

2025 Fall - STAT3612 Final Project

30-day All-Cause Hospital Readmission Prediction

Final Report

Leo Welma

Date: 30.11.2025

Contents

1 Introduction	3
2 Data Exploration	4
3 Baseline Models	6
3.1 Track 1: EHR-only	6
3.2 Track 2: Multimodal EHR	6
4 Methods	7
4.1 Standardization and Encoding	7
4.1.1 Ordinal Encoder	7
4.1.2 One-hot Encoder	7
4.2 Classification	7
4.2.1 Logistic regression	7
4.2.2 SVM	7
4.2.3 Ensemble Trees	7
4.2.4 Neural Networks	8
4.3 Model Development & Fusion	8
4.4 Feature Engineering	9
5 Results	10
5.1 Track 1: EHR-only	10
5.1.1 GRU	11
5.1.2 LSTM	11
5.2 Track 2: Multimodal EHR	11
6 Discussion	12
7 Appendix	14
7.1 Chronic Medical Conditions	14
7.2 Cancer-Related Conditions	15
7.3 Specific and Combined Conditions	16

1 Introduction

Reducing avoidable hospital readmission's is a primary goal for healthcare systems worldwide. In the United States, readmission rates serve as critical indicators of healthcare quality. Unplanned readmission's within a short time frame after discharge, typically occurring within 30 days, impose substantial burdens on both patients and healthcare systems. Machine learning, with its capacity to learn from extensive datasets encompassing diverse modalities, including patient clinical and imaging data, has demonstrated remarkable progress in various fields. This project aims to leverage machine learning techniques to predict 30-day hospital readmission's, contributing to improved patient care and healthcare system efficiency.

The aim is to predict 30-day, all-cause hospital readmission's at discharge using the MIMIC-IV v1.0 dataset. Performance was evaluated by AUROC, and all model selection was based on the provided validation split. The core dataset consisted of a day-by-day electronic health record (EHR) with 171 features covering demographics, comorbidities, laboratory test results, and medications. For Track 1 (EHR-only) a sequence-level summary statistics (such as last-day values and temporal aggregates) and compared different models was created, including logistic regression, support vector machines, random forests, XGBoost, and CatBoost. A CatBoost model trained on these summary features achieved the best validation performance (approximately 0.79 AUROC) and served as the final Track 1 model.

EHR data is not equally exhaustive for all features. Time-invariant features, such as demographic and comorbidity variables, remain constant throughout an admission and are therefore available every day. Time-varying predictors, such as lab tests and medications, are only performed or administered on certain days. Consequently, these features often have a value of zero. While this makes the data easy to store and feed into models, it also means the model must learn the difference between "true zero" and "not measured" and deal with a lot of sparsity, which can make learning stable temporal patterns more difficult.

For longitudinal modeling a GRU and LSTM based sequential models with daily EHR sequences using padded batches was performed with class-weighted binary cross-entropy loss and early stopping. These models achieved validation AUROCs of approximately 0.77, but consistently overfitted and did not outperform the CatBoost baseline.

For Track 2 (multimodal), the CatBoost framework was extended by combining EHR features with clinical notes and chest radiograph features from the MIMIC-CXR-JPG v2.0 dataset. Incorporating text improved discrimination ($\text{AUROC} = 0.82$), whereas adding image features via simple early fusion provided no additional benefit and increased overfitting.

Kaggle predictions for both tracks were submitted, documented in form of a full pipeline in a Jupyter notebook, and summarized in this final report.

2 Data Exploration

Each admission is organized as a daily EHR sequence, where each day is represented by 171 features. These features are categorized as follows:

- Demographics: Age, gender, and ethnicity (see Figure 2.).
- ICD labels (binary coded) that classify diseases, symptoms, and injuries. (91 unique features)
- Laboratory test variables (binary coded) that indicate whether lab tests were performed on that particular day. (36 unique features)
- Continuous medication variables that represent the quantity of administered prescriptions. (41 unique features)

The data consists of 64 time-varying features (mainly labs/meds) and 107 time-invariant features (mainly ICD labels). Length of stay is right-skewed with a mode near 8 days and a long tail of prolonged admissions (Figure 1.). Predefined train/validation/test splits for training, model selection, and final evaluation were used (Table 1.). Beyond descriptive, risk differences were computed to screen and compare variables. Risk differences are absolute-effect measures from survival analysis, representing the change in 30-day readmission probability between two groups by assessing the final day of admission. For binary indicators, the contrast is $x = 1$ and $x = 0$. Continuous variables were contrasted by the top and bottom quintiles as proxies for "high" and "low". The risk differences serve as simple guidelines to help with feature selection and model tuning (Figure 3.).

Split set	Total rows	Patients	Admissions	Re-admissions (< 30 days)	Prevalence (in %)
train	49451	6625	8234	1449	17.60
valid	16721	2208	2788	481	17.25
test	16293	2208	2741	no ground truth	

Table 1. Split table for the training, validation and testing sets

The prediction target is a binary indicator of 30-day, all-cause hospital readmission following discharge. For each MIMIC-IV admission, the patient was followed forward in time, and the stay was labeled as positive if an admission occurred within 30 days of the discharge date and as negative if an admission did not occur. In the training set, readmission prevalence was approximately 17%. This class imbalance motivated the use of class-weighted loss functions, (GRU/LSTM) and appropriate `scale_pos_weight` settings (XGBoost and CatBoost), which prevent the models from exploiting class imbalance by mostly predicting the majority class.

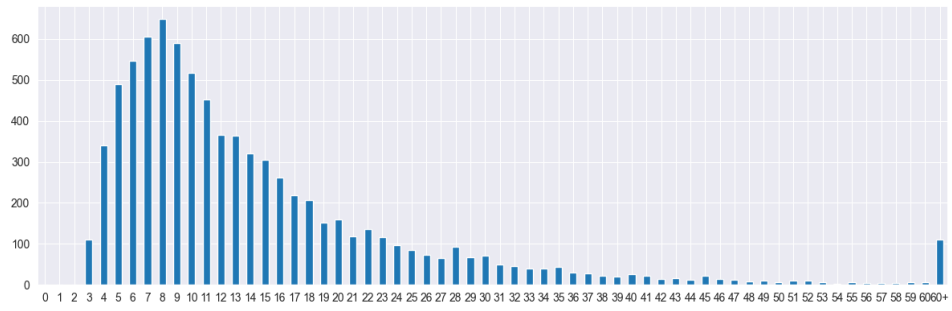


Figure 1: Distribution of admission length

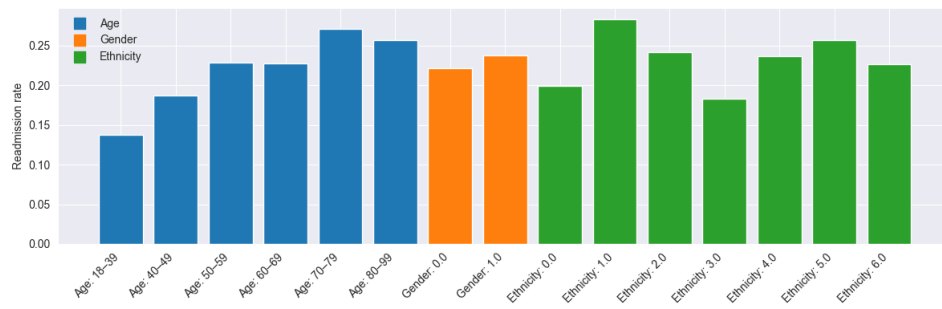


Figure 2: Readmission by Demographics [age, gender, ethnicity]

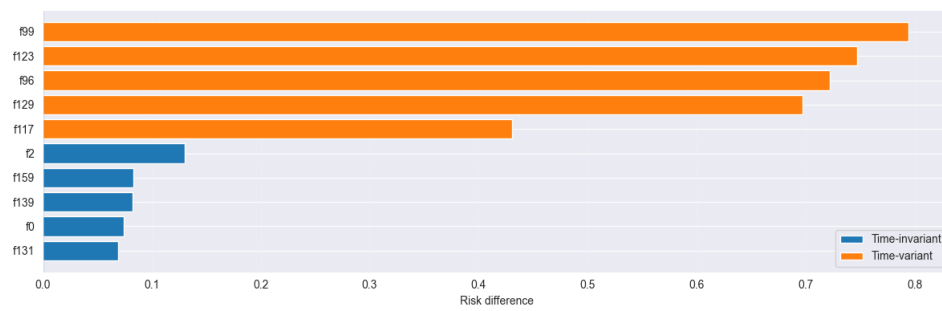


Figure 3: Risk differences for the features with the largest effects

3 Baseline Models

The EHR pickle file (`ehr_preprocessed_seq_by_day_cat_embedding.pkl`) was successfully loaded, and design matrices were constructed for the training, validation, and test sets. Design matrices were constructed for the training, validation, and test sets, indexed by unique identifiers. Parsing routines were implemented, and the `scale_pos_weight` parameter was computed to address class imbalance. Hyperparameter tuning for the CatBoost and XGBoost model was conducted using Optuna with cross-validation. The final CatBoost model was trained using the AUC metric with early stopping, and all relevant procedures were logged.

3.1 Track 1: EHR-only

For the Track 1 baseline, XGBoost and CatBoost models were explored. These models were trained using sequence-level summary statistics such as mean and standard deviation engineered from the sequential EHR data. As shown in Table 2, the CatBoost model achieved the highest validation AUROC. While this result is promising, the model exhibited significant overfitting. It is likely that these static summary features, while useful, fail to capture the complex temporal dynamics within the patient sequences, leading to the model overfitting on the training data.

Model	Train AUROC	Validation AUROC
XGBoost	0.8633	0.7838
CatBoost	0.9299	0.7898

Table 2. Baseline model metrics (EHR-only)

Future experiments should focus on expanded cross-validation and continued parameter exploration. Despite the simplicity of the current feature engineering, it yielded substantial performance improvements, highlighting the importance of further feature enhancement.

3.2 Track 2: multimodal EHR + images + notes

For the Track 2 baseline, a different feature engineering strategy was adopted to manage the high dimensionality of multimodal data, which made the Track 1 "sequence-level summary" approach computationally prohibitive. The EHR features were simplified to include only data from the last day of admission. For the text modality (`notes.csv`), Latent Semantic Analysis (LSA) was applied, reducing TF-IDF vectors (5000 features) to 128 components via TruncatedSVD. The pre-extracted image features (`image_features.zip`) were used directly. These three feature sets (EHR-last day, LSA-text, and image) were then concatenated and used to train a CatBoost model.

The results, presented in Table 3, reveal a key finding: the full multimodal baseline was outperformed by the **EHR + Text** model. This strongly suggests that while text features are highly predictive, the pre-extracted image features, when simply concatenated, introduce noise and degrade performance. This hypothesis is further supported by the low AUROC of the **EHR + Image** model.

	Train AUROC	Validation AUROC
EHR + Text	–	0.8156
EHR + Image	–	0.7071
EHR + Text + Image (Baseline)	0.9621	0.8087

Table 3. Baseline model metrics (EHR + images + notes)

The Track 2 baseline models also exhibit severe overfitting, particularly the full multimodal model, underscoring the need for further tuning and regularization. The current simple feature concatenation, while yielding a crucial insight—that text features are highly predictive, while image features (in their current form) are detrimental—is clearly suboptimal. This highlights the importance of further enhancement, moving beyond simple concatenation to explore more sophisticated multimodal fusion techniques in future experiments.

4 Methods

4.1 Standardization and Encoding

For some models such as logistic regression and support vector machine, it is necessary to standardize continuous features and encode categorical features.

Standardization is a preprocessing technique that transforms a continuous feature to have a mean of 0 and standard deviation of 1 by subtracting the mean and then divided by the standard deviation so as to ensure all features with the same scale. In the EHR dataset, most of the medication features are continuous.

Encoding is a preprocessing technique that transforms a categorical feature as shown below:

4.1.1 1. Ordinal encoder:

transforms a categorical feature to a new feature of integers (0 to n-1 for n categories). In the EHR dataset, ordinal encoder was applied to `readmitted_within_30days` to transform it by mapping “True” to 1 and “False” to 0.

4.1.2 2. One-hot encoder:

transforms each categorical feature of n categories into n binary features, with one of them 1 and others 0. In the EHR dataset, “ethnicity” is a categorical feature of 6 categories and thus one-hot encoder was applied to it.

4.2 Machine learning methods for classification

4.2.1 Logistic regression with LASSO regularization

A linear classification method that estimates probabilities using a logistic function. LASSO (L1) regularization adds a penalty term to the model to drive the coefficients of some features to zero. This effectively performs automatic feature selection and prevents overfitting.

4.2.2 Nonlinear Gaussian kernel support vector machine

A classifier that finds the optimal decision boundary by mapping data to a highdimensional space using the Gaussian (RBF) kernel. It captures complex nonlinear patterns and creates decision boundaries based on similarity measures between data points.

4.2.3 Ensemble tree-based models

Ensemble tree-based models, such as XGBoost and CatBoost, performed exceptionally well on the EHR data, outperforming simpler models, like logistic regression and random forests, by a large margin. These models can capture nonlinear relationships and interactions between variables with relatively little feature engineering. This makes them strong baselines and the best overall model for predicting 30-day readmissions.

- **Random Forest:** An ensemble method that builds multiple decision trees using bootstrap sampling and random feature selection. It combines predictions through averaging to reduce variance and prevent overfitting.

- **Gradient Boosting Machine (XGBoost):** A boosting algorithm that builds trees sequentially with each new tree correcting errors made by the previous ones. It optimizes computational performance and includes regularization to control model complexity.
- **CatBoost:** A gradient boosting implementation specifically designed to handle categorical features natively without extensive preprocessing. It uses ordered boosting to reduce overfitting and can provide good performance with minimal hyperparameter tuning.
- **AdaBoost:** An adaptive boosting algorithm that combines multiple weak classifiers by giving higher weights to misclassified instances in each iteration. It focuses on difficult cases to build a strong ensemble classifier.

4.2.4 Neural Network models

Flexible models which consist of interconnected layers of neurons. Through multiple hidden layers and activation functions, it learns hierarchical representations and approximate complex nonlinear relationships in data.

- A **GRU** is a type of recurrent neural network (RNN) designed to model trends and temporal dependencies in sequential data [Mienye et al., 2024]. Compared to fellow recurrent neural network LSTM, GRU offers a simpler gating mechanism with only an update and a reset gate, making it computationally more efficient.
- **Long short-term memory (LSTM)** networks are a type of recurrent neural network designed to model sequential data by maintaining information over time through gated memory cells. In this project, an LSTM was used to process day-level electronic health record (EHR) sequences for each admission. This allowed the model to learn how changes in labs, medications, and comorbidities during a hospital stay relate to the risk of readmission within 30 days.

4.3 Model Development and Fusion Strategy

For Track 1, a range of classical classification methods on engineered EHR features were explored, including logistic regression with LASSO regularization, decision trees, random forests, gradient boosted trees, XGBoost, and CatBoost. Using sequence-level summary statistics (such as last-day values and temporal aggregates). Ensemble tree-based models, in particular CatBoost, consistently outperformed simpler baselines and served as the strongest EHR-only model. In addition GRU- and LSTM-based models were implemented that directly consumed the daily EHR sequences, but these recurrent models did not surpass the best CatBoost configuration and tended to overfit after a few epochs.

For Track 2, the CatBoost was extended to a framework which incorporates additional modalities. Multimodal feature sets were constructed by concatenating last-day EHR features with text-derived representations from deidentified clinical notes and with pre-extracted image features from chest radiographs. This early-fusion strategy was evaluated in several variants:

- **EHR only:** CatBoost trained on EHR features as in Track 1.
- **EHR + Text:** EHR features combined with TF-IDF + truncated SVD representations of clinical notes.
- **EHR + Image:** EHR features combined with pre-extracted chest radiograph features.
- **EHR + Text + Image:** joint early fusion of all three modalities.

This systematic comparison showed that adding text features substantially improved performance over EHR alone, whereas the inclusion of image features in the current form did

not yield further gains and often increased overfitting. More advanced fusion strategies (for example late fusion or attention-based architectures) were therefore identified as promising directions for future work but were not implemented within the scope of this project.

As a first step in the modeling pipeline, classification methods and feature engineering choices were evaluated to see which provided the strongest predictive performance on readmission, and then built sequential and multimodal extensions on top of the best-performing tree-based baseline.

4.4 Feature Engineering

Feature engineering was essential in making the high-dimensional EHR data usable for a variety of models. For diagnostic information, the use of raw International Classification of Diseases (ICD) diagnosis indicators was compared with the grouping of ICD codes into the Elixhauser Comorbidity Index (ECI). The ECI categorizes comorbidities into 31 binary indicators derived from ICD codes, and it was originally proposed to predict hospital resource use [Elixhauser et al., 1998]. Replacing raw ICD features with ECI categories improved the performance of several classical models and provided a more compact, clinically interpretable representation of comorbidity burden.¹

Different feature groups were treated separately beyond comorbidities. Demographic variables and comorbidity indicators were handled as time-invariant features within each admission. Laboratory test and medication features, which vary over time and are often sparse, were used in two ways. For tree-based and other tabular models, simple, sequence-level summary statistics were derived, such as last-day values, as well as basic aggregates, such as means and standard deviations, over the stay. These statistics were used to obtain fixed-length feature vectors per admission. For the sequential GRU and LSTM models, the preprocessed day-level feature sequences supplied in the data were used directly, after standardising each feature using global training-set means and standard deviations.

For the multimodal track, additional features for text and imaging data were engineered. Deidentified clinical notes were transformed into high-dimensional TF-IDF vectors, which were then projected onto a lower-dimensional latent semantic space using truncated singular value decomposition (SVD). This yielded dense text representations that could be concatenated with EHR features. Chest radiograph information was incorporated with the provided pre-extracted image feature vectors. These vectors were connected with EHR features and, in some models, text features for early-fusion CatBoost classifiers. More advanced feature engineering, such as learned embeddings from clinical language models (like ClinicalBERT) or targeted feature selection on image components, was identified as a promising approach, though it was reserved for future research.

¹A description of the 31 ECI categories and their abbreviations is given in the Appendix.

5 Results

Table 4 shows AUROC for both training and testing datasets when fitting the above-mentioned models with and without ECI categorization of ICD. It is observed that ensemble tree-based models are generally better than other models. Moreover, there is not much improvement in prediction performance after ECI categorization.

Since the tree-based models demonstrated better prediction performance compared to other models, models along this direction were mainly used alongside feature engineering to further improve model performance.

Model	Without ECI categorization		With ECI categorization	
	Training AUROC	Validation AUROC	Training AUROC	Validation AUROC
Logistic regression	0.708	0.668	0.700	0.665
SVM	0.869	0.676	0.863	0.651
Random forest	0.821	0.713	0.807	0.705
XGBoost	0.797	0.694	0.783	0.702
CatBoost	0.823	0.711	0.815	0.702
AdaBoost	0.726	0.690	0.720	0.685
Neural network	0.830	0.653	0.744	0.665

Table 4. Training and validation AUROC for each model

5.1 Track 1: EHR-only

After promising results from the Catboost models, recurrent neural networks were explored. While deep learning architectures theoretically offer the advantage of capturing complex temporal dependencies, the empirical results demonstrated a consistent performance ceiling of ~ 0.77 AUROC for the recurrent models, even after implementing advanced techniques such as weighted loss for class imbalance, bidirectional layers, and explicit feature augmentation. Close but not close enough to the CatBoost baseline of ~ 0.80 AUROC. Further more developing a LSTM-CatBoost ensemble model where a prediction is averaged between the two models with weights was performed. The ensemble of LSTM and CatBoost failed to surpass the single-model performance of CatBoost, indicating that the weaker sequence model introduced more noise than signal.

Model	Validation AUROC
CatBoost (Baseline)	0.7982
Simple GRU (Weighted Loss)	0.7701
LSTM	0.7681
Ensemble (CatBoost + LSTM)	0.7589

Table 5. Performance comparison of Track 1 models (EHR-only)

Therefore the CatBoost model was prioritized for the final submission, over the more complex but less effective deep learning alternatives.

5.1.1 GRU model implementation

The Track 1 baseline, which relied on static summary statistics, showed significant overfitting and failed to capture temporal dynamics. To address this limitation, the primary plan is to develop longitudinal models, starting with a Gated Recurrent Unit (GRU), which is specifically designed to model sequences. Mienye et al. [2024] suggests that the performance of GRU is comparable to LSTM in many cases with a considerably faster training time with less computation required. The authors, however, point out that performance is task dependent, meaning that some datasets benefit from the complexity of LSTM while others perform well on a simpler GRU. Choi et al. [2016] employed a GRU, citing their simpler architecture and similar performance as LSTM, on a sequential EHR dataset. Outperforming their baseline models and achieving a recall@30 of 79%, meaning that 79% of the true diagnosis or medication codes for a patient’s next visit were contained within the model’s top 30 predicted codes.

Given these findings, exploring GRU seems like an attractive option for the project. They seem to be able to capture the temporal patterns in the sequential EHR data while maintaining a simpler and more cost-effective architecture compared to LSTM.

5.1.2 LSTM model implementation

A standalone LSTM using EHR sequences instead of summary features will be also implemented. For each admission, the LSTM is fed a variable-length sequence of daily feature vectors (171 features per day) extracted from the preprocessed data. Standardized data of all features using means and standard deviations computed on the training set is used. Padded mini-batches so that sequences of different lengths could be also processed together.

The LSTM model consisted of a single LSTM layer with 64 hidden units, a dropout layer, and a final linear layer. The final layer produced one logit per admission, which is converted into a readmission probability using a sigmoid function:

$$z_i = \mathbf{w}^\top \mathbf{h}_i + b, \quad \hat{p}_i = \sigma(z_i) = \frac{1}{1 + e^{-z_i}}.$$

The model was trained with a class-weighted binary cross-entropy loss:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N [w_1 y_i \log(\hat{p}_i) + w_0 (1 - y_i) \log(1 - \hat{p}_i)].$$

The model is trained using the Adam optimizer with a batch size of 64 and a class-weighted binary cross-entropy loss function to account for the smaller number of readmitted patients. The LSTM achieved a validation AUROC of approximately 0.76, showed clear signs of overfitting after a few epochs, and did not outperform the optimal CatBoost model using EHR summary features.

5.2 Track 2: multimodal EHR + images + notes

For Track 2, the EHR-only models is expanded by incorporating information from de-identified clinical notes and chest radiographs. They represent each admission with EHR features from the last day, TF-IDF + truncated SVD (LSA) embeddings of the associated clinical notes, and pre-extracted image features from MIMIC-CXR. These feature sets are concatenated and are used as inputs for CatBoost classifiers under several configurations: EHR only, EHR + text, EHR + image, and EHR + text + image. Adding text features improved performance over using EHR alone. However, image features, in their current form, did not improve performance and tended to increase overfitting when combined with EHR and text.

6 Discussion

In this project, we compared a wide variety of modeling strategies for predicting 30-day, all-cause hospital readmissions. We used longitudinal EHR data, as well as additional text and imaging information in Track 2. The results from Track 1 consistently showed that ensemble tree-based models, particularly CatBoost, outperformed classical classifiers. CatBoost achieved a validation AUROC of about 0.80 using simple sequence-level summary statistics (last-day values and basic aggregates) as input features, outperforming logistic regression, SVMs, random forests, AdaBoost, and a shallow neural network. These results suggest that for this task, most useful information can be captured using non-linear tree ensembles applied to well-engineered features rather than more complex architectures.

Driven by the sequential nature of the data, GRU- and LSTM-based models were developed that operate directly on daily EHR sequences. These recurrent models used padded batches and class-weighted binary cross-entropy loss to address the approximately 17% readmission prevalence. GRU and LSTM both improved quickly during the first few epochs, reaching validation AUROCs around 0.77. However, further training led to overfitting: training loss decreased while validation performance plateaued and then declined. Therefore, the models did not surpass the CatBoost baseline. An ensemble that averaged CatBoost and LSTM predictions performed worse than CatBoost alone. This suggests that the sequence models mainly contributed noise rather than a complementary signal. One possible explanation is that many features are static or sparse over time, limiting the additional temporal information that RNNs can exploit compared to summary statistics.

For Track 2, the CatBoost was adapted to a framework for a multimodal setting by combining electronic health record (EHR) features with text representations derived from clinical notes combined with pre-extracted chest radiograph features. This revealed a consistent trend: incorporating text features enhanced discrimination relative to EHR alone (validation AUROC approximately 0.82), whereas including image features in their current form diminished performance and exacerbated overfitting. The EHR + Image model performed substantially worse than the EHR + Text model. Additionally, the full EHR + Text + Image model did not outperform the EHR + Text model, despite a very high training AUROC.

References

- Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2):423–443, 2018.
- Marianne Chetcuti. Predicting hospital readmission from electronic health record data using a machine learning approach. Master’s thesis, University of Malta, 2024.
- Edward Choi, Mohammad Taha Bahadori, Andy Schuetz, Walter F Stewart, and Jimeng Sun. Doctor ai: Predicting clinical events via recurrent neural networks. In *Machine learning for healthcare conference*, pages 301–318. PMLR, 2016.
- Anne Elixhauser, Claudia Steiner, D Robert Harris, and Rosanna M Coffey. Comorbidity measures for use with administrative data. *Medical care*, 36(1):8–27, 1998.
- Ibomoiye Domor Mienye, Theo G Swart, and George Obaido. Recurrent neural networks: A comprehensive review of architectures, variants, and applications. *Information*, 15(9): 517, 2024.

7 Appendix

7.1 Chronic Medical Conditions

1. aids: AIDS/HIV
 - With immunocompromising condition.
2. cpd: Chronic Pulmonary Disease
 - A group of diseases affecting the lungs and airways, such as COPD, asthma, emphysema, and bronchiectasis.
3. diabc: Diabetes without chronic complications
 - Diabetes that is controlled and has not yet led to major issues like neuropathy, retinopathy, or nephropathy.
4. diabunc: Diabetes with chronic complications
 - Diabetes that has caused end-organ damage (e.g., kidney disease, eye damage, nerve problems).
5. hypc: Hypertension with complications
 - High blood pressure that has caused damage to organs like the heart, kidneys, or brain.
6. hypothy: Hypothyroidism
 - An underactive thyroid gland.
7. ld: Liver Disease
 - Mild to moderate liver conditions, such as cirrhosis, chronic hepatitis, or liver fibrosis (this is separate from severe liver failure).
8. pud: Peptic Ulcer Disease (excluding bleeding)
 - Ulcers in the stomach or duodenum.
9. rf: Renal Failure
 - Chronic kidney disease, requiring dialysis or significantly reduced kidney function.
10. valv: Valvular Disease
 - Disorders of the heart valves (e.g., stenosis or regurgitation).
11. chf: Congestive Heart Failure
 - A condition where the heart doesn't pump blood as well as it should.
12. pcd: Pulmonary Circulation Disorders
 - Conditions affecting the blood vessels in the lungs, such as pulmonary hypertension or pulmonary embolism.
13. pvd: Peripheral Vascular Disease
 - Circulatory problems causing narrowed blood vessels outside the heart and brain, often in the legs.
14. para: Paralysis

- Neurological disorders resulting in paralysis, such as hemiplegia, paraplegia, or quadriplegia.
- 15. ond: Other Neurological Disorders
 - Chronic neurological conditions like Parkinson's disease, multiple sclerosis, or epilepsy (excluding paralysis and stroke).
- 16. rheumd: Rheumatoid Arthritis / Collagen Vascular Diseases
 - Autoimmune and inflammatory disorders such as rheumatoid arthritis, lupus, and vasculitis.
- 17. coag: Coagulopathy
 - Bleeding or clotting disorders, such as hemophilia or thrombocytopenia.
- 18. obes: Obesity
 - Diagnosed obesity, typically a Body Mass Index (BMI) of 30 or greater.
- 19. wloss: Weight Loss
 - Severe, recent, and unintentional weight loss, often indicative of malnutrition or a severe chronic illness like cancer.
- 20. fed: Fluid and Electrolyte Disorders
 - Imbalances like dehydration, hyponatremia, or acidosis. This is a strong predictor of mortality.
- 21. blane: Blood Loss Anemia
 - Anemia caused by acute or chronic blood loss.
- 22. dane: Deficiency Anemia
 - Anemia due to nutritional deficiencies, such as iron, B12, or folate deficiency.
- 23. alcohol: Alcohol Abuse
 - Harmful patterns of alcohol use, including dependence and associated complications.
- 24. drug: Drug Abuse
 - Substance use disorders involving illicit or prescription drugs.
- 25. psycho: Psychoses
 - Severe mental disorders, such as schizophrenia, delusional disorders, or other psychotic conditions.
- 26. depre: Depression
 - Major depression and other depressive disorders.

7.2 Cancer-Related Conditions

1. lymph: Lymphoma
 - Cancers of the lymphatic system, including Hodgkin's and non-Hodgkin's lymphoma.
2. metacanc: Metastatic Cancer
 - Cancer that has spread from its original site to other parts of the body.
3. solidtum: Solid Tumor without Metastasis
 - A localized malignant tumor (e.g., breast, lung, colon cancer) that has not yet spread.

7.3 Specific and Combined Conditions

1. hypunc: Hypertension, Uncomplicated
 - High blood pressure without documented end-organ damage.
2. carit: Cardiac Arrhythmias
 - Abnormal heart rhythms, such as atrial fibrillation, flutter, or conduction disorders.