

## Milestone 1 Report:

# Problem Statement, Data Overview, Data Acquisition & Processing, Initial Exploratory Analysis

### Problem Statement:

For my second capstone, my goal was to predict campaign success on Kickstarter by leveraging NLP techniques.

Potential parties that could be interested in this project include:

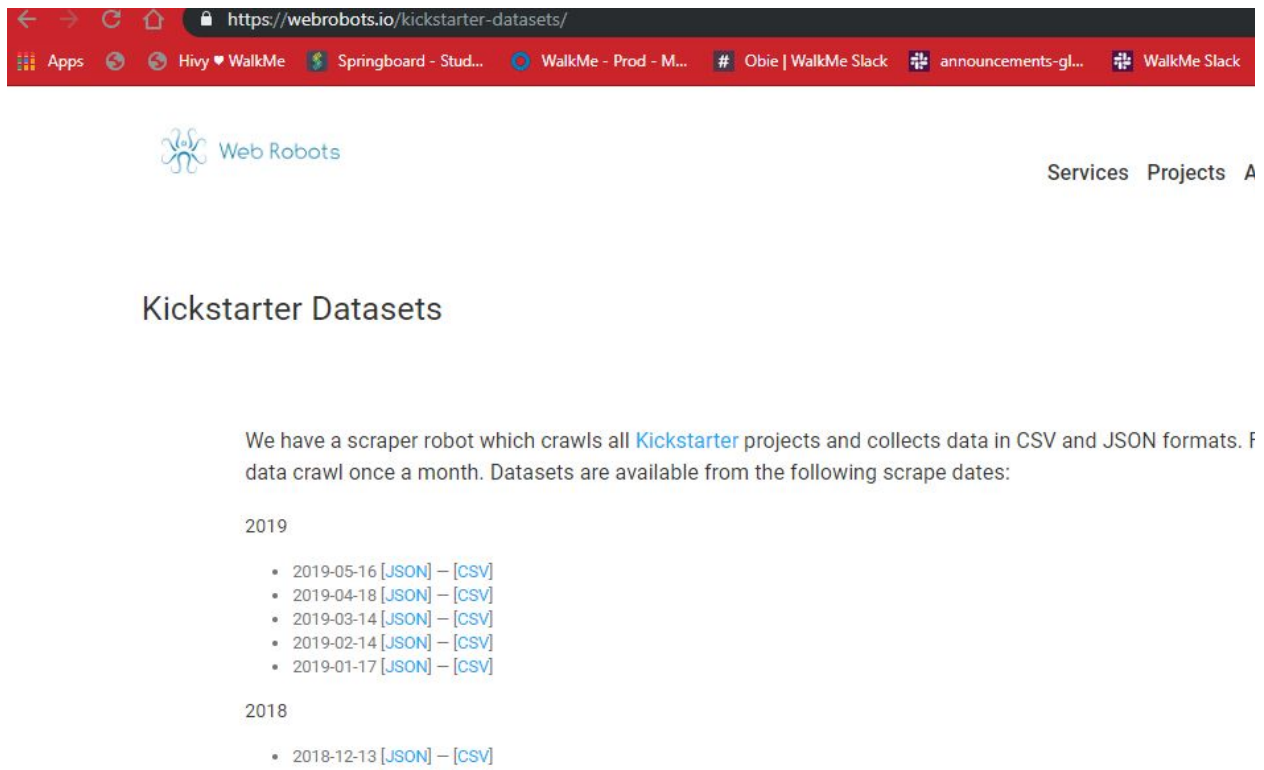
1. Potential creators wanting to understand:
  - 1) What is their realistic chance of being successful on Kickstarter
  - 2) What does the competitive landscape look like for similar categories, etc
  - 3) How they can better message their products.
2. Knock-off manufacturers wanting to understand who they should be ripping off
3. Competitive crowdfunding sites wanting to understand:
  - 1) The supply & demand of crowdfunding sites
  - 2) The products being moved (or not)

### Data Overview:

Data for the project came from a series of csv files (hosted by Webrobots) which were produced by scraping the Kickstarter site once a month (<https://webrobots.io/kickstarter-datasets/>).

Each scrape could produce anywhere from 25-50+ csv files which were stored in folders labeled by the month and day the data was scraped. Each row in a csv file

represents a single project and includes features like the creator, goal amount, country of origin and important text fields like “blurb” (a 1 sentence summary of the project), “category”, “title” and “creator”. Given the scrapes were performed once a month and the average campaign length was 30 days, the same campaign could be captured multiple times and a majority of campaigns could have been posted between scrapes.



A concern of using data from a point in time too far into the campaign is producing a model that results in overly high accuracy because of signal words like “success”, “fail” or even phrases like “reached our goal ahead of schedule!”.

Given we’re trying to predict a specific outcome (“successful” or “failed”) and we have labeled data, this would be a classification problem that could be addressed using a logistic regression, random forest classifier, or similar algorithms. Our goal is to use as many of the predictors as possible, given each file contains only about 32 features. We have a mix of data types including datetime columns (deadline, state\_changed\_at,

created\_at), numeric columns (backers, goal, usd\_pledged), text (blurb, title), and categorical (country, currency, state).

Deliverables will include a jupyter notebook, a summary paper, and a slide deck. The comprehensive list of possible features were:

- 'backers\_count', 'blurb', 'category',  
'converted\_pledged\_amount', 'country', 'created\_at', 'creator',  
'currency', 'currency\_symbol', 'currency\_trailing\_code',  
'current\_currency', 'deadline', 'dirname', 'disable\_communication',  
'friends', 'fx\_rate', 'goal', 'id', 'is\_backing', 'is\_starrable',  
'is\_starred', 'last\_update\_published\_at', 'launched\_at', 'location',  
'name', 'permissions', 'photo', 'pledged', 'profile', 'slug',  
'source\_url', 'spotlight', 'staff\_pick', 'state', 'state\_changed\_at',  
'static\_usd\_rate', 'unread\_messages\_count', 'unseen\_activity\_count',  
'urls', 'usd\_pledged', 'usd\_type'

Formula bar content: {"id":36,"name":"Classical Music","slug":"music/classical music","position":3,"parent\_id":14,"color":10878931,"urls":{"web":{"discover":"http://www.kickstarter.com/discover/categories/music/classical%20music"}}}

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
	backers_c	blurb	category	converted	country	created_a	creator	currency	currency	current_c	deadline	disable_c	frie	
2	25	Provocat	{"id":39,"r	1547	US	1.33E+09	{"id":4132	USD	\$	TRUE	USD	1.33E+09	FALSE	
3	5	Perfectly	{"id":313,'	152	CA	1.45E+09	{"id":8899	CAD	\$	TRUE	USD	1.47E+09	FALSE	
4	7	Live event	{"id":24,"r	235	US	1.46E+09	{"id":3840	USD	\$	TRUE	USD	1.47E+09	FALSE	
5	15	The 2nd A	{"id":54,"r	1710	US	1.41E+09	{"id":9458	USD	\$	TRUE	USD	1.42E+09	FALSE	
6	1	Hour-long	{"id":303,'	1	US	1.53E+09	{"id":1176	USD	\$	TRUE	USD	1.54E+09	FALSE	
7	170	Motion co	{"id":16,"r	7796	US	1.52E+09	{"id":7016	USD	\$	TRUE	USD	1.52E+09	FALSE	
8	182	"Hilariou	{"id":292,'	11021	GB	1.44E+09	{"id":9918	GBP	£	FALSE	USD	1.44E+09	FALSE	
9	22	The Robin	{"id":43,"r	938	US	1.34E+09	{"id":1641	USD	\$	TRUE	USD	1.35E+09	FALSE	
10	216	A versatil	{"id":337,'	35176	NZ	1.51E+09	{"id":1103	NZD	\$	TRUE	USD	1.52E+09	FALSE	
11	0	Murica Fe	{"id":24,"r	0	US	1.45E+09	{"id":1582	USD	\$	TRUE	USD	1.45E+09	FALSE	
12	25	A Survivor	{"id":30,"r	12180	US	1.48E+09	{"id":1511	USD	\$	TRUE	USD	1.48E+09	FALSE	
13	3	House of	{"id":303,'	7	SE	1.42E+09	{"id":8184	SEK	kr	TRUE	USD	1.42E+09	FALSE	
14	43	Music hap	{"id":36,"r	6645	US	1.4E+09	{"id":5971	USD	\$	TRUE	USD	1.4E+09	FALSE	
15	77	50 years a	{"id":30,"r	10107	GB	1.46E+09	{"id":1132	GBP	£	FALSE	USD	1.46E+09	FALSE	
16	3	The Interg	{"id":342,'	1097	US	1.54E+09	{"id":8709	USD	\$	TRUE	USD	1.55E+09	FALSE	
17	22	Assist pro	{"id":332,'	16276	US	1.43E+09	{"id":1168	USD	\$	TRUE	USD	1.43E+09	FALSE	
18	0	A video of	{"id":54,"r	0	US	1.37E+09	{"id":9381	USD	\$	TRUE	USD	1.38E+09	FALSE	
19	27	After mor	{"id":43,"r	1737	US	1.35E+09	{"id":1584	USD	\$	TRUE	USD	1.36E+09	FALSE	

In total the original data set was about 30.6 MB and consisted of projects scraped between Jan 28, 2016 and Feb 14, 2019. Some important

considerations needed to be handled given the data collection method.

- 1) Projects would be duplicated between scrapes and we wanted the earliest instance as opposed to the duplicates.
- 2) I was only interested in using data from projects within the early days of the campaign in order to avoid data leakage.
- 3) Many of the columns had data that was nested in a JSON format and needed to be split out.
  - a) Ex:

```
{"id":45,"name":"Art Books","slug":"publishing/art books","position":3,"parent_id":18,"color":14867664,"urls":{"web":{"discover":"http://www.kickstarter.com/discover/categories/publishing/art%20books"}}}
```

Because of the timing of the data collection, we won't be able to answer questions like:

- How did the campaigns perform over time?
- Are there differences in success factors between categories?
- Are there regional differences?
- How were edits handled within the same campaigns over the lifetime of the campaign?

## Data Acquisition & Processing:

In order to begin cleaning the data I needed to first compile the 30.6MB of data spread across 100+ csv files into a single source that could be queried. I decided to leverage the sqlite3 & os python packages to create a local SQLite database that could be loaded with data from the csv files, which would allow me to define the schema, query the records, and align the common columns.

```

53 def main():
54
55     sql_delete_table = """ DROP TABLE IF EXISTS maindump;"""
56
57     sql_create_maindump_table = """ CREATE TABLE IF NOT EXISTS maindump (
58         ROW_ID INTEGER PRIMARY KEY AUTOINCREMENT,
59         id integer,
60         backers_count integer,
61         dirname text,
62         blurb text,
63         category text,
64         converted_pledged_amount integer,
65         country text,
66         created_at integer,
67         creator text,
68         currency text,
69         currency_symbol text,
70         currency_trailing_code text,
71         current_currency text,
72         deadline text,
73         disable_communication text,
74         friends text,
75         fx_rate text,
76         goal text,
77         is_backing text,
78         is_starrable text,
79         is_starred text,
80         last_update_published_at text,
81         launched_at text,
82         location text,
83         name text,
84         permissions text,
85         photo text,
86         pledged text,
87         profile text,
88         slug text,
89         source_url text,
90         spotlight text,
91         staff_pick text,
92         state text,
93         state_changed_at text,
94         static_usd_rate text,
95         unread_messages_count text,
96         unseen_activity_count text,
97         urls text,
98         usd_pledged text,
99         usd_type text,
100        scrape_date text,
101        state_changed_at_clean text,
102        created_at_clean text,
103        deadline_clean text,
104        launched_at_clean text
105    ); """
106
107

```

```

10 import sqlite3
11 from sqlite3 import Error
12
13
14 def create_connection(db_file):
15     """ create a database connection to the SQLite database
16         specified by db_file
17     :param db_file: database file
18     :return: Connection object or None
19     """
20     try:
21         conn = sqlite3.connect(db_file)
22         return conn
23     except Error as e:
24         print(e)
25
26     return None
27
28
29
30 def create_table(conn, create_table_sql):
31     """ create a table from the create_table_sql statement
32     :param conn: Connection object
33     :param create_table_sql: a CREATE TABLE statement
34     :return:
35     """
36     try:
37         c = conn.cursor()
38         c.execute(create_table_sql)
39     except Error as e:
40         print(e)
41
42
43
44 def delete_table(conn, delete_table_sql):
45     try:
46         c = conn.cursor()
47         c.execute(delete_table_sql)
48     except Error as e:
49         print(e)
50
51

```



```

45 def main():
46     database = "G:\\My Drive\\sqlite\\db\\kickstarter.db"
47     path = "Users\\mikiko.b\\Documents\\Kickstarter\\data"
48     directory = os.path.join("c:\\", path)
49
50     # create a database connection
51     conn = create_connection(database)
52     with conn:
53         for dirpath, dirnames, filenames in os.walk(directory):
54             global dirname
55             dirname = dirpath.split(os.path.sep)[-1]
56             print("dirpath:", dirpath)
57             print("dirnames:", dirnames)
58             print("dirname:", dirname)
59             for file in filenames:
60                 with open(os.path.join(dirpath, file), encoding="utf8") as csvfile:
61                     # remove this line
62                     print("dirpath+file", os.path.join(dirpath, file))
63                     reader = csv.DictReader(csvfile)
64                     for row in reader:
65                         a_id=row.get('id',None)
66                         backers_count=row.get('backers_count',None)
67                         blurb=row.get('blurb',None)
68                         category=row.get('category',None)
69                         converted_pledged_amount=row.get('converted_pledged_amount',None)
70                         country=row.get('country',None)
71                         created_at=row.get('created_at',None)
72                         creator=row.get('creator',None)
73                         currency=row.get('currency',None)
74                         currency_symbol=row.get('currency_symbol',None)
75                         currency_trailing_code=row.get('currency_trailing_code',None)
76                         current_currency=row.get('current_currency',None)
77                         deadline=row.get('deadline',None)
78                         disable_communication=row.get('disable_communication',None)
79                         friends=row.get('friends',None)
80                         fx_rate=row.get('fx_rate',None)
81                         goal=row.get('goal',None)
82                         is_backing=row.get('is_backing',None)
83                         is_starrable=row.get('is_starrable',None)
84                         is_starred=row.get('is_starred',None)
85                         last_update_published_at=row.get('last_update_published_at',None)
86                         launched_at=row.get('launched_at',None)
87                         location=row.get('location',None)
88                         name=row.get('name',None)
89                         permissions=row.get('permissions',None)
90                         photo=row.get('photo',None)
91                         pledged=row.get('pledged',None)
92                         profile=row.get('profile',None)
93                         slug=row.get('slug',None)
94                         source_url=row.get('source_url',None)
95                         spotlight=row.get('spotlight',None)
96                         staff_pick=row.get('staff_pick',None)
97                         state=row.get('state',None)
98                         state_changed_at=row.get('state_changed_at',None)
99                         static_usd_rate=row.get('static_usd_rate',None)
100                         unread_messages_count=row.get('unread_messages_count',None)
101                         unseen_activity_count=row.get('unseen_activity_count',None)
102                         urls=row.get('urls',None)
103                         usd_pledged=row.get('usd_pledged',None)
104                         usd_type=row.get('usd_type',None)
105
106
107                     # now I'm going to clean the dates
108                     state_changed_at_clean = datetime.utcfromtimestamp(int(state_changed_at)).strftime('%Y-%m-%d %H:%M:%S')
109                     created_at_clean = datetime.utcfromtimestamp(int(created_at)).strftime('%Y-%m-%d %H:%M:%S')
110                     deadline_clean = datetime.utcfromtimestamp(int(deadline)).strftime('%Y-%m-%d %H:%M:%S')
111                     launched_at_clean = datetime.utcfromtimestamp(int(launched_at)).strftime('%Y-%m-%d %H:%M:%S')
112
113                     # create scrape date
114                     scrape_date = str(dirname)[12:22]
115
116                     #combine to create an insertable row - "project"
117                     #call create_project
118                     campaign = (a_id,backers_count,dirname,blurb,category,converted_pledged_amount,country,created_at,c
119                     create_campaign(conn, campaign)
120
121             print("added: ", dirname)
122
123
124
125
126
127
128
129 if __name__ == '__main__':
130     main()
131

```

After defining the table schema, I also converted the datetime columns (state\_changed\_at\_clean, created\_at\_clean, deadline\_clean, launched\_at\_clean) from a utc timestamp to a human readable format. Once all the projects were loaded into a single table, I then needed to stage the creation of the de-duped master data set.

```
In [20]: 1 project_outcomes.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1885073 entries, 0 to 1885072
Data columns (total 10 columns):
Unnamed: 0      int64
id              int64
slug            object
creator         object
created_at_clean object
launched_at_clean object
profile         object
country         object
state          object
state_changed_at_clean object
dtypes: int64(2), object(8)
memory usage: 143.8+ MB
```

Leveraging pandas read\_sql\_query function I first needed to get the list of unique projects :

```
projects_outcomes =
pd.read_sql_query("select id, slug,
creator, created_at_clean,
launched_at_clean, profile, country, state,
state_changed_at_clean from maindump
where state != 'live';", conn)
```

Specifically, I selected projects that were posted within less than 15 days prior to a scrape being run. This list would consist of the projects outcomes and would have the duplicates dropped from the pandas dataframe produced from running the query.

```
In [19]: 1 project_starting.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 478 entries, 0 to 477
Data columns (total 12 columns):
Unnamed: 0      478 non-null int64
id              478 non-null int64
slug            478 non-null object
name            478 non-null object
blurb           478 non-null object
category        478 non-null object
goal            478 non-null int64
launched_at_clean 478 non-null object
deadline_clean   478 non-null object
location        478 non-null object
scrape_date      478 non-null object
state           478 non-null object
dtypes: int64(3), object(9)
memory usage: 44.9+ KB
```

The starting data, that would consist of the predictors we'd want to use to create the predictive model, was created through the following query:

```
projects_starting_data =
pd.read_sql_query("select id, slug, name, blurb, category, goal, launched_at_clean,
deadline_clean, location, scrape_date, state from maindump where
(julianday(date(launched_at_clean)) - julianday(date(scrape_date))) < 15 and
(julianday(date(launched_at_clean)) - julianday(date(scrape_date))) > -1;", conn)
```

```
: 1 consolidated_table.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 478 entries, 0 to 1650
Data columns (total 21 columns):
Unnamed: 0_x    478 non-null int64
id_x            478 non-null int64
slug            478 non-null object
name            478 non-null object
blurb           478 non-null object
category        478 non-null object
```

The starting and outcomes data sets were then merged to produce a master data set, which only included 478 campaigns.

Additional cleaning including using regex to break out the nested columns, specifically “creator”, “category”, and “location” to isolate the creator name, category and country, city, state.

### Initial Exploratory Analysis:

When analyzing the data I first wanted to understand if there were potential issues with class imbalance between failed and successful campaigns and found that the ratios ended up being nearly equal.

I also wanted to understand whether the average campaign duration differed by outcome and average time to launch>

