

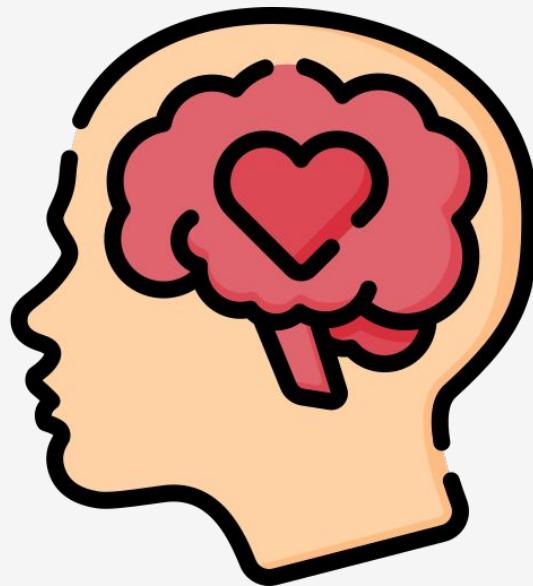


SC1015 Mini Project

Mental Health Predictor

REP2 group 1
Dion, Isaac, Ruiyun

Contents



01

Problem and Goals
Project overview
Datasets

02

Dataset EDA
Depression Dataset
Anxiety Dataset
Sentiment Analysis Dataset

03

Machine Learning
ANN (Depression) & results
SVM (Anxiety) & results
NLP model (Text based) & results

04

Recommendations
Future recommendations



01

Problem and Goals

1 in 8

people in the world live with a mental disorder

With the most common ones being:



Anxiety Disorder



Depression

Yet, approximately
50% goes
undiagnosed

Yet, approximately

Can we use machine learning to accurately predict the likelihood of depression and anxiety in University Students, given that many cases remain undiagnosed?

undiagnosed



MindScope

Mental Health Predictor app

General Information section

MindScope

Welcome to our Mental Health Predictor App, MindScope

This app is designed to support your mental well-being by helping you better understand potential indicators of anxiety and depression. By filling out a short questionnaire based on well-established predictive factors, you'll receive a personalized percentage estimate indicating the likelihood of experiencing symptoms related to anxiety or depression.

Please note: This tool is **not** a diagnostic instrument and should not replace professional mental health advice, diagnosis, or treatment. It is built on historical data and research to provide insight—not conclusions—into your mental health.

Your responses are handled with care and confidentiality. We hope this tool serves as a helpful starting point for reflection, awareness, and seeking further support if needed.

Gender

- Male
 Female

Age

Your answer _____

CGPA

Your answer _____

Work/Study Hours (per day)

Your answer _____

Next

Clear form

MindScope

Rating Section

For each of the following questions, please give a rating from 1-5, with 1 being the least intense and 5 being the most intense

How pressured do you feel by academics?

put 0 if not applicable

0 1 2 3 4 5

Low pressure

High pressure

How pressured do you feel by work?

put 0 if not applicable

0 1 2 3 4 5

High pressure

Low pressure

How satisfied are you with your studies?

put 0 if not applicable

0 1 2 3 4 5

Low satisfaction

High satisfaction

How satisfied are you with your job?

put 0 if not applicable

0 1 2 3 4 5

Low satisfaction

High satisfaction

How stressed are you by finance related matters

1 2 3 4 5

Low stress High stress

How satisfied are you by your social life?

1 2 3 4 5

Low satisfaction High satisfaction

Would you say that your future is secure?

1 2 3 4 5

Low security High security

How isolated would you say you are from friends and family?

1 2 3 4 5

Low isolation High isolation

Back

Next

Clear form

Rating section

Categorical section

MindScope

Categorical section

What is your sleep duration?

- <5 hours
- 5-6 hours
- 7-8 hours
- >8 hours

How would you describe your dietary habits?

- Unhealthy
- Moderate
- Healthy

Have you ever had suicidal thoughts?

- Yes
- No

Any family history of mental illness?

- Yes
- No

Back

Next

Clear form

Text based section

MindScope

Text section

Provide a brief write up about how your day was, include details and how you felt (about 150 words)

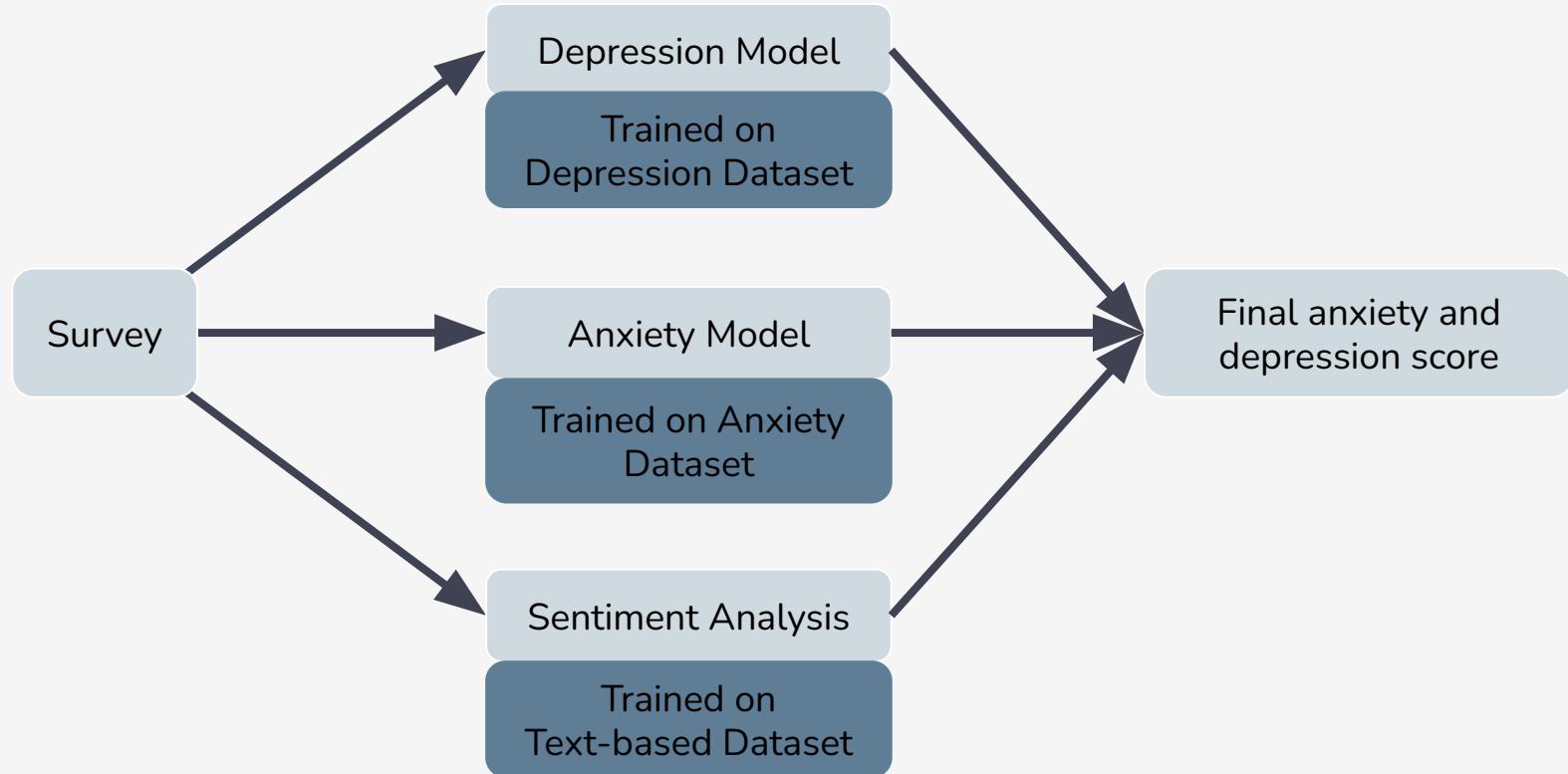
Your answer

Back

Submit

Clear form

MindScope Overview



Dataset:

Targeting University students and young adults

kaggle



02

Dataset EDA

Depression dataset: Overview



Depression dataset: Understand Data

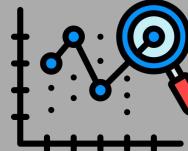
```
id                         int64
Gender                     object
Age                        float64
City                       object
Profession                  object
Academic Pressure          float64
Work Pressure               float64
CGPA                       float64
Study Satisfaction          float64
Job Satisfaction            float64
Sleep Duration              object
Dietary Habits              object
Degree                      object
Have you ever had suicidal thoughts ?    object
Work/Study Hours            float64
Financial Stress             object
Family History of Mental Illness        object
Depression                  int64
dtype: object
```

- Combination of numeric, categorical and binary columns
- 27901 values
- output column: Depression (binary)

Depression dataset: Correlation Analysis



Categorical columns
encoded

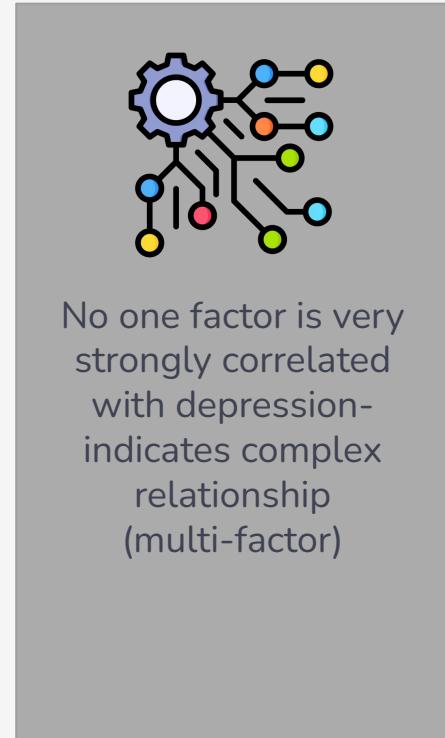
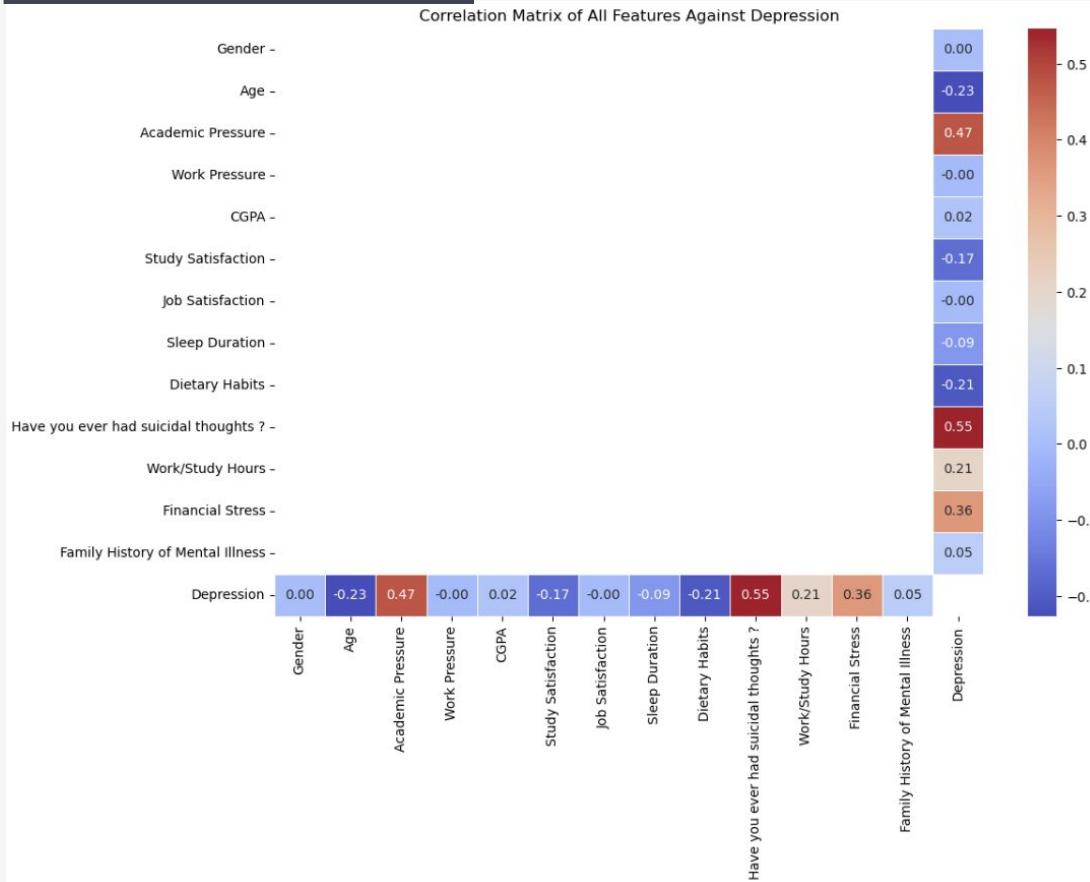


Numeric and
categorical columns
correlation with
depression calculated
using **point-biserial**
correlation

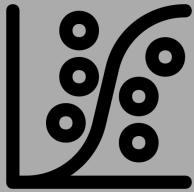
01
10

Binary columns
correlation with
depression calculated
with **phi-squared**
coefficient

Depression dataset: Correlation Analysis



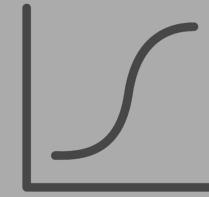
Depression dataset: Linearity Analysis



Regression plots
plotted with
`logistic=True`

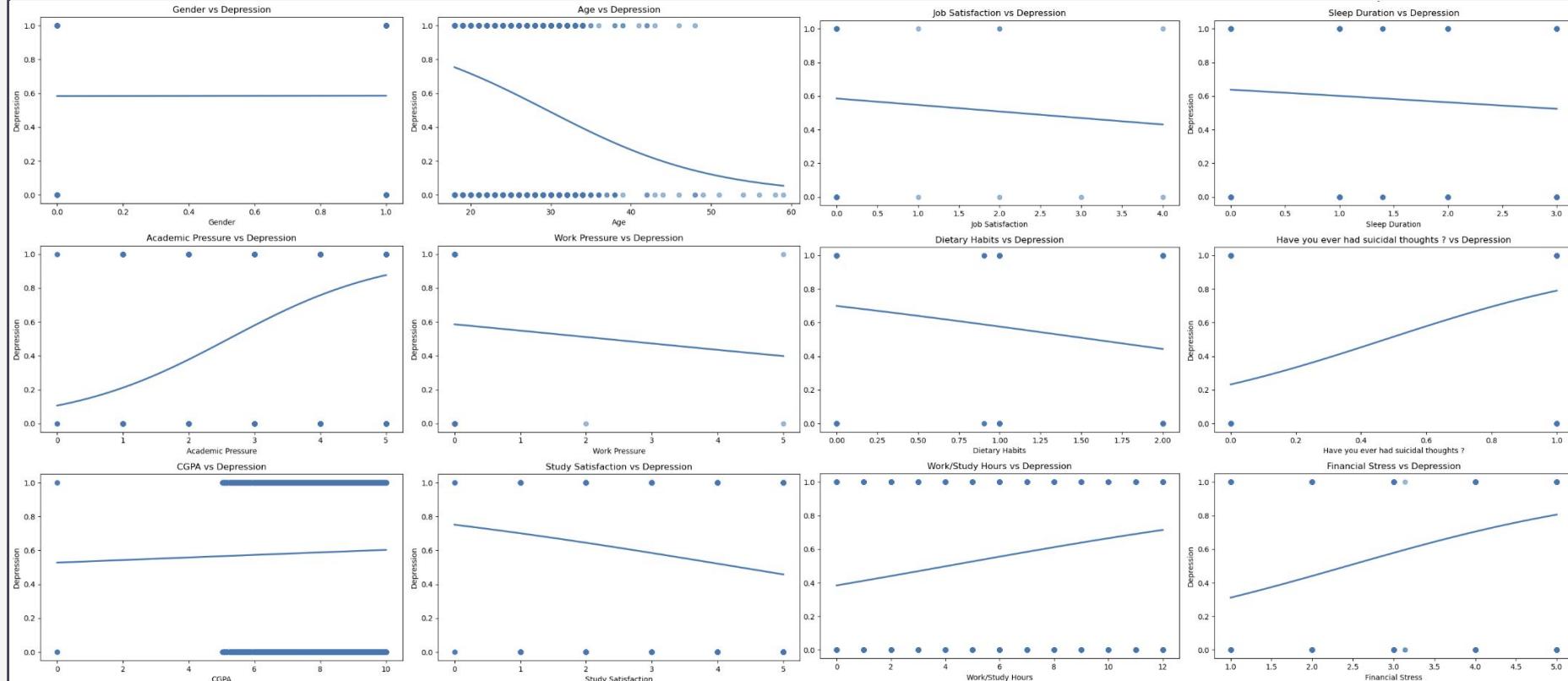


Visualise correlation
with depression
(output)

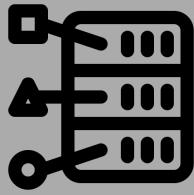


Sigmoid shapes indicate
good fit for logistic
regression model

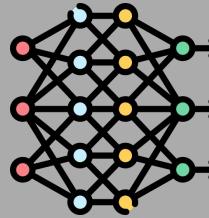
Depression dataset: Linearity Analysis



Depression dataset: Model Selection



- Vast amount of data (27901)
- individually low correlation with depression
- combination of flat, linear and sigmoidal logistic relationships



Select ANN
Individual features may look weak or flat (non linear or complex relationship), but **combination could be highly predictive**



ANN can detect subtle,
multi-feature interactions effectively

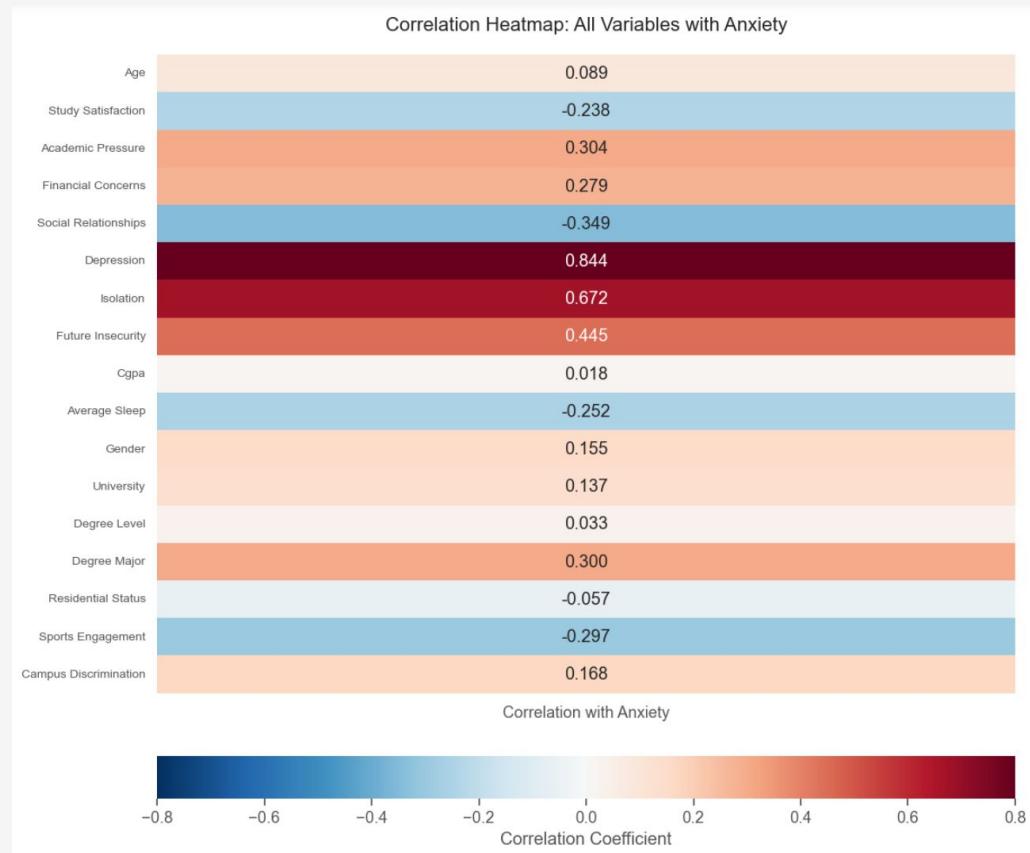
Anxiety dataset: Understand Data

Data Info:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 87 entries, 0 to 86
Data columns (total 21 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   gender          87 non-null      object  
 1   age              87 non-null      int64  
 2   university       87 non-null      object  
 3   degree_level    87 non-null      object  
 4   degree_major    87 non-null      object  
 5   academic_year   87 non-null      object  
 6   cgpa             87 non-null      object  
 7   residential_status 87 non-null    object  
 8   campus_discrimination 87 non-null    object  
 9   sports_engagement 87 non-null    object  
 10  average_sleep   87 non-null      object  
 11  study_satisfaction 87 non-null    int64  
 12  academic_workload 87 non-null    int64  
 13  academic_pressure 87 non-null    int64  
 14  financial_concerns 87 non-null    int64  
 15  social_relationships 87 non-null    int64  
 16  depression       87 non-null      int64  
 17  anxiety          87 non-null      int64  
 18  isolation         87 non-null      int64  
 19  future_insecurity 87 non-null      int64  
 20  stress_relief_activities 87 non-null    object  
dtypes: int64(10), object(11)
memory usage: 14.4+ KB
```

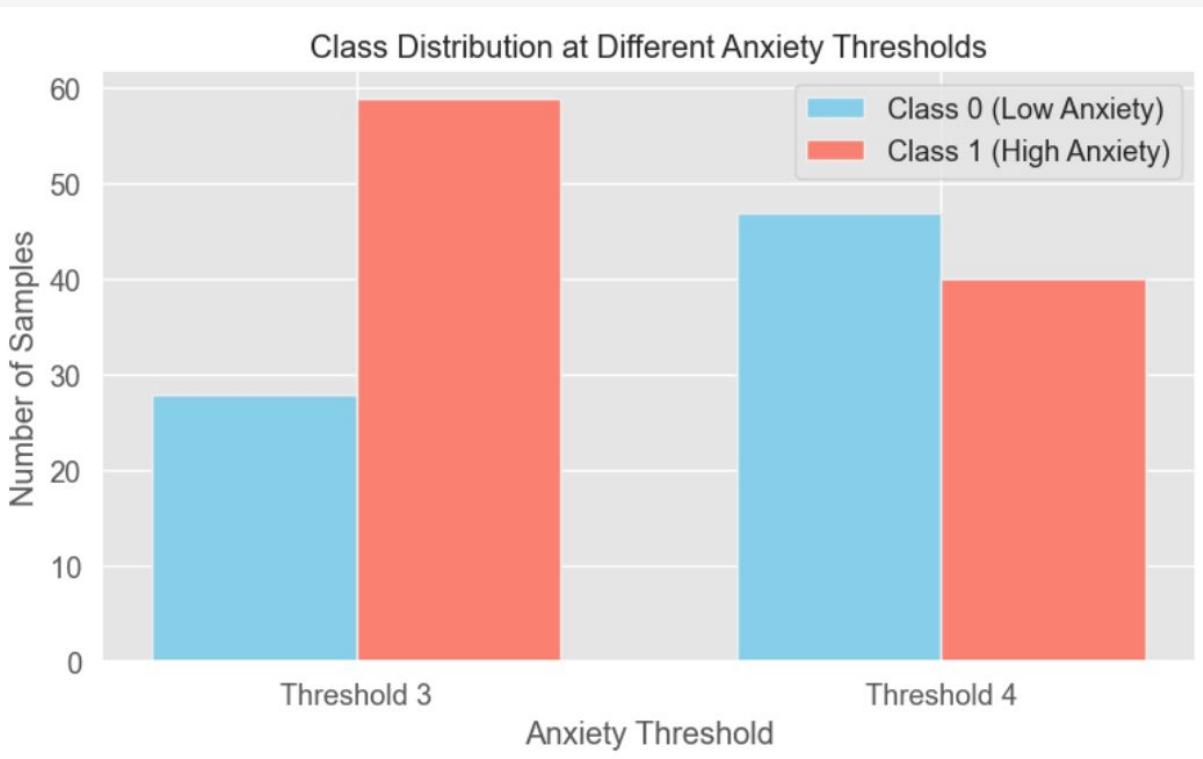
- Combination of numeric and categorical data
- 86 values

Anxiety dataset: Correlation Analysis



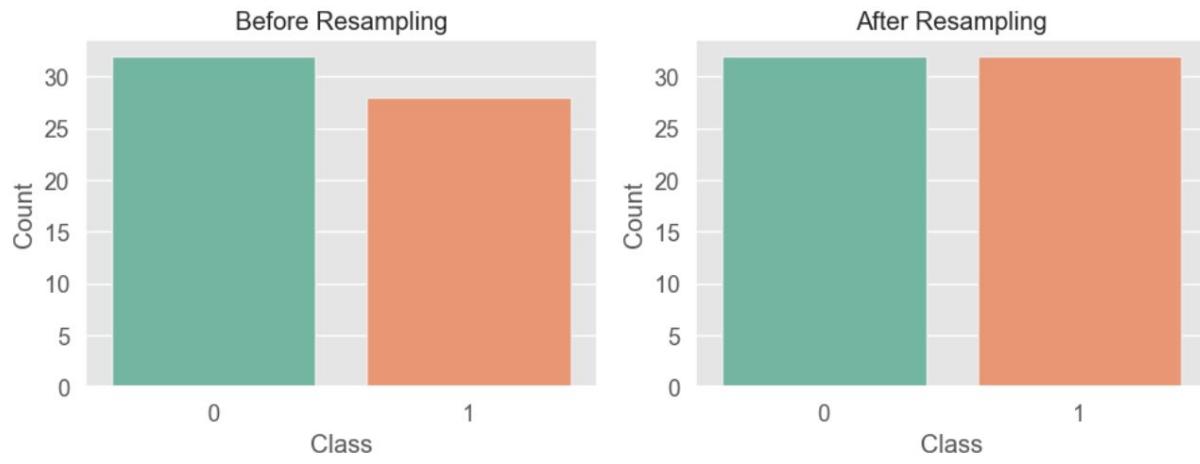
- Top 3 correlated variables with anxiety is
 1. Isolation
 2. Future Insecurity
 3. Social Relationship (negatively correlated)
- Excluding depression despite high correlation to avoid circular logic (since we predicting depression separately)

Anxiety dataset: Data Preparation



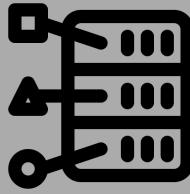
- Reclassifying Anxiety as High or Low (for binary output)
- Threshold of 4 gives more balanced distribution

Anxiety dataset: Data Preparation

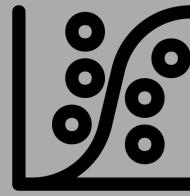


- Handled class imbalance after splitting
- Applied feature scaling
- Upsampled minority class

Anxiety dataset: Model Selection



- Small dataset (86 values)
- Binary Classification with Multivariate Input



Possible models:

- Logistic regression
- Random Forest
- SVM
- XGBoost



Selection Criteria:

- Performance on Validation set
- Resistance to overfitting

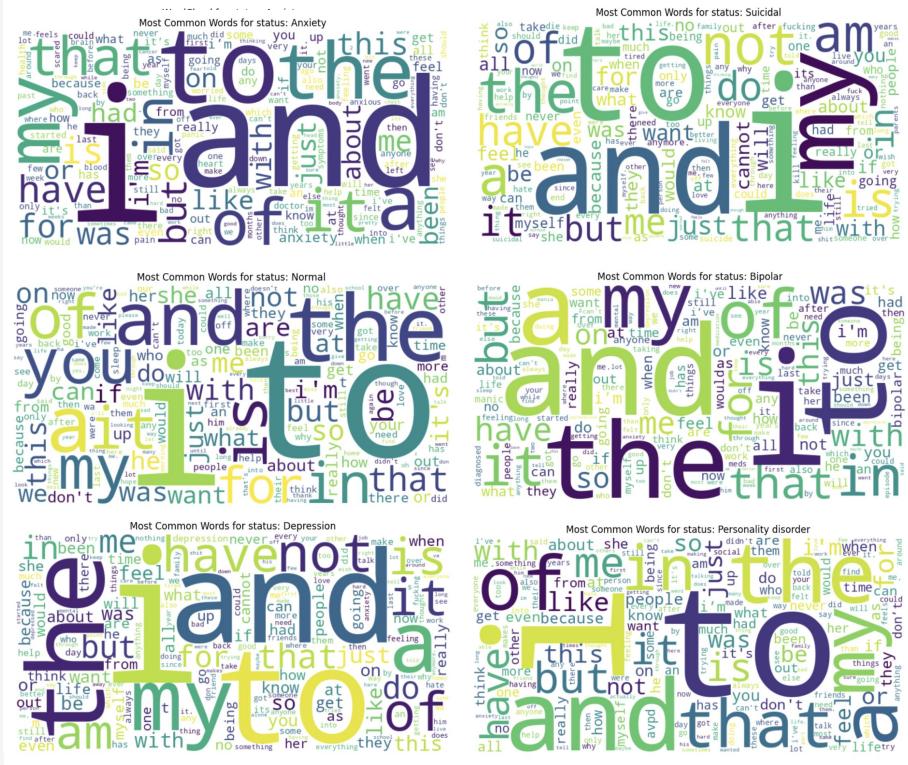
Sentiment Analysis Dataset: Understand Data

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 53043 entries, 0 to 53042
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   index        53043 non-null   int64  
 1   statement    52681 non-null   object  
 2   status        53043 non-null   object  
dtypes: int64(1), object(2)
memory usage: 1.2+ MB
```

```
Nulls per column:
index          0
statement     362
status          0
dtype: int64
Empty 'statement' entries: 0
```

- Combination of textual and categorical data and columns.
- 53043 values / rows
- Index column: Index
- Statement column: Contains string of text
- Status column: Depression, Anxiety, Personality Disorder, Bipolar, Suicidal or Normal

Sentiment Analysis dataset: Understand Data



First look:

- Many common words across different statuses like i, to, and, my, that, the.
- These words do not help much to distinct statuses
- More needs to be to clean the dataset to make each group of text more unique.

Sentiment Analysis dataset: Understand Data



TF-IDF:

term freq - inverse document freq

- Use TF-IDF to find more important words in each status
- Many common words across different statuses still present like 'feel', 'want', 'know'.
- More needs to be done to clean the dataset to make each group of text more unique.

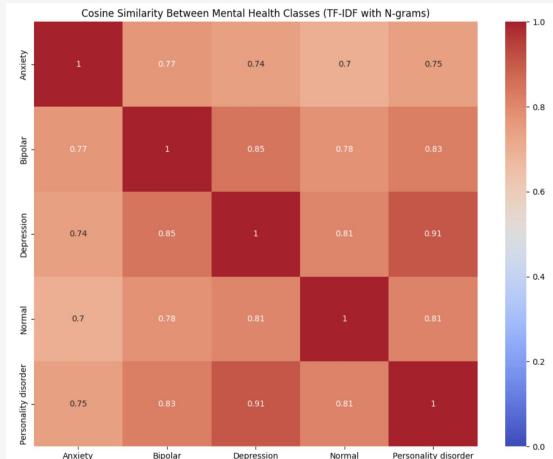
Sentiment Analysis dataset: Understand Data



TF-IDF + N-grams:
term freq - inverse document freq

- Use TF-IDF with N-grams to find more important words in each status
- We can see that the words across statuses are much more distinct.
- However, there are some that are groups of text that are still similar like Personality Disorder & Depression and Bipolar & Depression.

Sentiment Analysis dataset: Data Insights



- Common words used by people with Personality Disorder are significantly similar to words used by people with Depression (0.91). The next most similar pair would be words used by people with Bipolar and Depression (0.85).
- The dataset for people with Personality Disorder and Bipolar are the smallest at 1201 and 2877 statements.
- Might be better to drop these statements to minimise confusion for our model in discerning between mental health statuses.

Sentiment Analysis dataset: Insight

Validation

Status Classes Ranked by Similarity of Distinctive Words:

Rank 1: Normal and Personality disorder: 0.0546
Rank 2: Bipolar and Personality disorder: 0.0314
Rank 3: Normal and Bipolar: 0.0312
Rank 4: Depression and Personality disorder: 0.0225
Rank 5: Anxiety and Normal: 0.0000
Rank 6: Anxiety and Depression: 0.0000
Rank 7: Anxiety and Bipolar: 0.0000
Rank 8: Anxiety and Personality disorder: 0.0000
Rank 9: Normal and Depression: 0.0000
Rank 10: Depression and Bipolar: 0.0000

Cosine Similarity Matrix between Distinctive Word Sets:

	Anxiety	Normal	Depression	Bipolar	\
Anxiety	1.0	0.000000	0.000000	0.000000	
Normal	0.0	1.000000	0.000000	0.031231	
Depression	0.0	0.000000	1.000000	0.000000	
Bipolar	0.0	0.031231	0.000000	1.000000	
Personality disorder	0.0	0.054557	0.022511	0.031424	

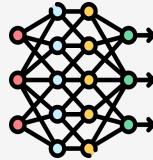
	Personality disorder
Anxiety	0.000000
Normal	0.054557
Depression	0.022511
Bipolar	0.031424
Personality disorder	1.000000

- Conducted log-likelihood analysis to identify the most distinct word in each class
- Cosine similarity for the Distinctive words in each class and ranked pairs by similarity
- Most distinct words used by people with Personality Disorder are very similar to normal people (Rank 1), Bipolar Disorder (Rank 2) and Depression (Rank 4).
- It can also be seen that the most distinct words used by Bipolar people are very similar to normal people (Rank 3).
- With such considerations, this confirms our insights and we will remove both the statements for Personality Disorder and Bipolar statuses.

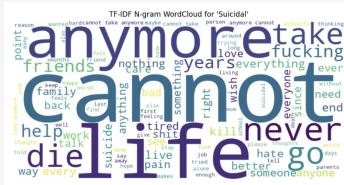
Sentiment Analysis dataset: Final Thoughts



Observations



Model Selection



- While Anxiety, Depression, and Normal classes have some distinct top TF-IDF words, their overall language usage remains highly similar (as shown by cosine similarity scores).
- To go beyond surface-level word usage, we should use modern NLP approach like Large Language Models. This enables the model to learn semantic representations of entire statements, capturing context, tone, and subtle cues. As a result, we can achieve more accurate and meaningful classification of mental health statuses.
- Words used by people afflicted with Suicidal tendencies will be flagged. Most language models have safeguards that do not allow these words. Hence we shall drop the Suicidal status.



03

Machine Learning

ANN (Depression dataset)

```
model = Sequential([
    Dense(16, input_dim=X.shape[1], activation='relu'),
    Dense(8, activation='relu'),
    Dense(1, activation='sigmoid') # sigmoid outputs a probability between 0 and 1
])

# Compile the model
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
```

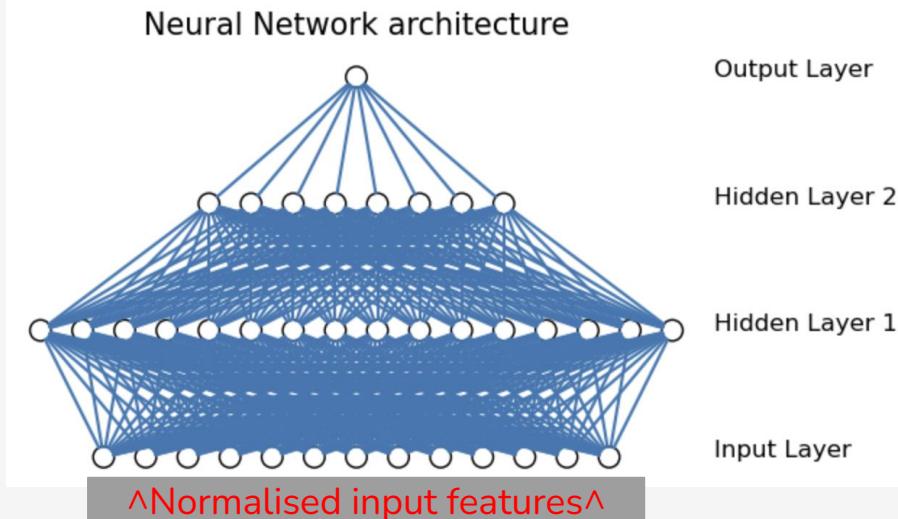
Model: "sequential_1"

Layer (type)	Output Shape	Param #
=====		
dense_3 (Dense)	(None, 16)	224
dense_4 (Dense)	(None, 8)	136
dense_5 (Dense)	(None, 1)	9
=====		

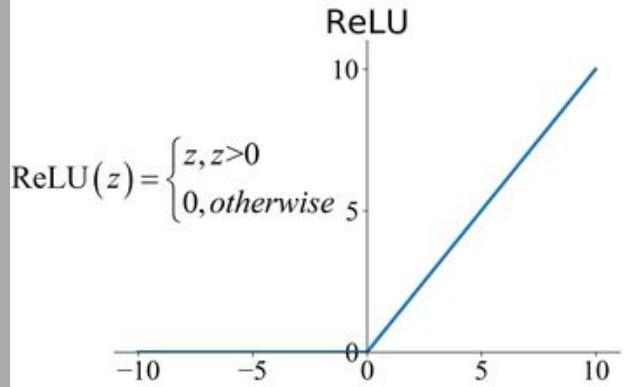
Total params: 369

Trainable params: 369

Non-trainable params: 0

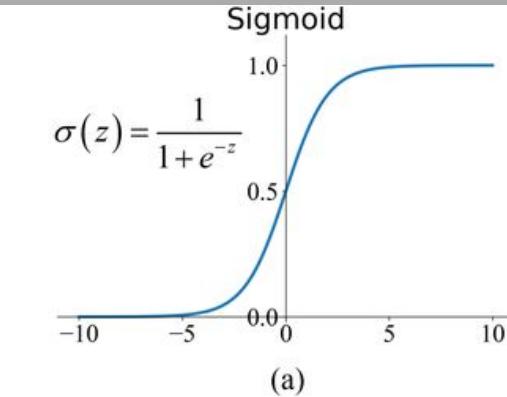


ANN (Depression dataset)



Hidden layers uses **RELU activation function**

- only allows x to pass through if $x>0$. 0 otherwise
- **capture piecewise linear patterns**
- **ignores flat features**, helps model to focus on meaningful signals



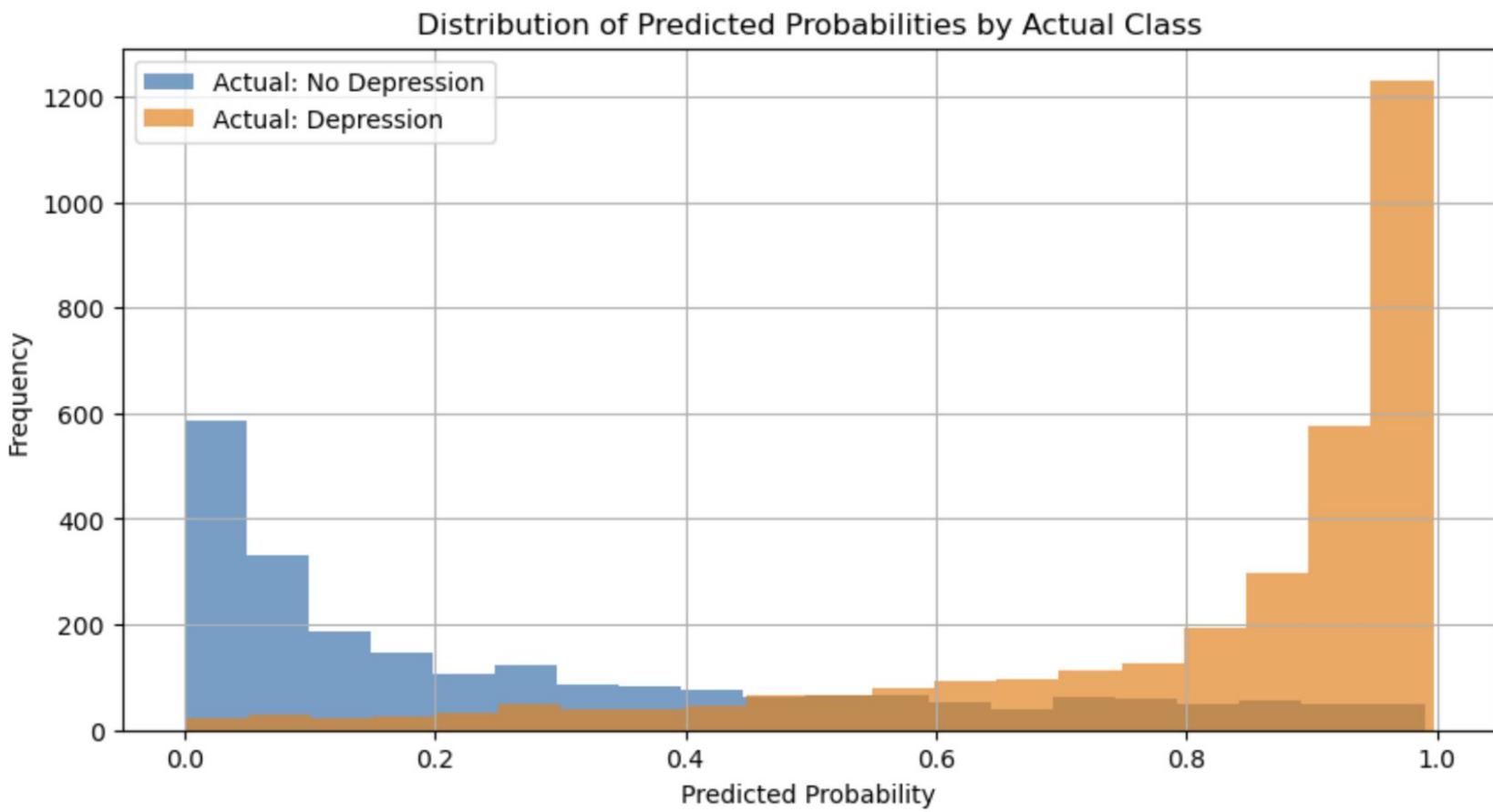
Output layer uses **Sigmoid activation function**

- sigmoid squashes final output into range [0,1]
- useful for probabilistic prediction



Loss calculated using **Binary cross entropy**- backpropagation to adjust weights using **adam optimiser**

ANN (Depression dataset)



ANN (Depression dataset)

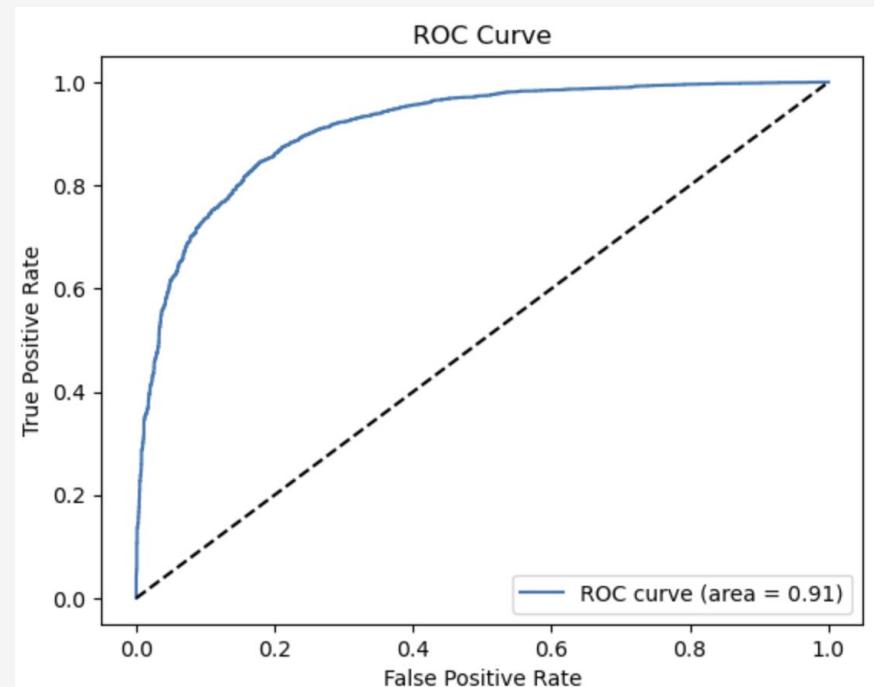
Classification Report:

	precision	recall	f1-score	support
0	0.83	0.77	0.80	2343
1	0.84	0.89	0.86	3238
accuracy			0.84	5581
macro avg	0.84	0.83	0.83	5581
weighted avg	0.84	0.84	0.84	5581

Confusion Matrix:

```
[[1800  543]
 [ 368 2870]]
```

ROC AUC Score: 0.91



Model Performance (Anxiety Dataset)

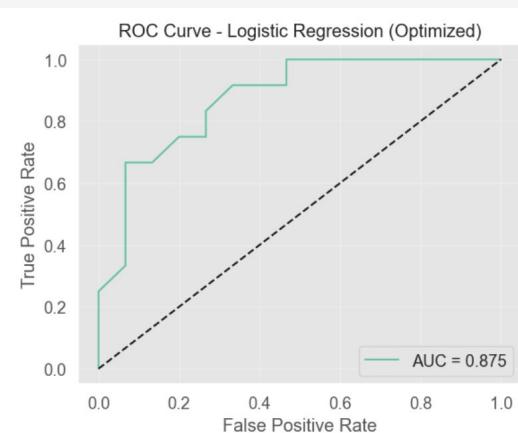
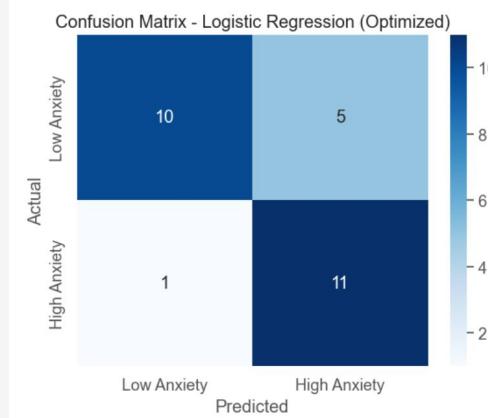
1. Logistic Regression

--- Logistic Regression (Optimized) ---

Optimal threshold: 0.434

Classification Report:

	precision	recall	f1-score	support
0	0.91	0.67	0.77	15
1	0.69	0.92	0.79	12
accuracy			0.78	27
macro avg	0.80	0.79	0.78	27
weighted avg	0.81	0.78	0.78	27



Model Performance (Anxiety Dataset)

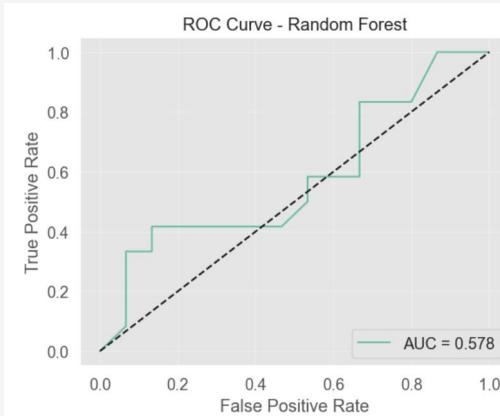
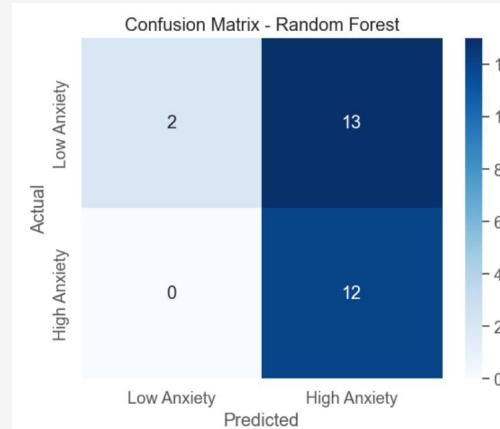
2. Random Forest

--- Random Forest ---

Optimal threshold: 0.020

Classification Report:

	precision	recall	f1-score	support
0	1.00	0.13	0.24	15
1	0.48	1.00	0.65	12
accuracy			0.52	27
macro avg	0.74	0.57	0.44	27
weighted avg	0.77	0.52	0.42	27



Model Performance (Anxiety Dataset)

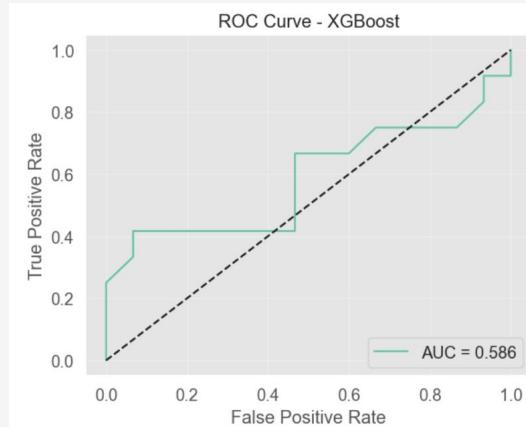
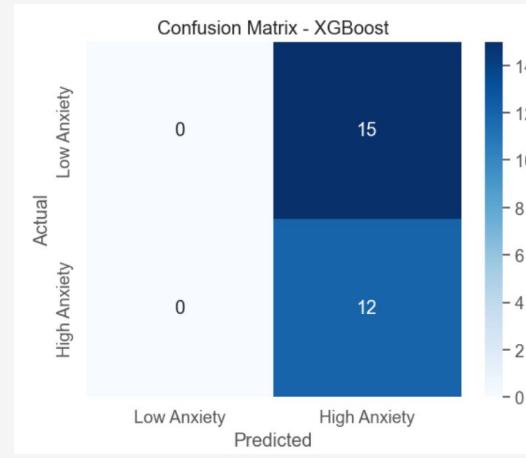
3. XGBoost

--- XGBoost ---

Optimal threshold: 0.000

Classification Report:

	precision	recall	f1-score	support
0	0.00	0.00	0.00	15
1	0.44	1.00	0.62	12
accuracy			0.44	27
macro avg	0.22	0.50	0.31	27
weighted avg	0.20	0.44	0.27	27



Model Performance (Anxiety Dataset)

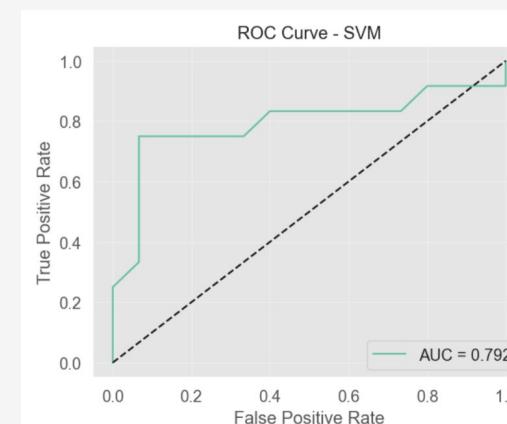
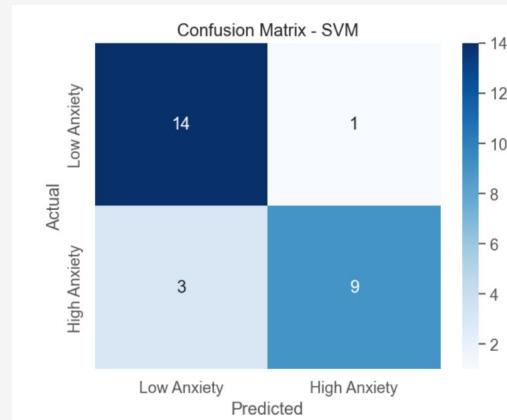
4. Best Model SVM

--- SVM ---

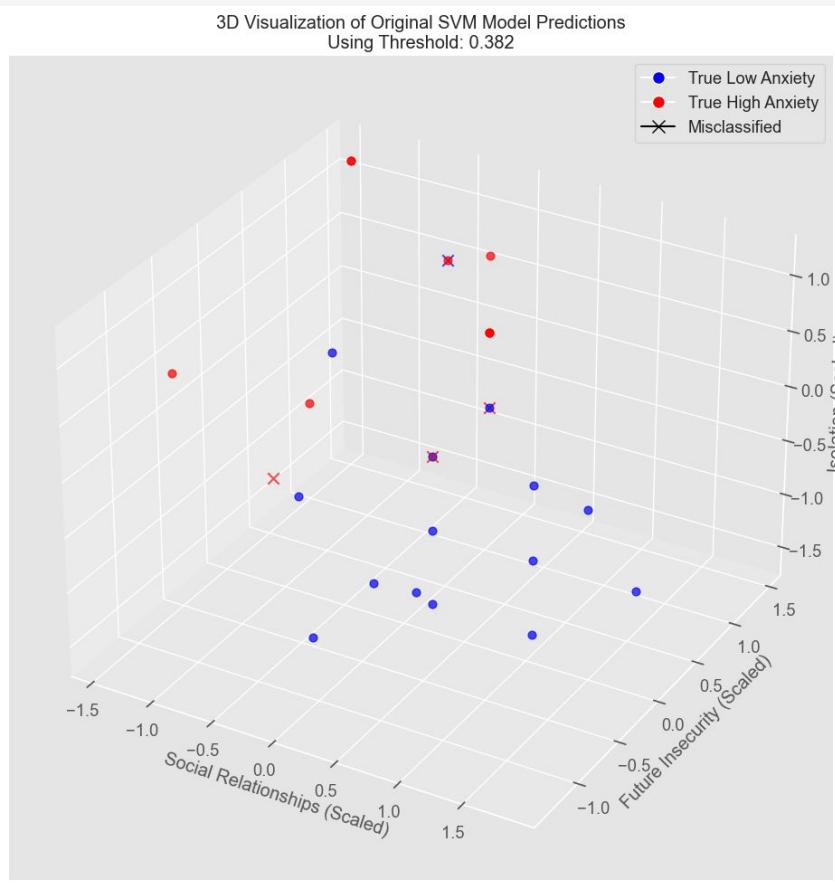
Optimal threshold: 0.382

Classification Report:

	precision	recall	f1-score	support
0	0.82	0.93	0.88	15
1	0.90	0.75	0.82	12
accuracy			0.85	27
macro avg	0.86	0.84	0.85	27
weighted avg	0.86	0.85	0.85	27



SVM (Anxiety Dataset)

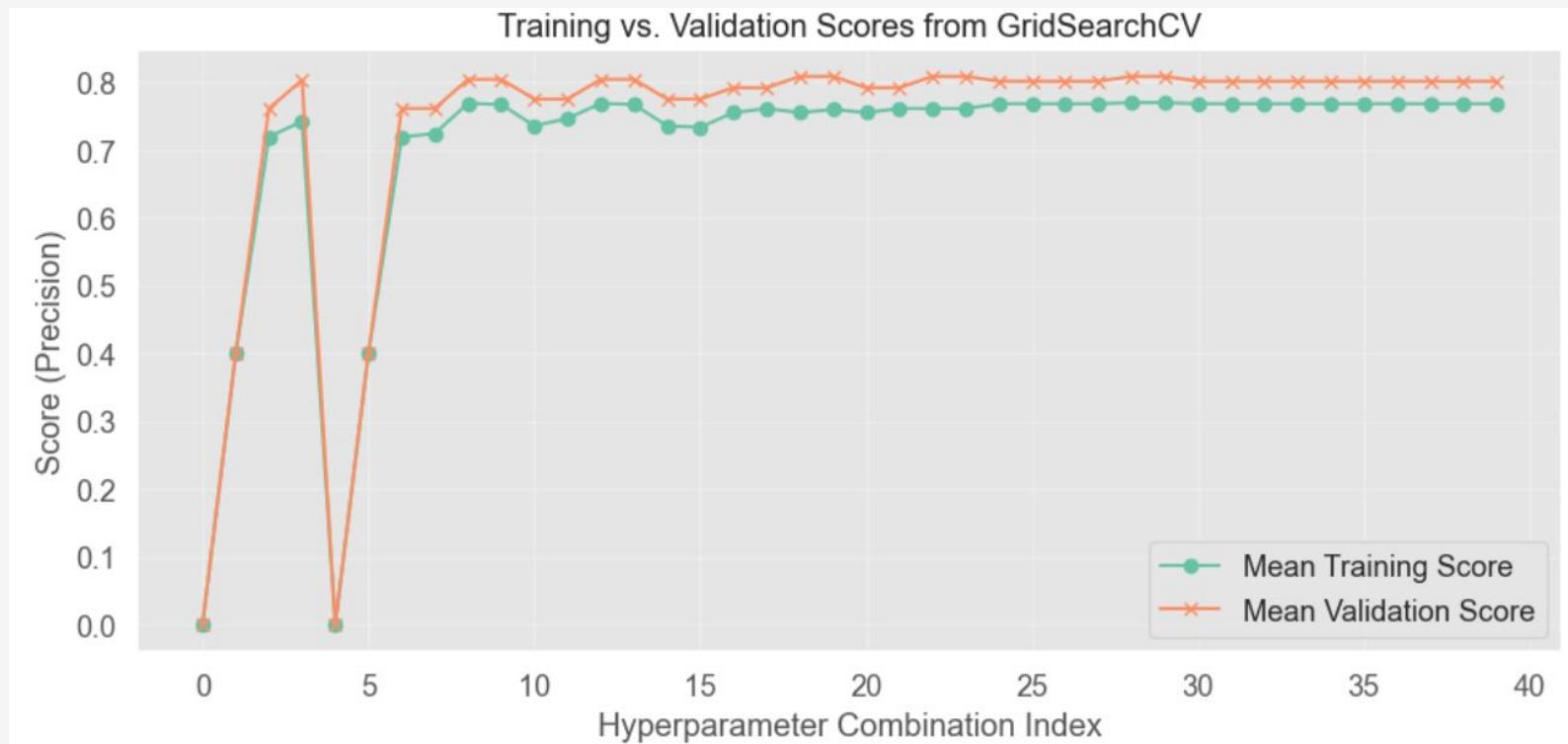


- Radial Basis Function (RBF) to create non-linear decision boundaries
- SVM finds the optimal hyperplane that maximises the margin between high and low anxiety cases

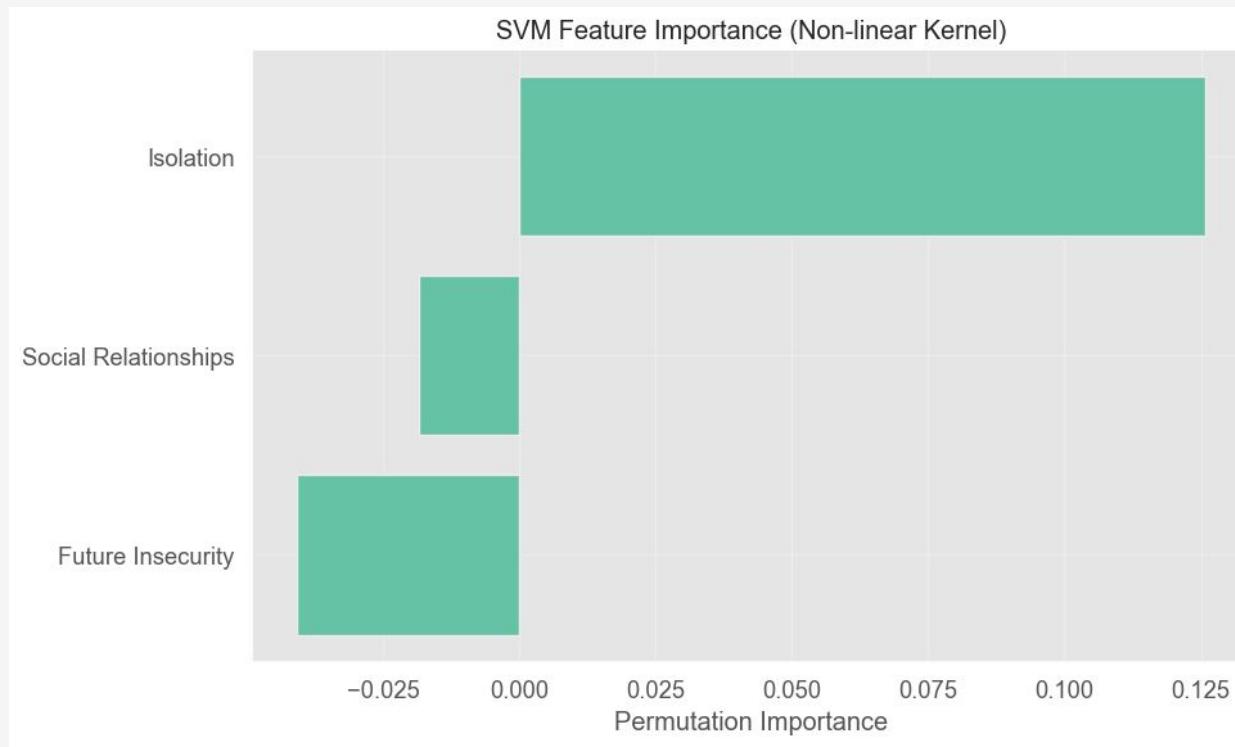
scaler = StandardScaler()

SVM with threshold 0.382

SVM Results (Anxiety Dataset)



SVM Data Insights (Anxiety Dataset)



Model Selection (Senti- Analysis Dataset)

Training VS Fine-Tuning a pre-trained Model



*Leverage on powerful,
pre-existing knowledge of
language*



*Use highly optimized
and proven model
architectures*



*Faster and less
computationally
intensive*

Fine-Tuning LLM (Sentiment Analysis Dataset)

```
from unsloth import FastLanguageModel

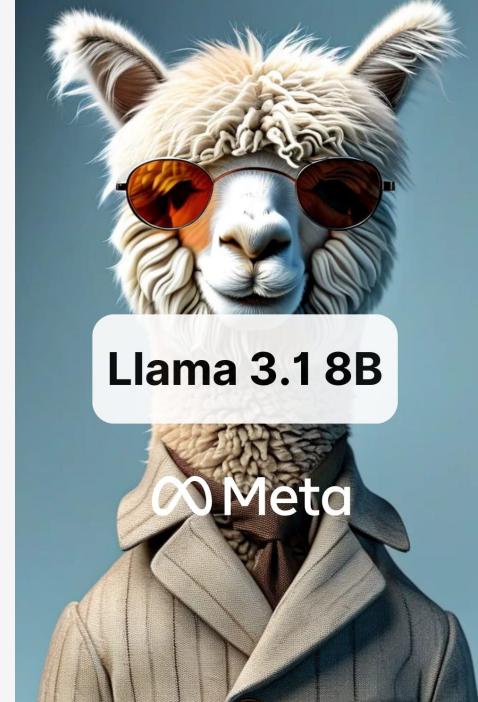
max_seq_length = 1024 # Choose any! We auto support RoPE Scaling internally!
dtype = None # None for auto detection. Float16 for Tesla T4, V100, Bfloat16 for Ampere+
load_in_4bit = True # Use 4bit quantization to reduce memory usage. Can be False.

model, tokenizer = FastLanguageModel.from_pretrained(
    model_name = "unsloth/Meta-Llama-3.1-8B-Instruct-bnb-4bit",
    max_seq_length = max_seq_length,
    dtype = dtype,
    load_in_4bit = load_in_4bit,
    token = HF_TOKEN, # use one if using gated models like meta-llama/Llama-2-7b-hf
)

# Set model configuration parameters similar to reference code
model.config.use_cache = False
model.config.pretraining_tp = 1

# Ensure pad token is set correctly
if tokenizer.pad_token_id is None:
    tokenizer.pad_token_id = tokenizer.eos_token_id
```

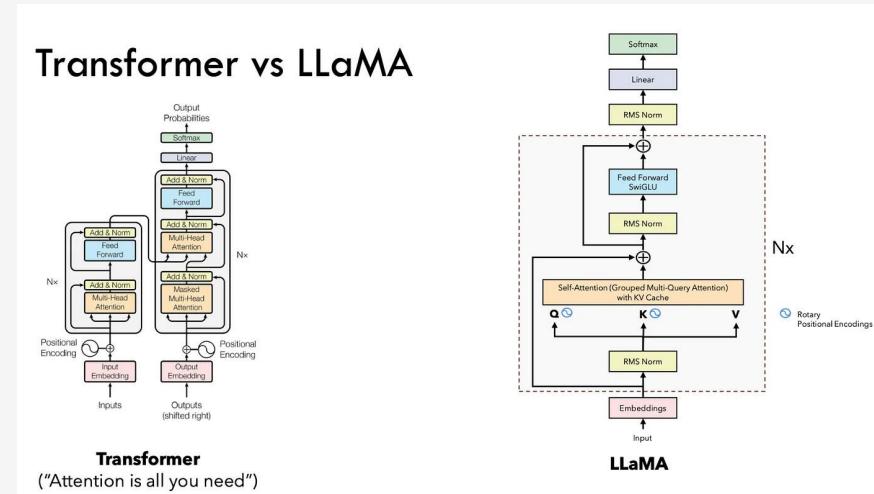
- Balance of computational efficiency and performance
- Suitable for deployment on normal consumer-grade hardware, but still delivers
- 4-bit quantized for memory efficiency



Model Overview

Model Differences

- Encoder-Decoder vs. Decoder-only
- LayerNorm vs. RMSNorm
- Absolute Sinusoidal vs. Rotary/RoPE
- Multi-Head Attention vs. Grouped-Query Attention
- Standard FFN + ReLU vs. SwiGLU



(**In comparison with standard transformer proposed in the landmark paper, “Attention is all you need”)

Model Performance (Senti- Analysis Dataset)

Pre-Fine Tuning

- Already rather accurate having an accuracy of around 73%
- However, it is not very accurate when classifying people with anxiety.
- Some overfitting may exist, maybe just classifying anyone that sounds NOK to be Depressed.

Accuracy: 0.730

Accuracy for label Normal: 0.746

Accuracy for label Depression: 0.853

Accuracy for label Anxiety: 0.591

Classification Report:

	precision	recall	f1-score	support
Normal	0.76	0.75	0.75	389
Depression	0.67	0.85	0.75	389
Anxiety	0.81	0.59	0.68	389
micro avg	0.73	0.73	0.73	1167
macro avg	0.75	0.73	0.73	1167
weighted avg	0.75	0.73	0.73	1167

Confusion Matrix:

```
[[290  82  17]
 [ 14 332  38]
 [ 76  83 230]]
```

Model Performance (Senti- Analysis Dataset)

Fine Tuning

- Trained the model for 1 epoch to save time
- Trained the model on the balanced dataset obtained from EDA.
- Split dataset into train-eval-test splits

```
Balanced class distribution:      Eval set class distribution:  
status  
Anxiety    3888          status  
Depression 3888          Depression   389  
Normal     3888          Normal       389  
Name: count, dtype: int64      Anxiety      388  
                                Name: count, dtype: int64  
  
Train set class distribution:    Test set class distribution:  
status  
Anxiety    3111          status  
Normal     3110          Normal       389  
Depression 3110          Anxiety      389  
Name: count, dtype: int64      Depression   389  
                                Name: count, dtype: int64
```

Run summary:

eval/loss	2.04277
eval/runtime	42.8131
eval/samples_per_second	27.235
eval/steps_per_second	6.82
total_flos	5.316330584776704e+16
train/epoch	0.99968
train/global_step	583
train/grad_norm	0.40559
train/learning_rate	0.0
train/loss	2.0802
train_loss	2.13006
train_runtime	3452.8734
train_samples_per_second	2.702
train_steps_per_second	0.169

Model Performance (Senti- Analysis Dataset)

Post-Fine Tuning

- Prediction accuracy increased significantly
- Achieves an average accuracy of 95.5% on prediction

Accuracy: 0.955

Accuracy for label Normal: 0.982

Accuracy for label Depression: 0.938

Accuracy for label Anxiety: 0.943

Classification Report:

	precision	recall	f1-score	support
Normal	0.96	0.98	0.97	389
Depression	0.94	0.94	0.94	389
Anxiety	0.97	0.94	0.95	389
micro avg	0.96	0.95	0.96	1167
macro avg	0.96	0.95	0.96	1167
weighted avg	0.96	0.95	0.96	1167

Confusion Matrix:

```
[[382  6  1]
 [ 9 365 12]
 [ 5 17 367]]
```

Model Performance (Senti- Analysis Dataset)

Output

```
text = "I'm trapped in a storm of emotions that I can't control, and it feels like no one understands the chaos inside me"
prompt = f"""Classify the text into Normal, Depression, Anxiety, and return the answer as the corresponding mental health disorder label.
text: {text}
label: """".strip()

pipe = pipeline(
    "text-generation",
    model=merged_model,
    tokenizer=tokenizer,
    device_map="auto",
)
outputs = pipe(prompt, max_new_tokens=2, do_sample=True, temperature=0.1)
print(outputs[0]["generated_text"].split("label: ")[-1].strip())
```

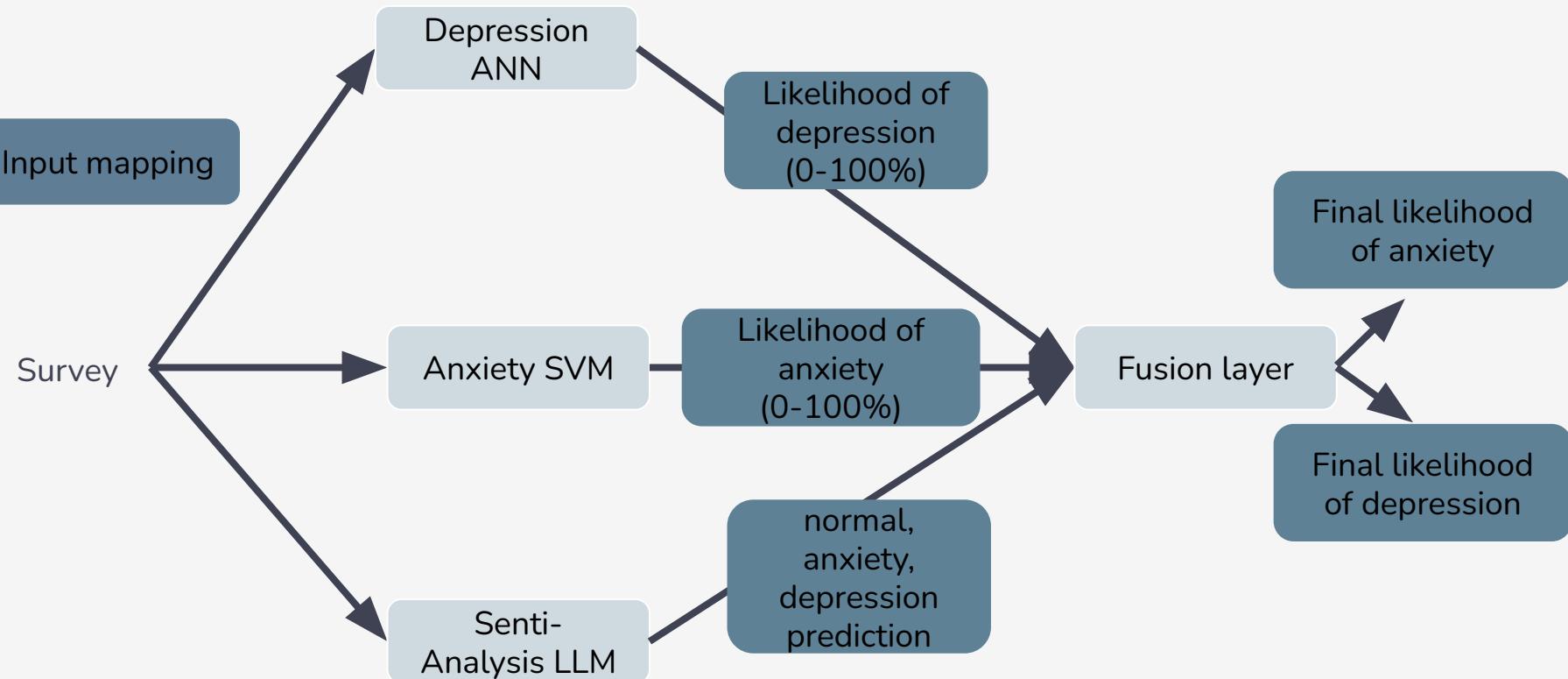
```
Device set to use cuda:0
The model 'PeftModelForCausalLM'
Depression
```

- Possible text input from a user and subsequent model output

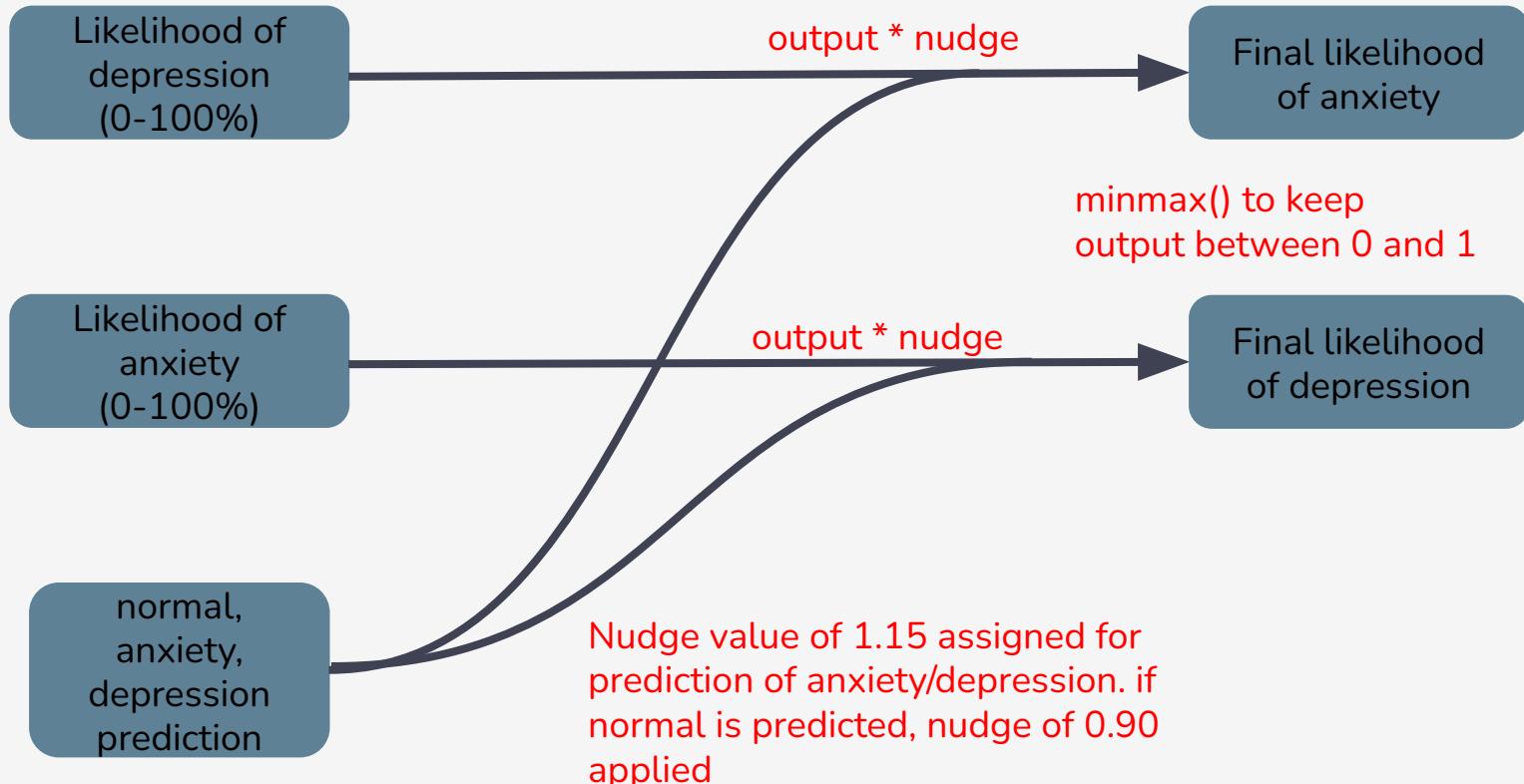


04 Recommendations

Recommendations: Output merging



Recommendations: fusion layer



Thank you!