

Unveiling the Power of Social Support in Healthcare



TEXAS

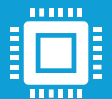
The University of Texas at Austin



Executive Summary: Unveiling the Power of Social Support in Healthcare



Our upcoming Self-Learning Tutorial (SLT) aims to shed light on the often-underestimated impact of social connections—specifically marriage and companionship—on patient health outcomes. This tutorial will provide a practical, step-by-step guide for anyone interested in exploring this vital link.



We'll start by accessing a vast, real-world medical database (MIMIC-III) stored in Google's cloud. From this extensive data, we'll carefully select and extract key information to build a compelling case for how social support influences health. This extracted data will then be prepared for in-depth analysis and compelling visual presentations, allowing us to easily see and understand important trends.

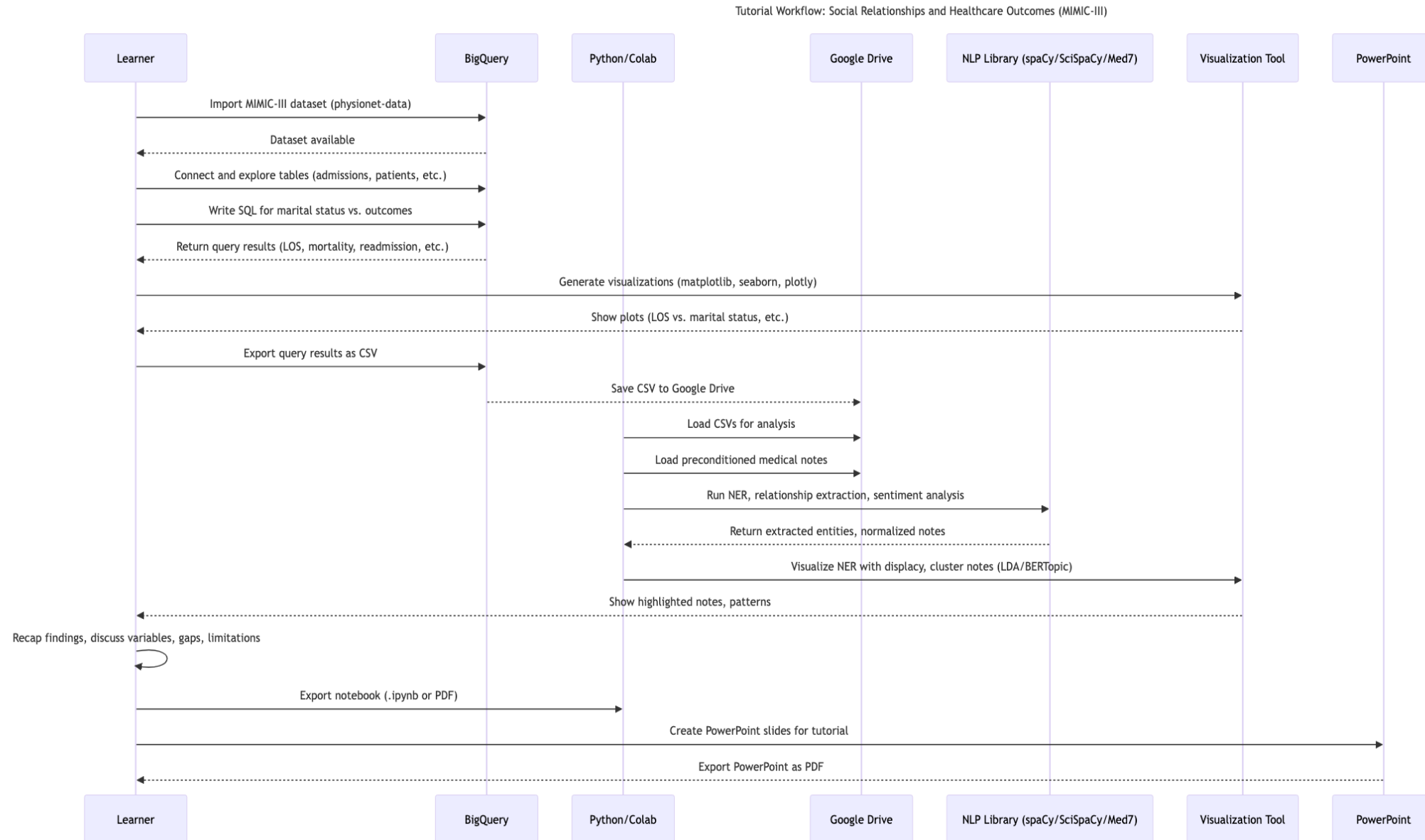


A significant part of the tutorial will involve using advanced text analysis tools to uncover crucial medical details hidden within large volumes of patient notes. This will help us pinpoint specific mentions related to social support and its context in healthcare records.



Ultimately, the SLT will synthesize all findings into clear conclusions, demonstrating the tangible benefits of social connections on health. The entire tutorial is designed to be highly hands-on, providing explicit instructions and code walk-throughs using various data analysis and text processing tools, ensuring that learners can easily follow along and grasp the concepts. Our goal is to empower individuals to understand and articulate the "unseen" yet critical role of social support in promoting better healthcare outcomes.

Synthesis and Workflow in Self-Learning Tutorial (SLT)



sequenceDiagram

title Tutorial Workflow: Social Relationships and Healthcare Outcomes (MIMIC-III)

- participant User as Learner
- participant BigQuery
- participant Python as Python/Colab
- participant Drive as Google Drive
- participant NLP as NLP Library (spaCy/SciSpaCy/Med7)
- participant Viz as Visualization Tool
- participant PPT as PowerPoint

```
%% 1. Data Import
User->>BigQuery: Import MIMIC-III dataset (physionet-data)
BigQuery-->>User: Dataset available

%% 2. Data Exploration & SQL Analysis
User->>BigQuery: Connect and explore tables (admissions, patients, etc.)
User->>BigQuery: Write SQL for marital status vs. outcomes
BigQuery-->>User: Return query results (LOS, mortality, readmission, etc.)

%% 3. Visualization of Findings
User->>Viz: Generate visualizations (matplotlib, seaborn, plotly)
Viz-->>User: Show plots (LOS vs. marital status, etc.)

%% 4. Data Export & Python Analysis
User->>BigQuery: Export query results as CSV
BigQuery-->>Drive: Save CSV to Google Drive
Python->>Drive: Load CSVs for analysis

%% 5. NLP with Medical Notes
Python->>Drive: Load preconditioned medical notes
Python->>NLP: Run NER, relationship extraction, sentiment analysis
NLP-->>Python: Return extracted entities, normalized notes

%% 6. Highlight Key Notes
Python->>Viz: Visualize NER with displacy, cluster notes (LDA/BERTopic)
Viz-->>User: Show highlighted notes, patterns

%% 7. Synthesis & Recap
User->>User: Recap findings, discuss variables, gaps, limitations

%% 8. Exportable Output
User->>Python: Export notebook (.ipynb or PDF)

%% 9. Create PowerPoint & Export as PDF
User->>PPT: Create PowerPoint slides for tutorial
PPT-->>User: Export PowerPoint as PDF
```



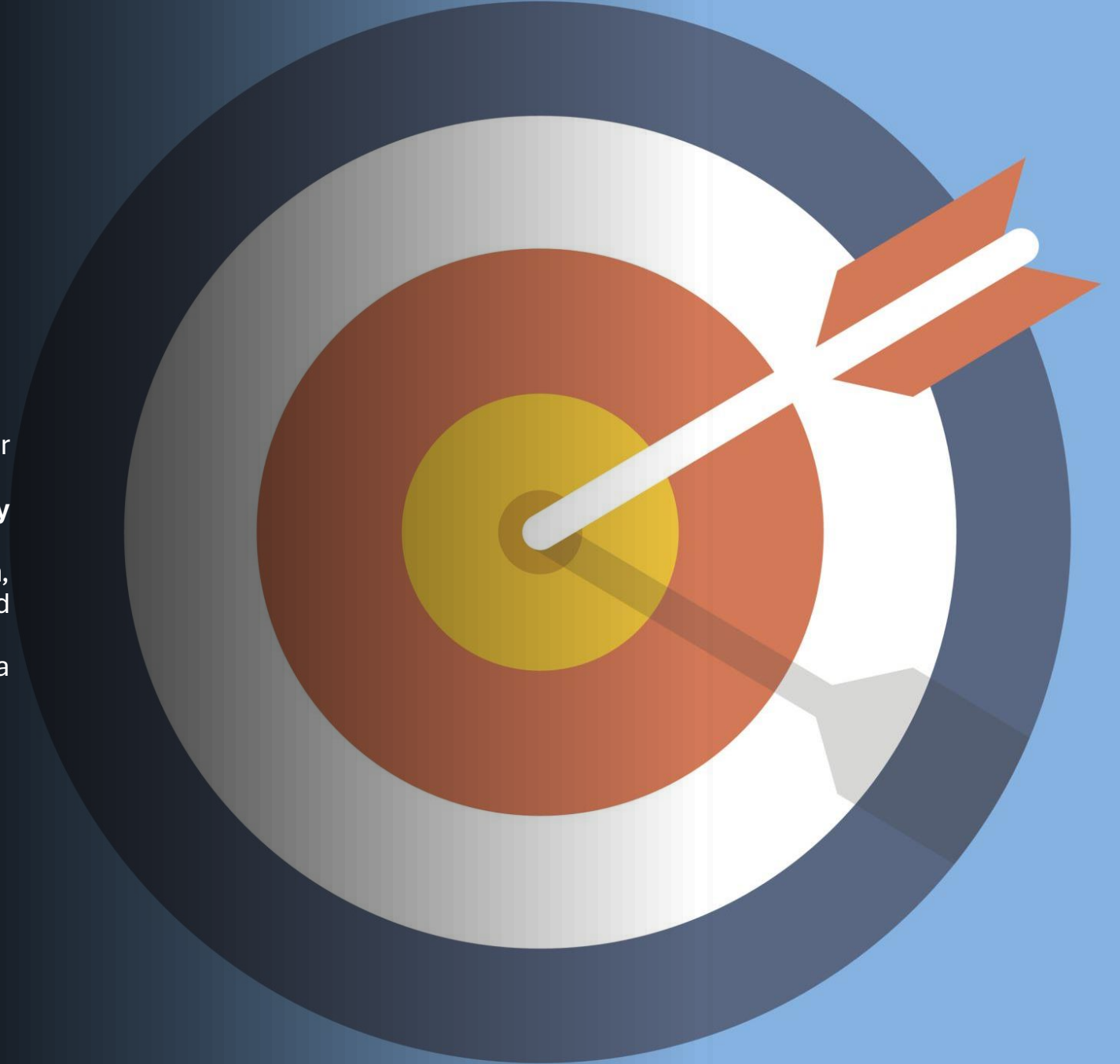
<https://mermaid.live/>

Synthesis and Workflow in Self-Learning Tutorial (SLT) UML Code

Tutorial Goal

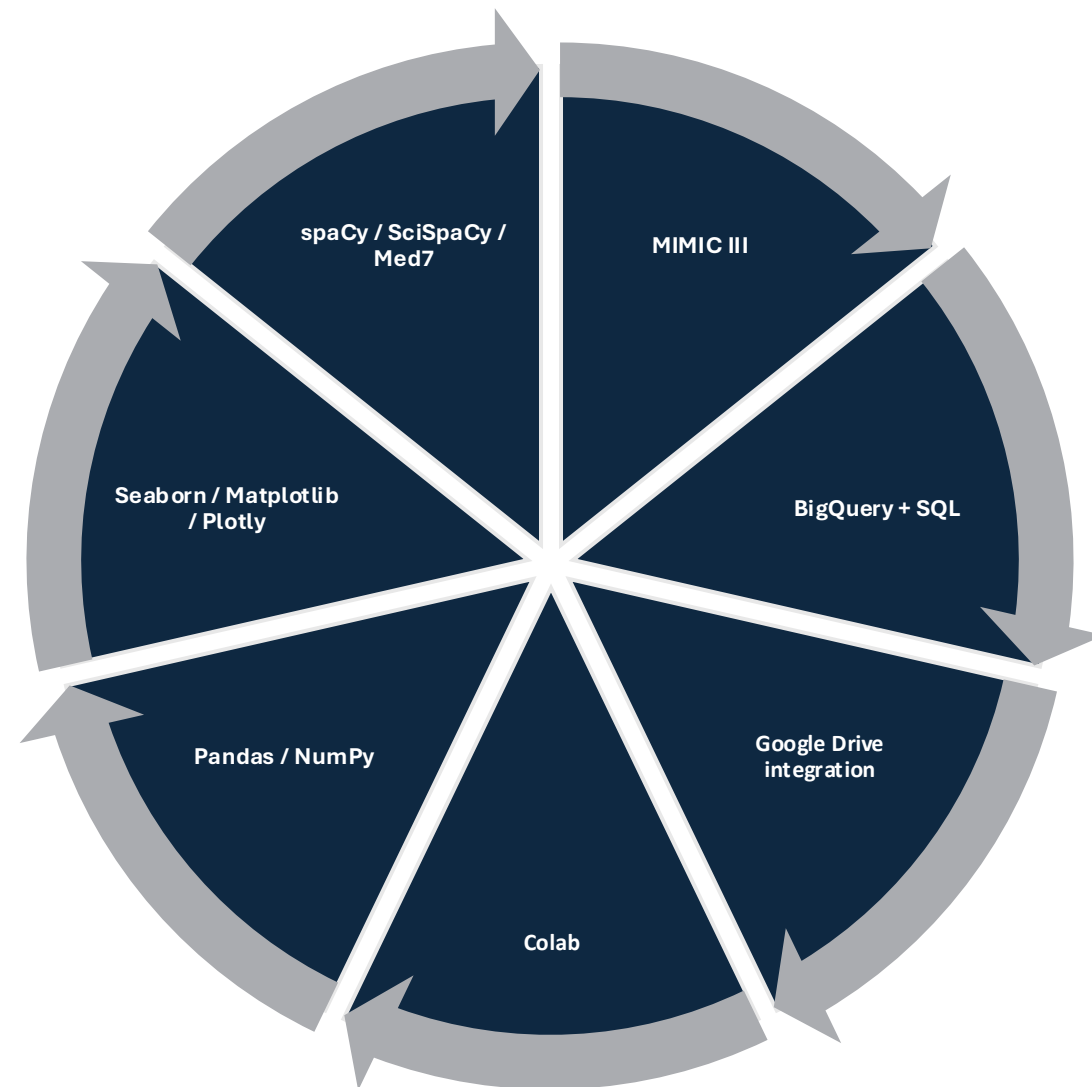
Learners will:

- Understand how social relationships (like marriage or companionship) correlate with healthcare outcomes.
- Gain hands-on experience querying **MIMIC-III** in **BigQuery**
Learn Clean, Tokenize and Lemmatization, NER
- Analyze structured data and medical notes using **Python**, **Pandas**, **NLP** (**spaCy/medSpaCy/SciSpaCy/Med7**), and **visualization tools**.
- Learn to normalize, visualize, and interpret clinical data with a health policy/public health lens.





Tools/Tech Stack



Collab to Bigquery

BigQuery Integration and Data Prep

- Connect to BigQuery
- Introduce key MIMIC-III tables: admissions, patients, diagnoses_icd, chartevents, noteevents, etc.
- Write SQLs to explore:
 - Marital status vs. outcomes
 - Length of stay, readmission, mortality, etc.
 - Compare outcomes across married vs. single vs. widowed/divorced patients

```
1 # 1. Install Necessary Libraries
2 # Connecting to BigQuery from Google Colab and pulling data is a straightforward process, primarily involving authentication and using the google-cloud-bigquery library.
3 # First, we'll need the google-cloud-bigquery library to interact with BigQuery and pandas to work with the data in a DataFrame.
4 !pip install google-cloud-bigquery pandas
```

```
1 # 2. Authenticate Your Colab Environment
2 # Google Colab provides a simple way to authenticate using your Google account. This is the most common method for interactive use in Colab
3 from google.colab import auth
4 auth.authenticate_user()
5 print('Authenticated')
```

↗ Authenticated

```
1 # 3. Connect Google Cloud Project ID
2 # We need to use the BigQuery client which Google Cloud Project your data resides in.
3 project_id = '<<your-project-id>>' # Replace with your actual Google Cloud Project ID
```

```
1 # 4. Initialize the BigQuery Client
2 # Create an instance of the BigQuery client. This object will manage your connection and allow you to send queries.
3 # |following block is commented here because we're consolidating all imports in one place which will be shown later.
4 from google.cloud import bigquery
5 client = bigquery.Client(project=project_id)
6 print('BigQuery client initialized.')
```

↗ BigQuery client initialized.



Mortality Counts and Average Age at Death by Marital Status

Mortality Counts and Average Age at Death by Marital Status

Goal:

To determine the number of deceased patients and their average age at death, grouped by marital status.

Implementation:

- Joins patients and admissions tables on subject_id.
- Filters for patients with a non-null date of death (p.dod) and non-null admission time.
- Calculates:
 - death_count: Number of deaths per marital status.
 - avg_age_at_death: Average age at death, in years, rounded to 1 decimal.
- Results are grouped by marital status and ordered by average age at death (descending).

```
1 # 5. Write Your SQL Query
2 # Craft your SQL query. Remember to specify the full path for your tables using the format resources.dataset_name.table_name.
3 # This Counts of patients who died within 7 days of hospital admission by marital status and age at death
4 query = """
5 SELECT
6     a.marital_status,
7     COUNT(*) AS death_count,
8     ROUND(AVG(DATE_DIFF(p.dod, p.dob, YEAR)), 1) AS avg_age_at_death
9 FROM
10    physionet-data.mimiciii_demo.patients AS p
11 JOIN
12    physionet-data.mimiciii_demo.admissions a
13    ON p.subject_id = a.subject_id
14 WHERE
15     p.dod IS NOT NULL
16     AND a.admittime IS NOT NULL
17     --AND p.dod <= DATE_ADD(a.admittime, INTERVAL 7 DAY)
18 GROUP BY
19     a.marital_status
20 ORDER BY
21     avg_age_at_death DESC;
22 """
23 print('SQL query defined.')
```

SQL query defined.

Outcome & Interpretation

```
1 # import pandas as pd
2
3 try:
4     df_bigquery_data = client.query(query).to_dataframe()
5     print('\nQuery executed successfully! Data loaded into pandas DataFrame.')
6     print('First 5 rows of the data:')
7     print(df_bigquery_data.head())
8 except Exception as e:
9     print(f'\nAn error occurred: {e}')
10    print("Please double-check:")
11    print(" - Your `project_id` is correct.")
12    print(" - The table names in your SQL query are fully qualified")
13    print(" - Your Google account has the necessary BigQuery permissions")
14    print(" - The BigQuery API is enabled for your Google Cloud Project.")
```



Query executed successfully! Data loaded into pandas DataFrame.
First 5 rows of the data:

	marital_status	death_count	avg_age_at_death
0	None	16	122.4
1	WIDOWED	15	97.7
2	MARRIED	60	81.9
3	UNKNOWN (DEFAULT)	5	79.6
4	DIVORCED	6	76.5

Interpretations:

- **Implausible Ages:** The group with `marital_status = None` has an average age at death of 122.4 years, which is biologically implausible and suggests data entry or calculation errors (e.g., missing or incorrect dates of birth/death).
- **Widowed:** High average age at death (97.7), which is plausible given demographic trends.
- **Married:** Largest group (60 deaths), average age at death 81.9.
- **Unknown/Default/None:** These categories indicate incomplete or improperly coded marital status.
- **Divorced:** Lowest average age at death among these groups (76.5).

Recommendations:

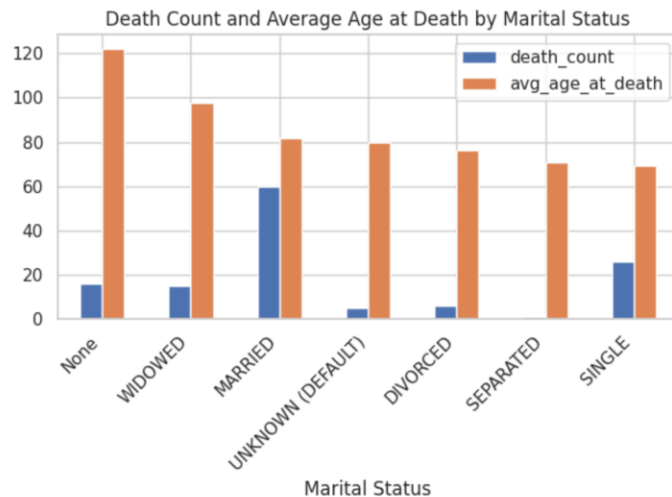
- The results are affected by data quality issues and the lack of a time-to-death filter. Data cleaning is recommended for more accurate insights.
- The presence of implausible ages and ambiguous marital status categories (None, UNKNOWN (DEFAULT)) indicates a need for data cleaning.
- The filter for deaths within 7 days of admission is commented out. As written, the query includes all deaths, not just those soon after admission.
- Clean data by removing implausible ages (e.g., age <0 or >110).
- Standardize or exclude ambiguous marital status.
- We can analyze early mortality, uncomment and use the 7-day filter.

Death Count by Marital Status (Bar Plot)

Goal:

Visualize the query results using a bar chart in pandas and matplotlib, which is ideal for comparing categorical data such as marital status against death counts or average age at death.

```
1 # import matplotlib.pyplot as plt
2
3 ax = df_bigquery_data.plot.bar(x='marital_status', y=['death_count', 'avg_age_at_death'])
4 plt.title('Death Count and Average Age at Death by Marital Status')
5 plt.xlabel('Marital Status')
6 plt.xticks(rotation=45, ha='right')
7 plt.tight_layout()
8 plt.show()
```



- Obtain The Data: We have successfully queried the data from Google BigQuery into a pandas DataFrame named `df_bigquery_data`.
- Run the Plotting Code: Once you have your `df_bigquery_data` DataFrame ready, execute the provided Python code in your environment (e.g., a Jupyter Notebook or Google Colab cell).
- Interpret the Plot: After running the code, a bar chart will be displayed. This chart visually represents the `death_count` and `avg_age_at_death` for each category of `marital_status`, allowing for quick comparison and analysis of these metrics across different marital groups.
- Conclusion: This result provides a summary of mortality counts and average age at death by marital status, but the results are affected by data quality issues and the lack of a time-to-death filter. Data cleaning is required for more accurate insights. We can perform this cleaning two ways:
 - 1. with SQL filter or
 - 2. Clean in Pandas after loading

Define BigQuery Helper Function

```
1 # Execute the cleaned SQL to retrieve data set
2 # import pandas as pd
3
4 df_bigquery_data = run_bq_query(query)
```



Query executed successfully! Data loaded into pandas DataFrame.
First 5 rows of the data:

	marital_status	death_count	avg_age_at_death
0	None	16	122.4
1	WIDOWED	15	97.7
2	MARRIED	60	81.9
3	UNKNOWN (DEFAULT)	5	79.6
4	DIVORCED	6	76.5



```
1 # Execute the SQL to retrieve data set
2 # import pandas as pd
3 def run_bq_query(query):
4     try:
5         df_bigquery_data = client.query(query).to_dataframe()
6         print('\nQuery executed successfully! Data loaded into pandas DataFrame.')
7         print('First 5 rows of the data:')
8         print(df_bigquery_data.head())
9         return df_bigquery_data
10    except Exception as e:
11        print(f'\nAn error occurred: {e}')
12        print("Please double-check:")
13        print("  - Your `project_id` is correct.")
14        print("  - The table names in your SQL query are fully qualified")
15        print("  - Your Google account has the necessary BigQuery permissions")
16        print("  - The BigQuery API is enabled for your Google Cloud Project.")
```

Goal:

Now we take the opportunity to improvise our client query to avoid writing redundant code and convert it to a Helper function.

Data Cleansing with SQL:

This SQL query is incorporating data cleaning steps for analysis

```
1 # SQL query for your analysis, incorporating data cleaning steps
2 query_cleaned = """
3 SELECT
4     CASE
5         WHEN a.marital_status IS NULL
6             OR a.marital_status IN ('UNKNOWN (DEFAULT)', 'None', '')
7             THEN 'Unknown'
8         ELSE a.marital_status
9     END AS marital_status_cleaned,
10    COUNT(*) AS death_count,
11    ROUND(AVG(DATE_DIFF(p.dod, p.dob, YEAR)), 1) AS avg_age_at_death
12 FROM
13     physionet-data.mimiciii_demo.patients AS p
14 JOIN
15     physionet-data.mimiciii_demo.admissions AS a
16     ON p.subject_id = a.subject_id
17 WHERE
18     p.dod IS NOT NULL
19     AND a.admittime IS NOT NULL
20     AND DATE_DIFF(p.dod, p.dob, YEAR) BETWEEN 0 AND 110
21     -- Uncomment the next line to restrict to deaths within 7 days of admission
22     -- AND p.dod <= DATE_ADD(a.admittime, INTERVAL 7 DAY)
23 GROUP BY
24     marital_status_cleaned
25 ORDER BY
26     avg_age_at_death DESC;
27 """
28 print('SQL query defined.')
29
```

SQL query defined.

Cleaned Query Results & Interpretations:

```
1 # Execute the cleaned SQL to retrieve data set
2 # import pandas as pd
3
4 marital_counts_cleaned = run_bq_query(query_cleaned)
```



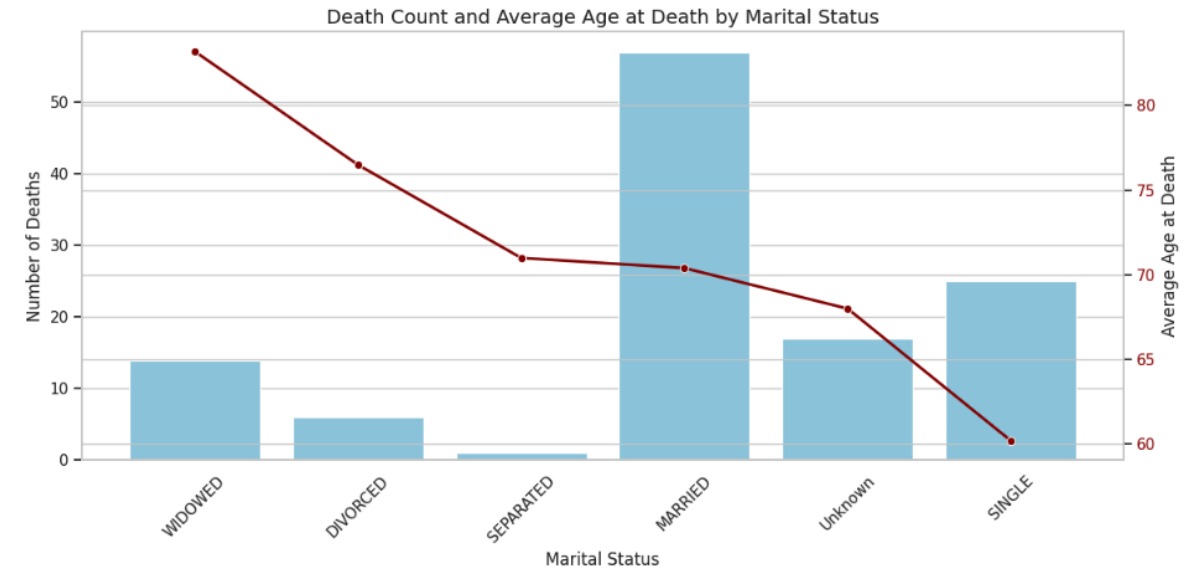
Query executed successfully! Data loaded into pandas DataFrame.
First 5 rows of the data:

	marital_status_cleaned	death_count	avg_age_at_death
0	WIDOWED	14	83.2
1	DIVORCED	6	76.5
2	SEPARATED	1	71.0
3	MARRIED	57	70.4
4	Unknown	17	68.0

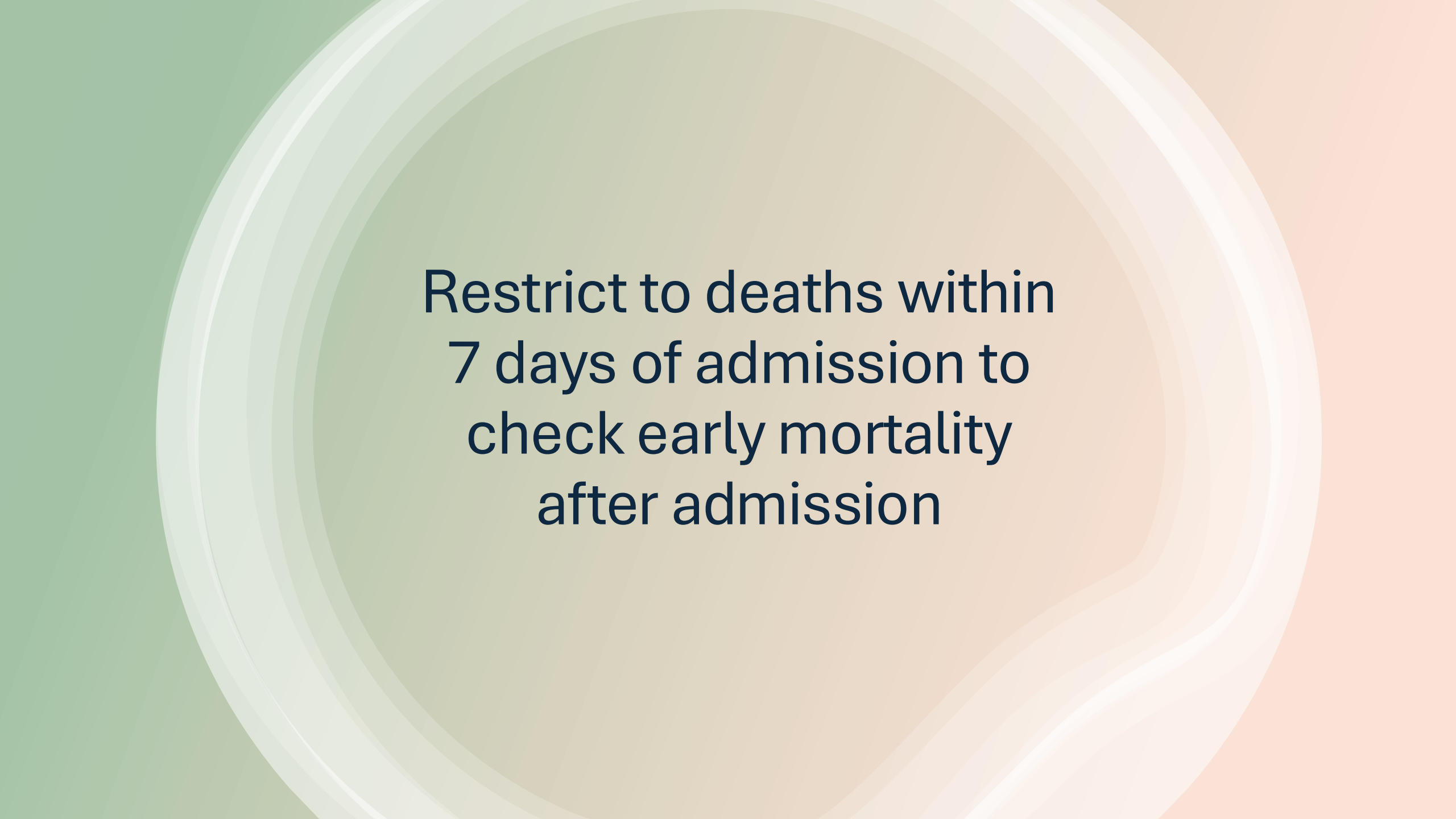
- Married individuals have the highest death count, but they also tend to die at a younger average age (~70) compared to widowed and divorced individuals. This could be due to a higher representation of married individuals in the overall dataset, and not necessarily poorer health outcomes.
- Widowed individuals show a lower death count but the highest average age at death (~83). This indicates longer life expectancy, potentially due to prior long-term companionship or spousal care before widowhood.
- Divorced and separated individuals fall in between - fewer deaths but also lower average age at death than the widowed. May suggest reduced social support after separation, affecting long-term health outcomes.
- Single individuals (from the extended version of the dataset you plotted) show lower life expectancy (~60) with a relatively high number of deaths, supporting the hypothesis that lack of companionship may contribute to poorer outcomes.
- Unknown marital status also skews younger, but is hard to interpret definitively due to lack of context.

Combined Visualization of Death Count and Average Age at Death by Marital Status

```
1 # Combined Bar + Line Plot: Death Count and Average Age at Death
2 fig, ax1 = plt.subplots(figsize=(12, 6))
3
4 # Bar plot for death count
5 sns.barplot(
6     x='marital_status_cleaned',
7     y='death_count',
8     data=marital_counts_cleaned,
9     ax=ax1,
10    color='skyblue'
11 )
12 ax1.set_ylabel('Number of Deaths', fontsize=12)
13 ax1.set_xlabel('Marital Status', fontsize=12)
14 ax1.set_title('Death Count and Average Age at Death by Marital Status', fontsize=14)
15 ax1.tick_params(axis='y')
16 plt.setp(ax1.get_xticklabels(), rotation=45)
17
18 # Line plot for average age at death on secondary y-axis
19 ax2 = ax1.twinx()
20 sns.lineplot(
21     x='marital_status_cleaned',
22     y='avg_age_at_death',
23     data=marital_counts_cleaned,
24     ax=ax2,
25     color='darkred',
26     marker='o',
27     linewidth=2
28 )
29 ax2.set_ylabel('Average Age at Death', fontsize=12)
30 ax2.tick_params(axis='y', labelcolor='darkred')
31
32 # Show plot
33 plt.tight_layout()
34 plt.show()
```



- **Create a figure and primary axis:** Set the canvas size to ensure clear display.
- **Draw the bar plot :** Use the primary axis to plot the number of deaths per marital status. Apply a consistent color for clarity (Sky blue). Add axis labels and a title. Rotate x-axis labels for readability.
- **Create a secondary y-axis:** Mirror the x-axis but allow for a different y-scale. This axis is used to plot average age at death.
- **Overlay the line plot:** Use the secondary axis to draw a line graph of average age at death. Add data markers for each marital status. Use a distinct color (dark red) to differentiate from the bar plot.
- **Label the secondary axis:** Add appropriate y-axis label and match it to the line color.
- **Finalize the layout :** Apply tight layout to avoid overlaps. Display the combined chart.



Restrict to deaths within
7 days of admission to
check early mortality
after admission

Mortality Counts and Average Age at Death by Marital Status within 7 days of admission

Implementation:

- Joins patients and admissions tables on subject_id.
- Filters for patients with a non-null date of death (p.dod) and non-null admission time.
- Calculates:
- death_count: Number of deaths per marital status.
- avg_age_at_death: Average age at death, in years, rounded to 1 decimal.
- This time Restrict to deaths within 7 days of admission to check early mortality after admission
- Results are grouped by marital status and ordered by average age at death (descending).

Goal:

To determine the number of deceased patients and their average age at death, grouped by marital status.

```
1 | # Restrict to deaths within 7 days of admission to check early mortality after admission
2
3 query_early_mortality = """
4 SELECT
5     CASE
6         WHEN a.marital_status IS NULL
7             OR a.marital_status IN ('UNKNOWN (DEFAULT)', 'None', '')
8             THEN 'Unknown'
9         ELSE a.marital_status
10    END AS marital_status_cleaned,
11    COUNT(*) AS death_count,
12    ROUND(AVG(DATE_DIFF(p.dod, p.dob, YEAR)), 1) AS avg_age_at_death
13 FROM
14     physionet-data.mimiciii_demo.patients AS p
15 JOIN
16     physionet-data.mimiciii_demo.admissions AS a
17   ON p.subject_id = a.subject_id
18 WHERE
19     p.dod IS NOT NULL
20     AND a.admittime IS NOT NULL
21     AND DATE_DIFF(p.dod, p.dob, YEAR) BETWEEN 0 AND 110
22     AND p.dod <= DATE_ADD(a.admittime, INTERVAL 7 DAY)
23 GROUP BY
24     marital_status_cleaned
25 ORDER BY
26     avg_age_at_death DESC;
27 """
28 print('SQL query defined.')
```

SQL query defined.

Outcome & Interpretation

```
1 # Execute the cleaned SQL to retrieve data set
2 # import pandas as pd
3
4 early_mortality_df = run_bq_query(query_early_mortality)
```



Query executed successfully! Data loaded into pandas DataFrame.
First 5 rows of the data:

	marital_status_cleaned	death_count	avg_age_at_death
0	WIDOWED	2	85.0
1	MARRIED	7	79.1
2	DIVORCED	3	72.3
3	Unknown	7	64.4
4	SINGLE	4	51.8

Interpretations:

- **WIDOWED** patients, although few in count, had the **highest average age at early death**, potentially due to **limited social support** post-loss of spouse.
- **MARRIED** patients had relatively **more early deaths** (7) but at a **higher average age (79.1)**, suggesting **supportive companionship may extend life** even in acute cases.
- **SINGLE** individuals had **lower average age at death (51.8)**—a striking indicator that **lack of support correlates with younger mortality**.
- **Unknown** marital status shows a comparable death count to married (7) but **much lower age at death**, which may indicate either missing documentation or neglected social support tracking.

Recommendations:

- Explore Further Stratification by Breaking down these stats by ICU type, gender, or comorbidity index.
- NLP Follow-Up to use unstructured clinical notes to detect loneliness, depression, or caregiver presence.

Early Mortality with in 7 days of admission (Bubble Plot)

Goal:

Visualize early mortality within 7 days of admission by marital status using a bubble plot, where the x-axis represents marital status, the y-axis shows average age at death, and bubble size encodes the number of deaths. This enables clear, multi-dimensional comparison of early mortality patterns across social relationship groups.

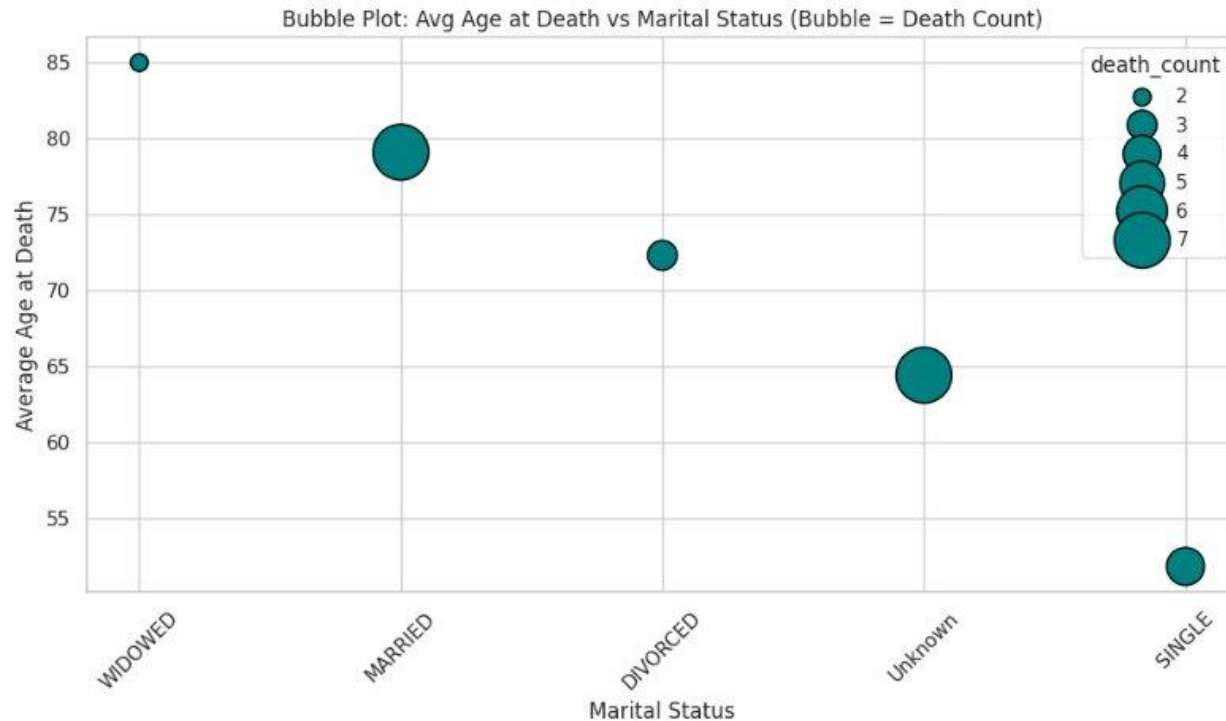
```
1 # Early Mortality with in 7 days of admission (Bubble Plot)
2
3 # Convert DataFrame to dictionary with lists as values
4 dict_stats = early_mortality_df.to_dict(orient='list')
5
6 # Bubble Plot: Death count (size), average age (y), marital status (x)
7 plt.figure(figsize=(10,6))
8 sns.scatterplot(
9     data=death_stats,
10    x='marital_status_cleaned',
11    y='avg_age_at_death',
12    size='death_count',
13    legend='brief',
14    sizes=(100, 1000),
15    color='teal',
16    edgecolor='black'
17 )
18 plt.title('Bubble Plot: Avg Age at Death vs Marital Status (Bubble = Death Count)')
19 plt.xlabel('Marital Status')
20 plt.ylabel('Average Age at Death')
21 plt.xticks(rotation=45)
22 plt.tight_layout()
23 plt.show()
```

Code Implementation:

- Convert DataFrame `early_mortality_df` to a dictionary with columns as keys and lists of column values as values, using the `to_dict` method with the 'list' orientation.
- Set up your plot area with a suitable figure size.
- Use Seaborn's `scatterplot` function to create a bubble plot:
- Set the x-axis to marital status.
- Set the y-axis to average age at death.
- Set the bubble size proportional to the death count.
- Customize the appearance with color, edge color, and legend.
- Add a title and axis labels, rotate the x-axis labels for readability, and use tight layout to prevent overlap.
- Display the plot.

This workflow will visually encode three variables—marital status, average age at death, and death count—into a single, interpretable bubble plot.

Visualization:



- **Conclusions:**

Married individuals likely benefit from stronger social and emotional support, which may delay mortality even during early critical illness. Those listed as single or divorced/separated often showed earlier average age of death—possibly due to loneliness, stress, or poor health adherence. Widowed Group: Though older (expected), their mortality can reflect loss of companionship impact.



Load files form Google
Drive

`df_4019_radiology_notes.csv` is a pre-filtered CSV file containing records with the ICD-9 code **4019** (Hypertension), specifically filtered for **Radiology** events. The source of this data is the **NOTEEVENTS** table from the **MIMIC-III** database.

```
1 # Steps to Load CSV from Google Drive in Google Colab
2 # df_4019_radiology_notes.csv' is a pre conditioned csv files with ICD9 code 4019, which is "Hypertension"
3 # and we have filtered for "Radiology" events
4 from google.colab import auth, drive, files
5
6 drive.mount('/content/drive')
7 auth.authenticate_user()
8
9 df_copy = pd.read_csv('/content/drive/My Drive/df_4019_radiology_notes.csv')
10 print(df_copy.head(3))
```

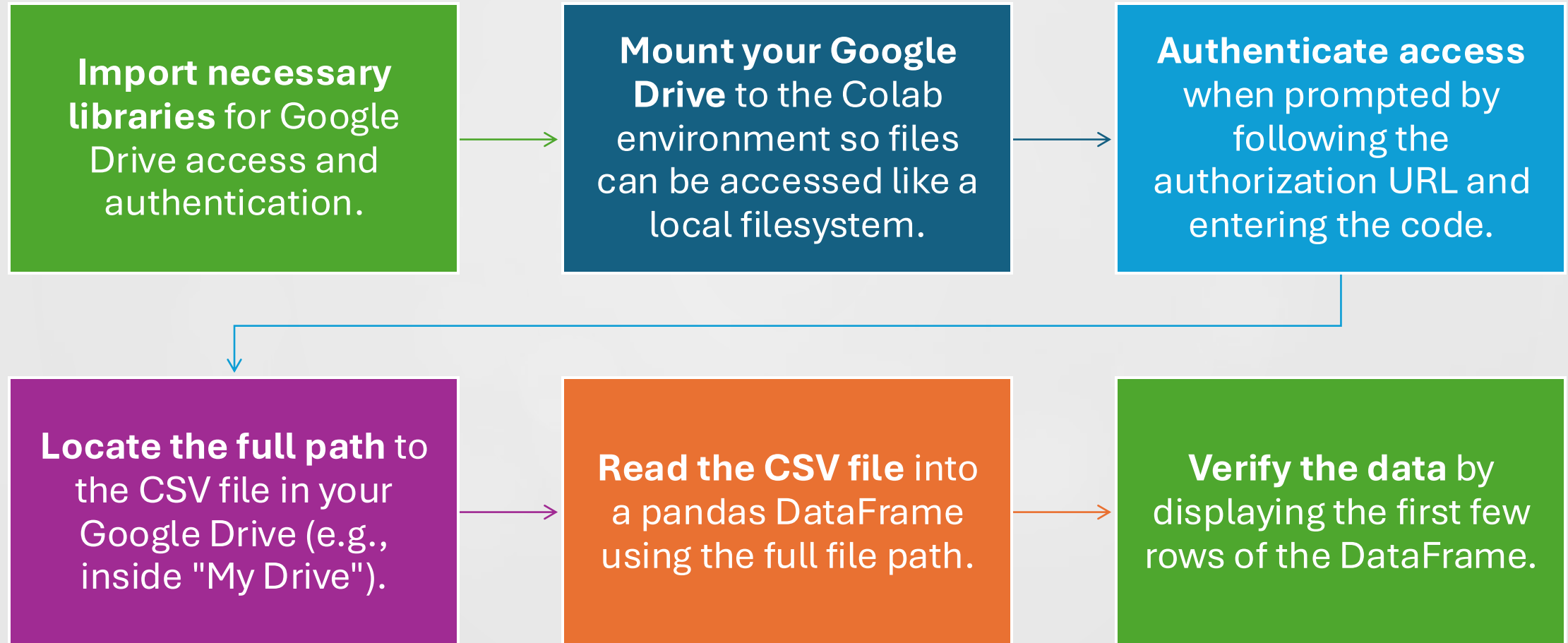
⇒ Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

	ROW_ID	SUBJECT_ID	HADM_ID	CHARTDATE	CHARTTIME	STORETIME	\
0	738405	93207	104110.0	2166-10-09	2166-10-09 20:24:00	NaN	
1	738661	24552	NaN	2153-01-12	2153-01-12 22:03:00	NaN	
2	738662	10390	NaN	2145-03-13	2145-03-13 02:58:00	NaN	

	CATEGORY	DESCRIPTION	CGID	ISERROR	\
0	Radiology	DISTINCT PROCEDURAL SERVICE	NaN	NaN	
1	Radiology	CHEST (PORTABLE AP)	NaN	NaN	
2	Radiology	CT PELVIS W/CONTRAST	NaN	NaN	

	TEXT
0	[**Last Name (LF) **],[**First Name3 (LF) **] ...
1	[**2153-1-12**] 10:03 PM\n CHEST (PORTABLE AP)...
2	[**2145-3-13**] 2:58 AM\n CT ABDOMEN W/CONTRAS...

Steps to Load a CSV from Google Drive in Google Colab



Clean and Tokenize



Clean and Tokenize Texts

```
1 # cleans and tokenize
2 # import spacy
3 # import re
4 # import string
5 # from collections import Counter
6
7 # import pandas as pd
8 # import numpy as np
9
10 # from spacy.pipeline import EntityRuler
11 # from spacy.tokens import DocBin
12
13 nlp = spacy.load("en_core_web_sm") # Use medSpaCy model here if installed
14 stop_words = nlp.Defaults.stop_words
15
16 def clean_and_tokenize(text,
17     .....: lowercase=True,
18     .....: remove_stopwords=True,
19     .....: remove_punctuation=True,
20     .....: min_token_len=2):
21
22     .....: if lowercase:
23     .....:     text = text.lower()
24
25     .....: doc = nlp(text)
26
27     .....: tokens = []
28     .....: for token in doc:
29     .....:     if remove_stopwords and token.text in stop_words:
30     .....:         continue
31     .....:     if remove_punctuation and token.text in string.punctuation:
32     .....:         continue
33     .....:     if len(token.text) < min_token_len:
34     .....:         continue
35     .....:     if token.is_space or token.is_digit:
36     .....:         continue
37     .....:     tokens.append(token.lemma_.strip())
38
39     .....: return " ".join(tokens)

```

```
1 df_copy['cleaned_tokens'] = df_copy["TEXT"].apply(lambda text: clean_and_tokenize(text))

```

Library Imports (Lines 1–13)

Imports essential NLP and data processing libraries like `spacy`, `re`, `string`, and `Counter`. `pandas` and `numpy` are included for dataframe operations and numerical tasks. Some optional SpaCy components (`EntityRuler`, `DocBin`) are commented out but can be used for advanced customization.

Loading SpaCy Language Model (Line 13)

Loads the `en_core_web_sm` model for English NLP tasks.

You can replace this with a more medically tuned model like `medSpaCy` if available.

Stop Words Setup (Line 14)

Retrieves default English stop words (e.g., “the,” “is,” “in”) from SpaCy for later filtering.

Token Cleaning Function: `clean_and_tokenize` This custom function processes raw text through several steps

Lowercase (Line 22) – Converts all text to lowercase (if enabled).

SpaCy Tokenization (Line 24) Breaks text into SpaCy Token objects.

Filtering Loop (Lines 27–36) Removes stop words, punctuation, very short tokens (e.g., length < 2), digits, and spaces. Only keeps the **lemma** (base form) of each remaining token, with whitespace stripped.

Join Clean Tokens (Line 39) – Converts the cleaned list of tokens back to a single string.

Apply Cleaning to DataFrame Column (Line 41)

Applies the `clean_and_tokenize()` function to each row in the `TEXT` column of `df_copy`. Stores the output in a new column named 'cleaned_tokens'



Benefits of This Approach

Tokenization with SpaCy

This uses a powerful natural language processing engine that understands the structure of language. It breaks down text into tokens, identifies their lemmas (base forms), and can recognize parts of speech and named entities.

Lowercasing & Stopword Removal

Converting text to lowercase ensures consistency, and removing common stop words (like “the,” “is,” “and”) reduces noise in the data. This helps in improving the clarity of models and visualizations.

Lemmatization

Lemmatizing words reduces them to their root forms—for example, “running” becomes “run.” This helps in grouping similar words together and provides more meaningful insights during analysis.

Punctuation & Digit Removal

Removing punctuation and numbers strips out irrelevant characters that don’t contribute to the semantic meaning of the text.

Minimum Token Length Filtering

This step helps eliminate short, often meaningless fragments (like “a” or “x”) which are not useful in text analysis.

Result: Cleaned Text

The final output is a cleaned, structured version of the original text. This dramatically improves the quality and accuracy of downstream tasks such as keyword extraction, topic modeling, and word embeddings.

Named Entity Recognition (NER)

NER helps you quickly inspect in recognizing and labeling entities in your notes, and visually highlights them for manual review or further analysis



```
1 # from spacy import displacy
2
3 nlp = spacy.load("en_core_web_sm")
4 # Loop through first 3 processed notes and render entities
5 for i in range(10):
6     text = df_copy['cleaned_tokens'][i]
7     doc = nlp(text)
8     displacy.render(doc, style="ent", jupyter=True)
9     print("*****")
```

if name3 ORG If eu 8:24 pm TIME cta head w&w recon cta neck w&w oc recon clip clip number radiology ct brain perfusion -59 distinct procedural service reason eval ich cva ORG contrast optiray amt hospital medical condition year old woman new onset dysarthria know reason examination eval ich cva contraindication iv contrast pfi report evidence

acute infarct intracranial hemorrhage brain evidence focal flow limit stenosis occlusion aneurysm great mm artery anterior posterior circulation head calcify plaque left distal common carotid artery cause approximately stenosis

10:03 pm TIME chest portable ap clip clip number radiology reason scv cvl PERSON check line ptx hospital medical condition year old woman gangrene rt foot PERSON rt fem tib reason examination scv cvl check line ptx final report clinical indication central venous catheter placement comparison previous study day early interval placement left

subclavian NORP central venous catheter terminate junction left brachiocephalic vein PERSON superior vena cava pneumothorax identify cardiac silhouette remain enlarge bilateral pleural effusion adjacent area increase lung opacification low lobe note previous mediansternotomy coronary bypass surgery impression central venous catheter terminate

confluence left brachiocephalic vein superior vena cava pneumothorax persistent bilateral pleural effusion adjacent lung opacity likely reflect atelectasis infectious process exclude proper clinical setting

2:58 TIME ct abdomen contrast ct pelvis contrast clip clip number radiology ct coronal sagittal obl reconstruction ct 150cc nonionic contrast LOC reason evaluate diverticulitis PRODUCT contrast optiray amt hospital medical condition year old woman diverticuli ORG recent admit diverticulitis PRODUCT low gi bleed -- ct show thickening stranding

paracolic fluid left descending prox sigmoid colon worrisome diverticulitis vs ischemic bowel vs hemorrhage bowel wall present recurrent lq pain diarrhea evaluate diverticulitis perforation reason examination evaluate diverticulitis final report history diverticulitis technique helically acquire ct image lung basis pubic symphysis contrast ct abdomen contrast comparison

prior ct scan date demonstrate lung basis mm PERSON calcify granuloma left base lung clear pleural effusion liver spleen pancrea adrenal kidney unremarkable see large cm gallbladder stone gallbladder appear unremarkable bowel upper abdomen appear unremarkable free fluid see ct pelvis contrast mark interval improvement thickened sigmoid colon see prior

study numerous diverticula minimal stranding sigmoid colon likely clear residual disease free fluid see kidney demonstrate bilateral symmetrical contrast excretion ureter see course bladder filling defect bladder appear unremarkable free fluid present focal fluid collection identify bone window bone demonstrate suspicious lytic blastic lesion impression interval

improvement patient sigmoid disease mild stranding sigmoid colon likely relate residua prior disease rule early diverticulitis

11:46 CARDINAL ct head w/o contrast clip clip number radiology reason cerebral edema central pontine myelinolysis mass hospital medical condition year old man liver transplant w/ persistent decrease mental status reason examination cerebral edema central pontine myelinolysis mass final report history year old DATE man status post liver transplant

persistent decrease mental status rule cerebral edema central pontine myelinolysis technique contiguous axial image head obtain iv contrast skull base vertex ct head w/o contrast sulci ventricle appropriate patient age intra extraaxial fluid collection masse shift midline structure observe doctor -while matter differentiation unremarkable impression evidence cerebra

edema

1:06 pm TIME chest portable ap clip clip number radiology reason chf hospital medical condition year old woman niddm cad aicd pace acute sob leg swelling reason examination chf final report indication year old DATE woman diabete cad present follow icd placement acute shortness breath leg swelling rule chf comparison upright ap chest radiograph interval

placement left side dual chamber icd lead locate right atrium right ventricle persistent cardiomegaly slight decrease extent upper zone redistribution pulmonary vascularity increase size large right side pleural effusion associate right middle low lobe collapse consolidation linear opacity left cp angle likely represent subsegmental atelectasis interval decrease extent left

pleural effusion soft tissue osseous structure unchanged impression new left side dual chamber icd lead appropriate position pneumothorax persistent cardiomegaly GPE interval decrease chf interval increase large right pleural effusion small left pleural effusion leave basilar subsegmental atelectasis right middle low lobe collapse consolidation

2:21 pm TIME chest portable ap clip clip number radiology reason acute sob hospital medical condition year old man pvd ischemic ile cad PERSON heparinwithacut sob low bp reason examination acute sob final report indication shortness breath low blood pressure comparison previous study early day DATE left picc line remain place terminate distal

superior vena cava patient mediansternotomy coronary bypass surgery heart enlarge stable size previously note mild congestive heart failure pattern demonstrate significantly change allow difference technique confluent area consolidation see lung pleural effusion evident left cp sulcus exclude impression significant change recent study early day DATE persistent

mild congestive heart failure pattern

4:00 pm TIME chest pa lat clip clip number radiology reason follow cpr pt pneumomia rml hospital medical condition year old man sob chest leave le dvt reason examination follow cpr pt pneumomia rml final report indication right middle lobe pneumonia chest view consolidation collapse right middle lobe anterior segment right low lobe significantly change minimal

linear atelectasis left lung base small right pleural effusion present left costophrenic angle clear pulmonary vascularity normal cardiomedialastinal silhouette normal impression consolidation collapse right middle lobe anterior segment right low lobe unchanged delayed follow time course week recommend evaluate long term change

8:33 TIME ct chest w/o contrast clip clip number radiology reason new pulmonary nodule right upper lobe hospital medical condition year old ORG woman copd exacerbation incidental finding solitary pulmonary nodule cpr complain occasional tightness area chest pack year hx smoking reason examination new pulmonary nodule right upper lobe final report

indication year old DATE female chronic obstructive pulmonary disease exacerbation new nodule right upper lobe chest radiograph positive smoking history helical ct thorax perform intravenous oral contrast administration mm PERSON collimation mm PERSON reconstruction interval image lung reveal evidence right upper lobe pulmonary nodule apparent

nodule chest radiograph show healing right-sided rib fracture evidence emphysema PERSON severe upper lobe mild severity remain portion lung left low lobe streaky peripheral opacity predominantly linear ORG configuration associate ground glass opacification scatter linear opacity see lingula periphery right upper lobe probably reflex scar central airways

appear patent segmental level bilaterally minimal small alway disease posteriorly right low lobe leave low lobe manifest small branching y- shape structure review osseous structure thorax reveal healing subacute fracture right fifth ORDINAL posterior rib assessment soft tissue structure thorax reveal lymph node precarinal area ORG slightly cm short axis

Code walkthrough and Utility of NER

- This code uses spaCy and its visualization tool displaCy to process and visually inspect named entities in clinical or textual notes. Here's a step-by-step explanation:
- Model Loading:
`nlp = spacy.load("en_core_web_sm")`
Loads the small English language model, which includes a pre-trained named entity recognizer (NER).
- Iterate Through Notes:
The for loop goes through the first 10 entries (`i in range(10)`) in the DataFrame column `cleaned_tokens`, which is assumed to contain preprocessed text notes.
- Process Text:
For each note, `doc = nlp(text)` processes the text with the spaCy pipeline, including tokenization, tagging, and NER.
- Visualize Entities:
`displacy.render(doc, style="ent", jupyter=True)`
Uses displaCy to highlight and label any named entities detected in the text (such as PERSON, DATE, ORG, etc.) directly in a Jupyter notebook cell.
 - The `style="ent"` argument specifies that named entities should be visualized.
 - The `jupyter=True` flag ensures the visualization appears inline in Jupyter notebooks.
- Separator:
`print("*****")`
Prints a line of asterisks to visually separate the output for each note.

Benefits of NER (Named Entity Recognition):

- Extracts key info from unstructured text (e.g., names, dates, places).
- Saves time by automating data tagging and classification.
- Enables analysis like trend detection, summarization, and search.
- Improves decision-making with structured insights from raw text.
- Supports automation in healthcare, finance, legal, and more.



Most Common Diagnoses by Marital Status in MIMIC-III Admissions

```
[98] 1 # Top most frequent diagnoses (by short title) grouped by marital status.
2
3 query_early_critical = """
4 SELECT
5     CASE
6         WHEN adm.marital_status IS NULL
7             OR adm.marital_status IN ('UNKNOWN (DEFAULT)', 'None', '')
8             THEN 'Unknown'
9         ELSE adm.marital_status
10    END AS marital_status_cleaned,
11    dd.short_title AS diagnosis_name,
12    COUNT(*) AS diagnosis_count
13 FROM
14     `physionet-data.mimiciii_demo.admissions` AS adm
15 JOIN
16     `physionet-data.mimiciii_demo.diagnoses_icd` AS d
17   ON adm.hadm_id = d.hadm_id
18 JOIN
19     `physionet-data.mimiciii_demo.d_icd_diagnoses` AS dd
20   ON d.icd9_code = dd.icd9_code
21 WHERE
22     dd.short_title IS NOT NULL
23 GROUP BY
24     marital_status_cleaned, diagnosis_name
25 ORDER BY
26     marital_status_cleaned, diagnosis_count DESC
27 LIMIT 10000;
28 """
```

```
1 query_early_critical_df = run_bq_query(query_early_critical)
```



Query executed successfully! Data loaded into pandas DataFrame.
First 5 rows of the data:

	marital_status_cleaned	diagnosis_name	diagnosis_count
0	DIVORCED	Hypertension NOS	4
1	DIVORCED	Atrial fibrillation	3
2	DIVORCED	CHF NOS	2
3	DIVORCED	Hypothyroidism NOS	2
4	DIVORCED	Ac posthemorrhag anemia	2

This process identifies the most frequently occurring medical diagnoses across different marital status categories in the MIMIC-III dataset. It helps examine potential health trends linked to social support. We get a detailed profile of which illnesses are most common for married, divorced, widowed, single, and unknown-status patients - providing insight into how social factors might influence disease burden.

Standardize Marital Status: Marital statuses like "UNKNOWN", "None", or empty values are grouped under a single label, "Unknown", to ensure consistent categorization.

Gather Diagnoses: Diagnoses for each patient are extracted using short, human-readable labels that describe the condition (like: "Hypertension").

Join Relevant Tables: Combines data from admissions, diagnosis events, and diagnosis dictionary to map each hospital stay to a diagnosis and corresponding marital status.

Count Frequencies: Calculates how many times each diagnosis occurs within each marital status group.

Group and Rank: Organizes the data by marital status and ranks diagnoses from most to least frequent.

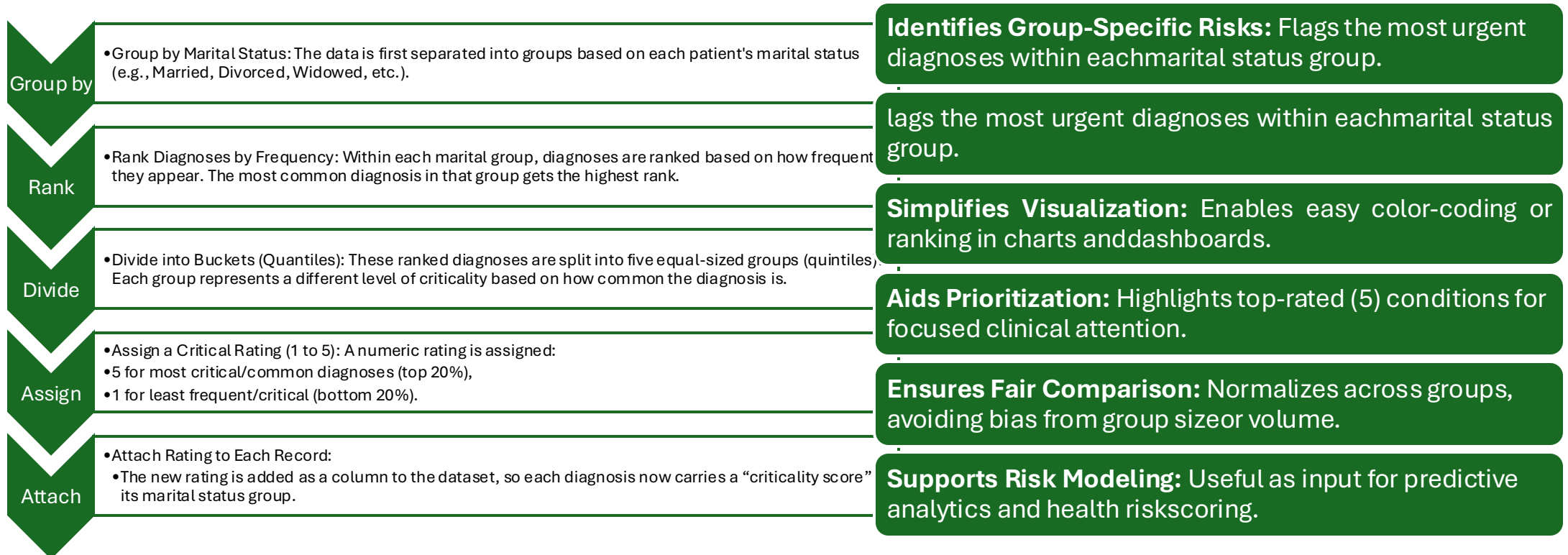
Filter Out Invalid Data: Any diagnosis without a proper title is excluded to maintain clean and interpretable results.

Limit Output:

Keeps the results manageable by returning only the top 10,000 records.

From Raw Counts to Risk Scores: Ranking Diagnoses Smarter

```
1 # Assign Critical Rating
2
3 # Assign critical rating based on frequency rank within each marital group
4 query_early_critical_df['critical_rating'] = query_early_critical_df.groupby('marital_status_cleaned')['diagnosis_count'] \
5     .transform(lambda x: pd.qcut(x.rank(method='first'), 5, labels=[1, 2, 3, 4, 5]))
6
7 # Preview
8 query_early_critical_df.head(1000)
9
```



Assign critical rating based on frequency rank within each marital group

1 of 1

to

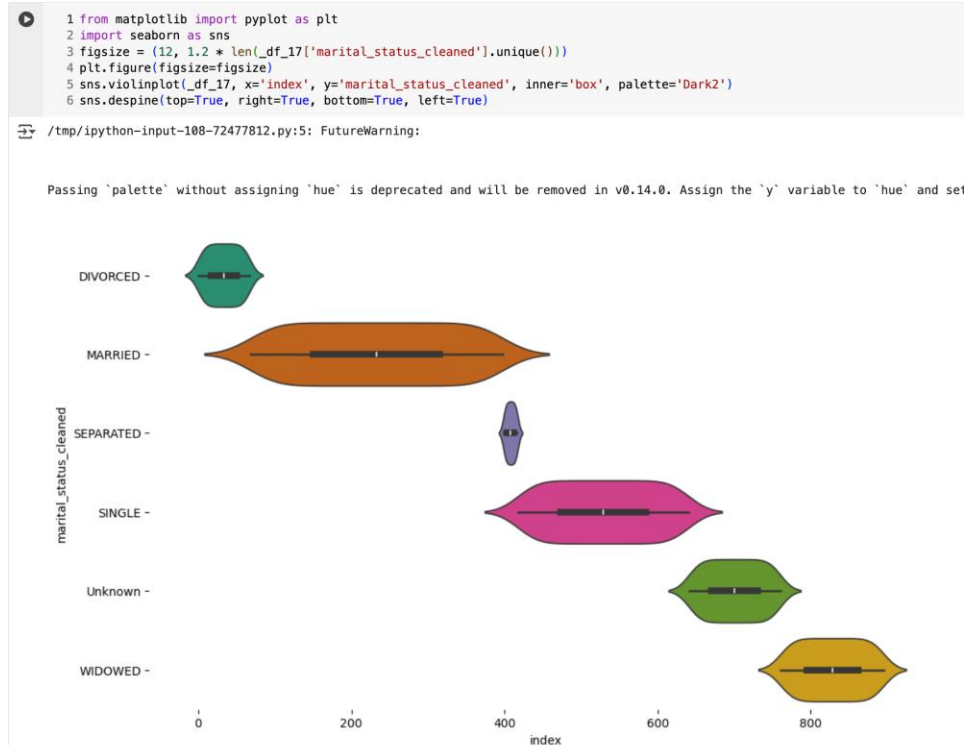
to

to

Search by all fields:

index	marital_status_cleaned	diagnosis_name	diagnosis_count	critical_rating
0	DIVORCED	Hypertension NOS	4	5
1	DIVORCED	Atrial fibrillation	3	5
2	DIVORCED	CHF NOS	2	5
3	DIVORCED	Hypothyroidism NOS	2	5
4	DIVORCED	Ac posthemorrhag anemia	2	5
5	DIVORCED	DMII w/o cmp nt st uncontr	2	5
6	DIVORCED	Delirium d/t other cond	2	5
7	DIVORCED	Severe sepsis	2	5
8	DIVORCED	Acute respiratory failure	2	5
9	DIVORCED	Cholangitis	1	1
10	DIVORCED	Shock w/o trauma NEC	1	1
11	DIVORCED	Emphysema NEC	1	1
12	DIVORCED	Anoxic brain damage	1	1
13	DIVORCED	Urin tract infection NOS	1	1
14	DIVORCED	Traum arthropathy-pelvis	1	1
15	DIVORCED	Hyperlipidemia NEC/NOS	1	1
16	DIVORCED	Alcoh dep NEC/NOS-unspec	1	1
17	DIVORCED	Septic shock	1	1
18	DIVORCED	Ac kidney fail, tubr neor	1	1
19	DIVORCED	Fx malar/maxillary-close	1	1
20	DIVORCED	Ac alcoholic hepatitis	1	1
21	DIVORCED	Adv eff opiates	1	1
22	DIVORCED	Hypopotassemia	1	1
23	DIVORCED	Idiopathic scoliosis	1	2
24	DIVORCED	Hx-prostatic malignancy	1	2
25	DIVORCED	Mitral valve disorder	1	2
26	DIVORCED	Divide colon w/o hmrhg	1	2
27	DIVORCED	Persistent vegtv state	1	2
28	DIVORCED	Paralysis agitans	1	2
29	DIVORCED	Acute pancreatitis	1	2
30	DIVORCED	Pneumonia, organism NOS	1	2
31	DIVORCED	Comp-ren dialys devlgrft	1	2
32	DIVORCED	Blood in stool	1	2
33	DIVORCED	Hx-circulatory dis NEC	1	2
34	DIVORCED	Septicemia NOS	1	2
35	DIVORCED	Hypotmolality	1	2
36	DIVORCED	Hyp kid NOS w or kid V	1	3
37	DIVORCED	Ac systolic hrt failure	1	3
38	DIVORCED	Ventricular fibrillation	1	3
39	DIVORCED	Subendo infarct, initial	1	3
40	DIVORCED	End stage renal disease	1	3
41	DIVORCED	Hematoma complic proc	1	3
42	DIVORCED	Cardiogenic shock	1	3
43	DIVORCED	Unc behav neo bone	1	3
44	DIVORCED	Cerv spondyl w myelopath	1	3
45	DIVORCED	Personal history of fall	1	3
46	DIVORCED	Status autm ord dfrtr	1	3
47	DIVORCED	Iatrogen pulm emb/infarc	1	3
48	DIVORCED	Other pulmonary insuff	1	3
49	DIVORCED	Iatrogenic hypotnsion NEC	1	3

Uncovering Hidden Health Patterns with Criticality Scores

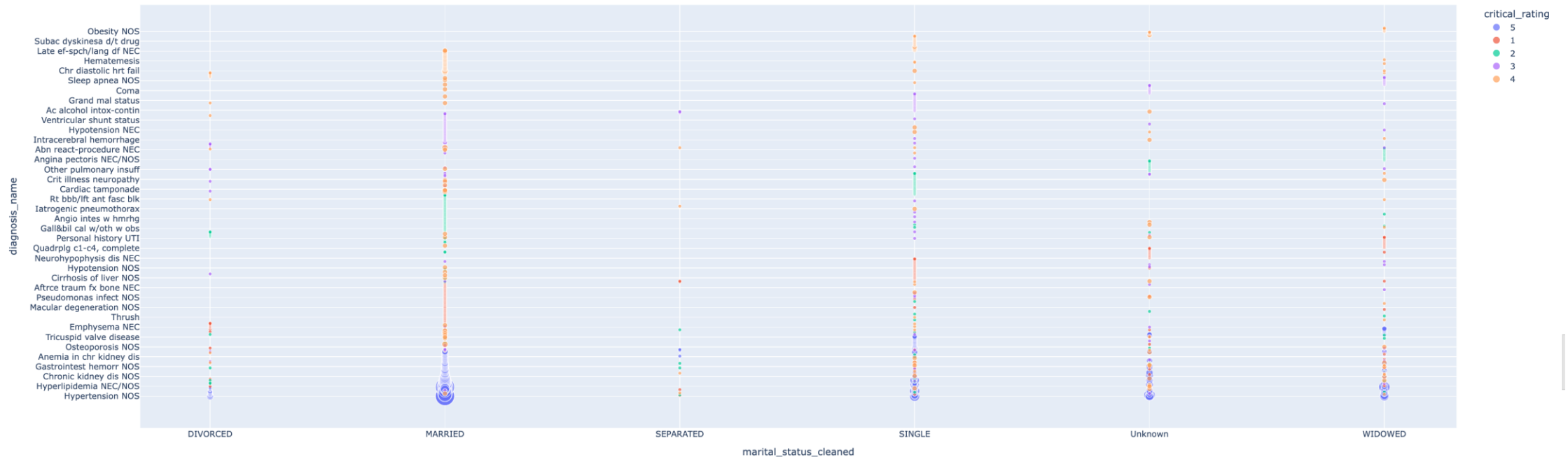


Intuitive, Interactive Insight into Critical Diagnoses Across Marital Status Groups

```
1 import plotly.express as px
2
3 fig = px.scatter(
4     query_early_critical_df,
5     x='marital_status_cleaned',
6     y='diagnosis_name',
7     size='diagnosis_count',
8     color='critical_rating',
9     color_continuous_scale='Reds',
10    title='Critical Diagnoses per Marital Status',
11    hover_data=['diagnosis_name', 'diagnosis_count', 'critical_rating']
12 )
13 fig.update_layout(height=800)
14 fig.show()
15
```



Critical Diagnoses per Marital Status





Benefits of This Visualization Approach

Highlights Group-Specific Health Risks: You can easily compare which diagnoses are most common and critical among different marital statuses.

Size and Color Encode Multiple Data Dimensions: Using both bubble size and color allows quick understanding of two variables (count and criticality) simultaneously without clutter.

Interactive Exploration: Hover details provide granular data without overwhelming the visual, making it suitable for presentations or clinical reviews.

Facilitates Prioritization: Identifying which diagnoses have the highest critical rating and frequency can help prioritize healthcare interventions or further analysis.

Clear Visual Segmentation: The categorical x-axis groups patients logically, making patterns and disparities easy to spot.

Customizable and Scalable: Plotly's interactive features allow zooming, panning, and easy modification, adapting well for dashboards or reports.

Leveraging NLP to Extract and Visualize Medical Entities from Diagnoses

para hyperglycemia

streptococcal septicemia

malfunc oth device graft

bipolar current PERSON nos

fx facial bone nec ORG close

fx surg nck ORG humerus clo

asthma nos

abn react surg proc nec ORG

polymyalgia GPE rheumatica

crnry athrscd PERSON natve vssl

esophageal reflux

pleural GPE effusion nos

choledochlith nos obst

atrial fibrillation

hypertension nos

acute kidney failure nos

hyperlipidemia nec ORG nos

severe anemia

hypothyroidism nos

ac posthemorrhag anemia ORG

dmii PERSON will cmp not st uncntd

delirium cond

severe sepsis

acute respiratory failure

cholangitis

shock w/o trauma nec ORG

emphysema PERSON nec ORG

anoxic brain damage

urin tract infection nos

traum arthropathy PERSON pelvis

hyperlipidemia nec ORG nos

alcoh dep PERSON nec ORG nos unspc

septic shock

ac ORG kidney ORG fail tubr necr

fx malar maxillary close

ac alcoholic ORG hepatitis

adv eff ORG opiate

hypopotassemia

hyposmolality

pressure ulcer low

neuropathy PERSON diabetes

depressive disorder nec ORG

pneumonia organism nos

dmii neuro nt st uncntrl PERSON

septic shock

acidosis

tracheostomy status

esophageal reflux

hyperosmolality

chronic kidney dis nos

pleural GPE effusion nos

food vomit pneumonitis

history tobacco use

secondary malign neo lung

obstructive sleep apnea

pneumonia organism nos

comp ren dialys dev grft PERSON

blood stool

hx circulatory dis nec ORG

septicemia nos

hyposmolality

hyp kid PERSON nos cr kid

ac systolic hrt ORG failure

ventricular fibrillation

subendo infarct initial ORG

end stage renal disease

hematoma complic GPE proc

cardiogenic shock

unc behav neo bone

cerv spondyl myelopath PERSON

personal history fall

status autm crd dfrlrr

Named Entity Recognition (NER) on Clinical Diagnoses Using spaCy

```
1
2 query_early_critical_df['cleaned_tokens'] = query_early_critical_df["diagnosis_name"].apply(lambda text: clean_and_tokenize(text))
3
4 nlp = spacy.load("en_core_web_sm")
5 # Loop through first 3 processed notes and render entities
6 for i in range(100):
7     text = query_early_critical_df['cleaned_tokens'][i]
8     doc = nlp(text)
9     displacy.render(doc, style="ent", jupyter=True)
10    # print("*****")
```

Code Walkthrough:

Text Cleaning and Tokenization:

The code applies a custom `clean_and_tokenize` function to the "diagnosis_name" column of the DataFrame, creating a new column "cleaned_tokens" containing the processed text.

Load spaCy Model:

Loads the pre-trained English language model (`en_core_web_sm`) to enable NLP tasks including Named Entity Recognition (NER).

Entity Extraction and Visualization:

Iterates over the first 100 cleaned diagnosis texts. For each, it:

- Processes the text with spaCy to generate a Doc object containing linguistic annotations.

- Uses `displacy.render` to visually highlight and label any named entities (such as diseases, drugs, or other biomedical terms) found in the text, directly in a Jupyter notebook cell for easy inspection.

Benefits of spaCy Modeling Approach

Transforms Unstructured Text into Structured Data:

Converts raw diagnosis text into tokenized and labeled components that computers can work with efficiently.

Highlights Clinically Relevant Information:

Automatically identifies and visualizes important medical entities, facilitating better understanding of diagnoses.

Supports Data Exploration and Quality Checks:

Visualizing entities helps spot inconsistencies or errors in the data early on.

Enables Downstream Analytics: Structured entities can be used in predictive models, clustering, or further natural language processing tasks.

Improves Interpretability: Makes complex clinical notes accessible for clinicians, researchers, or data scientists through clear visualization.

Scalable to Large Datasets: Automates entity recognition on many records, saving manual review time and effort.

The Unseen Role of Social Support: Marital Status and Healthcare Outcomes

While clinical factors are vital in healthcare outcomes, analysis of the MIMIC III dataset reveals that social determinants, especially marital status, play a crucial yet often overlooked role in patient health trajectories.

Research shows that married individuals tend to have better outcomes. This advantage is linked to stronger support systems that promote treatment adherence, provide emotional and practical help, and encourage earlier medical intervention. Such support can improve recovery rates and reduce mortality risk.

Our MIMIC III analysis highlights the vulnerability of single or isolated patients. Without immediate social support, these individuals may delay seeking care, face emotional challenges during health crises, and encounter higher risks of complications or death, particularly after major events like surgery.

Marital status acts as a proxy for broader social support networks in health research. When combined with age at death, it offers deeper insights. For example, a younger average age at death among single patients suggests an urgent need for targeted support and intervention. Conversely, older age at death may reflect natural disease progression, but social support remains a critical factor.

This study emphasizes that clinical data alone cannot capture the full picture. Ignoring social factors like marital status risks missing key influences on patient outcomes. Future research should integrate social determinants more systematically, improve data quality by addressing missing or null marital statuses, and adjust for variables like age, gender, illness severity, and comorbidities. Only then can we fully understand the complex interplay of clinical and social factors shaping health outcomes in datasets like MIMIC III.



<https://mermaid.live/>



[https://colab.research.google.com/drive/1l_tOX2wXK8LuZqx
mPWaWfBJmlc6wmZdE?usp=sharing](https://colab.research.google.com/drive/1l_tOX2wXK8LuZqxmPWaWfBJmlc6wmZdE?usp=sharing)



[SLT- Social support in HealthCare.pptx](#)