


```
1 import numpy as np          #NUMPY IS USED FOR NUMERICAL OPERATIONS
2 import pandas as pd # PANDAS FOR IMPORTING THE DATASET
3 from sklearn.model_selection import train_test_split    #TO SPLIT DATA INTO TRAINING AND TESTING DATA
4 from sklearn.tree import DecisionTreeClassifier    #FOR DECISION TREE CLASSIFIER
5 from sklearn.metrics import accuracy_score    #TO CHECK ACCURACY
6 import matplotlib.pyplot as plt    # FOR DATA VISUALIZATION PURPOSE
7 from sklearn import tree    # TO VISUALIZE THE TREE
8 import seaborn as sns    # FOR DATA VISUALIZATION
9
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.neighbors import KNeighborsClassifier
12 from sklearn.metrics import classification_report, confusion_matrix
```

LOADING THE DATASET:


```
1 df = pd.read_csv("/content/index.csv")    # using the Pandas library's read_csv function. The dataset is st
2 df
```



	date		datetime	cash_type		card	money	coffee_name
0	2024-03-01	2024-03-01	10:15:50.520	card	ANON-0000-0000-0001	38.70		Latte
1	2024-03-01	2024-03-01	12:19:22.539	card	ANON-0000-0000-0002	38.70		Hot Chocolate
2	2024-03-01	2024-03-01	12:20:18.089	card	ANON-0000-0000-0002	38.70		Hot Chocolate
3	2024-03-01	2024-03-01	13:46:33.006	card	ANON-0000-0000-0003	28.90		Americano
4	2024-03-01	2024-03-01	13:48:14.626	card	ANON-0000-0000-0004	38.70		Latte
...
1128	2024-07-31	2024-07-31	20:53:35.077	card	ANON-0000-0000-0443	23.02		Cortado
1129	2024-07-31	2024-07-31	20:59:25.013	card	ANON-0000-0000-0040	27.92		Americano with Milk
1130	2024-07-31	2024-07-31	21:26:26.000	card	ANON-0000-0000-0444	32.82		Latte
1131	2024-07-31	2024-07-31	21:54:11.824	card	ANON-0000-0000-0445	32.82		Latte
1132	2024-07-31	2024-07-31	21:55:16.570	card	ANON-0000-0000-0446	32.82		Latte


1133 rows × 6 columns

```
1 df.head(10)
```



	date		datetime	cash_type		card	money	coffee_name
0	2024-03-01	2024-03-01	10:15:50.520	card	ANON-0000-0000-0001	38.7		Latte
1	2024-03-01	2024-03-01	12:19:22.539	card	ANON-0000-0000-0002	38.7		Hot Chocolate
2	2024-03-01	2024-03-01	12:20:18.089	card	ANON-0000-0000-0002	38.7		Hot Chocolate
3	2024-03-01	2024-03-01	13:46:33.006	card	ANON-0000-0000-0003	28.9		Americano
4	2024-03-01	2024-03-01	13:48:14.626	card	ANON-0000-0000-0004	38.7		Latte
5	2024-03-01	2024-03-01	15:39:47.726	card	ANON-0000-0000-0005	33.8		Americano with Milk
6	2024-03-01	2024-03-01	16:19:02.756	card	ANON-0000-0000-0006	38.7		Hot Chocolate
7	2024-03-01	2024-03-01	18:39:03.580	card	ANON-0000-0000-0007	33.8		Americano with Milk
8	2024-03-01	2024-03-01	19:22:01.762	card	ANON-0000-0000-0008	38.7		Cocoa
9	2024-03-01	2024-03-01	19:23:15.887	card	ANON-0000-0000-0008	33.8		Americano with Milk

```
1 df.info()
```



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1133 entries, 0 to 1132
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  ---      -
0   date        1133 non-null   object
```

```

1  datetime    1133 non-null  object
2  cash_type   1133 non-null  object
3  card        1044 non-null  object
4  money       1133 non-null  float64
5  coffee_name 1133 non-null  object
dtypes: float64(1), object(5)
memory usage: 53.2+ KB

```

```
1 df.shape
```

```
(1133, 6)
```

```
1 df.describe()
```

```

money
count  1133.000000
mean    33.105808
std      5.035366
min     18.120000
25%     28.900000
50%     32.820000
75%     37.720000
max     40.000000

```

```
1 df.isnull().sum()
```

```

0
date      0
datetime  0
cash_type  0
card      89
money     0
coffee_name  0

```

```
1 df.columns
```

```
Index(['date', 'datetime', 'cash_type', 'card', 'money', 'coffee_name'], dtype='object')
```

```
1 df.dropna(inplace=True)
```


```
1 df.isnull().sum()
```

```

0
date      0
datetime  0
cash_type  0
card      0
money     0
coffee_name  0

```

```
1 df['date']
```




	date
0	2024-03-01
1	2024-03-01
2	2024-03-01
3	2024-03-01
4	2024-03-01
...	...
1128	2024-07-31
1129	2024-07-31
1130	2024-07-31
1131	2024-07-31
1132	2024-07-31

1044 rows × 1 columns



1 df.dtypes




	0
date	object
datetime	object
cash_type	object
card	object
money	float64
coffee_name	object



1 df['date'] = pd.to_datetime(df['date'])


1 df.dtypes



	0
date	datetime64[ns]
datetime	object
cash_type	object
card	object
money	float64
coffee_name	object



1 df.head()



	date	datetime	cash_type	card	money	coffee_name
0	2024-03-01	2024-03-01 10:15:50.520	card	ANON-0000-0000-0001	38.7	Latte
1	2024-03-01	2024-03-01 12:19:22.539	card	ANON-0000-0000-0002	38.7	Hot Chocolate
2	2024-03-01	2024-03-01 12:20:18.089	card	ANON-0000-0000-0002	38.7	Hot Chocolate
3	2024-03-01	2024-03-01 13:46:33.006	card	ANON-0000-0000-0003	28.9	Americano
4	2024-03-01	2024-03-01 13:48:14.626	card	ANON-0000-0000-0004	38.7	Latte



```
1 df.loc[:,['cash_type', 'card', 'coffee_name']].describe().T
```

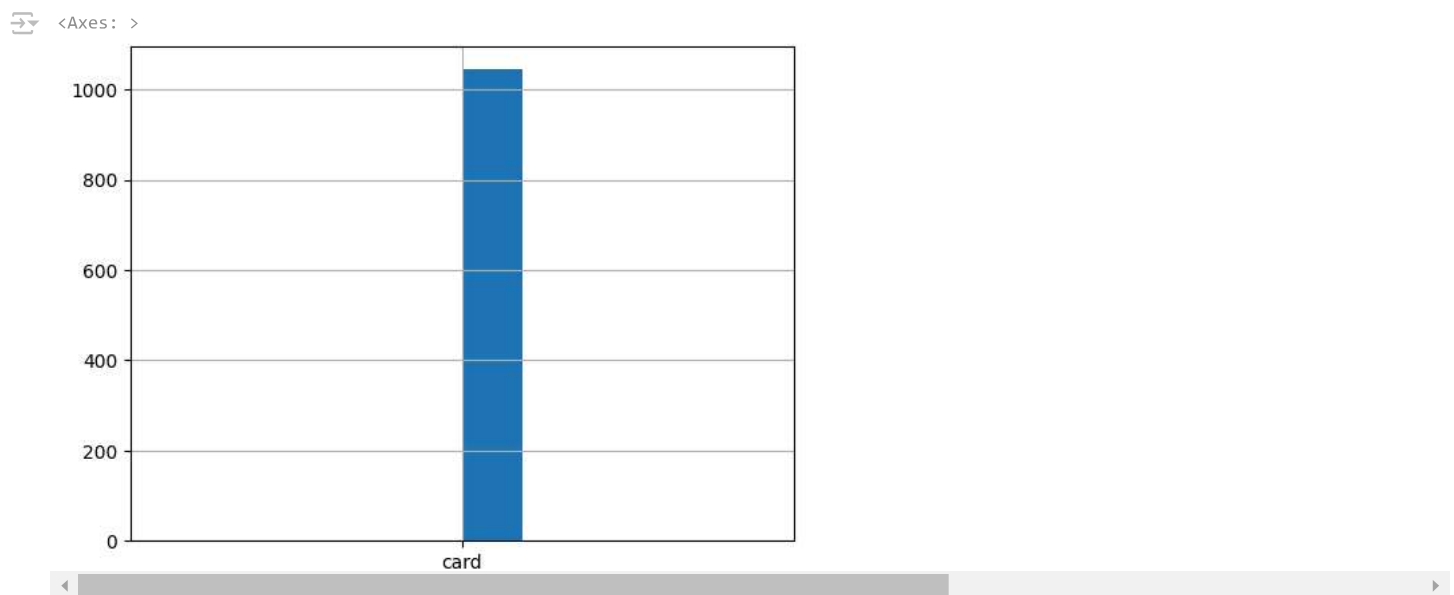
	count	unique	top	freq
cash_type	1044	1	card	1044
card	1044	446	ANON-0000-0000-0012	88
coffee_name	1044	8	Americano with Milk	253

```
1 Start coding or generate with AI.
```

- There are 1044 transactions in the data.
- 1 unique value of 'cash_type' and that is -> 'card'
- Americano with milk is the most popular product

```
1 Start coding or generate with AI.
```

```
1 df['cash_type'].hist()
```



```
1 Start coding or generate with AI.
```

Almost all transactions are made from card payment

```
1 Start coding or generate with AI.
```

```
1 df['coffee_name'].value_counts(normalize=True).sort_values(ascending=False).round(4)*100
```



proportion	
coffee_name	
Americano with Milk	24.23
Latte	20.88
Cappuccino	17.34
Americano	14.85
Cortado	9.00
Hot Chocolate	6.51
Espresso	4.21
Cocoa	2.97

- In this we can see the most popular products which are arranged in descending order
- Americano with Milk is at the top where as the second place belongs to Latte followed by other coffees

1 Start coding or generate with AI.

```
1 #Convert date and datetime to datetme format
2 df['date']=pd.to_datetime(df['date'])
3 df['datetime']=pd.to_datetime(df['datetime'])
4
5
6
```

```
1 #Create column of Month, Weekdays, and Hours
2 df['month']=df['date'].dt.strftime('%Y-%m')
3 df['day']=df['date'].dt.strftime('%w')
4 df['hour']=df['datetime'].dt.strftime('%H')
```

1 df.info()



```
<class 'pandas.core.frame.DataFrame'>
Index: 1044 entries, 0 to 1132
Data columns (total 9 columns):
#   Column      Non-Null Count  Dtype
---  -
0    date        1044 non-null   datetime64[ns]
1    datetime    1044 non-null   datetime64[ns]
2    cash_type    1044 non-null   object
3    card         1044 non-null   object
4    money        1044 non-null   float64
5    coffee_name  1044 non-null   object
6    month        1044 non-null   object
7    day          1044 non-null   object
8    hour         1044 non-null   object
dtypes: datetime64[ns](2), float64(1), object(6)
memory usage: 113.9+ KB
```

1 df.head()



	date	datetime	cash_type	card	money	coffee_name	month	day	hour
0	2024-03-01	2024-03-01 10:15:50.520	card	ANON-0000-0000-0001	38.7	Latte	2024-03	5	10
1	2024-03-01	2024-03-01 12:19:22.539	card	ANON-0000-0000-0002	38.7	Hot Chocolate	2024-03	5	12
2	2024-03-01	2024-03-01 12:20:18.089	card	ANON-0000-0000-0002	38.7	Hot Chocolate	2024-03	5	12
3	2024-03-01	2024-03-01 13:46:33.006	card	ANON-0000-0000-0003	28.9	Americano	2024-03	5	13
4	2024-03-01	2024-03-01 13:48:14.626	card	ANON-0000-0000-0004	38.7	Latte	2024-03	5	13

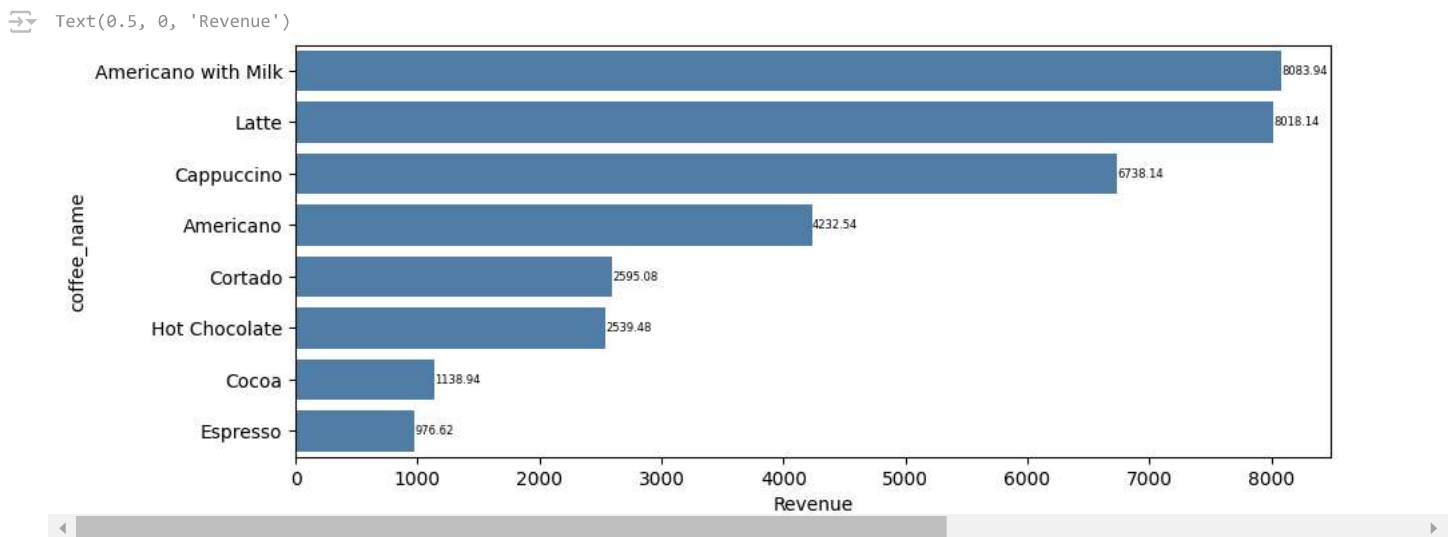
1 Start coding or generate with AI.

```
1 # Let's first check the overall revenue by products.
```

```
1 revenue = df.groupby(['coffee_name']).sum(['money']).sort_values(by='money', ascending=False)
2 revenue
```

	money
coffee_name	
Americano with Milk	8083.94
Latte	8018.14
Cappuccino	6738.14
Americano	4232.54
Cortado	2595.08
Hot Chocolate	2539.48
Cocoa	1138.94
Espresso	976.62

```
1 plt.figure(figsize=(10,4))
2 ax = sns.barplot(data=revenue, x='money', y='coffee_name', color='steelblue')
3 ax.bar_label(ax.containers[0], fontsize=6)
4 plt.xlabel("Revenue")
```



The product with the highest revenue is Americano with Milk whereas Espresso is at the bottom

```
1 Start coding or generate with AI.
```

```
1 monthly_sales = df.groupby(['coffee_name', 'month']).count()['date'].reset_index().rename(columns={'date': 'count'})
2 monthly_sales
```

☰

	coffee_name	month	Americano	Americano with Milk	Cappuccino	Cocoa	Cortado	Espresso	Hot Chocolate	Latte
0		2024-03	32	30	15	6	28	9	20	35
1		2024-04	33	38	36	4	16	4	10	27
2		2024-05	40	54	52	8	17	7	13	50
3		2024-06	14	66	46	4	19	10	14	50
4		2024-07	36	65	32	9	14	14	11	56

☰

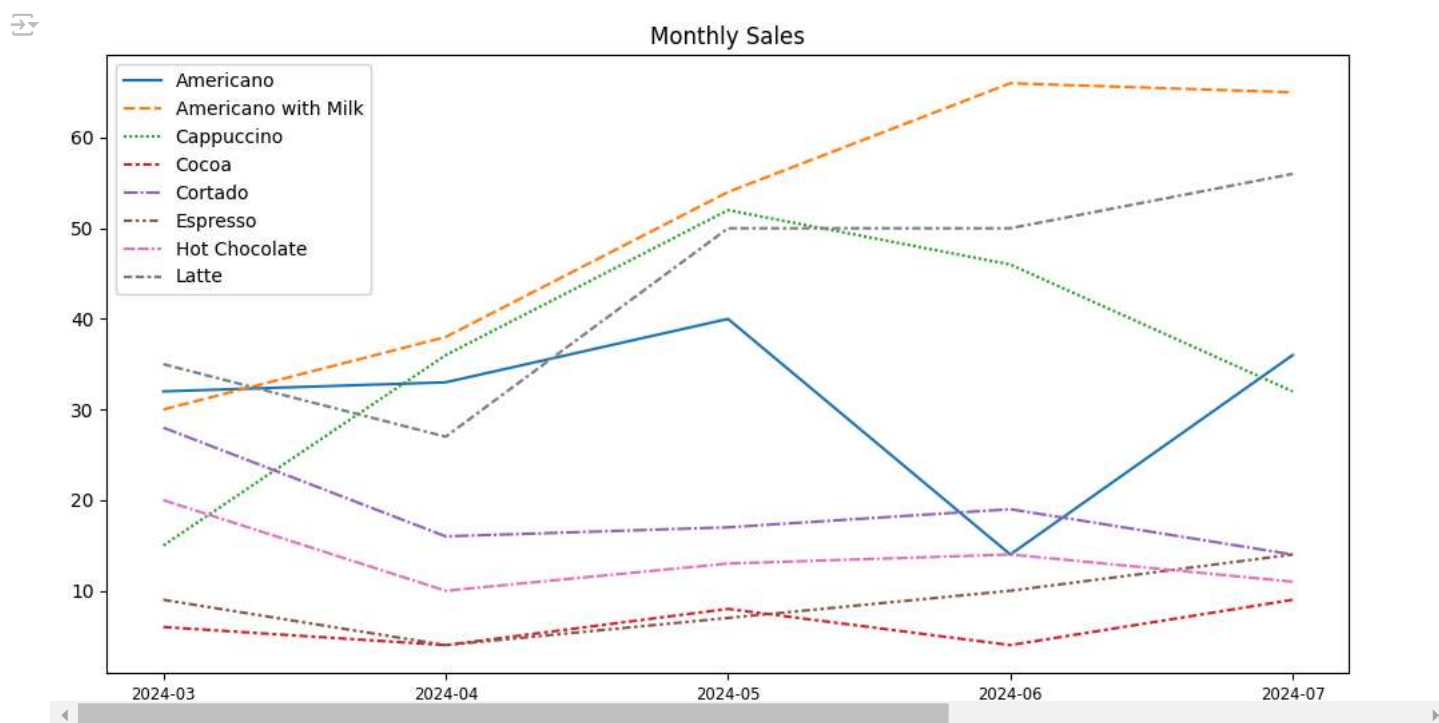
```
1 monthly_sales.describe().T.loc[:,['min', 'max']]
```

	min	max
coffee_name		
Americano	14.0	40.0
Americano with Milk	30.0	66.0
Cappuccino	15.0	52.0
Cocoa	4.0	9.0
Cortado	14.0	28.0
Espresso	4.0	14.0
Hot Chocolate	10.0	20.0
Latte	27.0	56.0

```

1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=monthly_sales)
3 plt.legend(loc='upper left')
4 plt.title("Monthly Sales")
5 plt.xticks(range(len(monthly_sales['month'])),monthly_sales['month'],size='small')
6 plt.show()

```



As shown in the line chart above, Americano with Milk and Latte, and Cappuccino are top selling coffee types, while Cocoa and Espresso have lowest sales. Additionally, Americano with Milk and Latte show an upward trending.

1 Start coding or generate with AI.

```

1 weekday_sales = df.groupby(['day']).count()['date'].reset_index().rename(columns={'date': 'count'})
2 weekday_sales

```

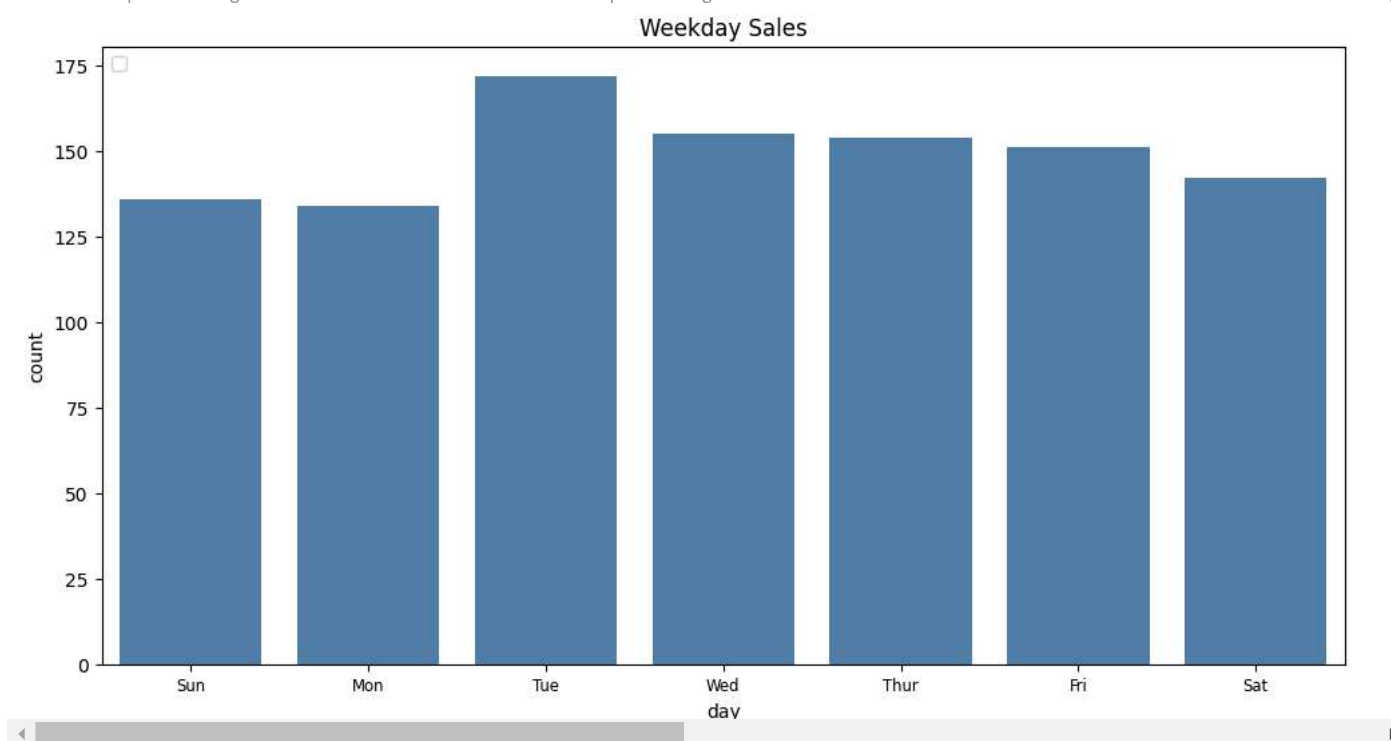
	day	count
0	0	136
1	1	134
2	2	172
3	3	155
4	4	154
5	5	151
6	6	142

```

1 plt.figure(figsize=(12,6))
2 sns.barplot(data=weekday_sales, x='day' , y='count', color='steelblue')
3 plt.legend(loc='upper left')
4 plt.title("Weekday Sales")
5 plt.xticks(range(len(weekday_sales['day'])),['Sun','Mon','Tue','Wed','Thur','Fri','Sat'],size='small')
6 plt.show()

```

WARNING:matplotlib.legend.No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored



The bar chart reveals that Tuesday has the highest sales of the week, while sales on the other days are relatively similar.

1 Start coding or generate with AI.

```

1 daily_sales = df.groupby(['coffee_name','date']).count()['datetime'].reset_index().reset_index().rename(cc
2 daily_sales

```




	coffee_name	date	Americano	Americano with Milk	Cappuccino	Cocoa	Cortado	Espresso	Hot Chocolate	Latte
0		2024-03-01	1.0	4.0	0.0	1.0	0.0	0.0	3.0	2.0
1		2024-03-02	3.0	3.0	0.0	0.0	0.0	0.0	0.0	0.0
2		2024-03-03	1.0	2.0	0.0	1.0	2.0	0.0	2.0	1.0
3		2024-03-04	0.0	1.0	0.0	0.0	0.0	1.0	0.0	2.0
4		2024-03-05	0.0	0.0	0.0	1.0	1.0	0.0	4.0	3.0
...	
145		2024-07-27	0.0	5.0	4.0	0.0	0.0	2.0	0.0	2.0
146		2024-07-28	0.0	1.0	0.0	0.0	0.0	1.0	0.0	1.0
147		2024-07-29	3.0	2.0	2.0	1.0	0.0	0.0	2.0	1.0
148		2024-07-30	2.0	12.0	2.0	0.0	3.0	2.0	0.0	3.0
149		2024-07-31	2.0	6.0	1.0	2.0	4.0	0.0	0.0	7.0

150 rows × 9 columns

```
1 daily_sales.iloc[:,1:].describe().T.loc[:,['min','max']]
```



	min	max
coffee_name		
Americano	0.0	4.0
Americano with Milk	0.0	12.0
Cappuccino	0.0	9.0
Cocoa	0.0	2.0
Cortado	0.0	4.0
Espresso	0.0	3.0
Hot Chocolate	0.0	4.0
Latte	0.0	7.0

```
1 Start coding or generate with AI.
```

```
1 hourly_sales = df.groupby(['hour']).count()['date'].reset_index().rename(columns={'date':'count'})
2 hourly_sales
```

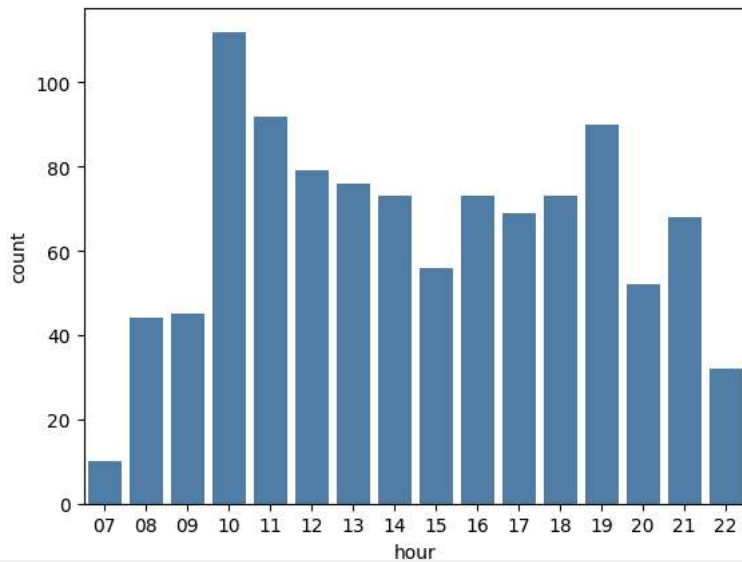


	hour	count
0	07	10
1	08	44
2	09	45
3	10	112
4	11	92
5	12	79
6	13	76
7	14	73
8	15	56
9	16	73
10	17	69
11	18	73
12	19	90
13	20	52
14	21	68
15	22	32

```
1 sns.barplot(data=hourly_sales,x='hour',y='count',color='steelblue')  
2
```




<Axes: xlabel='hour', ylabel='count'>




Overall,two peak hours within each day can be observed: 10:00am and 7:00pm.Then, let's check if any difference for different products.

1 Start coding or generate with AI.

```
1 hourly_sales_by_coffee = df.groupby(['hour','coffee_name']).count()['date'].reset_index().rename(columns={  
2 hourly_sales_by_coffee
```



coffee_name	hour	Americano	Americano with Milk	Cappuccino	Cocoa	Cortado	Espresso	Hot Chocolate	Latte
0	07	3.0	3.0	1.0	0.0	1.0	0.0	0.0	2.0
1	08	10.0	7.0	8.0	1.0	6.0	0.0	0.0	12.0
2	09	7.0	15.0	6.0	0.0	5.0	2.0	0.0	10.0
3	10	19.0	30.0	8.0	4.0	7.0	1.0	7.0	36.0
4	11	19.0	23.0	14.0	1.0	12.0	6.0	7.0	10.0
5	12	13.0	24.0	15.0	3.0	5.0	4.0	3.0	12.0
6	13	18.0	18.0	8.0	2.0	12.0	3.0	4.0	11.0
7	14	13.0	17.0	13.0	4.0	6.0	5.0	2.0	13.0
8	15	11.0	14.0	7.0	0.0	2.0	4.0	4.0	14.0
9	16	10.0	17.0	10.0	2.0	12.0	5.0	4.0	13.0
10	17	8.0	10.0	15.0	3.0	6.0	4.0	6.0	17.0
11	18	9.0	15.0	12.0	2.0	5.0	5.0	10.0	15.0
12	19	4.0	18.0	33.0	2.0	5.0	0.0	8.0	20.0
13	20	1.0	12.0	12.0	5.0	5.0	3.0	6.0	8.0
14	21	5.0	23.0	13.0	1.0	3.0	1.0	3.0	19.0
15	22	5.0	7.0	6.0	1.0	2.0	1.0	4.0	6.0



1 Start coding or [generate](#) with AI.

```

1 fig, axs = plt.subplots(2, 4, figsize=(20, 10))
2 # Flatten the array of subplots for easy iteration
3 axs = axs.flatten()
4 # Loop through each column in the DataFrame, skipping the 'Index' column
5 for i, column in enumerate(hourly_sales_by_coffee.columns[1:]):
6     axs[i].bar(hourly_sales_by_coffee['hour'], hourly_sales_by_coffee[column])
7     axs[i].set_title(f'{column}')
8     axs[i].set_xlabel('Hour')
9
10 #axs[i].set_ylabel('Sales')
11 plt.tight_layout()
12 # Show the plot
13 plt.show()

```

