




Capstone Project: NLP, Deep Learning, & Reddit

Eileen Hartnett
General Assembly


01

Introduction





Why look at Reddit text to study mental illness?





128



r/depression · Posted by u/Anders-94 8 hours ago 2

I've never been so depressed in my life.

I'm not making this because I need sympathy, I'm making this cause I think I need help.



3



r/AnorexiaNervosa · Posted by u/moonring_ 1 day ago

Really struggling with recovery and didn't know where else to turn

Recovery Related



Vote



r/schizophrenia · Posted by u/mowyourass 1 hour ago

I am experiencing almost all of these on the list, but I am scared of hospitals and my family's judgement. How do you know that this is the right decision?

Need Support

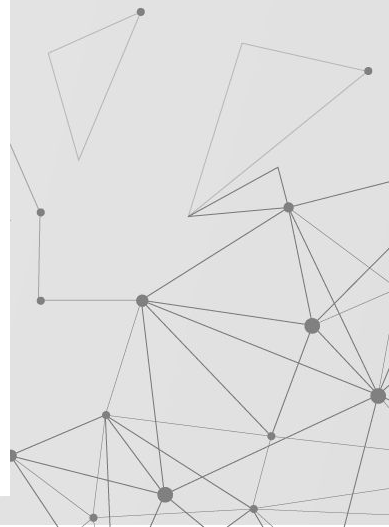


OPEN

A deep learning model for detecting mental illness from user content on social media

Jina Kim^{1,2}, Jieon Lee¹, Eunil Park^{1,3}✉ & Jinyoung Han³✉

Users of social media often share their feelings or emotional states through their posts. In this study, we developed a deep learning model to identify a user's mental state based on his/her posting information. To this end, we collected posts from mental health communities in *Reddit*. By analyzing and learning posting information written by users, our proposed model could accurately identify whether a user's post belongs to a specific mental disorder, including depression, anxiety, bipolar, borderline personality disorder, schizophrenia, and autism. We believe our model can help identify potential sufferers with mental illness based on their posts. This study further discusses the implication of our proposed model, which can serve as a supplementary tool for monitoring mental health states of individuals who frequently use social media.





INTRODUCTION

Motivation for this project.

01

PREPROCESSING

Web-scraping, cleaning, and pre-processing

02

ANALYSIS 1: XGBoost

Overview of the first angle of analysis using NLTK and XGBoost.

03

TABLE OF CONTENTS

04


ANALYSIS 2: CNN

Description of using Word2Vec and multiple CNNs.

05

CONCLUSIONS

A summary of what I've learned from taking this approach.

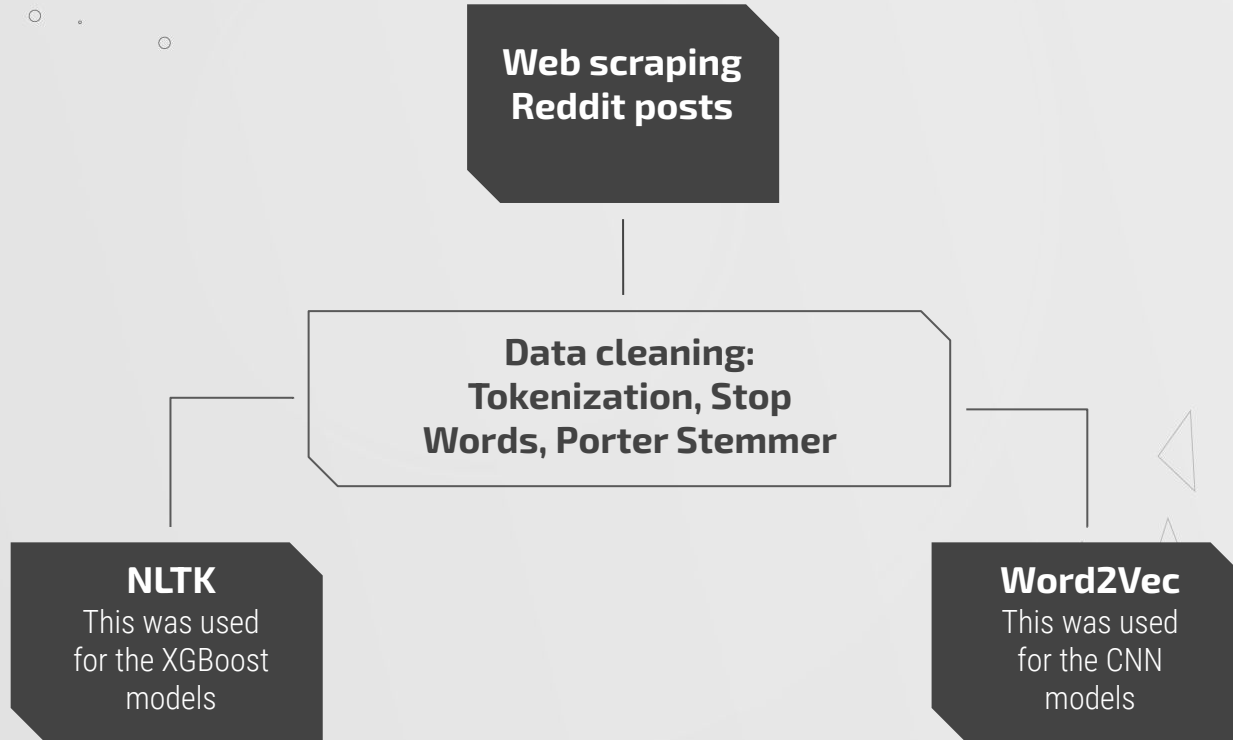




02

PREPROCESSING

Preprocessing steps



Pre-Processing in depth



Scrape Reddit Posts

r/mentalhealth
r/depression
r/Anxiety
/bipolar
r/BPD
r/schizophrenia
r/autism
r/AnorexiaNervosa
r/Bulimia



Clean Text

- Remove punctuation
- Stop words
- Numbers
- non-English characters
- Tokenize
- Porter-Stem



EDA

- 84,879 posts
- Roughly even amounts across subreddits
- "PreCovid" shutdown: 41,955
- "Post-Covid" shutdown: 42,924



NLTK



Word2Vec



03

Analysis 1: XGBoost

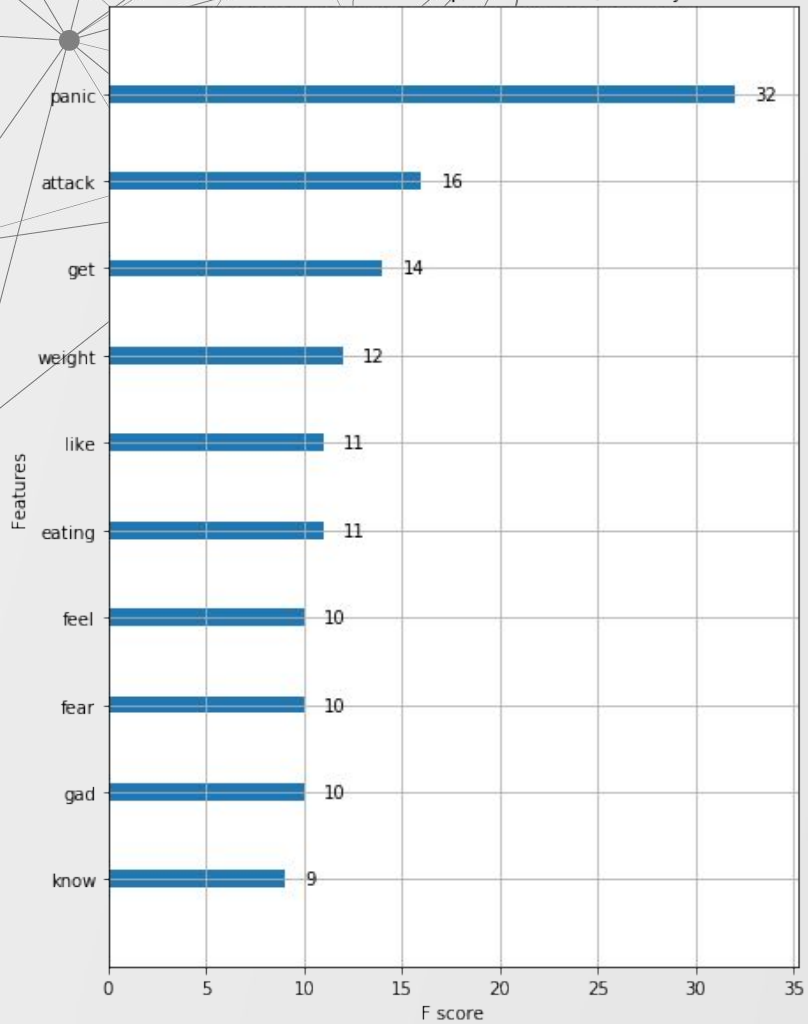
Kim et al. XGBoost Results

Channel	Class	F1-Score	Accuracy
r/depression	Depression	58.02	71.69
	Non-depression	78.65	
r/Anxiety	Anxiety	55.92	70.41
	Non-anxiety	77.73	
r/bipolar	Bipolar	53.59	85.53
	Non-bipolar	91.43	
r/BPD	BPD	46.43	85.14
	Non-BPD	91.37	
r/schizophrenia	Schizophrenia	40.97	86.72
	Non-schizophrenia	92.52	
r/autism	Autism	38.31	94.91
	Non-autism	97.35	

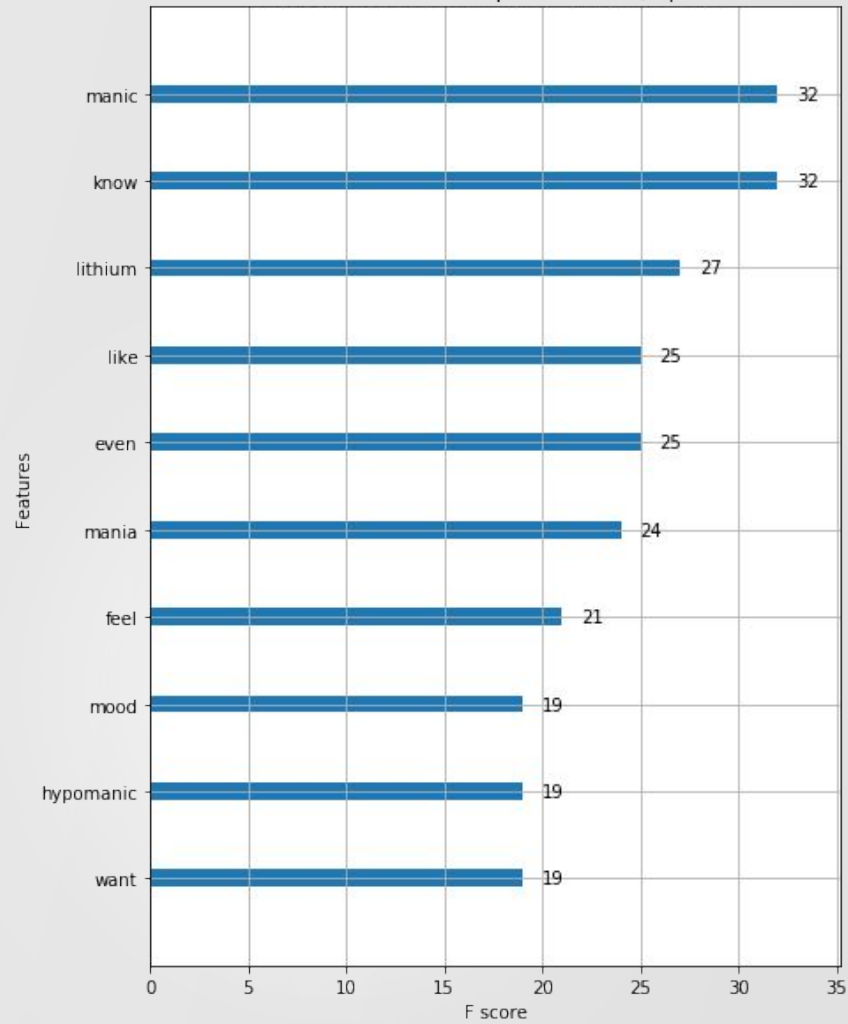
My XGBoost Results

Channel	Class	F1-Score	Accuracy
r/depression	Depression	23	74.24
	Non-depression	93	
r/Anxiety	Anxiety	17	65.52
	Non-anxiety	92	
r/bipolar	Bipolar	34	72.25
	Non-bipolar	82	
r/BPD	BPD	28	71.77
	Non-BPD	82	
r/schizophrenia	Schizophrenia	37	77.83
	Non-schizophrenia	87	
r/autism	Autism	48	83.15
	Non-autism	90	

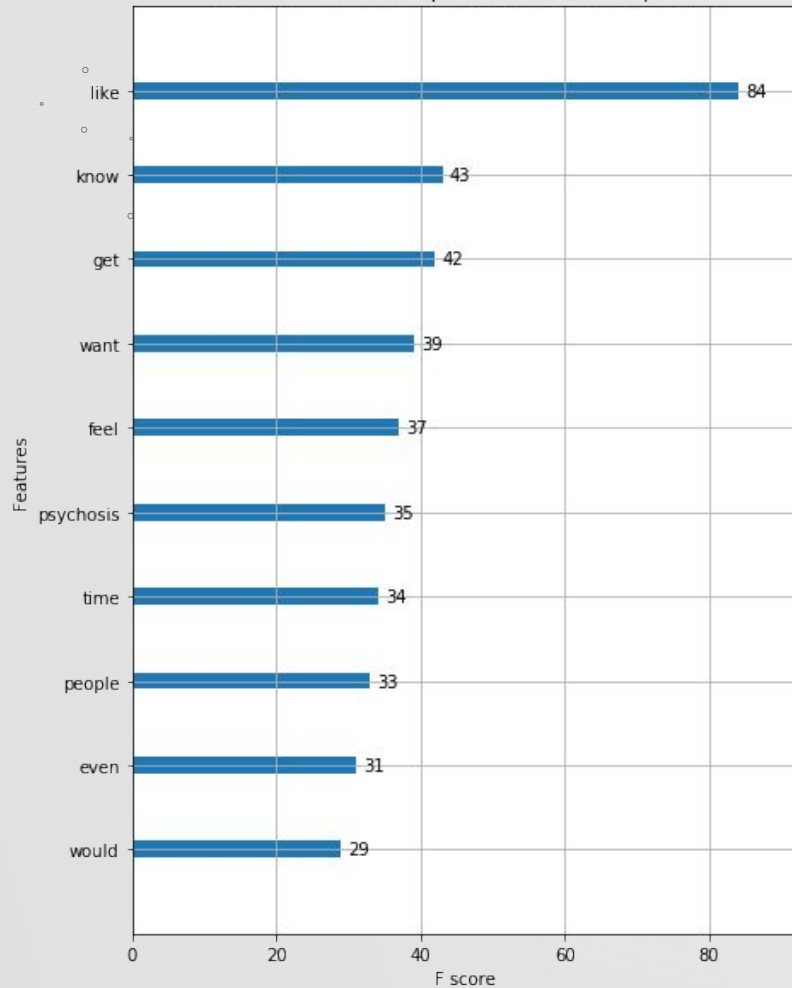
XGBOOST Feature Importance for r/Anxiety



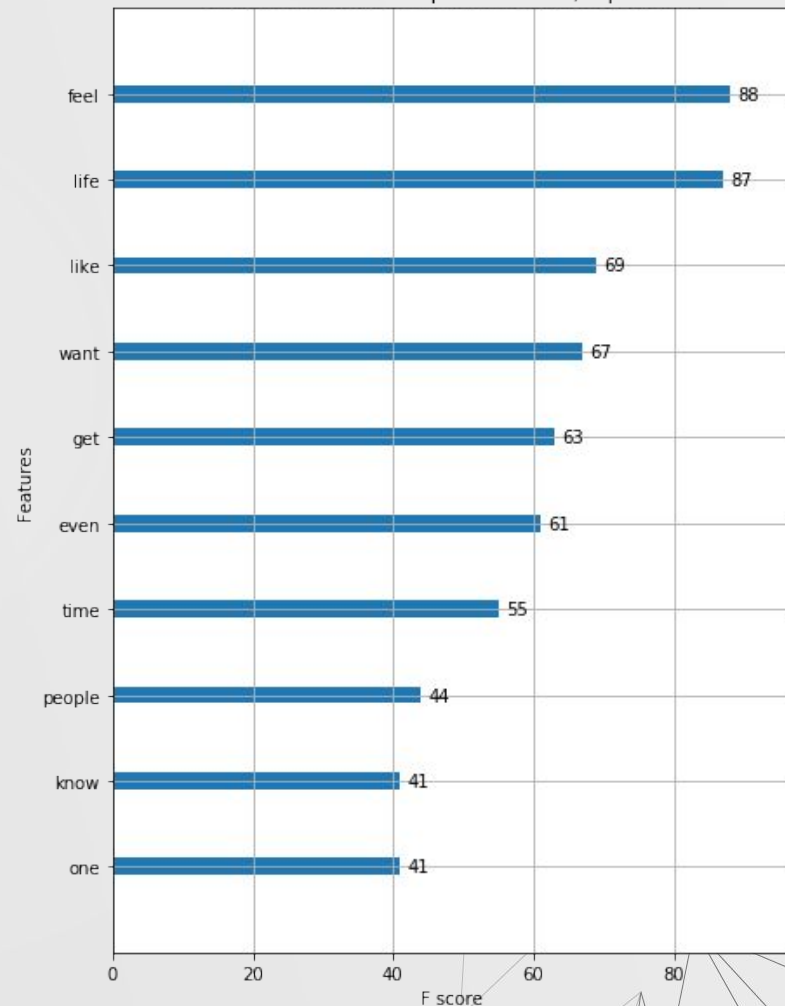
XGBOOST Feature Importance for r/bipolar



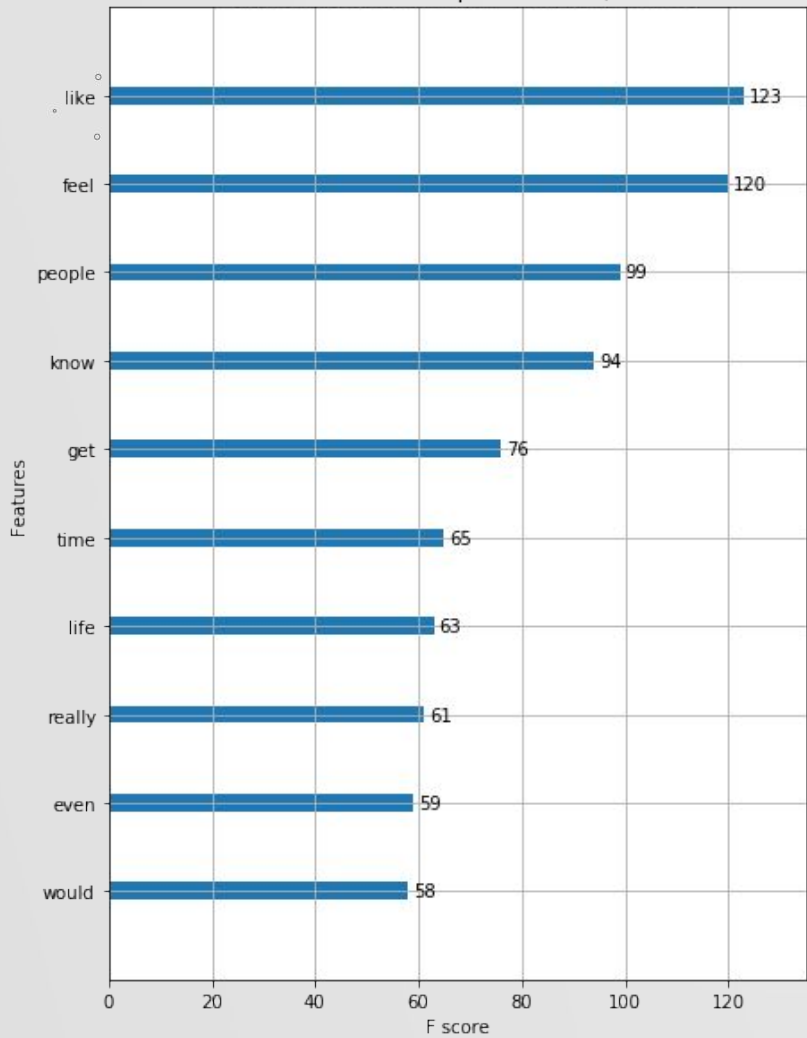
XGBOOST Feature Importance for r/schizophrenia



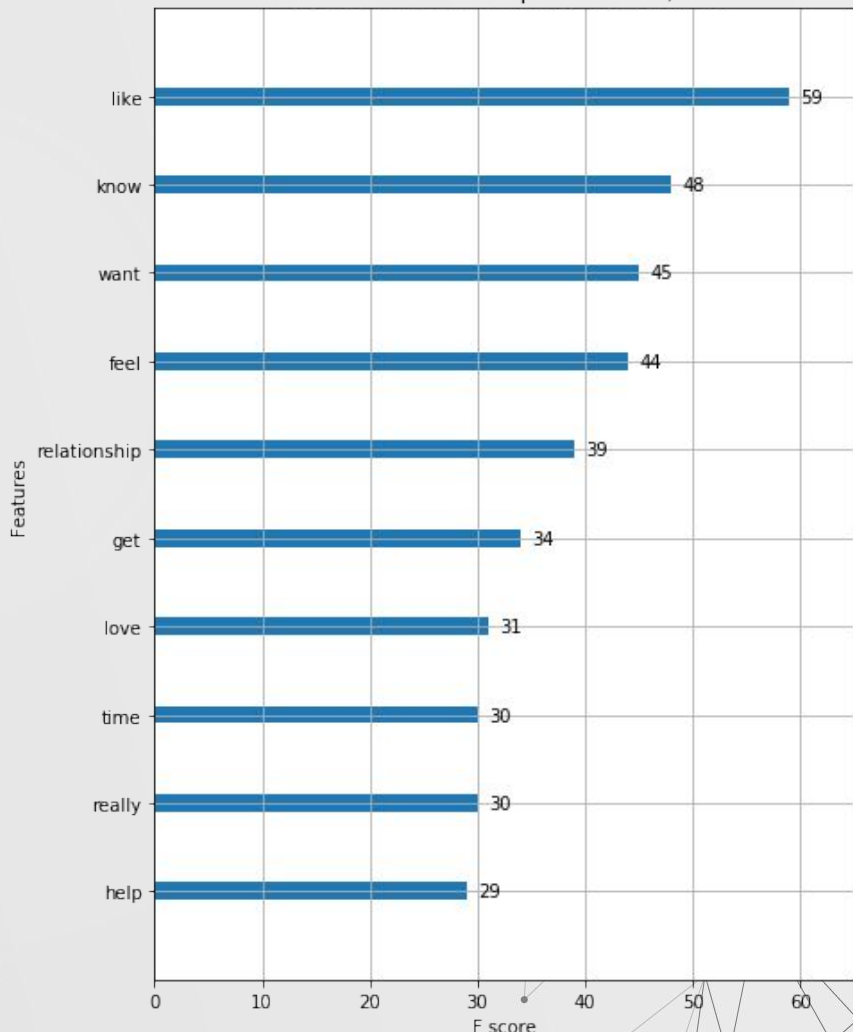
XGBOOST Feature Importance for r/depression



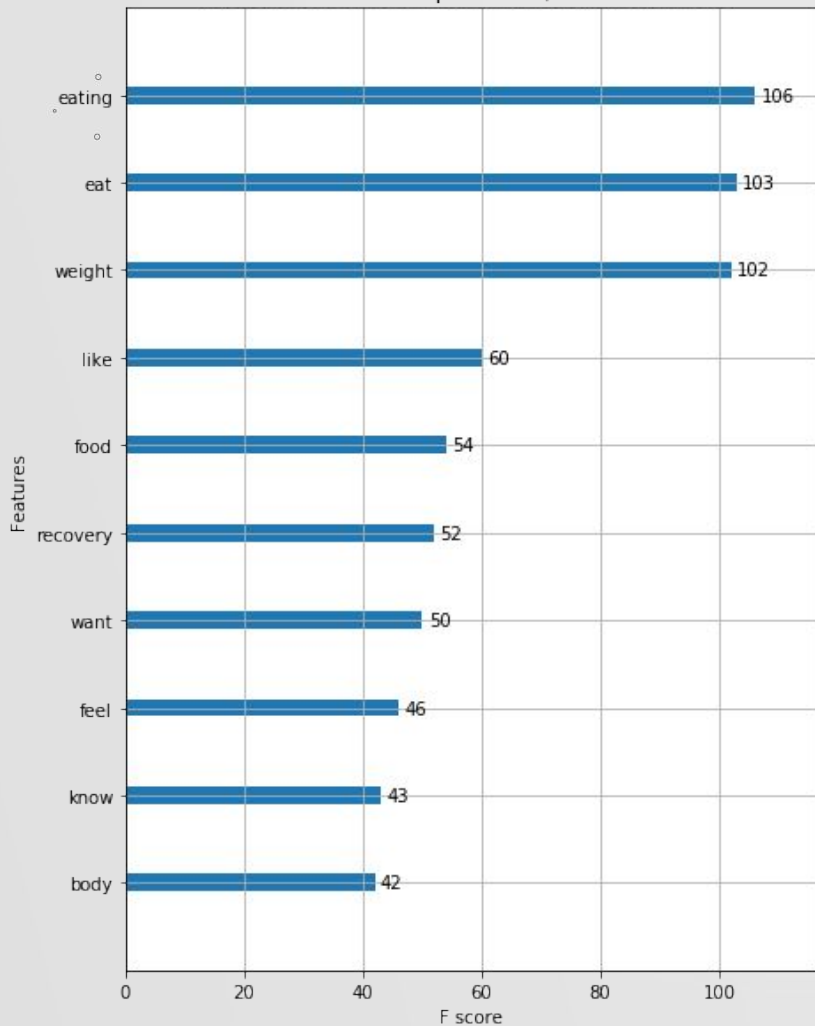
XGBOOST Feature Importance for r/Autism



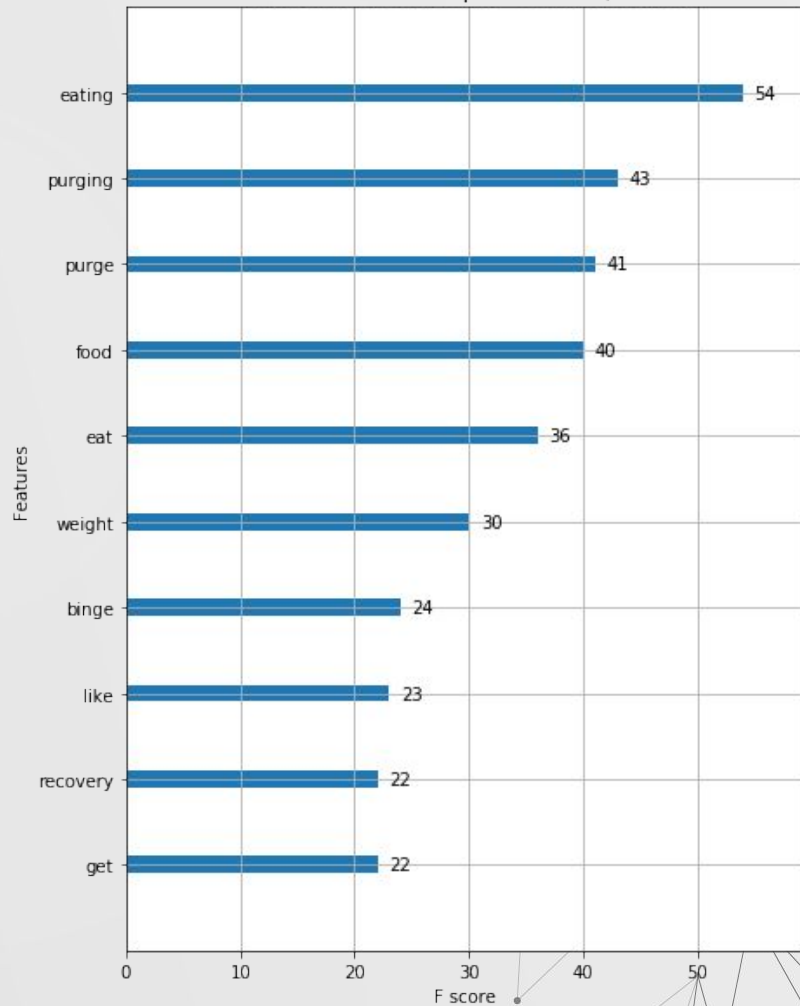
XGBOOST Feature Importance for r/BPD



XGBOOST Feature Importance r/AnorexiaNervosa



XGBOOST Feature Importance for r/bulimia

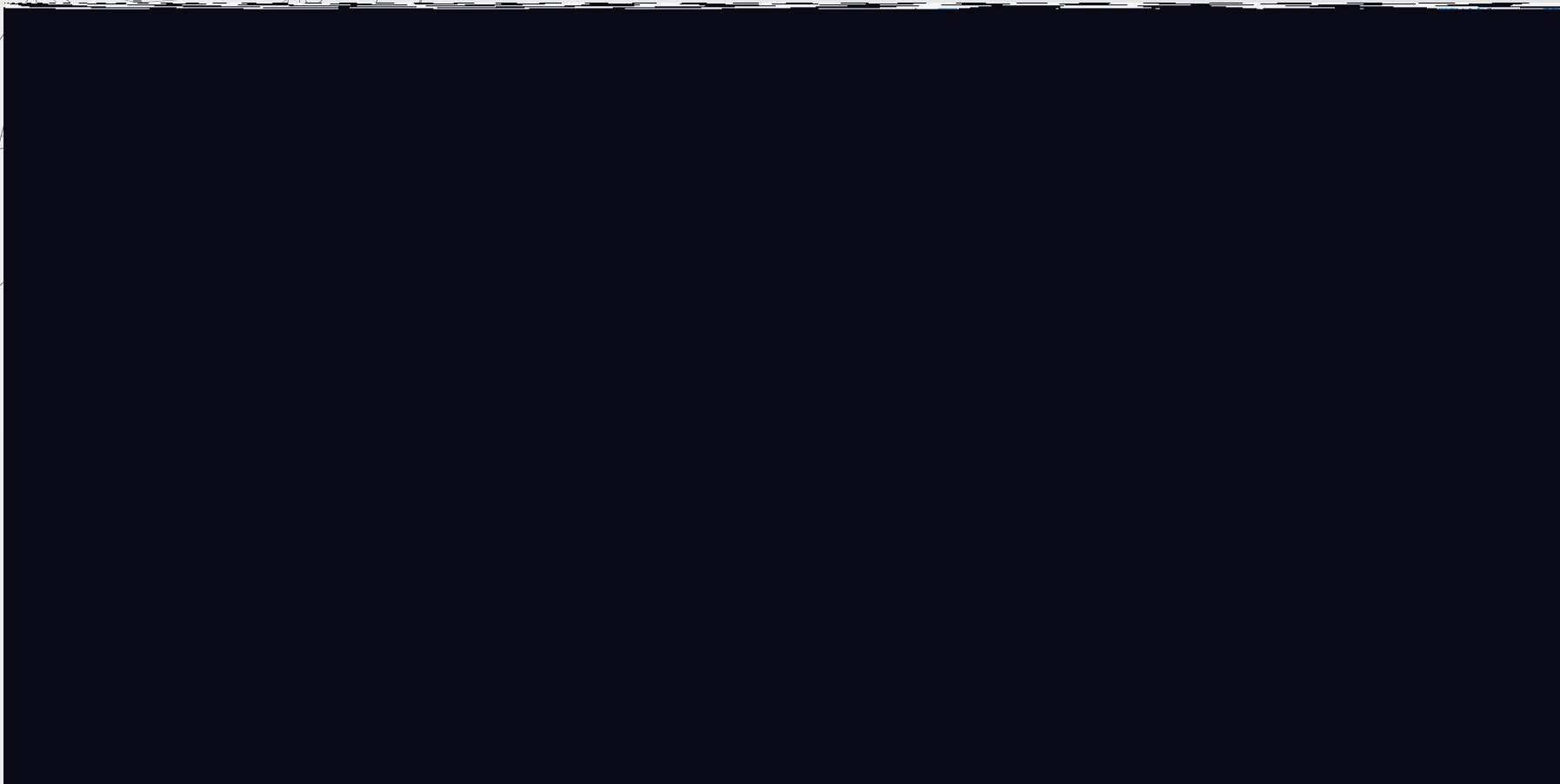


04

Analysis 2: Word2Vec & CNN



Word2Vec

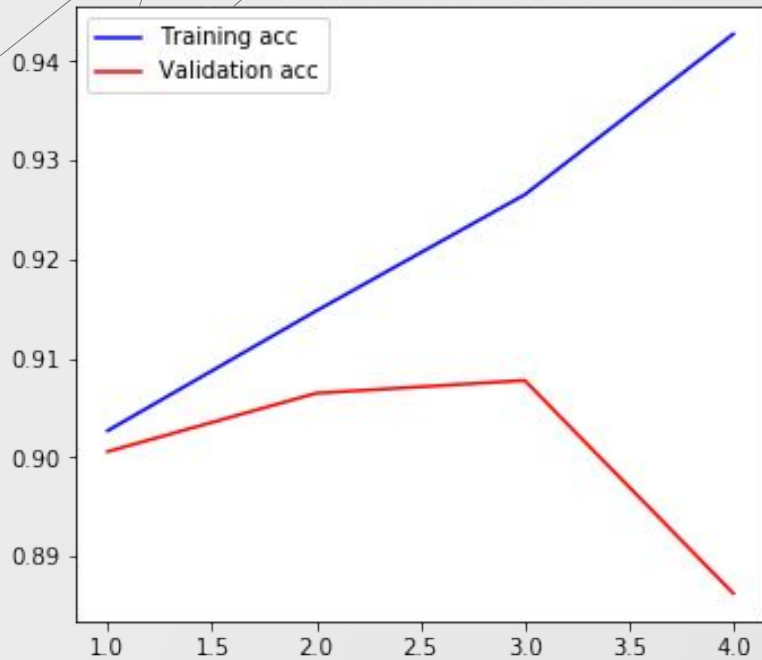


CNN Accuracy Results

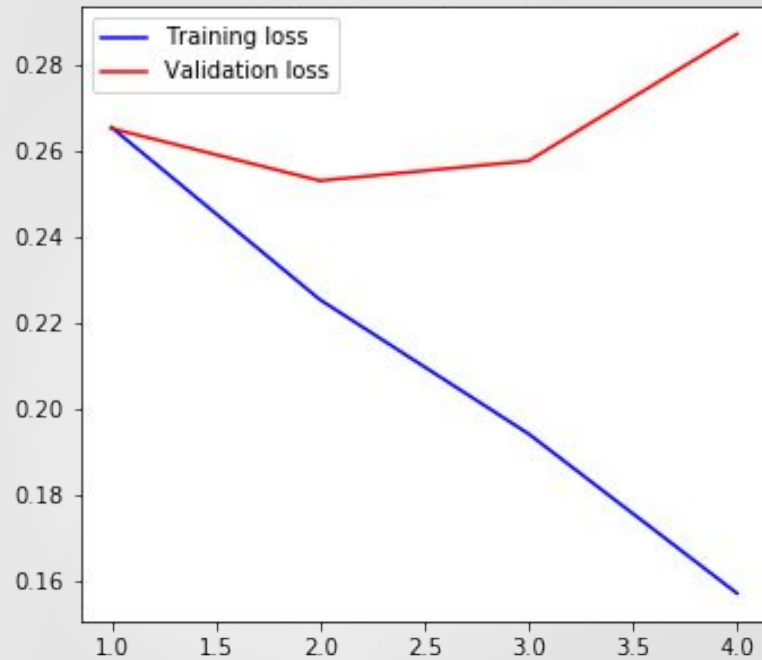
	r/Anxiety	r/Anorexia	r/Schizophrenia
Test Accuracy	82.76%	83.09%	79.14%
Train Accuracy	98.14%	98.15%	97.53%

r/Anxiety (classes imbalanced)

Training and validation accuracy

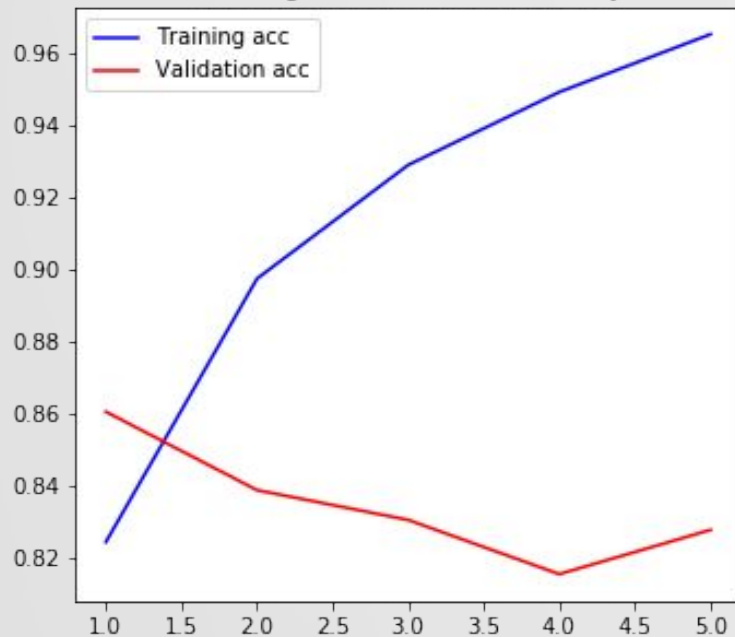


Training and validation loss

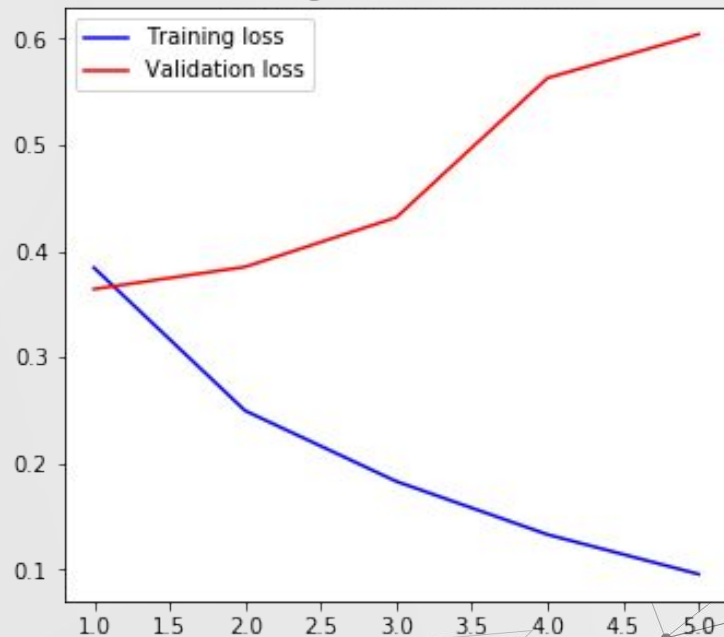


r/Anxiety (classes balanced with SMOTE)

Training and validation accuracy

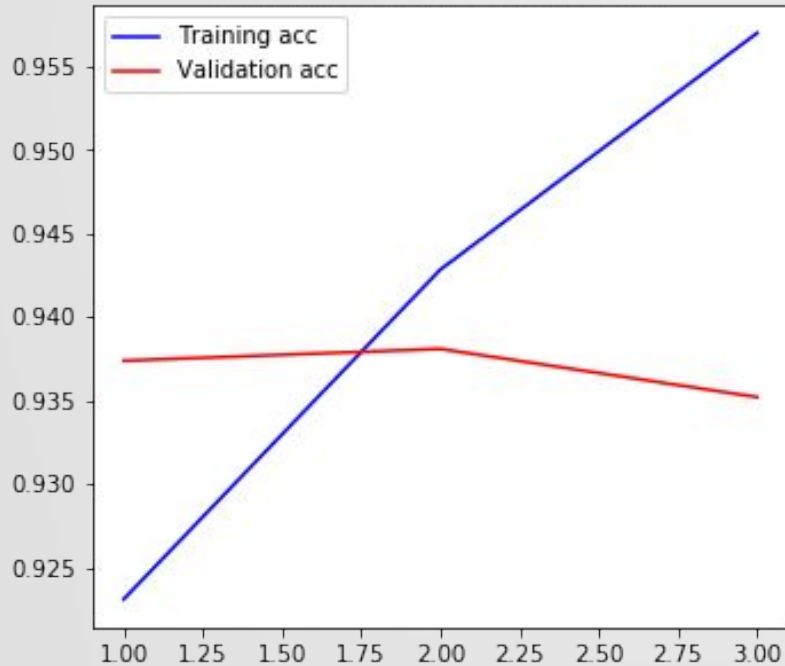


Training and validation loss

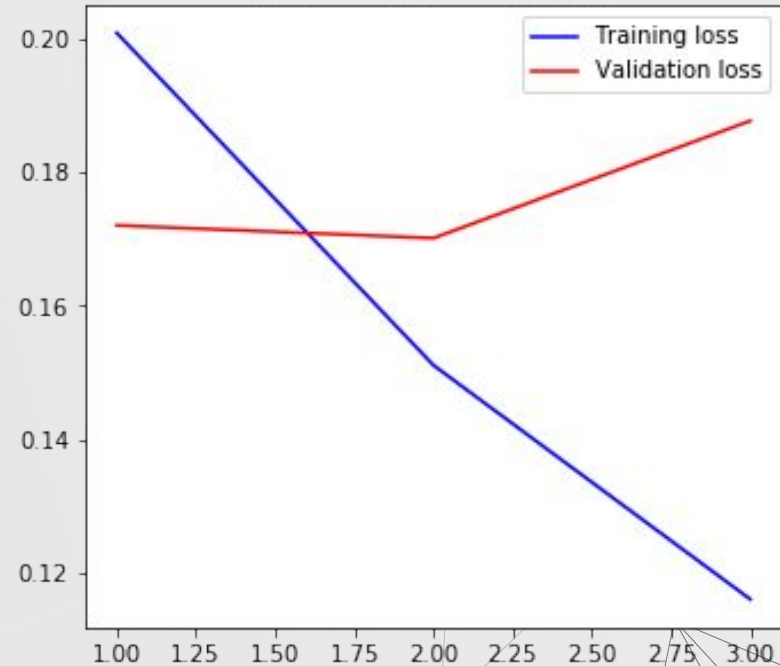


r/Schizophrenia (classes imbalanced)

Training and validation accuracy

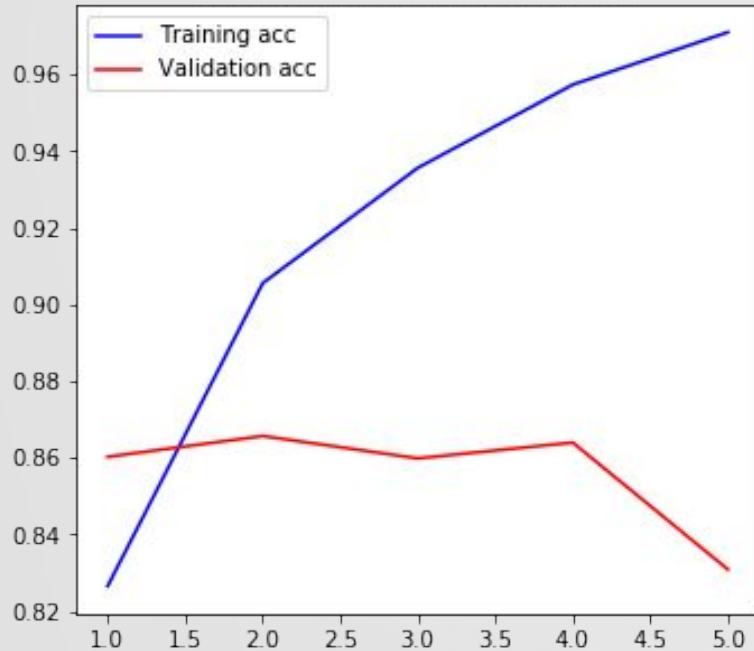


Training and validation loss

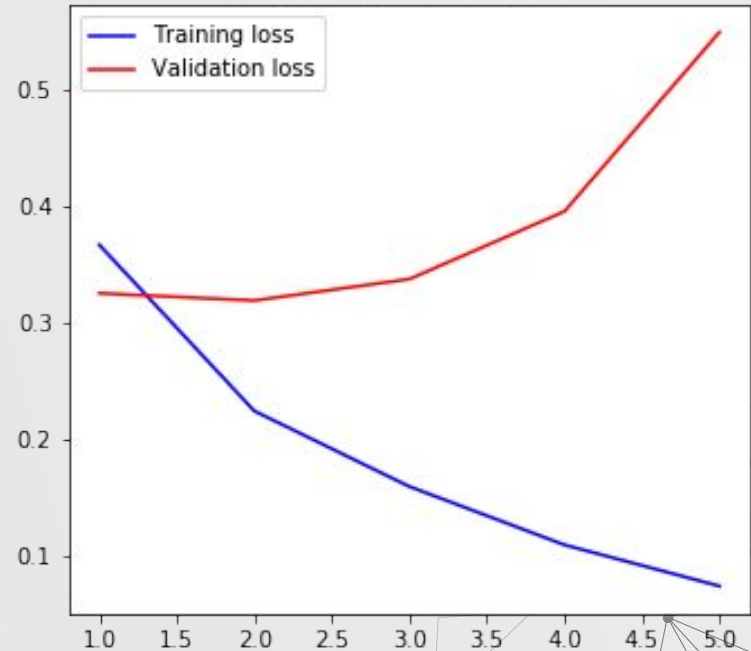


r/Schizophrenia (classes balanced with SMOTE)

Training and validation accuracy

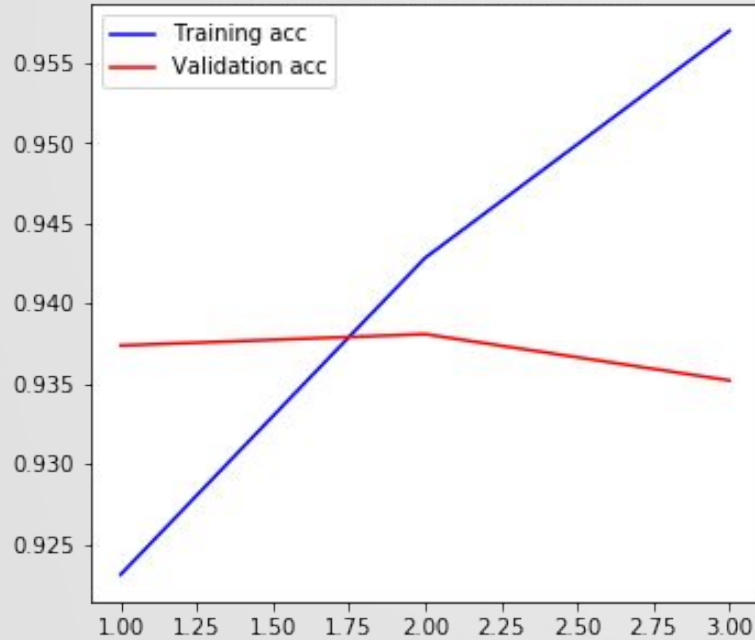


Training and validation loss

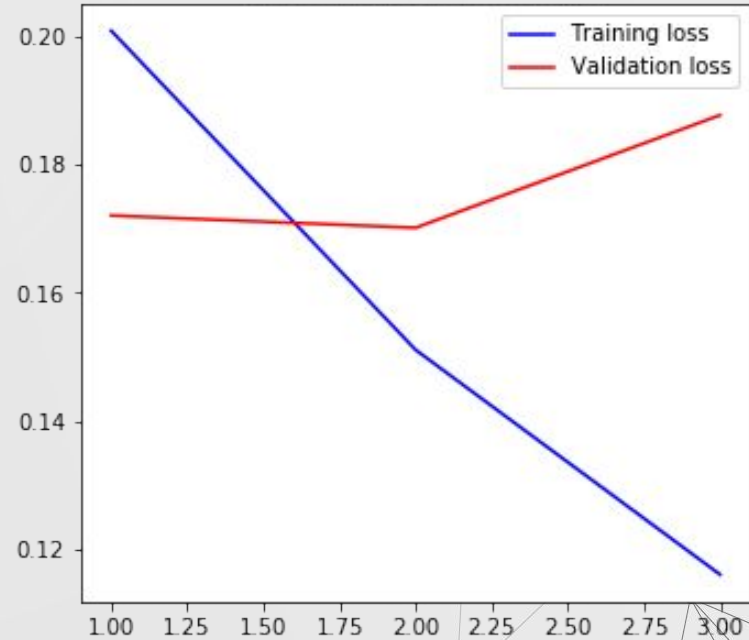


r/AnorexiaNervosa (classes imbalanced)

Training and validation accuracy

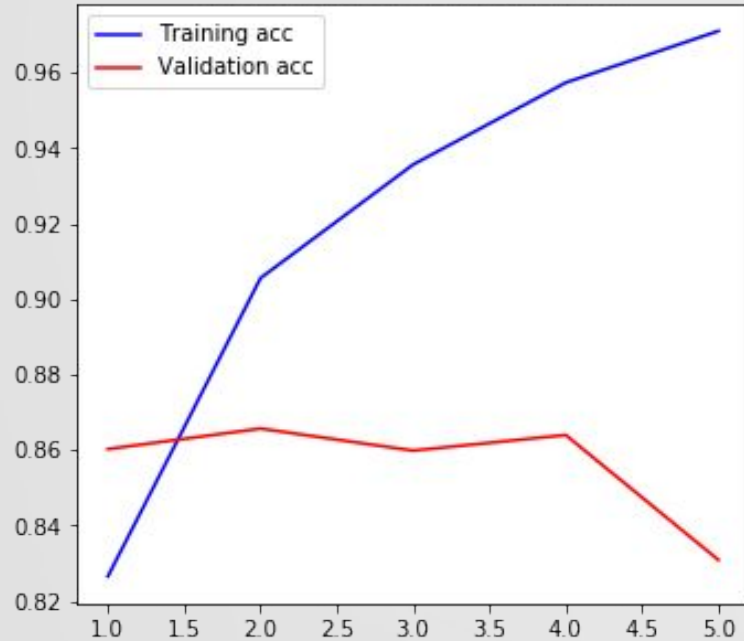


Training and validation loss

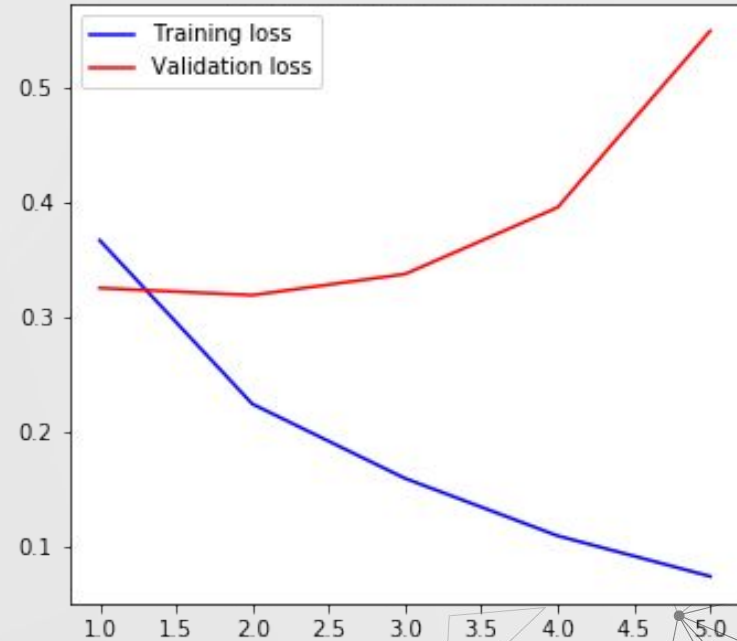


r/AnorexiaNervosa (classes balanced with SMOTE)

Training and validation accuracy



Training and validation loss



05

Conclusions



- We can obtain high accuracy scores with our neural networks, however more tuning is needed to ensure we are not overfitting.
- Read papers with a healthy dose of skepticism.
- This is an important research area and more work is needed on the topic.



The background of the slide features a complex, abstract geometric pattern. It consists of numerous thin, light gray lines that connect various points, creating a network-like structure. Some of these points are represented by small, solid dark gray circles. The overall effect is a modern, tech-inspired aesthetic. The text "THANK YOU!" is centered in the middle of the slide, rendered in a large, bold, dark gray sans-serif font.

THANK YOU!

CREDITS: This presentation template was created by **Slidesgo**, including icons by **Flaticon**, and infographics & images by **Freepik**.

Please keep this slide for attribution.

Questions?

