



Data Incubator Capstone Project Proposal

Michelle (Qin) Peng

Dataset Overview

- **Data source:**

Ebay API and Ebay WebScrapping

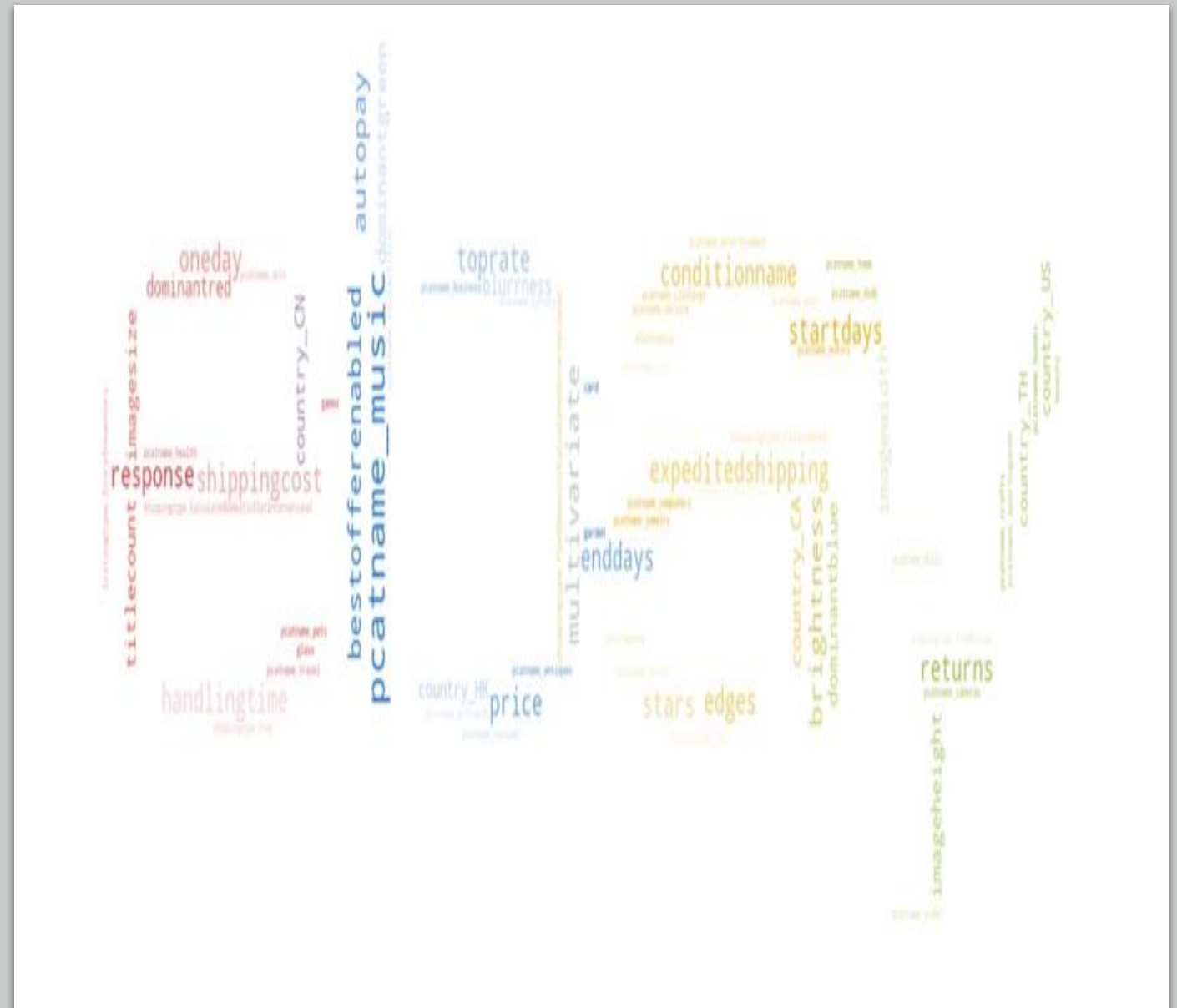
- **Number of Observations:**
708

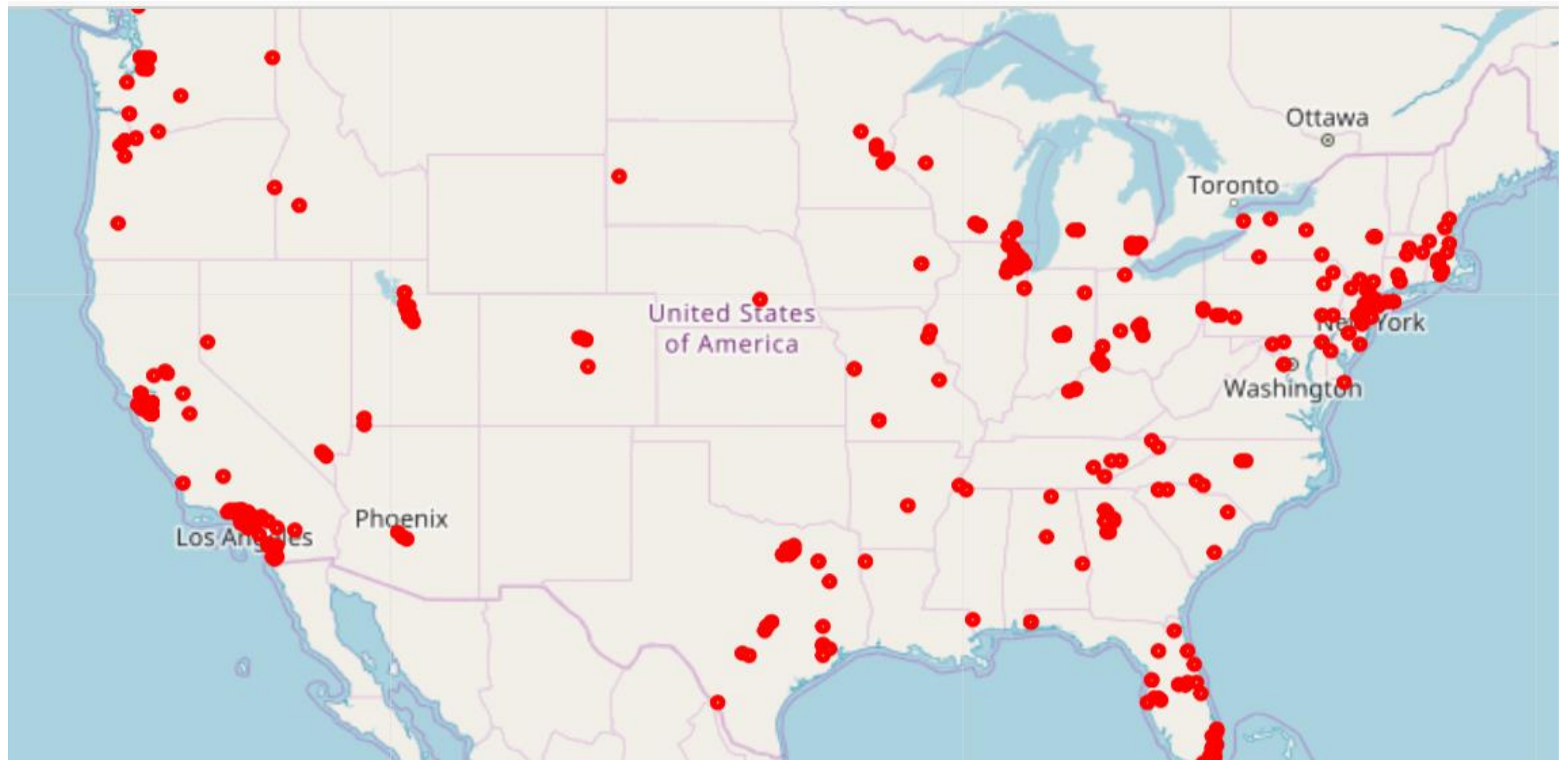
- **Response variable:**

Sales Percentage (How much is sold)

- **Features:**

31 (Categorical Variables (6), Numerical Variables (25))

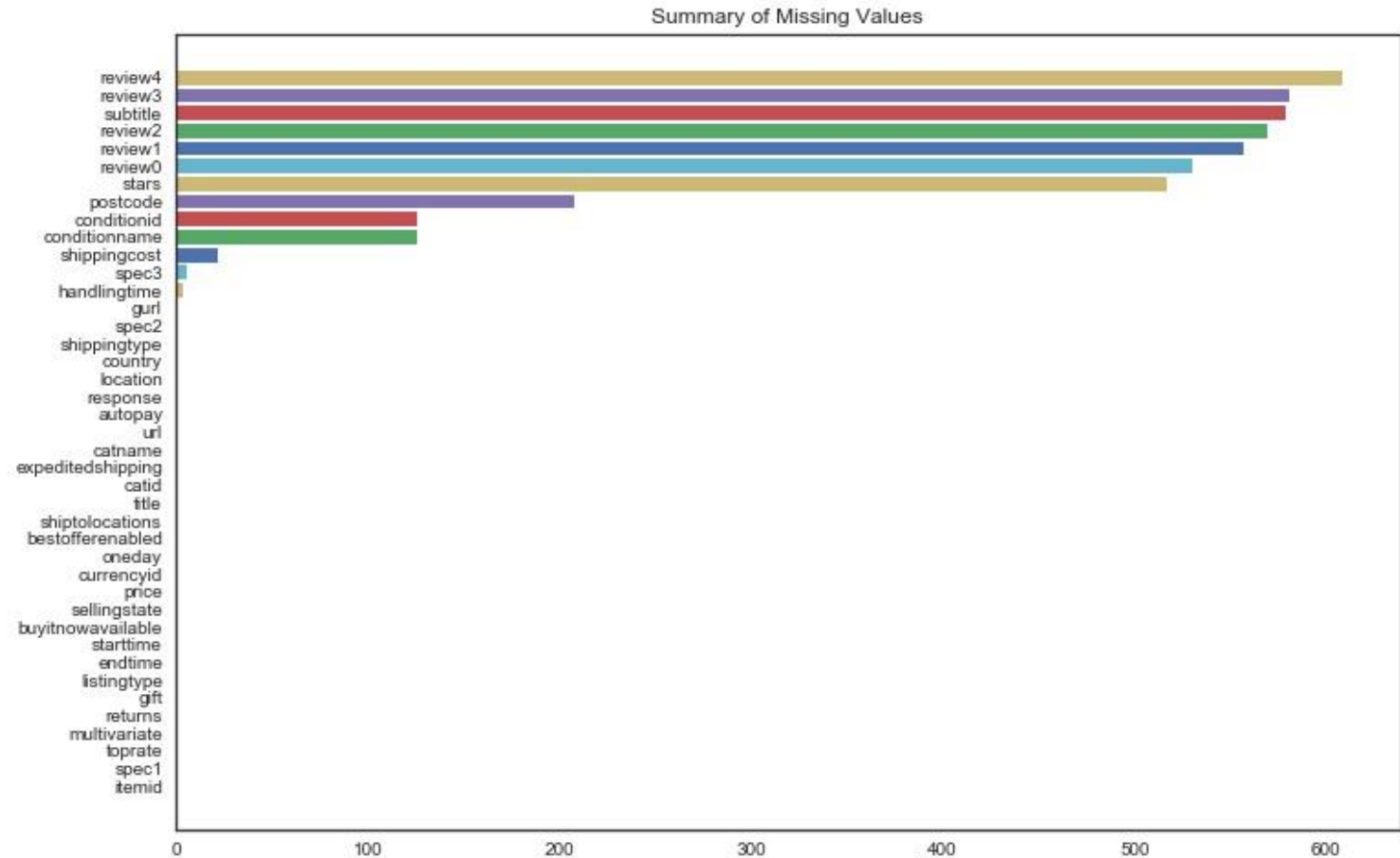


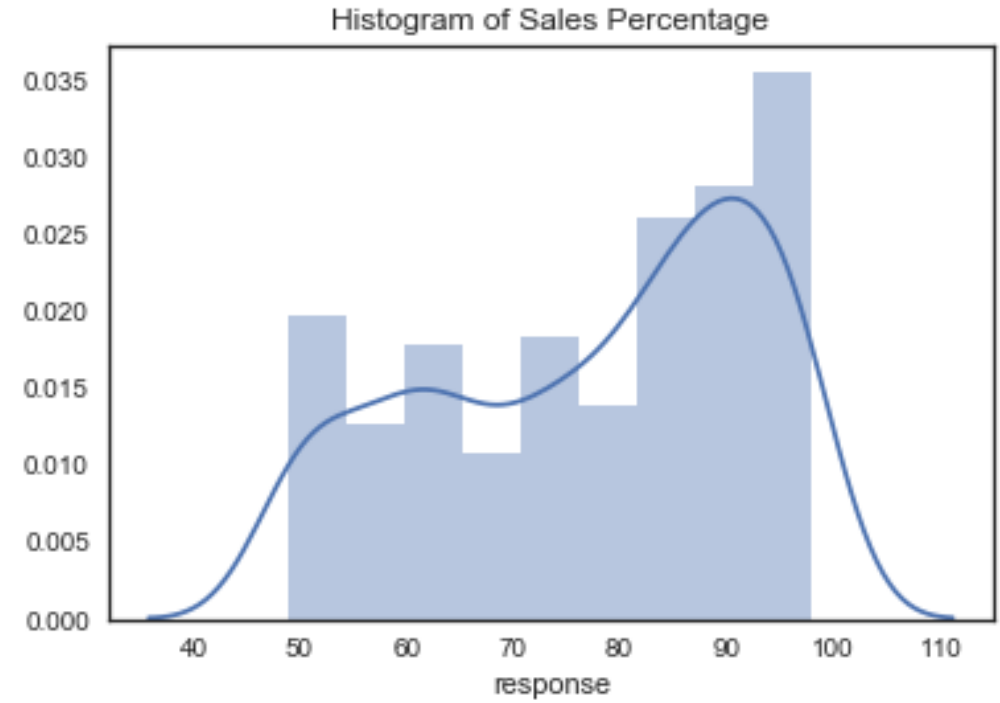
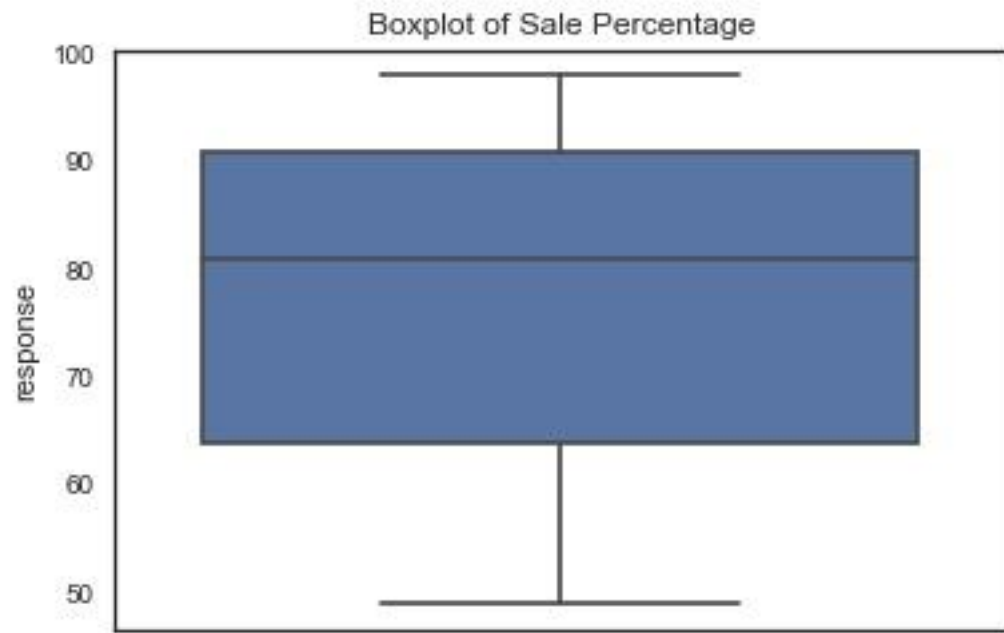




Exploratory Data Analysis

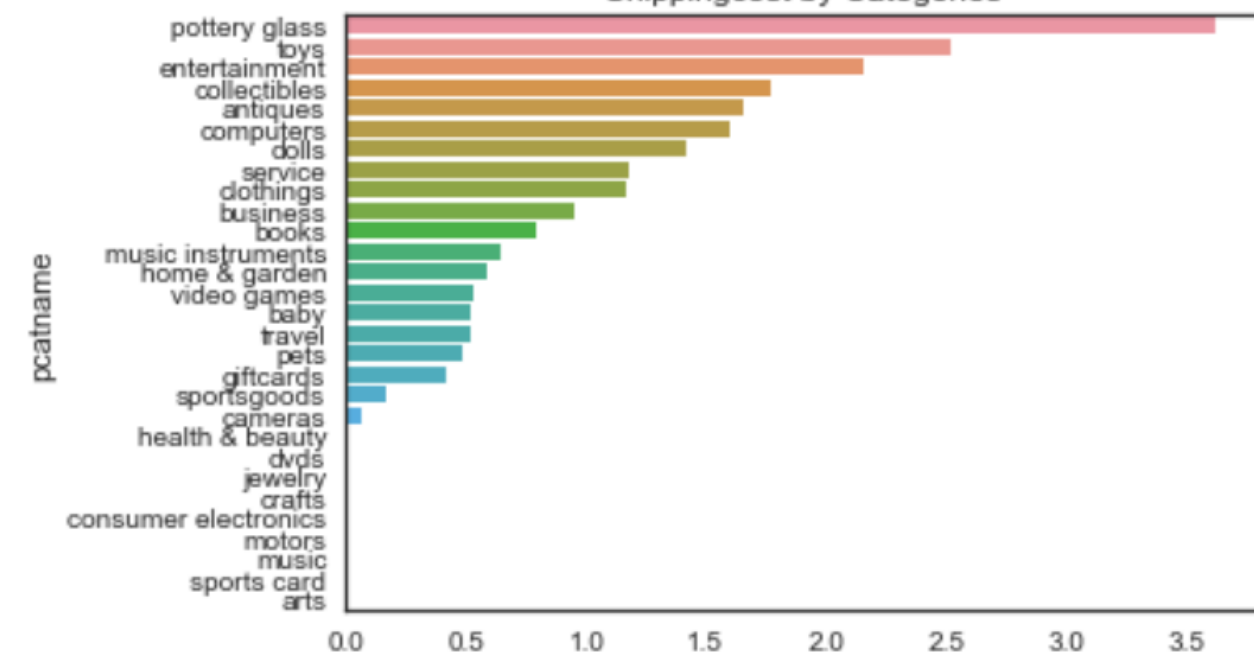
Missing Values



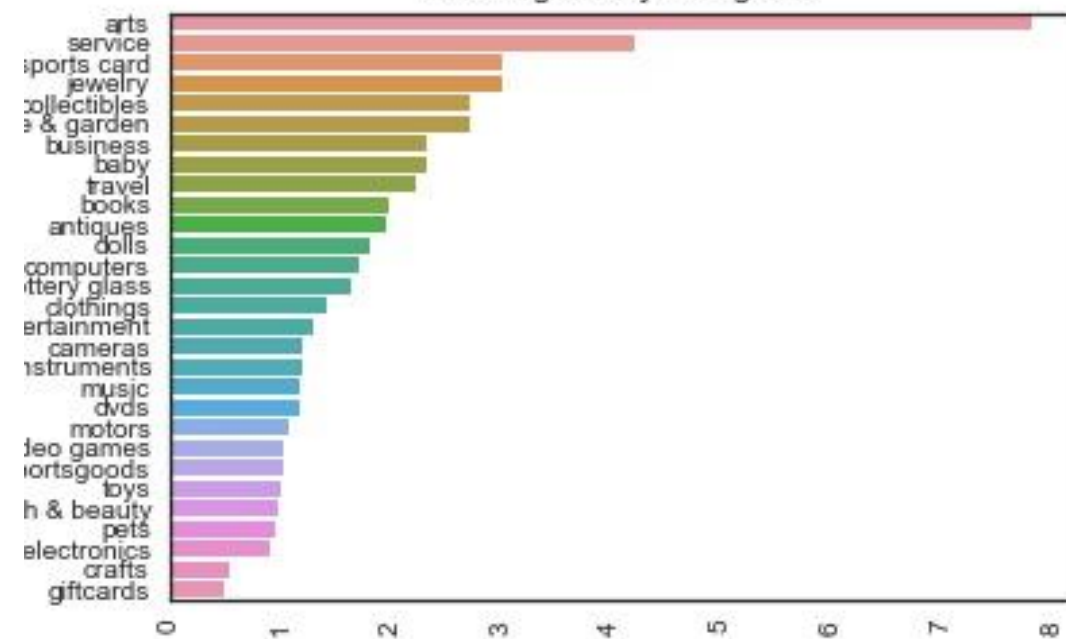


Response Variable: Sales Percentage

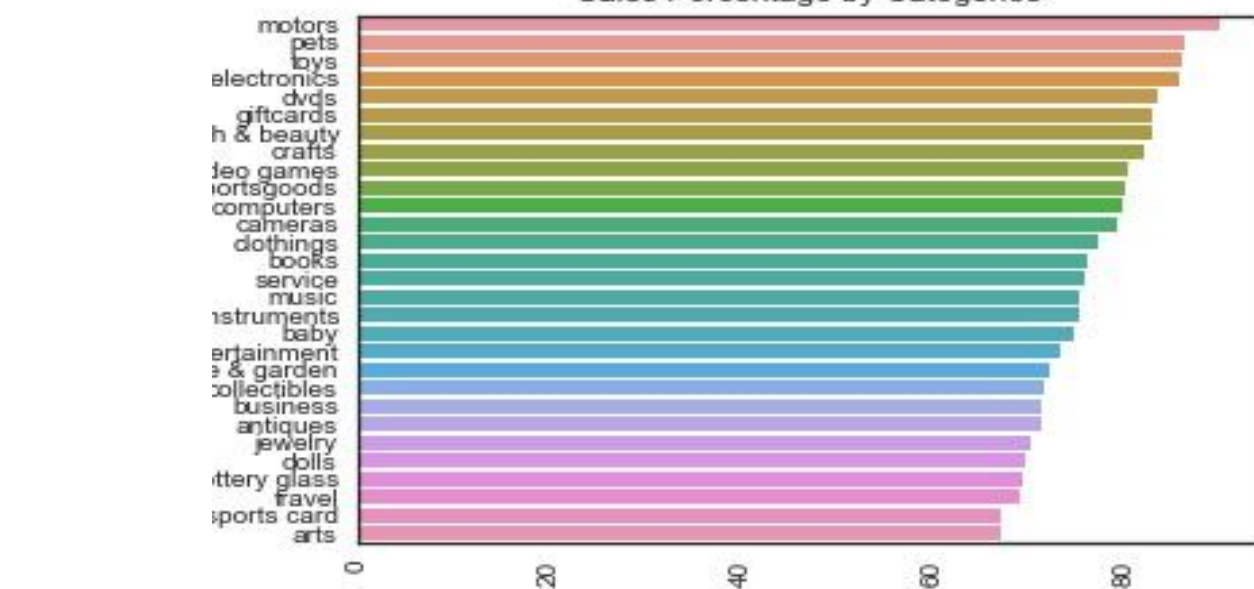
Shippingcost by Categories



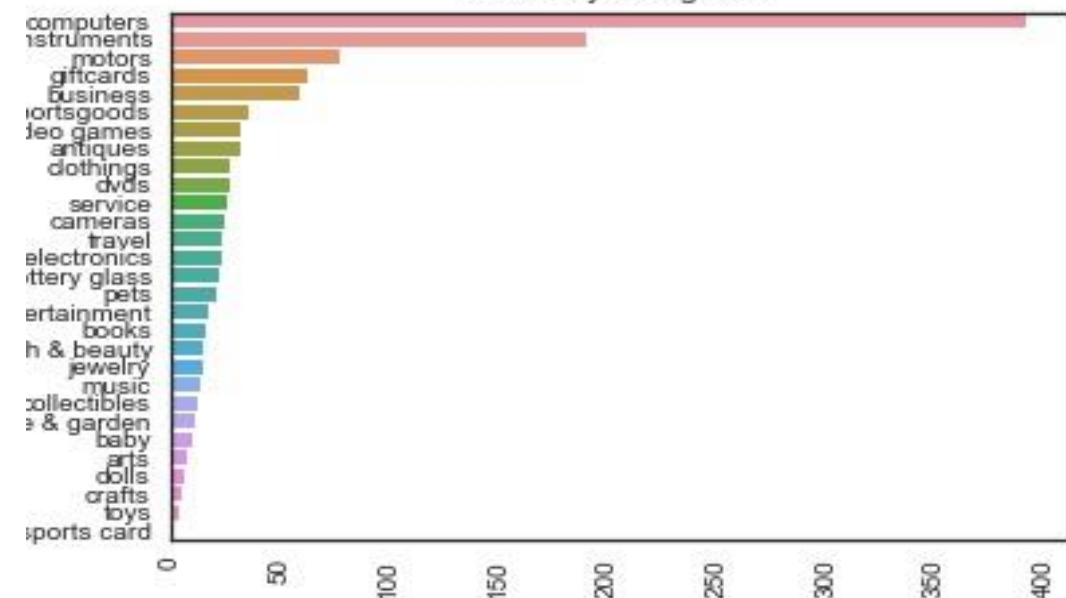
Handlingtime by Categories



Sales Percentage by Categories



Prices by Categories



A thin vertical black line is positioned to the left of the text.

Feature Engineering

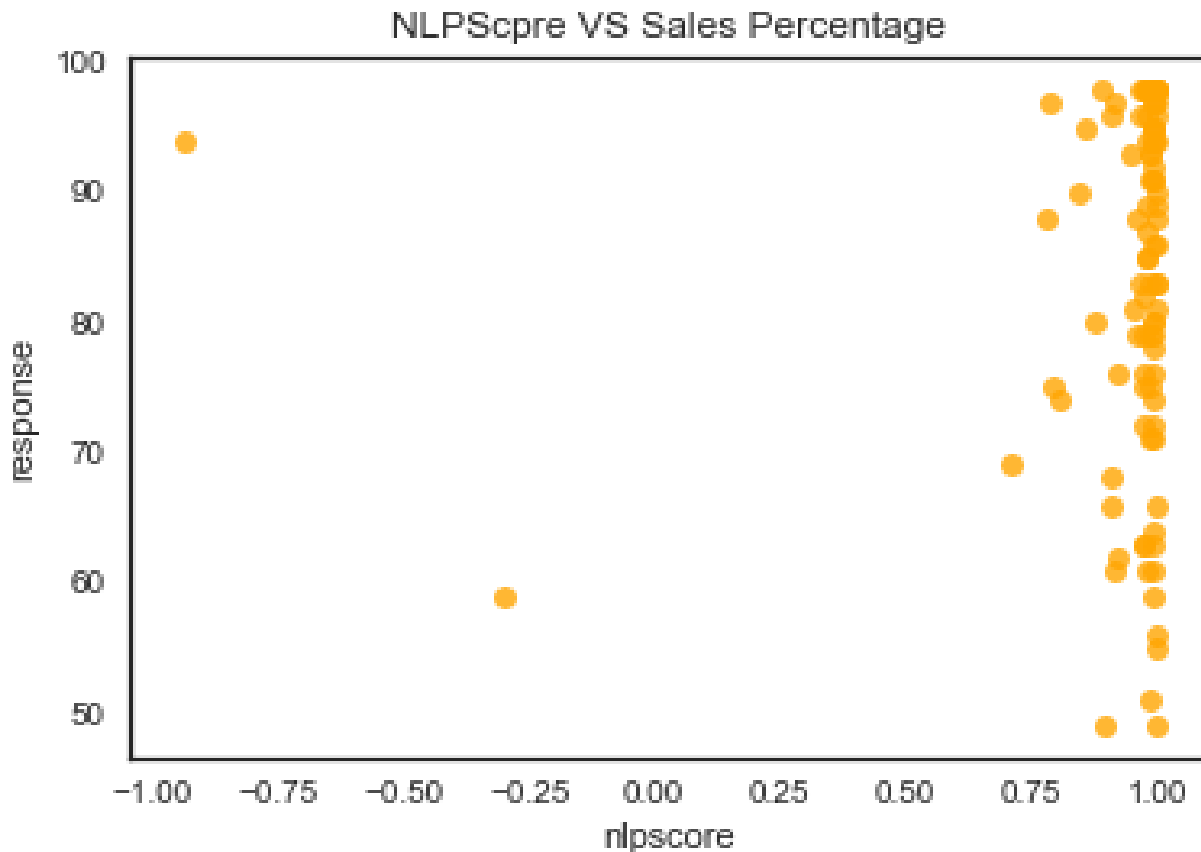
Image Variables:

Each picture is broken into Brightness, Blurness, Dominant Colors, Size



NLP Variables:

A sentiment score is calculated for each review using the VaderSentiment



"I\'m a longtime Shania fan but "Now" is a dud in my books.

It really saddens me to say that.

It\'s very different from what she has previously released

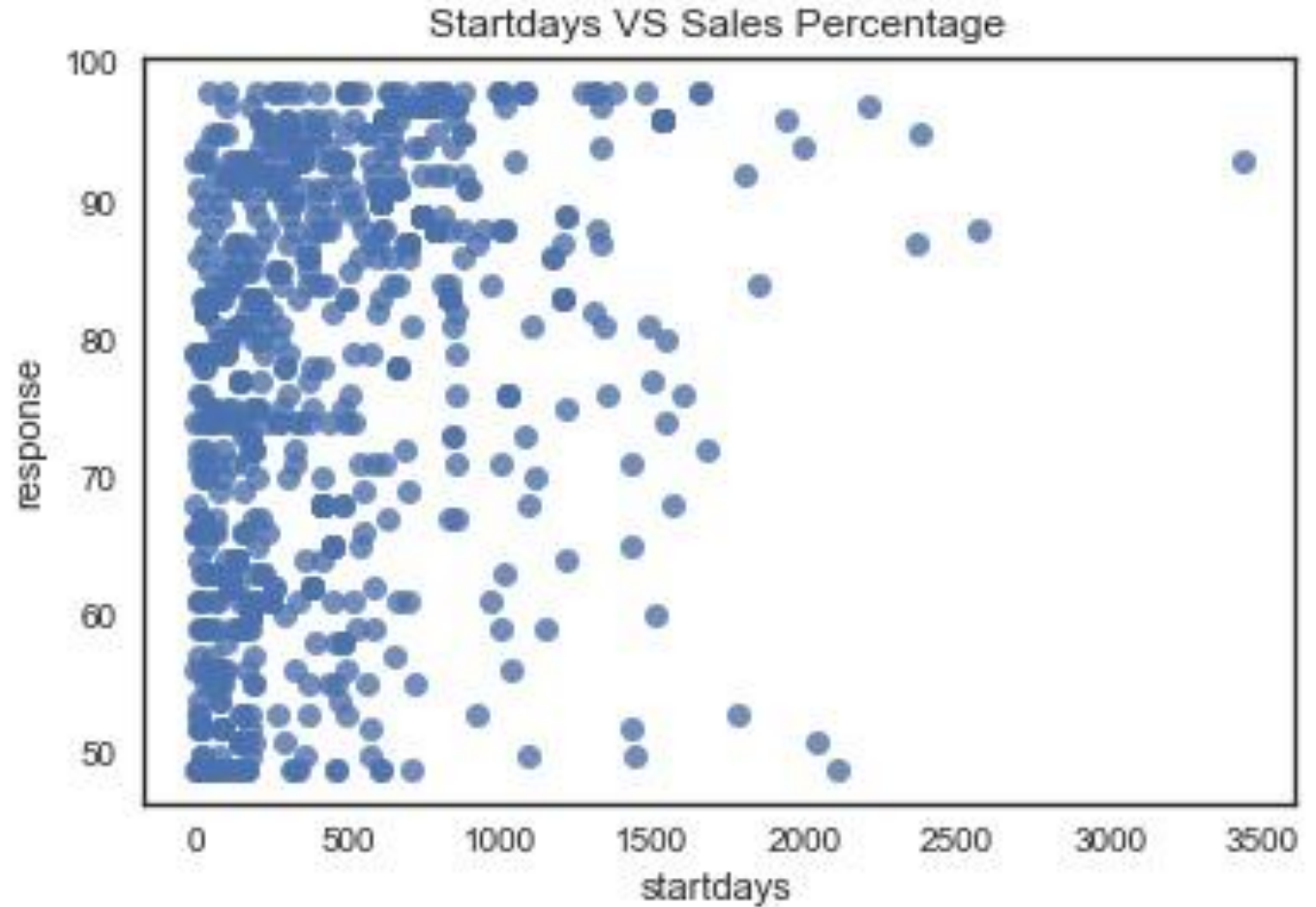
and I just cannot get into it.

If your expecting the sounds you\'ve come accustomed to with.."

↓
{'compound': 0.7781, 'neg': 0.08, 'neu': 0.788, 'pos': 0.132}

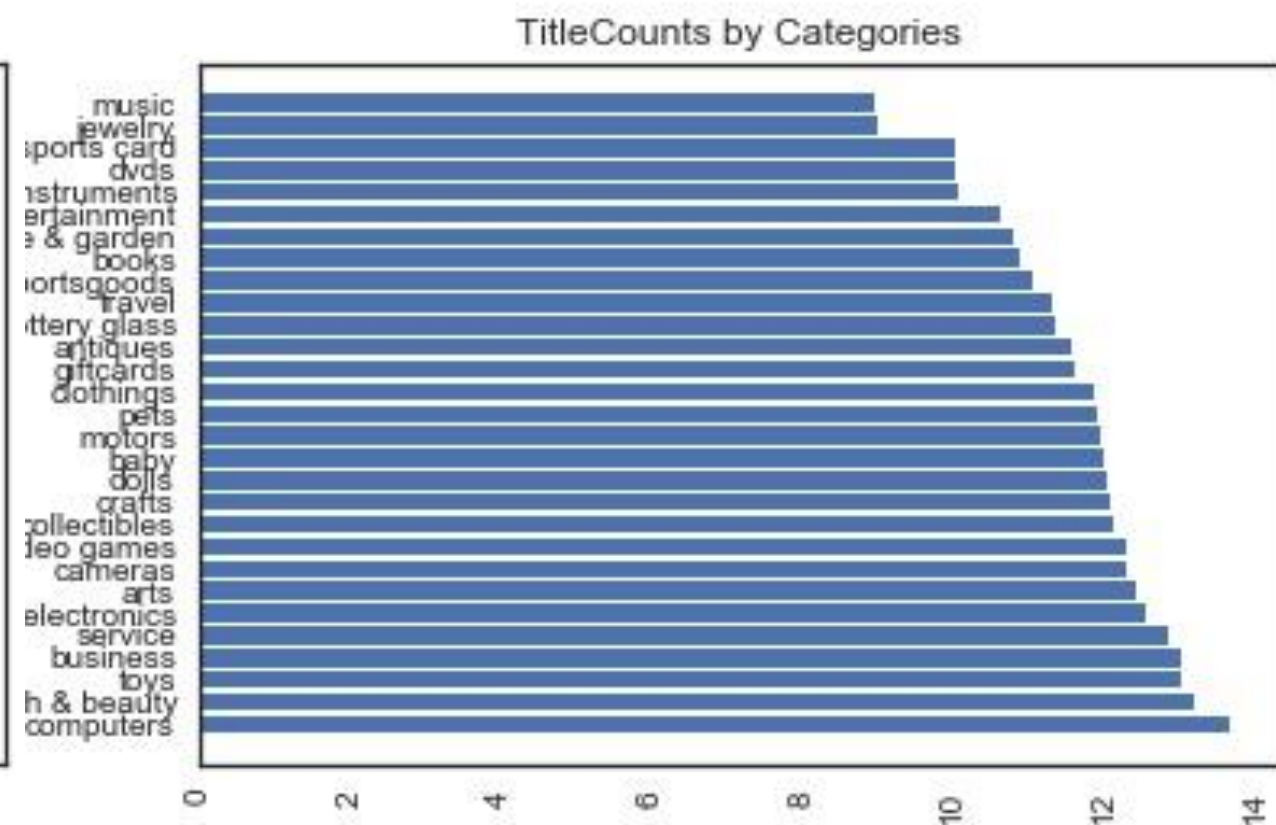
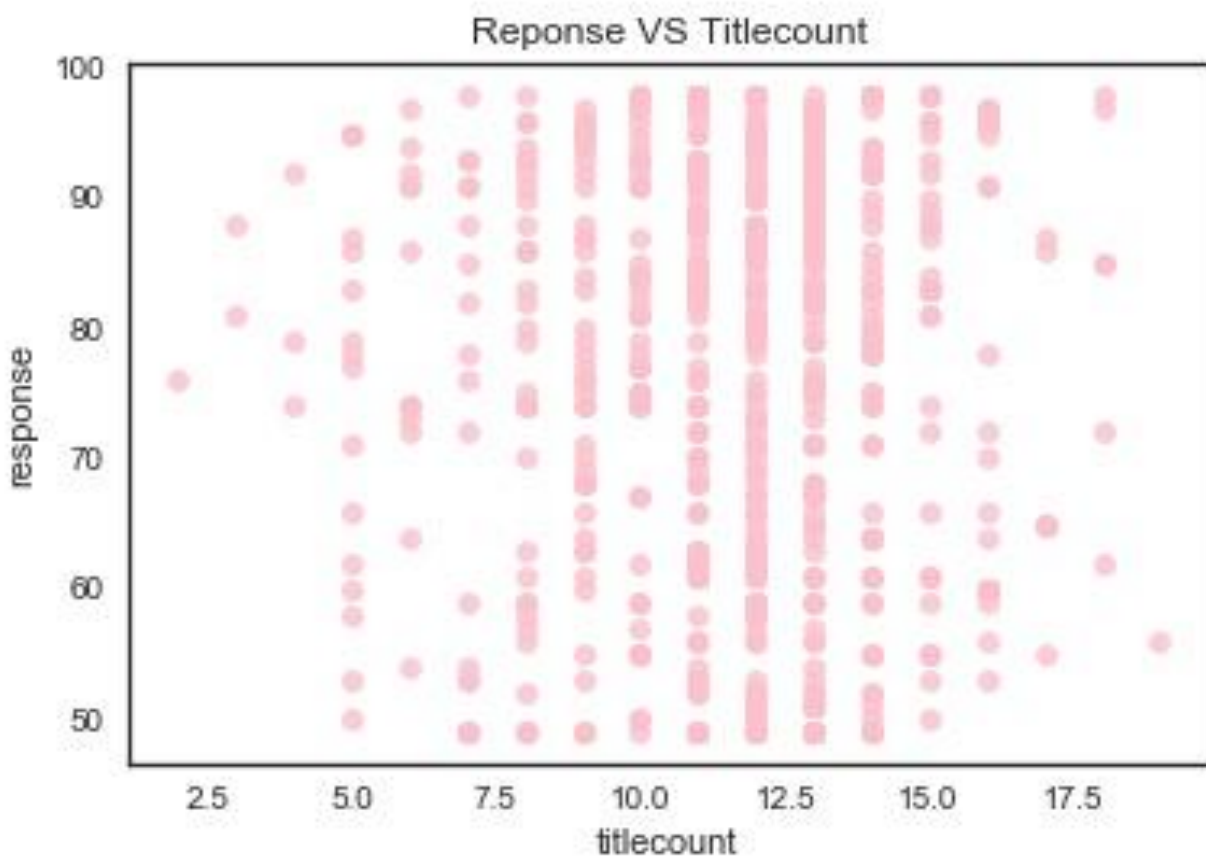
Date Variables

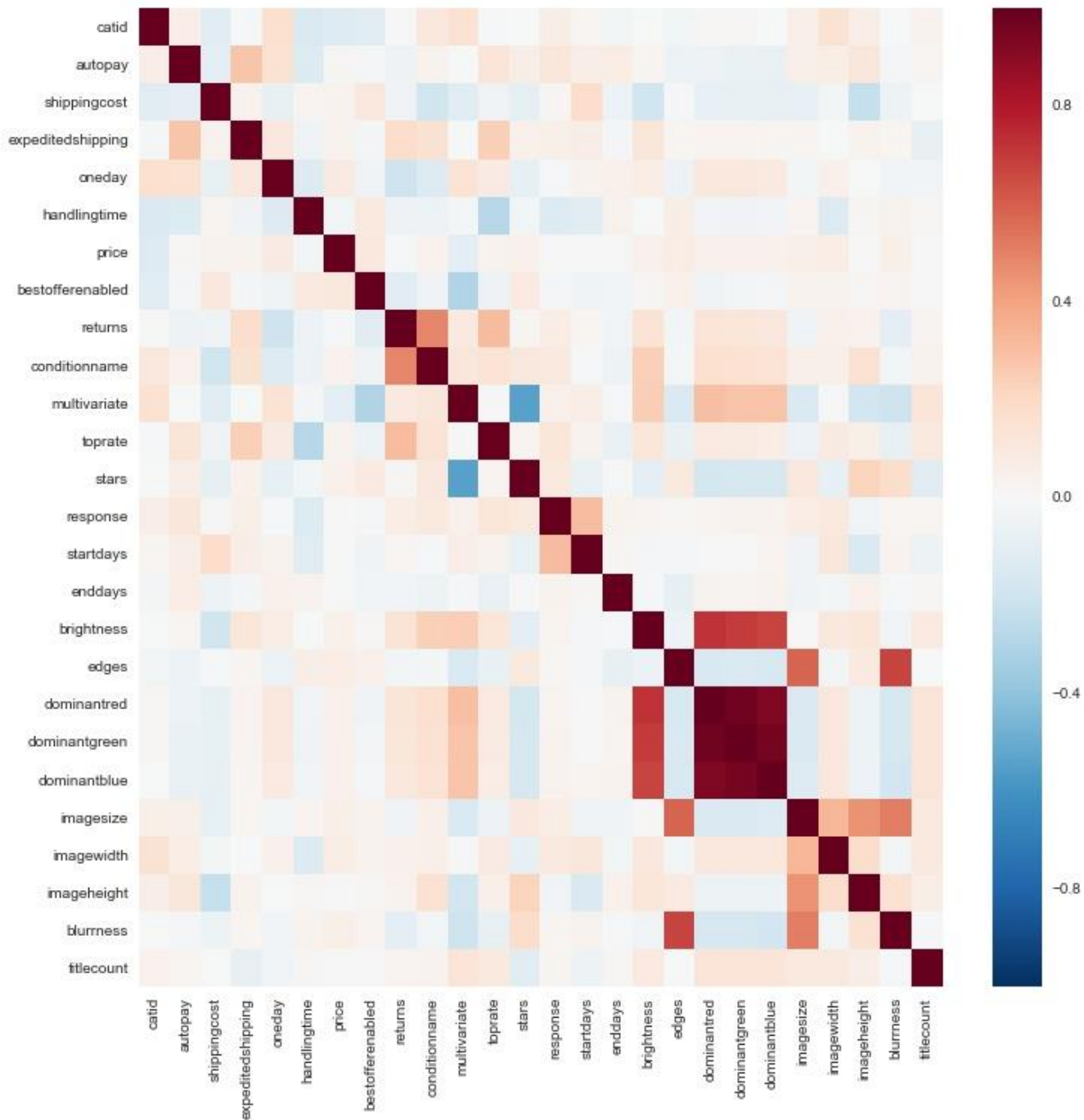
Turns starttime,
endtime into startdays,
enddays which means
how long have the
product has been
posted(startdays) and
how many days
left(enddays)



Title Variables

Construct the counts of words in the title.





Strong Correlation pairs:

1. Image colors
2. Image size VS Image widths & heights
3. Returns VS ConditionName
4. Number of Stars VS Multivariate Listing
5. Sales Percentage VS
6. Startdays

A thin vertical black line is positioned to the left of the title text.

Basic Model Fitting

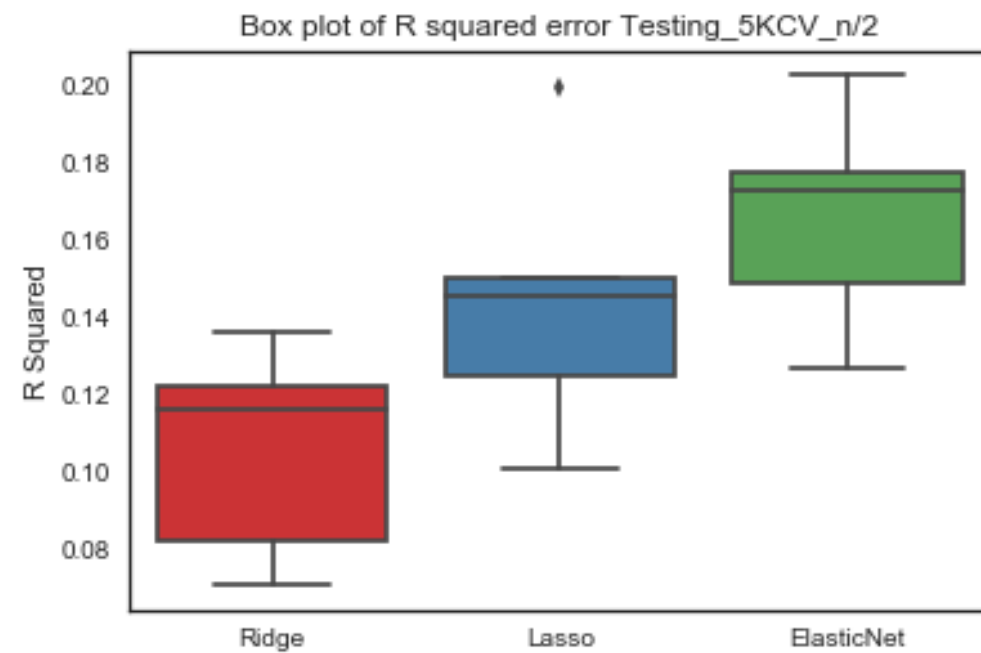
OLS Regression Results

```

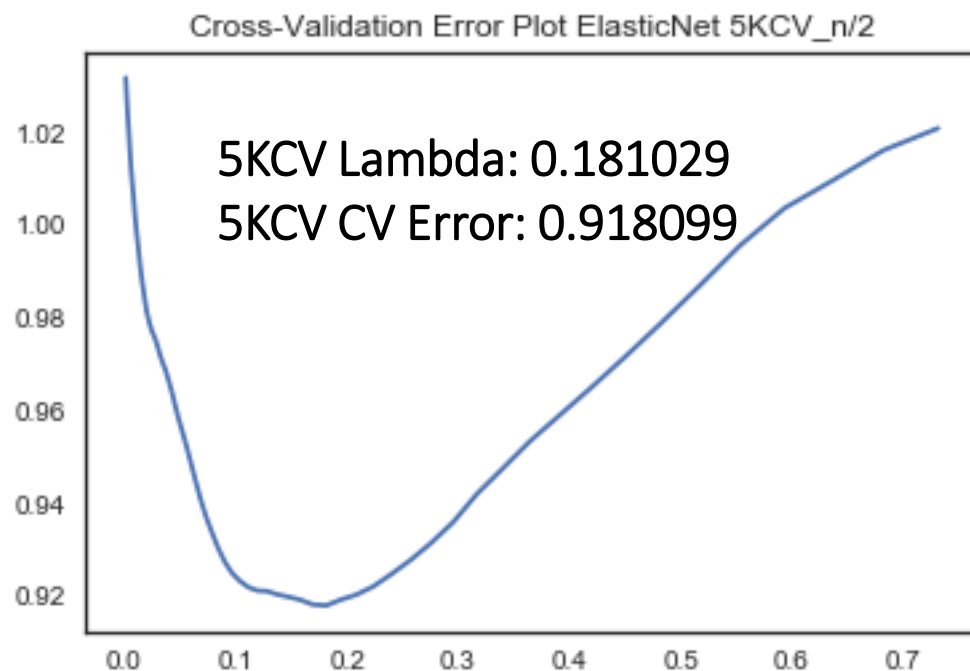
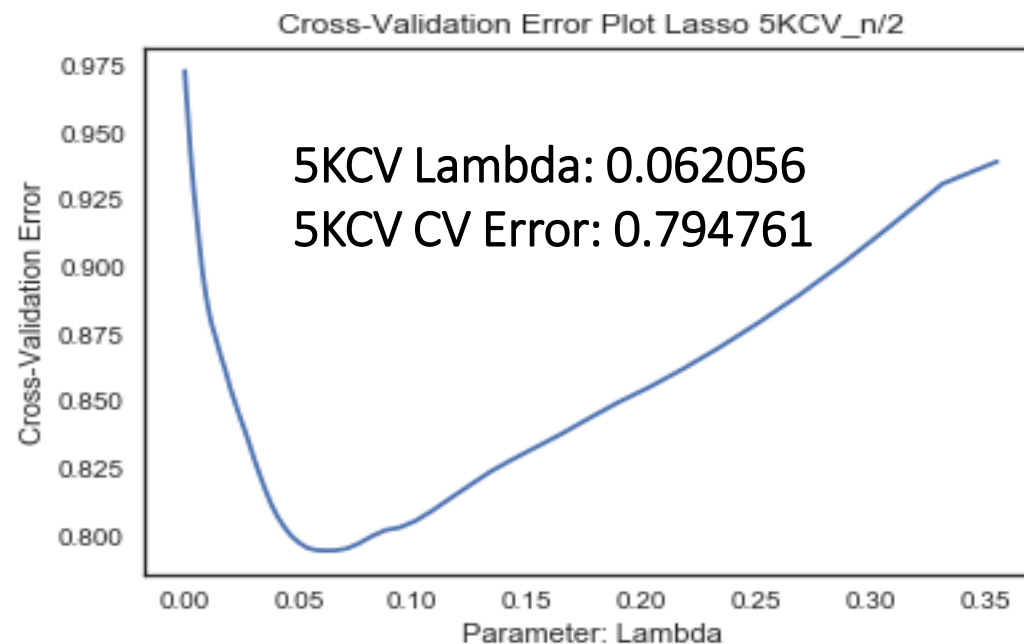
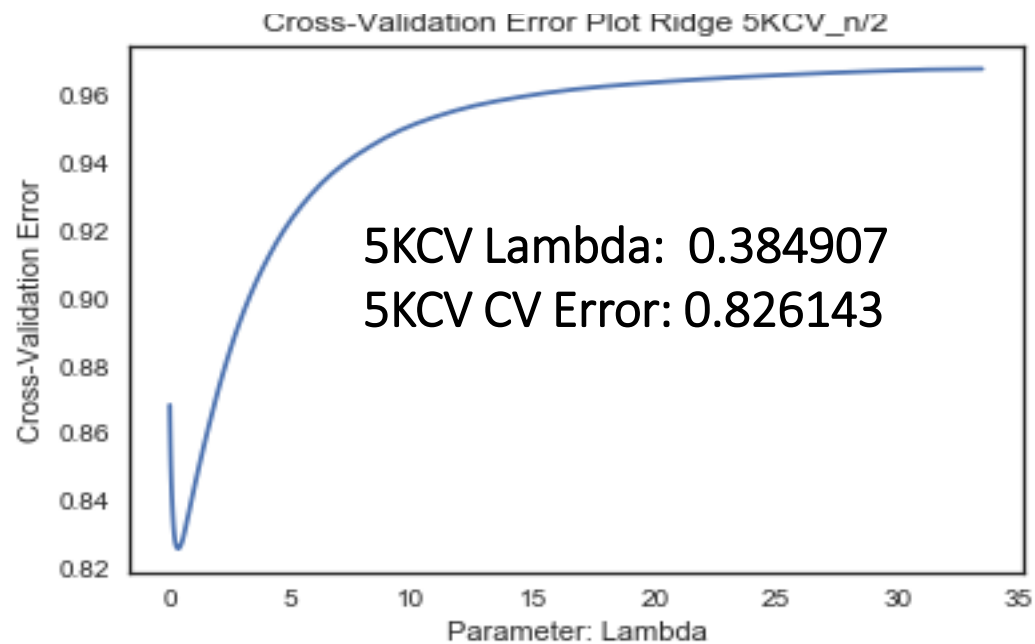
=====
Dep. Variable:          y      R-squared:          0.296
Model:                OLS     Adj. R-squared:       0.227
Method:             Least Squares   F-statistic:        4.276
Date:                Tue, 15 May 2018   Prob (F-statistic):  2.25e-21
Time:                  01:57:18   Log-Likelihood:     -856.22
No. Observations:      692     AIC:              1836.
Df Residuals:          630     BIC:              2118.
Df Model:               62
Covariance Type:       nonrobust
=====

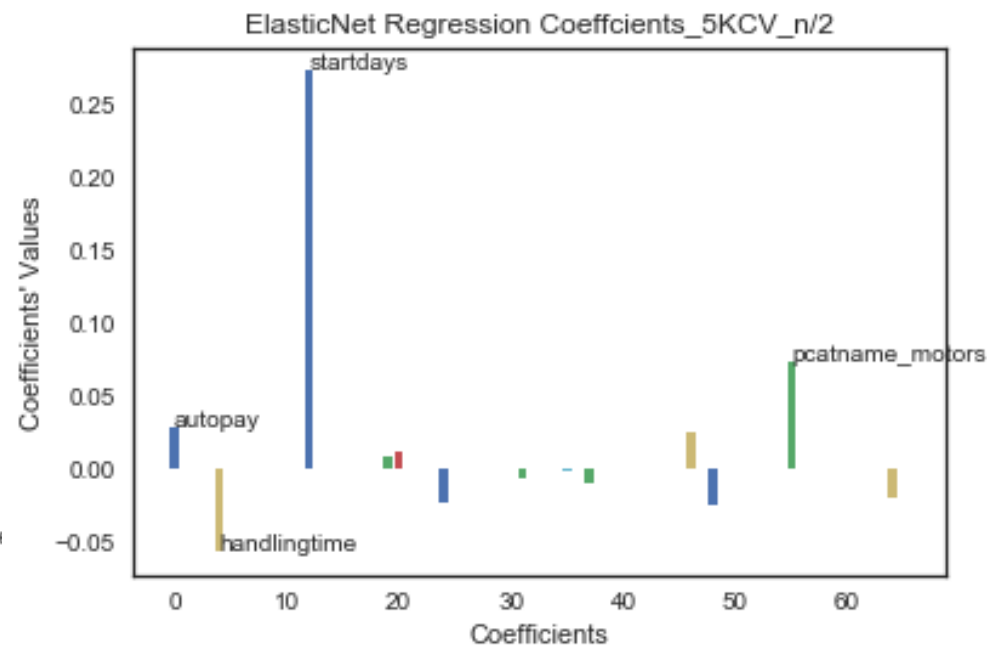
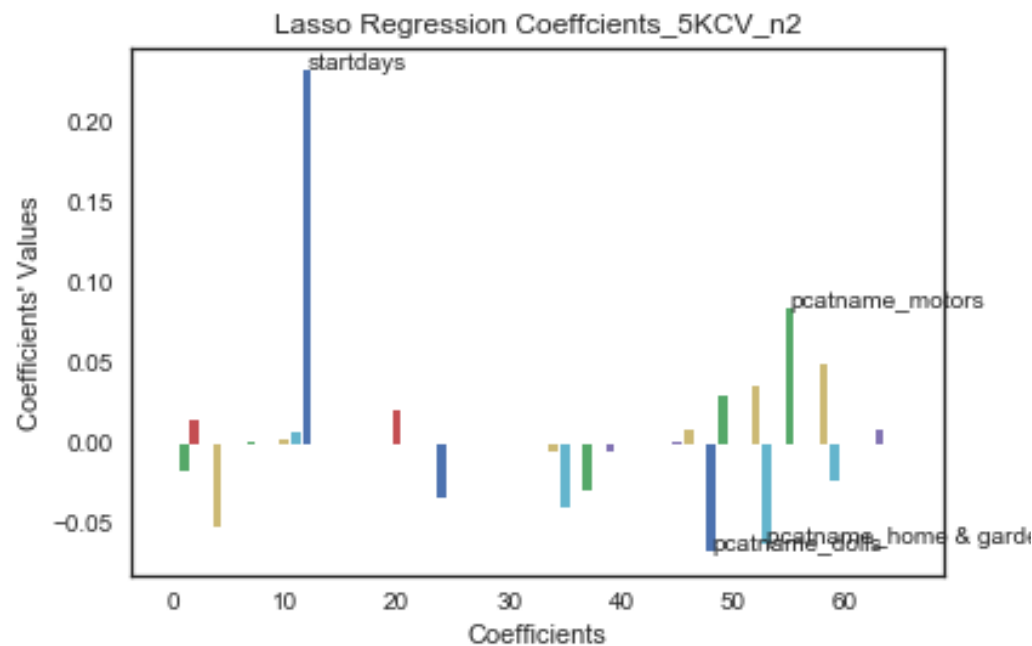
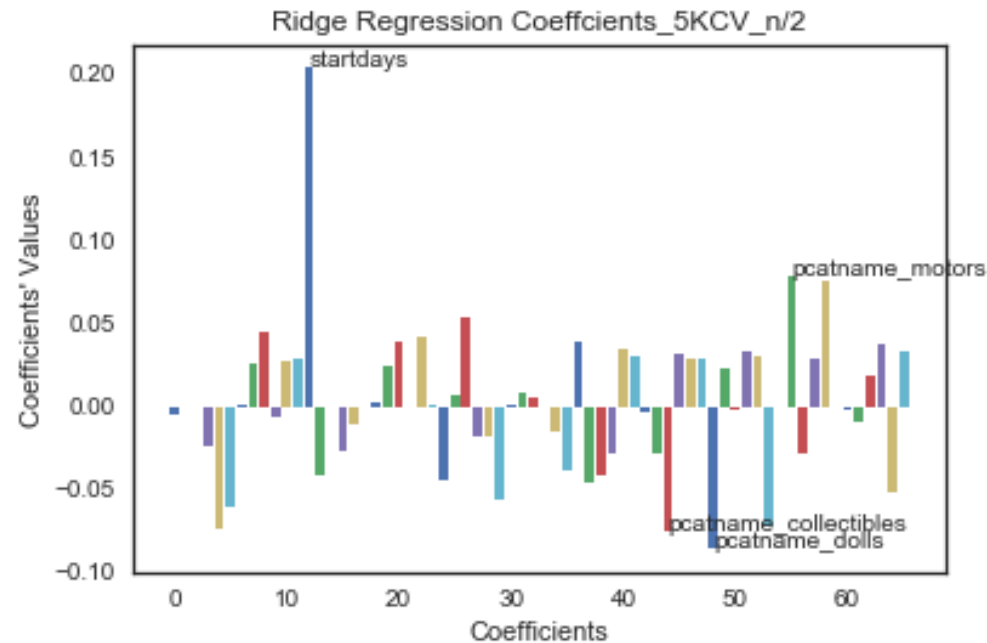
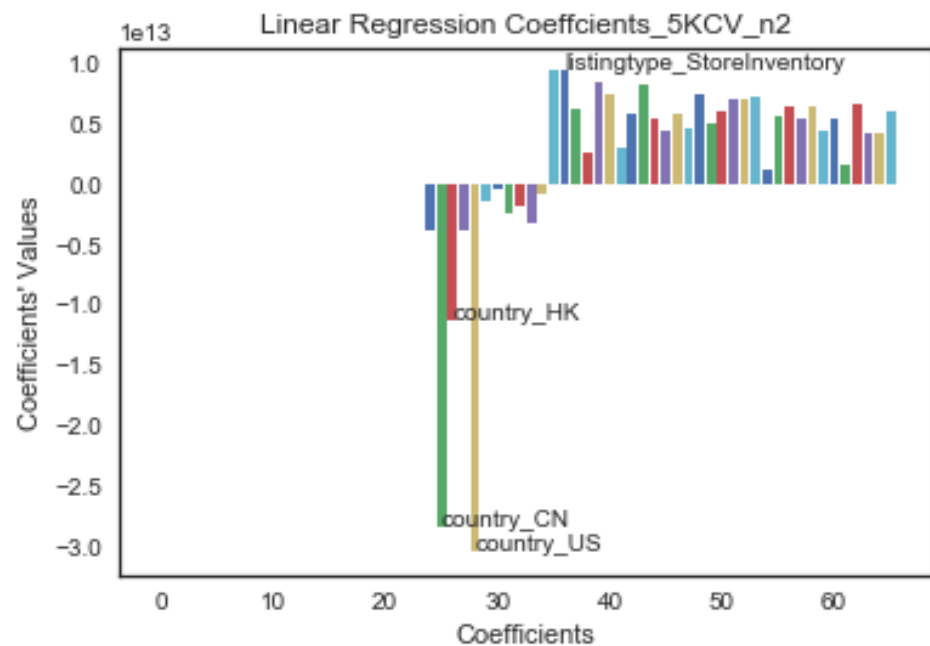
```

	coef	std err	t	P> t	[0.025	0.975]
x1	0.0026	0.042	0.062	0.951	-0.079	0.084
x2	0.0086	0.062	0.138	0.890	-0.114	0.131
x3	0.0023	0.040	0.058	0.954	-0.077	0.081
x4	-0.0408	0.041	-1.001	0.317	-0.121	0.039
x5	-0.0175	0.040	-0.440	0.660	-0.096	0.061
x6	-0.0534	0.041	-1.315	0.189	-0.133	0.026
x7	0.0155	0.038	0.408	0.684	-0.059	0.090
x8	0.1052	0.054	1.944	0.052	-0.001	0.211



0.384907





Thanks for watching!