

~\Downloads\final (1).py

```

1  """
2  Personalised Menstrual Tracking and Prediction Using Machine Learning
3  This project uses data science techniques to analyze and predict menstrual cycle patterns. It
4  follows the CRISP-DM methodology and applies Exploratory Data Analysis (EDA), statistical
5  analysis, and machine learning (Linear Regression) to forecast the next cycle start date.
6  """
7
8
9  from google.colab import drive
10 drive.mount('/content/drive')
11
12 import pandas as pd
13 import numpy as np
14 import seaborn as sns
15 import matplotlib.pyplot as plt
16 from scipy import stats
17
18 df = pd.read_excel('/content/drive/MyDrive/google colab files/Menstrual_cycle_tracking.xlsx',
19 na_values=[' ', '', 'NA', 'NaN'])
20 df
21
22
23 """Data Cleaning
24 Replace empty strings with NaN
25 Check missing values in each column
26 """
27
28 df = df.replace(r'^\s*$', pd.NA, regex=True)
29 missing_values = df.isnull().sum()
30 print(missing_values)
31 df.shape
32
33 df.info()
34
35 df.describe()
36
37 df.duplicated().sum() #no duplicates
38
39 columns_to_fill = ['MeanCycleLength', 'MeanMensesLength', 'MeanBleedingIntensity', 'Height', 'Weight', 'Age', 'BMI'] # add more if needed
40
41 df[columns_to_fill] = df.groupby('ClientID')[columns_to_fill].ffill()
42 df
43
44 df['Cycle_Length'] = pd.to_numeric(df['Cycle_Length'], errors='coerce')
45
46 df.to_csv("/content/drive/My Drive/google colab files/cleaned_menstrual_data.csv",
47 index=False) #saved cleaned dataset
48
49 df.shape

```

```
44
45 df.columns
46
47 df['ClientID'].nunique() # Unique Value Counts
48 df['PMS_intensity'].value_counts()
49
50 # Summary Statistics
51 print(df['Cycle_Length'].mean())
52 print(df['Cycle_Length'].median())
53 print(df['Cycle_Length'].std())
54 print(df['Cycle_Length'].min())
55 print(df['Cycle_Length'].max())
56 print(df['Cycle_Length'].mode())
57 Q1 = df['Cycle_Length'].quantile(0.25)
58 Q3 = df['Cycle_Length'].quantile(0.75)
59 IQR = Q3 - Q1
60 outliers = df[(df['Cycle_Length'] < Q1 - 1.5*IQR) | (df['Cycle_Length'] > Q3 + 1.5*IQR)]
61 print(f"Number of outliers: {len(outliers)}")
62
63 ***correlation and covariance**
64
65
66 df.corr(numeric_only=True)
67 df.cov(numeric_only=True)
68
69 df = pd.read_csv("/content/drive/My Drive/google colab files/cleaned_menstrual_data.csv")
70 df['Start_Date'] = pd.to_datetime(df['Start_Date'], errors='coerce')
71 df['start_month'] = df['Start_Date'].dt.month
72 df['start_weekday'] = df['Start_Date'].dt.dayofweek
73 df.to_csv("/content/drive/My Drive/google colab files/updated_menstrual_cycle_data.csv",
74           index=False)
75 df
76
77 df['start_month'].value_counts().sort_index()
78 df['start_weekday'].value_counts()
79
80 """Exploratory Data Analysis (EDA)
81 We explore key features through visualizations and analyze correlations among them.
82 """
83
84 sns.set(style="whitegrid", palette="pastel") # style settings for better visuals
85
86 plt.figure(figsize=(12, 6))
87
88 min_val = int(df['Cycle_Length'].min())
89 max_val = int(df['Cycle_Length'].max())
90 bins = np.arange(min_val - 0.5, max_val + 1.5, 1)
91
92 sns.histplot(df['Cycle_Length'], bins=bins, kde=True, color='lightcoral', edgecolor='black',
93             linewidth=1.2)
```

```

92
93 mean_val = df['Cycle_Length'].mean()
94 median_val = df['Cycle_Length'].median()
95 plt.axvline(mean_val, color='blue', linestyle='--', label=f'Mean: {mean_val:.1f}')
96 plt.axvline(median_val, color='green', linestyle='-.', label=f'Median: {median_val:.1f}')
97
98 plt.title("Distribution of Menstrual Cycle Lengths", fontsize=14)
99 plt.xlabel("Cycle Length (days)", fontsize=12)
100 plt.ylabel("Frequency", fontsize=12)
101 plt.xticks(np.arange(min_val, max_val + 1, 1))
102 plt.legend()
103 plt.grid(axis='y', linestyle='--', alpha=0.5)
104 plt.tight_layout()
105 plt.show()
106
107 """Top 10 Clients with Longest Average Cycle Length"""
108
109 avg_cycle = df.groupby('ClientID')['Cycle_Length'].mean()
110 avg_cycle_df = avg_cycle.reset_index()
111 avg_cycle_df.columns = ['ClientID', 'AverageCycleLength']
112
113 top10 = avg_cycle_df.sort_values(by='AverageCycleLength', ascending=False).head(10)
114
115 plt.figure(figsize=(10, 5))
116 sns.barplot(data=top10, x='ClientID', y='AverageCycleLength', hue="ClientID", legend=False,
117 palette='Blues_r')
118 plt.title("Top 10 Clients with Longest Average Cycle Length")
119 plt.xlabel("Client ID")
120 plt.ylabel("Average Cycle Length (days)")
121 plt.xticks(rotation=30)
122 plt.tight_layout()
123 plt.show()
124
125 client_id = 'nfp8122'
126 user_data = df[df['ClientID'] == client_id].sort_values('Start_Date')
127 plt.figure(figsize=(12, 5))
128 plt.plot(user_data['Start_Date'], user_data['Cycle_Length'], marker='o', linestyle='-',
129 color='purple') # Changed 'start_date' to 'StartOfPeriodDate'
130 plt.title(f"Cycle Length Changes Over Time for Client {client_id}")
131 plt.xlabel("Start Date")
132 plt.ylabel("Cycle Length (days)")
133 plt.grid(True)
134 plt.xticks(rotation=45)
135 plt.tight_layout()
136 plt.show()
137
138 """The top 15 users with the lowest standard deviation in cycle length were identified to
139 evaluate consistency. These users exhibit minimal fluctuations in their menstrual cycle
140 durations, indicating **high regularity**. Such patterns are useful for training initial
141 predictive models, as they represent relatively stable biological cycles."""

```

```
138 client_cycle_std = df.groupby('ClientID')['Cycle_Length'].std().dropna().sort_values()
139 top_clients = client_cycle_std.head(15)
140 colors = sns.color_palette("hls", len(top_clients))
141
142 plt.figure(figsize=(10, 5))
143 top_clients.plot(kind='bar', color=colors)
144 plt.title("Top 15 Most Consistent Users (Lowest Std Dev of Cycle Length)")
145 plt.xlabel("Client ID")
146 plt.ylabel("Cycle Length Std Dev")
147 plt.xticks(rotation=45)
148 plt.grid(axis='y', linestyle='--', alpha=0.5)
149 plt.tight_layout()
150 plt.show()
151
152 """BOX PLOT
153 This plot shows how the menstrual cycle length changes for the top 10 most active users in
154 your dataset.
155 Each vertical box is for one user (Client ID).
156 The line inside each box is the middle value (median) of that user's cycle lengths.
157 The box shows where most of their cycles fall (the middle 50%).
158 The lines outside the box show the full range (except for unusual values).
159 The dots outside the lines are outliers – cycles that were much shorter or longer than usual.
160 The plot shows that some users have very regular cycle lengths, while others have a lot of
161 variation. A few users had periods that were much shorter or longer than normal.
162 """
163
164 top_users = df['ClientID'].value_counts().head(10).index # top 10 users with the most cycle
165 records
166 top_df = df[df['ClientID'].isin(top_users)]
167
168 plt.figure(figsize=(12, 6))
169 sns.boxplot(x='ClientID', y='Cycle_Length', data=top_df, hue='ClientID', palette='Set2',
170 legend=False)
171 plt.title("Cycle Length Variation Across Top 10 Users")
172 plt.xlabel("Client ID")
173 plt.ylabel("Cycle Length (days)")
174 plt.grid(axis='y', linestyle='--', alpha=0.6)
175 plt.tight_layout()
176 plt.show()
177
178 """
179 Box Plot: Most Irregular Users (Highest Std Dev)
180 This plot shows the top 10 most irregular users based on cycle length variability.
181 Most of them show wide variation, and many have outliers, indicating inconsistent cycles.
182 This suggests that a one-size-fits-all prediction model may not be reliable and personalized
183 predictions could be more effective.
184 """
185
186 client_cycle_std = df.groupby('ClientID')['Cycle_Length'].std().dropna()
187
188 irregular_users = client_cycle_std.sort_values(ascending=False).head(10).index
```

```

184 irregular_df = df[df['ClientID'].isin(irregular_users)]
185
186 plt.figure(figsize=(12, 6)) #box plot
187 sns.boxplot(x='ClientID', y='Cycle_Length', data=irregular_df, hue='ClientID',
188 palette='Set3', legend=False)
189 plt.title("Cycle Length Variation Across Top 10 Most Irregular Users")
190 plt.xlabel("Client ID")
191 plt.ylabel("Cycle Length (days)")
192 plt.grid(axis='y', linestyle='--', alpha=0.6)
193 plt.tight_layout()
194 plt.show()
195
196 """HEATMAP
197 Does stress level increase or decrease cycle length?
198 Is PMS_intensity linked with bleeding intensity?
199 Are any features strongly negatively or positively correlated?
200 """
201 selected_cols = ['LengthofCycle', 'LengthofMenses', 'stress_level',
202 'PMS_intensity', 'MeanBleedingIntensity', 'TotalNumberofHighDays', 'Age']
203 filtered_df = df[selected_cols].dropna()
204 corr = filtered_df.corr()
205
206 plt.figure(figsize=(8, 6)) #heatmap
207 sns.heatmap(corr, annot=True, cmap='coolwarm', fmt='.2f', square=True, linewidths=0.5)
208 plt.title("Correlation Heatmap: Key Menstrual Features")
209 plt.xticks(rotation=45)
210 plt.yticks(rotation=0)
211 plt.tight_layout()
212 plt.show()
213
214 """Hypothesis **Test**
215 **✅ T-Test:** Does PMS Intensity Affect Cycle Length?
216 Test whether there's a significant difference in the average cycle length between:
217 Users with low PMS intensity (e.g., PMS ≤ 5)
218 Users with high PMS intensity (e.g., PMS > 5)
219 """
220 from scipy.stats import ttest_ind
221 low_pms = df[df['PMS_intensity'] <= 5]['LengthofCycle'].dropna()
222 high_pms = df[df['PMS_intensity'] > 5]['LengthofCycle'].dropna()
223 print("Mean cycle length (Low PMS):", low_pms.mean())
224 print("Mean cycle length (High PMS):", high_pms.mean())
225 t_stat, p_value = ttest_ind(low_pms, high_pms, equal_var=False) #ttest
226 print("T-statistic:", t_stat)
227 print("P-value:", p_value)
228 if p_value < 0.05:
229     print("✅ Result: Significant difference in cycle length between low and high PMS
230 groups.")
231 else:
232     print("❌ Result: No significant difference in cycle length between the two groups.")

```

```

232
233 """data suggests that PMS intensity doesn't significantly affect the length of the menstrual
cycle
234 ✅ **ANOVA Test**
235 To determine whether menstrual cycle length significantly varies across age groups.
236 Test Performed:
237 One-Way ANOVA
238 Groups: 20-29, 30-39, 40-49
239 Variable analyzed: LengthofCycle
240 """
241
242 bins = [20, 30, 40, 50]
243 labels = ['20-29', '30-39', '40-49']
244 df['AgeGroup'] = pd.cut(df['Age'], bins=bins, labels=labels, right=False)
245 group_counts = df.groupby('AgeGroup', observed=False)['LengthofCycle'].count()
246 valid_groups = group_counts[group_counts >= 2].index
247 df_valid = df[df['AgeGroup'].isin(valid_groups)]
248 print("Groups included in ANOVA:", df_valid['AgeGroup'].unique())
249 anova_groups = [group['LengthofCycle'].dropna() for name, group in
df_valid.groupby('AgeGroup', observed=False) if len(group) > 1]
250 f_stat, p_value = f_oneway(*anova_groups) # ANOVA
251
252 print("F-statistic:", f_stat) # result
253 print("P-value:", p_value)
254 if p_value < 0.05:
255     print("✅ Significant difference in cycle length across age groups.")
256 else:
257     print("❌ No significant difference in cycle length across age groups.")
258
259 """📊 Results:
260 F-statistic: 4.10
261 P-value: 0.0167
262 Cycle length is significantly affected by age. Users in different age groups may experience
different menstrual patterns.
263 ✅ **Shapiro-Wilk Test:**
264 To see if the LengthofCycle column is normally distributed, which is helpful before applying
regression or other parametric models.
265 """
266
267 from scipy.stats import shapiro
268 cycle_lengths = df['LengthofCycle'].dropna()
269 stat, p_value = shapiro(cycle_lengths)
270 print("Shapiro-Wilk Test Statistic:", stat)
271 print("P-value:", p_value)
272
273 # Interpretation
274 if p_value < 0.05:
275     print("❌ Data is NOT normally distributed (reject H0).")
276 else:
277     print("✅ Data IS normally distributed (fail to reject H0).")
278

```

```
279 """MODEL
280 We implemented and compared two regression models – Linear Regression (aligned with syllabus)
    and Random Forest (real-world extension) – to predict LengthofCycle using lifestyle and
    health features. After preprocessing and splitting the data, both models were trained and
    evaluated using R2 score, MAE, and RMSE. This step aligns with the CRISP-DM framework's
    modeling and evaluation phases.
281 """
282
283 from sklearn.model_selection import GroupShuffleSplit
284 from sklearn.ensemble import RandomForestRegressor
285 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
286
287 df = pd.read_csv("/content/drive/MyDrive/google colab files/updated_menstrual_cycle_data.csv")
288 df['Start_Date'] = pd.to_datetime(df['Start_Date'], errors='coerce')
289 df = df.sort_values(by=['ClientID', 'Start_Date'])
290 df['Next_Start_Date'] = df.groupby('ClientID')['Start_Date'].shift(-1)
291 df['Target_DaysToNextCycle'] = (df['Next_Start_Date'] - df['Start_Date']).dt.days
292 df = df.dropna(subset=['Target_DaysToNextCycle'])
293 mapping_dicts = {
294     'exercise': {'Cardio': 0, 'Strength': 1, 'Yoga': 2},
295     'diet_quality': {'Poor': 0, 'Average': 1, 'Good': 2},
296     'spotting': {'No': 0, 'Yes': 1},
297     'sexual_activity': {'No': 0, 'Yes': 1},
298     'sleep_quality': {'Poor': 0, 'Average': 1, 'Good': 2}
299 }
300 for col, mapping in mapping_dicts.items():
301     df[col] = df[col].map(mapping).fillna(-1).astype(int)
302 features = [
303     'LengthofCycle', 'LengthofMenses',
304     'MeanCycleLength', 'MeanMensesLength', 'BMI',
305     'PMS_intensity', 'EstimatedDayofOvulation',
306     'MeanBleedingIntensity', 'TotalDaysofFertility',
307     'TotalMensesScore', 'Age', 'sleep_quality',
308     'stress_level', 'exercise', 'diet_quality',
309     'spotting', 'sexual_activity'
310 ]
311 df = df.dropna(subset=features + ['Target_DaysToNextCycle'])
312 X = df[features]
313 y = df['Target_DaysToNextCycle']
314 groups = df['ClientID']
315
316 splitter = GroupShuffleSplit(test_size=0.2, n_splits=1, random_state=42)
317 train_idx, test_idx = next(splitter.split(X, y, groups=groups))
318
319 X_train, X_test = X.iloc[train_idx], X.iloc[test_idx]
320 y_train, y_test = y.iloc[train_idx], y.iloc[test_idx]
321
322 # ✅ Train model
323 rf_model = RandomForestRegressor(random_state=42)
324 rf_model.fit(X_train, y_train)
```

```

325 y_pred_rf = rf_model.predict(X_test)
326
327 # ✅ Evaluate model
328 mse_rf = mean_squared_error(y_test, y_pred_rf)
329 mae_rf = mean_absolute_error(y_test, y_pred_rf)
330 r2_rf = r2_score(y_test, y_pred_rf)
331
332 print("📊 Random Forest Evaluation:")
333 print(f"Mean Squared Error (MSE): {mse_rf:.2f}")
334 print(f"Mean Absolute Error (MAE): {mae_rf:.2f}")
335 print(f"R² Score: {r2_rf:.2f}")
336
337 """# ✅ Visualize
338 Actual vs Predicted Days (Random Forest
339 """
340
341 plt.figure(figsize=(8, 6))
342 sns.scatterplot(x=y_test, y=y_pred_rf, color='darkorange', alpha=0.6)
343 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], '--r', label='Perfect
344 Prediction Line')
345 plt.xlabel("Actual Days to Next Cycle")
346 plt.ylabel("Predicted Days")
347 plt.title("Actual vs Predicted Days (Random Forest)")
348 plt.legend()
349 plt.grid(True)
350 plt.tight_layout()
351 plt.show()
352
353 """✅ Predict next cycle for latest entry"""
354
355 target_client = 'nfp8237'
356 client_data = df[df['ClientID'] == target_client].sort_values('Start_Date')
357 latest_entry = client_data.iloc[[-1]]
358 last_start_date = latest_entry['Start_Date'].values[0]
359 latest_features = latest_entry[features]
360
361 predicted_days = rf_model.predict(latest_features)[0]
362
363 predicted_next_date = pd.to_datetime(last_start_date) +
364 pd.Timedelta(days=int(predicted_days))
365
366 # 📌 Output
367 print("🔮 Prediction for Client:", target_client)
368 print("📅 Last Known Cycle Start Date:", pd.to_datetime(last_start_date).date())
369 print(f"🕒 Predicted Days to Next Cycle: {int(predicted_days)} days")
370 print("📅 Predicted Next Cycle Start Date:", predicted_next_date.date())
371
372 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score

```



```
373
374 # y_test: actual values, y_pred_rf: predicted values
375 mae = mean_absolute_error(y_test, y_pred_rf)
376 mse = mean_squared_error(y_test, y_pred_rf)
377 rmse = np.sqrt(mse)
378 r2 = r2_score(y_test, y_pred_rf)
379
380 # Custom accuracy: within ±2 days
381 tolerance = 2
382 within_tolerance = np.abs(y_test - y_pred_rf) <= tolerance
383 custom_accuracy = within_tolerance.mean() * 100
384
385 print("📊 Evaluation Metrics:")
386 print(f"MAE(Mean Absolute Error): {mae:.2f} days")
387 print(f"MSE(Mean Squared Error): {mse:.2f}")
388 print(f"RMSE(Root Mean Squared Error): {rmse:.2f} days")
389 print(f"R² Score: {r2:.2f}")
390 print(f"Custom Accuracy (±{tolerance} days): {custom_accuracy:.2f}%")
```