# Social Media Impact on Young Female Mental Health

April 8, 2025

```python
[39]: #Load packages
      import seaborn as sns
      import pandas as pd
      import matplotlib.pyplot as plt
      from sklearn.ensemble import RandomForestRegressor
      from sklearn.model_selection import train_test_split
      from sklearn.metrics import accuracy_score

      import warnings
      warnings.filterwarnings('ignore')
```

```python
[40]: #Importing csv dataset
      df= pd.read_csv('smmh.csv')
```

```python
[41]: #Checking the datype
      df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 481 entries, 0 to 480
Data columns (total 21 columns):
 #   Column
Non-Null Count  Dtype
---  ------
--------------  -----
 0   Timestamp
481 non-null    object
 1   1. What is your age?
481 non-null    float64
 2   2. Gender
481 non-null    object
 3   3. Relationship Status
481 non-null    object
 4   4. Occupation Status
481 non-null    object
 5   5. What type of organizations are you affiliated with?
451 non-null    object
 6   6. Do you use social media?
481 non-null    object
```

```
 7   7. What social media platforms do you commonly use?
481 non-null    object
 8   8. What is the average time you spend on social media every day?
481 non-null    object
 9   9. How often do you find yourself using Social media without a specific
purpose?                                          481 non-null    int64
10  10. How often do you get distracted by Social media when you are busy doing
something?                                        481 non-null    int64
11  11. Do you feel restless if you haven't used Social media in a while?
481 non-null    int64
12  12. On a scale of 1 to 5, how easily distracted are you?
481 non-null    int64
13  13. On a scale of 1 to 5, how much are you bothered by worries?
481 non-null    int64
14  14. Do you find it difficult to concentrate on things?
481 non-null    int64
15  15. On a scale of 1-5, how often do you compare yourself to other
successful people through the use of social media?  481 non-null    int64
16  16. Following the previous question, how do you feel about these
comparisons, generally speaking?                  481 non-null    int64
17  17. How often do you look to seek validation from features of social media?
481 non-null    int64
18  18. How often do you feel depressed or down?
481 non-null    int64
19  19. On a scale of 1 to 5, how frequently does your interest in daily
activities fluctuate?                             481 non-null    int64
20  20. On a scale of 1 to 5, how often do you face issues regarding sleep?
481 non-null    int64
dtypes: float64(1), int64(12), object(8)
memory usage: 79.0+ KB
```

[42]: *#Renaming the columns*

```
df = df.rename(columns={'1. What is your age?':'Age', '2. Gender':'Gender', '3.␣
 ↪Relationship Status': 'Relationship status', '4. Occupation Status':␣
 ↪'Occupation', '5. What type of organizations are you affiliated with?':␣
 ↪'Organizations', '6. Do you use social media?':'Do you use social media?',␣
 ↪'7. What social media platforms do you commonly use?': 'Social media␣
 ↪platform', '8. What is the average time you spend on social media every day?␣
 ↪':'Time on social media', '9. How often do you find yourself using Social␣
 ↪media without a specific purpose?':'Using social media without a purpose',␣
 ↪'10. How often do you get distracted by Social media when you are busy doing␣
 ↪something?' : 'Distracted by social media', "11. Do you feel restless if you␣
 ↪haven't used Social media in a while?": 'Restless without social media','12.␣
 ↪On a scale of 1 to 5, how easily distracted are you?': 'Easily distracted',␣
 ↪'13. On a scale of 1 to 5, how much are you bothered by worries?':'Bothered␣
 ↪by worries', '14. Do you find it difficult to concentrate on things?':␣
 ↪'Difficult to concentrate', '15. On a scale of 1-5, how often do you compare␣
 ↪yourself to other successful people through the use of social media?' :␣
 ↪'Comparing yourself in social media', '16. Following the previous question,␣
 ↪how do you feel about these comparisons, generally speaking?': 'General␣
 ↪comparisons', '17. How often do you look to seek validation from features of␣
 ↪social media?':'Seeking validation through social media', '18. How often do␣
 ↪you feel depressed or down?':'Depressed or down', '19. On a scale of 1 to 5,␣
 ↪how frequently does your interest in daily activities fluctuate?': 'Interest␣
 ↪change in daily activities', '20. On a scale of 1 to 5, how often do you␣
 ↪face issues regarding sleep?':'Sleep issues'})
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 481 entries, 0 to 480
Data columns (total 21 columns):
 #   Column                                   Non-Null Count  Dtype
---  ------                                   --------------  -----
 0   Timestamp                                481 non-null    object
 1   Age                                      481 non-null    float64
 2   Gender                                   481 non-null    object
 3   Relationship status                      481 non-null    object
 4   Occupation                               481 non-null    object
 5   Organizations                            451 non-null    object
 6   Do you use social media?                 481 non-null    object
 7   Social media platform                    481 non-null    object
 8   Time on social media                     481 non-null    object
 9   Using social media without a purpose     481 non-null    int64
 10  Distracted by social media               481 non-null    int64
 11  Restless without social media            481 non-null    int64
 12  Easily distracted                        481 non-null    int64
 13  Bothered by worries                      481 non-null    int64
 14  Difficult to concentrate                 481 non-null    int64
 15  Comparing yourself in social media       481 non-null    int64
 16  General comparisons                      481 non-null    int64
```

```
17  Seeking validation through social media  481 non-null   int64
18  Depressed or down                        481 non-null   int64
19  Interest change in daily activities      481 non-null   int64
20  Sleep issues                             481 non-null   int64
dtypes: float64(1), int64(12), object(8)
memory usage: 79.0+ KB
```

[43]:
```
#Checking the head of the columns
df.head()
```

[43]:

|   | Timestamp | Age | Gender | Relationship status | Occupation | \ |
|---|---|---|---|---|---|---|
| 0 | 4/18/2022 19:18:47 | 21.0 | Male | In a relationship | University Student | |
| 1 | 4/18/2022 19:19:28 | 21.0 | Female | Single | University Student | |
| 2 | 4/18/2022 19:25:59 | 21.0 | Female | Single | University Student | |
| 3 | 4/18/2022 19:29:43 | 21.0 | Female | Single | University Student | |
| 4 | 4/18/2022 19:33:31 | 21.0 | Female | Single | University Student | |

|   | Organizations | Do you use social media? | \ |
|---|---|---|---|
| 0 | University | Yes | |
| 1 | University | Yes | |
| 2 | University | Yes | |
| 3 | University | Yes | |
| 4 | University | Yes | |

|   | Social media platform | Time on social media | \ |
|---|---|---|---|
| 0 | Facebook, Twitter, Instagram, YouTube, Discord… | Between 2 and 3 hours | |
| 1 | Facebook, Twitter, Instagram, YouTube, Discord… | More than 5 hours | |
| 2 | Facebook, Instagram, YouTube, Pinterest | Between 3 and 4 hours | |
| 3 | Facebook, Instagram | More than 5 hours | |
| 4 | Facebook, Instagram, YouTube | Between 2 and 3 hours | |

|   | Using social media without a purpose | … | Restless without social media | \ |
|---|---|---|---|---|
| 0 | 5 | … | 2 | |
| 1 | 4 | … | 2 | |
| 2 | 3 | … | 1 | |
| 3 | 4 | … | 1 | |
| 4 | 3 | … | 4 | |

|   | Easily distracted | Bothered by worries | Difficult to concentrate | \ |
|---|---|---|---|---|
| 0 | 5 | 2 | 5 | |
| 1 | 4 | 5 | 4 | |
| 2 | 2 | 5 | 4 | |
| 3 | 3 | 5 | 3 | |
| 4 | 4 | 5 | 5 | |

|   | Comparing yourself in social media | General comparisons | \ |
|---|---|---|---|
| 0 | 2 | 3 | |

|   | | |
|---|---|---|
| 1 | 5 | 1 |
| 2 | 3 | 3 |
| 3 | 5 | 1 |
| 4 | 3 | 3 |

|   | Seeking validation through social media | Depressed or down \ |
|---|---|---|
| 0 | 2 | 5 |
| 1 | 1 | 5 |
| 2 | 1 | 4 |
| 3 | 2 | 4 |
| 4 | 3 | 4 |

|   | Interest change in daily activities | Sleep issues |
|---|---|---|
| 0 | 4 | 5 |
| 1 | 4 | 5 |
| 2 | 2 | 5 |
| 3 | 3 | 2 |
| 4 | 4 | 1 |

[5 rows x 21 columns]

[44]:
```python
#Creating a copy of the dataset
df_clean = df.copy()
#Confirming the copy was created
df_clean.head(5)
```

[44]:

|   | Timestamp | Age | Gender | Relationship status | Occupation \ |
|---|---|---|---|---|---|
| 0 | 4/18/2022 19:18:47 | 21.0 | Male | In a relationship | University Student |
| 1 | 4/18/2022 19:19:28 | 21.0 | Female | Single | University Student |
| 2 | 4/18/2022 19:25:59 | 21.0 | Female | Single | University Student |
| 3 | 4/18/2022 19:29:43 | 21.0 | Female | Single | University Student |
| 4 | 4/18/2022 19:33:31 | 21.0 | Female | Single | University Student |

|   | Organizations | Do you use social media? \ |
|---|---|---|
| 0 | University | Yes |
| 1 | University | Yes |
| 2 | University | Yes |
| 3 | University | Yes |
| 4 | University | Yes |

|   | Social media platform | Time on social media \ |
|---|---|---|
| 0 | Facebook, Twitter, Instagram, YouTube, Discord… | Between 2 and 3 hours |
| 1 | Facebook, Twitter, Instagram, YouTube, Discord… | More than 5 hours |
| 2 | Facebook, Instagram, YouTube, Pinterest | Between 3 and 4 hours |
| 3 | Facebook, Instagram | More than 5 hours |
| 4 | Facebook, Instagram, YouTube | Between 2 and 3 hours |

```
      Using social media without a purpose  …  Restless without social media  \
0                                       5  …                               2
1                                       4  …                               2
2                                       3  …                               1
3                                       4  …                               1
4                                       3  …                               4

      Easily distracted  Bothered by worries  Difficult to concentrate  \
0                     5                    2                         5
1                     4                    5                         4
2                     2                    5                         4
3                     3                    5                         3
4                     4                    5                         5

      Comparing yourself in social media  General comparisons  \
0                                      2                    3
1                                      5                    1
2                                      3                    3
3                                      5                    1
4                                      3                    3

      Seeking validation through social media  Depressed or down  \
0                                            2                  5
1                                            1                  5
2                                            1                  4
3                                            2                  4
4                                            3                  4

      Interest change in daily activities  Sleep issues
0                                        4             5
1                                        4             5
2                                        2             5
3                                        3             2
4                                        4             1

[5 rows x 21 columns]
```

We will split column "7. What social media platforms do you commonly use" based on the commas for our analysis.

```
[45]: #Creating a dataframe with the split data
      data = {
          '7. What social media platforms do you commonly use?': [
              'Facebook, Instagram, Twitter',
              'Instagram, Youtube',
              'Twitter, Snapchat, TikTok',
              'Reddit, Facebook'
```

```
        ]
}

df_split = pd.DataFrame(data)

# Splitting column '7. What social media platforms do you commonly use?' into␣
 ↪multiple columns based on commas
all_platforms = df_split['7. What social media platforms do you commonly use?'].
 ↪astype(str).str.split(',', expand=True)

# Define the list of possible social media platforms
platforms = ['Facebook', 'Twitter', 'Instagram', 'Youtube', 'Snapchat',␣
 ↪'Discord', 'Reddit', 'Pinterest', 'TikTok']

# Create a column for each platform and set 1 if that platform is in the␣
 ↪response
for platform in platforms:
    df_split[platform] = df_split['7. What social media platforms do you␣
 ↪commonly use?'].apply(lambda x: 1 if platform in x else 0)

#Printing the results
df_split
```

[45]:
```
   7. What social media platforms do you commonly use?  Facebook  Twitter  \
0                      Facebook, Instagram, Twitter             1        1
1                                Instagram, Youtube             0        0
2                           Twitter, Snapchat, TikTok          0        1
3                                   Reddit, Facebook           1        0

   Instagram  Youtube  Snapchat  Discord  Reddit  Pinterest  TikTok
0          1        0         0        0       0          0       0
1          1        1         0        0       0          0       0
2          0        0         1        0       0          0       1
3          0        0         0        0       1          0       0
```

[46]:
```
#Calculating the percentage of missing values
missing = df_clean.isnull().sum().sort_values(ascending=False)/len(df_clean)*100
missing
```

[46]:
```
Organizations                           6.237006
Timestamp                               0.000000
Restless without social media           0.000000
Interest change in daily activities     0.000000
Depressed or down                       0.000000
Seeking validation through social media 0.000000
General comparisons                     0.000000
Comparing yourself in social media      0.000000
```

```
Difficult to concentrate              0.000000
Bothered by worries                   0.000000
Easily distracted                     0.000000
Distracted by social media            0.000000
Age                                   0.000000
Using social media without a purpose  0.000000
Time on social media                  0.000000
Social media platform                 0.000000
Do you use social media?              0.000000
Occupation                            0.000000
Relationship status                   0.000000
Gender                                0.000000
Sleep issues                          0.000000
dtype: float64
```

### 0.0.1 Now that the data is splitted by the social media, we will proceed to rename the columns

```python
[47]: #Renaming the social media columns
      df_split = df_split.rename(columns={'Facebook':'Social Media 1', 'Twitter':␣
       ↪'Social Media 2', 'Instagram':'Social Media 3', 'Youtube':'Social Media 4',␣
       ↪'Snapchat': 'Social Media 5', 'Discord': 'Social Media 6', 'Reddit':'Social␣
       ↪Media 7', 'Pinterest':'Social Media 8', 'TikTok':'Social Media 9'})
      #Removing column '7. What social media platforms do you commonly use?'
      df_split = df_split.drop(columns=['7. What social media platforms do you␣
       ↪commonly use?'])
      df_split.head(5)
```

```
[47]:    Social Media 1  Social Media 2  Social Media 3  Social Media 4  \
      0               1               1               1               0
      1               0               0               1               1
      2               0               1               0               0
      3               1               0               0               0

         Social Media 5  Social Media 6  Social Media 7  Social Media 8  \
      0               0               0               0               0
      1               0               0               0               0
      2               1               0               0               0
      3               0               0               1               0

         Social Media 9
      0               0
      1               0
      2               1
      3               0
```

### 0.0.2 Now we will proceed to merge our "df_split" columns to our dataset "df_clean"

```
[48]: df_clean = pd.concat([df, df_split], axis=1)
      df_clean
```

```
[48]:                  Timestamp   Age  Gender Relationship status  \
      0       4/18/2022 19:18:47  21.0    Male    In a relationship
      1       4/18/2022 19:19:28  21.0  Female               Single
      2       4/18/2022 19:25:59  21.0  Female               Single
      3       4/18/2022 19:29:43  21.0  Female               Single
      4       4/18/2022 19:33:31  21.0  Female               Single
      ..                     ...   ...     ...                  ...
      476     5/21/2022 23:38:28  24.0    Male               Single
      477      5/22/2022 0:01:05  26.0  Female              Married
      478     5/22/2022 10:29:21  29.0  Female              Married
      479     7/14/2022 19:33:47  21.0    Male               Single
      480    11/12/2022 13:16:50  53.0    Male              Married

                    Occupation            Organizations Do you use social media?  \
      0     University Student              University                       Yes
      1     University Student              University                       Yes
      2     University Student              University                       Yes
      3     University Student              University                       Yes
      4     University Student              University                       Yes
      ..                   ...                    ...                       ...
      476      Salaried Worker  University, Private                          Yes
      477      Salaried Worker              University                       Yes
      478      Salaried Worker              University                       Yes
      479   University Student              University                       Yes
      480      Salaried Worker                 Private                       Yes

                              Social media platform    Time on social media  \
      0     Facebook, Twitter, Instagram, YouTube, Discord…  Between 2 and 3 hours
      1     Facebook, Twitter, Instagram, YouTube, Discord…      More than 5 hours
      2             Facebook, Instagram, YouTube, Pinterest  Between 3 and 4 hours
      3                               Facebook, Instagram       More than 5 hours
      4                      Facebook, Instagram, YouTube  Between 2 and 3 hours
      ..                                            ...                    ...
      476                   Facebook, Instagram, YouTube  Between 2 and 3 hours
      477                             Facebook, YouTube  Between 1 and 2 hours
      478                             Facebook, YouTube  Between 2 and 3 hours
      479  Facebook, Twitter, Instagram, YouTube, Discord…  Between 2 and 3 hours
      480                             Facebook, YouTube      Less than an Hour

           Using social media without a purpose  …  Sleep issues  Social Media 1  \
      0                                        5  …             5             1.0
      1                                        4  …             5             0.0
```

```
2                                             3   …              5            0.0
3                                             4   …              2            1.0
4                                             3   …              1            NaN
..                                            …   …             …             …
476                                           3   …              4            NaN
477                                           2   …              1            NaN
478                                           3   …              2            NaN
479                                           2   …              4            NaN
480                                           2   …              3            NaN

     Social Media 2  Social Media 3  Social Media 4  Social Media 5  \
0               1.0             1.0             0.0             0.0
1               0.0             1.0             1.0             0.0
2               1.0             0.0             0.0             1.0
3               0.0             0.0             0.0             0.0
4               NaN             NaN             NaN             NaN
..               …               …               …               …
476             NaN             NaN             NaN             NaN
477             NaN             NaN             NaN             NaN
478             NaN             NaN             NaN             NaN
479             NaN             NaN             NaN             NaN
480             NaN             NaN             NaN             NaN

     Social Media 6  Social Media 7  Social Media 8  Social Media 9
0               0.0             0.0             0.0             0.0
1               0.0             0.0             0.0             0.0
2               0.0             0.0             0.0             1.0
3               0.0             1.0             0.0             0.0
4               NaN             NaN             NaN             NaN
..               …               …               …               …
476             NaN             NaN             NaN             NaN
477             NaN             NaN             NaN             NaN
478             NaN             NaN             NaN             NaN
479             NaN             NaN             NaN             NaN
480             NaN             NaN             NaN             NaN

[481 rows x 30 columns]
```

### 0.0.3 Since there weren't missing values from column question 7, and we introduce missing values by splitting that column, we will proceed to replace NaN with 0.

```
[49]:  #Fill Nan values with 0
```

```
df_clean[['Social Media 1','Social Media 2','Social Media 3','Social Media↵
    ↪4','Social Media 5','Social Media 6', 'Social Media 7', 'Social Media 8',↵
    ↪'Social Media 9']] = df_clean[['Social Media 1','Social Media 2','Social↵
    ↪Media 3','Social Media 4','Social Media 5','Social Media 6', 'Social Media↵
    ↪7', 'Social Media 8', 'Social Media 9']].fillna(0)
```

### 0.0.4 We will now look at missing values in our data.

```
[50]: #Calculating the percentage of missing values
      missing = df_clean.isnull().sum().sort_values(ascending=False)/len(df_clean)*100
      missing
```

```
[50]: Organizations                             6.237006
      Timestamp                                 0.000000
      General comparisons                       0.000000
      Social Media 8                            0.000000
      Social Media 7                            0.000000
      Social Media 6                            0.000000
      Social Media 5                            0.000000
      Social Media 4                            0.000000
      Social Media 3                            0.000000
      Social Media 2                            0.000000
      Social Media 1                            0.000000
      Sleep issues                              0.000000
      Interest change in daily activities       0.000000
      Depressed or down                         0.000000
      Seeking validation through social media   0.000000
      Comparing yourself in social media        0.000000
      Age                                       0.000000
      Difficult to concentrate                  0.000000
      Bothered by worries                       0.000000
      Easily distracted                         0.000000
      Restless without social media             0.000000
      Distracted by social media                0.000000
      Using social media without a purpose      0.000000
      Time on social media                      0.000000
      Social media platform                     0.000000
      Do you use social media?                  0.000000
      Occupation                                0.000000
      Relationship status                       0.000000
      Gender                                    0.000000
      Social Media 9                            0.000000
      dtype: float64
```

### 0.0.5 There is some missing data in the Organizations column, since that column is similar to Ocupation we decided to remove it. And we will proceed to remove unnecesary columns like Timestamp, and Social media platform

```
[51]: #Droping columns that are not needed
      df_clean = df_clean.drop(columns=['Timestamp', 'Organizations','Social media␣
       ↪platform'])
      df_clean.head(5)
```

```
[51]:      Age  Gender Relationship status          Occupation  \
      0  21.0    Male   In a relationship  University Student
      1  21.0  Female              Single  University Student
      2  21.0  Female              Single  University Student
      3  21.0  Female              Single  University Student
      4  21.0  Female              Single  University Student

        Do you use social media?   Time on social media  \
      0                      Yes  Between 2 and 3 hours
      1                      Yes      More than 5 hours
      2                      Yes  Between 3 and 4 hours
      3                      Yes      More than 5 hours
      4                      Yes  Between 2 and 3 hours

        Using social media without a purpose  Distracted by social media  \
      0                                    5                           3
      1                                    4                           3
      2                                    3                           2
      3                                    4                           2
      4                                    3                           5

        Restless without social media  Easily distracted  …  Sleep issues  \
      0                              2                  5  …             5
      1                              2                  4  …             5
      2                              1                  2  …             5
      3                              1                  3  …             2
      4                              4                  4  …             1

        Social Media 1  Social Media 2  Social Media 3  Social Media 4  \
      0             1.0             1.0             1.0             0.0
      1             0.0             0.0             1.0             1.0
      2             0.0             1.0             0.0             0.0
      3             1.0             0.0             0.0             0.0
      4             0.0             0.0             0.0             0.0

        Social Media 5  Social Media 6  Social Media 7  Social Media 8  \
      0             0.0             0.0             0.0             0.0
      1             0.0             0.0             0.0             0.0
```

```
2               1.0              0.0              0.0              0.0
3               0.0              0.0              1.0              0.0
4               0.0              0.0              0.0              0.0

     Social Media 9
0              0.0
1              0.0
2              1.0
3              0.0
4              0.0

[5 rows x 27 columns]
```

### 0.0.6 We will now look at describe() so we can have an idea if there are any outliers in the dataset.

```
[52]: #Describing the data
      df.describe()
```

```
[52]:            Age  Using social media without a purpose  \
      count  481.00000                            481.000000
      mean    26.13659                              3.553015
      std      9.91511                              1.096299
      min     13.00000                              1.000000
      25%     21.00000                              3.000000
      50%     22.00000                              4.000000
      75%     26.00000                              4.000000
      max     91.00000                              5.000000

             Distracted by social media  Restless without social media  \
      count                 481.000000                      481.000000
      mean                    3.320166                        2.588358
      std                     1.328137                        1.257059
      min                     1.000000                        1.000000
      25%                     2.000000                        2.000000
      50%                     3.000000                        2.000000
      75%                     4.000000                        3.000000
      max                     5.000000                        5.000000

             Easily distracted  Bothered by worries  Difficult to concentrate  \
      count         481.000000           481.000000                481.000000
      mean            3.349272             3.559252                  3.245322
      std             1.175552             1.283356                  1.347105
      min             1.000000             1.000000                  1.000000
      25%             3.000000             3.000000                  2.000000
      50%             3.000000             4.000000                  3.000000
      75%             4.000000             5.000000                  4.000000
```

```
max                          5.000000               5.000000                5.000000
```

```
        Comparing yourself in social media  General comparisons  \
count                        481.000000               481.000000
mean                           2.831601                 2.775468
std                            1.407835                 1.056479
min                            1.000000                 1.000000
25%                            2.000000                 2.000000
50%                            3.000000                 3.000000
75%                            4.000000                 3.000000
max                            5.000000                 5.000000

        Seeking validation through social media  Depressed or down  \
count                                481.000000         481.000000
mean                                   2.455301           3.255717
std                                    1.247739           1.313033
min                                    1.000000           1.000000
25%                                    1.000000           2.000000
50%                                    2.000000           3.000000
75%                                    3.000000           4.000000
max                                    5.000000           5.000000

        Interest change in daily activities  Sleep issues
count                          481.000000    481.000000
mean                             3.170478      3.201663
std                              1.256666      1.461619
min                              1.000000      1.000000
25%                              2.000000      2.000000
50%                              3.000000      3.000000
75%                              4.000000      5.000000
max                              5.000000      5.000000
```

### 0.0.7 The dataset doesn't seem to have outliers. We will proceed to transform to create a "Social Media Count" column.

```python
[53]: #Creating a total social media count column by summing the social media columns
      df_clean['Social Media Count']= df_clean['Social Media 1'] + df_clean['Social␣
       ↪Media 2'] + df_clean['Social Media 3'] + df_clean['Social Media 4'] +␣
       ↪df_clean['Social Media 5'] + df_clean['Social Media 6'] + df_clean['Social␣
       ↪Media 7'] + df_clean['Social Media 8'] + df_clean['Social Media 9']

      df_clean.head(5)
```

```
[53]:    Age  Gender Relationship status          Occupation  \
      0  21.0    Male   In a relationship  University Student
      1  21.0  Female              Single  University Student
      2  21.0  Female              Single  University Student
```

14

```
3  21.0  Female              Single  University Student
4  21.0  Female              Single  University Student

  Do you use social media?   Time on social media  \
0                      Yes  Between 2 and 3 hours
1                      Yes      More than 5 hours
2                      Yes  Between 3 and 4 hours
3                      Yes      More than 5 hours
4                      Yes  Between 2 and 3 hours

  Using social media without a purpose  Distracted by social media  \
0                                    5                           3
1                                    4                           3
2                                    3                           2
3                                    4                           2
4                                    3                           5

  Restless without social media  Easily distracted  …  Social Media 1  \
0                              2                  5  …             1.0
1                              2                  4  …             0.0
2                              1                  2  …             0.0
3                              1                  3  …             1.0
4                              4                  4  …             0.0

  Social Media 2  Social Media 3  Social Media 4  Social Media 5  \
0             1.0             1.0             0.0             0.0
1             0.0             1.0             1.0             0.0
2             1.0             0.0             0.0             1.0
3             0.0             0.0             0.0             0.0
4             0.0             0.0             0.0             0.0

  Social Media 6  Social Media 7  Social Media 8  Social Media 9  \
0             0.0             0.0             0.0             0.0
1             0.0             0.0             0.0             0.0
2             0.0             0.0             0.0             1.0
3             0.0             1.0             0.0             0.0
4             0.0             0.0             0.0             0.0

  Social Media Count
0                3.0
1                2.0
2                3.0
3                2.0
4                0.0

[5 rows x 28 columns]
```

### 0.0.8 We will group the non binary, unsure and there are others as 'others' since we are just interested to know data between male and female for this study.

```python
# Grouping non-binary, others, unsure into one category called 'Others'
df_clean['Gender'] = df_clean['Gender'].replace({'Nonbinary ':'Others', 'Non
  binary ':'Others', 'There are others???':'Others', 'Non-binary':'Others',
  'unsure ':'Others', 'NB':'Others', 'Trans':'Others'})
df_clean.head(5)
```

```
[54]:     Age  Gender Relationship status          Occupation  \
     0  21.0    Male   In a relationship  University Student
     1  21.0  Female              Single  University Student
     2  21.0  Female              Single  University Student
     3  21.0  Female              Single  University Student
     4  21.0  Female              Single  University Student

       Do you use social media?   Time on social media  \
     0                      Yes  Between 2 and 3 hours
     1                      Yes      More than 5 hours
     2                      Yes  Between 3 and 4 hours
     3                      Yes      More than 5 hours
     4                      Yes  Between 2 and 3 hours

       Using social media without a purpose  Distracted by social media  \
     0                                     5                           3
     1                                     4                           3
     2                                     3                           2
     3                                     4                           2
     4                                     3                           5

       Restless without social media  Easily distracted  …  Social Media 1  \
     0                              2                  5  …             1.0
     1                              2                  4  …             0.0
     2                              1                  2  …             0.0
     3                              1                  3  …             1.0
     4                              4                  4  …             0.0

       Social Media 2  Social Media 3  Social Media 4  Social Media 5  \
     0             1.0             1.0             0.0             0.0
     1             0.0             1.0             1.0             0.0
     2             1.0             0.0             0.0             1.0
     3             0.0             0.0             0.0             0.0
     4             0.0             0.0             0.0             0.0

       Social Media 6  Social Media 7  Social Media 8  Social Media 9  \
     0             0.0             0.0             0.0             0.0
     1             0.0             0.0             0.0             0.0
```

```
2              0.0              0.0              0.0              1.0
3              0.0              1.0              0.0              0.0
4              0.0              0.0              0.0              0.0

   Social Media Count
0                 3.0
1                 2.0
2                 3.0
3                 2.0
4                 0.0

[5 rows x 28 columns]
```

We will proceed to filter our data for young females (18-29)

```
[55]: young_females = df_clean[(df_clean['Age']>=18) & (df_clean['Age']<=29) &␣
      ↪(df_clean['Gender']=='Female')]
      young_females.head(5)
```

```
[55]:    Age  Gender Relationship status         Occupation  \
      1  21.0  Female              Single  University Student
      2  21.0  Female              Single  University Student
      3  21.0  Female              Single  University Student
      4  21.0  Female              Single  University Student
      5  22.0  Female              Single  University Student

        Do you use social media?  Time on social media  \
      1                      Yes     More than 5 hours
      2                      Yes  Between 3 and 4 hours
      3                      Yes     More than 5 hours
      4                      Yes  Between 2 and 3 hours
      5                      Yes  Between 2 and 3 hours

        Using social media without a purpose  Distracted by social media  \
      1                                     4                           3
      2                                     3                           2
      3                                     4                           2
      4                                     3                           5
      5                                     4                           4

        Restless without social media  Easily distracted  …  Social Media 1  \
      1                              2                  4  …             0.0
      2                              1                  2  …             0.0
      3                              1                  3  …             1.0
      4                              4                  4  …             0.0
      5                              2                  3  …             0.0

        Social Media 2  Social Media 3  Social Media 4  Social Media 5  \
```

```
   1          0.0          1.0          1.0          0.0
   2          1.0          0.0          0.0          1.0
   3          0.0          0.0          0.0          0.0
   4          0.0          0.0          0.0          0.0
   5          0.0          0.0          0.0          0.0

      Social Media 6  Social Media 7  Social Media 8  Social Media 9  \
   1          0.0             0.0             0.0             0.0
   2          0.0             0.0             0.0             1.0
   3          0.0             1.0             0.0             0.0
   4          0.0             0.0             0.0             0.0
   5          0.0             0.0             0.0             0.0

      Social Media Count
   1               2.0
   2               3.0
   3               2.0
   4               0.0
   5               0.0

   [5 rows x 28 columns]
```

[56]: *#Checking the data shape to ensure the data was filtered correctly*
      young_females.shape

[56]: (218, 28)

### 0.0.9 We will create a young_females_depressed dataframe so we can create different analysis charts.

[57]: *#Creating a dataframe with depressed young females.*
      young_females_depressed = young_females[young_females['Depressed or down'] >= 3]
      young_females_depressed

[57]:
```
        Age  Gender Relationship status          Occupation  \
   1    21.0  Female              Single  University Student
   2    21.0  Female              Single  University Student
   3    21.0  Female              Single  University Student
   4    21.0  Female              Single  University Student
   5    22.0  Female              Single  University Student
   ..    …     …                   …                   …
   462  28.0  Female             Married     Salaried Worker
   470  20.0  Female              Single  University Student
   471  20.0  Female              Single  University Student
   473  26.0  Female             Married  University Student
   477  26.0  Female             Married     Salaried Worker
```

```
     Do you use social media?   Time on social media  \
1                      Yes       More than 5 hours
2                      Yes  Between 3 and 4 hours
3                      Yes       More than 5 hours
4                      Yes  Between 2 and 3 hours
5                      Yes  Between 2 and 3 hours
..                     …                      …
462                    Yes  Between 2 and 3 hours
470                    Yes  Between 1 and 2 hours
471                    Yes       More than 5 hours
473                    Yes  Between 2 and 3 hours
477                    Yes  Between 1 and 2 hours

     Using social media without a purpose  Distracted by social media  \
1                                        4                           3
2                                        3                           2
3                                        4                           2
4                                        3                           5
5                                        4                           4
..                                       …                           …
462                                      5                           4
470                                      2                           2
471                                      5                           4
473                                      4                           4
477                                      2                           1

     Restless without social media  Easily distracted  …  Social Media 1  \
1                                 2                  4  …             0.0
2                                 1                  2  …             0.0
3                                 1                  3  …             1.0
4                                 4                  4  …             0.0
5                                 2                  3  …             0.0
..                                …                  … …               …
462                               1                  5  …             0.0
470                               1                  3  …             0.0
471                               2                  4  …             0.0
473                               3                  3  …             0.0
477                               2                  3  …             0.0

     Social Media 2  Social Media 3  Social Media 4  Social Media 5  \
1               0.0             1.0             1.0             0.0
2               1.0             0.0             0.0             1.0
3               0.0             0.0             0.0             0.0
4               0.0             0.0             0.0             0.0
5               0.0             0.0             0.0             0.0
..                …               …               …               …
462             0.0             0.0             0.0             0.0
```

```
470            0.0            0.0            0.0            0.0
471            0.0            0.0            0.0            0.0
473            0.0            0.0            0.0            0.0
477            0.0            0.0            0.0            0.0

     Social Media 6  Social Media 7  Social Media 8  Social Media 9  \
1               0.0             0.0             0.0             0.0
2               0.0             0.0             0.0             1.0
3               0.0             1.0             0.0             0.0
4               0.0             0.0             0.0             0.0
5               0.0             0.0             0.0             0.0
..              ...             ...             ...             ...
462             0.0             0.0             0.0             0.0
470             0.0             0.0             0.0             0.0
471             0.0             0.0             0.0             0.0
473             0.0             0.0             0.0             0.0
477             0.0             0.0             0.0             0.0

     Social Media Count
1                   2.0
2                   3.0
3                   2.0
4                   0.0
5                   0.0
..                  ...
462                 0.0
470                 0.0
471                 0.0
473                 0.0
477                 0.0

[175 rows x 28 columns]
```

### 0.0.10 We will keep young_females dataframe as a categorical dataset and we will transform to numerical, and be called young_females_numerical dataframe so we can create different analysis charts.

```python
[58]: young_females_numerical = young_females.copy()
```

```python
[59]: #Converting categorical data to numerical data.
      young_females_numerical['Time on social media'] = df_clean['Time on social␣
       ↪media'].replace({'Less than an Hour':0, 'Between 1 and 2 hours':1, 'Between␣
       ↪2 and 3 hours':2, 'Between 3 and 4 hours':3, 'Between 4 and 5 hours':4,␣
       ↪'More than 5 hours':5})
      young_females_numerical['Do you use social media?'] = df_clean['Do you use␣
       ↪social media?'].replace({'Yes':1,})
```

```
young_females_numerical['Occupation'] = df_clean['Occupation'].replace({'School␣
 ↪student':0, 'School Student':0, 'University Student':1, 'Salaried Worker':2,␣
 ↪'Retired':3 })
young_females_numerical['Relationship status'] = df_clean['Relationship␣
 ↪status'].replace({'Single':0, 'In a relationship':1, 'Married':2, 'Divorced':
 ↪3 })
young_females_numerical['Gender'] = df_clean['Gender'].replace({'Male':0,␣
 ↪'Female':1 })
young_females_numerical.head(5)
```

[59]:      Age Gender  Relationship status  Occupation Do you use social media?  \
    1  21.0      1                    0           1                         1
    2  21.0      1                    0           1                         1
    3  21.0      1                    0           1                         1
    4  21.0      1                    0           1                         1
    5  22.0      1                    0           1                         1

       Time on social media  Using social media without a purpose  \
    1                      5                                     4
    2                      3                                     3
    3                      5                                     4
    4                      2                                     3
    5                      2                                     4

       Distracted by social media  Restless without social media  \
    1                            3                               2
    2                            2                               1
    3                            2                               1
    4                            5                               4
    5                            4                               2

       Easily distracted  …  Social Media 1  Social Media 2  Social Media 3  \
    1                  4  …             0.0             0.0             1.0
    2                  2  …             0.0             1.0             0.0
    3                  3  …             1.0             0.0             0.0
    4                  4  …             0.0             0.0             0.0
    5                  3  …             0.0             0.0             0.0

       Social Media 4  Social Media 5  Social Media 6  Social Media 7  \
    1             1.0             0.0             0.0             0.0
    2             0.0             1.0             0.0             0.0
    3             0.0             0.0             0.0             1.0
    4             0.0             0.0             0.0             0.0
    5             0.0             0.0             0.0             0.0

       Social Media 8  Social Media 9  Social Media Count
    1             0.0             0.0                 2.0
```

```
2              0.0              1.0              3.0
3              0.0              0.0              2.0
4              0.0              0.0              0.0
5              0.0              0.0              0.0
```

[5 rows x 28 columns]

### 0.0.11 We will create a two dataframes from mental_health_columns and social_media_columns, and then get the statistics

```python
[60]: # Summary statistics for mental health variables
mental_health_columns = ['Depressed or down', 'Easily distracted', 'Bothered by␣
 ↪worries', 'Difficult to concentrate', 'Interest change in daily activities',␣
 ↪'Sleep issues']
social_media_columns = ['Using social media without a purpose','Time on social␣
 ↪media', 'Distracted by social media', 'Restless without social media',␣
 ↪'Seeking validation through social media']
# Calculate mean, median, and standard deviation for mental health and social␣
 ↪media usage
mental_health_stats = young_females_numerical[mental_health_columns].describe()
social_media_stats = young_females_numerical[social_media_columns].describe()
print("Mental Health Statistics:\n", mental_health_stats)
print("Social Media Usage Statistics:\n", social_media_stats)
```

```
Mental Health Statistics:
       Depressed or down  Easily distracted  Bothered by worries  \
count         218.000000         218.000000           218.000000
mean            3.573394           3.559633             3.853211
std             1.197782           1.114886             1.142603
min             1.000000           1.000000             1.000000
25%             3.000000           3.000000             3.000000
50%             4.000000           4.000000             4.000000
75%             5.000000           4.000000             5.000000
max             5.000000           5.000000             5.000000

       Difficult to concentrate  Interest change in daily activities  \
count                218.000000                           218.000000
mean                   3.472477                             3.463303
std                    1.222551                             1.164384
min                    1.000000                             1.000000
25%                    3.000000                             3.000000
50%                    4.000000                             4.000000
75%                    4.000000                             4.000000
max                    5.000000                             5.000000

       Sleep issues
count    218.000000
```

```
mean        3.302752
std         1.433661
min         1.000000
25%         2.000000
50%         4.000000
75%         5.000000
max         5.000000
Social Media Usage Statistics:
       Using social media without a purpose  Time on social media  \
count                           218.000000            218.000000
mean                              3.692661              3.412844
std                               1.030432              1.359120
min                               1.000000              0.000000
25%                               3.000000              2.000000
50%                               4.000000              3.000000
75%                               4.000000              5.000000
max                               5.000000              5.000000


       Distracted by social media  Restless without social media  \
count                  218.000000                     218.000000
mean                     3.504587                       2.788991
std                      1.267269                       1.295674
min                      1.000000                       1.000000
25%                      3.000000                       2.000000
50%                      4.000000                       3.000000
75%                      5.000000                       4.000000
max                      5.000000                       5.000000


       Seeking validation through social media
count                               218.000000
mean                                  2.587156
std                                   1.189129
min                                   1.000000
25%                                   2.000000
50%                                   3.000000
75%                                   3.000000
max                                   5.000000
```

```python
[61]:  # Correlation between mental health and social media variables
       correlation_matrix = young_females_numerical[mental_health_columns +
         ↪social_media_columns].corr()
       # Display correlation matrix
       import seaborn as sns
       import matplotlib.pyplot as plt
       # Plot heatmap of correlation matrix
       plt.figure(figsize=(8, 6))
       sns.heatmap(correlation_matrix, annot=True, fmt=".2f")
```

```
plt.title("Correlation between Mental Health and Social Media Usage")
plt.show()
```



Correlation between Mental Health and Social Media Usage

### 0.0.12 We will proceed to create feature importance based on the higher depression count.

```
[62]: #Split the data into training and testing sets
      X = young_females_numerical.drop(columns=['Depressed or down'])
      y_binary = (young_females_numerical['Depressed or down'] >= 3).astype(int)
```

```
[63]: #Creating the feature names
```

```
features_names= ['Age', 'Gender', 'Relationship status', 'Occupation',␣
 ↪'Organizations', 'Do you use social media?', 'Social media platform', 'Time␣
 ↪on social media', 'Using social media without a purpose', 'Distracted by␣
 ↪social media', 'Restless without social media', 'Easily distracted',␣
 ↪'Bothered by worries', 'Difficult to concentrate', 'Comparing yourself in␣
 ↪social media', 'General comparisons', 'Seeking validation through social␣
 ↪media', 'Depressed or down', 'Interest change in daily activities', 'Sleep␣
 ↪issues']
```

[64]:
```
#Creating feature importance plot
rf =RandomForestRegressor(n_estimators=100, random_state=42, n_jobs=-1)
rf.fit(X, y_binary)
feat_importances = pd.Series(rf.feature_importances_, index= X.columns)
feat_importances = feat_importances.sort_values(ascending=False)
feat_importances.plot(kind ='bar', figsize=(10,8), title='Feature Importances')
plt.show()
```

## Feature Importances



```
[65]: #Creating a "Is depressed" column
      df_clean['Is_depressed'] = df_clean['Depressed or down']>= 3

      #Grouping the data by gender, and calculating the percentage of depressed people
      depression_by_gender = df_clean.groupby('Gender')['Is_depressed'].mean()*100
      depression_by_gender
```

```
[65]: Gender
      Female    74.524715
      Male      64.454976
      Others    85.714286
      Name: Is_depressed, dtype: float64
```

### 0.0.13 Statistics comparing women vs men mental health are presenting a higher impact in women where anxiety is present 23% more compared to men, and depression is 50% higher. (Percentage of Depressed Individuals by Gender)

```
[66]: #Creating a plot for the depression
      # Plot the results
      plt.figure(figsize=(8, 5))
      depression_by_gender.plot(kind='bar', color=['blue', 'pink', 'orange'])

      # Add labels and title
      plt.title('Percentage of Depressed Individuals by Gender')
      plt.xlabel('Gender')
      plt.ylabel('Percentage (%)')
      plt.xticks(rotation=0)
      plt.show()
```

```
[67]:   # Gender vs Feeling Down
        plt.figure(figsize=(8, 5))
        age_filtered = df_clean[df_clean['Age'].between(18, 29)]
        sns.countplot(data=age_filtered, x="Depressed or down", hue="Gender",
                      palette="muted")
        plt.title("Feeling Depressed or Down by Gender (18-29)")
        plt.tight_layout()
        plt.show()
```

Feeling Depressed or Down by Gender (18–29)



### 0.0.14 According to Data Reportal, a person spends 2 hours and 30 minutes daily using social media. Our dataset showed that the majority of young females analyzed, around 70 of 218, spend More than 5 hours daily.

```
[68]:   #Time on social media
        plt.figure(figsize=(8, 5))
        sns.countplot(data=young_females, x="Time on social media",
                      order=young_females["Time on social media"].value_counts().index,
                      palette="Set2")
        plt.title("Young Females 18-29: Time Spent on Social Media")
        plt.xticks(rotation=90)
        plt.tight_layout()
        plt.show()
```

## Young Females 18-29: Time Spent on Social Media



```
[69]:  # Usage without purpose vs time
       plt.figure(figsize=(8, 5))
       sns.boxplot(data=young_females, x="Using social media without a purpose",␣
        ↪y="Time on social media",
                   palette="Set3")
       plt.title("Purpose Use vs Time Spent on Social Media (Young Females 18-29)")
       plt.tight_layout()
       plt.show()
```

## Purpose Use vs Time Spent on Social Media (Young Females 18-29)



```
[70]: #Creating a "Is depressed" column
      df_clean['Is_depressed'] = df_clean['Depressed or down']>= 3
```

```
[71]: #Creating an "old female" dataframe
      old_females = df_clean[(df_clean['Age']>29) & (df_clean['Gender']=='Female')]
      old_females.head(5)
```

```
[71]:      Age  Gender Relationship status      Occupation  \
      25  35.0  Female             Married  Salaried Worker
      42  56.0  Female             Married          Retired
      49  33.0  Female              Single  Salaried Worker
      87  32.0  Female             Married  Salaried Worker
      94  30.0  Female              Single  Salaried Worker

         Do you use social media?    Time on social media  \
      25                      Yes  Between 3 and 4 hours
      42                      Yes  Between 1 and 2 hours
      49                      Yes  Between 3 and 4 hours
      87                      Yes  Between 1 and 2 hours
      94                      Yes  Between 4 and 5 hours

         Using social media without a purpose  Distracted by social media  \
      25                                    4                           4
      42                                    1                           1
      49                                    3                           3
```

```
87                                                    4                        3
94                                                    4                        5

        Restless without social media  Easily distracted  …  Social Media 2  \
25                                   3                  3  …             0.0
42                                   1                  1  …             0.0
49                                   1                  3  …             0.0
87                                   2                  5  …             0.0
94                                   3                  2  …             0.0

        Social Media 3  Social Media 4  Social Media 5  Social Media 6  \
25                 0.0             0.0             0.0             0.0
42                 0.0             0.0             0.0             0.0
49                 0.0             0.0             0.0             0.0
87                 0.0             0.0             0.0             0.0
94                 0.0             0.0             0.0             0.0

        Social Media 7  Social Media 8  Social Media 9  Social Media Count  \
25                 0.0             0.0             0.0                 0.0
42                 0.0             0.0             0.0                 0.0
49                 0.0             0.0             0.0                 0.0
87                 0.0             0.0             0.0                 0.0
94                 0.0             0.0             0.0                 0.0

        Is_depressed
25              True
42             False
49             False
87             False
94              True

[5 rows x 29 columns]
```

```
[72]:  #Creating a "old female" dataframe
       old_females['Is_depressed'] = old_females['Depressed or down']>= 3
       young_females['Is_depressed'] = young_females['Depressed or down']>= 3
```

```
[73]:  # Calculate the depression percentage
       young_depression_percentage = young_females['Is_depressed'].mean() * 100
       old_depression_percentage = old_females['Is_depressed'].mean() * 100

       # Print the depression percentages
       print(f"Young Females Depression Percentage: {young_depression_percentage}%")
       print(f"Old Females Depression Percentage: {old_depression_percentage}%")
```

```
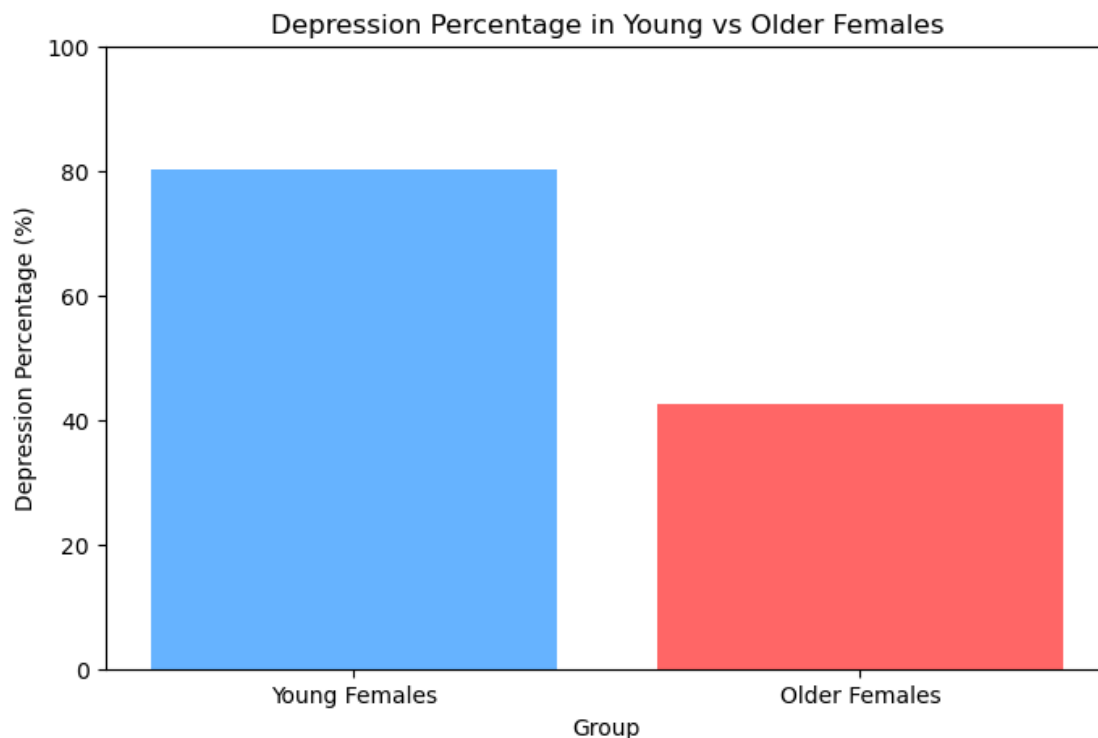Young Females Depression Percentage: 80.27522935779817%
Old Females Depression Percentage: 42.5%
```

**0.0.15** **The question is social media causing depression? Answer is yes. It is impacting on Young female more than older female. According to the Childmind Institute, Teenagers and young female who are spending more time on social media platforms are suffering 13 to 66 percent more with depression compare with the individuals spend less time on social media.**

**0.0.16** **In 2007, smartphones were introduced, limited availability and were expensive. By 2015, many smartphone companies came into the market. The competition is high, prices are affordable. Now a days, individuals are having phone in such a young age. To open an account in social media, age restriction is 12 years and older(many apps). Phone is a very easy source to access social media. However, social media became part of the life for young females' life.**

[74]:
```python
# Plot the depression percentages
plt.figure(figsize=(8, 5))
plt.bar(['Young Females', 'Older Females'], [young_depression_percentage,
 ↪old_depression_percentage], color=['#66b3ff', '#ff6666'])

# Add labels and title
plt.title('Depression Percentage in Young vs Older Females')
plt.ylabel('Depression Percentage (%)')
plt.xlabel('Group')
plt.ylim(0, 100)
plt.show()
```

**0.0.17 According to 19 th news – In young female brain process and emotions part develops faster than the critical thinking and judgement. They react very fast for smaller things. Most adolescent girls respond emotionally when they see something harsh like comments or content. Social media attracts many of us. It is very easy to attract young female in different ways. By nature females have sensitive nature. Single young female use social media for dating. When things are not going well, leads to anxiety and depression.**

**0.0.18 Cyberbullying: Many individuals are facing online bullying. However, young female are targeted. It could be anything by their body features, color, appearance etc. This issue is leading to anxiety and depression, anger and mental health issues.**

```python
[75]: #Creating a "relationship status" dataframe
      relation_status_depressed_female = young_females_depressed['Relationship
       ↪status'].value_counts()
      relation_status_percentage = 100 * relation_status_depressed_female /
       ↪relation_status_depressed_female.sum()
      plt.figure(figsize=(8, 5))
      ax = relation_status_depressed_female.plot(kind='bar', color=['#FFB6C1',
       ↪'#FF69B4', '#FF1493', '#C71585', '#8B008B'])  # Different tones of pink
      plt.title('Relationship Status of Depressed Young Females (18-29)')
      plt.xlabel('Relationship Status')
      plt.ylabel('Count')
      plt.xticks(rotation=90)
      for i, v in enumerate(relation_status_depressed_female):
          plt.text(i, v + 0.5, f'{relation_status_percentage[i]:.1f}%', ha='center',
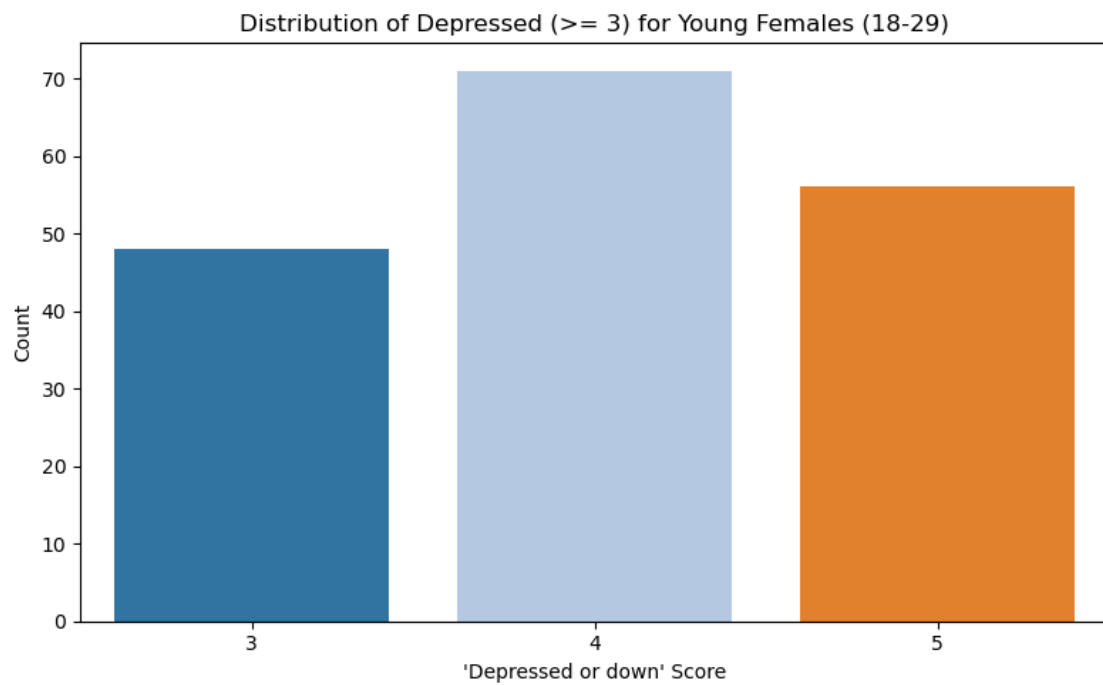       ↪va='bottom', fontsize=10, color='black')
      plt.show()
```

## Relationship Status of Depressed Young Females (18-29)



**0.0.19  As per my last week's research, addiction is one of cause for depression.**

**0.0.20  Addiction: Many apps make using social media easy and flexible. The younger generation always mentions that "I am very relaxed when I am on my phone." Individuals have no track of how many times per day they check the app. This can lead to increased screen time, which can affect sleeping and mental health issues. Making friends online leads to a lack of real-time interactions.**

**0.0.21  After school or work, they spend time on social media, missing play, and mingling with people. Over the time this leads to depression.**

```
[76]:  # Filter data for 'Depressed or down' >= 3
        depressed_data = young_females[young_females['Depressed or down'] >= 3]
        # Plot the data
        plt.figure(figsize=(8, 5))
        sns.countplot(data=depressed_data, x='Depressed or down', palette='tab20')
        plt.title("Distribution of Depressed (>= 3) for Young Females (18-29)")
        plt.xlabel("'Depressed or down' Score")
        plt.ylabel("Count")
```

```
plt.tight_layout()
plt.show()
```



Distribution of Depressed (>= 3) for Young Females (18-29)

Future Investigation: Appearance Comparisons/Body Image Concentration - Sleeping issues, depression or Young vs Old Females Worries