

REGRESSION TREE CREDIBILITY MODEL

LIQUN DIAO AND CHENGGUO WENG

Department of Statistics and Actuarial Science, University of Waterloo

Advances in Predictive Analytics Conference, Waterloo, Ontario
Dec 1, 2017

Overview

Statistical Method + Actuarial Model
Regression Trees Credibility Model

Outline

1. CREDIBILITY MODEL
2. BENEFIT OF PARTITIONING THE DATA SPACE
3. REGRESSION TREE CREDIBILITY MODEL
4. SIMULATION STUDIES
5. AN APPLICATION TO US MEDICARE DATA
6. CONCLUDING REMARKS

Bühlmann-Straub Credibility Model

- CREDIBILITY THEORY has become the paradigm for insurance experience rating and widely used by actuaries.

Bühlmann model (1967, 1969), and the Bühlmann-Straub model (1970).

- Consider a portfolio of I risks, where each individual risk i has n_i years of claim experiences, $i = 1, 2, \dots, I$.
- Let $Y_{i,j}$ denote the claim ratio of individual risk i in year j , and $m_{i,j}$ be the associated volume measure, also known as weight variable.
- Collect all the claims experience of individual risk i into a vector \mathbf{Y}_i , i.e., $\mathbf{Y}_i = (Y_{i,1}, \dots, Y_{i,n_i})^T$.
- The profile of individual risk i is characterized by θ_i , which is the realization of a random element Θ_i (usually either a random variable or a random vector).

- Assume that the following conditions are satisfied:

H01. Conditionally given $\Theta_i = \theta_i$, $\{Y_{i,j} : j = 1, 2, \dots, n_i\}$ are independent

$$\mathbb{E}[Y_{i,j} | \Theta_i = \theta_i] = \mu(\theta_i) \text{ and } \text{Var}[Y_{i,j} | \Theta_i = \theta_i] = \frac{\sigma^2(\theta_i)}{m_{i,j}}$$

for some unknown but deterministic functions $\mu(\cdot)$ and $\sigma^2(\cdot)$;

H02. The pairs $(\Theta_1, \mathbf{Y}_1), \dots, (\Theta_I, \mathbf{Y}_I)$ are independent, and $\{\Theta_1, \dots, \Theta_I\}$ are independent and identically distributed.

- Define structural parameters $\sigma^2 = \mathbb{E}[\sigma^2(\Theta_i)]$ and $\tau^2 = \text{Var}[\mu(\Theta_i)]$ for risks within the collective $\mathcal{I} := \{1, 2, \dots, I\}$, and denote

$$m_i = \sum_{j=1}^{n_i} m_{i,j}, \quad m = \sum_{i=1}^I m_i, \quad \bar{Y}_i = \sum_{j=1}^{n_i} \frac{m_{i,j}}{m_i} Y_{i,j}, \quad \bar{Y} = \sum_{i=1}^I \frac{m_i}{m} \bar{Y}_i,$$

- Let $\mu = \mathbb{E}[Y_{i,j}]$ denote the collective net premium.
- We are interested an estimator $\hat{\mu}(\Theta_i)$ of $\mu(\Theta_i)$, is called the correct individual premium of the individual risk (the fair risk premium), such that it makes the following quadratic loss as small as possible

$$E [(\hat{\mu}(\Theta_i) - \mu(\Theta_i))^2].$$

Inhomogeneous Credibility Premium

- The (inhomogeneous) credibility premium for an individual risk is defined as the best premium predictor P_i among the class

$$\left\{ Q : Q = a_0 + \sum_{i=1}^I \sum_{j=1}^{n_i} a_{i,j} Y_{i,j}, \quad a_0, a_{i,j} \in \mathbb{R} \right\}$$

to minimize the quadratic loss

$$\mathbb{E} \left[(Q_i - \mu(\Theta_i))^2 \right]$$

and its formula is given by

$$P_i = \alpha_i \bar{Y}_i + (1 - \alpha_i) \mu, \quad (1)$$

where

$$\alpha_i = \frac{m_i}{m_i + \sigma^2 / \tau^2} \quad (2)$$

is called *credibility factor* for an individual risk.

- The corresponding minimum quadratic loss for $\mathbb{E} \left[(Q_i - \mu(\Theta_i))^2 \right]$ is given by

$$L_i = \frac{\sigma^2}{m_i + \sigma^2/\tau^2}. \quad (3)$$

- The minimum quadratic loss for $\mathbb{E} \left[(Q_i - Y_{i,n+1})^2 \right]$ is given by

$$L_i = \frac{\sigma^2}{m_i + \sigma^2/\tau^2} + \sigma^2. \quad (4)$$

Homogeneous Credibility Premium

- The homogeneous credibility premium for an individual risk i from the collective \mathcal{I} is defined as the best premium predictor P_i among the class

$$\left\{ Q_i : Q_i = \sum_{i=1}^I \sum_{j=1}^{n_i} a_{i,j} Y_{i,j}, \mathbb{E}[Q_i] = \mathbb{E}[\mu(\Theta_i)], a_{i,j} \in \mathbb{R} \right\}$$

to minimize the quadratic loss

$$\mathbb{E} \left[(Q_i - \mu(\Theta_i))^2 \right]$$

and its formula is given by

$$P_i = \alpha_i \bar{Y}_i + (1 - \alpha_i) \bar{Y}, \quad (5)$$

where

$$\alpha_i = \frac{m_i}{m_i + \sigma^2/\tau^2} \quad (6)$$

is called *credibility factor* for individual risk i .

- The corresponding minimum quadratic loss for $\mathbb{E} \left[(Q_i - \mu(\Theta_i))^2 \right]$ is given by

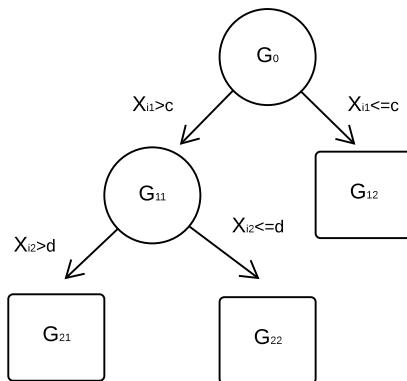
$$L_i = \tau^2(1 - \alpha_i) \left(1 + \frac{1 - \alpha_i}{\alpha_{\bullet}} \right), \quad (7)$$

where $\alpha_{\bullet} = \sum_{i=1}^I \alpha_i$.

- The minimum quadratic loss for $\mathbb{E} \left[(Q_i - Y_{i,n+1})^2 \right]$ is given by

$$L_i = \tau^2(1 - \alpha_i) \left(1 + \frac{1 - \alpha_i}{\alpha_{\bullet}} \right) + \sigma^2, \quad (8)$$

where $\alpha_{\bullet} = \sum_{i=1}^I \alpha_i$.



Better if we partition the collective?

The answer is **artificially YES** but **genuinely not necessary**

- For $n_i = n$ and $m_{i,j} = 1$ for and $i = 1, \dots, I$ and $j = 1, \dots, n$. Consider the loss function of homogenous premium and $\mu(\Theta_i)$.
 - Arbitrarily partition a given collective of individuals into two sub-collectives, and then apply credibility formula separately for each of the two resulting sub-collectives.
 - Let L_1 and L_2 be the total credibility losses we have for the two sub-collectives respectively, while L denotes the credibility loss with premium prediction apply with the whole collective without any partitioning.
 - We formally proved $L_1 + L_2 \leq L$.
- The above result relies on an **artificial assumption**:

We know the credibility factor α for each sub-collective".
- The premium prediction is given by

$$\hat{P}_i = \hat{\alpha}_i \bar{Y}_i + (1 - \hat{\alpha}_i) \bar{Y},$$

where $\hat{\alpha}_i$ is an estimator for α_i , $i = 1, \dots, I$.

Overview of Regression Trees

(Breiman, Friedman, Stone and Olshen, 1984)

- **CLASSIFICATION AND REGRESSION TREES (CART;** Breiman, Friedman, Stone and Olshen, 1984) and **RANDOM FORESTS (RF;** Breiman, 2001) are the most popular single-tree and ensemble recursive partitioning methods respectively.
- Any tree building process can be broadly described in three steps:
 - ① Choosing a criterion for making splitting decisions;
 - ② Generating a corresponding sequence of candidate trees;
 - ③ Selecting best candidate tree.
- It's common to allow each step to depend on a given loss function
- Prevailing software implementation of CART: `rpart` by R

Covariate-Dependent Model Setup

- Consider a portfolio of I risks numbered with $1, \dots, I$.
- Let $\mathbf{Y}_i = (Y_{i,1}, \dots, Y_{i,n_i})^T$ be the vector of claim ratios, $\mathbf{m}_i = (m_{i,1}, \dots, m_{i,n_i})$ be the corresponding weight vector, and $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,p})^T$ be the covariate vector associated with individual risk $i, i = 1, \dots, I$.
- The risk profile of each individual risk i is characterized by a scalar θ_i , which is a realization of a random element Θ_i .
- The following two conditions are further assumed:
 - H11. The triplets $(\Theta_1, \mathbf{Y}_1, \mathbf{X}_1), \dots, (\Theta_I, \mathbf{Y}_I, \mathbf{X}_I)$ are independent;
 - H12. Conditionally given $\Theta_i = \theta_i$ and $\mathbf{X}_i = \mathbf{x}_i$, the entries $Y_{i,j}, j = 1, \dots, n$, are independent with

$$\mathbb{E}[Y_{i,j} | \mathbf{X}_i = \mathbf{x}_i, \Theta_i = \theta_i] = \mu(\mathbf{x}_i, \theta_i) \text{ and}$$

$$\text{Var}[Y_{i,j} | \mathbf{X}_i = \mathbf{x}_i, \Theta_i = \theta_i] = \frac{\sigma^2(\mathbf{x}_i, \theta_i)}{m_{i,j}}$$

for some unknown but deterministic functions $\mu(\cdot, \cdot)$ and $\sigma^2(\cdot, \cdot)$.

- We approximate $\mu(\mathbf{X}_i, \Theta_i)$ and $\sigma^2(\mathbf{X}_i, \Theta_i)$:

$$\sum_{k=1}^K \mathbf{I}\{\mathbf{X}_i \in A_k\} \mu_{(k)}(\Theta_i), \quad \text{and} \quad \sum_{k=1}^K \mathbf{I}\{\mathbf{X}_i \in A_k\} \frac{\sigma_{(k)}^2(\Theta_i)}{m_{i,j}},$$

where

- $\{A_1, A_2, \dots, A_K\}$ is a partition of the covariate space
- $\mu_{(k)}(\Theta_i)$ and $\sigma_{(k)}^2(\Theta_i)$ respectively represent the net premium and the variance of an individual risk i from the k th sub-collective with a risk profile Θ_i , i.e.,

$$\mu_{(k)}(\theta_i) = \mathbb{E}(Y_{i,j} | \mathbf{X}_i \in A_k, \Theta_i = \theta_i) \quad \text{and}$$

$$\sigma_{(k)}^2(\theta_i) = \text{Var}(Y_{i,j} | \mathbf{X}_i \in A_k, \Theta_i = \theta_i).$$

- The condition of $\mathbf{X}_i \in A_k$ means that the individual risk i is classified into the k th sub-collective based on its covariate information.

Regression Tree Credibility Premium

- We target to find a “Good” partition $\{A_1, A_2, \dots, A_K\}$ and apply credibility formula for each sub-collective A_i separately.

- By “Good”, we should minimize the true prediction error

$$\mathbb{E} \left[\left(\mu(\Theta_i) - \hat{P}_i \right)^2 \right] \text{ as much as possible}$$

- **Credibility Regression Trees**

- Adopt one of the four credibility loss function with plugged-in estimates of structure parameters;
- Invent a heuristic **longitudinal cross-validation**.
- Credibility formula is applied for premium prediction for each terminal node separately, which is **Regression Tree Credibility Premium**.

Simulation Studies

- Consider a collective of $I = 300$ individual risks
- For each $i = 1, \dots, I$, independently simulate covariate vector

$$\mathbf{X}_i = (X_{i,1}, \dots, X_{i,p}) \stackrel{\text{i.i.d.}}{\sim} U\{1, \dots, 100\}^p$$

when $p = 10$ and $p = 50$.

- Balanced Claims Model: individual risks with $n = 5, 10$ and 20 years of claims experience, respectively. Unbalanced claims model is also explored.
- 1,000 independent samples

Simulation Scheme 1

Interaction effect

For each $i = 1, \dots, I$, independently simulate $\varepsilon_{i,j}$ from a given distribution function $F(\cdot)$, for $j = 1, \dots, n$, and define the n claims of individual risk i as

$$Y_{i,j} = e^{f(\mathbf{X}_i)} + \varepsilon_{i,j}, \quad j = 1, \dots, n, \quad (9)$$

where

$$f(\mathbf{X}_i) = 0.01 (X_{i,1} + 2X_{i,2} - X_{i,3} + 2\sqrt{X_{i,1}X_{i,3}} - \sqrt{X_{i,2}X_{i,4}}). \quad (10)$$

In our simulation, F takes one of the following three distributions:

- EXP(1.6487)
- LN (0, 1)
- PAR(3, 3.2974)

Simulation Scheme 2

Interactive effect+heterogeneous variance

For each $i = 1, \dots, I$, independently simulate $\varepsilon_{i,j}$ from a distribution function $F(\cdot; \mathbf{X}_i)$ which depends on the covariate \mathbf{X}_i of the individual risk i , for $j = 1, \dots, n$, and define the n claims of individual risk i as

$$Y_{i,j} = e^{f(\mathbf{X}_i)} + \varepsilon_{i,j}, \quad j = 1, \dots, n, \quad (11)$$

where $f(\mathbf{X}_i)$ is given by (10). We respectively consider three distinct distributions for $F(\cdot; \mathbf{X}_i)$:

$$(1) \text{EXP}(e^{\gamma(\mathbf{X}_i)/2}), \quad (2) \text{LN}(0, \gamma(\mathbf{X}_i)), \quad (3) \text{PAR}(3, 2e^{\gamma(\mathbf{X}_i)/2}), \quad (12)$$

where $\gamma(\mathbf{X}_i) = \frac{1}{102} |2X_{i,1} - X_{i,2} + \sqrt{X_{i,1}X_{i,2}}|$.

Simulation Scheme 3

Interactive effect + heterogeneous variance + **multiplication random effect**

For each $i = 1, \dots, I$,

- (1) independently simulate random effect variable Θ_i from the uniform distribution $U(0.9, 1.1)$
- (2) independently simulate $\varepsilon_{i,j}$ from a distribution function $F(\cdot; \mathbf{X}_i)$ which depends on the covariates \mathbf{X}_i associated with risk i , for $j = 1, \dots, n$, and
- (3) define the n claims of individual risk i as

$$Y_{i,j} = \Theta_i \left[e^{f(\mathbf{X}_i)} + \varepsilon_{i,j} \right], \quad j = 1, \dots, n, \quad (13)$$

where $f(\mathbf{X}_i)$ is given by (10).

We consider each of the three distributions described in Scheme 2 for the distribution $F(\cdot; \mathbf{X}_i)$.

Simulation Scheme 4

Interactive effect + heterogeneous variance + **complex random effect structure**

- For each $i = 1, \dots, I$, $\Theta_i = (\xi_{i,1}, \xi_{i,2})^T$,
 - (1) independently simulate random effect variables $\xi_{i,1}$ and $\xi_{i,2}$ from the uniform distribution $U(0.9, 1.1)$.
 - (2) independently simulate $\varepsilon_{i,j}$ from a distribution function $F(\cdot; \mathbf{X}_i, \xi_{i,2})$ for $j = 1, \dots, n$, and
 - (3) define the n claims of individual risk i as

$$Y_{i,j} = e^{\xi_{i,1} \cdot f(\mathbf{X}_i)} + \varepsilon_{i,j}, \quad j = 1, \dots, n, \quad (14)$$

where $f(\mathbf{X}_i)$ is defined in (10).

- We respectively consider the three distributions as defined in equation (12) in Scheme 2 **with $\gamma(\mathbf{X}_i)$ replaced by $\xi_{i,2} \cdot \gamma(\mathbf{X}_i)$** so that the distribution of $\varepsilon_{i,j}$ depends on the random effect variable $\xi_{i,2}$ in addition to covariate variable \mathbf{X}_i .

- **Data-driven covariate-dependent partitioning:** For each simulated sample, we grow and prune regression trees built using **four credibility loss functions** (R1-4) and L_2 (RL_2) loss function and the best tree is selected using longitudinal cross-validation.
- In addition, we consider **ad hoc covariate-dependent partitioning**, which is defined via the following notation:

$$\mathcal{R}(X_j) = \left\{ \{i \in \mathcal{I} : X_{i,j} \leq 50\}, \{i \in \mathcal{I} : X_{i,j} > 50\} \right\}, \quad j = 1, \dots, 5.$$

and

$$\mathcal{R}(X_{j_1}, X_{j_2}) = \left\{ \begin{aligned} &\{i \in \mathcal{I} : X_{i,j_1} \leq 50, X_{i,j_2} \leq 50\}, \{i \in \mathcal{I} : X_{i,j_1} \leq 50, X_{i,j_2} > 50\}, \\ &\{i \in \mathcal{I} : X_{i,j_1} > 50, X_{i,j_2} \leq 50\}, \{i \in \mathcal{I} : X_{i,j_1} > 50, X_{i,j_2} > 50\} \end{aligned} \right\},$$

for $j_1, j_2 = 1, \dots, 5$ with $j_1 \neq j_2$.

- We consider
 - $\mathcal{R}(X_2), \mathcal{R}(X_4),$
 - $\mathcal{R}(X_1, X_2, X_3), \mathcal{R}(X_1, X_2, X_4), \mathcal{R}(X_2, X_3, X_4), \mathcal{R}(X_1, X_3, X_4),$ and
 - $\mathcal{R}(X_1, X_2, X_3, X_4).$

- Prediction error for a given partitioning $\{\mathcal{I}_1, \dots, \mathcal{I}_K\}$:

$$\text{PE} = \frac{1}{I} \sum_{i=1}^I \sum_{k=1}^K \mathbf{I}\{\mathbf{X}_i \in \mathcal{I}_k\} \left(\pi_i^{(\text{H})^{(k)}} - \mu(\mathbf{X}_i, \Theta_i) \right)^2, \quad (15)$$

where

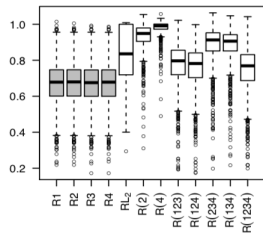
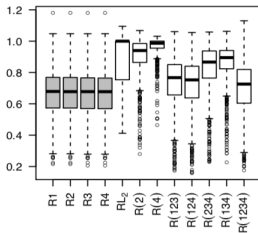
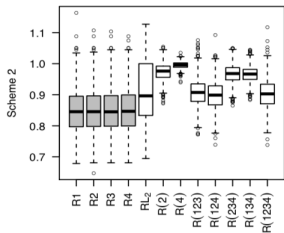
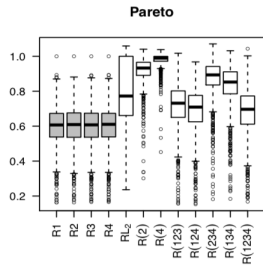
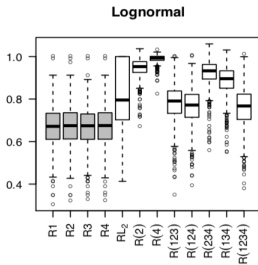
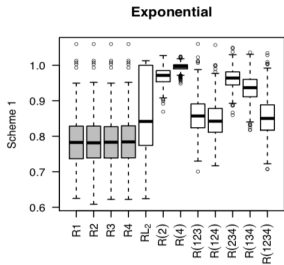
- $\pi_i^{(\text{H})^{(k)}}$ is the resulting premium prediction
- $\mu(\mathbf{X}_i, \Theta_i)$ is the true net premium
- The collective prediction error:

$$\text{PE}_0 = \frac{1}{I} \sum_{i=1}^I \left(P_i^{(\text{H})} - \mu(\mathbf{X}_i, \Theta_i) \right)^2,$$

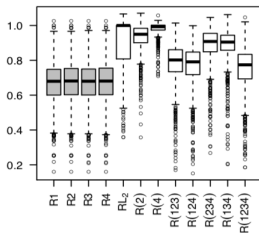
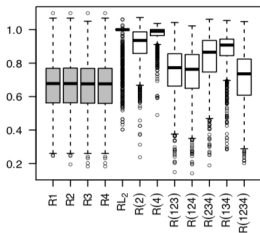
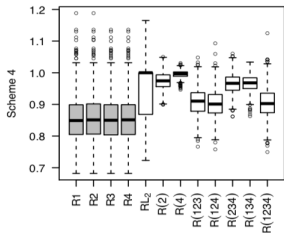
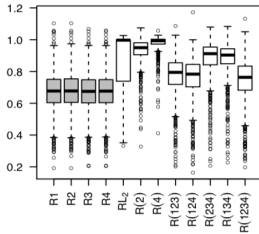
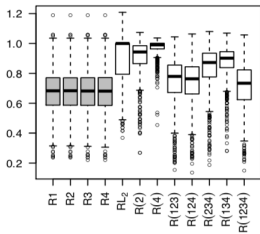
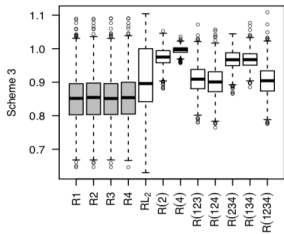
which does not use any covariate information and is anticipated to underperform compared to various kinds of covariate-dependent partitioning.

- The relative prediction error (RPE):

$$\mathbf{R} = \text{PE}/\text{PE}_0.$$



Concluding Remarks



Concluding Remarks

- We propose novel regression tree credibility (RTC) model, and bring machine learning techniques into the framework of credibility theory to enhance the prediction accuracy of credibility premium.
- In our proposed model, no ex ante analysis on the relationship between individual net premium and covariate variables is necessary, and the designed regression tree algorithm automatically selects influential covariate variables and informative cutting points to form a partition of data space, upon which a well-performed premium prediction rule can be consequently established.
- Our simulation studies and data analysis show that the proposed RTC model performs very well compared to no partitioning, ad-hoc partitioning and the L_2 loss based binary partition procedure.

- Although only the Classification and Regression Trees is introduced in this paper, it will be fruitful to pursue further research by considering other recursive partitioning methods, e.g.,
 - partDSA (partitioning deletion/substitution/addition algorithm) and
 - MARS (multivariate adaptive regression splines)
- It will be even more promising to consider the applications of ensemble algorithms, such as bagging, boosting, and random forests.
- It will be useful to develop an algorithm which can adopt time-dependent covaraites.
- It will be even more fruitful to consider their applications in various other insurance problems in addition to the premium rating, since many practical insurance problems amount to quantifying the relationship between insureds' claims and their demographic information.
- It is the authors' hope that the present paper will stimulate more actuarial applications of these machine learning techniques and eventually contribute to the development of **insurance predictive analytics** in general.