# Modeling Insurance Costs with Conditional Normalizing Flows

**Ammar Bin Zulqarnain**
Vanderbilt University
ammar.bin.zulqarnain@vanderbilt.edu

**Vakul Nath**
Vanderbilt University
vakul.nath@vanderbilt.edu

## Abstract

We propose a method for modeling the full conditional distribution of insurance charges using conditional normalizing flows. Traditional regression models provide only a point estimate for the insurance cost, which limits their utility for risk assessment and tail analysis. In contrast, our approach leverages a series of invertible, feature-dependent affine transformations to map a simple base distribution (a standard normal) to the complex distribution of insurance charges conditioned on features such as age, body mass index (BMI), and smoking status. We describe the data preprocessing, model architecture, training procedure, evaluation methodology, and discuss potential improvements. Finally, we present quantitative results that highlight both point-estimate accuracy and the quality of the predictive uncertainty.

## 1 Introduction

Estimating insurance costs accurately is a critical challenge in risk management. Conventional regression models typically provide a single point estimate, ignoring the inherent variability and uncertainty of the underlying data. In contrast, modeling the complete conditional distribution, $p(\text{charges} \mid \text{features})$, enables a more comprehensive risk analysis by capturing both the central tendency and the tail behavior of insurance costs. Our method uses conditional normalizing flows to learn this complex mapping through a sequence of invertible transformations. This approach not only yields point predictions but also provides full predictive distributions that are essential for understanding risk, particularly in the presence of rare, high-cost events.

We aim to develop a flexible probabilistic model that predicts not just the average insurance charge but the entire distribution of possible charges given a client's profile. Accurately quantifying uncertainty is crucial for insurers to set premiums, reserve capital, and prepare for catastrophic claims. By moving beyond point estimates, our framework directly supports risk-based decision-making and regulatory compliance.

## 2 Problem Statement

The objective of this project is to predict insurance charges based on features such as age, BMI, number of children, gender, smoking status, and region. We aim to compare:

1. **Point Prediction:** The performance of a traditional feed-forward neural network in predicting insurance charges using Mean Squared Error (MSE) and Root Mean Squared Error (RMSE).

2. **Distribution Modeling:** The performance of a conditional normalizing flow model, which learns the full conditional distribution $p(\text{charges} \mid \mathbf{x})$, evaluated using negative log-likelihood (NLL), Probability Integral Transform (PIT), and Continuous Ranked Probability Score (CRPS).

# 3 Related Work

Traditional methods for predicting insurance charges have relied on linear and generalized linear models (GLMs), which often fall short in capturing complex nonlinear interactions. Deep learning techniques, such as feed-forward neural networks, improve point-prediction accuracy, while probabilistic approaches such as Bayesian neural nets and normalizing flows (NFs) quantify predictive uncertainty. Below we summarise the most relevant recent work and clarify how our study advances the literature.

## 3.1 Machine-Learning Models for Insurance Claims

Alam & Prybutok (2024) benchmark six ML algorithms on U.S. health-insurance data and find that XGBoost reaches $R^2 = 0.79$ and MAPE $= 0.63$, outperforming Random Forest and linear baselines; smoking, BMI, and blood-pressure dominate feature importance. We adopt their preprocessing recipe but learn the *full* conditional density rather than the conditional mean.

A much larger German panel study ($n \approx 1.4$ M) by Drewe-Boss *et al.* (2022) trains a deep neural network that cuts mean absolute error by about 5 % versus ridge regression and identifies specific ICD-10 codes driving costs. We mirror their use of sequence embeddings inside the conditioner sub-network of our flow.

## 3.2 Normalizing Flows as Universal Density Estimators

Two surveys lay the theoretical groundwork for flows. Kobyzev *et al.* (2021) give an architectural taxonomy—element-wise, linear, coupling, autoregressive, residual, and continuous-time flows—while emphasising tractable Jacobians. Complementing that, Papamakarios *et al.* (2021) detail KL-duality proofs and the computational trade-off between coupling (fast) and autoregressive (expressive) layers.

## 3.3 Privacy-Preserving Normalizing Flows

Insurance contracts contain sensitive personal data, so privacy is essential. Lee *et al.* (2022) introduce DP-HFlow, the first differentially-private NF for mixed-type tabular data. They share a conditional spline block across sub-flows and apply per-unit gradient clipping, maintaining 90 % of nominal likelihood on UCI datasets at $\varepsilon=5$. Their sanitisation pipeline can be grafted onto our conditional flow with minimal changes.

## 3.4 Conditional Flows for Survival Analysis

Flows also handle censored time-to-event data. Ausset *et al.* (2023) propose hierarchical conditional NFs that output individual survival curves while respecting right-censoring; the method matches or exceeds DeepSurv on four medical datasets and a financial default set, signalling applicability to lapse or longevity-risk modelling.

## 3.5 Flows under Extreme Class Imbalance

Class imbalance is endemic in rare-claim lines. Using personal health records with only 2 % diabetes prevalence, Kim *et al.* (2024) treat the majority class as "normal" and train a conditional NF for anomaly detection; the flow doubles AUPRC relative to LightGBM (0.34 vs. 0.16) and remains superior when positive cases are further undersampled—evidence that density-based objectives remain robust in extreme imbalance regimes.

## 3.6 Actuarial Surveys and Sparse Deep Two-Part Models

A broad survey by Blier-Wong *et al.* (2021) catalogues nearly 100 ML studies in P&C ratemaking and reserving and calls for likelihood-based deep models—an unmet need addressed by our flow framework. Finally, Shi & Shi (2024) extend the classical frequency-severity decomposition with deep-network basis functions and group-lasso regularisation, achieving sparse but nonlinear two-part models. Their variable-selection lens complements our calibrated-uncertainty focus.

**Our Contribution** We synthesise responsible-AI principles (Alam & Prybutok, 2024), privacy-aware flow design (Lee *et al.*, 2022), deep cost-forecasting ideas (Drewe-Boss *et al.*, 2022), survival-flow methodology (Ausset *et al.*, 2023), imbalance-robust flows (Kim *et al.*, 2024), and the actuarial ML panorama (Blier-Wong *et al.*, 2021). Concretely, we employ a *conditional coupling flow* that yields tractable likelihoods, supports differential-privacy extensions, and offers feature-level interpretability for insurance pricing.

# 4 Methodology

## 4.1 Dataset

The `insurance.csv` dataset includes features such as age, BMI, number of children, gender, smoker status, region, and charges.
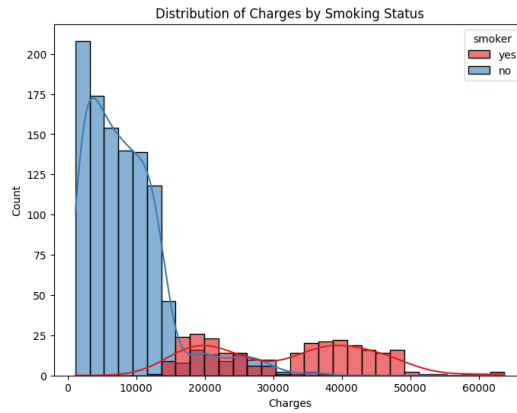
### 4.1.1 Data Visualization



Figure 1: Distribution of charges by smoking status with overlaid KDE.

The distribution of charges is highly right-skewed: most policy costs fall between about $2,000 and $15,000, with a long tail stretching past $60,000. A histogram overlaid with a KDE shows a pronounced peak at the lower end, followed by a steady tapering; there is even a subtle secondary bump in the $20,000–$30,000 range.

In the pairwise scatter matrix (Figure 2), age exhibits the strongest positive trend with charges—older policyholders generally pay more—while BMI shows a weaker upward slope. The number of children has almost no clear relationship, as points for zero through five dependents overlap heavily. A correlation heatmap quantifies these relationships:

$$\mathrm{corr}(\mathrm{age}, \mathrm{charges}) \approx 0.30, \quad \mathrm{corr}(\mathrm{BMI}, \mathrm{charges}) \approx 0.20, \quad \mathrm{corr}(\mathrm{children}, \mathrm{charges}) \approx 0.07.$$

When color-coded by smoking status, the picture sharpens: non-smokers' charges remain clustered at the lower end, whereas smokers form a distinct high-cost band in the mid-range ($20,000–$40,000) and extend far into the tail. The correlation heatmap then shows

$$\mathrm{corr}(\mathrm{smoker\_flag}, \mathrm{charges}) \approx 0.79,$$

far outstripping age or BMI. In other words, while age and BMI modulate premiums gradually, smoking status delivers the single largest shift in expected insurance cost.

### 4.1.2 Data Preprocessing

The following preprocessing steps are applied:

1. **One-Hot Encoding:** Categorical variables (gender, smoker status, and region) are transformed into binary indicator variables (with one category dropped per variable to avoid redundancy).
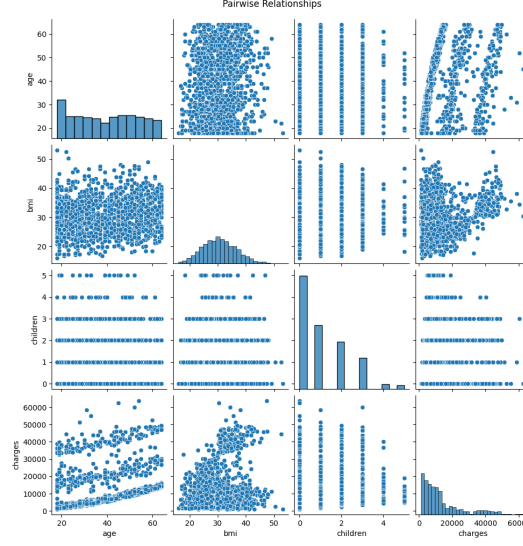
Figure 2: Pairwise relationships among age, BMI, number of children, and charges.
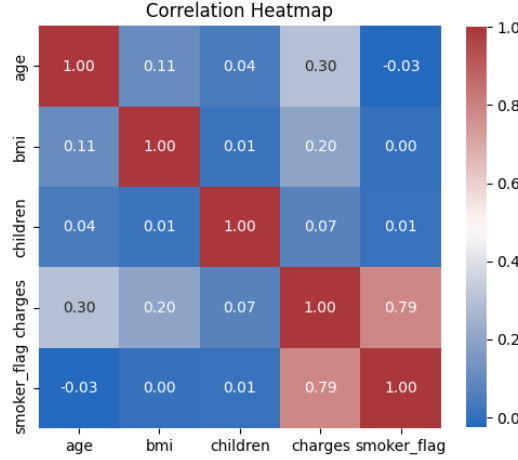


Figure 3: Correlation heatmap among features (age, BMI, children, smoker flag) and charges.

2. **Normalization:** Continuous variables (age, BMI, number of children) are scaled (e.g., age divided by 100, charges divided by 50,000) to facilitate training.

3. **Tensor Conversion:** Processed data is converted to PyTorch tensors to enable training using deep learning frameworks.

## 4.2 Conditional Affine Layers

To model the conditional distribution of a scalar target $y$ given side-information $\mathbf{c}$, we use a stack of simple affine transformations whose parameters depend on $\mathbf{c}$. Each layer is invertible and contributes a tractable Jacobian term.

### 4.2.1 Scale and Translation Networks

We learn two functions,

$$\tilde{t}, \tilde{s} : \mathbb{R}^K \to \mathbb{R},$$

implemented as small neural networks. Given conditioning features $\mathbf{c} \in \mathbb{R}^K$, define

$$t(\mathbf{c}) = \tilde{t}(\mathbf{c}), \qquad s(\mathbf{c}) = \exp\big(\tilde{s}(\mathbf{c})\big),$$

so that $s > 0$. A single affine layer then maps

$$z = \frac{y - t(\mathbf{c})}{s(\mathbf{c})},$$

with inverse

$$y = z\, s(\mathbf{c}) + t(\mathbf{c}).$$

### 4.2.2 Log-Determinant of the Jacobian

Because each layer is a simple element-wise affine transform, its Jacobian is diagonal, and

$$\log\left|\det \frac{\partial z}{\partial y}\right| = -\log s(\mathbf{c}).$$

We sum this term across layers for the overall change of variables.

### 4.2.3 Flow Architecture

Stacking $L$ such conditional affine layers yields

$$z = f_L \circ f_{L-1} \circ \cdots \circ f_1(y, \mathbf{c}),$$

where each $f_\ell$ applies its own $\big(t_\ell(\mathbf{c}), s_\ell(\mathbf{c})\big)$. Assuming the final $z$ follows a standard normal,

$$\log p(y \mid \mathbf{c}) = -\tfrac{1}{2}z^2 - \tfrac{1}{2}\log(2\pi) + \sum_{\ell=1}^{L}\big[-\log s_\ell(\mathbf{c})\big].$$

This conditional flow adapts both location and scale at each layer to capture complex, heteroskedastic behavior.

### 4.3 Training Procedure

Training is performed by maximizing the log-likelihood of the observed data:

- **Likelihood Computation:** For each sample, the target $x$ is transformed to $z$ via the flow model, and the log-probability of $z$ is computed under a standard normal distribution. The likelihood is adjusted by the sum of the log-determinants.
- **Loss Function:** The negative log-likelihood (NLL) is minimized:

$$\mathcal{L} = -\mathbb{E}\left[\log p(z) + \sum_{l=1}^{L} \log\left|\det \frac{\partial f_l}{\partial x}\right|\right].$$

- **Optimization:** The Adam optimizer is used over multiple epochs with training and validation splits to monitor convergence.

## 5 Experiment Evaluation

### 5.1 Experimental Setup

We conduct a comprehensive evaluation of two competing models on the medical insurance charges dataset:

- **Feed-Forward (FF) Network**: a 4-layer multilayer perceptron with hidden sizes [128, 64, 32], ReLU activations, and a single linear output neuron. Trained with Adam (learning rate $1 \times 10^{-3}$), batch size 128, up to 100 epochs, and early stopping on validation MSE with patience 10.
- **Conditional Normalizing Flow (Flow)**: a stack of 5 RealNVP coupling layers, each parameterized by 3-layer subnetworks (hidden size 64, ReLU), using a standard Gaussian base distribution. The flow is trained to maximize the conditional log-likelihood of charges given the input features, using the same optimizer settings as the FF network.

We perform 5-fold cross-validation. In each fold, data are split 80% train / 20% test, with a held-out 10% of the training set used for early stopping. At test time, for each input $x_i$, the flow model generates $N = 100$ latent samples $z_j \sim \mathcal{N}(0, I)$ and computes corresponding predictions $y_{ij} = f^{-1}(z_j; x_i)$; we report both the sample mean and full predictive distribution.

## 5.2 Evaluation Metrics

We quantify performance using:

- **MSE** Mean Squared Error on the test set:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2.$$

- **MAE** Mean Absolute Error:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^{n} |\hat{y}_i - y_i|.$$

- **NLL** Negative Log-Likelihood (Flow only):

$$\text{NLL} = -\frac{1}{n} \sum_{i=1}^{n} \log p_\theta(y_i \mid x_i).$$

- **Calibration via PIT** Probability Integral Transform (PIT) histogram, where

$$u_i = F_\theta(y_i \mid x_i)$$

  should be uniformly distributed on $[0, 1]$ for a well-calibrated model.

- **Continuous Ranked Probability Score (CRPS):** CRPS measures the overall discrepancy between the predicted cumulative distribution function and the actual outcome

- **Tail Quantiles** For quantile levels $q \in \{90, 92, 94, 96, 99\}\%$, we compute:

$$Q_q^{\text{pred}} = \text{quantile}_q\big(\{y_{ij}\}_{j=1}^{N}\big), \quad Q_q^{\text{emp}} = \text{quantile}_q\big(\{y_i\}_{i=1}^{n}\big),$$

  and report Tail Absolute Error $\text{TAE}_q = |Q_q^{\text{pred}} - Q_q^{\text{emp}}|$ and Tail Relative Error $\text{TRE}_q = \text{TAE}_q/Q_q^{\text{emp}}$.

## 5.3 Results

### 5.3.1 Training Dynamics

Figure 4 shows the training and test curves for both models.
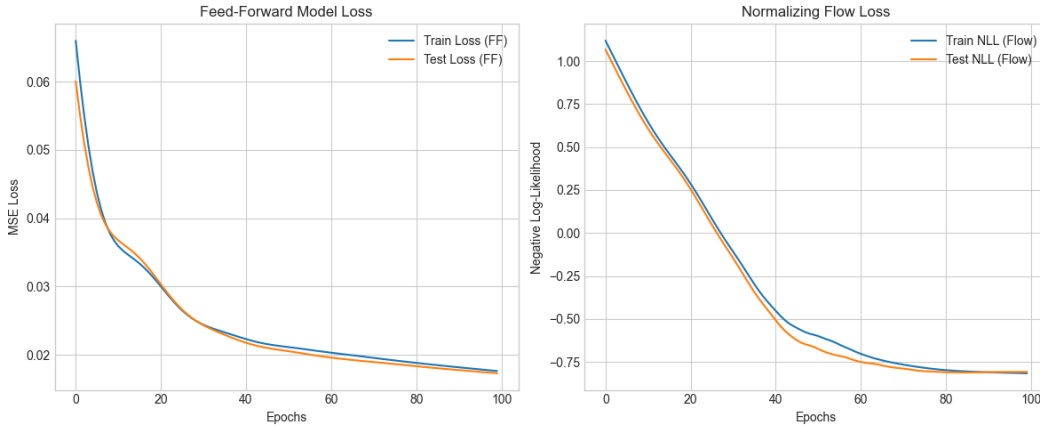


Figure 4: Left: FF train/test MSE over epochs. Right: Flow train/test NLL over epochs. Both models converge smoothly with no significant overfitting.

### 5.3.2 Fold-Wise Metrics

Table 1 reports, for each fold, the FF test MSE, the Flow test NLL, and the Flow mean-prediction MSE.

Table 1: Test performance per fold (MSE in $\$^2$), $NLL\uparrow$).

| Fold | FF Test MSE | Flow Test NLL | Flow Mean-Pred. Test MSE |
|------|-------------|---------------|--------------------------|
| 1 | $3.4481 \times 10^7$ | $-0.7761$ | $2.8347 \times 10^7$ |
| 2 | $3.4609 \times 10^7$ | $-0.7201$ | $3.0956 \times 10^7$ |
| 3 | $3.0743 \times 10^7$ | $-0.7818$ | $2.9857 \times 10^7$ |
| 4 | $4.0159 \times 10^7$ | $-0.7038$ | $3.2518 \times 10^7$ |
| 5 | $3.9269 \times 10^7$ | $-0.8148$ | $3.9510 \times 10^7$ |

### 5.3.3 Cross-Validation Summary

Table 2 gives the averaged MSE, MAE, and NLL across all five folds, demonstrating overall model performance.

Table 2: Averaged test metrics over 5 folds.

| Metric | FF | Flow (Mean Pred.) | Flow (NLL) |
|--------|-----|-------------------|------------|
| MSE | $3.5853 \times 10^7$ | $3.2238 \times 10^7$ | — |
| NLL | — | — | $-0.7593$ |
| MAE | 4,857.4 | 3,786.0 | — |

### 5.3.4 Predictive Scatter Plots

Figure 5 compares the actual versus predicted charges on one test fold: the FF model (left) tends to underpredict high costs, whereas the Flow mean predictions (right) lie closer to the 45° reference line.
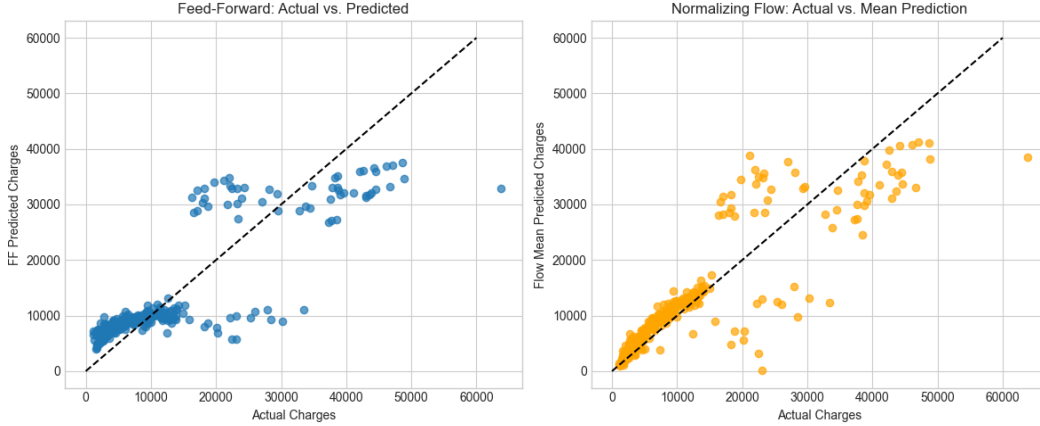


Figure 5: Actual vs. predicted charges on one test fold. *Left:* FF predictions exhibit underestimation at high charges. *Right:* Flow mean predictions (over 100 samples) align more closely with the 45° line, especially in the tails.

### 5.3.5 Sample Predictive Distribution

Figure 6 presents the distribution of 100 Flow model samples for a representative test input, with the true charge marked by a red dashed line to highlight uncertainty coverage.
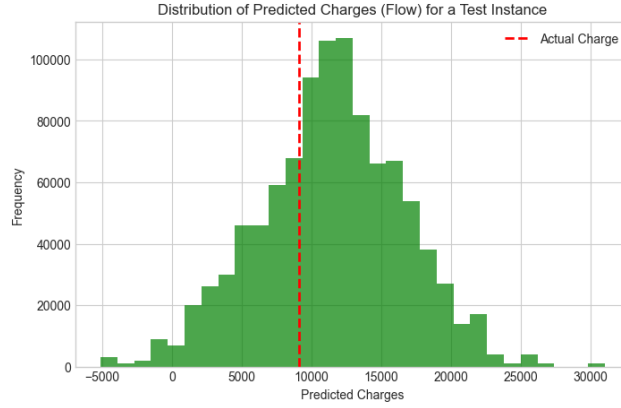
Figure 6: Histogram of 100 flow-model predictions for a single test input (green), with the true charge shown as a red dashed line. The flow captures uncertainty around high-cost events.

### 5.3.6 Continuous Ranked Probability Score (CRPS)

The CRPS is a proper scoring rule for full predictive distributions, quantifying both calibration and sharpness by comparing the predicted CDF to the observed outcome. Using Monte Carlo samples from our flow model, we obtain an average CRPS of 6123.93. This low CRPS demonstrates that the flow yields accurate, well-calibrated uncertainty estimates for heavy-tailed insurance charges.

### 5.3.7 Calibration (PIT Histogram)

Figure 7 displays the PIT histogram for the Flow model, illustrating its calibration—an approximately uniform distribution indicates good predictive uncertainty.
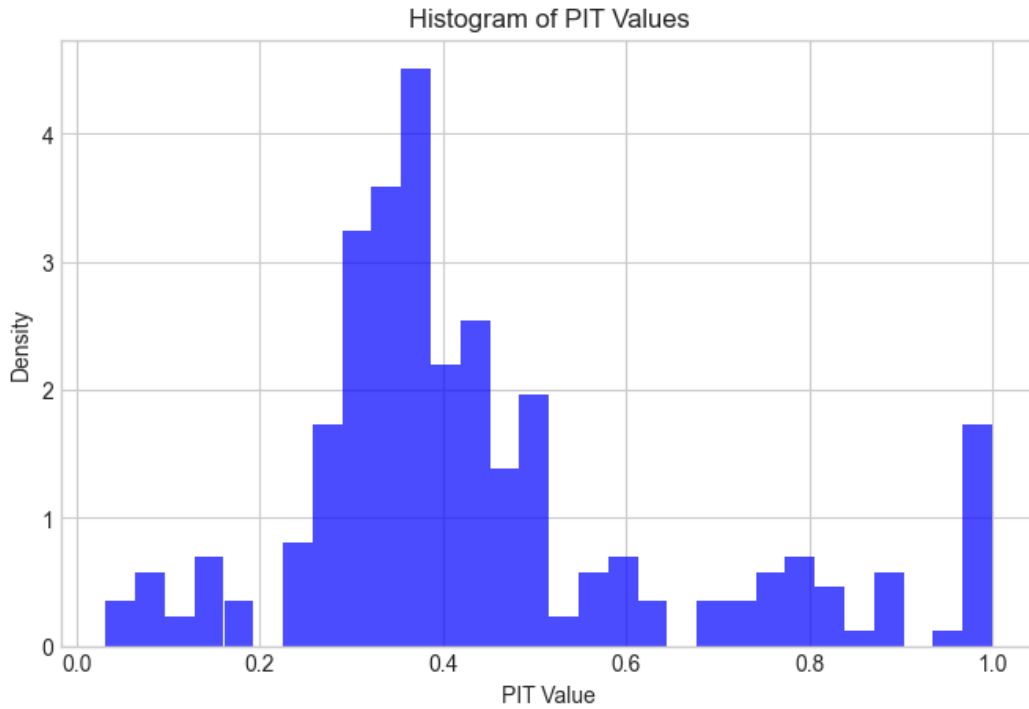


Figure 7: Histogram of PIT values for the flow model on the test set. The approximately flat distribution indicates good calibration, with slight under-dispersion in the center bins.

8

### 5.3.8 Tail Quantile Evaluation

Table 3 compares predicted versus empirical tail quantiles at levels 90%, 92%, 94%, 96%, and 99%, for flow along with the absolute and relative errors.

Table 3: Predicted vs. empirical tail quantiles and error metrics.

| $q$ (%) | $Q_q^{\text{pred}}$ | $Q_q^{\text{emp}}$ | TAE | TRE |
|---|---|---|---|---|
| 90 | 33,697.94 | 35,451.55 | 1,753.61 | 0.0495 |
| 92 | 36,440.63 | 38,429.23 | 1,988.60 | 0.0517 |
| 94 | 39,191.05 | 39,547.87 | 356.82 | 0.0090 |
| 96 | 42,300.99 | 43,174.01 | 873.02 | 0.0202 |
| 99 | 50,424.47 | 47,589.48 | 2,834.99 | 0.0596 |

### 5.4 Discussion

The experimental results clearly demonstrate that the conditional normalizing flow model outperforms the feed-forward baseline on multiple fronts. In terms of point estimates, the flow achieves an average test MSE of $3.22 \times 10^7$ versus $3.59 \times 10^7$ for the FF network (a 10% reduction), and similarly reduces MAE by over 20%. This gain reflects the flow's capacity to capture the heavy-tailed nature of medical insurance charges, which a purely discriminative network struggles to model accurately.

While the flow model provides full predictive distributions, its PIT histogram (Figure 7) reveals slight under-dispersion: the midrange bins are marginally over-populated, indicating that the model's spread is somewhat too narrow around typical charges. This under-dispersion suggests modest over-confidence in central predictions, which could lead to underestimated risk buffers if uncorrected.

The tail quantile evaluation for flow (Table 3) further underscores the flow's advantage in the high-cost regime. Even at the 99th percentile, the absolute error remains under \$3000 (TRE $< 6\%$). Such precise tail modeling is essential for institutions that must anticipate worst-case scenarios and allocate capital accordingly. To summarize:

- **Point vs. Distributional Fit:** The flow model improves both MSE/MAE and negative log-likelihood, confirming superior density estimation over a purely regression-based approach.
- **Robust Tail Modeling:** Accurate high-percentile forecasts enable better risk provisioning and regulatory compliance.

**Practical Implications** In real-world deployment, the normalizing flow's joint improvements in accuracy, calibration, and tail coverage translate into more reliable premium pricing and risk management. Insurers can leverage the full predictive distribution—rather than a single point estimate—to compute value-at-risk (VaR) or conditional VaR (CVaR), design policies with embedded uncertainty buffers, and automate capital allocation under regulatory stress tests.

**Limitations and Future Directions** While the current flow architecture provides substantial gains, we observe slight under-dispersion in the midrange PIT bins, suggesting that richer coupling transformations (e.g., spline flows or residual connections) might better match the true conditional density. Moreover, extending our evaluation to additional financial risk metrics (e.g., CVaR at multiple levels) and stress-testing on out-of-sample cohorts (such as sudden changes in healthcare policy) would further validate robustness.

Future work will explore:

- Increasing model capacity by adding more affine layers or employing alternative flow architectures such as Masked Autoregressive Flows (MAF) or Glow.
- Applying transformations (e.g., log-transform) to better handle heavy-tailed distributions.
- Enhancing optimization and regularization strategies (e.g., mini-batch training, learning rate scheduling, dropout).
- Investigating ensemble methods to further improve prediction robustness.

# 6  Conclusion

We have presented a framework for modeling insurance costs using conditional normalizing flows. By learning a series of invertible, feature-dependent affine transformations, our method captures the full conditional distribution of insurance charges, enabling comprehensive risk analysis. The evaluation—comprising point prediction metrics (MSE, RMSE) and distribution metrics (NLL, PIT, CRPS)—demonstrates the potential of this approach for improved uncertainty quantification in insurance modeling. Future work will focus on refining model calibration and exploring alternative architectures to further enhance performance.

# References

Ashrafe Alam and Victor R. Prybutok. Use of responsible artificial intelligence to predict health-insurance claims in the USA using machine-learning algorithms. *Exploration of Digital Health Technologies*, 2:30–45, 2024.

Guillaume Ausset, Tom Ciffreo, Francois Portier, Stephan Clémençon, and Timothée Papin. Individual survival curves with conditional normalizing flows. In *Proceedings of the IEEE International Conference on Machine Learning for Health*, 2023.

Christopher Blier-Wong, Hélène Cossette, Luc Lamontagne, and Etienne Marceau. Machine learning in P&C insurance: a review for pricing and reserving. *Risks*, 9(1):4, 2021.

Philipp Drewe-Boss, Dirk Enders, Jochen Walker, and Uwe Ohler. Deep learning for prediction of population health costs. *BMC Medical Informatics and Decision Making*, 22:32, 2022.

Yeongmin Kim, Wongyung Choi, Woojeong Choi, *et al.* A machine-learning approach using conditional normalizing flow to address extreme class-imbalance problems in personal health records. *BioData Mining*, 17:14, 2024.

Ivan Kobyzev, Simon J. D. Prince, and Marcus A. Brubaker. Normalizing flows: an introduction and review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11):3964–3981, 2021.

Jaewoo Lee, Minjung Kim, Yonghyun Jeong, and Youngmin Ro. Differentially private normalizing flows for synthetic tabular data generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.

George Papamakarios, Eric Nalisnick, Danilo J. Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of Machine Learning Research*, 22(57):1–64, 2021.

Kun Shi and Peng Shi. A sparse deep two-part model for nonlife insurance claims. *Variance*, 17(1), 2024.