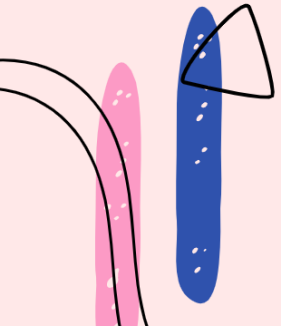
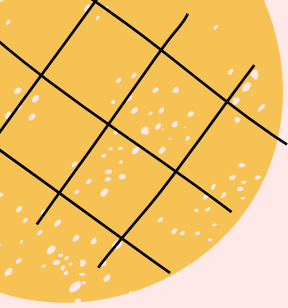




# User Profile Analysis

Team 11



# OUR TEAM



**Vince  
Pan**

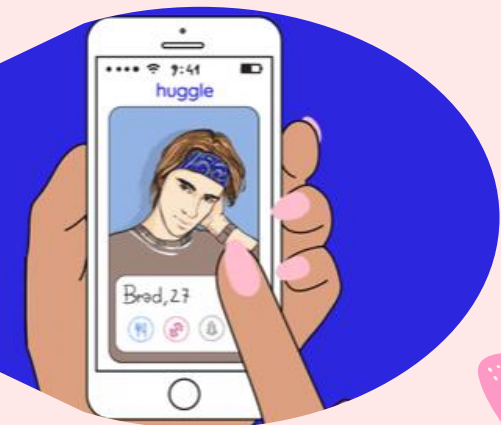


**Phoenix  
Wang**



**Peter  
Mankiewich**

# AGENDA



1

## Project Overview

What is our goal? Where did we get the data?

2

## Exploratory Data Analysis

What did we find out about the users?

3

## Text Analysis & Cluster Analysis

How did the users answer the open-ended questions?  
What types of users can you expect to see?

4

## Insights & Challenges

How can OkCupid benefit from our insights?



## PROJECT OBJECTIVE

We aimed to understand the types of users that are on the online dating app OkCupid. By creating distinct clusters of users, we can help people make better decisions while swiping, leading to better results for match seekers.



## **ABOUT THE DATASET**

The dataset was taken from Kaggle.  
Each row represents an anonymous user.  
It contains 59,946 rows and 31 columns,  
including structured personal information and  
written answers to 10 open-ended questions.



02

# EXPLORATORY DATA ANALYSIS

# DATA CLEANING & PREPROCESSING




## Missing Values

- Replace with mean
- Use logical answers
- Fill in "unknown"

## Categorical Features

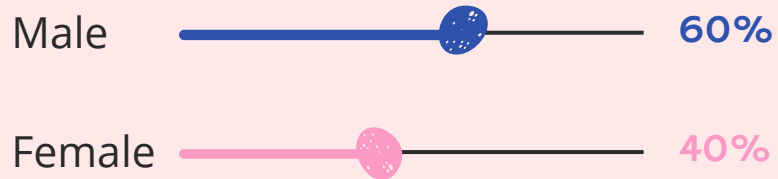
- Group values
- Convert to ordinal
- Create dummy variables

## Textual Data

- Replace meaningless content (e.g., links)
  - Concatenate strings
- 

# USERS OVERVIEW

## GENDER



## STRICTLY FOLLOWING DIET



Halal  
(23.4%)

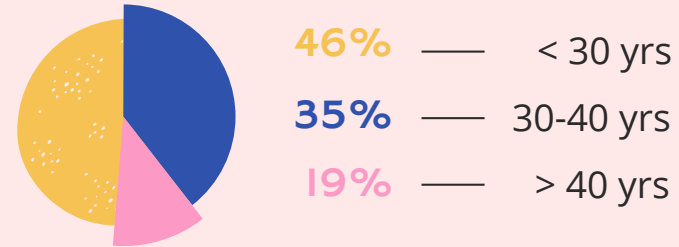


Veggie  
(19.4%)



Kosher  
(15.7%)

## AGE



## TAKING RELIGION SERIOUSLY







# INTERESTING FACTS

## Relationship Status

3.4% are in a relationship  
0.5% are married

## Sexual Orientation

9.3% are homosexual  
4.6% are bisexual

1

3

2

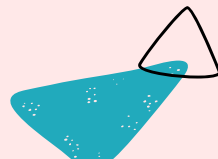
4

## Education Level

52.5% - college/university  
3.5% - law or med school

## Language Skills

Majority of Asian, Hispanic, and  
Middle Eastern users are bilingual



The background is a solid light pink color. It is decorated with several abstract shapes: a large pink circle with white speckles in the upper center containing the number '03'; a smaller yellow circle with white speckles overlapping the pink circle; a blue teardrop shape with white speckles on the left; a blue circle with white speckles at the top right; a pink vertical pill shape with white speckles on the right; a yellow pill shape with white speckles at the bottom right; and a small black-outlined triangle at the bottom left.

**03**

# **TEXT ANALYSIS**

# ESSAY QUESTIONS



*About me...*

*The six things I could never do without...*

*The most private thing I am willing to admit...*



**What do these responses tell us about a user,  
and their level of commitment?**



**How can we use this information to make intelligent matches?**



# NAMED-ENTITY RECOGNITION

hello! i enjoy traveling, watching movies, and hanging out with friends. my # 1 **MONEY** rule for traveling is to go to places where where i can drink the tap water. =p i've been to japan **GPE** , washington dc **GPE** , new york, london **GPE** , paris **GPE** , china **GPE** (an exception to rule # 1 **MONEY** ) and many other places. for the year **DATE** , i haven't decided yet. maybe toronto and montreal. my taste in movies (and television) fall into three **CARDINAL** genres: comedies, action, and drama. it doesn't matter if the film is in english or sub-titled (i can't stand dubbed), but it has to be well-written. aside from the traveling and movies, i also enjoy broadway shows ( 4 **CARDINAL** shows in 5 nights **DATE** in vegas a couple years ago **DATE** ), bay area sports (except the sharks), online gaming, good food, and keeping up with current events. my friends would describe me as easy-going, goofy, thoughtful, consistent, honest, and dependable. i can get serious when i have to be. overall, i like to think that i lead my life by doing what is "right".\_i'm a civil engineer, who enjoys helping the citizens of san san francisco\_- looking at things objectively - getting things done (but others may disagree) =p - remaining calm in any situation\_i'm **ORG** quiet until i get used to the environment (but isn't that normal?).\_last book: "game change". movies: bourne series, action, smart comedies (in contrast to slapstick), must have a decent plot shows: local sports, ncis, good wife, top chef, currently exploring japanese **NORP** dramas music: koit... food: korean bbq, japanese **NORP** , chinese **NORP** , thai **NORP** , italian, burgers, almost anything\_- iphone - friends and family - internet - bay area sports - humor - point & shoot camera\_aside from work, how to improve my home.\_out enjoying friendly conversation over dinner.\_please let me think about this more.\_we have similar interests.



# HOBBY ANALYSIS

- What are users interested in, and are they serious about the app?
- Benefits/problems identified
  - Marks important and relevant interests
  - The algorithm could be improved for more accurate identification

'acting'	'travel'	'traveling'	'traveling'
'acting'	'arts'	'drama'	'games'
'gaming'	'internet'	'shooting'	'sports'
'traveling'			





# SENTIMENT & SUBJECTIVITY

## Step 1

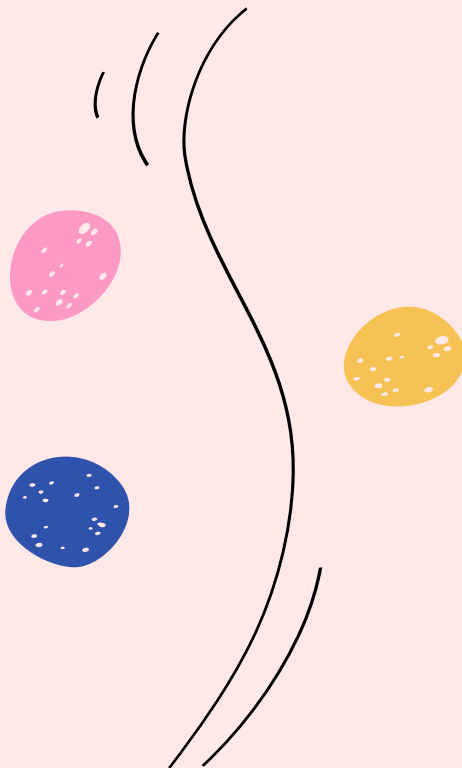
Apply TextBlob to each column to gather sentiment & subjectivity scores

## Step 2

Calculate the average sentiment & subjectivity scores

## Step 3

Create two columns in the dataset to save each user's scores



The background is a solid light pink color. It is decorated with several abstract shapes: a large pink circle with white speckles in the upper center, a smaller yellow circle with white speckles to its left, a blue teardrop shape with white speckles on the left, a blue circle with white speckles at the top right, a pink vertical pill shape with white speckles on the right, and a yellow pill shape with white speckles at the bottom right. There are also thin black outlines of a triangle at the bottom left and a line on the right side.

04

# CLUSTER ANALYSIS



# PCA/t-SNE & K-Means

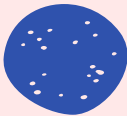
## Step 1

Standardize numeric features  
in the dataset

## Step 2

Implement **PCA** & get components  
for 95% explained variance ratio

Implement **t-SNE** & get back 2  
output columns/variables

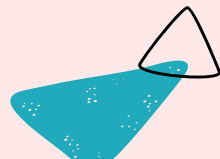


## Step 3

Record the inertia and silhouette  
scores for different K's

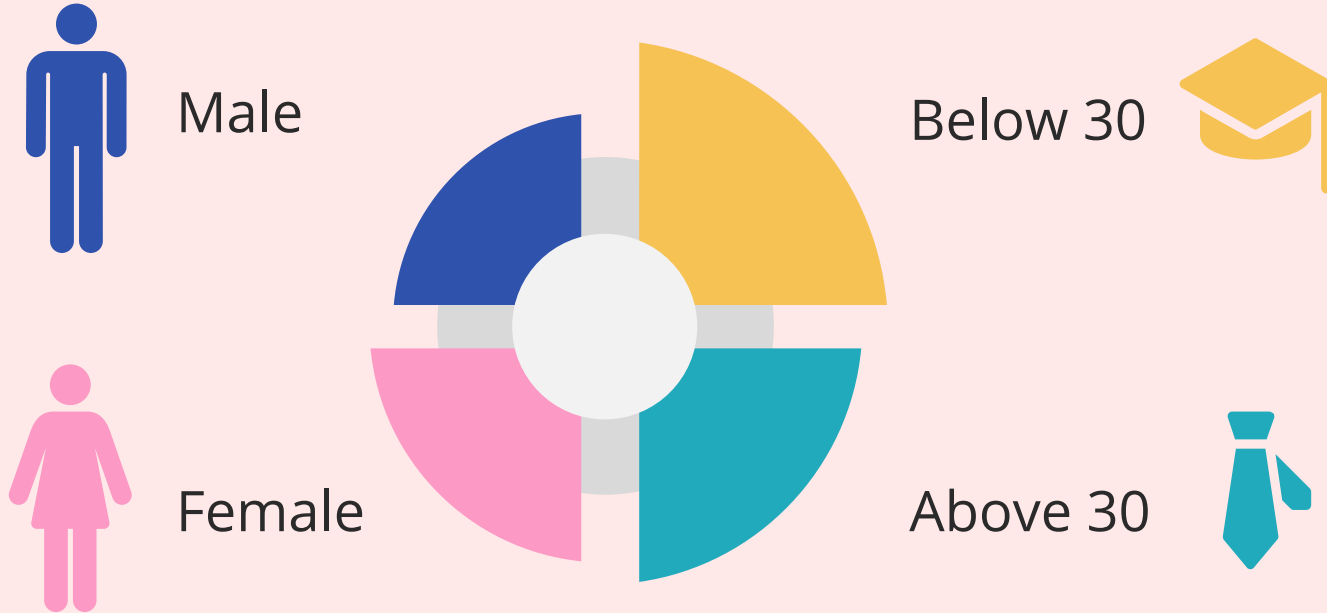
## Step 4

Select the optimal K value for  
K-Means & examine each  
cluster's descriptive statistics

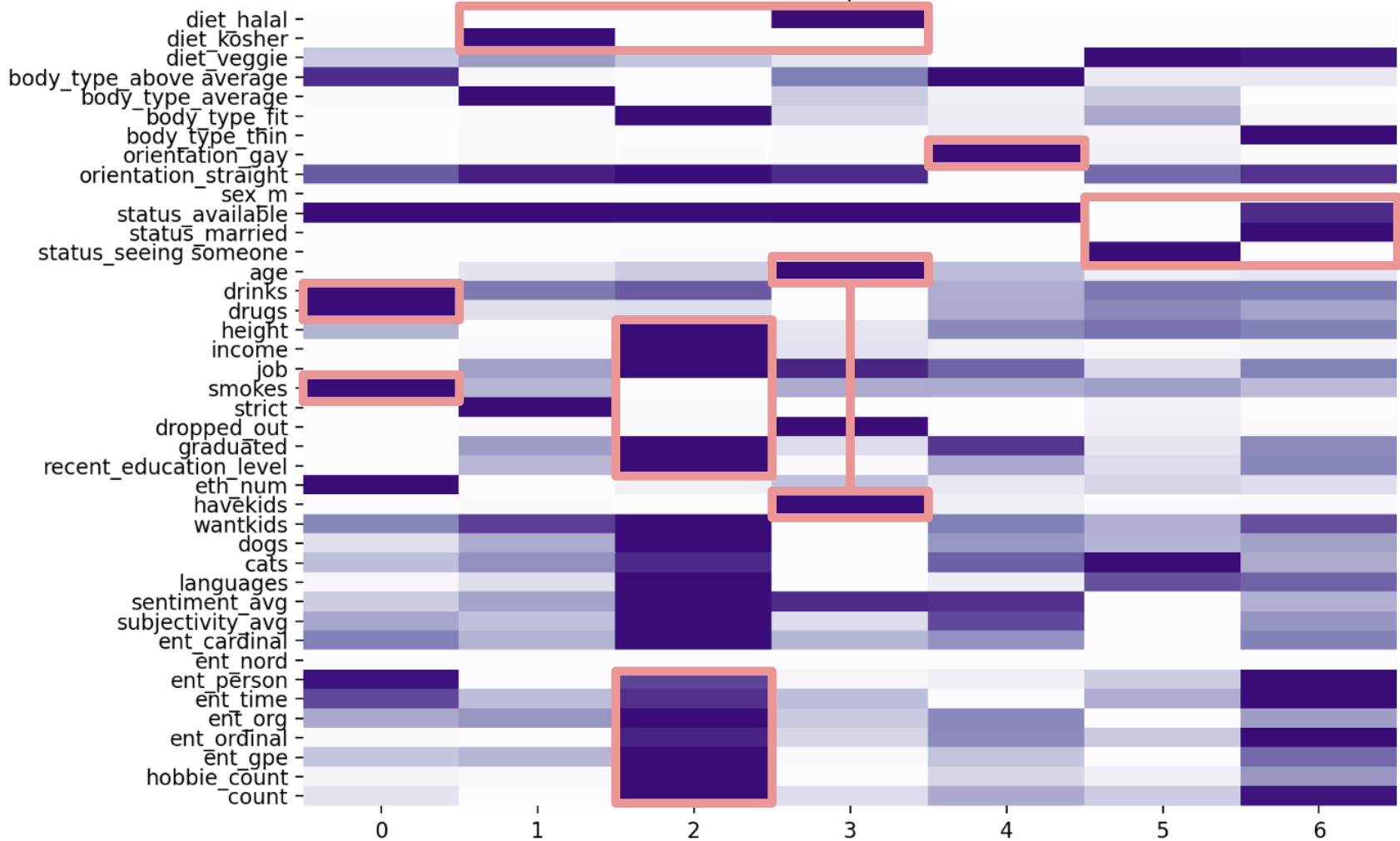




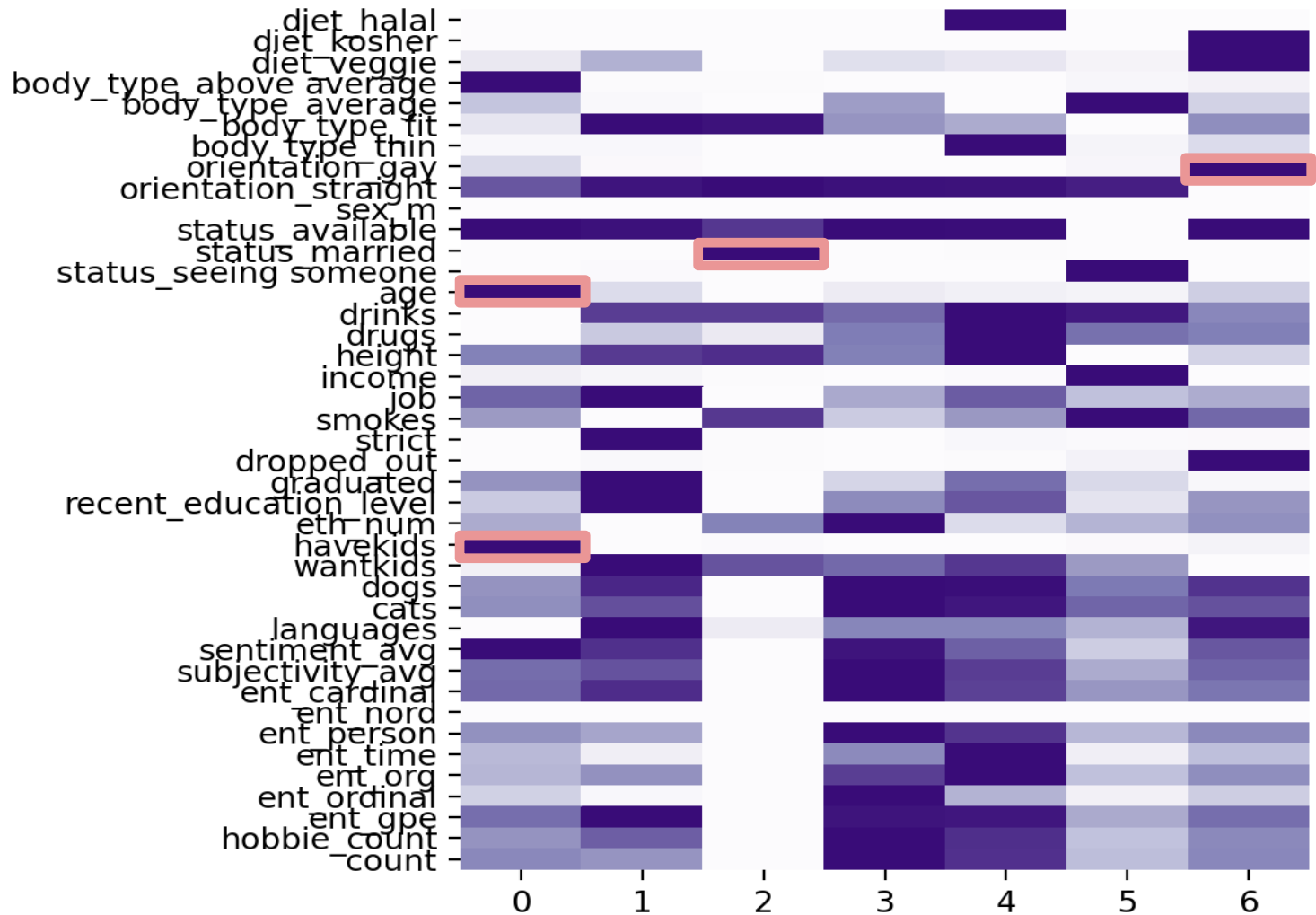
# K-MEANS CLUSTERING



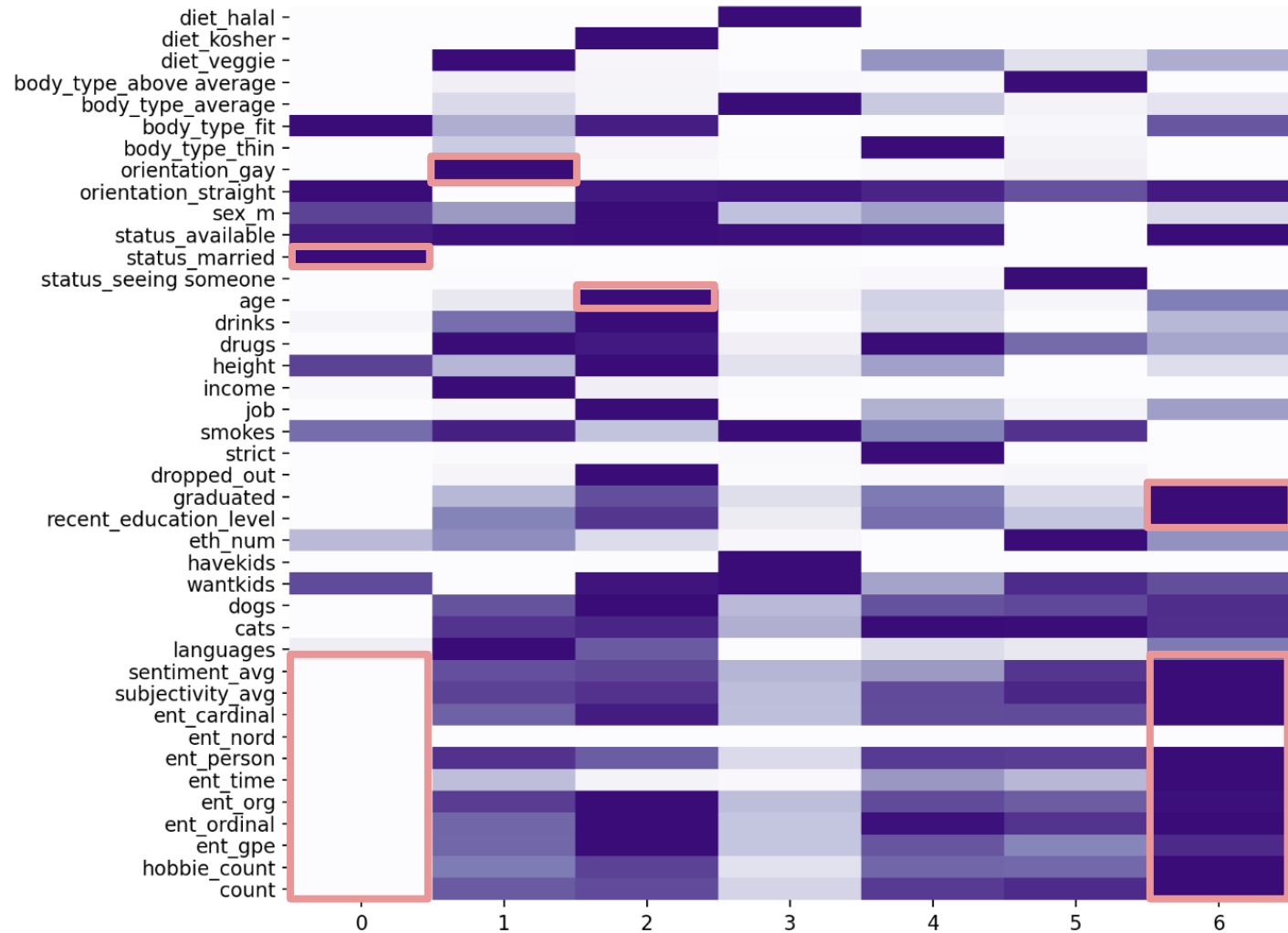
## Female Descriptive Statistics



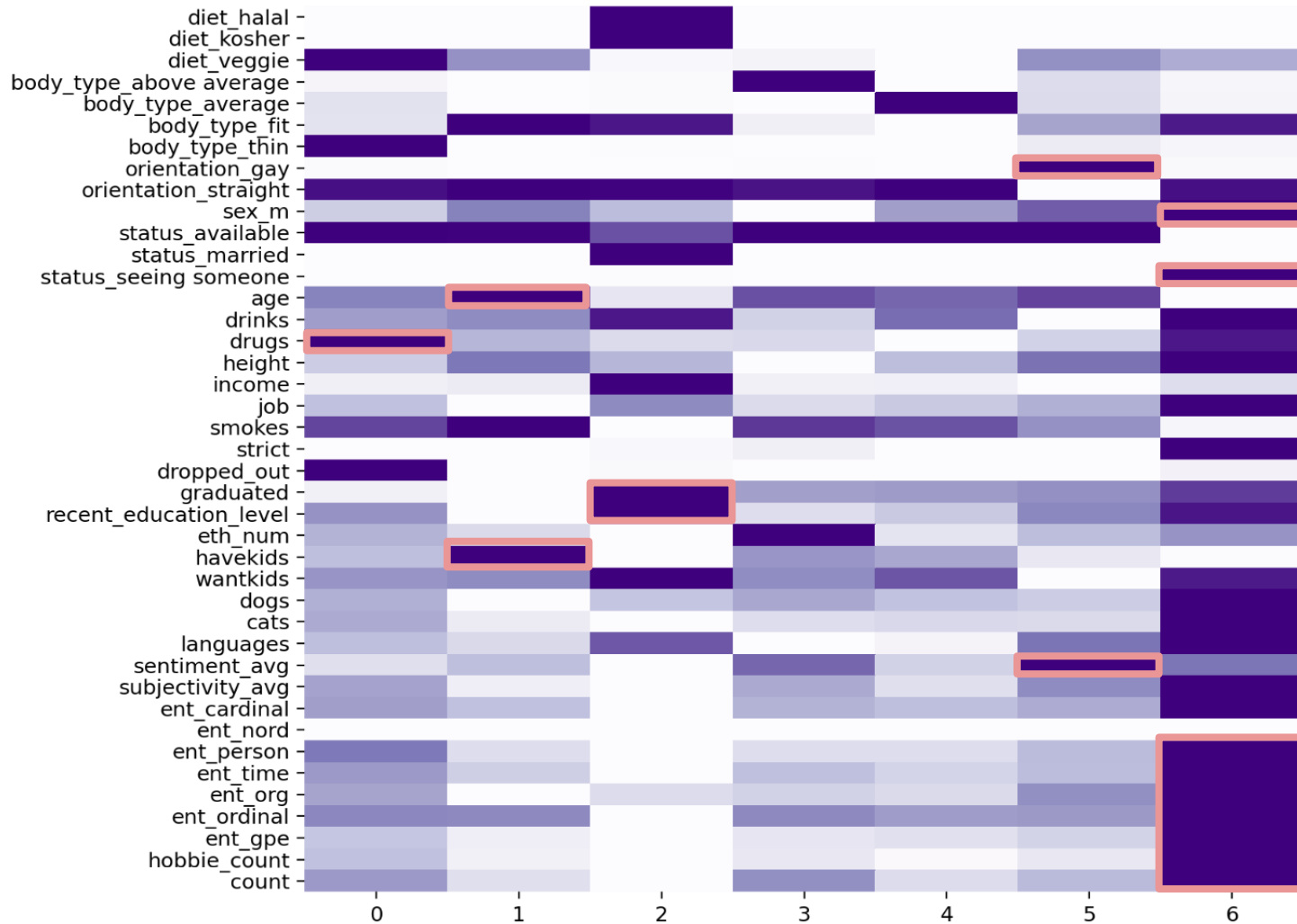
# Male Descriptive Statistics



AgeBelow30 Descriptive Statistics



AgeAbove30 Descriptive Statistics





05

# INSIGHTS & CHALLENGES



# BUSINESS VALUE FOR



## Internal User Rating

- Detect scam accounts
- Identify serious match seekers based on written answers

## Recommendation System

- Recommend profiles to users based on their swiping behaviors
- Recommend activities based on the users' characteristics

## Filter Optimization

- Combine structured information and essays to create intelligent filters
- Generate prompts to inspire profile updates



# CHALLENGES

## Missing Values

Users chose to not fill out certain fields when creating their profiles.

1

## Computational Cost

The computational cost of creating a distance matrix is large.

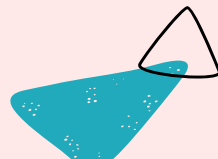
3

2

## Mixed Data Types

Future work:

- Explore alternative techniques, e.g., multiple correspondence analysis, factor analysis, etc.





**THANK YOU  
&  
HAPPY SWIPING**

