# Prediction of YouTube Video Trend for BigFame Co.

**Group 16** 

Qiaoling Chen: e0983208

Lyu Xinyu: e0983260

Xie Siteng: e0983376

Yashna Gogineni: e0983387

Zhao Rudan: e0983390





### Content



Biz & Data

Models

Conclusion

**Background** 

Data Processing

**Data Insight** 

**Topic Modeling** 

**Time Series** 

Clustering

Sentiment Analysis **Limitations** 

**Business Insight** 

Project Management

# **BIGFAME** — Company

- BigFame:Multi-ChannelNetwork (MCN) Co.
- Org. Structure:

Marketing Department

Through various channels, activities, to promote videos and activities, business negotiations and cooperation with advertisers and so on

Content Department

Responsible for the creative design and development of short video and live content Frequently analyze and study popular video works on various platforms, and keep unwith the hot spots

Operation Department

Responsible for the overall operation planning of the company, analyze and summarize all operational indicators, and formulate reasonable operational objectives and plans.

4 Host Management Department

Be responsible for the discovery, recruitment and training of influencers to create star Internet celebrity anchors.

- 5 HR Department
- 6 Administration Department
- 7 Financial Department
- 8 Technology Department

# Challenges

Revenue generate from Key Opinion
 Leader(KOL) in different categories

 Challenges we can deal with by this data analysis project:

#### **Host Management**

Lack of quality talent.

#### **Content Creation**

Lack of quality content ideas.

The content is seriously homogenized

### **Competitive Pressures**

KOL erosion.

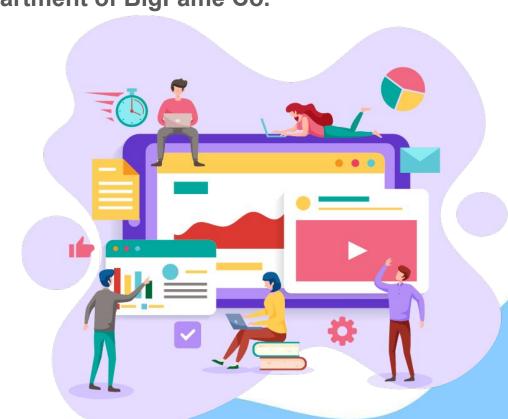
The problem of professionalization and contract.

#### **Platform Pressure**

Content operation investment in different categories of platforms

### RESPECT — Team

- Respect:
  - Data analysis team for Content Operations Department of BigFame Co.
- **Team members:** 
  - **Qiaoling Chen**
  - Lyu Xinyu
  - Xie Siteng
  - Yashna Gogineni
  - **Zhao Rudan**







# **Business Objectives**



Become a much more strategic generator of value and potentially also a powerful differentiator



Looking at video popularity trends and make well preparation for the hot category next period



Make agile, flexible, and resilient operational decisions based on full range of clean data

# **Technical Objectives**

- Classifying the sentiment
- Machine learning model such as clustering to identify different viewers group

3 Strategic

- Mathematical model such as TOPSIS to calculate videos' popular
- Time series model to predict the future popularity trend for all categories

Predictive

- Preparing data including integration
- Generating tags for each video based on the title and comments
- Building an interactive dashboard

1 Operational



# Data Pre-pocessing

### **Data Pre-processing**



### Data

#### Remove

Missing value • Duplicated value

**Outliers** 

**Generate Pop Score:** 

**TOPSIS Method Based on Entropy Weight** 

**Normalization** 

```
# Def: Standardize the matrix
def Y_ij(data1):
    for i in data1.columns:
      for j in range(n+1):
          if i == str(f'X{j}negative'): # negative
              data1[i]=(np.max(data1[i])-data1[i])/(np.max(data1[i])-np.min(data1[i]))
              data1[i]=(data1[i]-np.min(data1[i]))/(np.max(data1[i])-np.min(data1[i]))
   return data1
```

```
# Def: Caculate the entropy
# Build empty matrix
None ij = [[None] * n for i in range(m)]
def E j(data2):
    data2 = np.array(data2)
    E = np.array(None ij)
    for i in range(m):
        for j in range(n):
            if data2[i][j] == 0:
               e_ij = 0.0
                p_ij = data2[i][j] / data2.sum(axis=0)[j]
                e_{ij} = (-1 / np.log(m)) * p_{ij} * np.log(p_{ij})
            E[i][j] = e_ij
    E j=E.sum(axis=0)
    return E j
```

```
Likes
                                                                     Likes
                                                                                       0.301295
1600
                            Comment
                                                                     Comments
                                                                                       0.294011
1400
1200
                                        Views
        1600
                                                                    Views
                                                                                       0.404695
1000
        1400
800
        1200
                   1600
600
        1000
                   1400
         800
                   1200
         600
```

200

150

100

50

Pop Score

```
# Calculate the optimum and the worst
Imax j = Z ij.max(axis=0)
Imin j = Z ij.min(axis=0)
Dmax ij = np.array(None ij)
Dmin ij = np.array(None ij)
for i in range(m):
    for j in range(n):
        Dmax ij[i][j] = (Imax j[j] - Z ij[i][j]) ** 2
        Dmin_{ij}[i][j] = (Imin_{j}[j] - Z_{ij}[i][j]) ** 2
```

```
# Calculate the pop score
Dmax i = Dmax ij.sum(axis=1)**0.5
Dmin i = Dmin ij.sum(axis=1)**0.5
S i = Dmin i/(Dmin i + Dmax i)
S i= pd.Series(S i, index=video data.index, name='pop score')
```

200

#### Add **Feature**

### **Data Pre-processing**



#### Replace

- Contractions
- Capital → Lowercase

#### Remove

- Non ASCII Words
- Punctuation
- Stopwords

#### Lemmatize

#### **Tokens**

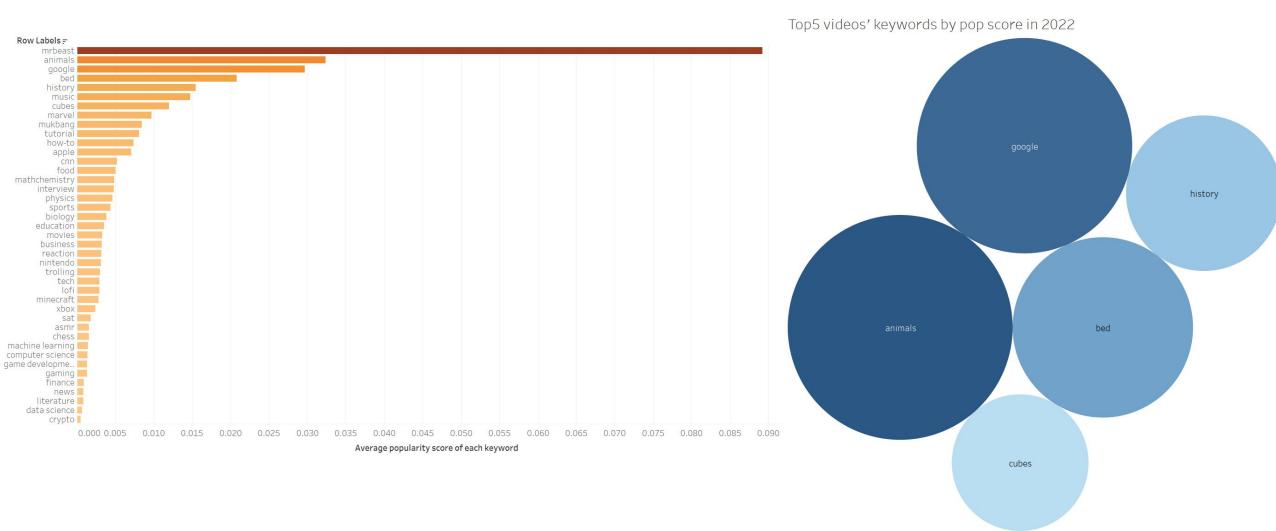
```
def replace contractions(text):
    """Replace contractions in string of text"""
   return contractions, fix(text)
def remove_non_ascii(words):
    """Remove non-ASCII characters from list of tokenized words"""
   new_words = []
   for word in words:
        new_word = unicodedata.normalize('NFKD', word).encode('ascii', 'ignore').decode('utf-8', 'ignore')
       new_words.append(new_word)
   return new words
                                                          def remove punctuation (words):
                                                               """Remove punctuation from list of tokenized words"""
def to lowercase (words):
                                                              new words = []
                                                               for word in words:
    """Convert all characters to lowercase from list of to
                                                                  new_word = re. sub(r'[^\w\s]', '', word)
   new words = []
                                                                  new_word = re. sub(r'(.)\1+', r'\1\1', new_word)
   for word in words:
                                                                   if new word != '':
                                                                      new_words.append(new_word)
       new word = word. lower()
                                                               return new words
       new_words.append(new_word)
   return new_words
                                                          def remove_stopwords(words):
                                                               """Remove stop words from list of tokenized words"""
                                                               new words = []
                                                               sw = stopwords, words('english')
                                                               for word in words:
                                                                   if word not in sw:
                                                                      new words, append (word)
                                                              return new_words
                                                          def lemmatize verbs(words):
                                                               """Lemmatize verbs in list of tokenized words"""
                                                              lemmatizer = WordNetLemmatizer()
                                                               1emmas = []
                                                               for word in words:
                                                                  lemma = lemmatizer.lemmatize(word, pos='v')
                                                                  1emmas. append (1emma)
                                                               return lemmas
```



# Data Insight

## **Data Insight**







# **Topic Modeling**

### **Topic Modeling**

# **Preparation for Clustering & Sentiment Analysis**

- Each video has 10 comments
- Topic Extraction for Each Video

```
Textrank summa Keywords
```

```
from summa import keywords

def generate_key(comm):

    TR_keywords = keywords.keywords(comm, scores=True)
    com_key=[]
    for x,y in TR_keywords:
        com_key.append(x)
    return com_key

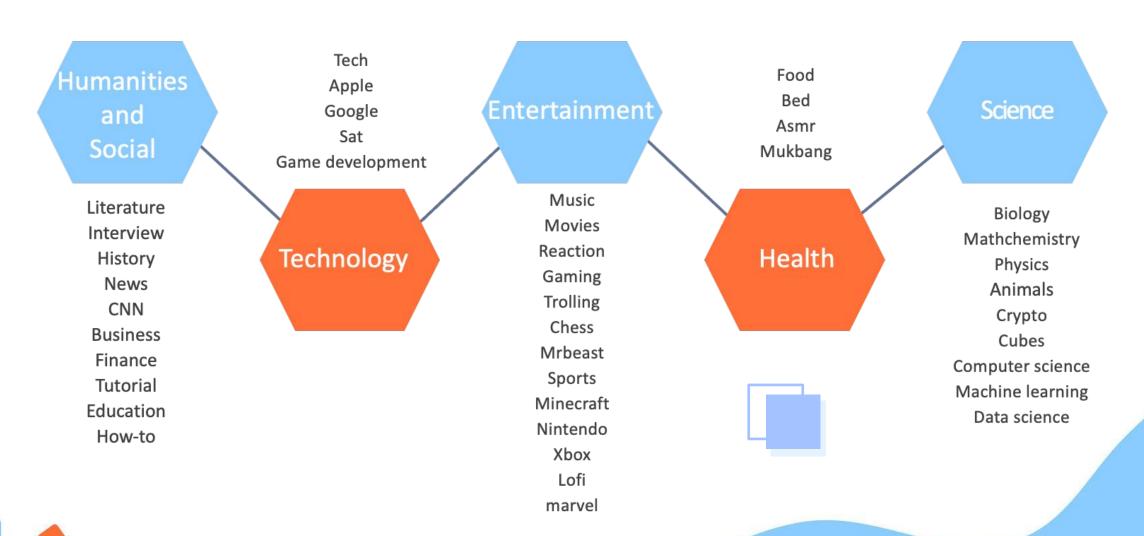
df_comment['Comment_keyword']=df_comment['Comment_tag'].apply(generate_key)
# for i range(len(df_comment)):
# df_comment['Comment'][i]=df_comment['Comment'][i]
```

```
df_comment. to_csv('data/comment_keyword_true_tag.csv')
df_comment['Comment_keyword'][0:10]
```

```
Video ID
--ZIOdSbbNU
               [eating, eat, eats, fry, feel, smile, man, suf...
               [russia, til, n, wo, realize, ukraine, analyza...
--hxd1Cr0ag
--ixiTvpG8g
               [college, colleges, tuition, regulated, regula...
-64r1hcxtV4
               [life, inspiration, inspiring, inspire, health...
               [thankful, thank, thanks, time, help, helped, ...
-6IgkG5yZfo
-7hzaGya86g
               [good, like, mikey, dinner, food, eating, dump...
-8TnsjDRXUE
               [kid, amazing, know, knowing, guy, time, work,...
-9h jdvULDvc
               [reading, prediction, th, nd, wa, writing, tes...
-Cr69sGnrk8
               [doe, u, fan, thank, thanks, sccoby, support, ...
-D4S6Tpn044
                                                       getbonus
```

Name: Comment\_keyword, dtype: object





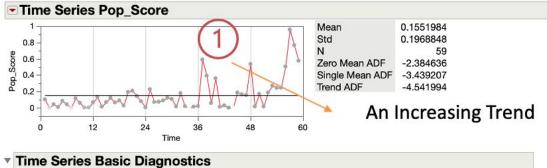
### ENTERTAINMANT

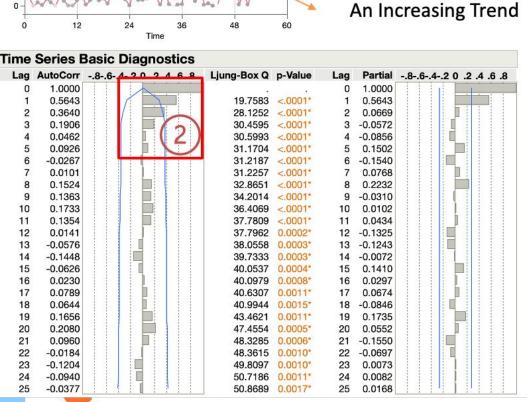
## **Time Series Analysis**



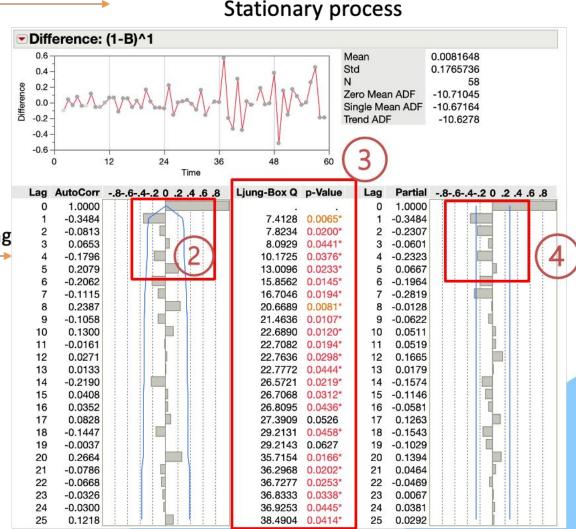
#### Non- stationary process







One Regular Differencing



### ENTERTAINMANT

# Time Series Analysis

Model	DF	Variance	AIC	SBC	RSquare	-2LogLH	MAPE	MAE
ARIMA(1, 1, 1)	55	0.024106	-45.44431	-39.26298	0.404	-51.44431	692.7482	0.105634
AR(1)	57	0.026241	-44.93511	-40.78003	0.345	-48.93511	564.0437	0.10742
IMA(1, 1)	56	0.02596	-44.82172	-40.70083	0.362	-48.82172	616.0846	0.108038
ARMA(1, 1)	56	0.025985	-44.41087	-38.17826	0.362	-50.41087	557.495	0.105578
AR(2)	56	0.026079	-44.25493	-38.02232	0.36	-50.25493	556.4575	0.105483
IMA(1, 2)	55	0.02572	-44.20452	-38.02319	0.379	-50.20452	651.6136	0.104614
ARIMA(2, 1, 1)	54	0.024432	-43.93216	-35.69039	0.407	-51.93216	679.9363	0.105215
Seasonal ARIMA(1, 1, 1)(0, 0, 1)12	54	0.024561	-43.45321	-35.21144	0.404	-51.45321	693.1448	0.10542
Seasonal ARIMA(1, 1, 1)(1, 0, 0)12	54	0.024559	-43.45152	-35.20975	0.404	-51.45152	693.0725	0.105457
ARIMA(1, 1, 2)	54	0.026167	-42.35606	-34.11429	0.38	-50.35606	609.5492	0.105772

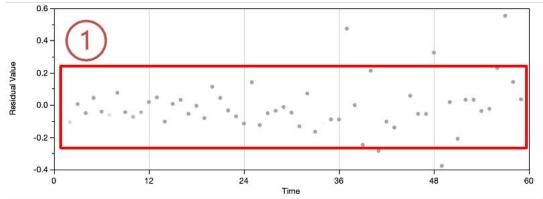
- Diagnostic Checking---Model Selection Criteria
- ✓ Akaike Information Criterion (AIC) = -2 ln(L) + 2k
- ✓ Schwartz Bayesian Criterion (SBC) = -2 In(L) + k In(n)

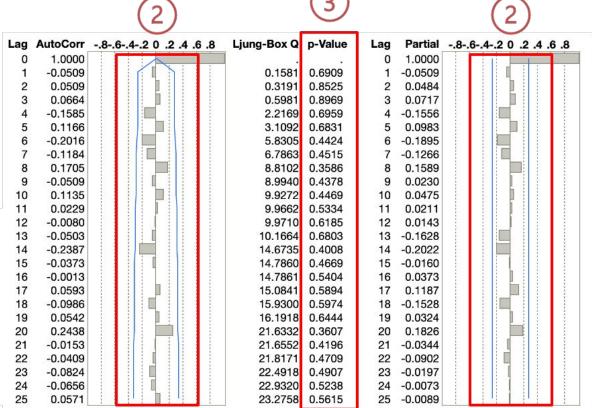


•



- ✓ Check the Residual Values
- ✓ Check the ACF and PACF
- ✓ Check the Result of Protmanteau Test

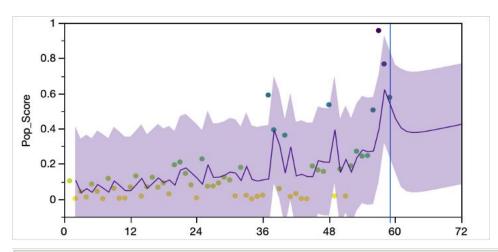




### ENTERTAINMANT

# Time Series Analysis





- ARIMA(1, 1, 1) Model Summary
- ✓ Stable and Invertible
- ✓ Significance of Parameters (at 95% confidence level)

$$(1 - \phi_1 B)(1 - B)y_t = \phi_0 + (1 - \theta_1 B)\varepsilon_t$$

### Result Analysis

✓ In general, the popularity of the topic of ENTERTAINMENT shows an upward trend.

#### ■ Model: ARIMA(1, 1, 1)

Model Summar	У
--------------	---

55 Stable 1.32581033 Invertible Yes Sum of Squared Errors 0.02410564 Variance Estimate Standard Deviation 0.15525992 Akaike's 'A' Information Criterion -45.444309 Schwarz's Bayesian Criterion -39.26298 **RSquare** 0.40369567 RSquare Adj 0.38201188 MAPE 692.74816 MAE 0.10563359

-51.444309

2

#### ▼ Parameter Estimates

-2LogLikelihood

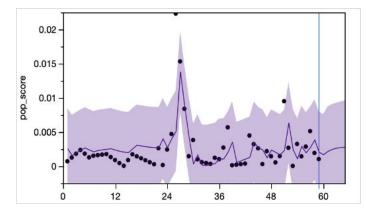
Term	Lag	Estimate	Std Error	t Ratio	Prob> t	Constant Estimate	Mu
AR1	1	0.4617515	0.1225077	3.77	0.0004*	0.00341949	0.006353
MA1	1	1.0000000	0.0735173	13.60	<.0001*	**************************************	
Intercept	0	0.0063530	0.0020585	3.09	0.0032*		



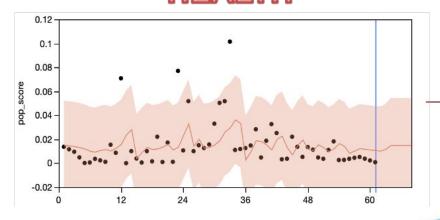
### **TECHNOLOGY**

- Stationary Process
- ARMA(2, 1)

- Stable trend but with fluctuates in September and October from 2019.
- Focus on the technology innovation, especially topics like APPLE.







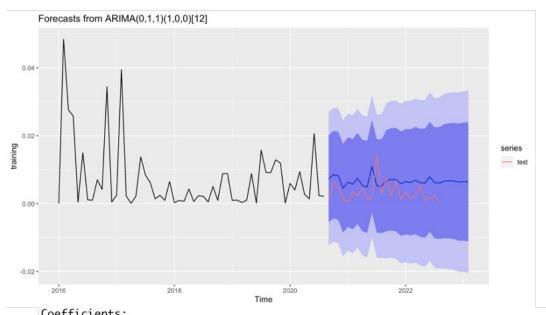
 Encourage innovation to attract more followers.

 Showing a stable trend without seasonal variation.

- Stationary Process
- MA(2)



### SCIENCE



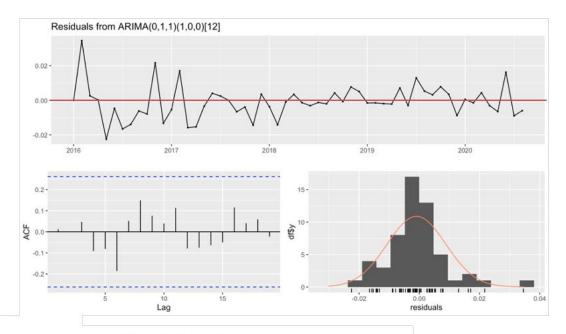


ma1 sar1 -0.8757 0.3137 0.0683 0.1695

sigma^2 = 9.978e-05: log likelihood = 175.04 AIC=-344.08 AICc=-343.61 BIC=-338.06

Training set error measures:

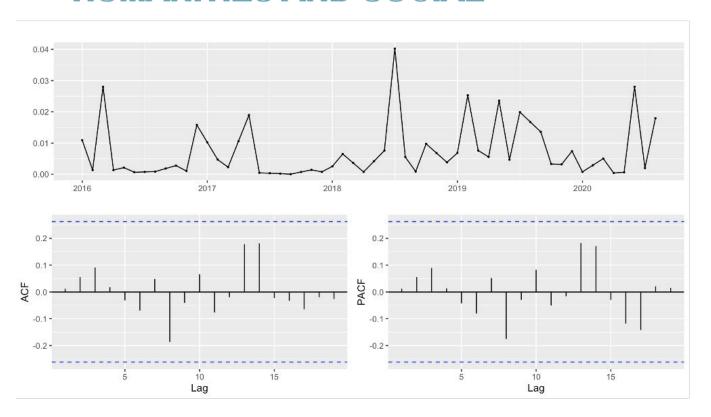
**RMSE** MAE MASE ACF1 Training set -0.0009385169 0.009717594 0.006951047 -764.5449 803.7078 0.8912433 0.01250267



Ljung-Box test

data: Residuals from ARIMA(0,1,1)(1,0,0)[12] $Q^* = 6.4604$ , df = 9, p-value = 0.6931

### HUMANITIES AND SOCIAL

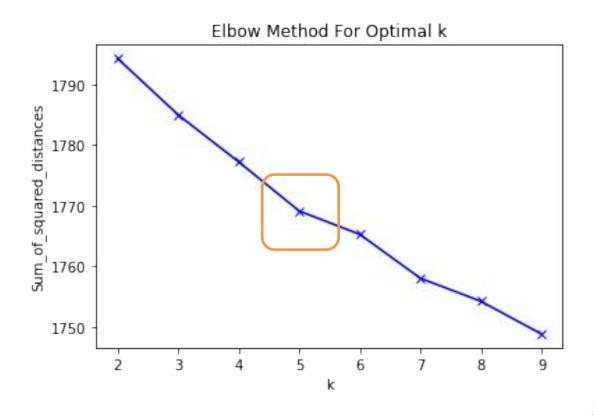


White Noise Series



# **Elbow Method**

determining the number of clusters in a data set



- luster 0: game play playing wa like player gaming played look make
- luster 1: song que music la love listening en se listen mix
- luster 2: love loved watching like watch wa look make man great
- luster 3: thank thanks learning learn wa help helpful work great like
- luster 4: wa like people time ha thing make great know guy
  - 1. GAME
  - 2. MUSIC
  - 3. LEARNING
  - 4. TECHNICAL
  - 5. EATING

#### Cluster 0 LR accuracy: 0.8888888888888888 Dummy classifier accuracy: 0.8739316239316239

tually making pol played new team :onsole - <sup>hit</sup>chess rele looking development gaming .ay goal excited hardware monitor

# 

```
música
     listen
            que son
                       listening
      mais
como
      esse
                        com
        music
                         amo
playlist
                      por
            um
               não
                         para
   lofimuito
                               tan
```

# Loving Watched respect asmr cooking feel animal food makeup movie eats channel sub look of eating eat respect asmr cooking feel animal food makeup movie eats channel talk eat respect asmr channel watching talk eat respect asmr channel man watching talk eat respect asmr channel man watching talk eat respect asmr channel watching talk eat respect asmr channel watching talk eat respect asmr cooking feel animal food makeup watching talk eat respect asmr channel watching talk eat respect asmr channel watching talk eat respect asmr channel watching watching talk eat respect asmr channel watching talk eat respect asmr channel watching talk eat respect asmr channel watching watching talk eat respect as man watching talk eat

#### Cluster 3 LR accuracy: 0.9145299145299145 Dummy classifier accuracy: 0.811965811965812

helpful question class with the light work of learned started the learn thanks learning explained helped explained teach thanks

#### Cluster 4 LR accuracy: 0.9294871794871795 Dummy classifier accuracy: 0.5769230769230769

by cubing trading proud live level by level by level by story achievement people cubing achievement achievement apple right amarket apple right apple google





#### Problem Statement

The aim of the problem is to predict the sentiment of the comment

Positive - 2

Neutral - 1

Negative - 0

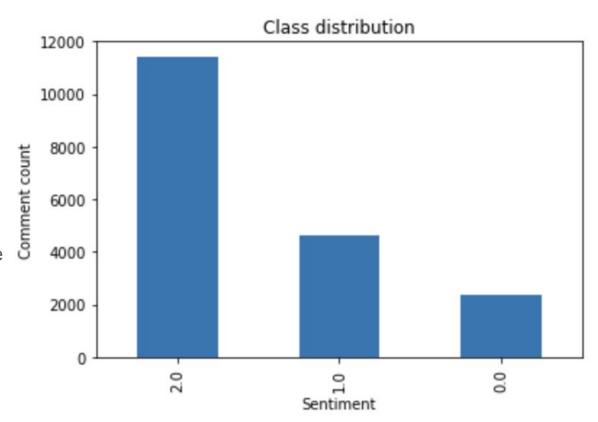
#### Imbalanced Class distribution

The bar chart shows the class distribution of the comments in the data set

#### · F1 score

Due to the imbalanced distribution,we chose F1 score as a evaluation method

$$\frac{2(P*R)}{P+R}$$



### **Word Cloud**



SENTIMENT: 2.0

- Sentiment 2 / Positive Sentiment
- This has words like 'Thank', 'Love', 'Great', 'much'
- Sentiment 1 / Neutral Sentiment
- This has words like 'see', 'think'
- Sentiment 0 / Negative Sentiment

This has words like 'never', 'need', 'even'



SENTIMENT: 1.0



SENTIMENT: 0.0

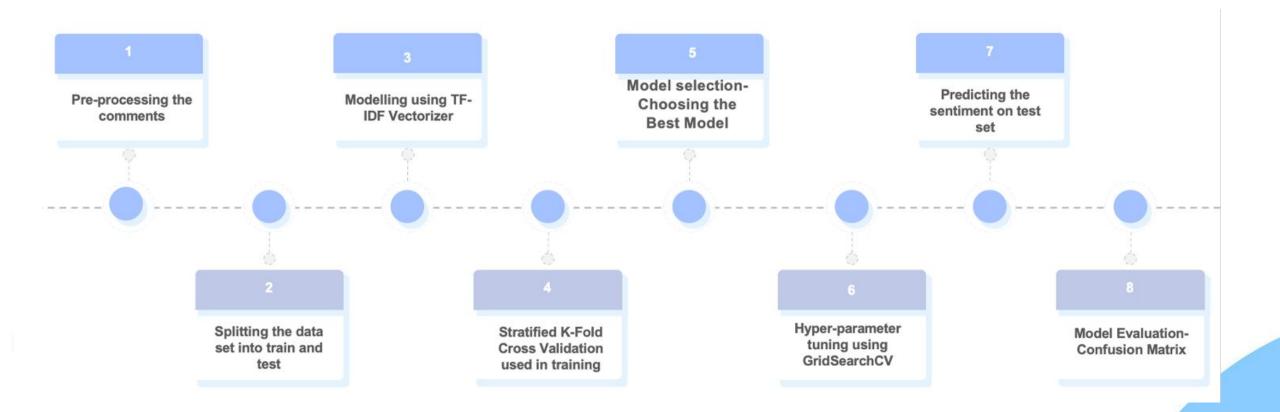






# Steps followed for Predictive Analytics Sentiment Prediction



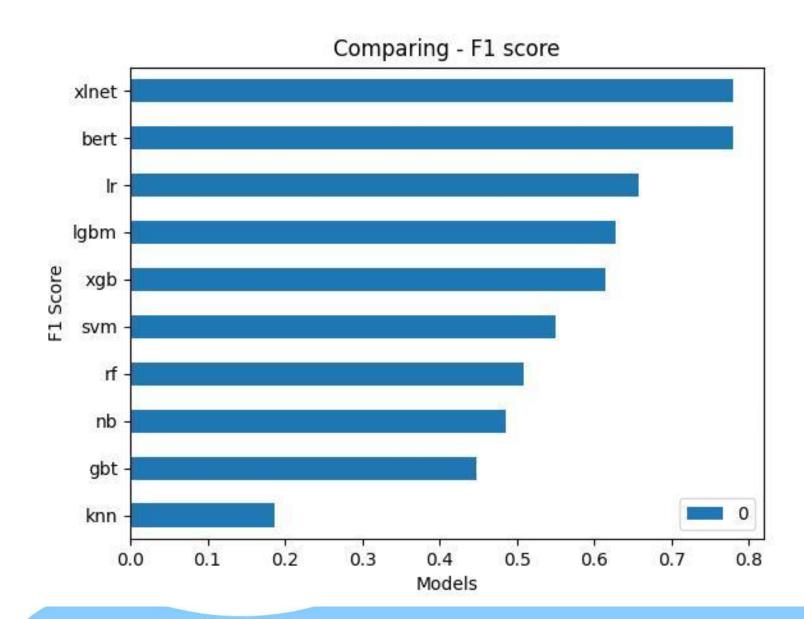


Machine Learning:

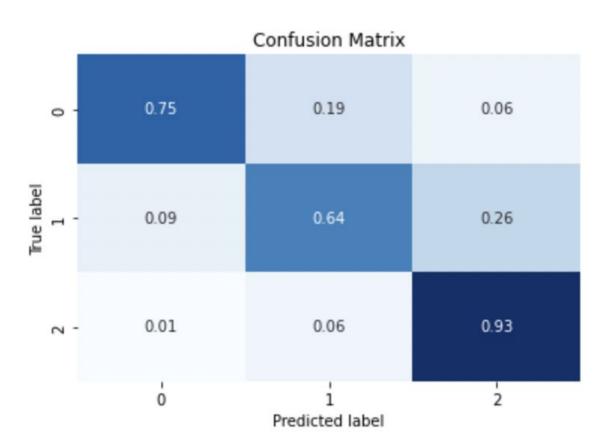
**Best: 65%** 

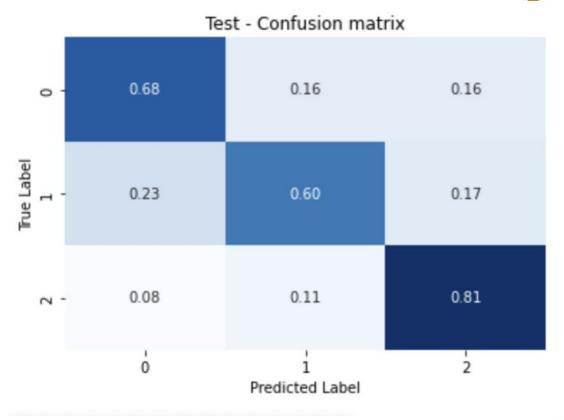
DeepLearning:

**Best:78%** 









**Deep Learning Model** 

Logistic Regression

Cost= Device\_price \* time\_consuming

	BERT	LR
Device_price	GPU= \$15/h	CPU= \$0.15/h
Time_consuming	60s	32ms
Cost	0.25	0.00001
F1-score(Reward)	0.78	0.67





- The classification Report on Train and Test set using Logistic Regression model
- •The model performs well both on training and testing data

Classification	on report on			
	precision	recall	f1-score	support
0.0	0.79	0.99	0.88	1753
1.0	0.75	0.93	0.83	3479
2.0	0.98	0.84	0.91	8574
accuracy			0.88	13806
macro avg	0.84	0.92	0.87	13806
weighted avg	0.90	0.88	0.88	13806
Classification	on report on	test set		
	precision	recall	f1-score	support
0.0	0.45	0.68	0.54	585
0.0 1.0	0.45 0.62	0.68 0.60	0.54 0.61	585 1159
1.0	0.62	0.60	0.61	1159
1.0	0.62	0.60	0.61	1159
1.0	0.62	0.60	0.61 0.84	1159 2858

CPU times: user 18 s, sys: 3.73 s, total: 21.7 s

Wall time: 3.58 s

### Limitation

1. Sentiment analysis have long tail class distribution

2. Some results of ARIMA have extremely low AIC or SBC

3. The iterative method of clustering makes results only locally optimal.



# **Business Insight**

- 1. Business Suggestions
- 2. Business Value

### **Business Suggestion**



# Focus on vloggers in Entertainment and Technology area.

- Help them catch the headlines
- Invest more on generating more creative and technical videos

### **Clustering Model**

# Differentiate viewers and catch up the attractive point they value on

- Games Experience improvement
- Music Relax
- Learning Math & Science
- Tech Financial Technique, Tech Market
- Eating Kind; animal

### **Business Value**



### Sentiment Analysis Model — LR Model

Lower Cost

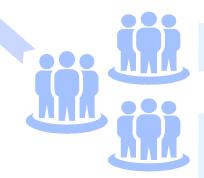
Higher Accuracy

Strong Explanatory



**Differentiate the KOLs** 

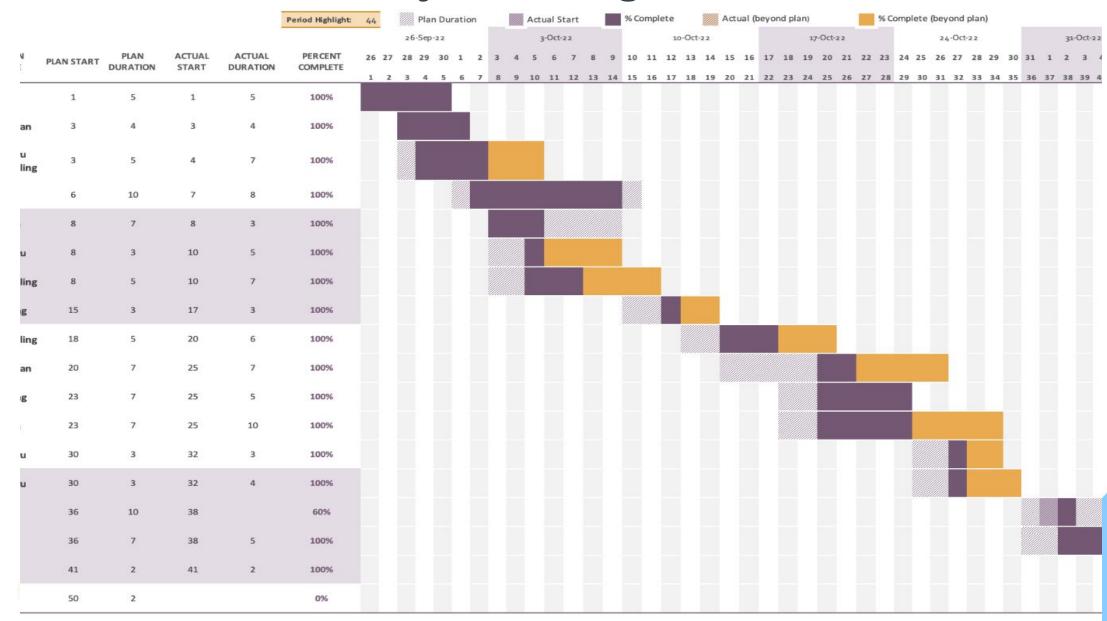
Large Scale of Comments



Follow the trend

Catch up the hottest topic

### **Project Management**







# Thank You