# Groupons: Understand and Group

By: Nicholas Maloof

# Introduction

# What is Groupon?

A website that publishes coupons known as "groupon's" from different venders that allow you to receive a discounted price for certain activities, food purchases, or etc.



Limited Time Remaining!

Up to 53% Off

★★★★★ 3,846 Ratings

Two Hours of All-You-Can-Eat Sushi, Sashimi, and Teriyaki Dinner with Two Beers, Wine, or Sake

Over 25,000 bought

~~$60~~ **$28**

53% off

1

Buy

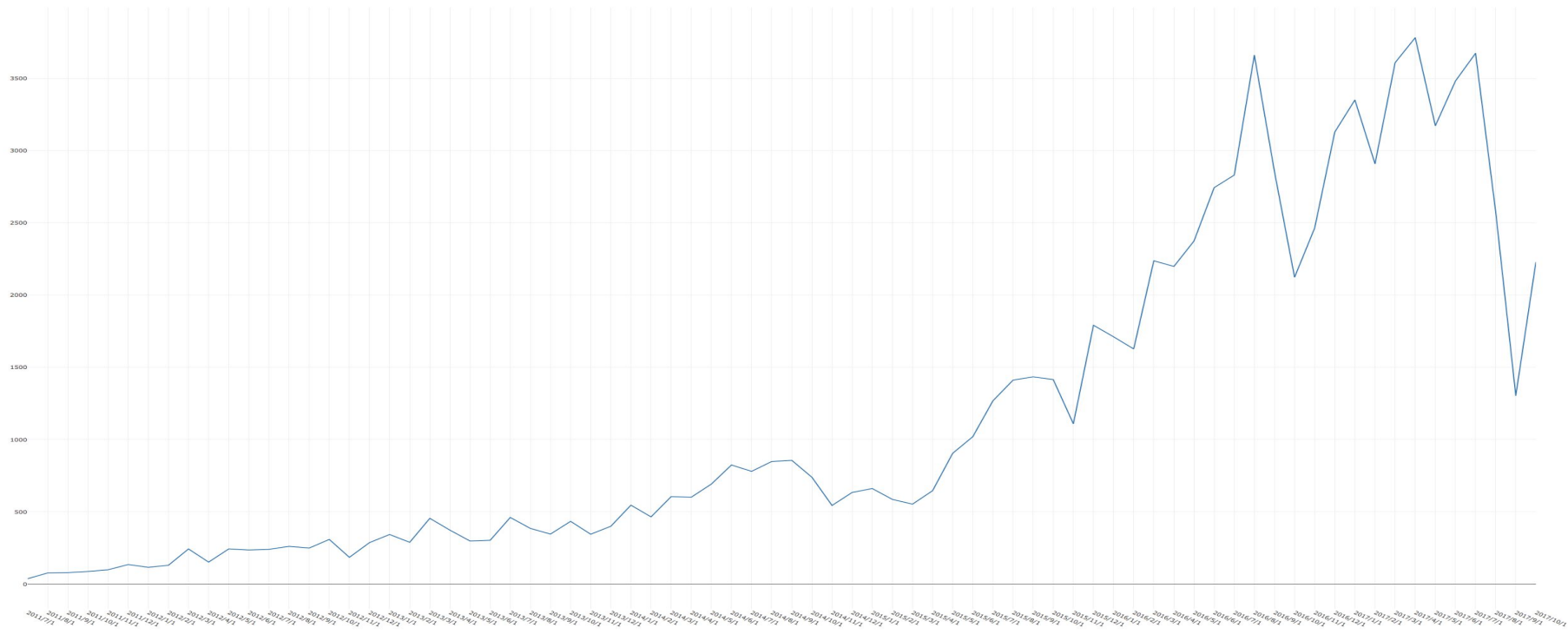🎁 Give as a Gift

# Scraping

# Scrapy the Fast
## and
# Selenium the Slow

- Scrapy
  - Scraped all the available local groupon data present at the time.
  - Scraped everything important from the individual groupon
  - Compiled a list a urls
- Selenium
  - Using urls, went into each groupon and clicked on the read reviews button
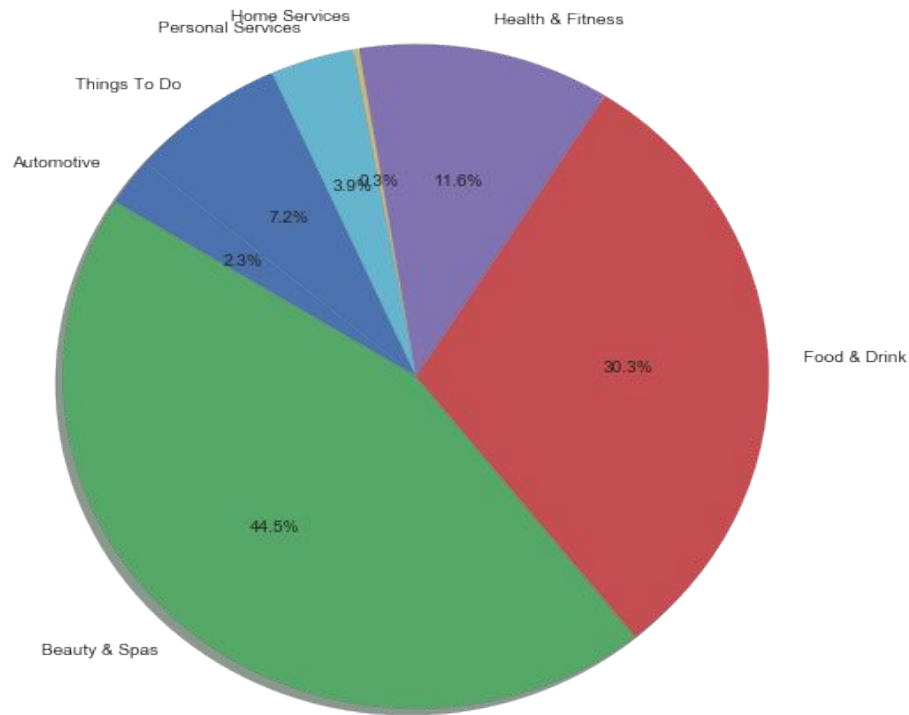  - Scraped every review and continued along the pages until the "end"
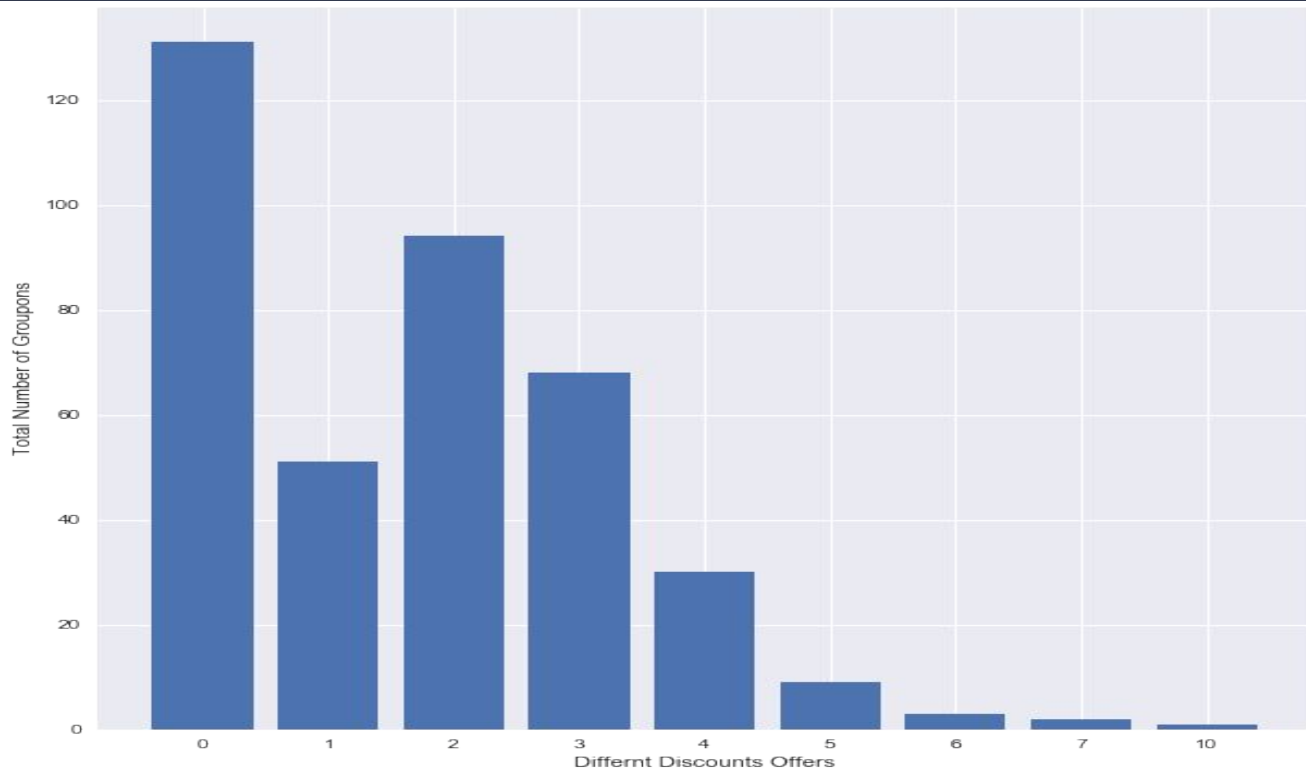
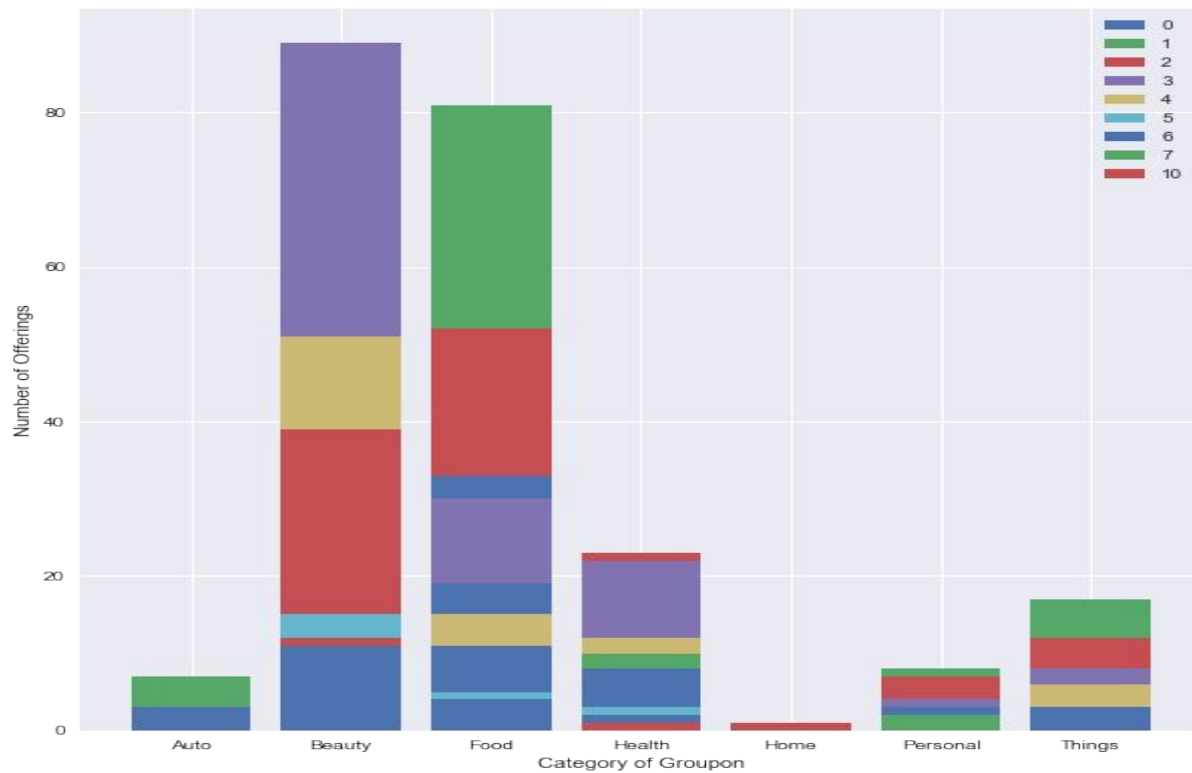# Data Analysis
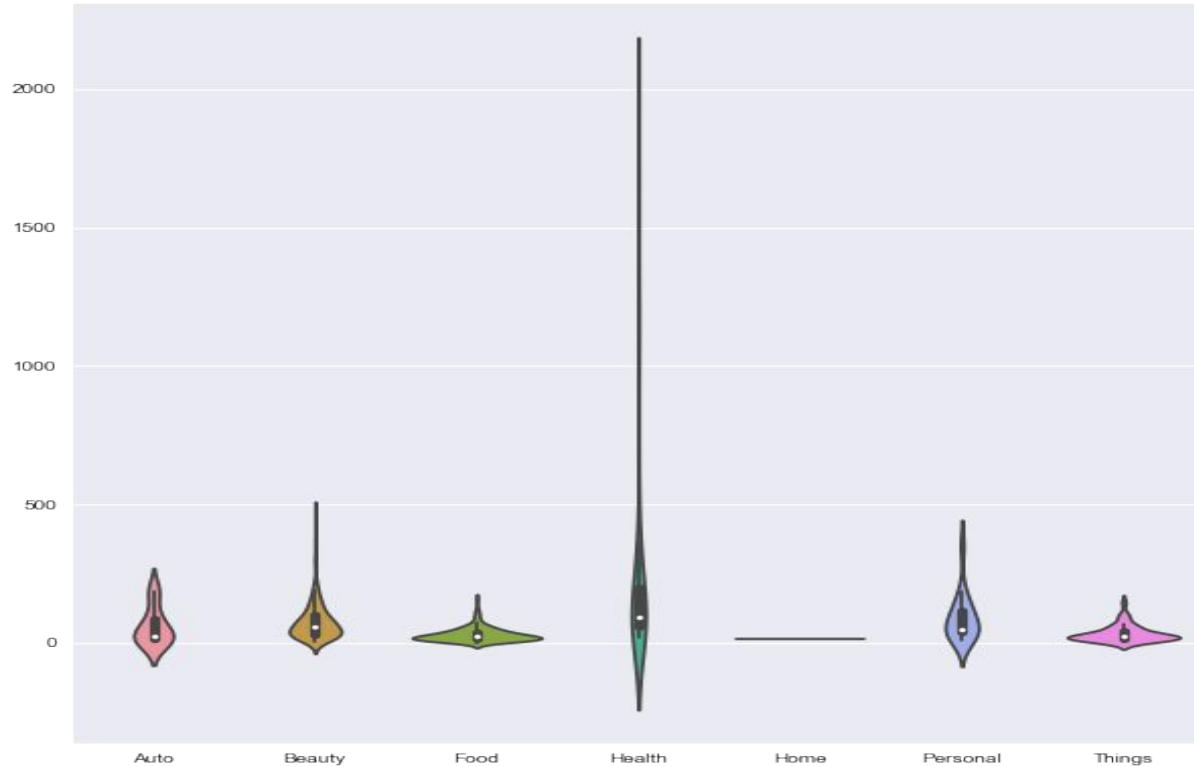
# More People are Using Groupon

# Distribution of Categories

# Groupon's Have Differing Deal Offers

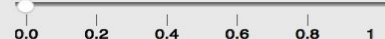# Offering Counts vs. Category

# Distribution of the Savings

Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

1

3

2

Marginal topic distribution

2%

5%

10%

Selected Topic: 1    Previous Topic    Next Topic    Clear Topic

Slide to adjust relevance metric:(2)
λ = 0

0.0   0.2   0.4   0.6   0.8   1

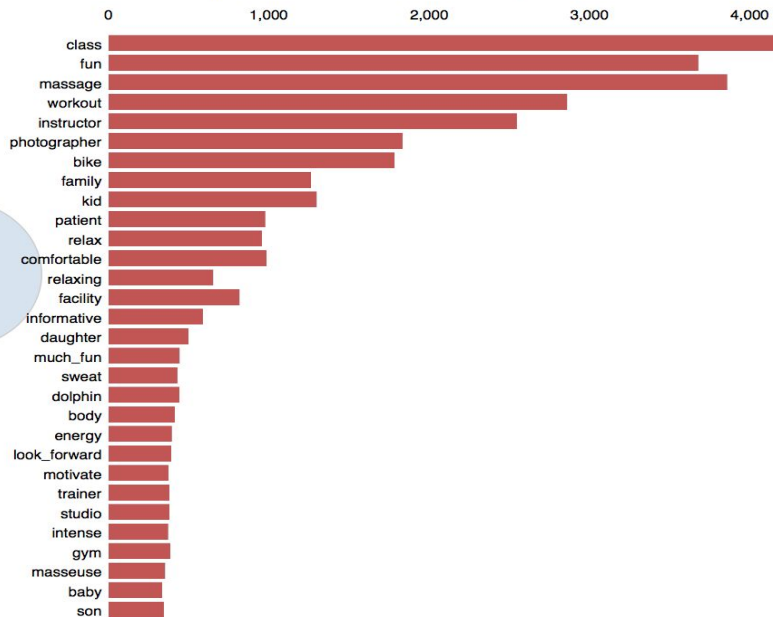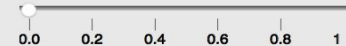Top-30 Most Relevant Terms for Topic 1 (48.6% of tokens)

0          1,000         2,000         3,000         4,000

use
bus
tell
pay
buy
tip
money
charge
hair
2
purchase
appointment
car
ticket
rude
lash
$
3
tax
arrive
line
long
horrible
voucher
extra
early
pick_up
terrible
phone
redeem

Overall term frequency
Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (2014)