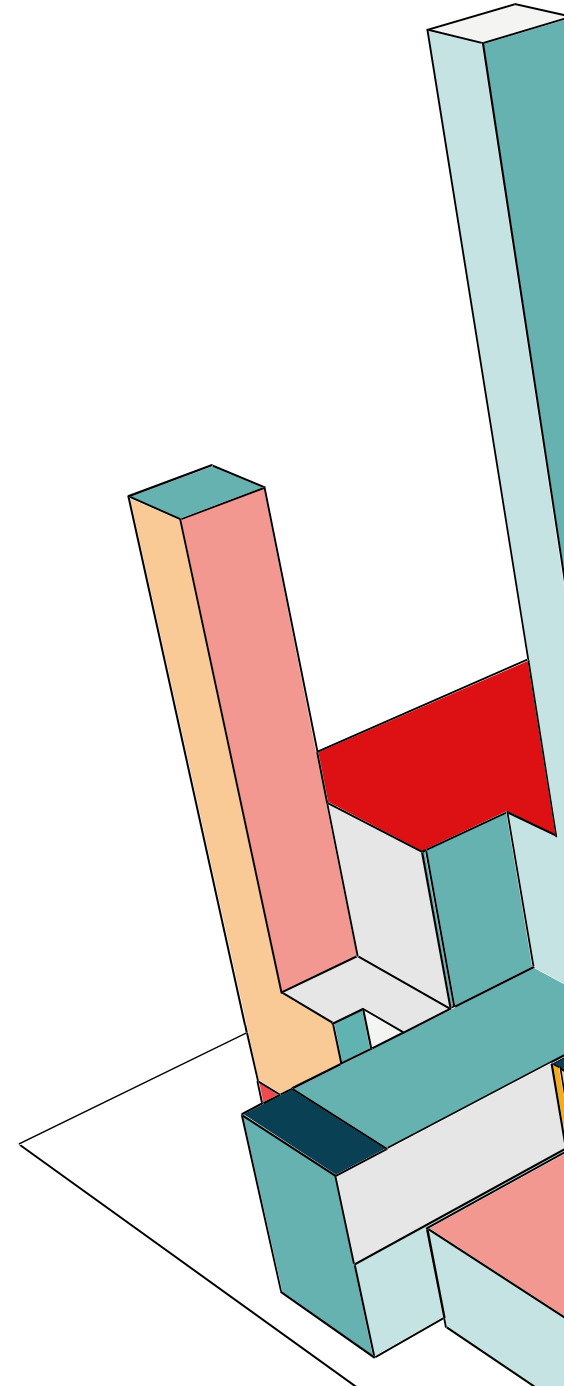# RARE DISEASE DETECTION FROM STANDARD LAB TESTING

Maina Musa

University of Texas at Austin

April 2025

# SCHEMA

- Introduction
- Methodology
- Results
- Demo
- Future Directions

# INTRODUCTION: FIGHTING FOR THOSE WITH RARE DISEASES

# INTRODUCTION: CLOSER TO HOME
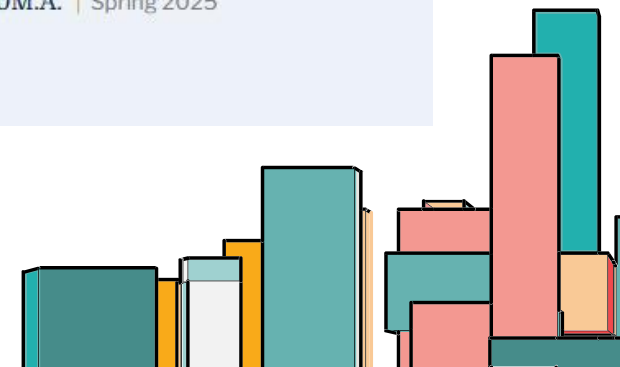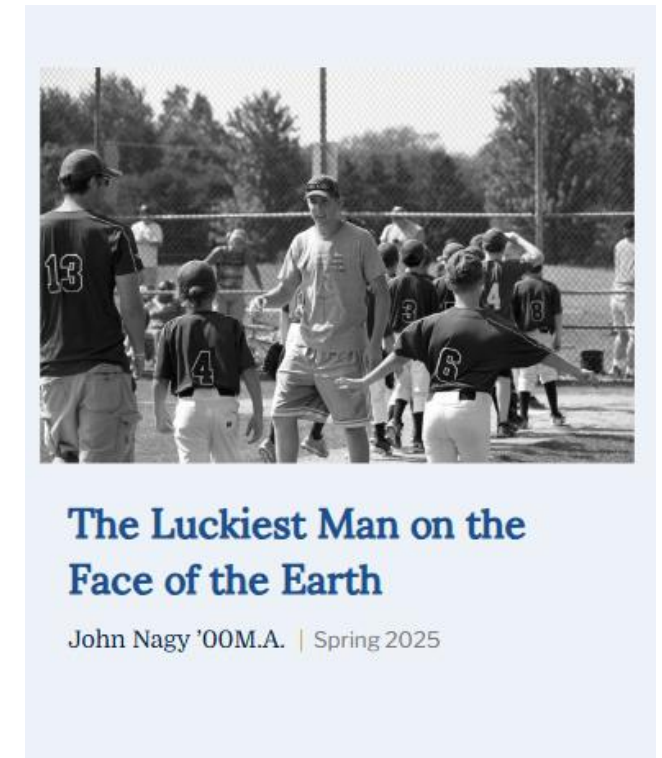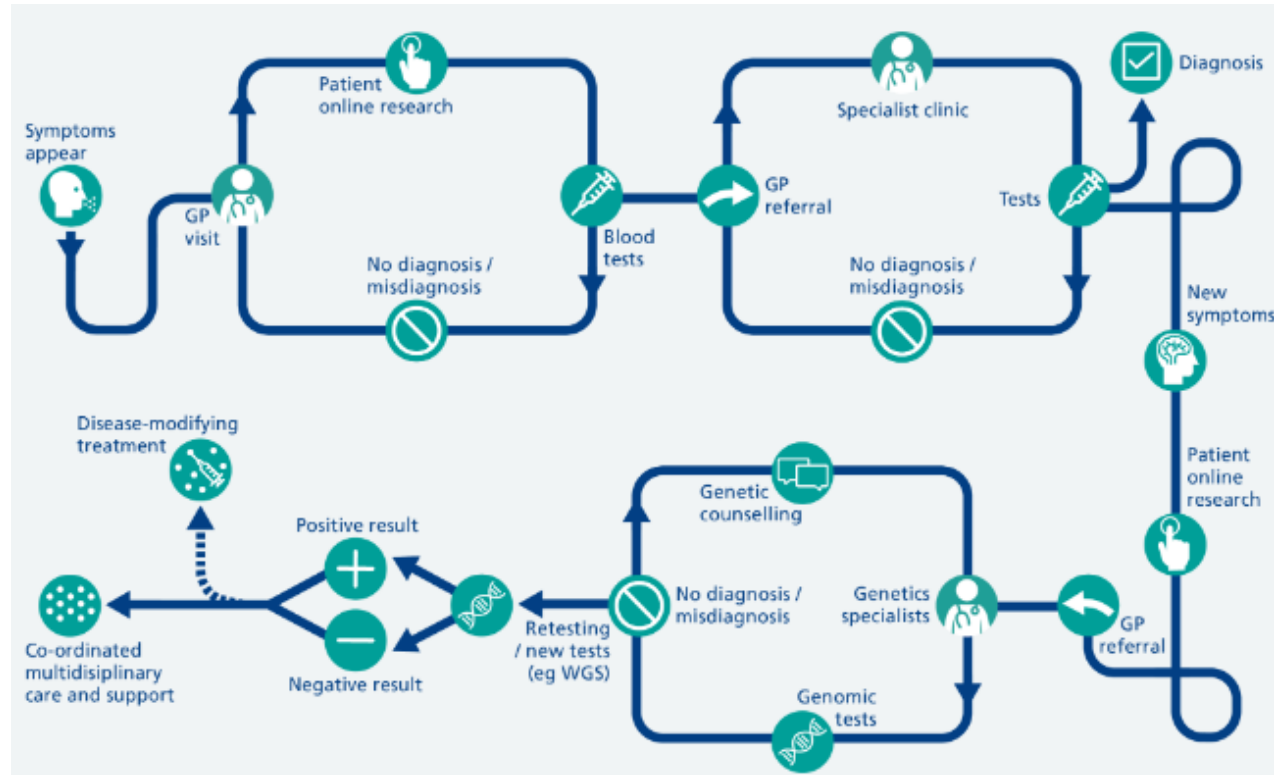


The Luckiest Man on the Face of the Earth

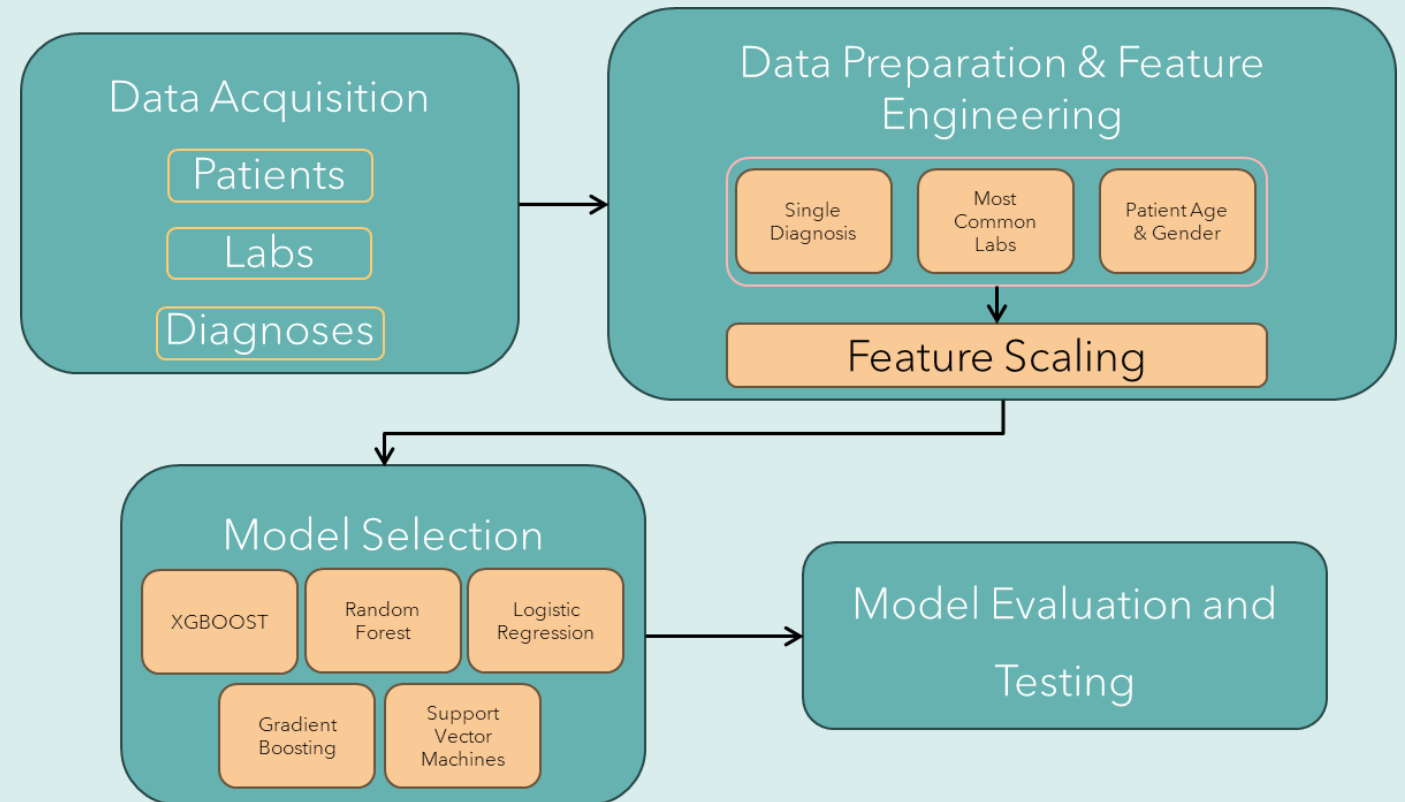John Nagy '00 M.A. | Spring 2025

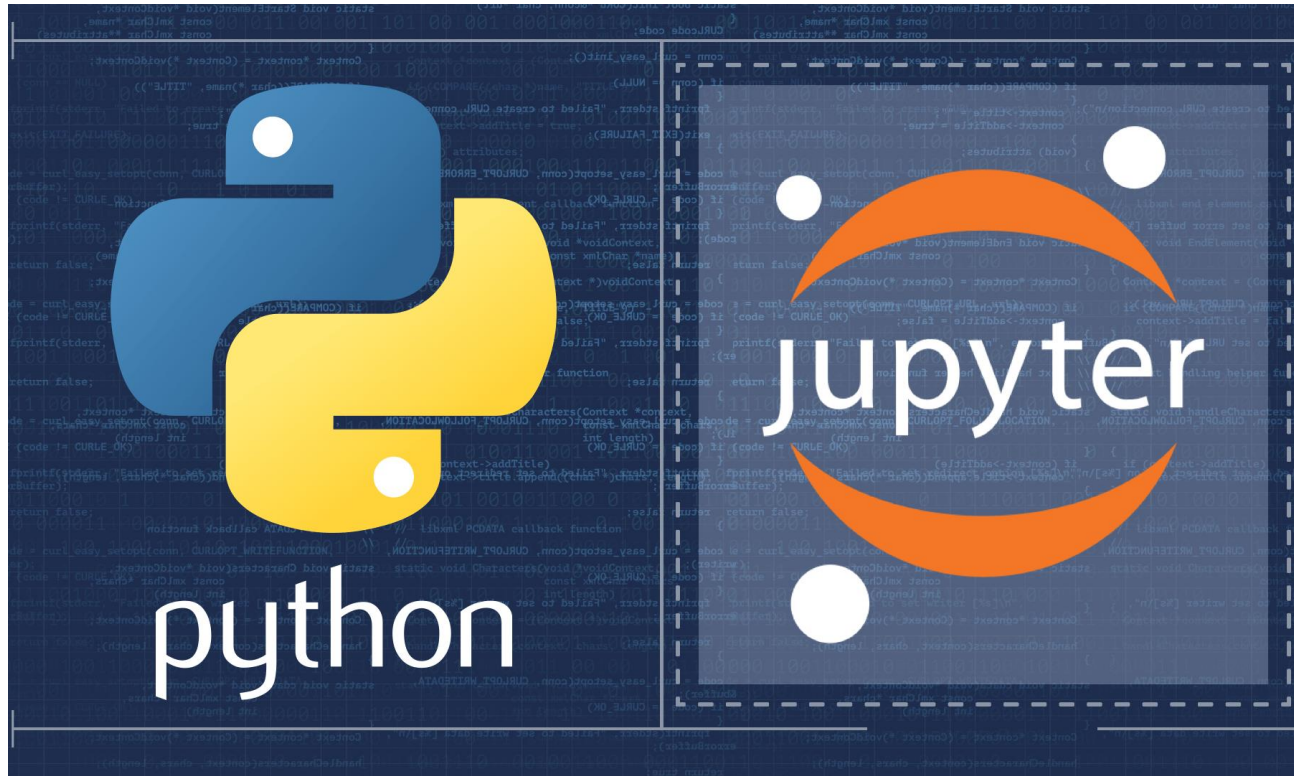Diagnostic odyssey in rare disease – Knowledge Hub

# METHODOLOGY

Steps to procure data, set up experiments, and evaluate results
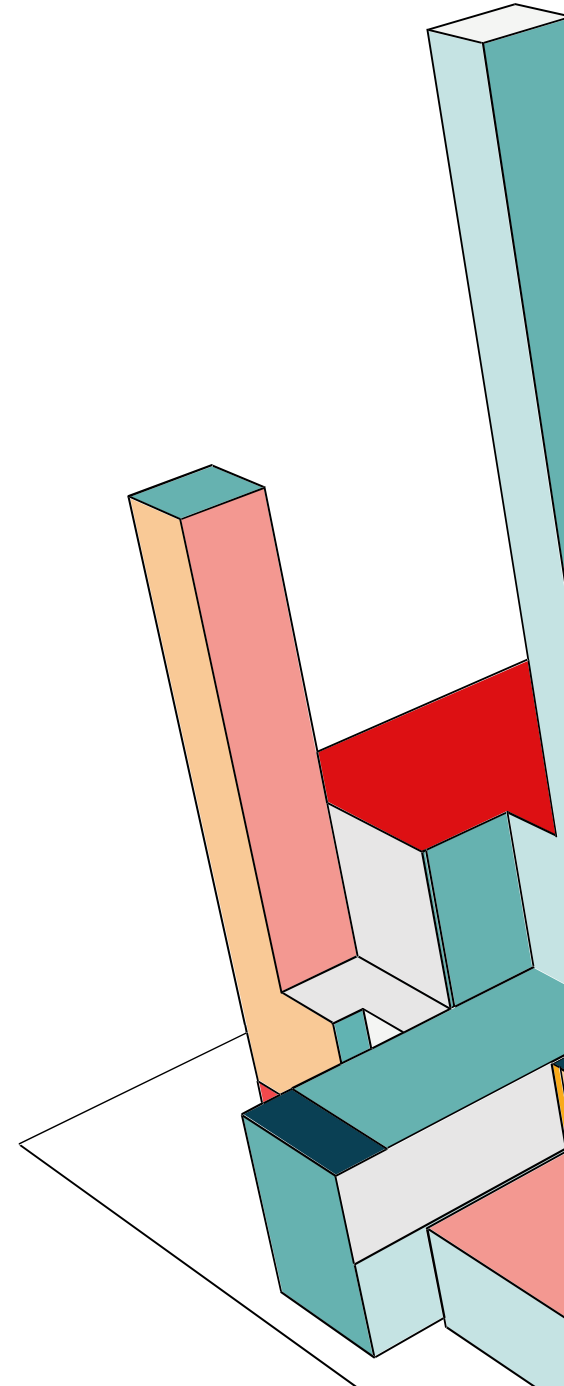
# METHODOLOGY: DATA ACQUISITION



- Used the MIMIC-III Clinical Database, provided by PhysioNet as the model's data source

- Comprehensive retrieval of tables including patients, lab events, lab items, diagnosis ICD, and diagnosis items

- Local copies of the gzipped CSV files were downloaded and imported into a Python notebook via the pandas package

# METHODOLOGY: DATA PREPARATION & FEATURE ENGINEERING

- Due to the unavailability of resources listing ICD-9 codes for diseases considered rare (less than 200,000 diagnoses), rare diseases were defined as those appearing only once in the dataset

- A threshold was set to identify the most common labs, defined as the top 10 most frequent labs

- Gender data was extracted from the patients table and mapped to a binary column (1 or 0)

- Age was calculated by subtracting the chart time from the lab admission date from the patient's date of birth

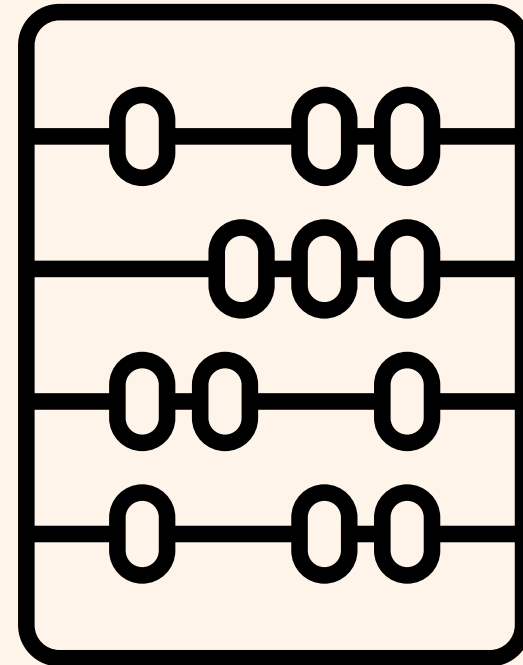- The MinMaxScaler from scikit-learn was utilized to standardize the data to a coherent scale

# METHODOLOGY: DATA IMBALANCE & MODEL SELECTION

- First used SMOTE then ADASYN to help deal with target variable imbalances
- Used grid search on multiple algorithm types

| Algorithm | Result |
| --- | --- |
| RandomForestClassifier | Poor |
| LogisticRegression | Very Poor |
| SVC | Very Poor |
| GradientBoostingClassifier | Very Poor |
| XGBClassifier | Crash Machine |

# RESULTS

# BEST RESULTS: RANDOM FOREST

Classification Report:

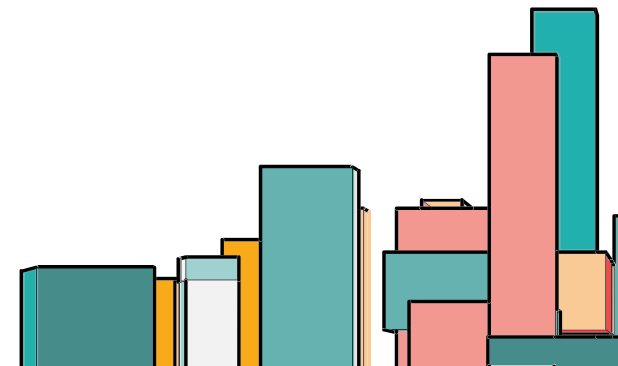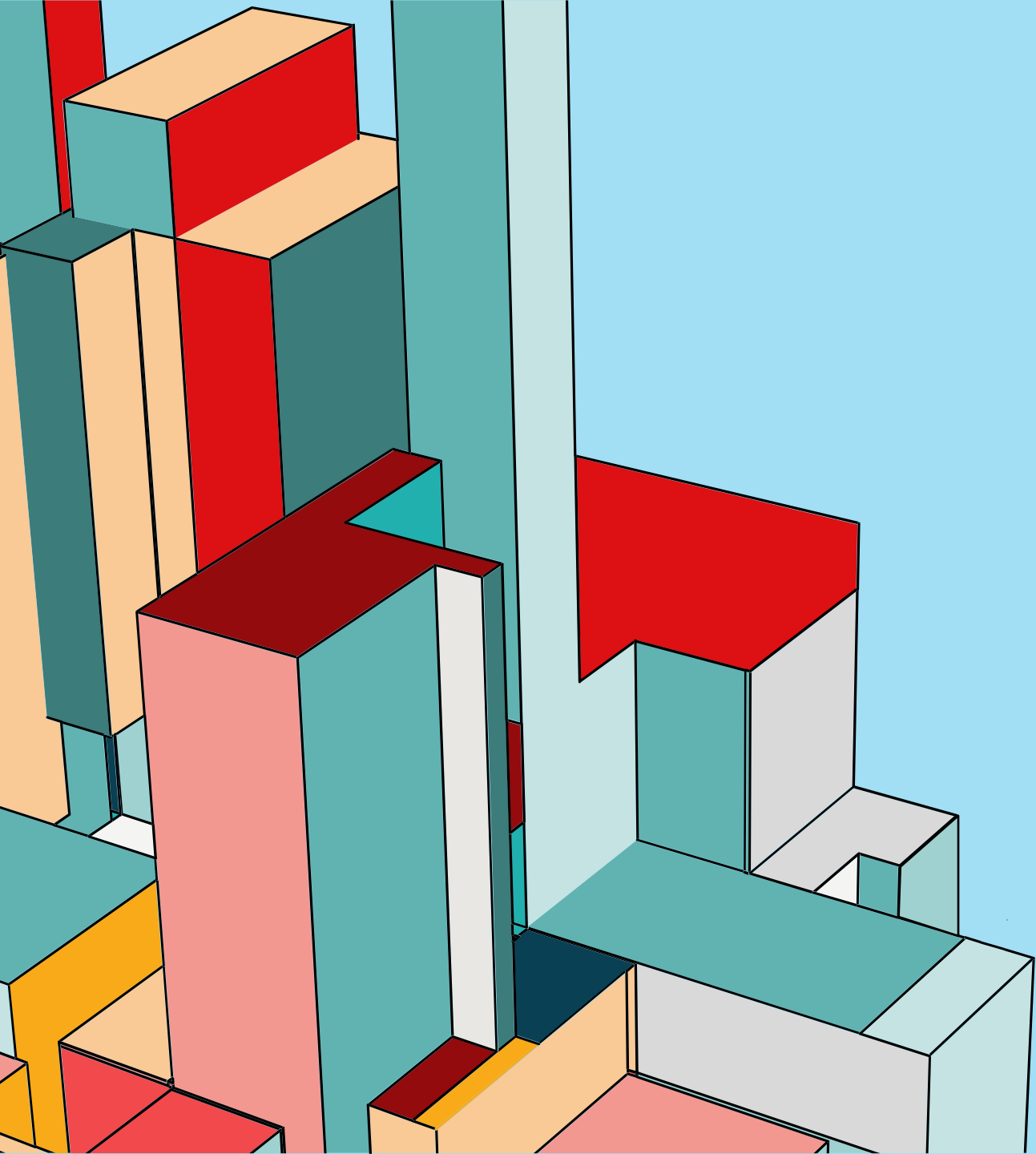|  | Precision | Recall | F1-score | Support |
|---|---|---|---|---|
| 0.0 | 0.98 | 0.88 | 0.93 | 11523 |
| 1.0 | 0.05 | 0.27 | 0.09 | 273 |
| | | | | |
| Accuracy | | | 0.87 | 11796 |
| Macro Avg | 0.52 | 0.58 | 0.51 | 11796 |
| Weighted Avg | 0.96 | 0.87 | 0.91 | 11796 |

Accuracy Score:

0.869616819261

Best results were obtained with a Random Forest Classifier with the hyperparams:

- 250 estimators

- 14 max depth

- 5 minimum samples split
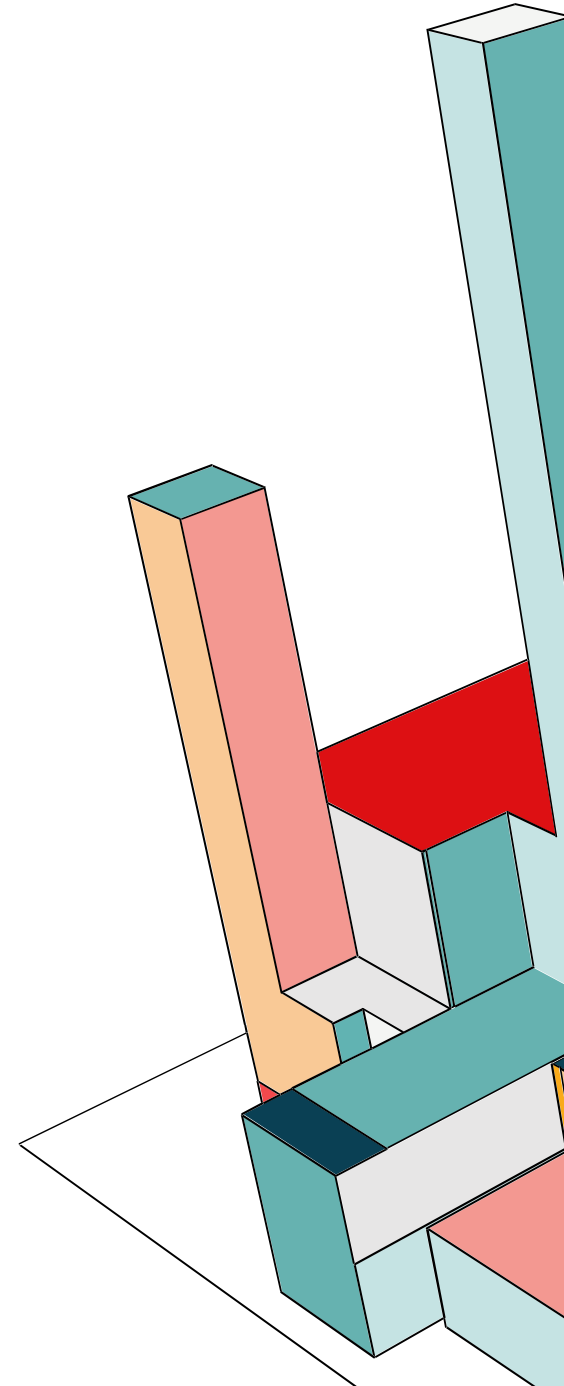
# DEMO

# DEMO: MODEL DATA & ALGORITHMS

# DEMO: MODEL SELECTION

```python
best_models = {}
for model_name, model_info in models_and_parameters.items():
    print(f"Running GridSearchCV for {model_name}...")
    grid_search = GridSearchCV(
        estimator=model_info["model"],
        param_grid=model_info["params"],
        scoring="accuracy",
        cv=5,  # 5-fold cross-validation
        n_jobs=1
    )
    grid_search.fit(X_train_adasyn, y_train_adasyn)
    grid_search = run_model(model_name, model_info, X_train_adasyn, y_train_adasyn)
    best_models[model_name] = grid_search.best_estimator_
    print(f"Best parameters for {model_name}: {grid_search.best_params_}")
    print(f"Best cross-validation accuracy for {model_name}: {grid_search.best_score_:.4f}")

print("\nEvaluating Best Models on Test Data:")
for model_name, model in best_models.items():
    print(f"\nEvaluating {model_name}...")
    y_pred = model.predict(X_test)
    print("Confusion Matrix:")
    print(confusion_matrix(y_test, y_pred))
    print("\nClassification Report:")
    print(classification_report(y_test, y_pred))
    print("\nAccuracy Score:")
    print(accuracy_score(y_test, y_pred))
```

# FUTURE DIRECTIONS

Further ideation

# FUTURE DIRECTIONS

- Exploring different thresholds for defining 'most common' tests beyond the top 10.

- Identifying the true list of ICD-9 codes for rare diseases.

- Incorporating additional demographic information to enrich the dataset.

- Including common chart data.

- Utilizing more powerful computational resources for modeling.

- Applying deep learning techniques to the data.

- Conducting a random search followed by a grid search to refine hyperparameters.

# THANK YOU

Maina Musa

University of Texas at Austin

April 2025

# APPENDIX

References
Flowcharts
Diagrams

# REFERENCES

[1] RareDisease.net. "Getting a Rare Disease Diagnosis." Available: https://raredisease.net/diagnosis.

[2] J. Benito-Lozano, G. Arias-Merino, M. Gómez-Martínez, A. Ancochea-Díaz, A. Aparicio-García, M. Posada de la Paz, and V. Alonso-Ferreira. "Diagnostic Process in Rare Diseases: Determinants Associated with Diagnostic Delay." Int. J. Environ. Res. Public Health, vol. 19, no. 11, p. 6456, 2022. Available: https://www.mdpi.com/1660-4601/19/11/6456.

[3] ScienceDirect. "Search for low-mass resonances decaying into two jets and produced in association with a photon using pp collisions at s=13 TeV with the ATLAS detector." Phys. Lett. B, vol. 795, pp. 56-75, 2019. Available: https://www.sciencedirect.com/science/article/pii/S0012369218300643.
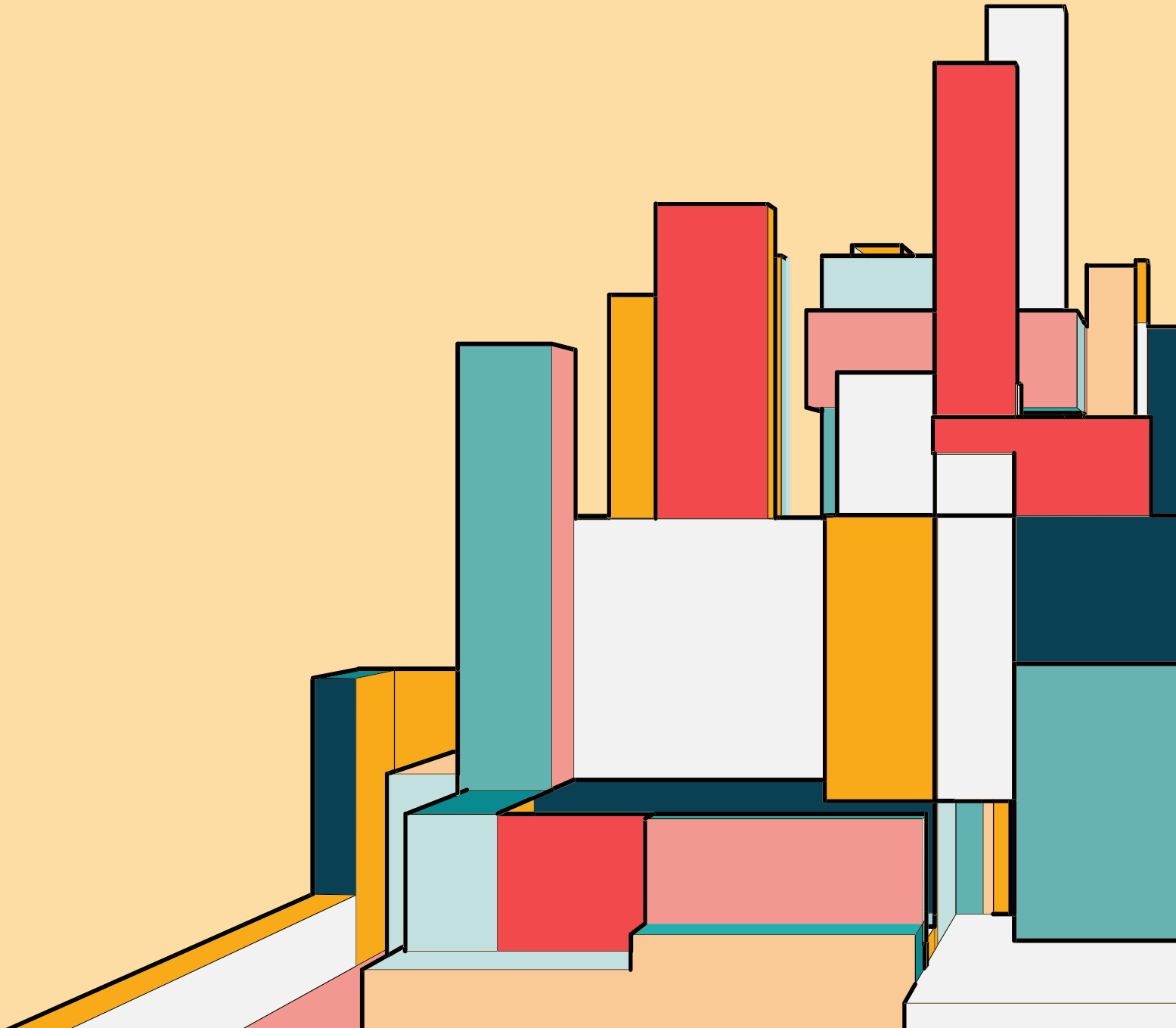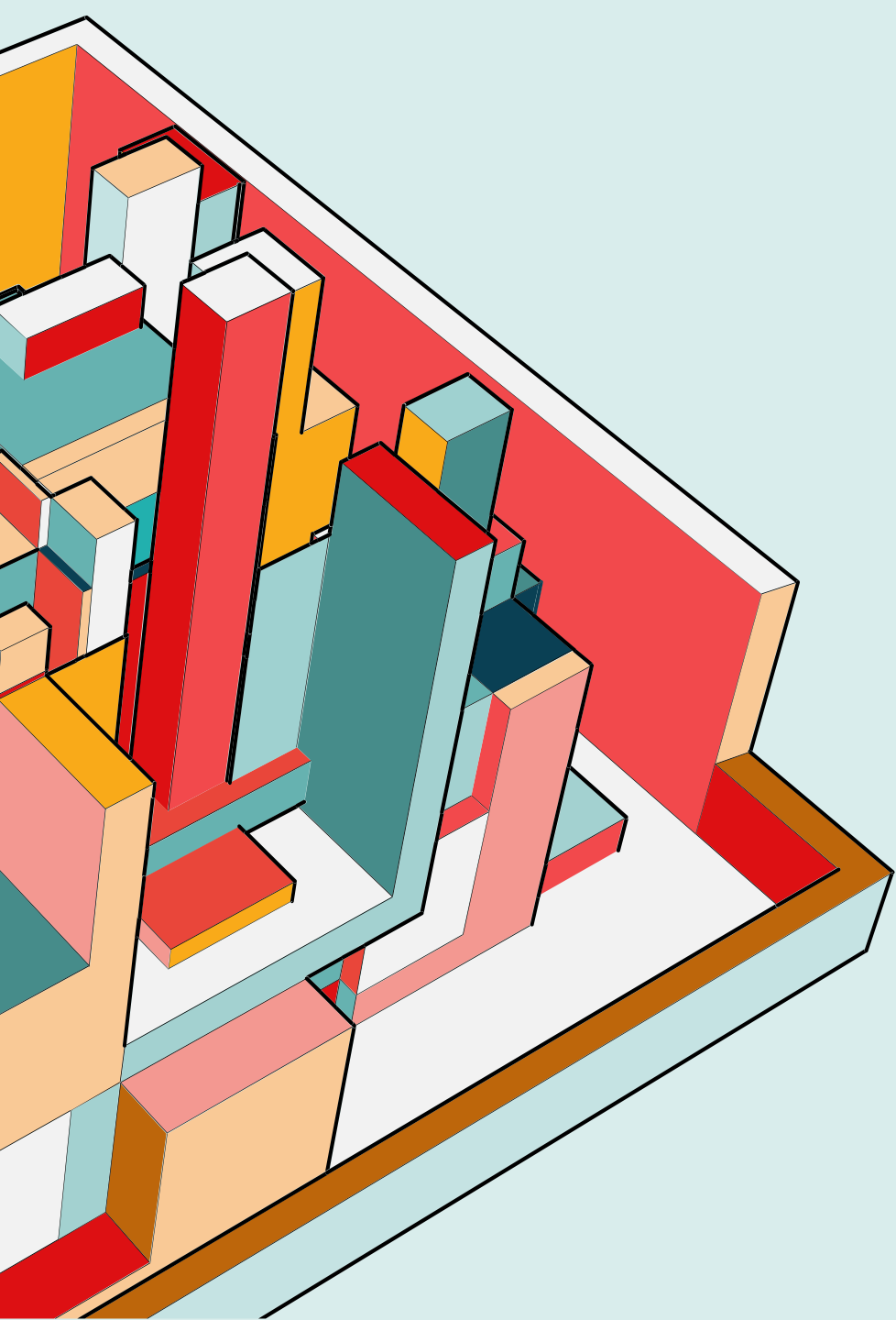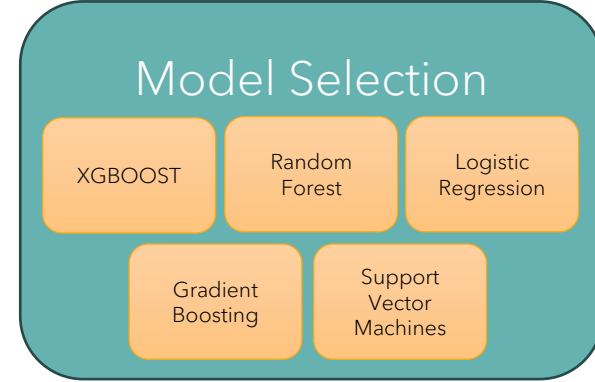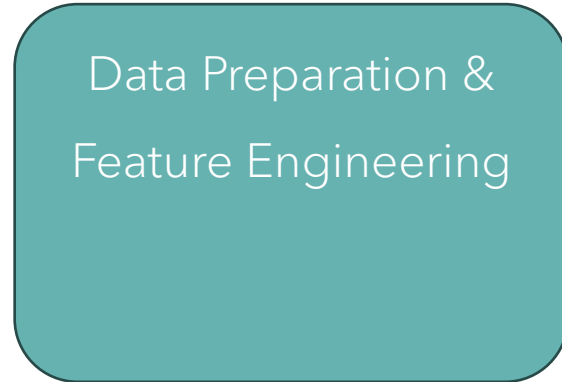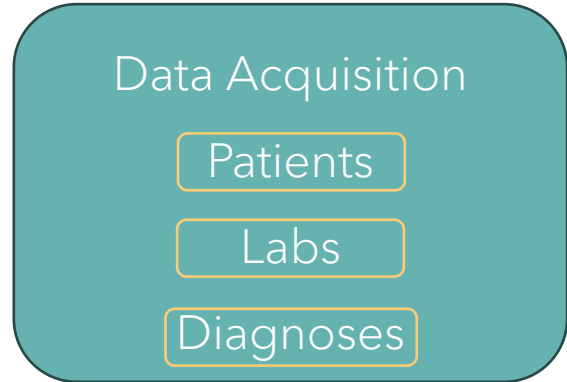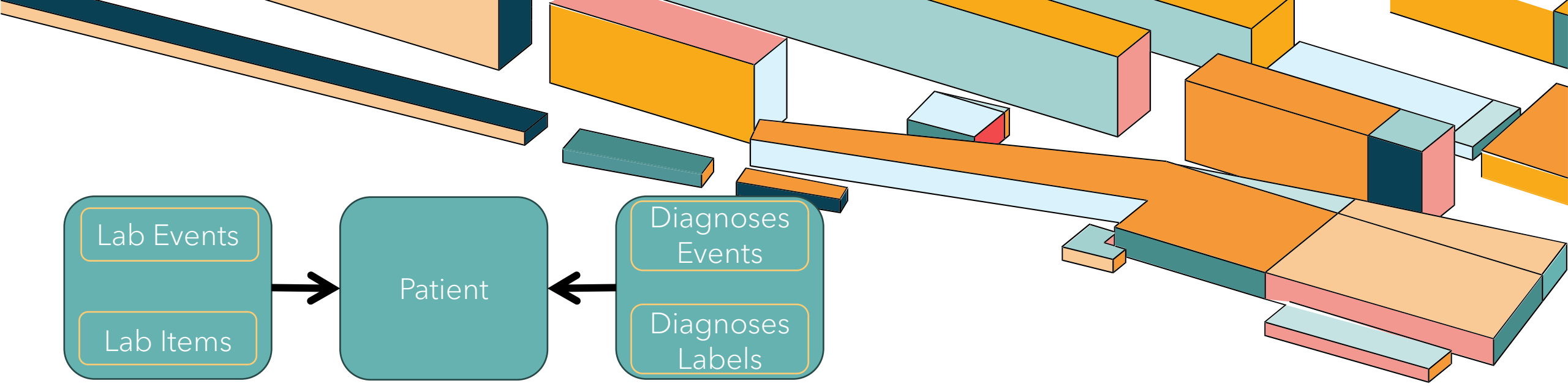
[4] R. Giugliani, S. Castillo Taucher, S. Hafez, J. B. Oliveira, M. Rico-Restrepo, P. Rozenfeld, I. Zarante, and C. Gonzaga-Jauregui. "Opportunities and challenges for newborn screening and early diagnosis of rare diseases in Latin America." Front. Genet., vol. 13, 2022. Available: https://www.frontiersin.org/journals/genetics/articles/10.3389/fgene.2022.1053559/full.

[5] P. Kováč, P. Jackuliak, A. Bražinová, I. Varga, M. Aláč, M. Smatana, D. Lovich, and A. Thurzo. "Artificial Intelligence-Driven Facial Image Analysis for the Early Detection of Rare Diseases: Legal, Ethical, Forensic, and Cybersecurity Considerations." AI, vol. 5, no. 3, pp. 990-1010, 2024. Available: https://www.mdpi.com/2673-2688/5/3/49.

[6] M. Kaasgaard, K. Grebosz-Haring, C. Davies, G. Musgrave, J. Shriraam, J. M. McCrary, and S. Clift. "Is it premature to formulate recommendations for policy and practice, based on culture and health research? A robust critique of the CultureForHealth (2022) report." Front. Public Health, vol. 12, 2024. Available: https://www.frontiersin.org/journals/public-health/articles/10.3389/fpubh.2024.1373649/full.

[7] FDA. "Rare Diseases at FDA." Available: https://www.fda.gov/patients/rare-diseases-fda.

[8] Johnson, T. Pollard, and R. Mark. "MIMIC-III Clinical Database v1.4." PhysioNet, 2016. Available: https://physionet.org/content/mimiciii/1.4/.

[9] McMullen, Tara Hunt "Fighting for Those With Rare Diseases" Available: https://fightingfor.nd.edu/2024/fighting-for-those-with-rare-diseases/

Lab Events

Lab Items

Patient

Diagnoses Events

Diagnoses Labels

Data Acquisition

Patients

Labs

Diagnoses

Data Preparation & Feature Engineering

Model Selection

XGBOOST

Random Forest

Logistic Regression

Gradient Boosting

Support Vector Machines

| Hyperparameter | Values |
| --- | --- |
| n_estimators | [200, 250] |
| max_depth | [10,12,14] |
| min_samples_split | [5, 8] |