# Project 1: DMP system

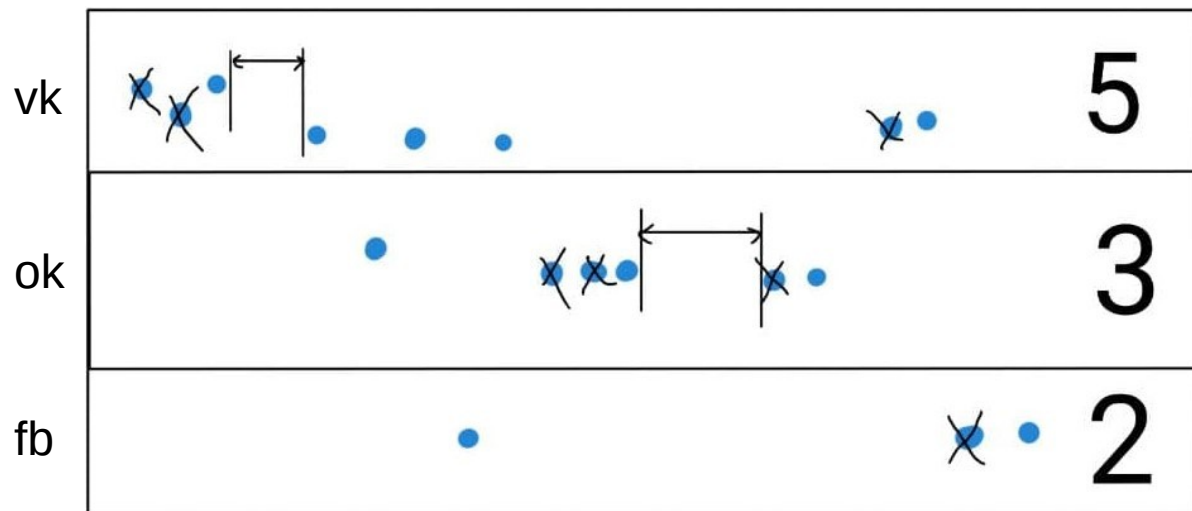Mukukenova Victoria
Prutko Alexander

# Contents

- Features
  - Feature cleaning
  - Feature engineering
- Model / Conclusion

# Features: feature cleaning

- url → domain
- unique visitor counts
- remove rare domains
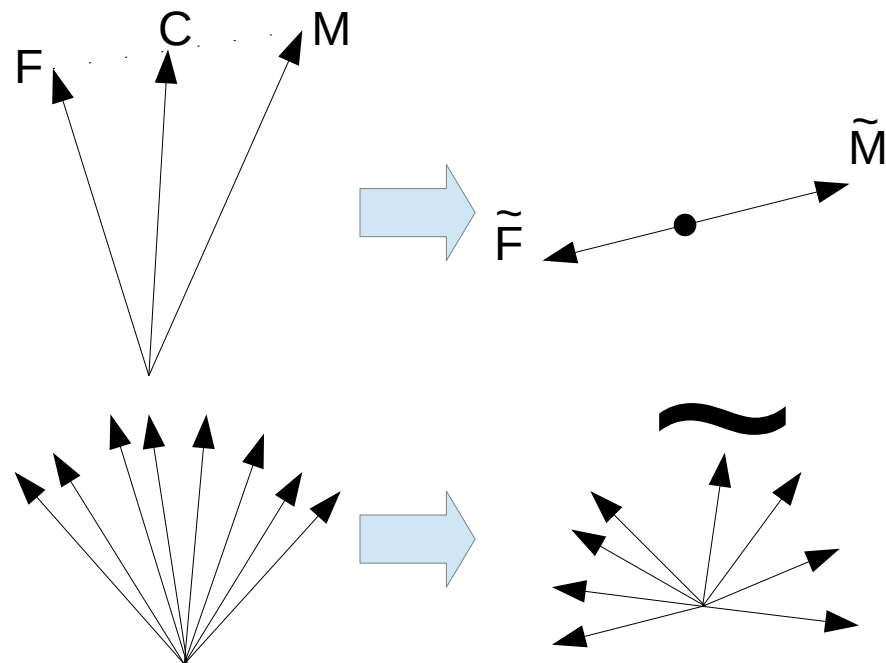- remove common domains (tf-idf stop words)

# Features: timestamps

- timestamp → [day of week, weekend, day, hour] visit counts
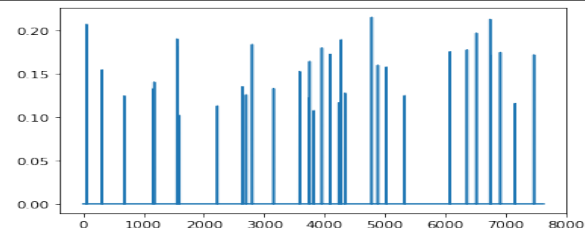
- timestamp → visit sessions
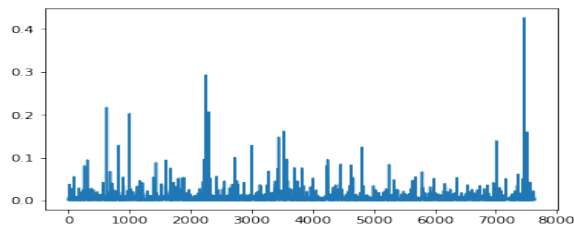
# Features: tf-idf and vectors

- domain string → tf-idf vector
- mean {Female, Male} vector
- mean Common vector
- subtraction of Common vector
- cosine similarity between object vectors and F,M-vectors
- common domains from common vector → stop words
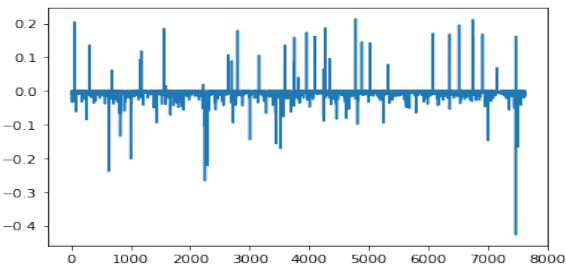- the same with age vectors and category vectors
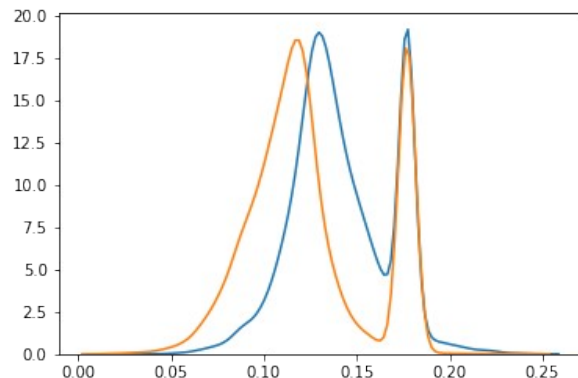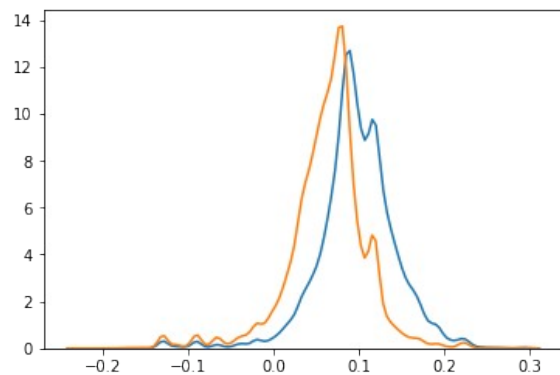
# Features: cosine similarity
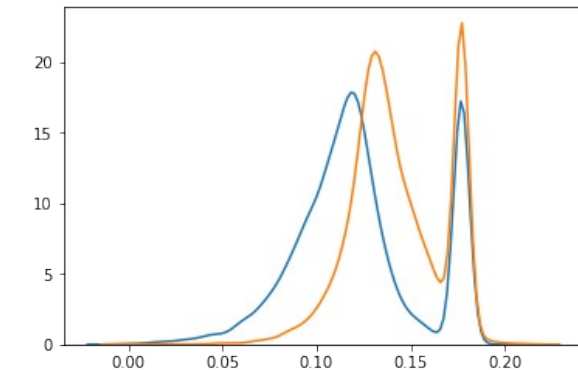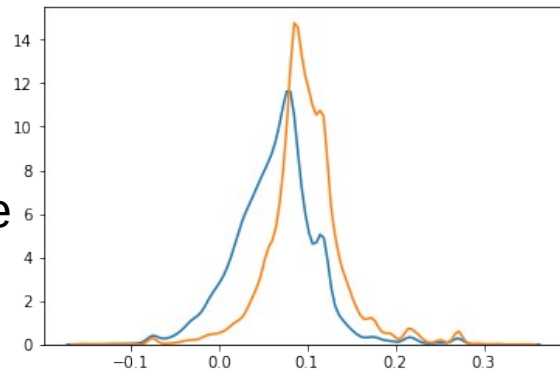


vector

common vector

vector after subtraction

cos sim to male

cos sim to female

all domains

all domains - stop words

# Model / Conclusion

- CatBoost
  - loss_function: MultiClassOneVsAll
  - max_depth: 2
  - subsample: 0.7
  - colsample_bylevel: 0.6
  - learning_rate: 0.05
  - n_estimators: 200

```
              Train                        Test
F0 08.9235%   189 / 2118     F0 07.9427%    61 / 768
F1 63.6293%  3254 / 5114     F1 57.4240%   963 / 1677
F2 21.1250%   676 / 3200     F2 16.7134%   179 / 1071
F3 17.6471%   351 / 1989     F3 12.6645%    77 / 608
F4 00.0000%     0 / 673      F4 00.0000%     0 / 222
M0 00.0000%     0 / 1507     M0 00.0000%     0 / 505
M1 59.3100%  3851 / 6493     M1 57.6622%  1253 / 2173
M2 26.8127%  1028 / 3834     M2 24.7809%   311 / 1255
M3 00.0000%     0 / 1592     M3 00.0000%     0 / 555
M4 00.0000%     0 / 583      M4 00.0000%     0 / 201
------------------------     -----------------------
-- 34.4943% 9349 / 27103     -- 31.4776% 2844 / 9035
```

# Questions?



No questions you have

# Features: cosine similarity