

Capstone Project Session 1

Derrick Mills

Student

Michigan Engineering Professional Education AIML – Cohort 3

Aura: Intelligent Prediction and Recommendation Engine for Digital Health Analytics

Introduction

The development of Aura, an intelligent prediction and recommendation engine for healthcare-focused digital marketing, requires a robust and reliable data pipeline capable of handling complex, heterogeneous datasets. The NSMES1988 dataset serves as the initial input for establishing the fundamental data ingestion, preprocessing, and optimization procedures that will support Aura's analytical, visualization, and machine-learning components.

The present report details the activities completed during Session 1, including data importation, structural evaluation, preliminary preprocessing, type refinement, memory footprint analysis, and exportation of optimized data artifacts. These steps establish a clean, validated, and computationally efficient foundation for subsequent modeling and system integration.

Methods:

Data Importation Environment

All tasks were conducted within Google Colab, using Python's scientific computing stack (pandas, NumPy). The dataset was imported from the file path /NSMES1988.csv. Colab provides a cloud-hosted execution environment that supports reproducible, scalable computational workflows suitable for healthcare analytics.

Initial Structural Inspection

Upon loading, the dataset was examined to characterize its dimensionality, variable composition, data types, and overall integrity. The inspection process involved:

- Evaluating dataset shape
- Reviewing variable data types
- Displaying summary statistics for numerical variables
- Assessing missingness patterns
- Identifying anomalous or non-standard values
- Characterizing categorical vs. continuous variables

Data Cleanliness Assessment

A complete analysis was performed using `df.isna().sum()`. Logical consistency checks were conducted to identify possible encoding abnormalities or out-of-range values. Special focus was given to socioeconomic variables such as income and clinical utilization variables that influence predictive modeling.

Type Optimization and Memory Reduction

Given Aura's requirement to scale to larger datasets in future phases, memory efficiency was prioritized. The following transformations were applied:

1. Removal of redundant variables

- Unnamed: 0, an auto-generated index column, was removed.

2. Categorical encoding

Variables representing discrete states (e.g., health, adl, region, gender, married, employment, insurance, medicaid) were cast to pandas category types.

3. Numeric downcasting

- Integer columns were downcast from int64 to smaller integer types (int16 or int8), where appropriate.
- Floating-point columns (age, income) were downcast from float64 to float32.

Memory usage was measured before and after optimization using `df.info(memory_usage="deep")`.

Data Exportation

To support downstream processing pipelines, two export formats were generated:

- The original dataset was exported as JSON (NSMES1988.json), enabling API-friendly ingestion.
- The cleaned, optimized dataset was exported as CSV (NSMES1988new.csv) for analytical and modeling tasks.

Results:

Dataset Structure and Composition

The dataset contains 4,406 observations and 19 variables, encompassing:

- **Healthcare utilization** - visits, non-physician visits, physician visits, emergency services, hospital admissions
- **Health status indicators** - chronic conditions, functional limitations (ADL), self-reported health

- **Demographic factors** - age (expressed in decades), gender, marital status
- **Socioeconomic indicators** - education, income (in units of \$10,000), employment status
- **Insurance variables** - private insurance, Medicaid coverage
- **Geographic region** - urban, rural, or regional classifications

The dataset exhibits no missing values, indicating complete records across all variables.

Observed Anomalies and Considerations

Two key observations emerged from the structural assessment:

1. **Age Encoding**
The age variable ranges from approximately 6.6 to 10.9, indicating that age is measured in decades rather than years. This transformation will require reconversion or careful interpretative handling in later modeling stages.
2. **Negative Income Values**
Income ranges from 1.0 to 54.8 (i.e., -\$10,000 to \$548,000). Negative values may reflect special coding, reporting errors, or debt and should be evaluated through business rules in subsequent preprocessing sessions.

Memory Utilization Before and After Optimization

- Original dataset memory footprint - 2.4 MB
- Optimized dataset memory footprint - 123 KB

This represents a 95% reduction in memory usage, achieved through categorical encoding and numeric downcasting. Although this dataset is relatively small, the optimization approach establishes a scalable pattern suitable for large-volume healthcare datasets expected in later stages of Aura's development.

Exported Artifacts

Two export files were successfully generated:

- NSMES1988.json - original dataset
- NSMES1988new.csv - cleaned and optimized dataset

These files are now available within the Colab working directory for downstream analysis.

Discussion

The preprocessing performed in Session 1 establishes a rigorous and efficient foundation for integrating the NSMES1988 dataset into Aura's analytical framework. The absence of missing

data simplifies early modeling tasks; the presence of encoded or transformed values (e.g., age in decades, negative income values) warrants domain-specific treatment in later sessions.

The dramatic memory footprint reduction demonstrates the importance of dtype optimization in scalable analytics architectures. As Aura evolves to accommodate larger, multi-source, and potentially real-time healthcare datasets, such optimizations will play a pivotal role in supporting low-latency processing, cost-efficient cloud workloads, and smooth visualization rendering.

Furthermore, the structured export of both original and cleaned datasets ensures compatibility with diverse components of the Aura pipeline, including RESTful ingestion endpoints, visualization engines, and machine-learning model trainers.

Conclusion

Session 1 accomplished the foundational tasks of data importation, structural assessment, data type optimization, memory analysis, and exportation of both original and cleaned datasets. The dataset is now prepared for more advanced phases of processing, including exploratory data analysis (Session 2), feature engineering, and predictive model development within the Aura engine.

The dataset's completeness, coupled with the systematic preprocessing conducted, provides a strong basis for subsequent analytical and modeling activities central to Aura's mission of delivering intelligent, data-driven insights for healthcare marketers.

Recommendations for Future Work

To prepare for predictive modeling and visualization in subsequent sessions, the following steps are recommended:

1. **Address anomalous income values** through imputation, capping, or transformation based on domain policy.
2. **Normalize key continuous variables**, particularly income and age (if converted back to years).
3. **Create a correlation matrix and univariate distributions** to identify predictive signal strength.
4. **Explore visit utilization clusters** (e.g., high-utilizers vs. low-utilizers) as potential segmentation insights.
5. **Establish metadata documentation** for variable semantics and business rules to support auditability and reproducibility.

File Edit View Insert Runtime Tools Help

Commands + Code + Text | Run all

```
import numpy as np
import pandas as pd

# *** CONFIG ***
CSV_PATH = "/NHEES1988.csv" # update if your file is elsewhere, e.g. "NHEES1988.csv"

def load_data(csv_path: str) -> pd.DataFrame:
    """Load the NHEES1988 dataset from CSV."""
    df = pd.read_csv(csv_path)
    return df

def basic_inspection(df: pd.DataFrame) -> None:
    """Print basic info and cleanliness checks."""
    print("\n*** BASIC INSPECTION ***")
    print(f"Shape (rows, columns): {df.shape}\n")
    print("Column dtypes:")
    print(df.dtypes, "\n")
    print("First 5 rows:")
    print(df.head(), "\n")
    print("Summary statistics (numeric):")
    print(df.describe(), "\n")
    print("Missing values per column:")
    print(df.isnull().sum(), "\n")

    if "age" in df.columns:
        print("Age stats:")
        print(df["age"].describe(), "\n")

    if "income" in df.columns:
        print("Income stats:")
        print(df["income"].describe(), "\n")

def memory_report(df: pd.DataFrame, title: str) -> None:
    """Print DataFrame.info with deep memory usage."""
    print(f"\n*** {title} MEMORY REPORT: {title} ***")
    df.info(memory_usage="deep")

def optimize_dtypes(df: pd.DataFrame) -> pd.DataFrame:
    """
    Optimize dtypes:
    - drop artificial index column if present
    - convert object columns to category
    - downcast integer and float columns
    """
    df_opt = df.copy()

    # Drop the 'Unnamed: 0' index column if present
    if "Unnamed: 0" in df_opt.columns:
```

File Edit View Insert Runtime Tools Help

Commands + Code + Text | Run all

```
    df_opt = df.copy()

    # Drop the 'Unnamed: 0' index column if present
    if "Unnamed: 0" in df_opt.columns:
        df_opt = df_opt.drop(columns=["Unnamed: 0"])

    # Categorical columns from the data dictionary
    cat_cols = [
        "health",
        "adl",
        "region",
        "gender",
        "marital",
        "employed",
        "insurance",
        "medicalid",
    ]

    for col in cat_cols:
        if col in df_opt.columns:
            df_opt[col] = df_opt[col].astype("category")

    # Downcast integer columns
    int_cols = df_opt.select_dtypes(include=["int64"]).columns
    df_opt[int_cols] = df_opt[int_cols].apply(
        pd.to_numeric, downcast="integer"
    )

    # Downcast float columns
    float_cols = df_opt.select_dtypes(include=["float64"]).columns
    df_opt[float_cols] = df_opt[float_cols].apply(
        pd.to_numeric, downcast="float"
    )

    return df_opt

def export_files(df_original: pd.DataFrame, df_clean: pd.DataFrame) -> None:
    """
    Export:
    - original data to JSON (NHEES1988.json)
    - cleaned / optimized data to CSV (NHEES1988new.csv)
    """
    print("\n*** EXPORTING FILES ***")

    # JSON export of original data frame
    json_path = "NHEES1988.json"
    df_original.to_json(json_path, orient="records")
    print(f"Exported Original Dataframe to JSON: {json_path}")

    # CSV export of cleaned / optimized data frame
    csv_new_path = "NHEES1988new.csv"
    df_clean.to_csv(csv_new_path, index=False)
    print(f"Exported Cleaned Dataframe to CSV: {csv_new_path}")

    # *** MAIN EXECUTION IN COLAB ***

```

The screenshot shows a Google Colab notebook titled "aura_session1.ipynb". The code cell contains the following Python script:

```
df_opt[col] = df_opt[col].astype("category")
# Downcast integer columns
int_cols = df_opt.select_dtypes(include=["int64"]).columns
df_opt[int_cols] = df_opt[int_cols].apply(
    pd.to_numeric, downcast="integer"
)

# Downcast float columns
float_cols = df_opt.select_dtypes(include=["float64"]).columns
df_opt[float_cols] = df_opt[float_cols].apply(
    pd.to_numeric, downcast="float"
)

return df_opt

def export_files(df_original: pd.DataFrame, df_clean: pd.DataFrame) -> None:
    """
    Export:
        - original data to JSON (NMES1988.json)
        - cleaned / optimized data to CSV (NMES1988new.csv)
    """
    print("==>>> EXPORTING FILES ==><==")

    # JSON export of original dataframe
    json_path = "NMES1988.json"
    df_original.to_json(json_path, orient="records")
    print(f"Exported original DataFrame to JSON: {json_path}")

    # CSV export of cleaned / optimized dataframe
    csv_new_path = "NMES1988new.csv"
    df_clean.to_csv(csv_new_path, index=False)
    print(f"Exported cleaned DataFrame to CSV: {csv_new_path}")

# === MAIN EXECUTION IN COLAB ===
df = load_data(CSV_PATH)

# 1. Basic Inspection & cleanliness check
basic_inspection(df)

# 2. Memory usage before optimization
memory_report(df, title="Original DataFrame")

# 3. Optimize dtype and clean structure
df_opt = optimize_dtypes(df)

# 4. Memory usage after optimization
memory_report(df_opt, title="Optimized DataFrame")

# 5. Export JSON and new CSV back to Colab workspace
export_files(df, df_opt)

print("\nNone. File now in the current working directory:")
!ls -lh
```

aura_session1.ipynb

```
File Edit View Insert Runtime Tools Help
Innbounds + Code + Text Run all
...
<ipython-input-1>
-- count 4406.000000 4406.000000 4406.000000 4406.000000 4406.000000 \
mean 0.263544 0.295968 1.541988 7.462480 16.209000
std 0.155559 0.173232 0.200000 0.287756
min 0.000000 0.000000 0.000000 0.000000 0.000000
25% 0.000000 0.000000 1.000000 6.900000 8.000000
50% 0.000000 0.000000 1.000000 7.300000 11.000000
75% 0.000000 0.000000 2.000000 12.000000 12.000000
max 12.000000 8.000000 8.000000 10.000000 15.000000

Income
count 4406.000000
mean 2.527132
std 2.924648
min -1.012500
25% 0.912159
50% 1.698158
75% 3.172858
max 54.853100

Missing values per column:
Unamed: 0 0
visits 0
nvisits 0
ovisits 0
novidits 0
nemergency 0
hospital 0
health 0
chronic 0
adl 0
region 0
age 0
gender 0
married 0
school 0
income 0
employed 0
insurance 0
medicaid 0
dtype: int64

Age stats:
count 4406.000000
mean 2.402480
std 2.033485
min 6.600000
25% 6.900000
50% 7.200000
75% 7.800000
max 10.000000
Name: age, dtype: float64

Income stats:
count 4406.000000
mean 2.527132
std 2.924648
min -1.012500
...
<ipython-input-2>
How can I install Python libraries? Load data from Google Drive Show an example of training:
[Code] [Text]
```

aura_session1.ipynb

```
File Edit View Insert Runtime Tools Help
Innbounds + Code + Text Run all
...
<ipython-input-1>
-- Income stats:
count 4406.000000
mean 2.527132
std 2.924648
min -1.012500
25% 0.912159
50% 1.698158
75% 3.172858
max 54.853100
Name: income, dtype: float64

*** MEMORY REPORT: Original DataFrame ***
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4406 entries, 0 to 4405
Data columns (total 19 columns):
 # Column Non-Null Count Dtype 
 --- ... -----
 0 emergency: 4406 non-null int64
 1 visits: 4406 non-null int64
 2 nvisits: 4406 non-null int64
 3 ovisits: 4406 non-null int64
 4 novidits: 4406 non-null int64
 5 emergency: 4406 non-null int64
 6 hospital: 4406 non-null int64
 7 health: 4406 non-null object 
 8 employed: 4406 non-null int64
 9 adl: 4406 non-null object 
 10 region: 4406 non-null object 
 11 age: 4406 non-null float64
 12 gender: 4406 non-null object 
 13 married: 4406 non-null object 
 14 school: 4406 non-null int64
 15 income: 4406 non-null float64
 16 employed: 4406 non-null object 
 17 insurance: 4406 non-null object 
 18 medicaid: 4406 non-null object 
 19 gender_1: 4406 non-null category
dtypes: float64(2), int64(9), object(8)
memory usage: 2.2 M

*** MEMORY REPORT: Optimized DataFrame ***
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4406 entries, 0 to 4405
Data columns (total 18 columns):
 # Column Non-Null Count Dtype 
 --- ... -----
 0 visits: 4406 non-null int8
 1 nvisits: 4406 non-null int8
 2 ovisits: 4406 non-null int16
 3 novidits: 4406 non-null int8
 4 emergency: 4406 non-null int8
 5 hospital: 4406 non-null int 
 6 health: 4406 non-null category
 7 chronic: 4406 non-null int 
 8 adl: 4406 non-null category
 9 region: 4406 non-null category
 10 age: 4406 non-null float32
 11 gender: 4406 non-null category
 12 gender_1: 4406 non-null category
 13 gender_2: 4406 non-null category
 14 gender_3: 4406 non-null category
 15 gender_4: 4406 non-null category
 16 gender_5: 4406 non-null category
 17 gender_6: 4406 non-null category
 18 gender_7: 4406 non-null category
 19 gender_8: 4406 non-null category
dtypes: int8(10), int16(1), int32(1), float32(1)
memory usage: 1.2 M
...
<ipython-input-2>
How can I install Python libraries? Load data from Google Drive Show an example of training:
[Code] [Text]
```

aura_session1.ipynb

File Edit View Insert Runtime Tools Help

commands + Code + Text Run all

```
1 visits    4400 non-null int64
2 ovvisits   4400 non-null int64
3 novists    4400 non-null int64
4 novists    4400 non-null int64
5 emergency  4400 non-null int64
6 hospital   4400 non-null int64
7 health     4400 non-null object
8 chronic    4400 non-null int8
9 adl       4400 non-null object
10 region    4400 non-null object
11 age      4400 non-null float64
12 gender    4400 non-null object
13 married   4400 non-null object
14 school   4400 non-null int64
15 income    4400 non-null float64
16 employed  4400 non-null object
17 insurance 4400 non-null object
18 medicaid  4400 non-null object
dtypes: float64(2), int64(9), object(8)
memory usage: 2.2 MB

*** MEMORY REPORT: Optimized DataFrame ***
<class 'pandas.core.frame.DataFrame'>
Rows: 4400 Columns: 18
Data columns (total 18 columns):
 #   Column   Non-Null Count  Dtype  
--- 
 0   visits    4400 non-null  int64  
 1   ovvisits   4400 non-null  int64  
 2   novists    4400 non-null  int64  
 3   novists    4400 non-null  int64  
 4   emergency  4400 non-null  int64  
 5   hospital   4400 non-null  int64  
 6   health     4400 non-null  category
 7   chronic    4400 non-null  int8   
 8   adl        4400 non-null  object  
 9   region    4400 non-null  category
 10  age        4400 non-null  float32 
 11  gender    4400 non-null  category
 12  married   4400 non-null  category
 13  gender    4400 non-null  int8   
 14  income    4400 non-null  float32 
 15  employed  4400 non-null  category
 16  insurance 4400 non-null  category
 17  medicaid  4400 non-null  category
dtypes: category(8), float32(2), int64(2), int8(6)
memory usage: 113.0 KB

*** EXPORTING FILES ***
Exported original DataFrame to JSON: NSMES1988.json
Exported cleaned DataFrame to CSV: NSMES1988new.csv

Done. Files now in the current working directory:
total 1.5M
drwx----- 5 root root  4.0K Nov 29 16:41 drive
-rw-r--r--  1 root root  1.2M Nov 29 16:46 NSMES1988.json
-rw-r--r--  1 root root 299K Nov 29 16:46 NSMES1988new.csv
drwxr-xr-x 1 root root  4.0K Nov 28 14:38 sample_data
```