



Project 3

Newsgroups NLP Classification

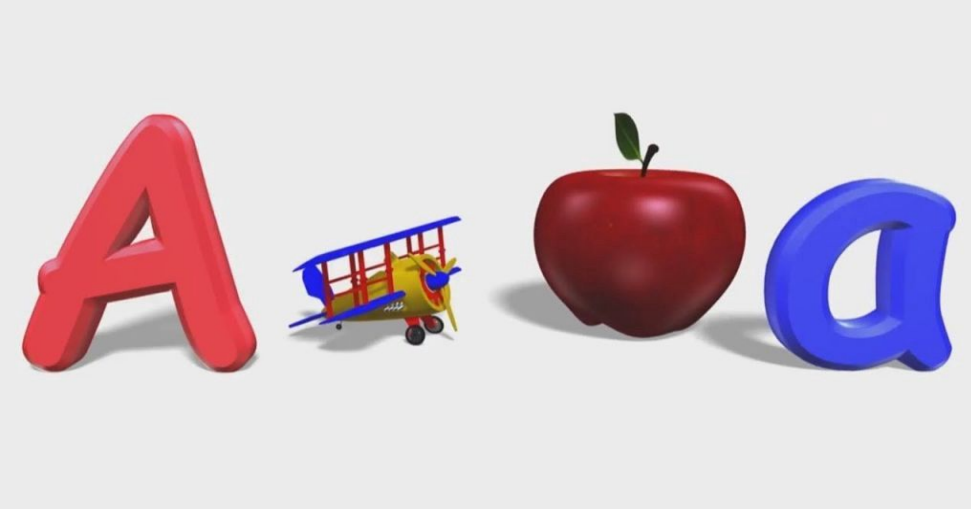
What are we trying to accomplish?

- Automatically Import fresh articles into newsgroup discussions

Why is this important?

- Increase usage and relevance of newsgroups

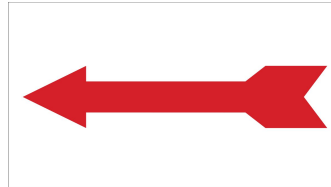




How it Works:

- Use NLP model to cluster existing newsgroup data
- Many iterations to establish some cluster separation
- TRY to associate Topic Clusters with Newsgroups

Build a Classifier using these associations



Dataset

18,846



NewsGroups

- comp.graphics
- comp.os.ms-windows.misc
- comp.sys.ibm.pc.hardware
- comp.sys.mac.hardware
- comp.windows.x
- rec.autos
- rec.motorcycles
- rec.sport.baseball
- rec.sport.hockey
- sci.crypt
- sci.electronics
- sci.med
- sci.space
- misc.forsale
- talk.politics.misc
- talk.politics.guns
- talk.politics.mideast
- talk.religion.misc
- alt.atheism
- soc.religion.christian

A lot of
these are
very
similar!

Bag of Words

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



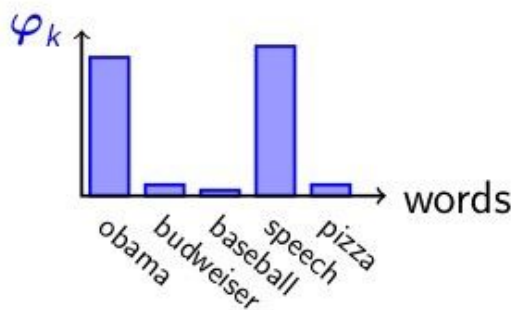
it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

Latent Dirichlet Allocation

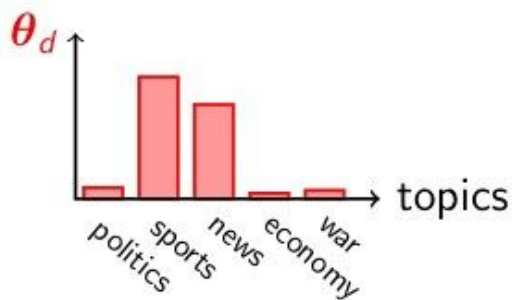
LDA discovers topics into a collection of documents.

LDA tags each document with topics.

Topic k



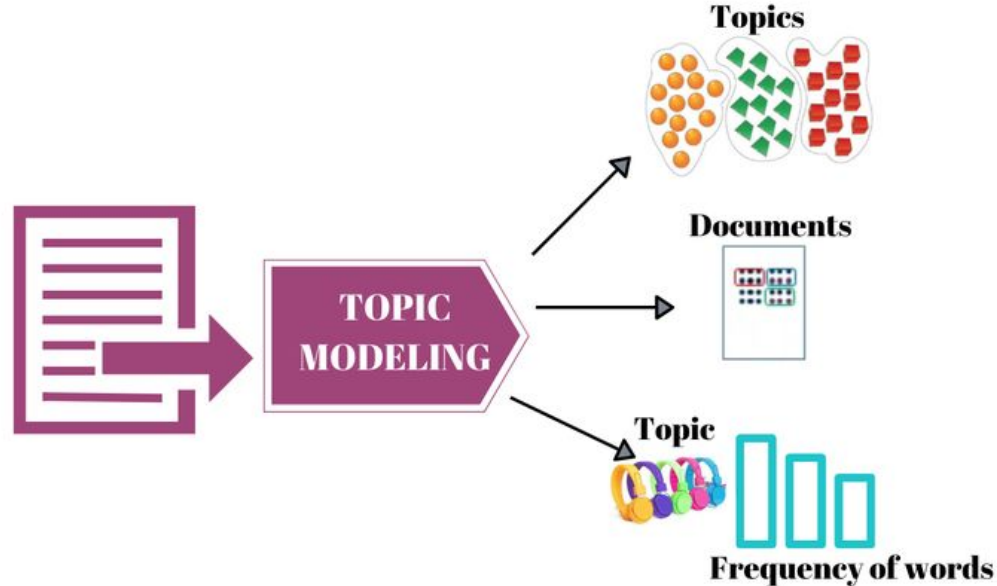
Document d



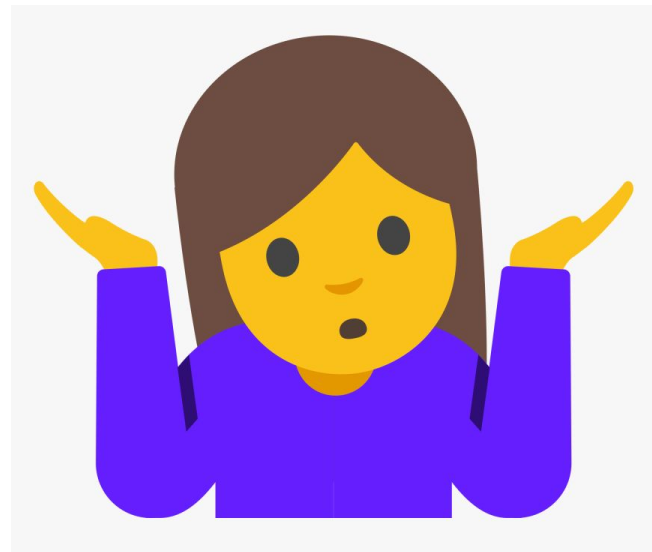
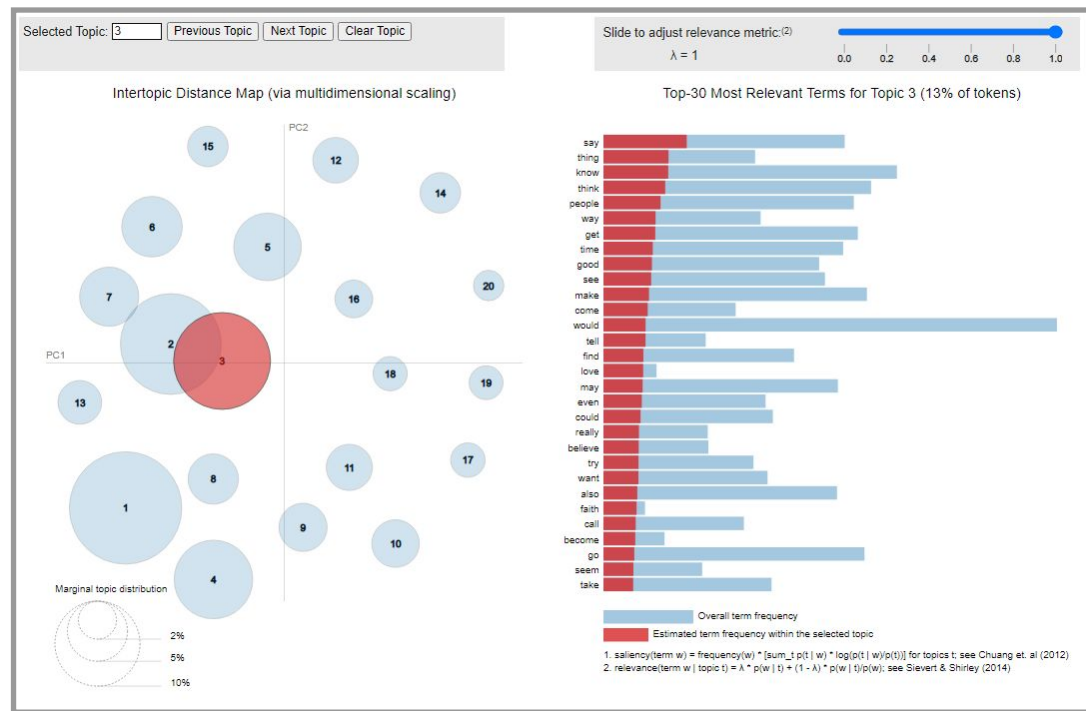
- Uses Bag of Words to identify topics in document
- Not a linear model, LDA came from research in genetics
- Assumption is that documents contain several topics

Technology and Model

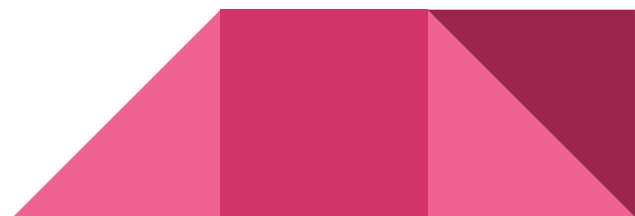
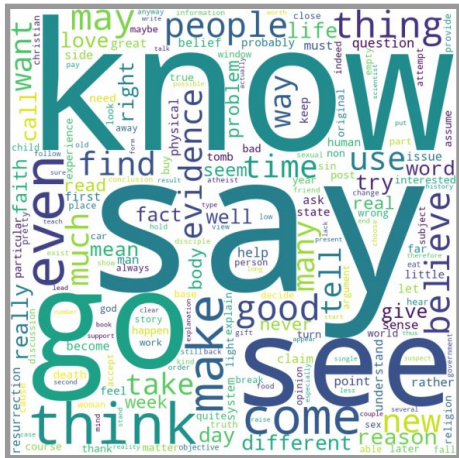
PYCARET



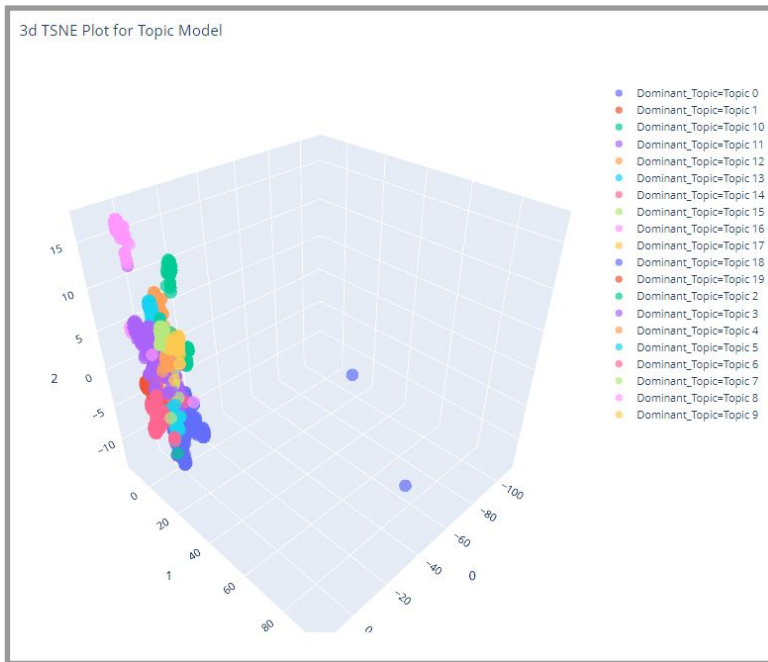
Experiment 1 - Topic Model



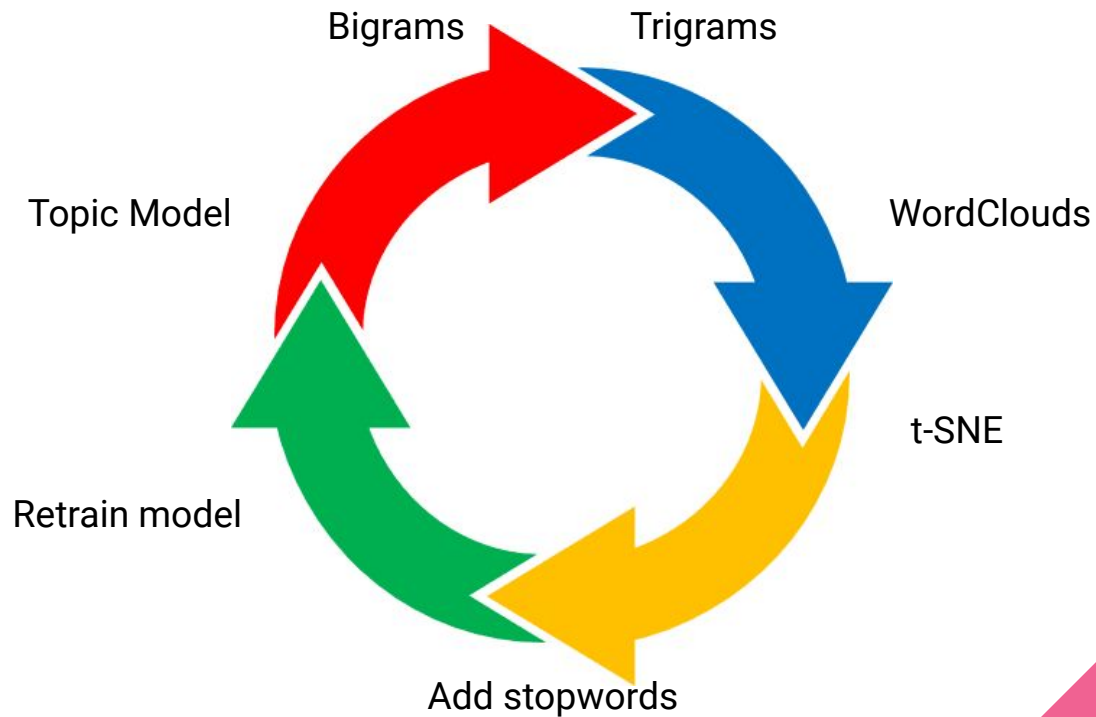
WordCloud



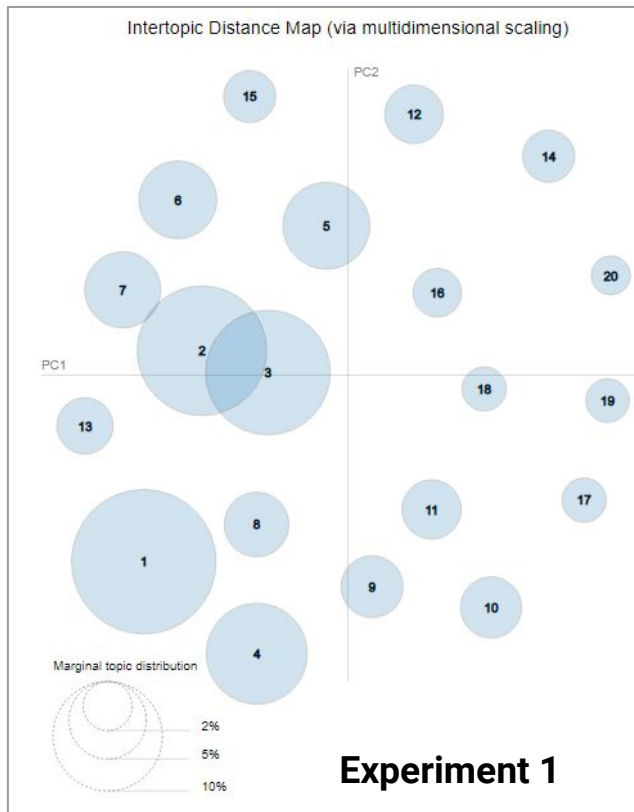
t-distributed Stochastic Neighbor Embedding (t-SNE)



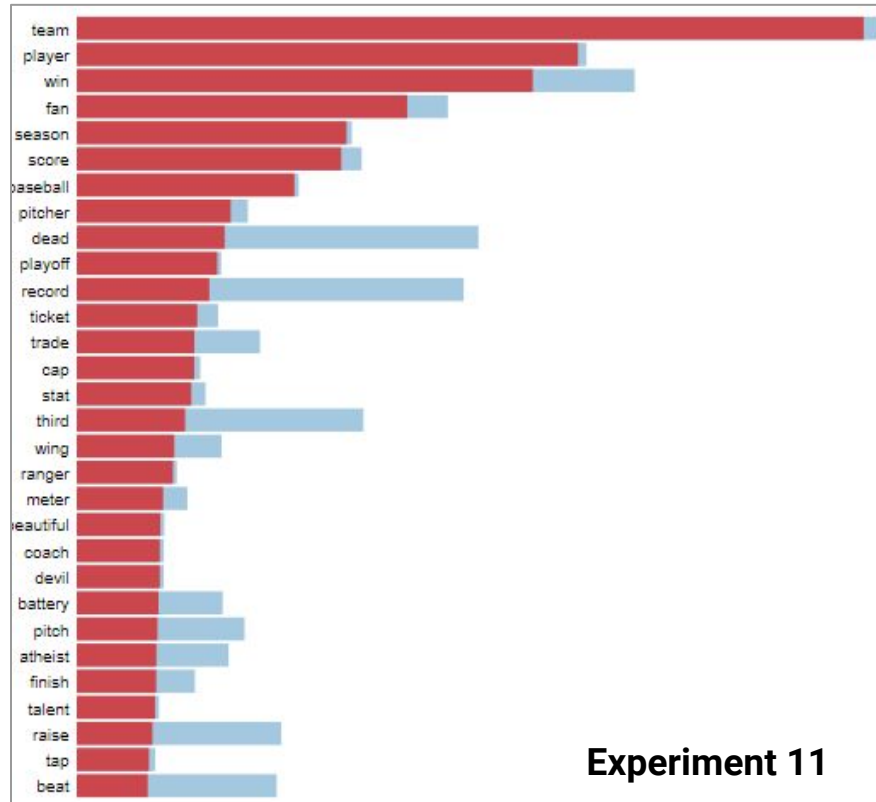
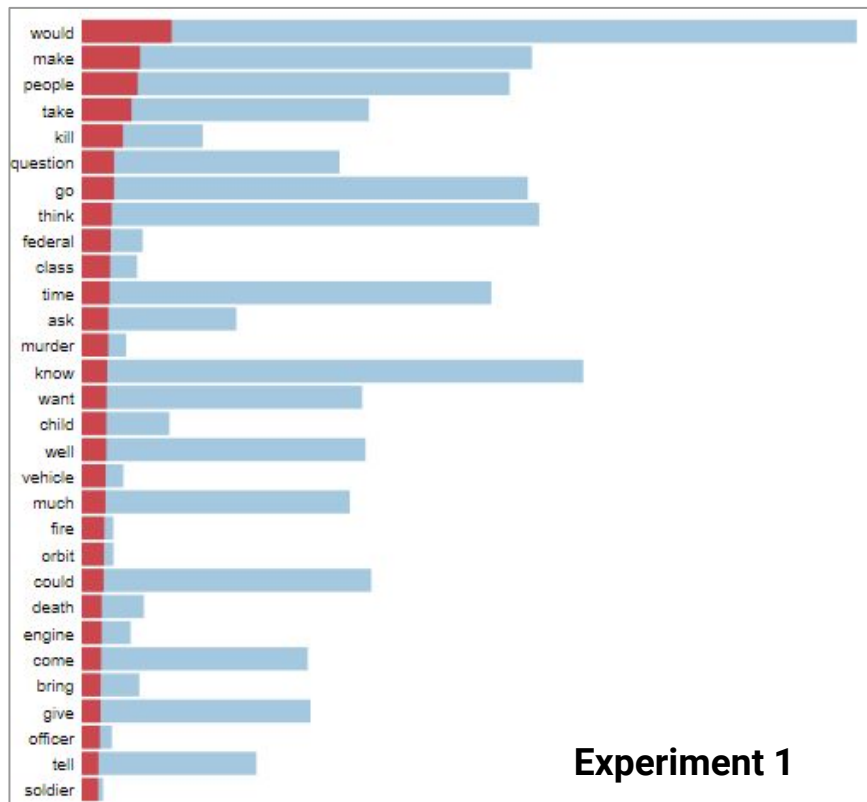
Experiment Methodology



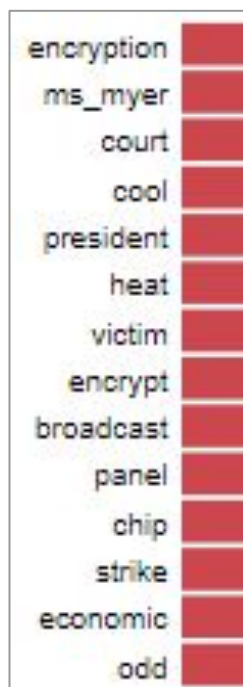
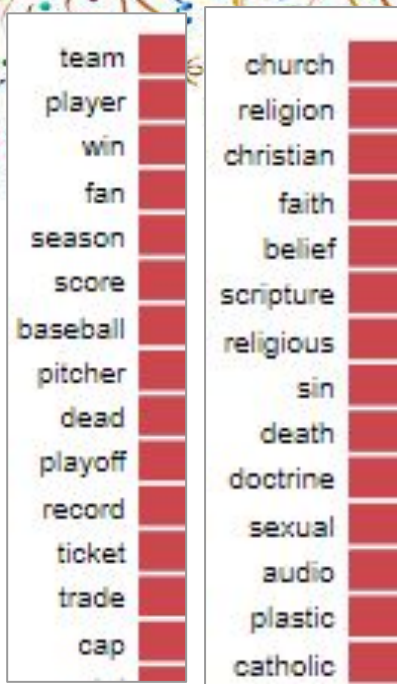
Topic Model - Comparison IDM



Topic Model - Term Comparison



Themes



Bigrams



loss tie
division finish
dead dead
coach corner
camera camera
battery camera
film cassette
innocent murder
ice stat
hate team
bike storage
stat player
goal assist
score goal
rod brind_amour
pitcher pitch
fan fan
score team
team player
team win



Trigrams

injury prone loose
suffer rush overweight
adirondack skipjack ranger
vdc tip conductor
win loss tie
team win team
team score team
player ridiculous absurd
skipjack ranger devil
record win loss
ranger devil oiler
player flyer season
loss tie team
islander adirondack wing
idle conductor vdc
ice team coach
cop fed listen
conductor vdc tip
canuck goal jet
adirondack wing adirondack
team record win
team player team
player ice team
win team win



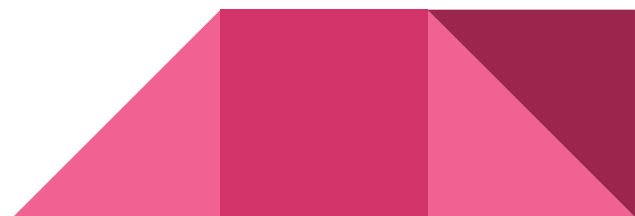
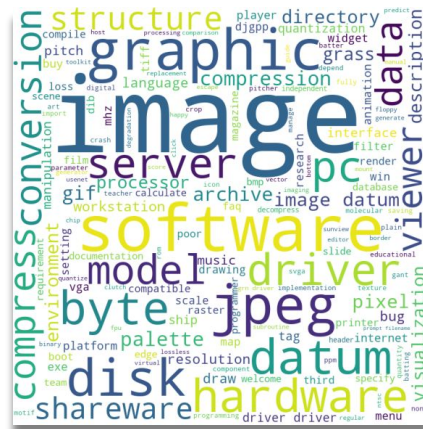
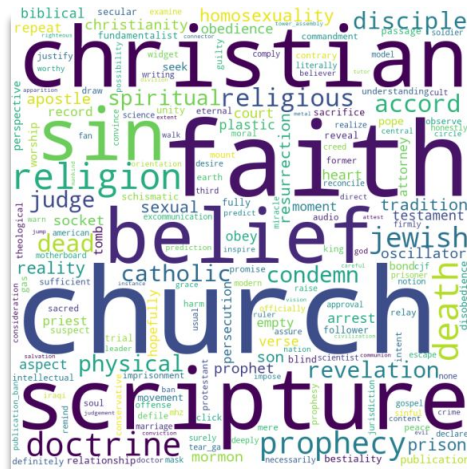
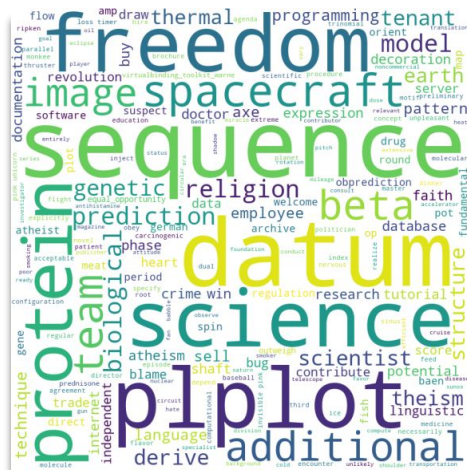
Trigrams - not all converge



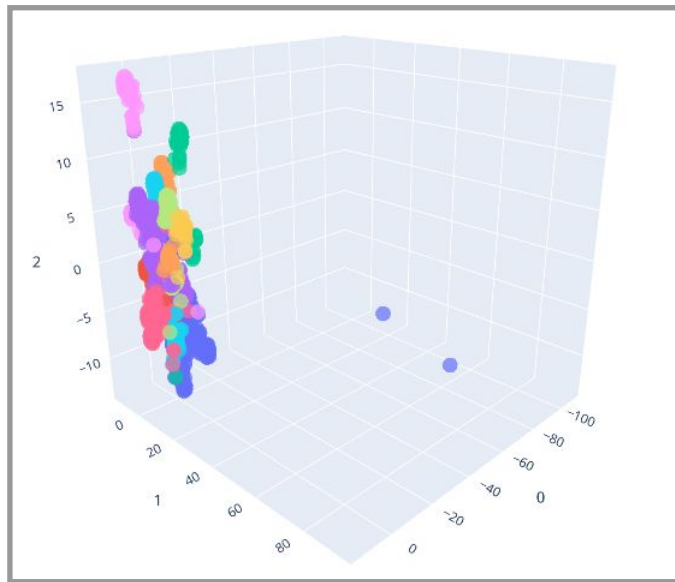
relevant datum science
roll heart soul
round round quiet
science datum image
theism theism favor
touch possession obsession
genetic algorithm genetic
drag dirty hall
danced zombie nervous
coupon coupon coupon
bam boom touch
hall oate bam
hooter nervous danced
plant stir journey
oate bam boom
monkee monkee mediate
journey raise suzanne
identify relevant datum
invisible pink unicorn



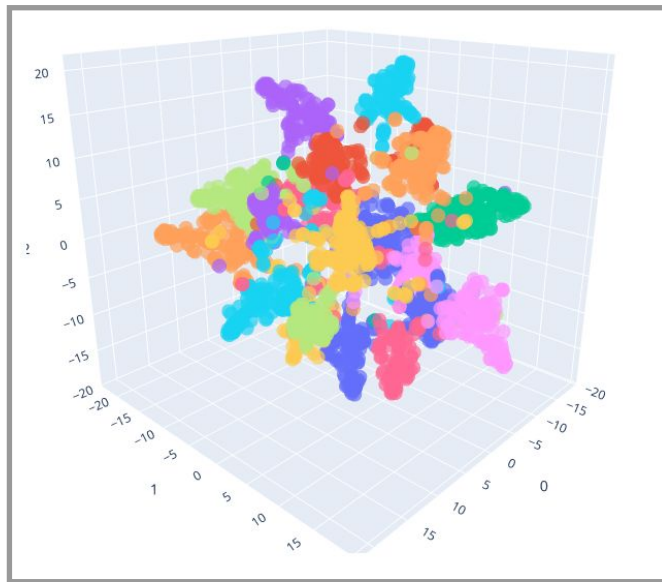
Wordcloud



t-SNE - Model Improvement



Experiment 1



Experiment 11

Ethics



Conclusion





Thank You