

Federal State Autonomous Educational Institution of High Professional Education National Research
University «Higher School of Economics»

Faculty of Computer Science School of Data Analysis and Artificial Intelligence

Report on the course ”Modern Methods of Data Analysis”.

«Data Sciences» Master program
Lecturer: Boris G. Mirkin.

Authors:

Cherny Artem
Ivanova Polina
Shvechikov Pavel

Dataset: Online News Popularity Data Set.

URL: <http://archive.ics.uci.edu/ml/datasets/Online+News+Popularity>

Explanation of choice

This dataset is based on recent articles, published in 2013–2015 years, so the data is actual up to now. Moreover, it is composed of different types of variables.

Data Set Information

- The articles were published by Mashable (www.mashable.com) and their content as the rights to reproduce it belongs to them. Hence, this dataset does not share the original content but some statistics associated with it. The original content could be publicly accessed and retrieved using the provided urls.
- Acquisition date: January 8, 2015
- The estimated relative performance values were estimated by the authors using a Random Forest classifier and a rolling windows as assessment method. See their article for more details on how the relative performance values were set.
- From initial data-set we chose 34 attributes and 10000 instances (instances were chosen randomly).

Attribute Information:

1. `url`: URL of the article (non-predictive)
2. `timedelta`: Days between the article publication and the dataset acquisition (non-predictive)
3. `n_tokens_title`: Number of words in the title
4. `n_tokens_content`: Number of words in the content
5. `n_unique_tokens`: Rate of unique words in the content
6. `num_hrefs`: Number of links
7. `num_self_hrefs`: Number of links to other articles published by Mashable
8. `num_imgs`: Number of images
9. `num_videos`: Number of videos
10. `average_token_length`: Average length of the words in the content
11. `num_keywords`: Number of keywords in the metadata
12. `data_channel_is_lifestyle`: Is data channel 'Lifestyle'?
13. `data_channel_is_entertainment`: Is data channel 'Entertainment'?
14. `data_channel_is_bus`: Is data channel 'Business'?
15. `data_channel_is_socmed`: Is data channel 'Social Media'?
16. `data_channel_is_tech`: Is data channel 'Tech'?
17. `data_channel_is_world`: Is data channel 'World'?
18. `self_reference_avg_sharess`: Avg. shares of referenced articles in Mashable
19. `weekday_is_monday`: Was the article published on a Monday?

20. `weekday_is_tuesday`: Was the article published on a Tuesday?
21. `weekday_is_wednesday`: Was the article published on a Wednesday?
22. `weekday_is_thursday`: Was the article published on a Thursday?
23. `weekday_is_friday`: Was the article published on a Friday?
24. `weekday_is_saturday`: Was the article published on a Saturday?
25. `weekday_is_sunday`: Was the article published on a Sunday?
26. `global_sentiment_polarity`: Text sentiment polarity
27. `global_rate_positive_words`: Rate of positive words in the content
28. `global_rate_negative_words`: Rate of negative words in the content
29. `rate_positive_words`: Rate of positive words among non-neutral
30. `avg_positive_polarity`: Avg. polarity of positive words
31. `avg_negative_polarity`: Avg. polarity of negative words
32. `title_sentiment_polarity`: Title polarity
33. `shares`: Number of shares (target)

1 Assignment 1

In the first task we consider the attribute called `shares`, which is the number of article shares in various social networks. Let construct a histogram and boxplot of chosen attribute (see Figure ??). From the histogram we can see, that the chosen attribute probably has a lognormal distribution, so we construct a new feature `log(shares)`. Histogram and boxplot for this new feature `log(shares)` are presented on Figure ??.

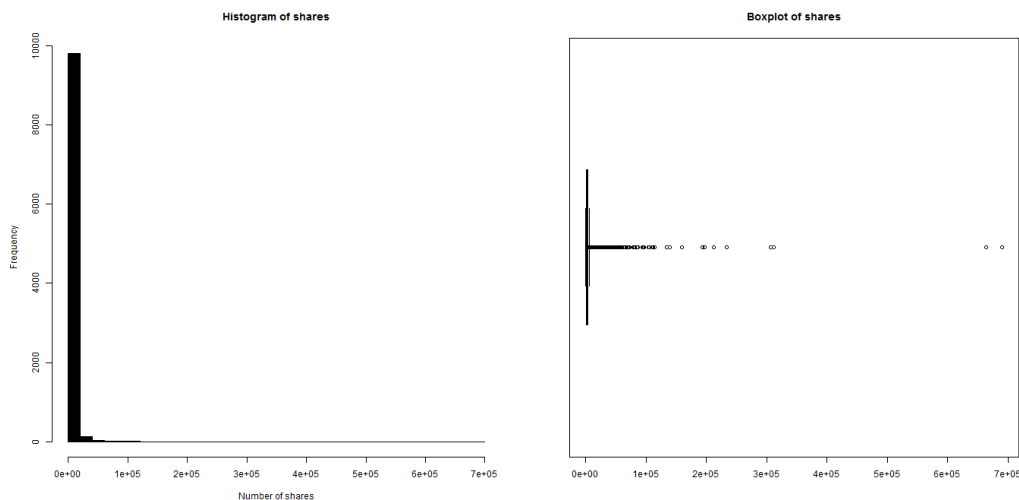


Figure 1

Below we will apply methods to the feature `log(shares)` instead of `shares`. As we can see in the Table ??, sample mean, median and mode of the `log(shares)` agree closely with each other, indicating that distribution is similar to symmetric. If we look at the same characteristics for the `shares`, we can see that the mean is significantly greater than the median and the mode. That could be easily explained with the histogram of `shares` (see Figure ??), which shows that the majority of entities have less than 20 000, and very few have more than 600 000 shares.

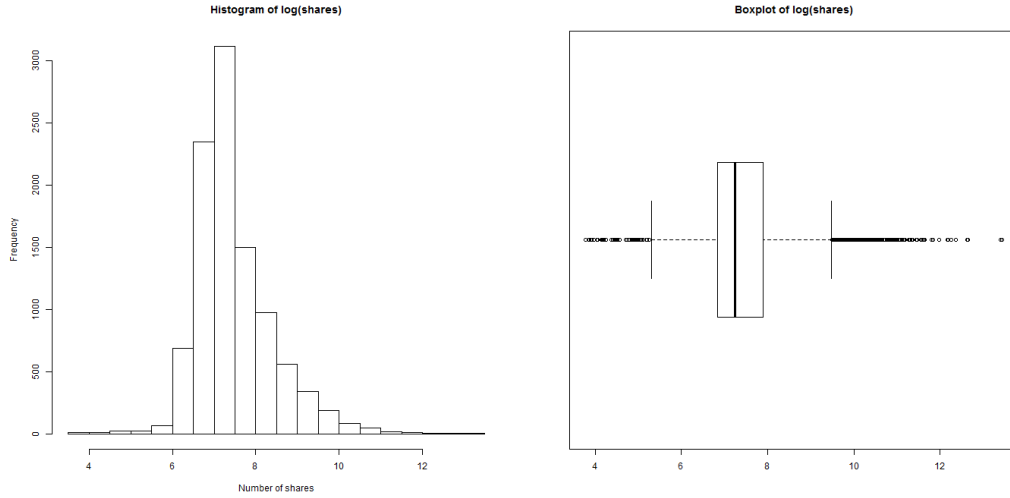


Figure 2

Table 1: Mean, median and mode for **shares** and **log(shares)**

shares	Mean	Median	Mode
	3374.5	1400	1100
log(shares)	Mean	Median	Mode
	7.45	7.24	7.0

Confidence intervals for the mean

The task is to find three confidence intervals (CI) for the mean of **log(shares)**. To do this, we make $N = 5000$ trials each of which consists of sampling with replacement from initial set of **log(shares)** and estimating the mean of population using that sampled data. The histogram of estimated means for the feature **log(shares)** is presented on the Figure ??, and computed 95% confidence intervals are shown in the Table 2 . The distribution of the means of **log(shares)** is very similar to normal distribution, so pivotal and non-pivotal intervals are similar too. It is worth to mention that statistic confidence interval is much more wider compared with any of the others.

Bootstrapping the mode and the median

The more the distribution resembles the power law distribution, the more appropriate is to choose median of the distribution as the center value. That is because the median is very stable against outliers. And pivotal or non-pivotal bootstrap methods can be applied to medians.

In case of mode it is hard to decide when the bootstrap technique is appropriate. The mode, in some sense, is not a smooth functional of the distribution. So the result will be most likely uninterpretable.

Table 2: 95% confidence intervals (CI's) for mean of **log(shares)**

Mean	7.45
Statistic CI	(5.94; 8.97)
Pivotal CI	(7.35; 7.56)
Nonpivotal CI	(7.33; 7.58)

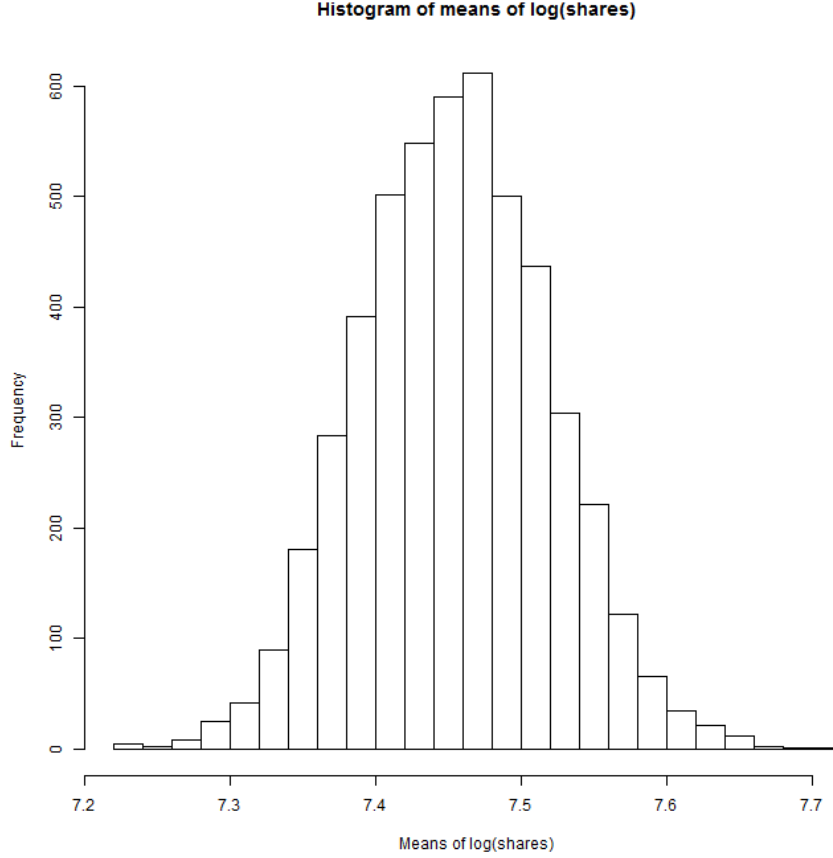


Figure 3: 30-bin histogram of means of $\log(\text{shares})$.

The $\log(\text{shares})$ is approximately distributed normally, so the values of mean, median and mode are close to each other. Therefore in this case bootstrap is likely to be a reliable method for computing confidence intervals of median and mode.

Histograms of sample medians and modes are presented in Figure ??, and respective confidence intervals are presented in Table ?. The distribution of the modes is far from the normal random variable distribution, so pivotal confidence interval could be incorrect. From the Table ?, as one could notice, it is evident that value of the mode is close to the left border of non-pivotal confidence interval.

It is much more interesting to compute the confidence interval for median on initial scale, that is not the median of logarithmic feature, but the median of initial **shares** feature. Since the distribution of **shares** is more similar to the power type distribution than to the normal one, it is better to choose the median as the central value due to the properties of median that were explained above. In fact, we can use either pivotal or non-pivotal approaches to estimate a median because of the next theorem.

Theorem (Median Theorem, [?]). *Let a sample of size $n = 2m + 1$ with n large be taken from an infinite population with a density function $f(\bar{x})$ that is nonzero at the population median $\tilde{\mu}$ and continuously differentiable in a neighborhood of $\tilde{\mu}$. The sampling distribution of the median is approximately normal with mean $\tilde{\mu}$ and variance $\frac{1}{8f(\tilde{\mu})^2m}$.*

The histogram of sample medians of **shares** is presented at Figure ?? (a), and computed confidence intervals are shown in the Table ?. The mode's distribution is obviously far from normal (so we don't examine pivotal CI), and, as well as in the case of $\log(\text{shares})$ modes, the value of population mode is close to the left border of non-pivotal confidence interval.

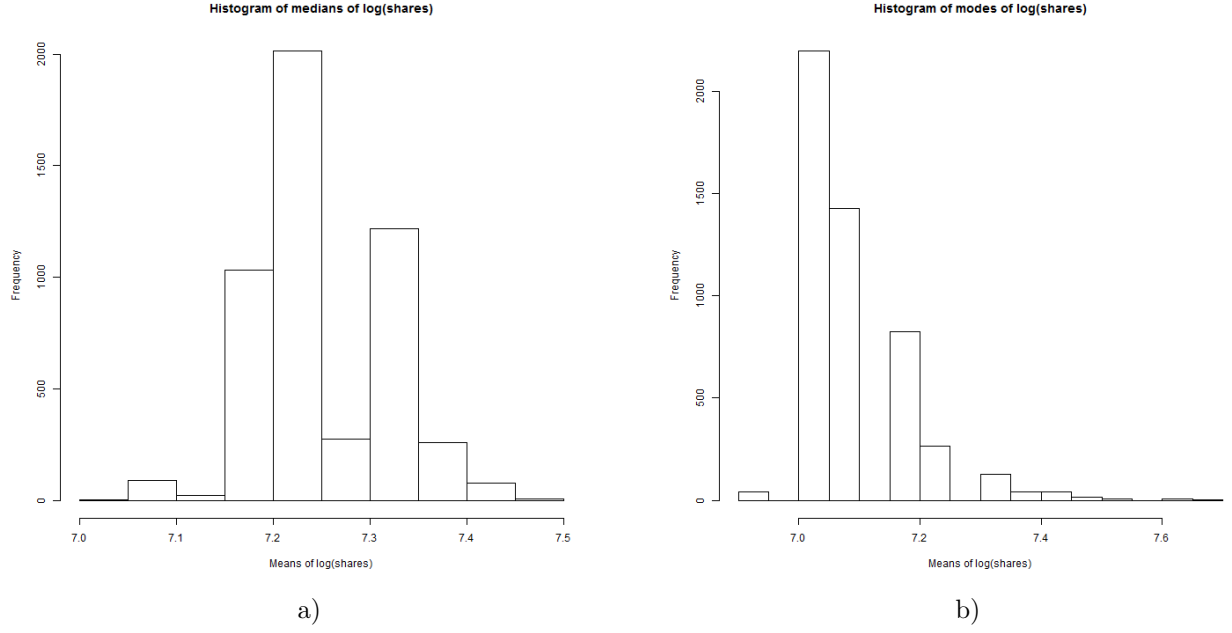


Figure 4: Histograms of medians (a) and modes (b) of $\log(\text{shares})$.

Table 3: 95% confidence intervals (CI's) for median and mode of $\log(\text{shares})$ and of shares

	$\log(\text{shares})$		shares	
	Median	Mode	Median	Mode
Value	7.24	7.0	1400	1100
Pivotal CI	(7.14; 7.36)	—	(1257.76; 1571.5)	—
Nonpivotal CI	(7.17; 7.38)	(7.0; 7.3)	(1250; 1600)	(1100; 1500)

Partitioning the population into two groups

We split our $\log(\text{shares})$ data according to the day, when the article was firstly published: weekday or weekend. In our dataset we have seven dummy variables, indicating the day of publishing a news: `weekday_is_monday`, `weekday_is_tuesday`, `weekday_is_wednesday`, `weekday_is_thursday`, `weekday_is_friday`, `weekday_is_saturday`, `weekday_is_sunday`. If we take entities, for which `weekday_is_saturday` or `weekday_is_sunday` are equal to 1, we will end up with the class `published_on_weekend`. All of the other entities will be considered as belonging to the class `published_on_workday`.

Histograms of the sample means in each of the classes are shown in figure ???. Each of the histograms closely resembles the density of normal distribution, therefore pivotal and non-pivotal bootstrap methods should compute the similar confidence intervals (CI's). 95% intervals for mean in each of the two groups are presented in Table ???.

The CI of the mean in both classes do not intersect with each other, so we can claim with 95% confidence that two means in these groups are different.

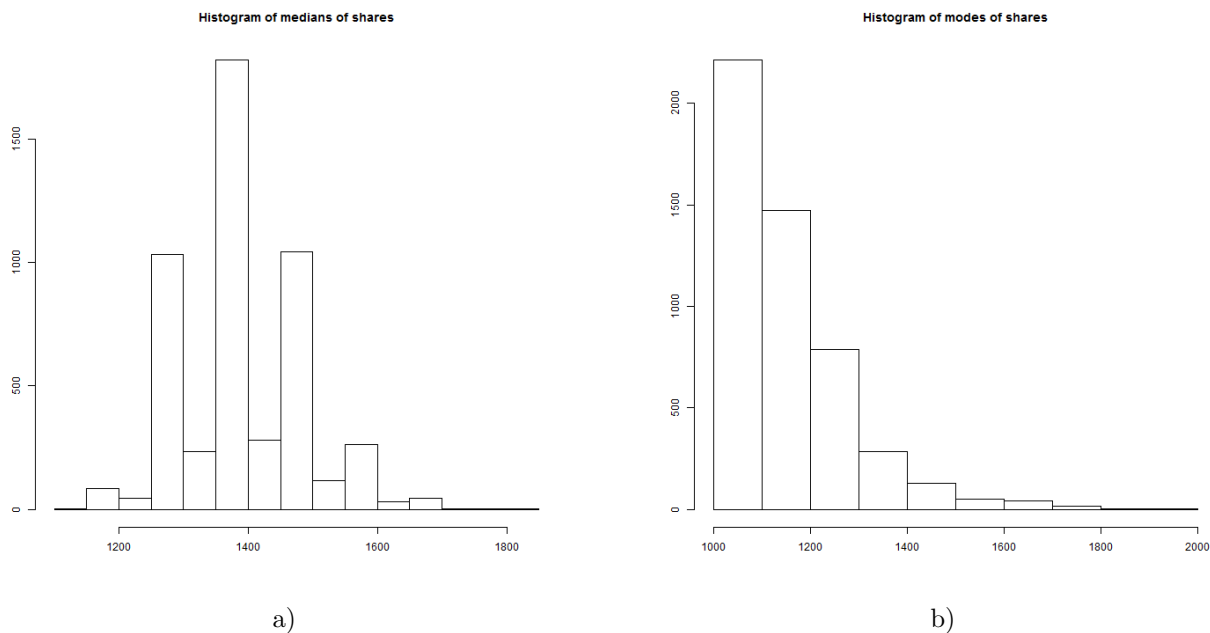


Figure 5: Histograms of sample medians (a) and sample modes (b) of `shares`.

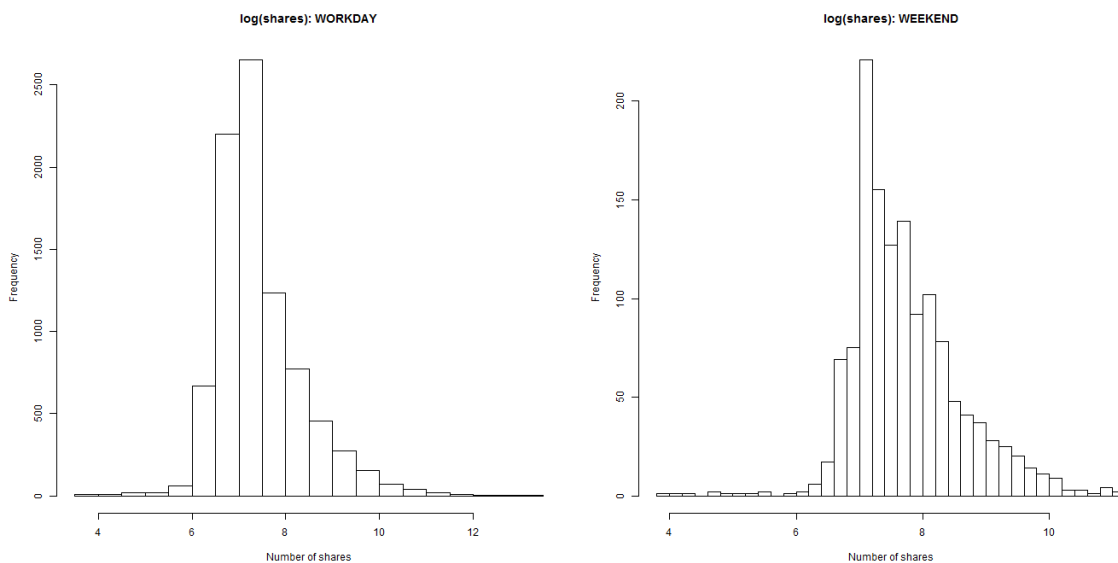


Figure 6: 30-bin histograms of $\log(\text{shares})$, grouped by the day, when the article was published: workday (left) or weekend (right)

2 Assignment 2. Linear regression.

2.1 Selection of features

First of all, let's take a look at all continuous features' scatterplot to identify which of them are linear dependent. Some features have nearly log-normal distributions, so, for more accurate and reliable linear regression we will logarithm these features.

As we can see from the scatterplot, the majority of pairs are not linear dependent. Fortunately, `global_sentiment_polarity` and `rate_positive_words` are linear dependent and we can easily understand why: they measure practically the same characteristic. The first feature is normalized from 0 to 1 (in our data from 0 to 0.7), the second one takes

Table 4: Confidence intervals of mean of $\log(\text{shares})$ feature, grouped by a day, when the article was published: weekday or weekend

	Workday	Weekend
Number of variables	8660	1340
Mean	7.41	7.73
Pivotal CI	(7.3; 7.5)	(7.63; 7.83)
Non-pivotal CI	(7.28; 7.54)	(7.6; 7.85)

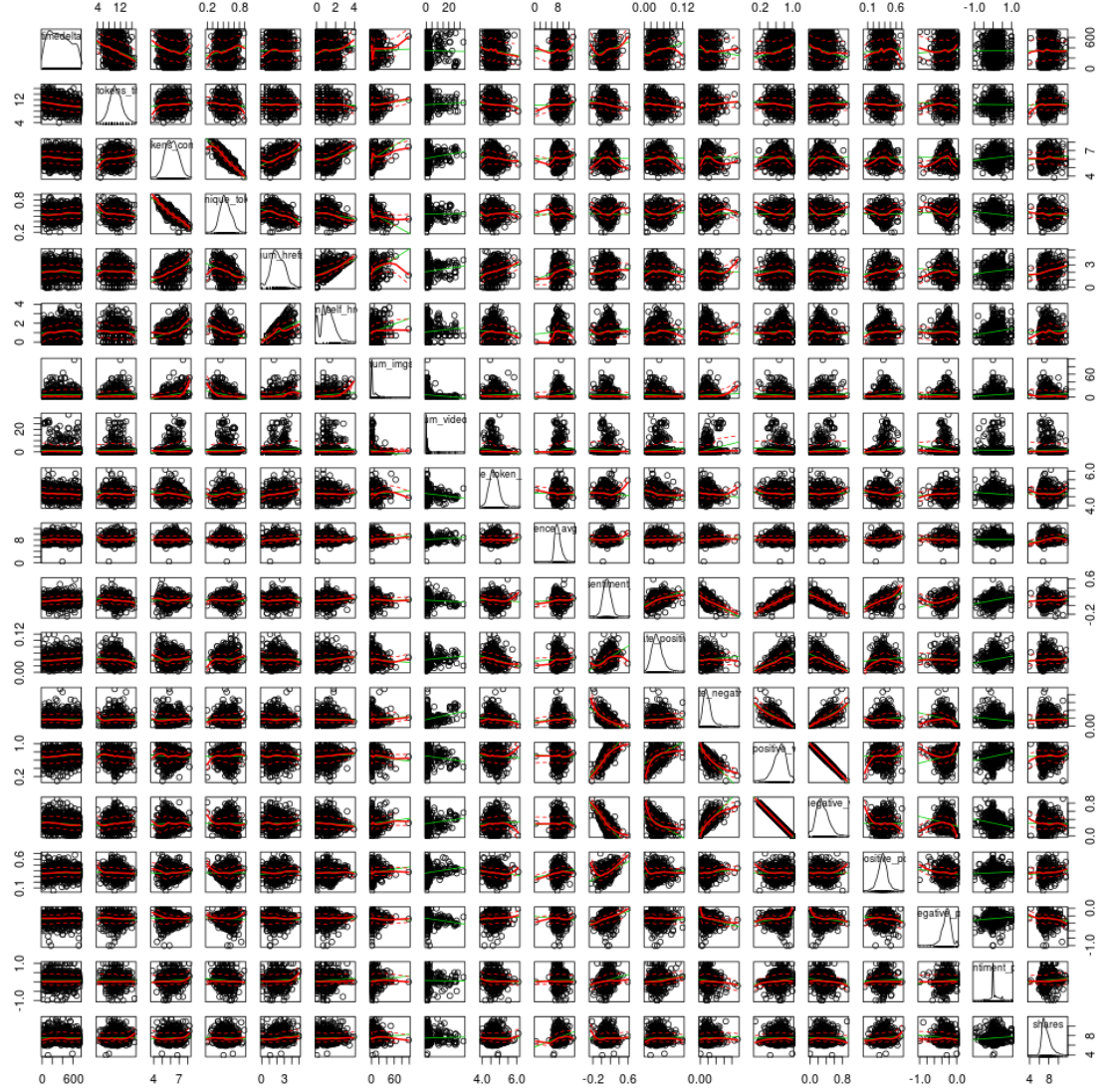


Figure 7: Scatterplot matrix of all considered continuous features

values from 0 to 1.

We will predict sentiment polarity over positive words rate. As we have rather heterogeneous data, let's make sure we won't be able to do our regression better with the help of grouping by **Channel** (Figure ??, a).

All channels look very similar, and we decided to consider only technical channel (just to reduce the sample size). After all these actions our scatterplot looks like at (Figure ??, b). Further at this section we will call the

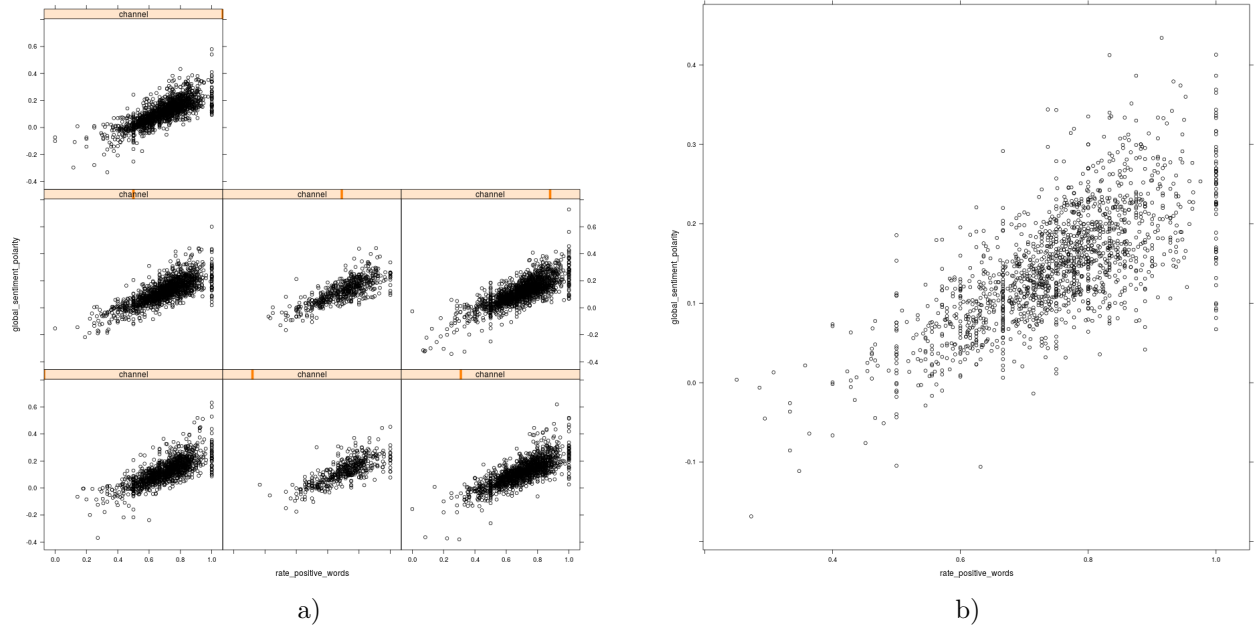


Figure 8: Grouped by channel (a) and only technical channel (b) dependence between `global_sentiment_polarity` and `rate_positive_words`.

predicted and the prediction features just `global_sentiment_polarity` and `rate_positive_words`, implying we work with only one channel.

2.2 Model of linear regression

Using basic functions in R, we have built a linear regression with slope equals 0.4341 and intercept equals -0.1788. The results of the regression you can see in Figure ??.

The slope is significantly positive (p-value equals 0) and it's not surprising: the more positive words are in the article, the more text is of positive polarity.

2.3 Correlation and determinacy coefficients

The correlation equals 0.7170736 and the coefficient of determination equals 0.5142 (adjusted is 0.5139). As we know from the definition of the coefficient of determination, R^2 measures of how well the regression line approximates the real data points and equals the ratio of explained variance. It's believed in practice that $R^2 > 0.5$ is acceptable, but not enough accurate.

Particularly the value of 0.5139 means that about 51% of variability between the two `global_sentiment_polarity` and `rate_positive_words` is captured by the linear model built with linear regression and the remaining 49% of variability still remains unaccounted for.

In another words the value of determination coefficient R^2 shows the rate of decrease of the variance of `global_sentiment_polarity` after its linear relation to `rate_positive_words` has been taken into account by the regression.

2.4 Bootstrap

We have conducted 5000 bootstrap trials to estimate 95% confidence intervals of slope, intercept and correlation coefficient. The results are summarized by the histograms shown in the figures ??-??.

It can be easily seen that histograms are pretty similar to the normal type of distribution. But let us prove that.

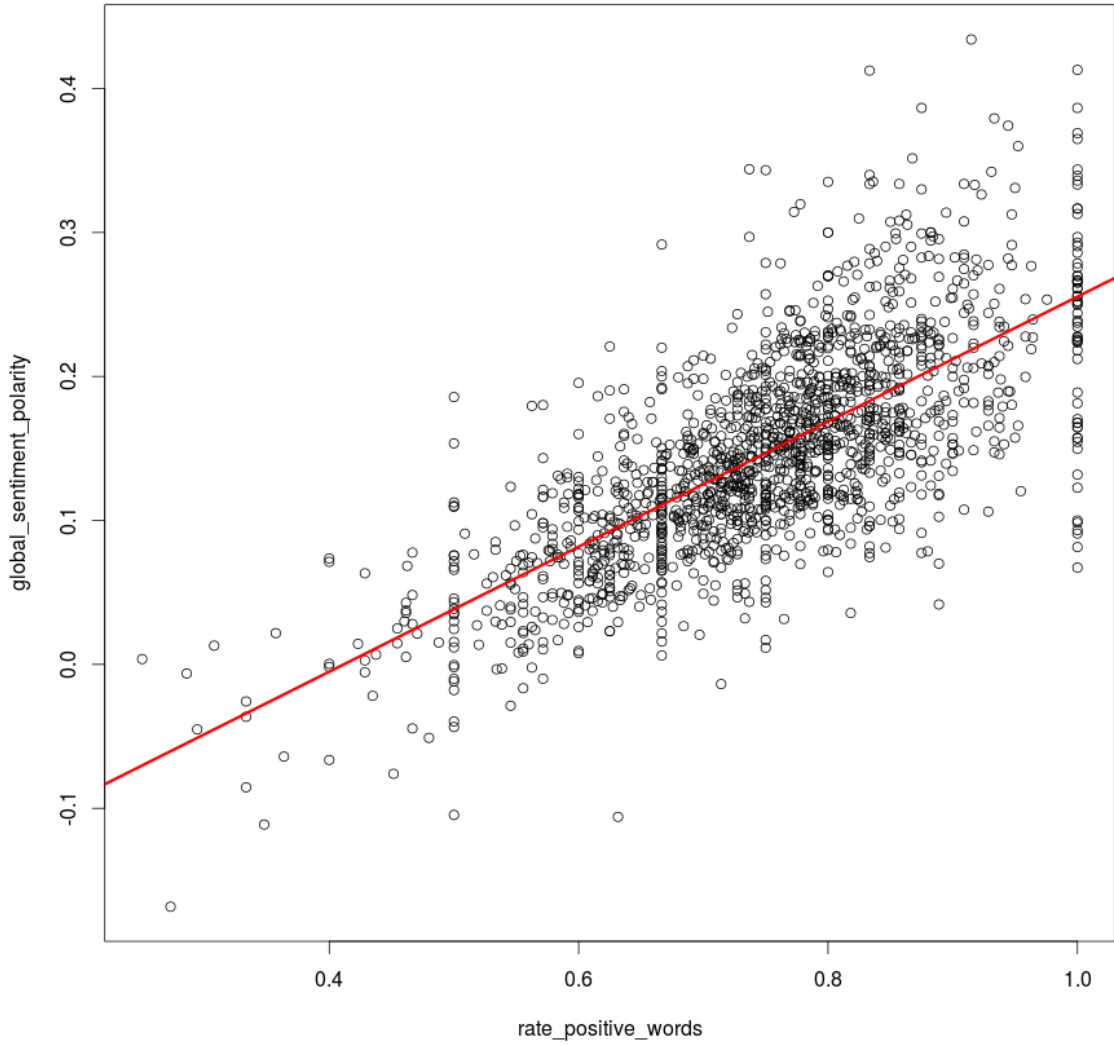


Figure 9

We have performed the Shapiro-Wilk normality test with intercept, slope and correlation coefficient bootstrap distributions. The results shown in the Table ?? are a little bit striking: with good p-value it is trustworthy that intercept and slope bootstrap samples are obtained from the normal distribution, but we could not say so about correlation coefficient samples data.

So for the correlation coefficient we compute CI using ranked quantiles. The aforesaid results we obtained are shown in Table ?. It is worth to mention that the difference in estimating correlation coefficient CI using the assumption of normality and without such is not very huge, but we think it is a good idea not only test normality of the data by it's visualization but also with statistical tests.

Table 5: Shapiro-Wilk normality test p-value

	95% p-value
Intercept	0.2279
Slope	0.244
Correlation	5.621e-08

Table 6: 95% confidence intervals (CI's) for intercept, slope, correlation coefficient based on bootstrap technique

	95% CI
Intercept	(-0.209; -0.148)
Slope	(0.392; 0.476)
Correlation (normal)	(0.668; 0.765)
Correlation (percentile)	(0.666, 0.763)

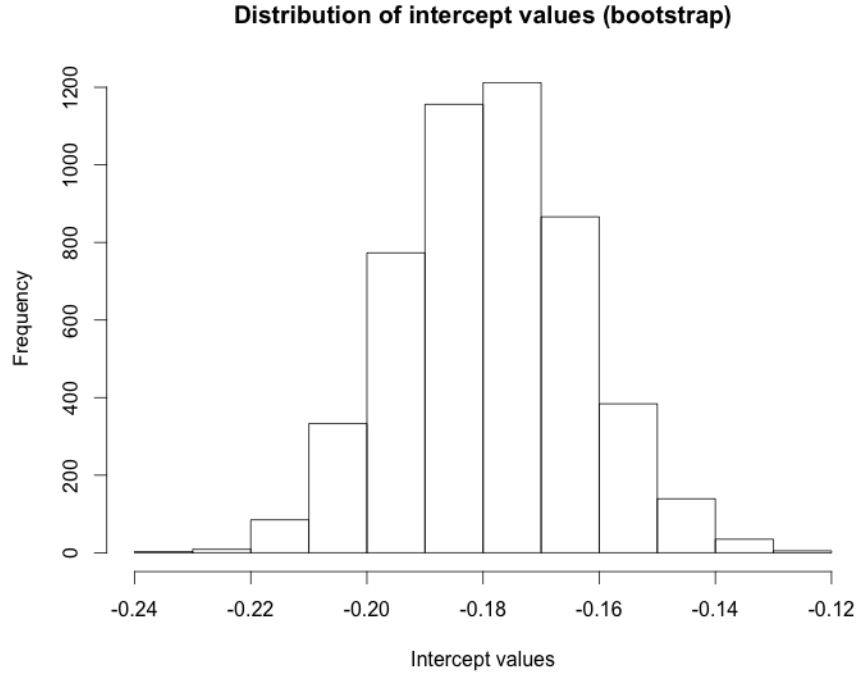


Figure 10: Distribution of intercept value of linear regression models built with each pair of sampled (global_sentiment_polarity) and (rate_positive_words) using bootsrap

2.5 Average relative error

Recall average relative error (ARE) and coefficient of determination (R^2) definitions:

$$\text{ARE} = \frac{1}{N} \sum_{i=1}^N \left| \frac{y_i - \hat{y}_i}{y_i} \right|,$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

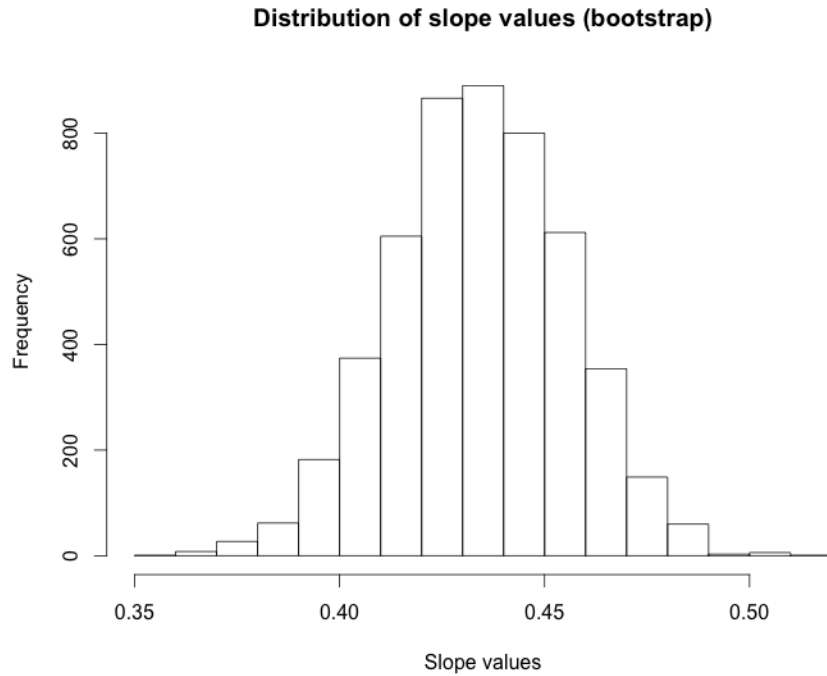


Figure 11: Distribution of slope value of linear regression models built with each pair of sampled (`global_sentiment_polarity`) and (`rate_positive_words`) using bootrsap

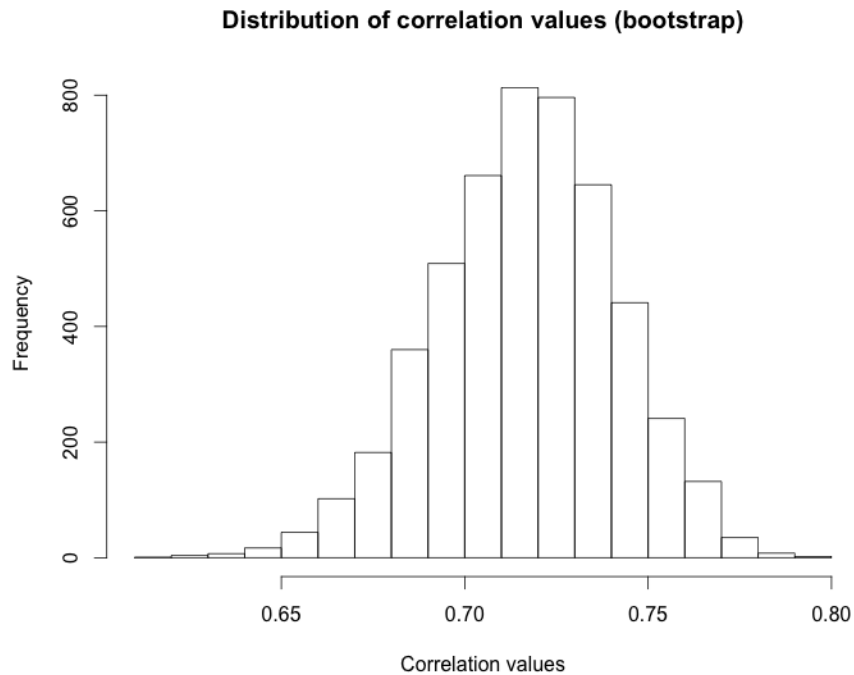


Figure 12: Distribution of correlation values between dependent variable (`global_sentiment_polarity`) and regressor (`rate_positive_words`)

```
mean( abs((y.feature - model$fitted.values)/y.feature) * 100) # in %

## [1] 56.3725

summary_regression$r.squared * 100 # in %

## [1] 51.41945
```

As we can see, considered values are reasonably close. But we should note that ARE is sensitive to the addition of a constant to all of the y_i while R^2 is not. That is, we could obtain any arbitrary value of ARE keeping the coefficient of determination constant by adding some constants to y_i .

It was suprise to us that ARE can be greater than 1 (which is the case if y_i is much less than $y_i - \hat{y}_i$).

According to all these facts one could conclude that comparing ARE and R^2 is meaningless without additional assumptions.

2.6 Nature-inspired algorithm

We will use nature-inspired algorithm to compute the parameters of linear regression that minimize the absolute relative error. If the target values of regression are y_i and predicted by linear regression values are \hat{y}_i , the absolute relative error is:

$$\frac{1}{N} \sum_{i=1}^N \left| \frac{y_i - \hat{y}_i}{y_i} \right|. \quad (1)$$

We have implemented algorithm similar to the one, which is described in [?], but we will minimize the function `delta(coefficients, x, y)`, which looks like

```
delta <- function(coefficients, x,y){
  a = coefficients[1]; b = coefficients[2]
  yp <- a*x + b
  esq <- mean( abs((y - yp)/y) )
}
```

We also need the function that compute permissible limits for coefficients. Let look at two values of target $y_i = a \cdot x_i + b$ and $y_j = a \cdot x_j + b$ ($i \neq j$) and express a, b in terms of x and y :

$$a_{ij} = \frac{y_j - y_i}{x_j - x_i}, \quad b_{ij} = \frac{y_i x_j - y_j x_i}{x_j - x_i}.$$

And then we calculate max and min of coefficients a and b among all pairs (x_i, y_i) and (x_j, y_j) .

Using the nature-inspired approach (we have implemented it in R-function `nlr`) we obtained the following values of slope (a) and intercept (b) and value of relative error:

```
# The regression model: y.feature = a * x.feature + b
# Coefficients of slope and intercept respectively:
model.nlr <- nlr(x.feature, y.feature)
model.nlr

## [1] 0.3610875 -0.1570541

# Value of relative error:
eps.nlr <- y.feature - model.nlr[1]*x.feature - model.nlr[2]
mean( abs(eps.nlr / y.feature) ) * 100

## [1] 50.04851
```

Comparing the values of two relative errors, we can see that nature-inspired approach reduce its value by a few percent.

In R-language there is a package **genalg** with function **rbga** that implement this approach. The results obtained with function **rbga** are very close to the results described above:

```
# The regression model: y.feature = a * x.feature + b
# Calculate the permissible limits for a and b
bound <- ddr(x.feature, y.feature)
bounds.min <- c(bound[[1]][1],bound[[2]][1]) # (a.min, b.min)
bounds.max <- c(bound[[1]][2],bound[[2]][2]) # (a.max, b.max)

rbga.res <- rbga(bounds.min, bounds.max, popSize = 30, iters = 5000,
  evalFunc = function(coefs) delta(coefs, x.feature, y.feature))

# Results (we need to take a look at "Best Solution")
cat(summary(rbga.res))

## GA Settings
##   Type                = floats chromosome
##   Population size      = 30
##   Number of Generations = 5000
##   Elitism              = 6
##   Mutation Chance      = 0.3333333333333333
##
## Search Domain
##   Var 1 = [-6.674999999946,6.47188552188199]
##   Var 2 = [-5.27885519411197,5.97499999995278]
##
## GA Results
##   Best Solution : 0.348751377798375 -0.14645601290741

# Value of relative error:
delta(c(0.3477, -0.1463), x.feature, y.feature) * 100

## [1] 49.92909
```

3 Assignment 3.

3.1 Selection and building nominal features