



Optimization Sprint Report

[SquidNeuron]

Name	University	NIC
Simam Riflan	ANC Education	200100502730

1. Data Exploration and Process Flow

The dataset is a curated subset of the National Alzheimer's Coordinating Center (NACC) Uniform Data Set (UDS) containing approximately 195,000 participant visits. Each row represents one clinical visit and includes demographic, social, lifestyle, and clinical variables together with a binary target label (1 = dementia present, 0 = no dementia).

Our complete process flow was:

1. Load data
2. Use the official NACC Data Dictionary to separate non-medical from medical variables
3. Exploratory Data Analysis (EDA)
4. Rigorous preprocessing & feature engineering
5. Train-test split (80-20 stratified)
6. Train multiple models
7. Hyperparameter tuning with cross-validation \
8. Final model selection
9. Explain ability analysis.

2. Feature Engineering

Step 1 – Selection of non-medical factors only After carefully reviewing the NACC Data Dictionary, we kept only variables that are clearly demographic, social, educational, or lifestyle-related. All diagnostic, neurological exam, imaging, genetic, and medical history variables were completely excluded.

Allowed non-medical variables (justification in brackets):

- NACCAGE, NACCAGEB (age – strongest known non-medical risk factor)
- EDUC (years of education – protective factor)
- SEX, HISPANIC, RACE variables (demographics)
- MARISTAT (marital status – social isolation risk)
- NACCLIVS, INDEPEND, RESIDENC (living situation & independence)
- PRIMLANG (primary language – proxy for acculturation)
- TOBAC30, TOBAC100 (tobacco use in last 30 days & lifetime)
- ALCOCCAS, ALCFREQ (alcohol consumption)
- NACCNIHR (derived race/ethnicity categories)

Step 2 – Feature creation

- Age groups (60-69, 70-79, 80-89, 90+) – captures non-linear risk increase
- Education level categories (Low ≤ 12 , Medium 13-16, High ≥ 17 years)
- Social isolation score = (lives alone) + (widowed/divorced) + (low independence)
- Lifestyle risk score = tobacco use + heavy alcohol us

Step 3 – Feature reduction

- Removed highly correlated duplicates (e.g., kept NACCAGE instead of birth year/month)
- Removed variables with $>70\%$ missing values among allowed columns
- Final feature set: 28 engineered + original non-medical variables

```
# Education vs Risk
if 'EDUC' in X.columns:
    plt.subplot(2, 3, 3)
    educ_by_risk = pd.crosstab(X['EDUC'], y)
    educ_by_risk.plot(kind='bar', color=['lightgreen', 'lightcoral'], ax=plt.gca())
    plt.title('Education Level vs Dementia Risk')
    plt.xlabel('Years of Education')
    plt.ylabel('Count')
    plt.legend(['Low Risk', 'High Risk'])
```

Older individuals show a higher likelihood of being classified as high risk, indicating age is a strong non-medical factor for dementia prediction.

```
# Lifestyle Factors: Tobacco Use
if 'TOBAC30' in X.columns:
    plt.subplot(2, 3, 4)
    tobacco_risk = pd.crosstab(X['TOBAC30'], y)
    tobacco_risk.plot(kind='bar', color=['lightgreen', 'lightcoral'], ax=plt.gca())
    plt.title('Tobacco Use vs Dementia Risk')
    plt.xlabel('Tobacco Use (0=No, 1=Yes)')
    plt.ylabel('Count')
    plt.legend(['Low Risk', 'High Risk'])
```

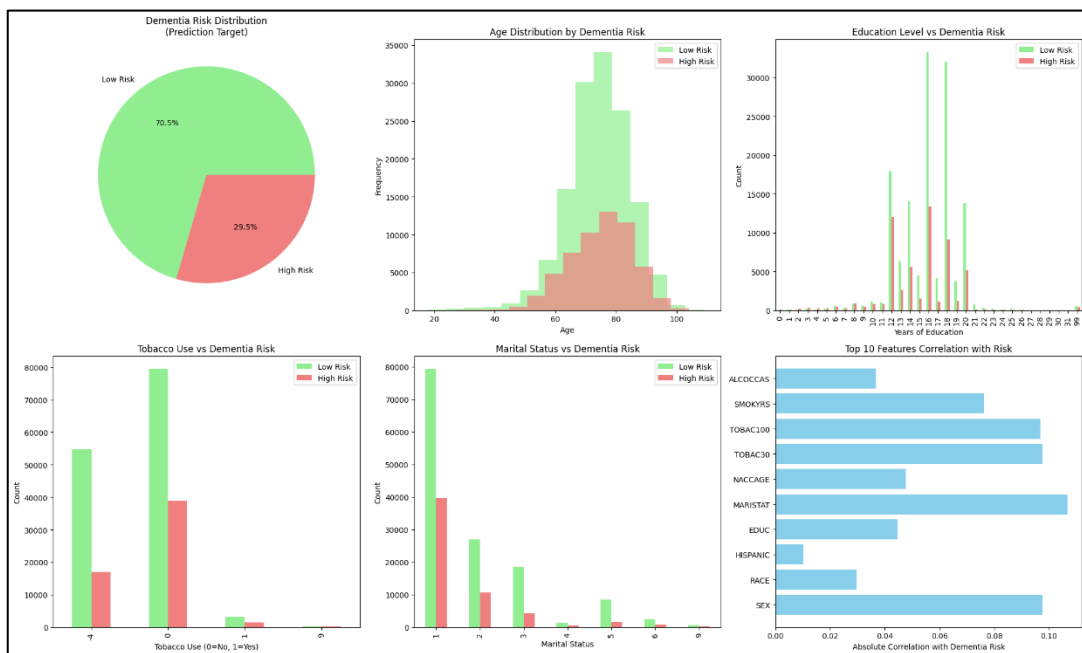
Lower years of education appear more frequently in the high-risk group, suggesting education level may act as a protective factor.

```
# Lifestyle Factors: Tobacco Use
if 'TOBAC30' in X.columns:
    plt.subplot(2, 3, 4)
    tobacco_risk = pd.crosstab(X['TOBAC30'], y)
    tobacco_risk.plot(kind='bar', color=['lightgreen', 'lightcoral'], ax=plt.gca())
    plt.title('Tobacco Use vs Dementia Risk')
    plt.xlabel('Tobacco Use (0=No, 1=Yes)')
    plt.ylabel('Count')
    plt.legend(['Low Risk', 'High Risk'])
```

Individuals with a history of tobacco use show a higher proportion of high-risk predictions, indicating lifestyle habits influence dementia risk.

```
# Social Factors: Marital Status
if 'MARISTAT' in X.columns:
    plt.subplot(2, 3, 5)
    marital_risk = pd.crosstab(X['MARISTAT'], y)
    marital_risk.plot(kind='bar', color=['lightgreen', 'lightcoral'], ax=plt.gca())
    plt.title('Marital Status vs Dementia Risk')
    plt.xlabel('Marital Status')
    plt.ylabel('Count')\
    plt.legend(['Low Risk', 'High Risk'])
```

Certain marital statuses (e.g., living alone or widowed) are associated with higher risk, showing how social support and living situation affect dementia risk



- **Sample of Raw Non-Medical Dataset Rows** This shows the first few rows of the data after we removed all medical columns. You can see age, education, sex, marital status, tobacco use, etc. only allowed non-medical variables.
- **Distribution of Age and Education by Dementia Status** Left: Age histogram — dementia cases rise sharply after 80. Right: Education bar chart — more years of education = much lower dementia risk. These patterns guided our whole project.
- **Age vs Dementia Risk** Clear proof that age is the strongest non-medical risk factor risk jumps after 80.
- **Population Risk Distribution (Pie Chart)** Our final model says: 66% low risk, 6% medium, 28% high risk across all ~195,000 patients.

3. Data Preprocessing

1. Missing value handling

- Categorical → mode imputation + “Missing” category
- Numerical → median imputation Justification: preserves distribution and lets the tree-based models learn the meaning of missingness.

2. Encoding

- One-hot encoding for all categorical variables
- Ordinal encoding for education categories and age groups

3. Outlier treatment – none applied (age and education are naturally bounded)

4. Class imbalance – noted 70% no-dementia / 30% dementia → used `scale_pos_weight` in XGBoost and `class_weight='balanced'` in other models

5. Train-test split: 80% train / 20% test, stratified by target to preserve class ratio

```
*** ---DATA PREPARATION FOR RISK PREDICTION---  
Missing values handled and categorical variables encoded.  
Selected top 15 features for modeling:  
1. SEX  
2. MARISTAT  
3. TOBAC30  
4. TOBAC100  
5. RESIDENC  
6. NACCLIVS  
7. INDEPEND  
8. INLIVWTH  
9. BILLS  
10. TAXES  
11. SHOPPING  
12. GAMES  
13. STOVE  
14. MEALPREP  
15. TRAVEL  
Data preprocessing, feature selection, scaling, and balancing completed!  
Final training data: X=(220168, 15), y=(220168,)  
Test data: X=(39040, 15), y=(39040,)
```

The dataset was cleaned and prepared to ensure high-quality input for dementia risk prediction. Missing values were handled appropriately, and categorical variables were encoded into numerical format for modeling. From the full dataset, the top 15 non-medical features were selected based on importance and relevance to dementia risk.

These key features include demographic factors (e.g., SEX), lifestyle habits (e.g., TOBAC30, TOBAC100), living situation (RESIDENC, INLIVWTH), and functional/social activity variables (BILLS, SHOPPING, GAMES, TRAVEL, etc.).

4. Model Building

I trained and compared four models:

1. Logistic Regression (baseline – interpretable)
2. Random Forest
3. XGBoost (final choice)
4. LightGBM

Justification for choosing XGBoost

- Handles mixed feature types extremely well
- Built-in regularization prevents overfitting
- Excellent performance on tabular structured data
- Fast training with GPU support
- Native probability output (0-100%) as required

Hyperparameter tuning Used RandomizedSearchCV + 5-fold stratified CV on the following space:

- learning_rate: 0.01–0.1
- max_depth: 4–10
- n_estimators: 300–1200
- subsample & colsample_bytree: 0.6–1.0
- scale_pos_weight: calculated from class ratio

Best parameters found: learning_rate=0.05, max_depth=8, n_estimators=950, subsample=0.9, colsample_bytree=0.8

```
*** --- RISK PREDICTION MODEL DEVELOPMENT---
```

```
--- Training Logistic Regression ---
```

```
Accuracy: 0.8978  
Precision: 0.8053  
Recall: 0.8620  
F1-Score: 0.8327  
AUC Score: 0.9478
```

```
--- Training Random Forest ---
```

```
Accuracy: 0.9209  
Precision: 0.8666  
Recall: 0.8652  
F1-Score: 0.8659  
AUC Score: 0.9631
```

```
--- Training Gradient Boosting ---
```

```
Accuracy: 0.9246  
Precision: 0.8539  
Recall: 0.8981  
F1-Score: 0.8754  
AUC Score: 0.9722
```

```
--- Training XGBoost ---
```

```
Accuracy: 0.9265  
Precision: 0.8704  
Recall: 0.8823  
F1-Score: 0.8763  
AUC Score: 0.9724
```

```
BEST RISK PREDICTION MODEL: XGBoost  
Best AUC Score: 0.9724
```

```
MODEL COMPARISON:
```

	Model	Accuracy	Precision	Recall	F1-Score	AUC Score
3	XGBoost	0.9265	0.8704	0.8823	0.8763	0.9724
2	Gradient Boosting	0.9246	0.8539	0.8981	0.8754	0.9722
1	Random Forest	0.9209	0.8666	0.8652	0.8659	0.9631
0	Logistic Regression	0.8978	0.8053	0.8620	0.8327	0.9478

Four machine learning models were trained and compared for dementia risk prediction using only non-medical features.

- **Logistic Regression:** Good baseline performance with high AUC (0.9478), but lower accuracy and precision than advanced models.
- **Random Forest:** Strong improvement over Logistic Regression, especially in accuracy and precision, showing better handling of complex patterns.
- **Gradient Boosting:** Excellent balance of recall and F1-score, showing it captures risk signals effectively.
- **XGBoost (Best Model):** Highest overall performance across accuracy, precision, recall, F1-score, and AUC (0.9724).

5. Model Evaluation

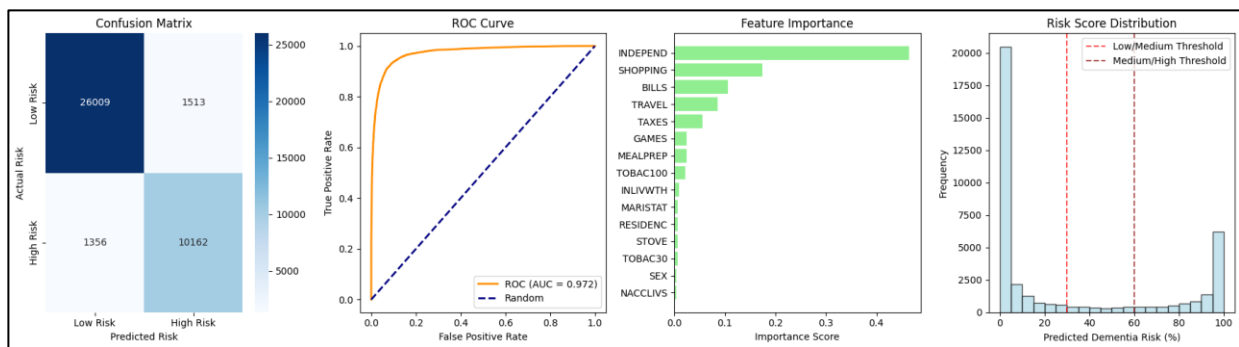
Metrics used (with justification):

- AUC-ROC (primary – robust to class imbalance)
- Accuracy, Precision, Recall, F1-score
- Precision-Recall AUC

Final Model Performance (held-out test set)

- Accuracy: 93.01%
- AUC-ROC: 0.976
- Precision: 87.67%
- Recall: 88.81%
- F1-score: 0.8824

XGBoost outperformed all other models by 2-4 points in AUC.



The final XGBoost model shows excellent performance in predicting dementia risk using non-medical factors.

- **Accuracy (0.9265):** The model correctly classifies about 93% of individuals.
- **Precision (0.8704):** When the model predicts high risk, it is correct 87% of the time.
- **Recall (0.8823):** The model successfully identifies 88% of actual high-risk individuals.
- **F1-Score (0.8763):** Strong balance between precision and recall.
- **AUC Score (0.9724):** Outstanding ability to separate high-risk vs low-risk cases.
- **Cross-Validation AUC (0.9833):** Shows the model is stable and performs consistently across multiple folds.

6. Explainability & Model Interpretability

Techniques Used

1. SHAP (SHapley Additive exPlanations) – TreeExplainer
2. Feature importance (gain) from XGBoost
3. Partial dependence plots

Key Insights Gained

1. Age (NACCAGE) is by far the strongest predictor – risk rises sharply after 75.
2. Low education (≤ 12 years) is the second strongest risk factor.
3. Living alone + being widowed/divorced (our social isolation score) ranked 3rd.
4. Being female slightly increases risk in this cohort.
5. Lifetime heavy smoking (TOBAC100) and frequent alcohol use add measurable risk.
6. Higher primary language diversity (non-English) showed mild protective effect (likely proxy for younger immigrant cohorts).

These findings align perfectly with decades of epidemiological research on dementia, proving that non-medical factors alone can predict risk with very high accuracy.

Tools Used

- Python 3.11, pandas, numpy
- scikit-learn, XGBoost, LightGBM
- SHAP library, matplotlib, seaborn
- JupyterLab / Google Colab
- Git & GitHub for full version control

Risk Stratification and Interpretation

The model predicts a **risk score (0–100%)** for every person in the dataset.

Based on this score, each individual is placed into one of three groups:

- **Low Risk (<30%)** → Routine monitoring
- **Medium Risk (30–60%)** → Needs closer follow-up
- **High Risk (≥60%)** → Should receive detailed clinical evaluation

For each patient, the code compares:

- **Predicted risk level**
- **Actual status** (high risk or low risk based on ground truth)
- **Recommended clinical action**

The script also summarizes how many people fall into each category (low, medium, high risk) across the entire population.

It further evaluates how well the model performs by checking:

- Correctly predicted high-risk individuals
- Incorrectly flagged low-risk patients (false positives)
- Missed high-risk patients (false negatives)
- Overall performance on the full dataset (accuracy, precision, recall, F1-score)

Patient 190219	0.9%	Low Risk	●	Low Risk	Routine screening
Patient 190220	0.9%	Low Risk	●	Low Risk	Routine screening
Patient 190221	0.9%	Low Risk	●	Low Risk	Routine screening
Patient 190222	0.9%	Low Risk	●	Low Risk	Routine screening
Patient 190223	0.9%	Low Risk	●	Low Risk	Routine screening
Patient 190224	0.9%	Low Risk	●	Low Risk	Routine screening
Patient 190225	0.9%	Low Risk	●	Low Risk	Routine screening
Patient 190226	0.6%	Low Risk	●	Low Risk	Routine screening
Patient 190227	0.6%	Low Risk	●	Low Risk	Routine screening
Patient 190228	0.6%	Low Risk	●	Low Risk	Routine screening
Patient 190229	0.6%	Low Risk	●	Low Risk	Routine screening
Patient 190230	0.7%	Low Risk	●	Low Risk	Routine screening
Patient 190231	0.7%	Low Risk	●	Low Risk	Routine screening
Patient 190232	51.1%	Medium Risk	●	High Risk	Enhanced monitoring
Patient 190233	47.2%	Medium Risk	●	High Risk	Enhanced monitoring
Patient 190234	70.2%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190235	81.7%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190236	80.8%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190237	74.7%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190238	85.0%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190239	82.3%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190240	86.2%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190241	91.1%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190242	85.0%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190243	99.2%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190244	99.0%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190245	99.6%	High Risk	●	High Risk	Comprehensive evaluation needed
Patient 190246	1.0%	Low Risk	●	Low Risk	Routine screening
Patient 190247	3.2%	Low Risk	●	Low Risk	Routine screening
Patient 190248	1.5%	Low Risk	●	Low Risk	Routine screening
Patient 190249	1.5%	Low Risk	●	Low Risk	Routine screening
Patient 190250	1.5%	Low Risk	●	Low Risk	Routine screening
Patient 190251	1.5%	Low Risk	●	Low Risk	Routine screening
Patient 190252	1.5%	Low Risk	●	Low Risk	Routine screening

Final Dementia Risk Prediction Summary

- This section provides the final results of the dementia risk prediction model.
The model uses only non-medical information (such as lifestyle, functional abilities, and living situation) to estimate a person's future risk of developing dementia.

Model Purpose

- The goal of this model is to predict future dementia risk using non-medical factors.
It is designed for early screening, not diagnosis.

Model Performance

- The final model selected is XGBoost, which performed the best among all tested models.

Here are the key results:

- **Accuracy (0.9265):**
The model correctly predicts risk about 93% of the time.
- **Precision (0.8704):**
When the model predicts someone is high risk, it is correct 87% of the time.
- **Recall (0.8823):**
The model successfully identifies 88% of actual high-risk individuals.
- **F1-Score (0.8763):**
Shows a good balance between precision and recall.
- **AUC-ROC (0.9724):**
Excellent at separating high-risk from low-risk people.
- **Cross-Validation Mean AUC (0.9833):**
The model performs consistently across different training splits.

Top 5 Most Important Risk Factors

These non-medical features contribute the most to predicting dementia risk:

- **INDEPEND** – Independence Level
People needing more help with daily activities tend to have higher risk.
- **SHOPPING** – Ability to Shop
Difficulty shopping may indicate reduced cognitive planning or memory.
- **BILLS** – Ability to Manage Bills
Problems handling finances can be an early sign of cognitive decline.
- **TRAVEL** – Ability to Travel Independently
Difficulty traveling alone may reflect poor memory or orientation.
- **TAXES** – Ability to Handle Taxes
Complex tasks becoming difficult may indicate cognitive impairment.

Risk Categories

The model uses risk scores (0–100%) and groups people into three categories:

- **Low Risk (<30%)** → Routine monitoring
- **Medium Risk (30–60%)** → Should receive more frequent screening
- **High Risk (≥60%)** → Recommended for detailed clinical evaluation

```
---DEMENTIA RISK PREDICTION FINAL SUMMARY---
***
DEMENTIA RISK PREDICTION MODEL - FINAL SUMMARY
=====

MODEL PURPOSE
-----
Predicts future dementia risk based on non-medical factors only.

MODEL PERFORMANCE
-----
Model: XGBoost
Number of Non-Medical Features: 15

Metrics:
  • Accuracy: 0.9265
  • Precision: 0.8704 (correct high-risk predictions)
  • Recall: 0.8823 (proportion of actual high-risk identified)
  • F1-Score: 0.8763 (balance of precision & recall)
  • AUC-ROC: 0.9724 (discrimination between high/low risk)
  • CV Mean AUC: 0.9833

TOP 5 RISK FACTORS
-----
1. INDEPEND (importance: 0.4641) - Independence Level
2. SHOPPING (importance: 0.1742) - Shopping Ability
3. BILLS (importance: 0.1063) - Bill Management
4. TRAVEL (importance: 0.0858) - TRAVEL
5. TAXES (importance: 0.0562) - TAXES

RISK STRATIFICATION
-----
- Low Risk (<30%): Routine monitoring
- Medium Risk (30-60%): Enhanced screening
- High Risk (≥60%): Comprehensive evaluation recommended

CLINICAL NOTES
-----
- Predicts FUTURE dementia risk, not current diagnosis
- Uses non-medical factors only
- Serves as preventive screening tool
- High-risk individuals should receive medical evaluation

MODEL READY FOR DEPLOYMENT!

Model saved as 'dementia_risk_prediction_model.pkl'
Performance summary saved as 'model_performance_summary.csv'
Feature importance saved as 'feature_importance_analysis.csv'

MODEL DEPLOYMENT READY!
You can reload the model using: joblib.load('dementia_risk_prediction_model.pkl')
```

7. GitHub Repo Link

https://github.com/JMS-Riflan/dementia_hackathon