**sol.py**

```python
# ============================================================
# Aadhaar Analytics — Full Graph Generation Code (Attractive)
# Covers: univariate, bivariate, trivariate, multivariate visuals
# Uses: pandas, numpy, matplotlib, scikit-learn
# NOTE: No seaborn, and no manual color setting.
# ============================================================

import os
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans

# ------------------------
# CONFIG
# ------------------------
ENROL_PATH = "api_data_aadhar_enrolment_0_1006029.csv"
DEMO_PATH  = "api_data_aadhar_demographic_0_2071700.csv"
BIO_PATH   = "api_data_aadhar_biometric_0_1861108.csv"

SAVE_FIGS = False
OUTDIR = "/mnt/data/aadhaar_figs"
os.makedirs(OUTDIR, exist_ok=True)

plt.rcParams["figure.figsize"] = (12, 5)
plt.rcParams["axes.grid"] = True

def savefig(name: str):
    if SAVE_FIGS:
        fp = os.path.join(OUTDIR, name)
        plt.savefig(fp, dpi=220, bbox_inches="tight")
        print("Saved:", fp)

# ------------------------
# LOAD + CLEAN
# ------------------------
def load_uidai_csv(path: str) -> pd.DataFrame:
    df = pd.read_csv(path)
    if "Unnamed: 0" in df.columns:
        df = df.drop(columns=["Unnamed: 0"])
    df.columns = [c.strip() for c in df.columns]
    if "date" not in df.columns:
        raise ValueError(f"Missing 'date' column in {path}")
    df["date"] = pd.to_datetime(df["date"], dayfirst=True, errors="coerce")
    # strip common geo fields if present
    for c in ["state", "district", "pincode"]:
        if c in df.columns:
            df[c] = df[c].astype(str).str.strip()
```

```python
52        return df
53
54   enrol = load_uidai_csv(ENROL_PATH)
55   demo  = load_uidai_csv(DEMO_PATH)
56   bio   = load_uidai_csv(BIO_PATH)
57
58   # ------------------------
59   # FEATURE ENGINEERING
60   # ------------------------
61   # totals
62   enrol["total_enrol"] = enrol[["age_0_5", "age_5_17",
     "age_18_greater"]].apply(pd.to_numeric, errors="coerce").fillna(0).sum(axis=1)
63
64   # demo/bio columns vary a bit; sum all columns that start with demo_ / bio_
65   demo_cols = [c for c in demo.columns if c.startswith("demo_")]
66   bio_cols  = [c for c in bio.columns  if c.startswith("bio_")]
67
68   demo["total_demo"] = demo[demo_cols].apply(pd.to_numeric,
     errors="coerce").fillna(0).sum(axis=1)
69   bio["total_bio"]   = bio[bio_cols].apply(pd.to_numeric,
     errors="coerce").fillna(0).sum(axis=1)
70
71   eps = 1e-9
72
73   # ------------------------
74   # INDIA DAILY TABLE
75   # ------------------------
76   daily = (
77       enrol.groupby("date")["total_enrol"].sum().to_frame("Enrolments")
78       .join(demo.groupby("date")["total_demo"].sum())
79       .join(bio.groupby("date")["total_bio"].sum())
80       .fillna(0)
81   )
82   daily.rename(columns={"total_demo": "Demographic Updates", "total_bio": "Biometric
     Updates"}, inplace=True)
83   daily["Total Updates"] = daily["Demographic Updates"] + daily["Biometric Updates"]
84   daily["Total Activity"] = daily["Enrolments"] + daily["Total Updates"]
85   daily["Enrol_7d_avg"] = daily["Enrolments"].rolling(7, min_periods=3).mean()
86   daily["Upd_7d_avg"]   = daily["Total Updates"].rolling(7, min_periods=3).mean()
87   daily["Act_7d_avg"]   = daily["Total Activity"].rolling(7, min_periods=3).mean()
88   daily["is_month_start"] = daily.index.is_month_start
89
90   # ============================================================
91   # 1) Time-series: Daily enrolments + demo + bio updates
92   # (supports: update-dominant lifecycle insight)
93   # ============================================================
94   plt.figure()
95   plt.plot(daily.index, daily["Enrolments"], label="Enrolments")
96   plt.plot(daily.index, daily["Demographic Updates"], label="Demographic Updates")
97   plt.plot(daily.index, daily["Biometric Updates"], label="Biometric Updates")
98   plt.title("India — Daily Aadhaar Enrolments and Updates")
99   plt.xlabel("Date")
100  plt.ylabel("Transactions")
101  plt.legend()
```

```python
102  plt.tight_layout()
103  savefig("01_india_daily_enrol_demo_bio.png")
104  plt.show()
105
106  # ================================================================
107  # 2) Time-series: Enrolments vs Total Updates vs Total Activity
108  # ================================================================
109  plt.figure()
110  plt.plot(daily.index, daily["Enrolments"], label="Enrolments")
111  plt.plot(daily.index, daily["Total Updates"], label="Total Updates")
112  plt.plot(daily.index, daily["Total Activity"], label="Total Activity")
113  plt.title("India — Enrolments vs Updates vs Total Activity")
114  plt.xlabel("Date")
115  plt.ylabel("Transactions")
116  plt.legend()
117  plt.tight_layout()
118  savefig("02_india_enrol_vs_updates_vs_activity.png")
119  plt.show()
120
121  # ================================================================
122  # 3) Rolling averages (trend vs noise)
123  # ================================================================
124  plt.figure()
125  plt.plot(daily.index, daily["Enrol_7d_avg"], label="Enrolments (7d avg)")
126  plt.plot(daily.index, daily["Upd_7d_avg"], label="Updates (7d avg)")
127  plt.plot(daily.index, daily["Act_7d_avg"], label="Activity (7d avg)")
128  plt.title("India — 7-Day Rolling Averages (Trend)")
129  plt.xlabel("Date")
130  plt.ylabel("Transactions")
131  plt.legend()
132  plt.tight_layout()
133  savefig("03_india_rolling_7d.png")
134  plt.show()
135
136  # ================================================================
137  # 4) Annotated: Month-start enrolment spikes
138  # ================================================================
139  plt.figure()
140  plt.plot(daily.index, daily["Enrolments"], label="Daily Enrolments")
141  ms = daily[daily["is_month_start"]]
142  plt.scatter(ms.index, ms["Enrolments"], label="1st of Month")
143  plt.title("India — Month-start Enrolment Spikes (Administrative Cycle)")
144  plt.xlabel("Date")
145  plt.ylabel("Enrolments")
146  plt.legend()
147  plt.tight_layout()
148  savefig("04_month_start_spikes.png")
149  plt.show()
150
151  # ================================================================
152  # 5) Anomaly plot (global z-score)
153  # (supports: anomaly periods for operational interpretation)
154  # ================================================================
155  s = daily["Enrolments"].astype(float)
```

```python
156  mu, sigma = s.mean(), s.std(ddof=0)
157  sigma = sigma if (sigma and not np.isnan(sigma)) else 1.0
158  z = (s - mu) / sigma
159  anom = z.abs() >= 3.0
160
161  plt.figure()
162  plt.plot(daily.index, s, label="Daily Enrolments")
163  plt.scatter(daily.index[anom], s[anom], label="Anomaly (|z|≥3)")
164  plt.title("India — Enrolment Anomalies ")
165  plt.xlabel("Date")
166  plt.ylabel("Enrolments")
167  plt.legend()
168  plt.tight_layout()
169  savefig("05_anomaly_global_z.png")
170  plt.show()
171
172  # ============================================================
173  # 6) Age-wise enrolments (bar)
174  # ============================================================
175  age_totals = enrol[["age_0_5", "age_5_17", "age_18_greater"]].apply(pd.to_numeric,
         errors="coerce").fillna(0).sum()
176  plt.figure()
177  plt.bar(["0-5", "5-17", "18+"], age_totals.values)
178  plt.title("India — Age-wise Aadhaar Enrolments")
179  plt.xlabel("Age Group")
180  plt.ylabel("Total Enrolments")
181  plt.tight_layout()
182  savefig("06_agewise_enrol_bar.png")
183  plt.show()
184
185  # ============================================================
186  # 7) Age-wise share (pie)
187  # ============================================================
188  plt.figure()
189  plt.pie(age_totals.values, labels=["0-5", "5-17", "18+"], autopct="%1.1f%%")
190  plt.title("India — Age-wise Aadhaar Enrolment Share")
191  plt.tight_layout()
192  savefig("07_agewise_enrol_pie.png")
193  plt.show()
194
195  # ============================================================
196  # STATE-LEVEL TABLES (for regional variation, burden, volatility)
197  # ============================================================
198  state = (
199      enrol.groupby("state")["total_enrol"].sum().to_frame("Enrol")
200      .join(demo.groupby("state")["total_demo"].sum())
201      .join(bio.groupby("state")["total_bio"].sum())
202      .fillna(0)
203  )
204  state.rename(columns={"total_demo": "Demo", "total_bio": "Bio"}, inplace=True)
205  state["Updates"] = state["Demo"] + state["Bio"]
206  state["Total Activity"] = state["Enrol"] + state["Updates"]
207  state["Update Burden Index"] = state["Updates"] / (state["Enrol"] + eps)
208
```

```python
209   # ================================================================
210   # 8) Top states by total enrolments (bar ranking)
211   # ================================================================
212   top_enrol_states = state.sort_values("Enrol", ascending=False).head(10)
213   plt.figure()
214   plt.bar(top_enrol_states.index, top_enrol_states["Enrol"])
215   plt.title("Top 10 States — Total Aadhaar Enrolments")
216   plt.xlabel("State")
217   plt.ylabel("Enrolments")
218   plt.xticks(rotation=45, ha="right")
219   plt.tight_layout()
220   savefig("08_top10_states_enrol.png")
221   plt.show()
222
223   # ================================================================
224   # 9) Top states: share of total Aadhaar activity (pie)
225   # ================================================================
226   top5_activity = state.sort_values("Total Activity", ascending=False).head(5)
227   plt.figure()
228   plt.pie(top5_activity["Total Activity"].values, labels=top5_activity.index,
          autopct="%1.1f%%")
229   plt.title("Top 5 States — Share of Total Aadhaar Activity")
230   plt.tight_layout()
231   savefig("09_top5_state_activity_share.png")
232   plt.show()
233
234   # ================================================================
235   # 10) Update Burden Index (bar ranking)
236   # ================================================================
237   top_burden = state.sort_values("Update Burden Index", ascending=False).head(10)
238   plt.figure()
239   plt.bar(top_burden.index, top_burden["Update Burden Index"])
240   plt.title("Top 10 States — Update Burden Index (Updates / Enrolments)")
241   plt.xlabel("State")
242   plt.ylabel("Updates per Enrolment")
243   plt.xticks(rotation=45, ha="right")
244   plt.tight_layout()
245   savefig("10_top10_update_burden.png")
246   plt.show()
247
248   # ================================================================
249   # 11) Volatility Score (std/mean of DAILY activity at state level)
250   # ================================================================
251   # build state-day activity (enrol + updates) to compute volatility
252   state_day_enrol = enrol.groupby(["state", "date"])["total_enrol"].sum().reset_index()
253   state_day_demo  = demo.groupby(["state", "date"])["total_demo"].sum().reset_index()
254   state_day_bio   = bio.groupby(["state", "date"])["total_bio"].sum().reset_index()
255
256   state_day = (
257       state_day_enrol.merge(state_day_demo, on=["state","date"], how="left")
258                      .merge(state_day_bio,  on=["state","date"], how="left")
259                      .fillna(0)
260   )
```

```python
261   state_day["activity"] = state_day["total_enrol"] + state_day["total_demo"] +
      state_day["total_bio"]
262
263   vol = state_day.groupby("state")["activity"].agg(["mean","std"]).fillna(0)
264   vol["Volatility Score"] = vol["std"] / (vol["mean"] + eps)
265
266   top_vol = vol.sort_values("Volatility Score", ascending=False).head(10)
267   plt.figure()
268   plt.bar(top_vol.index, top_vol["Volatility Score"])
269   plt.title("Top 10 States — Volatility Score (Std/Mean of Daily Activity)")
270   plt.xlabel("State")
271   plt.ylabel("Volatility Score")
272   plt.xticks(rotation=45, ha="right")
273   plt.tight_layout()
274   savefig("11_top10_volatility.png")
275   plt.show()
276
277   # ============================================================
278   # 12) Stacked bars: Top states × age group (trivariate)
279   # ============================================================
280   state_age = enrol.groupby("state")[["age_0_5","age_5_17","age_18_greater"]].apply(
281       lambda x: x.apply(pd.to_numeric, errors="coerce").fillna(0).sum()
282   )
283   state_age["total"] = state_age.sum(axis=1)
284   state_age = state_age.sort_values("total", ascending=False).head(10)
285
286   x = np.arange(len(state_age.index))
287   plt.figure(figsize=(12,5))
288   plt.bar(x, state_age["age_0_5"].values, label="0-5")
289   plt.bar(x, state_age["age_5_17"].values, bottom=state_age["age_0_5"].values, label="5-17")
290   plt.bar(x, state_age["age_18_greater"].values, bottom=
      (state_age["age_0_5"]+state_age["age_5_17"]).values, label="18+")
291   plt.title("Top 10 States — Enrolment by Age Group ")
292   plt.xlabel("State")
293   plt.ylabel("Enrolments")
294   plt.xticks(x, state_age.index, rotation=45, ha="right")
295   plt.legend()
296   plt.tight_layout()
297   savefig("12_stacked_age_top10.png")
298   plt.show()
299
300   # ============================================================
301   # 13) Trivariate heatmap: State × Month × Update/Enrol ratio
302   # (use Total Updates (demo+bio) / Enrolments)
303   # ============================================================
304   enrol_m = enrol.copy()
305   demo_m  = demo.copy()
306   bio_m   = bio.copy()
307
308   enrol_m["month"] = enrol_m["date"].dt.to_period("M").astype(str)
309   demo_m["month"]  = demo_m["date"].dt.to_period("M").astype(str)
310   bio_m["month"]   = bio_m["date"].dt.to_period("M").astype(str)
311
312   em = enrol_m.groupby(["state","month"])["total_enrol"].sum().reset_index()
```

```python
313   dm = demo_m.groupby(["state","month"])["total_demo"].sum().reset_index()
314   bm = bio_m.groupby(["state","month"])["total_bio"].sum().reset_index()
315
316   sm = em.merge(dm, on=["state","month"], how="left").merge(bm, on=["state","month"],
      how="left").fillna(0)
317   sm["updates"] = sm["total_demo"] + sm["total_bio"]
318   sm["ratio"] = sm["updates"] / (sm["total_enrol"] + eps)
319
320   top_states = sm.groupby("state")
      ["total_enrol"].sum().sort_values(ascending=False).head(12).index
321   pivot = sm[sm["state"].isin(top_states)].pivot(index="state", columns="month",
      values="ratio").fillna(0)
322
323   plt.figure(figsize=(14,6))
324   plt.imshow(pivot.values, aspect="auto")
325   plt.title("Heatmap — State × Month Update/Enrolment Ratio (Top States)")
326   plt.xlabel("Month")
327   plt.ylabel("State")
328   plt.xticks(range(len(pivot.columns)), pivot.columns, rotation=45, ha="right")
329   plt.yticks(range(len(pivot.index)), pivot.index)
330   plt.colorbar(label="Updates per Enrolment")
331   plt.tight_layout()
332   savefig("13_heatmap_state_month_ratio.png")
333   plt.show()
334
335   # ==============================================================
336   # 14) District-level analysis: Top districts by enrolments (bar)
337   # (supports: rural/semi-urban participation)
338   # NOTE: uses districts available in dataset; interpret with context.
339   # ==============================================================
340   if "district" in enrol.columns:
341       dist = enrol.groupby(["state","district"])["total_enrol"].sum().reset_index()
342       # take top 15 overall districts
343       dist_top = dist.sort_values("total_enrol", ascending=False).head(15)
344       labels = dist_top["state"] + " — " + dist_top["district"]
345
346       plt.figure(figsize=(12,6))
347       plt.bar(labels, dist_top["total_enrol"].values)
348       plt.title("Top 15 Districts — Total Aadhaar Enrolments")
349       plt.xlabel("District")
350       plt.ylabel("Enrolments")
351       plt.xticks(rotation=60, ha="right")
352       plt.tight_layout()
353       savefig("14_top15_districts_enrol.png")
354       plt.show()
355
356
```