# "Strax" — Advancements in Topic Modeling
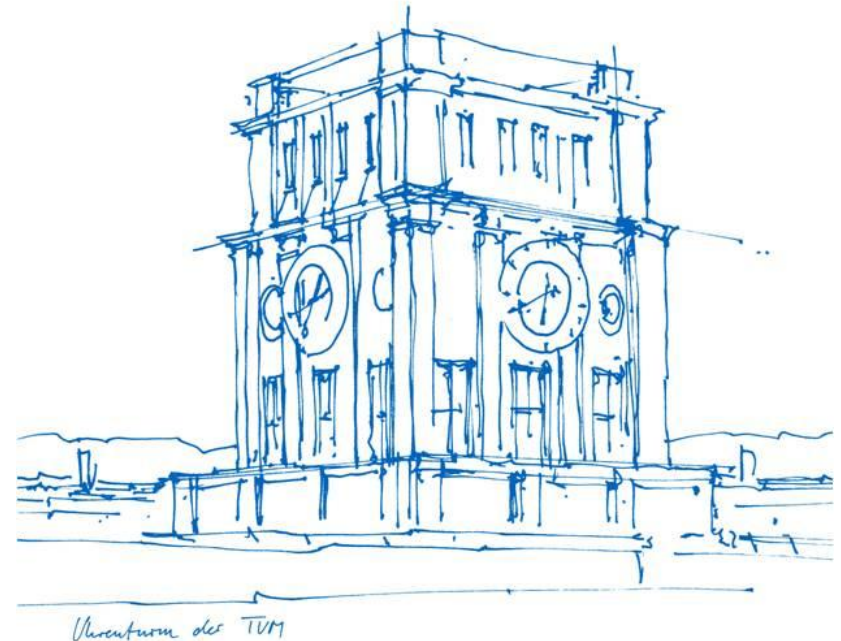
Technische Universität München

Fakultät für Informatik

NLP Lab Course, SS22

14.07.2022
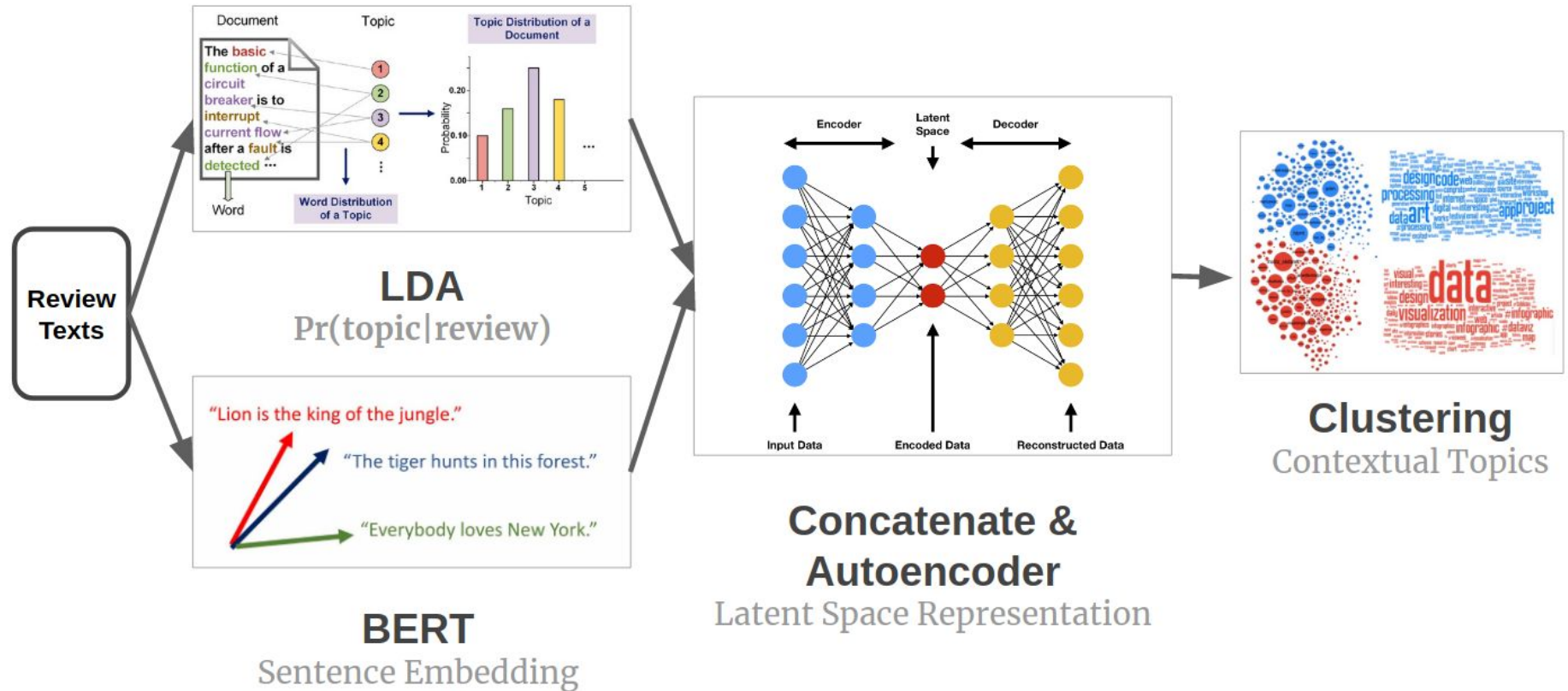
Berk Sudan, Ferdinand Kapl, Yuyin Lang

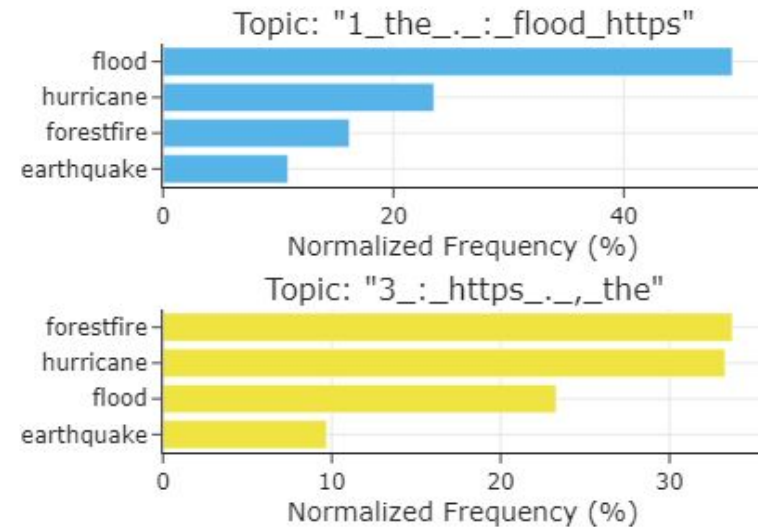Uhrenturm der TUM

# Overview

- LDA-BERT
- Visualization

- Yahoo Dataset & Visualizer

- Current Status

# LDA-BERT



LDA
Pr(topic|review)

BERT
Sentence Embedding

Concatenate &
Autoencoder
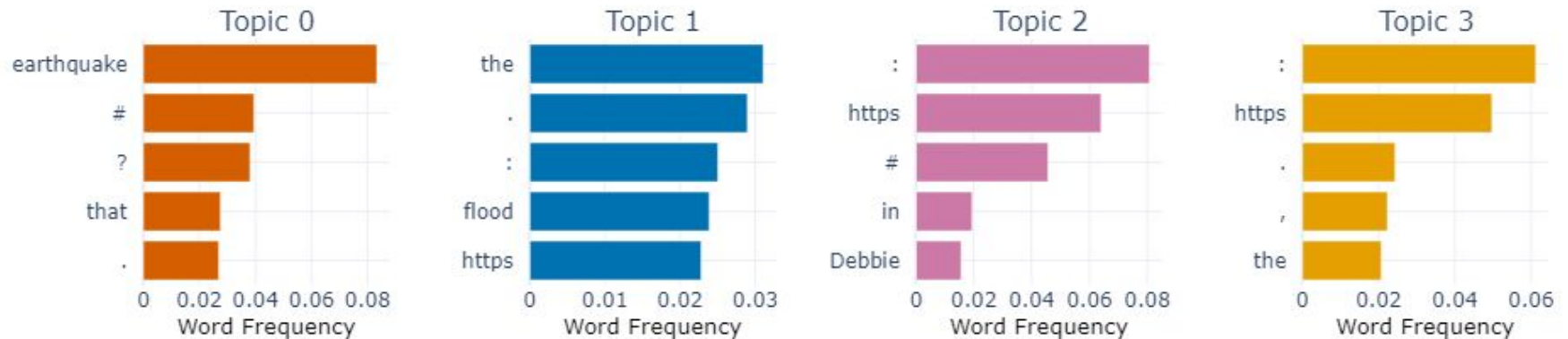Latent Space Representation

Clustering
Contextual Topics

# LDA-BERT: Current Results



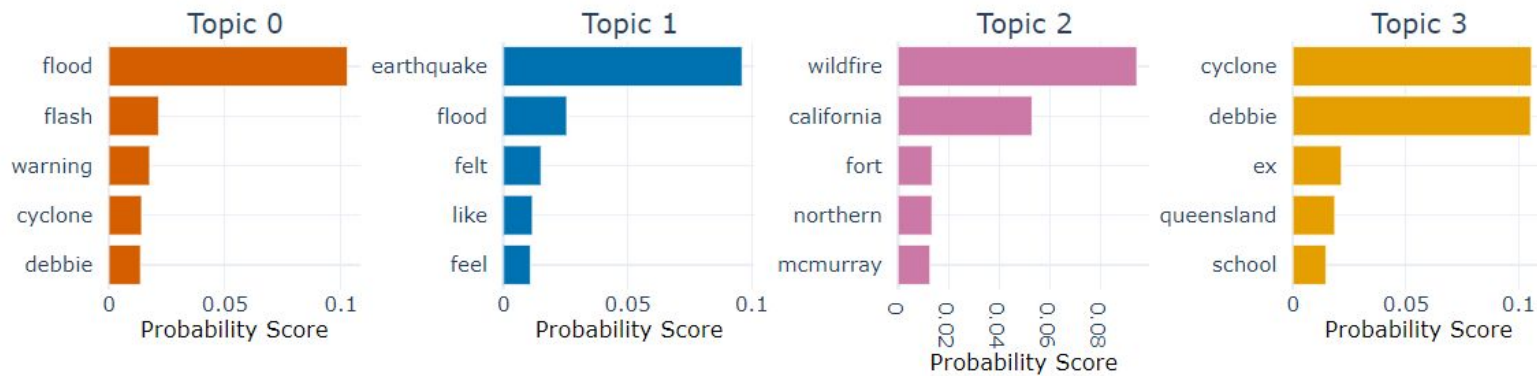Labels per Topic for algorithm="lda-bert", run_id="1657804665"

# LDA-BERT: Current Results



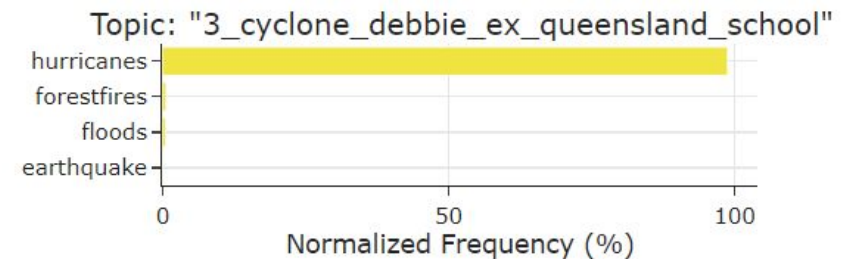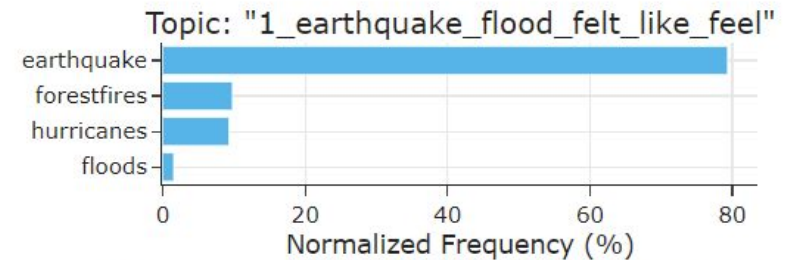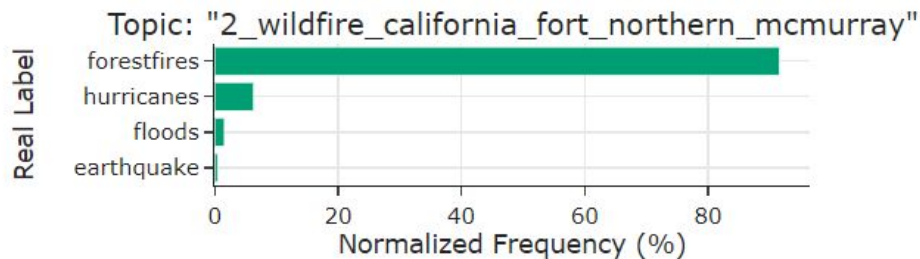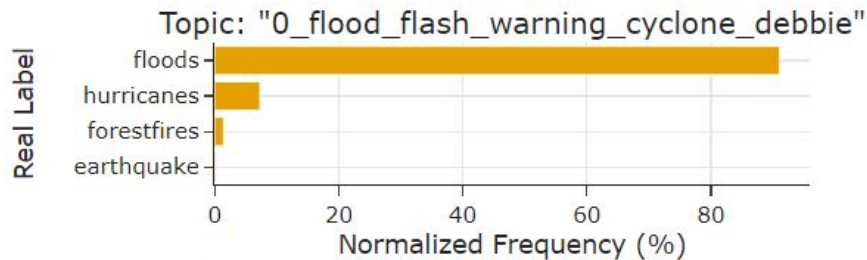Topic Word Scores for algorithm="lda-bert" and run_id="1657804665"

# Visualization of LDA and NMF



**Topic Word Scores for algorithm="NMF" and run_id="1657397615.511795"**

# Visualization of LDA and NMF



els per Topic for algorithm="NMF", run_id="1657397615.511795"

# Visualization of LDA and NMF (for comparison)

LDA:

```
Topic 0:
hurricanes      1253
floods           449
forestfires      347
earthquake        44
Name: Real Label, dtype: int64
-------------------------------
Topic 1:
forestfires     1136
floods           498
hurricanes       265
earthquake       116
Name: Real Label, dtype: int64
-------------------------------
Topic 2:
floods           555
hurricanes       233
forestfires      230
earthquake       140
Name: Real Label, dtype: int64
-------------------------------
Topic 3:
earthquake      1195
floods           474
forestfires      191
hurricanes       180
Name: Real Label, dtype: int64
-------------------------------
```

NMF:

```
Topic 0:
floods          1910
hurricanes       153
forestfires       30
earthquake         4
Name: Real Label, dtype: int64
-------------------------------
Topic 1:
earthquake      1481
forestfires      183
hurricanes       174
floods            29
Name: Real Label, dtype: int64
-------------------------------
Topic 2:
forestfires     1682
hurricanes       116
floods            29
earthquake         9
Name: Real Label, dtype: int64
-------------------------------
Topic 3:
hurricanes      1488
forestfires        9
floods             8
earthquake         1
Name: Real Label, dtype: int64
-------------------------------
```

# Glance at optional algorithms

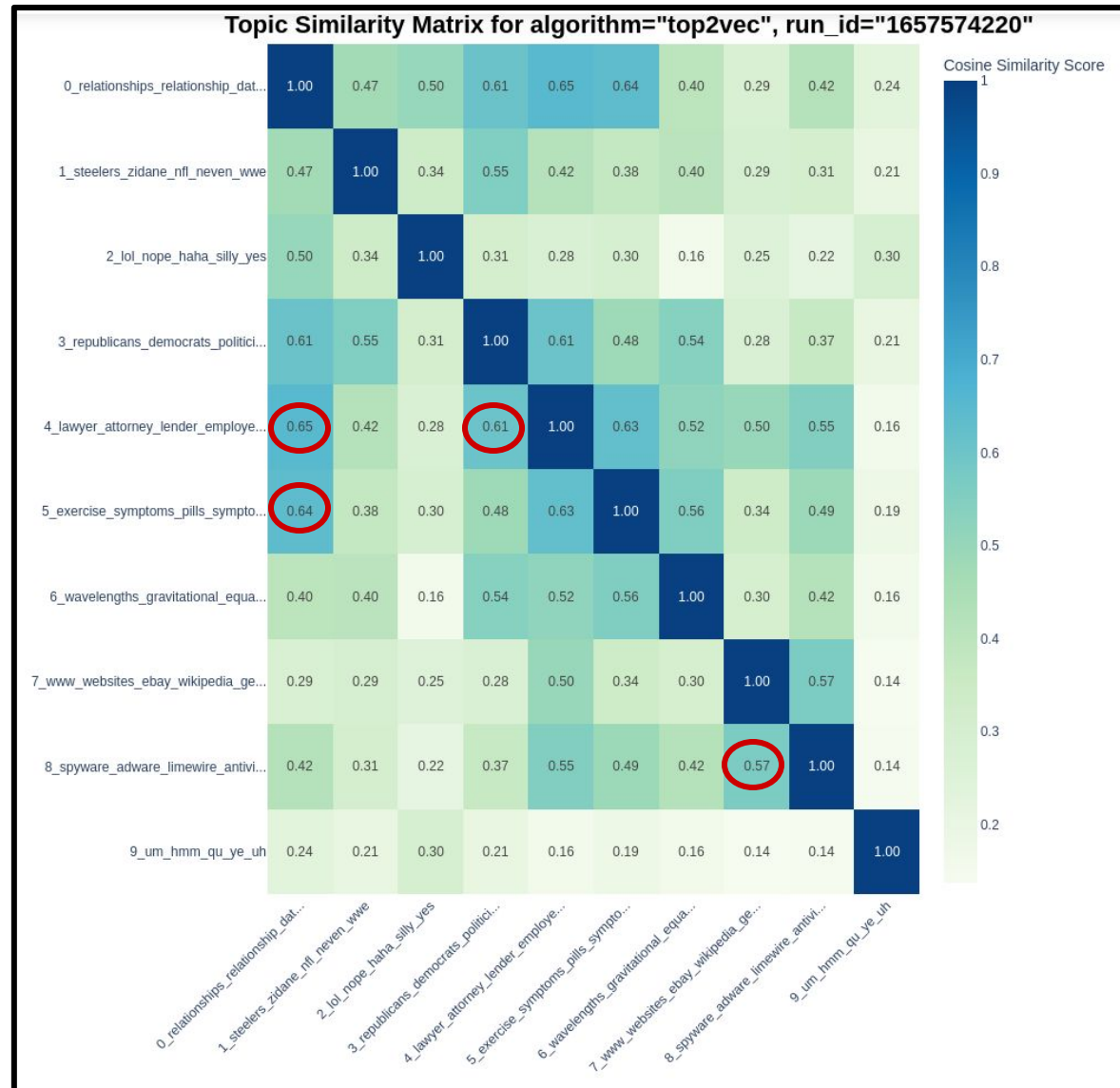| Algorithms | Special for | Code availability |
|---|---|---|
| LSA | - | Basic code |
| PLSA | - | Basic code in Enstop |
| BTM | Short text | No python |
| DMM (GSDMM) | Short text | Basic code |
| WNTM | Short text | No python |
| CTM | - | OCTIS, need to read source code for more information |

# Yahoo Dataset - Characteristics

- A Long Text Dataset

- # of Classes: 10

- # of Instances: 60K (selected randomly 300K+)

- Instance / (Real Label): 6K

- Excerpt from the data:

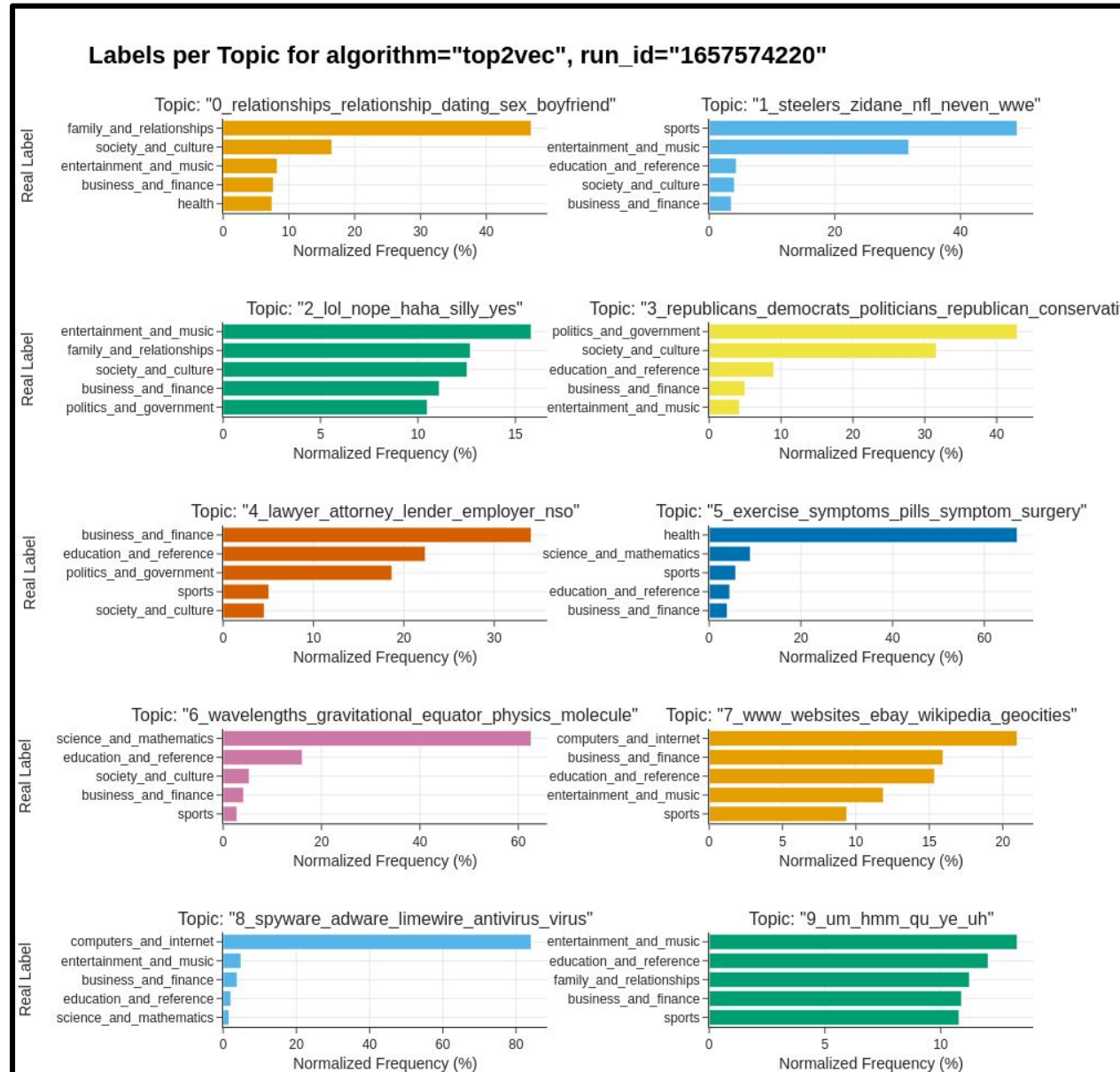| | Document ID | Document | Real Label | Assigned Topic Num | Assignment Score |
|---|---|---|---|---|---|
| 0 | 25095 | Being a man, I would say just go up to him and tell him how you feel. There is nothing I hate more than a guessing game. Just put it all out there, and you will get a response, even if it isn't one you want. Holding that in just makes things worse later on for you, and him. I like it when a girl comes to me to tell me she wants me. Even if im not interested in that girl, i admire that. | family_and_relationships | 0 | 0.700406 |
| 1 | 29164 | it all depends on how well you know him. find out everything about him, but do it discretely. only then will you what to do, and it wouldn't hurt to find out from his friends if he likes you as well. cuz there is no point to baring your heart to him, if he doesn't feel the same way about you, now is there? just sit back and relax and do your homework on the guy, it will be that much better, trust me. | family_and_relationships | 0 | 0.683313 |
| 2 | 29843 | excuse the other folks' rude answers. it sounds like you are very worried about this. do yourself a favor and relax. do you think maybe he's just telling you that to see if you'll tell him how much you care about him? or do you know for a fact his family may be planning to move there? could be just a head game, i don't k now you guys' ages. so many things we worry about never even happen, relax. if in fact he does move, and you find yourself with a broken heart, just try to keep really busy with your friends and family. if you two are meant to be together it will happen t hat way. and if not, you will fall in love with someone else and forget all about him if it'snot meant to be. good luck | family_and_relationships | 0 | 0.674734 |

# Yahoo Dataset - Topic Word Scores

# Yahoo Dataset - Topic Similarity Matrix (New)

# Yahoo Dataset - Labels per Topic (New)



Labels per Topic for algorithm="top2vec", run_id="1657574220"

# Current Status & Next

**DONE:**
- <u>ALGs:</u> LDA, NMF, Top2Vec, BERTopic, LDA-BERT
- <u>Visualizations:</u> Topic words barcharts, Real label - Topic distribution, Umap embedding plots, Similarity Heatmap

**NEXT:**
- <u>Preprocessing</u>
- <u>Evaluation Module:</u>
  - Coherence Metrics: C_V Coherence & NPMI Coherence
  - Diversity Metrics: Unique Words Topic Diversity & Kullback-Liebler Divergence
  - Supervised Metrics: Rand Index

# Questions Time!!