

# Credit Sesame Data Analysis and Findings

Hari Rajan, Paul Giraud, Audreya Metz, and Calvin Ma

September 27 2018

## Dataset

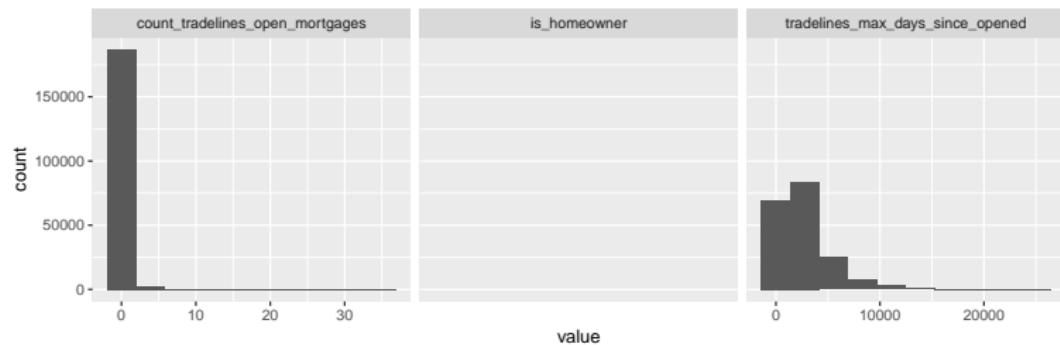
- ▶ The data was from Credit Sesame, a company that calculates credit scores to determine options for credit cards, mortgage rates and loans.
- ▶ This dataset included:
  - ▶ User Demographics
  - ▶ First Session Information
  - ▶ 30-day User Engagement Data

## Our Research Direction

- ▶ We wanted to give Credit Sesame an insightful description of their user base
- ▶ We believe that they would have an interest in the credit score of their users, since this helps them to predict the best loans, mortgages, etc.

# Exploring the Dataset

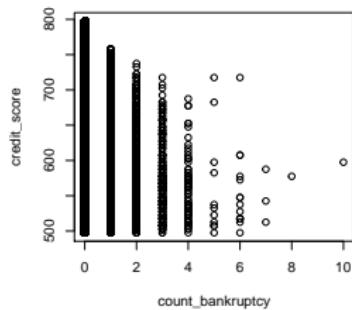
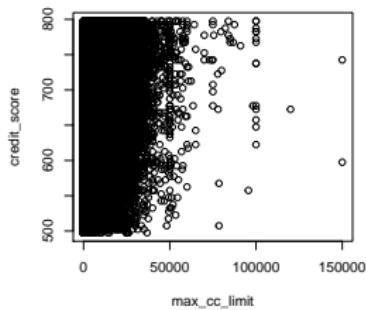
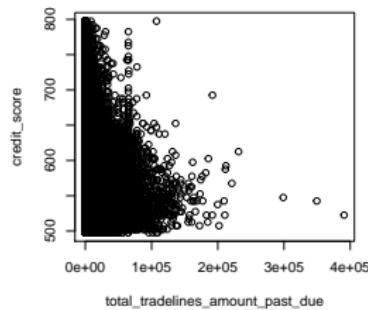
- ▶ We started out EDA by looking at the distribution of the variables in the dataset
- ▶ We found a significant right skew in the majority of the non-categorical data



This was important to keep in mind as we proceeded in our analysis as the distribution of each variable can impact the results of our analysis methods.

# Exploring the Dataset

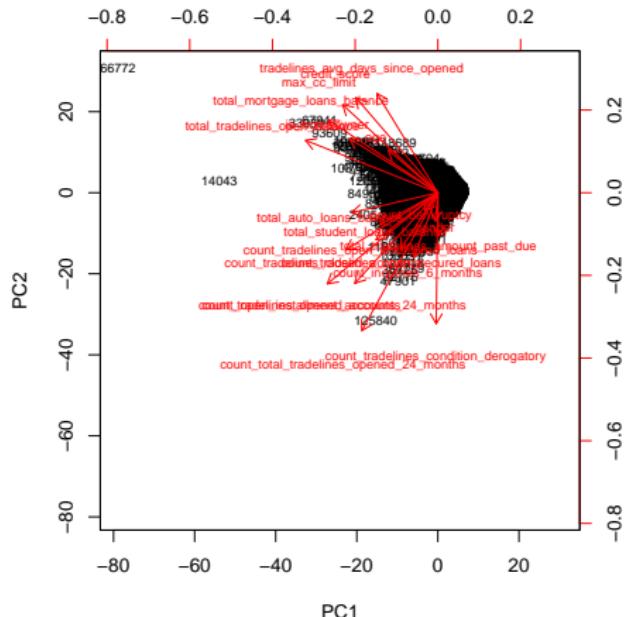
We also wanted to see if there was any obvious relationship between credit score and the others features in the user demographics.



Several of the features had a clear linear relationship with credit score, which motivated a linear regression to determine which features had the greatest impact on credit score. Other features didn't have such an obvious relationship.

# Exploring the Dataset

We ended our EDA using PCA to determine any relationships between data points. Due to the high-dimensionality, the biplot is hard to interpret, but we can still obtain some valuable information from the plot.



## Challenges of the Dataset

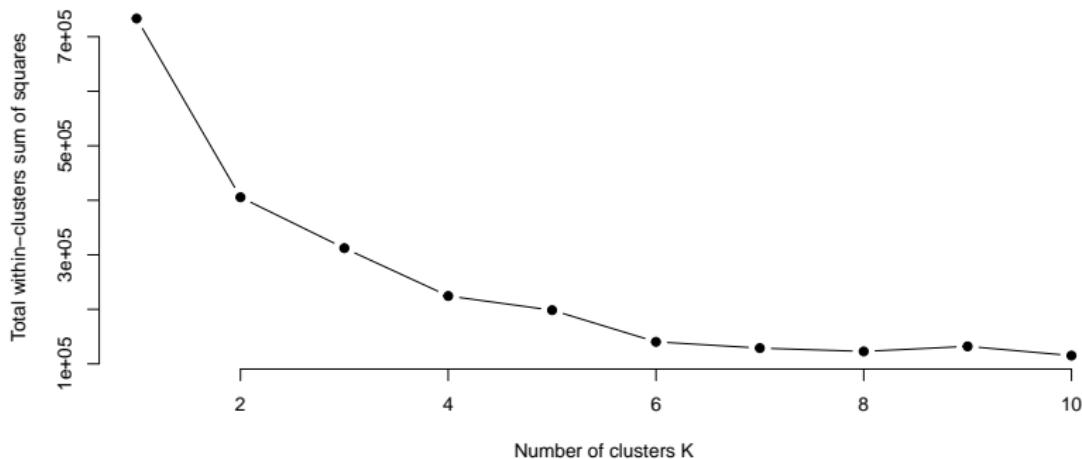
- ▶ The skewed distribution of our variables was the biggest challenge for determining how to analyze this data.
- ▶ Ultimately we did not let the skew effect our methods but were careful to consider it in our analysis.

# Methods

- ▶ Linear Regression
  - ▶ We determined which features to use in the regression using the plots and PCA
  - ▶ Only achieved an R-squared of .35, which is too low to draw meaningful conclusions
  - ▶ Could use other regression methods such as Lasso in the future
- ▶ K-means clustering
  - ▶ Cluster users based on most frequent page type
  - ▶ Compare within cluster user demographics with those of other clusters
  - ▶ Find differences between users

## K-means

- We use k-means because it is the quicker than hierarchical clustering, but we first need to define k. We look at the plot of k vs. WSS.



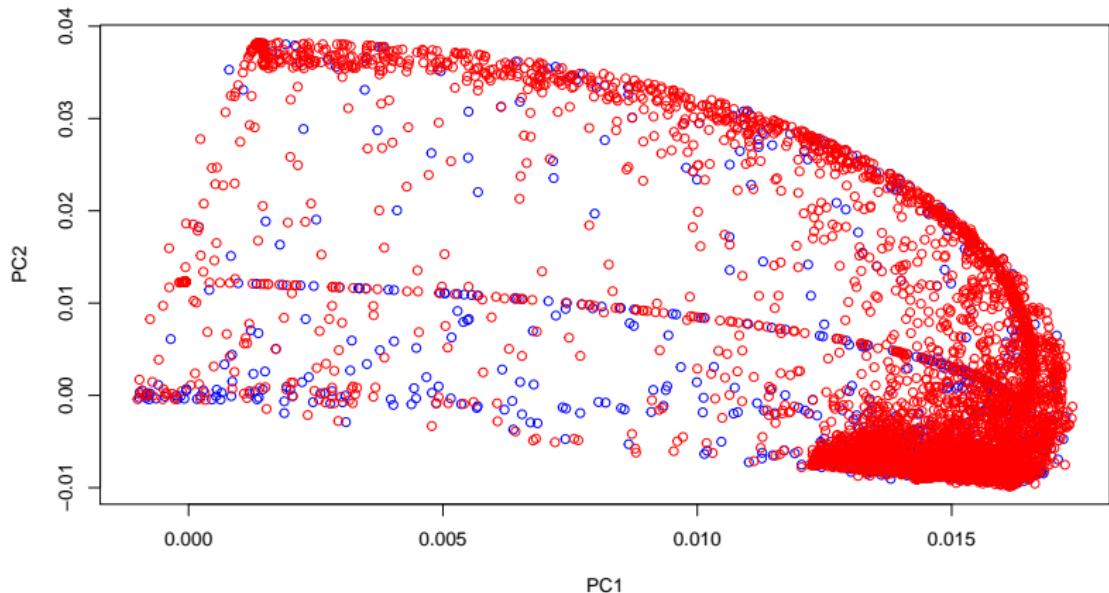
We use the Elbow method to determine the optimum number of k.

## Results

- ▶ After determining the number of clusters to use and accordingly performing the k-means clustering method, we noticed that in particular, there were two clusters which had relatively high values for `click_count_credit_card` and `click_count_personal_loan`.

```
##      click_count_credit_card click_count_personal_loan
## 1              0.2266731          2.40265451
## 2              0.0000000          0.01661813
## 3             10.5436328          0.21913381
## 4              3.9949445          0.12004829
## 5              1.2651661          0.03271203
## 6              0.7094105          6.74250259
```

## Results (continued)



There is a clear difference in the distribution of red dots (users clicking on credit pages) vs. distribution of blue dots (users clicking on loans).

## Results (continued)

The table below shows a comparison of the features that were relatively different between users who clicked credit and users who clicked loans.

	click_credit	click_loan
## open_installment_acc	2.663848	0.9634536
## unsecured_open_tradelines	3731.251586	1408.5340338
## open_cc_amount_past_due	1520.512685	4264.2048881
## tradelines_open_balance	1834.285412	4873.6174052
## auto_loans_balance	10737.634249	12049.6852444
## credit_score	587.521142	650.9730471

# Conclusions

- ▶ Our regression model was inconclusive
- ▶ The k-means clustering confirmed the demographics of people using Credit Sesame's website differentiate across features
- ▶ There is a distinction between clientele that uses the credit card features and loan features
- ▶ Future Work:
  - ▶ We could create improved regression models using techniques such as Lasso
  - ▶ Hierarchical clustering could have provided more insight into potential variable groups