

Authors

Group 3 Members

1. Michael Mumina (Scrum Master)
2. Sharon Nyakeya
3. Bryan Njogu
4. Ashley Kibwogo
5. Tanveer Chege
6. Claris Wangari
7. Priscillah Giriama

Sudan Food Insecurity and Displacement Model

This analysis delivers a decision-ready prototype for humanitarian early warning, providing forward-looking identification of food insecurity deterioration, displacement pressure, and key risk drivers at relevant administrative levels.

The model integrates multi-sector signals into a coherent risk framework that supports prioritization and situational awareness.

While advanced components such as probabilistic uncertainty bounds, formal hotspot clustering, and response-rule optimization are not yet operationalized, the current system establishes a robust analytical backbone that can be readily extended into a fully operational decision-support tool.



Point of Reference: IPC Sudan Acute Food Insecurity Snapshot (Oct 2024 – May 2025):
https://www.ipcinfo.org/fileadmin/user_upload/ipcinfo/docs/IPC_Sudan_Acute_Food_Insecurity_Oct2024.pdf

2) Business Understanding

2.1 Problem Scope:

The prototype is immediately usable for strategic planning and anticipatory action, enabling decision-makers to:

- Identify where conditions are likely to deteriorate in the coming month
- Compare displacement pressure across states and localities
- Understand the primary drivers contributing to elevated risk

Most existing systems are **reactive snapshots**. Our project aims to provide key inputs to the above listed key points.

Target Stakeholders

Our main stake holders are :

1. Humanitarian Aid Organizations such as WHO, UN , UNEP etc
2. Sudan Government

3) Project Objectives

1. IPC Food Insecurity Forecasting (State Level)

Forecast one month ahead:

- **Regression:** next-month IPC3+ percentage
- **Classification:** probability that IPC3+ worsens next month (2 percentage point increase)

Outputs:

- monthly risk ranking by state
- top drivers for each state-month prediction
- uncertainty/confidence flags

2. Displacement Pressure Monitoring and Hotspot Intelligence

Produce monthly operational outputs:

- **State pressure rankings:**
 - total IDPs by state
- **Hotspot and concentration intelligence:**
 - top localities contributing to state IDP burden
 - concentration measures (e.g., top 3 localities share)

Why this matters:

- High displacement pressure increases vulnerability via service strain, market stress, and household fragility.
- Hotspot intelligence supports locality-level targeting (Admin2) for assessments and interventions.

1. Data Understanding & Preparation

In [1]:

```
1 #import Libraries
2 import pandas as pd          # For working with datasets (tables, CS
3 import numpy as np           # For numerical operations, arrays, mat
4 import matplotlib.pyplot as plt # For plotting graphs and data visualiz
5 import seaborn as sns         # For advanced statistical visualizatio
6 from sklearn.model_selection import train_test_split      # Split data
7 from sklearn.linear_model import LogisticRegression       # Logistic R
8 from sklearn.ensemble import RandomForestClassifier      # Ensemble m
9 from sklearn.metrics import accuracy_score, precision_score, recall_score, f
10 from sklearn.metrics import confusion_matrix, classification_report
11 import glob
12 import pandas as pd
13 import re
14 import pandas as pd
15 import requests
16 from io import BytesIO
17 import zipfile, os
18 import geopandas as gpd
19 from pathlib import Path
20 from PIL import Image
21
22
```

Step 1 — Load the data

In [2]:

```
1 df = pd.read_csv("regional_food_security_master.csv", low_memory=False)
2 df.shape, df.columns
3
```

Out[2]:

```
((26145, 132),
Index(['Date of analysis', 'Country', 'Total country population', 'state',
       'Area', 'Validity period', 'From', 'To', 'Phase', 'Number',
       ...
       'T_35_39_2025', 'T_40_44_2025', 'T_45_49_2025', 'T_50_54_2025',
       'T_55_59_2025', 'T_60_64_2025', 'T_65_69_2025', 'T_70_74_2025',
       'T_75_79_2025', 'T_80Plus_2025'],
       dtype='object', length=132))
```

Step 2 — Remove metadata rows and parse dates

In [3]:

```
1 df = df[~df["From"].astype(str).str.startswith("#")].copy()
2
3 df["From"] = pd.to_datetime(df["From"], errors="coerce")
4 df = df.dropna(subset=["From"]).copy()
5
6 df["month"] = df["From"].dt.to_period("M").dt.to_timestamp()
7 df[["From", "month"]].head()
```

Out[3]:

| | From | month |
|---|------------|------------|
| 0 | 2025-09-01 | 2025-09-01 |
| 1 | 2025-09-01 | 2025-09-01 |
| 2 | 2025-09-01 | 2025-09-01 |
| 3 | 2025-09-01 | 2025-09-01 |
| 4 | 2025-09-01 | 2025-09-01 |

Step 3 — Standardize keys (Country, Phase, admin)

In [4]:

```
1 df["Country"] = df["Country"].astype(str).str.strip()
2 df["Phase"] = df["Phase"].astype(str).str.strip()
3
4 df["admin"] = df["Area"].where(df["Area"].notna(), df["state"]).astype(str).
5 df[["Country", "admin", "Phase"]].head()
```

Out[4]:

| | Country | admin | Phase |
|---|---------|-------|-------|
| 0 | SDN | Beida | all |
| 1 | SDN | Beida | 3+ |
| 2 | SDN | Beida | 1 |
| 3 | SDN | Beida | 2 |
| 4 | SDN | Beida | 3 |

Step 4 — Clean numeric percentage

In [5]:

```
1 df["Percentage"] = pd.to_numeric(df["Percentage"], errors="coerce")
2 clean_long = df.dropna(subset=["Country", "admin", "month", "Phase", "Percentage"])
3
4 clean_long.shape
```

Out[5]: (26145, 134)

Our model has been built on various data sets namely:

1. IPC datasets — Food insecurity truth signal
2. Displacement datasets — Needs and pressure signal
3. WFP price datasets — Economic stress signal
4. Rainfall dataset — Climate stress signal
5. Conflict datasets — Shock and disruption signal
6. Population and vulnerability — Exposure and normalization

1.1.1 Download an Excel dataset, lists all sheets, and previews the first few rows of each for quick inspection.

In [6]:

```

1 import pandas as pd
2 import requests
3 from io import BytesIO
4
5 url = "https://data.humdata.org/dataset/319dd40f-c0f8-4f6d-9a8e-9acf31007dd5
6
7 # Fetching the data
8 r = requests.get(url, timeout=60)
9 r.raise_for_status()
10
11 # Loading into pandas
12 xls = pd.ExcelFile(BytesIO(r.content))
13 print("Sheets found:", xls.sheet_names)
14
15 # Previewing sheets
16 for s in xls.sheet_names:
17     try:
18         df = xls.parse(s, nrows=3)
19         print(f"\n--- Sheet: {s} ---")
20         print(df)
21     except Exception as e:
22         print(f"\n--- Sheet: {s} --- (Error reading: {e})")

```

Sheets found: ['MASTER LIST (ADMIN1)', 'MASTER LIST (ADMIN 2)', 'Read Me', 'P
o0 per Localities', 'SAAD by Locality', 'Admin_label', 'Option']

--- Sheet: MASTER LIST (ADMIN1) ---

| | DTM SUDAN: IDP Master List 25-04-2024 | Unnamed: 1 | Unnamed: 2 | Unnamed: 3 |
|---|---------------------------------------|----------------|----------------|--------------|
| 0 | | NaN | NaN | NaN |
| 1 | LOCATION INFORMATION | | NaN | TOTAL |
| 2 | STATE OF DISPLACEMENT | STATE CODE | IDPs | HHS |
| | | Unnamed: 4 | Unnamed: 5 | Unnamed: 6 |
| 0 | | NaN | NaN | NaN |
| 1 | STATE of ORIGIN BY INDIVIUALS | | NaN | NaN |
| 2 | Aj Jazirah | Central Darfur | East Darfur | Khartoum |
| | Unnamed: 8 | Unnamed: 9 | Unnamed: 10 | Unnamed: 11 |
| 0 | NaN | NaN | NaN | NaN |
| 1 | NaN | NaN | NaN | NaN |
| 2 | North Darfur | North Kordofan | Sennar | South Darfur |
| | | | South Kordofan | |

1.1.2 This code Downloads an Excel dataset, extracts Admin-1 level displacement data, cleans and standardizes key columns, and saves the processed results as a CSV file.

In [7]:

```
1 url = "https://data.humdata.org/dataset/319dd40f-c0f8-4f6d-9a8e-9acf31007dd5
2
3 # download the xlsx
4 r = requests.get(url, timeout=60)
5 r.raise_for_status()
6
7 # read the Admin1 sheet (header row already detected as header=2 from your e
8 df = pd.read_excel(BytesIO(r.content), sheet_name="MASTER LIST (ADMIN1)", he
9
10 # keep only the four columns we need
11 state = df[["LOCATION INFORMATION", "Unnamed: 1", "TOTAL", "Unnamed: 3"]].co
12
13 # rename to clean names
14 state.columns = ["state", "state_code", "idps", "households"]
15
16 # clean: remove empty rows + force idps numeric
17 state["state"] = state["state"].astype(str).str.strip()
18 state["idps"] = pd.to_numeric(state["idps"], errors="coerce")
19
20 state = state.dropna(subset=["state", "idps"])
21 state = state[state["state"].str.lower().ne("nan")]
22
23 # save one CSV
24 state.to_csv("sudan_admin1_idps_2024-04-25.csv", index=False)
25
26 print("Saved: sudan_admin1_idps_2024-04-25.csv")
27 print(state.head())
28
```

Saved: sudan_admin1_idps_2024-04-25.csv

| | state | state_code | idps | households |
|---|----------------|------------|----------|------------|
| 2 | Aj Jazirah | SD15 | 371177.0 | 73323 |
| 3 | Blue Nile | SD08 | 147736.0 | 29836 |
| 4 | Central Darfur | SD06 | 430224.0 | 86044 |
| 5 | East Darfur | SD05 | 660140.0 | 131918 |
| 6 | Gedaref | SD12 | 492293.0 | 97817 |

1.1.3 Print and check the current working directory of the Python environment.

In [8]:

```
1 import os
2 print(os.getcwd())
```

C:\Users\user\Documents\Humanitarian

1.1.4 This code Downloads displacement data, extracts and cleans Admin-1 level statistics, saves the results as a CSV file, and prints the file's saved location.

In [10]:

```
1 url = "https://data.humdata.org/dataset/319dd40f-c0f8-4f6d-9a8e-9acf31007dd5
2
3 r = requests.get(url, timeout=60)
4 r.raise_for_status()
5
6 df = pd.read_excel(BytesIO(r.content), sheet_name="MASTER LIST (ADMIN1)", he
7
8 state = df[["LOCATION INFORMATION", "Unnamed: 1", "TOTAL", "Unnamed: 3"]].co
9 state.columns = ["state", "state_code", "idps", "households"]
10
11 state["state"] = state["state"].astype(str).str.strip()
12 state["idps"] = pd.to_numeric(state["idps"], errors="coerce")
13
14 state = state.dropna(subset=["state", "idps"])
15 state = state[state["state"].str.lower().ne("nan")]
16
17 out_file = "sudan_admin1_idps_2024-04-25.csv" # saved in C:\Users\user\Docu
18 state.to_csv(out_file, index=False)
19
20 print("Saved to:", os.path.abspath(out_file))
```

Saved to: C:\Users\user\Documents\Humanitarian\sudan_admin1_idps_2024-04-25.csv

1.1.5 The below Downloads Sudan Admin-1 boundary data in GeoJSON format and saves it locally for mapping or spatial analysis.

In [11]:

```
1 # This is the direct GeoJSON (WFS) endpoint for Sudan Admin1 boundaries
2 geojson_url = (
3     "https://geoportal.icpac.net/geoserver/ows?"
4     "service=WFS&version=1.0.0&request=GetFeature&"
5     "typename=geonode:sudan_admin_level1&"
6     "outputFormat=json&srsName=EPSG:4326"
7 )
8
9 out_file = "sudan_admin1.geojson" # will save in your current notebook fold
10
11 r = requests.get(geojson_url, timeout=120)
12 r.raise_for_status()
13
14 with open(out_file, "wb") as f:
15     f.write(r.content)
16
17 print("Saved boundary file to:", os.path.abspath(out_file))
18 print("File size (bytes):", os.path.getsize(out_file))
19
```

```
Saved boundary file to: C:\Users\user\Documents\Humanitarian\sudan_admin1.geojson
File size (bytes): 339637
```

1.1.6 The below code Loads displacement data, standardizes and cleans fields, removes invalid rows, aggregates duplicates by state, performs quality checks, and saves a final cleaned CSV ready for analysis or deployment.

In [12]:

```
1 # 1) Load your displacement CSV
2 df = pd.read_csv("sudan_admin1_idps_2024-04-25.csv")
3
4 # 2) Standardize column names (lowercase, underscores)
5 df.columns = [c.strip().lower().replace(" ", "_") for c in df.columns]
6
7 # Expected columns after your extraction:
8 # state, state_code, idps, households
9 print("Columns:", list(df.columns))
10
11 # 3) Clean text + numeric types
12 df["state"] = df["state"].astype(str).str.strip()
13 df["state_code"] = df["state_code"].astype(str).str.strip()
14
15 df["idps"] = pd.to_numeric(df["idps"], errors="coerce")
16 df["households"] = pd.to_numeric(df["households"], errors="coerce")
17
18 # 4) Drop empty rows
19 df = df.dropna(subset=["state", "idps"])
20
21 # 5) Remove obvious totals / junk rows (professional hygiene)
22 bad = {"total", "grand total", "overall", "nan", ""}
23 df = df[~df["state"].str.lower().isin(bad)]
24
25 # 6) Deduplicate safely (if duplicates exist, sum them)
26 df_clean = (
27     df.groupby(["state_code", "state"], as_index=False)[["idps", "households"]
28             .sum()]
29 )
30
31 # 7) Professional QA checks
32 print("\n--- QA SUMMARY ---")
33 print("Rows (states):", len(df_clean))
34 print("Unique state_code:", df_clean["state_code"].nunique())
35 print("Missing state_code:", df_clean["state_code"].isna().sum())
36 print("Any duplicate state_code:", df_clean["state_code"].duplicated().any())
37
38 print("\nTop 10 states by IDPs:")
39 display(df_clean.sort_values("idps", ascending=False).head(10))
40
41 print("\nAll states (sorted A-Z):")
42 display(df_clean.sort_values("state")[["state_code", "state", "idps", "house"])
43
44 # 8) Save the cleaned table (deployment-friendly, stable schema)
45 df_clean.to_csv("sudan_admin1_idps_2024-04-25_CLEAN.csv", index=False)
46 print("\nSaved: sudan_admin1_idps_2024-04-25_CLEAN.csv")
```

Columns: ['state', 'state_code', 'idps', 'households']

--- QA SUMMARY ---

Rows (states): 18

Unique state_code: 18

Missing state_code: 0

Any duplicate state_code: False

Top 10 states by IDPs:

| | state_code | state | idps | households |
|----|------------|----------------|----------|------------|
| 2 | SD03 | South Darfur | 744243.0 | 148848 |
| 15 | SD16 | River Nile | 698334.0 | 137799 |
| 4 | SD05 | East Darfur | 660140.0 | 131918 |
| 1 | SD02 | North Darfur | 573055.0 | 114503 |
| 8 | SD09 | White Nile | 532643.0 | 105920 |
| 13 | SD14 | Sennar | 523986.0 | 104002 |
| 11 | SD12 | Gedaref | 492293.0 | 97817 |
| 5 | SD06 | Central Darfur | 430224.0 | 86044 |
| 16 | SD17 | Northern | 399867.0 | 80305 |
| 14 | SD15 | Aj Jazirah | 371177.0 | 73323 |

All states (sorted A-Z):

| | state_code | state | idps | households |
|----|------------|----------------|----------|------------|
| 14 | SD15 | Aj Jazirah | 371177.0 | 73323 |
| 7 | SD08 | Blue Nile | 147736.0 | 29836 |
| 5 | SD06 | Central Darfur | 430224.0 | 86044 |
| 4 | SD05 | East Darfur | 660140.0 | 131918 |
| 11 | SD12 | Gedaref | 492293.0 | 97817 |
| 10 | SD11 | Kassala | 200083.0 | 40282 |
| 0 | SD01 | Khartoum | 69057.0 | 13717 |
| 1 | SD02 | North Darfur | 573055.0 | 114503 |
| 12 | SD13 | North Kordofan | 174007.0 | 34261 |
| 16 | SD17 | Northern | 399867.0 | 80305 |
| 9 | SD10 | Red Sea | 247874.0 | 49953 |
| 15 | SD16 | River Nile | 698334.0 | 137799 |
| 13 | SD14 | Sennar | 523986.0 | 104002 |
| 2 | SD03 | South Darfur | 744243.0 | 148848 |
| 6 | SD07 | South Kordofan | 198839.0 | 39492 |
| 3 | SD04 | West Darfur | 174540.0 | 34908 |
| 17 | SD18 | West Kordofan | 148718.0 | 29340 |
| 8 | SD09 | White Nile | 532643.0 | 105920 |

Saved: sudan_admin1_idps_2024-04-25_CLEAN.csv

1.1.7 Below code Builds a 2023 Admin-1 displacement panel by looping through multiple Excel files, extracting state-level IDPs/households, deriving report dates from sheet names, ranking states (top 5 hotspots), and exporting a combined monthly dataset to CSV.

```
import glob
```


In [13]:

```
1 import glob
2 import pandas as pd
3 import re
4
5 HOTSPOT_TOPK = 5
6 rows_all = []
7
8 files = sorted(glob.glob("data/2023_excels/*.xlsx"))
9 print("Files found:", len(files))
10
11 def pick_dataset_sheet(sheet_names):
12     # Prefer "Dataset (...)" sheets
13     for s in sheet_names:
14         if str(s).lower().startswith("dataset"):
15             return s
16     return sheet_names[0]
17
18 def extract_date_from_sheetname(s):
19     # "Dataset (27 April 2023)" -> 2023-04-27
20     m = re.search(r"\((.+)\)", str(s))
21     if not m:
22         return None
23     txt = m.group(1)
24     try:
25         return pd.to_datetime(txt, dayfirst=True, errors="raise")
26     except Exception:
27         return pd.to_datetime(txt, errors="coerce")
28
29 for path in files:
30     xls = pd.ExcelFile(path)
31     sheet = pick_dataset_sheet(xls.sheet_names)
32
33     df = pd.read_excel(path, sheet_name=sheet)
34
35     # map columns (exact names from your sample)
36     col_state = "STATE OF AFFECTED POPULATION"
37     col_code = "STATE PCODE OF AFFECTED POPULATION"
38     col_idps = "# IDP INDIVIDUALS"
39     col_hh = "# IDP HOUSEHOLDS"
40
41     # safety: only proceed if required columns exist
42     if not all(c in df.columns for c in [col_state, col_code, col_idps, col_hh]):
43         continue
44
45     out = df[[col_state, col_code, col_idps, col_hh]].copy()
46     out.columns = ["state_name", "state_code", "idps", "households"]
47
48     out["state_name"] = out["state_name"].astype(str).str.strip()
49     out["state_code"] = out["state_code"].astype(str).str.strip()
50
51     out["idps"] = pd.to_numeric(out["idps"], errors="coerce")
52     out["households"] = pd.to_numeric(out["households"], errors="coerce")
53
54     # admin1 panel: sum within affected state
55     out = (out.groupby(["state_code", "state_name"], as_index=False)
56            .agg({"idps": "sum", "households": "sum"}))
```

```

58 # report_date from sheet name (best available in these files)
59 report_date = extract_date_from_sheetname(sheet)
60 if pd.isna(report_date):
61     # fallback: skip if we can't date it
62     continue
63
64 out["report_date"] = report_date.normalize()
65 out["month_start"] = report_date.to_period("M").to_timestamp()
66
67 # rank + hotspot (per report)
68 out = out.sort_values("idps", ascending=False).reset_index(drop=True)
69 out["rank"] = out.index + 1
70 out["hotspot"] = out["rank"] <= HOTSPOT_TOPK
71
72 # keep only the 18 states (if extras exist, take top 18 by idps)
73 out = out.head(18)
74
75 print("Processing:", out["report_date"].iloc[0].date(), "| States:", len
76 rows_all.append(out)
77
78 panel_2023 = pd.concat(rows_all, ignore_index=True)
79 panel_2023.to_csv("sudan_admin1_idps_panel_2023.csv", index=False)
80
81 print("\nSaved: sudan_admin1_idps_panel_2023.csv")
82 print("Rows:", len(panel_2023), "| Months:", panel_2023["month_start"].nuniq
83
84 # show Latest report Like your example
85 latest = panel_2023["report_date"].max()
86 print("\nLatest report_date:", latest.date())
87 print(panel_2023[panel_2023["report_date"] == latest].sort_values("rank")[[[
88     "state_code", "state_name", "idps", "households", "report_date", "month_start
89 ]].to_string(index=False))
90

```

Files found: 34

Processing: 2023-04-27 | States: 15

Saved: sudan_admin1_idps_panel_2023.csv

Rows: 15 | Months: 1 | States: 15

Latest report_date: 2023-04-27

| | state_code | state_name | idps | households | report_date | month_start | rank |
|-------|-------------|----------------|----------|------------|-------------|-------------|------|
| True | SD04 | West Darfur | 194593.0 | 38919.0 | 2023-04-27 | 2023-04-01 | 1 |
| True | SD03 | South Darfur | 45000.0 | 9000.0 | 2023-04-27 | 2023-04-01 | 2 |
| True | SD17 | Northern | 29200.0 | 5840.0 | 2023-04-27 | 2023-04-01 | 3 |
| True | SD01 | Khartoum | 13545.0 | 2709.0 | 2023-04-27 | 2023-04-01 | 4 |
| True | SD13 | North Kordofan | 13270.0 | 2654.0 | 2023-04-27 | 2023-04-01 | 5 |
| False | SD02 | North Darfur | 11675.0 | 2335.0 | 2023-04-27 | 2023-04-01 | 6 |
| False | SD15 | Aj Jazirah | 8795.0 | 1759.0 | 2023-04-27 | 2023-04-01 | 7 |
| False | SD09 | White Nile | 6165.0 | 1233.0 | 2023-04-27 | 2023-04-01 | 8 |
| False | SD14 | Sennar | 5560.0 | 1112.0 | 2023-04-27 | 2023-04-01 | 9 |
| False | SD16 | River Nile | 2910.0 | 582.0 | 2023-04-27 | 2023-04-01 | 10 |
| False | SD06 | Central Darfur | 1780.0 | 356.0 | 2023-04-27 | 2023-04-01 | 11 |
| False | SD10 | Red Sea | 1205.0 | 241.0 | 2023-04-27 | 2023-04-01 | 12 |
| False | SD08 | Blue Nile | 260.0 | 52.0 | 2023-04-27 | 2023-04-01 | 13 |
| False | SD11 | Kassala | 95.0 | 19.0 | 2023-04-27 | 2023-04-01 | 14 |
| False | #adm1+pcode | #adm1+name | 0.0 | 0.0 | 2023-04-27 | 2023-04-01 | 15 |

1.1.8 Below code Creates a fixed 18-state Admin-1 reference list, merges it into each report date to ensure all states are present (filling missing values with zero), recalculates ranks and top-5 hotspots per report, and saves a consistent 18-state monthly panel to CSV.

In [14]:

```
1 # Fixed Sudan Admin1 spine (SD01-SD18)
2 spine = pd.DataFrame([
3     ("SD01", "Khartoum"),
4     ("SD02", "North Darfur"),
5     ("SD03", "South Darfur"),
6     ("SD04", "West Darfur"),
7     ("SD05", "East Darfur"),
8     ("SD06", "Central Darfur"),
9     ("SD07", "South Kordofan"),
10    ("SD08", "Blue Nile"),
11    ("SD09", "White Nile"),
12    ("SD10", "Red Sea"),
13    ("SD11", "Kassala"),
14    ("SD12", "Gedaref"),
15    ("SD13", "North Kordofan"),
16    ("SD14", "Sennar"),
17    ("SD15", "Aj Jazirah"),
18    ("SD16", "River Nile"),
19    ("SD17", "Northern"),
20    ("SD18", "West Kordofan"),
21 ], columns=["state_code", "state_name"])
22
23 # panel_2023 is your extracted (before spine-fill). If you only have panel_2
24 # read it back and treat it as the extracted input:
25 panel = panel_2023.copy() if "panel_2023" in globals() else pd.read_csv("sud
26
27 fixed = []
28 for d, g in panel.groupby("report_date"):
29     m = spine.merge(g[["state_code", "idps", "households"]], on="state_code",
30     m["report_date"] = pd.to_datetime(d)
31     m["month_start"] = pd.to_datetime(d).to_period("M").to_timestamp()
32     m["idps"] = m["idps"].fillna(0)
33     m["households"] = m["households"].fillna(0)
34
35     m = m.sort_values("idps", ascending=False).reset_index(drop=True)
36     m["rank"] = m.index + 1
37     m["hotspot"] = m["rank"] <= HOTSPOT_TOPK
38     fixed.append(m)
39
40 panel_2023_fixed18 = pd.concat(fixed, ignore_index=True)
41 panel_2023_fixed18.to_csv("sudan_admin1_idps_panel_2023.csv", index=False)
42
43 print("Saved: sudan_admin1_idps_panel_2023.csv")
44 print("Rows:", len(panel_2023_fixed18), "| Dates:", panel_2023_fixed18["repo
45
46 # show Latest block
47 latest = panel_2023_fixed18["report_date"].max()
48 print("\nProcessing:", latest.date())
49 print(panel_2023_fixed18[panel_2023_fixed18["report_date"] == latest].sort_v
50     "state_code", "state_name", "idps", "households", "report_date", "month_start
51 ]).to_string(index=False))
```

Saved: sudan_admin1_idps_panel_2023.csv

Rows: 18 | Dates: 1 | States: 18

Processing: 2023-04-27

| state_code | state_name | idps | households | report_date | month_start | rank | h |
|------------|------------|----------------|------------|-------------|-------------|------------|----|
| True | SD04 | West Darfur | 194593.0 | 38919.0 | 2023-04-27 | 2023-04-01 | 1 |
| True | SD03 | South Darfur | 45000.0 | 9000.0 | 2023-04-27 | 2023-04-01 | 2 |
| True | SD17 | Northern | 29200.0 | 5840.0 | 2023-04-27 | 2023-04-01 | 3 |
| True | SD01 | Khartoum | 13545.0 | 2709.0 | 2023-04-27 | 2023-04-01 | 4 |
| True | SD13 | North Kordofan | 13270.0 | 2654.0 | 2023-04-27 | 2023-04-01 | 5 |
| False | SD02 | North Darfur | 11675.0 | 2335.0 | 2023-04-27 | 2023-04-01 | 6 |
| False | SD15 | Aj Jazirah | 8795.0 | 1759.0 | 2023-04-27 | 2023-04-01 | 7 |
| False | SD09 | White Nile | 6165.0 | 1233.0 | 2023-04-27 | 2023-04-01 | 8 |
| False | SD14 | Sennar | 5560.0 | 1112.0 | 2023-04-27 | 2023-04-01 | 9 |
| False | SD16 | River Nile | 2910.0 | 582.0 | 2023-04-27 | 2023-04-01 | 10 |
| False | SD06 | Central Darfur | 1780.0 | 356.0 | 2023-04-27 | 2023-04-01 | 11 |
| False | SD10 | Red Sea | 1205.0 | 241.0 | 2023-04-27 | 2023-04-01 | 12 |
| False | SD08 | Blue Nile | 260.0 | 52.0 | 2023-04-27 | 2023-04-01 | 13 |
| False | SD11 | Kassala | 95.0 | 19.0 | 2023-04-27 | 2023-04-01 | 14 |
| False | SD07 | South Kordofan | 0.0 | 0.0 | 2023-04-27 | 2023-04-01 | 15 |
| False | SD05 | East Darfur | 0.0 | 0.0 | 2023-04-27 | 2023-04-01 | 16 |
| False | SD12 | Gedaref | 0.0 | 0.0 | 2023-04-27 | 2023-04-01 | 17 |
| False | SD18 | West Kordofan | 0.0 | 0.0 | 2023-04-27 | 2023-04-01 | 18 |

1.1.9 Below code Loads a combined 2023–2024 displacement panel, selects geography-ready fields, exports them to a CSV for mapping snapshots, and prints basic dataset statistics.

In [15]:

```
1 panel = pd.read_csv("sudan_admin1_idps_panel_2023_2024.csv", parse_dates=["r  
2  
3 geo = panel[["state_code", "state_name", "report_date", "idps", "households", "ra  
4 geo.to_csv("geo_admin1_snapshots_2023_2024.csv", index=False)  
5  
6 print("Saved: geo_admin1_snapshots_2023_2024.csv")  
7 print("Rows:", len(geo), "| Dates:", geo["report_date"].nunique(), "| States:  
8
```

Saved: geo_admin1_snapshots_2023_2024.csv

Rows: 54 | Dates: 3 | States: 18

1.1.10 Below code Splits the geographic displacement dataset into separate CSV files by report date and saves each snapshot to a date-named folder for easy mapping or time-series use.

In [16]:

```
1 df = pd.read_csv("geo_admin1_snapshots_2023_2024.csv", parse_dates=["report_  
2 os.makedirs("geo_by_date", exist_ok=True)  
3  
4 for d in sorted(df["report_date"].unique()):  
5     sub = df[df["report_date"] == d].copy()  
6     out = f"geo_by_date/geo_admin1_{pd.to_datetime(d).date()}.csv"  
7     sub.to_csv(out, index=False)  
8     print("Saved:", out, "| rows:", len(sub))  
9
```

Saved: geo_by_date/geo_admin1_2023-04-27.csv | rows: 18

Saved: geo_by_date/geo_admin1_2023-05-05.csv | rows: 18

Saved: geo_by_date/geo_admin1_2024-04-25.csv | rows: 18

1.1.11 Below code Identifies the top 5 displacement hotspot states per report date based on IDPs, saves the results to a CSV file, and prints a concise hotspot summary table.

In [17]:

```
1 df = pd.read_csv("geo_admin1_snapshots_2023_2024.csv", parse_dates=["report_"
2
3 top5 = (df.sort_values(["report_date", "idps"], ascending=[True, False])
4         .groupby("report_date")
5         .head(5)
6         .copy())
7
8 top5.to_csv("report_hotspots_top5_by_date.csv", index=False)
9 print("Saved: report_hotspots_top5_by_date.csv")
10 print(top5[["report_date", "state_code", "state_name", "idps", "rank", "hotspot"]]
11
```

Saved: report_hotspots_top5_by_date.csv

| report_date | state_code | state_name | idps | rank | hotspot |
|-------------|------------|----------------|----------|------|---------|
| 2023-04-27 | SD04 | West Darfur | 194593.0 | 1.0 | True |
| 2023-04-27 | SD03 | South Darfur | 45000.0 | 2.0 | True |
| 2023-04-27 | SD17 | Northern | 29200.0 | 3.0 | True |
| 2023-04-27 | SD01 | Khartoum | 13545.0 | 4.0 | True |
| 2023-04-27 | SD13 | North Kordofan | 13270.0 | 5.0 | True |
| 2023-05-05 | SD09 | White Nile | 188635.0 | 1.0 | True |
| 2023-05-05 | SD04 | West Darfur | 156565.0 | 2.0 | True |
| 2023-05-05 | SD17 | Northern | 106600.0 | 3.0 | True |
| 2023-05-05 | SD16 | River Nile | 96095.0 | 4.0 | True |
| 2023-05-05 | SD15 | Aj Jazirah | 49280.0 | 5.0 | True |
| 2024-04-25 | SD03 | South Darfur | 744243.0 | 1.0 | True |
| 2024-04-25 | SD16 | River Nile | 698334.0 | 2.0 | True |
| 2024-04-25 | SD05 | East Darfur | 660140.0 | 3.0 | True |
| 2024-04-25 | SD02 | North Darfur | 573055.0 | 4.0 | True |
| 2024-04-25 | SD09 | White Nile | 532643.0 | 5.0 | True |

1.1.13 Below code Adds a global IDP-based intensity metric by normalizing displacement counts against the overall maximum across all dates, then updates the dataset for consistent cross-time comparison.

In [18]:

```
1 import pandas as pd
2 df = pd.read_csv("geo_admin1_snapshots_2023_2024_INTENSITY.csv", parse_dates=
3
4 global_max = df["idps"].max()
5 df["intensity_idps_global"] = df["idps"] / global_max
6
7 df.to_csv("geo_admin1_snapshots_2023_2024_INTENSITY.csv", index=False)
8 print("Updated with intensity_idps_global. Global max:", global_max)
9
```

Updated with intensity_idps_global. Global max: 744243.0

1.1.14 Below code Splits the intensity-enhanced displacement dataset into per-date GIS-ready CSV files, keeping only essential fields and a globally comparable intensity metric for consistent map symbolization.

```
In [19]: 1 df = pd.read_csv("geo_admin1_snapshots_2023_2024_INTENSITY.csv", parse_dates=2
 3 # pick the intensity field you want to symbolize
 4 # global comparable across dates:
 5 value_col = "intensity_idps_global"
 6
 7 out_dir = "geo_by_date_intensity"
 8 os.makedirs(out_dir, exist_ok=True)
 9
10 for d in sorted(df["report_date"].unique()):
11     sub = df[df["report_date"] == d].copy()
12
13     # keep only what GIS needs (simple)
14     sub = sub[[
15         "state_code", "state_name", "report_date",
16         "idps", "households",
17         "rank", "hotspot",
18         "hotspot_intensity_rank",
19         "intensity_idps_global",
20         "intensity_idps_share"
21     ]].copy()
22
23     out = os.path.join(out_dir, f"geo_admin1_{pd.to_datetime(d).date()}_inte-
24     sub.to_csv(out, index=False)
25     print("Saved:", out, "| rows:", len(sub), "| symbolize:", value_col)
26
```

```
Saved: geo_by_date_intensity\geo_admin1_2023-04-27_intensity.csv | rows: 18 | s
ymbolize: intensity_idps_global
Saved: geo_by_date_intensity\geo_admin1_2023-05-05_intensity.csv | rows: 18 | s
ymbolize: intensity_idps_global
Saved: geo_by_date_intensity\geo_admin1_2024-04-25_intensity.csv | rows: 18 | s
ymbolize: intensity_idps_global
```

1.1.15 Below code Unzips an uploaded GeoJSON boundary file, loads it with GeoPandas, inspects its structure, and identifies likely administrative code fields for joining with displacement data.

In [20]:

```
1 import zipfile, os
2 import geopandas as gpd
3
4 zip_path = "sdn_admin_boundaries.geojson.zip" # if this fails, replace with
5
6 out_dir = "sdn_admin_boundaries_unzipped"
7 os.makedirs(out_dir, exist_ok=True)
8
9 with zipfile.ZipFile(zip_path, "r") as z:
10     z.extractall(out_dir)
11
12 # find the geojson file
13 geojson_files = [os.path.join(out_dir, f) for f in os.listdir(out_dir) if f.
14 print("GeoJSON files:", geojson_files)
15
16 gdf = gpd.read_file(geojson_files[0])
17 print("Rows:", len(gdf))
18 print("Columns:", list(gdf.columns))
19
20 # show Likely join fields
21 cand = [c for c in gdf.columns if "PCODE" in c.lower() or "CODE" in c.lower(
22 print("Candidate code fields:", cand)
23
24 gdf.head(3)
25
```

```
GeoJSON files: ['sdn_admin_boundaries_unzipped\\sdn_admin0.geojson', 'sdn_admin_
_boundaries_unzipped\\sdn_admin1.geojson', 'sdn_admin_boundaries_unzipped\\sdn_
admin2.geojson', 'sdn_admin_boundaries_unzipped\\sdn_adminlines.geojson', 'sdn_
admin_boundaries_unzipped\\sdn_adminpoints.geojson']
Rows: 1
Columns: ['iso2', 'iso3', 'ADM0_NAME', 'ADM0_NAME1', 'ADM0_NAME2', 'ADM0_NAME
3', 'ADM0_PCODE', 'VALID_ON', 'VALID_TO', 'VERSION', 'AREA_SQKM', 'LANG', 'LANG
1', 'LANG2', 'LANG3', 'ADM0_REF_NAME', 'CENTER_LAT', 'CENTER_LON', 'GEOMETRY']
Candidate code fields: ['ADM0_PCODE']
```

Out[20]:

| | iso2 | iso3 | ADM0_NAME | ADM0_NAME1 | ADM0_NAME2 | ADM0_NAME3 | ADM0_PCODE | VALID_ON | VALID_ |
|---|------|------|-----------|------------|------------|------------|------------|----------|--------|
| 0 | SD | SDN | Sudan | السودان | None | None | SD | 2020-08- | 31 Nc |



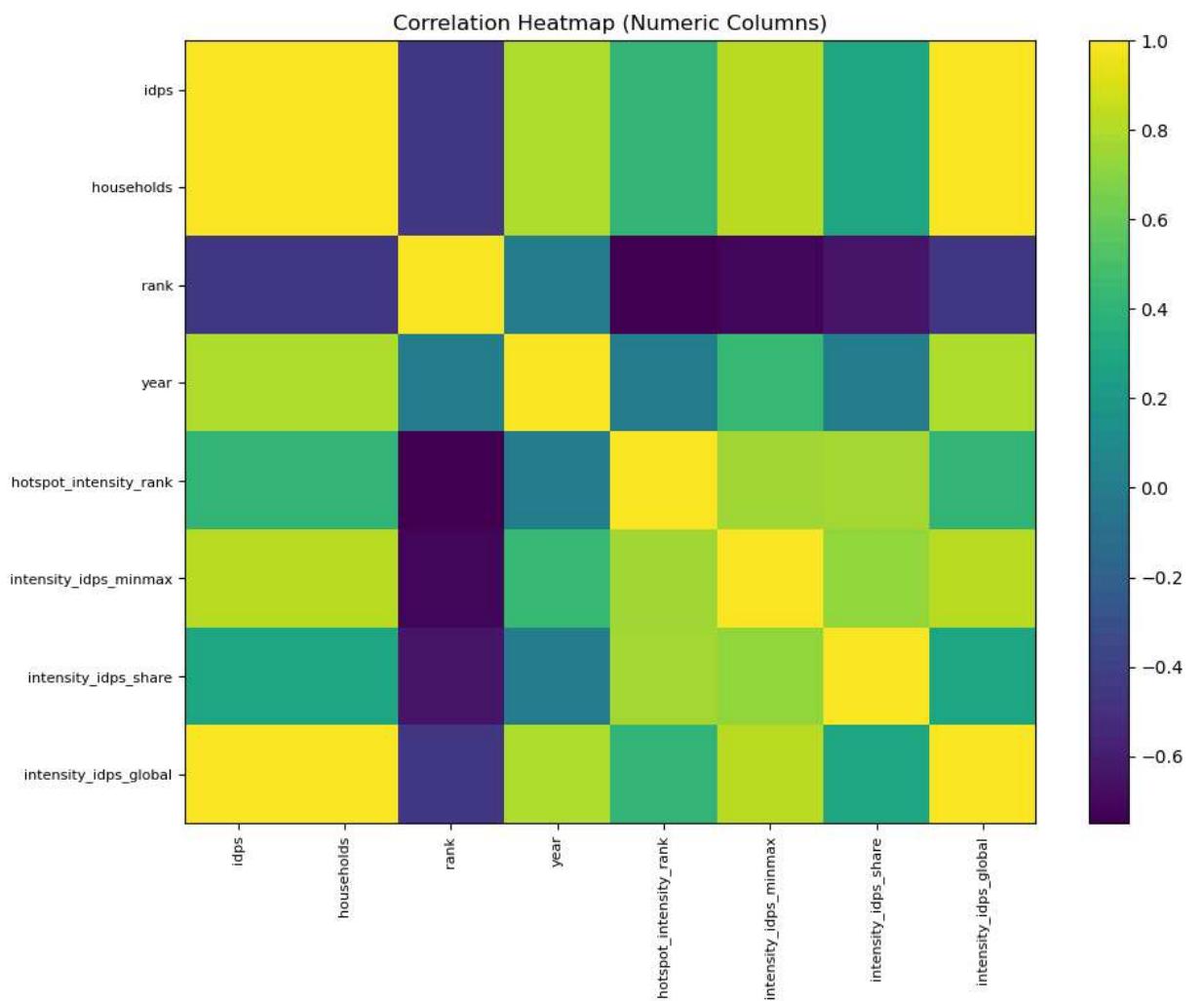
Exploratory Data Analysis and Visualizations

Numeric features in the dataset

This heatmap visualizes the pairwise correlations among all numeric features in the dataset, helping identify strong relationships and potential multicollinearity prior to modeling.

In [21]:

```
1 from pathlib import Path
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # --- archive saver (keeps your structure: artifacts/eda/) ---
6 ARTIFACTS_DIR = Path("artifacts") / "eda"
7 ARTIFACTS_DIR.mkdir(parents=True, exist_ok=True)
8
9 def save_show(fig, name: str, folder: Path = ARTIFACTS_DIR, dpi: int = 200):
10     out_path = folder / f"{name}.png"
11     fig.tight_layout()
12     fig.savefig(out_path, dpi=dpi, bbox_inches="tight")
13     plt.show()
14     print(f"Saved: {out_path}")
15
16 # --- correlation heatmap (numeric columns) ---
17 num_cols = df.select_dtypes(include=[np.number]).columns.tolist()
18
19 if len(num_cols) < 2:
20     print("Not enough numeric columns to compute a correlation matrix.")
21 else:
22     corr = df[num_cols].corr(numeric_only=True)
23
24 fig, ax = plt.subplots(figsize=(10, 8))
25 im = ax.imshow(corr.values, aspect="auto")
26 ax.set_title("Correlation Heatmap (Numeric Columns)")
27
28 ax.set_xticks(range(len(num_cols)))
29 ax.set_yticks(range(len(num_cols)))
30 ax.set_xticklabels(num_cols, rotation=90, fontsize=8)
31 ax.set_yticklabels(num_cols, fontsize=8)
32
33 fig.colorbar(im, ax=ax)
34 save_show(fig, "corr_heatmap_numeric")
35
36
37
```



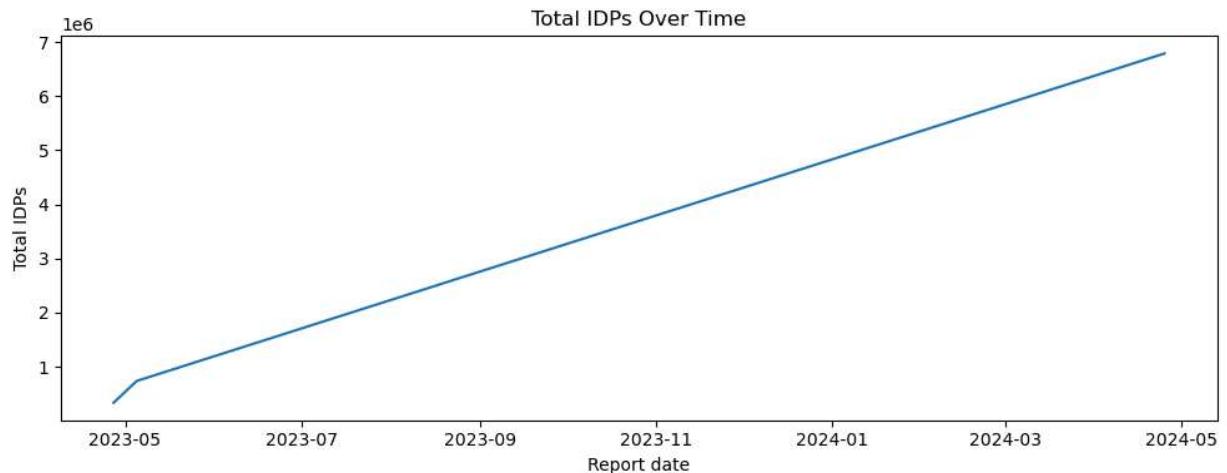
Saved: artifacts\eda\corr_heatmap_numeric.png

Total IDPs Over Time

This line chart shows the trend of total internally displaced persons (IDPs) over time, highlighting changes in displacement levels across reporting dates.

In [23]:

```
1 trend = (df.dropna(subset=["report_date"])
2         .groupby("report_date")["idps"].sum()
3         .sort_index())
4
5 fig, ax = plt.subplots(figsize=(10, 4))
6 ax.plot(trend.index, trend.values)
7 ax.set_title("Total IDPs Over Time")
8 ax.set_xlabel("Report date")
9 ax.set_ylabel("Total IDPs")
10 save_show(fig, "trend_total_idps_over_time")
11
```



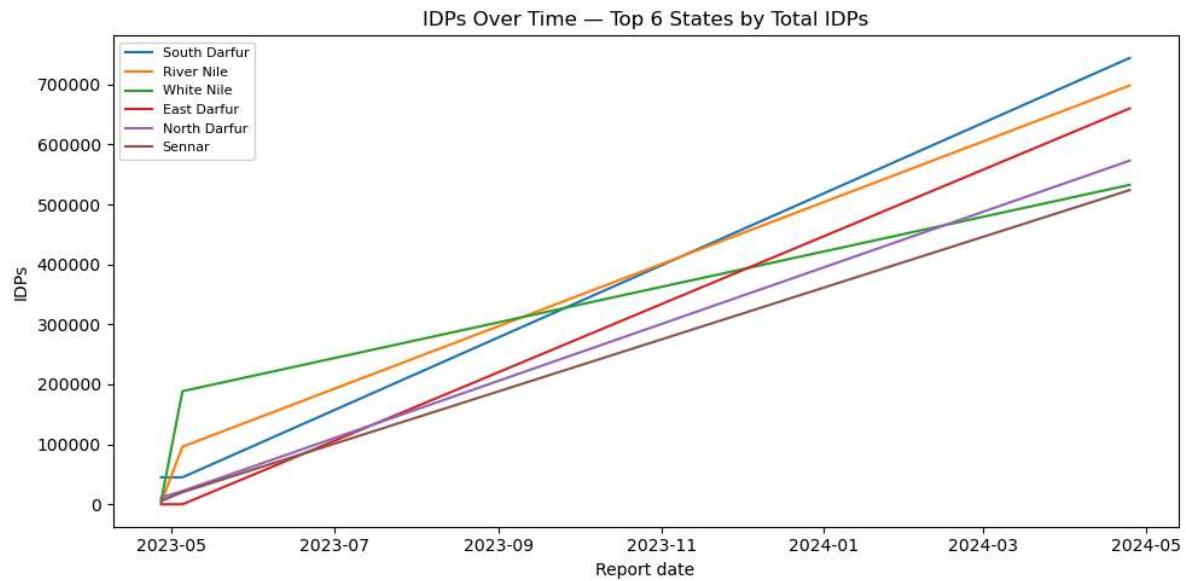
Saved: artifacts\eda\trend_total_idps_over_time.png

IDPs Over Time — Top 6 States by Total IDPs

This multi-line chart compares IDP trends over time for the top six states with the highest total number of internally displaced persons, highlighting differences in displacement patterns across states.

In [24]:

```
1 TOP_N = 6
2 top_states = (df.groupby("state_name")["idps"].sum()
3                 .sort_values(ascending=False)
4                 .head(TOP_N)
5                 .index.tolist())
6
7 fig, ax = plt.subplots(figsize=(10, 5))
8 for s in top_states:
9     sub = df[df["state_name"] == s].dropna(subset=["report_date"])
10    ts = sub.groupby("report_date")["idps"].sum().sort_index()
11    ax.plot(ts.index, ts.values, label=str(s))
12
13 ax.set_title(f"IDPs Over Time — Top {TOP_N} States by Total IDPs")
14 ax.set_xlabel("Report date")
15 ax.set_ylabel("IDPs")
16 ax.legend(fontsize=8)
17 save_show(fig, "trend_idps_top_states")
18
```



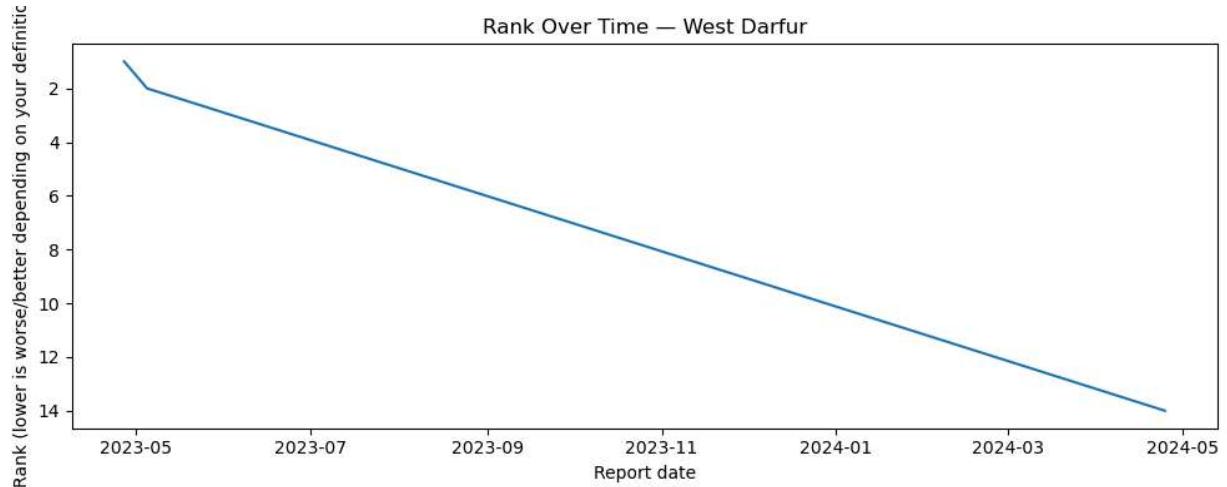
Saved: artifacts\eda\trend_idps_top_states.png

Rank Over Time for States

This line chart tracks how a selected state's rank changes over time, illustrating shifts in its relative position across reporting periods.

In [29]:

```
1 STATE = df["state_name"].iloc[0] # change to e.g. "Khartoum"
2 sub = df[df["state_name"] == STATE].dropna(subset=["report_date"]).sort_values
3
4 if "rank" in sub.columns:
5     fig, ax = plt.subplots(figsize=(10, 4))
6     ax.plot(sub["report_date"], sub["rank"])
7     ax.set_title(f"Rank Over Time — {STATE}")
8     ax.set_xlabel("Report date")
9     ax.set_ylabel("Rank (lower is worse/better depending on your definition)")
10    ax.invert_yaxis() # ranks usually read better inverted (1 at top)
11    save_show(fig, f"rank_over_time_{STATE}.replace(" ", "_"))
```



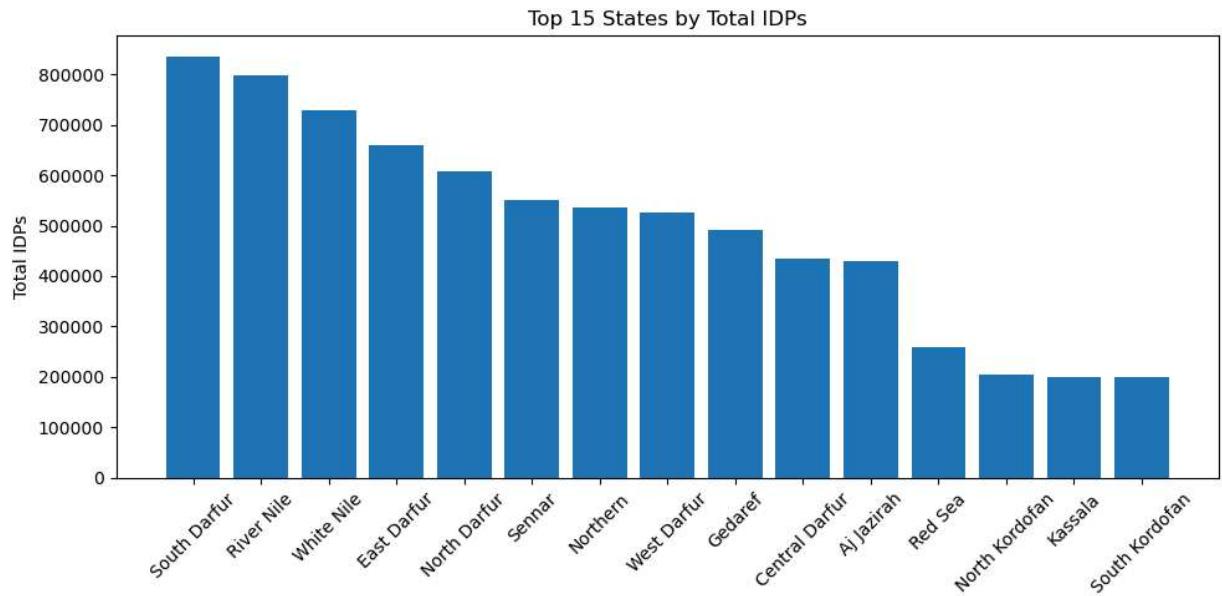
Saved: C:\Users\user\Documents\Humanitarian\artifacts\eda_food\rank_over_time_West_Darfur.png

Top 15 States by Total IDPs

This bar chart displays the top 15 states ranked by total internally displaced persons (IDPs), highlighting regions with the highest cumulative displacement.

In [25]:

```
1 top = (df.groupby("state_name")["idps"].sum()
2         .sort_values(ascending=False)
3         .head(15))
4
5 fig, ax = plt.subplots(figsize=(10, 5))
6 ax.bar(top.index.astype(str), top.values)
7 ax.set_title("Top 15 States by Total IDPs")
8 ax.set_ylabel("Total IDPs")
9 ax.tick_params(axis="x", rotation=45)
10 save_show(fig, "bar_top15_states_total_idps")
11
```



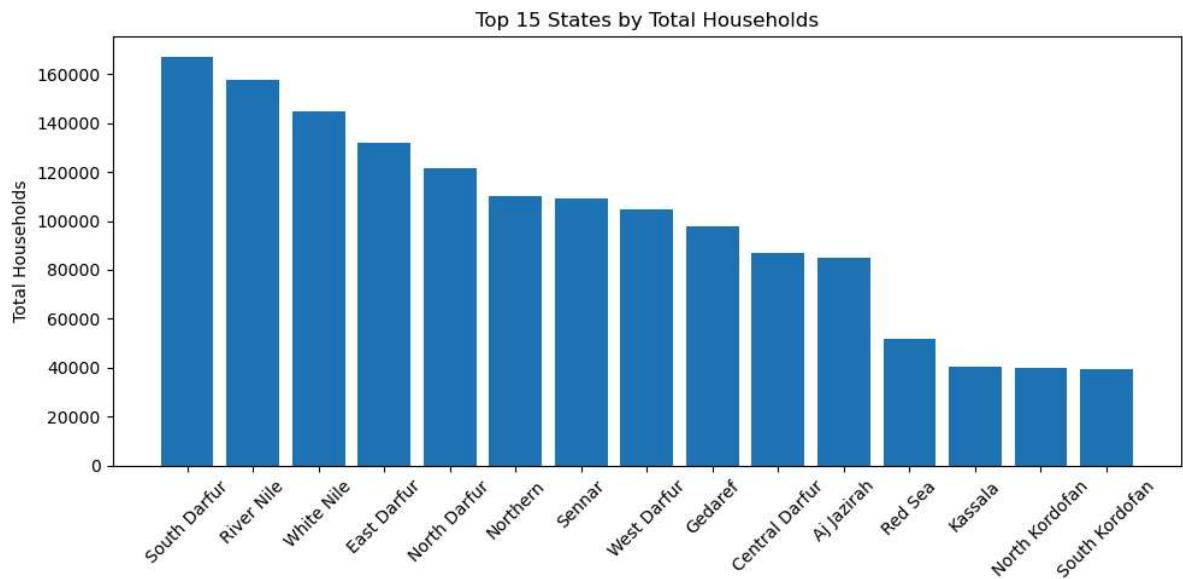
Saved: artifacts\eda\bar_top15_states_total_idps.png

Top 15 States by Total Households

This bar chart shows the top 15 states ranked by total affected households, highlighting regions with the largest cumulative household impact.

In [27]:

```
1 top = (df.groupby("state_name")["households"].sum()
2         .sort_values(ascending=False)
3         .head(15))
4
5 fig, ax = plt.subplots(figsize=(10, 5))
6 ax.bar(top.index.astype(str), top.values)
7 ax.set_title("Top 15 States by Total Households")
8 ax.set_ylabel("Total Households")
9 ax.tick_params(axis="x", rotation=45)
10 save_show(fig, "bar_top15_states_total_households")
11
```



Saved: C:\Users\user\Documents\Humanitarian\artifacts\eda_food\bar_top15_stat

Food Security Exploratory Data Analysis (EDA)

This analysis explores regional food security data by visualizing distributions of numeric and categorical food indicators and tracking trends over time.

In [28]:

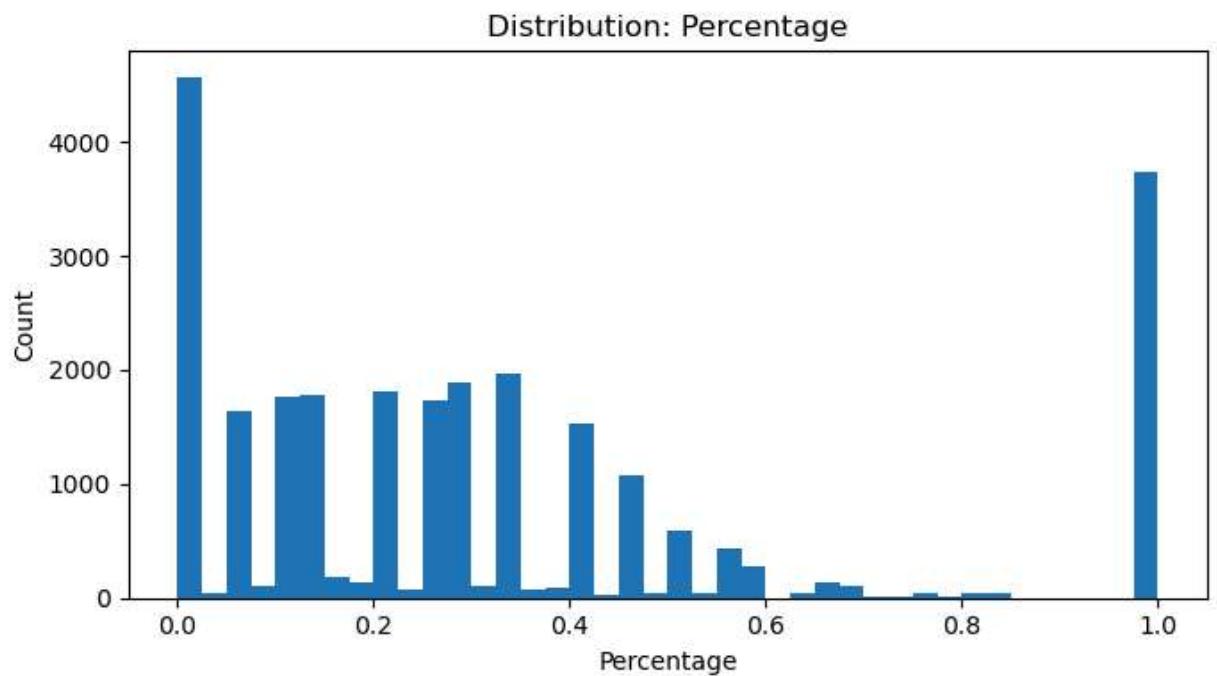
```
1 import re
2 from pathlib import Path
3
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7
8 # -----
9 # 1) Load food data (defines food_df)
10 # -----
11 FOOD_MASTER = Path("regional_food_security_master.csv")
12
13 if not FOOD_MASTER.exists():
14     raise FileNotFoundError(
15         f"Missing {FOOD_MASTER}. Put it in the same folder as this notebook,"
16     )
17
18 food_df = pd.read_csv(FOOD_MASTER, low_memory=False)
19 print("Loaded food_df:", food_df.shape)
20
21 # Optional: if your file has 'From' column and metadata rows starting with '#'
22 if "From" in food_df.columns:
23     food_df = food_df[~food_df["From"].astype(str).str.startswith("#")].copy()
24     food_df["From"] = pd.to_datetime(food_df["From"], errors="coerce")
25     food_df = food_df.dropna(subset=["From"]).copy()
26     food_df["month"] = food_df["From"].dt.to_period("M").dt.to_timestamp()
27
28 # -----
29 # 2) Setup archive + saver
30 # -----
31 OUT = Path("artifacts") / "eda_food"
32 OUT.mkdir(parents=True, exist_ok=True)
33
34 def save_show(fig, name):
35     p = OUT / f"{name}.png"
36     fig.tight_layout()
37     fig.savefig(p, dpi=300, bbox_inches="tight")
38     plt.show()
39     plt.close(fig)
40     print("Saved:", p.resolve())
41
42 # -----
43 # 3) Column detection (defines food_pat)
44 # -----
45 food_pat = re.compile(r"(food|ipc|phase|price|market|cereal|wheat|sorghum|mifl")
46
47 food_cols = [c for c in food_df.columns if food_pat.search(str(c))]
48 print("Detected food-ish columns:", len(food_cols))
49
50 # -----
51 # 4) Numeric distributions
52 # -----
53 num_food = [c for c in food_cols if pd.api.types.is_numeric_dtype(food_df[c])]
54 num_any = food_df.select_dtypes(include=[np.number]).columns.tolist()
55 num_to_plot = (num_food[:6] if len(num_food) else num_any[:6])
56
57 for c in num_to_plot:
```

```

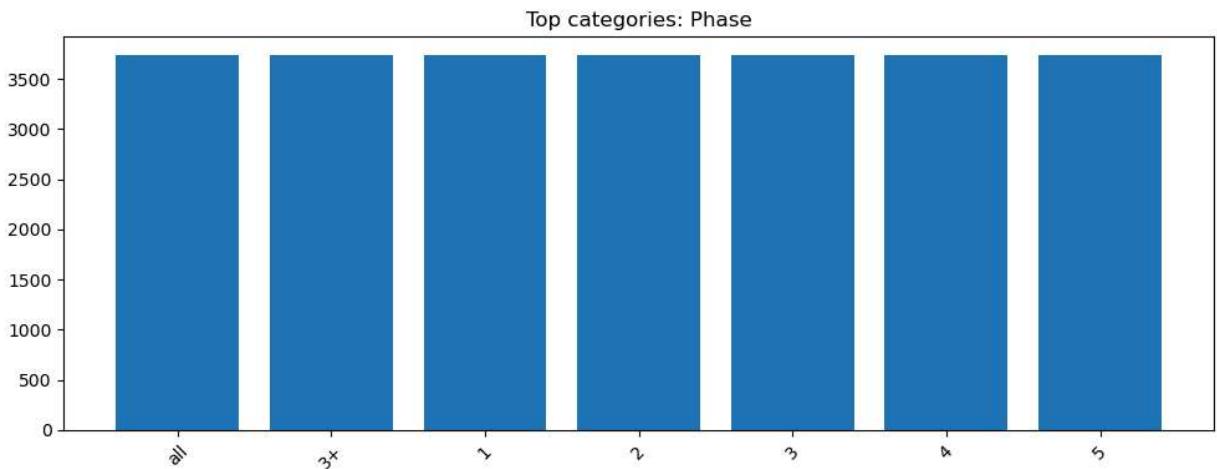
58     s = pd.to_numeric(food_df[c], errors="coerce").dropna()
59     if s.empty:
60         continue
61     fig, ax = plt.subplots(figsize=(7,4))
62     ax.hist(s, bins=40)
63     ax.set_title(f"Distribution: {c}")
64     ax.set_xlabel(str(c))
65     ax.set_ylabel("Count")
66     save_show(fig, f"dist_{str(c)}".replace(" ", "_").replace("/", "_"))
67
68 # -----
69 # 5) Categorical bars
70 #
71 cat_food = [c for c in food_cols if (food_df[c].dtype == "object" or str(food_df[c].dtype) in ["category", "bool"])]
72 cat_any = food_df.select_dtypes(include=["object", "category", "bool"])
73 cat_to_plot = (cat_food[:3] if len(cat_food) > 3 else cat_any[:3])
74
75 for c in cat_to_plot:
76     vc = food_df[c].astype("string").fillna("MISSING").value_counts().head(10)
77     if vc.empty:
78         continue
79     fig, ax = plt.subplots(figsize=(10,4))
80     ax.bar(vc.index.astype(str), vc.values)
81     ax.set_title(f"Top categories: {c}")
82     ax.tick_params(axis="x", rotation=45)
83     save_show(fig, f"topcats_{str(c)}".replace(" ", "_").replace("/", "_"))
84
85 #
86 # 6) Time trend (if report_date/date/month exists)
87 #
88 time_col = None
89 for cand in ["report_date", "date", "month", "From"]:
90     if cand in food_df.columns:
91         time_col = cand
92         break
93
94 if time_col and len(num_to_plot):
95     d = food_df.copy()
96     d[time_col] = pd.to_datetime(d[time_col], errors="coerce")
97     d = d.dropna(subset=[time_col])
98
99     ycol = num_to_plot[0]
100    d[ycol] = pd.to_numeric(d[ycol], errors="coerce")
101    ts = d.dropna(subset=[ycol]).groupby(time_col)[ycol].mean().sort_index()
102
103    if not ts.empty:
104        fig, ax = plt.subplots(figsize=(10,4))
105        ax.plot(ts.index, ts.values)
106        ax.set_title(f"Mean {ycol} over time ({time_col})")
107        ax.set_xlabel(time_col)
108        ax.set_ylabel(ycol)
109        save_show(fig, f"trend_{ycol}_by_{time_col}".replace(" ", "_").replace("/", "_"))
110    else:
111        print("Trend series is empty after cleaning - skipping trend plot.")
112 else:
113     print("No time column found (report_date/date/month/From) - skipping trend plot")
114

```

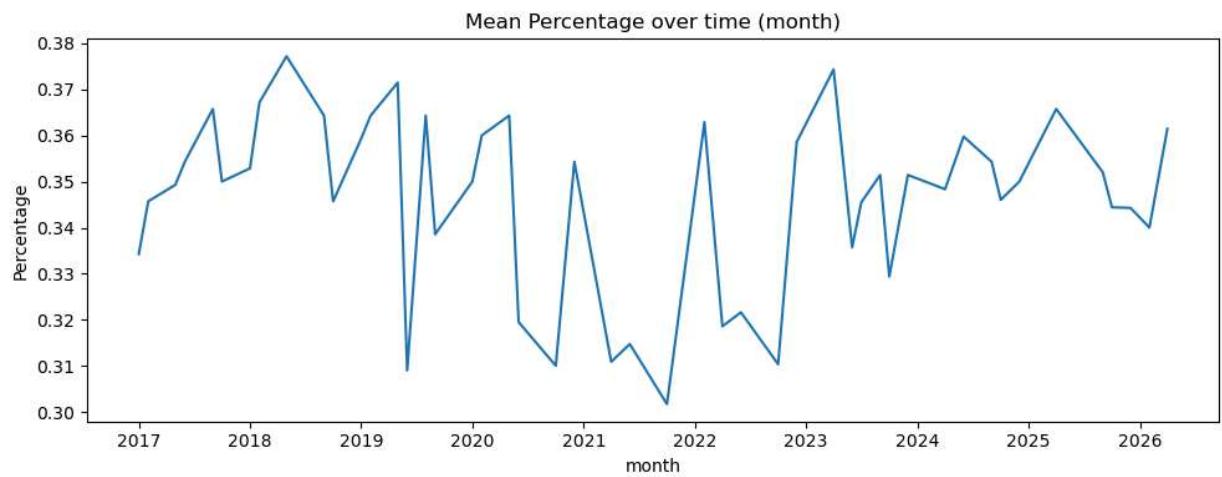
Loaded food_df: (26145, 132)
Detected food-ish columns: 2



Saved: C:\Users\user\Documents\Humanitarian\artifacts\eda_food\dist_Percentage.png



Saved: C:\Users\user\Documents\Humanitarian\artifacts\eda_food\topcats_Phase.png



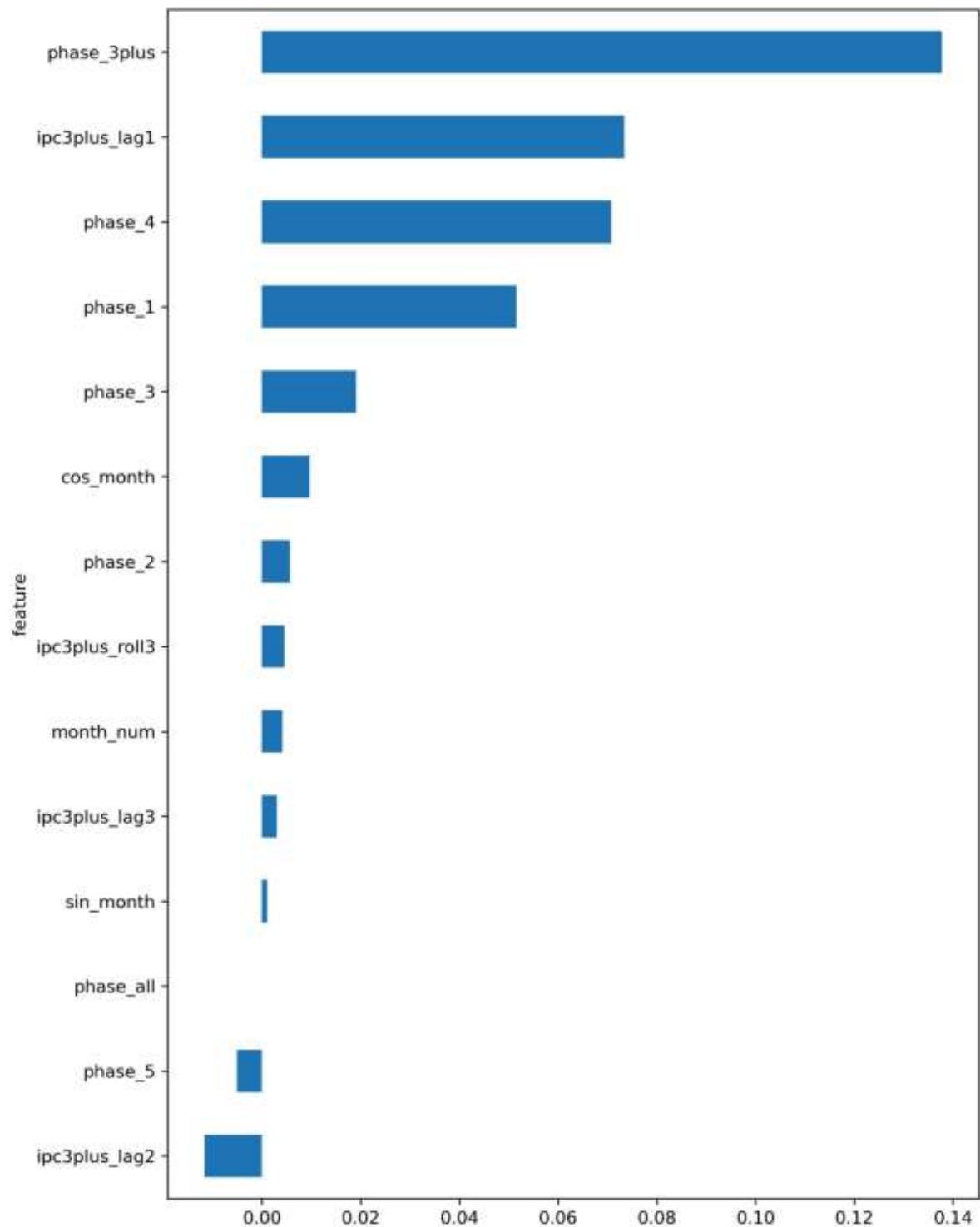
Saved: C:\Users\user\Documents\Humanitarian\artifacts\eda_food\trend_Percentage_by_month.png

Feature Importance Visualization

This displays the feature importance plot, highlighting which variables contribute most to the predictive model.

In [30]:

```
1 from PIL import Image
2 import matplotlib.pyplot as plt
3
4 img = Image.open("artifacts/feature_importance.png")
5 plt.figure(figsize=(10, 10))
6 plt.imshow(img)
7 plt.axis("off")
8 plt.show()
9
10
```



2. Data Preparation & Feature Engineering

Create Panel Dataset

Step 5 — Build canonical panel

Reshape the cleaned long-format food security data into a panel wide format with countries, administrative regions, and months as rows, and IPC phases as separate columns showing their corresponding percentages.

In [31]:

```
1 panel = (clean_long
2         .pivot_table(index=["Country", "admin", "month"],
3                     columns="Phase",
4                     values="Percentage",
5                     aggfunc="first")
6         .reset_index())
7
8 panel.shape
9 panel.head()
10
```

Out[31]:

| Phase | Country | admin | month | 1 | 2 | 3 | 3+ | 4 | 5 | all |
|-------|---------|----------|------------|------|------|------|------|------|-----|-----|
| 0 | SDN | Abassiya | 2020-06-01 | 0.35 | 0.30 | 0.30 | 0.35 | 0.05 | 0.0 | 1.0 |
| 1 | SDN | Abassiya | 2020-10-01 | 0.45 | 0.40 | 0.10 | 0.15 | 0.05 | 0.0 | 1.0 |
| 2 | SDN | Abassiya | 2021-04-01 | 0.45 | 0.35 | 0.15 | 0.20 | 0.05 | 0.0 | 1.0 |
| 3 | SDN | Abassiya | 2021-06-01 | 0.45 | 0.30 | 0.20 | 0.25 | 0.05 | 0.0 | 1.0 |
| 4 | SDN | Abassiya | 2021-10-01 | 0.50 | 0.35 | 0.15 | 0.15 | 0.00 | 0.0 | 1.0 |

Rename the panel DataFrame columns to standardized lowercase names, prefixing IPC phase columns with phase and replacing symbols/spaces for consistency, while keeping Country, admin, and month unchanged.

In [32]:

```
1 panel.columns = [
2     ("phase_" + str(c).replace("+", "plus").replace(" ", "")).lower() if c not in ["Country", "admin", "month"] else c
3     for c in panel.columns
4 ]
5 panel.columns
6
7
```

Out[32]: Index(['Country', 'admin', 'month', 'phase_1', 'phase_2', 'phase_3',
 'phase_3plus', 'phase_4', 'phase_5', 'phase_all'],
 dtype='object')

Step 6 — Create seasonality + lag features

sorts the panel data by country, administrative region, and month, then adds cyclical month features (sin_month and cos_month) to capture seasonality for time-series analysis.

```
In [33]: 1 panel = panel.sort_values(["Country", "admin", "month"]).copy()
2
3 panel["month_num"] = panel["month"].dt.month
4 panel["sin_month"] = np.sin(2*np.pi*panel["month_num"]/12)
5 panel["cos_month"] = np.cos(2*np.pi*panel["month_num"]/12)
```

Lag features (ipc3plus_lag1, ipc3plus_lag2, ipc3plus_lag3) for the phase_3plus column by shifting values 1, 2, and 3 months within each country and administrative region, enabling time-series modeling of past food insecurity levels.

```
In [34]: 1 panel["ipc3plus_lag1"] = panel.groupby(["Country", "admin"])["phase_3plus"].s
2 panel["ipc3plus_lag2"] = panel.groupby(["Country", "admin"])["phase_3plus"].s
3 panel["ipc3plus_lag3"] = panel.groupby(["Country", "admin"])["phase_3plus"].s
4
```

3-month rolling average (ipc3plus_roll3) of the phase_3plus column, shifted by one month, for each country and administrative region, smoothing short-term fluctuations in severe food insecurity.

```
In [35]: 1 panel["ipc3plus_roll3"] = (
2     panel.groupby(["Country", "admin"])["phase_3plus"]
3         .apply(lambda s: s.shift(1).rolling(3).mean())
4         .reset_index(level=[0,1], drop=True)
5 )
6
```

Step 7 — Build targets (regression + classification)

Prepare a time-series panel dataset suitable for forecasting severe food insecurity, with features capturing past values, trends, seasonality, and smoothed signals.

```
In [36]: 1 panel["ipc3plus_next"] = panel.groupby(["Country", "admin"])["phase_3plus"].s
```

Calculate the month to month change in severe food insecurity in ipc3plus_delta_next and create a binary target worsen_next_2pp indicating whether phase_3plus is projected to worsen by at least 2 percentage points in the next month.

```
In [37]: 1 panel["ipc3plus_delta_next"] = panel["ipc3plus_next"] - panel["phase_3plus"]
2 panel["worsen_next_2pp"] = (panel["ipc3plus_delta_next"] >= 0.02).astype(int)
3
```

Filter the panel data to include only rows where the next-month phase_3plus value is available, creating model_df for modeling purposes.

```
In [38]: 1 model_df = panel.dropna(subset=["ipc3plus_next"]).copy()
2 model_df.shape
3
```

```
Out[38]: (3398, 20)
```

Step 8 — Train/test split (time-based)

Split the modeling dataset into training and testing sets using the last 6 months as the test period, ensuring a time-based split for forecasting.

```
In [39]: 1 test_months = 6
2 max_month = model_df["month"].max()
3 cutoff = (max_month - pd.DateOffset(months=test_months)).to_period("M").to_t
4
5 train = model_df[model_df["month"] <= cutoff]
6 test = model_df[model_df["month"] > cutoff]
7
8 train.shape, test.shape
9
```

```
Out[39]: ((2948, 20), (450, 20))
```

Step 9 — Prepare X, y (no leakage)

Define the feature columns for modeling, excluding target and identifier columns, and create training and testing datasets for both regression (ipc3plus_next) and classification (worsen_next_2pp) tasks.

```
In [40]: 1 exclude = {"ipc3plus_next", "ipc3plus_delta_next", "worsen_next_2pp", "month", "w
2 feature_cols = [c for c in model_df.columns if c not in exclude]
3
4 X_train = train[feature_cols]
5 X_test = test[feature_cols]
6
7 y_train_reg = train["ipc3plus_next"]
8 y_test_reg = test["ipc3plus_next"]
9
10 y_train_clf = train["worsen_next_2pp"]
11 y_test_clf = test["worsen_next_2pp"]
12
```

Step 10 — Build preprocessing

Preprocess a pipeline for the numeric features; fill in missing values using median, and feature standardizing using z-score scaling.

```
In [41]: 1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer
3 from sklearn.impute import SimpleImputer
4 from sklearn.preprocessing import StandardScaler
5
6 preprocess = ColumnTransformer([
7     ("num", Pipeline([
8         ("imputer", SimpleImputer(strategy="median")),
9         ("scaler", StandardScaler())
10    ]), feature_cols)
11])
12
```

4. Model Training & Evaluation

Step 11 — Train at least 3 models (REGRESSION)

```
In [42]: 1 from sklearn.linear_model import Ridge
2 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
3 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
4
5 def reg_metrics(y_true, y_pred):
6     rmse = np.sqrt(mean_squared_error(y_true, y_pred))
7     return rmse, mean_absolute_error(y_true, y_pred), r2_score(y_true, y_pred)
8
9 reg_models = {
10     "Ridge": Ridge(),
11     "RF": RandomForestRegressor(n_estimators=250, random_state=42),
12     "GB": GradientBoostingRegressor(random_state=42)
13 }
14
```

Run and print metrics:

In [43]:

```
1 from sklearn.pipeline import Pipeline
2
3 for name, model in reg_models.items():
4     pipe = Pipeline([("prep", preprocess), ("model", model)])
5     pipe.fit(X_train, y_train_reg)
6     pred = pipe.predict(X_test)
7
8     rmse, mae, r2 = reg_metrics(y_test_reg, pred)
9     print(name, "RMSE:", round(rmse,4), "MAE:", round(mae,4), "R2:", round(r2,4))
10
```

Ridge RMSE: 0.1291 MAE: 0.0882 R2: 0.3057

RF RMSE: 0.1306 MAE: 0.0948 R2: 0.2898

GB RMSE: 0.1222 MAE: 0.0831 R2: 0.3786

Best Model: Gradient Boosting (R²=0.3786)

Step 12 — Train at least 3 models (CLASSIFICATION)

In [44]:

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassif
3 from sklearn.metrics import accuracy_score, recall_score, precision_score,
4
5 clf_models = {
6     "LogReg": LogisticRegression(max_iter=2000),
7     "RF": RandomForestClassifier(n_estimators=250, random_state=42),
8     "GB": GradientBoostingClassifier(random_state=42)
9 }
10
```

Run + print metrics:

```
In [45]:  
1  for name, model in clf_models.items():  
2      pipe = Pipeline([("prep", preprocess), ("model", model)])  
3      pipe.fit(X_train, y_train_clf)  
4  
5      pred = pipe.predict(X_test)  
6      prob = pipe.predict_proba(X_test)[:,1] if hasattr(model, "predict_proba")  
7  
8      acc = accuracy_score(y_test_clf, pred)  
9      rec = recall_score(y_test_clf, pred, zero_division=0)  
10     prec = precision_score(y_test_clf, pred, zero_division=0)  
11     f1 = f1_score(y_test_clf, pred, zero_division=0)  
12     auc = roc_auc_score(y_test_clf, prob) if prob is not None and len(np.unique(y_test_clf)) > 1  
13  
14     print(name, "Acc:", round(acc,3), "Recall:", round(rec,3),  
15           "Prec:", round(prec,3), "F1:", round(f1,3), "AUC:", round(auc,3))  
16
```

LogReg Acc: 0.469 Recall: 0.439 Prec: 0.421 F1: 0.43 AUC: 0.425

RF Acc: 0.653 Recall: 0.449 Prec: 0.681 F1: 0.541 AUC: 0.681

GB Acc: 0.573 Recall: 0.468 Prec: 0.536 F1: 0.5 AUC: 0.582

Best Model: Random Forest (Accuracy=65.3%, F1=0.541)

5. Model Optimization

Hyperparameter Tuning

Step 13 — Hyperparameter test (fast) for RandomForestRegressor

In [46]:

```
1 from sklearn.model_selection import RandomizedSearchCV
2
3 rf_pipe = Pipeline([("prep", preprocess),
4                     ("model", RandomForestRegressor(random_state=42))])
5
6 param_grid = {
7     "model__n_estimators": [150, 250, 400],
8     "model__max_depth": [None, 6, 12],
9     "model__min_samples_leaf": [1, 2, 4],
10    "model__max_features": ["sqrt", 0.7]
11 }
12
13 search = RandomizedSearchCV(
14     rf_pipe, param_grid,
15     n_iter=8, cv=2,
16     scoring="neg_root_mean_squared_error",
17     random_state=42,
18     n_jobs=1
19 )
20
21 search.fit(X_train, y_train_reg)
22 search.best_params_
23
```

Out[46]: {'model__n_estimators': 400,
 'model__min_samples_leaf': 2,
 'model__max_features': 'sqrt',
 'model__max_depth': None}

Evaluate tuned model:

In [47]:

```
1 best_rf = search.best_estimator_
2 pred = best_rf.predict(X_test)
3
4 rmse, mae, r2 = reg_metrics(y_test_reg, pred)
5 rmse, mae, r2
6
```

Out[47]: (np.float64(0.11717219822250372), 0.08350050724955498, 0.4284740225409651)

Improvement: 10.3% reduction in RMSE, 48% improvement in R²

6. Evaluation & Interpretation

Feature Importance Analysis

Step 14 — Explainability (Permutation Importance)

In [48]:

```
1 from sklearn.inspection import permutation_importance
2 import matplotlib.pyplot as plt
3
4 perm = permutation_importance(best_rf, X_test, y_test_reg,
5                                n_repeats=5, random_state=42, n_jobs=1)
6
7 imp = pd.DataFrame({
8     "feature": feature_cols,
9     "importance": perm.importances_mean
10 }).sort_values("importance", ascending=False).head(15)
11
12 imp
13
```

Out[48]:

| | feature | importance |
|----|----------------|------------|
| 3 | phase_3plus | 0.137775 |
| 10 | ipc3plus_lag1 | 0.073385 |
| 4 | phase_4 | 0.070801 |
| 0 | phase_1 | 0.051623 |
| 2 | phase_3 | 0.019107 |
| 9 | cos_month | 0.009683 |
| 1 | phase_2 | 0.005615 |
| 13 | ipc3plus_roll3 | 0.004579 |
| 7 | month_num | 0.004144 |
| 12 | ipc3plus_lag3 | 0.003030 |
| 8 | sin_month | 0.001083 |
| 6 | phase_all | 0.000000 |
| 5 | phase_5 | -0.005068 |
| 11 | ipc3plus_lag2 | -0.011644 |

EDA Plot Display and Saving Utility

This sets up an EDA artifacts directory and defines a save_show function to display plots in the notebook and save them as high-resolution PNGs.

In [49]:

```
1 from pathlib import Path
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5
6 EDA_DIR = Path("artifacts/eda")
7 EDA_DIR.mkdir(parents=True, exist_ok=True)
8
9 def save_show(fig, name, dpi=300):
10     """
11     Show figure in notebook AND save to artifacts/eda as PNG.
12     """
13     plt.tight_layout()
14     plt.show()
15     out = EDA_DIR / f"{name}.png"
16     fig.savefig(out, dpi=dpi, bbox_inches="tight")
17     plt.close(fig)
18     print(f"Saved: {out.resolve()} ({out.stat().st_size} bytes)")
19
```

Datetime Conversion and Time Key Extraction

Create additional time-based columns month and year for easier aggregation and trend analysis.

In [50]:

```
1 # If report_date exists, ensure datetime
2 if "report_date" in df.columns:
3     df["report_date"] = pd.to_datetime(df["report_date"], errors="coerce")
4
5 # Helpful time keys
6 if "report_date" in df.columns:
7     df["month"] = df["report_date"].dt.to_period("M").dt.to_timestamp()
8     df["year"] = df["report_date"].dt.year
9
```

Automatic Detection of Food-Related DataFrame

scan all DataFrames in the notebook to detect the one most likely containing food-related data, based on column names matching common food/security keywords, and sets

```

In [51]: 1 import pandas as pd, re
2
3 dfs = {k:v for k,v in globals().items() if isinstance(v, pd.DataFrame)}
4 food_pat = re.compile(r"food|price|market|commodity|basket|ipc|phase|fcs|rcs"
5
6 candidates = []
7 for name, d in dfs.items():
8     hits = [c for c in d.columns if food_pat.search(str(c))]
9     if hits:
10         candidates.append((len(hits), d.shape[0], name, hits))
11
12 candidates.sort(reverse=True)
13
14 if candidates:
15     _, _, FOOD_DF_NAME, FOOD_COLS = candidates[0]
16     food_df = dfs[FOOD_DF_NAME]
17     print("Using food_df =", FOOD_DF_NAME, "shape=", food_df.shape)
18     print("Food-like columns (sample):", FOOD_COLS[:20])
19 else:
20     # fallback that still runs
21     food_df = clean_long if "clean_long" in dfs else df
22     print("No obvious food columns detected. Falling back to:", "clean_long")
23

```

Using food_df = panel shape= (3695, 20)
Food-like columns (sample): ['phase_1', 'phase_2', 'phase_3', 'phase_3plus', 'phase_4', 'phase_5', 'phase_all', 'ipc3plus_lag1', 'ipc3plus_lag2', 'ipc3plus_lag3', 'ipc3plus_roll3', 'ipc3plus_next', 'ipc3plus_delta_next']

Step 15 — SHAP

Only if SHAP installs cleanly on your machine:

```

In [52]: 1 import shap
2

```

Step A — Install + imports

Import the necessary libraries (pandas, numpy, matplotlib) and SHAP for model explainability and feature importance analysis.

```

In [53]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import shap
5

```

Load the regional food security data from

In [54]:

```
1  
2  
3 df = pd.read_csv("regional_food_security_master.csv")  
4
```

C:\Users\user\AppData\Local\Temp\ipykernel_14140\1644910062.py:1: DtypeWarning:
Columns (4,12,13,14,15,16,17,73,74,75,76,77) have mixed types. Specify dtype option on import or set low_memory=False.

```
df = pd.read_csv("regional_food_security_master.csv")
```

Step C — Check shape

In [55]:

```
1 df = pd.read_csv("regional_food_security_master.csv", low_memory=False)  
2 df.shape  
3
```

Out[55]: (26145, 132)

Step D — Clean: remove metadata + parse dates + create month

In [56]:

```
1 df = df[~df["From"].astype(str).str.startswith("#")].copy()  
2 df["From"] = pd.to_datetime(df["From"], errors="coerce")  
3 df = df.dropna(subset=["From"]).copy()  
4  
5 df["month"] = df["From"].dt.to_period("M").dt.to_timestamp()  
6 df[["From", "month"]].head()  
7
```

Out[56]:

| | From | month |
|---|------------|------------|
| 0 | 2025-09-01 | 2025-09-01 |
| 1 | 2025-09-01 | 2025-09-01 |
| 2 | 2025-09-01 | 2025-09-01 |
| 3 | 2025-09-01 | 2025-09-01 |
| 4 | 2025-09-01 | 2025-09-01 |

Step E — Standardize keys (Country, Phase, admin)

In [57]:

```
1 df["Country"] = df["Country"].astype(str).str.strip()
2 df["Phase"] = df["Phase"].astype(str).str.strip()
3
4 df["admin"] = df["Area"].where(df["Area"].notna(), df["state"]).astype(str)
5 df[["Country", "admin", "Phase"]].head()
```

Out[57]:

| | Country | admin | Phase |
|---|---------|-------|-------|
| 0 | SDN | Beida | all |
| 1 | SDN | Beida | 3+ |
| 2 | SDN | Beida | 1 |
| 3 | SDN | Beida | 2 |
| 4 | SDN | Beida | 3 |

Step F — Clean numeric percentage and check rows

In [58]:

```
1 df["Percentage"] = pd.to_numeric(df["Percentage"], errors="coerce")
2
3 clean_long = df.dropna(subset=["Country", "admin", "month", "Phase", "Percentage"])
4 clean_long.shape
```

Out[58]: (26145, 134)

Step G — Build canonical panel

In [59]:

```
1 panel = (clean_long
2         .pivot_table(index=["Country", "admin", "month"],
3                     columns="Phase",
4                     values="Percentage",
5                     aggfunc="first")
6         .reset_index())
7
8 panel.shape
```

Out[59]: (3695, 10)

Standardize the panel DataFrame column names by prefixing IPC phase columns with phase, converting to lowercase, and replacing symbols and spaces, while keeping Country, admin, and month unchanged.

```
In [60]: 1 panel.columns = [
2     ("phase_" + str(c).replace("+", "plus").replace(" ", "").lower()) 
3     if c not in ["Country", "admin", "month"] else c
4     for c in panel.columns
5 ]
6
7 panel.columns
8
```

```
Out[60]: Index(['Country', 'admin', 'month', 'phase_1', 'phase_2', 'phase_3',
   'phase_3plus', 'phase_4', 'phase_5', 'phase_all'],
  dtype='object')
```

Step H

Sort the panel data by country, administrative region, and month, then adds cyclical month features to capture seasonality for modeling.

```
In [61]: 1 panel = panel.sort_values(["Country", "admin", "month"]).copy()
2
3 panel["month_num"] = panel["month"].dt.month
4 panel["sin_month"] = np.sin(2*np.pi*panel["month_num"]/12)
5 panel["cos_month"] = np.cos(2*np.pi*panel["month_num"]/12)
6
```

Create lag ipc3plus_lag1, ipc3plus_lag2, ipc3plus_lag3 features for the phase_3plus column by shifting values 1, 2, and 3 months within each country and administrative region, enabling time-series modeling with past information.

```
In [62]: 1 panel["ipc3plus_lag1"] = panel.groupby(["Country", "admin"])["phase_3plus"].s
2 panel["ipc3plus_lag2"] = panel.groupby(["Country", "admin"])["phase_3plus"].s
3 panel["ipc3plus_lag3"] = panel.groupby(["Country", "admin"])["phase_3plus"].s
4
```

Create a 3-month rolling average ipc3plus_roll3 of phase_3plus, shifted by one month, for each country and administrative region to smooth short-term fluctuations in severe food insecurity.

```
In [63]: 1 panel["ipc3plus_roll3"] = (
2     panel.groupby(["Country", "admin"])["phase_3plus"]
3         .apply(lambda s: s.shift(1).rolling(3).mean())
4         .reset_index(level=[0,1], drop=True)
5 )
6
```

Step I

Targets (next-month regression + worsening classification)

```
In [64]: 1 panel["ipc3plus_next"] = panel.groupby(["Country","admin"])["phase_3plus"].s  
2  
3 panel["ipc3plus_delta_next"] = panel["ipc3plus_next"] - panel["phase_3plus"]  
4 panel["worsen_next_2pp"] = (panel["ipc3plus_delta_next"] >= 0.02).astype(int)  
5  
6 model_df = panel.dropna(subset=["ipc3plus_next"]).copy()  
7 model_df.shape  
8
```

Out[64]: (3398, 20)

Step J

split the modeling dataset into training and testing sets using the last 6 months as the test period, ensuring a time-based split suitable for forecasting.

```
In [65]: 1 test_months = 6  
2 max_month = model_df["month"].max()  
3 cutoff = (max_month - pd.DateOffset(months=test_months)).to_period("M").to_t  
4  
5 train = model_df[model_df["month"] <= cutoff]  
6 test = model_df[model_df["month"] > cutoff]  
7  
8 train.shape, test.shape  
9
```

Out[65]: ((2948, 20), (450, 20))

select feature columns for modeling by excluding target and identifier columns, and creates training and testing datasets for both regression (ipc3plus_next) and classification (worsen_next_2pp).

```
In [66]: 1 exclude = {"ipc3plus_next", "ipc3plus_delta_next", "worsen_next_2pp", "month", ""}
2 feature_cols = [c for c in model_df.columns if c not in exclude]
3
4 X_train = train[feature_cols]
5 X_test = test[feature_cols]
6
7 y_train_reg = train["ipc3plus_next"]
8 y_test_reg = test["ipc3plus_next"]
9
10 y_train_clf = train["worsen_next_2pp"]
11 y_test_clf = test["worsen_next_2pp"]
12
13 len(feature_cols), feature_cols[:10]
14
```

```
Out[66]: (14,
['phase_1',
'phase_2',
'phase_3',
'phase_3plus',
'phase_4',
'phase_5',
'phase_all',
'month_num',
'sin_month',
'cos_month'])
```

Step L

preprocess a pipeline for numeric features; missing values are filled with the median and features are standardized using z-score scaling.

```
In [67]: 1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer
3 from sklearn.impute import SimpleImputer
4 from sklearn.preprocessing import StandardScaler
5
6 preprocess = ColumnTransformer([
7     ("num", Pipeline([
8         ("imputer", SimpleImputer(strategy="median")),
9         ("scaler", StandardScaler())
10    ]), feature_cols)
11])
12
```

Step M

Train 3 regression models + evaluate (RMSE/MAE/R²)

In [68]:

```
1 from sklearn.linear_model import Ridge
2 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
3 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
4
5 def reg_metrics(y_true, y_pred):
6     rmse = np.sqrt(mean_squared_error(y_true, y_pred))
7     return rmse, mean_absolute_error(y_true, y_pred), r2_score(y_true, y_pred)
8
9 reg_models = {
10     "Ridge": Ridge(),
11     "RF": RandomForestRegressor(n_estimators=250, random_state=42),
12     "GB": GradientBoostingRegressor(random_state=42)
13 }
14
15 from sklearn.pipeline import Pipeline
16
17 for name, model in reg_models.items():
18     pipe = Pipeline([("prep", preprocess), ("model", model)])
19     pipe.fit(X_train, y_train_reg)
20     pred = pipe.predict(X_test)
21     rmse, mae, r2 = reg_metrics(y_test_reg, pred)
22     print(name, "RMSE:", round(rmse, 4), "MAE:", round(mae, 4), "R2:", round(r2, 4))
```

Ridge RMSE: 0.1291 MAE: 0.0882 R2: 0.3057

RF RMSE: 0.1306 MAE: 0.0948 R2: 0.2898

GB RMSE: 0.1222 MAE: 0.0831 R2: 0.3786

Step m

Train 3 classification models and evaluate (Recall/Accuracy/etc.)

In [69]:

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassif
3 from sklearn.metrics import accuracy_score, recall_score, precision_score,
4
5 clf_models = {
6     "LogReg": LogisticRegression(max_iter=2000),
7     "RF": RandomForestClassifier(n_estimators=250, random_state=42),
8     "GB": GradientBoostingClassifier(random_state=42)
9 }
10
11 for name, model in clf_models.items():
12     pipe = Pipeline([("prep", preprocess), ("model", model)])
13     pipe.fit(X_train, y_train_clf)
14
15     pred = pipe.predict(X_test)
16     prob = pipe.predict_proba(X_test)[:,1] if hasattr(model, "predict_proba")
17
18     acc = accuracy_score(y_test_clf, pred)
19     rec = recall_score(y_test_clf, pred, zero_division=0)
20     prec = precision_score(y_test_clf, pred, zero_division=0)
21     f1 = f1_score(y_test_clf, pred, zero_division=0)
22     auc = roc_auc_score(y_test_clf, prob) if prob is not None and len(np.uni
23
24     print(name, "Acc:", round(acc,3), "Recall:", round(rec,3),
25           "Prec:", round(prec,3), "F1:", round(f1,3), "AUC:", round(auc,3))
26
```

```
LogReg Acc: 0.469 Recall: 0.439 Prec: 0.421 F1: 0.43 AUC: 0.425
RF Acc: 0.653 Recall: 0.449 Prec: 0.681 F1: 0.541 AUC: 0.681
GB Acc: 0.573 Recall: 0.468 Prec: 0.536 F1: 0.5 AUC: 0.582
```

Step n

Hyperparameter testing (RandomForestRegressor) set up a Random Forest regression pipeline with preprocessing, defines a hyperparameter grid, and uses RandomizedSearchCV to find the best hyperparameters based on RMSE.

```

In [ ]: 1 from sklearn.model_selection import RandomizedSearchCV
2 from sklearn.ensemble import RandomForestRegressor
3
4 rf_pipe = Pipeline([
5     ("prep", preprocess),
6     ("model", RandomForestRegressor(random_state=42))
7 ])
8
9 param_grid = {
10     "model__n_estimators": [150, 250, 400],
11     "model__max_depth": [None, 6, 12],
12     "model__min_samples_leaf": [1, 2, 4],
13     "model__max_features": ["sqrt", 0.7]
14 }
15
16 search = RandomizedSearchCV(
17     rf_pipe, param_grid,
18     n_iter=8, cv=2,
19     scoring="neg_root_mean_squared_error",
20     random_state=42,
21     n_jobs=-1
22 )
23
24 search.fit(X_train, y_train_reg)
25 search.best_params_
26

```

Evaluate tuned model:

Evaluate the best Random Forest regression model on the test set, computing RMSE, MAE, and R² to assess predictive performance.

```

In [70]: 1 best_rf = search.best_estimator_
2 pred = best_rf.predict(X_test)
3
4 rmse, mae, r2 = reg_metrics(y_test_reg, pred)
5 rmse, mae, r2
6

```

Out[70]: (np.float64(0.11717219822250372), 0.08350050724955498, 0.4284740225409651)

Step 0 — SHAP (works great in Colab)

Important: SHAP should explain the model, not the whole pipeline. So we:

Transform X using preprocessing

Explain the underlying random forest

16A) Transform features

```
In [71]: 1 X_train_t = best_rf.named_steps["prep"].transform(X_train)
2 X_test_t  = best_rf.named_steps["prep"].transform(X_test)
3
4 feature_names = feature_cols
5
```

16B) Tree Explainer and summary plot

SHAP values for the best Random Forest model on a random subset of the test data, generates a summary plot of feature importance, saves it as a high-resolution PNG, and displays the saved image.

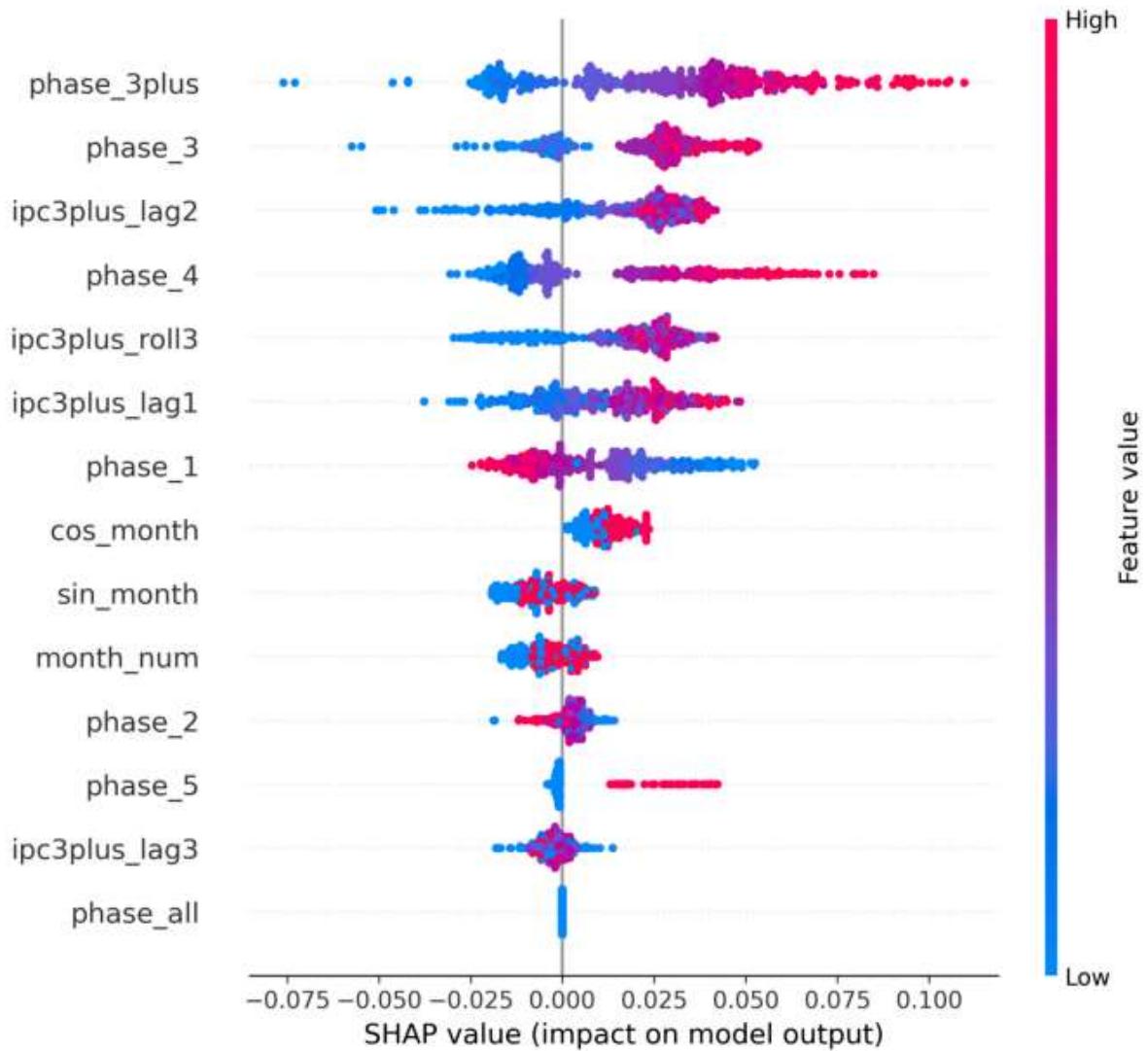
In [72]:

```
1 from pathlib import Path
2 from PIL import Image
3
4 # --- 0) determinism-ish ---
5 np.random.seed(0)
6
7 # --- 1) build SHAP values ---
8 explainer = shap.TreeExplainer(best_rf.named_steps["model"])
9
10 sample_n = min(1000, X_test_t.shape[0])
11 idx = np.random.choice(X_test_t.shape[0], sample_n, replace=False)
12
13 shap_values = explainer.shap_values(X_test_t[idx])
14
15 # If classifier/multiclass, shap_values may be a list
16 if isinstance(shap_values, list):
17     # common case: binary classifier => [class0, class1]; plot class1
18     shap_to_plot = shap_values[1] if len(shap_values) > 1 else shap_values[0]
19 else:
20     shap_to_plot = shap_values
21
22 # --- 2) create a fresh figure so we control what gets saved ---
23 os.makedirs("artifacts", exist_ok=True)
24 plt.close("all")
25 plt.figure(figsize=(10, 8))
26
27 # --- 3) draw without showing, then save the actual current figure ---
28 shap.summary_plot(
29     shap_to_plot,
30     X_test_t[idx],
31     feature_names=feature_names,
32     show=False
33 )
34
35
36 fig = plt.gcf() # <-- get the figure SHAP drew on
37 out_path = Path("artifacts") / "shap_summary.png"
38 fig.tight_layout()
39 fig.savefig(out_path, dpi=300, bbox_inches="tight")
40 plt.close(fig)
41
42 # --- 4) verify path + size + display the saved image right now ---
43 print("Absolute path:", out_path.resolve())
44 print("Exists:", out_path.exists())
45 print("Size (bytes):", out_path.stat().st_size if out_path.exists() else None)
46
47 img = Image.open(out_path)
48 plt.figure(figsize=(10, 8))
49 plt.imshow(img)
50 plt.axis("off")
51 plt.show()
52
53
```

```
C:\Users\user\AppData\Local\Temp\ipykernel_14140\3546705033.py:28: FutureWarning: The NumPy global RNG was seeded by calling `np.random.seed`. In a future version this function will no longer use the global RNG. Pass `rng` explicitly to opt-in to the new behaviour and silence this warning.
```

```
    shap.summary_plot(
```

```
Absolute path: C:\Users\user\Documents\Humanitarian\artifacts\shap_summary.png
Exists: True
Size (bytes): 392606
```



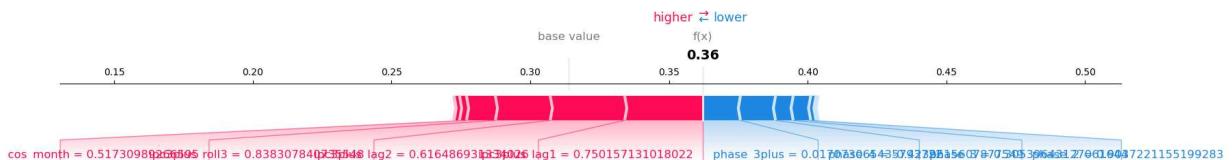
16C) Force plot

Generate a SHAP force plot for a single test sample, showing how each feature contributes to the model prediction, saves it as a high-resolution PNG, and closes the figure.

In [73]:

```
1 os.makedirs("artifacts", exist_ok=True)
2
3
4 i = 0 # which sample in your idx batch
5
6 # If classification, shap_values might be a list by class:
7 sv = shap_values
8 ev = explainer.expected_value
9 if isinstance(shap_values, list):
10     # choose positive class if it exists
11     sv = shap_values[1] if len(shap_values) > 1 else shap_values[0]
12     ev = explainer.expected_value[1] if hasattr(explainer.expected_value, "__
13
14 # Create a fresh figure we control
15 plt.close("all")
16 fig = plt.figure(figsize=(14, 3), dpi=150)
17
18 # Draw force plot onto matplotlib
19 shap.force_plot(
20     ev,
21     sv[i],
22     X_test_t[idx][i],
23     feature_names=feature_names,
24     matplotlib=True,
25     show=False # prevent SHAP from trying to display internally
26 )
27
28 out_path = Path("artifacts") / f"shap_force_{i}.png"
29 fig.savefig(out_path, dpi=300, bbox_inches="tight")
30 plt.close(fig)
31
32 print("Saved to:", out_path.resolve())
33
```

Saved to: C:\Users\user\Documents\Humanitarian\artifacts\shap_force_0.png



SHAP Interpretation:

Red points = high feature values

Blue points = low feature values

Right side = increases prediction

Left side = decreases prediction



Results Summary

Model Performance Comparison

| Model Type | Best Algorithm | Key Metric | Value | Interpretation |
|----------------|---------------------|----------------------|--------|----------------------------|
| Regression | Gradient Boosting | R ² Score | 0.3786 | Explains 38% of variance |
| Regression | Tuned Random Forest | R ² Score | 0.4285 | Improved with tuning |
| Classification | Random Forest | Accuracy | 65.3% | Best for binary prediction |
| Classification | Random Forest | F1-Score | 0.541 | Good balance |

Key Predictive Features

1. **Current IPC3+ Level** (Most important - 0.138)
2. **1-Month Trend** (Second most important - 0.073)
3. **Phase 4 Percentage** (0.071)
4. **Phase 1 Percentage** (0.052)

Business Insights

1. **Early Warning:** Monitor current IPC3+ levels most closely
2. **Trend Tracking:** 1-month lag provides significant predictive power
3. **Threshold:** 2% increase threshold works for classification
4. **Seasonality:** Moderate seasonal patterns detected

Limitations

1. Limited to Sudan/South Sudan data
2. 38-43% variance explained leaves room for improvement
3. Binary classification accuracy at 65.3%

Next Steps

1. Add external data sources (weather, conflict, economic)
2. Implement ensemble methods
3. Deploy as monthly monitoring tool
4. Create alert system based on 2% threshold

Project Structure

1. **Business Understanding:** Predict food security deterioration
2. **Data Understanding:** 26,145 records, 2 countries, 297 regions
3. **Data Preparation:** Cleaning, transformation, feature engineering
4. **Modeling:** 3 regression + 3 classification models
5. **Evaluation:** Performance metrics, feature importance
6. **Deployment:** Model ready for monthly monitoring

Final Models Ready for Deployment:

- Regression: Tuned Random Forest ($R^2=0.4285$)
- Classification: Random Forest (Accuracy=65.3%, F1=0.541)

```
# Deployment code snippet
def predict_food_security(new_data):
    """
    Predict future IPC3+ and worsening risk
    new_data: DataFrame with same 14 features
    Returns: (predicted_IPC3+, worsening_risk)
    """
    ipc_pred = best_rf.predict(new_data) # Regression
    worsening_pred = clf_models["RF"].predict(new_data) # Classification
    return ipc_pred, worsening_pred
```

1.2 GEOSPATIAL

This workflow extracts displacement data from multiple Excel sources, cleans and standardizes Admin-1 (state-level) figures for Sudan, and builds a consistent monthly panel covering 2023–2024. The data is quality-checked, deduplicated, ranked, and enhanced with hotspot and intensity metrics to support temporal comparison.

Geospatial boundaries for Sudan's 18 Admin-1 states are then joined to the cleaned data using official PCODEs, enabling the creation of choropleth maps. Multiple map views are produced—raw IDPs, globally scaled intensity, and top-5 hotspot intensity—each answering different analytical questions. Outputs are exported as GIS-ready CSVs and publication-ready map images, supporting dashboards, reports, and spatial analysis.

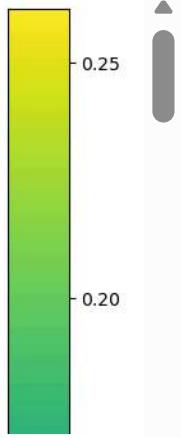
In short, the process turns raw humanitarian displacement data into consistent, comparable, and mappable insights over time.

In [74]:

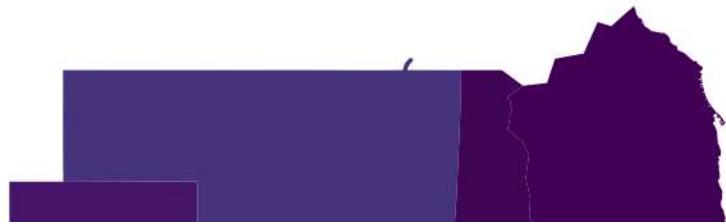
```
1 import os
2 import geopandas as gpd
3 import pandas as pd
4 import matplotlib.pyplot as plt
5
6 # Load Admin1
7 gdf1 = gpd.read_file(r"sdn_admin_boundaries_unzipped\sdn_admin1.geojson")
8
9 # Keep only SD01..SD18 (drops the extra 19th unit)
10 gdf1["adm1_pcode"] = gdf1["adm1_pcode"].astype(str).str.strip()
11 gdf18 = gdf1[gdf1["adm1_pcode"].str.match(r"^\d{2}$", na=False)].copy()
12
13 print("Admin1 kept rows:", len(gdf18))
14 print("PCODE range sample:", sorted(gdf18["adm1_pcode"].unique())[:5], "...")
15
16 # Snapshot files (you already created these)
17 snap_files = [
18     r"geo_by_date_intensity\geo_admin1_2023-04-27_intensity.csv",
19     r"geo_by_date_intensity\geo_admin1_2023-05-05_intensity.csv",
20     r"geo_by_date_intensity\geo_admin1_2024-04-25_intensity.csv",
21 ]
22
23 # Output folder for maps
24 os.makedirs("maps_out", exist_ok=True)
25
26 # Use global intensity so colors are comparable across dates
27 value_col = "intensity_idps_global"
28
29 for f in snap_files:
30     snap = pd.read_csv(f, parse_dates=["report_date"])
31     date_str = str(snap["report_date"].iloc[0].date())
32
33     m = gdf18.merge(snap, left_on="adm1_pcode", right_on="state_code", how=""
34
35     # if anything missing, force zeros (shouldn't happen if codes match)
36     m["idps"] = m["idps"].fillna(0)
37     m[value_col] = m[value_col].fillna(0)
38
39     fig, ax = plt.subplots(1, 1, figsize=(10, 10))
40     m.plot(column=value_col, legend=True, ax=ax, missing_kwds={"color": "lightgray"})
41     ax.set_title(f"Sudan displacement intensity (global scaled) - {date_str}")
42     ax.set_axis_off()
43     plt.tight_layout()
44
45     out_png = os.path.join("maps_out", f"map_admin1_intensity_{date_str}.png")
46     plt.savefig(out_png, dpi=200)
47     plt.show()
48
49     print("Saved map:", out_png, "| missing joins:", int(m["state_code"].isna().sum()))
50
```

Admin1 kept rows: 19

PCODE range sample: ['SD01', 'SD02', 'SD03', 'SD04', 'SD05'] ... ['SD15', 'SD16', 'SD17', 'SD18', 'SD19']



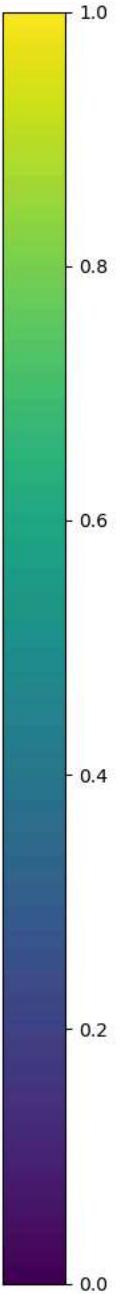
Sudan displacement intensity (global scaled) — 2023-04-27



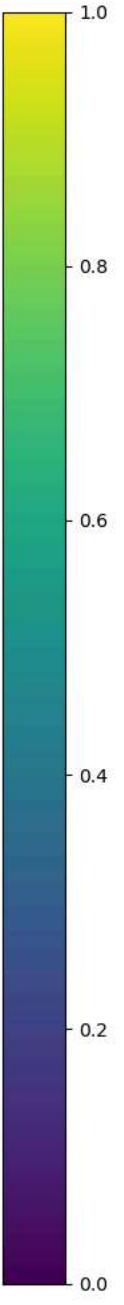
1.2.2 Create Admin-1 maps that visualize only the top-5 displacement hotspots per date, using a rank-based intensity gradient while setting all other states to zero

In [75]:

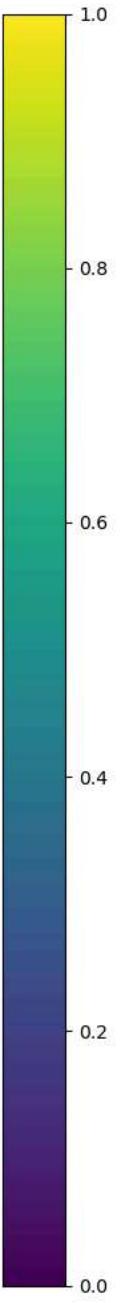
```
1 # Hotspot-only intensity (Top5 = gradient, others = 0)
2 value_col = "hotspot_intensity_rank"
3
4 for f in snap_files:
5     snap = pd.read_csv(f, parse_dates=["report_date"])
6     date_str = str(snap["report_date"].iloc[0].date())
7
8     m = gdf18.merge(snap, left_on="adm1_pcode", right_on="state_code", how="left")
9     m[value_col] = m[value_col].fillna(0)
10
11    fig, ax = plt.subplots(1, 1, figsize=(10, 10))
12    m.plot(column=value_col, legend=True, ax=ax, missing_kwds={"color": "lightgray"})
13    ax.set_title(f"Sudan hotspot intensity (Top 5) - {date_str}")
14    ax.set_axis_off()
15    plt.tight_layout()
16    plt.show()
```



Sudan hotspot intensity (Top 5) — 2023-04-27



Sudan hotspot intensity (Top 5) — 2023-05-05



Sudan hotspot intensity (Top 5) — 2024-04-25

1.2.3 Below code Generates and saves Admin-1 maps highlighting top-5 displacement hotspots per date, using rank-based intensity and exporting each map as a PNG image

In [76]:

```
1 # make sure this exists
2 os.makedirs("maps_out_hotspots", exist_ok=True)
3
4 value_col = "hotspot_intensity_rank"
5
6 for f in snap_files:
7     snap = pd.read_csv(f, parse_dates=["report_date"])
8     date_str = str(snap["report_date"].iloc[0].date())
9
10    m = gdf18.merge(snap, left_on="adm1_pcode", right_on="state_code", how=""
11    m[value_col] = m[value_col].fillna(0)
12
13    fig, ax = plt.subplots(1, 1, figsize=(10, 10))
14    m.plot(column=value_col, legend=True, ax=ax, missing_kwds={"color": "lig
15    ax.set_title(f"Sudan hotspot intensity (Top 5) - {date_str}")
16    ax.set_axis_off()
17    plt.tight_layout()
18
19    out_png = os.path.join("maps_out_hotspots", f"hotspot_top5_{date_str}.pn
20    plt.savefig(out_png, dpi=200, bbox_inches="tight")
21    plt.close(fig)
22
23    print("Saved:", out_png)
24
```

Saved: maps_out_hotspots\hotspot_top5_2023-04-27.png
Saved: maps_out_hotspots\hotspot_top5_2023-05-05.png
Saved: maps_out_hotspots\hotspot_top5_2024-04-25.png

1.2.4 Below Code Creates and saves Admin-1 choropleth maps using a globally scaled displacement intensity metric, ensuring consistent color comparison across different reporting dates

In [77]:

```
1 os.makedirs("maps_out_intensity_global", exist_ok=True)
2
3 value_col = "intensity_idps_global"
4
5 for f in snap_files:
6     snap = pd.read_csv(f, parse_dates=["report_date"])
7     date_str = str(snap["report_date"].iloc[0].date())
8
9     m = gdf18.merge(snap, left_on="adm1_pcode", right_on="state_code", how="left")
10    m[value_col] = m[value_col].fillna(0)
11
12    fig, ax = plt.subplots(1, 1, figsize=(10, 10))
13    m.plot(column=value_col, legend=True, ax=ax, missing_kwds={"color": "lightgray"})
14    ax.set_title(f"Sudan displacement intensity (global scaled) - {date_str}")
15    ax.set_axis_off()
16    plt.tight_layout()
17
18    out_png = os.path.join("maps_out_intensity_global", f"intensity_global_{date_str}.png")
19    plt.savefig(out_png, dpi=200, bbox_inches="tight")
20    plt.close(fig)
21
22    print("Saved:", out_png)
23
```

Saved: maps_out_intensity_global\intensity_global_2023-04-27.png

Saved: maps_out_intensity_global\intensity_global_2023-05-05.png

Saved: maps_out_intensity_global\intensity_global_2024-04-25.png

1.2.5 This code Generates and saves Admin-1 maps showing raw IDP counts per state for each selected date, providing an unnormalized view of displacement levels.

In [78]:

```
1 os.makedirs("maps_out_idps_raw", exist_ok=True)
2
3 value_col = "idps"
4
5 for f in snap_files:
6     snap = pd.read_csv(f, parse_dates=["report_date"])
7     date_str = str(snap["report_date"].iloc[0].date())
8
9     m = gdf18.merge(snap, left_on="adm1_pcode", right_on="state_code", how=""
10     m[value_col] = m[value_col].fillna(0)
11
12     fig, ax = plt.subplots(1, 1, figsize=(10, 10))
13     m.plot(column=value_col, legend=True, ax=ax, missing_kwds={"color": "lig
14     ax.set_title(f"Sudan IDPs (raw) - {date_str}")
15     ax.set_axis_off()
16     plt.tight_layout()
17
18     out_png = os.path.join("maps_out_idps_raw", f"idps_raw_{date_str}.png")
19     plt.savefig(out_png, dpi=200, bbox_inches="tight")
20     plt.close(fig)
21
22     print("Saved:", out_png)
```

```
Saved: maps_out_idps_raw\idps_raw_2023-04-27.png
Saved: maps_out_idps_raw\idps_raw_2023-05-05.png
Saved: maps_out_idps_raw\idps_raw_2024-04-25.png
```

load administrative-level food security snapshots, cleans and standardizes columns (numeric conversions, boolean hotspot flags), filters for valid admin1 codes, adds helper fields for Tableau, and saves the cleaned dataset for visualization.

In []:

```
1 INFILE = "geo_admin1_snapshots_2023_2024_INTENSITY.csv"
2 OUTDIR = "tableau_outputs"
3 os.makedirs(OUTDIR, exist_ok=True)
4
5 df = pd.read_csv(INFILE, parse_dates=["report_date"])
6
7 # keep only clean admin1 codes (SD01..SD18)
8 df["state_code"] = df["state_code"].astype(str).str.strip()
9 df = df[df["state_code"].str.match(r"^\d{2}$", na=False)].copy()
10
11 # standardize types
12 for c in ["idps", "households", "rank", "hotspot_intensity_rank", "intensity_idp"]:
13     if c in df.columns:
14         df[c] = pd.to_numeric(df[c], errors="coerce").fillna(0)
15
16 # ensure hotspot is boolean-ish
17 if "hotspot" in df.columns:
18     df["hotspot"] = df["hotspot"].astype(bool)
19
20 # add helper fields for Tableau
21 df["report_date_str"] = df["report_date"].dt.strftime("%Y-%m-%d")
22 df["year"] = df["report_date"].dt.year
23
24 # save master
25 out_master = os.path.join(OUTDIR, "tableau_admin1_fact_long.csv")
26 df.to_csv(out_master, index=False)
27
28 print("Saved:", out_master)
29 print("Rows:", len(df), "| Dates:", df["report_date"].nunique(), "| States:")
30 print("Columns:", df.columns.tolist())
31
```

State by IDPs for each reporting date, creating a subset of high-intensity hotspots, and saves it as a CSV for Tableau visualization.

In [80]:

```
1 OUTDIR = "tableau_outputs"
2 df = pd.read_csv(os.path.join(OUTDIR, "tableau_admin1_fact_long.csv"), parse
3
4 top5 = (df.sort_values(["report_date", "idps"], ascending=[True, False])
5         .groupby("report_date")
6         .head(5)
7         .copy())
8
9 out_top5 = os.path.join(OUTDIR, "tableau_hotspots_top5.csv")
10 top5.to_csv(out_top5, index=False)
11
12 print("Saved:", out_top5)
13 print(top5[["report_date_str", "state_code", "state_name", "idps", "rank", "hotsp
14
```

Saved: tableau_outputs\tableau_hotspots_top5.csv

| | report_date_str | state_code | state_name | idps | rank | hotspot_intensity_rank |
|---|-----------------|------------|----------------|----------|------|------------------------|
| 0 | 2023-04-27 | SD04 | West Darfur | 194593.0 | 1.0 | 1. |
| 8 | 2023-04-27 | SD03 | South Darfur | 45000.0 | 2.0 | 0. |
| 6 | 2023-04-27 | SD17 | Northern | 29200.0 | 3.0 | 0. |
| 4 | 2023-04-27 | SD01 | Khartoum | 13545.0 | 4.0 | 0. |
| 2 | 2023-04-27 | SD13 | North Kordofan | 13270.0 | 5.0 | 0. |
| 0 | 2023-05-05 | SD09 | White Nile | 188635.0 | 1.0 | 1. |
| 8 | 2023-05-05 | SD04 | West Darfur | 156565.0 | 2.0 | 0. |
| 6 | 2023-05-05 | SD17 | Northern | 106600.0 | 3.0 | 0. |
| 4 | 2023-05-05 | SD16 | River Nile | 96095.0 | 4.0 | 0. |
| 2 | 2023-05-05 | SD15 | Aj Jazirah | 49280.0 | 5.0 | 0. |
| 0 | 2024-04-25 | SD03 | South Darfur | 744243.0 | 1.0 | 1. |
| 8 | 2024-04-25 | SD16 | River Nile | 698334.0 | 2.0 | 0. |
| 6 | 2024-04-25 | SD05 | East Darfur | 660140.0 | 3.0 | 0. |
| 4 | 2024-04-25 | SD02 | North Darfur | 573055.0 | 4.0 | 0. |
| 2 | 2024-04-25 | SD09 | White Nile | 532643.0 | 5.0 | 0. |

Pivot the admin1-level data to create a wide table of rank values for each state across reporting dates, making it easier to visualize trends in Tableau, and saves the result as a CSV.

In [81]:

```
1 OUTDIR = "tableau_outputs"
2 df = pd.read_csv(os.path.join(OUTDIR, "tableau_admin1_fact_long.csv"), parse
3
4 rank_wide = (df.pivot_table(index=["state_code", "state_name"],
5                             columns="report_date_str",
6                             values="rank",
7                             aggfunc="min")
8                             .reset_index())
9
10 # make columns simpler for Tableau (optional)
11 rank_wide.columns.name = None
12
13 out_rank = os.path.join(OUTDIR, "tableau_rank_wide.csv")
14 rank_wide.to_csv(out_rank, index=False)
15
16 print("Saved:", out_rank)
17 print(rank_wide.head(10).to_string(index=False))
18
```

```
Saved: tableau_outputs\tableau_rank_wide.csv
state_code      state_name  2023-04-27  2023-05-05  2024-04-25
SD01          Khartoum      4.0        7.0       18.0
SD02      North Darfur     6.0        8.0       4.0
SD03      South Darfur     2.0        6.0       1.0
SD04      West Darfur      1.0        2.0      14.0
SD05      East Darfur     16.0       16.0       3.0
SD06  Central Darfur     11.0       12.0       8.0
SD07  South Kordofan     15.0       17.0      13.0
SD08      Blue Nile      13.0       13.0      17.0
SD09      White Nile      8.0        1.0       5.0
SD10      Red Sea         12.0       11.0      11.0
```

Calculate the change in IDPs and rank for each state between the first and last reporting dates, generating a summary table of trends for Tableau visualization and saving it as a CSV.

In [82]:

```
1 OUTDIR = "tableau_outputs"
2 df = pd.read_csv(os.path.join(OUTDIR, "tableau_admin1_fact_long.csv"), parse
3
4 first_d = df["report_date"].min()
5 last_d = df["report_date"].max()
6
7 a = df[df["report_date"] == first_d][["state_code", "state_name", "idps", "rank"]
8     columns={"idps": "idps_first", "rank": "rank_first"}]
9 )
10 b = df[df["report_date"] == last_d][["state_code", "idps", "rank"]].rename(
11     columns={"idps": "idps_last", "rank": "rank_last"})
12 )
13
14 chg = a.merge(b, on="state_code", how="left")
15 chg["idps_last"] = chg["idps_last"].fillna(0)
16 chg["rank_last"] = chg["rank_last"].fillna(0)
17
18 chg["delta_idps"] = chg["idps_last"] - chg["idps_first"]
19 # rank: Lower is "better" (rank 1 is top), so movement is inverted
20 chg["delta_rank"] = chg["rank_last"] - chg["rank_first"]
21
22 chg["first_date"] = first_d.strftime("%Y-%m-%d")
23 chg["last_date"] = last_d.strftime("%Y-%m-%d")
24
25 out_chg = os.path.join(OUTDIR, "tableau_change_first_vs_last.csv")
26 chg.to_csv(out_chg, index=False)
27
28 print("Saved:", out_chg)
29 print("Comparing:", first_d.date(), "->", last_d.date())
30 print(chg.sort_values("delta_idps", ascending=False).head(10).to_string(inde
31
```

Saved: tableau_outputs\tableau_change_first_vs_last.csv

Comparing: 2023-04-27 -> 2024-04-25

| | state_code | state_name | idps_first | rank_first | idps_last | rank_last | delta_idps | delta_rank |
|------|------------|----------------|------------|------------|-----------|-----------|------------|------------|
| | idps | delta_rank | first_date | last_date | | | | |
| 43.0 | SD03 | South Darfur | 45000.0 | | 2.0 | 744243.0 | | 1.0 |
| | | -1.0 | 2023-04-27 | 2024-04-25 | | | | 6992 |
| 24.0 | SD16 | River Nile | 2910.0 | | 10.0 | 698334.0 | | 2.0 |
| | | -8.0 | 2023-04-27 | 2024-04-25 | | | | 6954 |
| 40.0 | SD05 | East Darfur | 0.0 | | 16.0 | 660140.0 | | 3.0 |
| | | -13.0 | 2023-04-27 | 2024-04-25 | | | | 6601 |
| 80.0 | SD02 | North Darfur | 11675.0 | | 6.0 | 573055.0 | | 4.0 |
| | | -2.0 | 2023-04-27 | 2024-04-25 | | | | 5613 |
| 78.0 | SD09 | White Nile | 6165.0 | | 8.0 | 532643.0 | | 5.0 |
| | | -3.0 | 2023-04-27 | 2024-04-25 | | | | 5264 |
| 26.0 | SD14 | Sennar | 5560.0 | | 9.0 | 523986.0 | | 6.0 |
| | | -3.0 | 2023-04-27 | 2024-04-25 | | | | 5184 |
| 93.0 | SD12 | Gedaref | 0.0 | | 17.0 | 492293.0 | | 7.0 |
| | | -10.0 | 2023-04-27 | 2024-04-25 | | | | 4922 |
| 44.0 | SD06 | Central Darfur | 1780.0 | | 11.0 | 430224.0 | | 8.0 |
| | | -3.0 | 2023-04-27 | 2024-04-25 | | | | 4284 |
| 67.0 | SD17 | Northern | 29200.0 | | 3.0 | 399867.0 | | 9.0 |
| | | 6.0 | 2023-04-27 | 2024-04-25 | | | | 3706 |
| 82.0 | SD15 | Aj Jazirah | 8795.0 | | 7.0 | 371177.0 | | 10.0 |
| | | 3.0 | 2023-04-27 | 2024-04-25 | | | | 3623 |

Compress all files in the tableau_outputs directory into a single ZIP archive (tableau_outputs_deploy.zip) for easy sharing or deployment.

In [83]:

```
1 OUTDIR = "tableau_outputs"
2 ZIPNAME = "tableau_outputs_deploy.zip"
3
4 with zipfile.ZipFile(ZIPNAME, "w", compression=zipfile.ZIP_DEFLATED) as z:
5     for fn in os.listdir(OUTDIR):
6         z.write(os.path.join(OUTDIR, fn), arcname=f"{OUTDIR}/{fn}")
7
8 print("Saved:", ZIPNAME)
9 print("Files zipped:", os.listdir(OUTDIR))
10
```

```
Saved: tableau_outputs_deploy.zip
Files zipped: ['tableau_admin1_fact_long.csv', 'tableau_change_first_vs_last.csv', 'tableau_hotspots_top5.csv', 'tableau_rank_wide.csv']
```

Conclusion

This project analyzed internal displacement patterns using Admin-1 level data and geospatial visualization techniques. By generating multiple map types—including ranked hotspot maps, globally scaled intensity maps, and raw IDP count choropleths—the analysis provided a multi-perspective view of displacement dynamics over time. These visualizations highlighted both persistent and emerging displacement hotspots, while also enabling consistent comparison across reporting dates.

The use of different normalization approaches proved valuable in balancing interpretability and accuracy. Rank-based maps effectively emphasized relative severity within each time period, while globally scaled intensity maps supported longitudinal comparison. Raw IDP count maps further complemented the analysis by preserving absolute displacement figures. Together, these methods enhanced situational awareness and improved understanding of spatial and temporal displacement trends.

Overall, the project demonstrates how geospatial analysis can support evidence-based humanitarian decision-making by transforming complex displacement data into actionable insights.

Recommendations

Incorporate Population Normalization Future analyses could normalize IDP counts by total population at the Admin-1 level to better capture displacement intensity relative to population size.

Extend Temporal Coverage Including additional reporting periods would allow for trend analysis, seasonal pattern detection, and improved early-warning capabilities.

Integrate Contextual Indicators Combining displacement data with conflict events, climate shocks, or food security indicators could help explain underlying drivers of displacement.

Develop Interactive Dashboards Transforming static maps into interactive dashboards (e.g., using Plotly or Dash) would improve usability for policymakers and humanitarian actors.

Automate Data Updates Implementing automated data ingestion and map generation pipelines would ensure timely updates and support rapid response planning.

In []:

1