



## Trend Videos

Statistical Inference | Master in Statistics for Data Science

Javier Muñoz Flores & Luis Ángel Rodríguez García

17-05-2022

# 1 Introduction

## 1.1 Topic and motivation

The aim of this project is to analyze what are the most favorite YouTube videos and their categories associated. These trends are measured with the following quantifiers: number of accumulated views, number of accumulated likes, number of accumulated dislikes and number of accumulated comments. Therefore, the higher the value of any of these parameters, the more likely to be in the top of the trendiest videos.

## 1.2 Description of the dataset

Data comes from Kaggle separate by countries and we have merged all of them into a single csv file. As it is commented in the previous website, the information was obtained through the YouTube API.

The csv file we uploaded at the beginning of the course contains 375,942 rows and 10 variables. It describes the 200 top trending YouTube videos per day for some different countries: United States of America, Canada, Mexico, Japan, South Korea, India, Russia, United Kingdom, France and Germany. Not all countries have the same interval of time. YouTube videos do not have the same period of time between each other either. Let's see the data grouped, aggregating the dates for each group:

video_id	country_id	trending_dates
n1WpP7iowLc	CA	17.14.11 17.15.11 17.16.11 17.17.11
0dBikQ4Mz1M	CA	17.14.11 17.15.11 17.16.11 17.17.11 17.18.11
5qpjK5DgCt4	CA	17.14.11 17.15.11 17.16.11 17.17.11
d380meD0W0M	CA	17.14.11 17.15.11 17.16.11 17.17.11
2Vv-BfVoq4g	CA	17.14.11 17.15.11
0yIWz1XEeyc	CA	17.14.11

Then it is clear that the record dates depend on the `video_id` and the `country_id`. For example, in the original data the interval times for Canada, US, Germany, FR, UK, India, South Korea, Mexico and Russia are from 14-11-2017 to 14-06-2018 and for Japan from 07-02-2018 to 14-06-2018.

Therefore, due to the fact that time-series is not allowed for the purpose of this project, we are going to consider the last registered date time observation for a specific `video_id` and `country_id`.

Besides that, we are applying a couple of tweaks to have our data ready for the analysis. First, we are going to modify the type of a few variables:

- `trending_date` from character to Date class
- `category_id` from integer to factor class
- `country_id` from integer to factor class
- `comments_disables` from character to factor class

Secondly, in order to have available a continuous variable, we are going to change the value of the variable `views` to be measured in miles. Lastly, we are filling in the missing values for the variable `comments_disables` (setting 0 when there are comments and 1 when the number of comments is zero) and removing the variables that we are not going to use afterwards such as `trending_date` and `tags`.

```
data <- youtube %>%
  mutate(
    trending_date = as.Date(trending_date, '%y.%d.%m'),
    category_id = as.factor(category_id),
    country_id = as.factor(country_id),
    comments_disabled = factor(ifelse(is.na(comments_disabled), NA,
```

```

) %>%
group_by(video_id, country_id) %>%
slice(which.max(trending_date)) %>%
select(-c(trending_date, tags)) %>%
data.frame
data$comments_disabled[is.na(data$comments_disabled)] <-
  ifelse(data[is.na(data$comments_disabled),]$comment_count==0, 1, 0)
data$views = data$views/1000

```

Our final dataset contains 207,148 rows.

### 1.3 Description of the population

Observing some blogs about the YouTube's figures for 2018, we can say that the total number of YouTube videos approximately was around 8 billion videos. A significantly large population size comparing with the sample size of our data that is just 207,148 observations, then the population size is totally different from the sample size.

The total number of videos in YouTube at the time where the last observation was recorded vary based on the country. Considering the number of active users per country we could approximate the total videos for each market, however this value would not be accurate at all. Let's obtain the number of observations we have for each country.

country_id	n
CA	24427
DE	29627
FR	30581
GB	3272
IN	16307
JP	12912
KR	15876
MX	33513
RU	34282
US	6351

We could see in the table above that UK and US have a smaller number of observations, nevertheless it does not mean that they have less videos, it could be that they have a large number of videos repeated (consecutive days appeared in the 200 trendiest videos).

We would like to extrapolate our results (top 200 trending videos per day) to the whole population of YouTube videos. At the beginning just for the whole YouTube and later on segmented by country.

### 1.4 Description of the variables

The dataset has 8 dimensions, these variables are:

- **video\_id** (*qualitative nominal variable*) unique identifier for a YouTube video, it is a string of characters of length 11.
- **category\_id** (*qualitative nominal variable*) unique identifier for video's category. The mapping between categories and ids is the following one: Film & Animation (1), Autos & Vehicles (2), Music (10), Pets & Animals (15), Sports (17), Travel & Events (19), Gaming (20), People & Blogs (22), Comedy (23), Entertainment (24), News & Politics (25), Howto & Style (26), Education (27), Science & Technology (28), Nonprofits & Activism (29), Movies (30), Shows (43) and Trailers (44).

- **view** (*quantitative continuous variable*) number of views, measured in miles, up to the last record's date for a specific video and country.
- **likes** (*quantitative discrete variable*) number of likes up to the last record's date for a specific video and country
- **dislikes** (*quantitative discrete variable*) number of dislikes up to the last record's date for a specific video and country
- **comments\_count** (*quantitative discrete variable*) number of comments up to the last record's date for a specific video and country.
- **comments\_disabled** (*qualitative binary variable*) whether the comments are enabled or not, represented as 0 if it is false and 1 if it is true.
- **country\_id** (*qualitative nominal variable*) unique identifier for the country where the video is hosted. It consists a string of characters of length 2. The mapping between countries and ids is the following one: United States of America (US), Canada (CA), Mexico (MX), Japan (JP), South Korea (KR), India (IN), Russia (RU), United Kingdom (GB), France (FR) and Germany (DE).

As it is indicated in the guide project, we are going to reduce the number of factors for the categorical variables we have: **category\_id** and **country\_id**. For the first variable, we are going to merge those factors that are intimately related to, having:

- **Film**, that includes: Film & Animation (1), Comedy (23), Entertainment (24), Movies (30), Shows (43) and Trailers (44).
- **Music**, that is just Music (10).
- **Education**, that is formed by News & Politics (25), Howto & Style (26), Education (27), Science & Technology (28) and Nonprofits & Activism (29).
- **Leisure**, that includes: Autos & Vehicles (2), Pets & Animals (15), Sports (17), Travel & Events (19), Gaming (20) and People & Blogs (22).

The number of factors of the variable **country\_id** can be reduced using the continent associated instead of the country, therefore this variable will be renamed to **continent\_id**. Then, we will have three different ids:

- **America**, that contains United States of America (US), Canada (CA) and Mexico (MX).
- **Europe**, that includes Russia (RU), United Kingdom (GB), France (FR) and Germany (DE).
- **Asia**, that contains Japan (JP), South Korea (KR) and India (IN).

Notice that for this last case, we should add the value of the quantifiers for the same video and different countries in the same continent.

```
film_categories <- c(1, 23, 24, 30, 43, 44)
education_categories <- c(25, 26, 27, 28, 29)
leisure_categories <- c(2, 15, 17, 19, 29, 22)

new_category_ids <-
  ifelse(data$category_id %in% film_categories, "Film",
        ifelse(data$category_id %in% education_categories, "Education",
              ifelse(data$category_id %in% leisure_categories, "Leisure", "Music")))

new_continent_ids <-
  ifelse(data$country_id %in% c("US", "CA", "MX"), "America",
        ifelse(data$country_id %in% c("RU", "GB", "FR", "DE"), "Europe", "Asia"))
final_data <- data %>%
  rename(continent_id = country_id) %>%
  mutate(
    category_id = factor(new_category_ids),
```

```

continent_id = factor(new_continent_ids)
) %>%
group_by(video_id, continent_id, category_id) %>%
summarize(comments_disabled=comments_disabled[1],
          across(where(is.numeric), sum), .groups="keep") %>%
data.frame

```

After the previous transformation, we have a dataset with 195,582 observations.

In the table 1 we have summarized some important statistics related to the measures of centrality, variability and shape for the quantitative variables within the dataset.

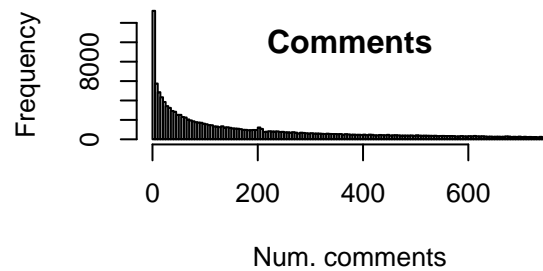
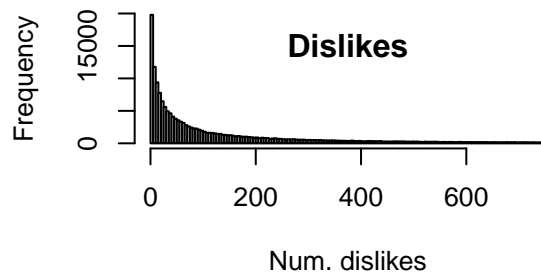
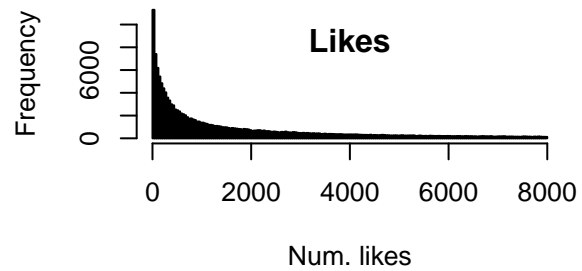
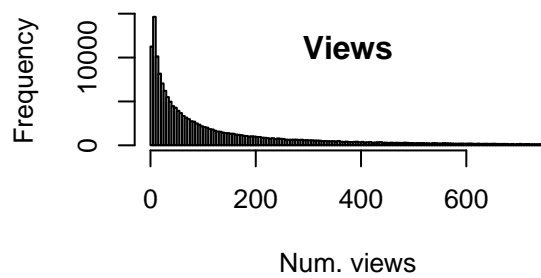
Table 1: Stats Table

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	sd	skewness	kurtosis
<b>views</b>	0.2	23.2	82.2	561.1	279.5	516003.4	4617.47	46.85	3290.39
<b>likes</b>	0	341	1529	17182	6420	19759060	149734.44	46.63	3537.62
<b>dislikes</b>	0	21	84	1037	329	5462288	20691.41	178.41	40459.19
<b>comment</b>	0	54	243	2141	919	4259927	23739.88	83.93	10519.16

Once we have analyzed the table above, we can conclude that all of the YouTube videos within the dataset at least contain 200 views, not occurring the same for the other three variables where data can have 0 likes, 0 dislikes or 0 comments.

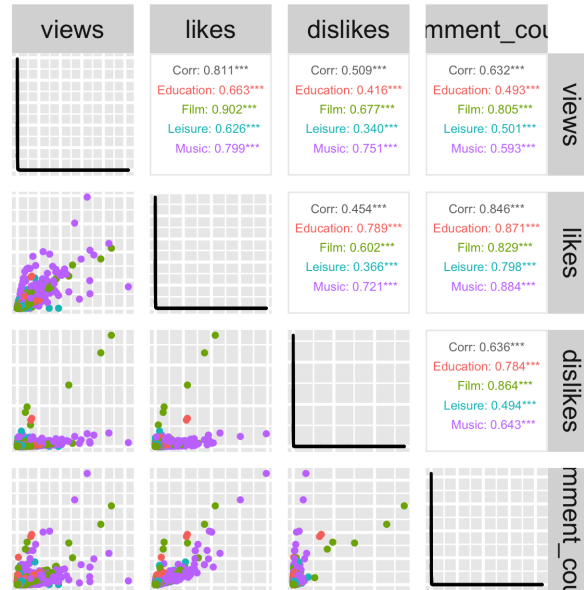
Taking a look to the measures of centrality we obtain an idea of how it is the variability, for example, the variability is much larger for the variable **likes** than the variable **views** since the difference between the mean with respect the minimum value and the maximum value of the variable **views** is considerably smaller than the same calculation for the variable **likes**. In order to verify this assumption we just need to check the value of the standard deviation (a measure of variability), for the case we have just commented the value for the variable **views** is 4617.47 in contrast with the value of the variable **likes** that is 149734.44, which is considerably much larger.

The other other two statistics provide a measure of the shape of the distribution. In particular, the skewness gives an idea of how much shifted the distribution is with respect a Normal distribution. In our case, the table 1 shows that all variables have a positive skewness, meaning that those distributions are right-skewed. It is worthy to say that this insight can be approximately acquired looking at the mean and the 1st and the 3rd quantiles. If the mean is more or less in the middle between the quantiles, then it is more likely to have a skewness close to 0. In our case, the four variables have the mean outside the interval between the 1st and the third quantile. Lastly, it is the kurtosis statistic that measures how heavy are the tails of the distribution. As it is occurred in the skewness, the Normal distribution is taken as the golden standard and the kurtosis is compared with respect this distribution. Therefore, for those variables -**views**, **likes**, **dislikes** and **comments\_count**- with a value greater than 3, the distribution is *leptokurtic* (heavier tails than a Normal distribution). Let's check all of this visually thanks to the following plot.



In the plot above, we can see clearly the characteristic that was pointed by the kurtosis value, distributions are *leptokurtic*, having heavier tails. It is noticeable as well the fact that there is not any variable that follows a Normal distribution. We will see more about this in the following point.

After describing some distribution details of each of the variables, now it is time to see if they are related to each other. In the following pairs chart we can see the relationship between each quantitative variable with others.



All variables are positive correlated, having that **views** and **likes** are highly correlated (0.811) since the higher of views the more likely that video is liked by the viewers. This logic works also for the number of comments and the number of likes, that's the reason the correlation is high as well (0.846). Moreover, we have shown the correlation split by the different category ids and it is clear that this correlation value depends directly on the category of the video.

Once we have analyzed the quantitative variables, we can move forward and describe a little bit the qualitative variables. For doing this, we are using the frequency tables showed below.

category_id	Freq	continent_id	Freq	comments_disabled	Freq
Education	43023	America	60202	0	189632
Film	79334	Asia	44386	1	5950
Leisure	55391	Europe	90994		
Music	17834				

The three categorical variables are unbalanced, then there are not the same number of observations for each factor. As we could expect, there are more videos with the comments enabled. It is curious that there are more videos in Europe or America than in Asia, the reason might make sense since YouTube (belonging to Google) is a western company. In case of the number of videos by category, the one that has more videos is the **Film** category, followed by **Leisure**, **Education** and lastly **Music**.

## 2 Model selection

### 2.1 Selecting a continuous random variable

### 2.2 Estimate the model parameters

#### 2.2.1 Method of moments

#### 2.2.2 Maximum likelihood

## 3 One-sample inference

### 3.1 Estimators of population mean

In this section, we use again the continuous variable `views` in order to calculate two different estimators of the population mean. The estimators selected are the following ones:

$$\hat{\mu}_1 = \bar{X} \quad \text{and} \quad \hat{\mu}_2 = \frac{X_1 + X_n}{2}$$

Estimator	Value of estimations
Estimator 1	561.135807835077
Estimator 2	258001.778

Now, let's analyze the properties of both estimators. First, we will see if any of them is unbiased.

•

$$E[\hat{\mu}_1] = E[\bar{X}] = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n E[X_i] = E[X]$$

•

$$E[\hat{\mu}_2] = E\left[\frac{X_1 + X_n}{2}\right] = \frac{1}{2}E[X_1 + X_n] = E[X]$$

Hence both estimators are unbiased. Now, let's calculate the variance of each one.

•

$$V(\hat{\mu}_1) = V(\bar{X}) = V\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n V(X_i) = \frac{V(X)}{n}$$

•

$$V(\hat{\mu}_2) = V\left(\frac{X_1 + X_n}{2}\right) = \frac{1}{4}V(X_1 + X_n) = \frac{1}{2}V(X)$$

It can be seen that the variance of  $\hat{\mu}_2$  is higher than the variance of  $\hat{\mu}_1$ , as the latter one depends on  $n$  (the number of observations) which in our case is much larger than 2. Therefore, the estimator with smaller variance is  $\hat{\mu}_1$ .

Other property we can analyse is the consistency in squared mean. We know that if the bias and the variance tend to zero as  $n \rightarrow \infty$ . Therefore, since both estimators are unbiased but  $\hat{\mu}_2$  does not tend to zero in the limit, so the only consistent estimator is  $\hat{\mu}_1$ .



### 3.2 Estimation of the error of the estimators

As we have mentioned previously, the estimators selected are unbiased. Thus, we must estimate the coefficient of variation (CV). The coefficient of variation (CV) is defined as the ratio of the standard deviation  $\sigma$  to the mean  $\mu$ ,  $CV = \frac{\sigma}{\mu}$ . We calculate the CV for both estimators:

$$\begin{aligned} \bullet \quad CV_{\hat{\mu}_1} &= \frac{\sqrt{Var(\hat{\mu}_1)}}{E[\hat{\mu}_1]} = \frac{\sqrt{Var(X)}}{\sqrt{n}E[X]} = \\ \bullet \quad CV_{\hat{\mu}_1} &= \frac{\sqrt{Var(\hat{\mu}_1)}}{E[\hat{\mu}_1]} = \frac{\sqrt{Var(X)}}{\sqrt{2}E[X]} = \end{aligned}$$

### 3.3 Selection of a qualitative/categorical variable

Now, we consider a categorical variable instead of a continuous one. We are interested in estimating the proportion of videos (in the population) that belong to a specific category. In particular, we want to know what proportion of videos are related with the music, i.e. the videos that belong to the category *Music* of the variable `category_id`.

Since the variable `category_id` has four categories, it follows a multinomial distribution  $Y \sim M(n, \mathbf{p})$  being  $\mathbf{p}$  the vector of probabilities of belong to each category, i.e.  $\mathbf{p} = \{\pi_1, \pi_2, \pi_3, \pi_4\} = \{Education, Film, Leisure, Music\}$ . We can write  $Y = \sum_{i=1}^n X_i$ , where  $X_i$  follows a Categorical distribution, a generalization of the bernoulli distribution. We want to estimate the proportion  $\pi_4 = P(X = Music)$  by computing the expression:

$$\hat{\pi}_4 = \frac{1}{n} \sum_{i=1}^n X_i$$

```
## [1] "The estimated proportion is:"
```

```
## [1] 0.09118426
```

Thus, the 9% of the population are videos related with music.

### 3.4 Estimation of the variance of the estimator of proportion

Here we need to calculate the variance of the estimator  $\hat{\pi}_4$ :

$$V(\hat{\pi}_4) = V\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n V(X_i) = \frac{1}{n} V(X)$$

Thus, since we know that  $V(X)$  is equal to  $\pi_k(1 - \pi_k)$ , since the categorical distribution is a generalized form of the Bernoulli distribution, we can estimate the variance of  $\hat{\pi}_4$  as follows:

$$\hat{\sigma}_{\hat{\pi}_4} = \frac{1}{n} \hat{\pi}_k (1 - \hat{\pi}_k)$$