

File - Users\jaytonschmeckle\Local\Repos\Math-CSCE-440-Final-Project\MATH440-Project\Notebooks\01-DataPreparation-CSV

```
1  #%% md
2  # Data Preparation Notebook
3
4  ## Purpose
5  This notebook is dedicated to preparing all
  datasets for subsequent analysis. The focus is on
  cleaning, structuring, and standardizing data to
  ensure consistency across analyses.
6
7  ## Datasets Overview
8  This section lists all the datasets involved, such
  as:
9  - Russian Losses as Documented by 3rd Party
10 - Ukrainian Losses as Reported by Ukrainian State
11 - [Add others as applicable]
12
13 ## Tools and Libraries
14 This notebook utilizes Python libraries including
  pandas for data manipulation and NumPy for
  numerical operations.
15 #%% md
16 # Setup
17
18 ## Import Libraries
19 Here, we import all necessary Python libraries
  needed for data preparation tasks.
20
21
22 #%%
23 # Import necessary libraries
24 import pandas as pd
25 import numpy as np
26 from sqlalchemy import create_engine
27 # TODO: Import any additional libraries needed
28
29 #%% md
30 # Data Loading
31
32 ## Load Data
33 Each dataset is loaded from its respective source.
  Detailed instructions and code for loading each
```

Page 1 of 7

File - Users\jaytonschmeckle\Local\Repos\Math-CSCE-440-Final-Project\MATH440-Project\Notebooks\01-DataPreparation-CSV

```
33 specific dataset are provided below.
34
35 #%%
36 import pandas as pd
37
38 # Load the first dataset (Documented Losses for
  both Ukraine and Russia)
39 file_path_documented = '/Users/jaytonschmeckle/
  Local/Repos/Math-CSCE-440-Final-Project/MATH440-
  Project/Data/Russia-Ukraine Equipment Losses -
  Original.csv'
40 df_documented = pd.read_csv(file_path_documented)
41 documented_headers = df_documented.columns.tolist()
42
43 # Load the second dataset (Russian Losses as
  Claimed by Ukraine)
44 file_path_claimed = '/Users/jaytonschmeckle/Local/
  Repos/Math-CSCE-440-Final-Project/MATH440-Project/
  Data/Russian Losses-ClaimedbyUkraine.csv'
45 df_claimed = pd.read_csv(file_path_claimed)
46 claimed_headers = df_claimed.columns.tolist()
47
48 # Display headers from both datasets with explicit
  descriptions
49 print("Headers from Documented Losses (both Ukraine
  and Russia):")
50 print(documented_headers)
51 print("\nHeaders from Russian Losses as Claimed by
  Ukraine:")
52 print(claimed_headers)
53
54 #%% md
55 # Data Structuring
56
57 ## Reshape Data
58 Data is organized into a consistent format across
  all datasets for easier analysis.
59
60 ## Merge/Concatenate
61 Data from different sources is combined if
  necessary to create a unified dataset for analysis.
```

Page 2 of 7

File - Users\jaytonschmeckle\Local\Repos\Math-CSCE-440-Final-Project\MATH440-Project\Notebooks\01-DataPreparation-CSV

```
62
63 #%%
64 # Initialize a new DataFrame for standardized data
65 df_standardized = pd.DataFrame()
66
67 # Date Column
68 df_standardized['Date'] = df_documented['Date']
69
70 # Tanks
71 df_standardized['RUS-Tanks-Documented'] =
  df_documented['Russia_Tanks']
72 df_standardized['RUS-Tanks-Claimed-UKR'] =
  df_claimed['tanks']
73
74 # Aircraft
75 df_standardized['RUS-Aircraft-Documented'] =
  df_documented['Russia_Aircraft']
76 df_standardized['RUS-Aircraft-Claimed-UKR'] =
  df_claimed['planes'] + df_claimed['helicopters']
  + df_claimed['uav']
77
78 # Artillery
79 df_standardized['RUS-Artillery-Documented'] =
  df_documented['Russia_Artillery']
80 df_standardized['RUS-Artillery-Claimed-UKR'] =
  df_claimed['artillery_systems'] + df_claimed['mlrs']
  + df_claimed['missiles']
81
82 # Air Defense
83 df_standardized['RUS-AirDefense-Documented'] =
  df_documented['Russia_Antiair']
84 df_standardized['RUS-AirDefense-Claimed-UKR'] =
  df_claimed['air_defense_equipment']
85
86 # Vehicles
87 df_standardized['RUS-Vehicles-Documented'] =
  df_documented['Russia_Vehicles']
88 df_standardized['RUS-Vehicles-Claimed-UKR'] =
  df_claimed['cars_and_tank_trucks']
89
90 # Personnel
```

Page 3 of 7

File - Users\jaytonschmeckle\Local\Repos\Math-CSCE-440-Final-Project\MATH440-Project\Notebooks\01-DataPreparation-CSV

```
91 df_standardized['RUS-Personnel-Documented'] =
  df_documented['Russia_Infantry']
92 df_standardized['RUS-Personnel-Claimed-UKR'] =
  df_claimed['personnel']
93
94 #%% md
95 # Output
96
97 ## Save Clean Data
98 The cleaned and structured data is saved to new
  files or databases for easy access in subsequent
  analysis phases.
99
100 #%%
101 # Save the standardized data to a new CSV file
102 df_standardized.to_csv('standardized_data.csv',
  index=False)
103
104 # Or analyze directly within the notebook
105 print(df_standardized.head())
106
107 #%% md
108
109 #%% md
110 ## Data Cleaning Function
111
112 This function is designed to streamline the
  process of checking and cleaning data for any
  selected category of losses (e.g., tanks, aircraft
  ). It handles missing values, corrects data types
  , and ensures that the dataset is ready for
  further analysis. By parameterizing the function,
  we can apply the same cleaning logic to different
  data categories efficiently.
113
114 #%%
115 import pandas as pd
116 import numpy as np
117 import matplotlib.pyplot as plt
118 from sklearn.metrics import mean_absolute_error,
  mean_squared_error
```

Page 4 of 7

File - Users\jaytonschmeckle\Local\Repos\Math-CSCE-440-Final-Project\MATH440-Project\Notebooks\01-DataPreparation-CSV

```
119 from scipy.optimize import newton
120 from scipy.interpolate import CubicSpline
121 import seaborn as sns
122
123 def clean_data(category):
124     """Cleans data for the specified category of
  losses.
125
126     Args:
127         category (str): The category to clean, e.g
  ., 'Tanks', 'Aircraft'.
128     """
129     documented_col = f'RUS-{category}-Documented'
130     claimed_col = f'RUS-{category}-Claimed-UKR'
131
132     # Check for missing values
133     print(f"Missing values before cleaning for {
  category}:")
134     print(df_standardized[[documented_col,
  claimed_col]].isnull().sum())
135
136     # Handle missing values by forward filling
137     df_standardized[documented_col].fillna(method=
  'ffill', inplace=True)
138     df_standardized[claimed_col].fillna(method='
  ffill', inplace=True)
139
140     # Convert data types to integer
141     df_standardized[documented_col] =
  df_standardized[documented_col].astype(int)
142     df_standardized[claimed_col] = df_standardized
  [claimed_col].astype(int)
143
144     # Verify changes
145     print(f"Data types after conversion for {
  category}:")
146     print(df_standardized[[documented_col,
  claimed_col]].dtypes)
147     print("\nSample data after cleaning:")
148     print(df_standardized[[documented_col,
  claimed_col]].head())
```

Page 5 of 7

File - Users\jaytonschmeckle\Local\Repos\Math-CSCE-440-Final-Project\MATH440-Project\Notebooks\01-DataPreparation-CSV

```
149
150 clean_data('Tanks') # Modify 'Tanks' to any other
  category needed
151
152 #%% md
153 ## Verifying Changes Across Categories
154
155 After applying the cleaning function to different
  categories, this section is intended for final
  verification. You can run diagnostics, visualize
  data, or apply further checks to ensure that all
  categories have been cleaned and are ready for in-
  depth analysis.
156
157 #%%
158 # Example of how to check data across multiple
  categories
159 categories = ['Tanks', 'Aircraft', 'Artillery', '
  AirDefense', 'Vehicles', 'Personnel']
160 for category in categories:
161     documented_col = f'RUS-{category}-Documented'
162     claimed_col = f'RUS-{category}-Claimed-UKR'
163     print(f"{category} - Missing values check:")
164     print(df_standardized[[documented_col,
  claimed_col]].isnull().sum())
165     print(f"{category} - Data type check:")
166     print(df_standardized[[documented_col,
  claimed_col]].dtypes)
167     print()
168
169 #%% md
170 ## Error Analysis
171
172 Calculate the Mean Absolute Error (MAE) and Mean
  Squared Error (MSE) between the documented and
  claimed data to quantify discrepancies.
173
174 #%%
175 def calculate_errors(documented, claimed):
176     mae = mean_absolute_error(documented, claimed)
177     mse = mean_squared_error(documented, claimed)
```

Page 6 of 7

```

178     return mae, mse
179
180 documented_tanks = df_standardized['RUS-Tanks-
    Documented']
181 claimed_tanks = df_standardized['RUS-Tanks-Claimed
    -UKR']
182 tanks_mae, tanks_mse = calculate_errors(
    documented_tanks, claimed_tanks)
183
184 print(f"Mean Absolute Error for Tanks: {tanks_mae}
    ")
185 print(f"Mean Squared Error for Tanks: {tanks_mse}"
    )
186
187 """ md
188 ## Linear Regression Analysis
189
190 Model the relationship between documented and
    claimed tank losses using linear regression to
    assess predictability and consistency of the
    reported data.
191
192 """
193 x = documented_tanks.values.reshape(-1, 1)
194 y = claimed_tanks.values.reshape(-1, 1)
195 A = np.vstack([x.T, np.ones(len(x))]).T
196 m, c = np.linalg.lstsq(A, y, rcond=None)[0]
197
198 plt.figure(figsize=(20, 10))
199 plt.plot(x, y, 'o', label='Original data',
    markersize=10)
200 plt.plot(x, m*x + c, 'r', label='Fitted line')
201 plt.xlabel('Documented Tank Losses')
202 plt.ylabel('Claimed Tank Losses by Ukraine')
203 plt.title('Linear Regression Analysis of Tank
    Losses')
204 plt.legend()
205 plt.show()
206
207 print(f"Slope: {m[0]}, Intencept: {c[0]}")
208

```