175, 0.6954545378684998]

## Make predictions os validation data

```
y_valid
```

```
array([1, 0, 1, ..., 0, 1, 1])
```

```
yhat = model.predict(X_valid)
print(yhat)
```

```
[[0.41861227]
 [0.2769578 ]
 [0.5777548 ]
 [0.29545367]
 [0.5777548 ]
 [0.45967266]]
```

```
y_valid
```

```
array([1, 0, 1, ..., 0, 1, 1])
```

```
y_class = np.where(yhat > 0.5, 1, 0)
```

```
y_class = y_class.squeeze()
```

```
y_class
```

```
array([0, 0, 1, ..., 0, 1, 0])
```

## Predict Labels on Test Set

```
df_test = pd.read_csv('/content/playground/test.csv')
print(df_test.columns)
print(len(df_test.columns))
print(df_test.head())
```

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	Cabin Embarked	3	Holliday, Daniel	male	19.0	0	0	24745	69.01
1	NaN	S							
1	100001	3	Nguyen, Lorraine	female	53.0	0	0	13264	5.81
2	NaN	S							
2	100002	1	Harris, Heather	female	19.0	0	0	25990	38.91
15315	100003	C							
3	NaN	S							
4	100004	1	Larsen, Eric	male	25.0	0	0	314011	12.93
22515	100004	C	Cleary, Sarah	female	17.0	0	2	26203	26.89

```
df_test["Sex"] = pd.Categorical(df_test["Sex"])
df_test["Sex"] = np.array(df_test.Sex.cat.codes)
df_test["Embarked"] = pd.Categorical(df_test["Embarked"])
df_test["Embarked"] = np.array(df_test.Embarked.cat.codes)
```

```
# Create Dataset
X = np.array(df_test[["Sex", "Embarked", "Pclass"]])
```

```
# Imputation of nan values
df_train['Embarked'] = df_train['Embarked'].replace(-1, np.nan)
from sklearn.impute import KNNImputer
imputer = KNNImputer(n_neighbors=2, weights='uniform')
X = imputer.fit_transform(X)
```

```
df_test.isnull().any()
```

```
# Create Test Dataset
X_test = np.array(df_test[["Sex", "Embarked", "Pclass"]])
Pids = np.array(df_test["PassengerId"])
```

```
Y_hat = model.predict(X_test)
y_class = np.where(Y_hat > 0.5, 1, 0)
```

```
df_classif = pd.DataFrame(data=[Pids, Y_class.squeeze()]).T
df_classif.columns = ["PassengerId", "Survived"]
```

```
df_classif.head()
```

```
df_classif.to_csv("classifications.csv", index=False)
```

## Draft

```
# optional model
# create an instance of a neural network:
k_model = Sequential()
n = 3
# the first hidden layer must have input dimensions:
k_model.add(Dense(10, activation='relu',
                  kernel_regularizer=tf.keras.regularizers.l2(0.01),
                  input_dim=n))
# additional hidden layers are optional:
k_model.add(Dense(20, activation='relu',
                  kernel_regularizer=tf.keras.regularizers.l2(0.01)))
# the output layer- a binary classifier w/ sigmoid activation:
k_model.add(Dense(1, activation='sigmoid',
                  kernel_regularizer=tf.keras.regularizers.l2(0.01)))

# Compile the model with Adam optimizer:
k_model.compile(optimizer=Adam(lr=1e-1),
                loss='binary_crossentropy',
                metrics=['accuracy'])

# Define a learning rate decay method:
# lr_decay = ReduceLRonPlateau(monitor='loss',
#                               patience=1,
#                               verbose=0,
#                               factor=0.5,
#                               min_lr=1e-7)

# Train the model:
k_model.fit(X_train,
            y_train,
            epochs=20,
            callbacks=[lr_decay],
            verbose=1,
            validation_data=(X_valid, y_valid)
            )
```