SportsStats

Capstone Project

LI RUI

Table Of Contents

Preparing Proposal

Developing Project Proposal

Preparing Proposal

Which client did you select and why?

• I choose the SportsStats client. I selected this project since I have strong knowledge in sports and have interests in working in sports analytics potentially.

 With an analysis of the dataset, I can find patterns and hidden insights for players, hidden insights.

Steps used to import and Clean data

• For importing the data, I used pandas to read in both csv files

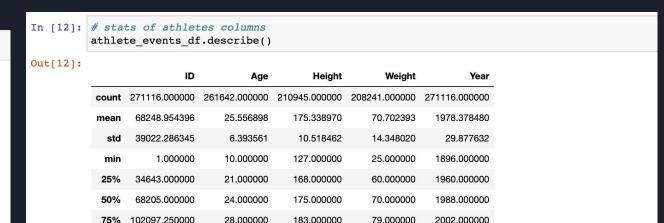
 Since there were many NAN values, I didn't do much cleaning since that would falsify the data!

Initial Exploration of Data and display screenshots of work

```
[11]: # info of athlete df columns
      athlete events df.info()
      <class 'pandas.core.frame.DataFrame'>
      RangeIndex: 271116 entries, 0 to 271115
      Data columns (total 15 columns):
           Column Non-Null Count
           TD
                  271116 non-null int64
                  271116 non-null object
           Sex
                  271116 non-null object
                  261642 non-null float64
           Age
          Height 210945 non-null float64
           Weight 208241 non-null float64
           Team
                  271116 non-null object
           NOC
                  271116 non-null object
                  271116 non-null object
                  271116 non-null int64
           Season 271116 non-null object
           City
                  271116 non-null
                                   object
          Sport
                  271116 non-null
                                   object
          Event
                  271116 non-null
                                   object
          Medal
                  39783 non-null
                                   object
      dtypes: float64(3), int64(2), object(10)
      memory usage: 31.0+ MB
```

Observation

- Fairly evenly distributed data types amongst the columns.
- · Object type is the most ubiquitous
- · Can see some missing values in data right away



226.000000

Observation

135571.000000

97.000000

Nothing too crazy about stats of athletes df columns, although I am curious of the 97 years old who participated

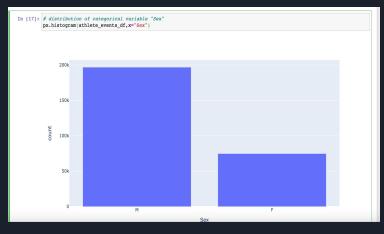
214.000000

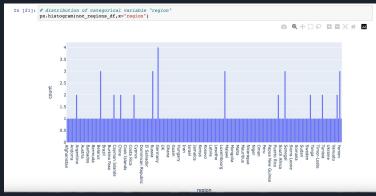
2016.000000

Initial Exploration of Data and display screenshots

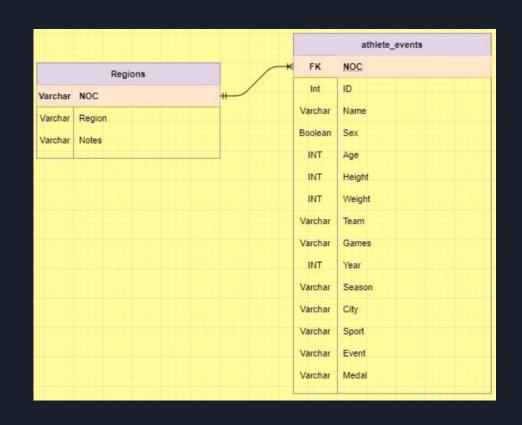
of work

```
In [14]: # making sure only two sexes
          athlete events df["Sex"].unique()
Out[14]: array(['M', 'F'], dtype=object)
In [15]: # seeing distribution of sexes
          athlete events df["Sex"].value counts(normalize=True)
Out[15]: M
               0.725129
               0.274871
          Name: Sex, dtype: float64
          Observation
           • There are more male athletes, with approximately 75% being male and 25% being female
         # distribution of categorical variable "Sex"
          px.histogram(athlete events df, x="Sex")
```





ERD



Developing Project Proposal

Possible questions

- When was the first season ever conducted
- Disparity between sexes
- Which country performed the best
- Age distribution of participants
- Which country had highest number of Players

Description

- My Project targets getting past Sports patterns and analysing it.
- Get to know more insights on the data, such as when the first event was organized and in which city/country.
- This analysis will not only help Sports Coaches to identify patterns and records, but it will also help the SportsStats firm aid in their clients' decision-making.
- My audience for the projects would not be limited to Coaches/Trainers but also players who will be able to see their records/performance in past events.

Possible Hypothesis

- Women have higher number of Medals.
- Year > 1956 will have the Highest number of Events.
- More people have participated in Football.
- 3 Teams will have medals > 40.
- People with Age > 40 have received medal in any of the events.