

# Kickstarter

## RWHD Final Project

Flower Yang, Qingwei Zhang, Keturah Jones  
04/30/2019



# Kickstarter Datasets

We have a scraper robot which crawls all [Kickstarter](#) projects and collects data in CSV and JSON formats. From March 2016 we run this data crawl once a month. Datasets are available from the following scrape dates:

2019

- 2019-03-14 [[JSON](#)] - [[CSV](#)]
- 2019-02-14 [[JSON](#)] - [[CSV](#)]
- 2019-01-17 [[JSON](#)] - [[CSV](#)]

2018

- 2018-12-13 [[JSON](#)] - [[CSV](#)]
- 2018-11-15 [[JSON](#)] - [[CSV](#)]
- 2018-10-18 [[JSON](#)] - [[CSV](#)]
- 2018-09-13 [[JSON](#)] - [[CSV](#)]
- 2018-08-16 [[JSON](#)] - [[CSV](#)]
- 2018-07-12 [[JSON](#)] - [[CSV](#)]
- 2018-06-14 [[JSON](#)] - [[CSV](#)]
- 2018-05-17 [[JSON](#)] - [[CSV](#)]
- 2018-04-12 [[JSON](#)] - [[CSV](#)]
- 2018-03-15 [[JSON](#)] - [[CSV](#)]
- 2018-02-15 [[JSON](#)] - [[CSV](#)]
- 2018-01-12 [[JSON](#)] - [[CSV](#)]

2017

Kickstarter - Excel

ome Insert Draw Page Layout Formulas Data Review View Help Tell me

Calibri 11

B I U A A

Font Alignment Styles

Conditional Formatting Format as Table Cell Styles

Cells Editing Ideas

backers\_count

	B	C	D	E	F	G	H	I	J	K	L
	_c blurb	category	converted	country	created_a	creator	currency	currency_	currency_	current_c	deadlin
1	Monsters,	["urls":{"v	20	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
37	Nano Art v	["urls":{"v	1974	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
81	Video and	["urls":{"v	4845	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
95	Finally, A	["urls":{"v	2948	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
10	If you like	["urls":{"v	522	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
0	WELCOME	["urls":{"v	0	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
0	Standing c	["urls":{"v	0	US	1.34E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
46	Wearable	["urls":{"v	1701	US	1.34E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
18	I will be d	["urls":{"v	350	US	1.34E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
3	I will be p	["urls":{"v	31	US	1.34E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
0	im an LA b	["urls":{"v	0	US	1.34E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
26	I am makin	["urls":{"v	2507	US	1.34E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
0	A unique	["urls":{"v	0	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
0	A portable	["urls":{"v	0	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
1	Limited Ec	["urls":{"v	25	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
1	Bringing fi	["urls":{"v	5	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
2	Digital Pai	["urls":{"v	120	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
0	An except	["urls":{"v	0	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
0	Springzeâ	["urls":{"v	0	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
0	1000's of r	["urls":{"v	0	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
30	We're brir	["urls":{"v	1627	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
3	Using con:	["urls":{"v	131	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
4	Classics fr	["urls":{"v	46	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+
4	The Royal	["urls":{"v	185	US	1.33E+09	["urls":{"v	USD	\$	TRUE	USD	1.34E+
0	Hello eve	["urls":{"v	0	US	1.32E+09	["urls":{"v	USD	\$	TRUE	USD	1.32E+
0	Turn thos	["urls":{"v	0	US	1.32E+09	["urls":{"v	USD	\$	TRUE	USD	1.33E+

Kickstarter



# Data Description

Data from 2018 is used in this project. Each month contains approximately 48 to 50 csv files.

## Examples of Variables:

1. ID : project ID, consist of numbers
2. Name: Project name
3. Blurb: Project Description
4. Goal: amount of the money goal to succeed the project
5. USD\_Pledged : Actual amount in USD reached from the goal
6. State: success, failed, canceled, or live
7. Country & location
8. Deadline: Time of the deadline
9. Launching time
10. Backers: Number of supporters
11. Category: 15 categories of project



# Data Pre-processing

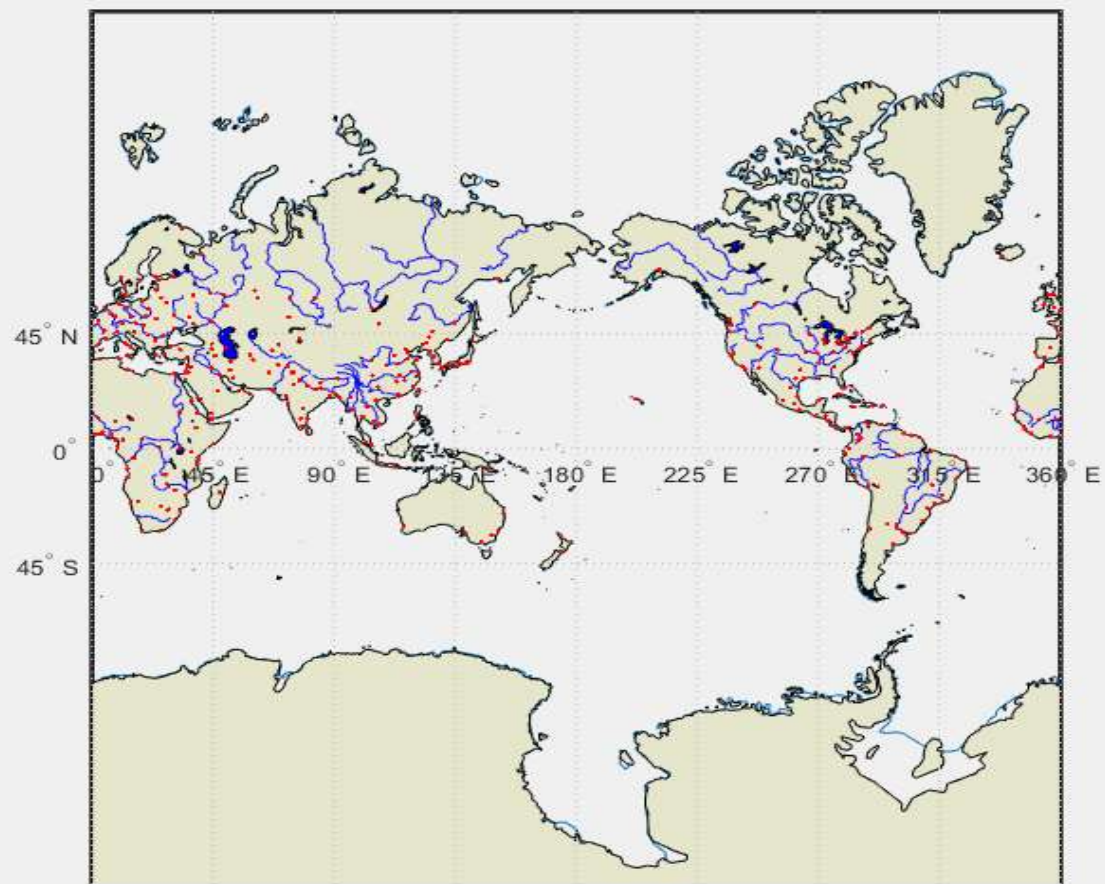
37 columns (variables) in original dataset, of which we keep for analysis: backers count; blurb; category; country; deadline; goal; id; launched at; location; name; spotlight; staff pick; state; state changed at; usd pledged

Adding numeric labels (1-15) to represent the 15 categories of projects: art; comics; crafts; dance; design; fashion; film & video; food; games; journalism; music; photography; publishing; technology; theater

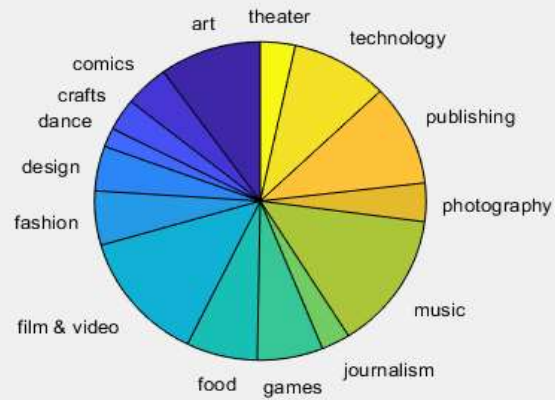
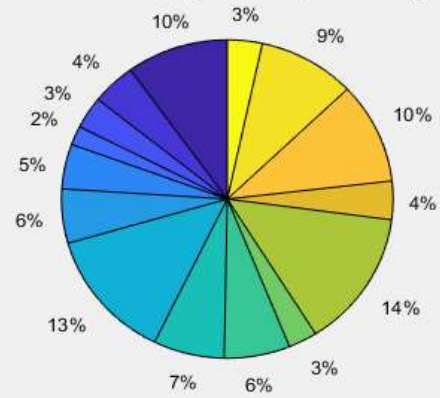


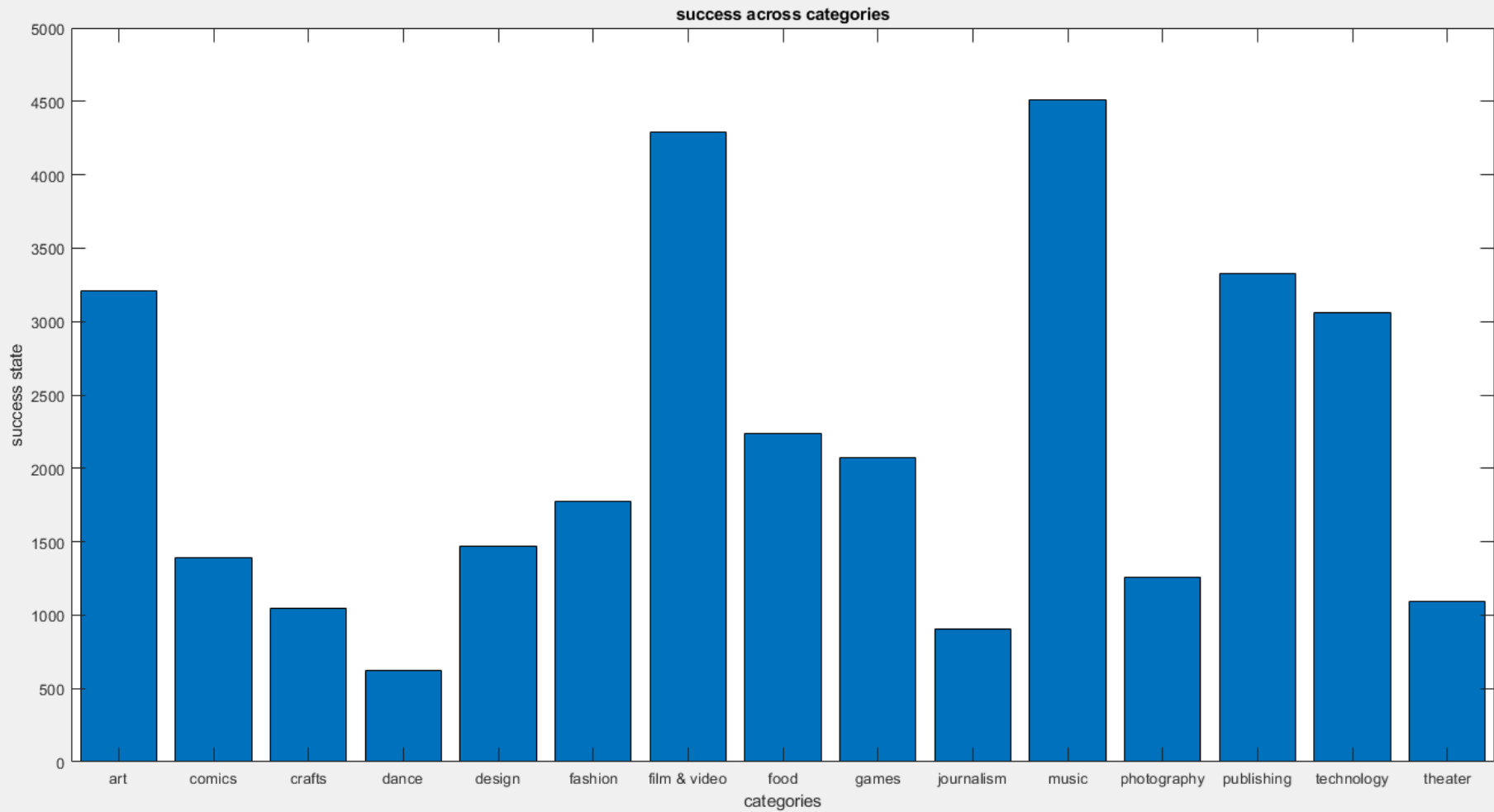
# Descriptive Statistics

- 1) Geo analysis -- The whole world
- 2) Histogram of success across categories (defining success using general ('state'))
- 3) relationship between "spotlighted"/"staff picked" and being successful
- 4) relationship between goal \$\$ and being successful
- 5) the number of projects in each category for every month of every year (a line plot with time on the x-axis).



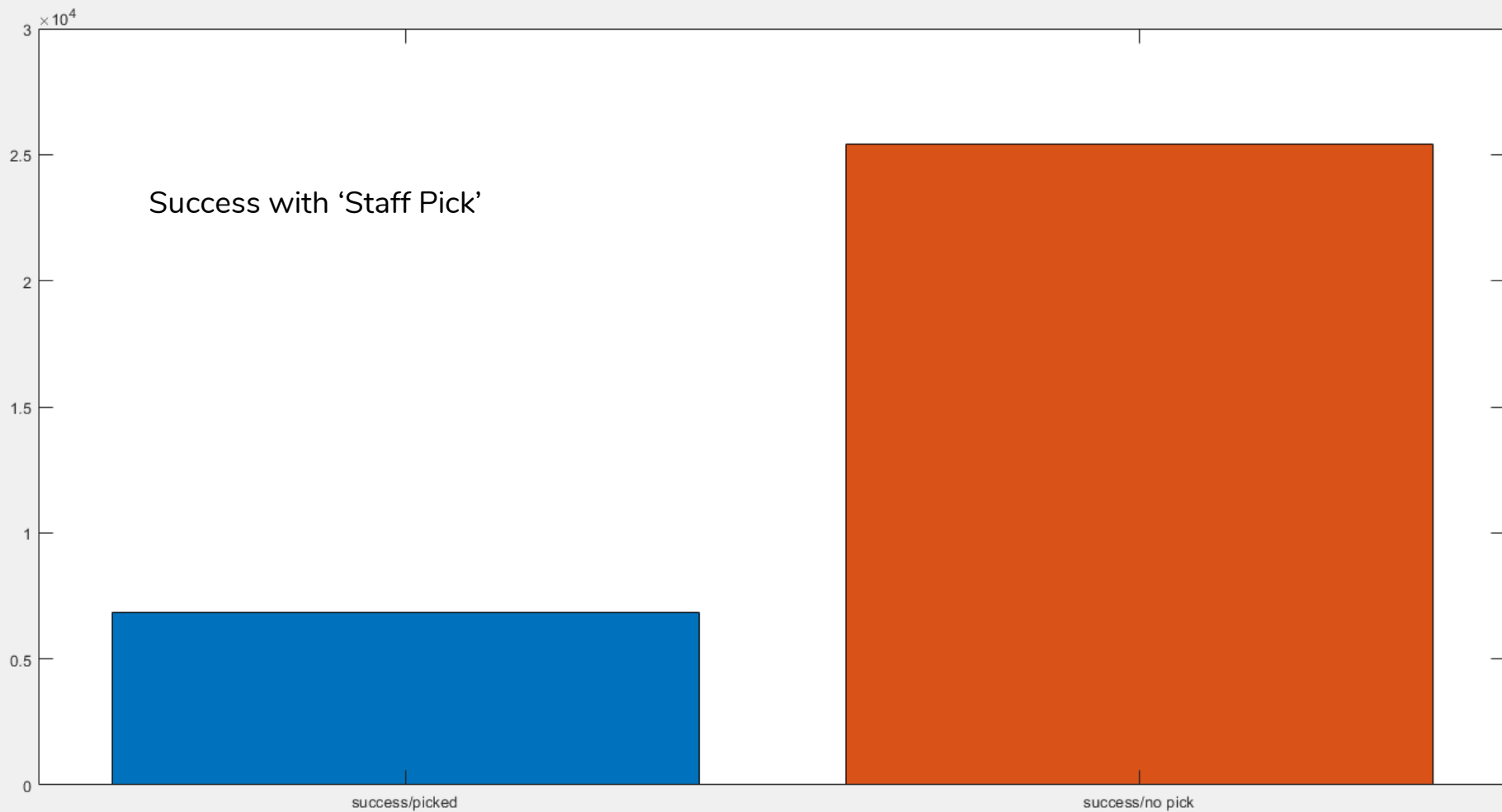
success across categories in percentage

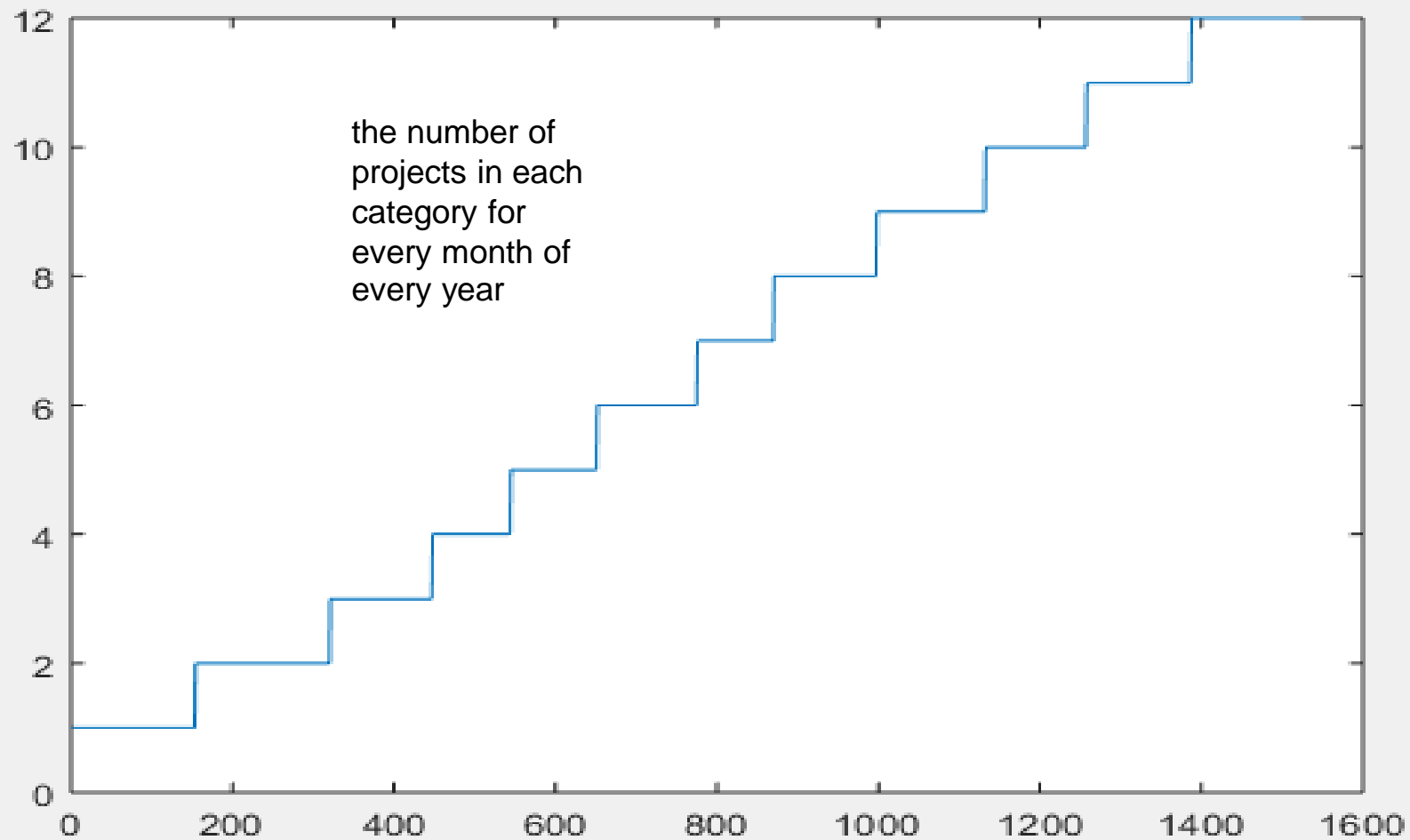






Success with 'Staff Pick'



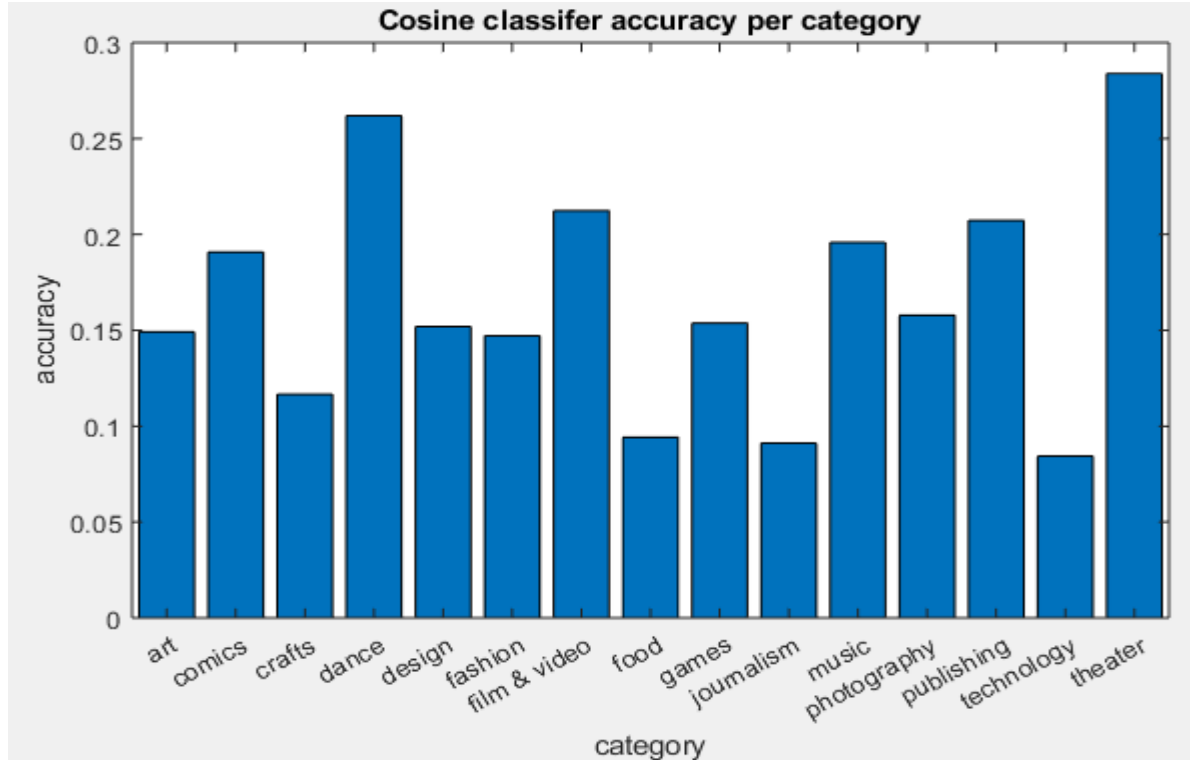




- 1) failed/cancelled projects
- 2) successful projects
- 3) wildly successful projects:  
exceed expected outcome  
by at least 20%



Classification using Cosine Classifier: generate the average word embedding vector for each class. Use 11 months of data as input (training data) to train a classifier, which is then used to predict the class of the remaining month.

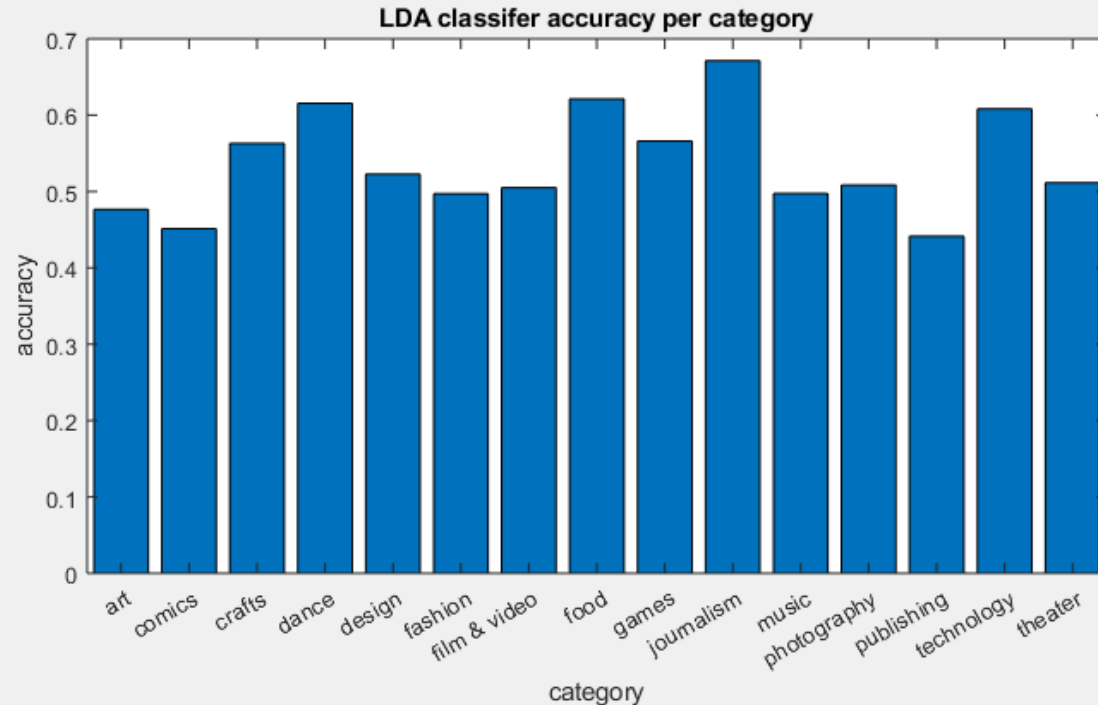


Overall mean accuracy  
across categories:

0.1883



Classification using LDA classifier: use 11 months of data as training data and the last 1 month as test data.



Overall mean accuracy  
across categories:

0.48



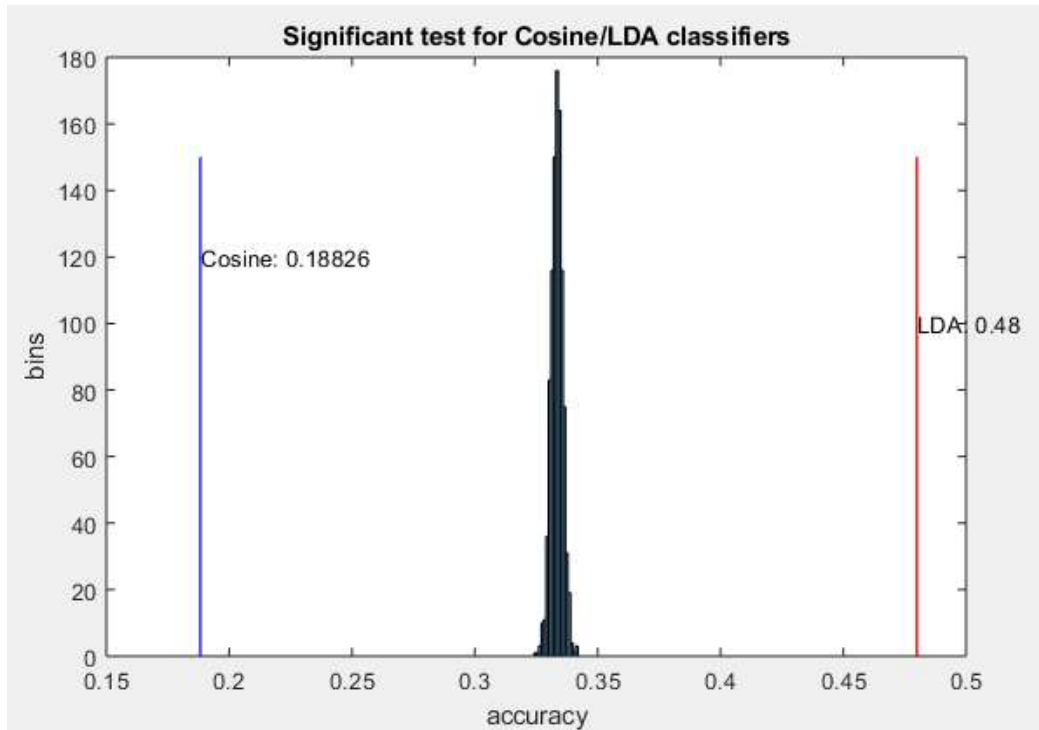
## Significance Test

Cosine classifier:

- Mean accuracy: 0.188
- P-value: 1

LDA classifier:

- Mean accuracy: 0.48
- P-value:  $8.5822e-15$





## **Detailed Analysis 2:**

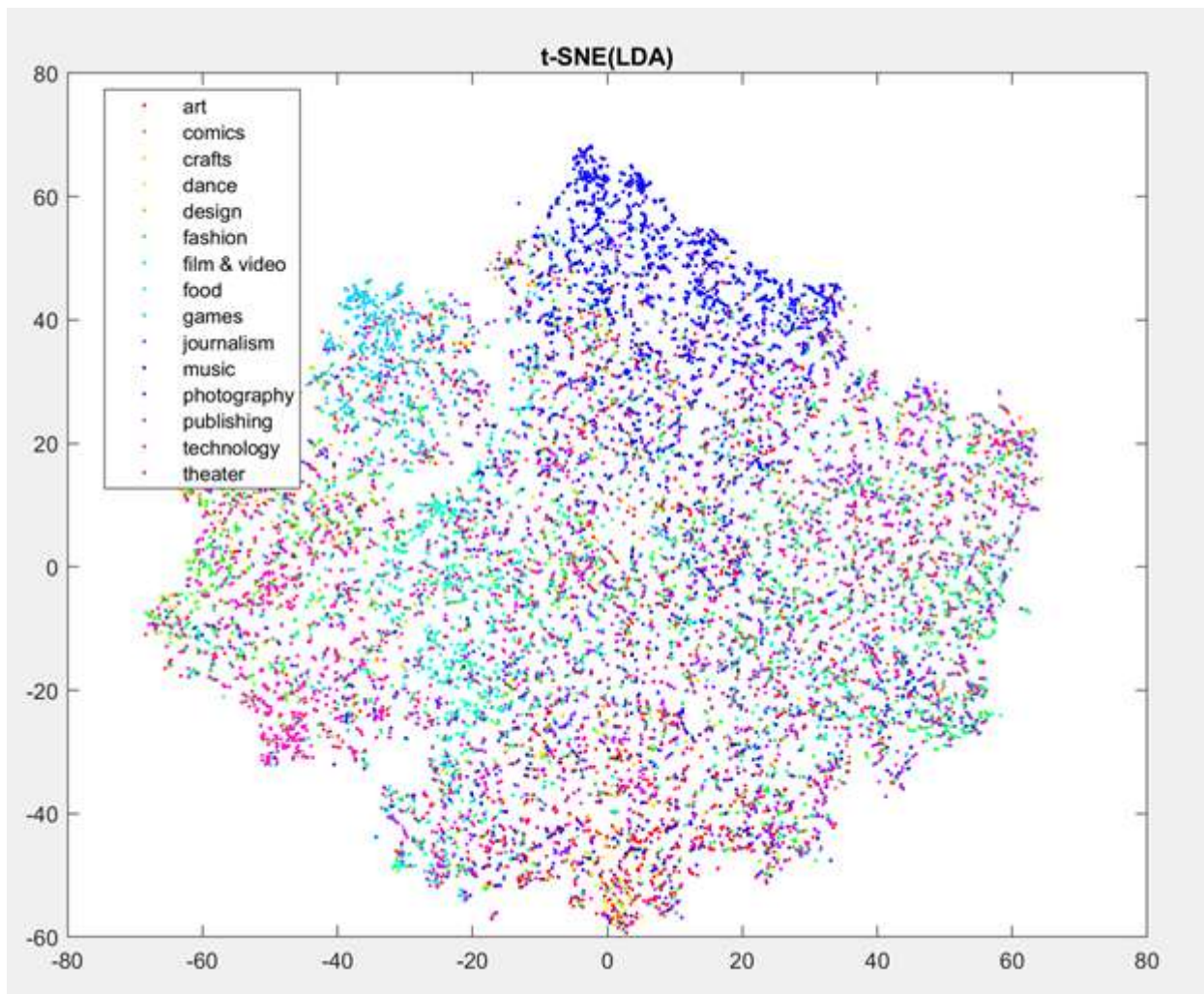
### **How well can t-SNE, LDA and word2vec recover the categories from the blurbs?**

- Trained an LDA model on the blurbs, using the same number of topics as there are categories (fifteen overall).
- For word2vec, calculated the average word embedding vector for each blurb
- Used a subset of the dataset to train LDA model and word2vec because full dataset was too large



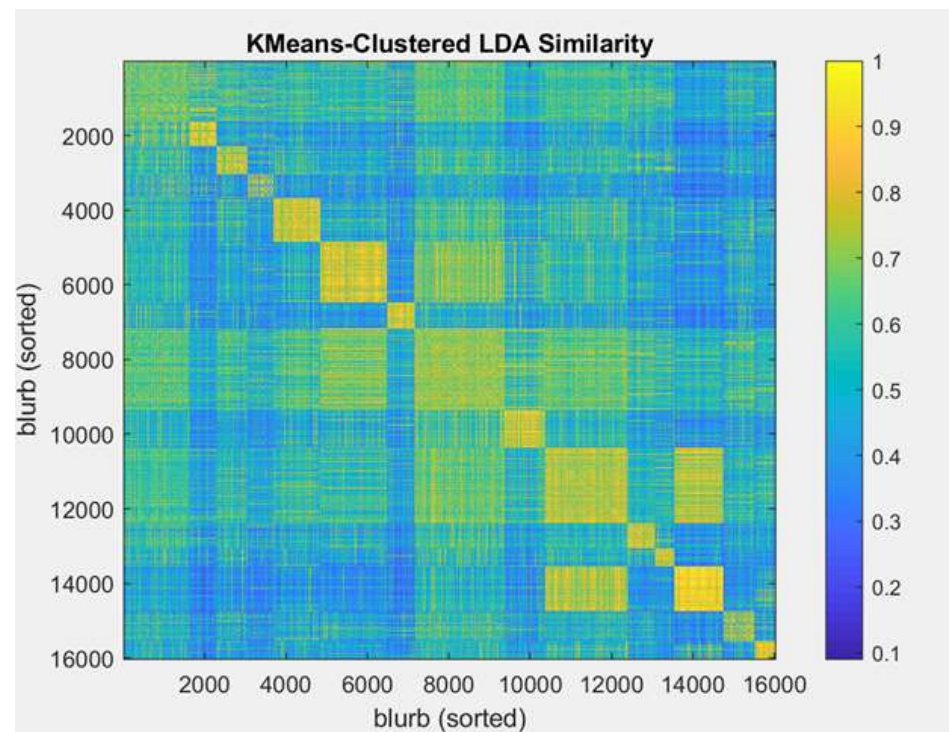
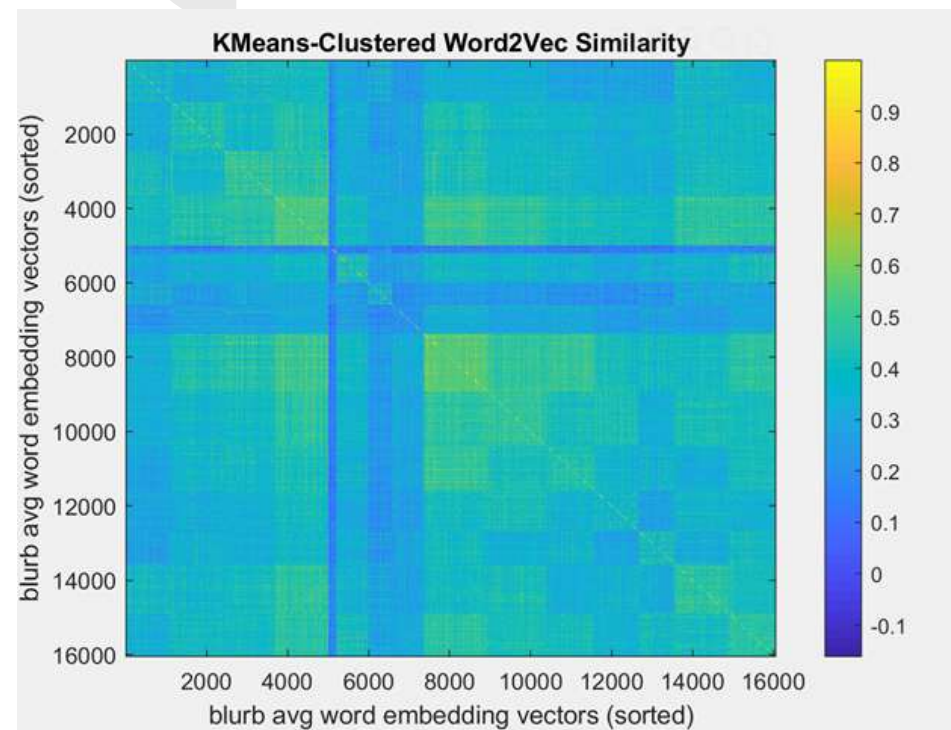
# t-SNE

Produced t-SNE to  
reduce  
dimensionality of  
the blurb x topic  
mixture matrix of  
LDA





# KMeans clustering

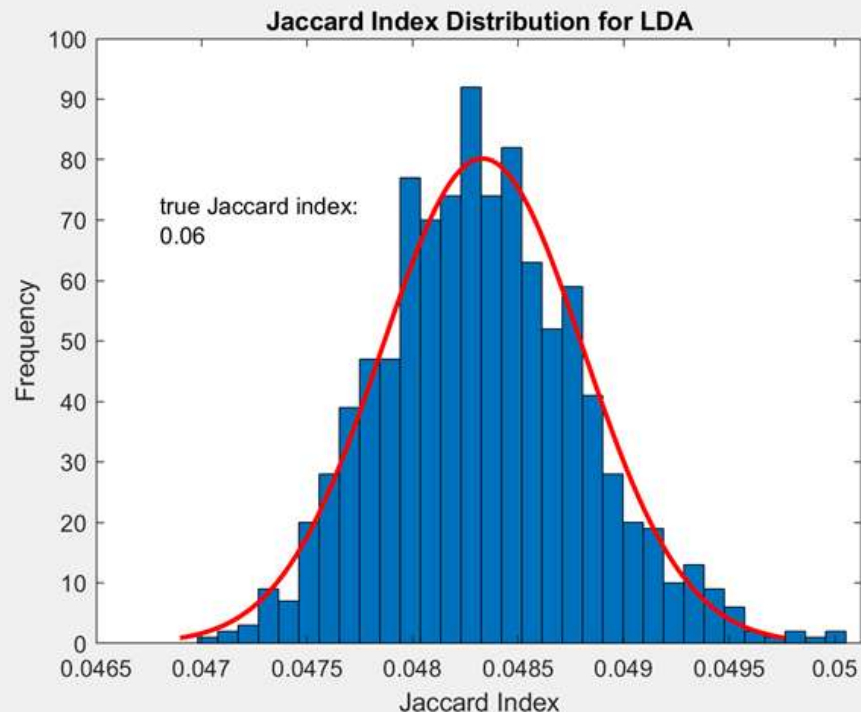
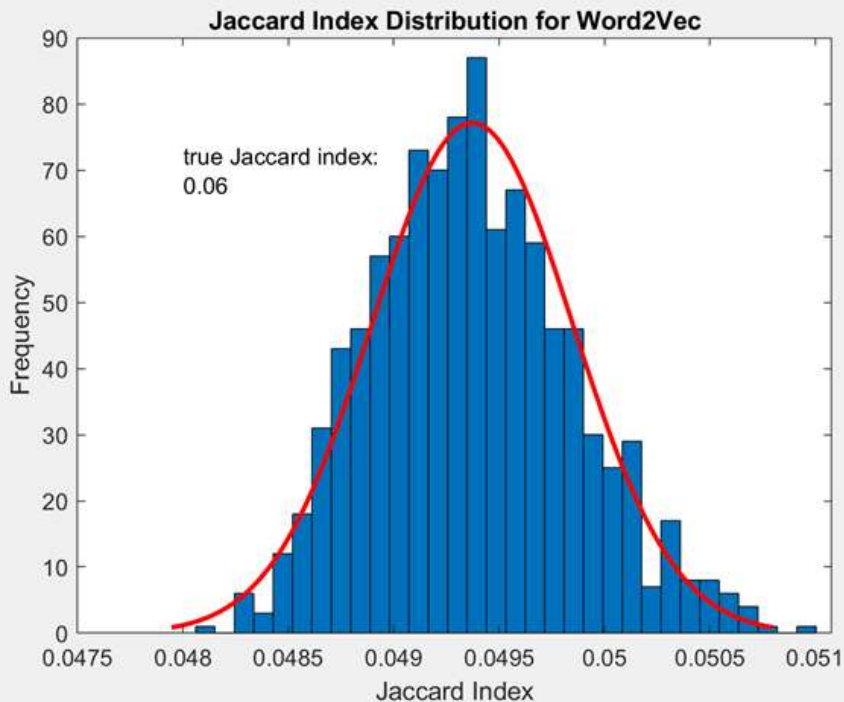




# Jaccard Index

- Calculated true overlap of actual categories and predicted categories (from kmeans) for both LDA and word2vec
- Randomized actual categories and calculated overlap of that result and predicted categories 1000 times.

# Jaccard Index Distribution Histograms





## **Solo Analysis: more classifiers**

Use the time difference between deadline and launching time to classify success status

- 1) Calculate the mean time diff. for the three classes;
- 2) Use 11 months of data to train the LDA classifier;
- 3) Use the remaining one month of data as test data;
- 4) Repeat;
- 5) Calculate accuracy.

Accuracy = 0.4119

P-value =  $6.8294e-12$



## **Solo Analysis: more classifiers**

Use the number of backers to classify success status

- 1) Calculate the mean number of backers for the three classes;
- 2) Use 11 months of data to train the LDA classifier;
- 3) Use the remaining one month of data as test data;
- 4) Repeat;
- 5) Calculate accuracy.

Accuracy = 0.4247

P-value =  $6.9781e-12$

# Solo Analysis: Significance

