# YouTube Analyzer

Insun Lee
Kim Nguyen
Chao Zheng

# Background and Motivation

- Extract meaningful information from a large set of dataset
- Project could be applicable on other big data problem
  - Ex. Amazon's shopping list
    - Top k queries: Find the top k categories/items in which sold to customers; top k most popular item
    - Range queries : find all items in cateogries X with price within range [t1,t2];
    - User identification in recommendation patterns
    - Using PageRank algorithms. With high PageRank score, that means that item is related to other items in graph, thus has a high influence.

# Problem Formulation

| | |
|---|---|
| Video ID | 11 digit string. Unique for each videos |
| uploader | Video uploader's username |
| updated_date | Integer number of days |
| category | String type. Decided by video uploader |
| length | Integer type representing video length |
| views | Integer type. Number of views |
| rate | Float number of video rate |
| rating | Integer type of ratings |
| comments | Integer type of comments |
| related_videos_IDs | IDs of related videos |

- Input & Output
  - Input : .txt format



  - Output : .txt format



Top 10 Rated Videos

# Algorithm

MapReduce (Hadoop)

- MapReduce is the heart of Hadoop
- Hadoop is a highly scalable storage platform designed to process very large data sets across hundreds to thousands of computing nodes that operate in parallel. It provides a cost-effective storage solution for large data volumes with no format requirements.
- Map - converting data set to another set, broken down into tuples (key/values pairs)
- Reduce - Taking the output of Map as input and combines data tuples into smaller set of tuples.

# Algorithm

```java
import java.io.IOException;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;



public class CategoryMapper extends Mapper <LongWritable,Text,Text,IntWritable>{
```

```java
import java.io.IOException;
import java.util.Iterator;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;



public class CategoryReducer extends Reducer<Text,IntWritable,Text,IntWritable>
{
```

# Experimental Study output

## Top 10 Viewed Videos:

```
[wless-user-███████5045:out0.█████████eng$ hadoop fs -cat part-r-00000|sort -n
-k2 -r|head -n10
17/12/03 17:22:21 WARN util.NativeCodeLoader: Unable to load native-hadoop libra
ry for your platform... using builtin-java classes where applicable
Rg6463aqOyA        668112
bRPeEVpHiI8        331333
gdTkR2VbBbI        94775
vZKbUYl13UY        83176
uTPfTfgT-Vo        82363
8ud8Mcmxo1M        79464
G7R634P1sDo        78172
7D6aRxltTuA        76366
OnIlPmgL5XI        76026
GaOH0TWfkiE        74424
```

## Top 10 Rated Videos:

```
qcfLRxFq-cY        5
n2tFFtmty-U        5
RN_BuZWGPGw        5
R0hTG3X7sKw        5
Pwo6Yf-sXoY        5
HbOKRT_YRDc        5
Cis7jshNJbs        5
CXY-H7LI6sU        5
8Fg2SyfqCJ0        5
RknupcBUXHo        2
```

## Top 10 Categories:

```
36 People & Blogs
34 Film & Animation
32 Music          1
31 Entertainment
25 Comedy         1
15 News & Politics
11 Sports         1
 7 Pets & Animals
 6 Howto & DIY
 3 Gadgets & Games
 1 Travel & Places
 1 Autos & Vehicles
```

# PAGERANK ALGORITHM

PageRank

- rank websites in their search engine results
- Measure the importance of website pages



| steps | A | B | C | D | E |
|-------|-------|-------|-------|-------|------|
| 0 | 0.2 | 0.2 | 0.2 | 0.2 | 0.2 |
| 1 | 0.05 | 0.25 | 0.1 | 0.25 | 0.35 |
| 2 | 0.025 | 0.075 | 0.125 | 0.375 | 0.4 |

# PAGERANK ALGORITHM

| Top | Video ID | PageRank Score |
|---|---|---|
| 1 | rkvEuAtErwQ | 0.15230242746613892 |
| 2 | skn2nNHH8co | 0.15168760879025242 |
| 3 | wvPOOXdlO8U | 0.151531893028099 |
| 4 | 5c9OBqgmjzE | 0.15147045107345972 |
| 5 | CCz1kmfqL7g | 0.1514101936466493 |
| 6 | U4bk0UH1poQ | 0.15136395778938208 |
| 7 | uAIIMI0BUQ8 | 0.15134156041956306 |
| 8 | OIkUJnXLU74 | 0.15130953057948002 |
| 9 | nMl6B9m1rDw | 0.15122525021758051 |
| 10 | RG2kMfInkFA | 0.15121963744246494 |

# Sample Recommendation

- Generate Graph based on
  - Video ID, Category, and rate

- Steps
  - Initially select the category
  - Clean the dataset based on selection
  - Select top 10 video of that category and set it as the main
  - Related video will be sub-node
  - Distance will be the rate of each videos

# Conclusion & Future work

This project enables us to analyze very large data set on trending topics and interests of people through social networking forum.

We have done data analysis on a single YouTube dataset. Future work can be done by doing the same analysis on weekly gathered data to identify the current trending topics and area of interest of the mass.

Future work also includes extracting more important information using other attributes in dataset.

# Related Work

Youtube Graph Network Model and Analysis

- figure out which videos the viewer can potentially watch
  - Probability computation, PageRank algorithm

Understanding User Behavior in Large-Scale Video-on-Demand Systems

- Make Analysis on User's behavior with national media company
  - User Access Pattern over time (Daily, Hourly, Weekly)
    - Used Poisson Distribution
  - Changes of user interest over time
    - When new videos were uploaded, changes of user interest on existing videos

QUESTION?