Understanding Happiness Using NLP

By: Madia & Manal

As Project 4 of SDAIA Data Science Bootcamp (T5)



Introduction Police Of the Control o

Being **happy** is one of the most fundamental requirements of a living being. Since the beginning of the human civilization, man has also been trying to develop new technologies, make new tools and improve his lifestyle for the sole purpose of attaining happiness. However, in its race of scientific endeavor and pursuit of money and luxuries, man is hardly aware of what exactly constitutes happiness.

Project Goals: Topic modelling to understand what makes Amazon's workers happy for the previous 24 hours

PROCESS & METHODOLOGY





Data

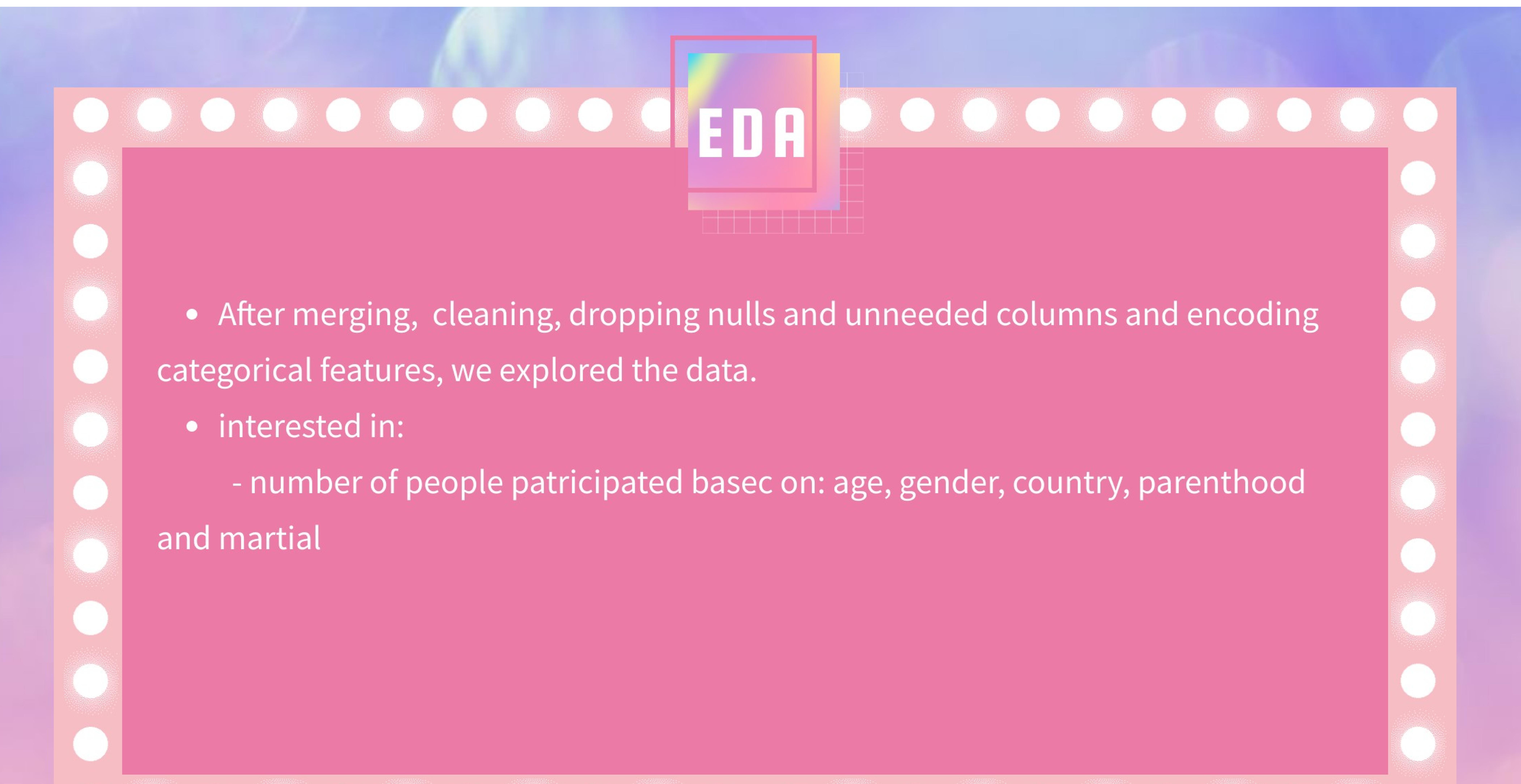
Data source: kaggle [https://www.kaggle.com/ritresearch/happydb]

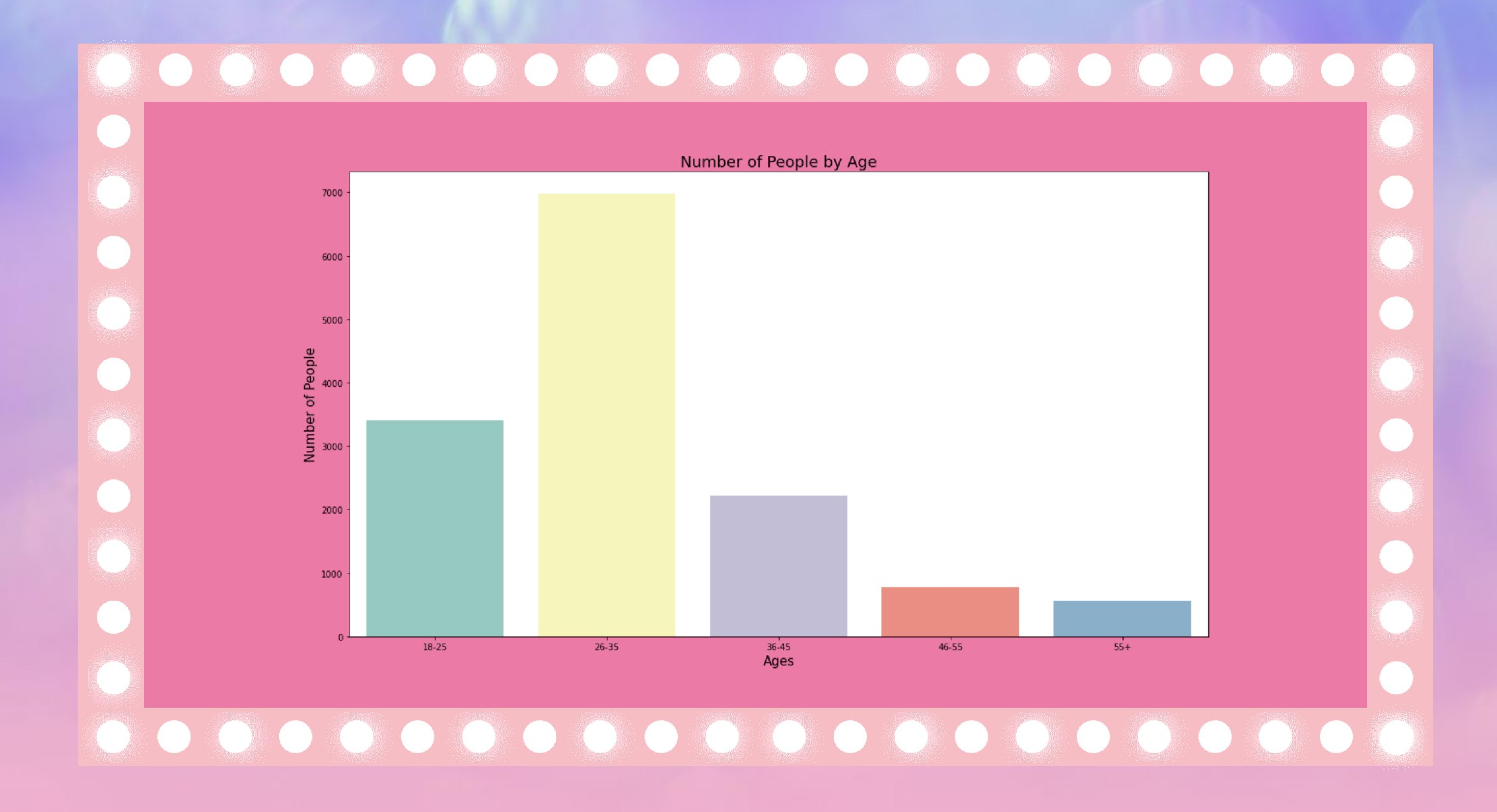
2 data files:

cleaned_hm.csv: the cleaned-up corpus of 100,000 crowd-sourced happy moments.

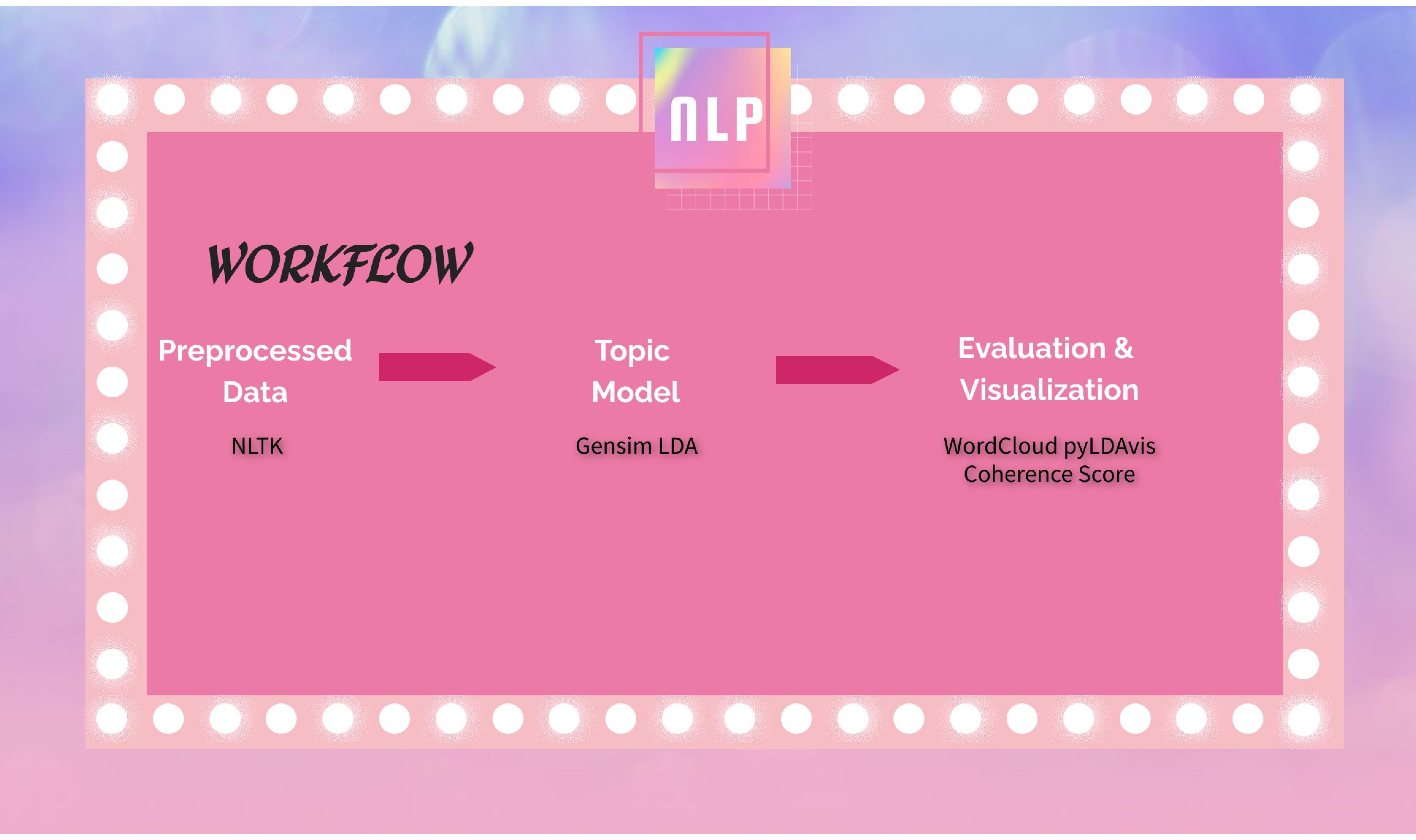
• demographic.csv: the demographic information of the worker who provided the moment.

100,535 observations & 14 features









Topics in LDA model

```
Topic: 0
Words: 0.172*"home" + 0.074*"money" + 0.070*"birthday" + 0.056*"house" + 0.049*"boss" + 0.047*"vacation" + 0.040*"girl" + 0.029*"mturk" + 0.024*"saving" + 0.022*"plan"

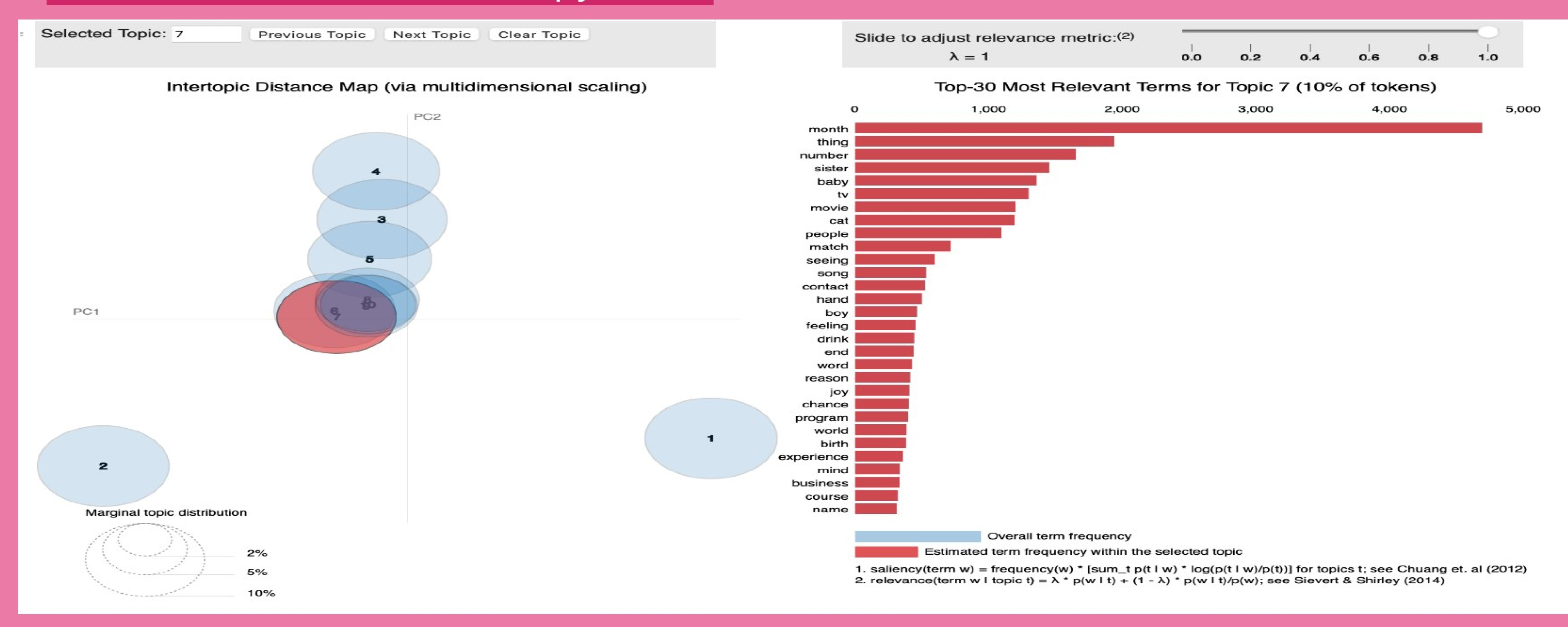
Topic: 1
Words: 0.146*"night" + 0.126*"morning" + 0.096*"wife" + 0.054*"girlfriend" + 0.054*"brother" + 0.039*"call" + 0.031*"trip" + 0.030*"sleep" + 0.027*"breakfast" + 0.025*"pizza"

Topic: 2
Words: 0.097*"game" + 0.044*"video" + 0.036*"company" + 0.034*"project" + 0.024*"team" + 0.023*"help" + 0.022*"student" + 0.021*"health" + 0.020*"music" + 0.020*"need"

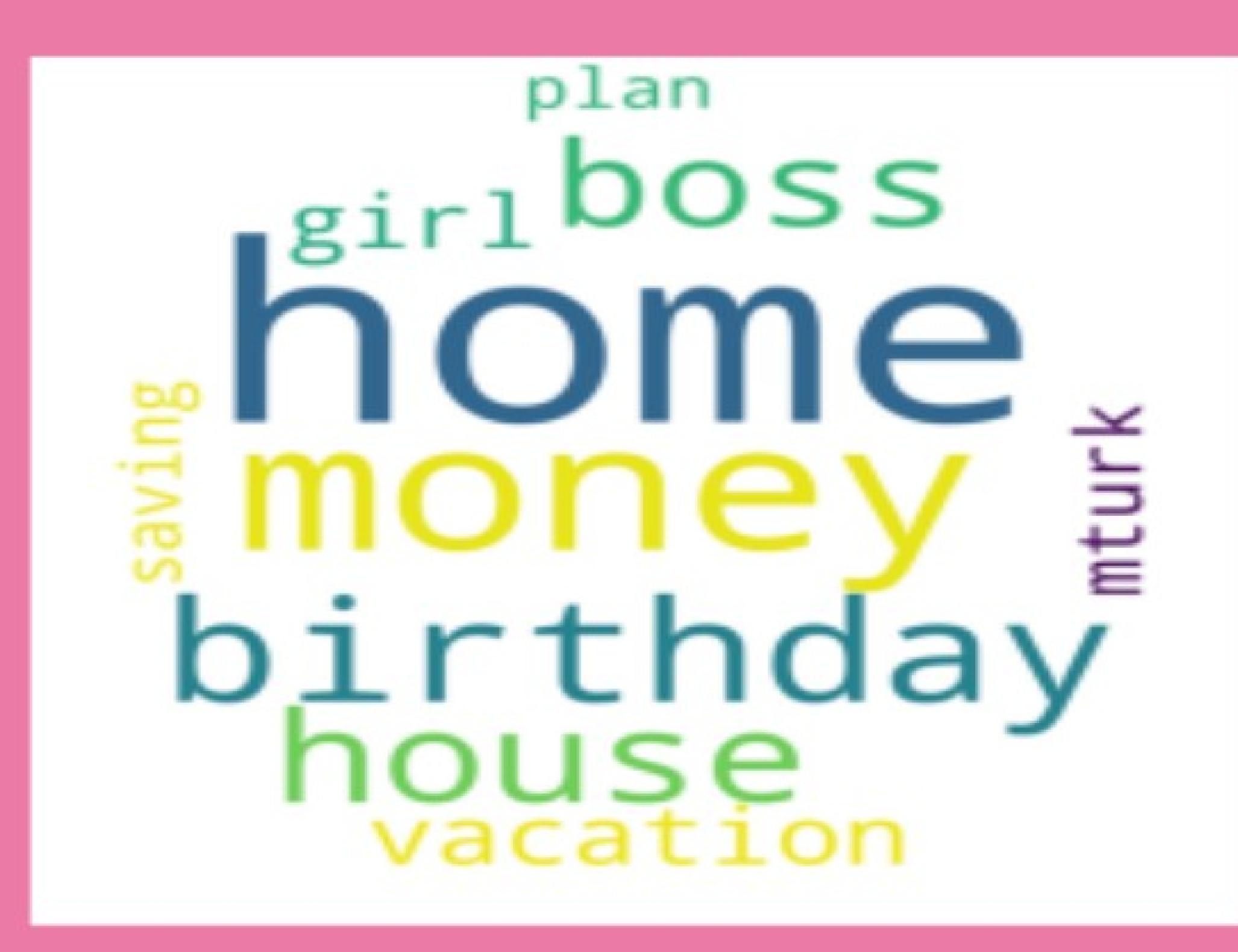
Topic: 3
Words: 0.196*"time" + 0.092*"daughter" + 0.070*"event" + 0.046*"lunch" + 0.043*"boyfriend" + 0.030*"kid" + 0.030*"child" + 0.025*"show" + 0.021*"bed" + 0.020*"way"

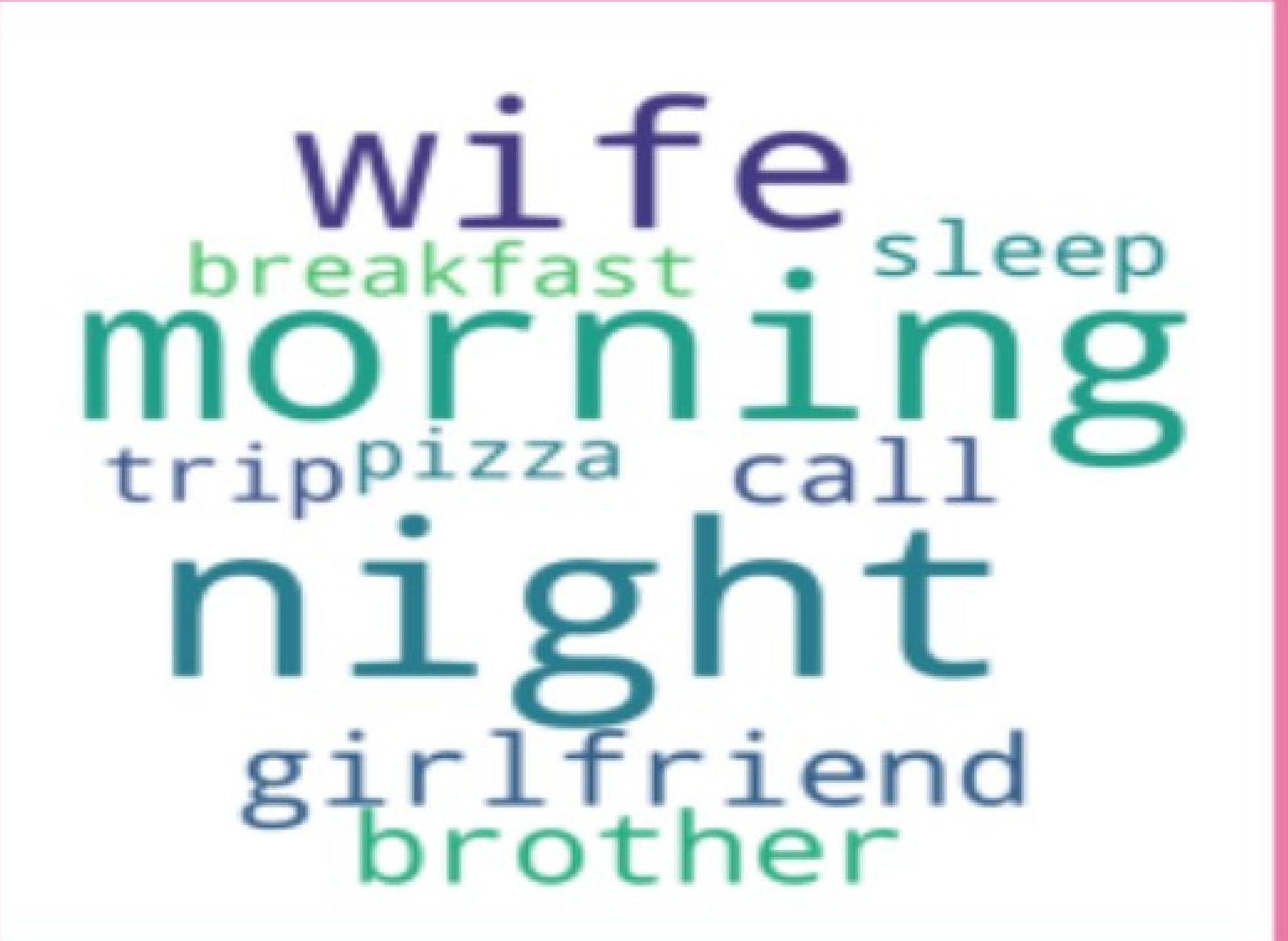
Topic: 4
Words: 0.209*"work" + 0.127*"today" + 0.091*"job" + 0.073*"family" + 0.020*"class" + 0.018*"gift" + 0.014*"store" + 0.014*"interview" + 0.012*"computer" + 0.012*"bonus"
```

TOPICS VISUALIZATION with pyLDAvis

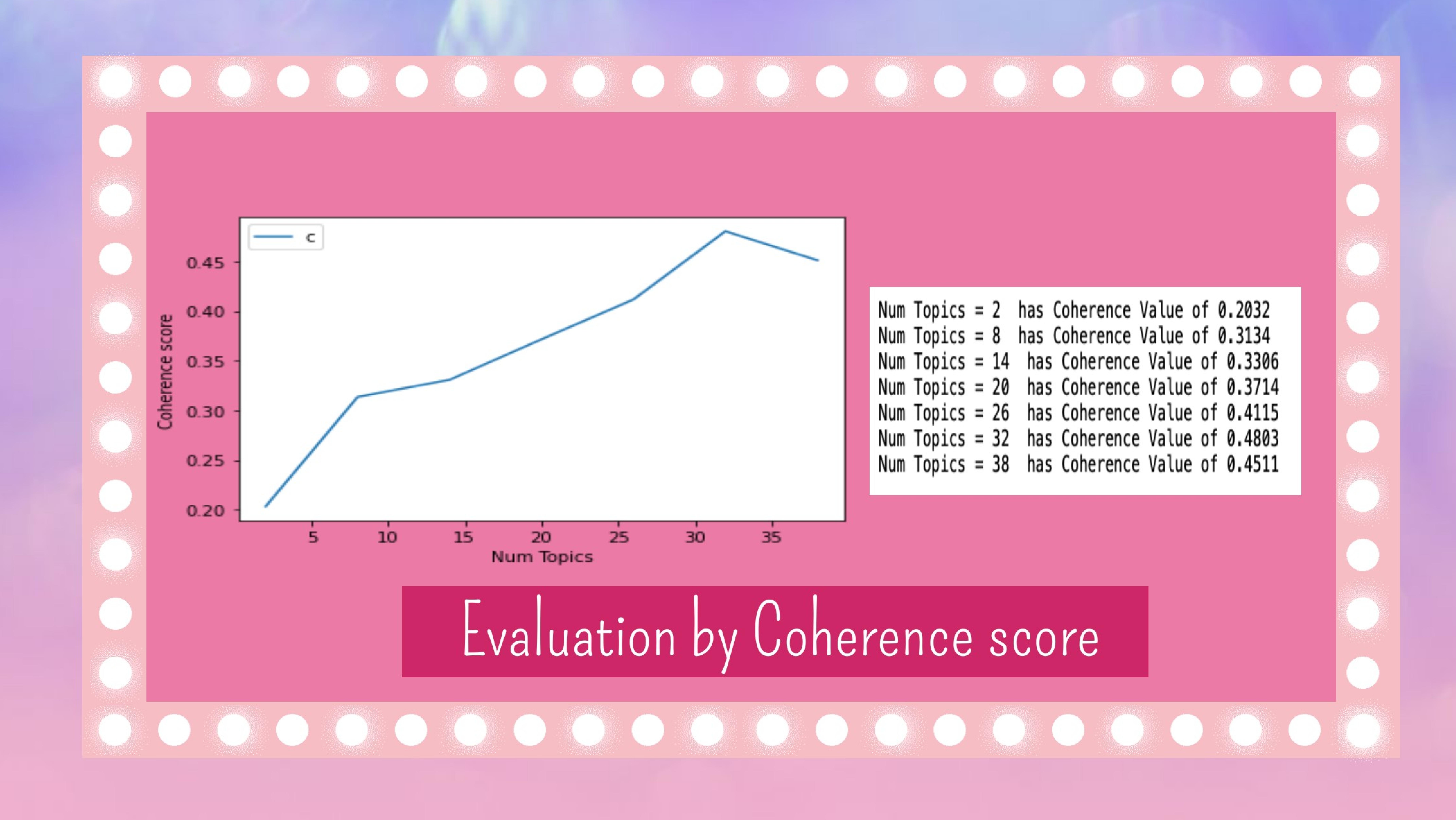


WordClouds based on LDA model









Topic classification based on LDA model + Domain Knowledge

Topics Classification:

Topic 0: spending
Topic 1: relaxing

Topic 2: entertainment

Topic 3: Social Activities

Topic 4: Achievement

The optimal topics number here is 32 topics with coherence score = 0.48

Conclusion

