29/11/2024

# Exploratory Telegram Data Analysis

Data collection & exploration of behavioural patterns

Project by Larysa Vasylenko
for Computational Social Science course;
Applied Mathematics, Bachelor's 3rd year

Teacher: Andrew Kurochkin

# Table of Contents

1. Introduction
2. Data collection
3. Exploration process
4. Results (shortly)
5. Further work
6. Link on the GitLab repo with the code

# Introduction

about the project

This data analysis looks at information taken from personal Telegram activity to find interesting **patterns in communication** and gain insights related to social behavior (for example Social Loafing, Groupthink, Synesthesia, homosociality). By studying messages, conversations, and interactions in private and group chats, the goal is to understand how people communicate, what they prefer, and explore any social trends that might appear.

# Data collection

To collect the data for this project, I used the Python library **Telethon** along with a pre-existing script available on GitHub. Initially, I faced challenges during the data extraction process. My first attempt yielded a dataset of **649 583 messages**, which was significantly limited due to the interruption of the code execution. (The code ran for almost a day and in the middle of the night its execution was interrupted due to network problems.)

At that stage, I decided to proceed with the data I had for the task titled **"Homework #3: Download and Understand a Dataset."**

However, it soon became evident that the dataset was insufficient for a comprehensive analysis. Determined to resolve the issue, I decided to solve the problem: I found the dialogues that were not loaded and uploaded them.

In the end, I was able to collect a much larger and more complete dataset, **comprising 1 262 176 messages**!

# Data collection

final data statistics

number of messages

1 262 176

dataset size in MB (.csv file)

376.9 MB

# Exploration process

My analysis began with reviewing the columns available in the dataframe and cleaning the data. During this process, I developed an understanding of the information and statistics I could visualize.

As a result, I had the following features at my disposal:

1. Date
2. Sender ID
3. Receiver ID
4. Message Content
5. Message Type
6. Duration
7. Reactions
8. Dialog ID
9. Forwarded From ID

10. Forwarded From Name
11. Dialog Type
12. Phone Number
13. Participant Information:
 Username
 User ID
 First Name
 Last Name

# Exploration process

**27 different aspects investigated!**

## Table of Contents

### Telegram activity analysis

### Phone numbers analysis

### Emoji usage

### Dataframe merging

### Event oriented

### Communication session duration analysis

### Sleep duration based on Telegram activity

### Message length analysis

### Telegram group chat analysis

### Telegram channel analysis

### Analysis of categorized chats by gender
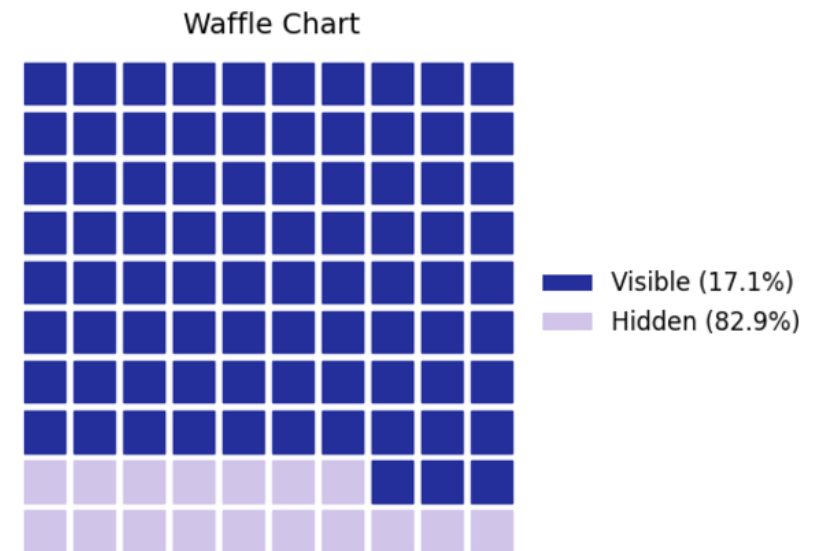
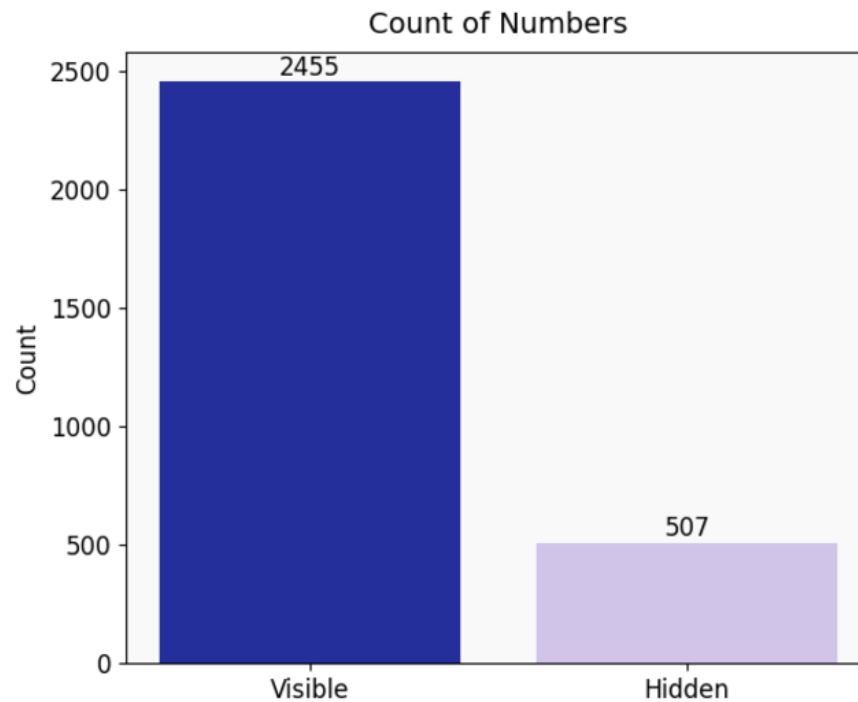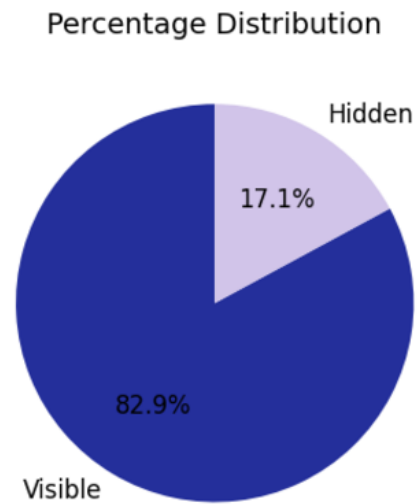### NLP (for Analyzing Communication Patterns)

### TF-IDF

# Exploration process & Results (shortly)

# Phone Numbers

# "Every sixth person makes their phone number visible on Telegram."



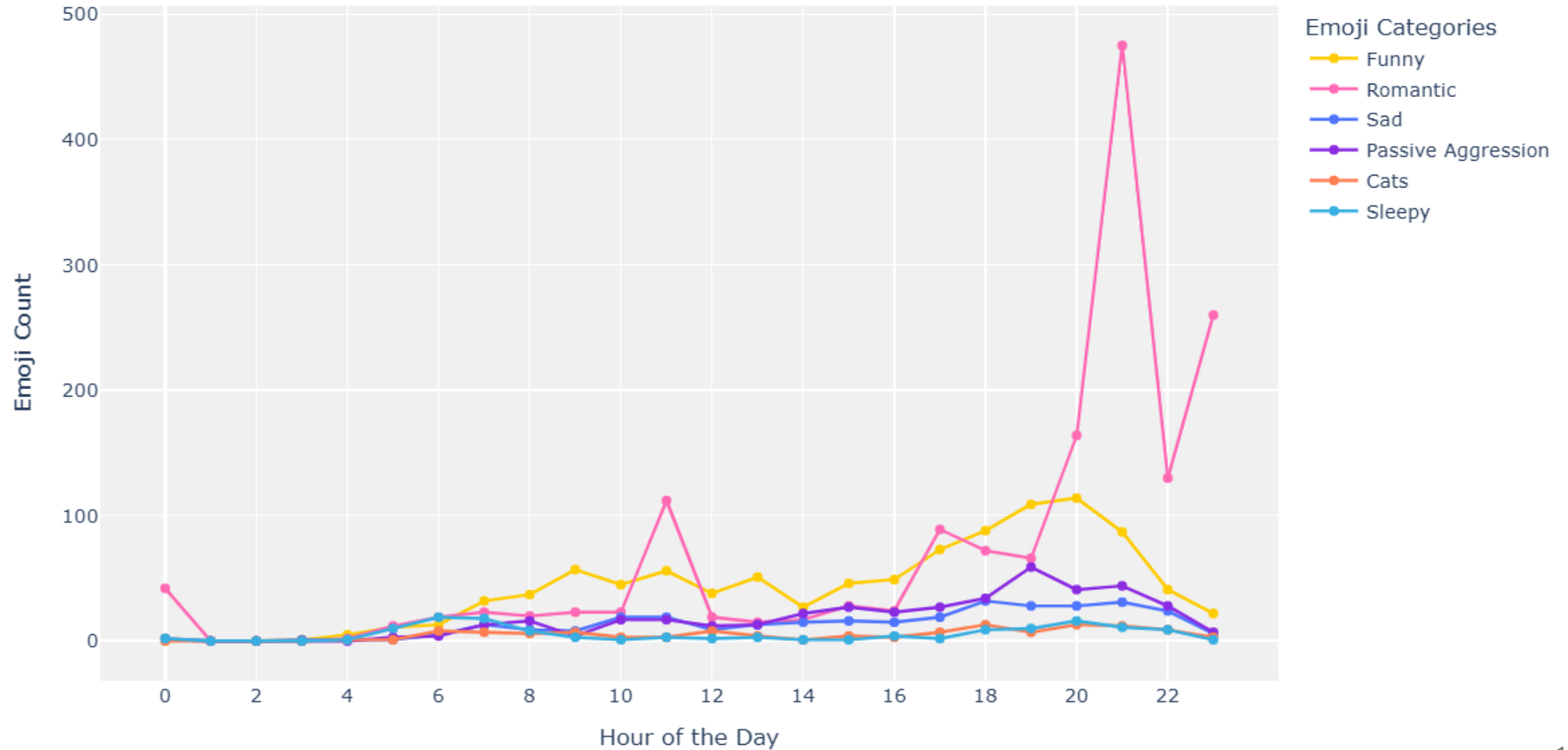Analysis of the tendency to share phone numbers on Telegram

Percentage Distribution

Count of Numbers

Waffle Chart

Visible (17.1%)
Hidden (82.9%)

🤞 Emoji

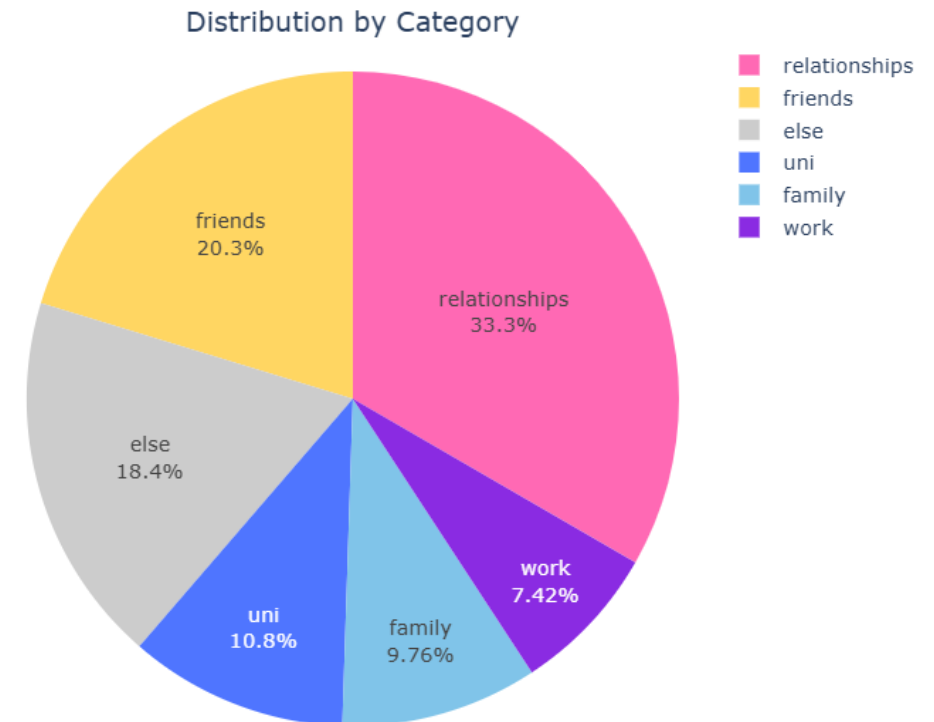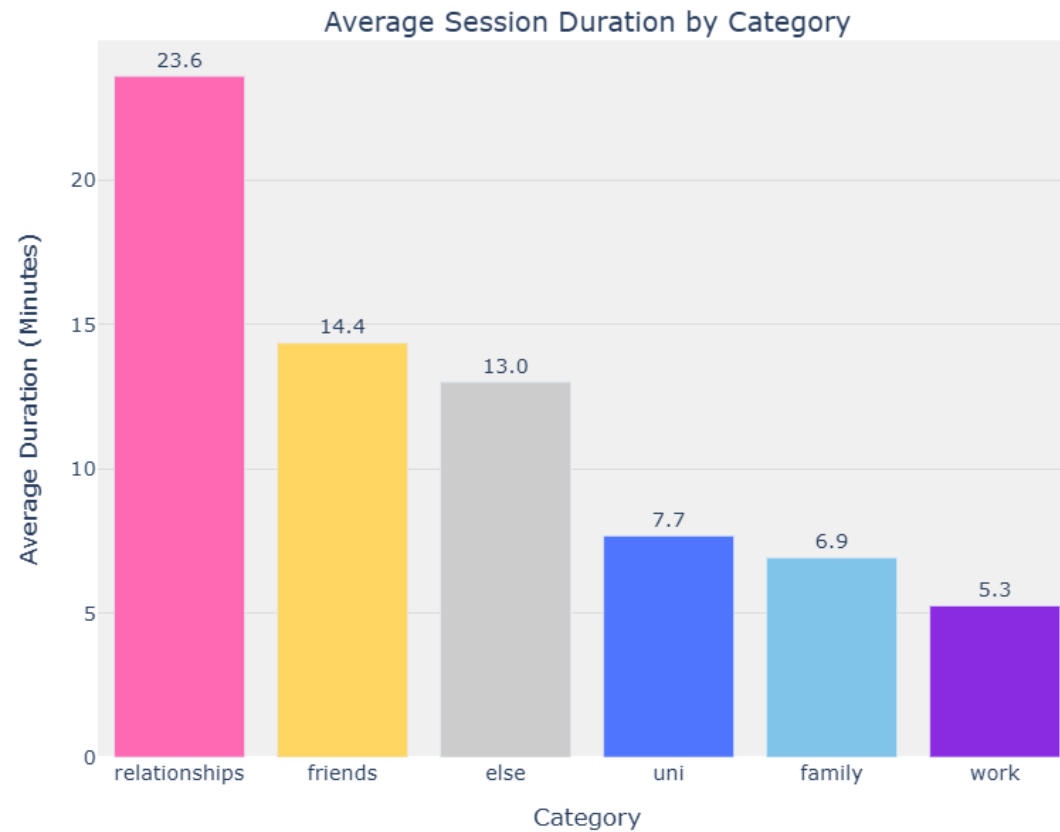# How does my mood change throughout the day?

# Session Duration

# With which category of people do I have the longest communication session?



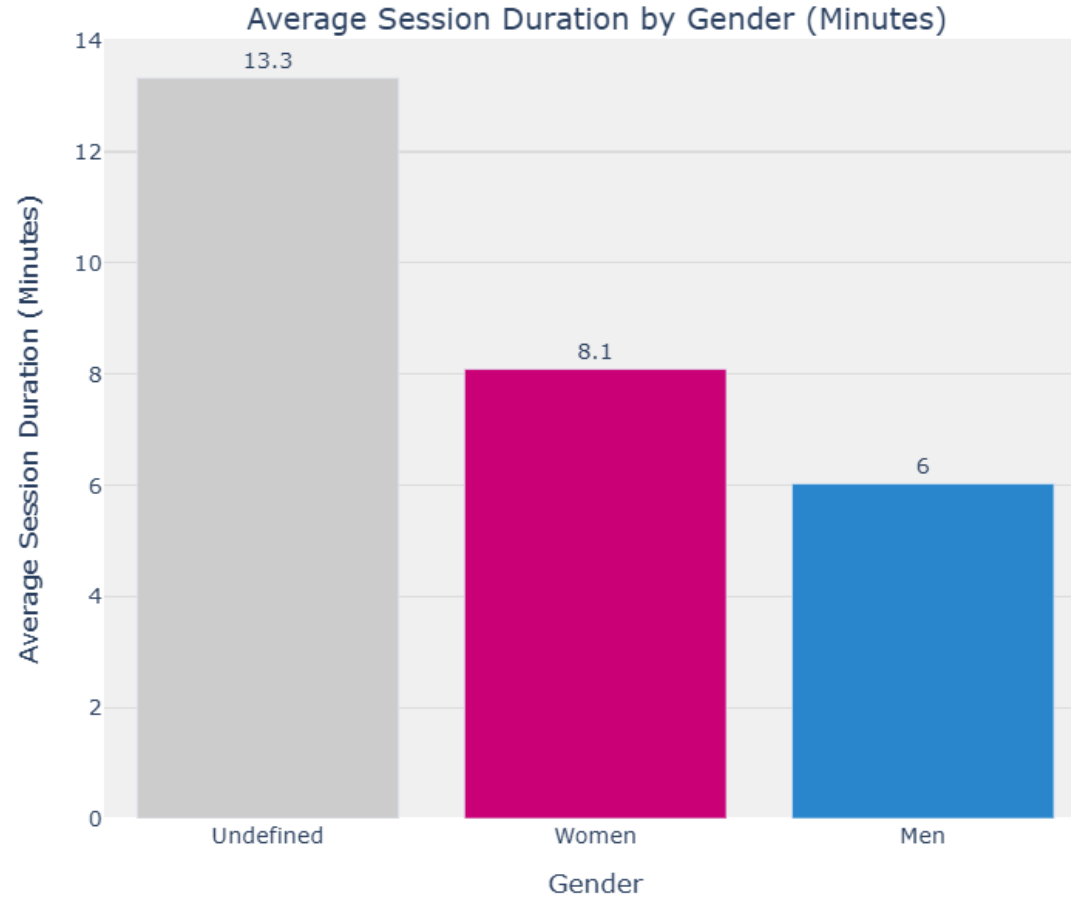Average Session Duration Analysis

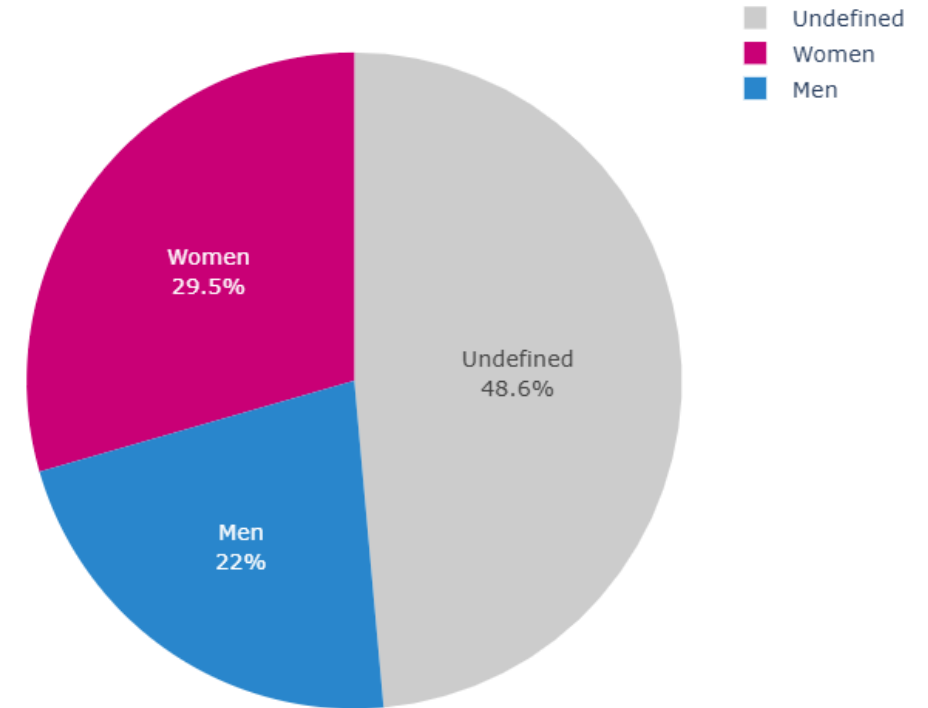How long do I usually sleep?

Average: 9.88 hours

# Which gender do I have the longest conversations with?

Average Session Duration by Gender



Average Session Duration by Gender (Minutes)
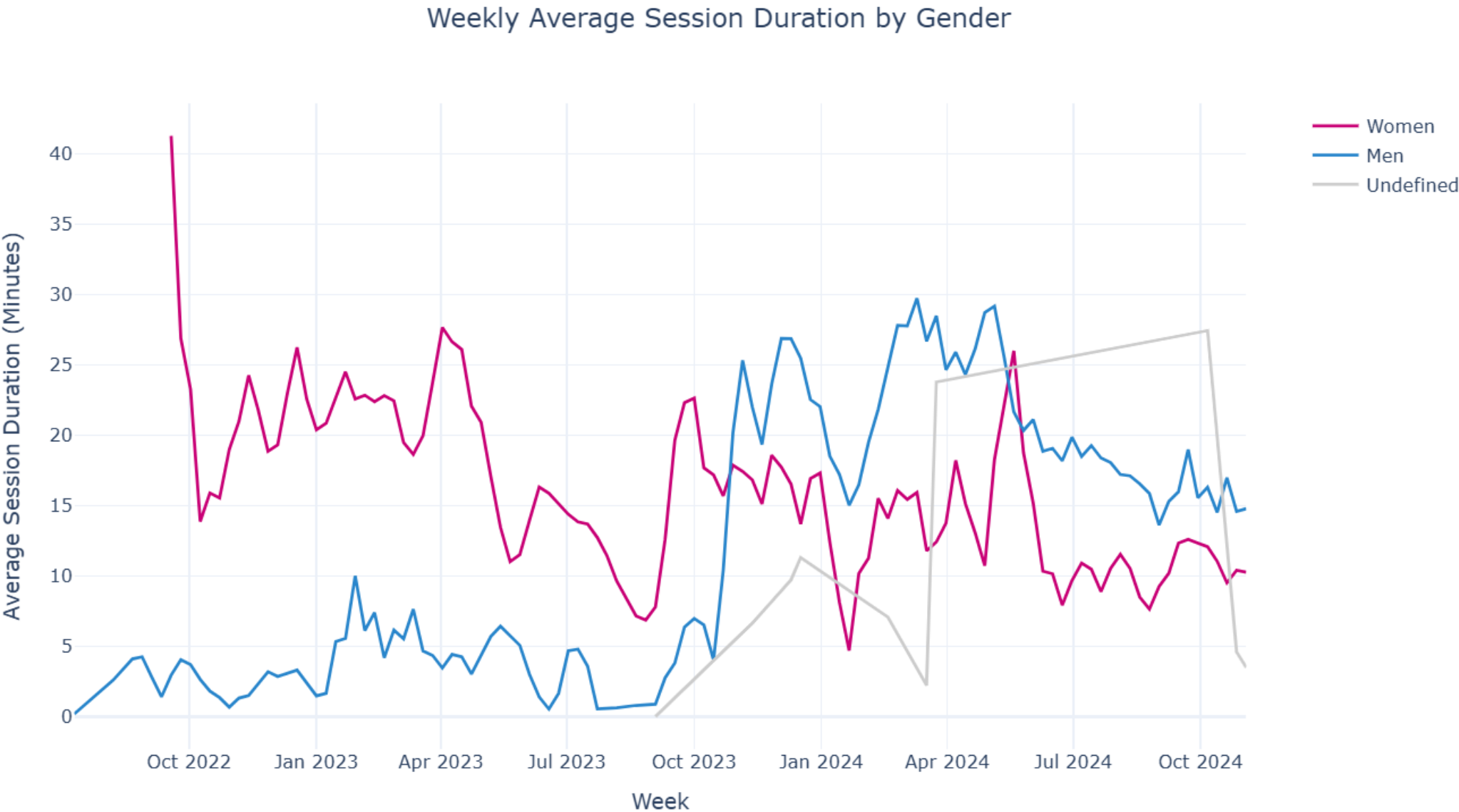
Gender Distribution (Percentage)

Legend: Undefined, Women, Men

Bar chart values: Undefined 13.3, Women 8.1, Men 6

Pie chart values: Undefined 48.6%, Women 29.5%, Men 22%

# How has this changed over time?



Weekly Average Session Duration by Gender
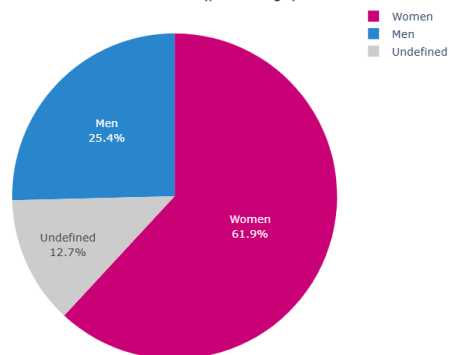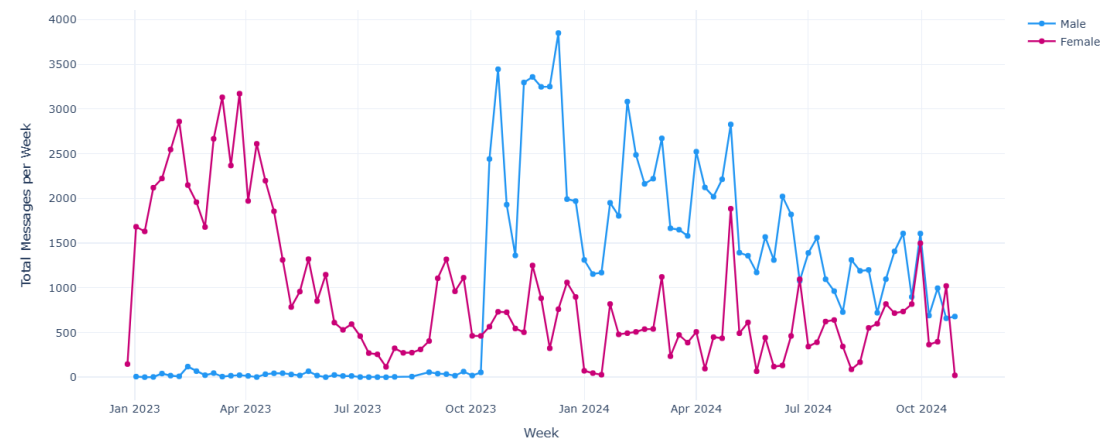
# Genders

Gender Distribution in Private Dialogs

Gender Distribution (counts)

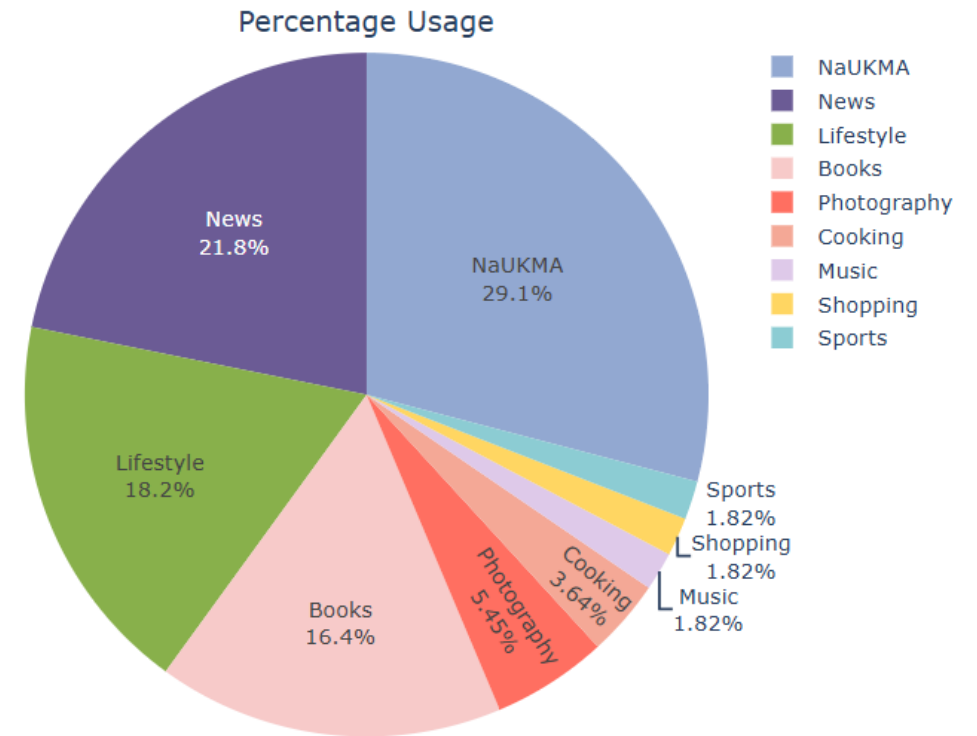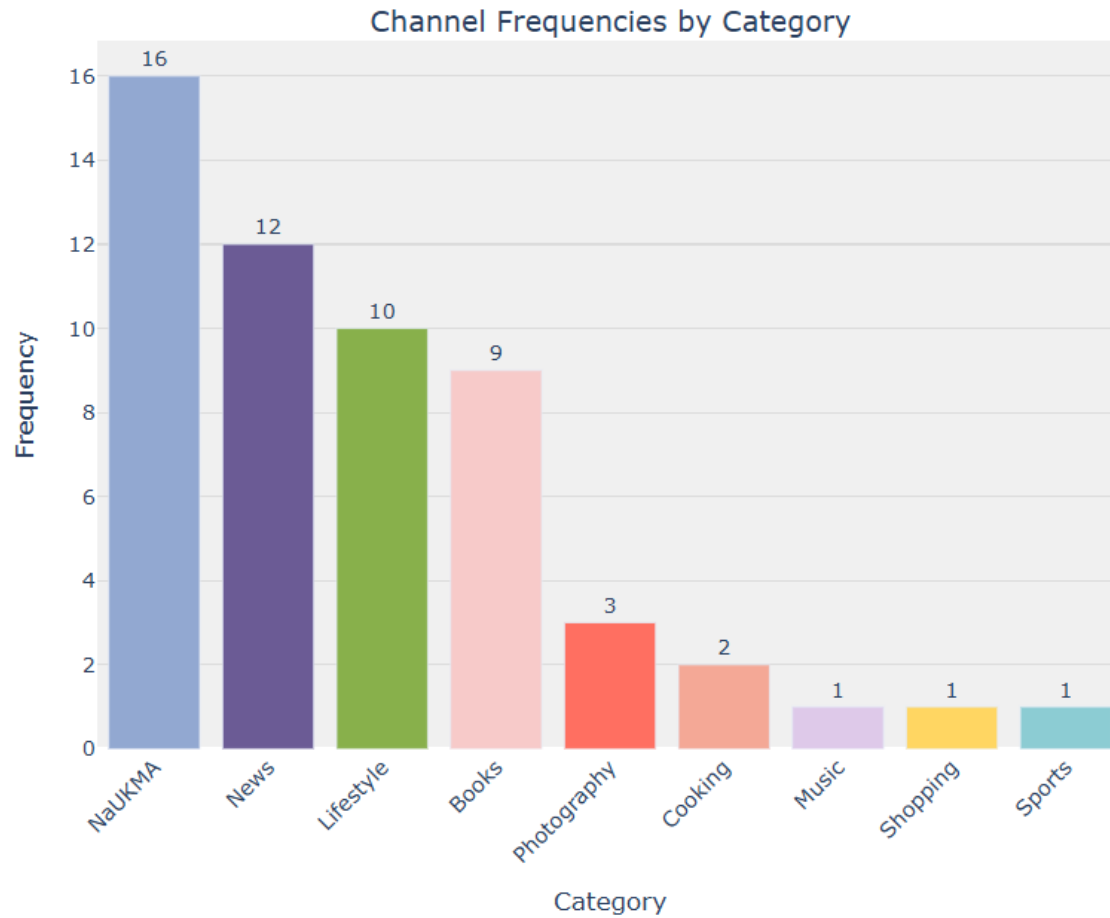Gender Distribution (percentage)

Weekly Total Messages in Private Dialogs by Gender

# Which Telegram channels do I mostly read?



Channel Analysis by Category

# Further work

1. Sentiment analysis

2. Deeper research into personal channels to increase activity there, for example, the channel of the NaUKMA student organization.

3. Using other available data, such as on reactions to messages, and building a model that would analyze the next message I want to send and predict the reaction of society. Or what message I will receive in response to mine!

4. Explore more with NLP

# GitHub
### repo with the code

To learn more about the insights I was able to find from my personal Telegram data, I invite you to the project repository!

Link: https://github.com/issaravas/Exploratory-Telegram-Data-Analysis/tree/main

Thanks for attention!
I will be happy to answer all questions about the project!