

Intro Analysis of Possible attributes contributing towards mental health in occupaing health

Python (Open Sourcing Mental Illness (OSMI) Survey on Health factors like :- family history, age, discuss\_CoWorker, discuss\_employer, sought treatment with mental health disorders.

Data Cleaning :- Column names were too long & this restricted importing the data into MySQL. renamed the column names into shorter abbreviations

Data was mostly categorical, converted to numerical for better visualization

### (OSMI) Open Sourcing Mental Survey

Total 123 Variables or 23 variables aligned with mental health were selected  
Data import to MySQL using

Shared data base 'I501sp21grp3-db' (user)

#### Markups

Data cleaning was done in SQL

Initially the column names were too big and of changed using SQL Queries

ALTER Table MHT\_Data

After removing changing the column names

Deleted null values

Deleted From MHT-Data

WHERE COALESCE (all the variables) is null

Then

UPDATE the Gender using SQL Query

(SET)

life cis woman, trans feminine, woman,  
cis female, transgender, female, female cis, female-ish,  
my sex is female to female using update

UPDATE MHDATA SET GENDER = 'male', where gender =  
cis heteromale, cis male, dude, m, male, male cis,  
male-ish, androgynous, man,

After Data cleaning

mysql → python (using mysql connector)

↳ Correlation matrix (positive correlation between

MHD - currently & family history

Negative correlation (disun coworker disun-employer,  
socgh treatment)

& heatmap

Exploratory data analysis

Normally distributed, outliers etc

Feature analysis

By taking all the columns as features required

↳ mostly → disun employer, disun coworker,  
features count foatoant, family history

Steps to identify the most significant features in the

dataset ↳ Step 1:

\* Independent features as 'X'  
attribute to be  
y → MHD - currently (predicted)

Step 2

logistic regression (categorical outcomes)  
variable

To perform logistic regression of feature  
Selection is key

Apply logistic regression on all features showed accuracy of 67.7%

Step-1 Scikit-learn → feature importance

Independent features as 'x' and MHD as as 'y' (attribute to be predicted) then performed logistic regression using all the features as independent variables, then the result showed accuracy of 67.7% (ROC Curve (Sensitivity vs Specificity) TPR (True Positive Rate) FPR (False Positive Rate))

Step-2 Then plotted feature importance to analyse most correlated & significant features in the dataset in python (the 'O' value shows no effect whereas - or +ve very high valued variable shows significant feature importance). Looking at the values features like Discm-employee, discm-coworker, Doubt- treatment & family history has more factors.

Step-3: Then performed logistic regression taking top 4 features as independent as 'x' variable & MHD as 'y'. I used scikit-learn package to import machine learning models. ~~Sklearn~~ Linear model imported logistic regression.

Step-4: Assign the input variable (x) with 4 independent features & y → MHD currently.

Training the model using fit() (train-test) 25% of data as test set and 75% train set

Logistic regression model showed it 68.7% & simultaneously plotted ROC curve

ROC (Receivers operator characteristic) (How well the model works)  
Family history → on higher feature importance

### Confusion matrix

to analyze the performance of the Ada-Boost model

Precision =  $\frac{\text{True positive}}{\text{True positive} + \text{False positive}}$

### Data visualization

#### Scatter plot - plotting

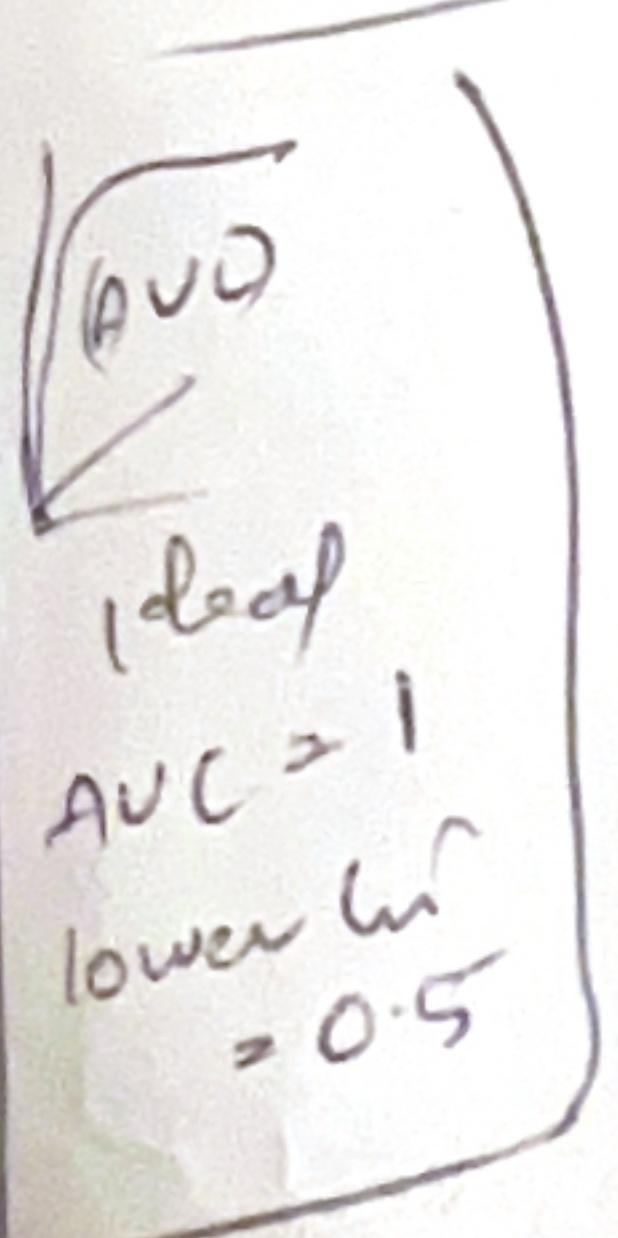


↳ dep → with int → sought treatment  
middle aged people seek treatment. So if treatment

#### multiple bar chart

If pivot table

dep = MHD anxiety, discuss\_CoWorker\_Egden  
not willing to discuss it with their  
coworker



#### Tableau

chloropleth maps (color mapping technique)

K → MHD Circles, areas in the world suffering  
from MHD.

because dependent variable is categorical in nature

#### man Whitney U test (nonparametric)

for hypothesis testing

P value < 0.05

Illy

why training  
and testing  
for machine  
learning

used to estimate the performance of  
machine learning algorithms for making  
predictions (understanding evaluation of algorithm  
performance) and how model performs with  
the new data.