

# Concert Ticket Price Prediction

Sarah Alabdulwahab & Asma Althakafi



# RAZORGATOR TICKETS

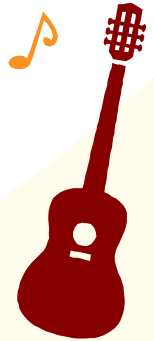
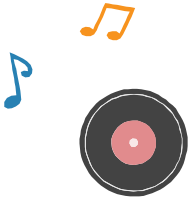
Online ticket **reselling** platform for sports, theater and concert tickets, and vacation packages for sporting events.





# GOAL

Predict the price of concert tickets in USA



# WEB SCRAPING PROCESS

## Step 1

Beautiful Soup &  
Selenium

## Step 3

Collecting the concert  
links for each artist

## Step 2

Collecting artists

## Step 4

Collecting the tickets  
of each concert



# WEB SCRAPED DATASET

## Artist

Individual artist or band

## Data

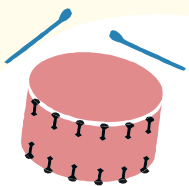
Venue, city, state, date,  
and time

## Level

Section and row

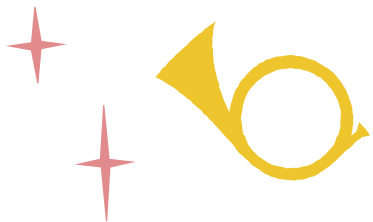
## Price

Ticket price in US Dollars



# 180,094 TICKETS

Took THIRTEEN Hours!!!



# ADDITIONAL DATASET



Average and median salaries of each state in USA

- From Wikipedia

# CLEANING



Drop duplicated data



Remove festivals from artist



Unify the level feature



Remove outliers



# FEATURE ENGINEERING



Extracting the venue, city, state, date, and time



Extracting the day, month, and year from date



Adding a price “class” feature:

- 0 for cheap, 1 for expensive

# AFTER CLEANING & FEATURE ENGINEERING

13 Features and 54,234 tickets

- Artist
- Level
- Venue
- City
- State
- Time
- Day
- Month
- Year
- Median Salary
- Average Salary
- Price
- Price Class

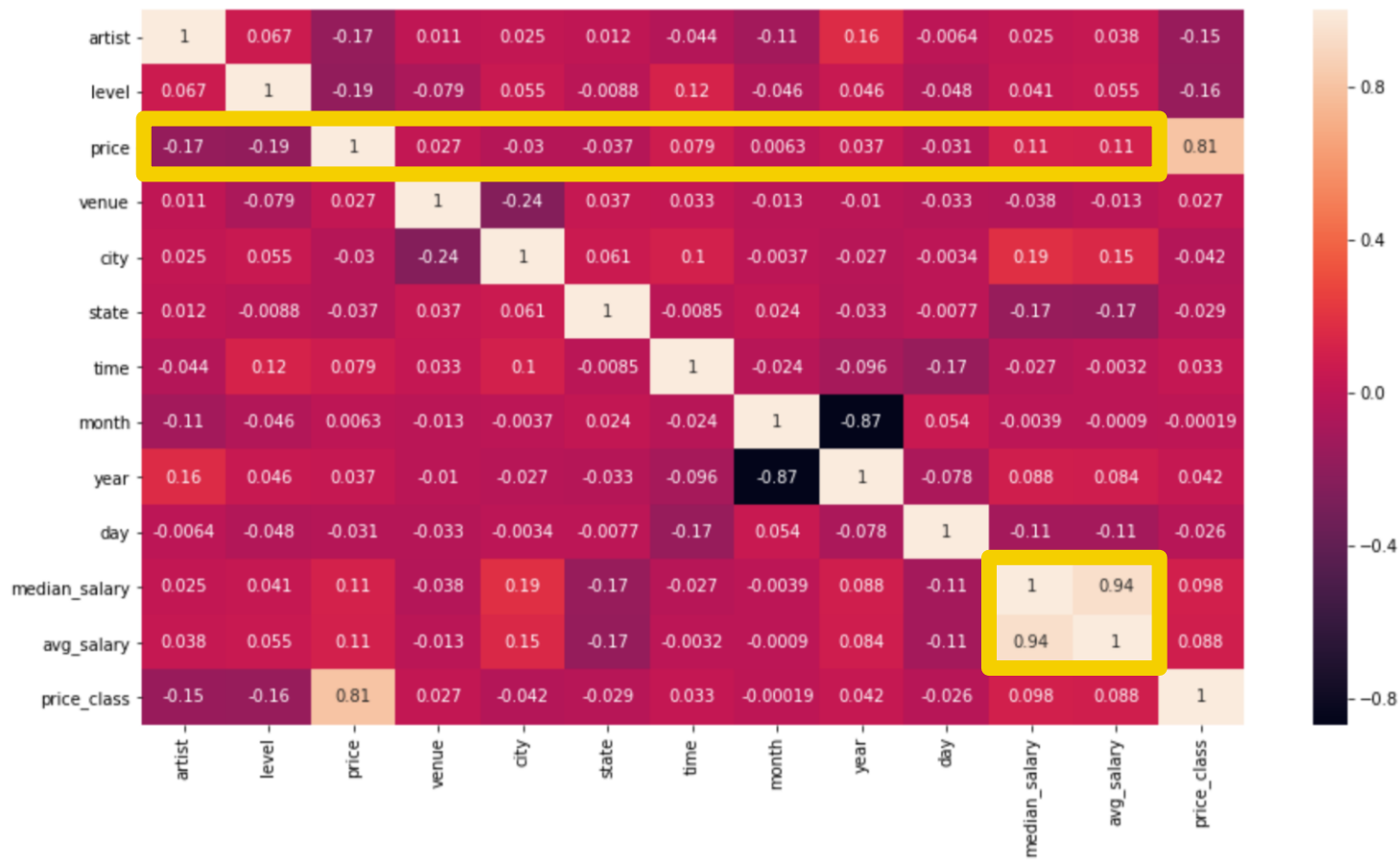
# VISUALIZATIONS

Moving to Tableau...

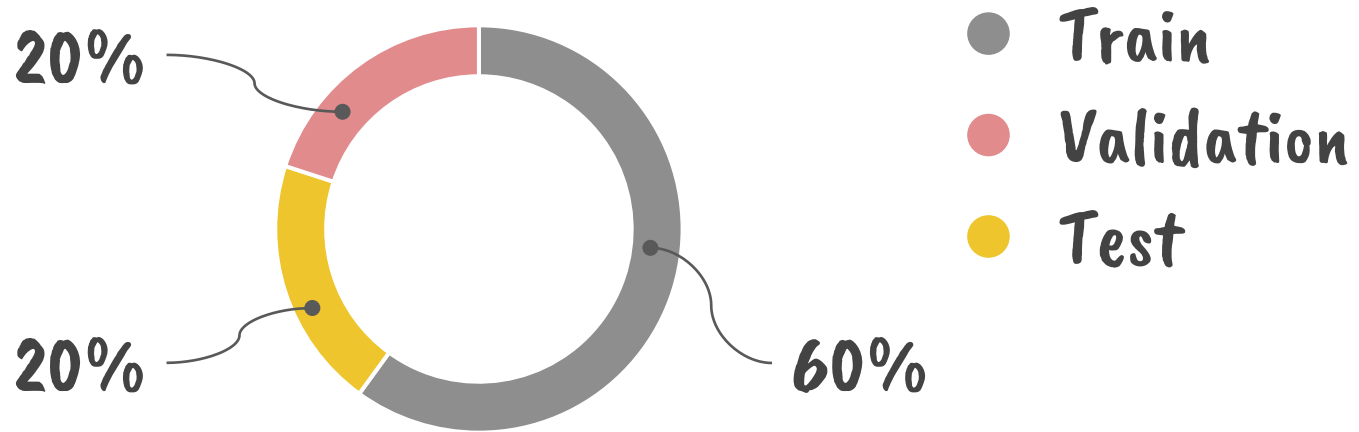




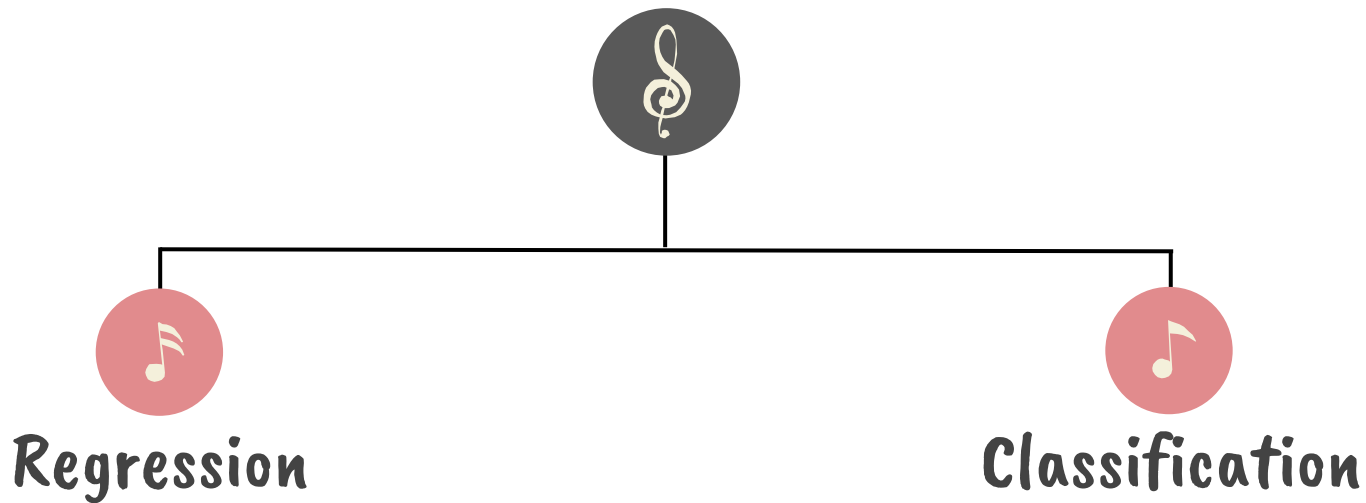
# CORRELATION BETWEEN THE FEATURES



# SPLITTING THE DATA



# MODELING



# LINEAR REGRESSION RESULTS

R Squared	MAE
0.103485	93.756173

Very bad...

# REGULARIZATION RESULTS

	R Squared	MAE
Lasso	0.102097	92.454616
Ridge	0.102108	92.453806
Elastic Net	0.102110	92.453659



# REGULARIZATION RESULTS

	R Squared	MAE
Lasso	0.102097	92.454616
Ridge	0.102108	92.453806
Elastic Net	0.102110	92.453659

# REGRESSION RESULTS

	R Squared	MAE
Linear Regression	0.107720	94.688636
Polynomial Features	0.421972	69.884868
Decision Tree	0.741464	44.595658
Ada Boost	0.709937	49.418598
Random Forest	0.745291	44.589948

# REGRESSION RESULTS

	R Squared	MAE
Linear Regression	0.107720	94.688636
Polynomial Features	0.421972	69.884868
Decision Tree	0.741464	44.595658
Ada Boost	0.709937	49.418598
Random Forest	0.745291	44.589948

# RANDOM FOREST RESULTS

On validation set:

R Squared	MAE
0.745291	44.589948

On test set:

R Squared	MAE
0.762172	42.128088

# CLASSIFICATION RESULTS

	Accuracy	F1
Logistic Regression	0.633631	0.602560
K Neighbors	0.849359	0.848807
Bagging	0.859132	0.858819
Decision Tree	0.863833	0.863353
Ada Boost	0.864294	0.863967
Random Forest	0.860791	0.860482

# CLASSIFICATION RESULTS

	Accuracy	F1
Logistic Regression	0.633631	0.602560
K Neighbors	0.849359	0.848807
Bagging	0.859132	0.858819
Decision Tree	0.863833	0.863353
Ada Boost	0.864294	0.863967
Random Forest	0.860791	0.860482

# ADA BOOST RESULTS

On validation set:

Accuracy	F1
0.864294	0.863967

On test set:

Accuracy	F1
0.867520	0.867361

# CONCLUSION

- Linear Regression was not suitable for this data
- The best regression model is Random Forest
- The best classification model is Ada Boost



THANK YOU!!  
Any questions?

