

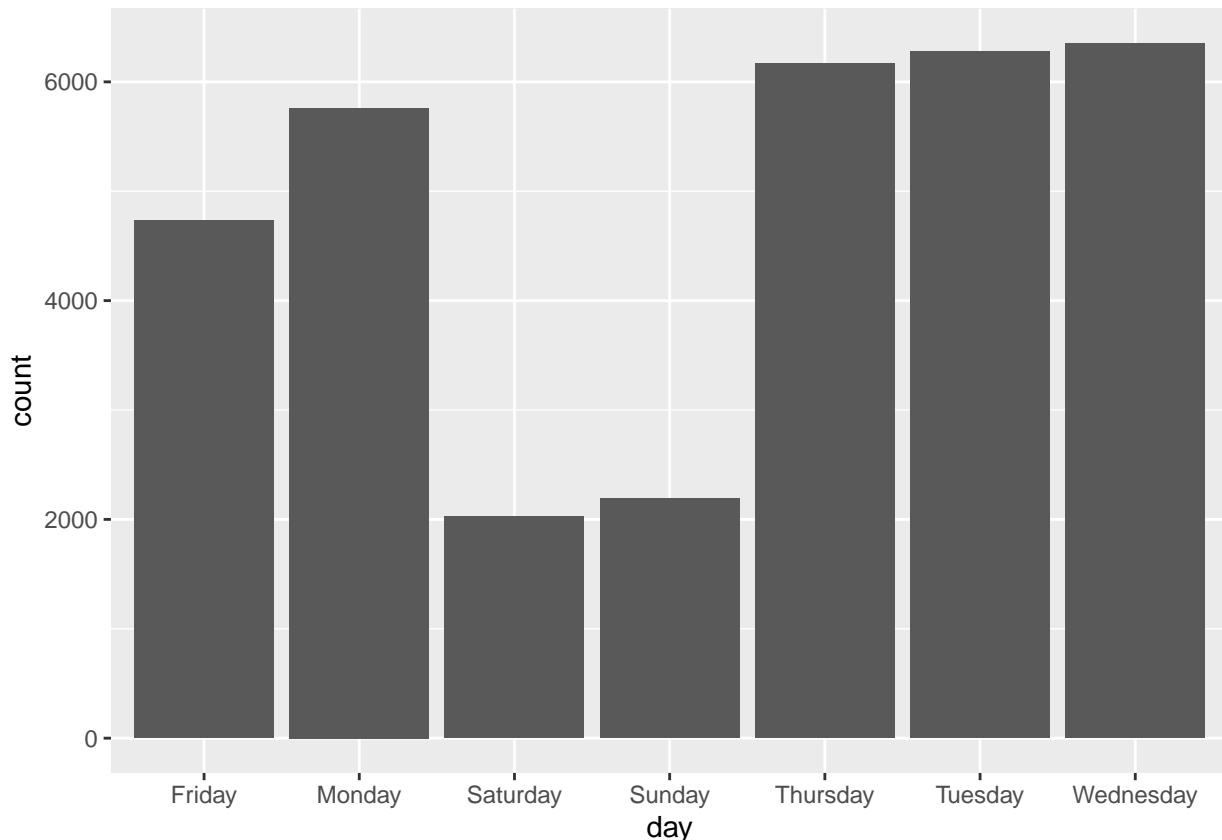
online-news-popularity

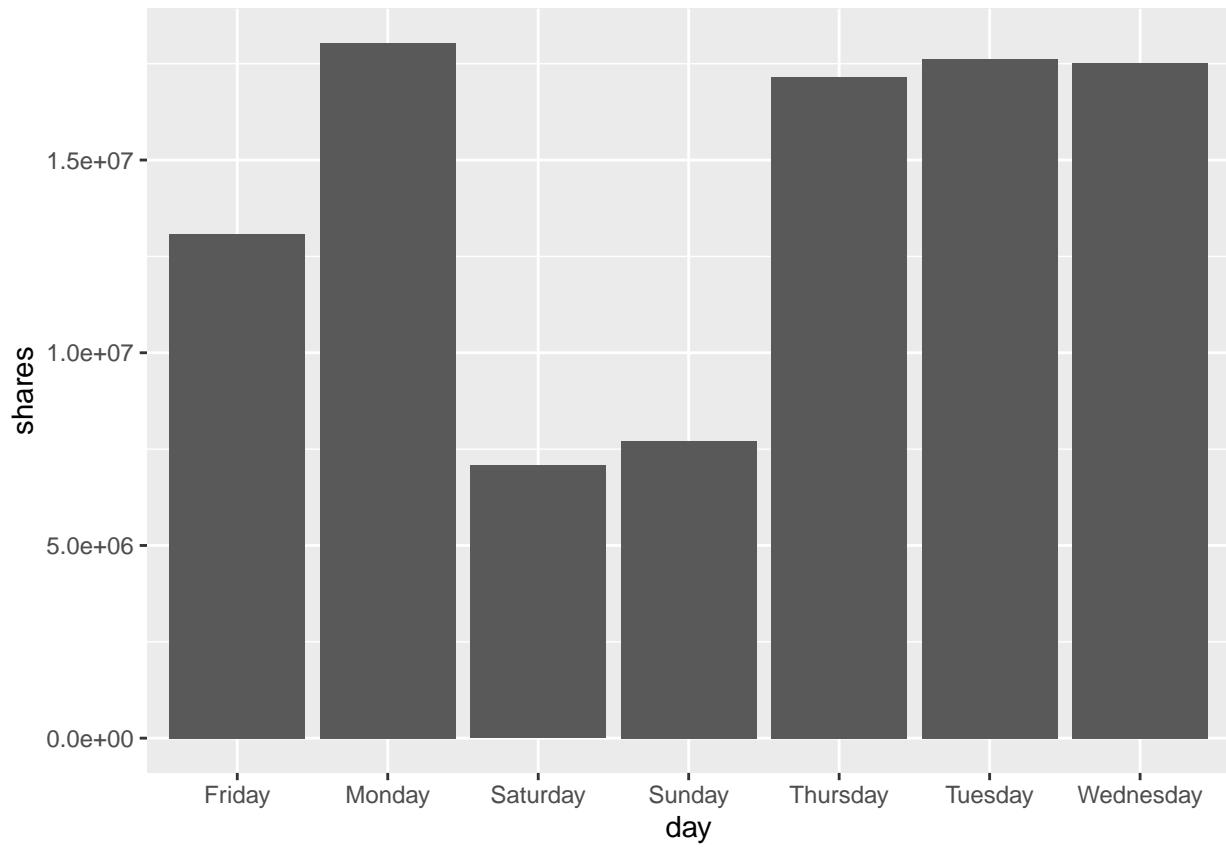
Sneh Gurdasani, Akshit Jain, Farhan Ansari, Sagar Singh

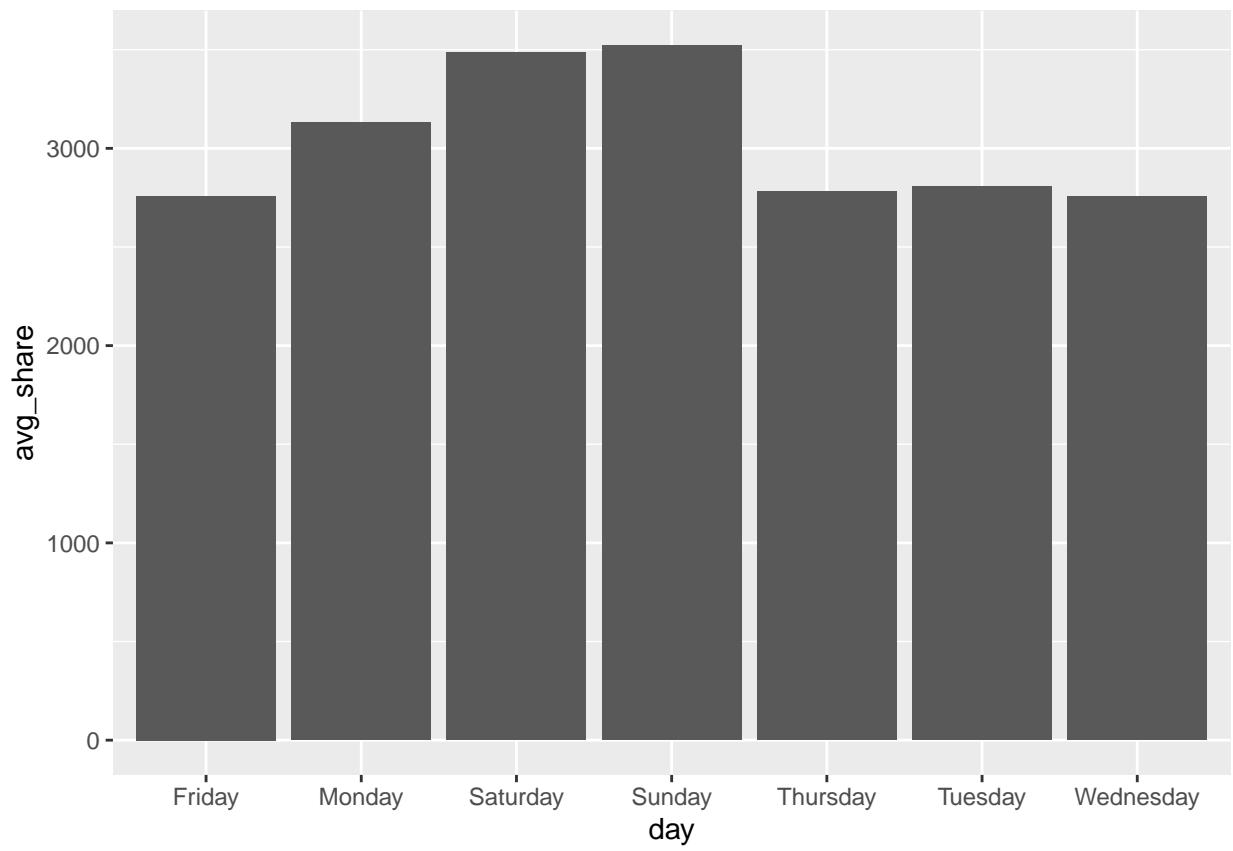
11/11/2019

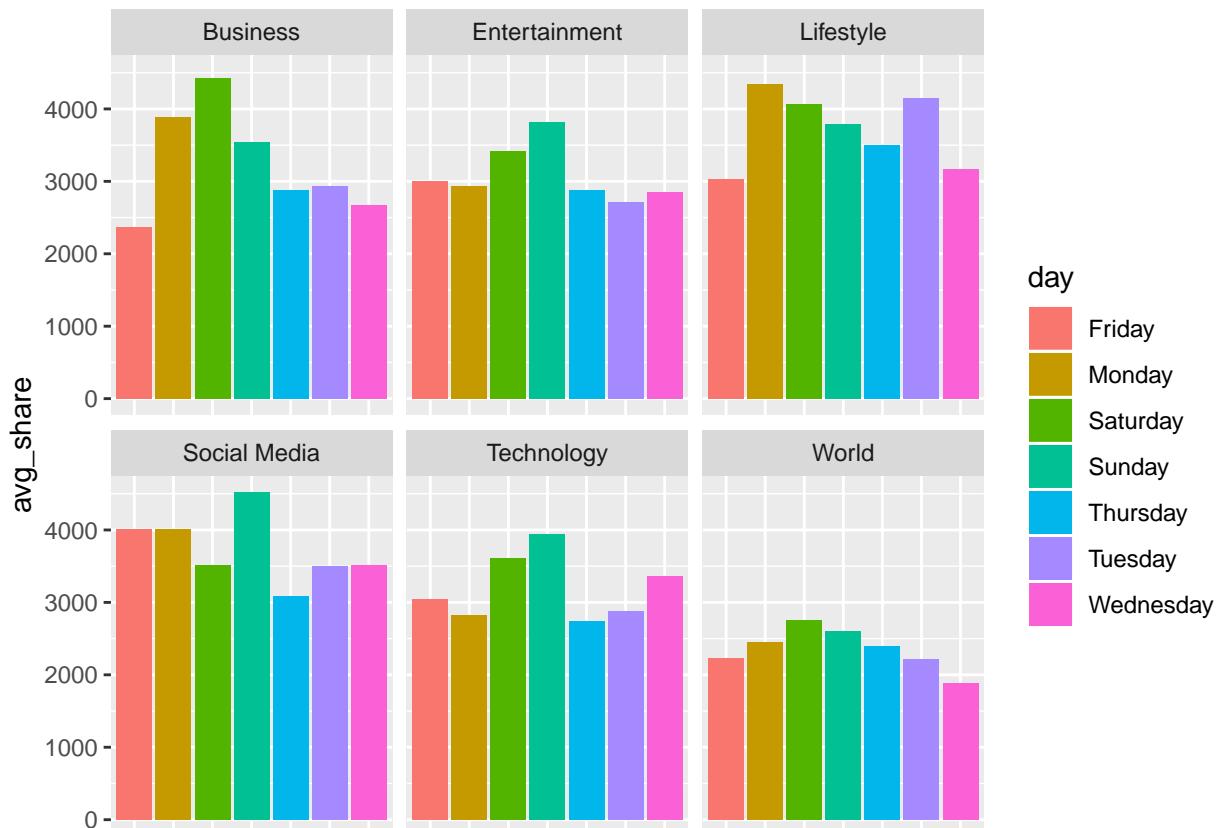
Load and Tidy Dataset

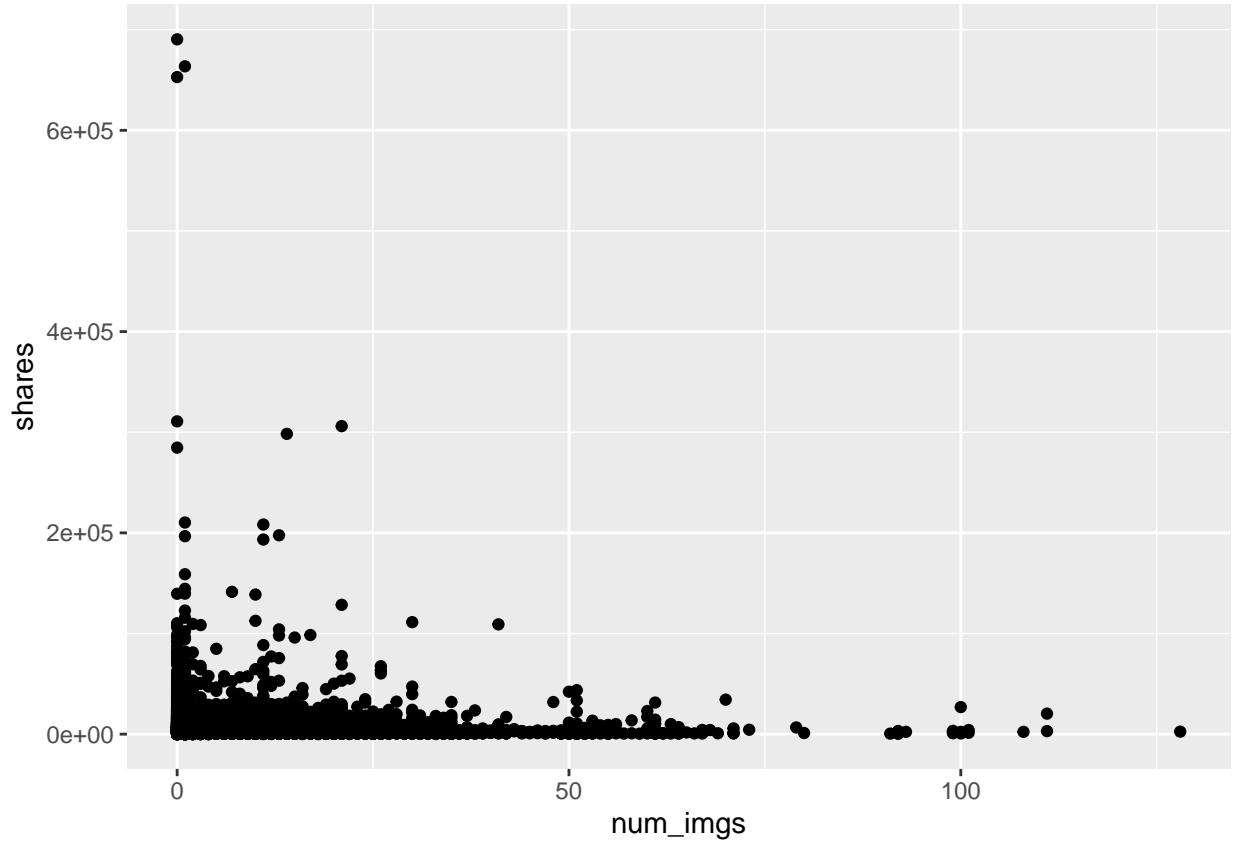
```
df <- read_csv('OnlineNewsPopularity/OnlineNewsPopularity.csv', col_types = cols())
list_of_dfs <- get_tidy_data(df)
cluster_numeric_df <- list_of_dfs[[1]]
classification_df <- list_of_dfs[[2]]  
  
visualize_summary_plots(classification_df)
```



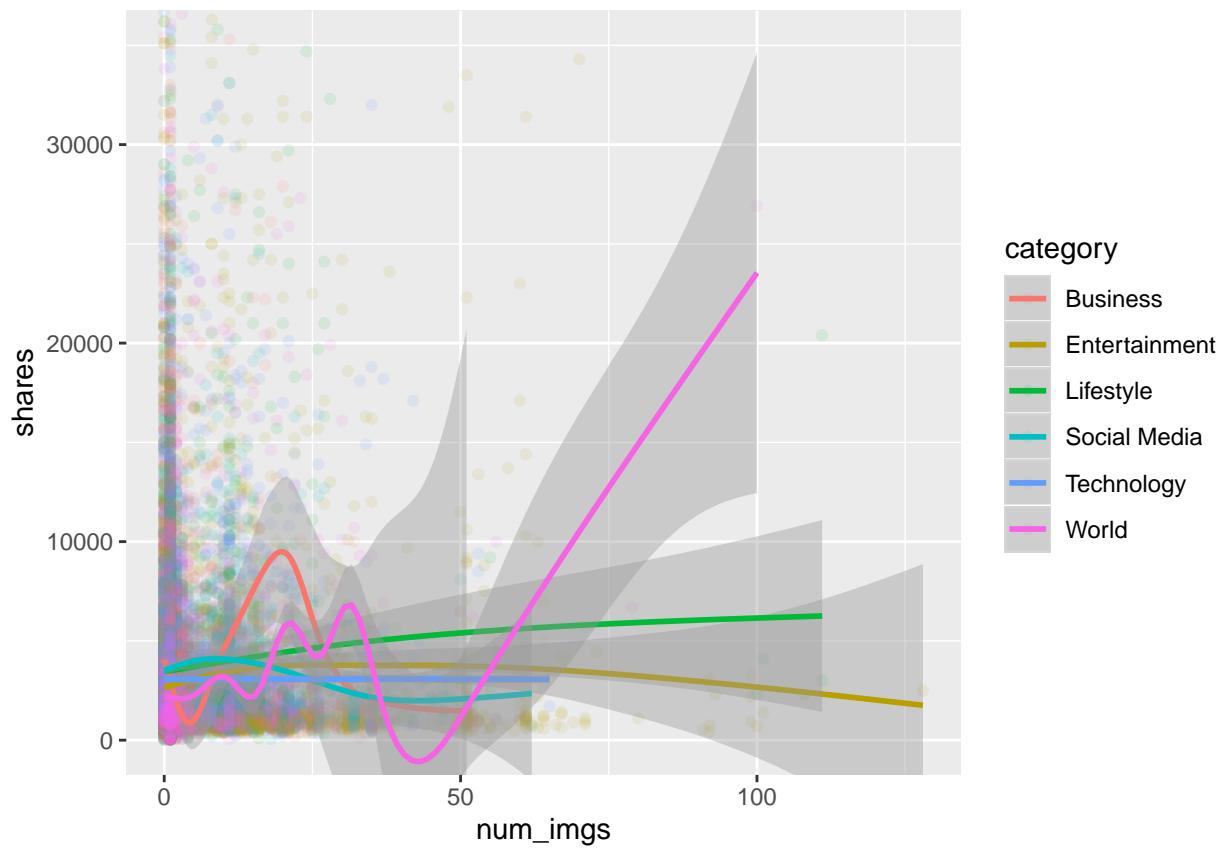


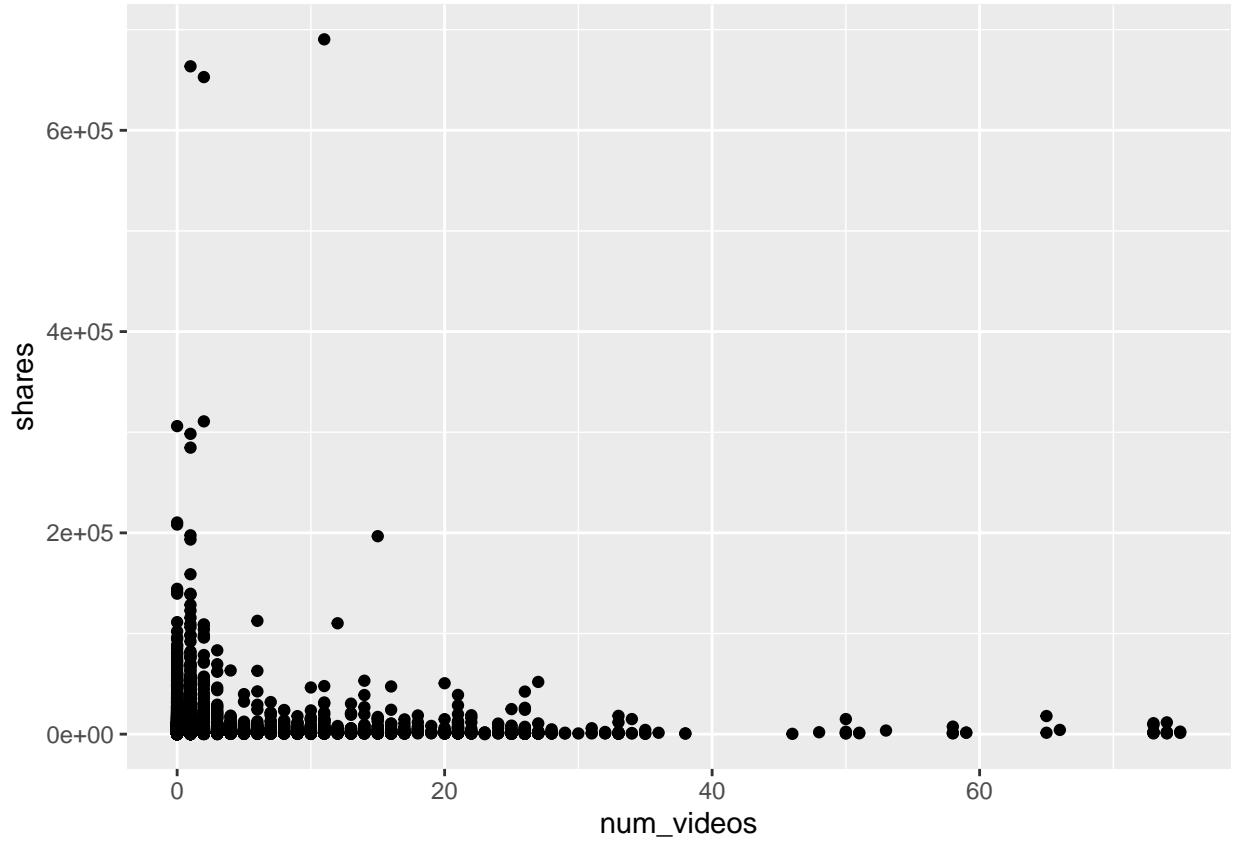




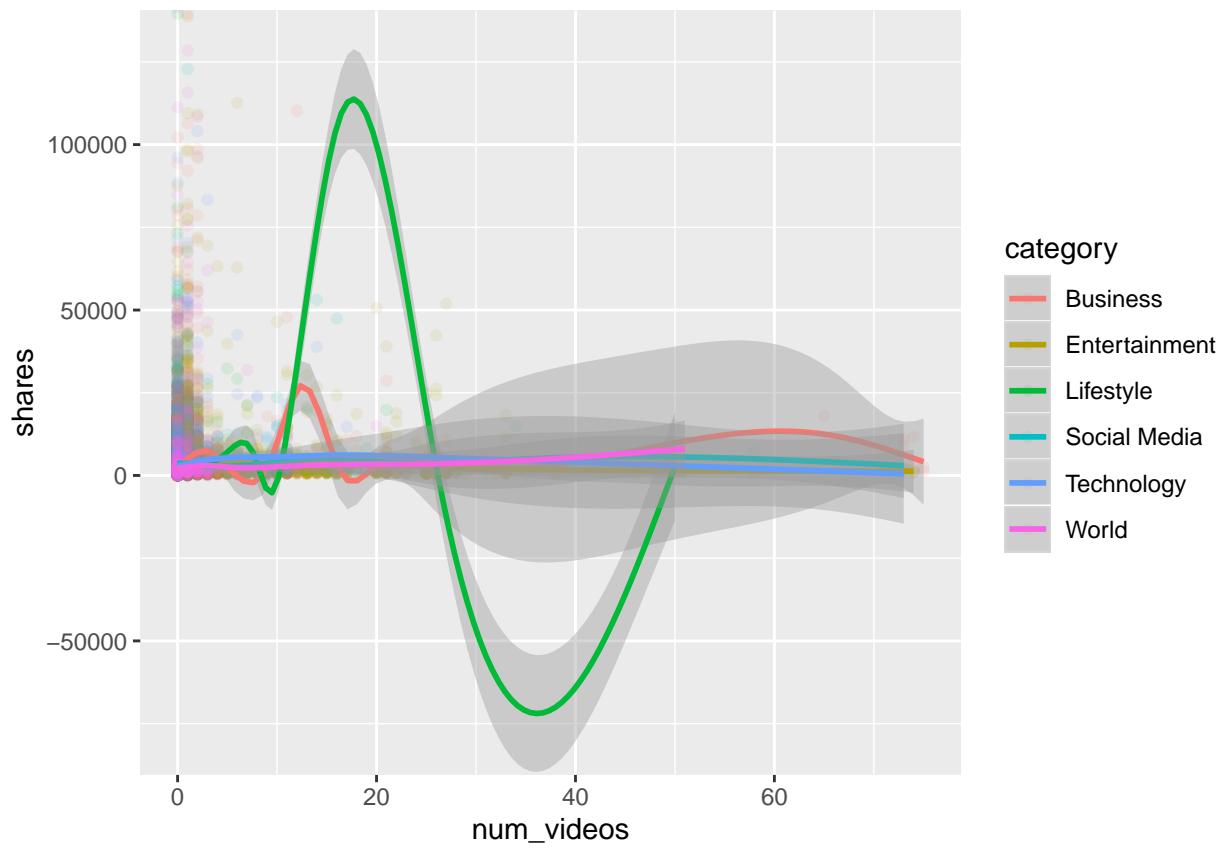


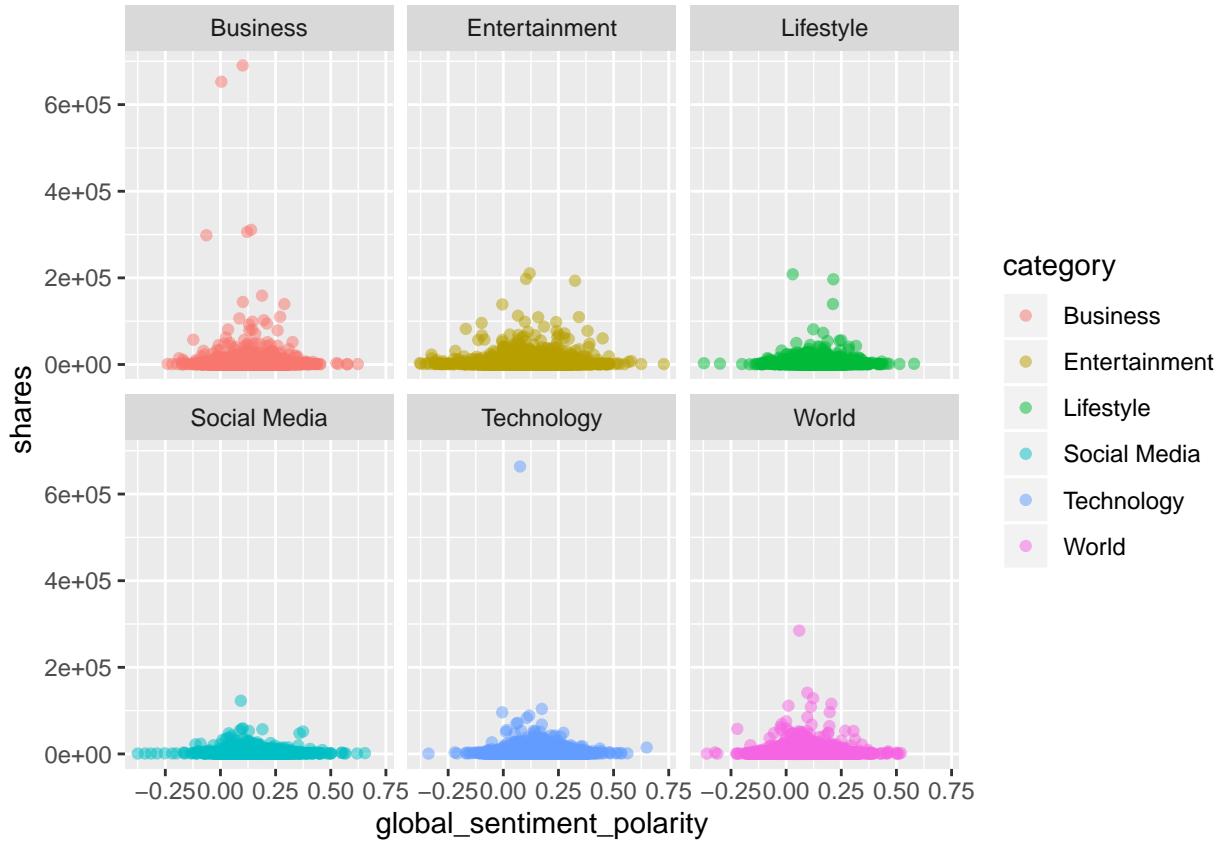
```
## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```

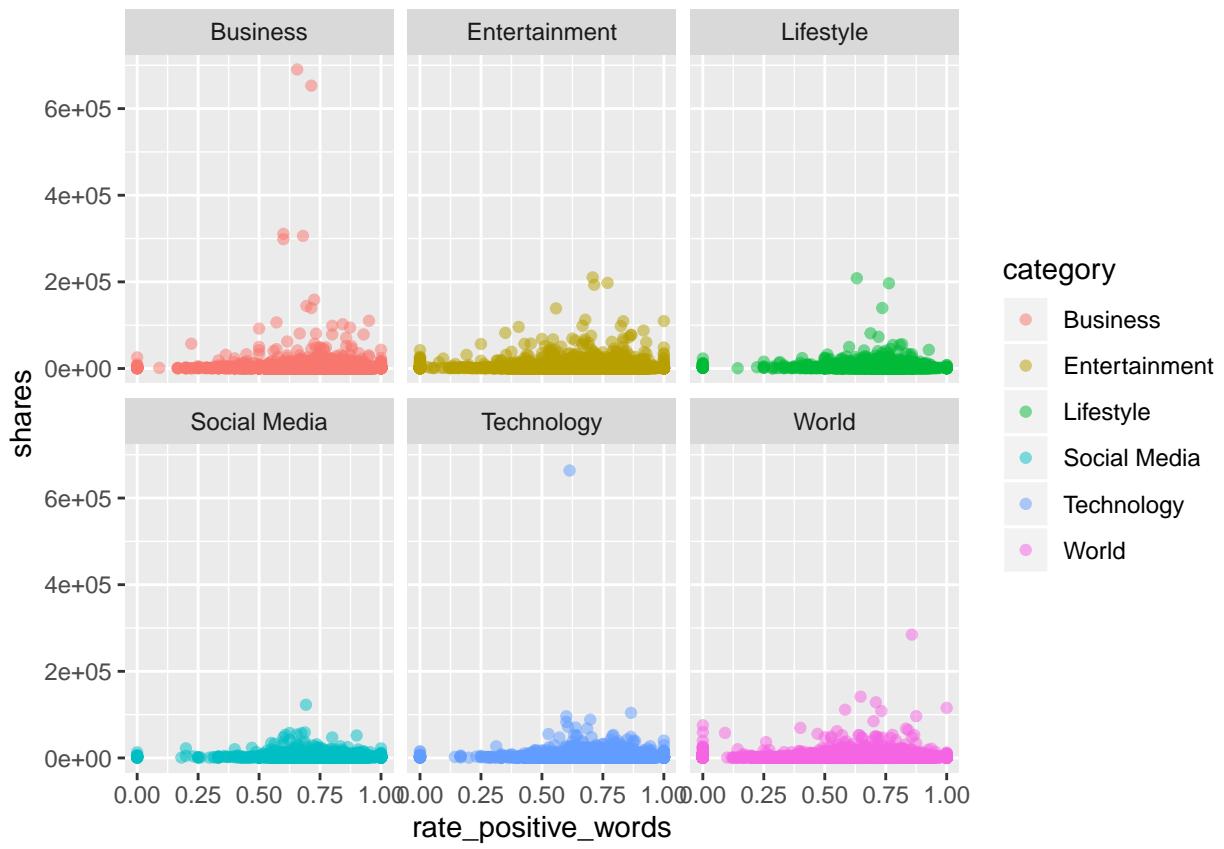


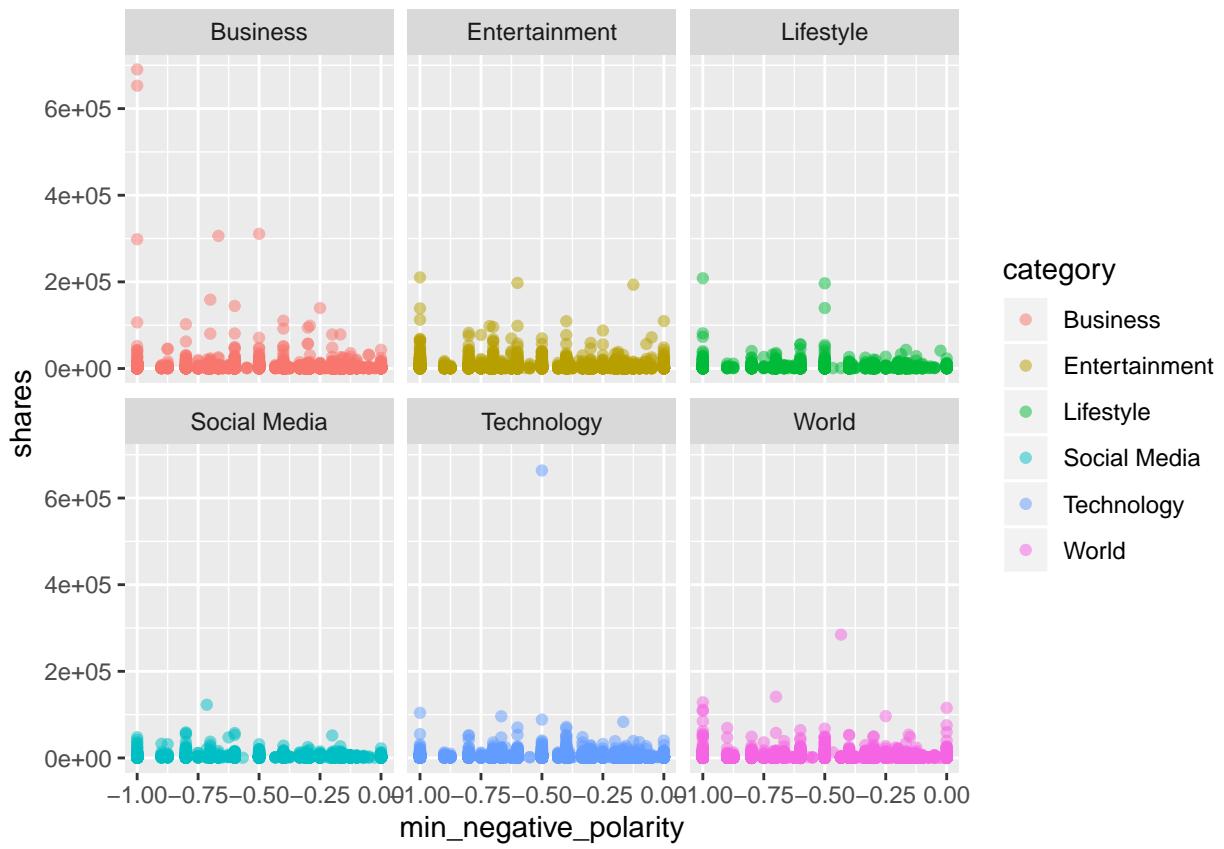


```
## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```

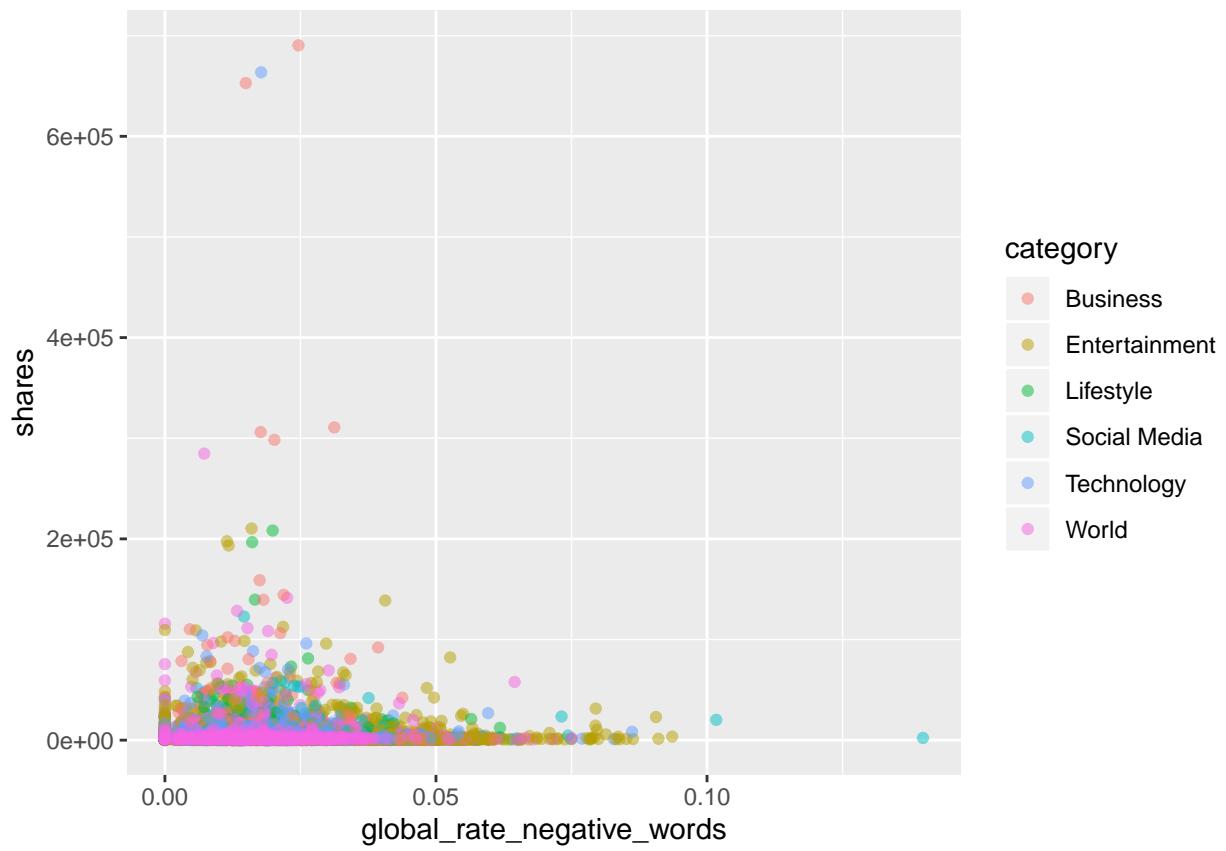


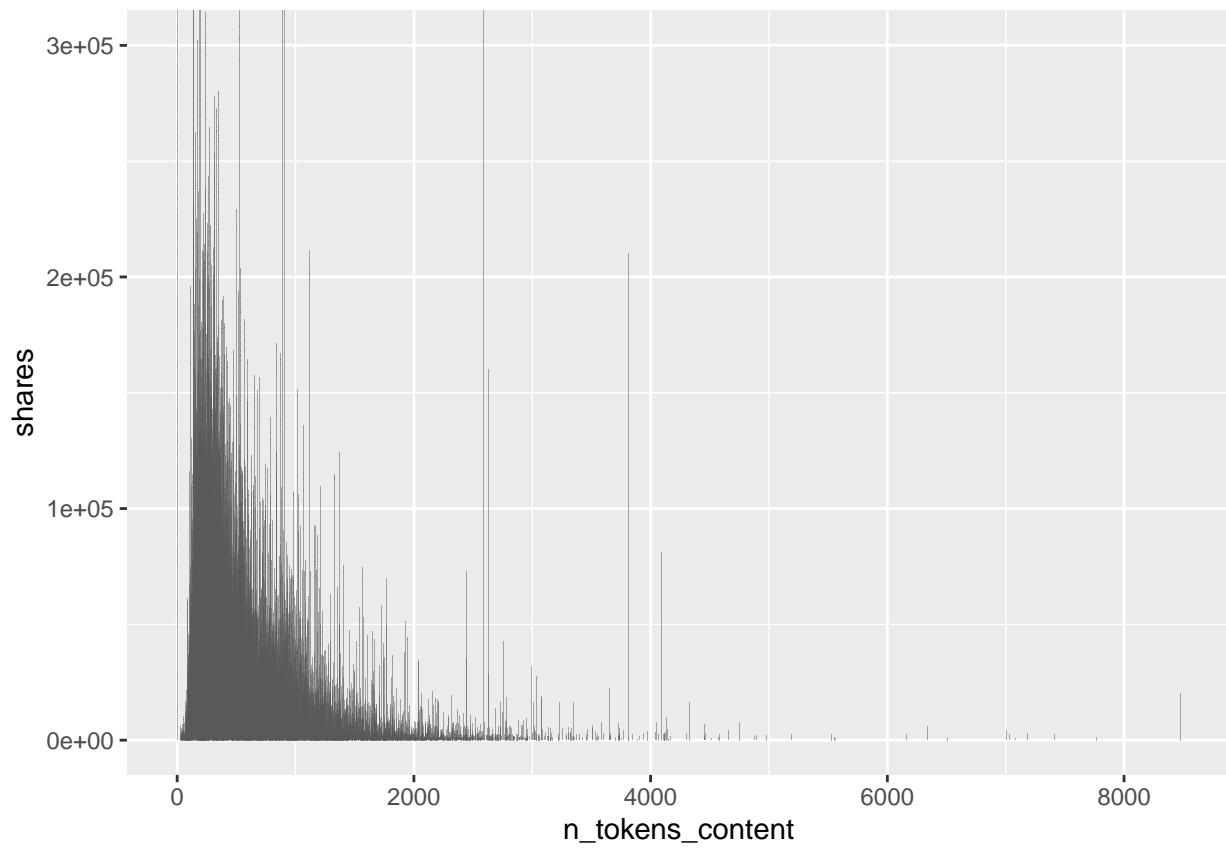


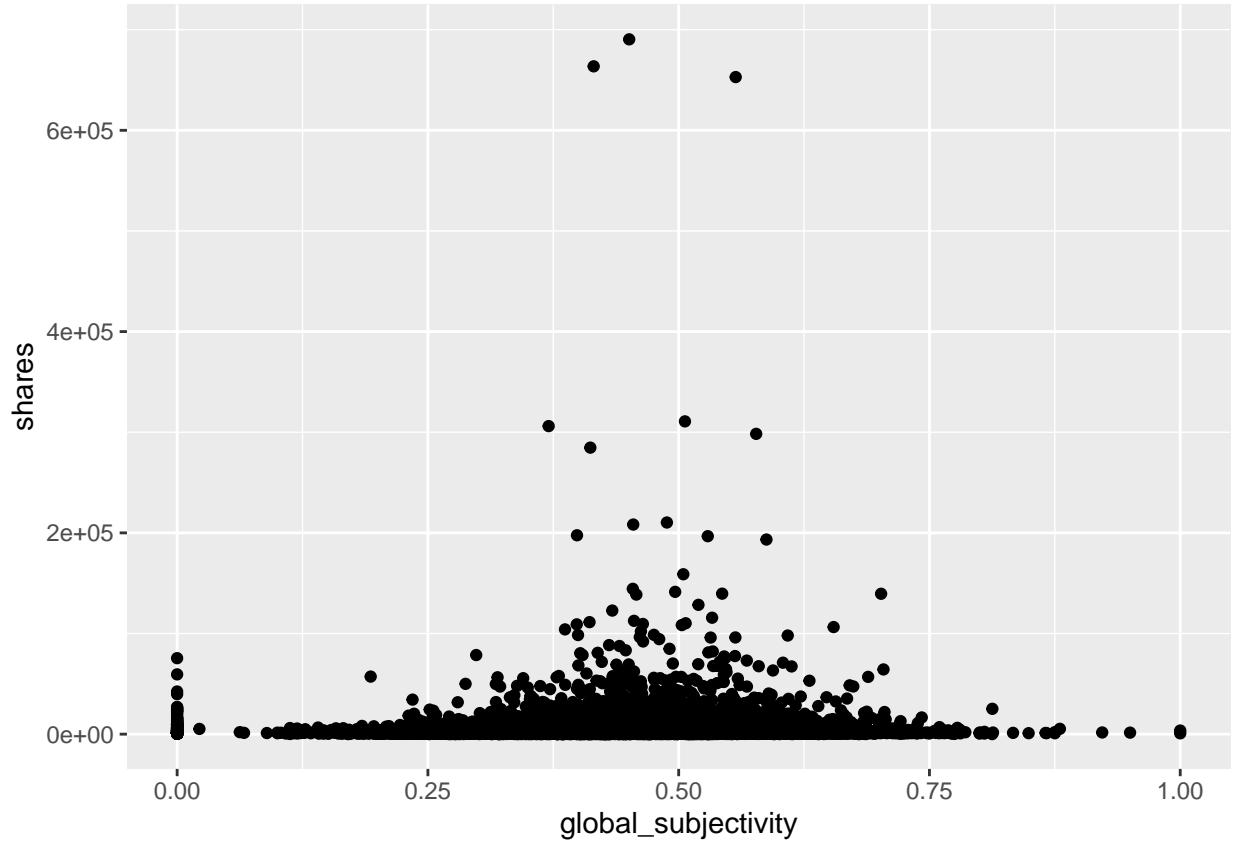


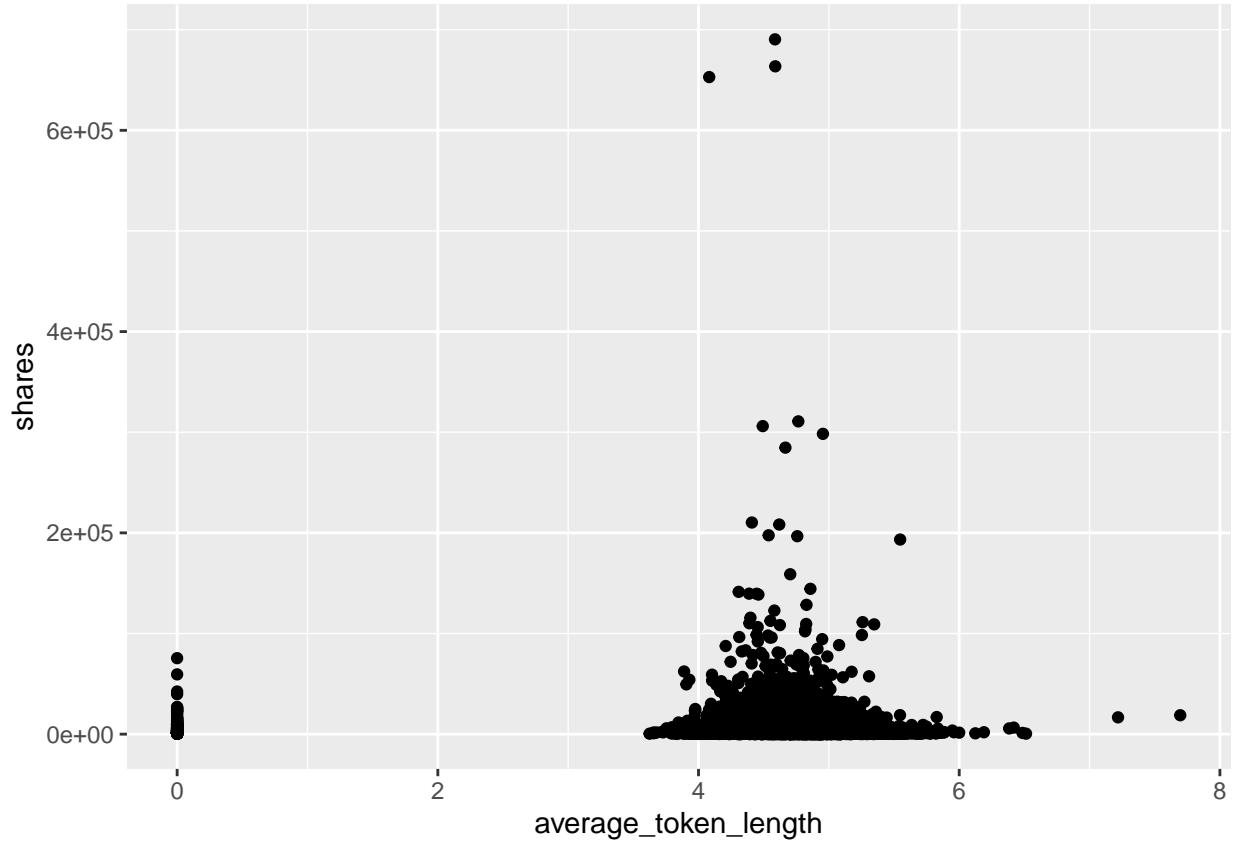


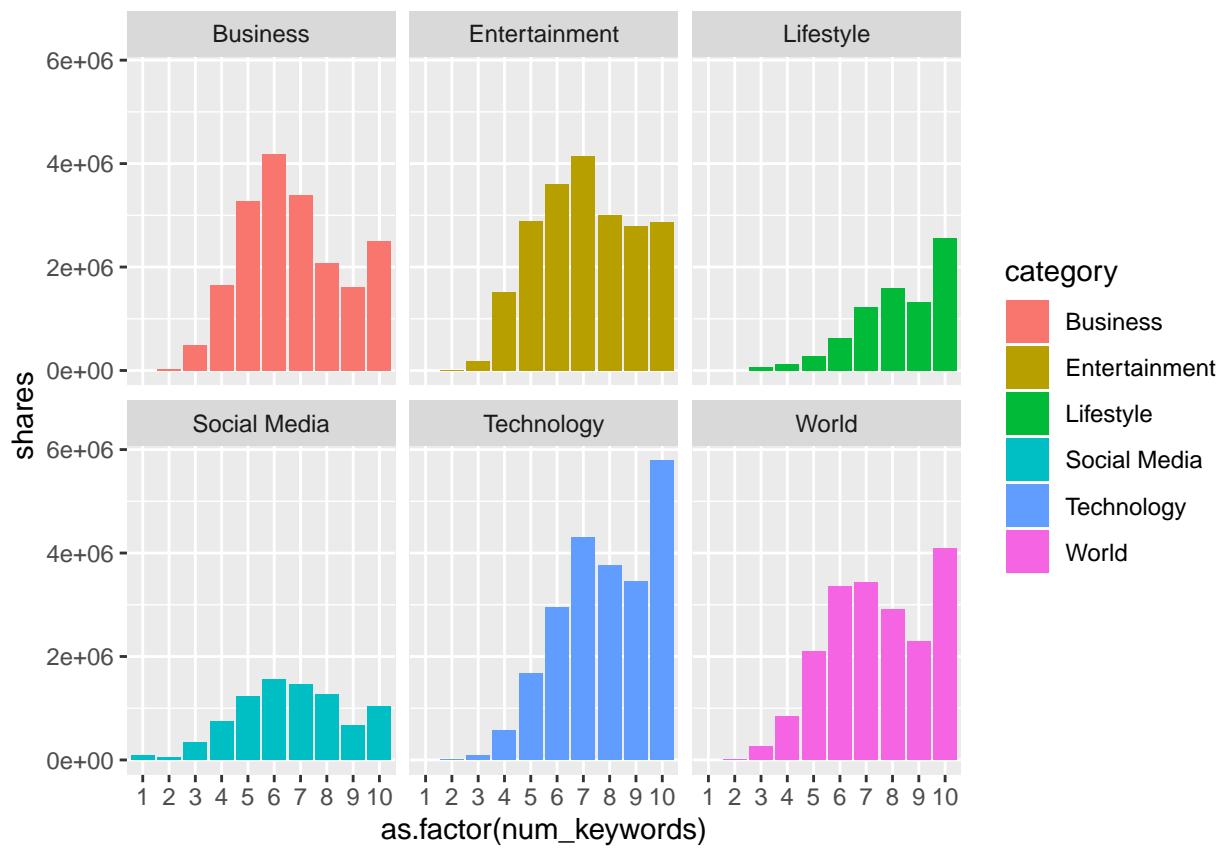
```
## Warning: Ignoring unknown parameters: binwidth, bins, pad
```

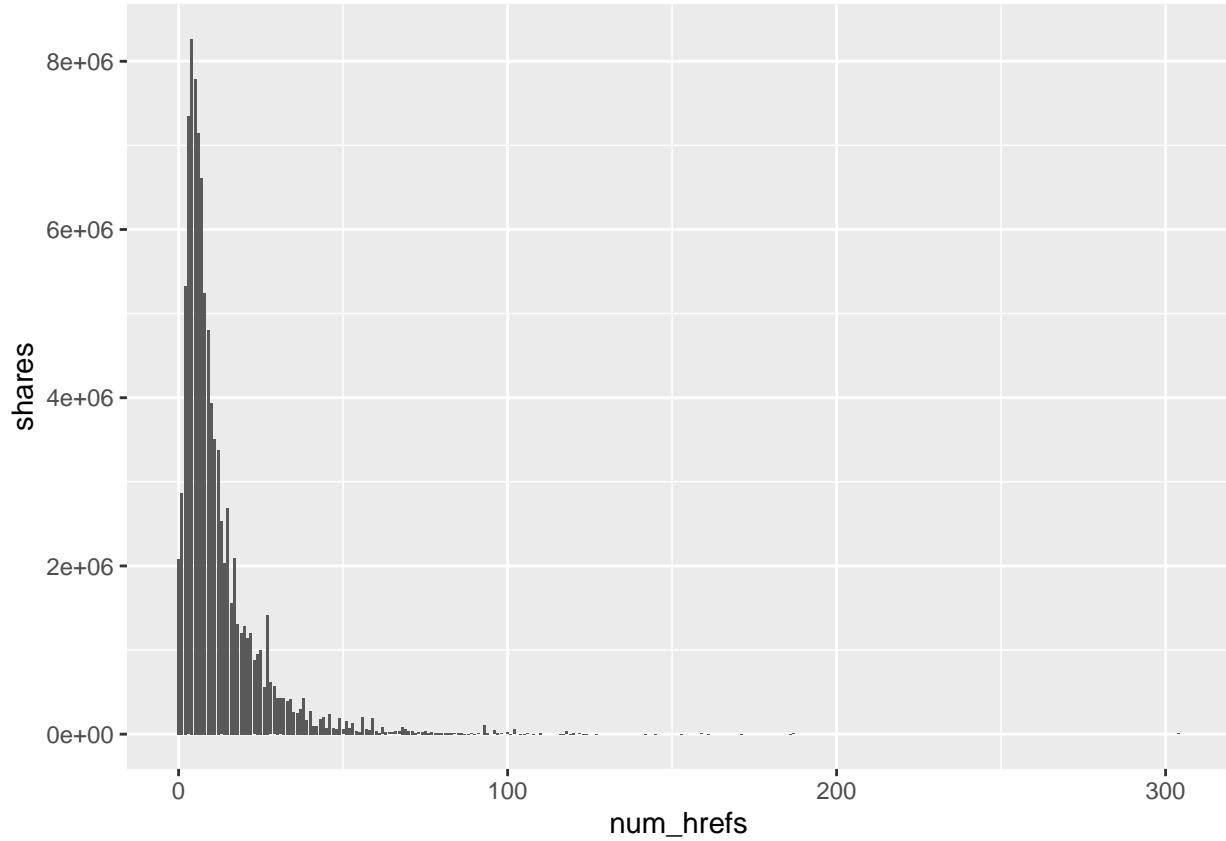


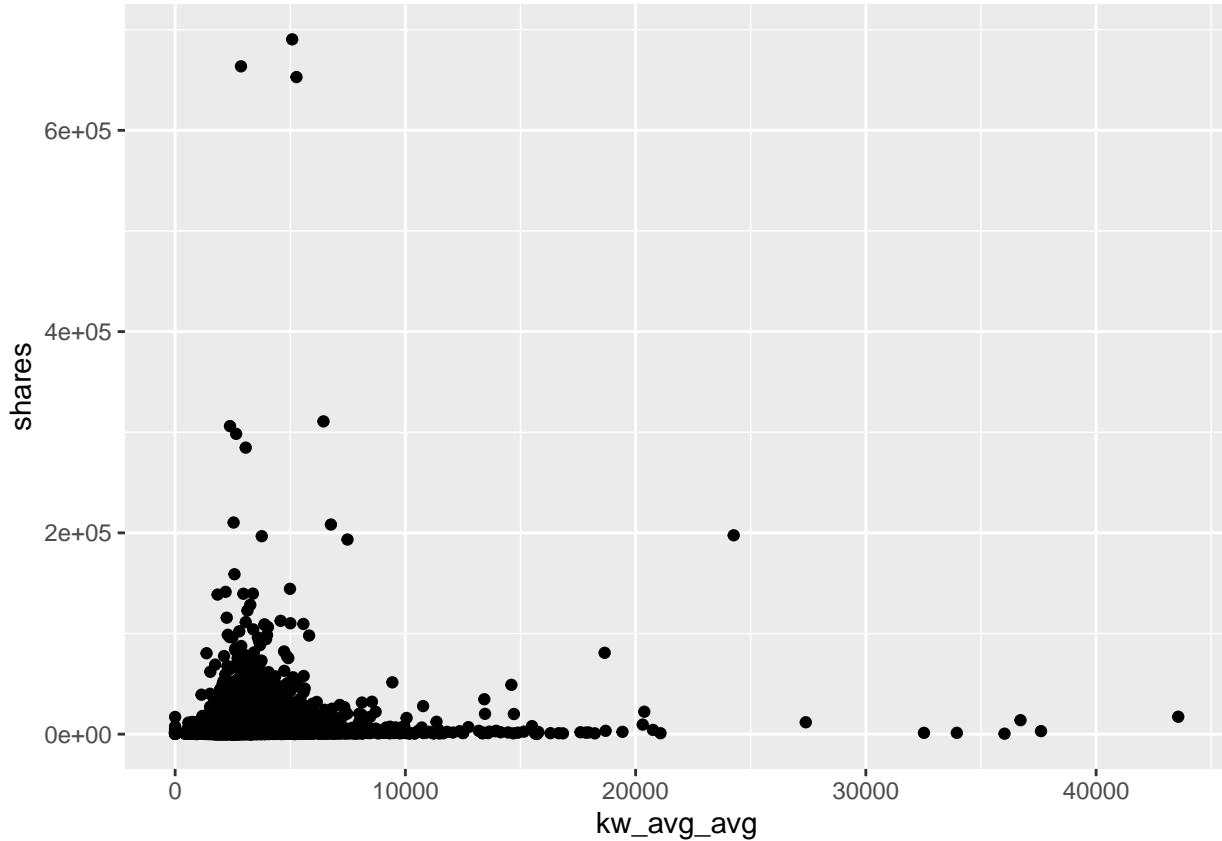












```
correlation_matrix(cluster_numeric_df)
```

```
## Selecting by shares
```

```
##          correlation_with      shares
## 1                  shares 1.0000000
## 2          kw_avg_avg 0.10002583
## 3          kw_max_avg 0.06586856
## 4 self_reference_avg_shares 0.06145828
## 5 self_reference_min_shares 0.05825928
## 6 self_reference_max_shares 0.04895216
## 7          kw_avg_min 0.04316214
## 8          kw_max_min 0.04307623
## 9          num_hrefs 0.04190182
## 10 global_subjectivity 0.03510054
## 11          LDA_03 0.03288953
## 12          kw_min_avg 0.02912447
## 13 n_tokens_content 0.02912325
## 14          num_imgs 0.02877061
## 15          num_keywords 0.02554536
```

```
prepare_data_classification(classification_df)
```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	1	930	1400	2929	2500	690400