# Suicide Prevention
DDS Projects Showcase
May 11, 2023

**\*Trigger Warning\***
topics include suicide, mental health issues, drug use, and crime

Team Lead: Zairan Xiang
Team Members: Courtney Cheung, Vicky Li, Lana Murray, Josh Puray

# Presentation Overview

- Background/Significance
- Dataset & Our Question
- Pre-Processing & EDA
- Models and Analysis
- Limitations/Challenges
- Conclusion

# Background/Significance

# Background/Significance

## Stats (World Health Organization - 2021)

- More than 700,000 people die due to suicide every year.

- Suicide is the fourth leading cause of death among 15-29 year-olds.

- 77% of global suicides occur in low and middle income countries.

## Significance

- By analyzing relevant features and attempting to predict suicide risk, we aim to spread awareness and emphasize that suicide is preventable.

- Suicide risk assessment is used by schools, healthcare services, and mental health professionals

- Important to understand what contributes to suicide risk

- Wanted to explore the relationship between drug use and mental health issues

- The better we understand risk, the faster we can provide quality mental health and harm reduction services

# Dataset & Research Question

# Our Dataset

- The National Survey on Drug Use and Health (NSDUH) Series

- Years 2010-2014

- Over 3000 features and 50,000 observations for each year

- Detailed documentation codebook with 900+ pages

- Needed to reduce number of features

# Main Feature Categories (52 in totals)

- Mental Health

- Demographics

- Crime

- Drug Use

- Religiosity

# Research Question

Which features are most predictive of suicide risk?

# Pre-Processing & EDA

# Roadmap for Pre-Processing

Before merging the 2010-2014 dataframes, we had to:

1.  Only use features shared by each dataframe

2.  EDA by topic category

3.  Delete features with too much missingness and impute missingness

4.  Find the most meaningful features

5.  Rename columns for clarity and replace values in dataset

6.  normalize all independent variables

# Generating Suicide Risk Feature (Predictive Feature)

Suicide Risk is calculated using the following variables

- Think:     Have you ever thought about

    suicide in the past year

- Plan:     Have you ever planned suicide

- Attempt:  Have you ever attempted suicide
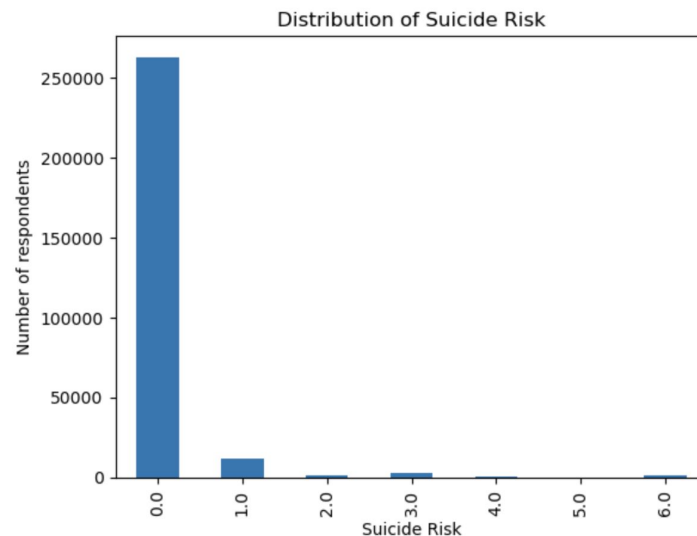
Suicide Risk = Think + 2*Plan + 3*Attempt



**Figure 1: Distribution of Suicide Risk**

# EDA

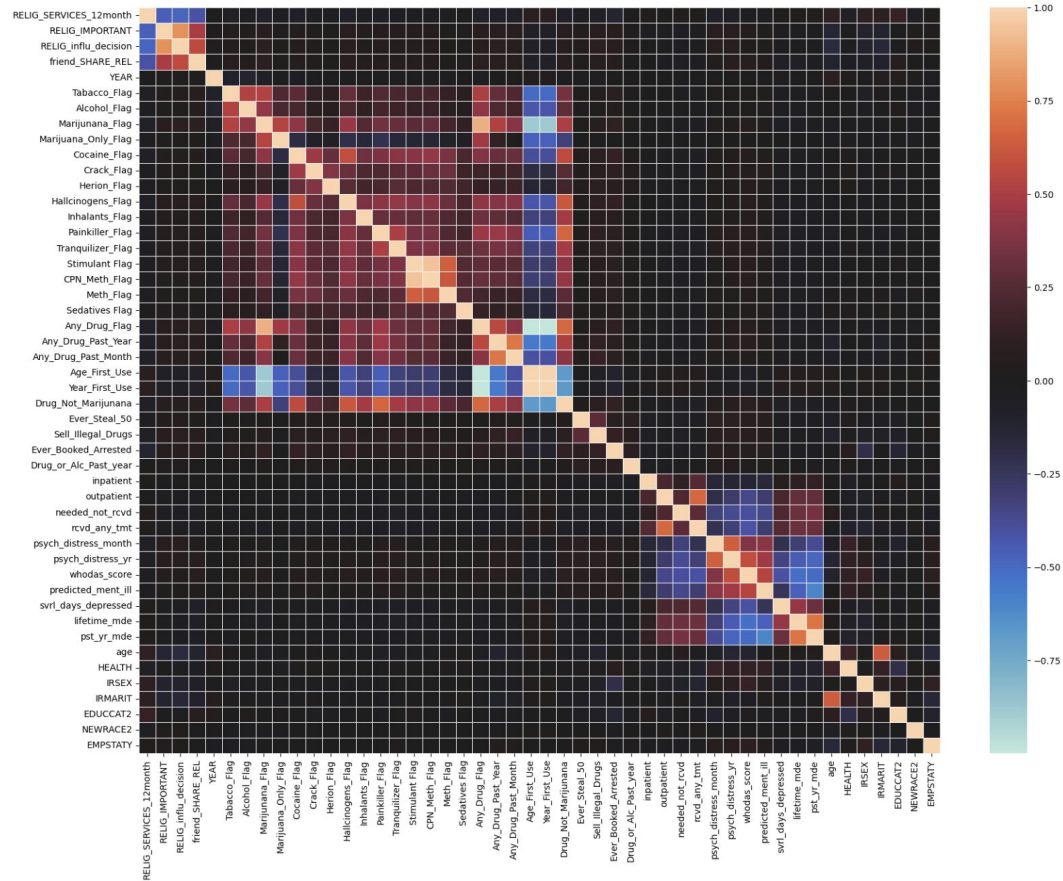# Checking for Multicollinearity



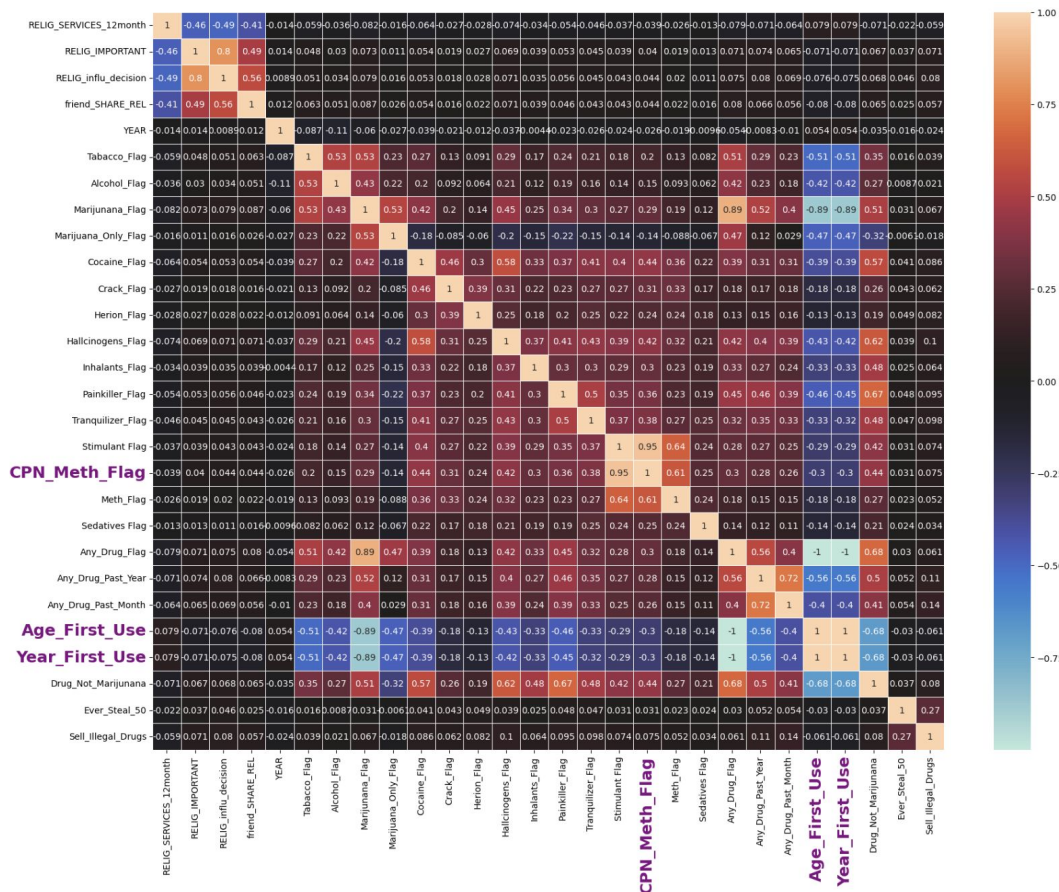**Figure 2A: Correlation Heatmap**

# Correlation Heatmap of subset 1



Figure 1B: Corr Heatmap of subset 1
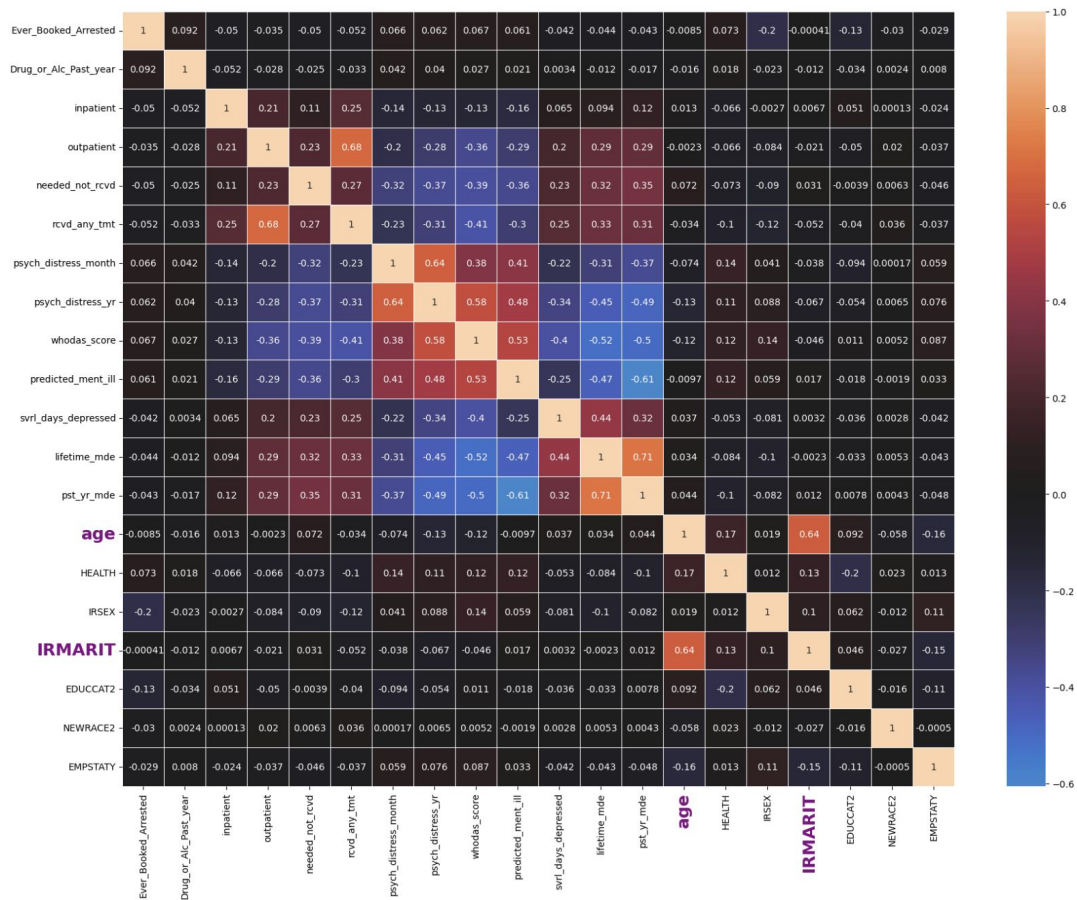
# Correlation Heatmap of subset 2
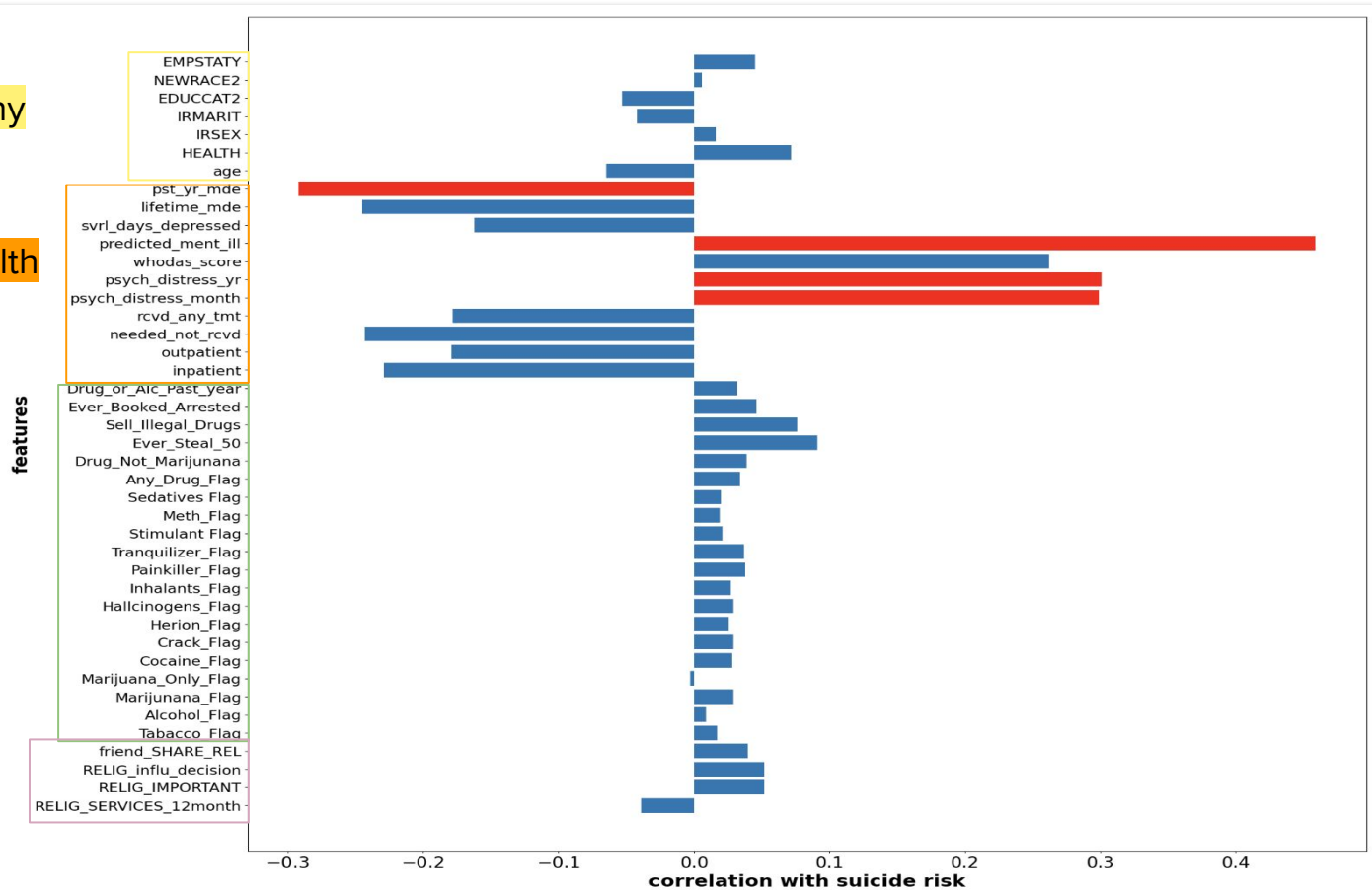


Figure 1B: Corr Heatmap of subset 2

# Linear Correlation with Suicide Risk

# Linear Relationship Between Features and Suicide Risk

**The most related features to the risk of committing suicide are also related to mental health/emotions.**
The top 5 are:

Binary:

- Serious Mental Illness (SMI) indicator: 0.459,
- Experiencing Serious Psychological Distress in Last Year: 0.301,
- Experiencing Serious Psychological Distress in Last Month: 0.299,
- Experiencing Major Depressive Episode in Past Year (For Adult): -0.292

Continuous:

- Level of Difficulty in Performing Daily Activities due to Problems with Mental Health: 0.262

# Discovery From the Most Correlated Features

- ~62% of individuals predicted to have **mental illness** (SMI) have a suicide risk > 0, with near 25% having a risk >= 3, 10% having a highest risk of 6.

- ~13% of individuals experiencing Major Depressive Episode in the last year have a risk >= 3

- ~10% of individuals experiencing Serious Psychological Distress in the last year have a risk >= 3

- A risk >= 3 means one has at least thought about and planed for committing suicide

# Hypothesis

We realized the risk of suicide is most correlated to features associated with recent mental health, so we are interested in

- whether having a good mental health is predicted to be able to help decrease the risk of someone committing suicide, and
- to what extent does it help

# Our Models

# K-Modes Clustering

K-modes is used for clustering categorical variables

Choosing Optimal K = 3

EDA:

- Generalize why one cluster would have higher suicide risk
- Create distribution graphs comparing the clusters



**Figure 3: Elbow Method**

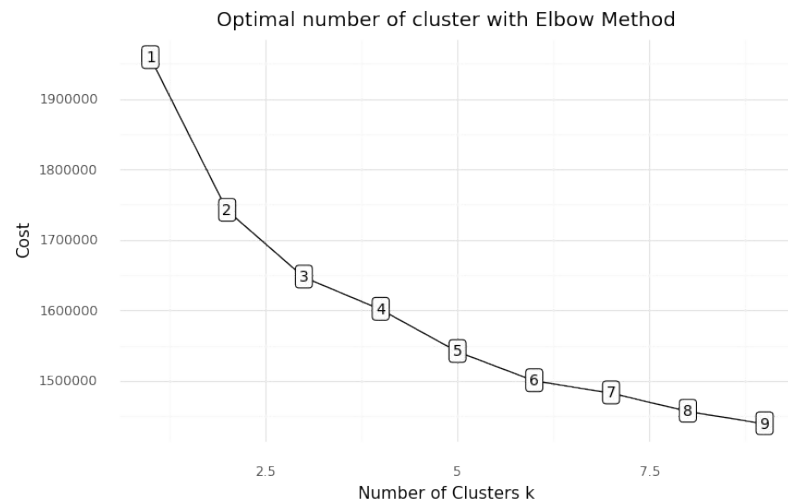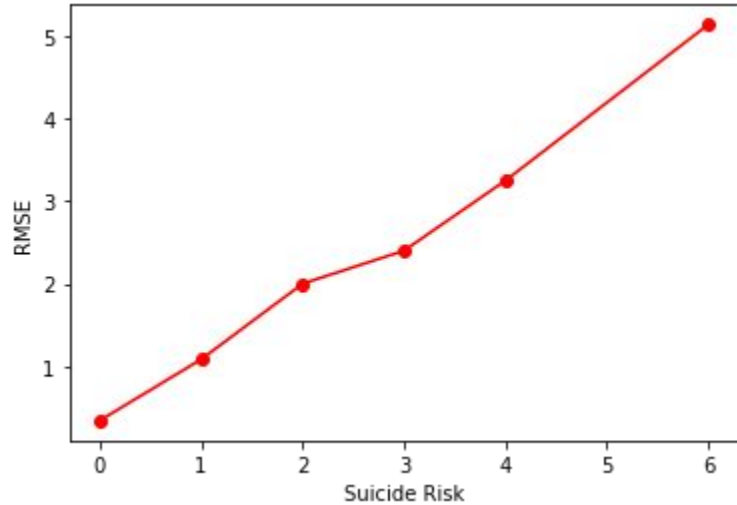| Cluster | suic_thought_pst_12month | suic_thought_pst_yr | suic_plan | suic_attempt |
|---|---|---|---|---|
| 0 | 0.959128 | 0.041670 | 0.012158 | 0.006088 |
| 1 | 0.953736 | 0.056209 | 0.017458 | 0.009477 |
| 2 | 0.928960 | 0.071269 | 0.022687 | 0.012462 |

# K-modes clustering

Permutation Testing

- **Null Hypothesis:** Both samples are drawn from the same population thus the observed TVD **is due to chance alone**.

- **Alternative Hypothesis:** Samples are drawn from two different populations thus the TVD **is not due to chance alone**, and there are other factors at hand.

Results: Significant difference in distributions of suicide risk, but hard to tell which features are most important.

# Regression Model Baseline
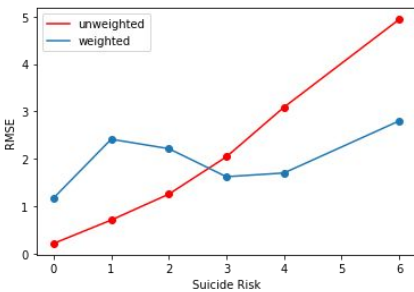
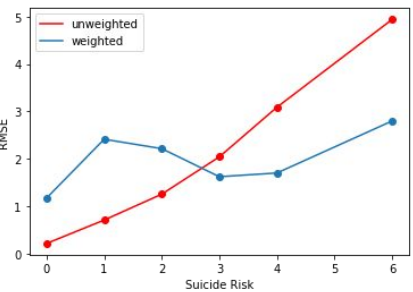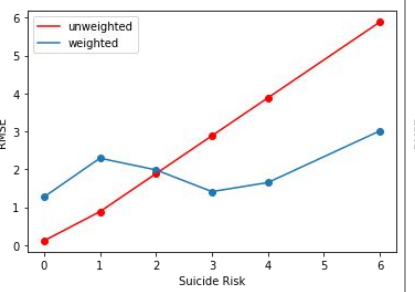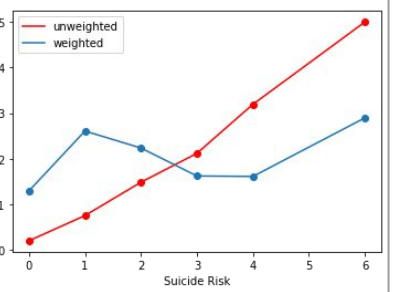K Neighbors Regressor

Overall RMSE: 0.6286

# Regression Model Criteria

Better to have better prediction performance (low RMSE) for high risk respondents at the expense of performance for low risk.

Idea: add sample weights

Weight = $2$^suicide_risk

|  | Least Squares | Ridge | LASSO | Random Forest |
|---|---|---|---|---|
| RMSE overall | 0.5208 | 0.5209 | 0.6042 | 0.5253 |
| RMSE weighted | 1.8405 | 1.8405 | 1.9511 | 1.9381 |
| RMSE by risk |  |  |  |  |

# Important Features

7 important features common in the least squares and ridge regression models

1. Predicted_ment_ill: predicted to have a serious mental illness
2. Inpatient: received inpatient mental health treatment in the past year
3. Psych_distress_yr: past year serious psychological distress indicator
4. Ever_Steal_50: ever stolen or tried to steal anything > $50
5. Drug_or_Alc_Past_year: used drugs or alcohol in the past year
6. Psych_distress_month: past month serious psychological distress indicator
7. EDUCCAT2: education level

*note: some features are still highly correlated with each other, and can be removed in further analysis

# Limitations/Challenges

# Challenges

Limitations include large dataset with many features, took a lot of time to understand and look through the dataset, reformulating the hypothesis question, reducing down our features, finding features that weren't highly correlated

Some features have many values representing missingness and we had to skim through each one. Some binary features use inconsistent values representing yes/no.

# Conclusion

# Conclusion

The features we found that impact suicide include not only mental health variables, but variables pertaining to education level, drug and alcohol use, and crime.

Although our modeling is not robust, it gives a good idea of the different signs that may indicate higher sucide risk, as well as the process others may use in risk modeling. It is important to understand the bias that will be present in risk modeling, as well as the socioeconomic and environmental factors that contribute to risk. Ideally, we would want to address the causes that lead to declining mental health, and learn how to promote and maintain positive mental health for all people.

# Thank you for your attention!