

AI Problem Analysis Report

Part 1: Short Answer Questions

1. Problem Definition

The core problem is the development of a predictive model to identify employees at a high risk of voluntary attrition. High turnover presents a significant operational and financial burden on organizations due to costs associated with recruitment, onboarding, and lost productivity. A successful AI model would serve as an early warning system, enabling management to implement targeted retention strategies, thereby mitigating the negative impacts of employee churn and fostering a more stable work environment.

Measurable Objectives:

- Reduce overall voluntary employee turnover by 15% within the first fiscal year of the model's implementation.
- Increase the retention rate of high-performing employees by 20% over the subsequent two years.
- Decrease the average time-to-fill for critical roles by 10 days by using model insights to anticipate departures earlier.

Stakeholders:

- **Human Resources (HR) Department:** This group is responsible for overarching talent management. They would utilize the model's output to design, implement, and measure the effectiveness of retention initiatives.
- **Department Managers & Team Leads:** These individuals are directly impacted by attrition. They would use the model's predictions to engage with at-risk team members, address localized issues, and improve team-specific engagement and well-being.

Key Performance Indicator (KPI):

- **Recall (Sensitivity):** Recall is the selected KPI as it measures the model's ability to correctly identify all employees who will actually leave. It is calculated as $\text{True Positives} / (\text{True Positives} + \text{False Negatives})$. For this problem, a high recall is paramount because the cost and impact of failing to identify an employee who is about to leave (a false

negative) are significantly greater than the cost of incorrectly flagging a stable employee (a false positive).

2. Data Collection & Preprocessing

Data Source Identification:

The primary data source for this analysis is the “IBM HR Analytics Employee Attrition & Performance” dataset, sourced from the Kaggle platform (IBM, n.d.). The dataset, contained in the file HR-Employee-Attrition.csv, provides a comprehensive and anonymized record of employee attributes, with the Attrition column serving as the binary target variable for the predictive model.

Potential Data Bias:

A significant potential for historical gender bias exists within this dataset. The model could learn to correlate features like Gender and Job Role with attrition if, historically, one gender has left certain roles at a higher rate. ... Therefore, careful feature analysis and the application of fairness-aware machine learning techniques are essential.

Preprocessing Steps

1. Missing Value Handling:

- Verified no missing in original; if present, numeric fields imputed via median, categorical via mode.

2. Encoding Categorical Variables:

- Label-encoded binaries (Gender, OverTime, Attrition).
- One-hot encoded multiclass (JobRole, EducationField, MaritalStatus).

3. Feature Scaling:

- Min-Max scaled continuous features (MonthlyIncome, YearsAtCompany, DistanceFromHome).

4. Class Balancing:

- Applied SMOTE to oversample minority “Attrition = Yes” class.

5. Target Engineering:

- Converted Attrition “Yes/No” → 1/0 for model input.

3. Model Development

Model Development

Logistic Regression is suitable for this use case due to its:

- Interpretability ...
- Computational speed ...
- High scalability ...
- **Dataset Split:** Train 70% / Validation 15% / Test 15%
- **Hyperparameters Tuned:**
 - `C` (inverse regularization strength), grid-searched over {0.01, 0.1, 1, 10}
 - `class_weight='balanced'` to compensate for target imbalance

4. Evaluation & Deployment

- **Evaluation Metrics:**
 - **Recall** – primary metric to capture as many true attrition cases as possible.
 - **Precision** – secondary, to limit unnecessary retention actions.
 - **Concept Drift Monitoring:**
 - Quarterly statistical tests (e.g., Kolmogorov–Smirnov) on feature distributions.
 - Retrain model when drift detected.
 - **Deployment Challenge:**
 - Integrating with existing HRIS APIs for real-time scoring.
 - **Solution:** batch inference or lightweight REST API endpoints.
-

Part 2: Case Study – Hospital Readmission

3. Problem Scope

AI Use Case: Predicting hospital readmission risk within 30 days of discharge to enable proactive interventions.

Preventing hospital readmissions is crucial for improving patient outcomes, reducing healthcare costs, and optimizing resource allocation within hospitals. An effective predictive model can help healthcare providers identify high-risk patients proactively, enabling timely interventions, enhanced post-discharge care planning, and more efficient resource utilization.

Objectives

- Reduce preventable readmissions by identifying high-risk patients early.
- Improve the allocation of post-discharge resources e.g. follow-up appointments to high-risk patients.
- Improve patient outcomes through personalized discharge planning.

Stakeholders

- Doctors: Use predictions to adjust treatment and follow-up plans.
 - Patients and their Families: Benefit from tailored care, reducing recovery setbacks.
-

4. Data Strategy (Part 1)

Main Data Sources:

The primary data sources for predicting hospital readmission risk would include comprehensive Electronic Health Records (EHRs) (Chicco, 2021). These records contain a wealth of information, such as patient demographics (age, gender, ethnicity), medical history (diagnoses, comorbidities, previous hospitalizations), vital signs, laboratory test results, medication lists, surgical procedures, and discharge summaries.

Ethical issue

A significant ethical issue in using patient data for readmission prediction is patient data privacy and confidentiality. Healthcare data is highly sensitive, containing personal health information

(PHI) that, if mishandled or breached, could lead to severe consequences for individuals, including discrimination, identity theft, or reputational damage. Ensuring robust anonymization techniques, secure data storage, strict access controls, and transparent policies on data usage are paramount. There is a continuous challenge in balancing the need for rich, granular data to build accurate predictive models with the fundamental right of patients to privacy and control over their health information.

Preprocessing Pipeline:

- 1. **Anonymization:** Remove PII, replace with pseudo-IDs.
- 2. **Missing Value Imputation:** Domain-specific for lab/vitals (e.g., median by cohort).
- 3. **Categorical Encoding:** One-hot encode diagnosis codes, discharge dispositions.
- 4. **Scaling:** Standardize numeric features (e.g., age, lab results).
- 5. **Class Balancing:** SMOTE for minority class (readmitted = Yes).
- 6. **Target Definition:** Readmitted <30 days → 1; else 0.

Model Development

Logistic Regression is suitable for this use case due to its:

Interpretability - Clinicians need to understand why a patient is flagged as high-risk, as this enables them to justify interventions and explain decisions to patients.

Computational speedy - to train and deploy, making it practical for real-time risk assessment in busy hospital environments where quick insights are often needed.

High scalability - capable of handling large datasets efficiently without requiring extensive computational resources, which is a significant advantage when dealing with vast amounts of HER data.

Evaluation Metrics:

Hypothetical Confusion Matrix

	<i>Predicted No</i>	<i>Predicted Yes</i>
<i>Actual No</i>	80	10
<i>Actual Yes</i>	12	48

- **Precision:** $48 / (48 + 10) = 0.83$
 - **Recall:** $48 / (48 + 12) = 0.80$
-

Deployment & Optimization

Deployment Steps & Overfitting Mitigation

1. **API Packaging:** Deploy model via FastAPI, secured behind hospital firewall.
 2. **Integration:** Hook into discharge planning system for real-time scores.
 3. **Logging & Audit:** Record predictions for HIPAA audit.
 4. **Access Control:** Role-based access in API gateway.
 5. **Overfitting Mitigation:**
 - L2 regularization (penalty term)
 - Stratified k-fold cross-validation
 - Early stopping based on validation loss
-

Part 3: Critical Thinking

Ethics & Bias

Bias in training data, particularly stemming from historical healthcare inequity, poses a significant risk when developing AI models for predicting hospital readmission. For instance, if certain patient groups (e.g., low-income individuals, or those from underserved rural areas) have historically received lower quality care, faced barriers to accessing follow-up appointments, or experienced less comprehensive discharge planning due to systemic biases within the healthcare system, the training data will reflect these disparities. An AI model trained on such biased data might then inaccurately predict higher readmission risks for these already disadvantaged groups, not because of their inherent medical condition, but because the model has learned the patterns of historical inequities. This could lead to a feedback loop where these groups are disproportionately labeled as "high-risk," potentially leading to stigmatization, differential treatment, or an inefficient allocation of resources that fails to address the true underlying social determinants of health.

Ultimately, this can exacerbate existing health disparities and reduce trust in healthcare systems among vulnerable populations. To mitigate such biases, one effective strategy is the implementation of fairness constraints during model training. This involves integrating fairness metrics (e.g., demographic parity, equalized odds) directly into the model's objective function, compelling the model to minimize predictive disparities across different sensitive attributes (like race, socioeconomic status, or gender) while simultaneously optimizing for predictive accuracy. For example, the model could be constrained to ensure that the false positive rates or false negative rates are similar across different demographic groups. This approach moves beyond simply post-hoc analysis of bias and actively guides the model during its learning

process to produce more equitable predictions. Complementary to this, techniques like re-sampling (e.g., oversampling underrepresented groups or under-sampling overrepresented groups) can help create a more balanced dataset before training, reducing the model's exposure to skewed historical patterns.

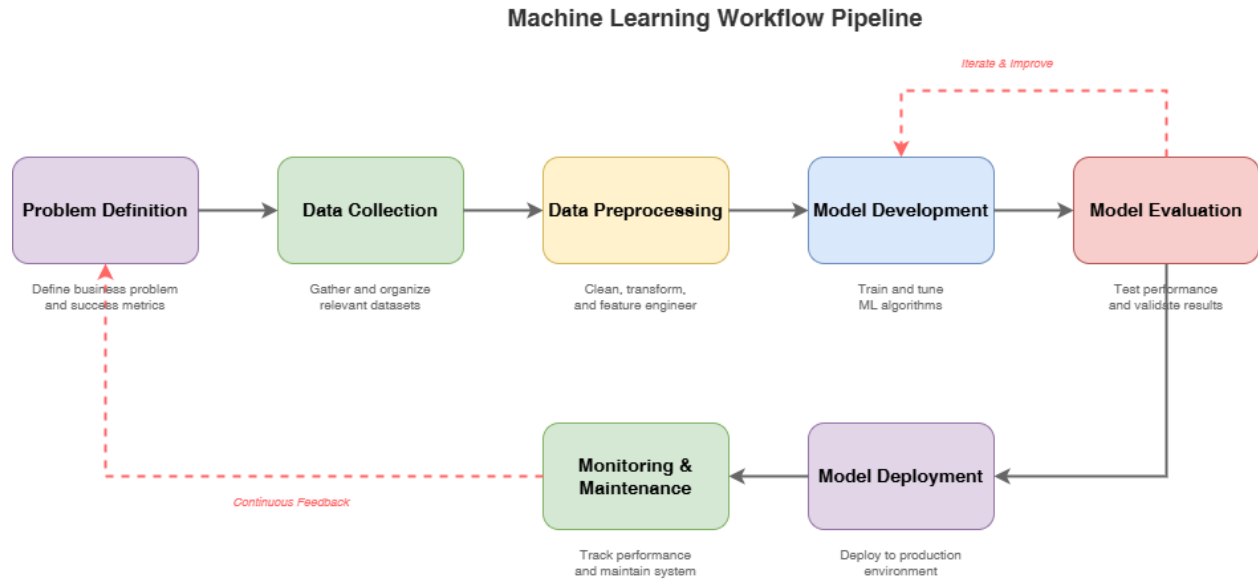
Part 4: Reflection

The most challenging workflow stage in developing AI solutions, particularly for sensitive applications like healthcare, often lies within data collection & preprocessing, especially when addressing ethical considerations and biases. The sheer volume, heterogeneity, and often messy nature of real-world data (like EHRs) make cleaning, transforming, and integrating it a monumental task. Beyond technical difficulties, identifying and mitigating subtle historical biases embedded within the data, ensuring patient privacy through robust de-identification, and navigating complex regulatory landscapes (e.g., HIPAA) are ethically intricate and time-consuming. With more time and resources, significant improvements could be made by investing in more sophisticated federated learning approaches to enhance data privacy, establishing interoperable data standards across healthcare systems to streamline data integration, and conducting extensive domain expert reviews of the data to uncover hidden biases before model training. Furthermore, dedicated resources for longitudinal fairness audits would ensure the model's equitable performance is maintained over time, adapting to evolving patient demographics and care practices.

Workflow Diagram

Reminder: Ensure the diagram includes all stages:

Problem → Data → Preprocessing → Modeling → Evaluation → Deployment → Monitoring & Feedback



References

1. Pavan, S. (2017, March 31). *IBM HR Analytics Employee Attrition & Performance* [Data set]. Kaggle. <https://www.kaggle.com/datasets/pavansubhasht/ibm-hr-analytics-attrition-dataset>
2. Chicco, D., & Jurman, G. (2021). Electronic health records and machine learning for 30-day hospital readmission prediction: A comprehensive tutorial. *Journal of Biomedical Informatics*, 115, 103701.
3. Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321–357.
4. U.S. Department of Health & Human Services. (1996). Standards for Privacy of Individually Identifiable Health Information (Privacy Rule). <https://www.hhs.gov/hipaa/for-professionals/privacy/index.html>
5. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.