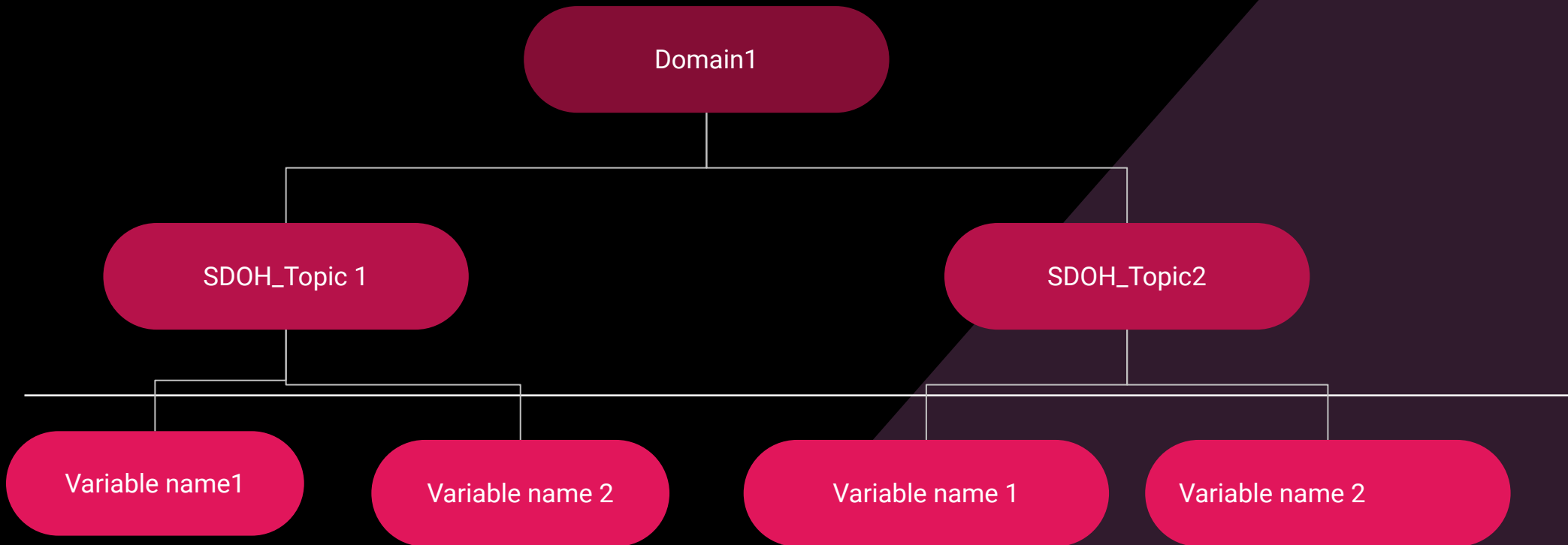# EDA AND PATTERNS OF MISSINGNESS IN SDOH for the EMORY CXR

Team 5 - CXR Dataset

Ali Aslam
Anudeep Errabelly
Enamul Hoq
Zeph Kaffey
David Nyarko
Chiratidzo Sanyika
Veera Venkata Satyavathi
Enzo Ferrante

# SOCIAL DETERMINANTS OF HEALTH (SDOH)

- Social factors allocated according to zip code
- 44 sources
- 8 categories (domains)

# SO MANY IDEAS

Merge SDOH, metadata, findings

**+**

Emeddings

# But....



- We discovered there are a lot of missing data in the SDOH dataset

# We decided to redirect our plan and focused on

Characterizing the fingerprint of "missingness" in the SDOH Dataset

**+**

Creating embeddings to allow for easier analysis

# Why focus on EDA?

- You cannot start running analysis on a dataset that is not well understood
- Patterns of missingness could lead to underrepresenting some groups and therefore sample bias
- Embeddings allow for easier analysis such as model training, and clustering identifies potential shortcuts for classification models
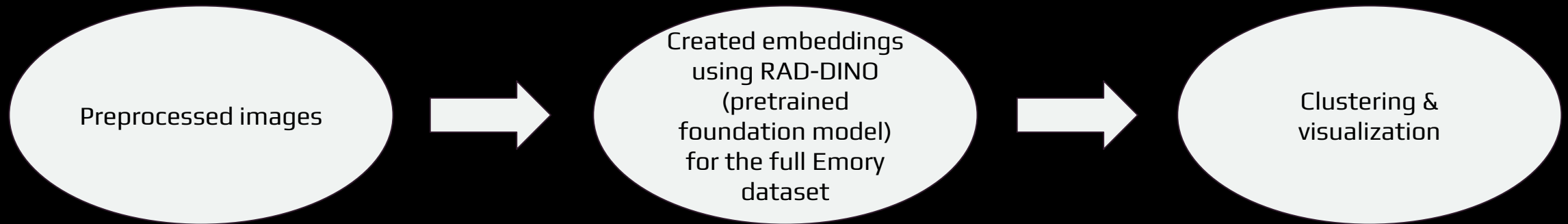
# Hypothesis

- Observe patterns of missingness in the SDOH table for different demographic group (sex and race)
- SDOH (like Gini index) vary strongly for different diseases, races, and sexes
- The embeddings will clusterize by demographics (sex and race) and SDOH
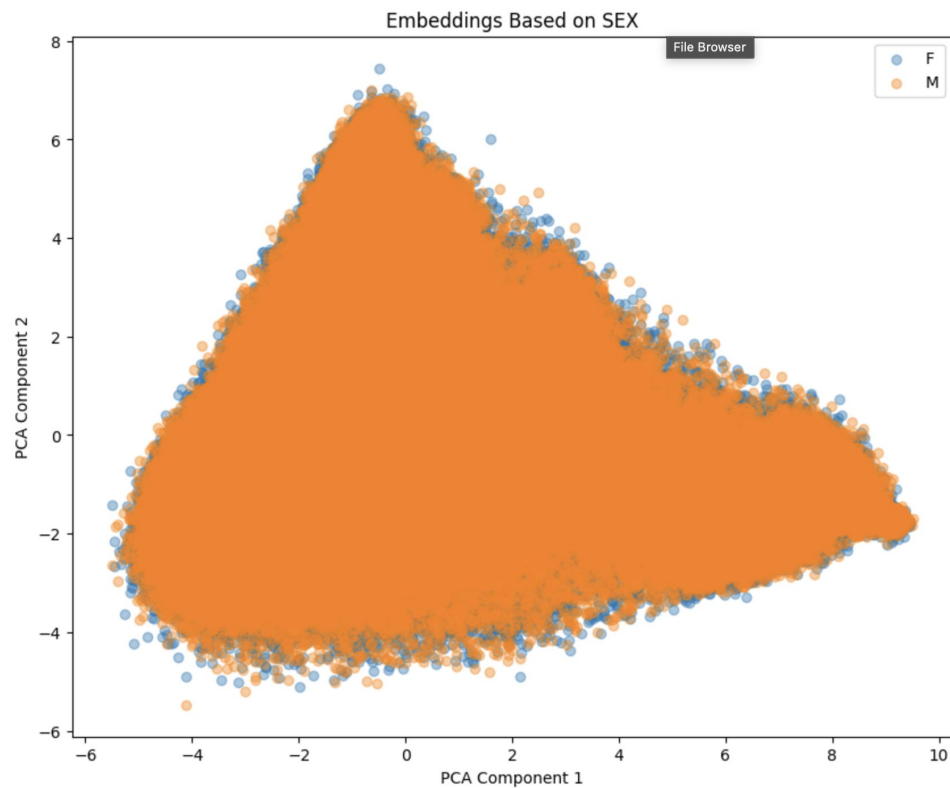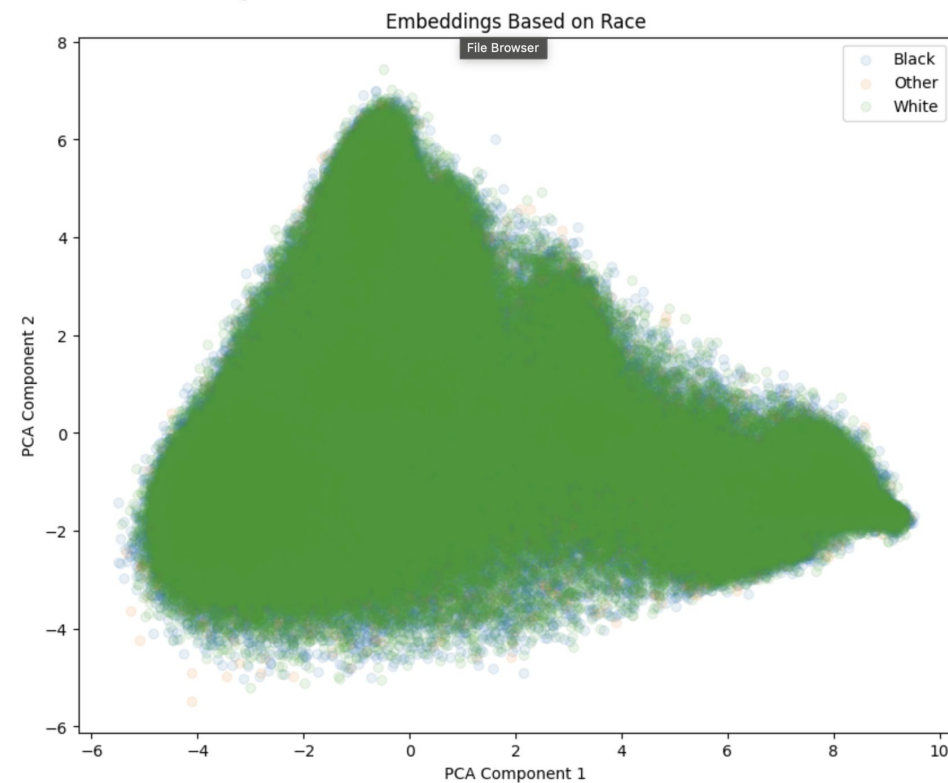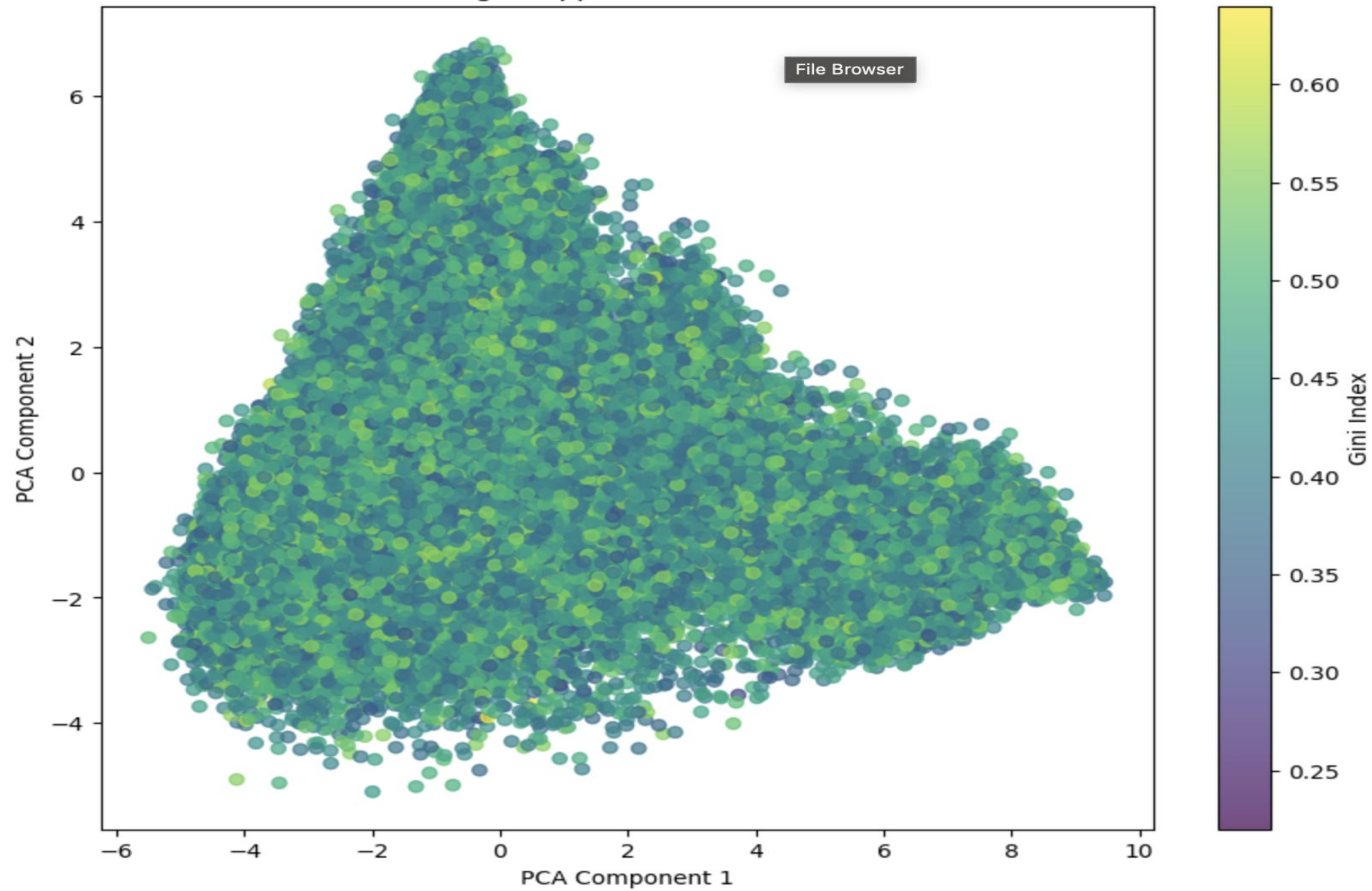
# Methods - SDOH

Merged SDOH, metadata and findings and split by demographic

→

% missingness by demographic by variable name

→

Plotted missingness by domain by demographic (sex, race)

# Methods - Embeddings

Preprocessed images → Created embeddings using RAD-DINO (pretrained foundation model) for the full Emory dataset → Clustering & visualization

# RESULTS: Embeddings

Embeddings Mapped Based on Gini Index

Percentage of Missingness in Each Domain by Race
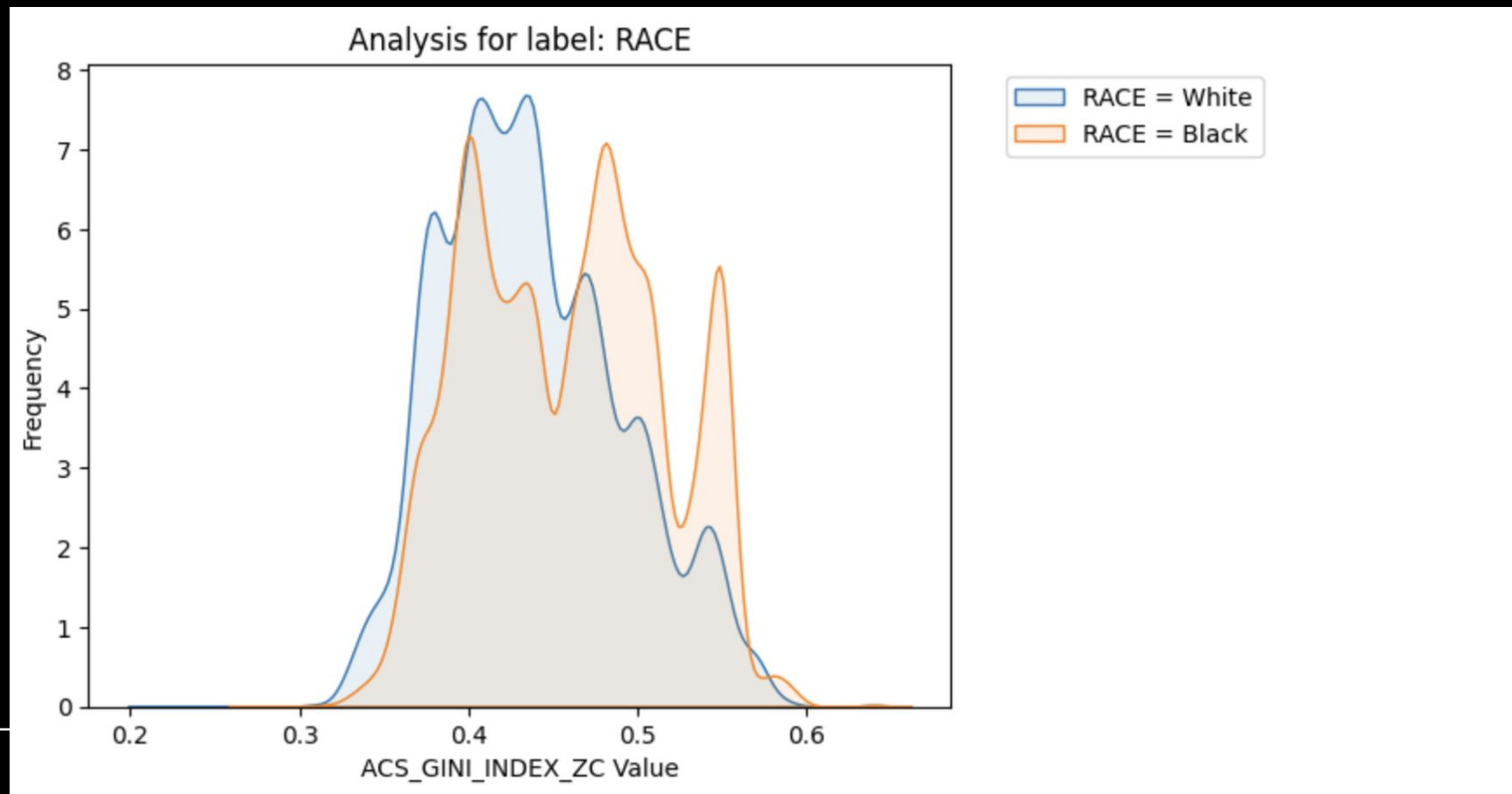
Percentage of Missingness in Each Domain by Sex

# RESULTS: Example of subpopulation analysis for the SDOH per Race

# Conclusion/Relevance

- Lay the framework for preparing & clustering embeddings, and interactive graphs to see who is and is not represented in the SDOH dataset
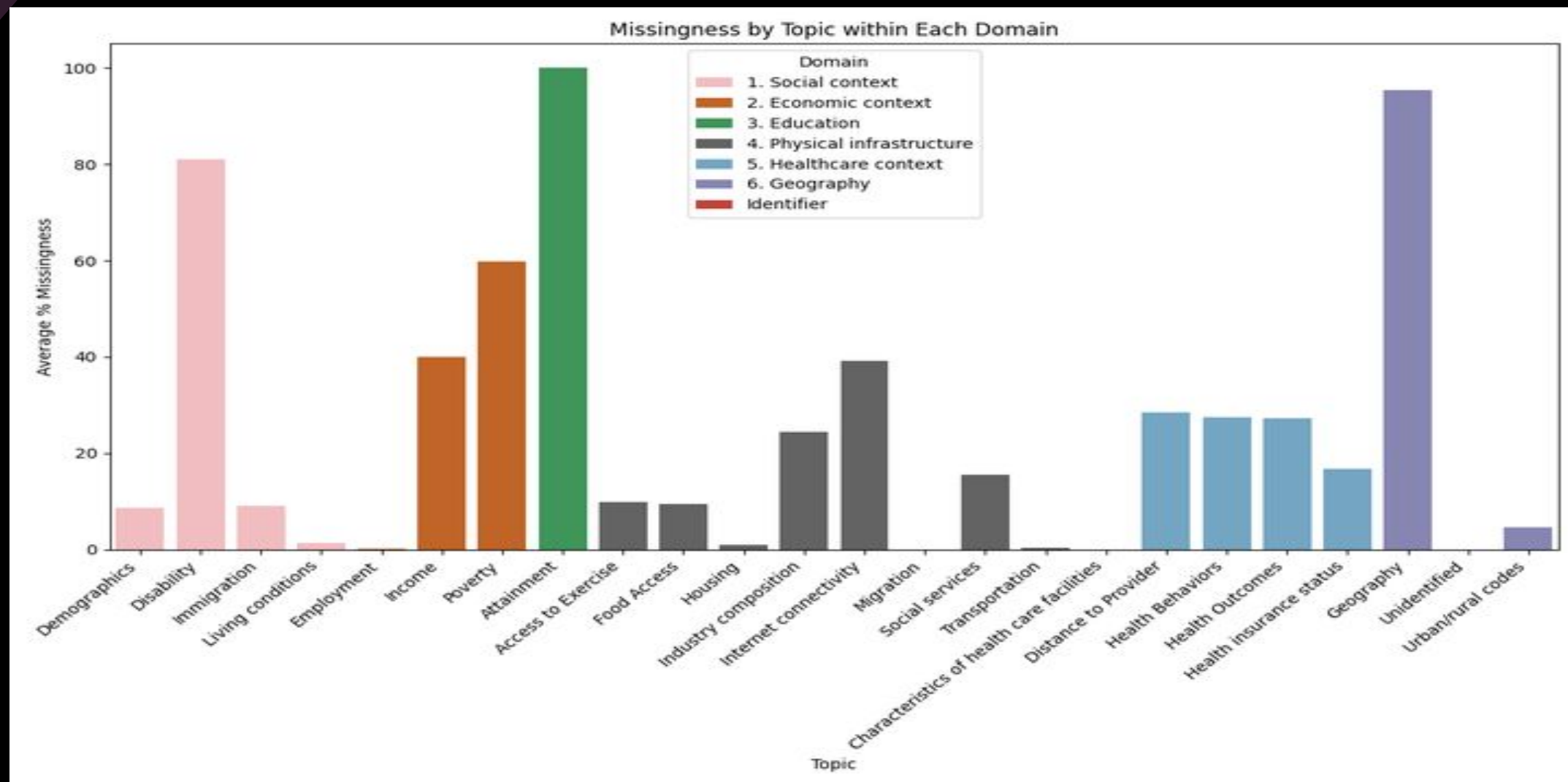
MEET THE TEAM

# Extra Slides

# RESULTS



Missingness by Topic within Each Domain

# RESULTS



Overall Missingness by Domain