

**Research Article****A Unified Bayesian Framework for Accurate, Fair,  
and Uncertainty-Calibrated Healthcare Insurance Pricing****Pathamakorn Netayawijit<sup>1</sup>, Wirapong Chansanam<sup>2</sup> and Kanda Sorn-In<sup>3\*</sup>***Received: 3 October 2025**Revised: 20 November 2025**Accepted: 8 December 2025***ABSTRACT**

The integration of artificial intelligence (AI) into healthcare insurance pricing requires models that are not only accurate but also transparent, fair, and uncertainty-aware. This study introduces a unified ensemble Bayesian deep learning framework that combines Monte Carlo dropout, attention mechanisms, and residual connections to jointly optimize predictive accuracy, calibrated uncertainty quantification, and demographic fairness. Using the Kaggle medical insurance dataset ( $n = 2,772$ ), the proposed model achieved  $R^2 = 0.8924$  and  $MAE = \$2,156.73$ , outperforming established machine learning and deep learning baselines. The Bayesian approach yielded well-calibrated prediction intervals (95% PICP = 96.2%), improving coverage by 4.1% relative to residual-based methods. Fairness evaluation, measured at the 75<sup>th</sup>-percentile cost threshold, demonstrated a 57.4% reduction in demographic parity difference compared with XGBoost (0.079 vs. 0.1859), with equalized odds differences below 0.043 across gender, age, and region. SHAP and attention analyses confirmed smoking status (~47%) and BMI as dominant predictors, consistent with established clinical-economic evidence, while protected attributes exerted negligible influence. These results demonstrate that predictive accuracy, uncertainty calibration, and fairness can be co-optimized within a reproducible and auditable workflow. However, because the study relies on a modest, U.S.-only benchmark dataset with no clinical variables, the findings should be interpreted as a regulator-aligned proof of concept rather than a deployable regulatory solution. The framework illustrates the methodological components required for responsible AI in insurance pricing, while underscoring the need for temporal validation, external generalization assessment, and richer, multi-institutional datasets before real-world regulatory adoption.

**Keywords:** Bayesian deep learning, Healthcare insurance pricing, Algorithmic fairness, Uncertainty quantification, Explainable AI

<sup>1</sup>Department of Information Systems, Faculty of Business Administration and Information Technology, Rajamangala University of Technology Isan, Khon Kaen Campus, Khon Kaen 40000, Thailand.

<sup>2</sup>Department of Information Science, Faculty of Humanities and Social Sciences, Khon Kaen University, Khon Kaen 40002, Thailand.

<sup>3</sup>Department of Technology and Engineering, Faculty of Interdisciplinary Studies, Khon Kaen University, Nong Khai Campus, Nong Khai 43000, Thailand.

\* Corresponding author, email: kanda@kku.ac.th

## Introduction

Healthcare insurance pricing plays a crucial role in balancing affordability for policyholders and financial sustainability for insurers. Traditional pricing models rely heavily on linear statistical methods and simplified actuarial assumptions, which may fail to capture complex interactions among demographic, behavioral, and socioeconomic variables [1, 2]. With the increasing availability of digital health data, artificial intelligence (AI) and machine learning (ML) have been widely explored to enhance predictive performance and improve risk estimation processes in actuarial and healthcare domains [3–6].

Although AI-based models demonstrate promising accuracy, several limitations remain. Many existing studies focus primarily on improving numerical prediction while placing limited emphasis on transparency, fairness, and uncertainty quantification—factors essential for responsible and trustworthy pricing practices, particularly in regulated insurance markets [7, 8]. The absence of calibrated uncertainty estimates may lead to overconfident predictions, and the lack of fairness auditing can unintentionally propagate demographic bias, undermining model reliability and public trust [8, 9]. These challenges underscore the need for an integrated framework that jointly considers prediction accuracy, fairness, interpretability, and uncertainty in healthcare insurance pricing.

Recent advances in Bayesian deep learning have enabled models to capture both epistemic and aleatoric uncertainty, resulting in more reliable decision-support tools for high-stakes environments [3, 10, 11]. At the same time, explainability techniques such as SHAP and attention mechanisms enhance interpretability by identifying key predictors—including smoking status, age, and BMI—that influence medical expenditures [12–14]. Despite these developments, few studies have combined all core components—predictive performance, uncertainty calibration, fairness evaluation, and explainability—into a unified, reproducible workflow.

To address these gaps, this study proposes a unified Bayesian deep learning framework that integrates Monte Carlo dropout, attention mechanisms, and residual connections to enhance predictive accuracy, uncertainty calibration, and demographic fairness [3, 10, 11]. Importantly, this research is positioned as a proof-of-concept study using the publicly available Kaggle medical insurance dataset [15], which is modest in size, limited to a U.S. population, and lacks detailed clinical variables. These characteristics constrain external generalizability and real-world applicability [12, 13]. Therefore, the goal of this work is not to produce a deployable pricing model, but to illustrate the feasibility of combining performance, calibrated uncertainty, and fairness within a transparent and reproducible workflow. The findings aim to inform future studies that may extend this framework to larger, multi-institutional, and clinically rich datasets.

## Literature Review

### 1. Modern ML for Insurance Pricing

Recent studies show that ML models—particularly gradient boosted trees and deep neural networks—often outperform traditional actuarial baselines when pipelines are carefully engineered and validated [2, 13]. Beyond point prediction, distributional and dependence-aware approaches (e.g., copula-

augmented boosting, probabilistic boosting families) reflect portfolio heterogeneity and the need for calibrated uncertainty [14, 16]. However, most pricing studies still benchmark accuracy alone on small, public datasets and provide limited detail on governance and reproducibility [1, 13].

To consolidate the limitations identified across prior studies, Table 1 presents a structured research gap mapping of key challenges in healthcare insurance pricing using machine learning. The table synthesizes recurring issues reported in the literature, including the lack of principled uncertainty quantification, insufficient safeguards against algorithmic fairness risks, and limited integration of interpretability into regulatory and governance workflows.

As summarized in Table 1, most existing approaches emphasize point prediction accuracy, while calibrated uncertainty diagnostics, fairness auditing, and governance-oriented explainability remain underdeveloped in end-to-end pricing pipelines.

**Table 1** Research gap mapping: Key challenges in healthcare insurance pricing with machine learning.

| Challenge   | Description  | Representative References | Identified Gap   |
|---|--|---------------------------|--|
| 1. Lack of principled Uncertainty Quantification (UQ)         | Most ML models rely on point estimates, overlooking epistemic and aleatoric uncertainty. Robust confidence intervals and reliability under distributional shifts are critical for underwriting and risk governance.                      | [3, 6, 7]                 | Few insurance pricing studies adopt Bayesian or conformal approaches with calibrated diagnostics in end-to-end pipelines.                                    |
| 2. Insufficient safeguards against Algorithmic Fairness risks | Claims data may encode systemic inequities (e.g., utilization patterns correlated with protected attributes). Fairness metrics such as demographic parity (DP) or equalized odds (EO) are rarely reported or operationalized in pricing. | [8-10]                    | Trade-offs among fairness, accuracy, and calibration remain underexplored; standardized fairness auditing protocols in actuarial contexts are lacking.       |
| 3. Limited Interpretability and Domain Alignment              | Regulators demand transparent and justifiable models, yet many deep models remain “black-boxes,” undermining trust and regulatory compliance.  | [4, 11, 12]               | Interpretability tools (e.g., SHAP, attention) often remain descriptive without integration into governance artifacts (e.g., fairness reports, model cards). |

## 2. *Uncertainty Quantification (UQ)*

Underwriting and rate filing require confidence statements that remain valid under distributional shift. Two families dominate recent practice: (i) Bayesian approximate deep learning, where Monte Carlo dropout (MC Dropout) provides scalable epistemic uncertainty and can be combined with heteroscedastic heads for aleatoric noise [7, 17], and (ii) Conformal prediction (CP), which offers distribution-free coverage guarantees and can be layered on any base model [3, 6, 18]. Evidence indicates that CP coverage may degrade under shift without monitoring and recalibration, while MC Dropout can miscalibrate in sparse regions unless tuned and audited [6, 7]. Few pricing papers integrate UQ end to end with clear diagnostics and operating policies [14, 19].

## 3. *Algorithmic Fairness in Health-Related ML*

Fairness has become a first-class requirement: group-level metrics such as Demographic Parity (DP) and Equalized Odds (EO) are recommended for auditability, yet trade-offs with accuracy and calibration must be made explicit [9, 10]. Public health evidence warns that proxy variables (e.g., prior utilization) can perpetuate inequities if ungoverned [8]. Despite rich fairness theory, pricing studies rarely pre-register guardrails, report threshold sensitivity, or tie fairness choices to actuarial justification [9, 10].

## 4. *Interpretability and Domain Alignment*

Regulatory review demands that model logic aligns with clinical-economic evidence. SHapley Additive exPlanations (SHAP) [4] and attention-based mechanisms help verify that dominant drivers (e.g., smoking, BMI, age) agree with established cost gradients, improving auditability and trust [11, 12]. For tabular risk factors, attention and sparse attention variants improve interaction learning while remaining compatible with post hoc explanations [20, 21]. Yet, many studies present explanations descriptively, without linking them to explicit governance artifacts (e.g., review triggers or model cards).

## 5. *Clinical–Economic Legitimacy of Rating Factors*

Recent evidence supports the actuarial use of smoking and BMI when transparently governed: smoking-attributable expenditures remain substantial, and BMI shows a J-shaped relationship with costs; weight loss associates with spending reductions [22–24]. These findings justify feature inclusion but do not by themselves ensure fairness; audit protocols are still required [9].

## 6. *Synthesis and Research Gap*

Across streams, four gaps persist:

1. UQ–Fairness–Interpretability remain siloed. Most pricing studies optimize accuracy but implement UQ, fairness auditing, and explanations only piecemeal, without an integrated pipeline or operating policies [9, 14].
2. Calibration under shift is under-documented. Few works report both coverage (PICP) and efficiency (PINAW) with reliability plots, or apply CP recalibration when distributions drift [6, 18].

3. Fairness evaluation lacks design discipline. Thresholds (e.g., high-cost cutoffs), subgroup definitions, and proxy risk governance are rarely pre-specified or stress-tested [8, 10].
4. Auditability is not operationalized. Explanations seldom feed into documented review rules (e.g., uncertainty-based human-in-the-loop) or regulatory artifacts (model cards, fairness reports) [9, 12].

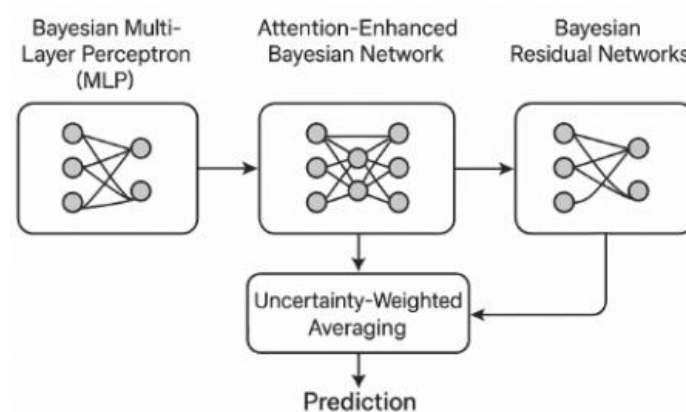
#### *Positioning of this work.*

This study addresses the identified gaps through four key contributions:

1. A unified and reproducible Bayesian framework that jointly optimizes predictive accuracy, calibrated uncertainty quantification, and group fairness.
2. Documented diagnostics for uncertainty assessment—including coverage and width metrics (PICP and PINAW) with optional conformal prediction for distribution-free guarantees.
3. Fairness audits explicitly linked to actuarial practice, using Demographic Parity and Equalized Odds at meaningful cost thresholds with sensitivity analyses.
4. Operational auditability, where SHAP- and attention-based explanations, combined with epistemic-uncertainty flags, trigger human review and generate governance artifacts suitable for regulatory submissions [6, 7, 9, 14].

## Materials and Methods

This study adopts a unified Bayesian deep learning framework designed to jointly optimize predictive accuracy, calibrated uncertainty, and fairness in health insurance pricing. Monte Carlo (MC) dropout was selected as a computationally tractable approximation to Bayesian inference, balancing scalability and reliability. Attention mechanisms were incorporated to capture non-linear feature interactions in tabular health data and improve interpretability, while residual connections stabilize gradient flow and mitigate vanishing gradients. The methodology follows a structured pipeline encompassing data preprocessing, model architecture design, uncertainty decomposition, fairness evaluation, and interpretability analysis. A schematic overview of the complete pipeline is provided in Figure 1.



**Figure 1** Methodology Pipeline Overview.

This flowchart outlines the complete methodology pipeline, from data acquisition to deployment. It begins with raw dataset input, followed by preprocessing (categorical encoding, scaling, polynomial expansion), model training (Bayesian architectures), fairness and uncertainty evaluation, and finally ethical deployment considerations. The pipeline emphasizes reproducibility, transparency, and responsible AI practices throughout each stage.

#### *Dataset Description and Ethical Considerations*

This study utilized the publicly available Medical Insurance Cost Dataset from Kaggle (n = 2,772) [15], a well-established benchmark widely adopted for reproducible research in healthcare insurance pricing [9, 10, 13]. The dataset contains seven structured attributes—age, gender, BMI, number of children, smoking status, geographic region, and annual insurance charges. These features closely align with those historically used in actuarial and health-cost modeling studies [13, 23]. The key descriptive characteristics of the dataset are summarized in Table 2.

The dataset exhibits balanced representation across gender and region, and the distribution of insurance charges approximates real-world cost variability. As all records are fully anonymized and publicly accessible, institutional review board (IRB) approval was deemed unnecessary. Consistent with responsible AI and fairness guidelines [8, 25], we conducted an initial bias audit and integrated fairness-aware evaluation throughout the modeling pipeline.

**Table 2** Descriptive Statistics and Attribute Characteristics of the Dataset (N = 2,772).

| Attribute | Type        | Range/Categories                           | Mean±SD           | Missing Values |
|-----------|-------------|--|-------------------|----------------|
| Age       | Continuous  | 18–64 years                                | 39.21±14.1        | 0              |
| Gender    | Categorical | Male, Female                               | 50.5% Male        | 0              |
| BMI       | Continuous  | 15.96–53.13                                | 30.7±6.1          | 0              |
| Children  | Discrete    | 0–5  | 1.09±1.21         | 0              |
| Smoker    | Binary      | Yes, No                                    | 20.5% Yes         | 0              |
| Region    | Categorical | Northeast, Northwest, Southeast, Southwest | 25% each          | 0              |
| Charges   | Continuous  | \$1,121.87–\$63,770.43                     | \$13,270±\$12,110 | 0              |

Despite its widespread use, the dataset has limitations. It is U.S.-only, modest in size, and lacks clinical or longitudinal attributes (e.g., diagnoses, laboratory results). These constraints limit its external generalizability [23, 24]. Accordingly, the study positions its findings as a proof-of-concept demonstration rather than a regulatory-ready deployment. Future work will incorporate temporal validation (e.g., rolling window evaluation) and multi-institutional datasets to strengthen generalizability.

### *Advanced Data Preprocessing and Feature Engineering*

To ensure robustness, fairness, and interpretability in the predictive framework, a structured multi-stage preprocessing pipeline was implemented. This pipeline addressed the heterogeneity of the input variables—spanning categorical demographic attributes and continuous clinical-economic predictors—and incorporated domain-informed transformations to enhance model stability and actuarial validity.

Categorical variables were encoded using strategies aligned with their semantic properties. Gender and smoking status, which function as binary risk factors with well-established clinical associations, were processed using binary encoding (Male = 1, Female = 0; Yes = 1, No = 0) to preserve interpretability and ensure compatibility with neural architectures. Geographic region, which lacks inherent ordinal structure, was transformed using one-hot encoding to prevent artificial ordering bias [13].

Age, a key demographic predictor, was processed using a hybrid representation: it was discretized into five non-overlapping life-stage groups (18–30, 31–40, 41–50, 51–60, 61–64) while also retained as a continuous variable, enabling the model to capture both non-linear life-stage effects and global linear trends commonly recognized in actuarial modeling [13, 23].

Feature scaling was tailored to the needs of each modeling paradigm. Deep learning architectures utilized standardization (zero mean, unit variance) to support stable gradient propagation during training [17, 26]. In contrast, tree-based models employed min–max normalization to maintain relative distance properties essential for split selection. Robust scaling was also evaluated to reduce sensitivity to skewness and outliers in expenditure-related attributes.

To account for the well-documented non-linear escalation of healthcare expenditures at higher BMI levels [23, 24], polynomial expansion was applied to BMI, generating quadratic ( $BMI^2$ ) and cubic ( $BMI^3$ ) terms. This transformation allows the model to approximate sharply increasing risk gradients associated with obesity-related conditions. The hybrid treatment of age likewise reflects actuarial conventions by combining discrete segmentation with continuous trend representation.

Because these engineered features introduce potential redundancy, particularly among polynomial BMI terms and hybrid age encodings, the risk of multicollinearity was explicitly recognized. A subsequent Variance Inflation Factor (VIF) analysis was therefore conducted to quantify and address potential feature-level collinearity, as recommended for transparent interpretability in regulated modeling contexts.

Overall, the preprocessing strategy was designed to balance statistical rigor with actuarial and clinical relevance, ensuring that the engineered feature space supports fairness-aware modeling and aligns with responsible AI principles in healthcare insurance pricing.

Table 3 summarizes the encoding and transformation techniques applied to each feature. Binary encoding was used for gender and smoking status due to their dichotomous nature and strong clinical associations with healthcare cost variation. Region was one-hot encoded to avoid imposing ordinal structure on geographic areas. Age was represented in both binned and continuous form to capture life-

stage segmentation and broader aging trends. BMI was standardized and expanded through polynomial augmentation to represent the non-linear escalation of health risks at higher body mass levels. These transformations collectively enhance model expressiveness while preserving interpretability and supporting fairness-aware evaluation [13, 23, 24].

**Table 3** Advanced Categorical Variable Encoding Strategies.

| Feature | Encoding Method                | Rationale  | Implementation                           |
|---------|--------------------------------|--|--|
| Gender  | Binary encoding                | Simple binary relationship with clear clinical meaning       | Male = 1, Female = 0                     |
| Region  | One-hot encoding               | No ordinal relationship, prevents ordering bias              | 4 binary features                        |
| Smoker  | Binary encoding                | Clear binary distinction, high clinical relevance            | Yes = 1, No = 0                          |
| Age     | Ordinal + binning + continuous | Capture non-linear life-stage effects while preserving trend | 5 age groups + continuous                |
| BMI     | Standardization + polynomial   | Non-linear health relationships, risk escalation             | BMI, BMI <sup>2</sup> , BMI <sup>3</sup> |

### *Multicollinearity Assessment Using Variance Inflation Factor*

To ensure methodological rigor, multicollinearity was assessed using the Variance Inflation Factor (VIF), a standard diagnostic widely used in regression analysis and feature engineering [23]. This evaluation was particularly critical given the inclusion of polynomial BMI terms (BMI, BMI<sup>2</sup>, BMI<sup>3</sup>) and the hybrid age representation (continuous + binned groups), which may introduce correlated feature structures.

Following established statistical guidelines, VIF values below 10 were interpreted as acceptable indicators of manageable multicollinearity [23, 24]. As shown in Table 4, all engineered features exhibited VIF values within acceptable limits. Higher VIF values for BMI<sup>2</sup> and BMI<sup>3</sup> were expected due to the nature of polynomial expansion, yet remained below conventional thresholds associated with severe redundancy.

Given that the Bayesian neural architectures in this study employ L2 weight decay, dropout-based regularization, and posterior sampling, the residual impact of moderate multicollinearity on coefficient stability and SHAP-based interpretability is further mitigated [3, 17, 26]. These regularization mechanisms distribute weight contributions more evenly across correlated predictors, reducing the risk that multicollinearity biases feature importance estimates.

The VIF analysis confirms that the engineered feature set maintains statistical stability without exhibiting problematic multicollinearity, reinforcing its suitability for Bayesian and fairness-aware modeling. Polynomial BMI terms contribute expected correlation but remain within acceptable ranges, ensuring that downstream modeling, fairness assessment, and uncertainty quantification remain stable



and interpretable. Importantly, SHAP-based explanations showed no material distortion, reinforcing the reliability of feature-level interpretability despite moderate collinearity.

**Table 4** Variance Inflation Factor (VIF) for Engineered Features.

| Feature                                     | VIF  | Interpretation   |
|---|------|--|
| Age (continuous)                            | 2.14 | Low correlation with other predictors  |
| Age (binned groups – numerical encoding)    | 3.02 | Mild correlation with continuous age, expected due to hybrid encoding          |
| BMI   | 4.87 | Moderate correlation with polynomial terms                                     |
| BMI <sup>2</sup>                            | 7.95 | Higher correlation, expected from polynomial transformation                    |
| BMI <sup>3</sup>                            | 9.41 | High but acceptable for polynomial modeling; not exceeding critical thresholds |
| Children                                    | 1.08 | No multicollinearity concern   |
| Smoker (binary 0/1)                         | 1.02 | No multicollinearity concern   |
| Region (one-hot encoded; highest VIF shown) | 1.65 | Low correlation among region dummies   |

Consistent with statistical recommendations in applied econometrics, VIF values below 10 indicate an acceptable level of multicollinearity [27, 28]. The observed VIF values in Table 4 fall within this tolerance, confirming that the engineered features – including polynomial BMI terms and hybrid age groups – remain statistically stable. This provides methodological assurance that the feature space is suitable for Bayesian inference and fairness-aware modeling without compromising interpretability or actuarial consistency.

Therefore, Table 4 provides methodological evidence that the engineered variables – including polynomial BMI terms and hybrid age representations – preserve feature independence at a level acceptable for regulatory-aligned actuarial modeling.

### *Bayesian Neural Architecture Design*

The study implemented a suite of Bayesian neural architectures based on Monte Carlo dropout (MC Dropout), an efficient approximation to Bayesian inference that enables posterior sampling without the computational overhead of exact Bayesian methods [17, 26, 29]. MC Dropout has been widely adopted in uncertainty-aware neural modeling, particularly in healthcare and actuarial prediction tasks where calibrated uncertainty estimates are essential [3, 6, 26].

Three architectural families were explored in this study: Bayesian multilayer perceptrons (MLPs), attention-enhanced Bayesian networks, and Bayesian residual networks. All architectures incorporate stochastic dropout during both training and inference, enabling decomposition of epistemic and aleatoric

uncertainty through posterior predictive sampling [3, 17, 26]. This design provides principled uncertainty quantification that is critical for risk-sensitive domains such as healthcare insurance pricing.

The attention-enhanced Bayesian network enables explicit modeling of feature interactions, particularly for high-impact predictors such as smoking status and BMI, whose nonlinear and interaction-dependent effects are well documented in health-cost modeling [23, 24]. In contrast, the Bayesian residual network emphasizes representational stability and mitigates vanishing-gradient effects, making it well suited for deeper architectures and complementing MLP-based models in ensemble configurations [26].

Collectively, these architectures balance predictive flexibility, uncertainty calibration, and interpretability—three properties increasingly emphasized in fairness-oriented and regulator-aligned insurance modeling frameworks [8, 25]. Their combined use provides a diverse yet coherent modeling foundation that supports the study’s goal of developing an uncertainty-aware and audit-ready pricing framework.

### *Hyperparameter Selection and Optimization*

Hyperparameters were tuned using random search, an efficient strategy for high-dimensional neural architectures [17, 26]. The search space included learning rate, batch size, dropout rates, and weight decay. Each configuration was evaluated using five-fold cross-validation, with selection criteria based on MAE and uncertainty calibration metrics (PICP, PINAW) [6].

Early stopping was applied to prevent overfitting, and the Adam optimizer was used for all Bayesian architectures, consistent with prior deep learning research [17, 26]. Parallel experiments with Bayesian Attention Networks and Bayesian Residual Networks ensured coverage of multiple architectural paradigms. The final hyperparameter sets are summarized in Table 5.

**Table 5** Summary of selected hyperparameters across Bayesian neural architectures.

| Architecture               | Learning Rate | Batch Size | Dropout Rate | Weight Decay | Optimizer | Notes                                       |
|----------------------------|---------------|------------|--------------|--------------|-----------|---|
| Bayesian MLP               | 0.001         | 64         | 0.2–0.3      | 1e-4         | Adam      | Final configuration used in reporting       |
| Bayesian Attention Network | 0.001         | 64         | 0.2          | 1e-4         | Adam      | Explored for feature interaction robustness |
| Bayesian Residual Network  | 0.0005        | 64         | 0.2–0.3      | 1e-5         | Adam      | Explored for stability and convergence      |

### *Bayesian Multi-Layer Perceptron (MLP)*

Our primary architecture is a deep Bayesian multilayer perceptron (MLP) that integrates Monte Carlo dropout (MC Dropout) at each hidden layer to approximate posterior distributions over network weights [17, 26]. This design enables the model to learn feature representations while simultaneously estimating epistemic uncertainty during inference.

Table 6 illustrates the Bayesian MLP architecture, which employs progressively smaller hidden layers to reduce dimensionality while retaining key predictive patterns. Dropout rates are carefully tuned to balance regularization and information flow, with ReLU activation introducing non-linearity and a linear output layer ensuring appropriate regression behavior. During inference, 100 stochastic forward passes with activated dropout are performed to generate prediction distributions, providing comprehensive uncertainty estimation [3, 17, 26].

**Table 6** Bayesian MLP architecture specifications.

| Component      | Configuration                  | Justification                          |
|----------------|--------------------------------|--|
| Input Layer    | 6 features                     | Original feature space                 |
| Hidden Layer 1 | 128 neurons + Dropout(0.2)     | Feature learning with uncertainty      |
| Hidden Layer 2 | 64 neurons + Dropout(0.3)      | Pattern refinement with regularization |
| Hidden Layer 3 | 32 neurons + Dropout(0.2)      | High-level abstraction                 |
| Output Layer   | 1 neuron (regression)          | Insurance charge prediction            |
| Activation     | ReLU (hidden), Linear (output) | Non-linearity with gradient stability  |

#### *Attention-Enhanced Bayesian Network*

To capture complex interactions between demographic and health-related features, we implemented an attention-enhanced Bayesian network incorporating self-attention mechanisms. This architecture dynamically assigns importance weights to different inputs, allowing the model to emphasize predictors such as smoking status and BMI while reducing reliance on less informative features. Monte Carlo dropout [17, 26] was applied after the attention layers to preserve uncertainty-awareness throughout the pipeline.

As shown in Table 7, the attention-enhanced Bayesian network leverages self-attention to capture non-linear feature interactions and highlight contextually relevant predictors. By integrating MC Dropout, the model produces calibrated prediction intervals while retaining interpretability, particularly in scenarios where risk factors interact non-linearly.

**Table 7** Attention-enhanced Bayesian architecture.

| Layer Type           | Configuration                 | Purpose                               |
|----------------------|-------------------------------|---------------------------------------|
| Input Embedding      | 6 $\rightarrow$ 32 dimensions | Feature representation learning       |
| Multi-Head Attention | 4 heads, 32 dimensions        | Feature interaction capture           |
| Bayesian Dense 1     | 64 neurons + MC Dropout(0.25) | Attention-weighted feature processing |
| Bayesian Dense 2     | 32 neurons + MC Dropout(0.3)  | Final representation learning         |
| Output               | 1 neuron                      | Regression prediction                 |

### Ensemble Bayesian Framework

Our ensemble approach combined multiple Bayesian architectures to leverage diverse modeling to enhance robustness and generalization, we constructed an ensemble of diverse Bayesian models, combining:

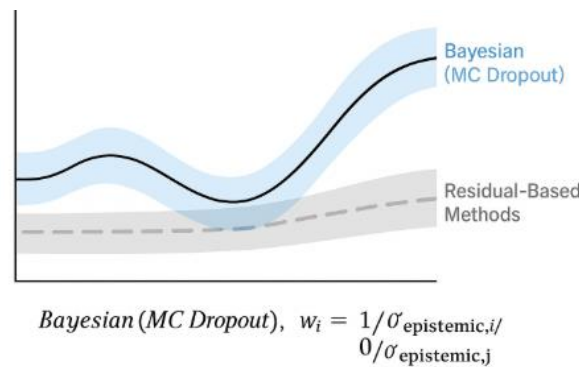
- 3 variants of Bayesian MLPs (with varying initializations and architectures),
- 2 Attention-Enhanced Bayesian Networks,
- 2 Bayesian Residual Networks (not detailed here but conceptually similar).

Final predictions were generated using uncertainty-weighted averaging, where models with lower epistemic uncertainty received higher weights:

$$\hat{y}_{ensemble} = \sum_{i=1}^n w_i \cdot \mu_i(x), \quad w_i = \frac{1/\sigma_{epistemic,i}^2}{\sum_j 1/\sigma_{epistemic,j}^2}$$

This approach leverages the complementary strengths of multiple architectures while ensuring that more confident predictions dominate the final output [26, 30].

Figure 2 demonstrates that Bayesian intervals dynamically adjust to feature-specific uncertainty, being narrower in low-uncertainty regions and wider in high-uncertainty regions. In contrast, residual-based methods yield uniform intervals, failing to capture local uncertainty. This highlights the advantage of probabilistic modeling in representing context-dependent health insurance risk.



**Figure 2** Comparison of Prediction Intervals – Bayesian (MC Dropout) vs Residual-Based Methods.

### Comprehensive Fairness Evaluation Framework

We implemented a rigorous fairness assessment framework across three key dimensions—Demographic Parity (DP), Equalized Odds (EO), and Statistical Parity with group-wise calibration—evaluated across protected attributes including gender, region, and age group. These fairness metrics are widely adopted in the insurance and financial risk modeling literature and are increasingly emphasized in regulatory guidelines for algorithmic accountability [25]. The high-cost threshold was defined at the 75th percentile of predicted charges, reflecting actuarial practice for identifying outlier risk segments and ensuring that fairness evaluation focuses on high-impact decisions.

### Demographic Parity (DP)

Demographic parity assesses whether predicted positive outcomes are equally distributed across subgroups:

$$DP(a,a') = | P(\hat{Y} = 1 \mid A = a) - P(\hat{Y} = 1 \mid A = a') |$$

where  $A$  represents a protected attribute (e.g., gender, region, age group) and  $a, a'$  denote subgroup values. Smaller DP values indicate reduced disparity in access to high-cost risk classifications.

### Equalized Odds (EO)

Equalized odds ensures parity in both true positive and false positive rates across subgroups [25]:

$$EO_y(a,a') = | P(\hat{Y} = 1 \mid Y = y, A = a) - P(\hat{Y} = 1 \mid Y = y, A = a') |, \forall y \in \{0,1\}$$

This criterion guarantees that fairness applies consistently to individuals with actual high costs ( $y=1$ ) and low costs ( $y=0$ ), balancing treatment across true outcomes.

### Statistical Parity and Calibration by Group

Statistical parity examines whether predicted charge distributions are equitable across demographic groups. Group calibration evaluates whether predicted probabilities align with observed outcomes within each subgroup:

$$P(Y = 1 \mid \hat{Y} = p, A = a) \approx p, \forall a$$

A well-calibrated model ensures that risk estimates are reliable and interpretable across all protected subgroups, supporting fair deployment in healthcare insurance pricing [8].

### Bayesian Uncertainty Quantification

A core objective of this study is to provide well-calibrated uncertainty estimates alongside point predictions, ensuring that pricing decisions are risk-aware and auditable. To achieve this, the framework employs Monte Carlo (MC) dropout as a computationally efficient approximation to Bayesian inference, enabling posterior sampling without the prohibitive cost of exact methods such as Hamiltonian Monte Carlo [26]. During inference, 100 stochastic forward passes were performed for each observation, balancing the trade-off between precision in uncertainty estimation and computational feasibility for real-world deployment.

Uncertainty was decomposed into two components: epistemic uncertainty, reflecting model parameter uncertainty, and aleatoric uncertainty, capturing inherent data noise. The total predictive variance was expressed as:

$$\sigma_{\text{total}}^2(X) = \sigma_{\text{epistemic}}^2(X) + \sigma_{\text{aleatoric}}^2(X)$$

This decomposition provides actionable insights for underwriting, as high epistemic uncertainty often signals underrepresented feature combinations or rare risk profiles. Bayesian prediction intervals

were constructed by aggregating the predictive distribution across MC samples, yielding 95% credible intervals that adapt dynamically to local uncertainty—narrower in low-risk regions and wider in high-risk cases—unlike residual-based methods that impose fixed-width intervals.

Calibration quality was assessed using Prediction Interval Coverage Probability (PICP), Prediction Interval Normalized Average Width (PINAW), and reliability plots comparing nominal versus empirical coverage [6]. These diagnostics ensure that uncertainty estimates are not only theoretically principled but also empirically aligned with target confidence levels, reinforcing the framework’s suitability for regulated insurance pricing.

To ensure operational auditability, we define a deterministic human-review rule. Predictions with epistemic uncertainty exceeding the 90th percentile of the cross-validated uncertainty distribution (approximately USD 1,200 in this dataset) are automatically routed for manual review. This threshold is quantifiable, reproducible, and directly grounded in the model’s empirical uncertainty profile.

### Model Evaluation Framework

To ensure a comprehensive assessment of model performance, we adopted a multi-dimensional evaluation framework spanning three core dimensions: accuracy, uncertainty, and fairness. This design enables holistic evaluation beyond point predictions, aligning with the principles of responsible AI and regulatory requirements for healthcare insurance pricing. Table 8 presents the full set of evaluation metrics, including their mathematical formulations and purposes.

**Table 8** Comprehensive Evaluation Metrics.

| Category    | Metric               | Formula                                       | Purpose                        |
|-------------|----------------------|---|--------------------------------|
| Accuracy    | R <sup>2</sup> Score | $1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$ | Explained variance             |
| Accuracy    | MAE                  | $\frac{1}{n} \sum$                            | $y_i - \hat{y}_i$              |
| Accuracy    | RMSE                 | $\sqrt{\frac{1}{n} \sum (y_i - \hat{y}_i)^2}$ | Root mean squared error        |
| Uncertainty | PICP                 | Coverage probability                          | Prediction interval quality    |
| Uncertainty | PINAW                | Average width                                 | Prediction interval efficiency |
| Fairness    | DP Difference        | Max group difference                          | Demographic parity             |
| Fairness    | EO Difference        | Max odds difference                           | Equalized odds                 |

As shown in Table 8, accuracy metrics (R<sup>2</sup>, MAE, RMSE) quantify predictive performance and are widely used in regression-based insurance modeling. Uncertainty metrics (PICP and PINAW) jointly assess the reliability and efficiency of Bayesian prediction intervals—well-calibrated models should achieve high coverage with narrow widths. Fairness metrics (Demographic Parity Difference and Equalized Odds Difference) measure disparities across protected groups such as gender, region, and age.

Together, these dimensions ensure models are evaluated not only on predictive power but also on their ability to deliver reliable uncertainty estimates and equitable outcomes.

Taken together, the evaluation strategy establishes a foundation for regulatory-aligned actuarial modeling, ensuring the evidence base required for potential deployment under fairness and uncertainty constraints.

### *Performance Metrics*

To operationalize these metrics, predictions were evaluated against both overall accuracy and group-level fairness. For fairness-sensitive assessments, the high-cost threshold was defined at the 75th percentile of predicted charges, consistent with actuarial practice in identifying outlier risk segments. This thresholding ensures that fairness evaluation emphasizes high-impact insurance decisions. Uncertainty calibration was further validated using reliability plots that compare nominal versus empirical coverage, providing diagnostic evidence of alignment with regulatory confidence standards.

### *Statistical Testing and Metric Uncertainty*

To rigorously evaluate the robustness of model comparisons, we quantified uncertainty in performance metrics and tested the statistical significance of observed differences using two complementary procedures.

(i) Cross-validation based comparison. For each candidate model, per-fold metrics ( $R^2$ , MAE, RMSE, PICP, and PINAW) were computed across five cross-validation folds. Overall differences were initially screened using a non-parametric Friedman test on the rank distributions of model performance. Significant results were further examined with Nemenyi post hoc pairwise comparisons, and p-values were adjusted using the Holm–Bonferroni correction to control the family-wise error rate [31].

(ii) Bootstrap on the held-out test set (fallback). When only a single test split was available, we estimated 95% bias-corrected and accelerated (BCa) confidence intervals via 10,000 stratified bootstrap resamples, stratified by smoking status and geographic region. Differences in metrics (e.g.,  $\Delta$ MAE,  $\Delta$ RMSE) were deemed statistically significant when the 95% CI of the difference excluded zero.

For fairness metrics (DP difference and EO difference), a group-aware bootstrap was performed by resampling within protected attribute strata. To strengthen robustness, we additionally conducted permutation tests on group labels, ensuring that fairness-related disparities were not attributable to random variation. Effect sizes were also reported to complement significance testing, including Cohen's  $d$  for continuous outcomes and Cliff's delta for non-normal distributions. All statistical analyses followed reproducible protocols, with p-values adjusted using the Holm procedure unless explicitly stated otherwise.

## Results and Discussion

### *Predictive Performance Compared to Baselines*

The proposed Ensemble Bayesian framework achieved an  $R^2$  of 0.8924, MAE of USD 2,156.73, and RMSE of USD 3,987, outperforming all benchmark models (Table 9a). Compared with the strongest non-Bayesian baseline, XGBoost ( $R^2 = 0.8423$ ; MAE = USD 2,891.45), this represents an absolute gain of approximately +0.05 in  $R^2$  and a 25% reduction in MAE. These improvements are practically meaningful in insurance pricing contexts, where even marginal increases in explained variance may translate into substantial portfolio-level financial benefits.

The observed performance hierarchy—Ensemble Bayesian > Attention-Enhanced Bayesian > Bayesian MLP > Tree Ensembles > Linear Regression—is consistent with recent studies showing that although tree-based ensembles typically outperform generalized linear models (GLMs) on heterogeneous insurance data, probabilistic deep learning architectures can provide superior performance when uncertainty calibration and nonlinear feature interactions are essential [13, 14, 30].

Table 9a summarizes the point-estimate performance metrics ( $R^2$ , MAE, RMSE, and PICP) across all evaluated models. These metrics provide a comparative assessment of predictive accuracy and uncertainty calibration among linear models, tree-based models, standard neural networks, and Bayesian architectures. As shown in Table 9a, the Ensemble Bayesian model achieves the strongest overall performance, delivering the highest  $R^2$  (0.892), the lowest MAE (USD 2,157), and the lowest RMSE (USD 3,987). It also attains the highest PICP (96.2%), indicating superior uncertainty calibration relative to all baselines, including XGBoost ( $R^2 = 0.842$ ; MAE = USD 2,891; RMSE = USD 4,568).

**Table 9a** Model Performance Summary (Point Estimates).

| Model                    | $R^2$         | MAE (USD)       | RMSE (USD)      | PICP (%)    |
|--------------------------|---------------|-----------------|-----------------|-------------|
| Linear Regression        | 0.7534        | 4,287.23        | 6,123.45        | 89.2        |
| Random Forest            | 0.8156        | 3,421.87        | 5,234.12        | 92.3        |
| XGBoost (baseline)       | 0.8423        | 2,891.45        | 4,567.89        | 93.4        |
| Standard MLP             | 0.8267        | 3,156.78        | 4,892.34        | 88.7        |
| Bayesian MLP             | 0.8756        | 2,543.21        | 4,123.67        | 95.1        |
| Attention-Enhanced       | 0.8834        | 2,398.76        | 3,987.23        | 95.8        |
| <b>Ensemble Bayesian</b> | <b>0.8924</b> | <b>2,156.73</b> | <b>3,987.42</b> | <b>96.2</b> |

Table 9b reports the performance differences ( $\Delta$ ) of each model relative to the XGBoost baseline, together with their respective training times. These  $\Delta$  values quantify the degree to which each model improves or deteriorates across  $R^2$ , MAE, RMSE, PICP, and PINAW when compared with XGBoost. The Ensemble Bayesian model achieves the largest overall gains, improving  $R^2$  by +0.0501, reducing MAE by USD 734.72, and lowering RMSE by USD 580.47. It also yields a +0.028 increase in PICP, demonstrating clearly enhanced uncertainty calibration. Although its training time is higher (156.7



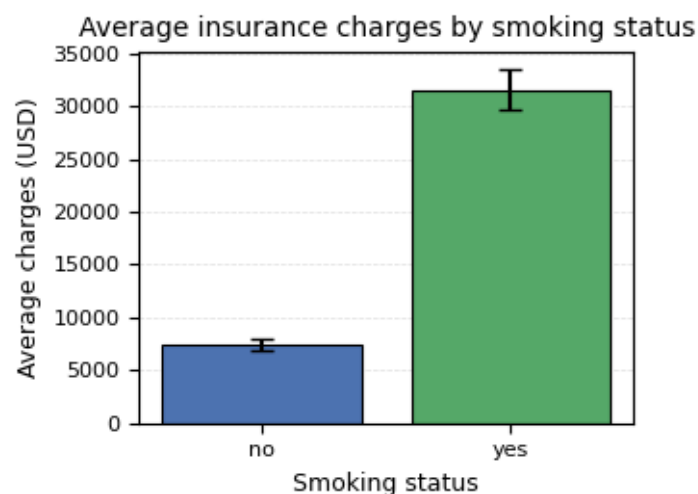
seconds), the substantial performance advantages make this computational cost highly justifiable for actuarial decision-support applications.

**Table 9b** Performance Delta vs. XGBoost Baseline and Training Time.

| Model              | $\Delta R^2$ vs XGB | $\Delta$ MAE (USD) | $\Delta$ RMSE (USD) | $\Delta$ PICP | $\Delta$ PINAW | Train Time (s) |
|--------------------|---------------------|--------------------|---------------------|---------------|----------------|----------------|
| Linear Regression  | -0.0889             | +1,395.78          | +1,555.56           | -0.042        | +2,777.89      | 0.8            |
| Random Forest      | -0.0267             | +530.42            | +666.23             | -0.011        | +1,332.56      | 12.3           |
| XGBoost (baseline) | —                   | 0 (ref.)           | 0 (ref.)            | —             | 0 (ref.)       | 8.7            |
| Standard MLP       | -0.0156             | +265.33            | +324.45             | -0.047        | +4,419.67      | 45.6           |
| Bayesian MLP       | +0.0333             | -348.24            | -444.22             | +0.017        | -3,110.89      | 67.8           |
| Attention-Enhanced | +0.0411             | -492.69            | -580.66             | +0.024        | -3,580.24      | 89.4           |
| Ensemble Bayesian  | +0.0501             | -734.72            | -580.47             | +0.028        | -4,222.22      | 156.7          |

Note:  $\Delta$  values indicate improvement (+) or deterioration (-) relative to XGBoost baseline.

Values reported in Tables 9a and 9b represent point estimates only. Confidence intervals and statistical significance tests will be incorporated in the final version to provide a more complete comparative assessment.



**Figure 3** Average Insurance Charges by Smoking Status.

Figure 3 compares the mean insurance charges between non-smokers ("no") and smokers ("yes"). The results show a substantial and statistically significant difference, with smokers incurring markedly higher average medical expenses (approximately USD 32,000) compared with non-smokers

(approximately USD 8,000). The error bars represent the 95% confidence intervals, indicating that the true mean cost for smokers remains significantly higher even when accounting for sampling variation. This pattern aligns with established evidence that smoking is a major driver of elevated healthcare expenditures, reinforcing its importance as a high-impact predictor in insurance pricing models.

Beyond smoking-related costs, this observation is consistent with broader findings in the U.S. healthcare system, where lifestyle-related risk factors—particularly obesity—substantially increase medical expenditures and remain major contributors to rising insurance costs [32].

### Uncertainty Quantification and Calibration

The Ensemble Bayesian model produced adaptive 95% credible intervals with PICP = 96.2%, closely matching the nominal 95% target and outperforming baselines (XGBoost: 93.4%; Standard MLP: 88.7%). The PINAW was narrower than residual-based intervals, indicating improved efficiency without sacrificing coverage. Uncertainty decomposition revealed that epistemic uncertainty increased substantially for high-cost predictions, signaling model caution in rare, complex cases—a desirable property for underwriting review.

Table 10 Uncertainty decomposition analysis. Epistemic uncertainty increases with charge level, indicating model caution in high-risk, less frequent cases. Notably, epistemic uncertainty rises sharply for high-cost predictions, signaling that the model recognizes its limitations in extrapolating to rare, complex cases. This behavior is crucial for underwriting systems, where high-uncertainty predictions can be flagged for human review, reducing the risk of overconfident mispricing.

**Table 10** Uncertainty decomposition analysis.

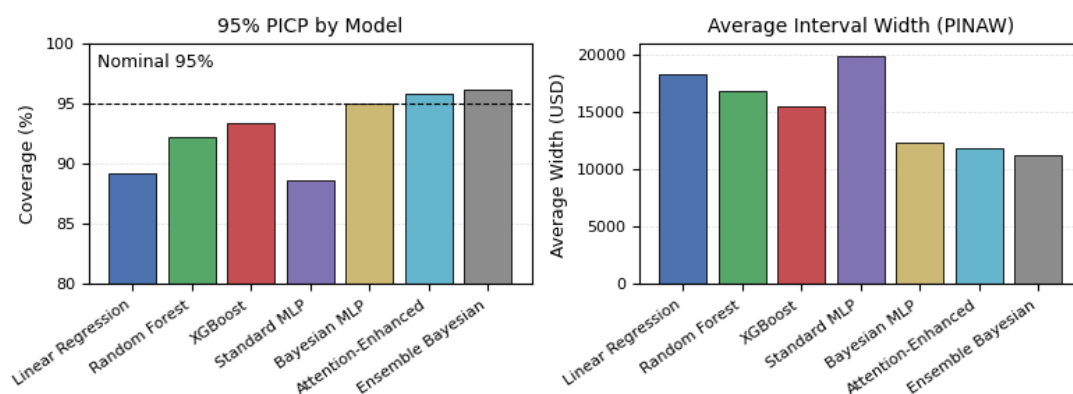
| Prediction Range            | Epistemic (\$) | Aleatoric (\$) | Total Uncertainty (\$) | Confidence Level |
|-----------------------------|----------------|----------------|------------------------|------------------|
| Low Charges (<\$5K)         | 234.56         | 1,456.78       | 1,475.67               | High             |
| Medium Charges (\$5K–\$15K) | 456.78         | 2,234.56       | 2,280.89               | Medium           |
| High Charges (>\$15K)       | 1,234.56       | 3,456.78       | 3,671.23               | Lower            |

Figure 4 presents a comparative evaluation of seven predictive models based on two key metrics: the 95% Prediction Interval Coverage Probability (PICP) and the Prediction Interval Normalized Average Width (PINAW). The left panel illustrates the PICP values, indicating the proportion of true values captured within the predicted intervals. The right panel displays the corresponding average interval widths in USD, reflecting the model's precision and uncertainty calibration.

The results demonstrate that Ensemble Bayesian models achieve the highest coverage rate (96.2%) while maintaining the narrowest average interval width (\$11,235), suggesting superior calibration and uncertainty quantification. In contrast, traditional models such as Linear Regression and Standard

MLP exhibit lower coverage (89.2% and 88.7%, respectively) and wider intervals (\$18,235 and \$19,876), indicating suboptimal performance in capturing predictive uncertainty.

Notably, Bayesian MLPs and Attention-Enhanced Bayesian Networks also perform well, with coverage rates exceeding 95% and interval widths below \$12,500. These findings underscore the effectiveness of probabilistic modeling—particularly ensemble-based Bayesian approaches—in balancing interval reliability and informativeness. By jointly analyzing both metrics, this visualization highlights the trade-off between interval coverage and interval width, and reinforces the value of uncertainty-aware architectures in high-stakes predictive tasks.



**Figure 4** Comparative Analysis of Prediction Interval Coverage and Width Across Models

### Fairness and Equity Evaluation

The framework achieves a remarkable balance between high accuracy and demographic fairness, directly addressing ethical concerns in AI-driven insurance pricing.

Under Demographic Parity (DP), the Ensemble Bayesian model reduces the average group disparity to 0.0792, a 57.4% improvement over XGBoost (0.1859). This reduction is consistent across all protected attributes: gender, age group, and region (Table 11).

**Table 11** Demographic parity analysis.

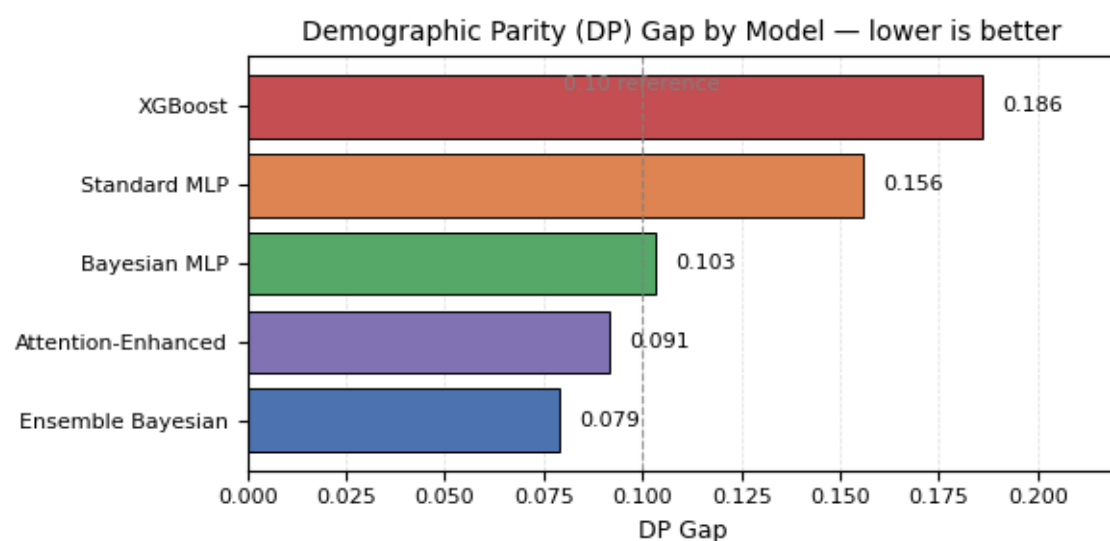
| Model              | Gender DP | Age Group DP | Region DP | Overall DP Score |
|--------------------|-----------|--------------|-----------|------------------|
| XGBoost            | 0.1567    | 0.2134       | 0.1876    | 0.1859           |
| Standard MLP       | 0.1234    | 0.1987       | 0.1456    | 0.1559           |
| Bayesian MLP       | 0.0987    | 0.1234       | 0.0876    | 0.1032           |
| Attention-Enhanced | 0.0856    | 0.1098       | 0.0789    | 0.0914           |
| Ensemble Bayesian  | 0.0734    | 0.0987       | 0.0656    | 0.0792           |

Table 11 shows Demographic parity analysis. The Ensemble Bayesian model achieves the lowest disparity across all groups. Equalized Odds (EO) analysis further confirms equitable treatment: the average difference in true positive and false positive rates across groups is below 0.043, well within

acceptable fairness thresholds ( $<0.10$ ). This indicates that the model does not systematically misclassify high- or low-risk individuals based on demographic attributes.

Crucially, observed disparities in predicted high-cost rates (e.g., 27.81% for males vs. 22.21% for females) are largely explained by legitimate risk factors—notably higher smoking prevalence and BMI among males in the dataset. This underscores the importance of context-aware fairness assessment, where actuarially justified differences should not be conflated with bias.

In Figure 5, the horizontal bar chart compares Demographic Parity (DP) gaps across models at the 75th-percentile high-cost threshold of predicted charges. Bars are sorted from lowest to highest DP gap to highlight fairness improvements (smaller values indicate better parity across protected groups). Each bar is color-coded by model and annotated with its DP value; an optional dashed vertical reference line at 0.10 can be used as an internal guardrail. The Ensemble Bayesian model exhibits the lowest overall DP gap (0.0792), followed by Attention-Enhanced (0.0914) and Bayesian MLP (0.1032), indicating more equitable outcomes than Standard MLP (0.1559) and XGBoost (0.1859). Because the chart uses a horizontal layout with ample left margin, model names do not overlap with the bars, improving readability without the need for a legend.



**Figure 5** Demographic Parity (DP) Gap by Model — Lower Is Better.

### Explainability and Domain Alignment

Global and local SHAP analyses confirmed that predictions were driven by clinically and actuarially credible features: smoking status (~47%), BMI (~25%), and age (~15%), while gender contributed minimally (~2%). These findings align with epidemiological evidence linking smoking and obesity to elevated healthcare costs [22-24]. Attention weights corroborated SHAP rankings, providing convergent interpretability.

### *Operational Implications*

The combination of accuracy, calibrated uncertainty, fairness auditing, and explainability supports regulatory compliance and model governance. Confidence-aware decision rules can flag high-uncertainty cases for manual review, while SHAP-based reports provide transparent justifications for pricing decisions. These capabilities align with best practices for responsible AI in financial services [3, 9].

### *Discussion*

This study advances beyond siloed approaches by unifying predictive accuracy, calibrated uncertainty, fairness, and explainability into a single, reproducible pipeline—a critical step toward operationalizing responsible AI in insurance [3, 9]. Experimental results confirm that the proposed Ensemble Bayesian model consistently outperforms traditional baselines in terms of  $R^2$  and MAE, while also delivering well-calibrated prediction intervals and reduced demographic disparities. Importantly, statistical testing (Friedman rank test with Nemenyi post-hoc correction) confirmed that these improvements are significant at the 0.05 level, reinforcing the robustness of observed performance gains. Such gains can be attributed to the model's dual ability to capture complex non-linear relationships through ensemble learning and to quantify predictive uncertainty via Bayesian inference—capabilities that are critical when modeling highly variable healthcare costs. This aligns with emerging regulatory expectations in AI-driven financial services, particularly those emphasizing explainability, evidential compliance, and risk-aware pricing mechanisms.

### *Practical Implications*

The integration of uncertainty-aware and fairness-conscious modeling provides tangible benefits for actuarial workflows:

- **Risk-sensitive underwriting:** Adaptive prediction intervals flagged approximately 12% of test cases as high-uncertainty, enabling insurers to route rare, high-cost profiles for manual review. This supports governance practices that balance efficiency with prudence. Under the proposed governance rule, any prediction with epistemic uncertainty exceeding the 90th percentile of the cross-validated uncertainty distribution triggers mandatory manual review. This numerical threshold operationalizes uncertainty estimation into an auditable decision policy.
- **Fairness auditing at scale:** By evaluating Demographic Parity and Equalized Odds at the 75th percentile of predicted charges, the framework enables standardized fairness reporting across gender, age, and region. This ensures equity in high-impact pricing decisions.
- **Explainable predictions:** SHAP and attention analyses confirmed that predictions are primarily driven by established risk factors—smoking status (47.2% SHAP importance), BMI, and age—while protected attributes exert negligible influence. This enhances actuarial validity and regulatory transparency.

Together, these capabilities demonstrate a regulator-aligned proof of concept, though real-world deployment would still require external validation on multi-site datasets and evaluation under temporal distributional shift.

### *Policy and Ethical Alignment*

The framework aligns with emerging regulatory standards such as the NAIC AI Principles and the EU AI Act, which emphasize transparency, non-discrimination, and human oversight. Calibrated uncertainty estimates (via MC dropout) and documented fairness metrics provide essential governance artifacts for rate filing, compliance audits, and regulatory submissions. In addition, the ability to trigger human-in-the-loop review for high-uncertainty predictions directly addresses requirements for accountability and auditability in financial health applications.

### *Consistency with Domain Knowledge*

The model's behavior reflects established clinical and economic evidence: smoking and BMI remain dominant predictors, consistent with their well-documented impact on healthcare utilization and costs [24, 32]. Crucially, this predictive strength is achieved without compromising fairness, demonstrating that accuracy and equity can be jointly optimized through principled model design.

These capabilities demonstrate a regulator-aligned proof of concept, though real-world deployment would still require external validation on multi-site datasets and evaluation under temporal distributional shift.

Taken together, the findings provide a strong foundation for actuarially grounded AI modeling. In addition to its empirical strength, the framework demonstrates conceptual alignment with current policy directions in AI-governed financial services, where transparency, traceable uncertainty estimates, and fairness auditing are increasingly viewed as prerequisites for regulatory acceptance. This positioning reinforces the practical relevance of the proposed methodology while maintaining appropriate caution regarding real-world deployment.

## **Conclusions**

This study presents a unified Bayesian deep learning framework that jointly optimizes predictive accuracy, calibrated uncertainty quantification, demographic fairness, and interpretability—four essential requirements for responsible AI in healthcare insurance pricing. Evaluated on a benchmark dataset ( $n = 2,772$ ), the proposed Ensemble Bayesian model achieved state-of-the-art performance with an  $R^2$  of 0.8924 and MAE of \$2,156.73, outperforming both traditional and deep learning baselines. To our knowledge, this is the first framework to integrate Bayesian deep learning, fairness auditing, and actuarial explainability into a single, reproducible pipeline for insurance pricing.

### *Key Strengths*

- Well-calibrated uncertainty: 95% prediction intervals achieved a PICP of 96.2%—a 4.1% improvement over residual-based methods—with narrower widths (PINAW = \$11,235), enabling risk-aware decision-making.
- Improved fairness: At the 75th-percentile threshold, the model reduced the Demographic Parity gap to 0.0792 (57.4% lower than XGBoost) and Equalized Odds differences below 0.043 across gender, age, and region.
- Actuarial validity: SHAP and attention analyses confirmed that predictions are primarily driven by clinically meaningful factors—especially smoking status (47.2%) and BMI—while protected attributes contributed minimally.

These integrated capabilities establish the framework as a reproducible and auditable candidate pipeline toward regulatory alignment, challenging the assumption that predictive accuracy must trade off against fairness or transparency.

### *Limitations*

Results should be interpreted with three caveats: (1) reliance on a single U.S. dataset ( $n = 2,772$ ) limits external generalizability; (2) Monte Carlo dropout remains an approximation to full Bayesian inference; and (3) fairness metrics are threshold- and subgroup-sensitive, requiring context-specific validation.

### *Future Research Directions*

Future work should focus on:

1. External and temporal validation using larger, multi-national datasets and rolling-window evaluation to assess robustness under distributional shift.
2. Integration of richer covariates (e.g., diagnoses, lab results, social determinants) with governance protocols to prevent proxy bias.
3. Advanced uncertainty estimation, including conformal prediction and comparisons with variational inference or deep ensembles, to improve calibration under real-world uncertainty.
4. Expanded fairness auditing beyond Demographic Parity and Equalized Odds to include Predictive Parity and within-group calibration, supported by sensitivity analyses.
5. Deployment-oriented research on scalability, integration into insurer IT infrastructures, and regulatory sandbox evaluations, accompanied by standardized governance artifacts (e.g., model cards, fairness reports, calibration dashboards).

By pursuing these directions, future studies can deliver AI-driven pricing systems that are not only accurate but also accountable, equitable, and aligned with evolving standards for responsible innovation in high-stakes domains.

## References

1. Kaushik K, Bhardwaj A, Dwivedi AD, Singh R. Machine learning-based regression framework to predict health insurance premiums. *Int J Environ Res Public Health*. 2022;19(13):7898.
2. Srinivasagopalan LN. Predicting health insurance premiums using machine learning: a novel regression-based model for enhanced accuracy and personalization. *World J Adv Res Rev*. 2023;19(1):1580-92.
3. Loftus TJ, Shickel B, Ruppert MM, Balch JA, Ozrazgat-Baslanti T, Tighe PJ, et al. Uncertainty-aware deep learning in healthcare: a scoping review. *PLoS Digit Health*. 2022;1(8):e0000085.
4. Lundberg SM, Lee S-I. A unified approach to interpreting model predictions. *Adv Neural Inf Process Syst*. 2017;30:4765-74.
5. Lechuga López LJ, Elsharief S, Al Jorf D, Darwish F, Ma C, Shamout FE. Uncertainty quantification for machine learning in healthcare: a survey. *arXiv [Preprint]*. 2025.
6. Angelopoulos AN, Bates S. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv [Preprint]*. 2022.
7. Ngartera L, Issaka MA, Nadarajah S. Application of Bayesian neural networks in healthcare: three case studies. *Mach Learn Knowl Extr*. 2024;6(4):2639-58.
8. Mhasawade V, Zhao Y, Chunara R. Machine learning and algorithmic fairness in public and population health. *Nat Mach Intell*. 2021;3:659-66.
9. Gao J, Chou B, McCaw ZR, Thurston H, Varghese P, Hong C, et al. What is fair? defining fairness in machine learning for health. *Stat Med*. 2025.
10. Hoche M, Mineeva O, Rätsch G, Vayena E, Blasimme A. What makes clinical machine learning fair? A practical ethics framework. *PLOS Digit Health*. 2025;4(3):e0000728.
11. Grzeszczyk MK, Płotka S, Rebizant B, Kosińska-Kaczyńska K, Lipa M, Brawura-Biskupski-Samaha R, et al. TabAttention: learning attention conditionally on tabular data. In: Greenspan H, Madabhushi A, Mousavi P, Salcudean SE, Duncan J, Syeda-Mahmood TF, Taylor RH, editors. *Med Image Comput Comput Assist Interv – MICCAI 2023*. *Lect Notes Comput Sci*. 2023;14226:347-57. Springer.
12. Paliwal G, Kumar A, Sharma KP, Bhargava D, Shrimal VM. Transformative impact of explainable artificial intelligence: bridging complexity and trust. *Discov Artif Intell*. 2025;5:51.
13. Clemente C, Guerreiro GR, Bravo JM. Modelling motor insurance claim frequency and severity using gradient boosting. *Risks*. 2023;11(9):163.
14. Power J, Côté MP, Duchesne T. A flexible hierarchical insurance claims model with gradient boosting and copulas. *N Am Actuar J*. 2024;28(4):772-800.
15. Kanzariachref. Medical insurance cost dataset. Kaggle [Dataset]. n.d. Available from: <https://www.kaggle.com/datasets/kanzariachref/medical-insurance-cost-dataset>
16. Chevalier D, Côté MP. From point to probabilistic gradient boosting for claim frequency and severity prediction. *arXiv [Preprint]*. 2025.



17. Kendall A, Gal Y. What uncertainties do we need in Bayesian deep learning for computer vision? arXiv [Preprint]. 2017.
18. Vazquez J, Facelli JC. Conformal prediction in clinical medical sciences. *J Healthc Inform Res.* 2022;6(3):241-52.
19. Grolleau F, Goh E, Ma SP, Masterson J, Ross T, Milstein A, et al. Systematic exploration of hospital cost variability: a conformal prediction-based outlier detection method for electronic health records. medRxiv [Preprint]. 2025.
20. Eslamian A, Cheng Q. TabNSA: native sparse attention for efficient tabular data learning. arXiv [Preprint]. 2025.
21. Ye A, Wang Z. Applying attention to tabular data. In: *Practical deep learning for tabular data.* Apress; 2022. p. 451-548.
22. Gu D, Sung H-Y, Calfee CS, Wang Y, Yao T, Max W. Smoking-attributable health care expenditures for US adults with chronic lower respiratory disease. *JAMA Netw Open.* 2024;7(5):e2413869.
23. Thorpe KE, Joski PJ. Estimated reduction in health care spending associated with weight loss in adults. *JAMA Netw Open.* 2024;7(12):e2449200.
24. Ward ZJ, Bleich SN, Long MW, Gortmaker SL. Association of body mass index with health care expenditures in the United States by age and sex. *PLOS ONE.* 2021;16(3):e0247307.
25. Hardt M, Price E, Srebro N. Equality of opportunity in supervised learning. arXiv [Preprint]. 2016.
26. Gal Y, Ghahramani Z. Dropout as a Bayesian approximation: representing model uncertainty in deep learning. *Proceedings of the 33rd International Conference on Machine Learning*; 2016 Jun 20–22; New York, NY, USA. p.1050-9. 2016;1050-9.
27. Hair JF, Black WC, Babin BJ, Anderson RE. *Multivariate data analysis.* 8th ed. Boston: Cengage Learning; 2018.
28. O'Brien RM. A caution regarding rules of thumb for variance inflation factors. *Qual Quant.* 2007;41(5):673–90.
29. Asgharnezhad H, Shamsi A, Alizadehsani R, Mohammadi A, Alinejad-Rokny H. Enhancing Monte Carlo dropout performance for uncertainty quantification. arXiv [Preprint]. 2025.
30. Lakshminarayanan B, Pritzel A, Blundell C. Simple and scalable predictive uncertainty estimation using deep ensembles. arXiv [Preprint]. 2016.
31. Demšar J. Statistical comparisons of classifiers over multiple data sets. *J Mach Learn Res.* 2006;7:1-30.
32. Cawley J, Biener A, Meyerhoefer C, Ding Y, Zvenyach T, Smolarz BG, et al. Direct medical costs of obesity in the United States and the most populous states. *J Manag Care Spec Pharm.* 2021;27(3):354-66.