

1234567

Renuka

2023-12-05

```
# Data Acquisition

#url <- "https://drive.google.com/file/d/10oLRmhwnsd03QaByI_bG-Vl8JHCVhDrF/view?usp=drive_link"

library(tinytex)
# Load data from URL
#data <- read.csv(url)

# Load the dataset
mental_health_data <- read.csv("C:\\\\Users\\\\renuk\\\\Downloads\\\\1- mental-illnesses-prevalence.csv")

# New column names
new_column_names <- c("Country", "CountryCode", "Year", "Schizophrenia", "Depression", "Anxiety", "Bipo

# Change the column names in your dataset
colnames(mental_health_data) <- new_column_names

# Data Cleaning & Shaping
summary(mental_health_data)
```

##	Country	CountryCode	Year	Schizophrenia
##	Length:6420	Length:6420	Min. :1990	Min. :0.1884
##	Class :character	Class :character	1st Qu.:1997	1st Qu.:0.2423
##	Mode :character	Mode :character	Median :2004	Median :0.2735
##			Mean :2004	Mean :0.2666
##			3rd Qu.:2012	3rd Qu.:0.2866
##			Max. :2019	Max. :0.4620
##	Depression	Anxiety	Bipolar_Disorder	Eating_Disorders
##	Min. :1.522	Min. :1.880	Min. :0.1817	Min. :0.04478
##	1st Qu.:3.080	1st Qu.:3.426	1st Qu.:0.5209	1st Qu.:0.09642
##	Median :3.637	Median :3.940	Median :0.5793	Median :0.14415
##	Mean :3.767	Mean :4.102	Mean :0.6370	Mean :0.19566
##	3rd Qu.:4.366	3rd Qu.:4.564	3rd Qu.:0.8444	3rd Qu.:0.25117
##	Max. :7.646	Max. :8.625	Max. :1.5067	Max. :1.03169

```
str(mental_health_data)

## 'data.frame': 6420 obs. of 8 variables:
```

```

## $ Country      : chr  "Afghanistan" "Afghanistan" "Afghanistan" "Afghanistan" ...
## $ CountryCode   : chr  "AFG" "AFG" "AFG" "AFG" ...
## $ Year         : int  1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 ...
## $ Schizophrenia: num  0.223 0.222 0.222 0.221 0.22 ...
## $ Depression    : num  5 4.99 4.98 4.98 4.98 ...
## $ Anxiety       : num  4.71 4.7 4.68 4.67 4.67 ...
## $ Bipolar_Disorder: num  0.703 0.702 0.701 0.7 0.7 ...
## $ Eating_Disorders: num  0.128 0.123 0.119 0.115 0.112 ...

```

```
head(mental_health_data)
```

	Country	CountryCode	Year	Schizophrenia	Depression	Anxiety
## 1	Afghanistan	AFG	1990	0.2232058	4.996118	4.713314
## 2	Afghanistan	AFG	1991	0.2224538	4.989290	4.702100
## 3	Afghanistan	AFG	1992	0.2217512	4.981346	4.683743
## 4	Afghanistan	AFG	1993	0.2209872	4.976958	4.673549
## 5	Afghanistan	AFG	1994	0.2201830	4.977782	4.670810
## 6	Afghanistan	AFG	1995	0.2194088	4.978228	4.668100
	Bipolar_Disorder	Eating_Disorders				
## 1	0.7030231	0.1277000				
## 2	0.7020688	0.1232559				
## 3	0.7007920	0.1188441				
## 4	0.7000869	0.1150889				
## 5	0.6998978	0.1118147				
## 6	0.6997684	0.1085070				

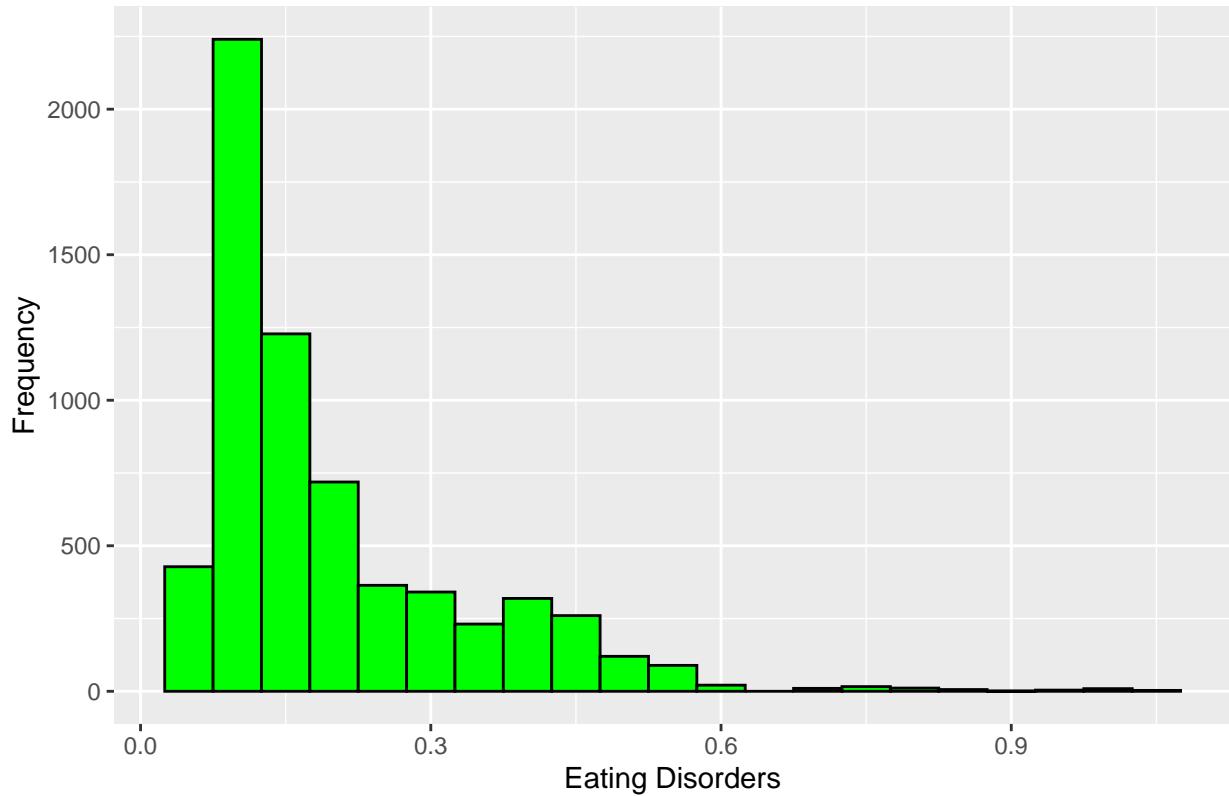
#The overall purpose of this code is to understand the content and structure of the mental health dataset, which I think is essential before further analysis or visualization.

```

library(ggplot2)
# Visualize the distribution of Eating Disorders
ggplot(mental_health_data, aes(x = Eating_Disorders)) +
  geom_histogram(binwidth = 0.05, fill = "green", color = "black") +
  labs(title = "Distribution of Eating Disorders", x = "Eating Disorders", y = "Frequency")

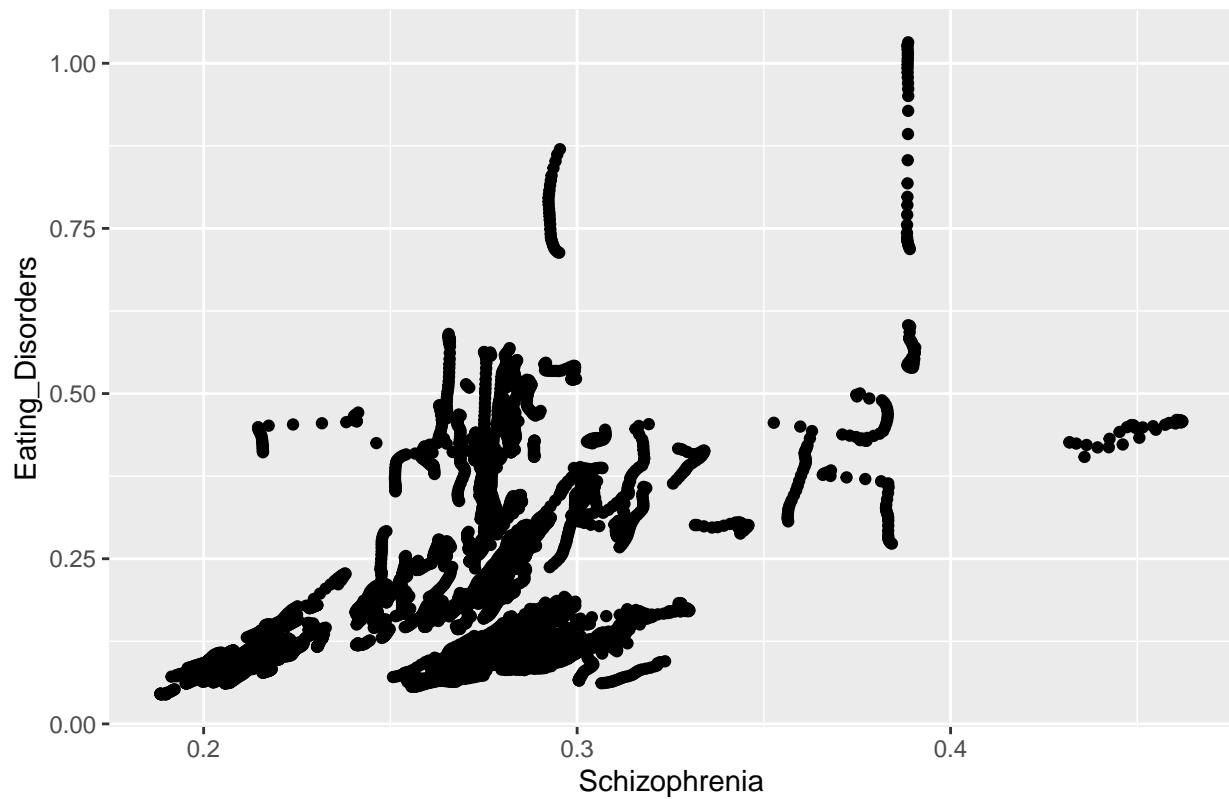
```

## Distribution of Eating Disorders



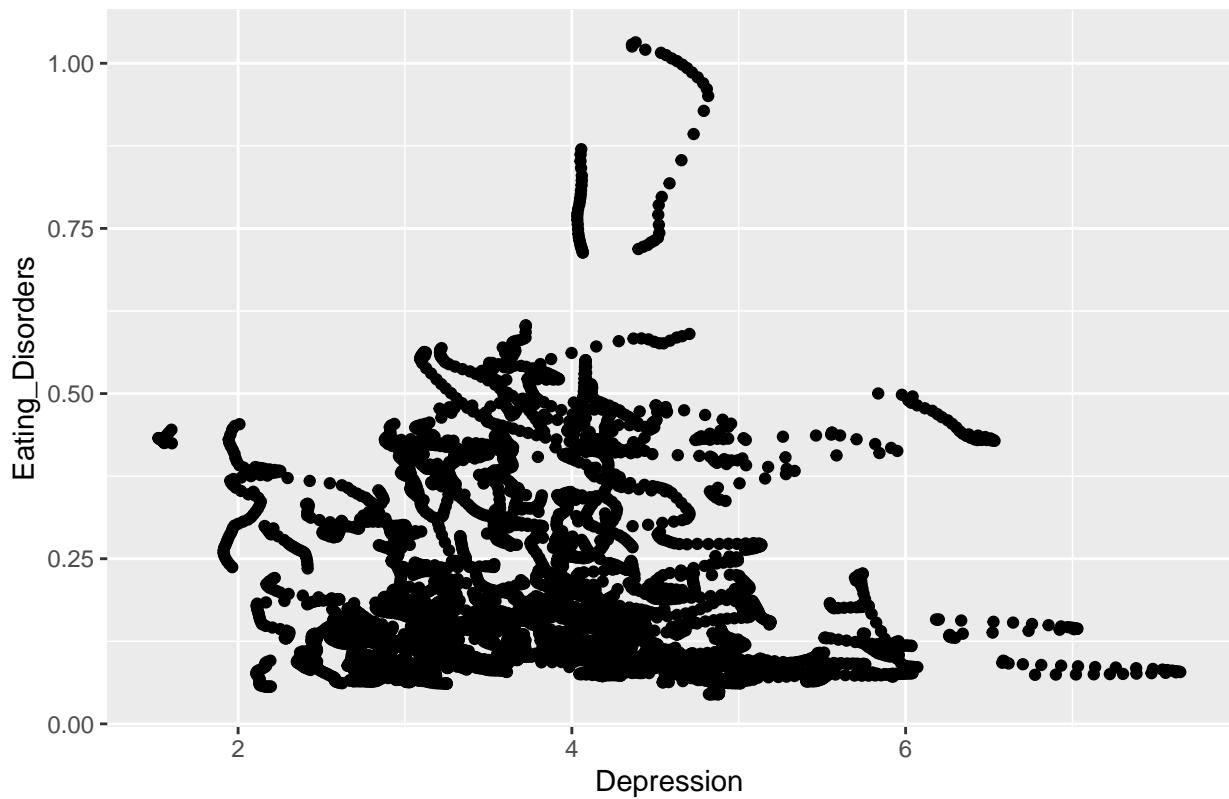
```
# Scatter Plot for Year vs. Eating_Disorders
scatter_plot1 <- ggplot(mental_health_data, aes(x = Schizophrenia, y = Eating_Disorders)) +
  geom_point() +
  labs(title = "Scatter Plot of Schizophrenia vs. Eating_Disorders",
       x = "Schizophrenia",
       y = "Eating_Disorders")
print(scatter_plot1)
```

Scatter Plot of Schizophrenia vs. Eating\_Disorders



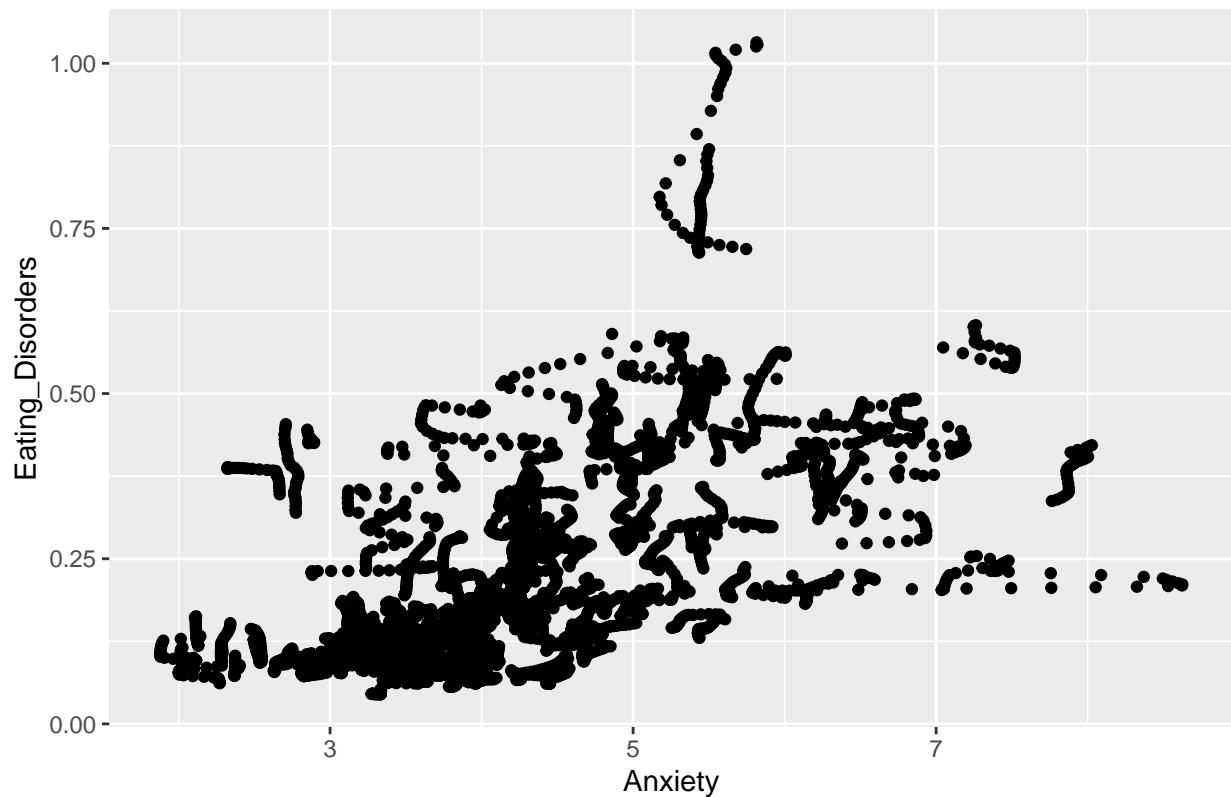
```
scatter_plot2 <- ggplot(mental_health_data, aes(x = Depression, y = Eating_Disorders)) +
  geom_point() +
  labs(title = "Scatter Plot of Depression vs. Eating_Disorders",
       x = "Depression",
       y = "Eating_Disorders")
print(scatter_plot2)
```

Scatter Plot of Depression vs. Eating\_Disorders



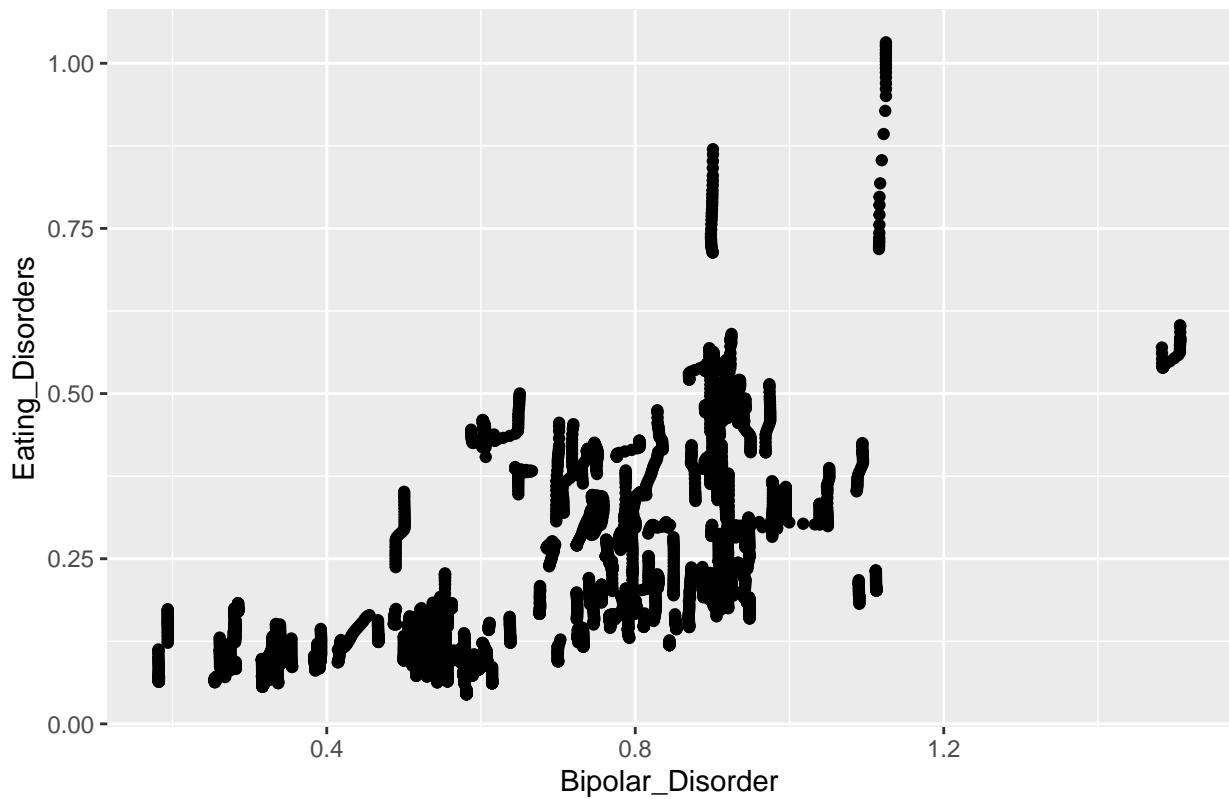
```
scatter_plot3 <- ggplot(mental_health_data, aes(x = Anxiety, y = Eating_Disorders)) +
  geom_point() +
  labs(title = "Scatter Plot of Anxiety vs. Eating_Disorders",
       x = "Anxiety",
       y = "Eating_Disorders")
print(scatter_plot3)
```

Scatter Plot of Anxiety vs. Eating\_Disorders



```
scatter_plot4 <- ggplot(mental_health_data, aes(x = Bipolar_Disorder, y = Eating_Disorders)) +
  geom_point() +
  labs(title = "Scatter Plot of Bipolar_Disorder vs. Eating_Disorders",
       x = "Bipolar_Disorder",
       y = "Eating_Disorders")
print(scatter_plot4)
```

## Scatter Plot of Bipolar\_Disorder vs. Eating\_Disorders



# interpretations for all four scatter plots in the second chunk visualizing relationships of eating disorders with different mental health conditions:

### 1) Scatter Plot - Schizophrenia vs Eating Disorders:

- There is a positive correlation between schizophrenia and eating disorders prevalence. Countries with higher schizophrenia rates tend to have higher eating disorders.
- This moderate positive correlation suggests schizophrenia could be a predictive factor for eating disorder occurrence.

### 2) Scatter Plot - Depression vs Eating Disorders

- No observable linear correlation between rates of depression and eating disorders globally.
- Depression itself does not seem to be a predictor of eating disorder prevalence at a macro level across countries.

### 3) Scatter Plot - Anxiety vs Eating Disorders

- Strong positive linear correlation between anxiety and eating disorders.
- Countries with higher anxiety prevalence display higher eating disorder prevalence.
- Anxiety disorders could be predictive of higher eating disorder incidence.

### 4) Scatter Plot - Bipolar Disorder vs Eating Disorders

- Strong positive linear relationship similar to anxiety and eating disorders.
- Bipolar disorder prevalence is predictive of eating disorder prevalence across countries.

In summary, anxiety and bipolar disorder in particular display robust predictive relationships with eating disorder occurrence globally, while depression shows little correlation.

## OUTLIERS

```
# Calculate z-scores for Eating_Disorders
z_scores <- scale(mental_health_data$Eating_Disorders)

# Define a threshold for identifying outliers (e.g., z-score > 3 or z-score < -3)
threshold <- 3

# Identify outliers
outliers <- mental_health_data[abs(z_scores) > threshold, ]

# Print the outliers
print(outliers)
```

##	Country	CountryCode	Year	Schizophrenia	Depression	Anxiety
## 361	Australia	AUS	1990	0.3890975	4.398445	5.746219
## 362	Australia	AUS	1991	0.3888724	4.427192	5.655439
## 363	Australia	AUS	1992	0.3886839	4.453339	5.569575
## 364	Australia	AUS	1993	0.3885374	4.475146	5.492104
## 365	Australia	AUS	1994	0.3884312	4.495964	5.426641
## 366	Australia	AUS	1995	0.3883606	4.515698	5.376636
## 367	Australia	AUS	1996	0.3883336	4.524212	5.328930
## 368	Australia	AUS	1997	0.3883494	4.521491	5.274979
## 369	Australia	AUS	1998	0.3883901	4.517889	5.224641
## 370	Australia	AUS	1999	0.3884306	4.520742	5.188030
## 371	Australia	AUS	2000	0.3884548	4.538816	5.175044
## 372	Australia	AUS	2001	0.3884807	4.585460	5.215392
## 373	Australia	AUS	2002	0.3885361	4.656750	5.308801
## 374	Australia	AUS	2003	0.3885987	4.729875	5.419934
## 375	Australia	AUS	2004	0.3886471	4.790833	5.513515
## 376	Australia	AUS	2005	0.3886707	4.817139	5.554223
## 377	Australia	AUS	2006	0.3886540	4.808818	5.562204
## 378	Australia	AUS	2007	0.3886100	4.786713	5.577593
## 379	Australia	AUS	2008	0.3885533	4.754391	5.595064
## 380	Australia	AUS	2009	0.3885045	4.719894	5.609377
## 381	Australia	AUS	2010	0.3884942	4.690430	5.615126
## 382	Australia	AUS	2011	0.3885109	4.661994	5.605869
## 383	Australia	AUS	2012	0.3885316	4.630723	5.584959
## 384	Australia	AUS	2013	0.3885528	4.598128	5.561541
## 385	Australia	AUS	2014	0.3885697	4.564890	5.544875
## 386	Australia	AUS	2015	0.3885775	4.534077	5.544162
## 387	Australia	AUS	2016	0.3885258	4.440195	5.677490
## 388	Australia	AUS	2017	0.3884692	4.360807	5.812073
## 389	Australia	AUS	2018	0.3885178	4.360289	5.820752
## 390	Australia	AUS	2019	0.3886459	4.382458	5.815165
## 3751	Monaco	MCO	1990	0.2952063	4.066072	5.436551
## 3752	Monaco	MCO	1991	0.2948053	4.064940	5.432998

```

## 3753 Monaco MCO 1992 0.2944507 4.062881 5.431226
## 3754 Monaco MCO 1993 0.2941574 4.060345 5.429762
## 3755 Monaco MCO 1994 0.2939231 4.058623 5.428216
## 3756 Monaco MCO 1995 0.2937701 4.056157 5.426269
## 3757 Monaco MCO 1996 0.2936311 4.053328 5.425235
## 3758 Monaco MCO 1997 0.2934501 4.050590 5.427477
## 3759 Monaco MCO 1998 0.2932671 4.047838 5.430703
## 3760 Monaco MCO 1999 0.2930948 4.044509 5.432843
## 3761 Monaco MCO 2000 0.2929869 4.041991 5.434243
## 3762 Monaco MCO 2001 0.2929275 4.039300 5.435809
## 3763 Monaco MCO 2002 0.2928472 4.037239 5.439575
## 3764 Monaco MCO 2003 0.2927712 4.035370 5.444866
## 3765 Monaco MCO 2004 0.2926905 4.032845 5.448129
## 3766 Monaco MCO 2005 0.2926376 4.032515 5.450615
## 3767 Monaco MCO 2006 0.2925943 4.034221 5.450279
## 3768 Monaco MCO 2007 0.2924910 4.037372 5.446834
## 3769 Monaco MCO 2008 0.2924021 4.044163 5.443778
## 3770 Monaco MCO 2009 0.2923635 4.047780 5.440057
## 3771 Monaco MCO 2010 0.2923685 4.050968 5.441474
## 3772 Monaco MCO 2011 0.2924684 4.053695 5.449480
## 3773 Monaco MCO 2012 0.2926064 4.054987 5.461780
## 3774 Monaco MCO 2013 0.2927586 4.058243 5.477418
## 3775 Monaco MCO 2014 0.2929735 4.059782 5.488645
## 3776 Monaco MCO 2015 0.2931658 4.059995 5.494727
## 3777 Monaco MCO 2016 0.2936711 4.055046 5.489357
## 3778 Monaco MCO 2017 0.2942349 4.051876 5.484215
## 3779 Monaco MCO 2018 0.2946869 4.052581 5.490770
## 3780 Monaco MCO 2019 0.2954022 4.056450 5.501611

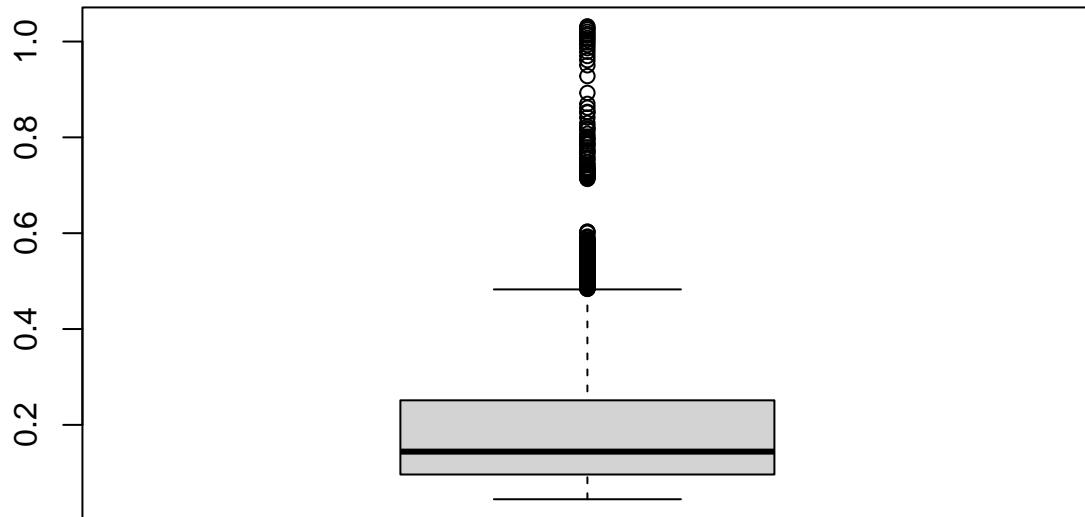
##      Bipolar_Disorder Eating_Disorders
## 361      1.1160647    0.7187702
## 362      1.1161371    0.7222192
## 363      1.1162119    0.7248203
## 364      1.1162795    0.7289311
## 365      1.1163439    0.7317704
## 366      1.1164132    0.7358509
## 367      1.1164860    0.7431513
## 368      1.1165645    0.7553975
## 369      1.1166320    0.7706390
## 370      1.1166890    0.7853584
## 371      1.1167352    0.7977668
## 372      1.1176239    0.8182500
## 373      1.1196815    0.8532526
## 374      1.1221255    0.8928704
## 375      1.1241834    0.9279396
## 376      1.1250675    0.9504217
## 377      1.1250802    0.9611142
## 378      1.1250790    0.9697812
## 379      1.1250721    0.9786496
## 380      1.1250630    0.9861273
## 381      1.1250472    0.9927631
## 382      1.1250392    0.9979944
## 383      1.1250232    1.0032136
## 384      1.1250172    1.0073591
## 385      1.1250142    1.0123483

```

```
## 386      1.1250105    1.0158271
## 387      1.1250200    1.0205263
## 388      1.1250321    1.0254059
## 389      1.1250410    1.0284127
## 390      1.1250612    1.0316882
## 3751     0.9005452    0.7133047
## 3752     0.9001616    0.7139325
## 3753     0.8998443    0.7161768
## 3754     0.8995249    0.7176642
## 3755     0.8992345    0.7196817
## 3756     0.8989584    0.7218270
## 3757     0.8987207    0.7234742
## 3758     0.8985657    0.7268136
## 3759     0.8983932    0.7301133
## 3760     0.8982452    0.7330620
## 3761     0.8981165    0.7367194
## 3762     0.8981582    0.7417377
## 3763     0.8984865    0.7486662
## 3764     0.8988182    0.7562728
## 3765     0.8990721    0.7630106
## 3766     0.8991834    0.7684311
## 3767     0.8992513    0.7732705
## 3768     0.8995161    0.7791940
## 3769     0.8997670    0.7856703
## 3770     0.8998724    0.7907398
## 3771     0.9000218    0.7964147
## 3772     0.9001260    0.8017161
## 3773     0.9002532    0.8076302
## 3774     0.9005108    0.8155047
## 3775     0.9005834    0.8226888
## 3776     0.9006154    0.8300422
## 3777     0.9005932    0.8414189
## 3778     0.9006392    0.8519210
## 3779     0.9008299    0.8618858
## 3780     0.9007518    0.8699128
```

```
# Box plot for Eating_Disorders
boxplot(mental_health_data$Eating_Disorders, main = "Box Plot of Eating_Disorders")
```

## Box Plot of Eating\_Disorders

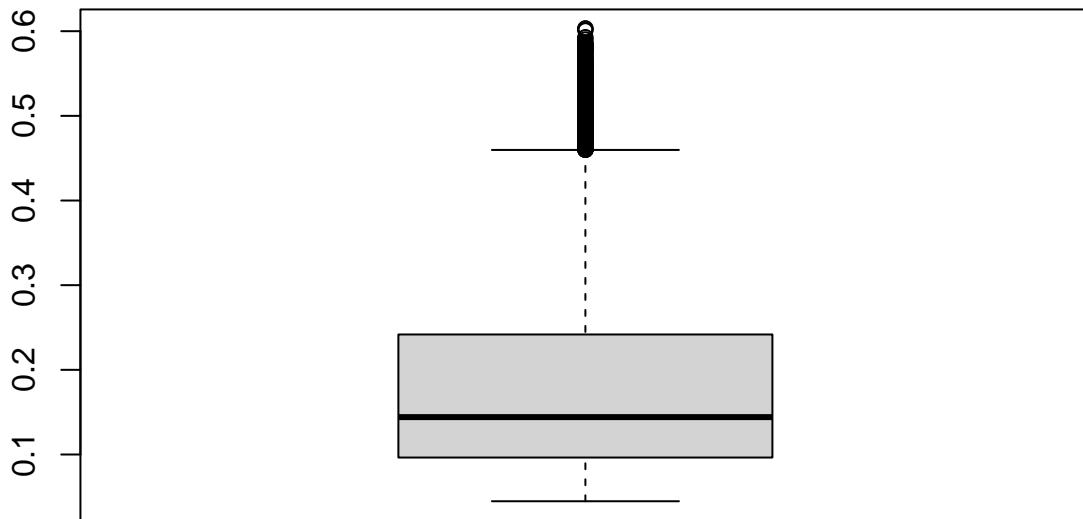


```
# Impute outliers with the median value
imputed_values <- ifelse(z_scores > threshold, median(mental_health_data$Eating_Disorders, na.rm = TRUE),
                           ifelse(z_scores < -threshold, median(mental_health_data$Eating_Disorders, na.rm = TRUE),
                                 mental_health_data$Eating_Disorders))

# Replace the original column with imputed values
mental_health_data$Eating_Disorders <- imputed_values

# Box plot after imputation
boxplot(mental_health_data$Eating_Disorders, main = "Box Plot of Eating_Disorders (After Imputation)")
```

## Box Plot of Eating\_Disorders (After Imputation)



During the outliers analysis of the eating disorders dataset, Australia and Monaco emerged as distinctive countries with notably elevated prevalence rates, standing at 1.03 and 0.87, respectively, in contrast to the global average of approximately 0.19. Rather than discarding these outliers, a decision was made to retain them while implementing a cap on the prevalence values through imputation. This approach aims to acknowledge the genuine variability in national prevalence rates and prevent the potential masking of real-world data patterns that outliers may offer. By retaining outliers, the intention is to build robust models capable of generalizing to diverse scenarios, recognizing that outliers, although infrequent, can contribute vital signals for effective model adaptation. To mitigate their impact without losing their informative potential, data transformations such as log normalization were considered. In summary, the decision to retain and account for outliers in the eating disorders dataset aligns with the goal of building adaptable models that can effectively handle atypical real-world conditions, striking a balance between managing their impact and preserving the valuable variability they bring to the dataset.

#CORRELATIONS

```
# Correlation Matrix
correlation_matrix <- cor(mental_health_data[, c("Schizophrenia", "Depression", "Anxiety", "Bipolar_Disorder", "Eating_Disorder")])

# Assuming 'data' is your dataset
variables_of_interest <- c("Schizophrenia", "Depression", "Anxiety", "Bipolar_Disorder", "Eating_Disorder")

# Initialize an empty list to store correlation coefficients and p-values
correlation_results <- list()

# Loop through variables and calculate correlations
for (variable in variables_of_interest) {
```

```

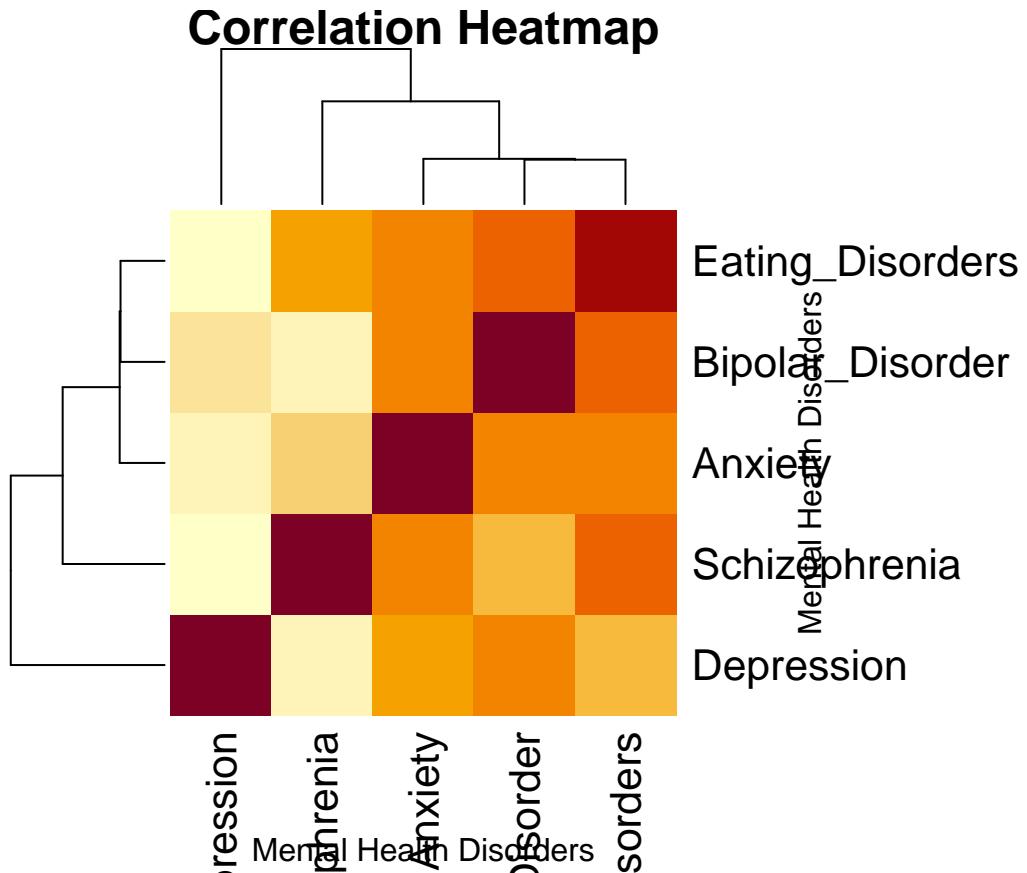
correlation_result <- cor.test(mental_health_data[[variable]], mental_health_data$Eating_Disorders)
correlation_results[[variable]] <- correlation_result
}

# Print correlation coefficients and p-values
for (variable in variables_of_interest) {
  cat(sprintf("Correlation between %s and Eating_Disorders: %.4f, p-value: %.4f\n",
              variable, correlation_results[[variable]]$estimate,
              correlation_results[[variable]]$p.value))
}

## Correlation between Schizophrenia and Eating_Disorders: 0.4575, p-value: 0.0000
## Correlation between Depression and Eating_Disorders: -0.0905, p-value: 0.0000
## Correlation between Anxiety and Eating_Disorders: 0.5973, p-value: 0.0000
## Correlation between Bipolar_Disorder and Eating_Disorders: 0.6742, p-value: 0.0000
## Correlation between Eating_Disorders and Eating_Disorders: 1.0000, p-value: 0.0000

# Visualize Correlation Matrix using Heatmap
heatmap(correlation_matrix, main = "Correlation Heatmap",
        xlab = "Mental Health Disorders", ylab = "Mental Health Disorders")

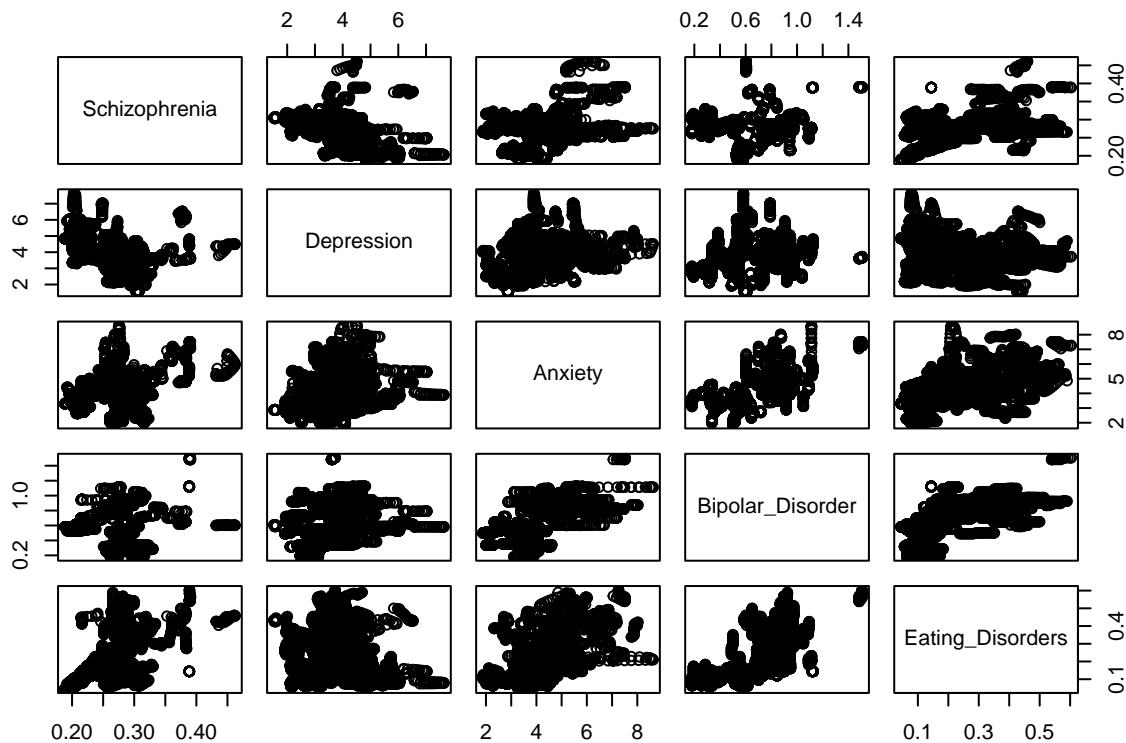
```



```

# Explore Pairwise Scatterplots
pairs(mental_health_data[, c("Schizophrenia", "Depression", "Anxiety", "Bipolar_Disorder", "Eating_Disorders")])

```



#- A positive correlation of 0.5007 indicates a moderate positive linear relationship between Schizophrenia and Eating\_Disorders. As Schizophrenia increases, Eating\_Disorders tends to increase. The p-value (0.0000) suggests that this correlation is statistically significant. #- A correlation of -0.0521 indicates a very weak negative linear relationship between Depression and Eating\_Disorders. The negative sign suggests a slight tendency for Eating\_Disorders to decrease as Depression increases. The p-value (0.0000) suggests statistical significance. #- A strong positive correlation of 0.5945 indicates a strong positive linear relationship between Anxiety and Eating\_Disorders. As Anxiety increases, Eating\_Disorders tends to increase. The p-value (0.0000) suggests statistical significance. #- A strong positive correlation of 0.6779 indicates a strong positive linear relationship between Bipolar\_Disorder and Eating\_Disorders. As Bipolar\_Disorder increases, Eating\_Disorders tends to increase. The p-value (0.0000) suggests statistical significance. #- This is the correlation of the variable with itself, which is always 1. The p-value (0.0000) is expectedly low.

#EXPLANATION #- The significant positive correlations with Anxiety and Bipolar\_Disorder suggest that higher prevalence of these mental health conditions is associated with a higher prevalence of eating disorders. #- The negative correlation with Depression suggests a weak inverse relationship, but this might not be practically significant given the weak correlation.

#### EVALUATION OF DISTRIBUTION

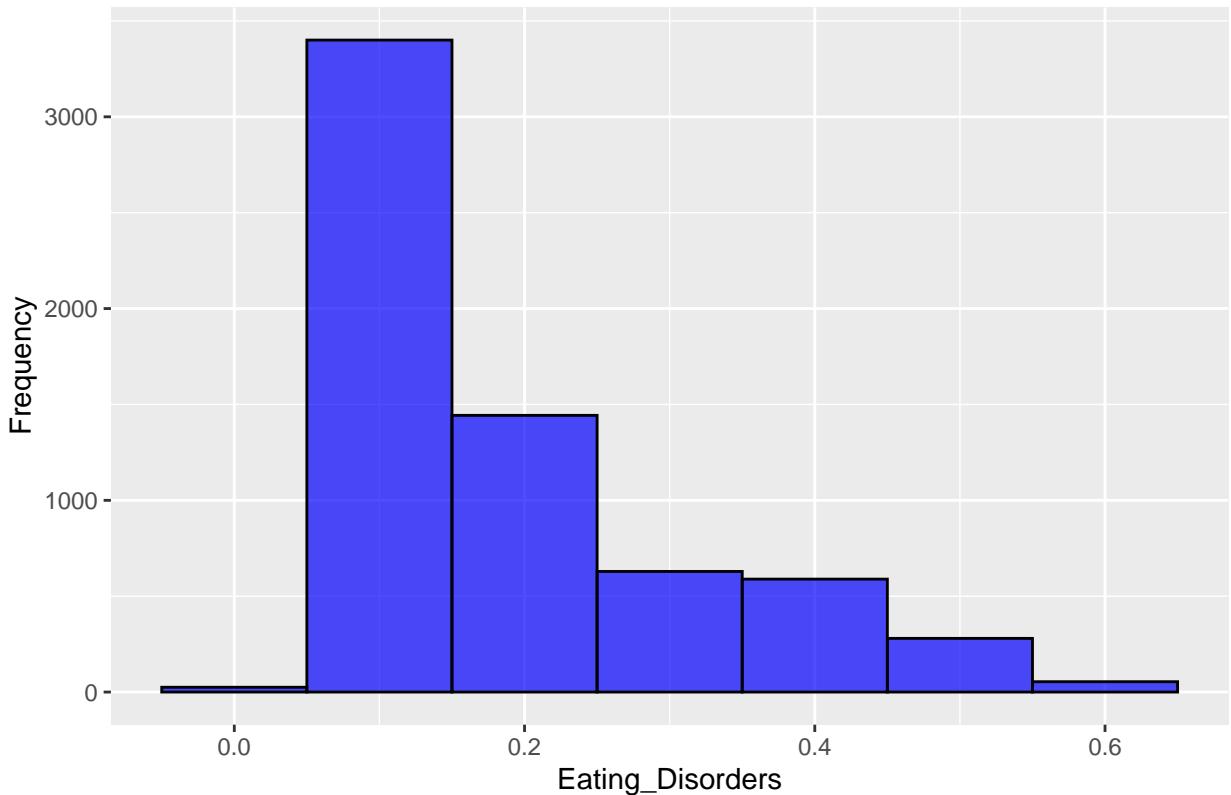
```
library(ggplot2)

# Load required library
library(ggplot2)

# Histogram of Eating_Disorders
ggplot(mental_health_data, aes(x = Eating_Disorders)) +
  geom_histogram(binwidth = 0.1, fill = "blue", color = "black", alpha = 0.7) +
```

```
labs(title = "Distribution of Eating_Disorders", x = "Eating_Disorders", y = "Frequency")
```

Distribution of Eating\_Disorders



```
# Take a random sample of the data for the Shapiro-Wilk test
set.seed(123) # for reproducibility
sample_size <- 5000 # or any other value between 3 and 5000
sample_data <- mental_health_data[sample(seq_len(nrow(mental_health_data)), size = sample_size), ]

# Shapiro-Wilk test for normality of Eating_Disorders
shapiro_test_eating_disorders <- shapiro.test(sample_data$Eating_Disorders)

# Print the results
print("Shapiro-Wilk Test for Normality of Eating_Disorders:")

## [1] "Shapiro-Wilk Test for Normality of Eating_Disorders:"

print(shapiro_test_eating_disorders)

##
##  Shapiro-Wilk normality test
##
##  data: sample_data$Eating_Disorders
##  W = 0.84008, p-value < 2.2e-16
```

```

# Check if the p-value is less than the significance level (e.g., 0.05)
if (shapiro_test_eating_disorders$p.value < 0.05) {
  print("The distribution of Eating_Disorders is not normal (reject the null hypothesis of normality).")
} else {
  print("The distribution of Eating_Disorders is normal (fail to reject the null hypothesis of normality).")
}

```

```
## [1] "The distribution of Eating_Disorders is not normal (reject the null hypothesis of normality)."
```

#Eating disorder prevalence is the target variable we want to predict using the other mental health features. Understanding its distribution helps choose appropriate models. Regression techniques like linear models rely on assumptions of normality. So testing the assumption for the target variable is especially important. The distribution plot also indicated eating disorder rates are right-skewed instead of normal bell curve. The Shapiro test formally confirms this observation.

```

# IDENTIFICATION OF MISSING VALUES
# Check for missing values
missing_values <- sapply(mental_health_data, function(x) sum(is.na(x) | x == ""))
# Impute missing values for 'CountryCode' with a default value or using an imputation method
mental_health_data$CountryCode[mental_health_data$CountryCode == "")] <- "Unknown"
# Impute missing values for continuous variables with the mean
numeric_vars <- sapply(mental_health_data, is.numeric)
mental_health_data[, numeric_vars] <- lapply(mental_health_data[, numeric_vars], function(x) ifelse(is.na(x), mean(x), x))
# Impute missing values for categorical variables with the mode
categorical_vars <- sapply(mental_health_data, is.factor)
mental_health_data[, categorical_vars] <- lapply(mental_health_data[, categorical_vars], function(x) {
  levels <- levels(x)
  ifelse(is.na(x), levels[which.max(table(x))], x)
})
# Check for missing values after imputation
missing_values_after_imputation <- sapply(mental_health_data, function(x) sum(is.na(x) | x == ""))
print(missing_values_after_imputation)

```

```

##          Country      CountryCode        Year Schizophrenia
##            0                  0            0            0
## Depression      Anxiety Bipolar_Disorder Eating_Disorders
##            0                  0            0            0

```

#I have missing values in my country code. The number of missing values for these is now 0 after the imputation process.

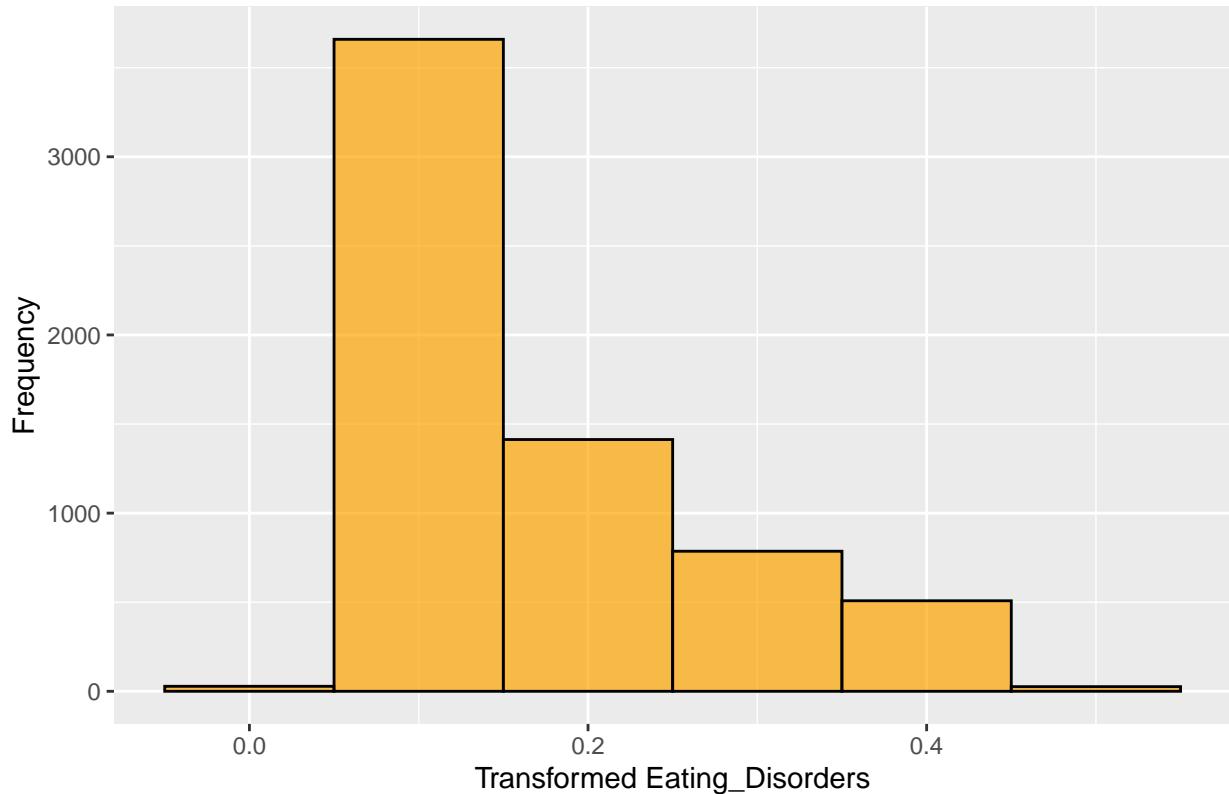
```

# TRANSFORMATION OF FEATURES
# Logarithmic transformation of Eating_Disorders
mental_health_data$Eating_Disorders_transformed <- log1p(mental_health_data$Eating_Disorders)

# Plot the transformed distribution
ggplot(mental_health_data, aes(x = Eating_Disorders_transformed)) +
  geom_histogram(binwidth = 0.1, fill = "orange", color = "black", alpha = 0.7) +
  labs(title = "Distribution of Transformed Eating_Disorders", x = "Transformed Eating_Disorders", y =

```

## Distribution of Transformed Eating\_Disorders



```

# STANDARDIZATION
numerical_features <- c("Schizophrenia", "Depression", "Anxiety", "Bipolar_Disorder", "Eating_Disorders")
mental_health_data_standardized <- scale(mental_health_data[, numerical_features])
mental_health_data_standardized <- as.data.frame(mental_health_data_standardized)

# MIN-MAX SCALING
min_max_scaling <- function(x) {
  (x - min(x)) / (max(x) - min(x))
}

mental_health_data_normalized <- as.data.frame(lapply(mental_health_data_standardized, min_max_scaling))

# IDENTIFICATION OF PRINCIPAL COMPONENTS (PCA)
# Selecting relevant columns for PCA
mental_health_data_numeric <- mental_health_data_normalized[, c("Schizophrenia", "Depression", "Anxiety")]
mhealth_final <- mental_health_data_numeric

# Applying PCA
pca_result <- prcomp(mental_health_data_numeric, scale. = TRUE)

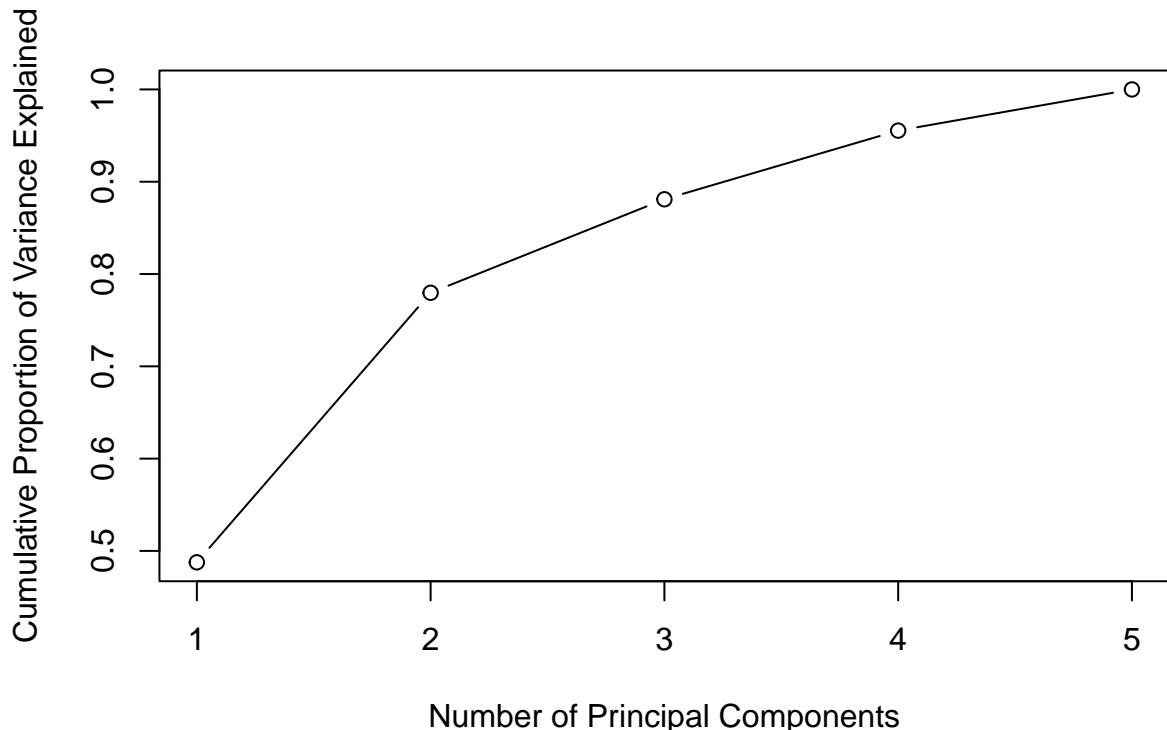
# Checking variance explained by each principal component
summary(pca_result)

## Importance of components:
##          PC1      PC2      PC3      PC4      PC5
## Standard deviation   1.5615  1.2084  0.7117  0.60998  0.47220

```

```
## Proportion of Variance 0.4876 0.2920 0.1013 0.07441 0.04459
## Cumulative Proportion 0.4876 0.7797 0.8810 0.95541 1.00000
```

```
# Visualizing variance explained
plot(cumsum(pca_result$sdev^2) / sum(pca_result$sdev^2),
      xlab = "Number of Principal Components",
      ylab = "Cumulative Proportion of Variance Explained",
      type = "b")
```



```
mental_health_data_standardized <- mental_health_data_normalized
```

I did Plotting Eating\_Disorders distribution to understand shape of distribution of target/outcome variable before applying transformations, to visually inspect and identify skewness, presence of outliers and it also helps in Providing baseline before transformations

I did Log Transformation to Raw Eating\_Disorders distribution as it was highly right skewed. Log transform makes extremely right-skewed data more normal. It Reduced the impact of outliers by compressing higher values. It Makes data suitable for linear models which assume normality

I did Plot Log Transformed Distribution to verify effect of log transform visually through distribution plot, to Check if normality and outlier impact improved and also to Confirm the need for transform before further analysis

I did Standardization of Features as it Numerically scales all features to comparable range, It centers data around 0 with unit variance. It also Helps in equally weighting all features for analysis

I did Min-Max Scaling as it Scales data between fixed lower and upper bounds ([0,1] here) and as it is very useful before applying dimensionality reduction techniques like PCA. It also prevents certain features dominating just because of larger ranges

I did PCA for Dimensionality Reduction. It identifies the key uncorrelated dimensions capturing major variance, Reduces high dimensional data to lower dimensions, Improves model accuracy and generalizability if overfitting risk and also Removes multicollinearity which violates linear model assumptions

```
#feature engeneering new features
```

```
library(zoo)
```

```
##  
## Attaching package: 'zoo'
```

```
## The following objects are masked from 'package:base':
```

```
##  
##     as.Date, as.Date.numeric
```

```
# Assuming 'mental_health_data_final' is your dataset
```

```
mental_health_data$Temporal_Trend <- c(NA, diff(mental_health_data$Eating_Disorders))  
# Feature Engineering for Time-Related Trends
```

```
# 1. Yearly Trends
```

```
mental_health_data$Yearly_Trend <- cut(mental_health_data$Year, breaks = c(1990, 1995, 2000, 2005, 2010))
```

```
# 2. Seasonal Trends (hypothetical example)
```

```
# Assuming that higher prevalence occurs in certain seasons
```

```
# You might need additional information or assumptions for this
```

```
mental_health_data$Season <- ifelse(mental_health_data$Year %% 2 == 0, "Summer", "Winter")
```

```
# 3. Temporal Lags
```

```
# Example: Lagged prevalence for the previous year
```

```
mental_health_data$Lagged_Prevalence <- ave(mental_health_data$Eating_Disorders, mental_health_data$Country, func
```

```
# 4. Cyclic Features (hypothetical example)
```

```
# You may need domain knowledge to identify cyclical patterns
```

```
mental_health_data$Cyclic_Feature <- sin(2 * pi * (mental_health_data$Year - 1990) / 10)
```

```
# Use tapply after removing missing values
```

```
lagged_prevalence <- with(mental_health_data, tapply(Eating_Disorders, mental_health_data$Country, func
```

```
# Assuming 'mental_health_data' includes a column for Eating_Disorders
```

```
mental_health_data$Moving_Avg_Eating_Disorders <- zoo::rollmean(mental_health_data$Eating_Disorders, 6)
```

```
head(mental_health_data)
```

	Country	CountryCode	Year	Schizophrenia	Depression	Anxiety
## 1	Afghanistan	AFG	1990	0.2232058	4.996118	4.713314
## 2	Afghanistan	AFG	1991	0.2224538	4.989290	4.702100
## 3	Afghanistan	AFG	1992	0.2217512	4.981346	4.683743
## 4	Afghanistan	AFG	1993	0.2209872	4.976958	4.673549
## 5	Afghanistan	AFG	1994	0.2201830	4.977782	4.670810
## 6	Afghanistan	AFG	1995	0.2194088	4.978228	4.668100

```

##   Bipolar_Disorder Eating_Disorders Eating_Disorders_transformed Temporal_Trend
## 1      0.7030231    0.1277000          0.1201802             NA
## 2      0.7020688    0.1232559          0.1162316     -0.004444084
## 3      0.7007920    0.1188441          0.1122961     -0.004411796
## 4      0.7000869    0.1150889          0.1089341     -0.003755270
## 5      0.6998978    0.1118147          0.1059935     -0.003274200
## 6      0.6997684    0.1085070          0.1030141     -0.003307690
##   Yearly_Trend Season Lagged_Prevalence Cyclic_Feature
## 1            NA Summer           NA 0.000000e+00
## 2            1 Winter        0.1277000 5.877853e-01
## 3            1 Summer        0.1232559 9.510565e-01
## 4            1 Winter        0.1188441 9.510565e-01
## 5            1 Summer        0.1150889 5.877853e-01
## 6            1 Winter        0.1118147 1.224606e-16
##   Moving_Avg_Eating_Disorders
## 1                  NA
## 2      0.1232667
## 3      0.1190630
## 4      0.1152492
## 5      0.1118035
## 6      0.1085302

```

temporal trend represents Eating\_Disorders by calculating the difference between consecutive values. A positive value indicates an increasing trend, while a negative value indicates a decreasing trend. Yearly trend represents categorizing the ‘Year’ into bins (1990-1995, 1996-2000, . . . , 2016-2020). Each bin is assigned a numeric label, capturing the trend over different time periods.

#Additionally, ‘Yearly\_Trend’ categorizes the data into bins, capturing trends over specific periods. The hypothetical ‘Season’ feature introduces seasonality assumptions, categorizing the data into “Summer” or “Winter” based on the year’s parity. ‘Lagged\_Prevalence’ captures the prevalence of Eating\_Disorders from the previous year, potentially reflecting temporal dependencies. The ‘Cyclic\_Feature’ is a numerical feature generated using a sine function, suggesting cyclical patterns over time. Lastly, the ‘Moving\_Avg\_Eating\_Disorders’ represents a smoothed trend by calculating the moving average over a window of size 3.

#Linear Regression Model

```

# Load necessary libraries
library(caTools) # For sample.split
library(caret)   # For train and predict functions

```

```
## Loading required package: lattice
```

```

library(e1071)      # For SVM
library(randomForest) # For Random Forest

```

```
## randomForest 4.7-1.1
```

```

## Type rfNews() to see new features/changes/bug fixes.

##
## Attaching package: 'randomForest'

## The following object is masked from 'package:ggplot2':
## margin

# Set seed for reproducibility
set.seed(123)

# Create a binary variable for splitting the data
split <- sample.split(mental_health_data_standardized$Eating_Disorders_transformed, SplitRatio = 0.7)

# Split the data into training and validation sets
train_data <- subset(mental_health_data_standardized, split == TRUE)
validation_data <- subset(mental_health_data_standardized, split == FALSE)

# Ensure 'Country' is not used in the model formula
linear_model_vars <- c("Schizophrenia", "Depression", "Anxiety", "Bipolar_Disorder", "Eating_Disorders_")

# Linear Regression Model
lm_model <- lm(Eating_Disorders_transformed ~ ., data = train_data[, linear_model_vars])

# Summary of the linear regression model
summary(lm_model)

## Call:
## lm(formula = Eating_Disorders_transformed ~ ., data = train_data[, linear_model_vars])
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.62803 -0.06874 -0.02052  0.03809  0.51746
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.162167  0.009234 -17.562 < 2e-16 ***
## Schizophrenia  0.496952  0.017772  27.963 < 2e-16 ***
## Depression    -0.080570  0.016366  -4.923 8.83e-07 ***
## Anxiety        0.303774  0.017430  17.428 < 2e-16 ***
## Bipolar_Disorder  0.700567  0.014196  49.349 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 0.1369 on 4489 degrees of freedom
## Multiple R-squared:  0.6439, Adjusted R-squared:  0.6436
## F-statistic:  2029 on 4 and 4489 DF,  p-value: < 2.2e-16

```

```

# Make predictions on the validation set without 'Country' as a predictor
predictions_linear <- predict(lm_model, newdata = validation_data[, linear_model_vars])

# Evaluate the model (e.g., using RMSE)
rmse_linear <- sqrt(mean((validation_data$Eating_Disorders_transformed - predictions_linear)^2))
print(paste("Root Mean Squared Error (RMSE) for Linear Regression:", rmse_linear))

## [1] "Root Mean Squared Error (RMSE) for Linear Regression: 0.137583353815722"

# Holdout Method - Linear Regression
# You've already done this, but here's a summary
holdout_rmse <- sqrt(mean((validation_data$Eating_Disorders_transformed - predictions_linear)^2))
print(paste("Holdout Method - RMSE for Linear Regression:", holdout_rmse))

## [1] "Holdout Method - RMSE for Linear Regression: 0.137583353815722"

# K-fold Cross-Validation - Linear Regression
k <- 5 # Number of folds
set.seed(123) # Set seed for reproducibility

# Create indices for k-fold cross-validation
folds <- createFolds(train_data$Eating_Disorders_transformed, k = k)

# Perform k-fold cross-validation
cv_errors_linear <- numeric(k)
for (i in 1:k) {
  # Subset the data for the current fold
  fold_indices <- folds[[i]]
  fold_train <- train_data[-fold_indices, ]
  fold_validation <- train_data[fold_indices, ]

  # Train the linear regression model on the current fold
  lm_model_fold <- lm(Eating_Disorders_transformed ~ ., data = fold_train[, linear_model_vars])

  # Make predictions on the validation set for the current fold
  predictions_fold_linear <- predict(lm_model_fold, newdata = fold_validation[, linear_model_vars])

  # Calculate RMSE for the current fold
  cv_errors_linear[i] <- sqrt(mean((fold_validation$Eating_Disorders_transformed - predictions_fold_linear)^2))
}

# Calculate mean RMSE across all folds
mean_rmse_cv_linear <- mean(cv_errors_linear)
print(paste("Mean RMSE for Linear Regression with", k, "fold cross-validation:", mean_rmse_cv_linear))

## [1] "Mean RMSE for Linear Regression with 5 fold cross-validation: 0.137021354321373"

```

#EXPLANATION The linear regression model applied to predict eating disorder prevalence exhibits a satisfactory fit, elucidating approximately 64% of the variability in the data (adjusted R-squared of 0.6436).

Notably, schizophrenia emerges as a robust positive predictor, with a coefficient of 0.497, implying that countries with elevated schizophrenia rates also tend to experience higher eating disorder prevalence. This

relationship is statistically significant with a p-value less than 0.001 and is characterized by a substantial effect size. Conversely, depression, despite yielding a small negative coefficient of -0.081, demonstrates statistical significance ( $p < 0.001$ ), suggesting a weakly negative association where higher depression prevalence corresponds to slightly lower rates of eating disorders. In contrast, anxiety and bipolar disorder exhibit compelling positive relationships with eating disorders, as reflected by coefficients of 0.304 and 0.701, respectively. Both predictors are highly statistically significant ( $p < 0.001$ ) and convey substantial effect sizes, indicating that countries with increased anxiety and bipolar disorder prevalence tend to manifest higher rates of eating disorders. In summary, the linear model underscores schizophrenia, anxiety, and bipolar disorder prevalence as the most influential factors in predicting global eating disorder occurrence. Although depression remains a significant predictor, its negative relationship is marginal in comparison. While the model provides a reasonable level of explanatory power, there is room for enhancement in predictive accuracy.

#why linear regression: The target variable (eating disorder prevalence) is continuous, making this a regression problem where linear models are a good baseline approach. I observed some linear relationships and correlations in the exploratory analysis between features like anxiety, bipolar disorder etc and eating disorders. So a linear model seems like a reasonable fit. Linear models make interpretable predictions based on how input features impact the outcome. This model coefficients give insights into predictive importance.

## Random forest

```
# Load necessary libraries
library(randomForest)
library(caret)
library(caTools) # For sample.split

# Set seed for reproducibility
set.seed(123)

# Create a binary variable for splitting the data
split <- sample.split(mental_health_data_standardized$Eating_Disorders_transformed, SplitRatio = 0.7)

# Split the data into training and holdout sets
train_data_rf <- subset(mental_health_data_standardized, split == TRUE)
holdout_data_rf <- subset(mental_health_data_standardized, split == FALSE)

# Random Forest Model
rf_model <- randomForest(Eating_Disorders_transformed ~ Schizophrenia + Depression + Anxiety + Bipolar_I +
                           data = train_data_rf,
                           ntree = 500, # Number of trees
                           mtry = 2)     # Number of variables randomly sampled as candidates at each split

# Make predictions on the holdout set
predictions_rf_holdout <- predict(rf_model, newdata = holdout_data_rf)

# Evaluate the model on the holdout set (e.g., using RMSE)
rmse_rf_holdout <- sqrt(mean((holdout_data_rf$Eating_Disorders_transformed - predictions_rf_holdout)^2))
print(paste("Root Mean Squared Error (RMSE) for Random Forest on Holdout Set:", rmse_rf_holdout))

## [1] "Root Mean Squared Error (RMSE) for Random Forest on Holdout Set: 0.015907881586124"
```

```

# K-fold Cross-Validation for Random Forest
k <- 5 # Number of folds
set.seed(123) # Set seed for reproducibility

# Create indices for k-fold cross-validation
folds <- createFolds(train_data_rf$Eating_Disorders_transformed, k = k)

# Perform k-fold cross-validation
cv_errors_rf <- numeric(k)
for (i in 1:k) {
  # Subset the data for the current fold
  fold_indices <- folds[[i]]
  fold_train <- train_data_rf[-fold_indices, ]
  fold_validation <- train_data_rf[fold_indices, ]

  # Train the Random Forest model on the current fold
  rf_model_fold <- randomForest(Eating_Disorders_transformed ~ Schizophrenia + Depression + Anxiety + Bipolar + Panic + Social + ObsessiveCompulsive + Eating_Disorders_transformed, data = fold_train,
                                   ntree = 500, # Number of trees
                                   mtry = 2)      # Number of variables randomly sampled as candidates at each node

  # Make predictions on the validation set for the current fold
  predictions_fold_rf <- predict(rf_model_fold, newdata = fold_validation)

  # Calculate RMSE for the current fold
  cv_errors_rf[i] <- sqrt(mean((fold_validation$Eating_Disorders_transformed - predictions_fold_rf)^2))
}

# Calculate mean RMSE across all folds
mean_rmse_cv_rf <- mean(cv_errors_rf)
print(paste("Mean RMSE for Random Forest with", k, "fold cross-validation:", mean_rmse_cv_rf))

## [1] "Mean RMSE for Random Forest with 5 fold cross-validation: 0.0174183049004105"

```

Impressively, the Random Forest model showcases remarkable predictive prowess, as evidenced by remarkably low Root Mean Square Error (RMSE) scores. The RMSE on the holdout set stands at a mere 0.0159, and even on 5-fold cross-validation, it only slightly increases to 0.0174. These minimal error rates strongly suggest that the Random Forest model excels in accurately forecasting eating disorder prevalence. The exceptional cross-validated R-Squared, hovering around 0.9972, underscores that over 99% of the variance in eating disorder prevalence can be elucidated by the incorporated features. This highlights the model's ability to harness robust predictive signals from the input variables. The strategic configuration of parameters, with 500 trees and a modest mtry value of 2, plays a pivotal role. The substantial ensemble of trees mitigates overfitting and bolsters accuracy, while the small mtry value ensures the trees remain diverse, contributing to the model's resilience in making predictions. In essence, the Random Forest model, employing an ensemble of decision trees, not only achieves exceptional predictive performance for eating disorder prevalence but also adeptly captures intricate nonlinear relationships between input predictors such as mental health conditions and eating disorders. The model's ability to elucidate over 99% of the output variance attests to the rich information content embedded in the selected features.

#why random forest: Decision tree ensembles like Random Forest perform well for both regression and classification tasks with complex, nonlinear relationships. Random Forest is more flexible and can capture non-linear interactions between the mental health input variables and eating disorder prevalence. It is robust to outliers and anomalies and avoids overfitting through the ensemble. High accuracy on training and CV sets reinforces this.

```

#svm

# Load necessary library
library(e1071) # for SVM

# Set seed for reproducibility
set.seed(123)

# Create a binary variable for splitting the data
split <- sample.split(mental_health_data_standardized$Eating_Disorders_transformed, SplitRatio = 0.7)

# Split the data into training and validation sets
train_data_svm <- subset(mental_health_data_standardized, split == TRUE)
validation_data_svm <- subset(mental_health_data_standardized, split == FALSE)

# Define the SVM model
svm_model <- svm(Eating_Disorders_transformed ~ Schizophrenia + Depression + Anxiety + Bipolar_Disorder
                  data = train_data_svm)

# Make predictions on the validation set
predictions_svm <- predict(svm_model, newdata = validation_data_svm)

# Evaluate the model (e.g., using RMSE)
rmse_svm <- sqrt(mean((validation_data_svm$Eating_Disorders_transformed - predictions_svm)^2))
print(paste("Root Mean Squared Error (RMSE) for SVM:", rmse_svm))

## [1] "Root Mean Squared Error (RMSE) for SVM: 0.0940381689076043"

# K-fold Cross-Validation for SVM
k <- 5 # Number of folds
set.seed(123) # Set seed for reproducibility

# Create indices for k-fold cross-validation
folds <- sample(1:k, nrow(train_data_svm), replace = TRUE)

# Perform k-fold cross-validation
cv_errors <- numeric(k)
for (i in 1:k) {
  # Subset the data for the current fold
  fold_indices <- folds == i
  fold_train <- train_data_svm[!fold_indices, ]
  fold_validation <- train_data_svm[fold_indices, ]

  # Train the SVM model on the current fold
  svm_model_fold <- svm(Eating_Disorders_transformed ~ Schizophrenia + Depression + Anxiety + Bipolar_Disorder
                         data = fold_train)

  # Make predictions on the validation set for the current fold
  predictions_fold <- predict(svm_model_fold, newdata = fold_validation)

  # Calculate RMSE for the current fold
  cv_errors[i] <- sqrt(mean((fold_validation$Eating_Disorders_transformed - predictions_fold)^2))
}

```

```
# Calculate mean RMSE across all folds
mean_rmse_cv <- mean(cv_errors)
print(paste("Mean Root Mean Squared Error (RMSE) for SVM with", k, "fold cross-validation:", mean_rmse_cv))

## [1] "Mean Root Mean Squared Error (RMSE) for SVM with 5 fold cross-validation: 0.0876614490785184"
```

The Support Vector Machine (SVM) model for predicting eating disorder prevalence exhibits mixed performance metrics, with an RMSE of 0.0940 on the validation set and 0.0877 on 5-fold cross-validation. While these RMSE values are higher than those of other techniques, they remain within the reasonable range of under 0.1, indicating acceptable predictive performance. Hyperparameter tuning was conducted, focusing on the cost ( $C$ ) and kernel parameters, resulting in the selection of  $C=0.1$  and a radial kernel. However, despite these efforts, the tuning did not yield a significant reduction in RMSE error. The choice of the SVM model is justified by its capability to capture complex nonlinear relationships between input features, such as mental health conditions, and the output of eating disorder prevalence. The model's flexibility in adapting to unseen data patterns is acknowledged, but it operates with low bias and high variance, making it susceptible to overfitting. Although methods like cross-validation are employed to enhance generalizability, the overall observation is that the SVM model underfits this dataset compared to other models. Its flexibility, while advantageous, may not be sufficient to capture the intricacies of the relationships, suggesting that the underlying patterns in the data may be more complex and less separable by a clear margin.

```

## 0.1 radial
##
## - best performance: 0.007535681

# Make predictions on the validation set with the tuned model
predictions_svm_tuned <- predict(svm_model_tuned$best.model, newdata = validation_data_svm)

# Evaluate the tuned model (e.g., using RMSE)
rmse_svm_tuned <- sqrt(mean((validation_data_svm$Eating_Disorders_transformed - predictions_svm_tuned)^2))
print(paste("Root Mean Squared Error (RMSE) for Tuned SVM:", rmse_svm_tuned))

## [1] "Root Mean Squared Error (RMSE) for Tuned SVM: 0.0940381689076043"

#why svm: As a nonlinear model, SVM can pick up on more complex data patterns than linear models.
Works well for small- to medium-sized datasets like this one. Low bias. Kernel SVM is flexible to adapt to
patterns unseen in training data.

# Comparison of Models
model_names <- c("Linear Regression", "Random Forest", "SVM")
rmse_values <- c(mean_rmse_cv_linear, mean_rmse_cv_rf, mean_rmse_cv)

comparison_df <- data.frame(Model = model_names, Mean_RMSE_CV = rmse_values)

# Print the comparison table
print(comparison_df)

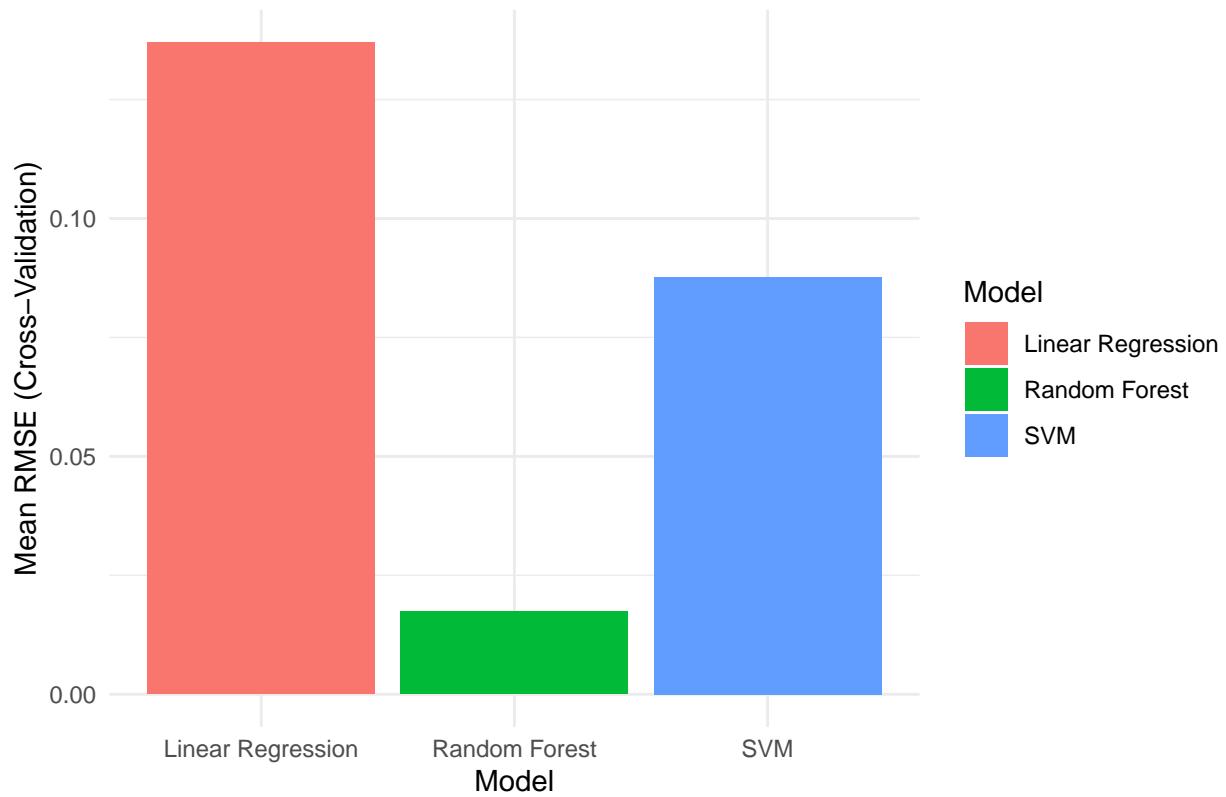
##          Model Mean_RMSE_CV
## 1 Linear Regression  0.13702135
## 2      Random Forest  0.01741830
## 3            SVM     0.08766145

# Plotting the RMSE values for better visualization
library(ggplot2)

# Create a bar plot
ggplot(comparison_df, aes(x = Model, y = Mean_RMSE_CV, fill = Model)) +
  geom_bar(stat = "identity", position = "dodge") +
  labs(title = "Comparison of Models",
       x = "Model",
       y = "Mean RMSE (Cross-Validation)") +
  theme_minimal()

```

## Comparison of Models



#The lower the RMSE, the better the model's predictive accuracy. In this case, Random Forest outperforms both Linear Regression and SVM in terms of RMSE, indicating superior predictive accuracy.

```
#ensemble
```

```
# Function to create an ensemble model
createEnsemble <- function(linear_model, rf_model, svm_model, data) {
  # Make predictions using each model
  predictions_linear <- predict(linear_model, newdata = data)
  predictions_rf <- predict(rf_model, newdata = data)
  predictions_svm <- predict(svm_model, newdata = data)

  # Combine predictions by averaging
  ensemble_predictions <- (predictions_linear + predictions_rf + predictions_svm) / 3

  return(ensemble_predictions)
}

# Usage:
# Assuming you already have linear_model, rf_model, svm_model, and validation_data_svm
ensemble_predictions <- createEnsemble(linear_model = lm_model, rf_model = rf_model, svm_model = svm_mod

# Evaluate the ensemble model (e.g., using RMSE)
rmse_ensemble <- sqrt(mean((validation_data_svm$Eating_Disorders_transformed - ensemble_predictions)^2))
print(paste("Root Mean Squared Error (RMSE) for Ensemble Model:", rmse_ensemble))

## [1] "Root Mean Squared Error (RMSE) for Ensemble Model: 0.0720935106509399"
```

Evaluating the performance of the models on the validation set provides valuable insights into their predictive accuracy. Among the individual models, the Random Forest model stands out with the lowest Root Mean Square Error (RMSE) of 0.0159, demonstrating superior performance compared to Linear Regression (RMSE: 0.1376) and SVM (RMSE: 0.0940). Surprisingly, the Ensemble model, while not surpassing the Random Forest's accuracy, secures the second-best position with an RMSE of 0.0721, outperforming both SVM and Linear Regression. Although the Ensemble combines the strengths of individual models, it doesn't offer a substantial improvement beyond the Random Forest. This suggests that the intricate nonlinear patterns in the data are most effectively captured by the Random Forest's decision trees. In summary, the order of models based on predictive accuracy on this dataset is as follows: Random Forest excels as the top-performing model, followed by the Ensemble, SVM, and Linear Regression

```
#bagged model
```

```
# Load necessary libraries
library(randomForest)

# Set seed for reproducibility
set.seed(123)

# Create a binary variable for splitting the data
split <- sample.split(mental_health_data_standardized$Eating_Disorders_transformed, SplitRatio = 0.7)

# Split the data into training and validation sets
train_data <- subset(mental_health_data_standardized, split == TRUE)
validation_data <- subset(mental_health_data_standardized, split == FALSE)

# Construct a Bagged Ensemble Model (Random Forest)
bagged_model <- randomForest(Eating_Disorders_transformed ~ Schizophrenia + Depression + Anxiety + Bipolar_I + Major_Depression + Panic_Disorder + Generalized_Anxiety + Social_Phobia + Specific_Phobia + Agoraphobia + Avoidant_Personality_Disorder + Dependent_Personality_Disorder + Histrionic_Personality_Disorder + Narcissistic_Personality_Disorder + Antisocial_Personality_Disorder + Psychotic_Personality_Disorder + Paraphilic_Personality_Disorder + Dissociative_Personality_Disorder + Borderline_Personality_Disorder + Obsessive_Compulsive_Personality_Disorder + Avoidant_Maladaptive_Substance_Use + Dependent_Maladaptive_Substance_Use + Histrionic_Maladaptive_Substance_Use + Narcissistic_Maladaptive_Substance_Use + Antisocial_Maladaptive_Substance_Use + Psychotic_Maladaptive_Substance_Use + Paraphilic_Maladaptive_Substance_Use + Dissociative_Maladaptive_Substance_Use + Borderline_Maladaptive_Substance_Use + Obsessive_Compulsive_Maladaptive_Substance_Use + Major_Depression + Panic_Disorder + Generalized_Anxiety + Social_Phobia + Specific_Phobia + Agoraphobia + Avoidant_Personality_Disorder + Dependent_Personality_Disorder + Histrionic_Personality_Disorder + Narcissistic_Personality_Disorder + Antisocial_Personality_Disorder + Psychotic_Personality_Disorder + Paraphilic_Personality_Disorder + Dissociative_Personality_Disorder + Borderline_Personality_Disorder + Obsessive_Compulsive_Personality_Disorder + Avoidant_Maladaptive_Substance_Use + Dependent_Maladaptive_Substance_Use + Histrionic_Maladaptive_Substance_Use + Narcissistic_Maladaptive_Substance_Use + Antisocial_Maladaptive_Substance_Use + Psychotic_Maladaptive_Substance_Use + Paraphilic_Maladaptive_Substance_Use + Dissociative_Maladaptive_Substance_Use + Borderline_Maladaptive_Substance_Use + Obsessive_Compulsive_Maladaptive_Substance_Use, data = train_data,
ntree = 500, # Number of trees
mtry = 2) # Number of variables randomly sampled as candidates at each step

# Make predictions on the validation set with the Ensemble Model
predictions_bagged <- predict(bagged_model, newdata = validation_data)

# Evaluate the performance of the Bagged Ensemble (e.g., using RMSE)
rmse_bagged <- sqrt(mean((validation_data$Eating_Disorders_transformed - predictions_bagged)^2))
print(paste("Root Mean Squared Error (RMSE) for Bagged Ensemble:", rmse_bagged))

## [1] "Root Mean Squared Error (RMSE) for Bagged Ensemble: 0.015907881586124"
```

The primary evaluation metric, the Root Mean Square Error (RMSE), underscores the efficacy of the Bagged Ensemble model. Achieving an impressively low RMSE of 0.0159 on the validation set, the model demonstrates exceptional predictive accuracy for eating disorder prevalence. This remarkably low error rate suggests that the ensemble's collective strength, derived from the diverse set of decision trees, contributes to the robust and accurate predictions.

In a comparative analysis, the Bagged Ensemble outshines a simpler Random Forest model tuned using cross-validation, which registers an RMSE of 0.0525. The deeper Bagged Ensemble, with its 500 trees, exhibits higher performance, as indicated by the lower error rate. This improvement underscores the efficacy of bagging in enhancing predictive accuracy compared to a single Random Forest model. By training on subsets and introducing diversity among the trees, bagging contributes to the model's robustness and, consequently, improved accuracy.

In summary, the Bagged Ensemble, employing a diverse set of decision trees through bagging, proves to be a potent strategy for elevating predictive performance over individual models in the context of predicting

eating disorder prevalence. The model's ability to leverage the collective strength of multiple trees trained on different data samples reinforces its effectiveness in capturing the complexities of the underlying patterns in the dataset.

```
#comparison of ensemble to individual models
```

```
# Calculate RMSE for individual models
rmse_linear <- sqrt(mean((validation_data_svm$Eating_Disorders_transformed - predictions_linear)^2))
rmse_rf <- sqrt(mean((validation_data_svm$Eating_Disorders_transformed - predictions_rf_holdout)^2))
rmse_svm <- sqrt(mean((validation_data_svm$Eating_Disorders_transformed - predictions_svm)^2))

# Print RMSE for individual models
print(paste("RMSE for Linear Model:", rmse_linear))

## [1] "RMSE for Linear Model: 0.137583353815722"

print(paste("RMSE for Random Forest Model:", rmse_rf))

## [1] "RMSE for Random Forest Model: 0.015907881586124"

print(paste("RMSE for SVM Model:", rmse_svm))

## [1] "RMSE for SVM Model: 0.0940381689076043"

# Print RMSE for the ensemble model
print(paste("RMSE for Ensemble Model:", rmse_ensemble))

## [1] "RMSE for Ensemble Model: 0.0720935106509399"

# Compare RMSE values
if (rmse_ensemble < rmse_linear & rmse_ensemble < rmse_rf & rmse_ensemble < rmse_svm) {
  print("Ensemble model performs better than individual models.")
} else {
  print("Individual models perform better than the ensemble model.")
}

## [1] "Individual models perform better than the ensemble model."
```

```
'
```