

Price Prediction using Regression

This is a tickets price monitoring system. It scrapes tickets pricing data periodically and stores it in a database. Ticket pricing changes based on demand and time, and there can be significant difference in price. We are creating this product mainly with ourselves in mind. Users can set up alarms using an email, choosing an origin and destination (cities), time (date and hour range picker) choosing a price reduction over mean price, etc.

Following is the description for columns in the dataset

- insert_date: date and time when the price was collected and written in the database
- origin: origin city
- destination: destination city
- start_date: train departure time
- end_date: train arrival time
- train_type: train service name
- price: price
- train_class: ticket class, tourist, business, etc.
- fare: ticket fare, round trip, etc

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
sns.set()
sns.set()
```

Task 1: Import Dataset and create a copy of that dataset

```
#write code here
data = pd.read_csv('data1.csv')
df = data.copy()
```

Task 2: Display first five rows

```
#write code here
df.head()
```

	Unnamed: 0	insert_date	origin	destination	start_date	end_date	train_type	price	train_class	fare
0	702	2019-04-19 05:37:35	PONFERRADA	MADRID	2019-06-02 16:00:00	2019-06-02 18:42:00	MD-AVE	66.50	Turista con enlace	Flexible
1	703	2019-04-19 05:37:35	PONFERRADA	MADRID	2019-06-02 17:15:00	2019-06-02 23:03:00	MD-AVE	34.65	Turista con +	Promo
2	704	2019-04-19 05:37:35	PONFERRADA	MADRID	2019-06-02 17:15:00	2019-06-02 23:10:00	MD-LD	38.95	Turista con enlace	Promo
3	705	2019-04-19 05:37:35	PONFERRADA	MADRID	2019-06-02 17:15:00	2019-06-02 22:14:00	MD-AVE	40.60	Turista con enlace	Promo +
4	706	2019-04-19 05:37:35	PONFERRADA	MADRID	2019-06-02 18:55:00	2019-06-02 23:03:00	ALVIA	27.90	Turista	Promo

Task 3: Drop 'Unnamed: 0' column

```
#write code here
df.drop('Unnamed: 0', axis=1, inplace=True)
```

Task 4: Check the number of rows and columns

```
#write code here
df.shape
```

```
(21509, 9)
```

Task 5: Check data types of all columns

```
#write code here
df.dtypes
```

```
insert_date    object
origin         object
destination    object
start_date     object
end_date       object
train_type     object
price          float64
train_class    object
fare           object
dtype: object
```

Task 6: Check summary statistics

```
#write code here
df.describe()
```

```
count    202321.000000
mean     56.723877
std      25.531787
min      16.600000
25%      28.350000
50%      53.400000
75%      76.300000
max      206.800000
```

Task 7: Check summary statistics of all columns, including object datatypes

```
df.describe(include='all')
```

	insert_date	origin	destination	start_date	end_date	train_type	price	train_class	fare
count	215909	215909	215909	215909	215909	215909	202321.000000	215266	215266
unique	30543	5	5	2231	2870	16	NaN	5	5
top	2019-05-09 05:37:35	MADRID	MADRID	2019-06-02 17:20:00	2019-06-02 23:03:00	AVE	NaN	Turista	Promo
freq	90	110440	105469	2089	1278	126577	NaN	164016	132085
mean	NaN	NaN	NaN	NaN	NaN	NaN	56.723877	NaN	NaN
std	NaN	NaN	NaN	NaN	NaN	NaN	25.531787	NaN	NaN
min	NaN	NaN	NaN	NaN	NaN	NaN	16.600000	NaN	NaN
25%	NaN	NaN	NaN	NaN	NaN	NaN	28.350000	NaN	NaN
50%	NaN	NaN	NaN	NaN	NaN	NaN	53.400000	NaN	NaN
75%	NaN	NaN	NaN	NaN	NaN	NaN	76.300000	NaN	NaN
max	NaN	NaN	NaN	NaN	NaN	NaN	206.800000	NaN	NaN

Question: Explain the summary statistics for the above data set

Answer: From this summary statistics we conclude that there are total 215909 set of rows from which some rows of price, train_class and fare are empty or contain null values. Also average price for ticket is 56 and max price for ticket is 209. There are 5 unique values for both origin and destination which means there are total 5 stations available in given dataset

Task 8: Check null values in dataset

```
#write code here
df.isnull().sum()
```

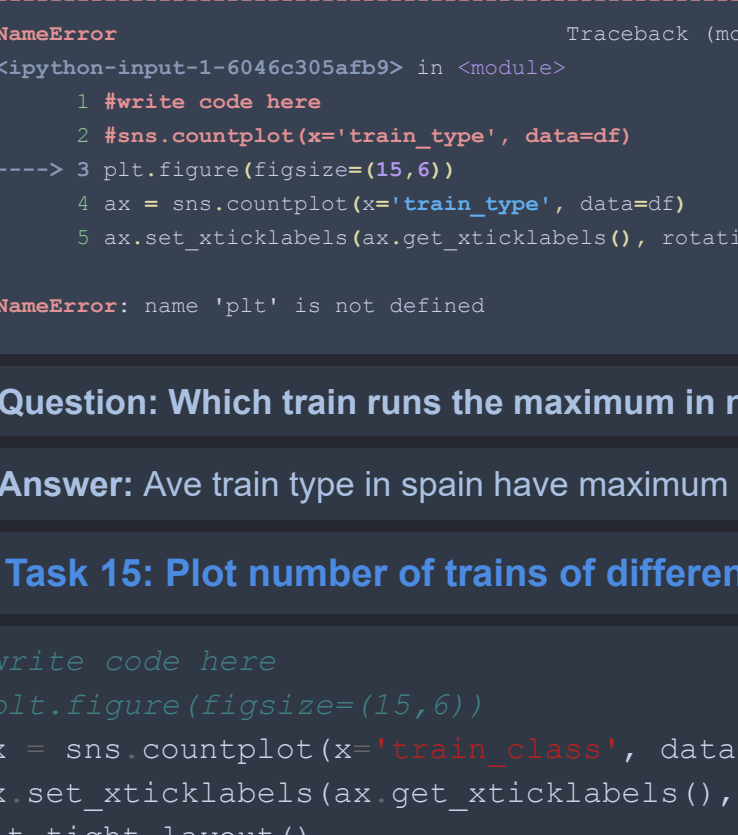
```
insert_date    0
origin         0
destination    0
start_date     0
end_date       0
train_type     0
price         13588
train_class    643
fare          643
dtype: int64
```

Task 9: Fill the Null values in the 'price' column.

Data in price column is right skewed so we replace null values with median

```
sns.distplot(df['price'])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x53f025e08>
```



```
df['price'].fillna(df['price'].median(),inplace=True)
```

Task 10: Drop the rows containing Null values in the attributes train_class and fare

```
#write code here
df.dropna(subset=['train_class','fare'],inplace=True)
```

Task 11: Drop 'insert_date'

```
#write code here
df.drop('insert_date', axis=1, inplace=True)
df.shape
```

```
(215266, 8)
```

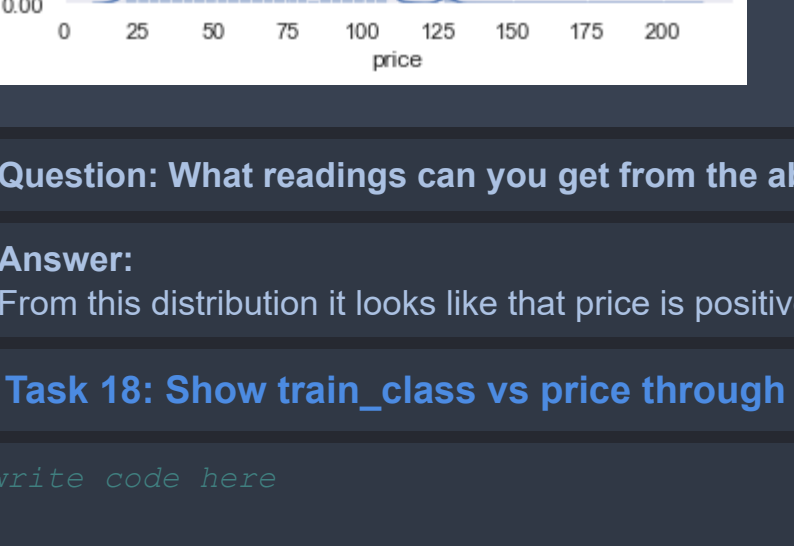
Check null values again in dataset

```
#write code here
df.isnull().sum()
```

```
origin         0
destination    0
start_date     0
end_date       0
train_type     0
price          0
train_class    0
fare           0
dtype: int64
```

Task 12: Plot number of people boarding from different stations

```
#write code here
sns.countplot(x='origin', data=df)
```

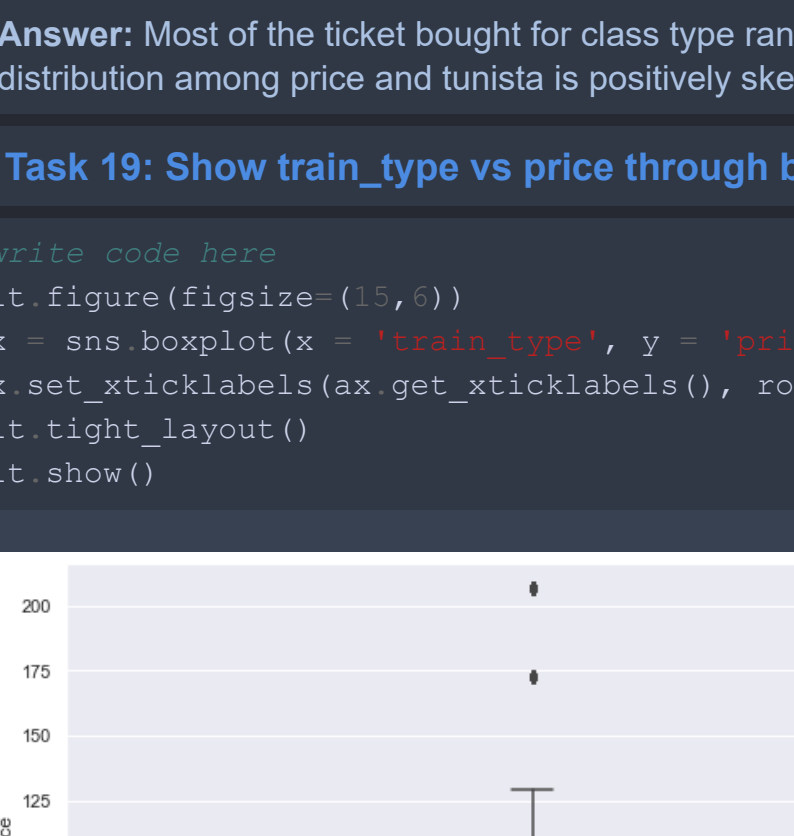


Question: What insights do you get from the above plot?

Answer: Most common station for boarding is madrid, around 150,000 or more people board from this station. Ponferrada have very less number of people boarding as compared to other stations. After madrid Barcelona have more number of boardings.

Task 13: Plot number of people for the destination stations

```
sns.countplot(x='destination', data=df)
```



Question: What insights do you get from the above graph?

Answer: Madrid is the most common station used for travelling there were also huge number of people boarding from this station and as well as huge number of people getting off this station and very less number of people are getting off from Ponferrada station

Task 14: Plot different types of train that runs in Spain

```
#write code here
#sns.countplot(x='train_type', data=df)
plt.figure(figsize=(15,5))
ax = sns.countplot(x='train_type', data=df)
ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha='right')
plt.tight_layout()
plt.show()
```

```
-----
NameError                                Traceback (most recent call last)
<ipython-input-1-6046305ab0> in <module>
      1 #write code here
      2 sns.countplot(x='train_type', data=df)
----> 3 plt.figure(figsize=(15,5))
      4 ax = sns.countplot(x='train_type', data=df)
      5 ax.set_xticklabels(ax.get_xticklabels(), rotation=40, ha='right')

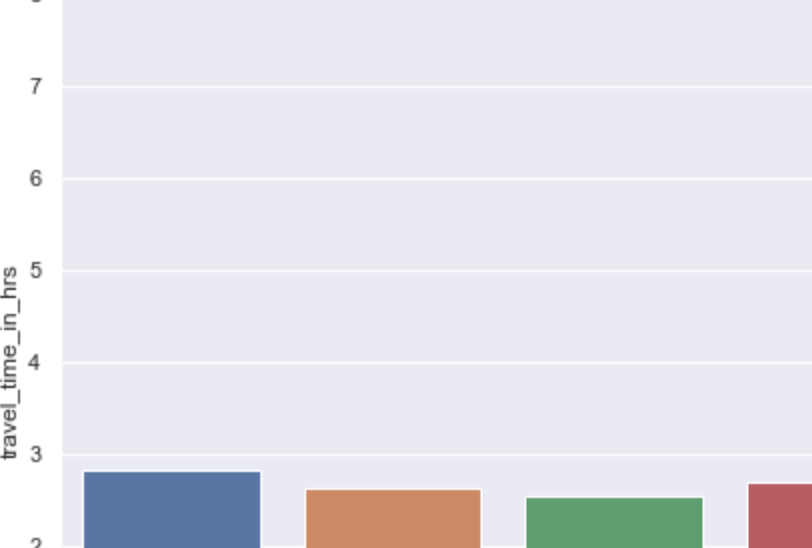
NameError: name 'plt' is not defined
```

Question: Which train runs the maximum in number as compared to other train types?

Answer: Ave train type in Spain have maximum number as compared to other train type

Task 15: Plot number of trains of different class

```
#write code here
plt.figure(figsize=(15,5))
ax = sns.countplot(x='train_class', data=df)
ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha='right')
plt.tight_layout()
plt.show()
```



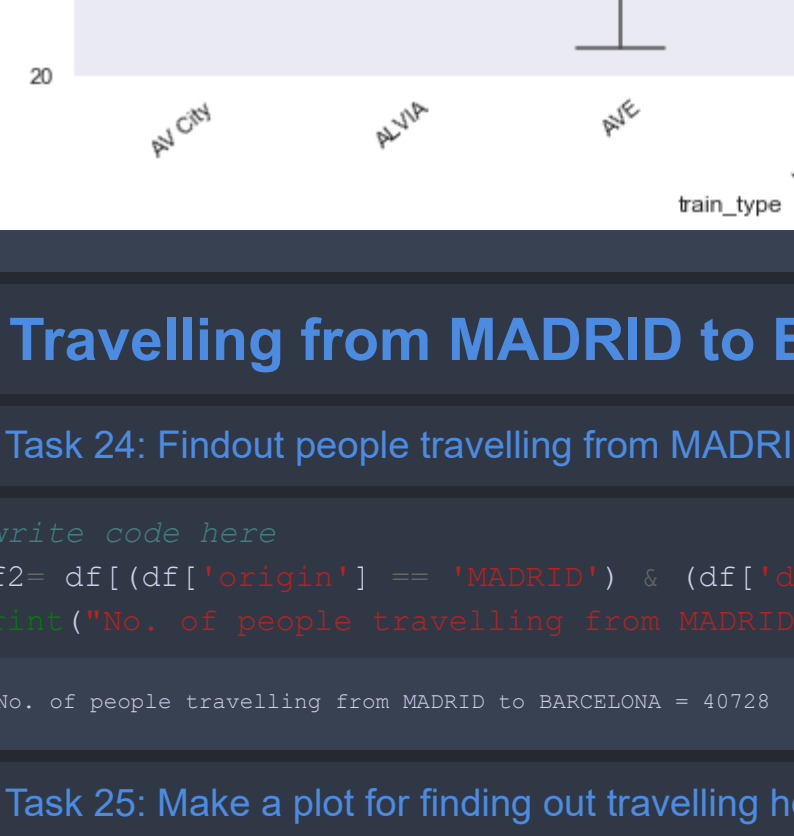
Question: Which the most common train class for traveling among people in general?

Answer: Turista is most common train class for traveling among people in Spain.

Task 16: Plot number of tickets bought from each category

```
#write code here
sns.countplot(x='fare', data=df)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x53f0d5d08>
```



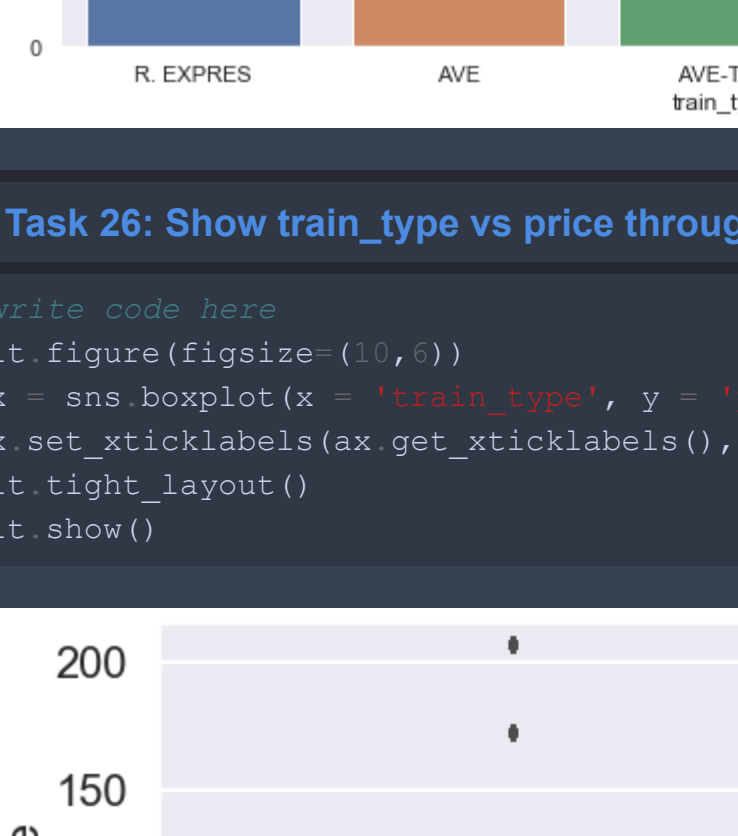
Question: Which the most common tickets are bought?

Answer: Promo tickets are the most common ticket bought among people

Task 17: Plot distribution of the ticket prices

```
#write code here
sns.distplot(df['price'])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x53f1b5e08>
```

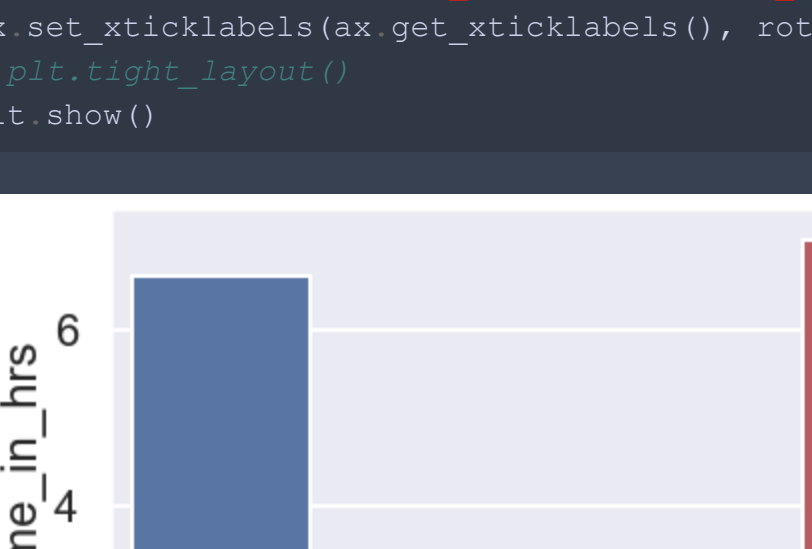


Question: What readings can you get from the above plot?

Answer: From this distribution it looks like that price is positively skewed

Task 18: Show train_class vs price through boxplot

```
#write code here
ax = sns.boxplot(x='train_class', y='price', data = df)
ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha='right')
plt.tight_layout()
plt.show()
```

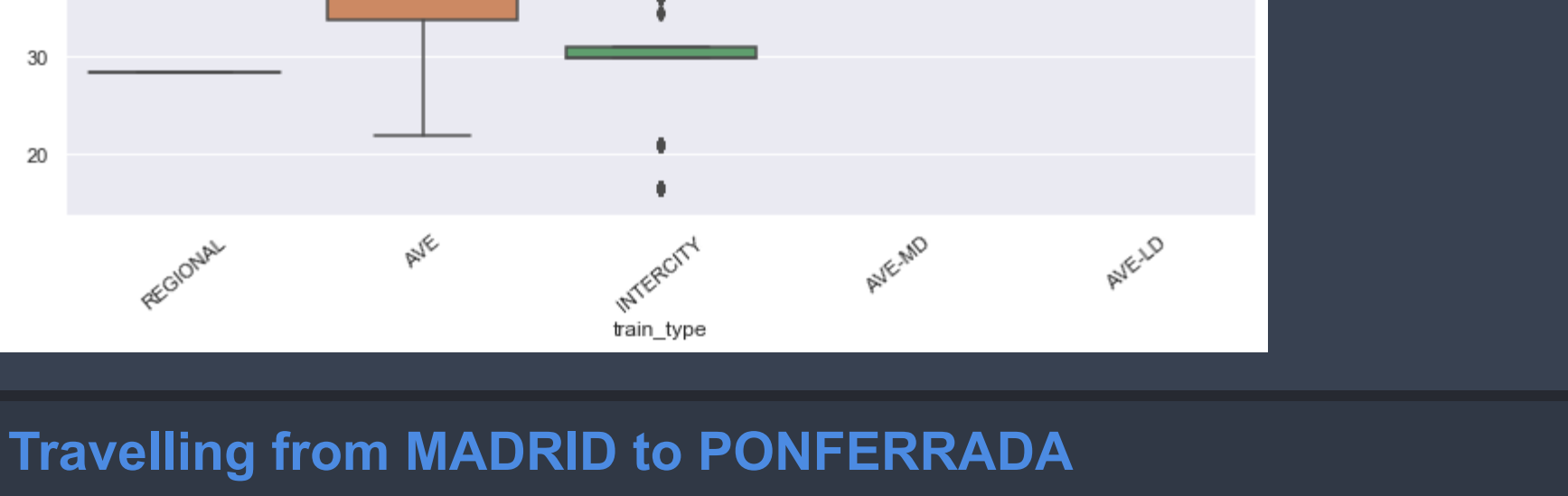


Question: What ticket trends can you find out by looking at the plot above?

Answer: Most of the ticket trends for class type range from 40 to 80. Premium class type have normal distribution among price and turista is positively skewed and turista plus is negatively skewed

Task 19: Show train_type vs price through boxplot

```
#write code here
plt.figure(figsize=(15,5))
ax = sns.boxplot(x='train_type', y='price', data = df)
ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha='center')
plt.tight_layout()
plt.show()
```



Question: Which type of trains cost more as compared to others?

Answer: Ave train type cost more as compared to other train type.

Feature Engineering

```
df = df.reset_index()
```

Finding the travel time between the place of origin and destination. We need to find out the travel time for each entry which can be obtained from the 'start_date' and 'end_date' column. Also if you see, these columns are in object type therefore datetimeFormat should be defined to perform the necessary operation of getting the required time.

Import datetime library

```
#write code here
import datetime
```

```
datetimeFormat = '%Y-%m-%d %H:%M:%S'
a,b = datetime.datetime.strptime(b, datetimeFormat) - datetime.datetime.strptime(a, datetimeFormat)
a = a + datetime.timedelta(seconds=3600.0)
```

```
df['travel_time_in_hrs'] = df.apply(lambda x:fun(x['start_date'],x['end_date']),axis=1)
```

Task 20: Remove redundant features

You need to remove features that are giving the related values as 'travel_time_in_hrs' Hint: Look for date related columns

```
#write code here
df.drop(['start_date','end_date'], axis=1, inplace=True)
```

We now need to find out the pricing from 'MADRID' to other destinations. We also need to find out time which each train requires for travelling.

Travelling from MADRID to SEVILLA

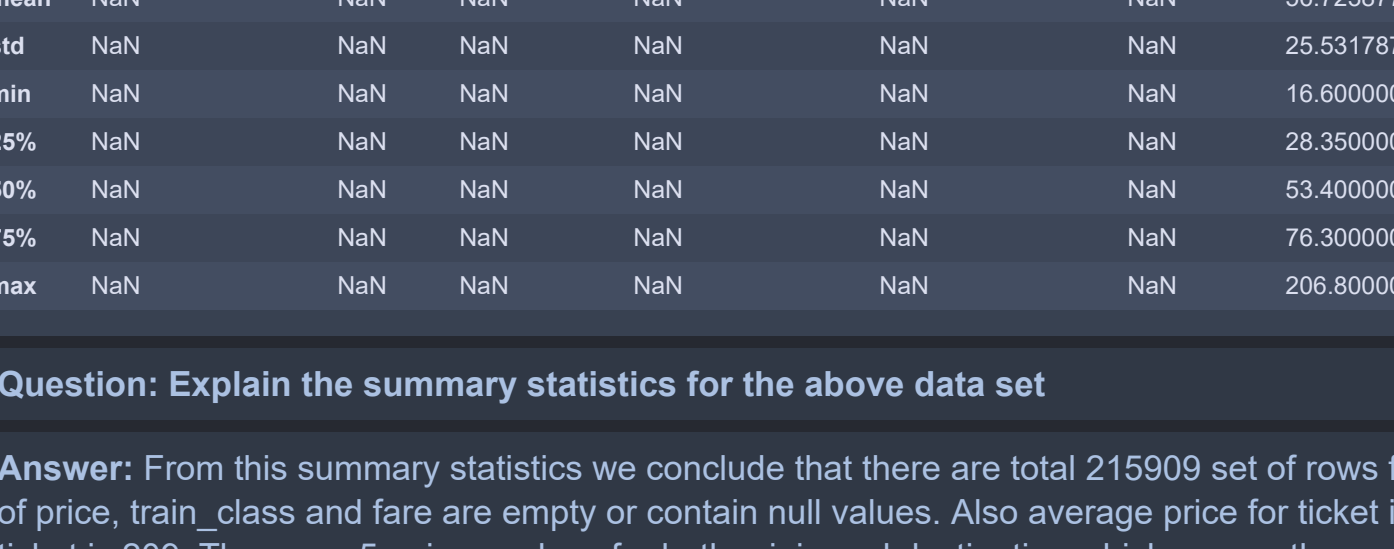
Task 21: Findout people travelling from MADRID to SEVILLA

```
#write code here
df1 = df[(df['origin'] == 'MADRID') & (df['destination'] == 'SEVILLA')]
print('No. of people travelling from MADRID to SEVILLA = '+str(df1.shape[0]))
```

```
No. of people travelling from MADRID to SEVILLA = 26361
```

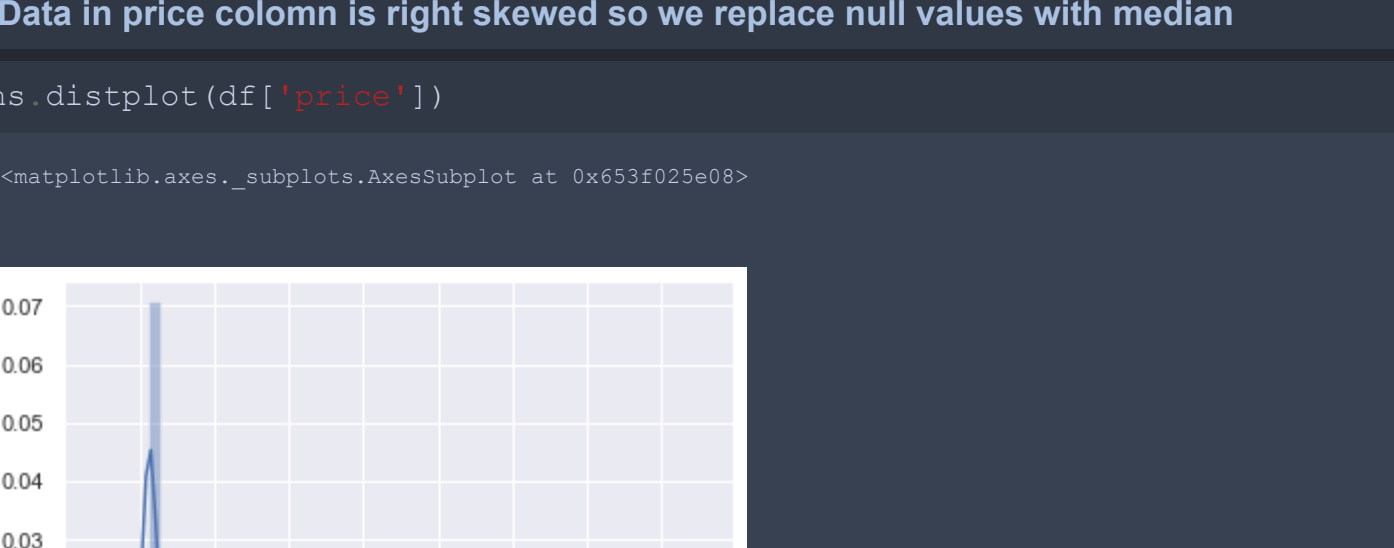
Task 22: Make a plot for finding out travelling hours for each train type

```
#write code here
plt.figure(figsize=(10,5))
ax = sns.barplot(x='train_type', y='travel_time_in_hrs', data=df1, ci=0)
ax.set_xticklabels(ax.get_xticklabels(), rotation=0, ha='center')
plt.tight_layout()
plt.show()
```



Task 23: Show train_type vs price through boxplot

```
#write code here
plt.figure(figsize=(10,5))
ax = sns.boxplot(x='train_type', y='price', data = df1)
ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha='center')
plt.tight_layout()
plt.show()
```



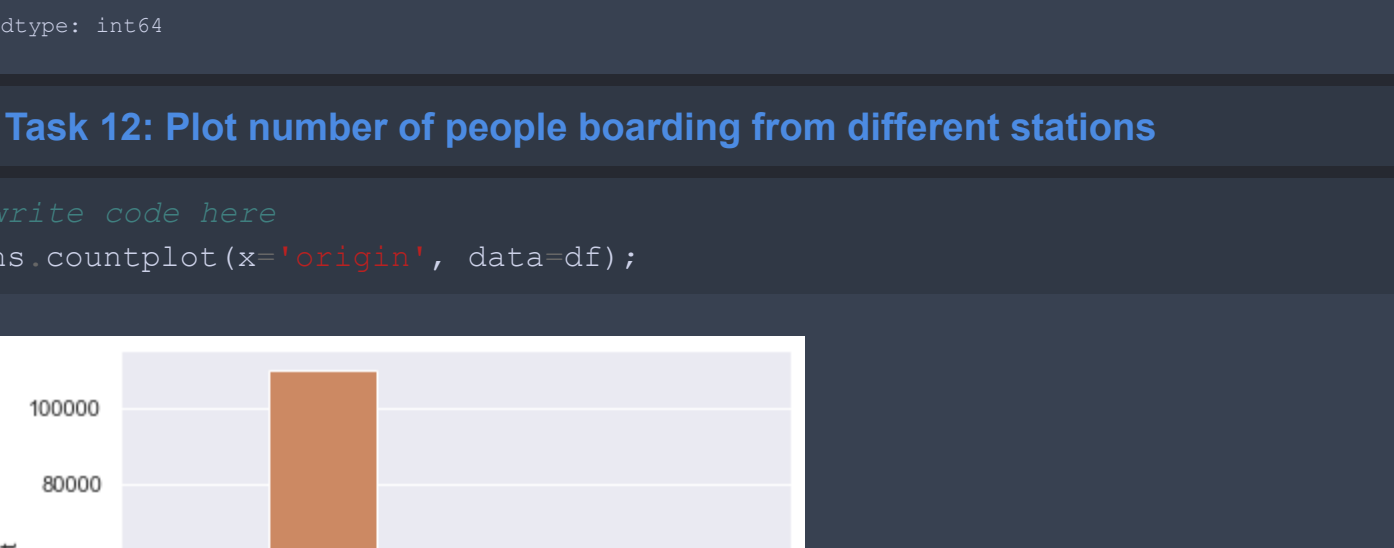
Task 24: Findout people travelling from MADRID to BARCELONA

```
#write code here
df2 = df[(df['origin'] == 'MADRID') & (df['destination'] == 'BARCELONA')]
print('No. of people travelling from MADRID to BARCELONA = '+str(df2.shape[0]))
```

```
No. of people travelling from MADRID to BARCELONA = 40728
```

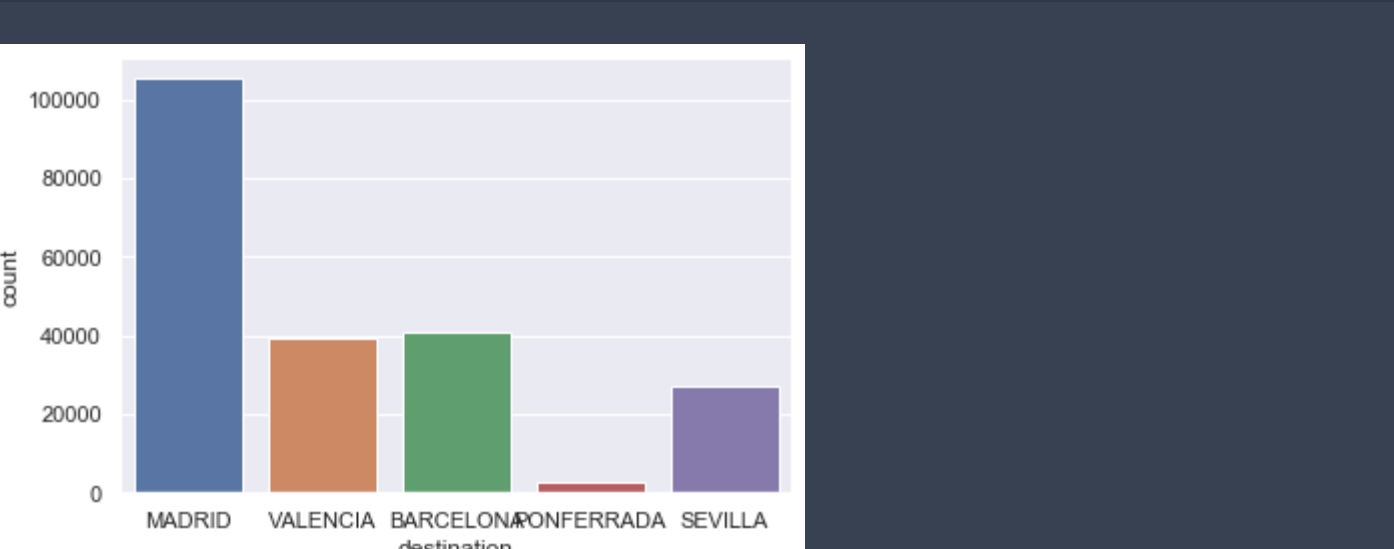
Task 25: Make a plot for finding out travelling hours for each train type

```
#write code here
plt.figure(figsize=(10,5))
ax = sns.barplot(x='train_type', y='travel_time_in_hrs', data=df2, ci=0)
ax.set_xticklabels(ax.get_xticklabels(), rotation=0, ha='center')
plt.tight_layout()
plt.show()
```



Task 26: Show train_type vs price through boxplot

```
#write code here
plt.figure(figsize=(10,5))
ax = sns.boxplot(x='train_type', y='price', data = df2)
ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha='center')
plt.tight_layout()
plt.show()
```



Travelling from MADRID to VALENCIA

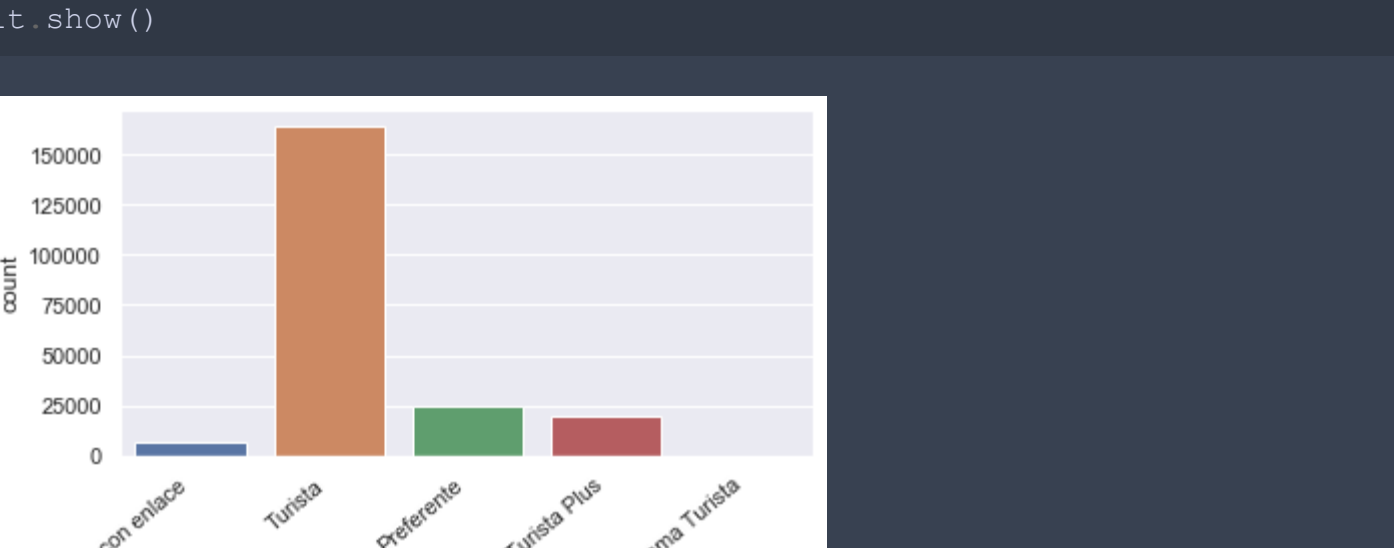
Task 27: Findout people travelling from MADRID to VALENCIA

```
#write code here
df3 = df[(df['origin'] == 'MADRID') & (df['destination'] == 'VALENCIA')]
print('No. of people travelling from MADRID to VALENCIA = '+str(df3.shape[0]))
```

```
No. of people travelling from MADRID to VALENCIA = 39447
```

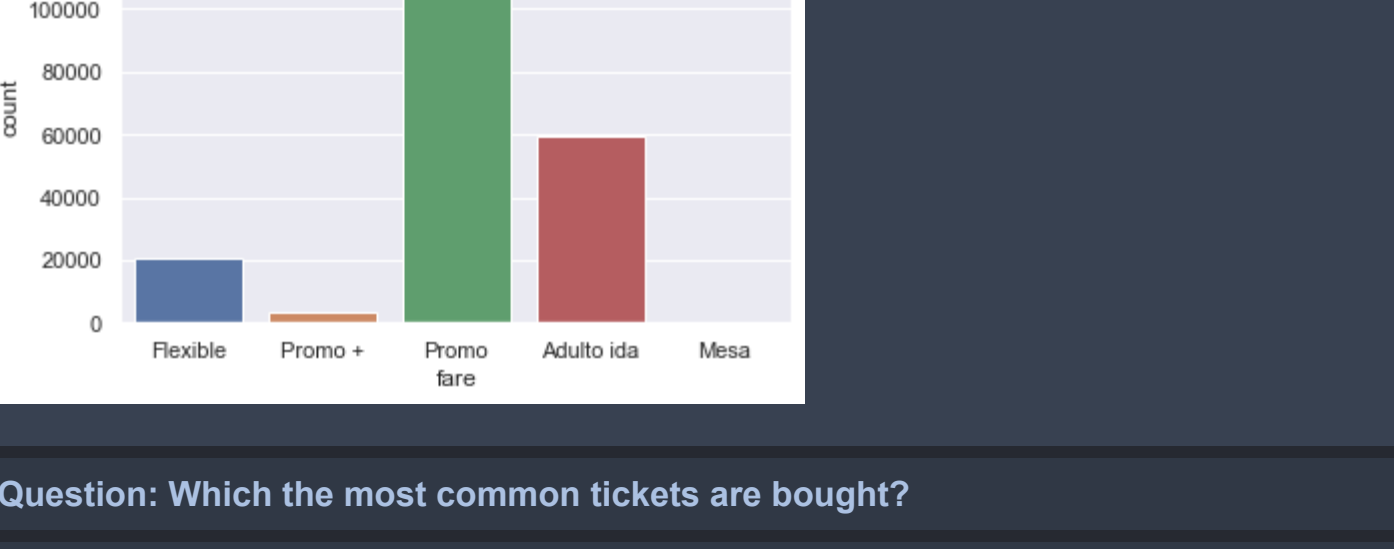
Task 28: Make a plot for finding out travelling hours for each train type

```
#write code here
plt.figure(figsize=(10,5))
ax = sns.barplot(x='train_type', y='travel_time_in_hrs', data=df3, ci=0)
ax.set_xticklabels(ax.get_xticklabels(), rotation=0, ha='center')
plt.tight_layout()
plt.show()
```



Task 29: Show train_type vs price through boxplot

```
#write code here
plt.figure(figsize=(10,5))
ax = sns.boxplot(x='train_type', y='price', data = df3)
ax.set_xticklabels(ax.get_xticklabels(), rotation=45, ha='center')
plt.tight_layout()
plt.show()
```



Travelling from MADRID to PONFERRADA

Task 30: Findout people travelling from MADRID to PONFERRADA

```
#write code here
df4 = df[(df['origin'] == 'MADRID') & (df['destination'] == 'PONFERRADA')]
print('No. of people travelling from MADRID to PONFERRADA = '+str(df4.shape[0]))
```

```
No. of people travelling from MADRID to PONFERRADA = 2839
```

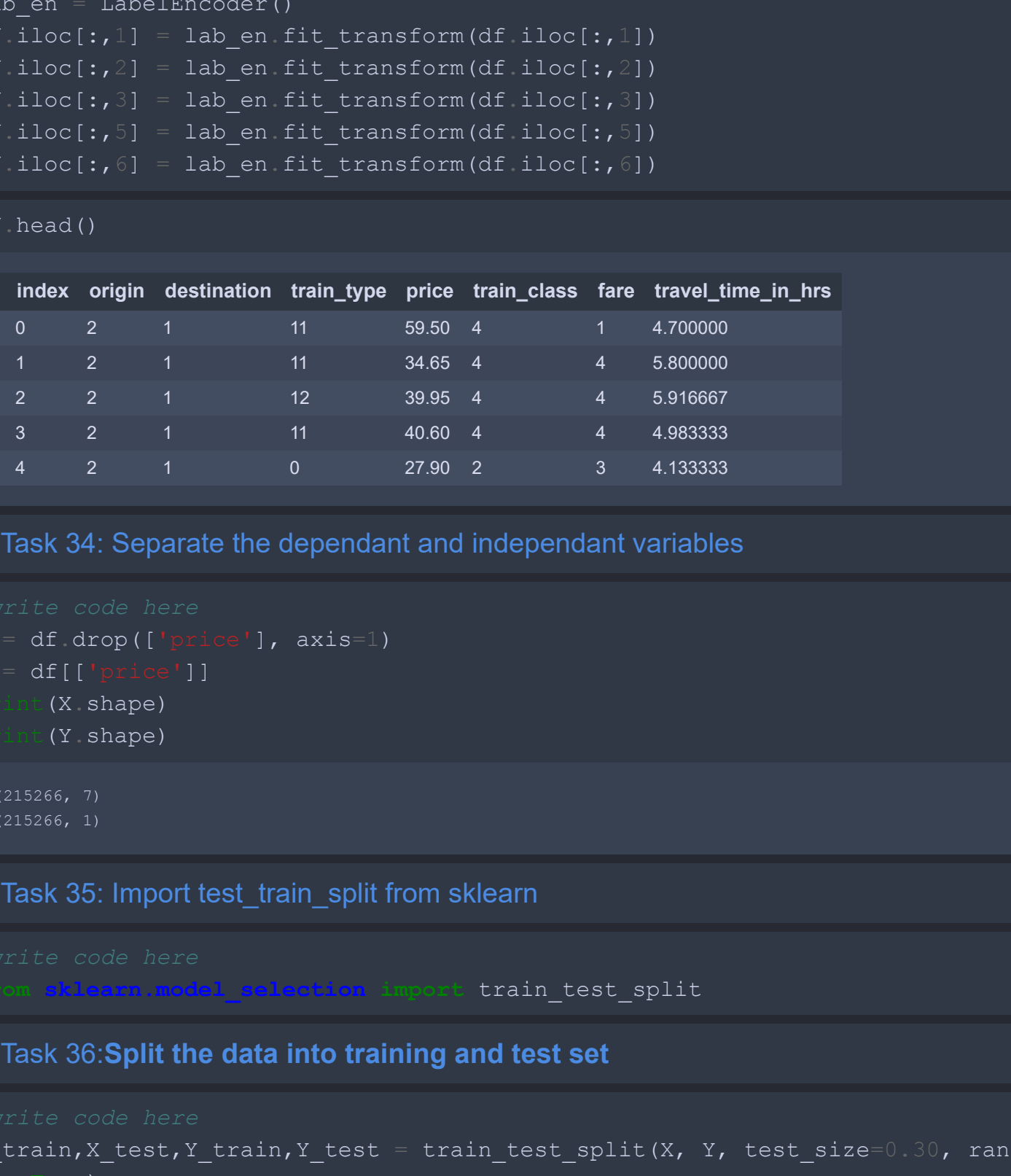
Task 31: Make a plot for finding out travelling hours for each train type


```
#write code here
plt.figure(figsize=(10,5))
ax = sns.barplot(x='train_type', y='travel_time_in_hrs', data=df4, ci=0)
ax.set_xticklabels(ax.get_xticklabels(), rotation=0, ha='center')
plt.tight_layout()
plt.show()
```



Task 32: Show train_type vs price through boxplot

```
#write code here
plt.figure(figsize=(10,5))
ax = sns.boxplot(x='train_type', y='price', data=df4)
ax.set_xticklabels(ax.get_xticklabels(), rotation=0, ha='center')
plt.tight_layout()
plt.show()
```



Applying Linear Regression

Task 33: Import LabelEncoder library from sklearn

```
#write code here
from sklearn.preprocessing import LabelEncoder
```

Data Encoding

```
lab_en = LabelEncoder()
df.iloc[:,1] = lab_en.fit_transform(df.iloc[:,1])
df.iloc[:,2] = lab_en.fit_transform(df.iloc[:,2])
df.iloc[:,3] = lab_en.fit_transform(df.iloc[:,3])
df.iloc[:,4] = lab_en.fit_transform(df.iloc[:,4])
df.iloc[:,5] = lab_en.fit_transform(df.iloc[:,5])
```

```
df.head()
```

	index	origin	destination	train_type	train_class	fare	travel_time_in_hrs
0	0	2	1	11	59.50	4	1.4700000
1	1	2	1	11	34.65	4	5.8000000
2	2	2	1	12	39.95	4	5.916667
3	3	2	1	11	40.80	4	4.983333
4	4	2	1	0	27.90	2	4.133333

Task 34: Separate the dependant and independant variables

```
#write code here
X = df.drop(['price'], axis=1)
Y = df[['price']]
print(X.shape)
print(Y.shape)
```

```
(155866, 7)
```

```
(155866, 1)
```

Task 35: Import test_train_split from sklearn

```
#write code here
from sklearn.model_selection import train_test_split
```

Task 36: Split the data into training and test set

```
#write code here
X_train,X_test,Y_train,Y_test = train_test_split(X, Y, test_size=0.30, random_state=0,shuf
file=False)
```

Task 37: Import LinearRegression library from sklearn

```
#write code here
from sklearn.linear_model import LinearRegression
```

Task 38: Make an object of LinearRegression() and train it using the training data set

```
#write code here
lr = LinearRegression()
lr.fit(X_train, Y_train)
```

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

Task 39: Find out the predictions using test data set.

```
#write code here
lr_predict = lr.predict(X_test)
```

Task 40: Find out the predictions using training data set.

```
#write code here
lr_predict_train = lr.predict(X_train)
```

Task 41: Import r2_score library form sklearn

```
#write code here
from sklearn.metrics import r2_score
```

Task 42: Find out the R2 Score for test data and print it.

```
#write code here
lr_r2_test= r2_score(Y_test,lr_predict)
print('r2 score of testing = ', lr_r2_test)
```

```
r2 score of testing 0.7290143663235672
```

Task 43: Find out the R2 Score for training data and print it.

```
#write code here
lr_r2_train = r2_score(Y_train,lr_predict_train)
print('r2 score of training = ', lr_r2_train)
```

```
r2 score of training 0.7238620495006879
```

Comaparing training and testing R2 scores

```
print('R2 score for Linear Regression Training Data is: ', lr_r2_train)
print('R2 score for Linear Regression Testing Data is: ', lr_r2_test)
```

```
R2 score for Linear Regression Training Data is: 0.7238620495006879
R2 score for Linear Regression Testing Data is: 0.7290143663235672
```

Applying Polynomial Regression

Task 44: Import PolynomialFeatures from sklearn

```
#write code here
from sklearn.preprocessing import PolynomialFeatures
```

Task 45: Make an object of default Polynomial Features

```
#write code here
poly_reg = PolynomialFeatures(degree=1)
```

Task 46: Transform the features to higher degree features.

```
#write code here
X_train_poly,X_test_poly = poly_reg.fit_transform(X_train),poly_reg.fit_transform(X_test)
```

```
print(X_train.shape)
print(X_train_poly.shape)
print(X_test.shape)
print(X_test_poly.shape)
```

```
(155866, 7)
```

```
(155866, 36)
```

```
(64590, 36)
```

Task 47: Fit the transformed features to Linear Regression

```
#write code here
poly_model = LinearRegression()
poly_model.fit(X_train_poly, Y_train)
```

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

Task 48: Find the predictions on the data set

```
#write code here
y_train_predicted,y_test_predict = poly_model.predict(X_train_poly),poly_model.predict(X_t
st_poly)
```

Task 49: Evaluate R2 score for training data set

```
#evaluating the model on training dataset
#write code here
r2_train = r2_score(Y_train, y_train_predicted)
```

Task 50: Evaluate R2 score for test data set

```
# evaluating the model on test dataset
#write code here
r2_test = r2_score(Y_test, y_test_predict)
```

Comaparing training and testing R2 scores

```
#write code here
print('The r2 score for training set is: ',r2_train)
print('The r2 score for testing set is: ',r2_test)
```

```
The r2 score for training set is: 0.819004180163414
The r2 score for testing set is: 0.8195478211895499
```

Task 51: Select the best model

Question: Which model gives the best result for price prediction? Find out the complexity using R2 score and give your answer

Hint: Use for loop for finding the best degree and model complexity for polynomial regression model

```
#write code here
r2_train=[]
r2_test=[]
for i in range(1,5):
    poly_reg = PolynomialFeatures(degree=i)

    X_tr_poly,X_tst_poly = poly_reg.fit_transform(X_train),poly_reg.fit_transform(X_test)
    poly_linearRegression()
    poly.fit(X_tr_poly, Y_train)

    y_tr_predicted,y_tst_predict = poly.predict(X_tr_poly),poly.predict(X_tst_poly)
    r2_train.append(r2_score(Y_train, y_tr_predicted))
    r2_test.append(r2_score(Y_test, y_tst_predict))

print('R2 Train', r2_train)
print('R2 Test', r2_test)
```

```
R2 Train [0.7238620495006879, 0.819004180163414, 0.8093183580432635, 0.716418397692478, 0.4956480399303697]
R2 Test [0.7290143663235678, 0.8195478211895499, 0.8109303933894887, 0.7189345889748313, 0.497889321205602315]
```

Plotting the model

```
plt.figure(figsize=(10,5))
sns.set_context('poster')
plt.subplot(1,2,1)
sns.lineplot(x=range(1,5), y=r2_train, label='Training');
plt.subplot(1,2,2)
sns.lineplot(x=range(1,5), y=r2_test, label='Testing');
```



Answer

Model with Polynomial degree 2 give the best optimum solution with accuracy of 81% and r2 score 0.81 as well as for training and testing data.