# KICKSTARTER Project Analytics

## A NLP Approach to Increase the Chance of Building a Successful Crowdfunding Project

**Dao Yong | Hanyu | Jing Jie | Jocille | Raymond | Soo Yeon**

**Date: 6 April 2021**

# PROJECT OVERVIEW – BACKGROUND

## About Kickstarter



- Global **crowdfunding** platform
- Fund **All-or-Nothing** model
- Collect **5%** of the successful funds raised.

## Business Problem



- **38.61%** projects fail, according to Kickstarter
- **$229,775,700 USD** opportunity cost

## Stakeholders



- New aspiring **project starters** -> successfully raise funds
- **Kickstarter management** -> collect fees from successful projects

2

# BUSINESS USE CASES

Fundraisers will be provided with the **sentiments of backers** towards project categories of their choice.

Fundraisers will be **recommended** with related and similar projects.

## OUR GOAL

**To use text mining to help future Kickstarter campaigns to increase their success rate**

3

# DATASET AND TOOLS



### Web Robots Kickstarter Datasets

- Data is crawled and published once a month
- Includes short project description, location, amount raised, goal and project status

### Crawled Data From Kickstarter

- Using Octoparse to crawl project comments from Kickstarter

### Our Final Datasets

- Combined through project names
- 605 Projects in final datasets
- 25158 comments

4

# SOLUTION OVERVIEW (1)

## 1. Topic Modelling and Classification

- Gather description & title
- To identify new key trending topics
- Optimally recommend categories for fundraisers
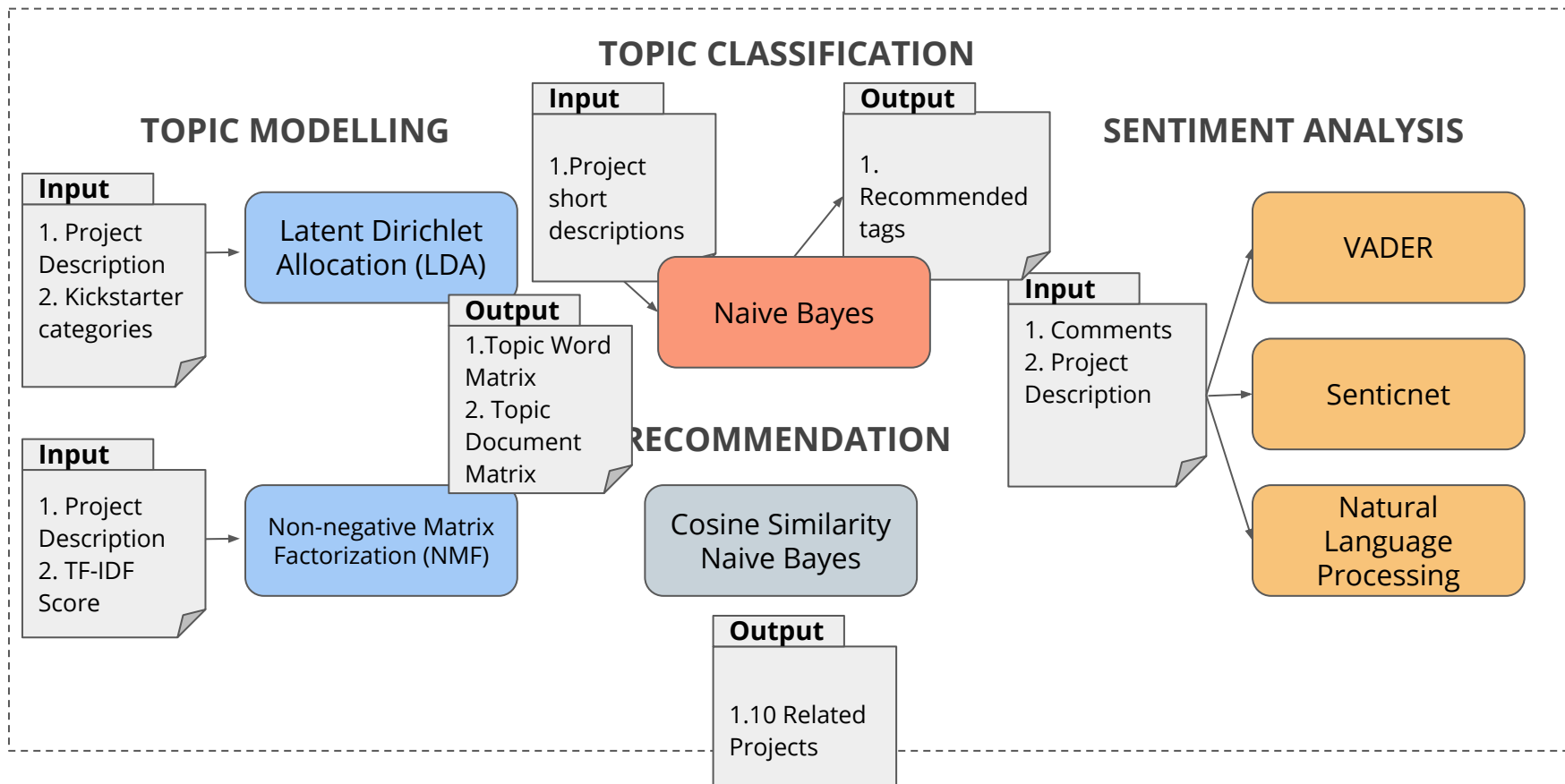
## 2. Sentiment Analysis

- To find out general attitude towards a campaign from comments
- To draw business insights from the different project categories in a top down manner

## 3. Content Based Recommender

- To **recommend** similar and related projects to fundraisers assist their research

5

# SOLUTION OVERVIEW (2)

## TOPIC CLASSIFICATION

### TOPIC MODELLING

**Input**
1. Project Description
2. Kickstarter categories

Latent Dirichlet Allocation (LDA)

**Input**
1. Project short descriptions

**Output**
1. Recommended tags

### SENTIMENT ANALYSIS

VADER

Naive Bayes

**Output**
1. Topic Word Matrix
2. Topic Document Matrix

**Input**
1. Comments
2. Project Description

Senticnet

**Input**
1. Project Description
2. TF-IDF Score

Non-negative Matrix Factorization (NMF)

RECOMMENDATION

Cosine Similarity Naive Bayes

Natural Language Processing

**Output**
1.10 Related Projects

6

# TASK 1 – TOPIC MODELLING

**PURPOSE**: To identify trending topics and keywords

**APPROACH: LDA**

- Unsupervised classification of documents

**RESULTS:**

- Coherence Score:  0.38

  Word intrusion: 0.27

**CHALLENGES:**
Choosing the number of topics and finding the right hyperparameters

**APPROACH: NMF**

- Application of TF-IDF vector to document and model the topics

**RESULTS:**

- Coherence Score: 0.42
- Word intrusion: 0.40

**CHALLENGES:**
Choosing the right number of topics and setting TF-IDF to work with NMF

7

# TASK 2 – PROJECT TOPIC CLASSIFICATION

**PURPOSE**:   Optimally recommend categories for fundraisers

**APPROACH:**

**Naive Bayes Classification Technique**

**Train:** Trained the model with 100 documents for each category

**Validate:** Validated with 20% of the volume of the trained model (Unique data)

**Test:** Tested with 10% of the volume of the trained model (Unique data)

**RESULT:**

**99.29%** accuracy with validation

**1.42%** error with testing (Prediction)

**CHALLENGES:**

- Non-english descriptions
- Wrong outputs in some cases

8

# TASK 3 – SENTIMENT ANALYSIS

**PURPOSE**: To find out the general attitude towards different groups of projects.

**APPROACH:**
1. **Comparison** of sentiment results on comments and projects descriptions from VADER and Senticnet.
   a. Word level - from model
   b. Sentence level - from model or average score
   c. Project level - average score
2. Natural Language Processing technique to find out the **most commonly used words** in different projects across locations, categories, status etc.
   a. Ngram

**RESULT:**
- **Visualisation** & **Insights**!

film & video,television



Sentiment Scores by Category          Unigram for All Comments

- **VADER** performs better!

**CHALLENGES:**
- Non-english comments and misspelling.
- Limitation on VADER.

9

# BONUS TASK – CONTENT BASED RECOMMENDER

**PURPOSE**: To find similar & related projects

**APPROACH:**
- Base on topic modeling, group the projects into different groups
- Use cosine similarity to find related projects
- Comparison between normal cosine similarity and topic based cosine similarity

```
Precision@k = (# of recommended items
@k that are relevant) / (# of
recommended items @k)
```

**RESULT:**

| Precision @ K | LDA + Cosine | 0.636 |
| --- | --- | --- |
| | NMF + Cosine | 0.472 |
| | Cosine | 0.916 |

**CHALLENGES:**
- Using the short description/category might not be enough
- Cosine similarity only matches the exact words
- Lexical ambiguity

10

# DEMO

| | Findings | Insights | Analysis |
|---|---|---|---|
| 1 | **NMF:**<br>Topic 9: film and horror appear frequently | ● Most kickstarter project ask for funds in areas of **films** and **horror**.<br>● Trending topic for people asking for fundings | Allows Kickstarters to **gauge** funding based on the topics |
| 2 | **Sentiment Analysis:**<br>● Technology-related projects receive most comments.<br>● Art/illustration raise the most amount. | ● Imbalance of public attention and who gets the fundings.<br>● **Technology-related projects** have the potentials to grow in the future. | Many people are looking forward to advanced technology, but taken aback by the chance of failure. |
| 3 | **Sentiment Analysis:**<br>Successful projects in Hong Kong and Germany do not receive as much comments | ● People from Hong Kong and Germany are likely to be non-native English speaker<br>● Project descriptions will not be translated. | **Auto-translation** of the project description can be implemented to increase the project exposure to non-english speaking countries. |
| 4 | **Sentiment Analysis:**<br>● Most successful project category is Art/Illustration.<br>● Card games such as Tarot and fictional characters are popular. | Audience and backers tend to look for products that are not easily found on mainstream markets on Kickstarters. | Fundraisers and Kickstarters should pay more attention to projects that are original, and in the field of **horror/weird fiction** and **fortune-telling.** |

12

# WHAT WORKS WELL AND WHAT DOES NOT

## GOOD

- **NMF** for Topic Modelling
- **Cosine Similarity** for Recommender
- **VADER** is able to find the most positive and negative comments

## BAD

- Recommender with **topic modelling** is still not optimal
- **Senticnet** did not perform well for our sentiment analysis

13

# LIMITATIONS

| Datasets | Unbalanced number of comments | Non-english and multilingual comments; Misspelling |
|---|---|---|

- Our dataset is a percentage of all the Kickstart projects

- Our approach is **scalable** once full dataset is obtained.



count of comments by project id top 50 projects

- Removed all the non-English comments using Langid model.

- Some sentiments are not considered.

14

# FUTURE WORKS

- Work on non-english comments for sentiment analysis

- Further tuning of hyperparameters

- Lexical ambiguity and dictionary for unknown words

# Beyond the class?

- **NMF**

  - Model used for document clustering

  - Non negative elements needed

  - Multiplication of 2 matrix (W & H) to derive the combine matrix (V)

- **Sentiment Analysis using VADER and Senticnet**

  - **VADER** - Valence Aware Dictionary for Sentiment Reasoning. It uses a combination of **sentiment lexicon** (phrases, emoticons, acronyms) and **grammatical rules**.

  - **Senticnet** - Concept-level sentiment analysis. It focuses on **semantic analysis** through the use of **semantic networks**.

16

# Thank You!

## QnA

# Why VADER Performs Better?

| Model | Most Positive Comments | Most Negative Comments |
|---|---|---|
| **VADER** | Congrats Sam, it is a great idea and I look forward to watching and if given the chance to participate I'll be there. I'm a truckin biker of the old school I make more time to ride these days as I have been a successful owner-operator in the trucking business. As a cold war vet from '69 to '71 I relish experiencing our great countries history and meeting the folks who hold secrets and stories not readily available to most. Best wishes and best of luck. | Fuck it, why not. But you better hit someone with that stick. |
| **Senticnet** | Nicely presented Kickstarter - Good Luck!<br><br>@ jerice50: Thank you very much. | Congratulations on funding! |

# How We Calculate Word Intrusion Score

| Topic | Tick if Intrurder (One Only) | | | | | | | Topic | Tick if Intrurder (One Only) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| app | | | | | | | | pin | | | | | | |
| power | | | | | | | | enamel | | | | | | |
| summer | ☑ | ☑ | ☑ | ☑ | ☑ | ☑ | | enamel pin | | | | | | |
| easy | | | | | | | | hard enamel | | | | | | |
| device | | | | | | | | hard | | | | | | ☑ |
| build | | | | | | | | pin inspri | | | | | | |
| smart | | | | | | | | theme | | | | | | |
| control | | | | | | | | cute | | | | | | |
| phone | | | | | | | | spooki | | | | | | |
| camera | | | | | | | | cabinet | ☑ | ☑ | ☑ | ☑ | ☑ | |
| screen | | | | | | | | set | | | | | | |

| | JJ | DY | Raymond | Hanyu | Soo Yeon | Jocille | | | JJ | DY | Raymond | Hanyu | Soo Yeon | Jocille |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Topic | Tick if Intrurder (One Only) | | | | | | | Topic | Tick if Intrurder (One Only) | | | | | |
| make | | | | | | | | card | | | | | | |
| great | | | | | | | | play | | | | | | |
| favorite | | | | | | | | play card | | | | | | |
| attractive | | | | | | | | recording | ☑ | ☑ | ☑ | ☑ | ☑ | ☑ |
| awesome | | | | | | | | deck | | | | | | |
| special | | | | | | | | deck play | | | | | | |
| machine | ☑ | ☑ | | ☑ | ☑ | ☑ | | card game | | | | | | |
| gift | | | | | | | | card deck | | | | | | |
| touch | | | | | | | | uspcc | | | | | | |
| christmas | | | | | | | | print uspcc | | | | | | |
| pen | | | ☑ | | | | | custom deck | | | | | | |

# FUII SOLUTION OVERVIEW

## TOPIC MODELLING

**Input**
1. Project Description
2. Kickstarter categories

→ Latent Dirichlet Allocation (LDA)

**Input**
1. Project Description
2. TF-IDF Score

→ Non-negative Matrix Factorization (NMF)

**output/input**
1. Extracted topics
2. Coherence score

+Chosen topic and project

## TOPIC CLASSIFICATION

**Input**
1. Project short descriptions

Naive Bayes

**Output**
1. Recommended tags

## RECOMMENDATION

Cosine Similarity Naive Bayes

**output**
Recommended projects similar to chosen project

## SENTIMENT ANALYSIS

**Input**
1. Comments
2. Project Description

VADER

Senticnet

Natural Language Processing

**Output**
Positive, negative and compound score

**Output**
Aggregated sentiment scores

**Output**
Wordclouds