# USED CARS DATASET..

By : Soaad Aljafr.

Reem Alghamdi.

# • Abstract:

- The goal of the project is to predict the prices of cars that will be displayed in the future, The prediction will be made from the data already in the database https://www.kaggle.com/austinreese/craigslist-carstrucks-data By training the machine with the previous data.

# Interdiction:

- Craigslist is an American classified advertisements website with sections devoted to jobs, housing, for sale, items wanted, services, community service, gigs, résumés, and discussion forums.

- Craigslist is the world's largest collection of used vehicles for sale, dataset which includes every used vehicle entry within the United States on Craigslist. it contains most all relevant information that provides on car sales including columns like price, condition, manufacturer, latitude/longitude, and 18 other categories.

# Design:

- Since 2008, guests and hosts have used Airbnb to expand on traveling possibilities and present more unique, personalized way of experiencing the world. This dataset describes the listing activity and metrics in NYC, NY for 2019.

# Data:

- The dataset includes 26 columns and more than 400,000 rows, The most important of them are price, color, year , odometer and condition.

- It differs between numeric and category, This information helped an adequate understanding of the database.

- Almost is category so we should convert to numeric.

# Num. null

| | 0 |
|---|---|
| region | 0 |
| price | 0 |
| year | 1205 |
| manufacturer | 17646 |
| model | 5277 |
| condition | 174104 |
| cylinders | 177678 |
| fuel | 3013 |
| odometer | 4400 |
| title_status | 8242 |
| transmission | 2556 |
| drive | 130567 |
| size | 306361 |
| type | 92858 |
| paint_color | 130203 |
| county | 426880 |
| state | 0 |
| posting_date | 68 |

# Object describe

| | COUNT | UNIQUE | TOP | FREQ |
|---|---|---|---|---|
| region | 426880 | 404 | columbus | 3608 |
| manufacturer | 409234 | 42 | ford | 70985 |
| model | 421603 | 29667 | f-150 | 8009 |
| condition | 252776 | 6 | good | 121456 |
| cylinders | 249202 | 8 | 6 cylinders | 94169 |
| fuel | 423867 | 5 | gas | 356209 |
| title_status | 418638 | 6 | clean | 405117 |
| transmission | 424324 | 3 | automatic | 336524 |
| drive | 296313 | 3 | 4wd | 131904 |
| type | 334022 | 13 | sedan | 87056 |
| paint_color | 296677 | 12 | white | 79285 |
| state | 426880 | 51 | ca | 50614 |
| posting_date | 426812 | 381536 | 2021-04-23T22:13:05-0400 | 12 |

# Algorithms:

- Cleaned the data for many steps

- First remove the null from year and fill null for mean in numeric columns and mode for categories.

- Then remove duplicate.
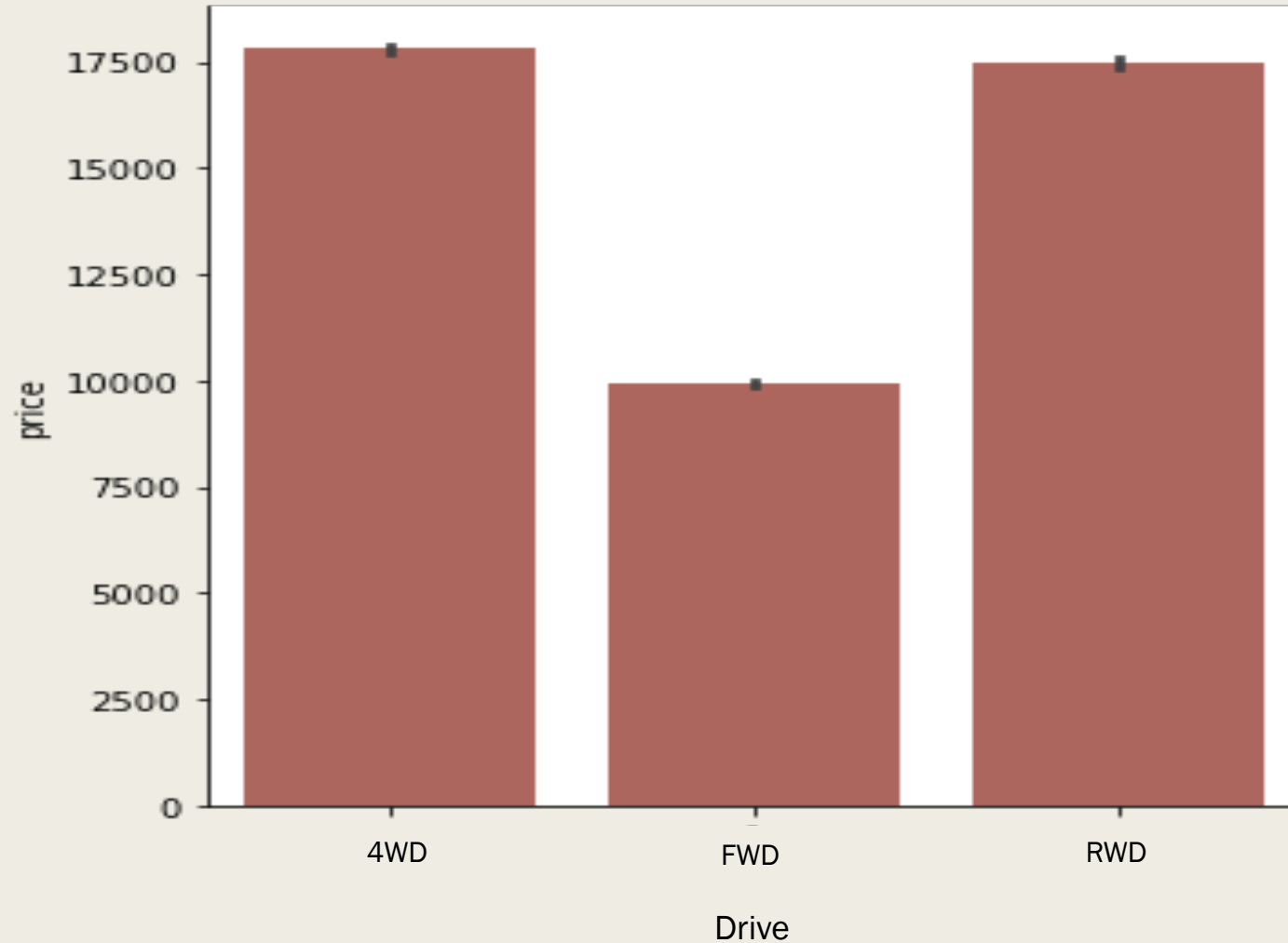
- Cheek if there any null or duplicate.

# Tools:

- The dataset includes 26 columns and more than 400,000 rows, The most important of them are price, color, year , odometer and condition.

- It differs between numeric and category, This information helped an adequate understanding of the database.

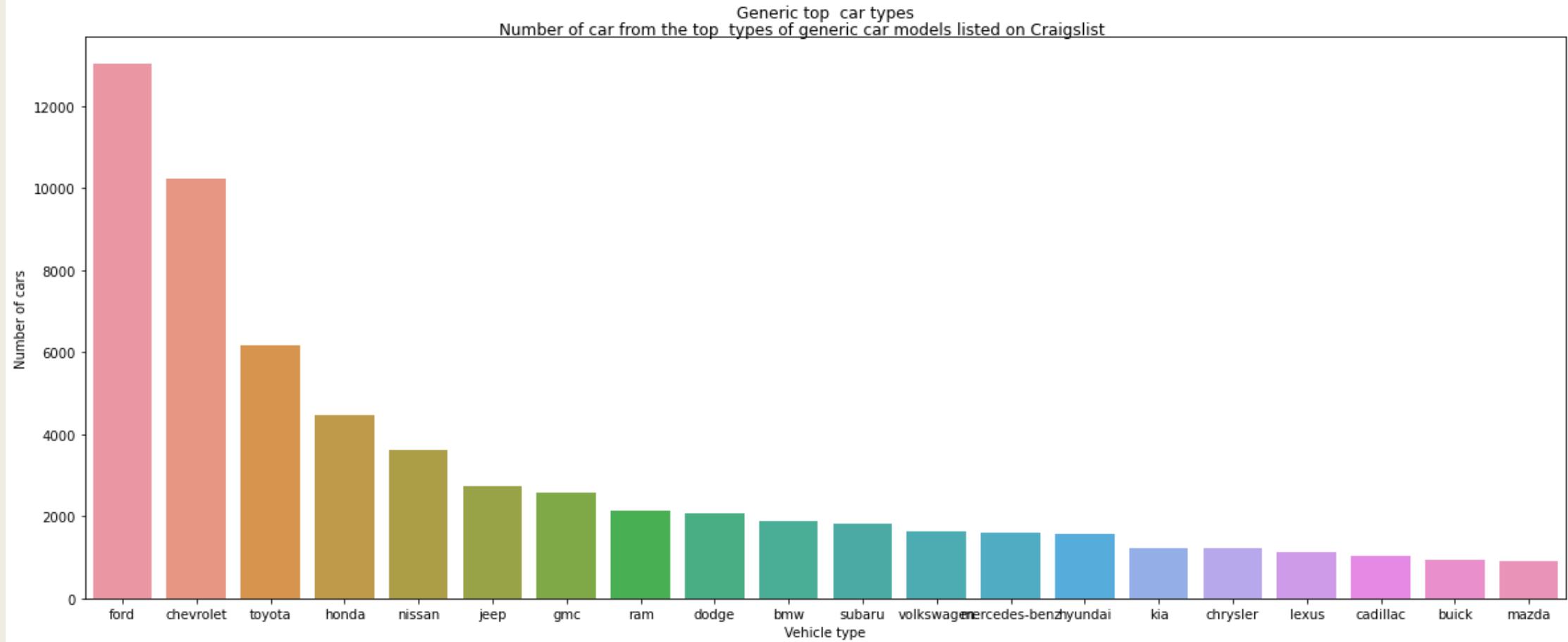- Almost is category so we should convert to numeric.

# Data Cleaning

- Read The Data from cvs

- Drop some columes

- Delet the Null values

- Fill Null values

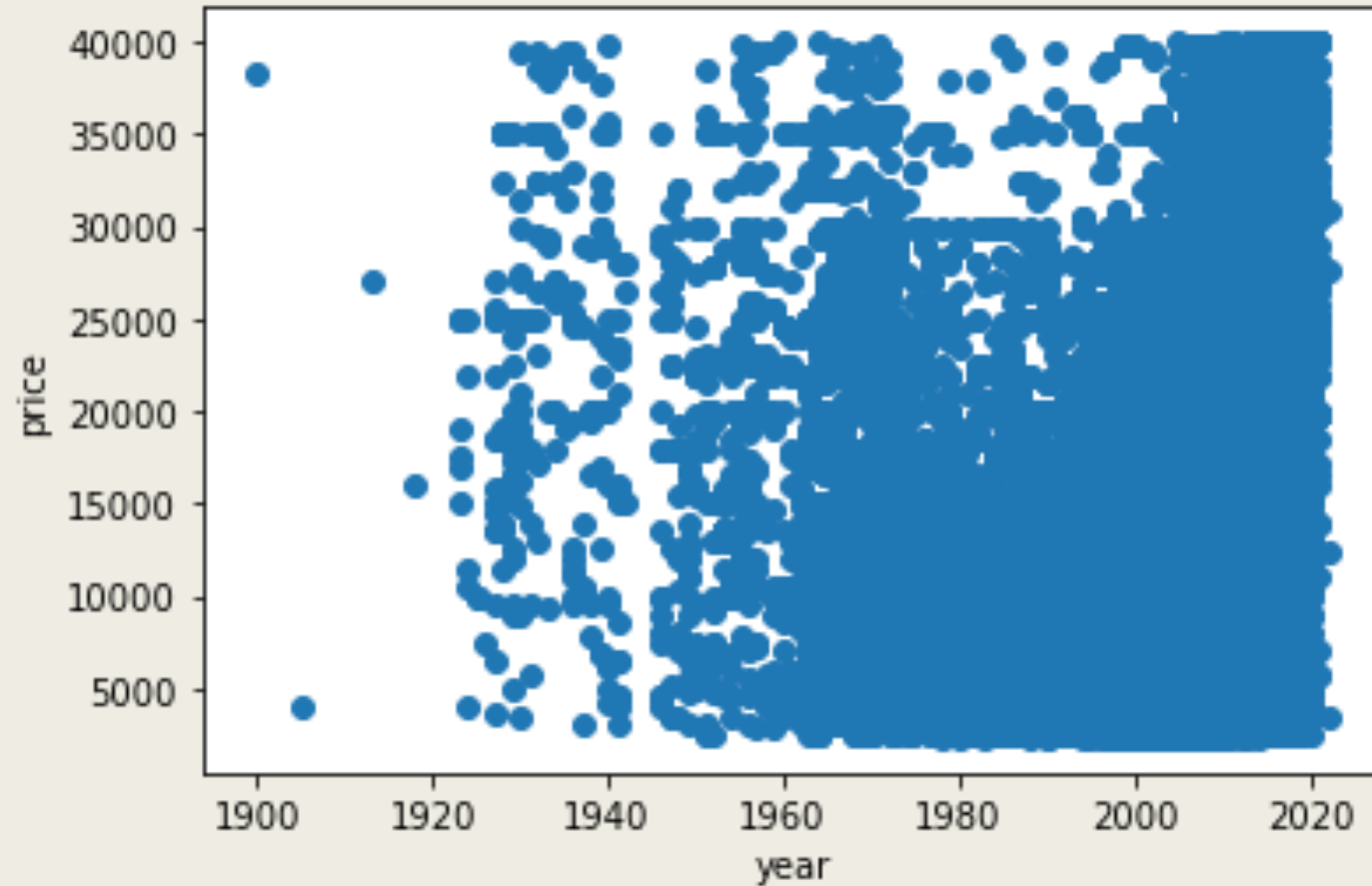- Change the object columes to numaric

# Bar plot between type and price

# Manufacturer



Generic top  car types
Number of car from the top  types of generic car models listed on Craigslist

# Scatler Years and Price

# Data Correlation

# Model and Result

| | Training | Validation | Test |
|---|---|---|---|
| Liner Regression | 0.377 | 0.358 | 0.339 |
| Polynomial | 0.801 | 0.875 | 0.845 |

# Thanks..