

DMW.p7jpyb

File Edit View Insert Runtime Tools Help

All changes saved

Comment

Share

RAM

Disk

+ Code

+ Text

Q

(x)

🔍

📄

21162101012_CBA_Yash_Lakhtariya

DMW Practical 7

From the given dataset, you need to create a classification model which will predict whether the particular person has survived the tragic incident or not. And to create the model, you first have to convert the dataset into appropriate form, e.g - remove the special characters, remove null/duplicate entries, etc. A decision tree should be constructed from the given data. More, to test the model if any new tuple assigned to it then, it should be able to classify that the given detail of the person has actually survived the incident or not.

[1] import pandas as pd
from bs4 import BeautifulSoup

[4] from google.colab import files
uploaded = files.upload()

[5] file = uploaded['titanicdata.htm']
soup = BeautifulSoup(file,'html.parser')

[6] table = soup.find('table')

[7] data=pd.read_html(str(table),flavor='bs4')[0]
data.head(10)

[8] data.info()

[9] import re
data['Name']=data['Name'].map(str).apply(lambda x:x.encode('utf-8').decode('ascii','ignore'))
data['Boat [Body]']=data['Boat [Body]'].map(str).apply(lambda x:x.encode('utf-8').decode('ascii','ignore'))

[10] data.head(10)
df=data

[11] def proc_age(value):
 if'm' in value:
 return float(re.findall(r'-(?d+\.?d+)',value)[0])/12
 else:
 return(float(value))

data['Age']= data['Age'].map(str).apply(proc_age)
data.head()

[12] data_trim = data[['Name', 'Age', 'Class/Dept', 'Boat [Body]']]
data_trim.head()

[13] def categories_person(value):
 if'PASSENGER' in value.upper():
 return 'Passenger'
 else:
 return 'Crew'

data_trim['Crew/Passenger']=data_trim['Class/Dept'].map(str).apply(categories_person)
data_trim.head()

[14] def categories_age(age):
 if age>18:
 return 'Adult'
 else:
 return 'Child'

data_trim['Adult/Child']=data_trim['Age'].apply(categories_age)
data_trim.head()

[15] def categories_gender(name):
 firstname= name[name.index(',')+2:]
 saluation = firstname.split(" ")[0]
 if saluation.upper() in ['MR','MASTER']:
 return 'Male'
 else:
 return 'Female'

data_trim['Gender']=data_trim['Name'].map(str).apply(categories_gender)
data_trim.head()

[16] def categories_survival(value):
 if value.strip()=='nan' or '[' in value:
 return 0
 else:
 return 1

data_trim['Survival']=data_trim['Boat [Body]'].map(str).apply(categories_survival)
data_trim.head(15)

Along with that try to analyze the given data by answering the below mentioned questions:

[19] from google.colab import files
datafile = files.upload()

[20] data=pd.read_csv('train.csv')
data.head()

[21] mean_age = data['Age'].mean()
data['Age'].fillna(mean_age, inplace=True)
data_trim = data[['Name', 'Age', 'Pclass', 'Survived']]
data_trim.head()

[22] from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score

[23] X = data_trim[['Age', 'Pclass']]
y = data_trim['Survived']

Train Decision Tree Classifier
clf = DecisionTreeClassifier()
clf.fit(X, y)

Test the model (predict survival for new data)
new_data = {
 "Name": ["New Passenger"],
 "Age": [30.0],
 "Pclass": [2]
}

new_data_trim = pd.DataFrame(new_data)
X_new = new_data_trim[['Age', 'Pclass']]
prediction = clf.predict(X_new)

if prediction[0] == 1:
 print("The passenger is predicted to have survived.")
else:
 print("The passenger is predicted to have not survived.")

The passenger is predicted to have not survived.

What is the survival percentage of Crew members ?

[24] crew_survival_percentage = data_trim['Survived'].mean() * 100
print("Survival percentage of Crew members:", crew_survival_percentage)

What is the survival percentage of first class passengers?

[25] first_class_survival_percentage = data_trim[data_trim['Pclass'] == 1]['Survived'].mean() * 100
print("Survival percentage of first class passengers:", first_class_survival_percentage.round(2))

Does age also play a role in survival ? How can you say that ?

[26] survivors_mean_age = data_trim[data_trim['Survived'] == 1]['Age'].mean()
non_survivors_mean_age = data_trim[data_trim['Survived'] == 0]['Age'].mean()
import matplotlib.pyplot as plt

plt.figure(figsize=(8, 5))
plt.hist(data_trim[data_trim['Survived'] == 1]['Age'], bins=10, alpha=0.5, label='Survived')
plt.hist(data_trim[data_trim['Survived'] == 0]['Age'], bins=10, alpha=0.5, label='Did not survive')
plt.xlabel('Age')
plt.ylabel('Count')
plt.title('Age Distribution of Survivors vs Non-survivors')
plt.legend()
plt.show()

Age Distribution of Survivors vs Non-survivors

correlation_coefficient = data_trim['Age'].corr(data_trim['Survived'])
print("Pearson correlation coefficient between Age and Survived:", correlation_coefficient)

Pearson correlation coefficient between Age and Survived: -0.06980851528714314

Here, the correlation of the age and survived is negative.

Hence, it can be said that more age means less chances of being survived and vice versa, but as it near to zero, age has negligible effect of survival chances

Colab paid products - Cancel contracts here

0s completed at 8:49 AM