

# Ian Perry

Capstone Project

# Project



# DATA SCIENCE WORKFLOW



# DataSet / Cleaning

Rank	Title	Category	User	Views	Upload	#	+/-	URL	Duration	Descript...	Keyword	Keyword	Keyword
2	Speed of	Science &	Science...	345732	2016-09-11	413	1160/108	<a href="#">http://...</a>	44:28	For more	yes	no	0.00%
3	The	Education	Real	22592	2016-11-15	171	213/11	<a href="#">http://...</a>	50:09	"You feel	yes	no	0.00%
4	[Blow	Science &	Discovery	550	2016-11-17	0	53/2	<a href="#">http://...</a>	0s	LIKE -----	yes	yes	8.50%
5	National	Science &	Star Staff	14277	2016-11-14	10	56/14	<a href="#">http://...</a>	44:03	The	yes	no	0.00%
6	Anonym...	Nonprofits	Anonym...	734223	2016-07-24	3073	10538/298	<a href="#">http://...</a>	58:00	Anonym...	yes	yes	2.01%
7	Living	Education	Real	122780	2016-11-11	483	829/64	<a href="#">http://...</a>	52:02	An	yes	yes	0.32%

# Successes

- Three laws
- Fascinating Topic
- A small difference
- NLTK

# Challenges

- The API
- Data Capture
- Duplication
- Skewness in Dislikes

# For example



This video is no longer available due to a  
copyright claim

Sorry about that.

# Top Docs

## Top Likes

- 3 - Beyonce
- 2 - Cyber Bully
- 1 - Suicide Forests in Japan.  
(Viewer Discretion)

## Top Dislikes

- 3 - Brexit
- 2 - Living as a baby
- 1 - Cult Leader  
Thinks He's Jesus

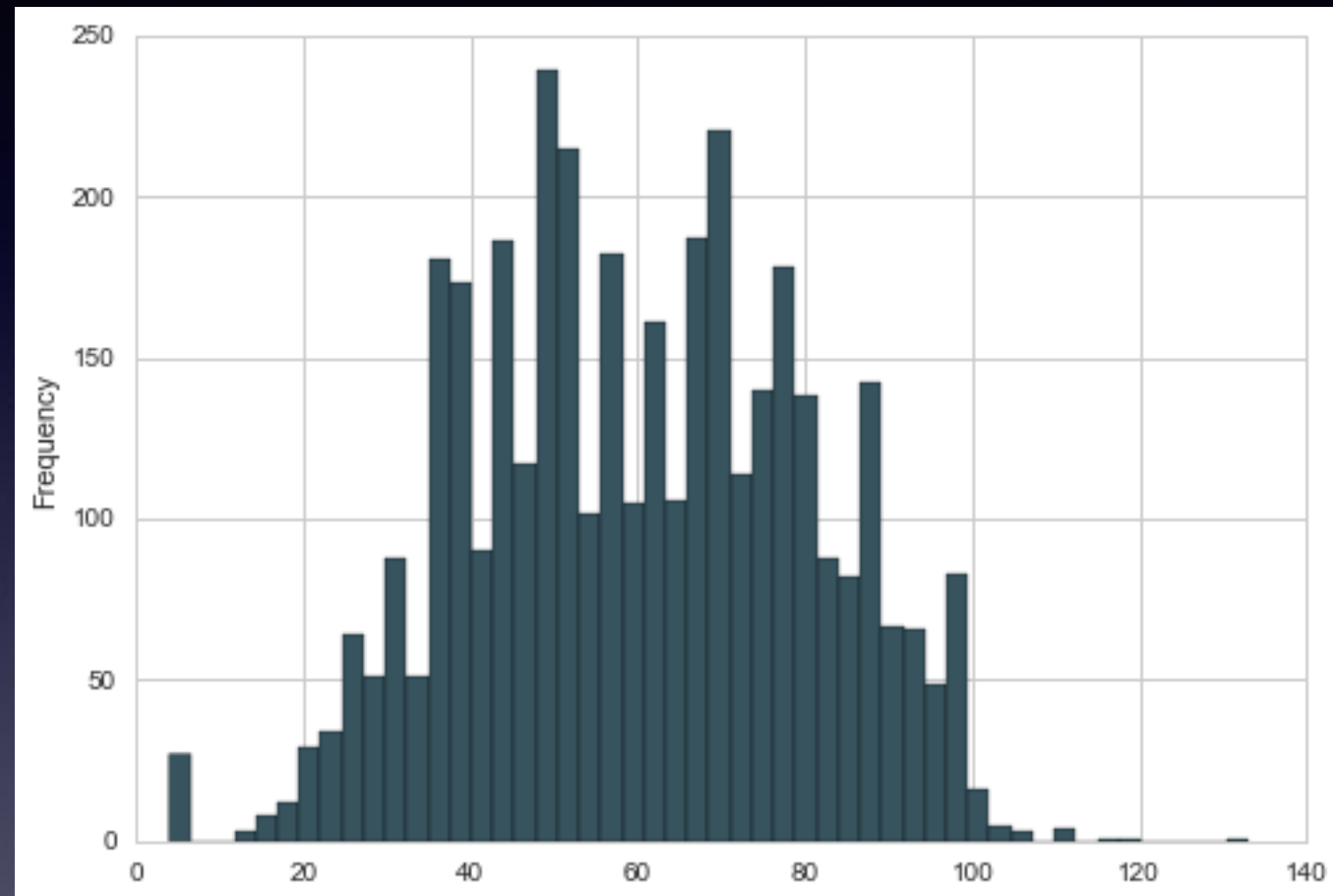


# Top Cats

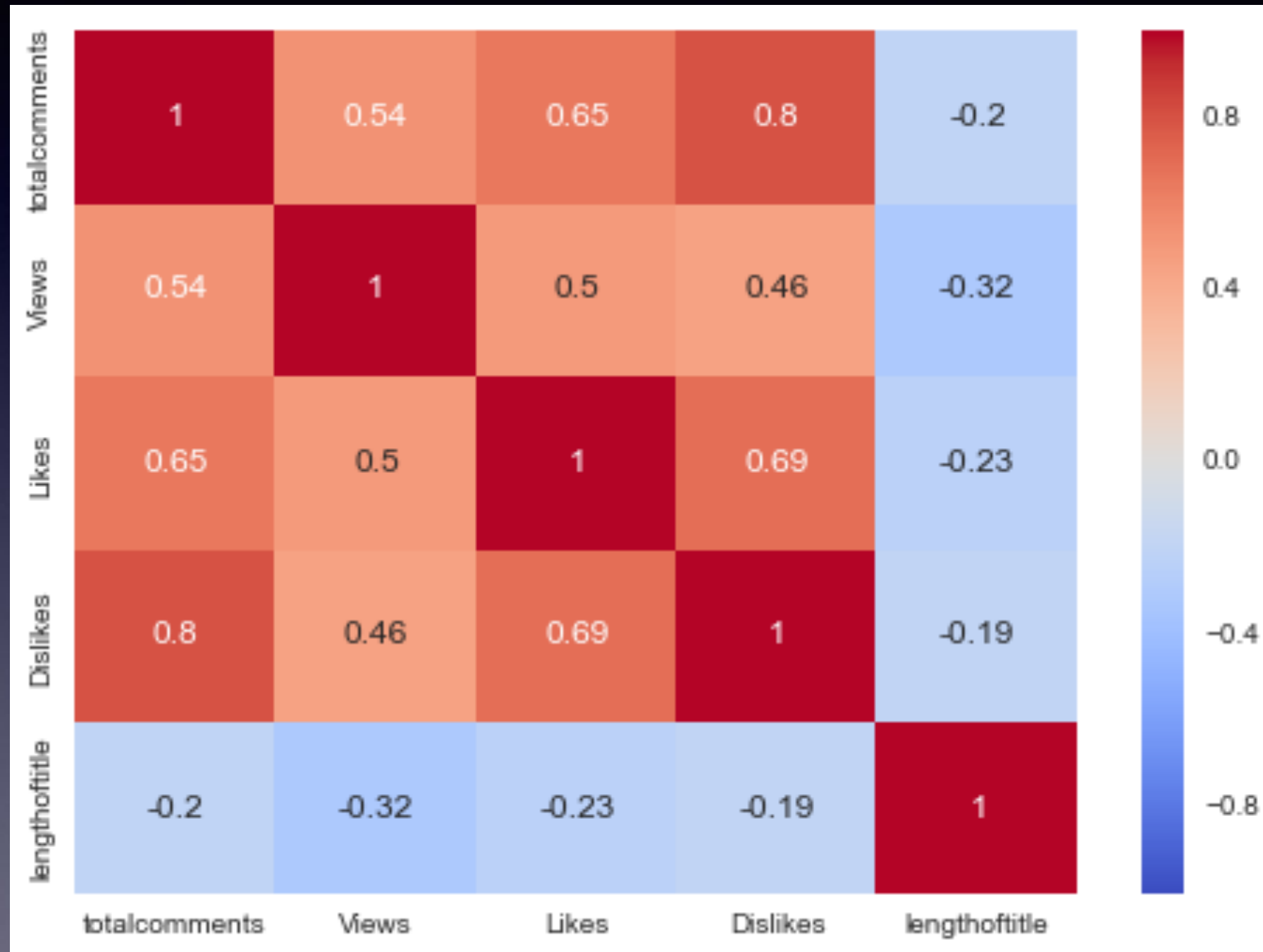
## **Top Categories**

- 3 - Education
- 2 - Film and animation
- 1 - People and Blogs

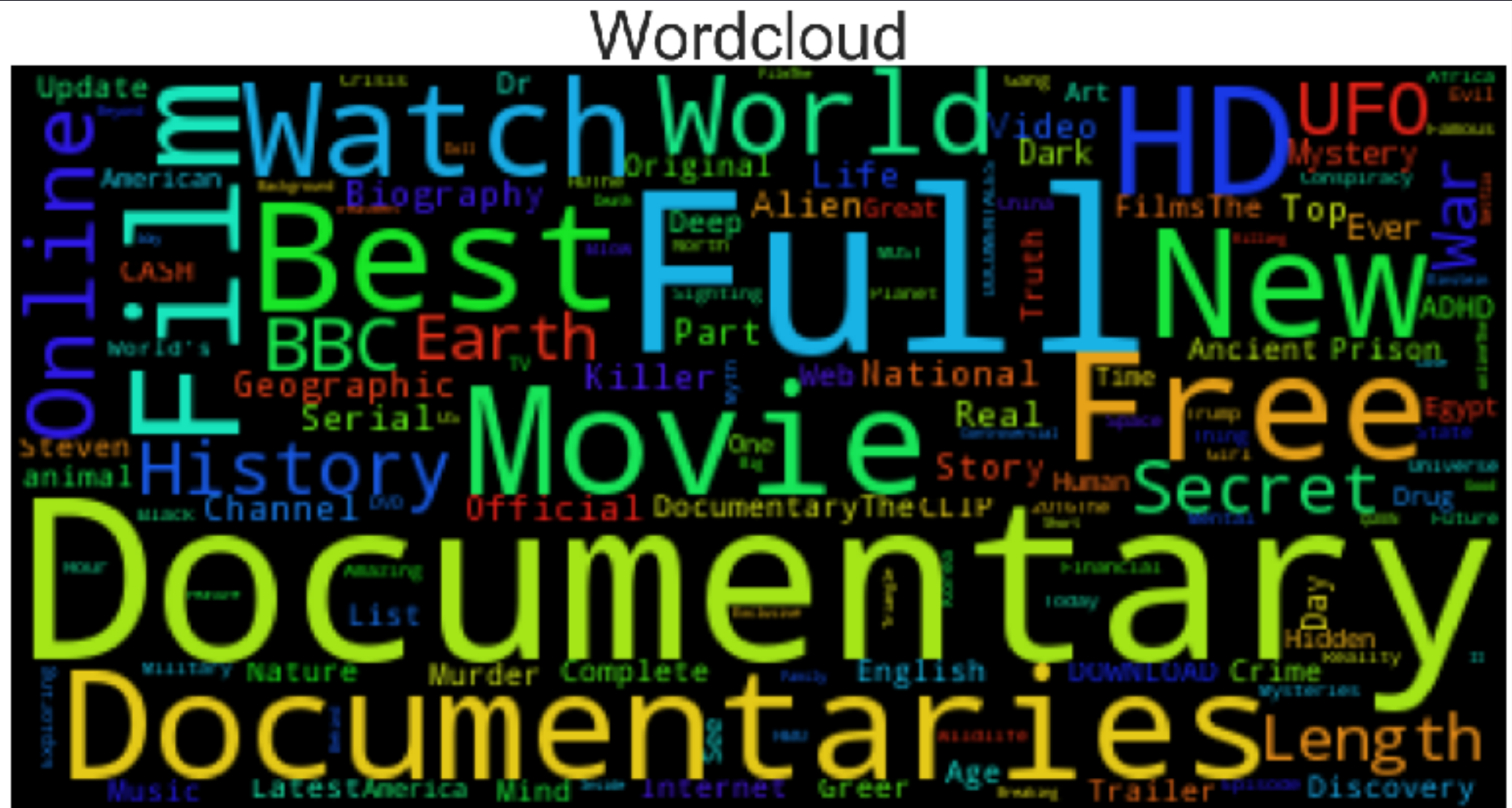
# Title Length



# Length of Title



# Wordcloud



# Modeling

	Bayes	LR	TFIDF
Accruacy	0.79	0.82	0.81
Precision	0.80	0.81	0.83
Recall	0.80	0.81	0.83
F1	0.80	0.81	0.83
Support	1143	1143	1258

Repeat

# Questions