

# Predict the Popularity of a TED Talk



Liu, Shijie | Dr. Liu, Xinlian | Hood College of Frederick Maryland

## Abstract



TED (Technology, Entertainment, and Design)
Talks are being posted by TED Conference
LLC for free on their website and YouTube
Channel under the slogan of "ideas worth
spreading". TED Talk offers a wide range of
topics within the research and practice of
science and culture, and often through
storytelling. It has a variety of presenters as

storytelling. It has a variety of presenters as well, such as Writer, Researcher, and Scientist. TED also gives people the right to raise their own local TEDx Event. How could TED filter and determine which talk shall be published? What an organizer could do to make his/her talk more popular? In this project, I would like to apply Exploratory Data Analysis on all the available attributes and help people to have a better understanding of what are the variables that could affect the popularity of a TED Talk. In addition, applying Natural Language Processing technique on the TED Talk transcript, and then train a classification model with algorithms that could predict the popularity of TED Talk.

**KEYWORDS**: Exploratory Data Analysis, Visualization, Classification, Natural Language Processing

# Techniques

### **Classification Algorithms:**

- Logistic Regression
- Random Forest Classifier
- Decision Tree Classifier
- Multinomial Naïve Bayes

### **Classification Model Evaluation Metrics:**

- Prediction Accuracy
- Confusion Matrix
- Classification ReportCross Validation

# Sample Data

### **TED Main Dataset- Sample Data:**

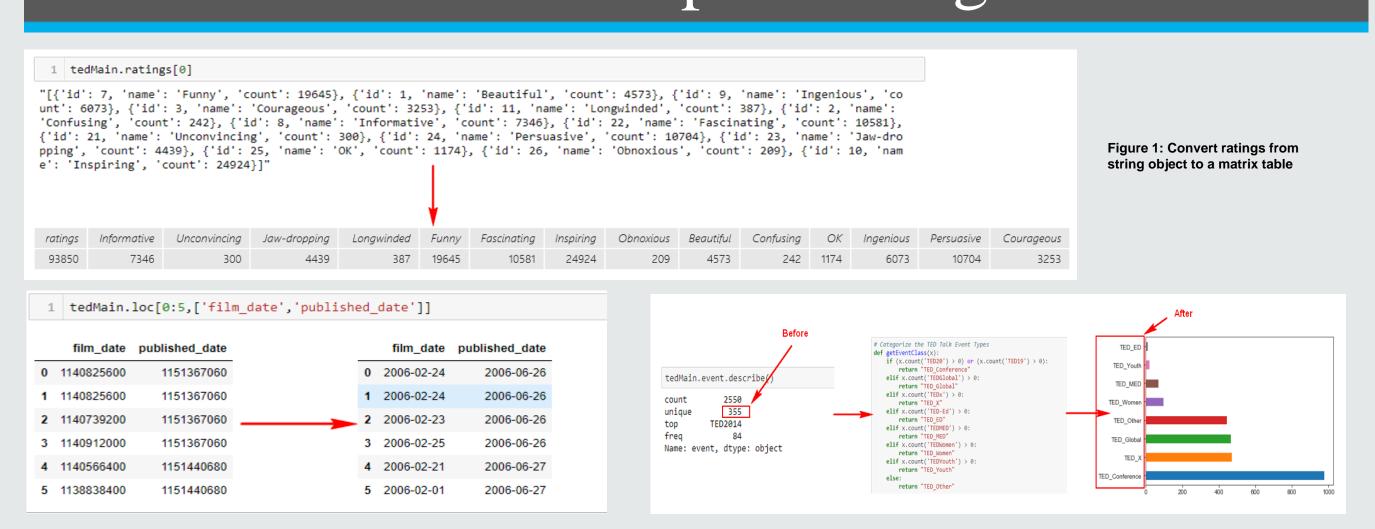
commen ts	descripti on	duration	event	film_dat e	languag es	main_spea ker	name	num_spea ker	published_ date	ratings	related_talk s	speaker_occu pation	tags	title	url	views
4553	Sir Ken Robinso n makes an entertain ing and pro	1164	TED200 6	1140825 600	60	Ken Robinson	Ken Robin son: Do schoo Is kill creati vity?	1	115136706 0	[{'id': 7, 'name' : 'Funny ', 'count' : 19645 }, {	[{'id': 865, 'hero': 'https://pe.t edcdn.com/ im	Author/educat or	['child ren', 'creat ivity', 'cultu re', 'danc e',	Do scho ols kill creat ivity ?	https:// www.te d.com/ talks/k en_rob inson_ says_s c	4722711 0

# TED Transcript Dataset- Sample Data:

Figure 2: Convert UNIX Timestamp to Human Date

ľ		
	transcript	url
	Good morning. How are you?(Laughter)It's been great, hasn't it? I've been blown away by the whole thing. In fact, I'm leaving.(Laughter)There have been three themes running through the conference which are relevant to what I want to talk about. One is the extraordinary evidence of human creativity in all of the presentations that we've had and in all of the people here. Just the variety of it and the range of it. The second is that it's put us in a place where we have no idea what's going to happen, in terms of the future. No idea how this may play out. I have an interest in education. Actually, what I find is everybody has an interest in education. Don't you? I find this very interesting. If you're at a dinner party, and you say you work in education — Actually, you're not often at dinner parties, frankly.(Laughter)If you work in education, you're not asked.(Laughter)And you're never asked back, curiously. That's strange to me. But if you are, and you say to somebody, you know, they say, "What do you do?" and you say you work in education, you can see the blood run from their face. They're like, "Oh my God," you know, "Why me?"(Laughter)"My one night out all week."(Laughter)But if you ask about their education, they pin you to the wall. Because it's one of those things that	https://www.ted.com/talks/ken_robinson_says_schools_kill_creativity

# Data Preprocessing



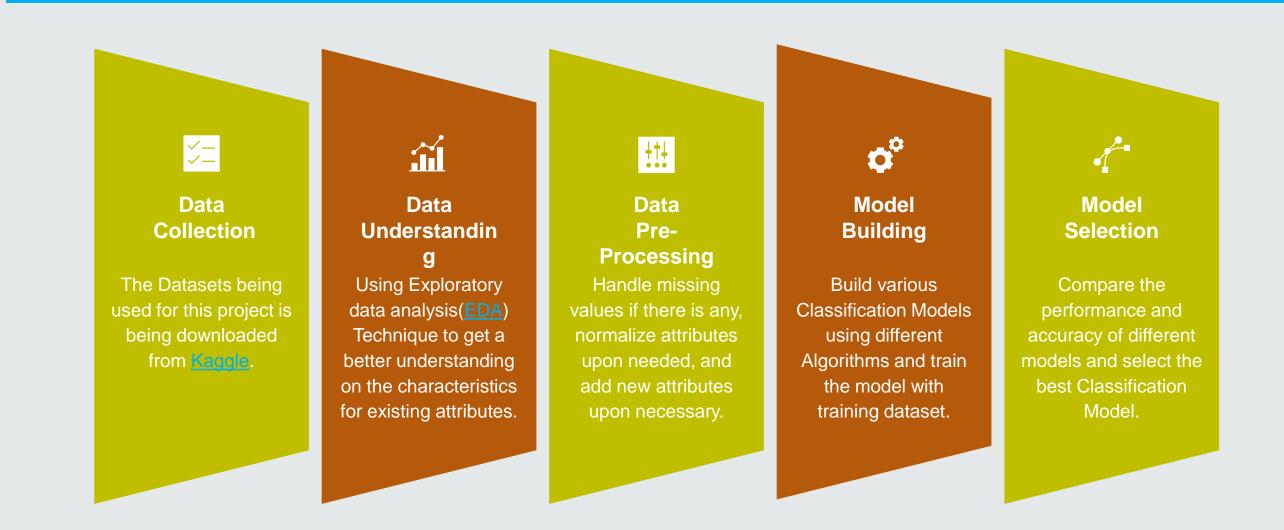
### Figure 3: Adding a Event Class Feature based on event

### Dataset Description **Feature Name** Data Type Description The number of first level comments made on the talk. comments A blurb of what the talk is about. description The duration of the talk in seconds. duration object The TED/TEDx event where the talk took place. event The Unix timestamp of the filming. film\_date The number of languages in which the talk is available. languages The first named speaker of the talk. object main\_speaker The official name of the TED Talk. Includes the title and the speaker. The number of speakers in the talk. num\_speaker published\_date The Unix timestamp for the publication of the talk on TED.com. object A stringified dictionary of the various ratings given to the talk. ratings object A list of dictionaries of recommended talks to watch next. related\_talks object speaker\_occupation The occupation of the main speaker. The themes associated with the talk. tags The title of the talk. The number of views on the talk. views int64 object The URL of the talk.

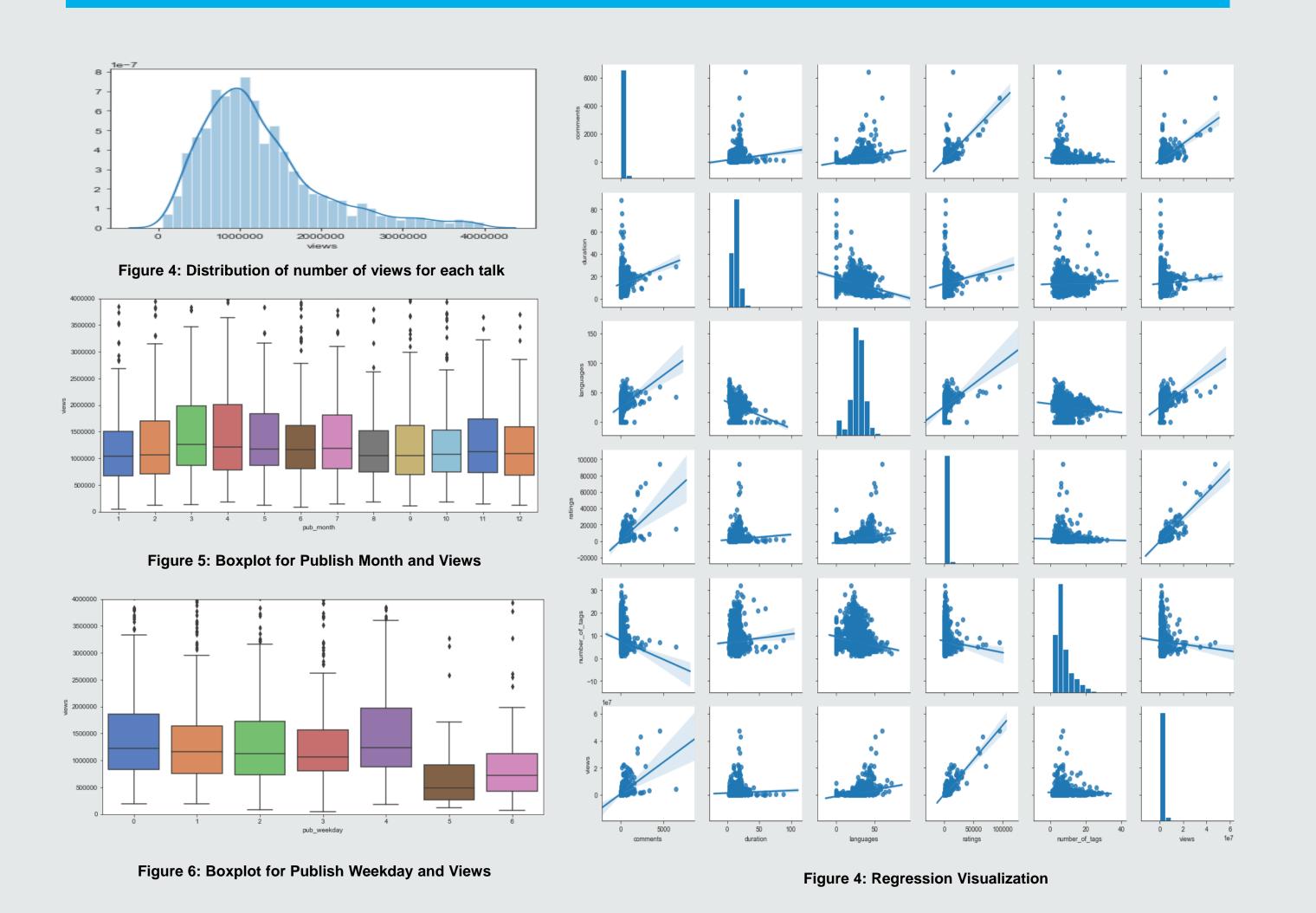
# Implementation

transcript

The official English transcript of the talk.



# Exploratory Data Analysis



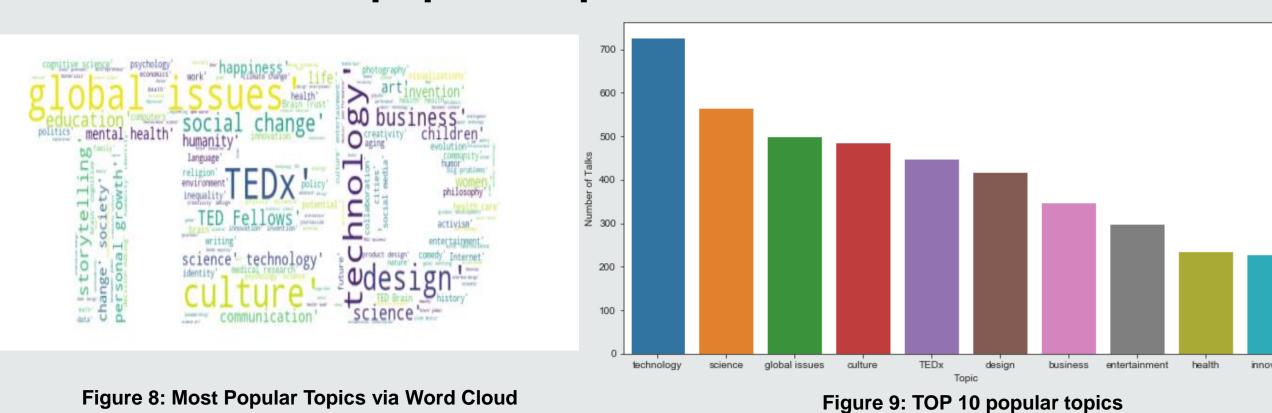
# Result

What makes a popular TED Talk?

# Important features: Ratings Comments Languages Duration Number of tags Published weekday rumber\_of\_attendences rumber\_of\_attendenc

Figure 7: Feature Importance of Random Forest Classifie

What are the most popular topics in TED Talk?



Predict the popularity of a TED Talk based on its transcript ?

• **Answer:** Prediction Model based on Multinomial Naïve Bayes [2] algorithm is selected as a better model that has Prediction Accuracy of 68%.

# Future Work

- Although this project has met its original objectives that it could identify those key features that contribute to the popularity of a TED talk and it could also predict the popularity based on transcript, but the features being used to train the classification model is limited and the prediction accuracy is lower than expected.
- Future work will be carried out to understand why the prediction accuracy is lower and do further pre-tuning and post-tuning to enhance the classification model. Further, understanding the reality of all features and take full advantage of all available features while training the classification model, such as, title, description, and tags are not being used for this project, but those features are actually important contributors for TED Talk popularity. Consequently, future work would also be carried out to apply Topic Modeling algorithms to classify the topic class to a categorical value, the proposed algorithm is *Latent Dirichlet allocation* (LDA). [3]

# Acknowledgement & Reference

- This project is a course project for Data Mining at Hood College of Frederick, Maryland, during the Fall Semester of 2018. Firstly, I would like to thank Dr. Liu, Xianlian for giving us the fantastic opportunity that we can select our project topic by ourselves and the great guidance given through the entire semester. In addition, I would like to give thanks to my classmates and friends for providing suggestions on this project. Finally, I would like to express my deepest thanks and sincere appreciation to my family, girlfriend, and colleagues for their love, understanding, and support.
  - [1] Wikipedia. 2018. Wikipedia: TED (conference). Retrieved from <a href="https://www.wikipedia.org">https://www.wikipedia.org</a>.
  - [2] L. Jiang, Z. Cai & D. Wang. 2010. Improving Naive Bayes for Classification, International Journal of Computers and Applications, 32:3, 328-332, DOI: 10.2316/Journal.202.2010.3.202-2747
- [3] Machine Learning Plus. 2018. Topic Modeling with Gensim (Python). Retrieved December 16, 2018 from <a href="https://www.machinelearningplus.com/nlp/topic-modeling-gensim-python/">https://www.machinelearningplus.com/nlp/topic-modeling-gensim-python/</a>