# A Global Analysis and Prediction of Food Insecurity and Hunger Crisis Outbreaks using Machine Learning

Bhavan K. Basavaraju, Namratha S. Kumar, Shreya Chikatmarla, Sowjanya Pamulapati

# Project Background

- In today's world, hunger and food scarcity are on the rise as more people in the world face the challenge of securing a meal on a regular basis, leading to malnutrition and food deprivation.

- According to Tarasuk (2001), food insecurity is a fundamental human need that is unaffected by a person's gender, colour, or ethnicity.

- Hernandez et al. (2017) report that 27% of adults, aged 36 on average, suffer food insecurity. This emphasizes the need for immediate monitoring and eradication.

- Previous studies have restricted its scope and narrowed the focus to particular regions.

- This area of study intends to improve global food intake estimates by addressing shortcomings, including a broader scope that includes post-pandemic scenarios, building upon the work of Martini et al. (2022).

# Executive Summary

- Understanding and addressing the global food issues, with a focus on malnutrition and food insecurity.

- Diverse data for effective algorithm designing and model development.

- Time-series biassed algorithms (LSTM, ARIMA, Random Forest, Prophet) are used for precise predictions.

- To evaluate algorithm performance and tweak as needed, metrics like as MAE, MSE, R2 Score, and Explained Variance Score are used.

- Integration of initial analysis, data visualisation, and a predictive model to provide an overview of the global hunger situation and its progression over time.

# Motivation

- The project is motivated by the global prevalence of those battling for a daily meal.

- Tarasuk (2001) and Hernandez et al. (2017) underline the severity of food insecurity. Emphasizing the urgency and need for study and deliberate elimination of a food crisis.

- The shortcomings of current prediction models, particularly missing out on post-pandemic scenarios such as COVID-19, further motivate our project.

- To address the global nature of food insecurity, extending outside existing geographical boundaries.

# Problem Statement

- Addressing global food hunger is a significant concern with consequences for both individuals and the nation.

- Current studies utilising machine learning approaches exhibit restrictions in terms of using diverse algorithms, efficiency in data procurement, and region-specific analysis.

- Our proposed approach intends to overcome these gaps by incorporating unique algorithms and using diversified data sources, with the objective of enhancing the precision of forecasts while offering actionable insights for successful mitigation strategies.

# Project Requirements

- Addressing the hunger crisis and food insecurity.

- Identifying factors and inclusion of post-pandemic data

- Diverse dataset collection: census (2010-2021), RASFF (2000-2023), OECD survey(1970-2022), IMF trade data (1980-present)

- Time-series biased algorithm selection: LSTM, ARIMA, Prophet, and Random Forest.

- Algorithmic workflow: preprocessing, splitting, configuration, training, and evaluation.

- Evaluation metrics: MSE, MAE, R2 Score, Explained Variance Score.

- Comparative study to determine the best algorithm performance.

- Organizing project workflow and task division for optimal collaboration.

# Technology & Literature Survey

| PAPER TITLE | AUTHORS | DATASET | MODELS | RESULTS |
|---|---|---|---|---|
| Food security and agricultural challenges in West-African rural communities: a machine learning analysis. | Ahn et al. (2022) | Food Security and Production Assessment of Ghana, Senegal, and Liberia | Random Forest Algorithm, Decision tree and Chi-Square Automatic Interaction Detection (CHAID) | The decision tree achieved an accuracy of 73.7% Classification for Senegal using decision trees was more precise than for Liberia and the misclassification rate was 20%. Random forest employed for Ghana had a 9% misclassification rate and 0.01 standard error indicating efficient model performance and precision. |
| Forecasting transitions in the state of food security with machine learning using transferable features | Westerveld et al. (2021) | food security transitions on a monthly basis in Ethiopia | "Extreme Gradient Boosting", "Random Forest", and "CatBoost" | Welch based T-test between the three algorithms is performed using normal distribution. It was found that with a t (198) =19.86 and t (198) =86.33 with p < 0.001 across comparison of F1 scores between XGboost and Catboost, XGBoost and Random Forest respectively. |
| Food security prediction from heterogeneous data combining machine and deep learning methods | Deléglise et al. (2020) | The datasets used in the study included time series data, meteorological data, population density data, World Bank economic data, and the Normalized Difference Vegetation | Random Forest (RF), Convolutional Neural Networks (CNNs), and Long-term and Short-term Memory (LSTM) models | The framework outperforms all competing methods, with model (b) surpassing model (a). The R̂2 values obtained for (0.469) and (0.434) are statistically significant, demonstrating the benefits of integrating various data science techniques. The WFP framework's results in Burkina Faso are relatively modest (0.34 for and 0.30 for), while Lentz et al.'s study shows even lower R̂2 values below 0.2 |
| Modeling and Forecasting of Food Security for Wheat in Egypt Through Year 2025 by Using time Series ARIMA Models | Negm et al. (2018) | Data on wheat production, consumption, and imports in Egypt obtained from sources | ARIMA | wheat consumption in Egypt is projected to increase by 18.54 million tons based on the ARIMA (0, 0, 1) model |

# Project Resource Requirements

| Hardware | Memory | Configuration | Purpose |
|---|---|---|---|
| Apple M2, macOS Sonoma Version 14.0 | 16 GB | 494.38 GB | ML Model |
| Apple M1 Pro, macOS Ventura Version 13.5.1 | 16 GB | 494.38 GB | ML Model |
| Apple M2, macOS Ventura Version 13.0 | 16 GB | 994.66 GB | Computation engine for running ML Models in Jupyter Notebook |

| Libraries/Packages | Purpose | Version |
|---|---|---|
| Python | Data Cleaning Data Preparation Data Analysis | 3.8 |
| Pandas | Data Manipulation and analysis Data Structures | 1.3.2 |
| Numpy | Fundamental Computing for Numerical data | 3.2 |
| Matplotlib | Data Visualization using plots and charts | 3.2 |
| Seaborn | Informative Statistical Graphics for Data Exploration | 0.12.0 |
| sci-kit-learn | Machine learning Model building | 1.1.2 |

**Hardware Requirements**

**Software Requirements**

**Tools and Licenses**

| Tools | Purpose | License |
|---|---|---|
| Jupyter NoteBook | Code Design and Development | Open Source |
| IntelliJ | Software Development | Open Source |
| GitHub | Create and storing our Project | Open Source |
| Draw.io | Aiding Design | Free |
| Zoom | Team Collaboration meeting | Free |
| Google Meet | Team Collaboration meetings | Free |
| Google Docs | Documentation | Free |
| Google Drive | Data Storage | Free |

# Project Management Tools & WBS

## WBS

WBS is a Project management tool for hierarchical structure based on scope assessment

It breaks down complex deliverables into simpler, manageable tasks and defines project objectives effectively.

Establishes milestones and divides deliverables into work packages and identifies dependencies for successful project delivery.

Integration of CRISP-DM methodology enhances work breakdown structure.

## Gantt Chart

They are essential project management tools that visually represent tasks, schedules, and dependencies of a project.

They provide a clear overview of the project schedule, help identify milestones and due dates, facilitate communication and collaboration among team members , and allow for tracking task progress

Creating a Gantt chart involves listing tasks, determining task durations and dependencies, and using software tools such as WBS - Gantt chart plugin in JIRA.

In our project, we follow the structure of Epics, User stories, Tasks, and, Sub Tasks which are all represented on the Gantt chart to track progress and ensure timely completion.

## PERT CHART

PERT is a project management tool for handling complex projects effectively and provides detailed explanations and addresses project complexities.

They help to determine task order, estimate task durations, identify critical paths, and manage project risks.

Critical path analysis is crucial in identifying tasks that can impact the overall project timeline.

They have limitations and can be challenging for large projects, but they remain a valuable tool for handling complex and risky projects.

# Work Breakdown Structure

# Gantt Chart

# PERT CHART



| Activity | Description | Predecessor/Dependency |
|---|---|---|
| 1 | Start Project | _ |
| 2 | Key factors identification | _ |
| 3 | Related works detailing | _ |
| 4 | Identify existing algorithms | _ |
| 5 | Finalized Idea | 1,3 |
| 6 | Data discovery | 5 |
| 7 | Exploratory data analysis | 6 |
| 8 | Accessing Data Quality | 7 |
| 9 | Cleaning on Data | 8 |
| 10 | Datasets consolidation | 9 |
| 11 | Selection of target variables/features | 10 |
| 12 | Model Design | 11 |
| 13 | Model Enhancement Techniques | 12 |
| 14 | Review the Model Design | 12 |
| 15 | Evaluation wrt algorithm based metrics | 12,14 |
| 16 | Determination of output | 15 |
| 17 | Model deployment | 16 |
| 18 | Final report submission | 17 |

# Project Workflow

# Data Collection

- Census information on food security(2010 to 2021)

  Datatypes - int64, float64

  Number of columns - 507

  Data Size - 485.4 MB


- RASSF(Ras association domain family) dataset EC (Europe Commision) data from 2017to present)

  Datatypes - object, float64

  Number of columns - 14

  Data Size - 546.6+ KB

- Gross domestic product (GDP) survey data - (OECD sourced survey data from 1970 to 2022)

  Datatypes - int64, float64, object

  Number of columns - 23

  Data Size - 12.0+ MB


- Commodity Terms of Trade dataset (IMF sourced circadian data from 1980 to present)

  Datatypes- int64, float64, object

  Number of columns - 592

  Data Size - 19.5+ MB

# Raw Dataset Samples

Census information on *Food security dataset*



RASSF dataset



Gross domestic product (GDP) survey data



Commodity Terms of Trade dataset

# Data Preprocessing: Exploratory Data Analysis



Sample HRPOOR feature analysis



Distribution of time series data



Sample HRFS12M1 feature analysis

# Data Preprocessing: Exploratory Data Analysis



HRPOOR features multivariate analysis



Correlation matrix

# Data Cleaning: Handling of Inconsistent Data

| | reference | category | type | subject | date |
|---|---|---|---|---|---|
| 0 | 2023.7463 | fruits and vegetables | food | Buprofezin in lemons from Türkiye | 2023-01-11 09:48:00 |
| 1 | 2023.7457 | dietetic foods, food supplements and fortified... | food | Novel Food in Food Supplement | 2023-10-31 18:27:17 |
| 2 | 2023.7456 | fruits and vegetables | food | Perchlorate in radish | 2023-10-31 17:55:46 |
| 3 | 2023.7455 | fruits and vegetables | food | Chlormequat in Nashi Pears | 2023-10-31 17:46:25 |
| 4 | 2023.7454 | soups, broths, sauces and condiments | food | Consignment Mayonaises not presented for veter... | 2023-10-31 17:32:54 |
| ... | ... | ... | ... | ... | ... |
| 4995 | 2022.5568 | crustaceans and products thereof | food | Ruptura de la cadena de frío en Litopenaeus va... | 2022-09-26 15:11:33 |
| 4996 | 2022.5567 | confectionery | food | Too high content of fumaric acid in gummy bear... | 2022-09-26 15:07:40 |
| 4997 | 2022.5566 | nuts, nut products and seeds | food | SALMONELLA IN SESAME SEEDS FROM TURKEY | 2022-09-26 14:25:06 |
| 4998 | 2022.5565 | nuts, nut products and seeds | food | SALMONELLA IN SESAME SEEDS FROM TURKEY | 2022-09-26 14:18:42 |
| 4999 | 2022.5564 | nuts, nut products and seeds | food | SALMONELLA IN SESAME SEEDS FROM TURKEY | 2022-09-26 14:12:57 |

Time in dd-m-y : hr format

```python
def parse_date(date_str):
    try:
        # Try one format
        return pd.to_datetime(date_str, format="%m/%d/%Y %H:%M")
    except ValueError:
        try:
            # Try another format
            return pd.to_datetime(date_str, format="%d-%m-%Y %H:%M:%S")
        except ValueError:
            # Handle any other format or errors
            return None
df2['date'] = df2['date'].apply(parse_date)
df2
```

Sample code snippet

| Year |
|---|
| 2023.0 |
| 2023.0 |
| 2023.0 |
| 2023.0 |
| 2023.0 |
| ... |
| 2022.0 |
| 2022.0 |
| 2022.0 |
| 2022.0 |
| 2022.0 |

Year extracted from time format

# Data Cleaning: Handling of Incomplete & Missing Data

```
clean_df=new_df.dropna()
```

```
clean_df
```

Null values count:

| | |
|---|---|
| Unnamed: 0 | 0 |
| LOCATION | 275831 |
| Country | 275831 |
| TRANSACT | 275831 |
| Transaction | 275831 |
| MEASURE | 275831 |
| Measure | 275831 |
| TIME | 275831 |
| Year | 275831 |
| Unit Code | 275831 |
| Unit | 275831 |
| PowerCode Code | 275831 |
| PowerCode | 275831 |
| Reference Period Code | 397267 |
| Reference Period | 397267 |
| Value | 275831 |
| Flag Codes | 497908 |
| Flags | 497908 |
| reference | 518979 |
| category | 518979 |
| type | 518979 |
| subject | 518979 |
| date | 518979 |
| notifying_country | 518979 |
| classification | 518979 |
| risk_decision | 518979 |
| distribution | 520547 |
| forAttention | 521408 |
| forFollowUp | 521667 |
| operator | 518986 |
| origin | 519105 |
| hazards | 520356 |
| Year.1 | 0 |
| HEFAMINC | 0 |
| HRPOOR | 0 |
| HRFS12CX | 0 |
| HRFS12M1 | 0 |
| HRFS30D1 | 0 |
| HESS1 | 0 |
| PREMPHRS | 0 |
| PREXPLF | 0 |
| PREMPNOT | 0 |
| HETSP3O | 0 |
| HESP6 | 0 |
| dtype: int64 | |
| 12650054 | |

**Null values count**

| | Unnamed: 0 | LOCATION | Country | TRANSACT | Transaction | MEASURE | Measure | TIME | Year | Unit Code | ... | HRPOOR | HRFS12CX | HRFS12M1 | HRFS30D1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | AUS | Australia | B1_GA | Gross domestic product (output approach) | C | Current prices | 2016.0 | 2016.0 | AUD | ... | -1.0 | -1.0 | -1.0 | -1.0 |
| **1** | 1 | AUS | Australia | B1_GA | Gross domestic product (output approach) | C | Current prices | 2017.0 | 2017.0 | AUD | ... | -1.0 | -1.0 | -1.0 | -1.0 |
| **2** | 2 | AUS | Australia | B1_GA | Gross domestic product (output approach) | C | Current prices | 2018.0 | 2018.0 | AUD | ... | -1.0 | -1.0 | -1.0 | -1.0 |
| **3** | 3 | AUS | Australia | B1_GA | Gross domestic product (output approach) | C | Current prices | 2019.0 | 2019.0 | AUD | ... | -1.0 | -1.0 | -1.0 | -1.0 |
| **4** | 4 | AUS | Australia | B1_GA | Gross domestic product (output approach) | C | Current prices | 2020.0 | 2020.0 | AUD | ... | -1.0 | -1.0 | -1.0 | -1.0 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **248142** | 248142 | TUR | Türkiye | P32S13 | Collective consumption expenditure of general ... | DOB | Deflator | 2018.0 | 2018.0 | IDX | ... | -1.0 | -1.0 | -1.0 | -1.0 |
| **248143** | 248143 | TUR | Türkiye | P32S13 | Collective consumption expenditure of general ... | DOB | Deflator | 2019.0 | 2019.0 | IDX | ... | -1.0 | -1.0 | -1.0 | -1.0 |
| **248144** | 248144 | TUR | Türkiye | P32S13 | Collective consumption expenditure of general ... | DOB | Deflator | 2020.0 | 2020.0 | IDX | ... | -1.0 | -1.0 | -1.0 | -1.0 |
| **248145** | 248145 | TUR | Türkiye | P32S13 | Collective consumption expenditure of general ... | DOB | Deflator | 2021.0 | 2021.0 | IDX | ... | -1.0 | -1.0 | -1.0 | -1.0 |
| **248146** | 248146 | TUR | Türkiye | P32S13 | Collective consumption expenditure of general ... | DOB | Deflator | 2022.0 | 2022.0 | IDX | ... | -1.0 | -1.0 | -1.0 | -1.0 |

248147 rows × 26 columns

**Sample data after clearing Null values**

```
<class 'pandas.core.frame.DataFrame'>
Index: 248147 entries, 0 to 248146
Data columns (total 26 columns):
 #   Column          Non-Null Count   Dtype
---  ------          --------------   -----
 0   Unnamed: 0      248147 non-null  int64
 1   LOCATION        248147 non-null  object
 2   Country         248147 non-null  object
 3   TRANSACT        248147 non-null  object
 4   Transaction     248147 non-null  object
 5   MEASURE         248147 non-null  object
 6   Measure         248147 non-null  object
 7   TIME            248147 non-null  float64
 8   Year            248147 non-null  float64
 9   Unit Code       248147 non-null  object
 10  Unit            248147 non-null  object
 11  PowerCode Code  248147 non-null  float64
 12  PowerCode       248147 non-null  object
 13  Value           248147 non-null  float64
 14  Year.1          248147 non-null  float64
 15  HEFAMINC        248147 non-null  float64
 16  HRPOOR          248147 non-null  float64
 17  HRFS12CX        248147 non-null  float64
 18  HRFS12M1        248147 non-null  float64
 19  HRFS30D1        248147 non-null  float64
 20  HESS1           248147 non-null  float64
 21  PREMPHRS        248147 non-null  float64
 22  PREXPLF         248147 non-null  float64
 23  PREMPNOT        248147 non-null  float64
 24  HETSP3O         248147 non-null  float64
 25  HESP6           248147 non-null  float64
dtypes: float64(16), int64(1), object(9)
memory usage: 51.1+ MB
```

**Dataset after clearing null values**

# Data Cleaning: Handling of Noisy Data

[7530 rows x 17 columns]

Row count for outliers

```python
z_scores = (numerical_df - numerical_df.mean()) / numerical_df.std()
print(z_scores)
outliers_df = numerical_df[~(z_scores < 3).all(axis=1)]

# Rows containing outliers
print("Rows with outliers:")
print(outliers_df)

# Rows after eliminating outliers
res_df = numerical_df[(z_scores < 3).all(axis=1)]
print("\nRows after eliminating outliers:")
print(res_df)
```

Sample code snippet

[240617 rows x 17 columns]

Total row count after removing outliers

To summarize:

- Exploratory data analysis is performed to understand how the data stands after merging the datasets which helps in understanding correlation between columns, dispersion of datatypes, etc.

- Data cleaning is performed in three phases,

- As we have many yearly insights to be drawn, the data is manipulated to a yearly format

- Next, we clean the data of all the null values

- Finally, noisy data is cleared using the z-scores from all the numerical columns

# Data Transformation: Label Encoding

```python
unique_values = sorted(final_clean_df['Country'].unique())

# Print the unique values
print(unique_values)
```

```
['Argentina', 'Australia', 'Austria', 'Belgium', 'Brazil', 'Bulgaria', 'Cameroon', 'Canada', 'Chile', "China (Peopl
e's Republic of)", 'Colombia', 'Costa Rica', 'Croatia', 'Cyprus', 'Czech Republic', 'Denmark', 'Estonia', 'Euro area
(19 countries)', 'European Union — 27 countries (from 01/02/2020)', 'Finland', 'France', 'Germany', 'Greece', 'Hungar
y', 'Iceland', 'India', 'Indonesia', 'Ireland', 'Israel', 'Italy', 'Japan', 'Korea', 'Latvia', 'Lithuania', 'Luxembou
rg', 'Malta', 'Mexico', 'Netherlands', 'New Zealand', 'Norway', 'OECD - Total', 'Poland', 'Portugal', 'Romania', 'Rus
sia', 'Saudi Arabia', 'Senegal', 'Slovak Republic', 'Slovenia', 'South Africa', 'Spain', 'Sweden', 'Switzerland', 'Tü
rkiye', 'United Kingdom', 'United States']
```

Sample code snippet to identify unique values in a column to be encoded

```python
from sklearn.preprocessing import LabelEncoder
country_encoder = LabelEncoder()
# Fit and transform the column
final_clean_df['Country_Column_encoded'] = country_encoder.fit_transform(final_clean_df['Country'])

# Print or use the encoded column
print(final_clean_df['Country_Column_encoded'])
```

```
0          1
1          1
2          1
3          1
4          1
          ..
240612    53
240613    53
240614    53
240615    53
240616    53
Name: Country_Column_encoded, Length: 240617, dtype: int32
```

Sample code snippet of label encoding a column

# Data Normalization

```python
import pandas as pd
from sklearn.preprocessing import MinMaxScaler, StandardScaler
num_df = final_clean_df.select_dtypes(include=['number'])

# Z-score Normalization
z_score_scaler = StandardScaler()
df_z_score_normalized = pd.DataFrame(z_score_scaler.fit_transform(num_df), columns=num_df.columns)
```

Sample code snippet to perform Z-score normalization

| | TIME | Year | PowerCode Code | Value | Year.1 | HEFAMINC | HRPOOR | HRFS12CX | HRFS12M1 | HRFS30D1 | ... | PREXPLF | PREMPNOT | HETSP3O | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | -1.472603 | 2016.0 | 0.487835 | -0.007763 | 0.988656 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |
| 1 | -0.962859 | 2017.0 | 0.487835 | -0.007583 | 0.988656 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |
| 2 | -0.453115 | 2018.0 | 0.487835 | -0.007359 | 0.988656 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |
| 3 | 0.056629 | 2019.0 | 0.487835 | -0.007289 | 0.988656 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |
| 4 | 0.566373 | 2020.0 | 0.487835 | -0.007072 | 0.988656 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 240612 | -0.453115 | 2018.0 | -2.251247 | -0.011545 | -1.011474 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |
| 240613 | 0.056629 | 2019.0 | -2.251247 | -0.011545 | -1.011474 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |
| 240614 | 0.566373 | 2020.0 | -2.251247 | -0.011545 | -1.011474 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |
| 240615 | 1.076117 | 2021.0 | -2.251247 | -0.011545 | -1.011474 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |
| 240616 | 1.585861 | 2022.0 | -2.251247 | -0.011545 | -1.011474 | -1.695627 | -1.134884 | -1.12064 | -1.071966 | -1.087396 | ... | -0.802531 | -1.108711 | -0.149102 | -0.: |

240617 rows × 22 columns

Sample look at dataset after normalization

# Data Regularization: L1 regularization

```
[14]: import numpy as np
      import pandas as pd
      from sklearn.model_selection import train_test_split
      from sklearn.linear_model import LogisticRegression
      from sklearn.preprocessing import StandardScaler

      numeric_columns = final_clean_df.select_dtypes(include=[np.number]).columns

      df_numeric = final_clean_df[numeric_columns]
      X = df_numeric.drop('HRFS30D1', axis=1)
      y = df_numeric['HRFS30D1']

      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

      # Standardize the data
      scaler = StandardScaler()
      X_train_scaled = scaler.fit_transform(X_train)
      X_test_scaled = scaler.transform(X_test)

      # Apply L2 regularization
      model = LogisticRegression(penalty='l2', solver='liblinear')
      model.fit(X_train_scaled, y_train)

      print(X_train_scaled)

      # Evaluate the model
      accuracy = model.score(X_test_scaled, y_test)
      print(f'Accuracy: {accuracy}')
```

Sample code snippet to perform L1 regularization

```
[[ 1.58656589  1.58656589  0.48719587 ... -0.52537944 -1.86778544
   0.12749048]
 [-0.96196616 -0.96196616  0.48719587 ... -0.52537944  1.24083575
   0.12749048]
 [ 0.56715307  0.56715307  0.48719587 ... -0.52537944 -1.97497928
   0.12749048]
 ...
 [ 0.56715307  0.56715307  0.48719587 ... -0.52537944 -1.33181627
   0.12749048]
 [ 0.05744666  0.05744666  0.48719587 ... -0.17354301  0.81206041
   0.12749048]
 [-0.45225975 -0.45225975  0.48719587 ... -1.40497053 -1.33181627
   0.12749048]]
```

Sample data after L1 regularization

Accuracy: 0.9742747901255091

Accuracy for L1 regularization

# Data Regularization: L2 regularization

```python
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler

numeric_columns = final_clean_df.select_dtypes(include=[np.number]).columns

df_numeric = final_clean_df[numeric_columns]
X = df_numeric.drop('HRFS30D1', axis=1)
y = df_numeric['HRFS30D1']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Standardize the data
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# Apply L2 regularization
model = LogisticRegression(penalty='l2', solver='liblinear')
model.fit(X_train_scaled, y_train)

print(X_train_scaled)

# Evaluate the model
accuracy = model.score(X_test_scaled, y_test)
print(f'Accuracy: {accuracy}')
```

Sample code snippet to perform L2 regularization

```
[[ 1.58656589  1.58656589  0.48719587 ... -0.52537944 -1.86778544
   0.12749048]
 [-0.96196616 -0.96196616  0.48719587 ... -0.52537944  1.24083575
   0.12749048]
 [ 0.56715307  0.56715307  0.48719587 ... -0.52537944 -1.97497928
   0.12749048]
 ...
 [ 0.56715307  0.56715307  0.48719587 ... -0.52537944 -1.33181627
   0.12749048]
 [ 0.05744666  0.05744666  0.48719587 ... -0.17354301  0.81206041
   0.12749048]
 [-0.45225975 -0.45225975  0.48719587 ... -1.40497053 -1.33181627
   0.12749048]]
```

Sample data after L2 regularization

Accuracy: 0.9741916715152522

Accuracy for L2 regularization

# Principal Component Analysis



Sample PC distribution



2D scatter plot



Explained variance with respect to PC



3D scatter plot

# Singular Value Decomposition

```
        Component 1  Component 2  Component 3
0          3.212763    -1.986734    -0.747473
1          3.204051    -1.267029    -0.714400
2          3.195338    -0.547324    -0.681326
3          3.186625     0.172381    -0.648253
4          3.177912     0.892086    -0.615180
...             ...          ...          ...
240612     3.213045    -0.615550    -0.568933
240613     3.204332     0.104156    -0.535859
240614     3.195619     0.823862    -0.502786
240615     3.186906     1.543567    -0.469712
240616     3.178193     2.263273    -0.436639

[240617 rows x 3 columns]
```



Sample SVD distribution

Explained variance with respect to PC

To summarize:

- PCA is usually used to work on linear relationships, thus there seems to be lower variance extraction from the 2 principal components

- SVD has shown efficiency in dealing with non-linear data, it better explains about data for non-linear relationships, works for broader range of applications, and hold numerical stability over PCA

- SVD aids working with time series data through extraction of relevant features, isolation of anomalous patterns or trends, etc. making it an efficient tool to work with time series-based algorithms and data.

# Train, Test, and Validation Datasets

```python
from sklearn.model_selection import train_test_split

features = transformed_df.drop('HRFS30D1', axis=1)
target = transformed_df['HRFS30D1']

# Split the data into training (80%) and temporary set (20%)
temp_features, test_features, temp_target, test_target = train_test_split(features, target, test_size=0.2, random_state=42)

# Further split the temporary set into training (60%) and validation (20%)
train_features, val_features, train_target, val_target = train_test_split(temp_features, temp_target, test_size=0.25, random_sta

print(f"Training set size: {len(train_features)}")
print(f"Validation set size: {len(val_features)}")
print(f"Test set size: {len(test_features)}")

print("\nSample from Training Set:")
print(train_features.head())

print("\nSample from Validation Set:")
print(val_features.head())

print("\nSample from Test Set:")
print(test_features.head())
```

Code snippet for splitting the dataset

```
Training set size: 144369
Validation set size: 48124
Test set size: 48124
```

Dataset size after splitting

# Model Development

## Long Short-Term Memory (LSTM)

- It is a recurrent neural network (RNN) architecture that is well-suited for processing and making predictions based on sequential data, such as time series data and natural language.

- It also has popular choices for various applications, including time series forecasting, language modeling, speech recognition, and more.

## ARIMA

- It is a well-known time series model for forecasting. It's divided into three parts: autoregression (AR), differencing (I), and moving average (MA).

- It is responsible for capturing a variety of standard temporal structures in time series data. It is a well-known time series model for forecasting. It's divided into three parts: autoregression (AR), differencing (I), and moving average (MA).

- It is responsible for capturing a variety of standard temporal structures in time series data.

## Prophet

- It is a process for forecasting time series data that is based on an additive model where non-linear trends are fit with yearly, weekly, and daily.

- It is designed to work best and strong with time series and historical data.

- It is available in both R and Python and used in many applications for providing forecasts.

## Random Forest

- It is an ensemble machine learning model because during training, it gives both features and targets to predict the data.

- It is commonly used for both regression and classification and is also trained on random subset of the data points and features at each node.

# Model Evaluation Metrics



Metrics
- MSE
- MAE
- R2 Score
- Explained Variance Score
- MedAE

**LSTM**
- MSE : 1.27
- MAE : 1.05
- R2 Score : -0.05
- Variance Score : -0.01
- MedAE : 0.96

**ARIMA**
- MSE : 1.17
- MAE : 1.02
- R2 Score : 0.00
- Variance Score : 0.00
- MedAE : 0.81

**Prophet**
- MSE : 1.17
- MAE : 1.02
- R2 Score : 0.00
- Variance Score : 0.00
- MedAE : 0.83

**Random Forest**
- MSE : 0.012
- MAE : 0.021
- R2 Score : 0.989
- Variance Score : 0.989
- MedAE : 0.00

Mean Squared Error (MSE)

Mean Absolute Error (MAE)

Median Absolute Error (MedAE)

# Model Justification

- The use of diverse machine learning algorithms, including Prophet, Long short-term memory (LSTM), Autoregressive Integrated Moving Average (ARIMA) demonstrates a comprehensive approach to model selection.

- The Random Forest model was chosen as the optimal model for this project due to its high performance, as evidenced by a 0.98 R2 Score and 0.98 Explained Variance Score.

- This indicates that the model can explain 98% of the variance in the data, suggesting a strong ability to make accurate predictions.

# Limitations

- The availability of comprehensive and up-to-date data, especially for post-pandemic scenarios, may pose challenges to the accuracy and reliability of the predictions.

- The model primarily focuses on key influencing factors and may not fully capture the complexity of socio-political, cultural, and environmental factors contributing to food insecurity.

- It provides a global analysis, but regional variations and unique challenges faced by different regions may not be fully captured.

- The interpretation of results and the implementation of policies to address food insecurity require additional considerations.

12/1/2023

# Summary

The project focuses on global analysis and prediction of food insecurity and hunger crisis outbreaks using machine learning techniques.

The goal is to address the constant rise in malnutrition and food deprivation by widening the scope of analysis from a regional to a global scale.

The proposed solution aims to incorporate a prediction-based approach that considers various factors, including economic shocks, extreme weather events, conflicts, and the impact of the COVID-19 pandemic.

# References

- Hernandez, D. C., Reesor, L., & Murillo, R. (2017). Food insecurity and adult overweight/obesity: Gender and race/ethnic disparities. Appetite, 117, 373–378. https://doi.org/10.1016/j.appet.2017.07.010

- Martini, G., Bracci, A., Riches, L., Jaiswal, S., Corea, M., Rivers, J., Husain, A., & Omodei, E. (2022). Machine learning can guide food security efforts when primary data are not available. Nature Food, 3(9), 716–728. https://doi.org/10.1038/s43016-022-00587-8

- Tarasuk, Valerie. (2001). Health Canada. (2002, December 4). ARCHIVED - Discussion Paper on Household and Individual Food Insecurity - Executive Summary. Canada.ca. https://www.canada.ca/en/health-canada/services/food-nutrition/healthy-eating/nutrition-policy-reports/executive-summary-discussion-paper-household-individual-food-insecurity.html

- Ahn, J., Briers, G. E., Baker, M., Price, E., Djebou, D. C. S., Strong, R. L., Piña, M., & Kibriya, S. (2022). Food security and agricultural challenges in West-African rural communities: a machine learning analysis. *International Journal of Food Properties*, *25*(1), 827–844. https://doi.org/10.1080/10942912.2022.2066124

- Westerveld, J. J., Van Den Homberg, M., Nobre, G. G., Van Den Berg, D. L., Teklesadik, A., & Stuit, S. (2021). Forecasting transitions in the state of food security with machine learning using transferable features. *Science of the Total Environment*, *786*, 147366. https://doi.org/10.1016/j.scitotenv.2021.147366

- Deléglise, H., Interdonato, R., Bégué, A., D'Hôtel, É. M., Teisseire, M., & Roche, M. (2022). Food security prediction from heterogeneous data combining machine and deep learning methods. Expert Systems With Applications, 190, 116189. https://doi.org/10.1016/j.eswa.2021.116189

- Negm, M. M. (2018). Modeling and Forecasting of Food Security for Wheat in Egypt Through Year 2025 by Using time Series ARIMA Models. *The Journal of Social Sciences Research*, *5*, 510–518. https://doi.org/10.32861/jssr.spi5.510.518

- *Commodity Terms of Trade dataset.* (2023). International Monetary Fund. https://data.imf.org/?sk=2cddccb8-0b59-43e9-b6a0-59210d5605d2

- *Food security dataset.* (2021). United States Census Bureau. https://www.census.gov/data/datasets/time-series/demo/cps/cps-supp_cps-repwgt/cps-food-security.html

- *Gross domestic product (GDP) survey data.* (2022). Organization for Economic Co-operation and Development. https://stats.oecd.org/

- *RASSF dataset.* (2022). European Commission. RASFF Database (Version 2.4.1) https://webgate.ec.europa.eu/rasff-window/screen/search

THANK YOU