

Problem for Covid - 19 Data Analysis

Project using Python

```

1 # covid_analysis.py
2 # Python 3.8+
3 # Required packages: pandas, numpy, matplotlib, seaborn
4 # Install if needed: pip install pandas numpy matplotlib seaborn
5
6 import pandas as pd
7 import numpy as np
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10
11 # ===== 1. Import dataset =====
12 url = "https://raw.githubusercontent.com/SR1608/Datasets/main/covid-data.csv"
13 df = pd.read_csv(url)
14
15 # ===== 2. High Level Data Understanding =====
16 print("== Shape (rows, columns) ==")
17 print(df.shape) # 2.a
18
19 print("\n== Data types ==")
20 print(df.dtypes) # 2.b
21
22 print("\n== Info() ==")
23 print(df.info()) # 2.b (info)pip
24
25 print("\n== Describe() ==")
26 print(df.describe(include='all')) # 2.c
27
28 # ===== 3. Low Level Data Understanding =====
29 # 3.a unique count in 'location'
30 if 'location' in df.columns:
31     print("\nUnique locations count:", df['location'].nunique())
32     print("Sample unique locations:", df['location'].unique()[:10])
33
34 # 3.b which continent has maximum frequency (value_counts)
35 if 'continent' in df.columns:
36     print("\nContinent frequency (value_counts):")
37     print(df['continent'].value_counts())
38     most_freq_continent = df['continent'].value_counts().idxmax()
39     print("Continent with maximum frequency:", most_freq_continent)
40
41 # 3.c max & mean of 'total_cases'
42 if 'total_cases' in df.columns:
43     print("\nTotal_cases max:", df['total_cases'].max())
44     print("Total_cases mean:", df['total_cases'].mean())
45
46 # 3.d quartiles of 'total_deaths'
47 if 'total_deaths' in df.columns:
48     print("\nTotal_deaths quantiles (25%,50%,75%):")
49     print(df['total_deaths'].quantile([0.25, 0.5, 0.75]))
50

```

```

51 # 3.e continent with maximum 'human_development_index'
52 if {'continent', 'human_development_index'}.issubset(df.columns):
53     # groupby continent and take idxmax or max
54     hdi_by_cont = df.groupby('continent')['human_development_index'].max().dropna()
55     print("\nMax human_development_index by continent:")
56     print(hdi_by_cont)
57     if not hdi_by_cont.empty:
58         continent_max_hdi = hdi_by_cont.idxmax()
59         print("Continent with maximum human_development_index:", continent_max_hdi)
60
61 # 3.f continent with minimum 'gdp_per_capita'
62 if {'continent', 'gdp_per_capita'}.issubset(df.columns):
63     gdp_by_cont = df.groupby('continent')['gdp_per_capita'].min().dropna()
64     print("\nMin gdp_per_capita by continent:")
65     print(gdp_by_cont)
66     if not gdp_by_cont.empty:
67         continent_min_gdp = gdp_by_cont.idxmin()
68         print("Continent with minimum gdp_per_capita:", continent_min_gdp)
69
70 # ===== 4. Filter dataframe to required columns and update df =====
71 cols_keep =
72     ['continent', 'location', 'date', 'total_cases', 'total_deaths', 'gdp_per_capita', 'human_development_index']
73 missing_cols = [c for c in cols_keep if c not in df.columns]
74 if missing_cols:
75     raise ValueError(f"The dataset is missing required columns: {missing_cols}")
76
77 df = df[cols_keep].copy()
78 print("\nFiltered dataframe shape:", df.shape)
79
80 # ===== 5. Data Cleaning =====
81 # 5.a Remove duplicates
82 before_dup = df.shape[0]
83 df = df.drop_duplicates()
84 after_dup = df.shape[0]
85 print(f"\nRemoved duplicates: {before_dup - after_dup}")
86
87 # 5.b Missing values in all columns
88 print("\nMissing values (per column):")
89 print(df.isnull().sum())
90
91 # 5.c Remove observations where continent is missing (dropna subset)
92 before_drop_cont = df.shape[0]
93 df = df.dropna(subset=['continent'])
94 after_drop_cont = df.shape[0]
95 print(f"\nRows removed due to missing continent: {before_drop_cont - after_drop_cont}")
96
97 # 5.d Fill all remaining missing values with 0
98 df = df.fillna(0)
99 print("\nMissing values after fillna(0):")
100 print(df.isnull().sum())
101
102 # ===== 6. Date time format =====

```

```

102 # 6.a convert date column to datetime
103 df['date'] = pd.to_datetime(df['date'], errors='coerce')
104
105 # if coercion produced NaT, you may want to check them:
106 n_invalid_dates = df['date'].isna().sum()
107 print(f"\nNumber of invalid/NaT dates after to_datetime: {n_invalid_dates}")
108
109 # 6.b Create new column 'month' extracted from date
110 df['month'] = df['date'].dt.month.fillna(0).astype(int) # 0 means invalid/missing date
111
112 # ===== 7. Data Aggregation =====
113 # 7.a Find max value in all columns using groupby on 'continent'
114 df_groupby = df.groupby('continent').max(numeric_only=False).reset_index()
115 # Note: groupby.max will take max for columns; for non-numeric like date it takes max
116 # lexicographically / latest date.
117 print("\n==== df_groupby (max per continent) ===")
118 print(df_groupby.head())
119
120 # ===== 8. Feature Engineering =====
121 # 8.a Create total_deaths_to_total_cases ratio
122 # Work on df_groupby (as requested further analysis to use df_groupby)
123 # Ensure numeric type:
124 df_groupby['total_cases'] = pd.to_numeric(df_groupby['total_cases'], errors='coerce').fillna(0)
125 df_groupby['total_deaths'] = pd.to_numeric(df_groupby['total_deaths'], errors='coerce').fillna(0)
126
127 # compute ratio safely (avoid division by zero)
128 df_groupby['total_deaths_to_total_cases'] = np.where(
129     df_groupby['total_cases'] == 0,
130     0,
131     df_groupby['total_deaths'] / df_groupby['total_cases']
132 )
133
134 print("\nFeature 'total_deaths_to_total_cases' added:")
135 print(df_groupby[['continent', 'total_deaths', 'total_cases', 'total_deaths_to_total_cases']])
136
137 # ===== 9. Data Visualization =====
138 sns.set(style='whitegrid')
139
140 # 9.a Univariate analysis on gdp_per_capita: histogram (seaborn)
141 plt.figure(figsize=(8,5))
142 # convert to numeric and drop zeros if you want to avoid 0-filled missing values dominating
143 gdp_series = pd.to_numeric(df_groupby['gdp_per_capita'], errors='coerce').dropna()
144 sns.histplot(gdp_series, kde=True)
145 plt.title('Distribution of gdp_per_capita (df_groupby)')
146 plt.xlabel('gdp_per_capita')
147 plt.tight_layout()
148 plt.show()
149
150 # 9.b Scatter plot of total_cases & gdp_per_capita (from df_groupby)
151 plt.figure(figsize=(8,5))
152 x = df_groupby['gdp_per_capita'].astype(float)
153 y = df_groupby['total_cases'].astype(float)

```

```
153 sns.scatterplot(x=x, y=y)
154 plt.xlabel('gdp_per_capita')
155 plt.ylabel('total_cases')
156 plt.title('Total Cases vs GDP per Capita (df_groupby)')
157 plt.tight_layout()
158 plt.show()
159
160 # 9.c Pairplot on df_groupby dataset
161 # For pairplot, choose numeric columns only to keep it readable
162 numeric_cols = df_groupby.select_dtypes(include=[np.number]).columns.tolist()
163 if len(numeric_cols) >= 2:
164     sns.pairplot(df_groupby[numeric_cols])
165     plt.suptitle('Pairplot of numeric features in df_groupby', y=1.02)
166     plt.show()
167 else:
168     print("Not enough numeric columns for pairplot.")
169
170 # 9.d Bar plot of 'continent' vs 'total_cases'
171 plt.figure(figsize=(10,6))
172 # Use seaborn catplot (kind='bar') as suggested:
173 sns.catplot(kind='bar', x='continent', y='total_cases', data=df_groupby, height=5, aspect=1.6)
174 plt.title('Total Cases by Continent (df_groupby)')
175 plt.tight_layout()
176 plt.show()
177
178 # ===== 10. Save df_groupby to local drive =====
179 out_filename = "df_groupby.csv"
180 df_groupby.to_csv(out_filename, index=False)
181 print(f"\nSaved df_groupby to: {out_filename}")
182
183 # Optionally, also save the cleaned and filtered df
184 df.to_csv("covid_cleaned_filtered.csv", index=False)
185 print("Saved cleaned filtered dataframe to: covid_cleaned_filtered.csv")
186
```