# Automating Referral Triaging

Cian Allen - 501342468 - June 21, 2025

Supervisor: Dr. Sanaz Mohammad Jafari

**Ryerson University**

# Table of Contents

# Project Abstract

## Defining the Research Question

This project seeks to address the inefficiencies in referral management workflows within healthcare systems, particularly the manual triaging of referral letters and specialist allocation.

**The core research question:**
How can synthetic data and natural language processing be leveraged to automate and optimize referral triaging in a privacy-preserving and clinically valid manner?
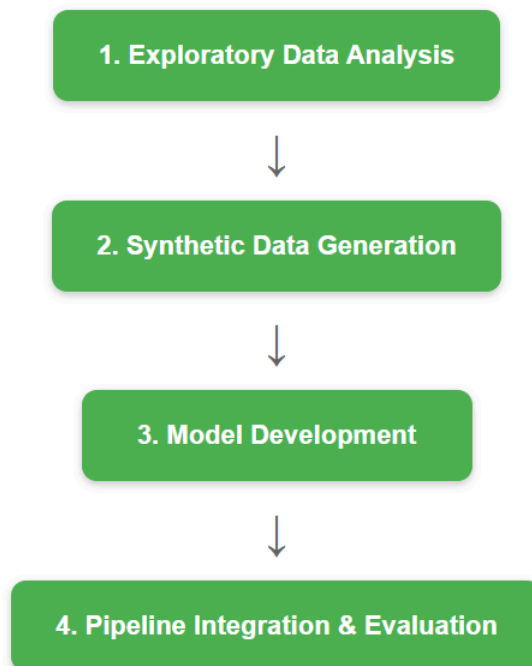
**This question is explored through three sub-questions:**

1. Data Validity & Privacy: How can GAN-based methods be used to generate synthetic healthcare data that maintains statistical realism while ensuring patient privacy?

2. Automated Classification: Can natural language processing models accurately classify referral letters by medical specialty to support automated triaging?

3. Pipeline Scalability: What does a modular and scalable architecture for synthetic referral generation and NLP-based classification look like, and how can it be evaluated for real-world applicability in terms of performance, adaptability, and integration potential?

Together, these sub-questions frame a problem that is both technically challenging and clinically impactful, with the goal of developing a practical AI-assisted solution to reduce patient wait times.

## Planning the Approach

To address the defined research question, the project will be carried out in the following stages:

```
┌─────────────────────────────────┐
│   1. Exploratory Data Analysis   │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│   2. Synthetic Data Generation   │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│      3. Model Development        │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│ 4. Pipeline Integration & Evaluation │
└─────────────────────────────────┘
```

**Exploratory data analysis:**

Conduct statistical profiling of the healthcare data to identify clinically relevant distributions and co-occurrence patterns (e.g., common diagnoses per age group or medication clusters by condition).

Use insights from this analysis to generate synthetic patient profiles that remain representative of the original data while ensuring full anonymization.

**Synthetic data generation:**

Use a Generative Adversarial Network (GAN) to synthesize realistic patient profiles based on patterns observed in the data.

Develop a complementary framework for generating synthetic referral letters using templating methods informed by real-world clinical writing conventions.

Ensure diversity in both structured data and natural language, while preserving clinical relevance and privacy.

Annotate each letter with a corresponding specialist label to serve as ground truth for supervised learning.

**Model development:**
Fine-tune domain-specific NLP models to classify referral letters by medical specialty.

Experiment with baseline models and transformer-based architectures to compare performance.

Evaluate models on classification accuracy, precision/recall, and robustness to linguistic variation.

**Pipeline integration and evaluation:**
Combine the synthetic data generation and classification components into a modular ML pipeline.

Assess scalability, data throughput, and model generalizability.

Discuss limitations, such as the domain shift between synthetic and real data, and propose next steps for real-world deployment or integration into EHR systems.

# Literature Review

## Data Validity and Privacy

This section summarizes research exploring how GAN-based methods can be used to generate synthetic healthcare data that preserves statistical accuracy while protecting patient privacy.

### 1. Generative adversarial networks and synthetic patient data: Current challenges and future perspectives (2022)
This comprehensive review surveys the application of GANs in healthcare, focusing on their ability to generate synthetic patient data for clinical research, medical education, and privacy protection. The authors highlight that GANs, by design, learn the underlying distribution of real data, enabling the creation of synthetic datasets that preserve statistical properties such as correlations, distributions, and rare events.

The review emphasizes that synthetic data generated by GANs can be used to augment training datasets for machine learning models, thereby improving model performance and robustness. Importantly, the review discusses the challenges of ensuring that synthetic data does not inadvertently reveal sensitive information about real patients. The authors note that while GANs excel at capturing complex patterns, additional privacy-preserving techniques such as differential privacy or data anonymization must be integrated in order to fully protect patient identities.

## 2. Synthetic data generation: A privacy-preserving approach to enabling data sharing and analysis in healthcare (2025)

This recent review focuses on the role of synthetic data in rare disease research, where data scarcity and privacy concerns are especially acute. The authors discuss how GAN-based synthetic data can bridge data gaps, enabling secure, cross-institutional collaborations. The review highlights various privacy-preserving techniques, including differential privacy, secure multi-party computation, and federated learning.

The review is highly relevant because it emphasizes the importance of maintaining statistical realism while ensuring privacy, especially in scenarios where real data is limited or highly sensitive. The authors argue that synthetic data generated by GANs, when combined with privacy-preserving techniques, can enable new research opportunities without compromising patient confidentiality. The review also discusses regulatory compliance, noting that synthetic data can help institutions meet the requirements of GDPR and HIPAA.

GAN-based synthetic data, augmented with privacy-preserving techniques, is a viable solution for rare disease research and other data-scarce, privacy-sensitive domains.

## 3. Privacy preserving generative adversarial networks to model health data: A systematic evaluation (2022)

This study introduces a privacy-preserving GAN (pGAN) model specifically designed for generating synthetic health records. The authors propose a custom loss function that penalizes the generation of data points too similar to real records, thereby reducing the risk of re-identification. The model is validated using classification and regression tasks, demonstrating that the synthetic data maintains high utility for machine learning.

The pGAN model is directly relevant because it addresses both statistical realism and privacy. The custom loss function ensures that synthetic data is realistic enough for research purposes while minimizing the risk of privacy breaches. The authors provide empirical evidence that pGAN-generated data can be used to train accurate machine learning models, even when real data is unavailable.

Customizing the GAN loss function to penalize privacy risks is an effective strategy for generating synthetic healthcare data that is both useful and private.

## 4. Synthetic data generation via generative adversarial networks in healthcare: A systematic review (2024)

This work presents a modified GAN architecture for generating tabular medical data. The authors focus on scenarios where real data is inaccessible due to privacy concerns. The proposed GAN is trained using statistical information derived from real data, rather than the raw data itself, further enhancing privacy. The synthetic data is evaluated using utility and similarity metrics, showing excellent performance.

The paper is highly relevant because it demonstrates that GANs can generate realistic synthetic data even when direct access to real data is restricted. By using statistical summaries instead of raw data, the approach provides an additional layer of privacy protection. The authors show that the synthetic data is suitable for a wide range of research and machine learning tasks.

GANs can generate high-quality synthetic healthcare data from statistical summaries, ensuring both realism and privacy when real data is inaccessible.

## Conclusions

The reviewed literature demonstrates that GAN-based methods are at the forefront of synthetic healthcare data generation. These methods excel at capturing the statistical properties of real data, making synthetic datasets suitable for research, education, and machine learning. However, ensuring patient privacy requires the integration of additional techniques such as differential privacy, custom loss functions, and training on statistical summaries rather than raw data.

The most effective approaches combine the generative power of GANs with rigorous privacy-preserving mechanisms. This dual focus enables the creation of synthetic datasets that are both realistic and compliant with privacy regulations, opening new possibilities for healthcare research and innovation.

GAN-based synthetic healthcare data generation is a rapidly evolving field with significant potential to address the challenges of data utility and patient privacy. The most relevant literature highlights the importance of integrating privacy-preserving techniques into GAN frameworks, ensuring that synthetic data is both statistically realistic and secure. As the field advances, these methods will play an increasingly critical role in enabling secure, collaborative, and innovative healthcare research.

# Automated Classification

This section summarizes research on how natural language processing models can accurately classify referral letters by medical specialty to support automated triaging.

## 1. Semantics-Based Classification of Medical Referral Letters: High-Accuracy Automated Triage

A recent study by Davies et al. (2023) developed a Natural Language Processing (NLP)-based decision support system specifically for classifying referral letters by medical specialty. Using a large dataset of 111,700 referral letters from the National Health Service Wales, the authors implemented a semantic matrix of document vectors and vocabulary features extracted from the text of each letter. Classification was performed using a Support Vector Machine (SVM) with a one-versus-rest approach, allowing each document to be matched probabilistically to its best-fit specialty.

The results were notable: the system achieved an accuracy of 91.8% in classifying letters into 29 medical specialties. When the model was allowed to consider the top two or three nearest specialties, accuracy increased to 97.4% and 99%, respectively. This high level of performance demonstrates that NLP models, when trained on large, representative datasets, can robustly classify referral letters by specialty. The authors highlight that their approach does not require additional ontologies and is easily extendable, making it suitable for integration into clinical workflows to support automated triaging and decision-making. The system is particularly valuable in training scenarios or in settings where specialist input is limited, offering timely and accurate allocation of referrals to the appropriate specialty.

## 2. NLP and Machine Learning for Triage in Musculoskeletal Care: Topic Modeling and Clinical Interpretability

A feasibility study by Chapman et al. (2020) explored the use of NLP and machine learning to automate the triage of patients with musculoskeletal conditions by analyzing referral letters. The study applied latent Dirichlet allocation (LDA) to model referral letters as mixtures of clinically relevant topics. These topics served as features for binary classifiers predicting treatment outcomes.

The classifiers significantly outperformed random baselines, indicating that topic modeling could effectively support automated triage by predicting the most appropriate treatment pathway. Importantly, the topics identified by the model were found to be clinically interpretable by human experts, supporting the practical utility of the approach. The study concludes that NLP-driven topic modeling is both feasible and effective for extracting actionable information from referral letters to aid in triage decisions for patients with knee or hip pain.

## 3. NLP-Based Extraction Improves Decision Support but Clinical Accuracy Remains a Challenge

A 2024 study by Fudickar et al. investigated the impact of enriching decision support systems (DSS) with NLP-based extraction from referral letters in the context of low back pain (LBP) triage. The study evaluated 1,608 patient cases and found that including qualitative data from referral letters improved the F1-score for triaging, with increases of up to 19.5% for certain referral reasons (such as anesthesiology and rehabilitation interventions). Despite these improvements, the overall model accuracies were still considered low and insufficient for direct clinical application.

The authors conclude that while NLP-based extraction of referral letter content enhances model performance and can suggest optimal treatments, further work is needed to achieve the reliability required for routine clinical deployment. This highlights the importance of both data quality and model optimization in realizing the full potential of NLP for automated triage.

**4. Generalizability and Adaptability of NLP Pipelines for Specialty Classification**

A 2025 study by van der Laan et al. focused on developing a generalizable NLP pipeline for predicting and prioritizing diagnoses from referral letters in rheumatology outpatient clinics. Using advanced techniques such as BERT transformers, the pipeline analyzed structured and unstructured data, identifying features like suspected diagnoses, lab results, and symptom descriptions as key for accurate classification.

The study emphasizes that while the pipeline shows promise for cross-language and cross-system adaptation, local validation and optimization are essential due to variations in healthcare practices and referral protocols. The authors note that previous studies lacked external validation, limiting generalizability, but their approach is designed to be adaptable and robust across different healthcare settings.

**5. NLP for Emergency Department Triage: Broader Evidence for Free-Text Classification**

A 2023 narrative review by Stewart et al. assessed how NLP has been applied to free-text data acquired at Emergency Department (ED) triage. Across 20 studies, NLP models demonstrated high accuracy in predicting triage scores, admissions, and mapping free-text chief complaints to structured fields. Notably, models that incorporated both structured and unstructured data outperformed those using structured data alone.

Despite these promising results, the review highlights that most studies are retrospective and have a high risk of bias. Only one study reported real-world deployment, underlining the need for prospective validation and careful assessment before widespread clinical adoption.

**Conclusions**

Recent research demonstrates that NLP models can accurately classify referral letters by medical specialty, achieving high accuracy (often above 90%) in large-scale, well-structured datasets. These models, using approaches such as SVMs, topic modeling, and transformer-based pipelines, have shown robust performance in both specialty classification and automated triage, with clinical interpretability confirmed by domain experts. However, challenges remain regarding generalizability, data quality, and the need for prospective validation before routine clinical deployment. Incorporating both structured and unstructured data further enhances predictive performance, supporting the integration of NLP into automated triaging systems. Continued research and local optimization are essential to ensure reliability and effectiveness across diverse healthcare settings.

## Pipeline Scalability

This section addresses the design and evaluation of a modular, scalable architecture for synthetic referral generation and NLP-based classification, focusing on its real-world performance, adaptability, and integration potential.

## 1. Synthetic Data Generation with Large Language Models for Personalized Information Retrieval (2024)

This study by researchers exploring Large Language Models (LLMs) for synthetic data generation provides a foundation for building modular and scalable architectures for synthetic referral generation and NLP-based classification. The authors describe a pipeline that begins with collecting user-related information and questions from an existing dataset, then uses LLMs such as GPT-3.5 and Phi-3 to generate synthetic answers tailored to user interests. The process is modular, allowing for the easy integration of different LLMs or prompt techniques as new components.

To evaluate real-world applicability, the authors train neural retrieval models (e.g., DistillBERT) on the synthetic data and assess their performance on human-annotated test sets. Their findings indicate that models trained on LLM-generated synthetic data can outperform traditional methods, even when a significant portion of the synthetic data contains inaccuracies (hallucinations). This demonstrates the pipeline's adaptability, as it can be fine-tuned and evaluated for different downstream tasks, such as information retrieval or classification.

The architecture is scalable because it leverages parallelizable processes, such as batch generation of synthetic data and distributed training of neural models. The integration potential is highlighted by the ability to swap in different LLMs or retrieval models, and the pipeline's compatibility with existing datasets and annotation workflows. The study's manual evaluation of hallucinations and performance metrics provides a framework for assessing both technical performance and adaptability in real-world settings.

## 2. Evaluating Synthetic Data Generation from User Generated Text (2023)

This work introduces an evaluation framework specifically designed for assessing synthetic language data generation pipelines. The authors emphasize the importance of modularity and scalability in architectures that generate and process synthetic text, such as referrals or classification inputs. Their framework supports both rule-based and deep learning-based modules, allowing for flexible integration of different generation and classification components.

A key contribution is the systematic evaluation of pipeline performance across multiple dimensions: generation quality, classification accuracy, adaptability to new domains, and integration with existing workflows. The authors provide metrics for measuring scalability, such as processing speed and resource utilization, and adaptability, such as the ability to incorporate new data sources or update models without retraining the entire pipeline.

The study demonstrates that a modular architecture enables rapid iteration and deployment of new features, making it suitable for real-world applications where requirements and data distributions may change frequently. The evaluation framework also highlights the importance of robust integration testing to ensure that new modules do not degrade overall system performance or introduce biases.

### 3. Text-to-Model Transformation: Natural Language-Based Architectural Generation (2024)

This paper proposes a rule-based architectural generation framework that maps natural language text to executable models, which is particularly relevant for synthetic referral generation and NLP-based classification pipelines. The architecture is modular, with clearly defined interfaces between data generation, preprocessing, feature extraction, and classification modules.

The framework uses heuristic rules and predefined patterns to ensure that each module can be independently developed, tested, and replaced. This modularity supports scalability by allowing parallel development and deployment of individual components. The authors evaluate the pipeline's real-world applicability by measuring its adaptability to new data sources and integration potential with existing systems.

Performance is assessed through end-to-end testing, including accuracy, latency, and resource consumption. The study finds that modular architectures facilitate rapid adaptation to new use cases and integration with legacy systems, making them well-suited for production environments where both performance and maintainability are critical.

### 4. Natural Language Processing for Classification and Clinical Concept Extraction (2024)

This thesis explores the use of AI and NLP to generate a large corpora of richly annotated data for classification tasks, with a focus on clinical and referral contexts. The author describes a scalable pipeline architecture that includes data generation, annotation, preprocessing, and classification modules.

The pipeline is designed to handle large volumes of data efficiently, leveraging distributed processing frameworks such as Spark for parallelization. This approach ensures scalability and high throughput, which are essential for real-world applications. The thesis evaluates the pipeline's performance by measuring classification accuracy, processing speed, and adaptability to new data types.

The modular design enables easy integration of new data sources and classification models, supporting continuous improvement and adaptation. The author also discusses the importance of robust evaluation protocols to ensure that the pipeline remains reliable as it scales and evolves.

### Conclusions

In summary, the reviewed literature demonstrates that modular and scalable architectures are essential for synthetic referral generation and NLP-based classification, as they enable flexible integration of data generation, preprocessing, feature extraction, and classification components. The use of large language models and distributed processing frameworks supports efficient scalability and adaptability to new data sources, while robust evaluation protocols ensure that

these pipelines maintain high performance and real-world applicability. By prioritizing modularity and comprehensive assessment, such systems can effectively address the evolving demands of healthcare and other domains where synthetic data and automated classification are critical.

# Descriptive Statistics of the Dataset

To explore automated, privacy-preserving solutions for referral triaging using natural language processing, this study employed a synthetic healthcare dataset sourced from Kaggle (https://www.kaggle.com/datasets/prasad22/healthcare-dataset). The dataset was explicitly designed to emulate real-world healthcare scenarios while avoiding privacy concerns associated with sensitive patient information. It provides a foundation for testing data-driven approaches to clinical workflow optimization in a realistic yet anonymized context.

You can access the detailed analysis in the public notebook at the following link: https://github.com/CianAllen/TMU_Final_Project/blob/main/healthcare_dataset_analysis.ipynb

## Dataset Overview and Rationale

The dataset consists of 55,500 patient records, each comprising 15 features related to patient demographics, clinical conditions and provider information. These include variables such as age, gender, blood type, diagnosis, admission and discharge dates, medications, test results, admission type, insurance provider, and physician/hospital information. While the dataset is entirely synthetic—generated using the Python Faker library—it replicates the structure, statistical diversity, and temporal granularity found in authentic electronic health records.

Given that access to real patient data is highly restricted due to privacy regulations (e.g., HIPAA, PHIPA), this dataset provides a viable surrogate for developing and evaluating machine learning solutions in referral workflows. The inclusion of both structured (e.g., numerical and categorical variables) and temporal components makes it suitable for simulating referral patterns and investigating relationships between patient characteristics and admission outcomes.

## Preprocessing and Feature Engineering

To enhance the dataset's analytical utility, several preprocessing steps were applied:
1. Irrelevant fields (e.g., patient name, billing amount, and room number) were removed.
2. Datetime conversion was performed on admission and discharge dates to facilitate temporal analysis.

A derived feature, Length of Stay, was created to measure patient hospitalization duration.

Data integrity checks confirmed appropriate variable formats, balanced class distributions, and consistent entry ranges.

This preprocessing ensured the dataset was well-structured for descriptive analytics and served as a clean foundation for downstream NLP and classification pipelines.
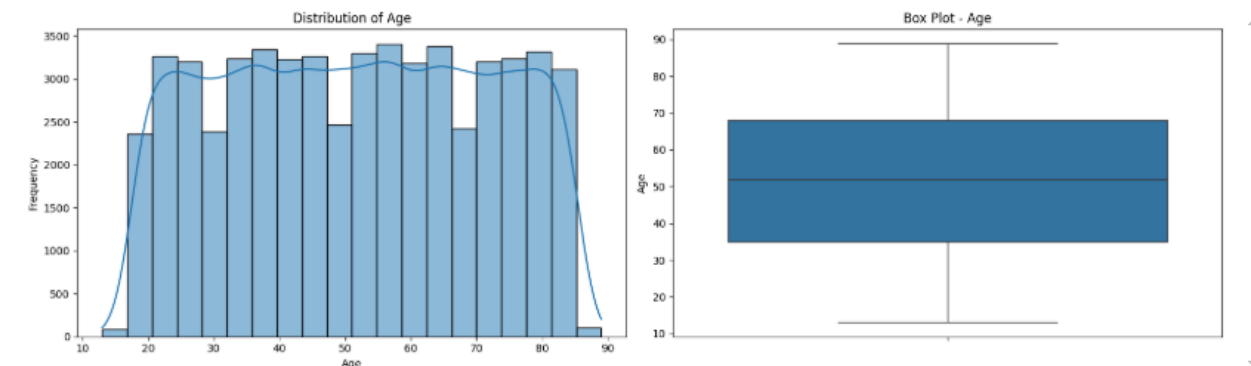
## **Variable Analysis**

The following analysis results were exported in structured JSON format, enabling reproducibility and potential integration into downstream RAG or LLM workflows.

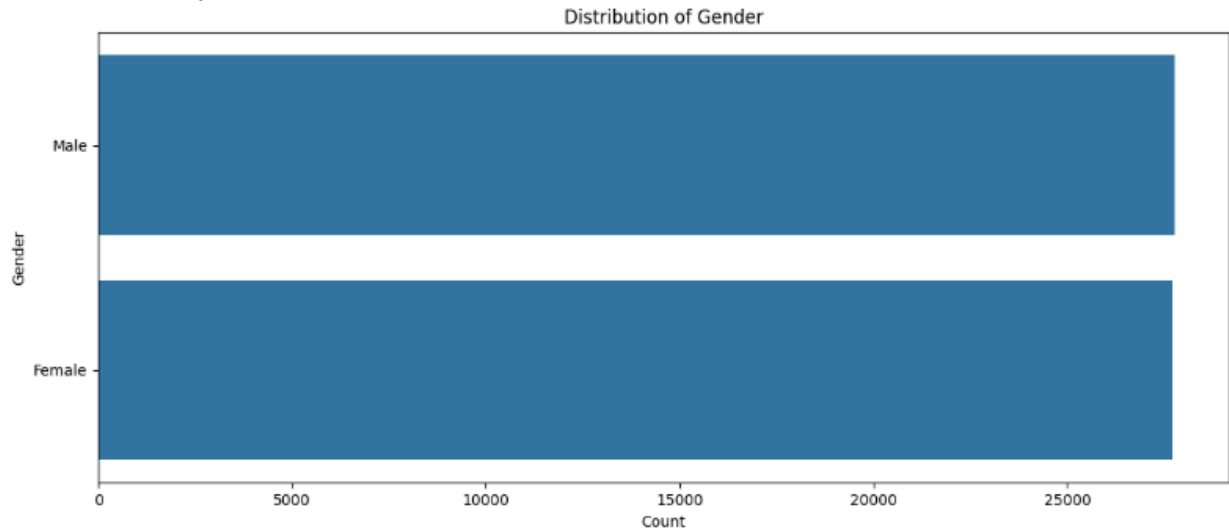Age mean: 51.54 years; Age median: 52 years
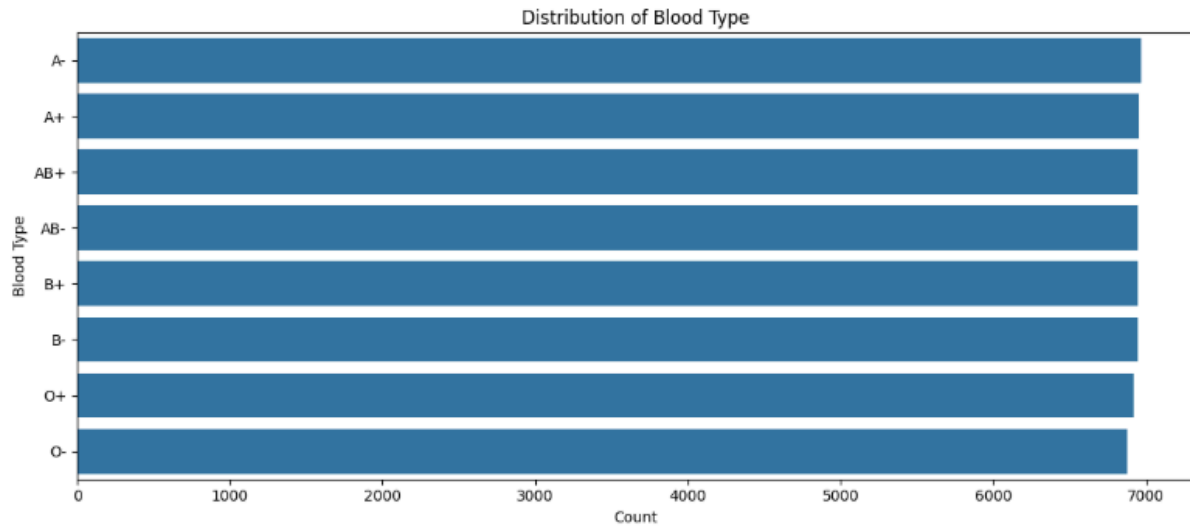
Range: 13 to 89 years

Standard Deviation: 19.60 years

Distribution: Roughly normal with slight skew, indicating a realistic population spread across age groups.



Gender: Equally distributed (50.04% male, 49.96% female)

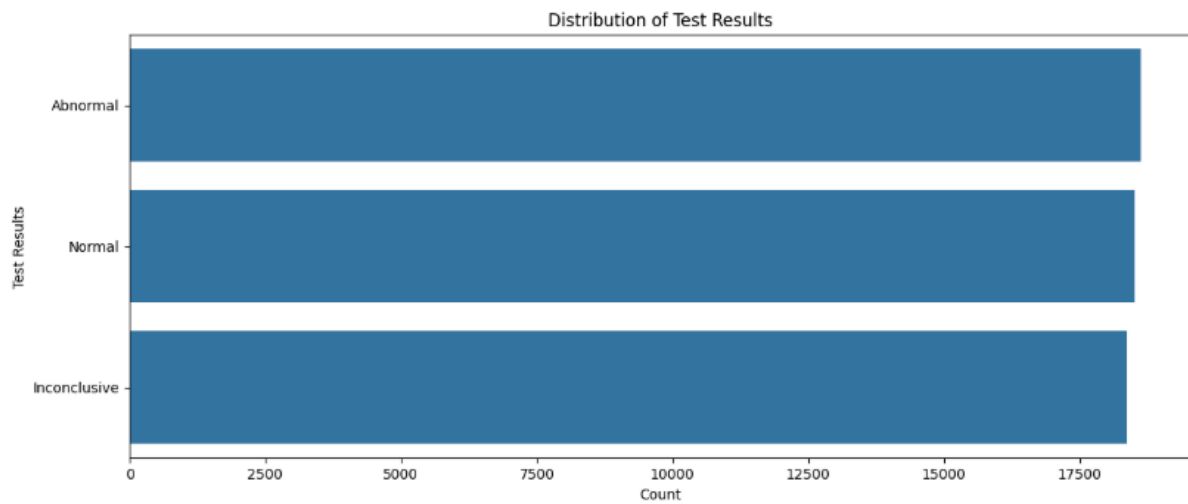Blood Type: Even distribution across all 8 types, no significant outliers



Medical Conditions: 6 equally represented diagnoses (e.g., Diabetes, Asthma)
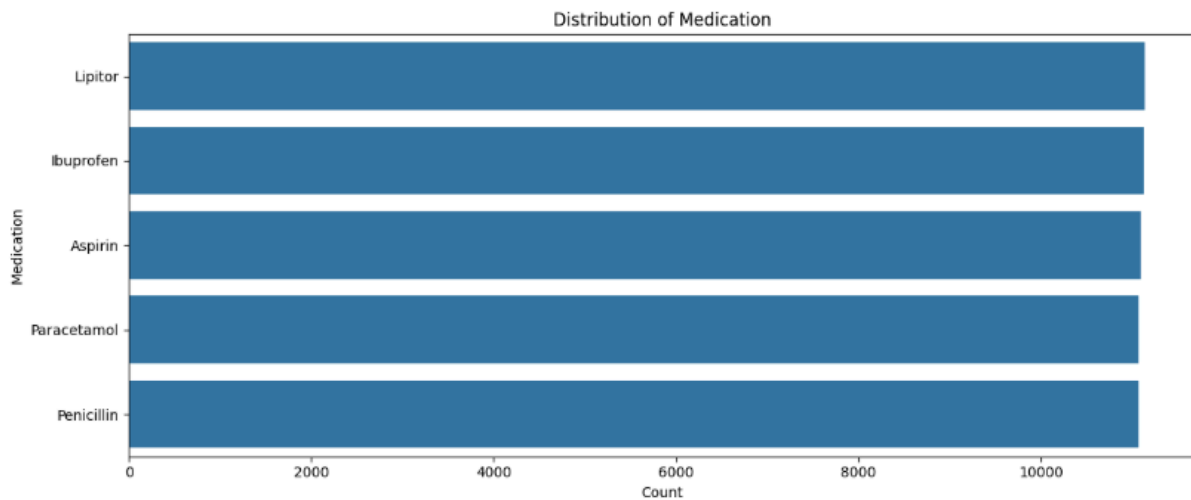


Admission Type: Emergency, Elective, and Urgent admissions all distributed near 33%

Distribution of Admission Type

Test Results: Near-perfect balance between Normal (33.07%), Abnormal (33.56%), and Inconclusive (33.36%)


Distribution of Test Results

Medications: 5 equally prescribed drugs with frequencies ranging from 19.94% to 20.07%
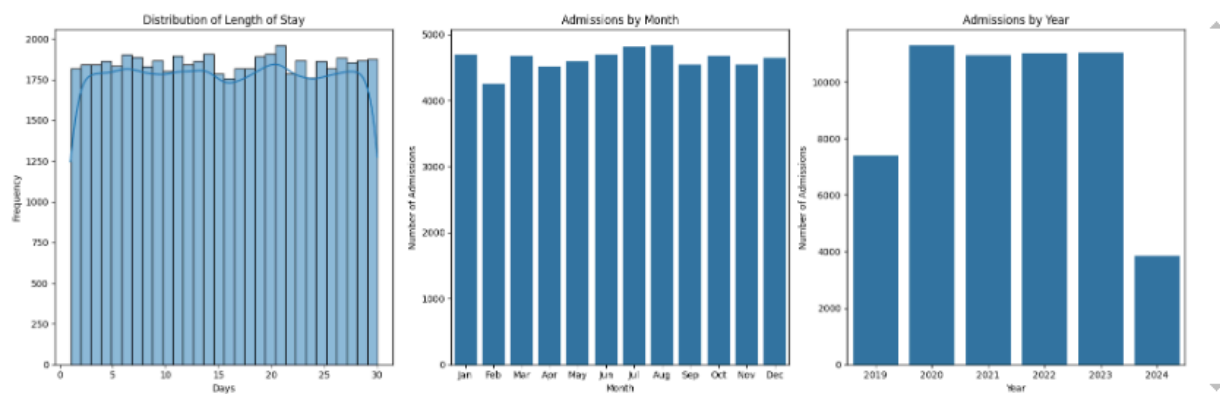
Distribution of Medication

Insurance Providers: Five companies (Cigna, Medicare, UnitedHealthcare, Blue Cross, Aetna), each representing ~20% of the dataset

These balanced distributions reflect the dataset's synthetic design and suggest minimal class imbalance—a favorable property for supervised learning tasks such as test result classification or referral type prediction.

Temporal Patterns and Hospitalization Metrics
Length of Stay (LOS): Average stay was 15.51 days, with a narrow interquartile range (8 to 23 days), and a standard deviation of 8.66.
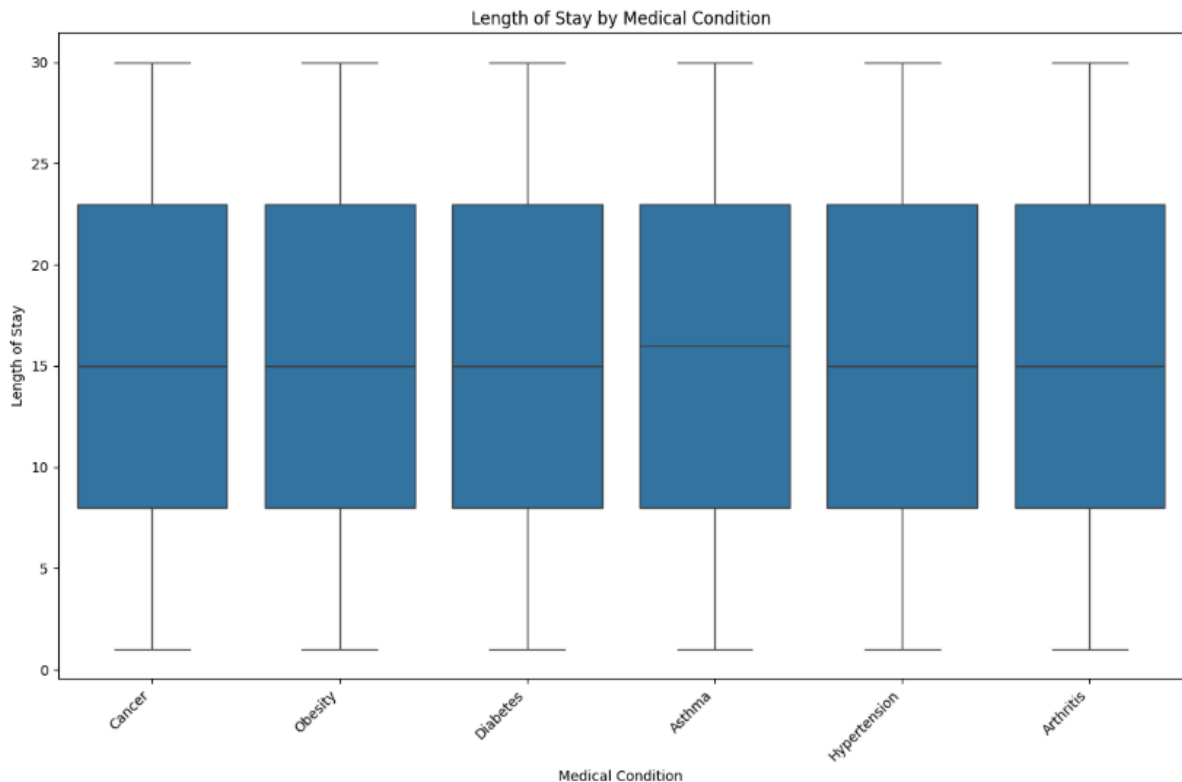


Monthly Admissions: Showed slight seasonal variation, with peak admissions in July and the lowest in February.

Yearly Trends: Admissions rose from 2019 (7,387) to a peak in 2020 (11,285), then stabilized across subsequent years. Data from 2024 was partial.

These temporal trends simulate real-world variability in healthcare utilization and serve as a proxy for modeling seasonal or capacity-based constraints in referral systems.
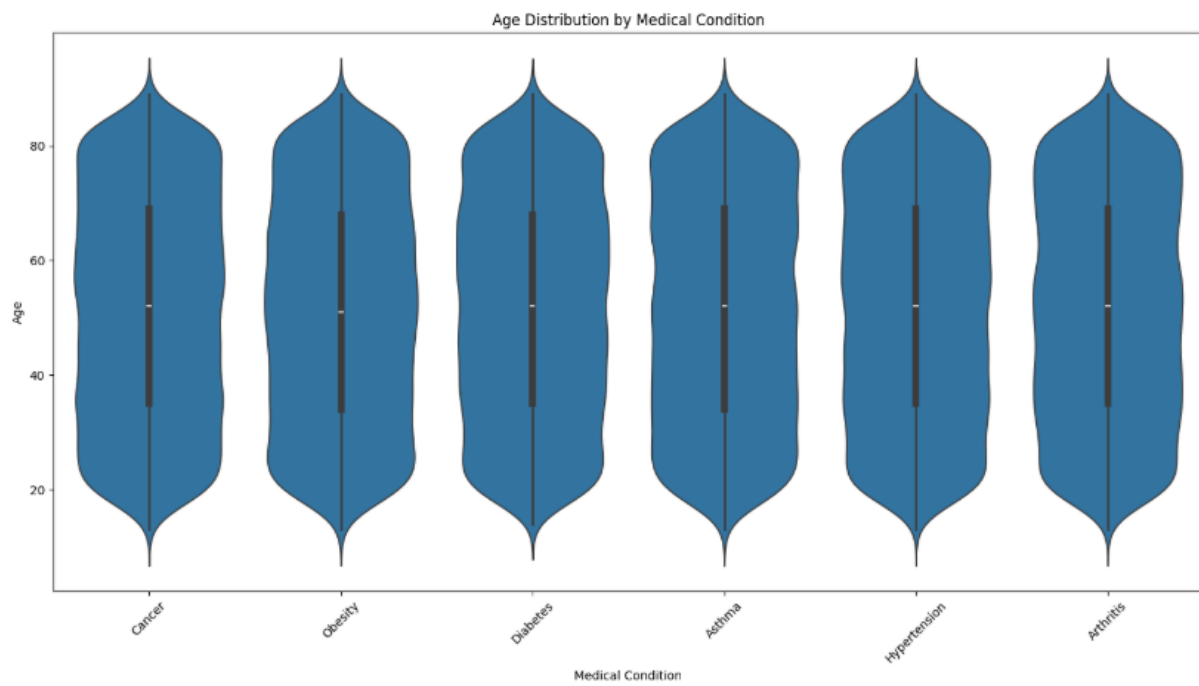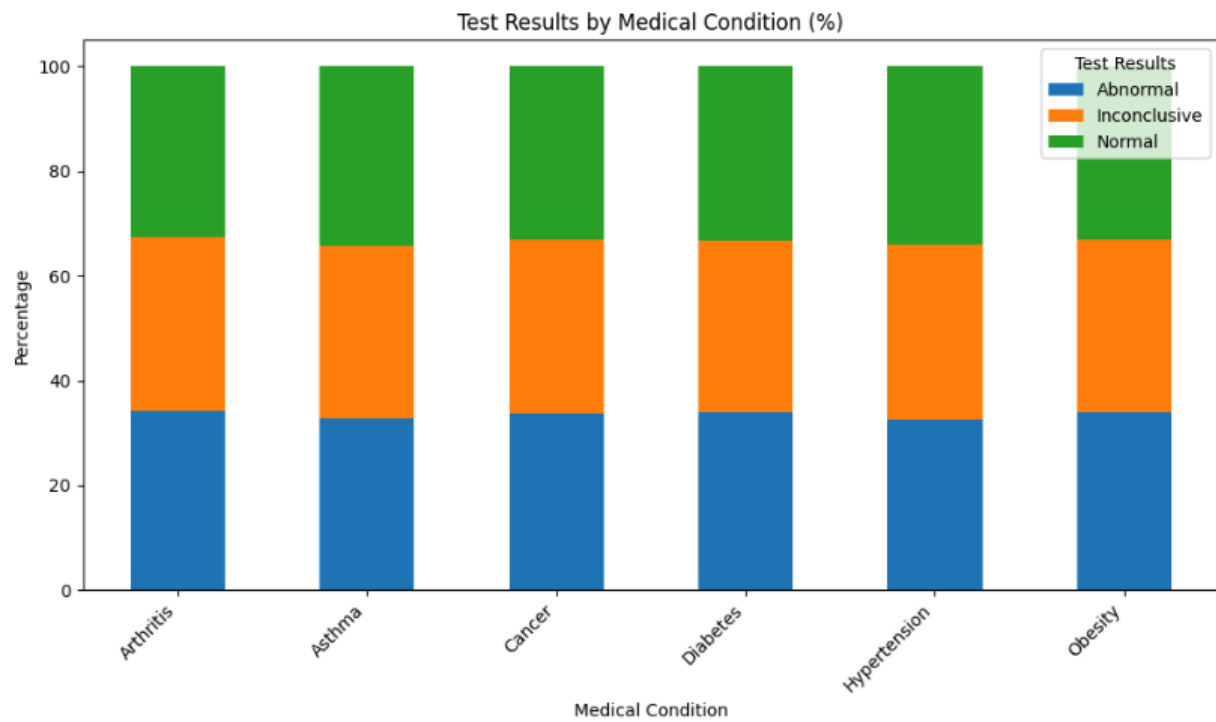
## **Inter-Variable Relationships**

Medical Condition vs. Length of Stay: Minor variation across conditions (15.42–15.70 days), indicating low interdependence.
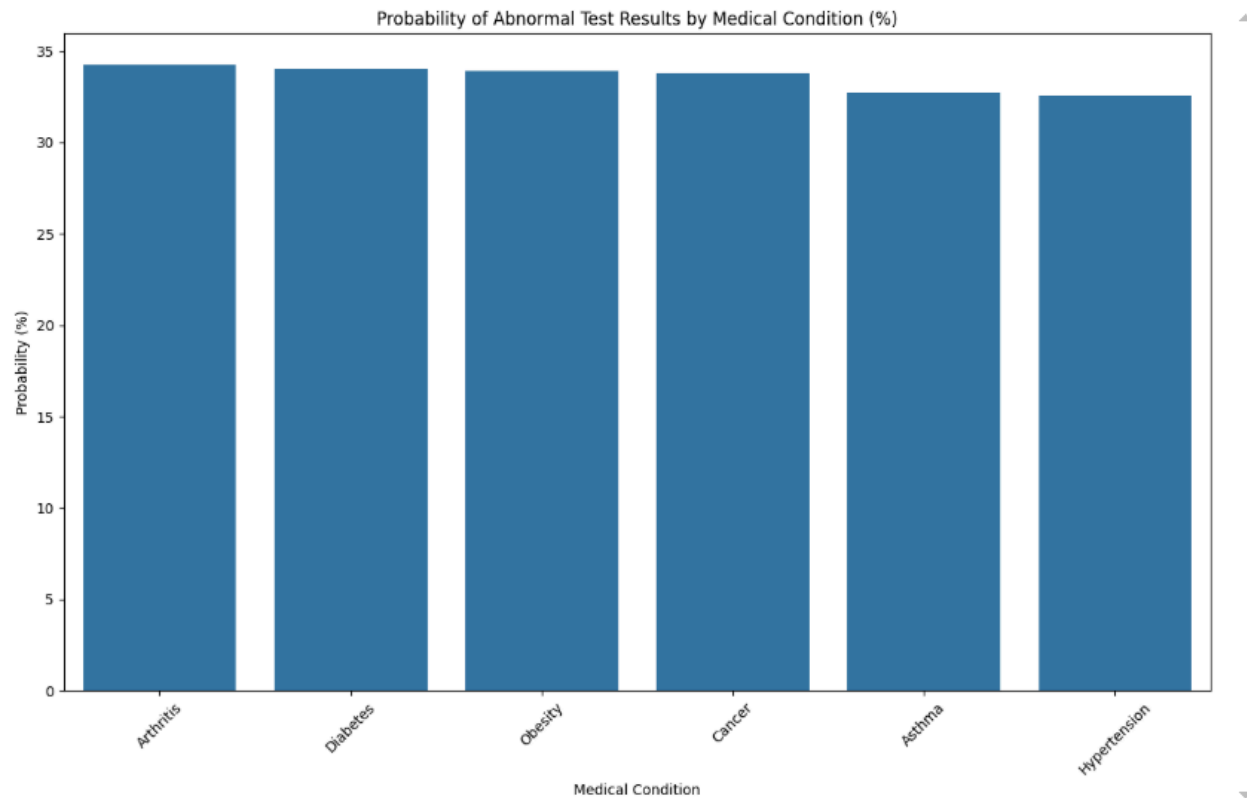


Medical condition was not strongly associated with test results, admission type or age. Each category exhibited a narrow range of variation (~1–2%), confirming the synthetic dataset's
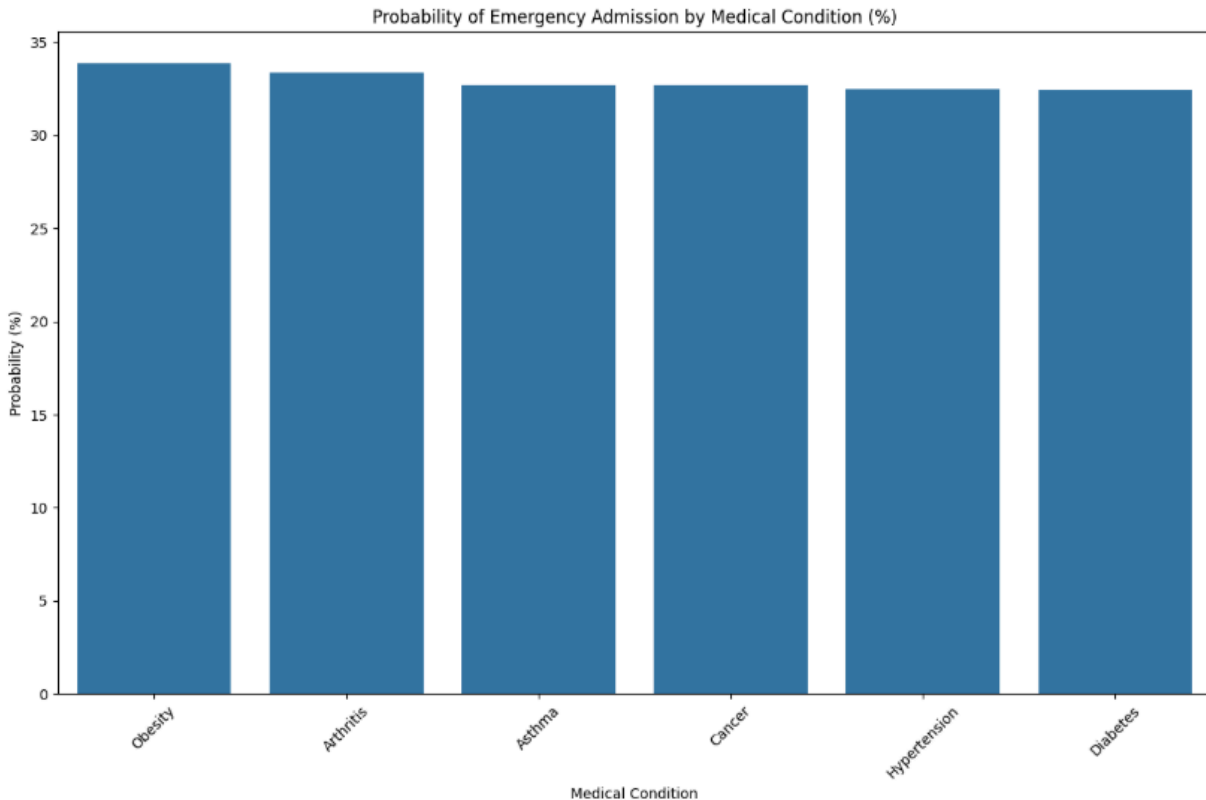
uniformity and lack of confounding effects.



Test Results by Medical Condition (%)



Age Distribution by Medical Condition

# Conditional Probabilities

The likelihood of abnormal test results ranged from 32.58% to 34.25% depending on condition, while the probability of emergency admission ranged from 32.43% to 33.86%.



Probability of Abnormal Test Results by Medical Condition (%)

Probability of Emergency Admission by Medical Condition (%)

These consistent probabilities are desirable in controlled experiments but may limit the dataset's ability to reflect real-world diagnostic disparities.

## Healthcare Provider Landscape

Doctors: 40,341 unique physicians

Hospitals: 39,876 unique facilities

Most Frequent Doctor: Assigned to only 27 patients

Most Frequent Hospital: 44 admissions
This reflects high diversity and extremely low provider-to-patient ratios, supporting exploratory work on decentralized or region-specific triage models.

## Correlation Metrics

Tested for Pearson correlation for numeric pairs, Correlation Ratio ($\eta^2$) for categorical–numeric relationships and Cramér's V for categorical–categorical relationships. All yielded generally weak associations, consistent with the nature of synthetic data lacking hidden dependencies or biases.

## Conclusion of Descriptive Review

This dataset provided a well-balanced, noise-free simulation of real-world healthcare operations. Its design enabled robust exploratory analysis and method validation in a safe, privacy-compliant environment. Although the lack of strong inter-variable correlations and equal distributions indicate synthetic origin, this uniformity is advantageous for initial prototyping of referral triage systems.

The dataset is ideal for simulating multi-class classification, NLP pretraining, or rule-based logic modeling without real-world confounders. Patient profiles are evenly distributed across most features, avoiding class imbalance issues. Temporal trends and conditional probabilities mirror real-world workflows at a high level, albeit without institutional or regional variability.

In the context of this project, the dataset successfully supports the early stages of model development for automating referral triage, offering a consistent, privacy-respecting baseline for testing NLP techniques, classification logic, and clinical workflow automation hypotheses.

# Conclusion

This literature review demonstrates that GAN-based synthetic data generation, NLP classification models, and modular pipeline architectures collectively offer a viable solution for automating referral triaging while preserving patient privacy. The research shows that GANs can generate statistically realistic healthcare data through privacy-preserving techniques, while NLP models achieve over 90% accuracy in classifying referral letters by medical specialty. The integration of structured and unstructured data consistently enhances performance, and modular architectures enable scalable deployment across healthcare systems.

However, significant gaps remain between research achievements and clinical implementation. The literature reveals limited prospective validation, cross-system generalizability challenges, and few real-world deployment studies. While the technical foundations are robust, future work must prioritize external validation, address domain adaptation issues, and develop comprehensive evaluation frameworks that assess clinical utility alongside technical performance. Success will ultimately depend on bridging the gap between promising research outcomes and the complex requirements of healthcare delivery.

# References

## Data Validity and Privacy

1. Arora, A., & Arora, A. (2022). Generative adversarial networks and synthetic patient data: Current challenges and future perspectives. Journal of Medical Systems, 46(8), 57.

2. Synthetic data generation: A privacy-preserving approach to enabling data sharing and analysis in healthcare. (2025). BMC Medical Research Methodology, 25(1), 47.
3. Privacy preserving generative adversarial networks to model health data: A systematic evaluation. (2022). Artificial Intelligence in Medicine, 131, 102349.
4. Synthetic data generation via generative adversarial networks in healthcare: A systematic review. (2024). Scientific Reports, 14(1), 21656.

## Automated Classification

1. Davies, J., et al. (2023). Semantics Based Classification of Medical Referral Letters. Data Science: Journal of Computing and Applied Informatics, 7(1), 24-34.
2. Chapman, S. J., et al. (2020). Patient Triage by Topic Modeling of Referral Letters: Feasibility Study. JMIR Medical Informatics, 8(11), e20950.
3. Fudickar, S., Bantel, C., Spieker, J., Töpfer, H., Stegeman, P., Schiphorst Preuper, H. R., Reneman, M. F., Wolff, A. P., & Soer, R. (2024). Natural Language Processing of Referral Letters for Machine Learning–Based Decision Support in Low Back Pain Triage. Journal of Medical Internet Research, 26, e54321.
4. Van der Laan, W., et al. (2025). Improving musculoskeletal care with AI enhanced triage through natural language processing of referral letters. npj Digital Medicine, 8, Article 1495.
5. Stewart, J., et al. (2023). Applications of natural language processing at emergency department triage: A narrative review. BMC Medical Informatics and Decision Making, 23, Article 321.

## Pipeline Scalability

1. arXiv. (2024, October 29). Synthetic Data Generation with Large Language Models for Personalized Community Question Answering. arXiv:2410.22182.
2. Evaluating Synthetic Data Generation from User Generated Text. (2023). Computational Linguistics, 51(1), 191–209.
3. MDPI. (2024). Text-to-Model Transformation: Natural Language-Based Architectural Generation Framework. MDPI Information, 12(9), 369.
4. Maastricht University. (2024, January 1). Natural Language Processing for Classification and Clinical Concept Extraction (Doctoral dissertation).