# Semantics of Kickstarter Projects Descriptions: Initial Results
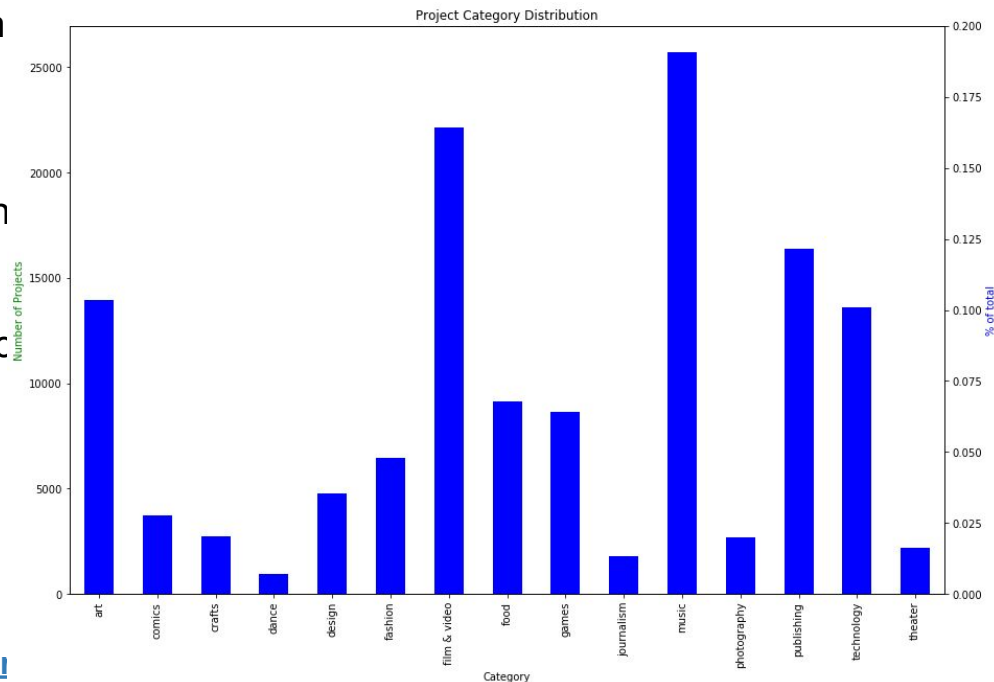
Mariya Petukhova
June 2017

# Objective

Conduct NLP analysis of project descriptions on Kickstarter.com to discover an underlying semantic structure

# Data

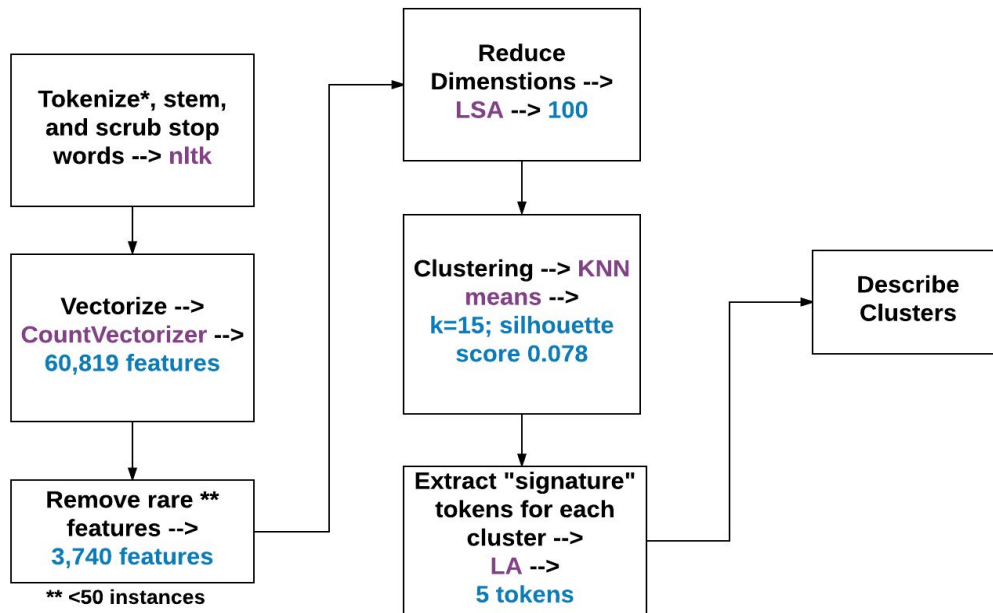Most recently scraped Kickstarter project data in JSON format available at
https://webrobots.io/kickstarter-datasets/ *

- ❏ 134,780 projects originated in US and plan for US market
  - ❏ Sampled down to 20,000 for
- ❏ 22 features, including text blurb with projec description
- ❏ Almost 1GB data

**\* Actual file name:**
**https://s3.amazonaws.com/weruns/forfun/Kickstarte**
**starter_2017-05-15T22_21_11_300Z.json.gz**



Project Category Distribution

3

# Methodology of Project Blurb Analysis

# Traditional Obscure Cultural Reference

"**William Shakespeare's vocabulary has been estimated by the experts at <span style="color:red">twelve thousand*</span> words. The vocabulary of the Mumbo Jumbo tribe amounts to three hundred words.**
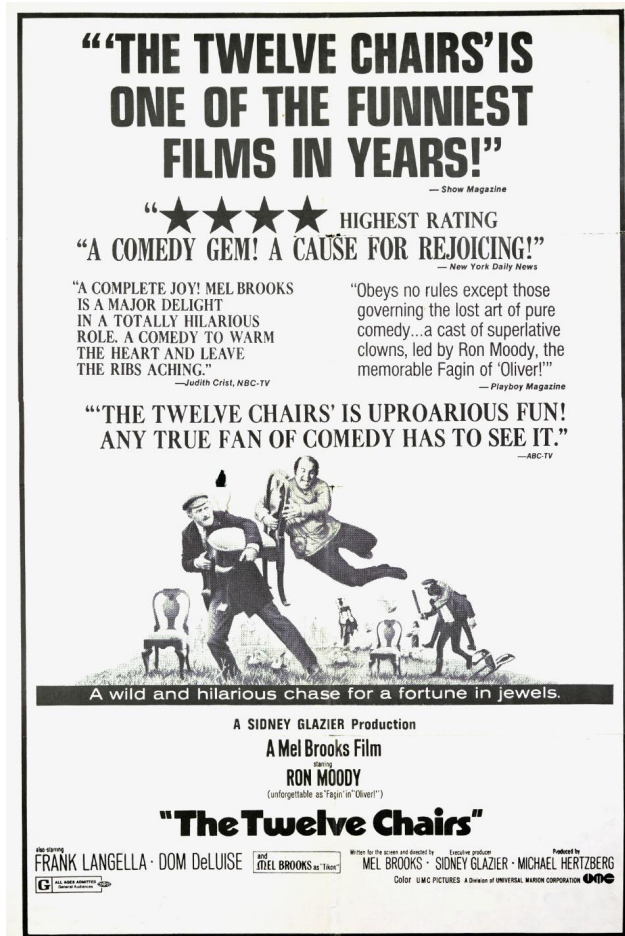
**Ellochka Schukin managed easily and fluently on thirty.**

Here are the words, phrases and interjections which she fastidiously picked from the great, rich and expressive Russian language:
1. You're being vulgar.
2. Ho-ho!
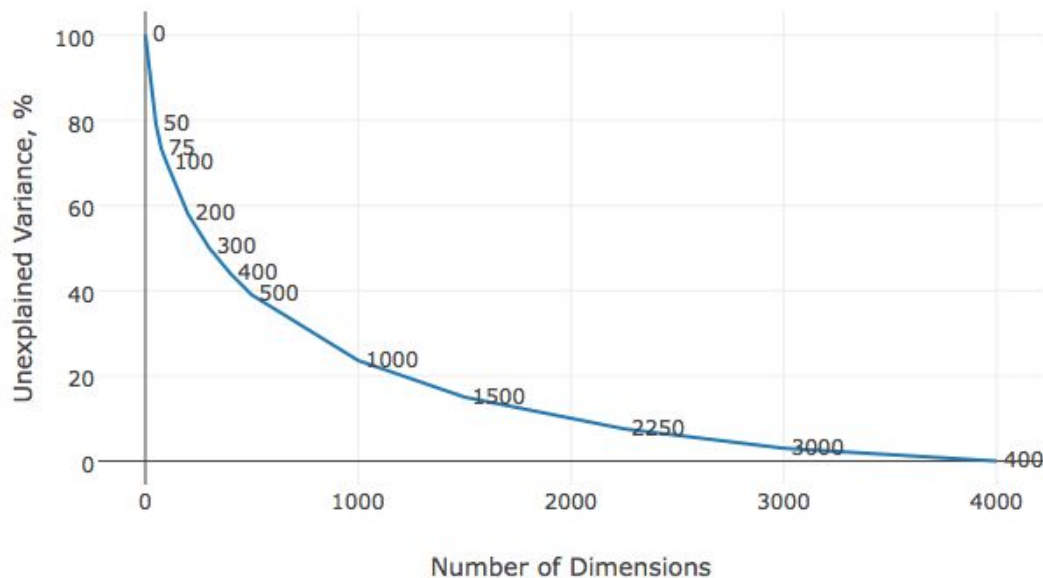[…]"

Ilf and Petrov
"The Twelve Chairs

\* This makes 60,819 tokens obtained from vectorization extremely high
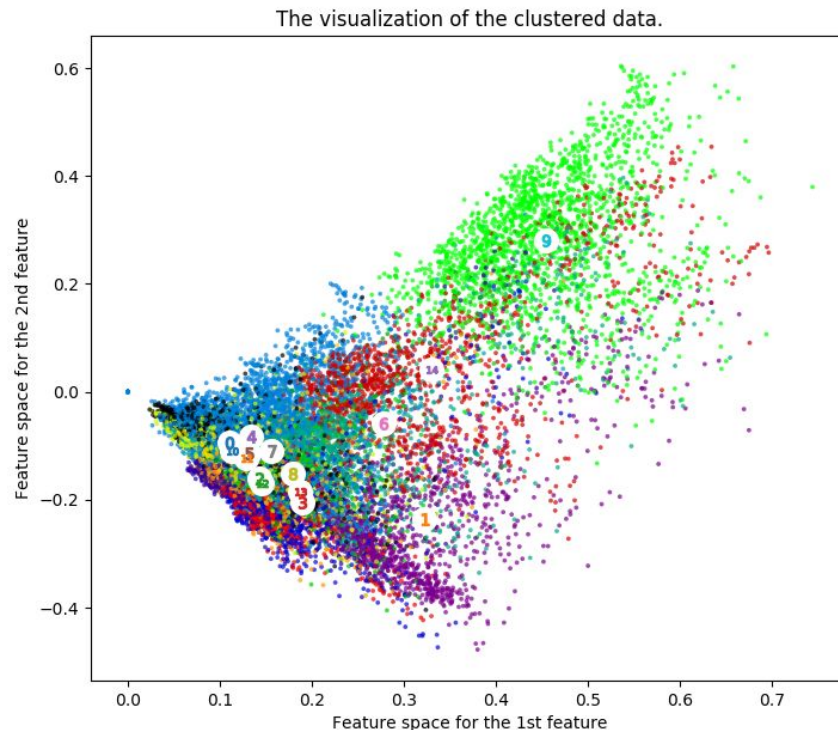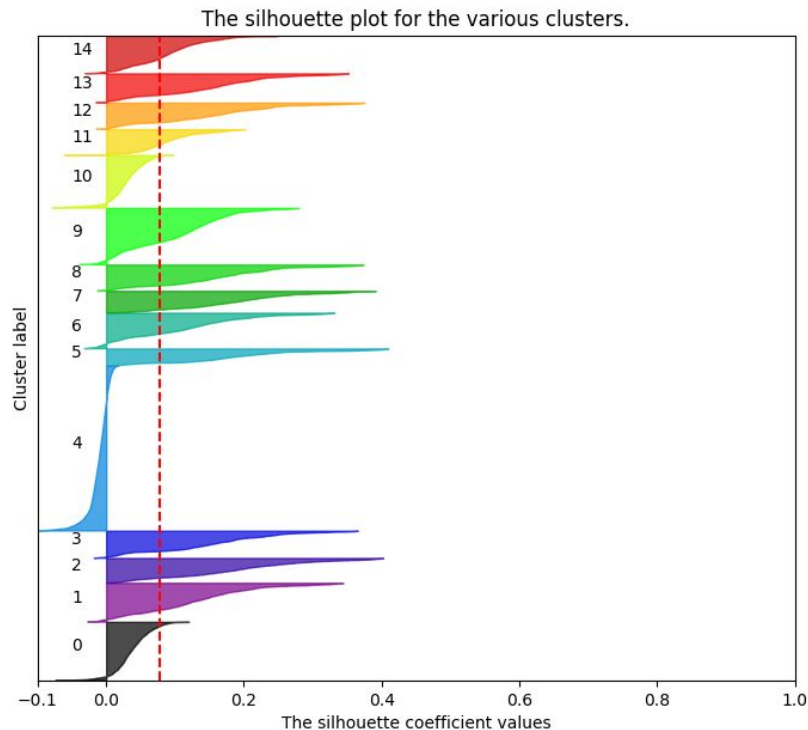
# Dimension Reduction

- ❏ 1,500 dimensions explain 85% of variation in the data
  - ❏ Curse of Dimensionality in action - very low silhouette scores
  - ❏
  - ❏
- ❏ Proceeded with 100 dimensions and 30% explainability

## Number of Dimensions vs. Unexplained Variance in the Data



Number of Dimensions (x-axis), Unexplained Variance, % (y-axis). Data points labeled: 0, 50, 75, 100, 200, 300, 400, 500, 1000, 1500, 2250, 3000, 400
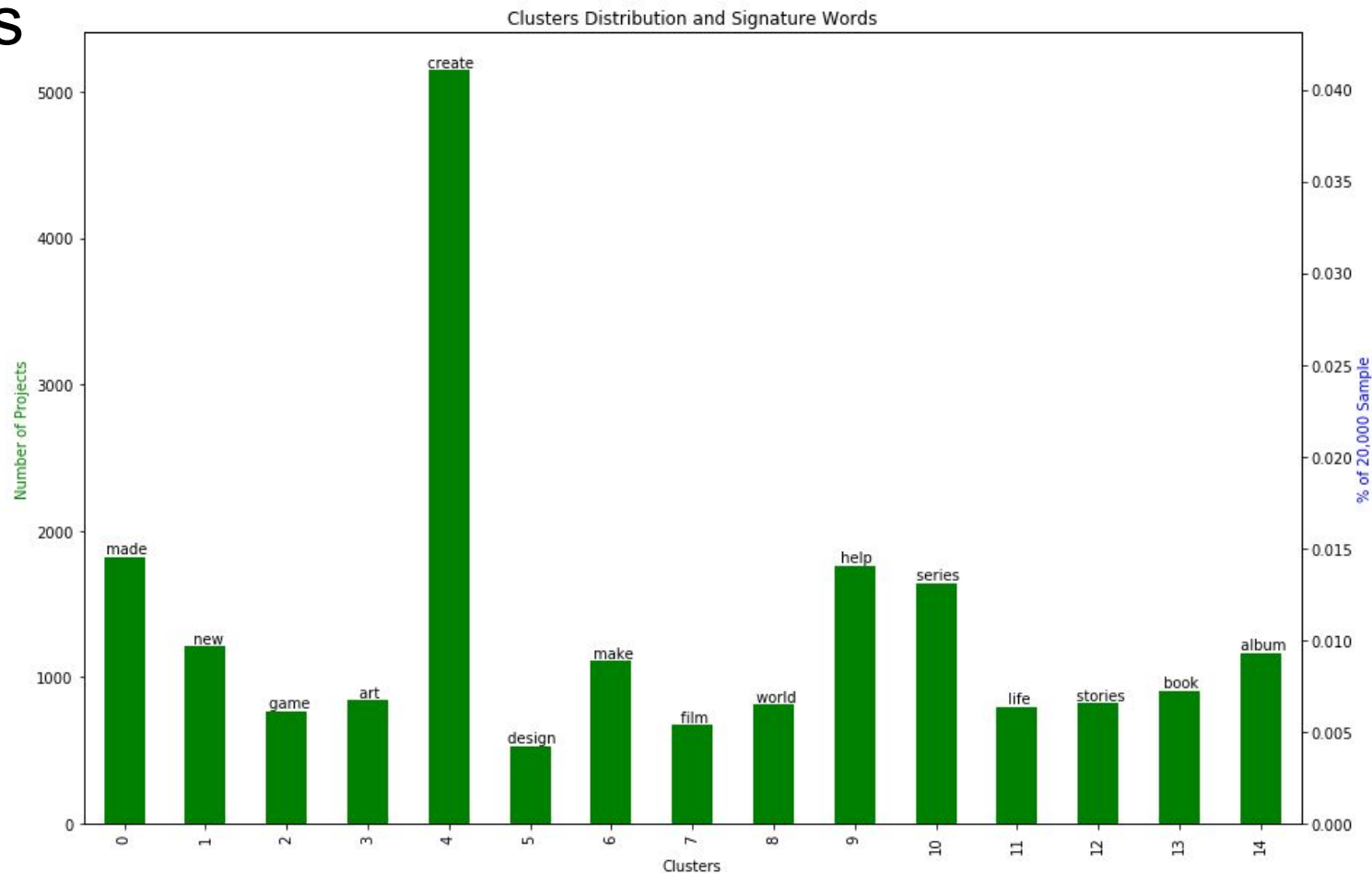
6

# Clustering: KNN Means with k=15

**Silhouette analysis for KMeans clustering on sample data with n_clusters = 15**



The average silhouette_score is : 0.078

# Clusters



Clusters Distribution and Signature Words

# The Most Separated (#9) and The Largest (#4) Segments

# Clusters Highlights

| Cluster ID | Category | Top % in Category | Signature Word | Fun Facts |
|:---:|:---:|:---:|:---:|:---:|
| 2 | Games | 77% | 'game' | The highest pledge, goal, staff picked, and backers numbers |
| 14 | Music | 96% | 'album' and 'record' | The highest success rate (65%) |
| 7 | Film & Video | 89% | 'film' | The 2nd most successful (60%) |
| 13 | Publishing | 68% | 'book' | The second lowest % live (1.05) |
| 2 | Technology Food | 30% 25% | 'made' & 'food' | The lowest success rate (31%) |

# Conclusion

- ❏ Only 5 out of 15 clusters closely follow project's Category, thus, more analysis needed to better understand the semantics behind groupings
- ❏ The standalone cluster 9 seems to be just asking for money (top 5 words are 'help', 'need','us', 'fund', 'get'
- ❏ Cluster 4 is a catch-it-all - the top words are generic, the size is the largest
- ❏ Game is king - the highest pledge and goal amounts, the highest number of staff-picked projects, and the highest backers count
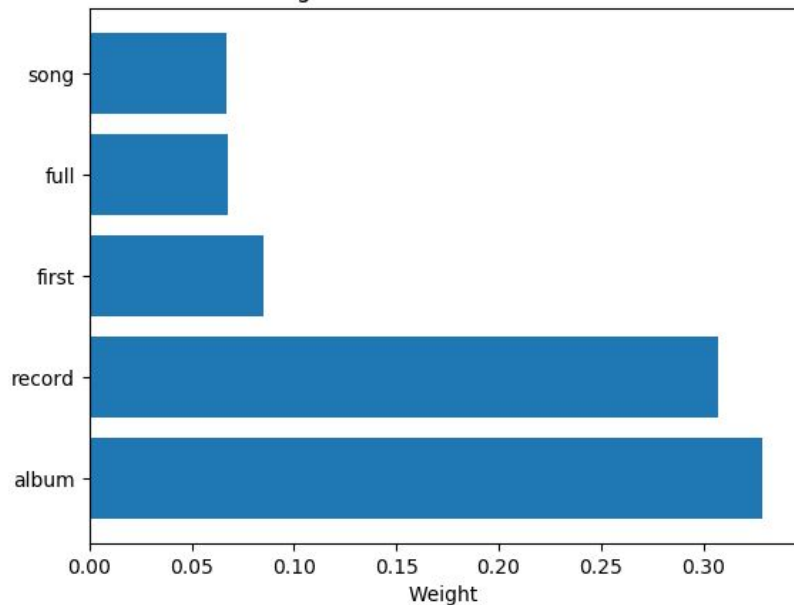
# Next Steps

- ❏ Use all 130K + observations
- ❏ Try other clustering techniques
- ❏ Build supervised classification model to see if new clusters or specific signature tokens can predict success of campaign
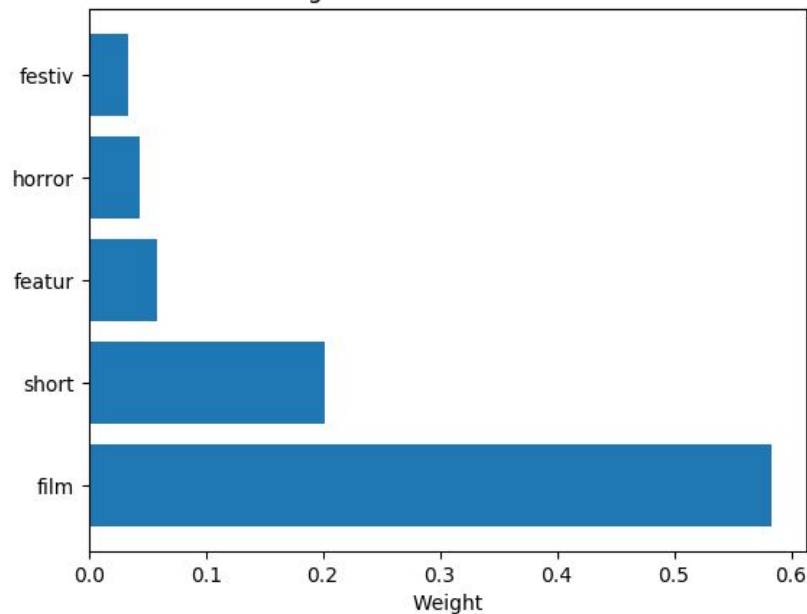
# Appendix

# Signature Words: Most Successful Clusters

# Signature Words: The Highest Pledge and Goal Amounts, Most supported by Staff, and Max Number of Backers Cluster



Signature Words: Cluster 2