

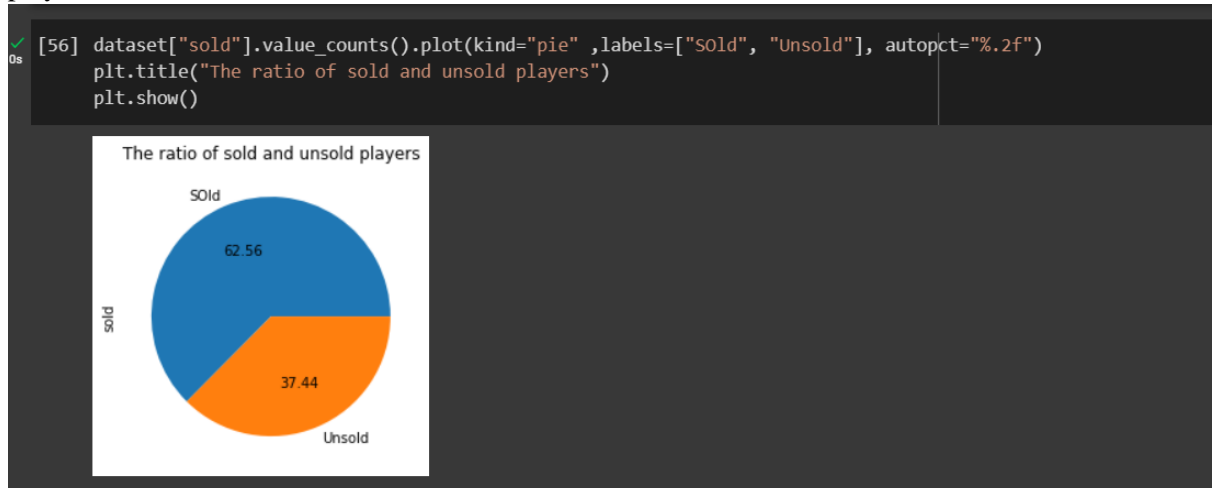
Decoding ml task 1

Name: Pasumarthy Sri Vishnu Vardhan

EDA

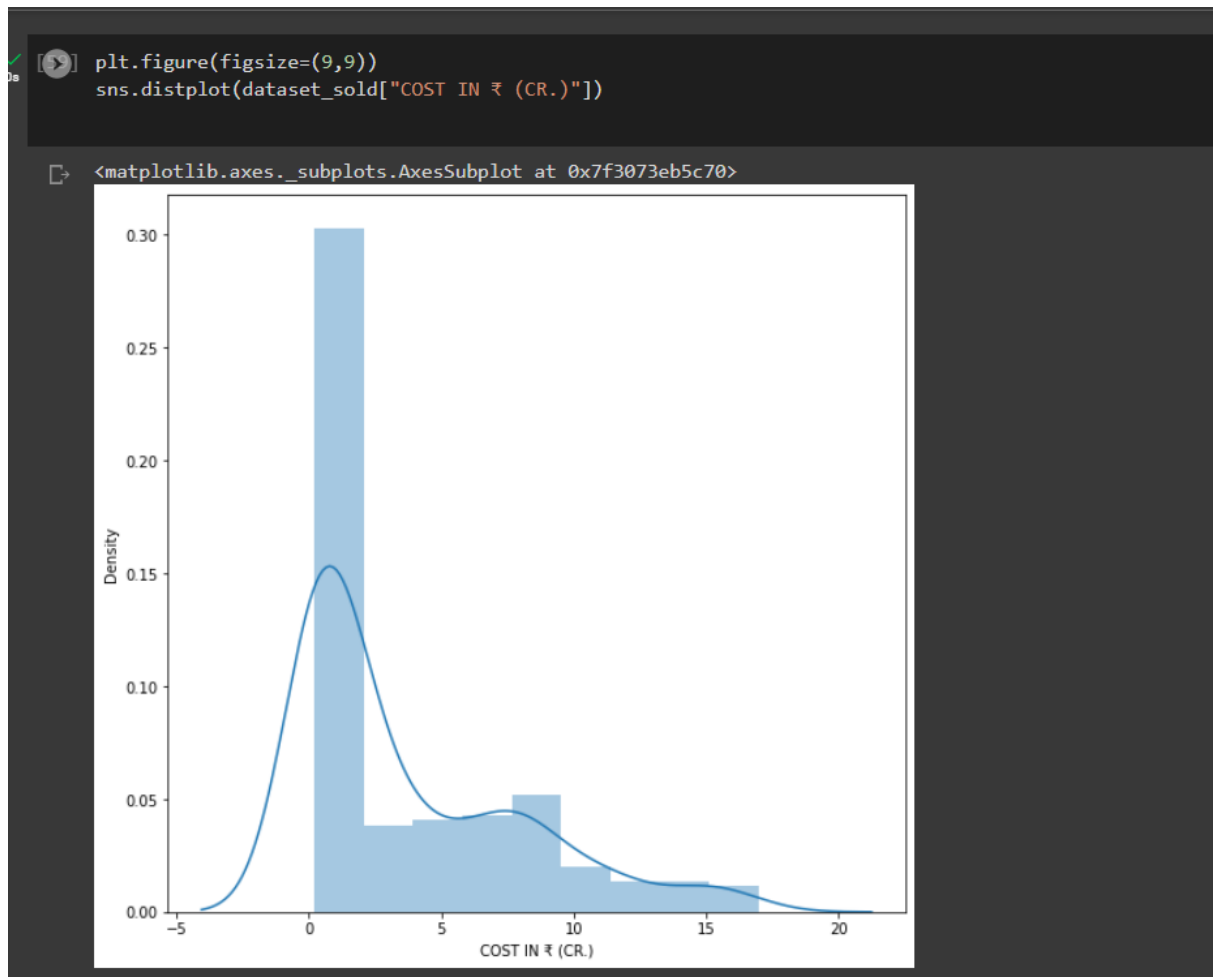
Dataset: IPL auction dataset

- 1) Importing necessary packages numpy, pandas, seaborn, matplotlib.pyplot, warnings
- 2) Reading the dataset using pandas
- 3) Understanding the dataset
 - dataset.head() : it gives you first 5 rows in the dataset
 - dataset.tail() : it gives you last 5 rows from the dataset
 - dataset.columns: it returns you column names
 - dataset.info(): it returns you null count and datatype of each column
- 4) By observing there are null values In Team column and other columns
I observed there are null values In Team column and other columns
- 4) By observing null values in dataset I split the dataset into sold players dataset and unsold players dataset



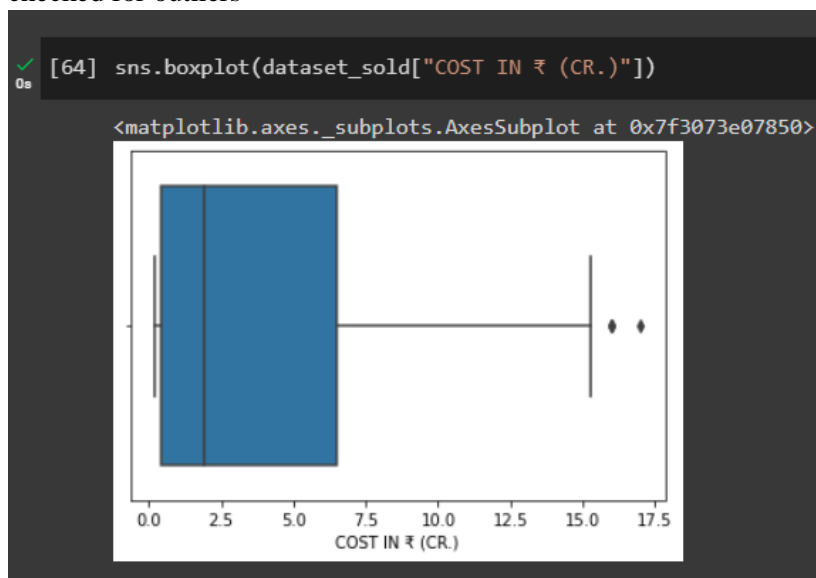
The above pie chart describes the percentage of sold and unsold players

- 5) Now I started exploratory data analysis on sold players
 - a. Let the dataset name for sold players be dataset_sold
 - b. Then I observed the price of each player starting from 0.2cr to 15 Crs
 - c. So I decided to plot a dist plot



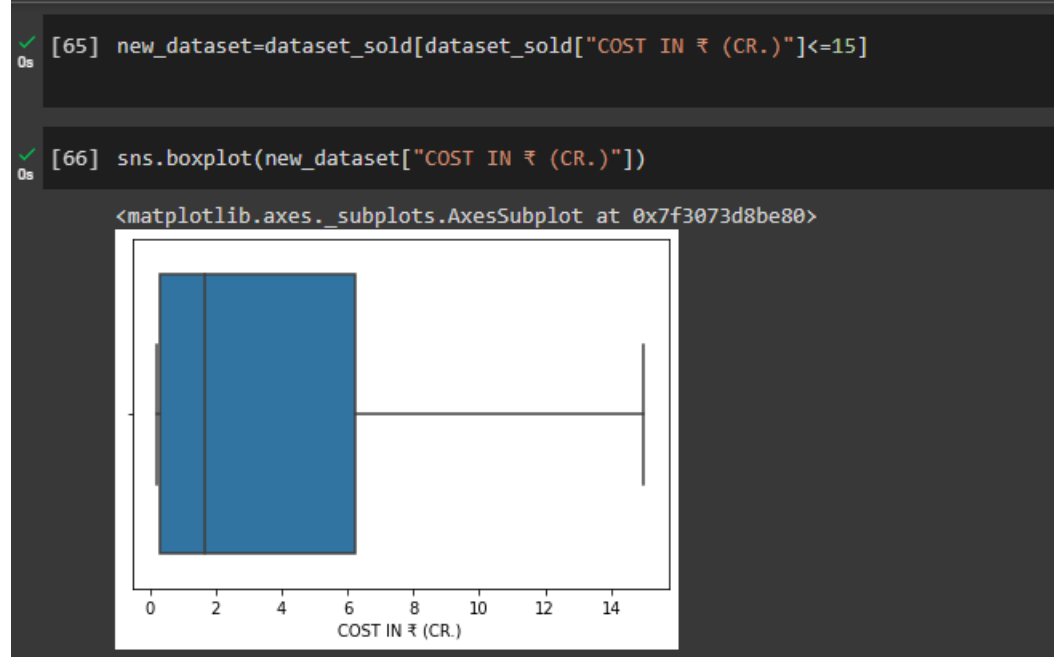
So the most of the players price at 0.2-3 cr

- d. Then I tried to calculate mean or average of each price of player, before going to checked for outliers

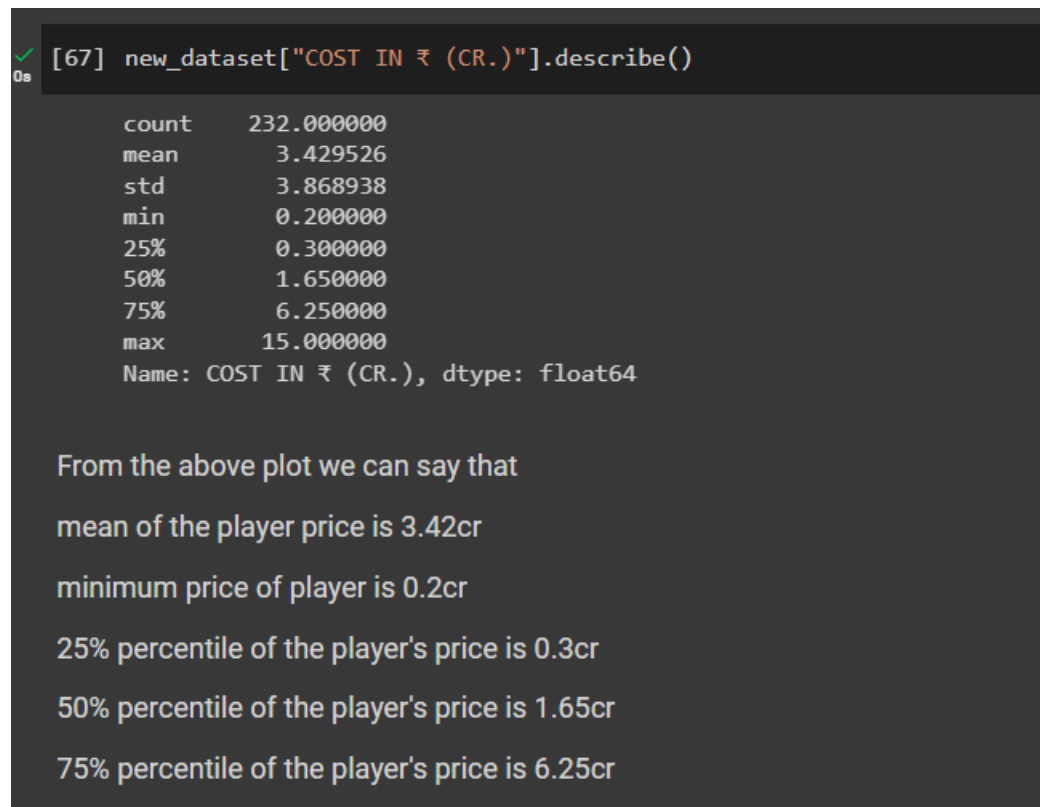


so I found an outliers and I removed them and calculated mean

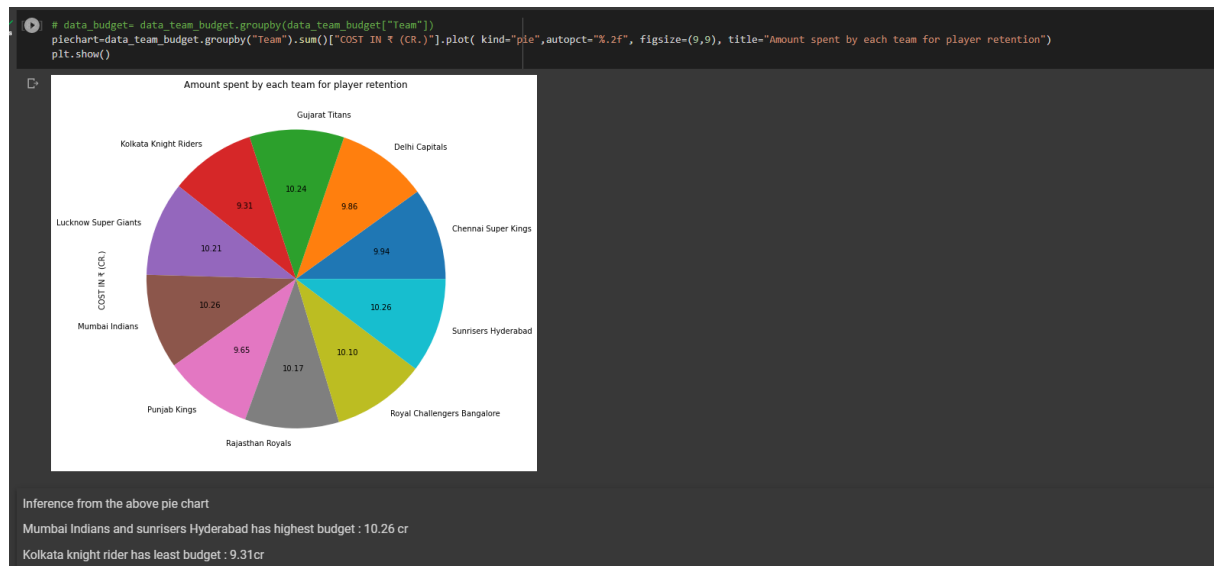
So the average price of each player is 3.42 crs



- e. The above chart is out come After removing outlier
- f. Dataset[‘COST IN \$(CR.)’].describe() : it returns a five tuple summary of the dataset or row



- g. Amount spent by each team for player retention



to get this above pie chart : First I grouped by teams column

- h. Next, I tried to figure out number of players in each type in a team

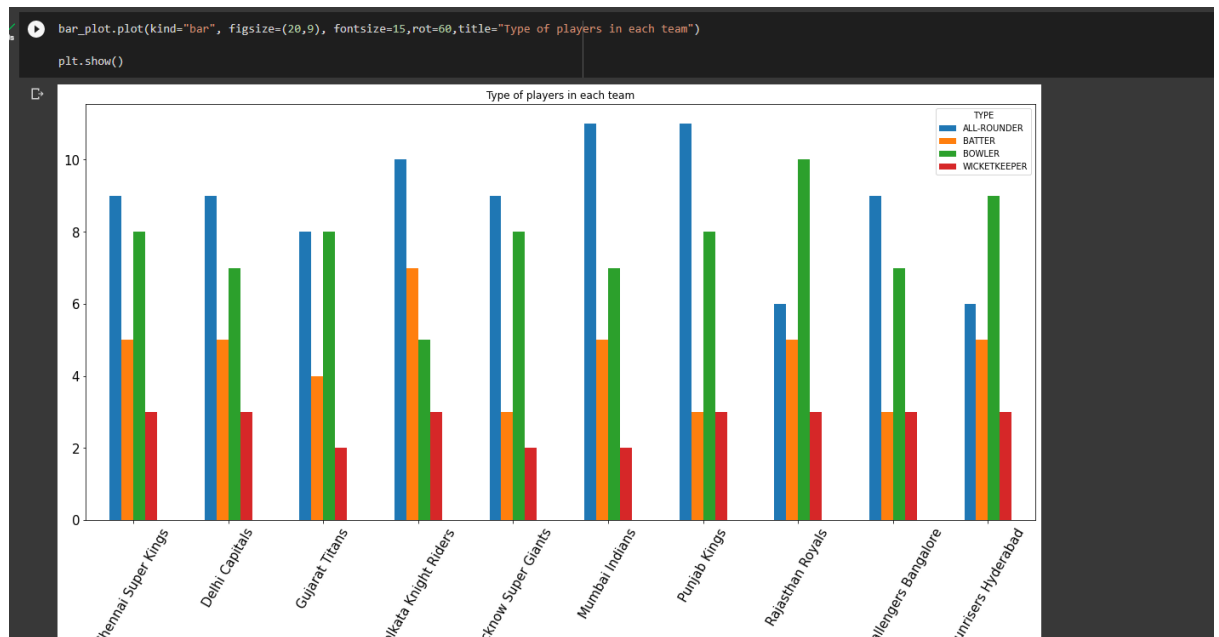
Type of players

```
[71] bar_plot=pd.crosstab(dataset_sold["Team"],dataset_sold["TYPE"])
```

```
[72] bar_plot.head()
```

| | TYPE | ALL-ROUNDER | BATTER | BOWLER | WICKETKEEPER |
|-----------------------|------|-------------|--------|--------|--------------|
| Team | | | | | |
| Chennai Super Kings | | 9 | 5 | 8 | 3 |
| Delhi Capitals | | 9 | 5 | 7 | 3 |
| Gujarat Titans | | 8 | 4 | 8 | 2 |
| Kolkata Knight Riders | | 10 | 7 | 5 | 3 |
| Lucknow Super Giants | | 9 | 3 | 8 | 2 |

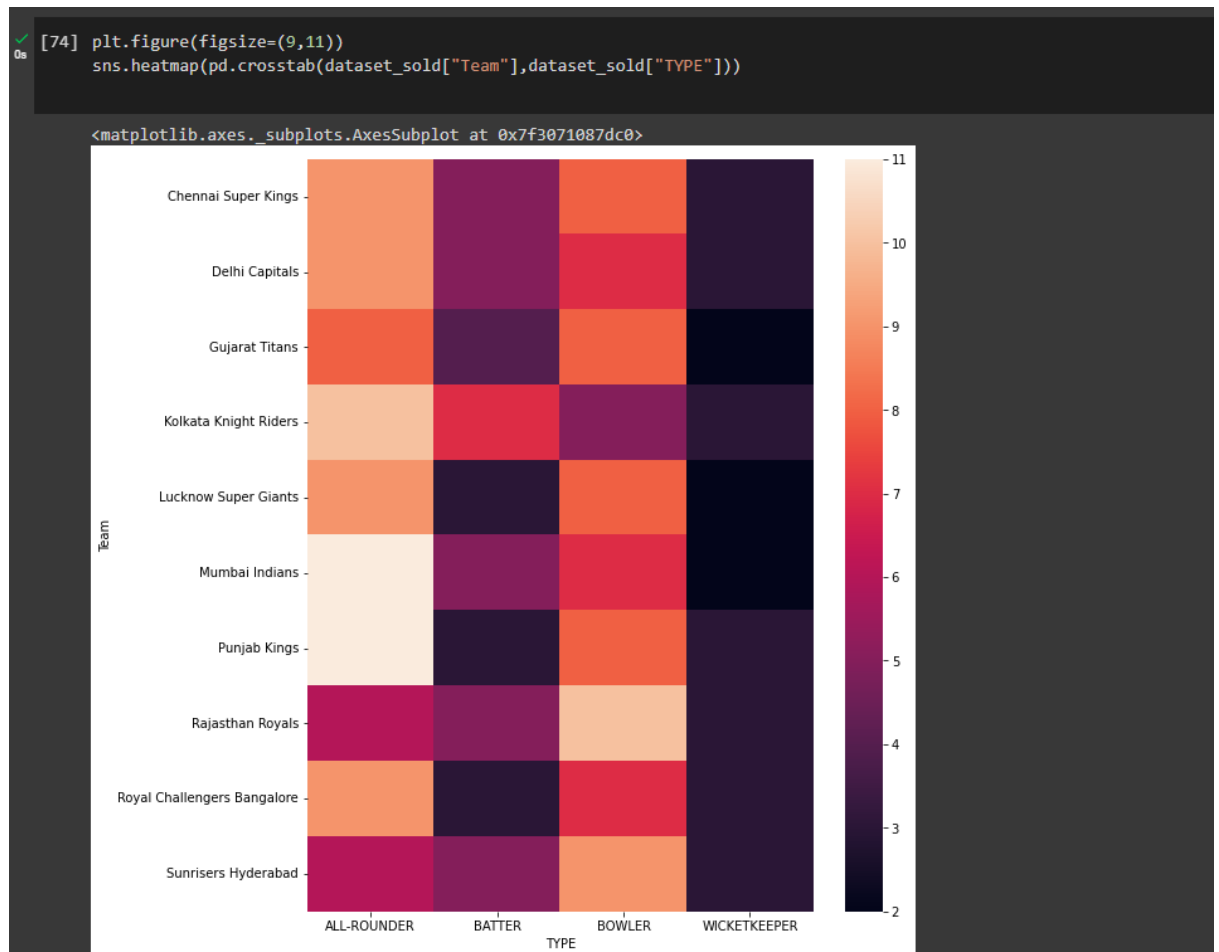
`crosstab("column1","column2")` : A crosstab is a table showing the relationship between two or more variables. Where the table only shows the relationship between two categorical variables, a crosstab is also known as a contingency table.



From the above chart we observe that

- 1) Mumbai Indians and Punjab Kings have more all rounder players
- 2) Rajasthan Royals has more number of bowlers
- 3) Kolkata Knight Riders has more number of batter

2) using heatmap



6) EDA on unsold data

```
[75] dataset_unsold=dataset[dataset["Team"]=="Unsold"]
```

```
[76] dataset_unsold.head()
```

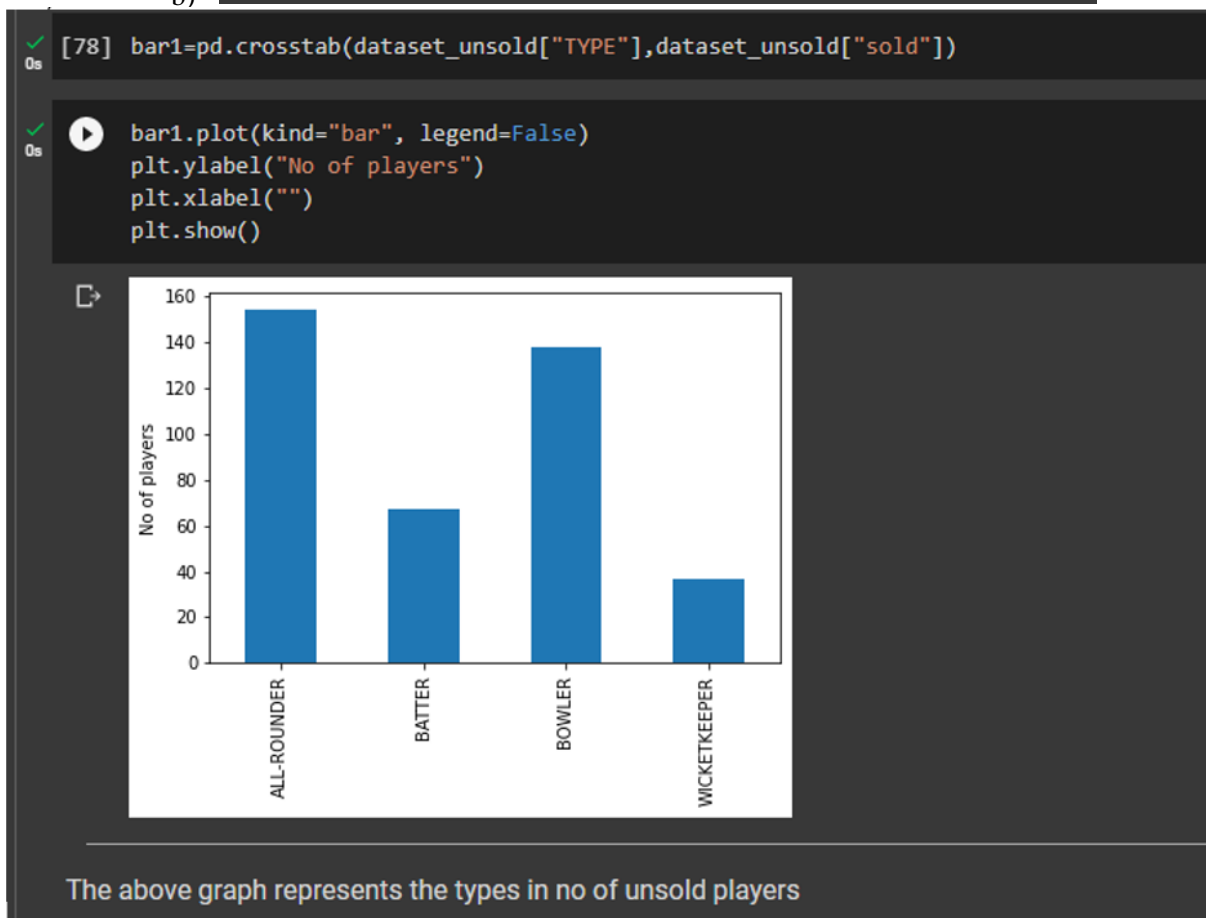
| | Player | Base Price | TYPE | COST IN ₹ (CR.) | Cost IN \$ (000) | 2021 Squad | Team | sold |
|-----|-----------------|------------|-------------|-----------------|------------------|------------|--------|------|
| 237 | Suresh Raina | 2 Cr | BATTER | NaN | NaN | CSK | Unsold | 0 |
| 238 | Steve Smith | 2 Cr | BATTER | NaN | NaN | DC | Unsold | 0 |
| 239 | Shakib Al Hasan | 2 Cr | ALL-ROUNDER | NaN | NaN | KKR | Unsold | 0 |
| 240 | Amit Mishra | 1.5 Cr | BOWLER | NaN | NaN | DC | Unsold | 0 |
| 241 | Adil Rashid | 2 Cr | BOWLER | NaN | NaN | PBKS | Unsold | 0 |

a)

```
[77] dataset_unsold.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 396 entries, 237 to 632
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  ---                ---
0   Player                 396 non-null   object
1   Base Price             396 non-null   object
2   TYPE                   396 non-null   object
3   COST IN ₹ (CR.)        0 non-null     float64
4   Cost IN $ (000)        0 non-null     float64
5   2021 Squad             40 non-null    object
6   Team                   396 non-null   object
7   sold                   396 non-null   int64
dtypes: float64(2), int64(1), object(5)
memory usage: 27.8+ KB
```

b)



c)

From the above bar chart observations:

1. No of All-rounder unsold :150
2. No of Batters unsold: 60
3. No of bowlers unsold: 130
4. No of wicketkeeper unsold: 40

The End