# The most important variables for prediction of "Premium upgrade" are "Proportion working time" and "Searches per day"
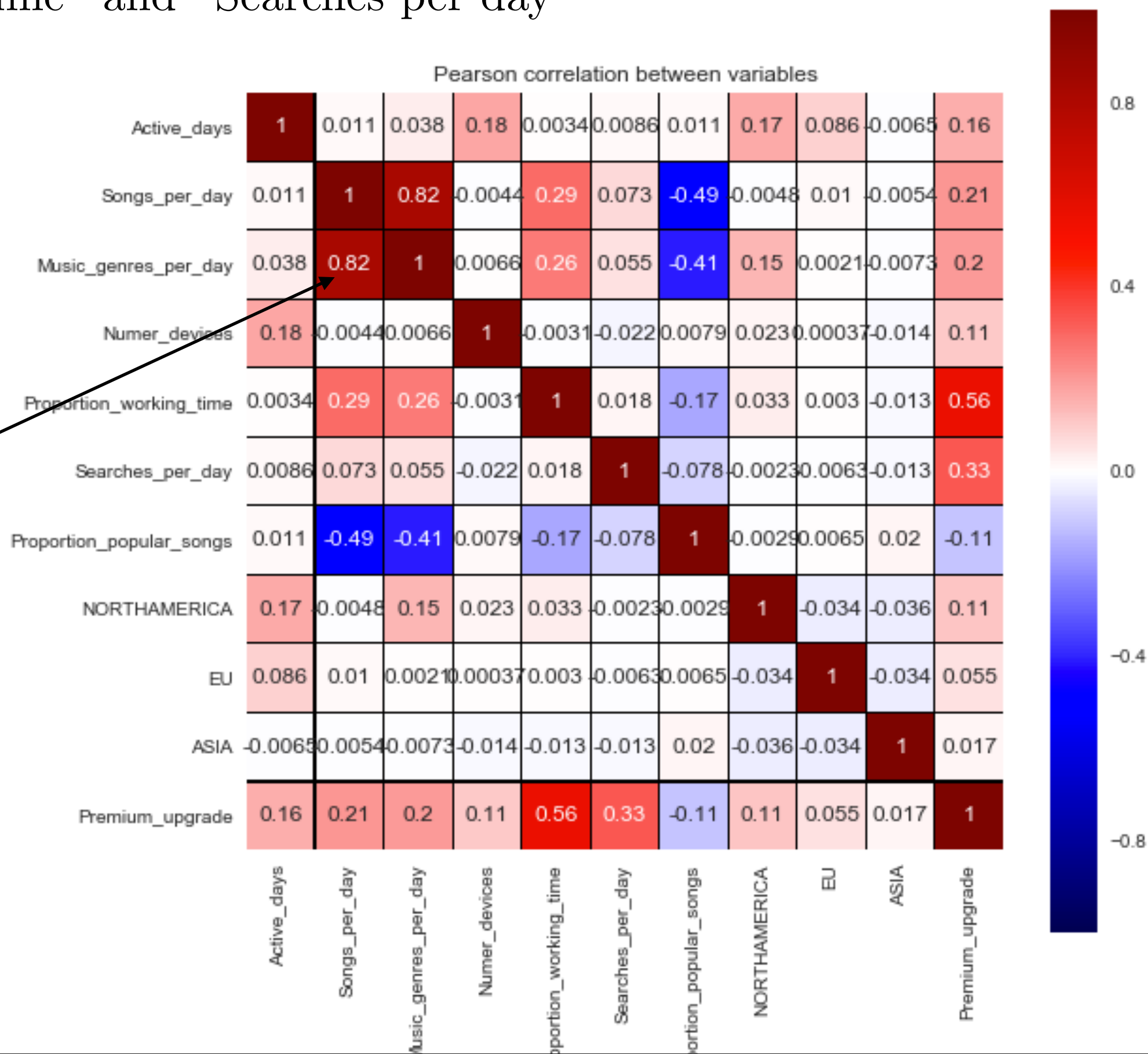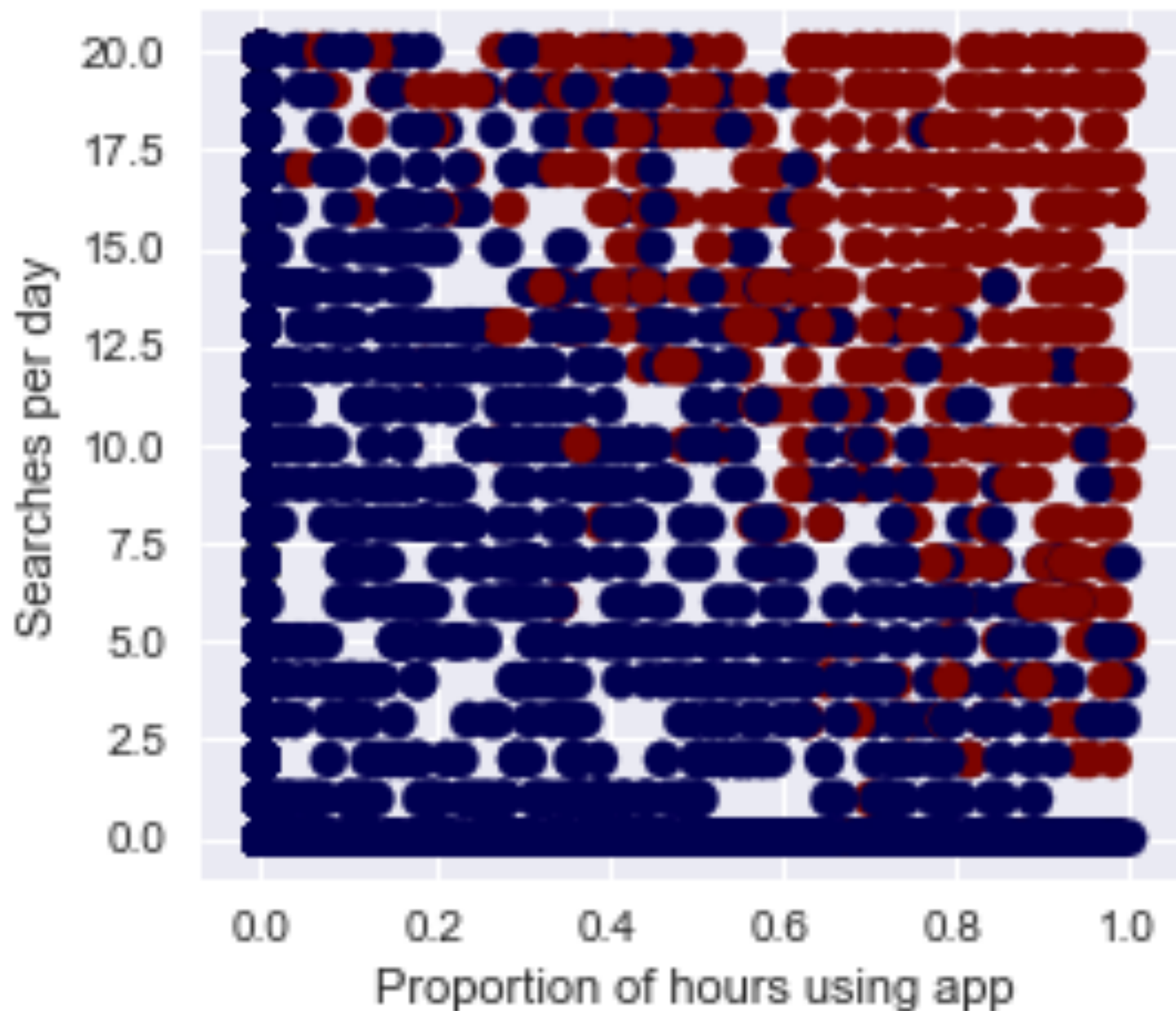


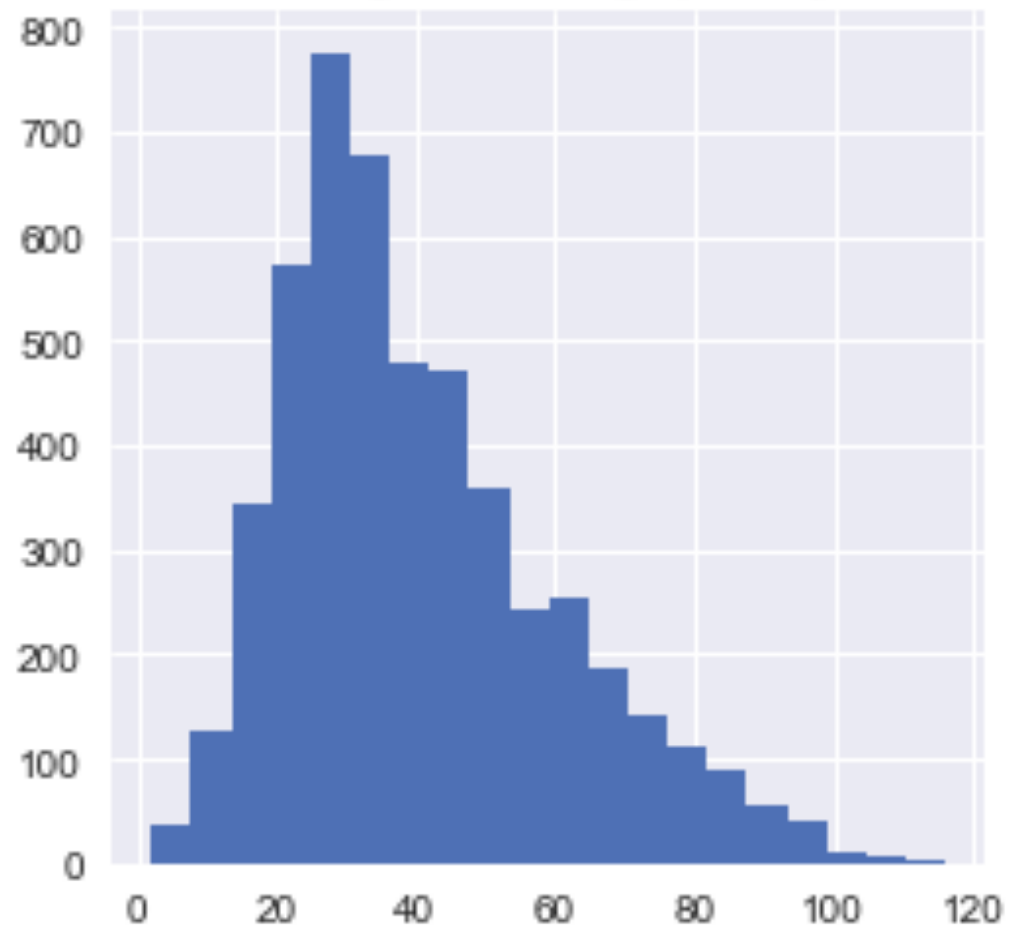Pearson correlation between variables

extremely correlated

Blues → non-premium users

Reds → premium users

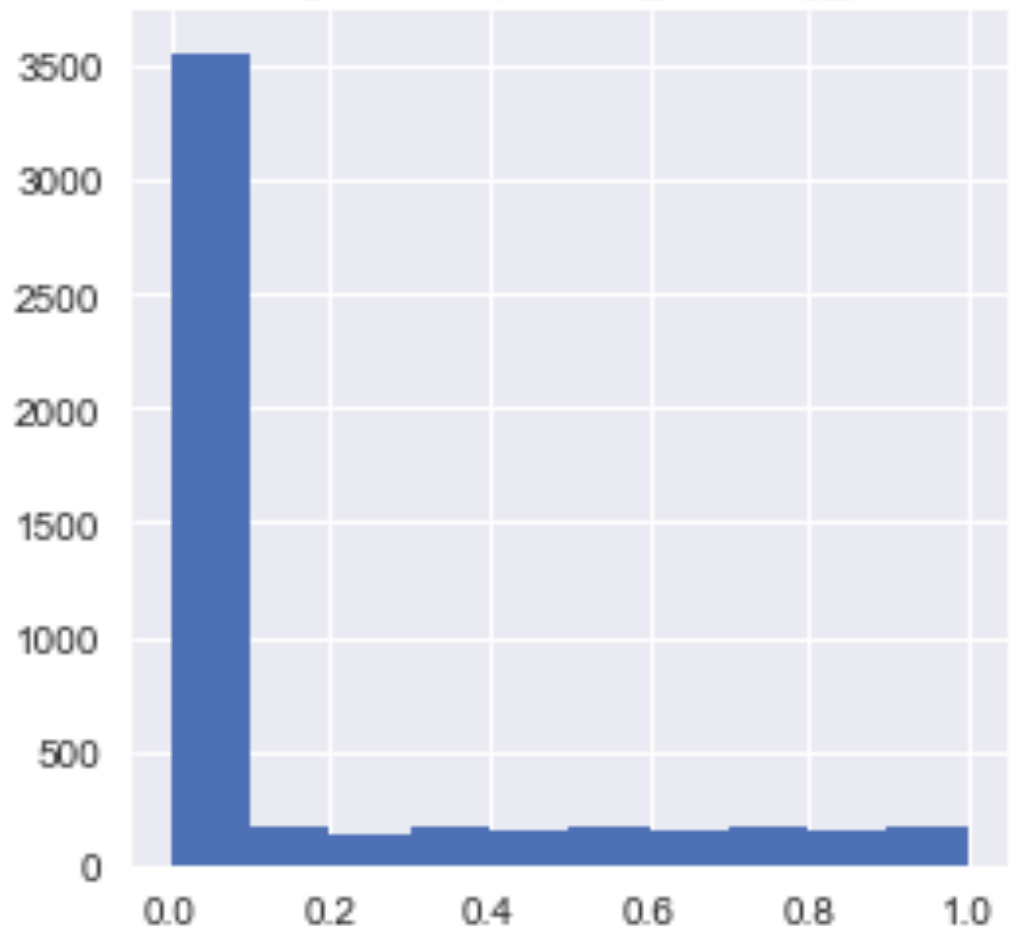Premium upgrade vs (Searches per day, Proportion of hours using app)

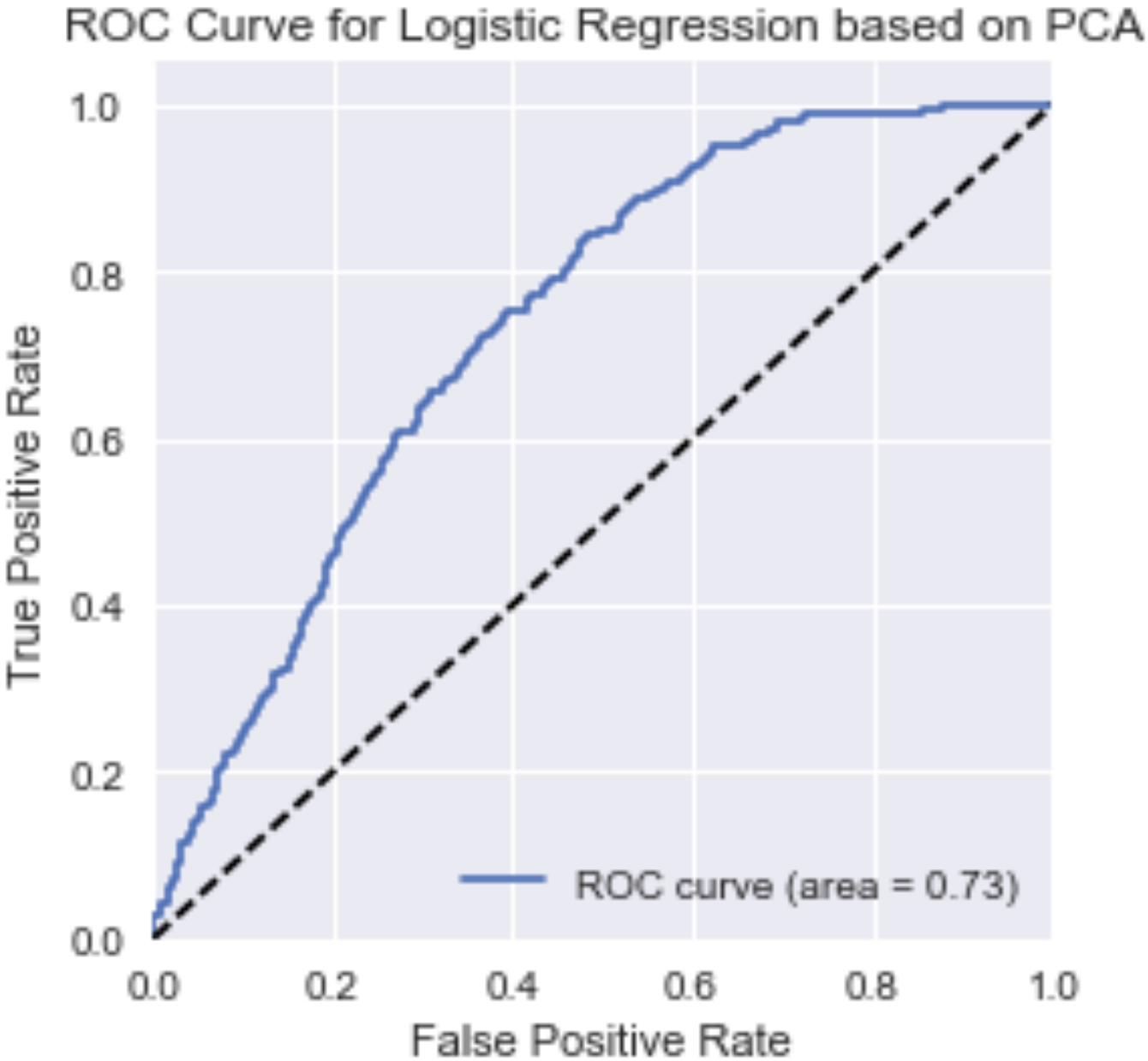# First two principal components explain more than 99% of the variance
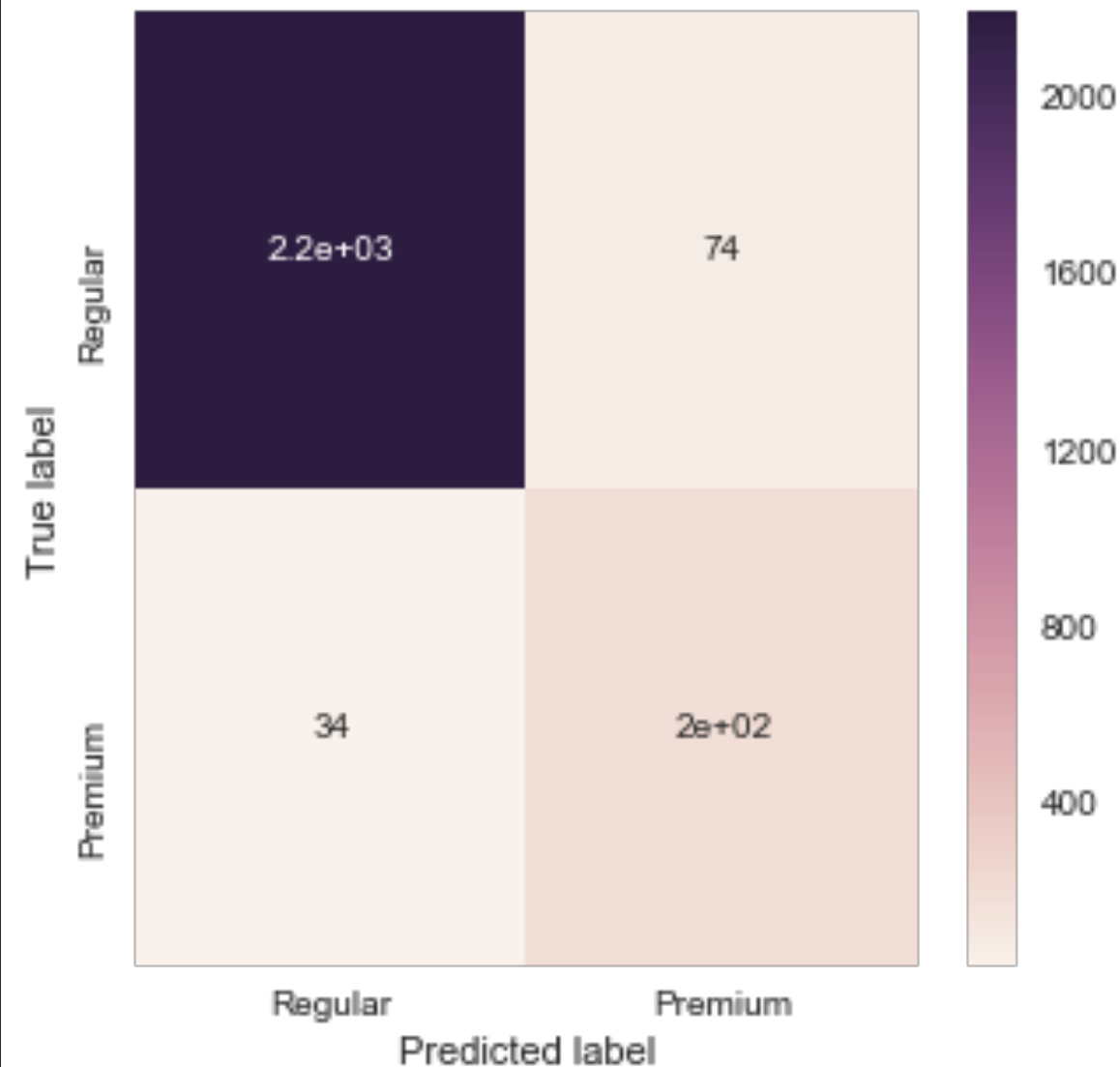
Not very good predictor for Logistic Regression



ROC Curve for Logistic Regression based on PCA



Confusion matrix for Logistic Regression Classifier based on PCA

# Logistic Regression works much better using all original variables



accuracy score = 0.957079823506
precision score = 0.728624535316
recall score = 0.852173913043

# kNN+ SVM + Decision Tree



Confusion matrix for Voting Classifier

accuracy score = 0.906137184116
precision score = 0.152416356877
recall score = 0.872340425532
cross_val_score = 0.905173909564

Best algorithm is xgboost

accuracy score = 0.973926995588
precision score = 0.810408921933
recall score = 0.939655172414
cross_val_score = 0.97321392557

# xgboost also gives us importance of each feature
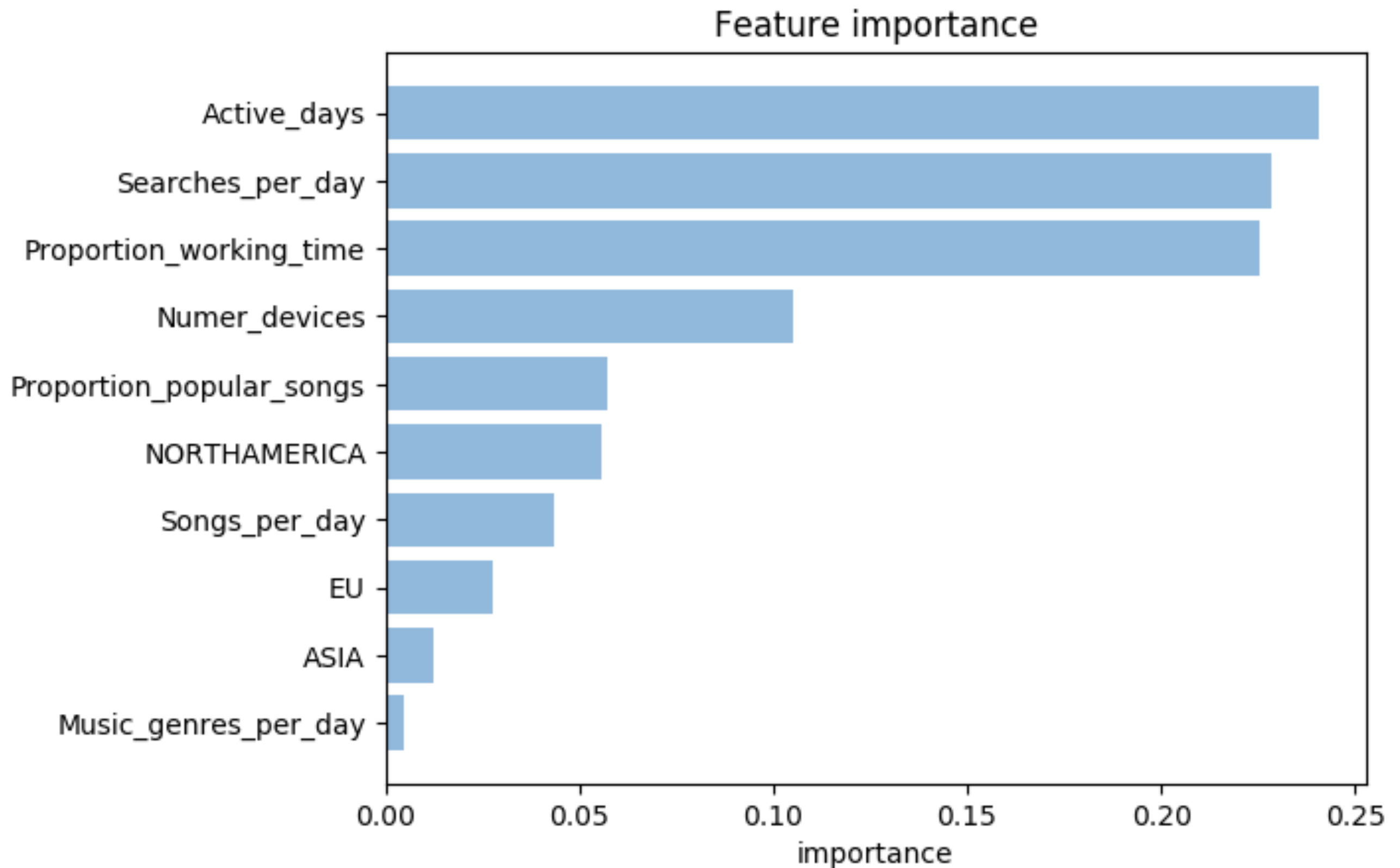


Feature importance

# Visualizing how Logit classifier works with most important two-variables



Logistic Regression performance

## Premium users

| | average |
|---|---|
| Active_days | 157.054745 |
| Songs_per_day | 51.841241 |
| Numer_devices | 2.666058 |
| Proportion_working_time | 0.615914 |
| Searches_per_day | 14.416058 |
| Proportion_popular_songs | 0.011327 |
| NORTHAMERICA | 0.093066 |
| EU | 0.058394 |
| ASIA | 0.043796 |

## Candidates for upgrade

| | average |
|---|---|
| Active_days | 122.868000 |
| Songs_per_day | 38.254000 |
| Numer_devices | 2.224000 |
| Proportion_working_time | 0.097731 |
| Searches_per_day | 7.494000 |
| Proportion_popular_songs | 0.017607 |
| NORTHAMERICA | 0.010000 |
| EU | 0.002000 |
| ASIA | 0.004000 |

selected from highest probabilities

**What kind of jobs do this people work at?. I would like more info on their workplace and job-status**



Histogram Proportion_working_time