

# What's there to a review?

...

Iris Borkovsky

# Action Plan

- Project scope
- Data summary
- Models
- Findings
- Appendix



# Project scope

- Can we predict if other users will find a particular review helpful?
- Topic modeling from reviews

Helpful

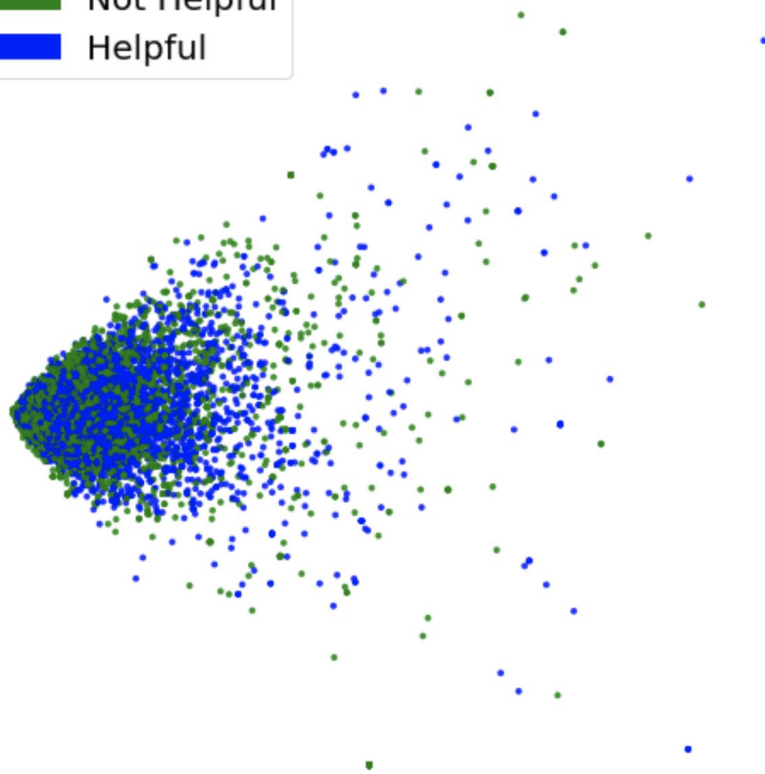
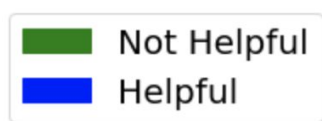
Not Helpful

[Comment](#)

[Report abuse](#)

# Data

- Working data set of ~20k, from ~500k Amazon food reviews
- Focus on helpfulness numerator and denominator to craft a classifier

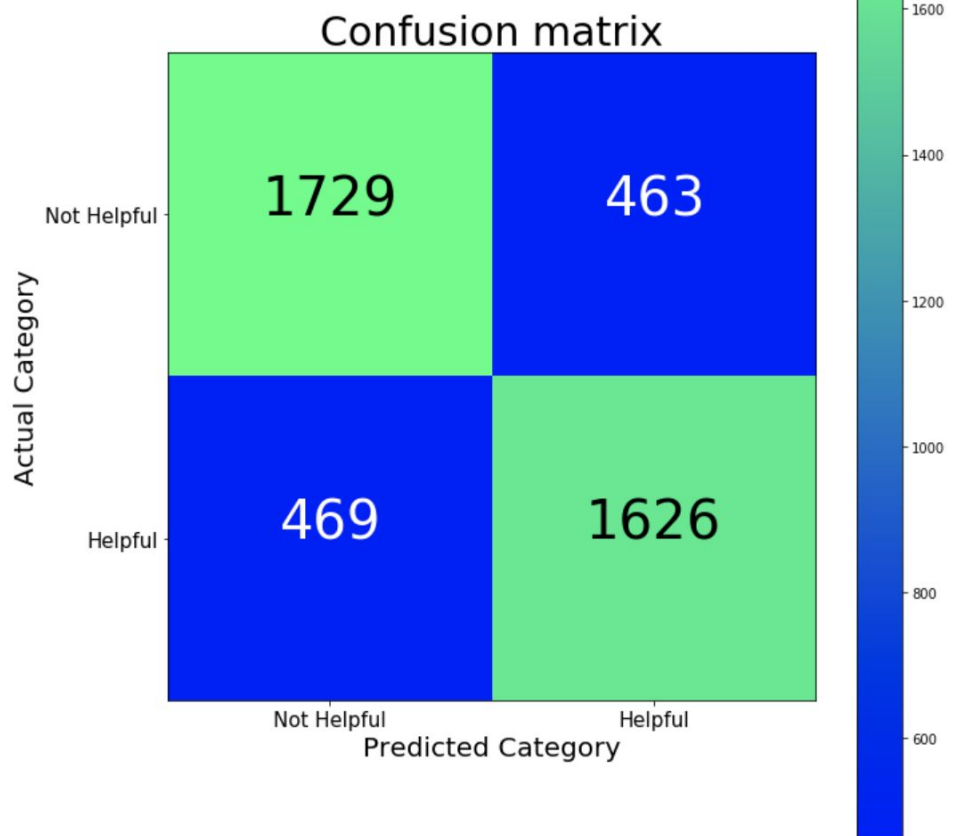


# Models

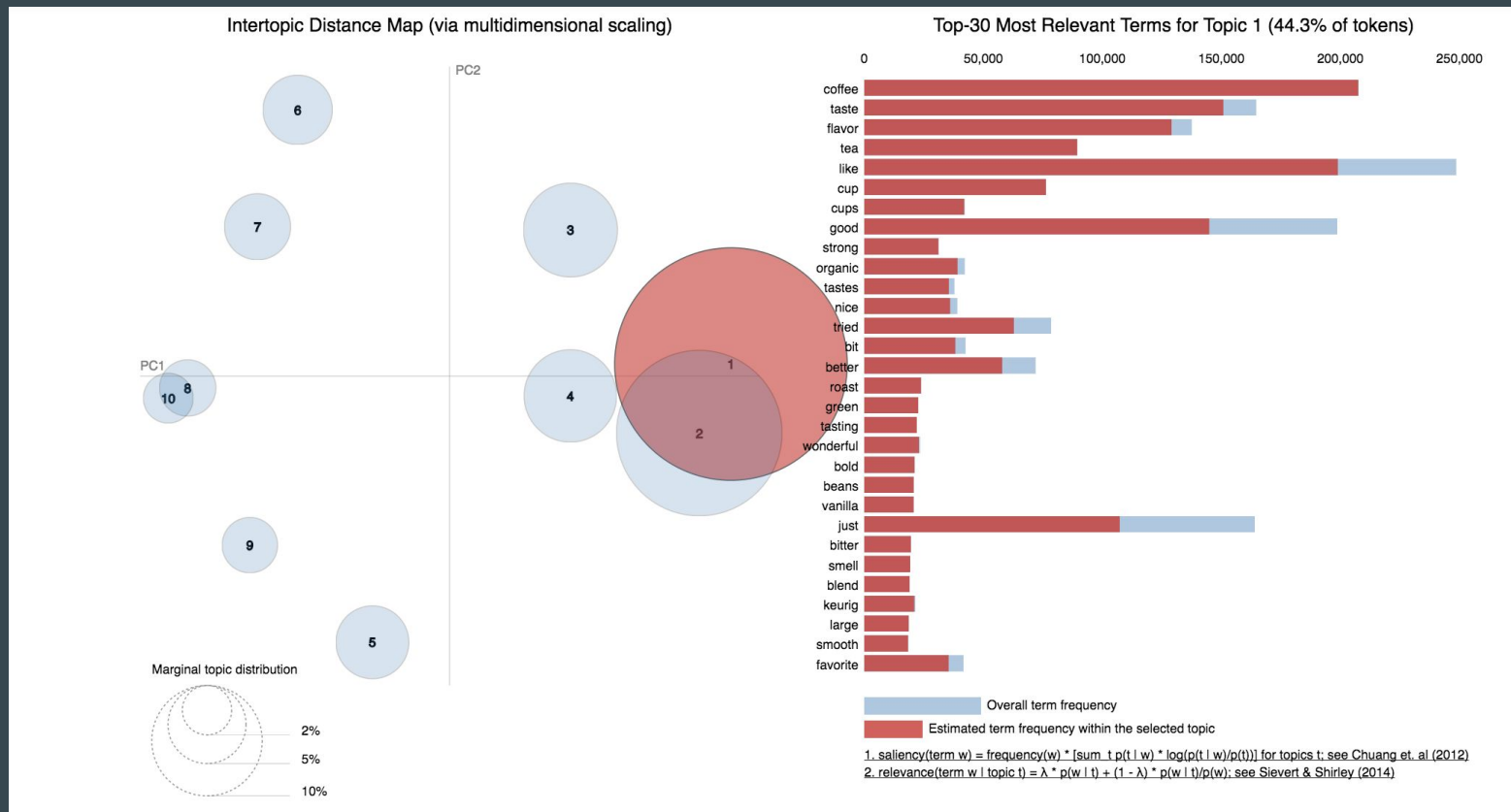
- Logistic Regression
- Latent Dirichlet Allocation
- K-means Clustering

# Findings

- Regression  
Accuracy score:  $\sim .78$



# Underlying categories in the data





# Next steps

- Semantic meaning of words (Word2Vec)
- Additional approaches to creating the original classifier (Helpful/ Not Helpful)
- How intent relates to helpfulness

# Thank you!



Iris Borkovsky

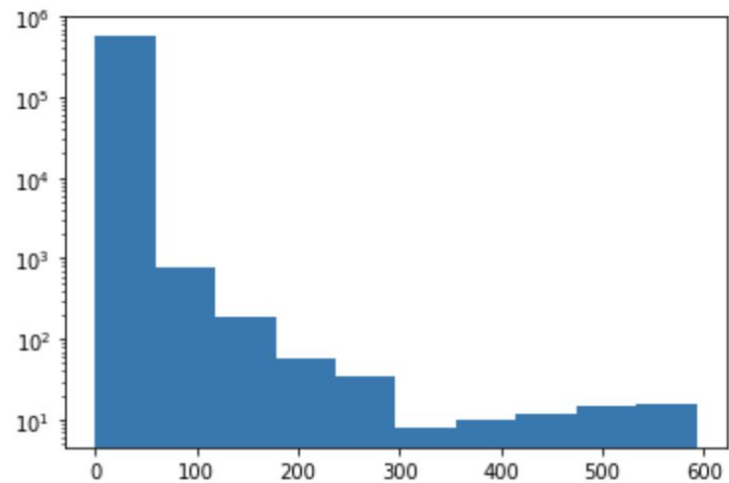
LinkedIn: <https://www.linkedin.com/in/iris-borkovsky>

GitHub: <https://github.com/borkovsky>

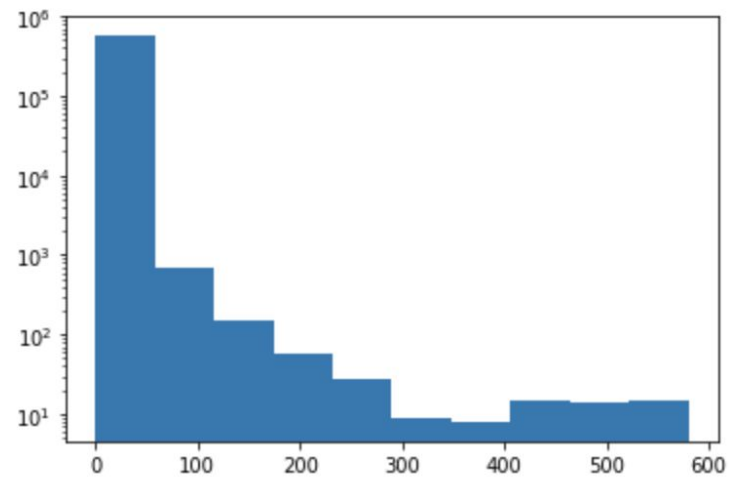
Blog: <https://medium.com/@borkovsky>

# Appendix

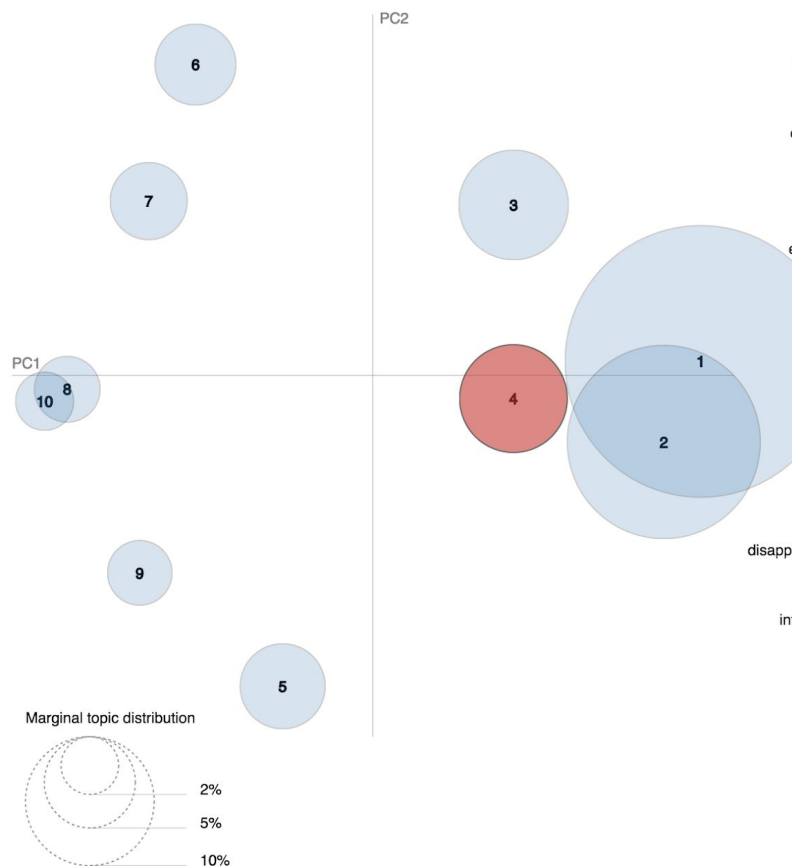
Helpfulness Denominator without outliers



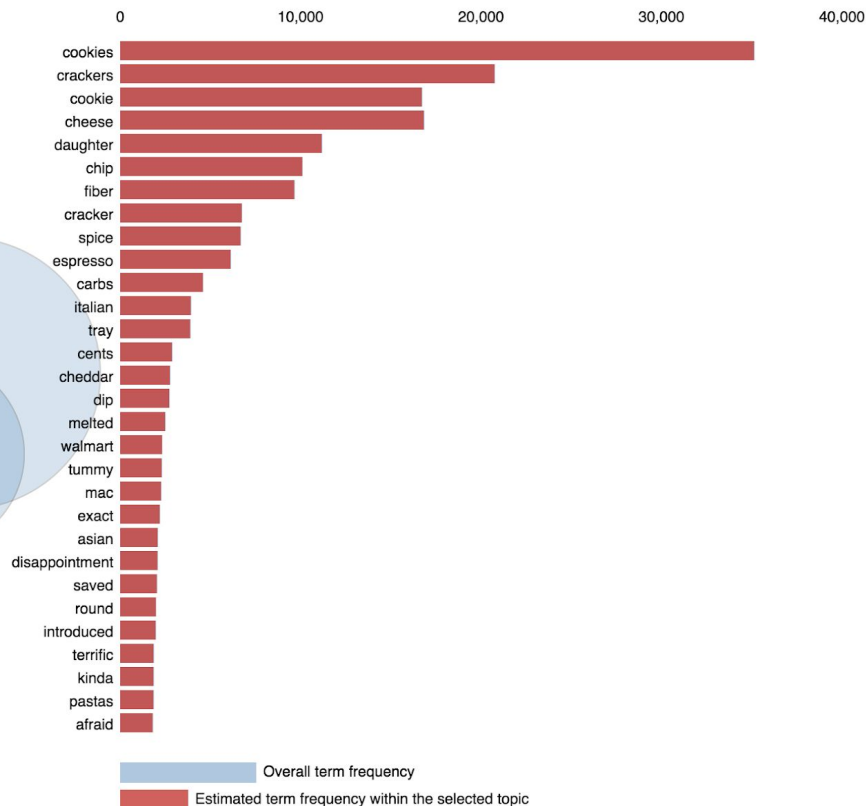
Helpfulness Numerator without outliers



Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 4 (7% of tokens)



1.  $saliency(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w) / p(t))]$  for topics  $t$ ; see Chuang et. al (2012)

2.  $relevance(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t) / p(w)$ ; see Sievert & Shirley (2014)

Selected Topic: 1

Previous Topic

Next Topic

Clear Topic

Slide to adjust relevance metric: (2)

 $\lambda = 0.58$ 

0.0

0.2

0.4

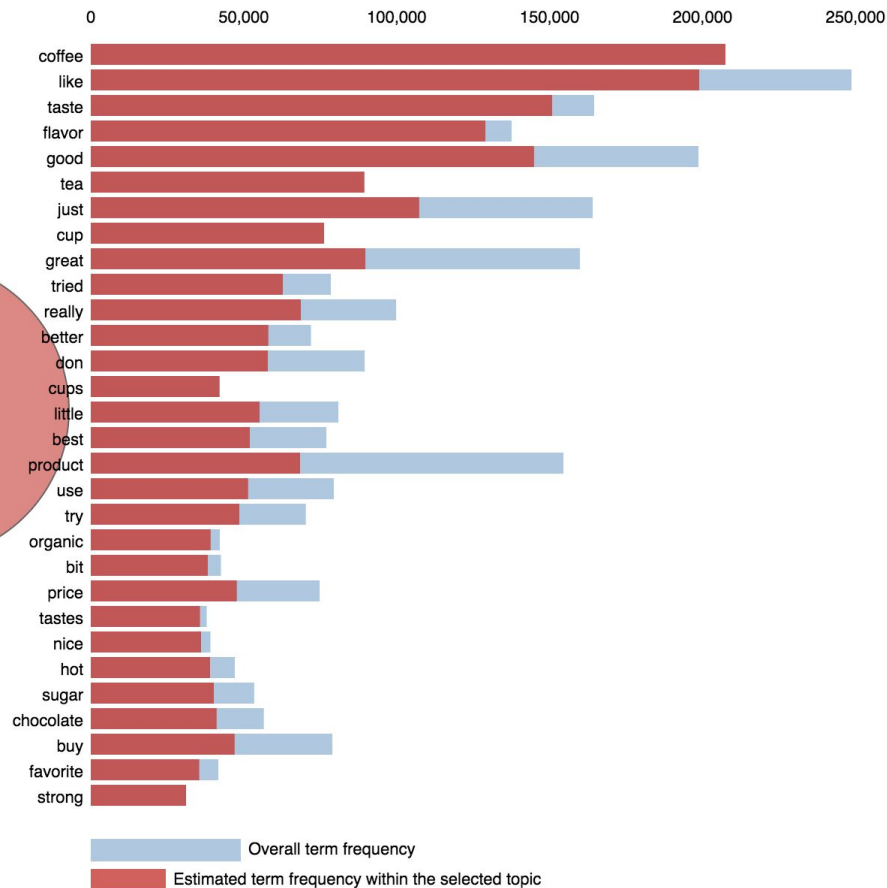
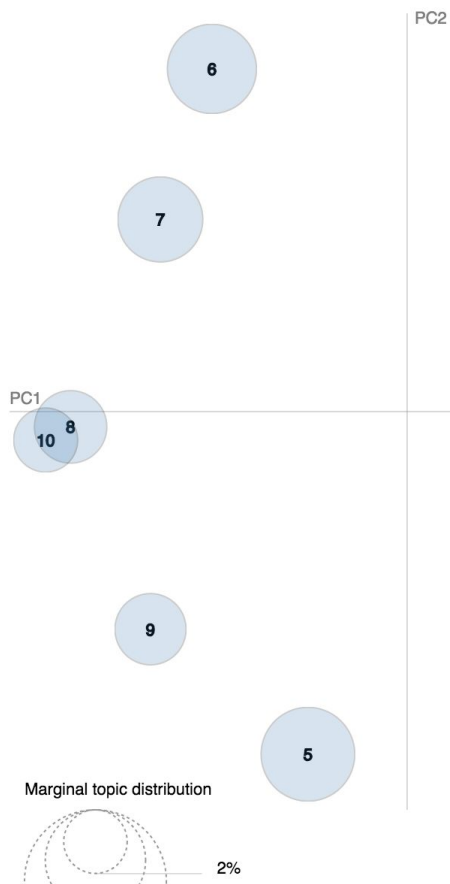
0.6

0.8

1

Intertopic Distance Map (via multidimensional scaling)

Top-30 Most Relevant Terms for Topic 1 (44.3% of tokens)



- There are 3,175,298 words total
- Vocabulary: 49,311