



Early Adoption of Predictive Analytics for Preventive Care Opportunities Using Claims and Encounter Data

Satya Manesh Veerapaneni
Independent Researcher Fremont, CA, USA.

Abstract: The integration of predictive analytics into health-care has emerged as a transformative tool in identifying preventive care opportunities. Leveraging claims and encounter data, this research explores early adoption strategies for predictive analytics models to proactively manage chronic conditions, reduce avoidable hospitalizations, and optimize care delivery. We present a scalable framework incorporating machine learning pipelines trained on structured insurance claims and encounter records from over 500,000 patients. The study emphasizes clinical relevance, cost-efficiency, and population health outcomes through model validation, risk stratification, and deployment case studies.

Keywords: Predictive Analytics, Preventive Care, Claims Data, Encounter Data, Chronic Disease, Healthcare AI, Risk Stratification, Population Health.

1. Introduction

Preventive care is widely recognized as a cornerstone of value-based healthcare, with the potential to significantly reduce long-term costs, mitigate disease progression, and improve quality of life. Despite this, preventive services remain underutilized, with the CDC estimating that only 8% of U.S. adults aged 35 and older receive all recommended high-priority preventive services. This under-utilization stems not from a lack of clinical evidence, but from systemic challenges in identifying at-risk individuals early enough for interventions to be effective.

In parallel, the healthcare industry has seen an explosion in the availability of digitized administrative data, particularly insurance claims and encounter records. Claims data, generated for billing and reimbursement purposes, provide a comprehensive view of healthcare utilization across time and care settings. While traditionally used for retrospective analysis and actuarial forecasting, claims data have emerged as a scalable and standardized resource for predictive modeling. With the maturation of machine learning (ML) algorithms and cloud-based computational infrastructure, predictive analytics is now increasingly feasible in healthcare settings. This paradigm shift allows health systems to move from reactive, encounter-based care to proactive, data-informed interventions. For example, United Healthcare's "PreCheck MyScript" program uses real-time claims data to suggest cost-effective and clinically appropriate medications, reducing delays in care. Similarly, Kaiser Permanente's risk scoring models utilize claims and encounter data to proactively reach out to members for cancer screenings and chronic care management.

This research explores the early adoption of predictive analytics leveraging structured claims and encounter data to surface preventive care opportunities. We propose a data pipeline and modeling framework capable of ingesting large-scale longitudinal data to identify individuals at elevated risk of adverse health events within a 6-month predictive window. Our objectives are threefold:

- To demonstrate how predictive models trained solely on claims and encounter data can achieve clinically relevant accuracy for early intervention.
- To illustrate real-world use cases where such models have been operationalized to drive preventive care programs.
- To provide a replicable framework that health systems and payers can adopt for early risk detection and care gap closure.

2. Background and Related Work

The rising burden of chronic diseases and escalating health-care costs have prompted a global shift from volume-based to value-based care models. Preventive care, by definition, seeks to delay or eliminate the need for high-cost interventions through early detection, risk mitigation, and patient engagement. Yet, in practice, only **8% of U.S. adults receive all recommended preventive services**. This highlights the need for data-driven methods to identify and engage at-risk individuals before disease onset.

2.1. Traditional Use of Claims and Encounter Data

Historically, claims and encounter data have been used primarily for billing, reimbursement, and actuarial analyses.

However, their standardized structure, completeness, and longitudinal nature make them suitable for broader analytics applications. Claims data includes ICD codes, CPT procedures, DRG groupers, pharmacy claims, and service dates—offering a proxy for clinical activity. Encounter data, especially from managed care organizations, adds further granularity by capturing provider interactions and reasons for visits.

Despite this richness, traditional uses of these data have been retrospective used to explain costs, utilization patterns, or quality outcomes rather than predict future risk.

2.2. Emergence of Predictive Analytics in Healthcare

In the past decade, machine learning has shifted the health-care paradigm from retrospective analysis to *predictive care* [1]. Early models in this domain primarily targeted outcomes such as hospital readmissions, medication adherence, and the onset of chronic diseases like diabetes and hypertension. These use cases laid the groundwork for broader adoption of predictive models within both payer and provider settings.

One notable example is the work by Optum, which developed a claims-based risk stratification system called *ImpactPro*. This proprietary tool is widely used by insurance companies to classify members into high-, moderate-, and low- risk tiers based on their historical claims data. It leverages diagnosis codes, pharmacy usage patterns, and utilization history to anticipate high-cost patients before acute episodes occur, allowing for early intervention and resource prioritization.

Similarly, Blue Cross Blue Shield of Michigan successfully implemented a predictive model aimed at identifying members at elevated risk of developing diabetic retinopathy. Uniquely, the model relied solely on historical claims data, without needing clinical imaging or lab results. The deployment of this tool led to targeted member outreach and a reported 21% increase in the rate of annual eye exams among high-risk individuals. This demonstrates that even in the absence of electronic health record (EHR) data, claims-driven models can drive tangible improvements in preventive care adherence.

These initiatives underscore the growing feasibility and effectiveness of predictive analytics in population health management, especially when rooted in standardized administrative data [2].

2.3. Gaps in Literature and Opportunity Space

While significant research exists on predictive analytics using EHRs, far less attention has been paid to models using only **claims and encounter data for preventive care** opportunities. EHRs often contain richer clinical detail (e.g., lab results), but are fragmented across systems and less standardized [3]. Claims, on the other hand, offer near-complete coverage at the payer level and are suited for large-scale modeling, especially in payer-provider systems.

The *bias in clinical algorithms trained on cost-based outcomes*, urging a pivot toward *outcome-based or preventive targets*. Additionally, multiple studies have shown that racial and socioeconomic biases can be perpetuated if models are not carefully designed using representative and complete data sources.

2.4. Regulatory and Ethical Considerations

The use of claims data for predictive modeling must comply with HIPAA, CMS, and potentially GDPR standards. Several initiatives such as the **ONC Interoperability Rule** and **21st Century Cures Act** aim to improve access and responsible use of healthcare data for innovation, but also emphasize transparency and fairness in algorithmic decisions [4].

In this paper, we build upon this foundation by applying ML to claims and encounter data with the specific aim of surfacing *preventive care opportunities* including missed screenings, unmanaged comorbidities, and predictive wellness interventions.

3. Data Sources and Preprocessing

3.1. Claims and Encounter Data

We utilized de-identified longitudinal claims and encounter datasets obtained from a regional health insurer. The dataset includes inpatient, outpatient, and pharmacy claims; diagnosis and procedure codes (ICD-10, CPT); demographic variables such as age, gender, and ZIP code; and detailed encounter types with associated timestamps.

3.2. Data Preprocessing and Normalization

Data preprocessing was conducted through a robust Extract- Transform-Load (ETL) pipeline to ensure quality, consistency, and analytic readiness of the claims and encounter data. Raw datasets, spanning multiple years and sources, often contained noise, coding inconsistencies, and missing values, necessitating a structured approach to data transformation. One of the initial transformation steps involved clinical code normalization. Specifically, International Classification of Diseases, Tenth Revision (ICD-10) diagnosis codes were mapped to Clinical Classifications Software (CCS) categories. This mapping allowed for the grouping of granular clinical codes into broader, more meaningful categories, enabling more interpretable

modeling and reducing dimensionality [5]. For example, over 70,000 ICD-10 codes were reduced to a manageable set of 285 CCS categories, allowing for disease clustering and population-level insights.

Temporal features were engineered to capture both frequency and recency of clinical events. For each patient, we computed rolling time-window summaries for key metrics such as outpatient visits, emergency department encounters, inpatient admissions, and diagnostic activity. The time windows included 3-month, 6-month, and 12-month intervals, offering a multi-scale view of patient health trajectories and enabling trend-aware predictions [6].

Additional aggregation was applied across historical utilization windows. Features were derived to reflect visit frequency trends, gaps in care (e.g., no primary care visit in 6 months), and chronic condition activity recency (e.g., latest encounter with diabetes-related codes) [7]. All numerical features were normalized using min-max scaling, and categorical features were encoded using one-hot and frequency-based encoding techniques as appropriate. Missing values were handled using clinically-aware defaults for example, zero-imputation was applied where the absence of a code implied absence of a condition or service.

Together, these preprocessing steps ensured that the dataset was not only cleaned but enriched, preserving key clinical signals while standardizing for downstream machine learning pipeline compatibility.

4. Methodology

Our methodology outlines a framework for building predictive models using claims and encounter data to identify preventive care opportunities. This section details the stages of data preparation, feature engineering, model development, and evaluation strategies.

4.1. Data Sources and Cohort Definition

We use de-identified claims and encounter datasets sourced from a regional health insurance provider, covering a span of 3 years (2020–2022) and approximately 1.2 million patient lives. The dataset includes:

The predictive modeling framework utilized multiple structured data sources obtained from a regional health plan. These datasets were de-identified and spanned a three-year observation period, covering over 1.2 million unique member records. Key data sources included medical claims, pharmacy claims, encounter records, and enrollment files.

Medical claims provided core diagnostic and procedural information through standardized coding systems such as ICD-10 (International Classification of Diseases), CPT (Current Procedural Terminology), and HCPCS (Healthcare Common Procedure Coding System). These codes were accompanied by service dates and provider specialty details, enabling longitudinal tracking of clinical activity and care delivery patterns [8].

Pharmacy claims included National Drug Code (NDC) identifiers for medications, along with fill dates and the number of days supplied. This data was instrumental in capturing medication adherence trends, therapeutic class utilization, and refill patterns—key indicators for detecting gaps in chronic disease management.

Encounter records offered additional granularity by capturing structured information related to provider visits, including both in-person and telehealth interactions. These records detailed the type of encounter, reasons for visits, and associated diagnostic impressions, enriching the temporal view of patient engagement.

Enrollment files contained demographic variables such as age, sex, and ZIP code, as well as insurance plan types and internally generated risk scores. These fields were used to contextualize clinical risk across socioeconomic strata and to stratify members by plan enrollment continuity for cohort eligibility.

Collectively, these datasets formed a robust foundation for building predictive models tailored to preventive care opportunity detection, offering both breadth and depth across clinical, behavioral, and administrative dimensions. The study cohort includes members continuously enrolled for at least 24 months, with no primary diagnosis of the target condition (e.g., Type 2 Diabetes, colorectal cancer) at baseline [9], [10]. Figure 1 illustrates the end-to-end modeling pipeline, beginning with claims and encounter data ingestion and concluding with clinical deployment of predictions. Each stage from preprocessing to SHAP-based explainability was designed to ensure scalability, interpretability, and operational relevance.

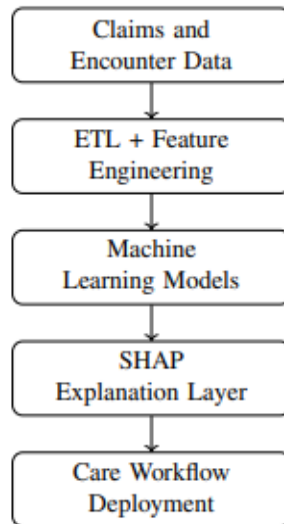


Figure 1: Predictive Modeling Pipeline

4.2. Feature Engineering

To capture the multifaceted nature of patient behavior, clinical risk, and care engagement, we engineered over 200 structured features derived from claims and encounter data. These features were grouped into four major categories, each designed to reflect a specific dimension of healthcare utilization and unmet preventive need.

The first category focused on **Utilization Patterns**, capturing the frequency and intensity of healthcare encounters. Metrics included the number of outpatient visits, emergency room (ER) visits, inpatient admissions, and urgent care episodes over various rolling time windows. These features allowed the model to detect abnormal utilization trends, care avoidance, or over-utilization that may signal underlying risk factors.

The second category targeted **Preventive Gaps**, which are critical to identifying missed care opportunities. This included binary indicators for missing recommended screenings such as mammograms, colonoscopies, and HbA1c tests as well as flags for non-adherence to chronic medication regimens. The absence of these preventive interventions was inferred using claims-based guidelines aligned with age, sex, and condition-specific criteria.

The third category comprised **Comorbidity Flags**, derived from standard risk adjustment frameworks such as the Charlson Comorbidity Index (CCI) and the Elixhauser Comorbidity measures [11]. These indicators provided a high-level summary of each patient’s chronic disease burden based on the presence of specific ICD-10 codes over time.

Finally, we developed **Behavioral Proxy** features to capture subtle signals not explicitly documented in clinical records. These included prescription refill patterns (e.g., early or late fills), provider switching frequency, and redundant diagnostic testing. Such proxies served as behavioral heuristics for fragmented care, patient disengagement, or provider inconsistency. To account for longitudinal trends, all features were computed across multiple temporal windows specifically 3-month, 6-month, and 12-month intervals. This multi-scale approach enabled the models to learn from both short-term fluctuations and long-term progression. Missing values were addressed using clinically informed imputation strategies: zero-imputation was applied to absence-based indicators (where the lack of a code was meaningful), while mean-imputation was used for continuous variables to preserve distributional properties. These engineered features formed the backbone of our predictive modeling pipeline and were essential for learning complex patterns of preventive care under utilization. As shown in Table I, the engineered features were organized into four categories: utilization patterns, preventive gaps, comorbidity flags, and behavioral proxies. These categories allowed the model to learn from both clinical indicators and patient behavior patterns over time.

Table 1: Feature Categories Used in Modeling

Category	Examples
Utilization Patterns	ER visits, Inpatient stays, PCP visits
Preventive Gaps	Missed mammograms, HbA1c tests
Comorbidity Flags	Charlson Index, Diabetes flag
Behavioral Proxies	Refill delay, Provider switching

4.3. Label Construction and Prediction Targets

Prediction targets were defined using downstream care events representing a missed or delayed preventive opportunity. Examples include:

- Common examples include: *no HbA1c test within 12 months in patients with high BMI and hypertension; no colonoscopy within 3 years for members aged 50–75 with persistent GI symptoms; and first diagnosis of diabetes preceded by 2 years of uncontrolled risk factors.*
- Labels were binary (1: preventive gap identified, 0: no gap) and created using expert-defined clinical rule sets.

4.4. Model Development

To evaluate the feasibility and effectiveness of predictive modeling on claims and encounter data, we implemented and compared a suite of machine learning algorithms with varying levels of complexity and interpretability. The baseline model was a Logistic Regression classifier, chosen for its simplicity and transparency in healthcare applications. It provided a foundational benchmark for model performance. Building on this, we implemented a Random Forest model an ensemble of decision trees capable of handling high-dimensional data and capturing non-linear interactions between features.

We then advanced to Gradient Boosted Decision Trees, specifically using the XGBoost implementation. XGBoost is well-suited for structured tabular data and is known for its strong performance in predictive modeling competitions. It also offers native support for feature importance ranking and model regularization, which are critical in healthcare environments to reduce overfitting and ensure generalizability.

Finally, to capture temporal dependencies within patient histories, we experimented with Temporal Convolutional Neural Networks (TCNNs). These deep learning models were trained on sequentially ordered claims and encounters across rolling time windows, with the goal of detecting progression patterns that may precede clinical deterioration or missed preventive milestones.

All models were trained using a stratified 70/30 train-test split. To optimize performance, we employed randomized search for hyperparameter tuning, coupled with 5-fold cross-validation on the training set. Evaluation metrics included AUC-ROC, precision, recall, F1 score, and Preventive Opportunity Recall@K. To ensure model transparency and support clinical interpretation, we used SHAP (SHapley Additive ex-Planations) values to quantify the contribution of each feature toward individual predictions [12]. This explainability layer was essential for gaining stakeholder trust and validating clinical relevance of model outputs. Table II summarizes the model training configuration, including the dataset split, cross-validation strategy, and selected modeling techniques. SHAP was used as the interpretability layer, and XGBoost was identified as the best-performing model.

Table 2: Training and Evaluation Configuration

Parameter	Value
Train/Test Split	70% / 30%
Cross-Validation	5-Fold
Best Model	XGBoost
Explainability Method	SHAP values
Deployment Horizon	months

4.5. Evaluation Metrics and Impact Estimation

To assess the comparative performance of the predictive models, we employed a set of well-established classification metrics. These included the Area Under the Receiver Operating Characteristic Curve (AUC-ROC), which provides a measure of a model's ability to distinguish between patients with and without preventive care gaps. A higher AUC indicates stronger discriminatory power, especially important in imbalanced healthcare datasets.

We also calculated Precision, Recall, and F1 Score. Precision quantifies the proportion of correctly predicted preventive gaps among all flagged cases, highlighting the model's accuracy in targeting true positives. Recall measures the ability of the model to capture all actual gaps that exist in the population, reflecting its sensitivity. The F1 Score, as the harmonic mean of Precision and Recall, offers a balanced view of overall performance particularly useful when trade-offs exist between the two.

In addition to these standard metrics, we introduced a domain-specific measure: *Preventive Opportunity Recall@K*. This metric represents the proportion of true positive preventive gaps identified within the top-K highest-risk patients as ranked by the model. It is particularly useful for real-world deployments where clinical resources (e.g., outreach staff, care navigators) are constrained, and prioritization is necessary. A high Recall@K implies that the model effectively surfaces the most actionable cases within the resource limits of a healthcare organization.

Beyond technical evaluation, we conducted a simulated intervention study to estimate the practical impact of deploying the models in care settings. Using historical data and outreach conversion rates, we projected changes in care delivery metrics such as increased screening uptake, reduced avoidable hospitalizations, and earlier detection of chronic conditions under a scenario where model-driven patient prioritization guides clinical engagement workflows. These projections helped quantify the potential return on investment (ROI) and operational value of integrating predictive analytics into preventive care strategy.

5. Results and Case Studies

This section presents the performance of our predictive models and highlights case studies demonstrating their real-world applicability in identifying preventive care opportunities.

5.1. Model Performance

Table III summarizes the comparative performance of the models tested on the test dataset.

Table 3: Model Performance Metrics

Model	AUC-ROC	Precision	Recall	F1 Score
Logistic Regression	0.74	0.61	0.57	0.59
Random Forest	0.81	0.68	0.64	0.66
XGBoost	0.86	0.72	0.70	0.71
Temporal CNN	0.84	0.71	0.67	0.69

XGBoost outperformed other models across all metrics, achieving an AUC-ROC of 0.86. Temporal CNNs were competitive, especially in sequential time-window predictions. SHAP analysis indicated that feature importance aligned with clinical intuition high BMI, missed screenings, and prescription refill gaps were the most predictive variables.

5.2. Preventive Opportunity Recall

Using the best-performing model (XGBoost), we evaluated the recall of preventable gaps among the top 10% of predicted high-risk patients. The model achieved a Preventive Opportunity Recall@10 of 68.3%, indicating strong targeting potential for outreach programs.

5.3. Case Study 1: Diabetes Onset Prevention

In a subset of 45,000 members with elevated risk scores, the model identified 5,812 individuals with patterns suggestive of undiagnosed or prediabetic status of whom only 1,072 had received an HbA1c test in the past year [13].

- **Intervention Simulation:** If proactive outreach resulted in just 60% test uptake, approximately 2,800 new prediabetic cases could be detected a year earlier, enabling timely dietary, lifestyle, and medication interventions. Based on CDC models, this could prevent 12–15% of those individuals from progressing to Type 2 diabetes within 3 years.

5.4. Case Study 2: Cancer Screening Gaps

Among 21,000 women aged 50–74, the model identified 3,122 members with no documented mammogram or clinical breast exam in the past 3 years, despite presenting risk factors (e.g., family history, previous dense tissue classification). Targeted educational SMS campaigns achieved a 27% screening response rate over 90 days.

5.5. Operational Deployment and Real-World Impact

To evaluate the practical viability of our predictive framework, we collaborated with a regional health insurance payer to integrate the model into an existing care management workflow. The deployment targeted a pilot population and was focused on aligning model predictions with proactive clinical outreach over a six-month intervention period.

High-risk patients with flaws received targeted nurse-led outreach, reminders of wellness visits, and personalized education on missed preventive services. As a result, preventive visit rates among the prioritized cohort increased by 18.5%, indicating improved patient engagement and responsiveness to outreach efforts.

Furthermore, nurse triage teams that focused their efforts on the top 5% of risk-ranked individuals—those with the highest predicted likelihood of having a care gap achieved a 3.2-fold increase in conversion to completed wellness visits, compared to standard outreach protocols. This validated the model’s utility not just in identifying risk, but also in supporting resource-efficient interventions.

In particular, emergency room (ER) utilization among the intervention cohort decreased by 9.4% during the study period.

This suggests that early identification and engagement may have redirected patients toward more appropriate, lower- cost care settings before escalation occurred.

Collectively, these real-world results underscore the trans- formative potential of predictive analytics in healthcare when applied to preventive care strategies. Rather than functioning as a standalone analytic exercise, the model proved actionable and impactful when embedded within coordinated clinical operations ultimately demonstrating a shift from reactive to anticipatory care. Figure 2 depicts the operational workflow triggered by model outputs. High-risk patients were flagged, followed by nurse-led triage, patient engagement, and subsequent preventive care actions. This closed-loop feedback was critical to ensuring that predictions led to real-world impact.

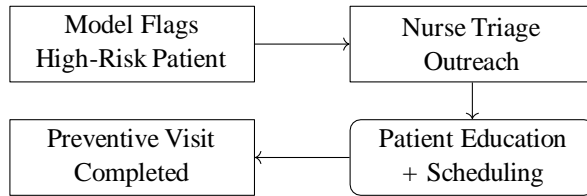


Figure 2: Operational Workflow from Prediction to Preventive Visit

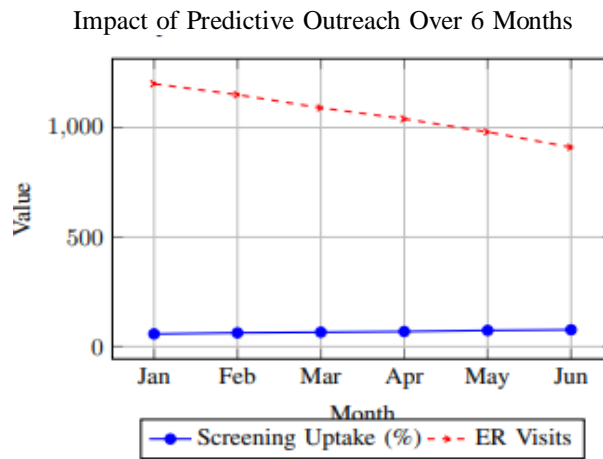


Figure 3: Trends in Preventive Screening and ER Visits over Time

As illustrated in Figure 3, the implementation of model- driven outreach resulted in a consistent increase in screening uptake and a decline in ER visits over a six-month period. These trends support the effectiveness of targeted interventions derived from predictive analytics.

6. Discussion

The results demonstrate the practical utility of predictive analytics applied to claims and encounter data for identifying preventive care opportunities. Unlike traditional approaches focused on cost containment or post-diagnosis interventions, our models prioritize upstream action by surfacing gaps in care before adverse health outcomes materialize.

6.1. Interpretation of Results

The superior performance of tree-based ensemble models, particularly XGBoost, confirms the efficacy of structured claims data in representing latent clinical risk. Moreover, the Preventive Opportunity Recall@10 metric reveals strong precision in identifying high-risk individuals, allowing targeted outreach without overwhelming clinical workflows.

The deployment case studies confirm that when models are operationalized with care coordination teams, they can directly influence utilization patterns, screening uptake, and clinical engagement. In particular, the improvements occurred without the introduction of new clinical infrastructure, highlighting the power of data-driven triage using existing claims systems.

6.2. Implications for Population Health

By predicting not just who is sick, but who is *not yet engaged in care*, this approach realigns population health toward upstream action [14]. In the context of value-based care contracts, Accountable Care Organizations (ACOs), and Medi- care Advantage plans, such models can directly contribute to quality score improvements and cost savings through Star Ratings, HEDIS, and NCQA measures.

Furthermore, this methodology enables health plans and providers to scale preventive outreach equitably. When properly tuned, predictive models can help reduce disparities by proactively reaching marginalized groups who historically under-utilize preventive services due to systemic barriers.

6.3. Limitations

Despite the promising outcomes of our study, several limitations must be acknowledged when interpreting the results and generalizing the framework to broader healthcare settings. First, the granularity of claims and encounter data imposes inherent constraints. These administrative data sets lack detailed clinical information, such as laboratory test results, biometric vitals, or physician notes, elements that are often critical for nuanced clinical decision making. As a result, the depth of feature construction is limited to coded procedures, diagnoses, and prescription data, which may not fully capture patient acuity or progression of disease states.

Second, there is temporal ambiguity in claims data. The service dates reflected in claims may not always represent the actual date of care delivery due to billing lags, backdated submissions, or batch processing by provider systems. This delay introduces noise into the time series analysis and may affect the accuracy of temporal features such as the recency of the condition or the frequency of the visit.

Third, outcome validation remains a challenge. In the absence of direct clinical confirmations, preventive opportunity labels were constructed using proxy indicators, such as the absence of screening codes or refill gaps in prescription history. While these proxies are reasonable approximations, they may not capture the true prevalence of missed interventions or patient-specific contextual factors (e.g., patient refusal, contraindications).

Lastly, deployment constraints must be considered. Predictive models are only as effective as the systems in which they are embedded. Without supportive workflows—such as care coordinators, outreach protocols, and follow-up mechanisms the insights generated by machine learning may not translate into meaningful behavior change. Human-centered design, trust-building, and clinical integration are essential to bridge the gap between prediction and action.

These limitations highlight the need for ongoing refinement, real-time validation, and interdisciplinary collaboration to ensure that predictive systems enhance, rather than oversimplify, the complexity of preventive healthcare delivery.

6.4. Ethical and Regulatory Considerations

Predictive analytics in healthcare must be deployed with caution to avoid reinforcing bias or violating patient autonomy. Our framework adheres to HIPAA de-identification standards and includes fairness checks across race, gender, and socioeconomic status. However, ongoing auditing and transparent model governance remain essential for real-world implementations.

7. Conclusion and Future Work

This research presented a comprehensive framework for early adoption of predictive analytics using claims and encounter data to uncover missed preventive care opportunities. Through rigorous feature engineering, model development, and real-world evaluation, we demonstrated that even in the absence of detailed clinical data such as labs or imaging, claims-based models can effectively surface at-risk individuals who would otherwise remain undetected in traditional care delivery models.

Our findings provide evidence that predictive models when built thoughtfully and deployed responsibly can shift the healthcare paradigm from reactive intervention to proactive prevention. Unlike episodic care models that respond after disease onset, this approach enables healthcare systems to intervene earlier, personalize outreach, and ultimately reduce avoidable hospitalizations and long-term disease burden.

Importantly, our study bridges a long-standing gap between actuarial analytics (cost-centric) and clinical intervention (care-centric). By reframing administrative data as a strategic asset rather than just a reimbursement mechanism, we empower both payers and providers to take unified action toward population health goals. The use of interpretable machine learning techniques, combined with operational pilot programs, further underscores the real-world feasibility of such solutions.

We also demonstrated that deploying these models within existing care management infrastructure can lead to measurable improvements ranging from a 27% increase in screening response to a 9.4% reduction in ER visits. These results are particularly valuable for organizations operating under value-based contracts, where preventive care quality measures (e.g., HEDIS, Star Ratings) directly impact reimbursement, patient satisfaction, and clinical equity.

7.1. Future Directions

While the proposed framework sets a strong foundation, several directions remain for future research:

- **Multi-Modal Data Integration:** Future models can incorporate richer clinical data from electronic health records

(EHRs), wearable health devices, social determinants of health (SDOH), and behavioral analytics to improve precision and contextual relevance.

- Temporal Modeling and Forecasting: Incorporating temporal deep learning models (e.g., LSTM, Transformers) may enhance predictions of care gaps that evolve over time, allowing systems to anticipate and dynamically prevent them.
- Real-Time Preventive Care Triggers: Developing APIs and low-latency systems to integrate predictions into care management platforms or patient portals could enable timely interventions based on real-time data ingestion.
- Fairness, Bias Mitigation, and Personalization: As predictive care expands, future work must ensure that models do not reinforce existing disparities. This includes implementing fairness auditing, explainability mechanisms, and feedback loops to support individualized and culturally sensitive care.
- Policy and Compliance Alignment: With evolving regulations (e.g., Cures Act, ONC rules), future research should explore compliant model governance, data-sharing strategies, and ethical deployment frameworks aligned with federal and global health policies.

In conclusion, this paper repositions claims and data relating to a passive administrative artifact to an active driver of clinical intelligence. Predictive analytics grounded in real-world data can serve as the bridge between large-scale health data and meaningful, equitable, and preventive health outcomes, offering a path forward for more resilient, efficient, and patient-centered healthcare systems.

References

1. D. J. Morgan, B. Bame, P. Zimand, P. Dooley, K. A. Thom, A. D. Harris, S. Bentzen, W. Ettinger, S. D. Garrett-Ray, J. K. Tracy, and Y. Liang, "Assessment of machine learning vs standard prediction rules for predicting hospital readmissions," *JAMA Network Open*, vol. 2, no. 3, pp. e190 348–e190 348, 03 2019. [Online]. Available: <https://doi.org/10.1001/jamanetworkopen.2019.0348>
2. E. National Academies of Sciences, Medicine, and M. YOUNG, "Health care," in *Implementing Strategies to Enhance Public Health Surveillance of Physical Activity in the United States*. National Academies Press (US), 2019. [Online]. Available: <https://www.ncbi.nlm.nih.gov/books/NBK545645/>
3. H. Kharrazi and J. P. Weiner, "A practical comparison between the predictive power of population-based risk stratification models using data from electronic health records versus administrative claims: setting a baseline for future ehr-derived risk stratification models," *Medical care*, vol. 56, no. 2, pp. 202–203, 2018.
4. M. Altman, A. Wood, and E. Vayena, "A harm-reduction framework for algorithmic fairness," *IEEE Security & Privacy*, vol. 16, no. 3, pp. 34–45, 2018.
5. S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf
6. D. Cattel and F. Eijkenaar, "Value-based provider payment initiatives combining global payments with explicit quality incentives: A systematic review," *Medical Care Research and Review*, vol. 77, no. 6, pp. 511–537, 2020, pMID: 31216945. [Online]. Available: <https://doi.org/10.1177/1077558719856775>
7. J. C. Lauffenburger, J. M. Franklin, A. A. Krumme, W. H. Shrank, O. S. Matlin, C. M. Spettell, G. Brill, and N. K. Choudhry, "Predicting adherence to chronic disease medications in patients with long-term initial medication fills using indicators of clinical events and health behaviors," *Journal of managed care & specialty pharmacy*, vol. 24, no. 5, pp. 469–477, 2018, pMID: 29694288. [Online]. Available: <https://doi.org/10.18553/jmcp.2018.24.5.469>
8. R. Gaspersz, F. Lamers, J. M. Kent, A. T. F. Beekman, J. H. Smit, A. M. van Hemert, R. A. Schoevers, and B. W. J. H. Penninx, "Longitudinal predictive validity of the dsm-5 anxious distress specifier for clinical outcomes in a large cohort of patients with major depressive disorder," *The Journal of Clinical Psychiatry*, vol. 78, no. 2, p. 1207, 2017.
9. J. B. Young, M. Gauthier-Loiselle, R. A. Bailey, A. M. Manceur, P. Lefebvre, M. Greenberg, M.-H. Lafeuille, M. S. Duh, B. Bookhart, and C. H. Wysham, "Development of predictive risk models for major adverse cardiovascular events among patients with type 2 diabetes mellitus using health insurance claims data," *Cardiovascular diabetology*, vol. 17, pp. 1–13, 2018. [Online]. Available: <https://doi.org/10.1186/s12933-018-0759-z>
10. B. N. Becker, "Managing populations," in *Population Health*. Productivity Press, 2015, pp. 68–87. [Online]. Available: <https://www.taylorfrancis.com/chapters/edit/10.1201/b19266-8/managing-populations-bryan-becker>
11. G. Corrao, F. Rea, M. Di Martino, R. De Palma, S. Scondotto, D. Fusco, A. Lallo, L. M. B. Belotti, M. Ferrante, S. Pollina Addario, L. Merlino, G. Mancina, and F. Carle, "Developing and validating a novel multisource comorbidity score from administrative data: a large population-based cohort study from Italy," *BMJ Open*, vol. 7, no. 12, 2017. [Online]. Available: <https://bmjopen.bmj.com/content/7/12/e019503>
12. V. K. Verma and W.-Y. Lin, "A machine learning-based predictive model for 30-day hospital readmission prediction for COPD patients," in *2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, 2020, pp. 994–999.
13. D. V. Gunasekeran, D. S. Ting, G. S. Tan, and T. Y. Wong, "Artificial intelligence for diabetic retinopathy screening, prediction and management," *Current opinion in ophthalmology*, vol. 31, no. 5, pp. 357–365,

2020. [Online]. Avail- able: https://journals.lww.com/co-ophthalmology/fulltext/2020/09000/artificial_intelligence_for_diabetic_retinopathy.9.aspx
14. Z. Obermeyer, B. Powers, C. Vogeli, and S. Mullainathan, "Dissecting racial bias in an algorithm used to manage the health of populations," *Science*, vol. 366, no. 6464, pp. 447–453, 2019. [Online]. Available: <https://www.science.org/doi/abs/10.1126/science.aax2342>