

# Accurate and Interpretable Machine Learning for Transparent Pricing of Health Insurance Plans

**Rohun Kshirsagar, Li-Yen Hsu, Vatshank Chaturvedi\*, Charles H. Greenberg, Matthew McClelland, Hilaf Hasson\*, Anushadevi Mohan, Wideet Shende, Nicolas P. Tilmans, Renzo Frigato, Min Guo, Ankit Chheda, Meredith Trotter, Shonket Ray, Arnold Lee, Miguel Alvarado**  
Lumiata Data Science and Engineering

rohun@lumiata.com

\*Work completed prior to joining Amazon

## Abstract

Health insurance companies cover half of the United States population through commercial employer-sponsored health plans and pay 1.2 trillion US dollars every year to cover medical expenses for their members. The actuary and underwriter roles at a health insurance company serve to assess which risks to take on and how to price those risks to ensure profitability of the organization. While Bayesian hierarchical models are the current standard in the industry to estimate risk, interest in machine learning as a way to improve upon these existing methods is increasing. Lumiata, a healthcare analytics company, ran a study with a large health insurance company in the United States. We evaluated the ability of machine learning models to predict the per member per month cost of employer groups in their next renewal period, especially those groups who will cost less than 95% of what an actuarial model predicts (groups with “concession opportunities”). We developed a sequence of two models, an individual patient-level and an employer-group-level model, to predict the annual per member per month allowed amount for employer groups, based on a population of 14 million patients. Our models performed 20% better than the insurance carrier’s existing pricing model, and identified 84% of the concession opportunities. This study demonstrates the application of a machine learning system to compute an accurate and fair price for health insurance products and analyzes how explainable machine learning models can exceed actuarial models’ predictive accuracy while maintaining interpretability.

## Introduction

The recent explosion of available electronic health record (EHR) and insurance claims data sets, coupled with the democratization of statistical learning algorithms, has set the stage for machine learning (ML) applications to fundamentally transform the healthcare industry. Employer-sponsored health insurance (ESI) currently covers 150 million Americans (?). With numerous subsidies in place to make ESI more affordable (??), it is by far the most popular option for obtaining health insurance in the United States (?). Since the passage of the Affordable Care Act (ACA) in 2009, premiums for single and family ESI plans have increased 50% and deductibles have doubled, making affordability of care

a major issue for many Americans. Thirty four percent of patients on ESI plans are reportedly unable to pay for an unexpected bill of \$500 and over 50% skip or postpone medical care and prescription fills due to cost (?). By making healthcare more affordable, it increases the chance patients will receive needed medical care and refill medications in a timely manner. Increasing affordability improves patient health outcomes and quality of life, reducing familial strain of medical debt, postponing major household spending, and obviating the need to hold multiple jobs (?).

However, ESI premiums are the dominating source of revenue for US-based health insurance companies, many of which are among the Fortune 500 companies. As such, accurate rate-setting is a crucial component of revenue and membership growth for insurance companies; setting inadequate rates can mean the difference between profitability and unprofitability (?). Traditionally, health insurance companies set rates using a combination of actuarial science and underwriting judgement. Actuaries apply statistical modeling to claims data to set premiums for each employer-group; underwriters use the actuary’s predicted rate alongside non-claims data (e.g. health questionnaires) to decide which groups to cover and what their rates will be. Accurate rate setting is essential to balance customer retention against business viability because the insurer needs to retain a group for several years before the account becomes profitable. Hence, insurers are willing to reduce rates in the near-term, in exchange for the chance of a longer term relationship (i.e. greater *persistence*) (?). Good financial standing allows insurers to focus on growing their business, influencing patient health outcomes, improving customer experience, and increasing efficiency (?). Reduced renewal premiums can align patients’ financial interests and the insurers strategic interests.

Within the ESI market, the <500 employer group segment (employer groups with fewer than 500 enrollees) is highly transactional, particularly during peak season (near January 1st of each year). The largest insurance carriers process tens of thousands of new and renewal business quotes. For instance, the carrier in this pilot study averaged 130 presale quotes and 70 renewal quotes per under-

writer. Some aspects unique to the <500 market are the: (1) Complex array of funding arrangements, including: fully-insured, level-funding Administrative Services Only (ASO, i.e. self-insured), pay-as-you-go and monthly shared-risk models (?) (2) high-risk, high-reward nature of this segment: the <500 market comprises one third of medical insurance customers but yields one half of earnings, providing opportunities to yield higher profit margins, and (3) lower persistency compared to larger clients (the average account length is 5 years for the <500 market vs >7 years for larger groups).

Due to the large volume of cases in the < 500 market, underwriters could benefit from an additional highly-accurate signal to increase the efficiency of their work. This signal would need to be easily interpretable and broadly applied for measurability. The unique circumstances of the <500 market provide an opportunity for improved renewal rate development, driving price optimization and increasing persistency in a risk-rich environment focused on strong, sustainable customer-base growth.

Although ML-based approaches on claims data are widely applied in clinical contexts (e.g. ? ?; ? ?; ? ?), few studies incorporate ML approaches for underwriting group health insurance. To our knowledge, there are no models of individual health risk used to form more competitive pricing for a group. Here, we show that ML approaches in medical underwriting can provide: (1) improved accuracy and broader applicability, facilitated by modeling cost and risk at the individual and group levels; and (2) interpretability for non-experts to use.

We jointly ran our study between Lumiata, a healthcare analytics company, and a large insurance company in the US, referred to here with the pseudonym “Delphi”.

## Materials and Methods

### Participants and Setting

We drew from a 14 million member population representative of Delphi’s entire fully-insured employer-group customer base (i.e. their *book of business*) from 2015 to 2017. We used this population to train and tune our models and then predicted the annual cost of groups in a separate “holdout” set. The holdout set consisted of 648 employer groups (referred to here simply as “group”) with renewal dates varying from 05/01/2016 to 04/01/2017. There were a total of 349,715 members still actively enrolled as of their respective groups renewal date. Using the holdout set, we evaluated our model by its performance of predicted cost incurred during the 12-month period starting from their group’s renewal date (the “projection period”).

In the holdout set, Delphi censored the data recorded four months prior to the renewal date for each group (the “blackout period”); group underwriting is usually done several months before the new contract year in order to create the renewal quotes presented to employer groups. The 12 months before the blackout period comprise the “experience



Figure 1: Holdout set “dynamic time-slicing”. This picture shows the way the holdout set groups are time-sliced according to their renewal dates, and thus have different lengths of historical data (shown as “Features”) and different projection periods (shown as “Targets”) with respect to each other. Lumiata was blinded from the information in blackout and projection periods before submitting predictions to Delphi.

period”, ending on the “slice date”. For example, if a group’s renewal date was 05/01/2017, the experience period was the entire 2016 calendar year and 05/01/2017 to 04/30/2018 was the cost projection period (Figure 1).

For this pilot study, we sliced the training data to select groups and members in a similar, but simpler way. Instead of dynamically slicing the groups, we imposed a fixed renewal date of 01/01/2017 and used only the data recorded until 08/31/2016 for feature extraction. Therefore, the members and groups for which we predicted cost were those eligible as of the 08/31/2016 slice date. After filtering the training data for eligibility, ~ 7.4 million patients were enrolled as of the slice date. In order to fine-tune and evaluate models, we split these eligible groups and members into “train”, “test”, and “evaluate” sets. The data were split using a 70:20:10 ratio (Figure 2).

### Data Sources

We built our models using medical, capitation and pharmacy claims, and lab and eligibility tables for Delphi’s patients and groups. The medical claims tables contained cost information, International Classification of Disease (ICD-9 and ICD-10) diagnostic codes, and Current Procedural Technology (CPT) procedure codes at the claim-level. The capitation tables contained only cost information. All tables reported each claim’s care setting: inpatient, outpatient, ancillary, emergency, primary care, or specialty care. The pharmacy claims tables contained National Drug Code (NDC) medication codes and cost for each drug prescribed to a patient and the written and fill dates for the drug prescriptions. The cost associated to claims were given as “allowed” amounts (the amount paid by the insurer plus the members cost share). Lab tables contained Logical Observation Identifiers Names and Codes (LOINC) lab test codes. Eligibility tables captured each patient’s health plans, enrollment time periods, and plan benefits.

### Models

**Actuarial Models** To compare model performance between Delphi’s actuarial models and Lumiata’s ML models, we followed the best practices of actuarial science. Actuarial models estimate group cost on a per member per month

<sup>2</sup>Independence Blue Cross: Large Group Underwriting Guidelines



Figure 2: “Splitting” and “time-slicing” for the training data. This figure shows how we split and sliced the data of groups and members enrolled on 08/31/2016 to train and validate our models. Groups were split using a 70:20:10 train:test:evaluate ratio.

(pmpm) basis. The normalization unit is called a “member month”, which is defined as one month of enrollment for one member. We used member months to normalize cost because predicted pmpm cost for a group translates to the monthly premium charged to each member in this group. Actuarial models estimate group-level cost, treating each member month of medical history as independent across members and within a member.

The following example shows the utility and limitation of this perspective. Consider three hypothetical groups X, Y, and Z, each having 100 members and costing \$1 million during the 2017 calendar year:

- Group X: each member is enrolled for 10 months - the pmpm cost equals \$1 million/(10 months  $\times$  100 members) = \$1,000 pmpm
- Group Y: each member is enrolled for five months - the pmpm cost equals \$1 million/(5 months  $\times$  100 members) = \$2,000 pmpm
- Group Z: each member is enrolled for five months and one member costs \$900,000, while the rest of the members cost \$100,000 total - the pmpm cost equals \$1 million/(5 months  $\times$  100 members) = \$2,000 pmpm

As a result of the pmpm formulation, two groups with the same cost but different member months will have a different pmpm cost (Group X vs Y or Z). If the group cost is highly concentrated on one individual, versus being evenly distributed amongst the members in the group, it does not necessarily reflect in the pmpm cost, (Group Y vs Z). In contrast, a member-level cost prediction model views Group Y and Z differently, therefore better modeling cost at the group level.

Actuarial predictive models used for quoting renewal business rely on dozens of rating factors (input variables) to

build a predicted rate for a given group. The factors rely on pre-computed demographic, medical trend, pharmacy trend and other actuarial coefficients derived from patients across a large (usually exogenous) population using regression-based methods. These rating factor “priors” are then used to assemble a predicted trend value for a particular group (called the “manual rate”,  $MR$ ). Historical claims for the group are used to create the “experience rate” ( $ER$ ), which is then blended with the manual rate to create the final prediction. The proportion of blending between the two quantities is called “credibility” ( $c$ , where  $0 \leq c \leq 1$ ; ??). Thus, the predicted total cost of a group can be expressed as:

$$\text{predicted cost} = c \cdot ER + (1 - c) \cdot MR \quad (1)$$

This formula is a linear Bayesian hierarchical model, and is the optimal linear least-squares solution for estimating the annual pmpm cost of an employer group, called the Buhlmann-Straub method (?). This approach is the industry standard, underlying most models in production at insurance companies; applications include pricing, plan design, and reserve setting (??). As a group becomes more credible, the actuarial model can rely more on the medical claims history of the group as an indicator of future expenses. In the absence of credibility, the safer bet is to rely on population-level cost estimations using only age and sex (the manual rate). Actuaries use a group’s member months to parameterize credibility (??). A larger group enrolled for a shorter period of time (e.g. 1000 members  $\times$  6 months = 6000 member months) can have the same number of member months as a smaller group enrolled for a longer period of time (e.g. 200 members  $\times$  30 months = 6000 member months), making them equally credible.

The two equations below are examples of the type of experience rating and manual rating models used by Delphi:

$$ER = (TC - TSC)(1 + AT)^{\frac{m}{12}} \cdot x_m x_b x_d + n_s x_p + BC_p(1 + AT_L)^{\frac{m}{12}} \cdot x_{ph} x_{gp} x_{dp} x_{ip} \cdot mm$$

$$MR = [BC_{med}(1 + AT_{med})^{\frac{m}{12}} \cdot x_{gm} x_{dm} x_{im} x_{udm} + BC_{cap}(1 + AT_{med})^{\frac{m}{12}} + BC_{ph}(1 + AT_{ph})^{\frac{m}{12}} \cdot x_{gph} x_{dph} x_{iph} x_{udph}] \cdot mm$$

The definitions of independent variables in these two equations are shown in Table 1. For example,  $x_d$  is the average of the group’s member’s age and sex factors, which are legal and widely used for large group underwriting in the United States. The experience rate is linear in terms of the total claims ( $TC$ ), and is combined with the manual rate and the group’s member months (Eqn. 1).

To improve accuracy, we modeled cost at both the individual and group levels. We used a sequence of two models: (1) Individual-level model predicting per month cost for a given member and (2) Group-level model predicting per member per month (pmpm) cost for a given group. Our approach contrasts with traditional actuarial methods which are heavily focused on group-level cost and lack individual-level information within each group.

| Variable                              | Definition  | Comment   |
|---------------------------------------|---|---|
| $AT$                                  | annual trend  | medical and pharmacy trends combined adjusted for the effects of pooling-point leveraging                               |
| $AT_L$                                | leveraged annual trend  |   |
| $AT_{med}$<br>$AT_{ph}$<br>$BC_{cap}$ | medical annual trend<br>pharmacy annual trend<br>base capitation claims | pmpm cost of medical claims during the experience period<br>pmpm cost of capitation claims during the experience period |
| $BC_{med}$                            | base medical claims   |   |
| $BC_p$                                | base pooled claims  | pmpm cost of pooled claims during the experience period<br>pmpm cost of pharmacy claims during the experience period    |
| $BC_{ph}$                             | base pharmacy claims  |   |
| $TC$                                  | total claims  | total cost during the experience period   |
| $TSC$                                 | total shock claims  | total cost of claims over the pooling level during the experience period  |
| $m$                                   | midpoint months   | number of months between the midpoints of experience and projection periods   |
| $mm$                                  | member months   | member months during the experience period  |
| $n_s$                                 | number of shock claims  | number of claims over the pooling level   |
| $x_b$                                 | benefit   | based on the experience period  |
| $x_d$                                 | demographic   |   |
| $x_{dm}$                              | demographic - medical   | based on census   |
| $x_{dp}$                              | demographic - pooling   | based on census   |
| $x_{dph}$                             | demographic - pharmacy  | based on census   |
| $x_{gm}$                              | geographic area - medical   | based on census   |
| $x_{gp}$                              | geographic area - pooling   | based on census   |
| $x_{gph}$                             | geographic area - pharmacy  | based on census   |
| $x_{im}$                              | industry - medical  | the cost threshold defining a shock claim (typically \$100,000)   |
| $x_{ip}$                              | industry - pooling  |   |
| $x_{iph}$                             | industry - pharmacy   |   |
| $x_m$                                 | maturation  |   |
| $x_p$                                 | pooling level   |   |
| $x_{ph}$                              | pharmacy load   | $1.2e^{-0.8S}$ , where $S$ is medical cost share determined from a table using medical cost share                       |
| $x_{udm}$                             | utilization dampening - medical   |   |
| $x_{udph}$                            | utilization dampening - pharmacy  |   |

Table 1: Definitions of all the independent variables in the manual rate and experience rate equations. The most recent census was used for lookups. All  $x$  variables are factors.

**Feature Engineering** We reshaped the claims and eligibility tables into longitudinal patient records per our proprietary data format, the “Lumiata Data Model” (LDM). From the LDM, we created member-level features (input variables) using information before the blackout period, based on techniques from the literature (e.g. ? ?; ? ?). Our demographic features were age and sex; all other features were time-dependent. Our time windows to compute a variety of features (e.g. diagnosis, medication, procedure, lab, and revenue codes, and cost and coverage) were: “last three months”, “last six months”, “last one year”, and “anytime” prior to the blackout date.

In addition to ICD-9, ICD-10, CPT, NDC, and LOINC codes, we transformed the codes into their grouped counterparts based on organ-type (SNOMED), condition categories (HCUP, CMS-HCC and HHS-HCC), drug molecule (RxNorm and ATC), and our proprietary clinical grouper (Lumiata disease code). We derived features from these codes by calculating the log count of every unique code in each time window, and the summary statistics of each system (total count, unique code count, minimum count, maximum count, mean count, etc) in the “anytime” window. The presence of revenue codes (binary) was computed in the “anytime” window.

We computed features using their observed lab interpretations or values from the LOINC codes. These include (1) log counts of interpretations, i.e. “high”, “low”, “abnormal”, and “normal”, (2) whether the value was increasing, decreasing or flat across time points for the same test (one-hot encoded), and (3) if the interpretation was fluctuating across time (binary). For (2), the calculation was based on a t-test’s  $p$ -value of a simple linear regression’s slope.

Cost features are the most powerful features to predict future cost. In addition to the cumulative allowed cost in different time windows, we computed cost attributed to different care settings. The length of coverage was computed for all time windows. In total, we constructed more than 5 million possible features per patient. The resulting feature matrix was very sparse; most of the columns had no values at all. We reduced the dimension of the feature space before training a model using feature selection techniques discussed below.

**Individual-Level Model** We regressed our first model on the allowed amount per month during the projection period for each member. We trained a gradient boosting tree that optimized the mean squared error (MSE) using the LightGBM package (?) in the Python programming language (?). We chose LightGBM over linear regression and XGBoost due to its better performance. To speed up training time and reduce over-fitting, we tested a variety of feature prevalence thresholds to reduce the features set to  $\leq 100,000$ . We defined the prevalence of a feature as the fraction of non-zero values for this feature across patients in the training set. The model implicitly selected useful features to split on; the typical number of features was  $\sim 4000$  (see Table 2).

We ran the model recursively, using a different threshold (thus a different number of features) in each iteration. Evaluating on the test set reached a performance plateau at

| feature type                     | value type | count |
|----------------------------------|------------|-------|
| age                              | numeric    | 1     |
| gender                           | binary     | 1     |
| cost                             | numeric    | 28    |
| coverage                         | numeric    | 4     |
| revenue code                     | binary     | 70    |
| ATC log count                    | numeric    | 31    |
| ATC summary stats                | numeric    | 7     |
| CMS-HCC log count                | numeric    | 11    |
| CPT log count                    | numeric    | 771   |
| CPT summary stats                | numeric    | 6     |
| HCUP log count                   | numeric    | 337   |
| HCUP summary stats               | numeric    | 7     |
| HHS-HCC log count                | numeric    | 6     |
| ICD10 log count                  | numeric    | 40    |
| ICD10 summary stats              | numeric    | 7     |
| ICD9 log count                   | numeric    | 12    |
| ICD9 summary stats               | numeric    | 6     |
| LOINC log count                  | numeric    | 115   |
| LOINC summary stats              | numeric    | 7     |
| LOINC interpretation log count   | numeric    | 117   |
| LOINC interpretation fluctuation | binary     | 4     |
| LOINC value trend                | binary     | 88    |
| NDC summary stats                | numeric    | 4     |
| RxNorm log count                 | numeric    | 230   |
| RxNorm summary stats             | numeric    | 6     |
| SNOMED log count                 | numeric    | 2082  |
| SNOMED summary stats             | numeric    | 7     |
| Lumiata disease code log count   | numeric    | 91    |

Table 2: A summary of features used by the individual-level model. Note that only the “anytime” window is used for revenue code and all the summary statistics. For LOINC value trend, each code corresponds to three one-hot encoded binary features, i.e., “increasing”, “decreasing”, “flat”, given a time window.

a threshold of  $\sim 0.001$ , which we used for all subsequent models. We used the test set data for hyper-parameter tuning and early-stopping. Using the trained model, we made cost predictions for all the individuals in the train, test, evaluate, and holdout sets.

Members enrolled as of 08/31/2016 were included in our train, test, and evaluate sets. Members who dropped out during the blackout period were filtered out before training; we did not train on a member’s data if this member was not enrolled in a group on 01/01/2017. We only omitted these members during training, not during inference.

**Aggregation** Our individual model predicts the per month cost in the projection period for the members who were enrolled in groups at the end of experience period. For our train, test, and evaluate sets, this date was 08/31/2016. For the holdout set, the slice date depended on the groups renewal date. We then aggregated these predictions based on the members active on this date to obtain the mean of member-level cost predictions for each group. This quantity has a unit of pmpm and became the input of our group-level

model described below.

**Group-Level Adjustment** The aggregated mean prediction of a group can be thought of as predicted pmpm cost. However, this quantity only considers the members enrolled at the end of the experience period and assumes enrollment remains constant throughout the projection period. In reality, a group can grow or shrink during the blackout and projection periods, affecting the true pmpm. We regressed a second model on the true group-level pmpm cost to adjust the aggregated predictions. We experimented with different group-level features in this model and use the following (per group): (i) mean cost of member-level predictions, (ii) mean member age, (iii) total number of member months for the group during the experience period, (iv) “growth” feature, defined as the change in the number of members during the experience period divided by the total number of member months, (v) average length of member coverage, (vi) fraction of experience period costs that were incurred during the final four months, before the blackout period (vii) fraction of high-cost members, defined as someone whose cost falls within the top 10% of all members in the training set. We then trained a LightGBM model that optimized pmpm Mean Absolute Error (MAE) and used the test set to perform hyper-parameter tuning and early-stopping. The mean of member-level cost predictions for each group highly correlates with the overall target pmpm cost for the group. With the additional features, the group-level model improved pmpm MAE by  $\sim 10\%$  compared to the individual-level model alone.

Note that a regression model optimizing MSE estimates the true mean of  $y$  at a given  $x$ , while optimizing MAE is equivalent to estimating the true median. Thus, in order to make the aggregated individual-level predictions unbiased estimators for the group pmpm costs, we must optimize MSE instead of MAE at the individual level, even though we are pursuing MAE at the group level.

Similar to the individual model, we trained the group model on groups still active during the projection period because only the groups that remained active will be included when evaluating the results. This operation was only done in model training and was not performed when doing inference using the trained model.

## Model Evaluation

**ML Metrics Evaluation** We adapted the standard metrics R-squared ( $R^2$ ), MAE, and Gini index (?) to predicted pmpm cost; evaluation metrics are measured by comparing “true pmpm cost” and “predicted pmpm cost”. Lumiata received claims data at the allowed amount level to preserve Delphi’s pharmaceutical and provider reimbursement rates, whereas Delphi’s production models predict the “paid amount” for employer groups. As a result, Delphi and Lumiata modeled two different types of cost, which are not directly comparable because the allowed amount is usually greater than the paid amount (though they are highly correlated; see ? ?). For business purposes, to make the predictions comparable, both teams agreed on two solutions. First, we computed “normalized” pmpm MAE, which equals



pmpm MAE divided by the “global pmpm cost”. The global pmpm cost is the total (allowed or paid) amount divided by the total number of member months across all groups. Second, we computed the trend versions of our predictions. The “allowed trend” was defined as the pmpm allowed amount in the projection period divided by the pmpm allowed amount in the experience period; the “paid trend” was defined similarly. Trend calculations are ubiquitous in actuarial science (?). ML and actuarial metrics used in this context provide strong directional support for one model over another, as we will see.

**Lift Plot and Concession Opportunities** While ML performance metrics like MAE and  $R^2$  are ubiquitous in the tech industry, they often lack direct connection to concrete business-level Key Performance Indicators (KPIs) in other fields (?). To address this gap, we computed a KPI called a “lift plot”, illustrating the implications of Lumiata’s model in a production context. The lift plot is made using the following steps: (1) compute Lumiata’s predicted allowed trend divided by Delphi’s predicted paid trend (the “trend ratio”), (2) rank groups according to the trend ratio, (3) segment the ranked groups into deciles, (4) compute the actual-to-expected (A/E) ratios within each decile, where “A” is the true paid amount of a given decile, and “E” is Delphi’s predicted paid amount of the same decile, (5) normalize the A/E of all the ten deciles by the global A/E of the holdout set, and (6) plot the “oracle” model (a model that perfectly predicts pmpm allowed trend). The predicted paid amount of a given decile is the sum of the predicted pmpm paid amount times the projected member months for all groups in this decile. The projected member months of the projection period is 12 months times the number of members at the end of the experience period.

A group is considered to be a “concession opportunity at the 5% level” if the group’s rates can be reduced by 5% while maintaining profitability (A/E remains  $< 1$ ). In other words, the true paid trend ratio (true pmpm paid trend divided by Delphi’s predicted pmpm paid trend) is  $< 0.95$ . We used precision and recall to evaluate concession opportunity identification performance (?). In order to identify concession opportunities at the 5% level, we applied a decision rule of  $< 1$  using a predicted “trend ratio” (Lumiata’s predicted pmpm allowed trend divided by Delphi’s predicted pmpm paid trend). That is, if Lumiata’s predicted pmpm allowed trend was lower than Delphi’s predicted pmpm paid trend for a group, we asserted that group was a concession opportunity at the 5% level. We used a decision rule of  $< 1$  instead of  $< 0.95$  in order to achieve a higher recall.

## Results

**Model Performance** Lumiata’s model was 20% better in normalized pmpm MAE, 26% better in pmpm  $R^2$  than Delphi’s model, and 2% lower in Gini index than Delphi’s model (Table 3). Lumiata correctly identified 84% of the groups in the holdout set that had concession opportunities at the 5% level.

Lumiata’s predicted pmpm allowed trend had a 65% precision and 84% recall in identifying concession opportuni-

| model   | MAE          | $R^2$        | Gini index   |
|---------|--------------|--------------|--------------|
| Delphi  | 0.239        | 0.265        | <b>0.600</b> |
| Lumiata | <b>0.192</b> | <b>0.334</b> | 0.588        |

Table 3: Performance metric results for Lumiata vs Delphi on a normalized pmpm basis. Bold indicates higher performance for that metric.



Figure 3: Flow-chart of the “stop-light” process implementing Lumiata’s model output. “UW” stands for underwriting.

ties at the 5% level, and a 56% precision and 85% recall in identifying concession opportunities at the 10% level. For comparison, an “oracle” model had an 84% precision and 96% recall at the 5% level, and a 73% precision and 97% recall at the 10% level. This similarity in precision and recall indicates Lumiata’s model was near optimal for predicting concession opportunities at the 5% level.

**Practical Application** Operationalizing Lumiata’s model relied on the “stop-light principle” to make the model output interpretable to non-data scientists. The decision process derived from Lumiata’s model is: (1) If Lumiata’s predicted trend is less than Delphi’s predicted trend, then the model suggests an underwriter give a concession of at least 5% on that group’s renewal quote (Green). (2) If Lumiata’s predicted trend is equal to or greater than Delphi’s predicted trend, no action is taken (Yellow or Red, see Figure 3).

The lift plot in Figure 4 shows the result of using Lumiata’s predicted pmpm allowed trend in this decision process. The average A/E per decile of Delphi’s model varies with the decile of Lumiata trend ratio (Lumiata’s predicted allowed trend divided by Delphi’s predicted paid trend) and true allowed trend ratio (true allowed trend divided by Delphi’s predicted paid trend), respectively. The A/E of the bottom five Lumiata trend ratio deciles are below 0.95, meaning Lumiata’s trend ratio can select groups for a rate drop of  $\geq 5\%$  with good accuracy. The median trend ratio was 0.89 which yielded higher precision than a decision rule of 1.0.

Had Delphi implemented our model for this pricing period, their underwriters could have dropped the renewal quote 5% or more for approximately half of the groups while retaining profitability. An “oracle” model (blue line in Figure 4) identified the same number of decile concession opportunities at the 5% level as our model (five deciles). In this application, our model acts as a decision-support tool to give



Figure 4: Lift plot for the holdout set using Lumiata (purple) and true allowed (blue) trend ratios. The trend ratio of a model is defined as its predicted trend divided by Delphi’s predicted trend. The solid black line indicates  $A/E = 1$ . The dotted black lines indicate where the  $A/E = 1.05$  or  $0.95$ . Concession opportunity groups are in deciles below the 0.95 dotted line.

efficient guidance to underwriters.

## Deployment

The goals of Lumiata’s pilot study with Delphi were to prove Lumiata’s (1) ML could improve over an established industry methodology, and (2) tech stack could deliver monthly predictions for Delphi’s renewal business groups. We designed a Kaggle-style competition between the two companies, with two holdout sets - a preliminary holdout set (658 groups) and a final holdout set (which we call “holdout set” throughout this paper - 648 groups). We had one chance to compare our predictions to Delphi’s on the final holdout set; the success or failure of the pilot study was predicated on whose model had the best group-level cost predictions.

**Quality Assurance** Due to several delicate calculations needed to assemble the predicted/true allowed trend predictions, we computed non-prediction fields to rule out non-data science confounding factors before running the results analysis. These fields included (per group): (i) number of members enrolled at the end of experience period, (ii) number of member months in the experience period, (iii) true allowed amount in the experience period, and (iv) predicted allowed amount in the projection period. After receiving the censored information, we computed the: (i) number of members at the beginning of projection period, (ii) number of member months in the projection period, and (iii) true allowed amount in the projection period. Due to a strong data model built off of FHIR and solid compute infrastructure, we were able to iterate and fix bugs quickly, until our calculations of the experience period non-prediction columns

matched Delphi’s within 5%. Completion of this analysis allowed quick and self-evident comparison of Lumiata’s and Delphi’s models performance metrics. Additionally, the non-prediction fields found their way into roll-out plan below.

**Data Challenges** For model comparison, we had to address discrepancies between Lumiata’s and Delphi’s claims data sets. Two major differences in the claims data were: (1) Delphi’s models used paid amount and Lumiata used allowed amount, and (2) Delphi used the “paid date”, but requested Lumiata use the “encounter date” for allowed amounts (Delphi felt “encounter date” was more appropriate for patient-level cost predictions). To solve these issues, we computed the “allowed trend” compared to the “paid trend” (Figure 4). Calculating the “allowed trend” required we compute three other quantities per group (“number of members at the end of experience period”, “member months in the experience period”, and “allowed amount in the experience period”), and the predicted allowed amount for the group. Accurately computing these quantities was more difficult than expected due to the consequences of dual paid/allowed conventions and time-dependent patient enrollment:

(i) Calculating the “allowed amount in the experience period” differs depending on whether “paid date” or “encounter date” is used. Using the paid vs encounter date segmented claims differently into a groups experience, black-out, and projection periods. For example, some claims were denied or paid claims could be reversed. Hence, claims data filtered on encounter date in the holdout set experience period were sometimes absent from the unblinded holdout set’s *experience* period.

(ii) The number of members enrolled in a group at the end of the experience period could change during the black-out and projection periods. For example, members can shift between different groups due to job or spousal health plan changes. Also, enrollment was updated on the 15th of the month, but we calculated “number of members at end of experience period” and “member months in the experience period” for the first of each month, making our enrollment data two weeks out of date.

**Roll-out Strategy** Delphi requested Lumiata provide monthly group-level cost predictions and concession opportunities for groups up for renewal within the next four months (Delphi renews groups throughout the year). Delphi asked for a seven day cost-prediction turn around time and fixed model feature set each quarter. In response, we built a platform that can create LDMs for 1 TB of claims data in under an hour and produce highly optimized group-level cost predictions in under four hours.

To streamline deployment, we proposed the following plan to Delphi: (1a) Each month, Delphi sends a moving three year snapshot of their claims and eligibility tables. (1b) Delphi sends their non-prediction field calculations and their paid trend group cost predictions for groups up for renewal. (2) Lumiata updates the existing patients LDMs and adds new patients from the eligibility file. (3) Lumiata updates the group cost prediction model, training using the prior 12 months as the projection period, four months prior as the blackout, and prior 14 months as the experience period. (4)

<sup>2</sup><https://www.kaggle.com/>

<sup>3</sup><https://wiki.hl7.org/FHIR>

Lumiata creates feature vectors including the new projection/blackout period claims information. (5) Lumiata applies the updated member- and group-level models to the claims data to produce group-level allowed trend predictions. (6) Using Delphi’s paid trend predictions from 1b and the “stop-light” principle (Figure 3), Lumiata recommends whether to drop the rate by 5% for each group. (7) Lumiata sends its non-prediction attributes and the allowed trend predictions to Delphi with one recommendation per group up for renewal. All non-prediction fields must agree within  $\leq 5\%$  between Delphi and Lumiata’s calculations. (8) Delphi verifies receipt of the data and the results are consumed by their actuarial and underwriting teams. All files are transferred to and from Lumiata’s platform (hosted on Google Cloud Platform) using sFTP.

**Transparent Rate Setting** Actuarial models have an essential property: they are “explainable” because a prediction can be decomposed into discrete multiplicative factors with an inherent interpretation. For example, “geographic area factor” = 0.9 means people in a particular zip code cost 10% less than the mean, so the base rate is adjusted (multiplied) by 0.9 for members from this zip code. This degree of explainability is crucial because actuaries need to file rates annually for individual/small group markets with the state insurance commissioner to ensure the factors used to produce the rate are compliant with legal guidelines. Furthermore, an underwriter needs to be able to explain how she arrived at a particular rate to a customer. A critique of ML models is that they lack explainability in terms of what input variables may have contributed to a particular prediction (?). However, explainable ML in healthcare is must-have, touching upon fundamental issues of bias, transparency, and reasonableness of ML model predictions.

Shapely values, a game theoretic algorithm, enable common ML algorithms to output feature weights specific to a prediction (?). We applied the Python package SHAP to our LightGBM models to yield member-level explanations. According to the algorithm, a model  $f$  and a member feature vector  $x$  admit an additive decomposition in terms of the mean value of the model  $\mathbf{E}[f(x)]$  and SHAP values  $\phi_i$  (interpreted as dollar value pmpm amounts) for the  $i$ th feature. For our member- and group-level models,  $f_1$  and  $f_2$  are:

$$f_1(x) = \mathbf{E}[f_1(x)] + \sum_{i=1}^{3996} \phi_{i,1} \quad (2)$$

$$f_2(x) = \mathbf{E}[f_2(x)] + \sum_{i=1}^7 \phi_{i,2} \quad (3)$$

Given our member and group model sequence, our group model (3) admits an additive decomposition in terms of member-level  $\phi_{i,1}$ . The mechanics in (3) are similar to an actuarial formula, in that SHAP values for a group’s members’ features additively adjust the mean pmpm cost of the model over all members and groups, while the actuarial factors for a particular group multiplicatively adjust the mean pmpm cost of their entire patient population. We often found that  $\leq 500$  features’ SHAP values account for 95% of a cost

prediction, for each group. However, the specific features involved varied by group.

The transparency afforded by the group-level SHAP values provides the opportunity to explain a rate adjustment to a customer in dollar pmpm amounts using the specific risk drivers for that group and modify a rate given by the ML model by the expected change in cost for specific drugs and services. For instance, if the price of Glipizide, a drug to treat type-2 diabetes, will drop for an insurer next year by 20%, imagine a world where the insurer can multiply the SHAP values corresponding to Glipizide-related pharmacy costs by 0.80 for all the members on Glipizide, thus lowering the projected rate. These mechanics would be similar to current actuarial methods, making them easy to implement. Furthermore, greater model transparency could increase patient adherence to prescribed medications. As drug prices rise and more patients purchase high-deductible plans, patients have to pay higher out of pocket costs for drug treatment, and patient medication adherence declines (?). Insurers could use the SHAP values from the patient-level model to identify drugs driving up projected cost for that group, and suggest the prescribing doctor offer a lower cost alternative drug with similar efficacy. This provides a win-win opportunity, lowering drug cost for the payer, and improving patient adherence to the cheaper drug through increased affordability.

## Discussion

Here, we demonstrate that: (1) ML approaches can significantly improve the accuracy and efficiency of group health insurance underwriting and (2) ML models can offer comparable interpretability to traditional actuarial methods. Our contributions provide clear direction for how to improve the efficiency and predictive performance of underwriting for employer-based insurance and how to lower the cost for members in groups of any size. Our ML-based approach improved MAE over actuarial models across the book of business: >500-member groups showed an improved performance.

Our model shows the most improvement over actuarial models in situations where the group size is  $\leq 500$  and/or the group claims experience is relatively short ( $< 8$  months). In these situations, the groups are not yet credible, so actuarial models perform sub-optimally. We believe the success of our model was due to modeling costs: (1) at the individual level; by contrast, actuarial methods aggregate medical history across group members, (2) with models that perform well with skewed distributions; the cost of care in an insurance population is often gamma distributed, making the ML method of gradient boosting trees, like LightGBM, highly effective (?), (3) using a model-agnostic approach to select features relevant for predictions, and (4) by combining individual- and group-level models to produce the final predictions.

We treated all members in the training set as if their group’s renewal was 01/01/2017, with a four month blackout period starting on 08/31/2016. However, patients often try to “fit in” healthcare services before their next renewal

<sup>4</sup><https://github.com/slundberg/shap>



period because they are likely near or above their plans deductible and hence will have these services fully covered by the insurance company. When patient records were “sliced” uniformly on 08/31/2016, there was a chance that this useful information would be lost and negatively impact the model performance. We found that dynamically “slicing” patient records according to their groups renewal date (instead of all patients on 08/31/2016), and training the models with this feature setup, the overall results on the holdout set remained roughly the same in terms of MAE and  $R^2$ . Model performance likely didn’t improve because claims data are inherently fuzzy with dates; claims can take variable amounts of time to get paid, and hence a model that tries to learn precise date information will add some but not meaningful value. To account for fuzzy date information, we aggregated features over 3 months, 6 months, etc (see Methods).

As ML models are increasingly compared to more traditional statistical techniques, the most appropriate study design and model evaluation metrics should be examined. For example, the holdout set data was “out of sample” (i.e. using patients/groups unseen before by the model) but not “out of time” (i.e. projecting costs for time periods subsequent to 2017). Furthermore, claims data are not time-stationary (e.g. new drugs and treatments will be developed), so the expected model performance may not be perfectly realized in practice, but the relative difference between the models should hold. Also, we obtained a slightly worse Gini index than Delphi, despite our much better  $R^2$ , MAE, and lift plot (Table 3 and Figure 4). This discrepancy can occur because the Gini index is a ranking-based metric, whereas a regression model minimizes the prediction errors. One difficulty is the Gini index is not a differentiable quantity. Future work should develop algorithms to address this problem. Data quality was crucial to our success. Alignment on non-prediction fields between Lumiata and Delphi ruled out errors in the data, pipeline, or output, improving communication across teams and increasing efficiency. These calculations must be automated for developing and productionalizing a medical underwriting ML application, due to the large size of data sets and rapid turnaround of results.

Due to its highly applied nature, some operational realities limit our study’s evaluation. A challenge for validating our predictions is the long feedback cycle (20 months). Also, not all of Lumiata’s concession recommendations could be granted due to a variety of quantitative and judgement factors under the underwriter and insurer’s discretion.

Additionally, we could not determine if our individual-level model was racially biased, because we did not receive patient ethnicity data. Avoiding racial bias is important as previous studies have found evidence of racial bias in commercial cost prediction models used for clinical management (?). Historically, poorer minorities under-utilized healthcare services due to mistrust of the system and confusion about how to navigate it (?). However, because our response variable is not a clinical outcome but a financial one, we think this effect on pricing may be less significant. More work will be needed to better understand the effect of pricing insurance more affordably for minority patients, predicated on their less frequent utilization of the healthcare system.

In practice, ML approaches can help insurers be more competitive, avoiding adverse risk. It can result in the design of more “exotic” funding arrangements due to better predictive power of patient health, following the industry trend towards capitated payments (?). Unlike previous ML models in healthcare (?), our model output is interpretable by a non-technical user, simplifying operationalization (Figure 3). A user does not need to understand the inner workings of our algorithm to apply our output as a multiplicative adjustment factor to their existing actuarial models and can output the most important group-specific risk factors.

## Conclusions

Machine learning on insurance claims data provides a powerful tool to improve the efficiency and affordability of plans and care offered to patients enrolled in employer-sponsored health plans. With more accurate rate-setting, health insurance companies can design nuanced plan attributes, reducing the cost of care for their members. Our ML model achieved 20% improved accuracy in absolute predictive performance over traditional actuarial methods and was able to identify over 80% of new concession opportunities available to Delphi. This allows underwriters to better price and retain <500 employer group customers. This study can be used by payers to give underwriters improved pricing guidance, retaining business and giving a better and more affordable experience to members.

## Acknowledgements

We thank our counterparts at Delphi for collaboratively working with us to validate our model against a production-grade model with real data. Additionally, we thank the following people for their contributions: Kim Branson, Dilawar Syed, Laika Kayani, Wil Yu, Shahab Hassani, Alexandra Pettet, Derek Gordon, Ash Damle, Thomas Watson, Leon Barovic, Diana Rypkema and our investors at Blue Cross Blue Shield Venture Fund, and Khosla Ventures.

Quasi officia doloribus nisi debitis perspicatis a modi quibusdam expedita consequuntur mollitia, aliquid itaque ab nostrum autem tempora odit velit rerum, repudiandae architecto neque ratione sapiente ducimus eum quis voluptatibus porro autem? Incidunt quo earum minus laudantium quis quasi rem, quo incidunt nihil? Iure sunt repellendus incidunt sint explicabo eos quos recusandae odio, corrupti odio doloribus fuga officiis, ipsam culpa inventore consequatur quae id velit asperiores ex quidem quam eos, distinctio vitae porro eum doloribus animi reprehenderit veritatis nulla delectus recusandae, earum quibusdam odit ex ratione? Eveniet minus quas eos quo libero asperiores recusandae sapiente, ullam adipisci veniam exercitationem nam doloremque consectetur repellat unde quaerat odit laboriosam, distinctio autem laboriosam, veritatis aliquam ex odio cumque nemo nesciunt non aliquid id consequuntur, aspernatur facilis fugiat obcaecati vitae eius minima? Labore animi doloribus fuga minima vel aliquid quod, dicta voluptas ipsam temporibus asperiores nulla id provident excepturi quaerat odit eveniet, doloremque dolorum ipsam excepturi magnam sapiente rerum. Dolorem explicabo placeat similique error, saepe aperiam facilis magni aliquam assumenda porro officia asperiores, earum delectus voluptatum voluptates ab