

Case Reserving in Non-Life Practice using Individual Data and Machine Learning

Marco Aleandri

Ph.D. Student in Actuarial Science

Dept. of Statistical Sciences

Univ. “La Sapienza”, Rome

Email: marco.aleandri@uniroma1.it

Abstract

This paper focuses on the potential of machine learning tools in micro-level reserving by using individual claim data, which is more and more available nowadays. This is especially relevant for non-life insurance, but it could also be useful for some specific life business branches.

After a brief introduction to the problem of reserve estimation in non-life, we will describe the algorithms behind some of the fundamental machine learning tools such as regression methods, naive Bayes, k -nearest neighbors, CARTs, and neural networks. All of them will be used to estimate closing delay and payment amount for individual claims of a specific automobile bodily injury claim dataset. Theoretically, these estimations represent the foundation for a triangle-free, machine-learning-based approach to non-life reserving.

Keywords: individual claims, case reserving, machine learning, generalized regression, naive Bayes, k -nearest neighbors, classification and regression trees, neural network

1. Introduction

Actuarial practice in non-life reserving is traditionally based on aggregate claims data structured in triangles. In fact, this has been proved (for instance, in [2] and [11]) to be an effective approach as long as we face high-probability low-impact claims such as those of motor insurance. Run-off triangles - like that in Figure 1 - are fundamentally based on the assumption that the reserve on future claim payments depend on the reporting year (or accident year) and closing delay only.

Actually, it is far to be true, unless the claims tend to be extremely homogeneous. And even in that case, triangle-based estimations will neglect relevant information about individual claims. This was necessary when actuaries had to use data in times of strong computational limits. Nowadays, this is no longer a major constraint. This is the reason why more and more studies promote micro-level reserving based on individual claims data, for instance [1], [7], [9], [12], [13], and [15]. To some extent, all of them assumes a rather fixed structural form for the timing or the amount of the payments. Unfortunately, such approaches are as rigid as any other parametric method, and cannot take into account all the available details for claims.

In the recent years, some researchers tried to introduce machine learning in actuarial practice by tackling classical problems with new techniques. In fact, such data-driven tools can ultimately overcome the issues explained so far in this section about more traditional approaches. Restricting ourselves to non-life practice, it is worth citing

[18] which applies classification trees to estimate number of future payments varying by accident year and reporting delay,

[19] which utilizes neural network to handle heterogeneity in data and improve Chain-Ladder reserving,

[20] which presents machine learning tools in non-life pricing.

This paper is especially inspired by [18], whose author uses decision trees to predict the number of payments by accident year and closing delay. On our side, we will try to extend those ideas by directly predicting claim amounts using a wider range of machine learning tools, in order to choose the most accurate one. However, a major difference with respect to [18] and, more generally, actuarial practice in non-life reserving, regards the concept of reserve we are going to refer to. Even if the results will be reported in the traditional

reporting year	closing delay					tail
	0	1	...	$n-1$	n	
$T-n$	$P_{T-n,0}$	$P_{T-n,1}$...	$P_{T-n,n-1}$	$P_{T-n,n}$	$\tilde{P}_{T-n,c}$
$T-(n-1)$	$P_{T-(n-1),0}$	$P_{T-(n-1),1}$...	$P_{T-(n-1),n-1}$	$\tilde{P}_{T-(n-1),n}$	$\tilde{P}_{T-(n-1),c}$
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots	\vdots
$T-1$	$P_{T-1,0}$	$P_{T-1,1}$...	$\tilde{P}_{T-1,n-1}$	$\tilde{P}_{T-1,n}$	$\tilde{P}_{T-1,c}$
T	$P_{T,0}$	$\tilde{P}_{T,1}$...	$\tilde{P}_{T,n-1}$	$\tilde{P}_{T,n}$	$\tilde{P}_{T,c}$

Fig. 1: Run-off triangle of paid amounts, and related reserve estimations (in red)

36 reporting-year-per-closing-delay format, our final goal is NOT the estimation
37 of the ultimate loss reserve as at year end. By contrast, we are interested
38 in estimating the so-called *case reserve*, that is, the final cost of each single
39 claim on its own.

40 After a brief introduction to the reserving problem in non-life insurance (see
41 Section 2) and the probabilistic model at the base of our analysis (see Section
42 3), we will recall some features of the machine learning tools we intend to use
43 (see Section 4). All of them will be applied to publicly available automobile
44 bodily injury claim data, in order to estimate individual case reserves. The
45 analysis is presented in Section 5. In line with the underlying model, the
46 results will be broken down into three steps:

- 47 1. closing delay estimation of individual claims through naive Bayes, k -
48 nearest neighbors, and classification tree (see Subsection 5.2)
- 49 2. payment amount estimation of individual claims through generalized
50 regression, regression tree, and neural network (see Subsection 5.3)
- 51 3. case reserve estimation of individual claims through some combinations
52 between tools in 1. and tools in 2. (see Subsection 5.4).

53 While no machine learning tool will be able to explain a relevant amount
54 of variance at the closing delay's step, payment amount estimations will be
55 quite accurate though. More specifically, as pointed out in Section 6, the es-
56 timations returned by decision trees and neural networks will be significantly
57 more accurate than those from generalized regression. Of course, this is not

58 to prove that machine learning can always outperform more traditional ap-
59 proaches. Nonetheless, it should provide actuaries with further techniques to
60 better estimate reserves when it comes with skewed and heavy-tailed claim
61 distributions.

62 2. Understanding claim timeline

63 Generally, a claim is triggered by an accident causing a damage covered
64 by the insurance contract. In an ideal world, the related benefit is paid as
65 soon as the accident occurs, but often this is not the case in non-life insur-
66 ance. In fact, a number of years may pass between the effective occurrence
67 and the final claim payment (or payments). This time gap represents the
68 reason why insurance companies must allocate reserve sufficient to cover any
69 future payments for outstanding loss liabilities.

70 Assume the premium is paid in t_0 for an insurance protection that is immedi-
71 ately effective for a period T . During that period, an accident occurs at time
72 $t_a < T$, the so-called *accident date*. Ideally, the accident is immediately re-
73 ported to the company, but for a number of reasons it may happen differently,
74 that is, the accident is reported at any time $t_r \geq t_a$, the so-called *reporting*
75 *date*. The difference $\Gamma := t_r - t_a$ is the *reporting delay*. If it is small, say
76 days, it does not really represent a problem to the company. However, if the
77 reporting delay extends for years, it generates an unknown outstanding loss
78 liability for the company, which is backed by the so-called *Incurred-But-Not-*
79 *Yet-Reported reserve* - or IBNYR reserve. Actually, the related claims are
80 not in the company's systems yet, so data-driven tools are hardly adaptable
81 to this problem. For this reason, we will not estimate the IBNYR reserve in
82 this paper - it will be shortly discussed in Section 6 as a topic for further
83 development.

84 As soon as the accident is reported, the company is able to collect informa-
85 tion about it, which represents the starting point for our analysis. Typically,
86 the claim cannot be settled immediately for a number of reasons, includ-
87 ing further investigation, new information, court decisions, and so on. As a
88 consequence, the claim is actually closed only at a future date $t_c \geq t_r$, the
89 so-called *closing date*. The difference $\Delta := t_c - t_r$ is the *closing delay*, which
90 may have very different features depending on the specific non-life business
91 involved as well as the claim severity. For standard claims, it might be very
92 small, like in health insurance contracts for the employees of a firm: the
93 company receives standard claim documentation from the policyholder, ap-

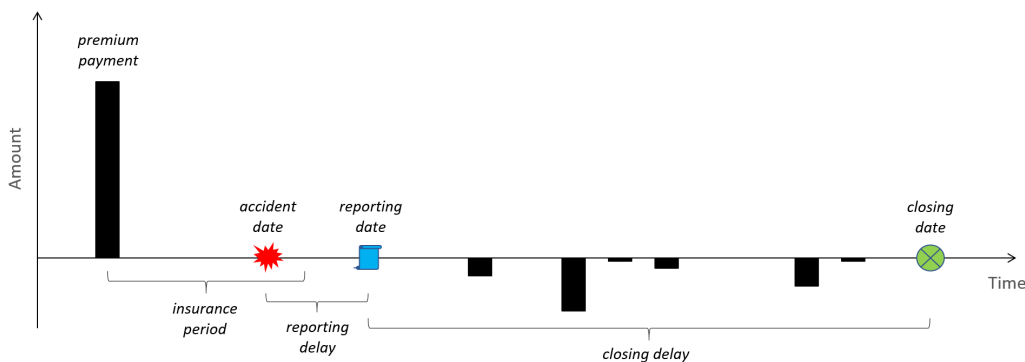


Fig. 2: Graphical representation of the claim timeline

94 prove it quickly, and refund him/her with one single payment. On the other
 95 hand, more severe claims often lead to more tortuous - and longer - closing
 96 delay: sometimes no payment is due, sometimes a final payment is due, and
 97 sometimes company's investigation justifies claim benefits all the way along
 98 and a series of cash-flows is correspondingly paid. At reporting date, the
 99 company must allocate reserve to cover the payments expected during the
 100 closing delay. Such a reserve is the so-called *Reported-But-Not-Yet-Settled*
 101 *reserve* - or RBNYS reserve. Given that it is allocated at reporting date,
 102 when some information about the claim is already known, we can use it to
 103 build our individual reserve estimation.
 104 The process described in this section is represented in a sort of timeline in
 105 Figure 2.

106 3. Assumptions and model

107 Since our valuation date is the reporting date, we can assume that all
 108 the information about the claim is known at that date, and thus it can
 109 be used to predict future payments. This is a first simplification, which
 110 may be unacceptable in some specific cases. In fact, one of the causes of
 111 the closing delay is the further investigation by the company, which could
 112 discover new information at a later date. For sake of simplicity, we will ignore
 113 this possibility.
 114 According to the explanation in Section 2, the company pays an amount to
 115 the policyholder as soon as it is justified by the aforementioned investigation,
 116 thus generating a series of cash-flows. The emerging of such cash-flows is

117 modeled in [18]. However, it would be a further complication for us, and the
 118 available data does not include those details. Therefore, we will assume one
 119 single, aggregate payment for each claim due at closing date.

120 Once these two assumptions are accepted, the model is a rather simple one.
 121 Assume that the closing delay Δ is a discrete variable measured in years, say
 122 $0, \dots, m$. That means: no payment is delayed more than m years after the
 123 reporting date. Moreover, consider a number of predictors x_1, \dots, x_n , that is,
 124 information about the policyholder, the claim, or any other relevant details.
 125 They are all available at reporting date, so they can be used to predict how
 126 likely the payment related to the claim i occurs after k years, that is, $\Delta_i = k$.
 127 Using a proper classification tool, formally represented by a function $p_k(\cdot)$ of
 128 the predictors, we will estimate $P(\Delta_i = k)$ for each admissible k :

$$\widehat{P}(\Delta_i = k) := p_k(x_{i1}, \dots, x_{in}), \quad \forall k \in [0, m]. \quad (1)$$

129 Now, we can still use the same predictors x_{i1}, \dots, x_{in} to estimate the payment
 130 amount due for the claim i ; additionally, we will also use the information
 131 about Δ_i . Using a proper regression tool, formally represented by a function
 132 $f(x_{i1}, \dots, x_{in} | \Delta_i)$ of the closing delay and the predictors, we will estimate
 133 the payment amount:

$$\widehat{C}_i(\Delta_i) := f(x_{i1}, \dots, x_{in} | \Delta_i). \quad (2)$$

134 providing us with an estimation conditioned to Δ_i , which is unknown at
 135 reporting date. To overcome this limit, we first calculate the following esti-
 136 mations:

$$\widehat{C}_i(k) := f(x_{i1}, \dots, x_{in} | k), \quad \forall k \in [0, m] \quad (3)$$

137 and then estimate the payment amount for the claim i as follows:

$$\begin{aligned} \widehat{C}_i &:= \sum_{k=0}^m \widehat{P}(\Delta_i = k) \widehat{C}_i(k) = \sum_{k=0}^m p_k(x_{i1}, \dots, x_{in}) f(x_{i1}, \dots, x_{in} | k) = \\ &= \sum_{k=0}^m p_k(\mathbf{x}_i) f(\mathbf{x}_i | k) \end{aligned} \quad (4)$$

138 which is an unconditioned estimation.

139 In the next sections, we will present some of the machine learning techniques
 140 that can be used to estimate p_k and f . We will separate them because of
 141 their different nature: while p_k estimates the probabilities related to the

142 categorical variable Δ_i , f estimates the claim amount $C_i(k)$. The former
143 refers to a classification problem, whereas the latter refers to a regression
144 problem.

145 4. Fundamental machine learning tools

146 The closing delay estimation is a classification problem, that is, the goal
147 is the prediction of how likely a claim will be closed - and the related amount
148 paid - after k years from the reporting date. By contrast, the claim amount
149 estimation is a regression problem since the target variable is numerical. Sim-
150 ilarly, the closing delay estimation would have been treated as a regression
151 problem too, if we had assumed k as a continuous time variable. However,
152 this would complicate the model, being inconsistent with the traditional as-
153 sumption of a discrete k with some upper limit m , just like in any triangle-
154 based reserve calculation exercise.

155 Although we will distinguish between tools used for closing delay estimation
156 and tools used for claims amount estimation as described at the end of Sec-
157 tion 1, this distinction is not strict. In fact, most of the fundamental machine
158 learning tools we will recall in the following subsections - generalized regres-
159 sion, naive Bayes, k -nearest neighbors, decision trees, and neural networks -
160 are flexible enough to be used for both classification and regression problems.

161 4.1. Generalized regression

162 Regression models are by far the most used fitting tools in industry.
163 They rely on a unique combination of theoretical solidity and practical in-
164 terpretability, which makes them one of the favorite tools among actuaries.
165 For the same reasons, regression models do not represent a powerful machine
166 learning tool. Unfortunately, the assumption set they require causes a sort
167 of rigidity that limits the ability of learning from data. Nonetheless, since we
168 are going to use - among others - a generalized regression model to predict
169 claim amounts (see Subsection 5.3), it is worth providing a brief description.
170 Very generally, we can assume that the target variable y is a function of some
171 predictors, that is, the explanatory variables of the i^{th} record:

$$y_i := \phi(x_{i1}, \dots, x_{in}), \quad \forall i. \quad (5)$$

172 If we can rely on an algorithm that estimates ϕ by building a new regular
173 function f , we will get an estimation of the target variable:

$$\hat{y}_i = \phi(x_{i1}, \dots, x_{in}) + u_i = f(x_{i1}, \dots, x_{in}), \quad \forall i \quad (6)$$

174 where u_i denotes the estimation error (notice the similarities between (3) and
 175 (6)). Quite intuitively, the function f will be somehow estimated in order
 176 to minimize the errors u_i , or, more precisely, a relevant combination of them
 177 representing the overall error of the model. However, the shape of f should
 178 be defined as an assumption. In *multiple linear regression*, for instance, its
 179 shape is linear:

$$\hat{y}_i = \beta_0 + \sum_{j=1}^n x_{ij}\beta_j, \quad \forall i \quad (7)$$

180 where β_j denotes the estimation parameter related to the j^{th} predictor. Re-
 181 member that the linearity refers to the coefficients β_0, \dots, β_n , not to the
 182 predictors x_{i1}, \dots, x_{in} . In other terms, we can use the predictors as they
 183 appear in the dataset as well as any transformation or interaction of them,
 184 linear or nonlinear: in any case, the model will still be linear. Defining

$$\mathbf{y} := \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}, \mathbf{X} := \begin{pmatrix} 1 & x_{11} & \dots & x_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & \dots & x_{Nn} \end{pmatrix}, \beta := \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_n \end{pmatrix}, \mathbf{u} := \begin{pmatrix} u_1 \\ \vdots \\ u_N \end{pmatrix} \quad (8)$$

185 the model (7) may be also expressed in a matricial form:

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{u} \quad (9)$$

186 and

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}. \quad (10)$$

187 The shape of f is not the only assumption of the model. Rather, the follow-
 188 ings should hold too:

- 189 1. $r(\mathbf{X}) = n$
- 190 2. $E[\mathbf{u}] = \mathbf{0}$
- 191 3. $Var[\mathbf{u}] = \sigma^2\mathbf{I}_N$.

192 In other words, the rank of \mathbf{X} is the maximum possible rank so that no re-
 193 dundant information is there (assumption 1.), the random variables of the
 194 estimation error are null on average (assumption 2.), and they are also inde-
 195 pendent and homoscedastic with variance σ^2 (assumption 3.). In particular,
 196 the assumption 2. directly implies the followings:

- 197 2a. $E[\mathbf{y}] = E[\mathbf{y}|\mathbf{X}] = \mathbf{X}\beta = \hat{\mathbf{y}}$

198 2b. $Var(\mathbf{y}) = Var(\mathbf{y}|\mathbf{X}) = \sigma^2\mathbf{I}_N$.

199 Once the assumptions have been verified, the vector of parameters β can be
200 estimated through the *least squares method*. It is based on the minimization
201 of the sum of squared error u_i^2 :

$$\begin{aligned} Q(\beta) &:= \sum_{i=1}^N u_i^2 = \sum_{i=1}^N (y_i - \hat{y}_i)^2 = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) = \\ &= \mathbf{y}^T \mathbf{y} - \beta^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\beta + \beta^T \mathbf{X}^T \mathbf{X}\beta = \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\beta + \beta^T \mathbf{X}^T \mathbf{X}\beta \end{aligned} \quad (11)$$

202 which is a quadratic function, thus its minimum is necessarily its only sta-
203 tionary point, that is, the solution of the following system:

$$\frac{\partial Q(\beta)}{\partial \beta} = \frac{\partial \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\beta + \beta^T \mathbf{X}^T \mathbf{X}\beta}{\partial \beta} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\beta = 0 \quad (12)$$

204 that is

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \quad (13)$$

205 Using (10), we also get

$$\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (14)$$

206 and using (9)

$$\hat{\mathbf{u}} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = [\mathbf{I}_N - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{y}. \quad (15)$$

207 Actually, we got a model now, that is, we are able to calculate $\hat{\mathbf{y}}$ by using
208 $\hat{\beta}$ in (10), as long as all the aforementioned assumptions hold. However, the
209 traditional multiple linear regression includes a further, crucial assumption:

$$\mathbf{u} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_N) \quad (16)$$

210 which implies

$$\mathbf{y} \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_N). \quad (17)$$

211 Thanks to the normality of the vector \mathbf{u} of the residuals, a whole range
212 of relevant statistics can be deduced from the model, including parameter
213 distributions, confidence interval, and so on. Among others, for instance, it
214 can be easily proved that $\hat{\beta}$ is normally distributed with mean β , and $\hat{\mathbf{y}}$ is
215 normally distributed with mean \mathbf{y} .

216 Once the model has been defined, we want an effective way to evaluate its
 217 performance. The theoretical framework helps us to define a rather intuitive
 218 but rigorous tool for performance evaluation. First, the *total deviance* is
 219 defined as

$$D_T(\mathbf{y}) := \sum_{i=1}^N (y_i - E[\mathbf{y}])^2 \quad (18)$$

220 and it represents a measure of the information contained into the data. The
 221 (18) may be further manipulate as follows:

$$\begin{aligned} D_T(\mathbf{y}) &= \sum_{i=1}^N (y_i - \hat{y}_i + \hat{y}_i - E[\mathbf{y}])^2 = \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \sum_{i=1}^N (\hat{y}_i - E[\mathbf{y}])^2 =: \\ &=: D_E(\mathbf{y}) + D_R(\mathbf{y}) \end{aligned} \quad (19)$$

222 where the *explained deviance* $D_E(\mathbf{y})$ measures the information explained by
 223 the model, while the *residual deviance* $D_R(\mathbf{y})$ measures the residual infor-
 224 mation that the model could not detect. Naturally, the greater $D_E(\mathbf{y})$, the
 225 smaller $D_R(\mathbf{y})$, the better the model. In fact, the *coefficient of determination*
 226 - known as *R-squared* and defined by the deviance measures - is the most
 227 common performance evaluation tool for multiple linear regression:

$$R^2(\mathbf{y}) := \frac{D_E(\mathbf{y})}{D_T(\mathbf{y})} = 1 - \frac{D_R(\mathbf{y})}{D_T(\mathbf{y})} \quad (20)$$

228 which is a value in $[0, 1]$. Once again, the greater $D_E(\mathbf{y})$, the greater $R^2(\mathbf{y})$,
 229 the better the model, but the R^2 has the additional advantage of the normal-
 230 ization. In practice, it means that the R^2 makes possible a fair performance
 231 comparison among different models.

232 Using (17), the model definition (9) may be written differently:

$$\mathbf{y} = E[\mathbf{y}] + \mathbf{u} \quad (21)$$

233 which is a direct consequence of the normality assumed. However, if we relax
 234 that assumption, it can be more generally assumed that there is a regular (i.e.,
 235 invertible and derivable) *link function* g - different to the identity function -
 236 mapping $E[\mathbf{y}]$ to $\mathbf{X}\beta$:

$$g(E[\mathbf{y}]) = \mathbf{X}\beta \quad (22)$$

237 so that, using (21),

$$\mathbf{y} = g^{-1}(\mathbf{X}\beta) + \mathbf{u}. \quad (23)$$

238 The main reason why such a *generalized regression* is very useful in practice is
 239 that g maps from some subset $X \subseteq \mathbb{R}$ to \mathbb{R} , that is, g^{-1} transforms the linear
 240 predictor $\mathbf{X}\beta \in \mathbb{R}$ to the target variable prediction in X . In fact, generalized
 241 regression is able to handle target variables defined in a specific subset of \mathbb{R} .
 242 This is not possible in multiple linear regression, unless we operate a proper
 243 transformation of the target variable itself (this is sometimes enough, but it
 244 introduces *transformation bias* into the model). For instance, if the target
 245 variable represents a positive amount, the quickest ways to use regression are

- 246 • convert amounts to logarithmic amounts, predict the latter through mul-
 247 tiple linear regression, and convert the predictions back by using the
 248 exponential function;
- 249 • choose a link function $g : (0, +\infty) \rightarrow \mathbb{R}$, and predict the amounts
 250 through the related generalized regression.

251 Actually, the choice of the link function is not direct, rather it is a conse-
 252 quence of the distribution we assume for the target variable. This is possible
 253 since generalized regression assumes that such a distribution belongs to the
 254 exponential family, whose density is:

$$f_e(y_i; \theta_i, \phi) := e^{\frac{y_i \theta_i - c(\theta_i)}{\phi} + h(y_i, \phi)} \quad (24)$$

255 where θ_i denotes the *canonical parameter* varying by observation, while ϕ
 256 denotes the *dispersion parameter* that is constant for all observations. In
 257 other words, the distribution constraint is not completely eliminated, rather
 258 it is “generalized” to a wider range of distributions - a distribution family.
 259 When it comes with prediction of amounts, whose distributions are often
 260 nonnegative and heavy-tailed, a common choice is the Gamma distribution,
 261 defined by a *shape parameter* $\varphi > 0$ and a *scale parameter* $\vartheta_i > 0$:

$$f_\Gamma(y_i; \vartheta_i, \varphi) := \frac{y_i^{\varphi-1} e^{-\frac{y_i}{\vartheta_i}}}{\vartheta_i^\varphi \Gamma(\varphi)}. \quad (25)$$

262 which can be easily rearranged as follows:

$$f_\Gamma(y_i; \vartheta_i, \varphi) = e^{-\frac{y_i}{\vartheta_i} - \varphi \ln \vartheta_i + (\varphi-1) \ln y_i - \ln \Gamma(\varphi)}. \quad (26)$$

263 If we define $\theta_i := -\frac{1}{\varphi \vartheta_i}$ and $\phi := \frac{1}{\varphi}$, then

$$f_\Gamma(y_i; \vartheta_i, \varphi) = e^{\frac{y_i \theta_i}{\phi} - \frac{1}{\phi} \ln \left(-\frac{1}{\varphi \vartheta_i}\right) + \left(\frac{1}{\phi} - 1\right) \ln y_i - \ln \Gamma\left(\frac{1}{\phi}\right)} =$$

$$\begin{aligned}
&= e^{\frac{y_i \theta_i}{\phi} - \frac{1}{\phi} \ln\left(-\frac{1}{\theta_i}\right) + \left(\frac{1}{\phi} - 1\right) \ln y_i + \varphi \ln \phi - \ln \Gamma\left(\frac{1}{\phi}\right)} = \\
&= e^{\frac{y_i \theta_i - \ln\left(-\frac{1}{\theta_i}\right)}{\phi} + \left(\frac{1}{\phi} - 1\right) \ln y_i - \frac{1}{\phi} \ln \phi - \ln \Gamma\left(\frac{1}{\phi}\right)} \quad (27)
\end{aligned}$$

264 which belongs to the exponential family in (24) with

$$c(\theta_i) := \ln\left(-\frac{1}{\theta_i}\right), \quad h(y_i, \phi) := \left(\frac{1}{\phi} - 1\right) \ln y_i - \frac{1}{\phi} \ln \phi - \ln \Gamma\left(\frac{1}{\phi}\right). \quad (28)$$

265 This is important because it implies that

$$E[y_i] = \frac{dc(\theta_i)}{d\theta_i} = \frac{d}{d\theta_i} \ln\left(-\frac{1}{\theta_i}\right) = -\frac{d}{d\theta_i} \ln \theta_i = -\frac{1}{\theta_i} = \varphi \vartheta_i \quad (29)$$

$$\text{Var}(y_i) = \phi \frac{d^2 c(\theta_i)}{d\theta_i^2} = -\phi \frac{d}{d\theta_i} \frac{1}{\theta_i} = \frac{\phi}{\theta_i^2} = \varphi \vartheta_i^2. \quad (30)$$

266 All in all, how should we define the link function starting from these consider-
267 ations? First, notice that the normal distribution belongs to the exponential
268 family too, because

$$f_N(y_i; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mu)^2}{2\sigma^2}} = \dots = e^{\frac{y_i \mu - \frac{\mu^2}{2}}{\sigma^2} - \frac{y_i^2}{2\sigma^2} - \ln \sqrt{2\pi\sigma^2}} \quad (31)$$

269 thus we should define $\theta_i := \mu$ and $c(\theta_i) := c(\mu) = \frac{\mu^2}{2}$, so that

$$E[y_i] = \frac{dc(\mu)}{d\mu} = \mu = \theta_i. \quad (32)$$

270 Actually, if the target variable is normally distributed, we will be back to the
271 multiple linear regression. In such a case, we may use (21) to write:

$$\mathbf{y} = E[\mathbf{y}] + \mathbf{u} = \boldsymbol{\theta} + \mathbf{u} \quad (33)$$

272 so a generalized regression model might be simply redefined as

$$\mathbf{y} = g^{-1}(\boldsymbol{\theta}) + \mathbf{u}. \quad (34)$$

273 In fact, the link function is now the function that maps $E[\mathbf{y}]$ to the vector $\boldsymbol{\theta}$
274 of the canonical parameters:

$$g(E[\mathbf{y}]) = \boldsymbol{\theta} \quad (35)$$

275 but if the target variable is distributed as a Gamma, we also know that
276 $E[y_i] = \varphi\vartheta_i$ and $\theta_i = -\frac{1}{\varphi\vartheta_i}$, so

$$g(\varphi\vartheta_i) = -\frac{1}{\varphi\vartheta_i} \quad (36)$$

277 which means that g may be simply defined as $g(x) := -\frac{1}{x}$. This is a somewhat
278 natural choice, coming directly from the theory, thus it is called *canonical link*
279 *function*. For generalized regression models, it is always possible to define g
280 in a canonical way, that is, imposing (35), and this choice implies a range
281 of desirable features in our model. However, others link functions might be
282 rather used, and sometimes there are good reasons to.

283 Unfortunately, parameters of generalized regression models cannot be explic-
284 itly calculated as in the multiple linear regression. In this case, we can only
285 look for maximum likelihood estimates by using numerical methods.

286 At the beginning of this subsection, we implicitly assumed to know the range
287 of explanatory variables x_1, \dots, x_n . Of course, we know the explanatory vari-
288 ables in the dataset, but how should we select them as x_1, \dots, x_n ? Because
289 of multicollinearity among potential explanatory variables, we cannot simply
290 run the regression on all of them, and then select only the most significant
291 ones based on their p-values. Rather, we should somehow select different
292 sets of explanatory variables and run the related regression: the model with
293 the highest R^2 will be selected. The different sets of explanatory variables
294 depend on the algorithm used to select them. There are mainly three popular
295 iterative search algorithms.

296 In *forward selection*, we start with no predictors, and then add them one by
297 one. Each added predictor is that (among all predictors) that has the largest
298 contribution to R^2 on top of the predictors that are already in it. The algo-
299 rithm stops when the contribution of additional predictors is not statistically
300 significant.

301 In *backward selection*, we start with all predictors, and then eliminate the
302 least useful one at each step according to statistical significance. The algo-
303 rithm stops when all the remaining predictors have significant contributions.

304 Finally, *stepwise selection* is like forward selection except that at each step
305 we consider dropping predictors that are no longer statistically significant,
306 as in backward selection. In the Subsection 5.3, our data will be regressed
307 through stepwise selection.

308 To finally conclude the section, some practical remarks are to be highlighted.

309 Whatever the specific regression model, its assumptions make it theoret-
 310 ically strong, but VERY weak from a practical perspective. Actually, We
 311 may condensate them with the following, practical premises:

- 312 • no allowance for nonlinear relationships
- 313 • no allowance for dependencies among predictors
- 314 • no allowance for outliers.

315 They somewhat simplify the original assumptions, but give a fair idea of what
 316 those assumptions really mean, i.e. we can hardly expect good performance
 317 from regression methods when information is far to be regular. Data is
 318 generally affected by missing values, redundancies, correlations, heavy tails,
 319 asymmetries, nonlinearities, and any other kind of distortion. In fact, as we
 320 will discuss in Subsection 5.1, our data is no exception.

321 4.2. Naive Bayes

322 This tool is surely one of the easiest machine learning techniques. It is
 323 a transformation of the well-known *Bayesian classifier*. Given the values for
 324 the predictor vector \mathbf{x}_i related to the claim i , the Bayes' theorem returns the
 325 (exact) Bayesian classifier:

$$\begin{aligned}
 p_k(\mathbf{x}_i) &= P(\Delta_i = k | \mathbf{x} = \mathbf{x}_i) = \\
 &= \frac{P(\Delta_i = k)P(\mathbf{x} = \mathbf{x}_i | \Delta_i = k)}{\sum_{h=0}^m P(\Delta_i = h)P(\mathbf{x} = \mathbf{x}_i | \Delta_i = h)}, \quad \forall k = 1, \dots, m. \quad (37)
 \end{aligned}$$

326 This approach is theoretically correct, but presents a fundamental limit. The
 327 predictor in (37) implicitly assumes that we can find a sufficient number of
 328 records in the sample sharing the same vector \mathbf{x}_i . Perhaps, this is reasonable
 329 when there are VERY few predictors in the dataset, otherwise it is completely
 330 impracticable.

331 A straight modification to (37) represents a very simple solution to such a
 332 problem. If we give up to the assumption that the best probability estimation
 333 is solely returned by those records matching the record to be classified, we
 334 will be able to use the whole dataset for the estimation. As a consequence
 335 of this assumption, the classifier in (37) changes as follows:

$$\begin{aligned}
 p_k(\mathbf{x}_i) &= P(\Delta_i = k | \mathbf{x} = \mathbf{x}_i) = \frac{P(\Delta_i = k)P(\mathbf{x} = \mathbf{x}_i | \Delta_i = k)}{\sum_{h=0}^m P(\Delta_i = h)P(\mathbf{x} = \mathbf{x}_i | \Delta_i = h)} = \\
 &= \frac{P(\Delta_i = k) \prod_{j=1}^n P(x_j = x_{ij} | \Delta_i = k)}{\sum_{h=0}^m P(\Delta_i = h) \prod_{j=1}^n P(x_j = x_{ij} | \Delta_i = h)}, \quad \forall k = 1, \dots, m. \quad (38)
 \end{aligned}$$

336 which is, indeed, the so-called *naive Bayes classifier*.
337 This approach is extremely simple to understand and to use. Moreover, it
338 presents no computational issues: it is just a formula to apply as it is, rather
339 than a complex algorithm. Unfortunately, this simplicity hides a major draw-
340 back, that is, the assumption of stochastic independence among predictors.
341 In effect, that is exactly the assumption which allows us to move from the
342 Bayesian classifier in (37) to the naive Bayes classifier in (38). Comparing
343 the two formulas, we can easily notice this point. However, this is seldom
344 the case.

345 Moreover, some studies (see, for instance, [10]) point out the strengths of such
346 a classifier when it comes with record ranking, but also its weaknesses when
347 it comes with probability estimation. In practice, naive Bayes classifier's
348 probability estimation can be very biased by the assumption of stochastic
349 independence. If the same bias is shared by each record, the classification
350 can be still good, but of course we cannot rely on the estimated probability.
351 This is the reason why this classifier often outperforms more sophisticated
352 classifiers as a classification tool, and it is still widely used in several fields
353 (see, for instance, the spam filtering case in [14]).

354 A last drawback is quite relevant. Actually, what if some predictor category
355 is not present in the training dataset (for instance, it could be very rare)?
356 In this case, $P(x_j = x_{ij} | \Delta_i = k) = 0$ for some j and k , thus $p_k(\mathbf{x}_i) = 0$,
357 which is clearly wrong. In fact, the naive Bayes classifier works well if each
358 and every category is well represented. That has a two-fold meaning. First,
359 the training dataset should be large enough to well represent each and every
360 category. Second, more importantly, numerical predictors are not admissible
361 at all by definition, that is, we can use it for closing delay estimation (which
362 is a classification problem), but not for claim amount estimation (which is a
363 regression model). For the latter, we need the methods described in the next
364 subsections.

365 4.3. *K-nearest neighbors*

366 The idea behind the k -nearest neighbors algorithm is very intuitive, but
367 it still guarantees a high level of adaptability to data. To score a new record,
368 the method relies on finding the most “similar” records - the so-called “neigh-
369 bors” - in the training dataset. In fact, this is a pure nonparametric method:
370 no assumption needs to be established, no parameter needs to be estimated,
371 no functional form needs to be assumed.

372 The sole issue regards the choice of a measure to calculate the “distance”

373 between two records, that is, their grade of similarity. The most popular
 374 measure is the *Euclidean distance*: given the vectors of predictors for two
 375 different records, say $\mathbf{x}_i = (x_{i1}, \dots, x_{in})$ and $\mathbf{x}_j = (x_{j1}, \dots, x_{jn})$, the Eu-
 376 clidean distance between them is defined as follows:

$$d(\mathbf{x}_i, \mathbf{x}_j) := \sqrt{\sum_{k=1}^n (x_{ik} - x_{jk})^2}, \quad \forall i, j. \quad (39)$$

377 It is worth noting that this definition of distance would give much more
 378 weight to higher scales than lower scales, so all the predictors should be first
 379 standardized before computing the (39). Otherwise, the k -nearest neighbors
 380 algorithm could result in extremely biased predictions.

381 Once we have chosen a distance measure, we can calculate the distance be-
 382 tween any pair of records, but we still need a rule to score new records. The
 383 simplest rule is: a new record is classified in the same category of its closest
 384 neighbor. Therefore, given the predictors of such a record, we will compute
 385 all the distances between it and the records in the training dataset. Among
 386 them, pick up that with the smallest distance to the new record: its category
 387 will be assigned to the new record itself. In effect, we have just applied a
 388 1-nearest neighbors algorithm.

389 Nonetheless, the approach might be easily generalized to any number k of
 390 neighbors. Instead of picking up the closest record only, pick up the k closest
 391 records, and assign the majority class among them to the new record. In
 392 practice, the usage of more neighbors tends to reduce misclassification error.
 393 If we use one neighbor only, it could be the case that a record is the closest
 394 one to the new record only by chance: there could be noise there, rather
 395 than information. To some extent, the greater the k , the lower the error, the
 396 greater the predictive power. On the other hand, however, if k is too high,
 397 we will miss out on the method's ability to capture the local structure in the
 398 data, that is, we will ignore information.

399 Anyway, the choice of k is straightforward. First, choose an upper limit K
 400 for k . Then, score each record in the *validation* dataset using the closest
 401 record in the training dataset (1-nearest neighbors algorithm), and calculate
 402 the validation error ϵ_1 . Likewise, score each record in the *validation* dataset
 403 using the two closest records in the training dataset (2-nearest neighbors
 404 algorithm), and calculate the validation error ϵ_2 . Repeat the process until
 405 the K -nearest neighbors algorithm, that is, the last admissible algorithm
 406 according to our upper limit. At the end, we will get the validation errors

407 $\epsilon_1, \dots, \epsilon_K$. Finally, pick up the k related to the smallest validation error, say
408 k_{min} , and use the k_{min} -nearest neighbors algorithm for scoring.
409 Remember that ϵ_k denotes the validation error of k -nearest neighbors algo-
410 rithm. Correspondingly, we could compute the training error as well, that is,
411 the error committed when scoring each record in the *training* (rather than
412 validation) dataset using the k closest records in the training dataset itself.
413 The training error, however, cannot be used for scoring because the smallest
414 training error always relates to the 1-nearest neighbors algorithm: whichever
415 the training record to be scored, its closest training record is obviously the
416 record itself!
417 Of course, the k -nearest neighbors classifier is as simple as the naive Bayes
418 classifier. As already discussed at the beginning of this section, its main ad-
419 vantage is the lack of parametric assumptions. Unfortunately, there are some
420 drawbacks. For instance, from a computational perspective, this algorithm
421 can be very expensive. Sometimes, data scientists try to reduce predictors
422 by using other, less expensive machine learning tools before applying the k -
423 nearest neighbors algorithm to their datasets. Dimension reduction may be
424 performed, among others, by classification trees, which is the topic of the
425 next subsection.

426 4.4. Classification and regression trees (CARTs)

427 Decision trees were used as a machine learning tool in [3] for the very
428 first time to segment a population by splitting up the dataset through bi-
429 nary rules. The algorithm is now referred to as *classification tree*. Since one
430 of our goals is the categorical classification among the closing delay classes,
431 this tool is a good candidate for us. By contrast, we should properly adapt it
432 to regression problems, if we want to use it to predict the claim amount. In
433 this case, we call it a *regression tree*. However, given that the basic algorithm
434 does not change, we can always refer to *classification and regression trees*, or
435 CARTs.

436 CARTs are based on recursive partitioning, which divides up the multidimensional space (that is, the dataset) of the explanatory variables into non-overlapping multidimensional rectangles. This division is accomplished recursively, i.e. operating on the results of the prior divisions. An example of CART is in Figure 3. First, one of the explanatory variables is selected, say $x_{k(0)}$ (the first *node* of the tree, so-called *root*), and a value of $x_{k(0)}$, say $s_{k(0)}$, is chosen to split the n -dimensional space into two parts: one part contains all the records with $x_{k(0)} \leq s_{k(0)}$, say $n(1, 1)$ records, while the other with all

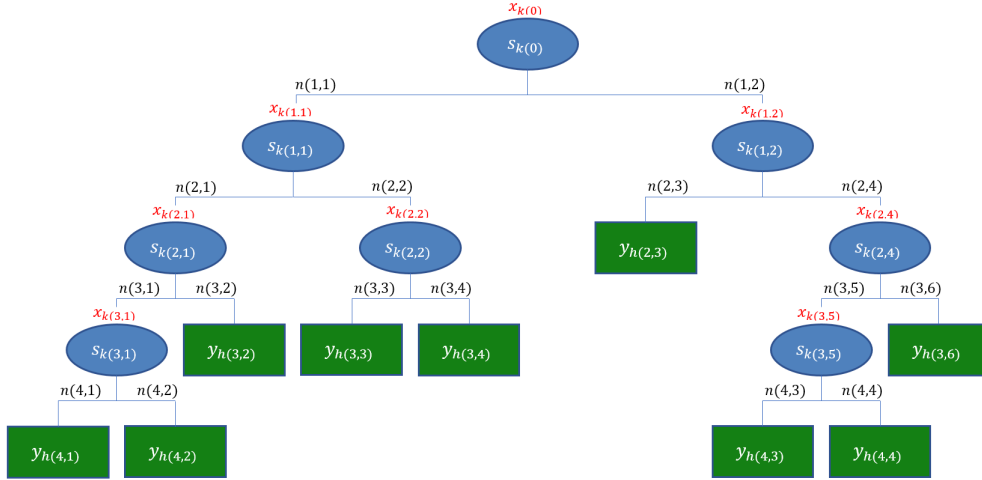


Fig. 3: Example of CART produced by recursive partitioning

444 the records with $x_{k(0)} > s_{k(0)}$, say $n(1,2)$ records. The two sub-datasets rep-
 445 resent the first level of the tree. Let's consider one of them: it could be either
 446 *pure* - i.e. it contains only records sharing the same value of the independent
 447 variable - or *impure*. In the first case, no further split is possible, so the
 448 sub-dataset will represent a *leaf* of the tree. In the second case, other splits
 449 are possible, so the sub-dataset will represent another node of the tree. In
 450 Figure 3, leaves are denoted by green rectangles, while nodes are denoted by
 451 blue circles. Unless both of the sub-datasets generated by the root node are
 452 pure, one of them (at least) will be divided in a similar manner by choosing a
 453 variable again, and a split value for the variable. In Figure 3, where both of
 454 the sub-datasets of the first level are impure, they represent the two nodes of
 455 the first level, and the initial dataset is further partitioned into four regions:

- 456 • the first one contains the $n(2,1)$ records with $x_{k(1,1)} \leq s_{k(1,1)}$, and it will
 457 represent a node for the next level using $x_{k(2,1)}$ as a splitting variable
- 458 • the second one contains the $n(2,2)$ records with $x_{k(1,1)} > s_{k(1,1)}$, and
 459 it will represent a node for the next level using $x_{k(2,2)}$ as a splitting
 460 variable
- 461 • the third one contains the $n(2,3)$ records with $x_{k(1,2)} \leq s_{k(1,2)}$, and it
 462 will represent a leaf for the next level containing all the records with
 463 target variable equal to $y_{h(2,3)}$

464 • the fourth one contains the $n(2, 4)$ records with $x_{k(1,2)} > s_{k(1,2)}$, and
 465 it will represent a node for the next level using $x_{k(2,4)}$ as a splitting
 466 variable.

467 Since some splits are still possible, the recursive partitioning goes on, getting
 468 smaller and smaller sub-datasets, either nodes or leaves. Sooner or later, we
 469 will have divided the whole dataset up into pure sub-datasets (of course, this
 470 is not always possible, as there may be records that belong to different classes
 471 but have exactly the same values for everyone of the predictor variables).

472 In the case of closing delay estimation, the dataset will be partitioned into
 473 sub-datasets which contain either claims closed in the reporting year (closing
 474 delay 0), or claims closed the year after (closing delay 1), or claims closed
 475 two years after (closing delay 2), ..., or claims closed m years after (closing
 476 delay m). In fact, the classification tree resulting from recursive partitioning
 477 is a *pure* tree: all the closing delay categories are perfectly separated.

478 The main problem of recursive partitioning is the choice of the splitting
 479 rule node by node, that is, the choice of $x_{k(\cdot,\cdot)}$ and $s_{k(\cdot,\cdot)}$ at each step of
 480 the algorithm. Assume to define an *impurity function* $i(A)$ as an impurity
 481 measure of some rectangle A , or its related node. A specific splitting rule
 482 on A results in two sub-rectangles A_L and A_R , which are generally impure,
 483 that is, $i(A_L)$ and $i(A_R)$ are both nonzero. Intuitively, we want to choose
 484 the splitting rule in order to minimize some combination of $i(A_L)$ and $i(A_R)$.
 485 The most natural choice is the function

$$I(A_L, A_R) := \frac{|A_L|}{|A|}i(A_L) + \frac{|A_R|}{|A|}i(A_R) \quad (40)$$

486 which is the average of the two impurity measures, weighted by the number
 487 of observations in each rectangle. By comparing the reduction in $I(A_L, A_R)$
 488 across all possible splits in all possible predictors, the next split is chosen.

489 What about the impurity function i ? In our application and in most of them,
 490 one uses the *Gini index* (as defined in [14]):

$$i(A) := 1 - \sum_{k=0}^m p_k^2(A) \quad (41)$$

491 where p_k is the proportion of records in rectangle A that are closed after k
 492 year from the reporting. However, other impurity measures are also widely

493 used, for example the *entropy index* (as defined in [14]):

$$E(A) := - \sum_{k=0}^m p_k(A) \log_2[p_k(A)]. \quad (42)$$

494 All in all, so far the algorithm is quite intuitive as well as its application in
495 classifying new records. For instance a new observation, whose explanatory
496 values are known, will be dropped down the tree until it reaches a leaf. So
497 the new observation will be simply classified on the base of the specific leaf's
498 classification.

499 As discussed at the beginning of the section, the algorithm can be easily
500 adapted to predict numerical variables. We only need some changes. First,
501 the response value assigned to a record in a leaf is determined by the av-
502 erage of the response variable among the records in that leaf (by contrast,
503 in a classification tree, it is determined by one of the possible category of
504 the response variable). Second, given that we cannot use discrete measure of
505 impurity such as the Gini index and the entropy index, the typical impurity
506 measure used in regression problems is the sum of squared deviations from
507 the mean, that is, the sum of squared errors.

508 By definition, recursive partitioning produces a tree which classify the records
509 without errors. Actually, we used a training dataset to train the classification
510 tree, which perfectly predict closing delay on that dataset. But what if we
511 use the same tree on the validation dataset? In general, the predicted values
512 on the validation dataset will result in a positive *misclassification error*. In
513 fact, the error cannot be zero on datasets other than the training dataset
514 itself. However, our full classification tree has a major drawback represented
515 by Figure 4. As it usually happens through the first splits on the validation
516 dataset, the full tree can still guarantee comparable misclassification errors
517 on the two datasets. However, as the number of splits increases, the full
518 tree starts *overfitting* the validation data: since it fully reflects the train-
519 ing dataset without distinguishing between “signal” and “noise”, the noisy
520 component cause too high misclassification error in the validation dataset.
521 Indeed, the typical consequence of overfitting is that, after some number of
522 splits, the misclassification error on the validation dataset stops decreasing
523 and starts increasing (in Figure 4, it occurs after ten splits). In the first ten
524 splits both the training and validation misclassification errors decrease, but
525 thereafter the full tree overfits the validation data.

526 Overfitting prevent us from using the full tree for predicting purposes, so we

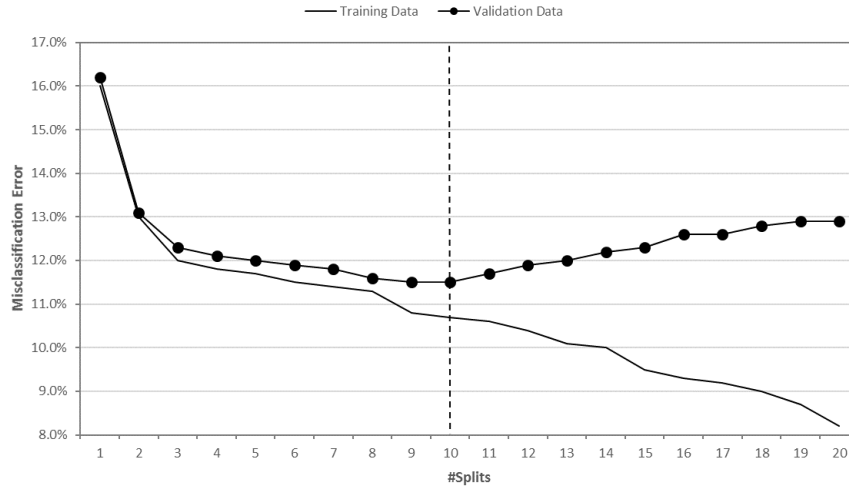


Fig. 4: Training Misclassification Error vs. Validation Misclassification Error

527 need to choose another tree, that is, some subtree of the full tree. There
 528 are several criteria to do that, but two major categories of methods can be
 529 distinguished:

- 530 • forward stopping-tree methods
- 531 • backward pruning-tree methods.

532 Some empirical forward methods can be easily implemented by setting condi-
 533 tions to the tree characteristics such as maximum number of splits, minimum
 534 number of records in a node, and minimum reduction in impurity. Unfortu-
 535 nately, these approaches are solely based on the tree complexity, rather than
 536 its predictive power.

537 Among the forward methods, the most natural choice is suggested by the
 538 Figure 4 itself, i.e. we can simply use the tree consisting of the the first n
 539 splits that do not induce overfitting ($n = 10$ in Figure 4). In other words, we
 540 let the full tree grow until the first step leading to a validation error higher
 541 than the previous step. If we have new observations to classify, they will be
 542 dropped down this subtree until they reach a leaf.

543 More complex methods have been developed as well, for instance, the so-
 544 called *chi-squared automatic in-training data* (CHAID). At each node, the
 545 algorithm select the predictor with the strongest association to the target

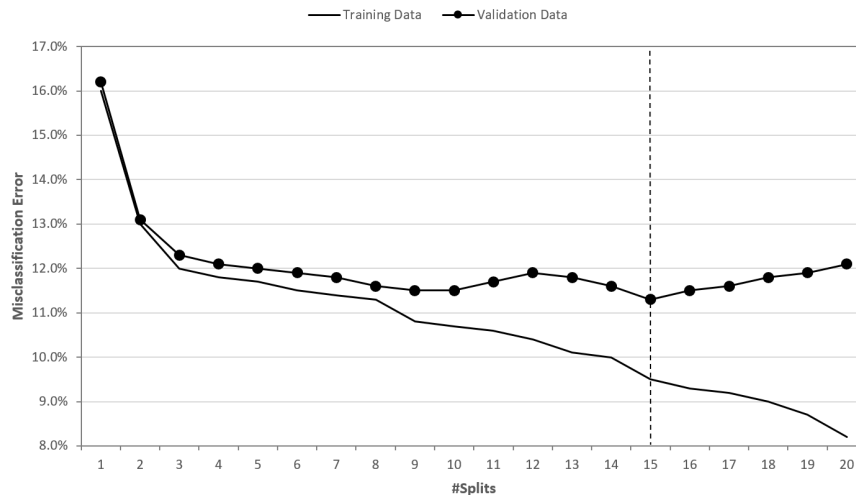


Fig. 5: Training Misclassification Error vs. Validation Misclassification Error

546 variable, measured by the p -value of the chi-squared test of independence. If
 547 such a p -value is low enough, the split of the node significantly improves its
 548 purity, so the algorithm will carry it out, and the growth of the tree goes on.
 549 Otherwise, the growth is stopped.

550 The alternative to stopping the growth of the tree is represented by prun-
 551 ing the tree, that is, “climb” the training full tree and “chop” the weakest
 552 branches until some conditions are met. Intuitively, pruning the tree is more
 553 computationally expensive, because we have to build the full tree anyway, and
 554 then work on that further. However, it has been proven to be more successful
 555 too, and it is not so surprising. In effect, the tree is pruned considering more
 556 information: not only what-if-we-keep-the-smaller-tree information, but even
 557 what-if-we-take-the-bigger-tree information.

558 A simple backward approach follows the idea of picking the tree with the
 559 smallest validation error, just like in one of the forward methods. However,
 560 the full tree is still built until the last leaf, and then we choose the subtree
 561 leading to the smallest validation error. Notice that we would still pick the
 562 subtree after 10 splits in Figure 4, but that is not always the case (see, for
 563 instance, Figure 5).

564 However, choosing the subtree according to the smallest validation error only
 565 means we completely ignore the complexity of the tree. For instance, take
 566 a look at Figure 5 once again. We are picking the 15-split subtree since it

567 leads to the smallest error, but we get a rather small error with the 10-split
568 subtree as well. Actually, we accept five splits more - a relevant increase in
569 complexity - for a little decrease in error. Somehow, it does not seem to be
570 the best choice. To consider this issue in pruning the tree, we may use the
571 so-called *cost complexity criterion* as described in [3]. For a tree with L leaves
572 and training error ϵ , the cost complexity of a tree T is defined as follows:

$$\gamma(T) := \alpha L(T) + \epsilon(T) \tag{43}$$

573 where α denotes a (nonnegative) penalty factor for the tree size. Notice
574 that, if $\alpha = 0$, there is no penalty for the tree size, and the best tree is
575 simply the full tree itself. On the other hand, the greater the α , the greater
576 the relevance of the tree size in pruning the tree. If α is great enough, the
577 training error is no longer relevant for the algorithm, and the tree is pruned
578 until the very first node, that is, the root. The algorithm therefore starts
579 with the n -level full tree, and compares its γ with the γ s of all the possible
580 $n - 1$ -level subtrees. The α gets increased little by little, until the γ of the
581 full tree exceeds that of one of those subtrees. Such a subtree is considered
582 as the best of its level, and the same procedure is repeated starting from it.
583 In fact, the algorithm finds the best subtree of each size based on the cost
584 complexity criterion: the lower γ , the better the subtree. In practice, we get
585 a sequence of “best subtrees” for their sizes, based on the training data only.
586 Finally, the so-called *best pruned tree* used for scoring will be the subtree
587 among them leading to the smallest validation error.

588 All in all, the full tree is useless for scoring purposes, rather it is just the
589 formal result of recursive partitioning. What is really useful for prediction is
590 the subtree extracted by using one of the several algorithms to stop growing
591 or pruning the full tree.

592 So far we described the fundamentals of CARTs. One of the reasons for their
593 popularity is that they are adaptable to a wide variety of applications, and
594 have been successfully used in many situations. In particular, if there is a
595 highly non-linear and complex relationship to describe, decision trees may
596 outperform regression models. Furthermore, CARTs do not require massive
597 data preparation, that is, they can handle non-standardized data, categorical
598 data, missing data, outliers, and so on. By contrast, we should standardize
599 the variables and take the natural logarithm of some numerical variables
600 before running standard linear regression, logistic regression, or even the k-
601 nearest neighbors method described in Subsection 4.3. Finally, trees provide
602 easily understandable classification rules (at least if they are not too large),

603 even easier than in regression.

604 An important advantage of CARTs is that no further selection algorithm is
605 necessary. As opposed to parametric regression methods, the process itself
606 selects the most relevant explanatory variables. We simply let the machine
607 learning tool run on the whole dataset, and the resulting tree will include
608 only some of the explanatory variables, which are the most significant on the
609 base of the impurity measure used to split the dataset.

610 4.5. Neural networks

611 Algorithms like k-nearest neighbors method or CARTs have the major
612 advantage of non-dependence on underlying structures or parameters. This
613 is what makes them perfect to handle - and, to some extent, discover - un-
614 known relationships in datasets. On the other hand, this flexibility makes
615 them extremely weak when data is not enough to train them properly. *Neural*
616 *networks* are trained by data too, but assume an underlying function that is
617 generally much more complex than the typical functions used in regression.
618 To some extent, we might consider neural networks as a trade-off between
619 pure nonparametric methods and traditional regressions.

620 A number of successful applications contributed to the great spread of the
621 neural network concept, including some relevant financial topics (see [16])
622 such as bankruptcy prediction, asset allocation, fraud detection, and cus-
623 tomer relationship management.

624 An example of neural network is in Figure 6. Actually, each neural net-
625 work has its own structure based on *neurons*, but all of them share the same
626 fundamental features:

- 627 • one *input layer* consisting of a number of neurons - one for each pre-
628 dictor
- 629 • one or more *hidden layers*, each of them consisting of its own neurons
- 630 • one *output layer* consisting of a number of neurons which returns the
631 predictions.

632 In particular, notice that the output layer consists of one neuron only if the
633 target variable is binary (i.e., the neural network predicts one probability
634 only) or numerical (i.e., the neural network predicts one numerical value
635 only).

636 Another important remark regards the number of hidden layers. One issue

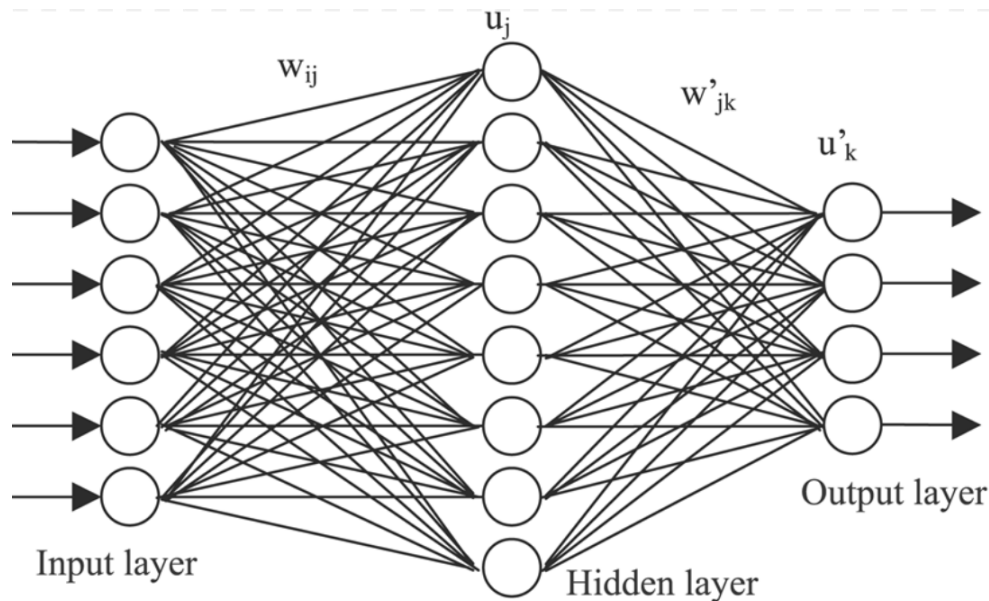


Fig. 6: Example of neural network

637 within this subject on which there is a consensus is the performance difference
 638 from adding additional hidden layers: the situations in which performance
 639 improves with a second (or third, etc.) hidden layer are very few. One hidden
 640 layer is sufficient for the large majority of problems, as stated in [14] as well:

641 *The most popular choice for the number of hidden layers is one.*
 642 *A single hidden layer is usually sufficient to capture even very*
 643 *complex relationships between the predictors.*

644 Therefore, we will solely consider one-hidden-layer neural networks, just like
 645 the example in Figure 6.

646 Furthermore, each input neuron is connected to each hidden neuron, and
 647 each hidden neuron is connected to each output neuron.

648 Neural network training is based on the calibration of the following parame-
 649 ters (see Figure 6):

- 650 • the weight parameters w_{ij} , one for each connection from the input layer
 651 to the hidden layer
- 652 • the bias parameters u_j , one for each hidden neuron

- 653 • the weight parameters w'_{jk} , one for each connection from the hidden
654 layer to the output layer
- 655 • the bias parameters u'_k , one for each output neuron.

656 The input layer, which knows the raw data of the predictors x_i , communicate
657 it to the hidden layer. Such an information, however, is weighted by the
658 connection itself, that is, by the related weight parameter. In other words,
659 the j^{th} hidden neuron receives the information $w_{ij}x_i$. Additionally, given
660 that it receives data from each and every input neuron (i.e., predictor) at
661 the same time, it aggregates information through some *hidden activation*
662 *function* $f(\mathbf{x}, \mathbf{w}_j, u_j)$ of the predictor values, the weights, and the bias. The
663 most popular hidden activation function is the traditional weighted average:

$$H_j := f(\mathbf{x}, \mathbf{w}_j, u_j) = u_j + \sum_i w_{ij}x_i, \quad \forall j. \quad (44)$$

664 Likewise, the output neurons receive weighted information from each hidden
665 neuron. In particular, the k^{th} output neuron receives the value $w'_{jk}H_j$ from
666 the j^{th} hidden neuron. In fact, the k^{th} output neuron receives information
667 from each and every hidden neuron, so it will manipulate it too. More specif-
668 ically, the output neurons use an *output activation function* $g(\mathbf{H}, \mathbf{w}'_k, u'_k)$ of
669 their own parameters. Of course, g is generally different to f , for instance, g
670 is often defined as a linear transformation of the logistic function:

$$\begin{aligned} O_k &:= g(\mathbf{H}, \mathbf{w}'_k, u'_k) = \text{logit}\left(u'_k + \sum_j w'_{jk}H_j\right) = \\ &= \frac{1}{1 + e^{-u'_k - \sum_j w'_{jk}H_j}}, \quad \forall k. \end{aligned} \quad (45)$$

671 Another common output activation function in artificial neural network is
672 the hyperbolic tangent function:

$$\begin{aligned} O_k &:= g(\mathbf{H}, \mathbf{w}'_k, u'_k) = \tanh\left(u'_k + \sum_j w'_{jk}H_j\right) = \\ &= \frac{e^{u'_k + \sum_j w'_{jk}H_j} - e^{-u'_k - \sum_j w'_{jk}H_j}}{e^{u'_k + \sum_j w'_{jk}H_j} + e^{-u'_k - \sum_j w'_{jk}H_j}}, \quad \forall k. \end{aligned} \quad (46)$$

673 Several other options are still possible, but O_k always represents the predic-
674 tion provided by the k^{th} output neuron. If the target variable is categorical,

675 O_k equals the probability of the k^{th} category related to a specific record.
 676 There is still an interesting observation to do. Assume that the target vari-
 677 able is binary, that is, the neural network gets only one output neuron that
 678 predicts the probability of “success”, whereas the number of hidden neuron
 679 is the same as the number of predictors. Moreover, assume the following
 680 parameters within the hidden layer:

$$u_i = 0, \quad w_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}, \quad \forall i, j \quad (47)$$

681 which actually means

$$H_j \equiv H_i = x_i, \quad \forall j. \quad (48)$$

682 Using (45) for the single output neuron, we get

$$\hat{P}(1) = O_1 := g(\mathbf{x}, \mathbf{w}'_k, u'_k) = \text{logit}\left(u'_1 + \sum_i w'_{i1} x_i\right) = \frac{1}{1 + e^{-u'_1 - \sum_i w'_{i1} x_i}} \quad (49)$$

683 which is equivalent to the functional form of a logistic regression’s prediction.
 684 However, we are going to see why it does not mean at all that this peculiar
 685 neural network equals the logistic regression in terms of predicted probabili-
 686 ties.

687 In fact, the main difference between neural networks and regressions lies in
 688 the way parameters are estimated. While regression methods rely on prede-
 689 termined target functions to minimize or maximize, neural networks actually
 690 “learn” from data. Estimating biases and weights is a consequence of such a
 691 learning process, whichever its algorithm is. Anyway, the most popular one
 692 is the so-called *back propagation*.

693 We denoted the prediction from the k^{th} node by O_k . Moreover, let’s define
 694 the prediction error related to the first record as

$$\epsilon_{1k} := O_{1k}(1 - O_{1k})(y_1 - O_{1k}) \quad (50)$$

695 where y_1 equals the actual value of the target variable. Given a global *learn-*
 696 *ing rate* $\lambda \in (0, 1)$ and some initialization values for the parameters of the
 697 network, weights and biases are updated as follows:

$$w_{ij} \longrightarrow w_{ij} + \lambda \epsilon_{1k}, \quad \forall i, j \quad (51)$$

$$u_j \longrightarrow u_j + \lambda \epsilon_{1k}, \quad \forall j \quad (52)$$

$$w'_{jk} \longrightarrow w'_{jk} + \lambda \epsilon_{1k}, \quad \forall j \quad (53)$$

$$u'_k \longrightarrow u'_k + \lambda \epsilon_{1k}. \quad (54)$$

698 After that, the second record goes through the network to get its own esti-
 699 mations O_{2k} , then the error is computed:

$$\epsilon_{2k} := O_{2k}(1 - O_{2k})(y_2 - O_{2k}) \quad (55)$$

700 and the parameters are updated once again starting from those of the previ-
 701 ous step:

$$w_{ij} \longrightarrow w_{ij} + \lambda\epsilon_{2k}, \quad \forall i, j \quad (56)$$

$$u_j \longrightarrow u_j + \lambda\epsilon_{2k}, \quad \forall j \quad (57)$$

$$w'_{jk} \longrightarrow w'_{jk} + \lambda\epsilon_{2k}, \quad \forall j \quad (58)$$

$$u'_k \longrightarrow u'_k + \lambda\epsilon_{2k}. \quad (59)$$

702 The computation is repeated for all the records, all the way through the
 703 training dataset: after the last observation, the first *epoch* is completed.
 704 Generally, a number of epochs is predefined, that is, the records are estimated
 705 several times until some tolerance on the significance of parameter updating
 706 is broken, or some threshold on the training error is finally met.

707 What we have just described is called *case updating*, but it is not the sole
 708 option. For instance, in *batch updating*, the whole training dataset is run
 709 through the network before each updating takes place. As a consequence,
 710 the parameters are updated on the base of the overall training error, or its
 711 average:

$$\bar{\epsilon}_k := \frac{1}{N} \sum_i \epsilon_{ik} \quad (60)$$

712 and

$$w_{ij} \longrightarrow w_{ij} + \lambda\bar{\epsilon}_k, \quad \forall i, j \quad (61)$$

$$u_j \longrightarrow u_j + \lambda\bar{\epsilon}_k, \quad \forall j \quad (62)$$

$$w'_{jk} \longrightarrow w'_{jk} + \lambda\bar{\epsilon}_k, \quad \forall j \quad (63)$$

$$u'_k \longrightarrow u'_k + \lambda\bar{\epsilon}_k. \quad (64)$$

713 In practice, case updating tends to be more accurate than batch updating,
 714 but it is also more computationally expensive, given that the parameters are
 715 updated N times per epoch rather than only once.

716 Neural networks can be very powerful, if their architecture is significant, that
 717 is, the number of hidden layers and hidden neurons is “right” - whatever it
 718 means. There are algorithms that automatically select this features, but

719 none of them seems clearly superior to a simple trial-and-error approach (see
720 [14]). Nonetheless, network architecture depends on the predictors too. Un-
721 fortunately, neural network are rigid in this sense: they cannot really choose
722 among predictors, as opposite to other methods such as stepwise regressions
723 and CARTs. Neural networks always use all the predictors given as input,
724 so they should be chosen very carefully by the data scientist, for instance, by
725 using a proper selection method. Clustering, principal component analysis,
726 and CARTs themselves are all suitable approaches.
727 The various forms of their architecture give neural network a unique flexibil-
728 ity in dealing with data. Potentially, they can recognize any type of pattern.
729 However, the architecture itself is the origin of their major drawback too,
730 that is, their black-box structure. While anyone can easily “read” a tree, or
731 interpret the parameters of a regression, this is generally impossible when
732 dealing with neural networks. Of course, knowing the transfer functions and
733 any parameters, we may write the ultimate function returning predictions,
734 but then we would probably find no meaning in that. Too many parameters
735 and too complex functions are often involved in neural network, and it must
736 be accepted as it is. The usual validation tools can be used to measure the
737 predictive power of a neural network, but unfortunately we cannot rely on
738 model interpretation.

739 **5. An application to automobile bodily injury claim data**

740 As we may have noticed in Section 4, machine learning tools are very
741 flexible, and could potentially improve many of the traditional processes in
742 actuarial practice. Non-life reserving is just one of them.
743 This section describes the path that has been followed to predict the closing
744 delay and claim amount on a publicly available motor insurance dataset.

745 *5.1. Data*

746 Data comes from an R package containing a number of datasets for ac-
747 tuarial and actuarial-affine applications (see [5]). It was also used in [4] as a
748 starting point for generalized regression analysis. As such, we can surely rely
749 on that without further validation. In [5], the dataset is called *Automobile*
750 *bodily injury claim dataset in Australia* (`ausautoBI8999`), with the following
751 description:

752 *This dataset contains information on 22036 settled personal in-*
753 *jury insurance claims in Australia. These claims arose from acci-*
754 *dents occurring from July 1989 through to January 1999. Claims*
755 *settled with zero payment are not included.*

756 Actually, it relates to a small part of the entire claim scope of a motor insur-
757 ance company, but it is probably the most interesting part to us. In effect,
758 bodily injury claims are the most expensive and long-lasting ones, of course
759 the most important to predict. In comparison to traditional triangle-based
760 reserving methods (more suitable for standard claim prediction), machine
761 learning can express its greatest potential for this kind of claims.

762 The dataset includes the following fields:

- 763 • `AccDate` for the claim accident date
- 764 • `ReportDate` for the claim reporting date
- 765 • `FinDate` for the claim closing date date
- 766 • `AccMth` for the claim accident month
- 767 • `ReportMth` for the claim reporting month
- 768 • `FinMth` for the claim closing month
- 769 • `OpTime` for the operational time
- 770 • `InjType1` for the injury severity of the first injured person
- 771 • `InjType2` for the injury severity of the second injured person
- 772 • `InjType3` for the injury severity of the third injured person
- 773 • `InjType4` for the injury severity of the fourth injured person
- 774 • `InjType5` for the injury severity of the fifth injured person
- 775 • `InjNb` for the number of injured persons (max 5)
- 776 • `Legal` for the legal representation (yes/no)
- 777 • `AggClaim` for the aggregate claim amount after closing.

category	score
<i>NA (no injury)</i>	0
<i>minor injury</i>	25
<i>small injury</i>	25
<i>medium injury</i>	50
<i>not recorded</i>	50
<i>high injury</i>	75
<i>severe injury</i>	75
<i>fatal injury</i>	100

Fig. 7: Codification of injury severity

778 The injury severity variables are categorical, but we converted them in nu-
779 merical between 0 and 100, in order to calculate an overall severity score
780 for each claim. The conversion is based on the rules in Figure 7, and the
781 overall score of a single claim is given by the sum of its own scores (no-
782 tice that it will be always positive because each claim caused one injury at
783 least). Therefore, we define the numerical variables `InjScore1`, `InjScore2`,
784 `InjScore3`, `InjScore4`, and `InjScore5`, in correspondence to the original
785 `InjType1`, `InjType2`, `InjType3`, `InjType4`, and `InjType5`. Additionally, we
786 get the overall severity score variable `InjScoreTot`.
787 Moreover, from the variables `AccDate`, `ReportDate`, and `FinDate`, we ex-
788 tract the year - `AccYr`, `ReportYr`, and `FinYr` - and use it to add two more
789 variables:

- 790 • `ReportTime` for the reporting delay defined as `ReportYr - AccYr`
- 791 • `FinTime` for the closing delay defined as `FinYr - ReportYr`.

792 One last additional variable that will be useful in multiple linear regression
793 is `LnAggClaim`, that is, the natural logarithm of the claim amount `AggClaim`.
794 Using `ReportYr` and `FinTime`, let's have a look at the run-off triangle of the
795 entire dataset in Figure 8 for the claim numbers, in Figure 9 for the claim
796 payments, and in Figure 10 for the claim average payments (a dozen of claims
797 relating to the highest amounts have been excluded). Unfortunately, some
798 reporting years and some closing delays include too few observations, as we
799 can easily observe in Figure 8. This is the reason why we will select a proper
800 training/validation dataset as well as a test dataset for the ultimate analysis.

<i>reporting year</i>	<i>closing delay</i>						
	0	1	2	3	4	5	6
1993	1.191	1.970	1.152	558	426	277	49
1994	874	1.531	870	709	443	65	0
1995	718	1.381	1.188	881	134	0	0
1996	529	1.511	1.314	174	0	0	0
1997	585	2.020	283	0	0	0	0
1998	807	352	0	0	0	0	0
1999	14	0	0	0	0	0	0

Fig. 8: Claim numbers

<i>reporting year</i>	<i>closing delay</i>						
	0	1	2	3	4	5	6
1993	26.038.265	77.236.474	74.285.051	58.268.484	48.169.135	38.347.205	8.444.935
1994	16.746.725	42.454.730	42.485.823	55.798.479	44.833.960	7.119.057	0
1995	6.076.308	20.958.156	41.895.020	50.580.334	8.286.925	0	0
1996	4.090.537	21.665.535	46.736.025	8.158.885	0	0	0
1997	3.787.223	33.675.352	8.543.935	0	0	0	0
1998	6.207.963	3.475.759	0	0	0	0	0
1999	36.599	0	0	0	0	0	0

Fig. 9: Claim amounts

801 More specifically, the reporting years 1993-1996 for the sole closing delays
802 from 0 to 3 (in bold in Figure 8, 9 and 10) and will be considered for training
803 and validation, while the reporting years 1997-1998 will be considered for
804 test using all the available closing delays, that is, from 0-2 for 1997 and 0-1
805 for 1998 (in Roman in Figure 8, 9 and 10). All the other observations (in
806 gray in Figure 8, 9 and 10) are not used in the analysis.
807 Figure 8 also demonstrate that there is no regular path by reporting year
808 or closing delay. For instance, the zero-delay claims tend to decrease by
809 reporting year, while the 3-year-delay claims have a rather different path.
810 Comparing reporting years, we see less differences, but still some relevant
811 discordances. For instance, compare the 1-year-delay claim average payment
812 to the 2-year-delay claim average payment in Figure 10: in 1993 the latter is
813 about 50% greater than the former, while in 1996 the latter is about 150%

<i>reporting year</i>	<i>closing delay</i>						
	0	1	2	3	4	5	6
1993	21.863	39.206	64.484	104.424	113.073	138.438	172.346
1994	19.161	27.730	48.834	78.700	101.205	109.524	0
1995	8.463	15.176	35.265	57.412	61.843	0	0
1996	7.733	14.339	35.568	46.890	0	0	0
1997	6.474	16.671	30.191	0	0	0	0
1998	7.693	9.874	0	0	0	0	0
1999	2.614	0	0	0	0	0	0

Fig. 10: Claim average amounts

814 greater than the former. Remember that triangle-based methods often as-
815 sume no differences in path among reporting years, but actually we would
816 lose some important information in this specific case.

817 In order to better include timing information, we consider both the report-
818 ing year and the closing delay as categorical variables. As a consequence,
819 we convert each of them in three binary variables (the forth one would be
820 redundant):

- 821 • `ReportYr` is converted to `ReportYr1` for 1994, `ReportYr2` for 1995, and
822 `ReportYr3` for 1996
- 823 • `FinTime` is converted to `FinTime1` for one year, `FinTime2` for two year,
824 and `FinTime3` for three year.

825 Similarly, we convert `InjNb` to `InjNb2` for two injuries, `InjNb3` for three
826 injuries, `InjNb4` for four injuries, and `InjNb5` for five injuries.

827 Finally, we get the following fields:

- 828 • `InjScoreTot` for the overall severity score
- 829 • `InjNb2` for two injures (1/0)
- 830 • `InjNb3` for three injures (1/0)
- 831 • `InjNb4` for four injures (1/0)
- 832 • `InjNb5` for five injures (1/0)
- 833 • `LegalBin` for the legal representation (1/0)

- 834 • `ReportTime` for the reporting delay
- 835 • `ReportYr1` for the reporting year 1994 (1/0)
- 836 • `ReportYr2` for the reporting year 1995 (1/0)
- 837 • `ReportYr3` for the reporting year 1996 (1/0)
- 838 • `FinTime` for the closing delay
- 839 • `FinTime1` for the 1-year closing delay (1/0)
- 840 • `FinTime2` for the 2-year closing delay (1/0)
- 841 • `FinTime3` for the 3-year closing delay (1/0)
- 842 • `AggClaim` for the aggregate claim amount after closing
- 843 • `LnAggClaim` for the logarithm of `AggClaim`

844 beyond an ID variable to identify the different records.

845 *5.2. Claim closing delay estimation*

846 In our framework, estimating the closing delay means using the predic-
 847 tors `InjScoreTot`, `InjNb2`, `InjNb3`, `InjNb4`, `InjNb5`, `LegalBin`, `ReportTime`,
 848 `ReportYr1`, `ReportYr2`, and `ReportYr3` as inputs for some machine learning
 849 tool to return an estimation of the probability that the claim is definitely
 850 closed after 0, 1, 2 and 3 years after the reporting. In other words, the target
 851 variable is `FinTime`. We use three methods: naive Bayes, k-nearest neigh-
 852 bors, and classification tree.

853 According to the description in Subsection 4.3, the number k of neighbors
 854 should correspond to the lowest validation error. As demonstrated in Figure
 855 14, it occurs when $k = 19$. It is worth noting the peculiar, opposite shape
 856 of the two lines. In Subsection 4.3, we noticed that the training dataset has
 857 actually nothing to gain from k-nearest neighbors prediction: in fact, the
 858 more the neighbors, the greater the error. This happens because the best
 859 prediction in the training dataset always occurs when $k = 1$ by definition.
 860 On the contrary, the validation error slightly decreases as k increases since
 861 the algorithm neglects more and more noisy information. So this is not sur-
 862 prising if the greatest validation error occurs when $k = 1$, just as the training
 863 error is minimum.

864 Likewise, we need to choose a proper classification tree on the base of the
865 validation error as well. In Figure 17, we observe that the lowest validation
866 error occurs after twelve splits (gray dashed line in Figure 17), but, as dis-
867 cussed in Subsection 4.4, this is not necessarily the best choice if we want
868 to take into account the complexity of the tree too. This is the reason why
869 we will use the best-pruned tree for scoring, that is, the tree consisting of
870 three splits (black dashed line in Figure 17). Furthermore, Figure 18, 19,
871 and 20 shows the full tree, the minimum error tree, and the best pruned tree
872 respectively.

873 Remember: we are not really interested in classifying record among one of the
874 four classes, rather we will directly used the estimated probabilities. How-
875 ever, trying to classifying them using the various methods is the easiest way
876 to compare their performances. Therefore, such an assessment is reported
877 in Figure 12-13, 15-16, and 21-22. There we can observe the related confu-
878 sion matrices together with the training and validation errors. All the errors
879 seems to swing around 60%. In other words: if we use the rule that assigns
880 a record the greatest predicted probability among the four possible category,
881 we will correctly predict around 40% closing delays.

882 Unfortunately, this represents a rather poor result, and it can be proved
883 through a straightforward remark. The training data is characterized by the
884 prior (empirical) probabilities per closing delay category as in Figure 11. The
885 most obvious prediction algorithm would classify all the observations in the
886 category that appears most often in the training dataset, that is, category
887 1 with 38,86%. In such a case, we would correctly predict 38,86% claims,
888 that is, an error of 61,14%. But it is just slightly higher than the overall
889 validation errors of the 19-nearest neighbors algorithm (see Figure 16) and
890 the best pruned classification tree (see Figure 22). Curiously, it is even lower
891 than the overall validation error of naive Bayes (see Figure 13). In other
892 word, it does not seem we gain any insight by using machine learning on this
893 dataset to predict closing delay. This is simply due to the low informative
894 value of the data itself. However, we are going to see in Subsection 5.3 the
895 great potential of machine learning in claim amount estimation as compared
896 to traditional regression methods.

closing delay	prior probability
0	20,01%
1	38,86%
2	27,00%
3	14,14%

Fig. 11: Prior probabilities using training data

<i>Actual Class</i>	<i>Predicted Class</i>				<i>Cases Number</i>	<i>Errors Number</i>	<i>Errors Percentage</i>
	0	1	2	3			
0	740	1.069	151	27	1.987	1.247	62,76%
1	870	2.210	689	90	3.859	1.649	42,73%
2	434	1.380	764	103	2.681	1.917	71,50%
3	242	730	281	151	1.404	1.253	89,25%
Total					9.931	6.066	61,08%

Fig. 12: Naive Bayes summary results using training data

<i>Actual Class</i>	<i>Predicted Class</i>				<i>Cases Number</i>	<i>Errors Number</i>	<i>Errors Percentage</i>
	0	1	2	3			
0	477	708	122	18	1.325	848	64,00%
1	536	1.454	490	54	2.534	1.080	42,62%
2	284	988	483	88	1.843	1.360	73,79%
3	153	445	226	94	918	824	89,76%
Total					6.620	4.112	62,11%

Fig. 13: Naive Bayes summary results using validation data

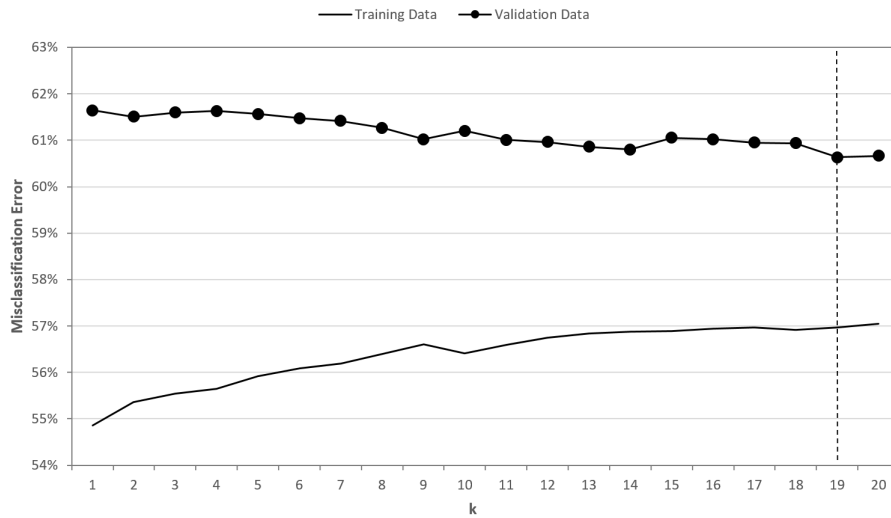


Fig. 14: K-nearest neighbors - training and validation error varying by k

Actual Class	Predicted Class				Cases Number	Errors Number	Errors Percentage
	0	1	2	3			
0	472	1.389	105	21	1.987	1.515	76,25%
1	367	3.140	286	66	3.859	719	18,63%
2	222	1.869	502	88	2.681	2.179	81,28%
3	115	938	191	160	1.404	1.244	88,60%
Total					9.931	5.657	56,96%

Fig. 15: K-nearest neighbors summary results using training data

Actual Class	Predicted Class				Cases Number	Errors Number	Errors Percentage
	0	1	2	3			
0	267	959	83	16	1.325	1.058	79,85%
1	269	1.990	224	51	2.534	544	21,47%
2	164	1.355	248	76	1.843	1.595	86,54%
3	75	614	128	101	918	817	89,00%
Total					6.620	4.014	60,63%

Fig. 16: K-nearest neighbors summary results using validation data

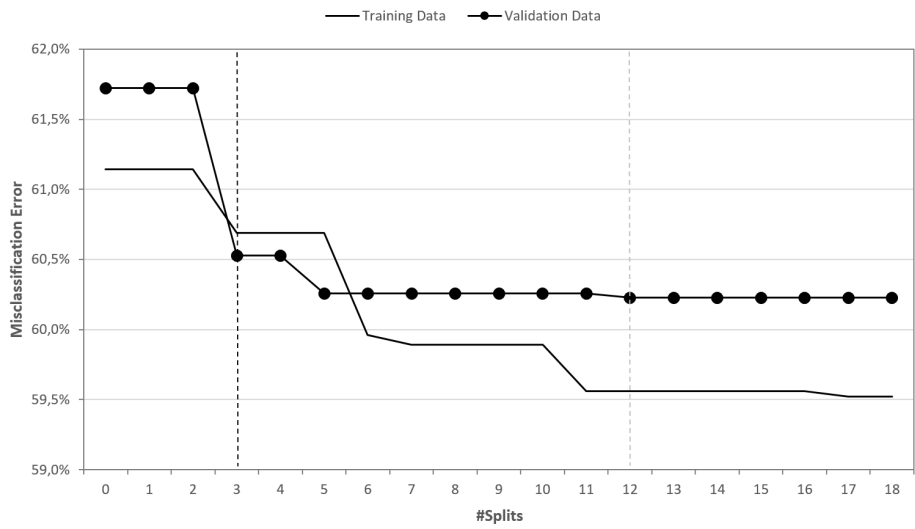


Fig. 17: Classification tree - training and validation error varying by number of splits

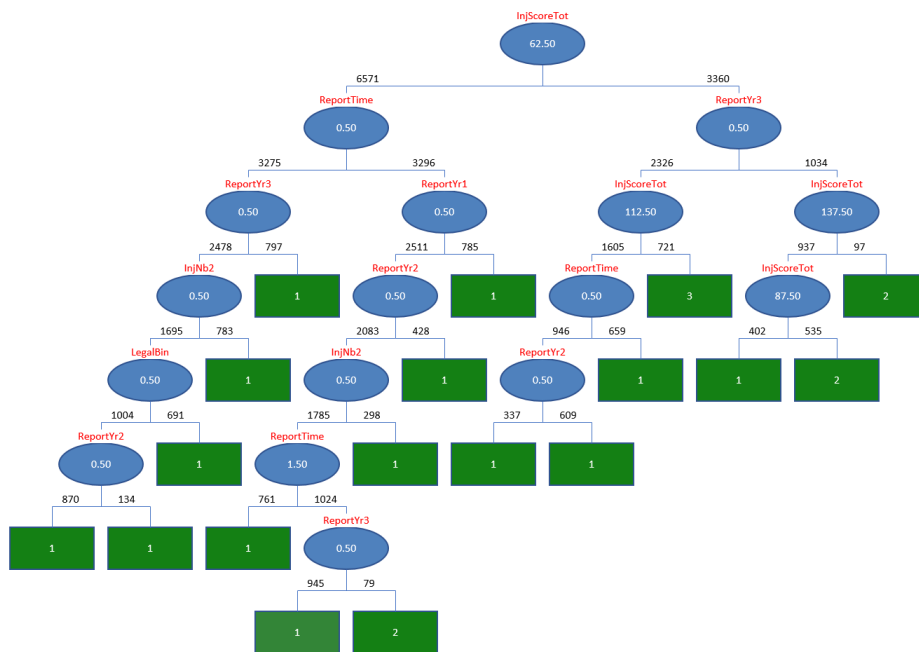


Fig. 18: Full classification tree

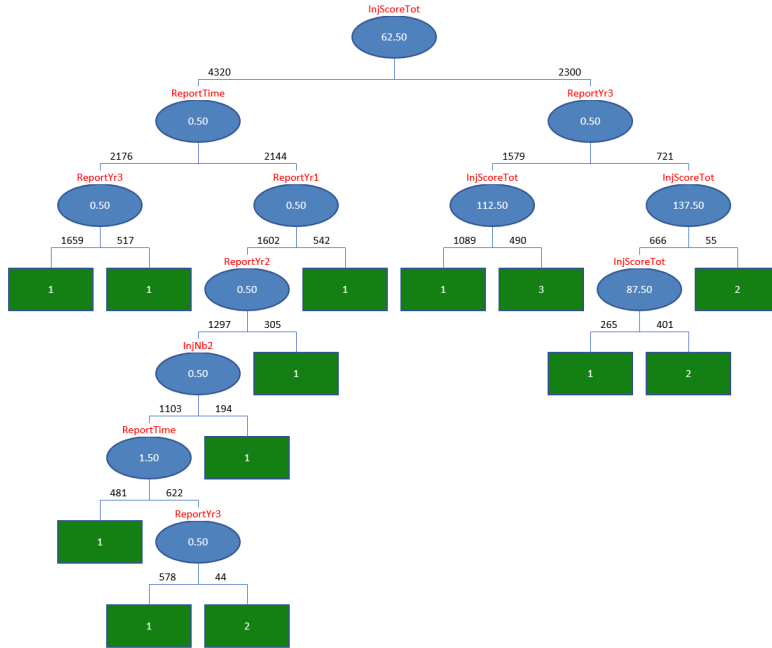


Fig. 19: Minimum error classification tree

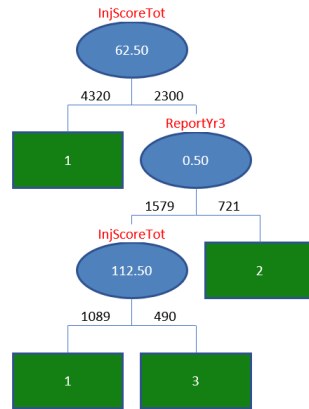


Fig. 20: Best pruned classification tree

<i>Actual Class</i>	<i>Predicted Class</i>				<i>Cases Number</i>	<i>Errors Number</i>	<i>Errors Percentage</i>
	0	1	2	3			
0	0	1.892	57	38	1.987	1.987	100,00%
1	0	3.412	258	189	3.859	447	11,58%
2	0	2.101	347	233	2.681	2.334	87,06%
3	0	1.094	49	261	1.404	1.143	81,41%
Total					9.931	5.911	59,52%

Fig. 21: Classification tree summary results using training data

<i>Actual Class</i>	<i>Predicted Class</i>				<i>Cases Number</i>	<i>Errors Number</i>	<i>Errors Percentage</i>
	0	1	2	3			
0	0	1.229	67	29	1.325	1.325	100,00%
1	0	2.121	294	119	2.534	413	16,30%
2	0	1.363	315	165	1.843	1.528	82,91%
3	0	696	45	177	918	741	80,72%
Total					6.620	4.007	60,53%

Fig. 22: Classification tree summary results using validation data

897 *5.3. Claim payment amount estimation*

898 In our framework, estimating the claim payment amount means using
899 the predictors previously used to estimate the closing delay (see Section 5.2)
900 together with the binary variables `FinTime1` for the 1-year closing delay,
901 `FinTime2` for the 2-year closing delay, and `FinTime3` for the 3-year closing
902 delay. They will represent the inputs for some machine learning tool to return
903 an estimation of `AggClaim`. We use three methods: generalized regression
904 with Gamma distribution, regression tree, and neural network.

905 When it comes with regression methods, we should pay attention to the fea-
906 tures of our predictors. First, it's important to avoid asymmetries and heavy
907 tails. This can be check by having a look at the skewness and kurtosis in
908 Figure 23. Only the target variable `AggClaim` got very high skewness (5,42)
909 and kurtosis (42,21), but this is simply due to the nature of the variable it-
910 self. A simple solution is represented by the natural logarithm of `AggClaim`,
911 `LnAggClaim`, which would permit us to use multiple linear regression (Figure
912 23 also reports skewness and kurtosis of `LnAggClaim`: both of them are very
913 low). However, it would also introduce a significant component of transfor-
914 mation bias (see Subsection 4.1), so we will not try this approach.

915 Another typical issue related to regression models is multicollinearity among
916 numerical variables. In our case, however, if all the correlations among the
917 binary variables of the same categorical variable are excluded, few remain-
918 ing correlations are significant. Therefore, we will run the gamma regression
919 including all the predictors, expecting that the highest correlations will be
920 automatically solved by the stepwise selection algorithm. Figure 24 shows
921 the related results (we also included multiple linear regression results, al-
922 though they will be not used for claim amount prediction). By comparing
923 Figure 23 and Figure 24, it is worth noting that the selected 9-coefficient
924 model includes the predictors characterized by the highest correlations with
925 `LnAggClaim`. Moreover, the exclusion of all the `InjNb` binary variables leads
926 to the exclusion of the most relevant multicollinearities, that is, those among
927 `InjNb` and `InjScoreTot` (see Figure 23). Actually, this is not surprising: to
928 some extent, in effect, the greater the number of injuries, the greater the
929 overall claim severity.

930 When it comes with regression trees, we need to do some further remarks
931 with respect to classification tree. Actually, the former predict numerical
932 variables, while the latter classify records among a range of categories. Given
933 that each claim relates to a different payment, a mere full tree would get as
934 many leaves as the number of records in the training dataset. Obviously,

	lnjScoreTot	lnjNb2	lnjNb3	lnjNb4	lnjNb5	LegalBin	ReportTime	ReportYr1	ReportYr2	ReportYr3	FinTime1	FinTime2	FinTime3	LnAggClaim	AggClaim
mean	56,99	0,19	0,13	0,08	0,09	0,63	0,82	0,24	0,25	0,21	0,39	0,27	0,14	9,59	35.857,34
st. deviation	37,98	0,39	0,34	0,27	0,29	0,48	1,14	0,43	0,43	0,41	0,49	0,45	0,35	1,42	67.073,67
skewness	1,36	1,58	2,15	3,14	2,79	-0,52	1,39	1,21	1,14	1,40	0,47	1,02	2,07	-0,53	5,42
kurtosis	2,16	0,49	2,63	7,86	5,79	-1,73	1,26	-0,53	-0,69	-0,04	-1,78	-0,97	2,29	1,48	42,21

	lnjScoreTot	lnjNb2	lnjNb3	lnjNb4	lnjNb5	LegalBin	ReportTime	ReportYr1	ReportYr2	ReportYr3	FinTime1	FinTime2	FinTime3	LnAggClaim
lnjScoreTot	100,0%													
lnjNb2	-5,4%	100,0%												
lnjNb3	22,3%	-19,0%	100,0%											
lnjNb4	36,9%	-14,1%	-11,5%	100,0%										
lnjNb5	68,6%	-15,6%	-12,6%	-9,4%	100,0%									
LegalBin	9,1%	0,1%	2,8%	4,1%	6,6%	100,0%								
ReportTime	-13,0%	-7,8%	-9,5%	-7,1%	-8,8%	8,3%	100,0%							
ReportYr1	-8,7%	1,0%	-4,2%	-4,4%	-6,2%	-29,4%	2,2%	100,0%						
ReportYr2	17,9%	6,4%	10,0%	8,9%	10,0%	3,2%	-20,3%	-32,7%	100,0%					
ReportYr3	17,8%	4,9%	8,2%	8,2%	11,0%	19,0%	-20,5%	-29,3%	-30,2%	100,0%				
FinTime1	-6,9%	1,7%	-0,1%	-0,6%	-6,7%	-0,1%	-2,1%	-0,2%	-6,5%	4,5%	100,0%			
FinTime2	13,0%	1,6%	4,3%	5,4%	9,2%	5,6%	-6,9%	-6,9%	1,5%	11,6%	-48,7%	100,0%		
FinTime3	14,6%	-0,7%	2,2%	3,2%	11,7%	-1,2%	-3,1%	6,1%	11,9%	-13,6%	-32,0%	-24,8%	100,0%	
LnAggClaim	26,8%	-1,4%	3,9%	7,2%	19,5%	20,6%	27,7%	0,9%	-5,1%	-10,4%	-11,9%	20,7%	28,1%	100,0%

Fig. 23: Descriptive statistics

#Coeffs	R ²	Adj. R ²	Intercept	InjScoreTot	InjNb2	InjNb3	InjNb4	InjNb5	LegalBin	ReportTime	ReportYr1	ReportYr2	ReportYr3	FinTime1	FinTime2	FinTime3
1	0,0%	0,0%	x													x
2	7,9%	7,9%	x													x
3	16,3%	16,3%	x													x
4	25,8%	25,8%	x													x
5	30,8%	30,8%	x													x
6	34,7%	34,7%	x													x
7	36,8%	36,8%	x													x
8	38,0%	38,0%	x													x
9	39,4%	39,3%	x	x	x	x	x	x	x	x	x	x	x	x	x	x

Stepwise Selection	
Coefficient	7,7530
Standard Error	0,0334
t-Statistic	232,26
P-Value	0,0000
Conf. Interval Lower	7,6876
Conf. Interval Upper	7,8185

Multiple Linear Regression									
Coefficient	7,7530	0,0090	0,5396	0,3511	-0,4377	-0,5941	0,7727	1,4149	1,8110
Standard Error	0,0334	0,0003	0,0239	0,0105	0,0294	0,0317	0,0310	0,0338	0,0399
t-Statistic	232,26	27,92	22,61	33,32	-14,90	-18,77	24,94	41,84	45,35
P-Value	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000
Conf. Interval Lower	7,6876	0,0083	0,4928	0,3305	-0,4952	-0,6562	0,7120	1,3486	1,7328
Conf. Interval Upper	7,8185	0,0096	0,5864	0,3718	-0,3801	-0,5321	0,8334	1,4812	1,8893

Gamma Regression									
Coefficient	8,5839	0,0084	0,3353	0,3985	-0,4863	-0,5905	0,5660	1,2402	1,6610
Standard Error	0,0478	0,0005	0,0342	0,0151	0,0421	0,0453	0,0444	0,0484	0,0572
t-Statistic	179,54	18,31	9,81	26,40	-11,56	-13,02	12,76	25,61	29,04
P-Value	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000	0,0000
Conf. Interval Lower	8,4849	0,0075	0,2677	0,3665	-0,5686	-0,6807	0,4782	1,1439	1,5476
Conf. Interval Upper	8,6840	0,0093	0,4027	0,4308	-0,4032	-0,4996	0,6532	1,3362	1,7752

Fig. 24: Regression summary after stepwise selection

935 this is not feasible, so we need an additional rule to stop the growth of the
936 full tree, and the prediction related to a specific leaf will equal the average
937 payment of the training records in that leaf. Generally, the rule is quite em-
938 pirical, for instance, a maximum number of tree levels, a maximum number
939 of nodes, and so on.

940 In our analysis, we choose a minimum number of records in any leaf, and
941 then predictions and errors from that tree are evaluated. The results are
942 summarized in Figure 25, not only training error, validation error, and over-
943 all error, but also a percentage weighted error and a percentage overall error.
944 The percentage weighted error is based on the sixteen percentage errors per
945 reporting year and closing delay weighted on the actual payment amount it-
946 self. Instead, the percentage overall error is simply the percentage difference
947 between the total of the payments in the dataset and the total of the related
948 estimates.

949 First, Figure 25 shows no overfitting, which is a quite important advantage
950 for any predictive model. In other words, training error and validation error
951 are very close regardless of the minimum records per leaf. Secondly, observe
952 the fluctuation of the percentage errors in the last to columns in Figure 25.
953 When very few records are required in each leaf, it seems that a lot of noise
954 affects the tree: errors are quite material, especially the percentage weighted
955 errors, which are greatest between 10 and 40 minimum records per leaf. After
956 all, the regression tree is predicting much better between 50 and 100: it will
957 be probably there where it performs at best. Nonetheless, further increase
958 in minimum records per leaf implies a new increase in error: in fact, if too
959 many records are required into each leaf, the regression tree will no longer
960 be able to detect information in data.

961 All in all, the chosen tree is highlighted in bold in Figure 25, that is, the tree
962 leading to the lowest percentage weighted record 3,67%. Further, although
963 it leads to the lowest average overall error, it does not lead to the lowest
964 percentage overall error (which is however very low).

965 Unfortunately, we cannot report the whole full tree and best pruned tree
966 since they are too big. Anyway, it is worth noting that all the predictors
967 have been used to build the tree, that is, each predictor is the splitting vari-
968 able for one node at least. As briefly mentioned in Subsection 4.5, standard
969 neural networks lack an embedded algorithm to select relevant predictors and
970 exclude irrelevant predictors, but we may refer to the predictors implicitly
971 selected by the regression tree itself. Since it has used all of them, we will
972 run the neural network on the full dataset.

<i>minimum records per leaf</i>	<i>average training error</i>	<i>average validation error</i>	<i>average overall error</i>	<i>percentage weighted error</i>	<i>percentage overall error</i>
10	28.262	28.554	28.379	16,63%	0,38%
20	28.393	28.622	28.485	16,74%	0,38%
30	27.819	28.020	27.899	14,28%	0,39%
40	27.819	28.020	27.899	14,28%	0,39%
50	26.695	27.528	27.028	3,67%	0,29%
60	26.982	27.477	27.180	5,06%	0,22%
70	27.088	27.326	27.183	5,54%	0,10%
80	27.154	27.393	27.250	5,49%	0,09%
90	27.185	27.331	27.243	6,10%	0,01%
100	27.273	27.385	27.318	6,30%	0,15%
110	27.339	27.457	27.386	6,83%	0,16%
120	27.525	27.568	27.542	9,14%	0,16%
130	27.603	27.596	27.600	9,47%	0,10%
140	27.647	27.679	27.660	9,83%	0,25%
150	27.729	27.753	27.739	9,85%	0,41%
160	27.632	27.680	27.651	8,43%	0,37%
170	27.652	27.662	27.656	8,45%	0,40%
180	27.743	27.718	27.733	9,09%	0,39%
190	27.764	27.741	27.755	9,38%	0,43%
200	27.806	27.797	27.802	9,60%	0,51%
210	27.830	27.777	27.809	9,59%	0,48%
220	27.842	27.788	27.820	9,81%	0,50%
230	27.826	27.766	27.802	9,67%	0,48%
240	27.952	27.866	27.918	9,93%	0,54%
250	27.971	27.906	27.945	10,18%	0,55%
260	27.971	27.913	27.948	10,96%	0,56%
270	28.000	28.116	28.046	11,02%	0,67%
280	28.012	28.139	28.063	11,30%	0,72%
290	28.032	28.171	28.088	11,50%	0,69%
300	28.196	28.430	28.289	13,86%	0,63%

Fig. 25: Regression tree performance varying by minimum number of records per leaf

<i>number of hidden neurons</i>	<i>average training error</i>	<i>average validation error</i>	<i>average overall error</i>	<i>percentage weighted error</i>	<i>percentage overall error</i>
1	30.420	30.068	30.279	16,23%	12,50%
2	29.634	29.621	29.629	15,20%	9,17%
3	29.706	29.838	29.758	13,96%	12,09%
4	28.858	28.947	28.894	8,96%	8,57%
5	29.028	29.285	29.131	11,26%	9,27%
6	28.734	28.828	28.771	10,08%	8,52%
7	28.848	29.338	29.044	11,13%	9,80%
8	28.068	28.607	28.284	8,30%	5,57%
9	28.609	29.273	28.875	8,89%	8,28%
10	28.615	29.411	28.933	10,66%	6,89%
11	28.318	29.056	28.613	10,04%	8,81%
12	28.553	29.229	28.824	10,60%	8,02%
13	28.376	28.975	28.616	9,81%	8,33%
14	28.375	29.368	28.772	9,82%	7,10%
15	28.261	28.900	28.516	8,82%	6,44%
16	28.480	29.020	28.696	9,64%	7,19%
17	28.402	29.299	28.761	10,07%	7,77%
18	28.188	28.939	28.488	9,00%	6,90%
19	28.321	28.831	28.525	8,64%	6,55%
20	28.044	28.890	28.382	9,06%	5,76%

Fig. 26: Neural network performance varying by number of hidden neurons

973 So far, we know the structure of the input layer: thirteen neurons for thirteen
974 predictors. As discussed in Subsection 4.5, we may also accept the assump-
975 tion of one single hidden layer. However, how many hidden neurons shall we
976 use? Intuitively, too few hidden neurons will not be able to detect informa-
977 tion, while too many hidden neurons will imply overfitting. Just like for the
978 regression tree, let's try various cases by gradually increasing the number of
979 hidden layers, say from 1 to 20. The results are summarized in Figure 26
980 (the fields have the same meaning as in Figure 25). Once again, training er-
981 ror and validation error are very close: the neural network is not overfitting.
982 The bad performance due to few hidden neurons - between 1 and 7 - is quite
983 clear. By adding more hidden neurons, the error tends to decrease slowly,
984 but it seems that it is not overfitting by 20 neurons yet: it will probably start
985 overfitting a bit later. The chosen neural network is that (in bold in Figure
986 25) leading to the lowest percentage weighted error, which also coincides to

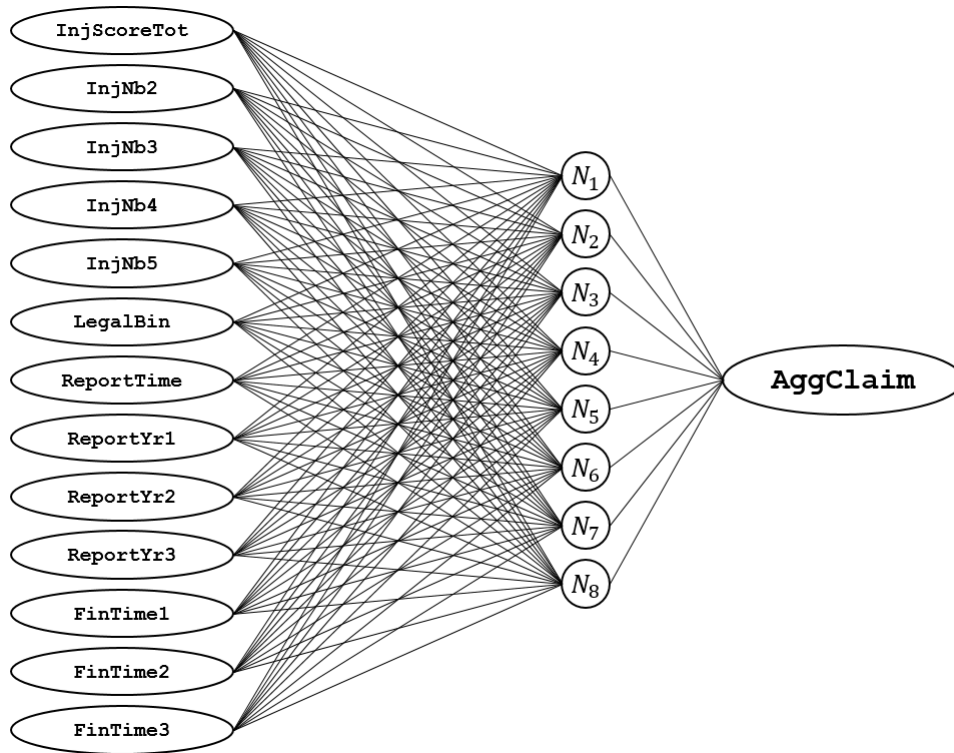


Fig. 27: 8-hidden-neuron neural network for claim amount prediction

987 the lowest average overall error and the lowest percentage overall error. A
 988 formal representation of this neural network is in Figure 27, while the related
 989 parameters are reported in Figure 28.

990 As a conclusion to this subsection, we will present some global results. Firstly,
 991 the convergence plots of the regression tree and the neural network in Fig-
 992 ure 29 and 30 respectively. More importantly, however, we should compare
 993 actual data and estimations - see Figure 31-34. Remember that the chosen
 994 regression tree and neural network got a weighted error of 3,67% and 8,30%
 995 respectively, so their good punctual estimations for the claim payments are
 996 not surprising at all. Likewise, the gamma regression shows good perfor-
 997 mance too, although there seems to be some overestimation by closing delay
 998 2 and 3.

<i>Bias</i>	1,6433	0,9728	-2,6430	3,4999	-0,4701	-0,8636	0,6062	1,2669	
<i>InjScoreTot</i>	-3,6164	1,6972	-1,7925	-7,9758	-1,0843	3,7490	-2,8758	-1,2362	
<i>InjNb2</i>	2,5818	1,4252	2,8352	-1,1509	-0,5093	-0,7649	-0,7377	-0,3781	
<i>InjNb3</i>	-1,7530	-2,0918	1,9940	-2,4436	1,1031	0,3761	1,7829	0,2160	
<i>InjNb4</i>	3,7145	1,6844	-0,2911	1,8491	-0,1809	-1,6364	0,5876	0,3556	
<i>InjNb5</i>	1,5213	-0,3456	2,8264	-1,0558	0,0054	-1,2920	1,6087	-1,5291	
<i>LegalBin</i>	-1,4256	-1,0866	0,2715	-0,2891	-1,0853	-0,0509	0,5342	-2,8530	
<i>ReportTime</i>	-6,4297	-3,9165	-0,2582	0,6026	-0,6748	-2,5407	-1,1021	-2,6102	
<i>ReportYr1</i>	4,7466	3,8977	0,5149	-2,6189	-0,4819	0,2611	2,5912	1,5195	
<i>ReportYr2</i>	-2,8519	-2,7222	0,8436	-1,0376	-0,3715	0,0108	1,1027	-1,6984	
<i>ReportYr3</i>	-1,7670	-2,6942	-0,4083	-2,1818	-1,1539	-0,0667	1,8782	0,1445	
<i>FinTime1</i>	1,0403	1,3477	-0,3577	-0,7600	0,1805	0,1251	0,8446	0,2697	
<i>FinTime2</i>	1,5695	1,8607	-0,6853	-0,7159	0,5422	-0,6917	1,0201	0,9820	
<i>FinTime3</i>	-1,4015	-0,6039	0,5690	-0,4504	-0,3393	-2,9462	-2,0530	1,1273	
	<i>Bias</i>	<i>Neuron 1</i>	<i>Neuron 2</i>	<i>Neuron 3</i>	<i>Neuron 4</i>	<i>Neuron 5</i>	<i>Neuron 6</i>	<i>Neuron 7</i>	<i>Neuron 8</i>
<i>Response</i>	0,3426	-2,8543	4,0754	-1,9584	-2,4236	1,7553	-3,2518	-2,0831	-1,4936

Fig. 28: 8-hidden-neuron neural network parameters

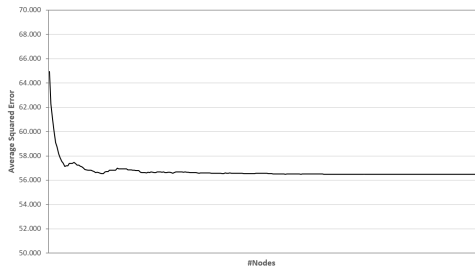


Fig. 29: Convergence of the error for the regression tree

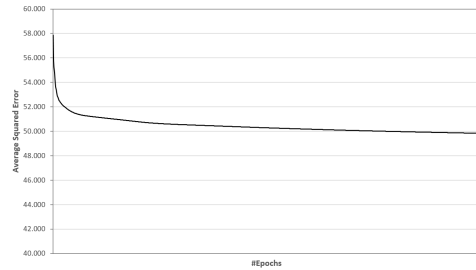


Fig. 30: Convergence of the error for the neural network

reporting year	closing delay			
	0	1	2	3
1993	26.038.265	77.236.474	74.285.051	58.268.484
1994	16.746.725	42.454.730	42.485.823	55.798.479
1995	6.076.308	20.958.156	41.895.020	50.580.334
1996	4.090.537	21.665.535	46.736.025	8.158.885

Fig. 31: Actual data

reporting year	closing delay			
	0	1	2	3
1993	23.942.277	72.096.426	85.372.988	68.135.627
1994	15.615.349	43.607.867	49.123.608	69.977.076
1995	6.294.610	21.590.823	41.199.048	51.380.998
1996	4.338.858	23.040.826	46.178.322	10.158.513

Fig. 32: Gamma regression

reporting year	closing delay			
	0	1	2	3
1993	27.777.729	74.715.426	74.860.543	53.671.653
1994	17.317.427	44.645.051	41.115.866	56.014.481
1995	7.819.221	22.443.511	40.888.916	51.092.375
1996	5.285.040	22.293.701	46.183.852	9.048.303

Fig. 33: Regression tree

reporting year	closing delay			
	0	1	2	3
1993	30.323.454	80.122.598	74.090.121	56.198.161
1994	21.876.621	49.791.426	44.957.040	51.693.100
1995	8.007.124	25.605.931	47.106.288	49.418.199
1996	5.358.636	27.349.566	46.177.586	8.455.383

Fig. 34: Neural network

999 *5.4. Claim reserve estimation as an ensemble*

1000 In machine learning, an *ensemble* is a complex machine learning algorithm
1001 consisting of a number of simpler machine learning tools. The combination
1002 may be very easy to implement (for instance, relating to the three predictive
1003 models in Subsection 5.3, an ensemble for the prediction of the claim amount
1004 could be the average of the three different predictions), or very difficult. In
1005 any case, the goal is the improvement of the predictive performance.
1006 Performance assessment for an ensemble typically requires more computa-
1007 tion than performance assessment for its constituents, so ensembles may be
1008 thought of as a way to compensate for poor learning algorithms by perform-
1009 ing a lot of extra computation. Therefore, fast algorithms such as CARTs
1010 are commonly used in some ensemble versions (for example *random forests*,
1011 *bagging trees*, and *boosting trees*).

1012 From this perspective, the estimation expressed by the (4) may be seen ex-
1013 actly this way, that is, as an ensemble resulting from the combination of
1014 a classification tool and a regression tool. Of course, considering all the
1015 tools described in Section 4, any of the classification methods may be com-
1016 bined with any of the regression methods. The choice will depend on the
1017 performance reported in Figure 35-38 by reporting year, and the overall per-
1018 formance reported in Figure 39.

1019 Globally, gamma regression significantly misestimates the claim payments,
1020 especially for the reporting year 1993 and 1994. More importantly, the ma-
1021 jor problem is that the performance is quite different across reporting years

Actual claim amount 235.828.275	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>	Δ%	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>
<i>random</i>	262.835.516	242.262.782	245.845.509	<i>random</i>	11,45%	2,73%	4,25%
<i>naive Bayes</i>	239.075.514	217.568.447	231.965.733	<i>naive Bayes</i>	1,38%	-7,74%	-1,64%
<i>k-nearest neighbors</i>	249.685.356	229.945.901	238.026.479	<i>k-nearest neighbors</i>	5,88%	-2,49%	0,93%
<i>classification tree</i>	258.083.300	233.422.075	243.846.286	<i>classification tree</i>	9,44%	-1,02%	3,40%

Fig. 35: Predictive performance for reporting year 1993

Actual claim amount 157.485.757	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>	Δ%	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>
<i>random</i>	174.026.704	155.206.357	167.166.597	<i>random</i>	10,50%	-1,45%	6,15%
<i>naive Bayes</i>	177.763.171	149.619.919	165.548.345	<i>naive Bayes</i>	12,88%	-4,99%	5,12%
<i>k-nearest neighbors</i>	182.202.506	159.124.199	169.743.174	<i>k-nearest neighbors</i>	15,69%	1,04%	7,78%
<i>classification tree</i>	182.470.513	156.827.235	170.113.669	<i>classification tree</i>	15,86%	-0,42%	8,02%

Fig. 36: Predictive performance for reporting year 1994

Actual claim amount 119.509.819	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>	Δ%	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>
<i>random</i>	103.147.343	107.212.659	114.257.199	<i>random</i>	-13,69%	-10,29%	-4,40%
<i>naive Bayes</i>	126.742.857	129.321.735	136.983.874	<i>naive Bayes</i>	6,05%	8,21%	14,62%
<i>k-nearest neighbors</i>	117.758.364	120.167.319	128.321.910	<i>k-nearest neighbors</i>	-1,47%	0,55%	7,37%
<i>classification tree</i>	116.673.918	119.593.669	126.714.629	<i>classification tree</i>	-2,37%	0,07%	6,03%

Fig. 37: Predictive performance for reporting year 1995

Actual claim amount 80.650.982	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>	Δ%	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>
<i>random</i>	83.862.742	82.642.642	86.413.193	<i>random</i>	3,98%	2,47%	7,14%
<i>naive Bayes</i>	87.835.102	86.321.719	91.700.340	<i>naive Bayes</i>	8,91%	7,03%	13,70%
<i>k-nearest neighbors</i>	83.617.031	81.500.277	87.200.645	<i>k-nearest neighbors</i>	3,68%	1,05%	8,12%
<i>classification tree</i>	83.901.388	82.788.649	87.006.662	<i>classification tree</i>	4,03%	2,65%	7,88%

Fig. 38: Predictive performance for reporting year 1996

Actual claim amount 593.474.833	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>	Δ%	<i>gamma regression</i>	<i>regression tree</i>	<i>neural network</i>
<i>random</i>	623.872.306	587.324.439	613.682.497	<i>random</i>	5,12%	-1,04%	3,40%
<i>naive Bayes</i>	631.416.644	582.831.820	626.198.292	<i>naive Bayes</i>	6,39%	-1,79%	5,51%
<i>k-nearest neighbors</i>	633.263.257	590.737.696	623.292.207	<i>k-nearest neighbors</i>	6,70%	-0,46%	5,02%
<i>classification tree</i>	641.129.119	592.631.628	627.681.248	<i>classification tree</i>	8,03%	-0,14%	5,76%

Fig. 39: Overall predictive performance

1022 regardless of the classification tool used. On the other hand, the best per-
1023 formances seem to be related to two ensembles:

1024 • k-nearest neighbors and regression tree (overall error -0,46%)

1025 • classification tree and regression tree (overall error -0,14%)

1026 and their error per reporting year is always lower than 3% in absolute value.
1027 Basically, if one of this two combinations is used to predict the payment as
1028 soon as the claim is reported (assuming that all the relevant information is
1029 immediately available), we will reserve an extremely accurate amount for any
1030 reporting year. That's sounds quite good, but it is not necessarily the best
1031 solution. Actually, a proper reserve should always account for some level of
1032 conservatism, as long as it is uniformly included into each allocation. For
1033 instance, Figure 35-38 shows that the two ensembles

1034 • k-nearest neighbors and neural network (overall error 5,02%)

1035 • classification tree and neural network (overall error 5,76%)

1036 slightly overestimate the claim payments. In particular, the overestimation is
1037 quite stable over the reporting years 1994-1996, that is, they overestimate in
1038 the same direction, with the same magnitude, regardless of the reporting year.
1039 If the company reserves using one of the two aforementioned ensembles, it will
1040 most probably overestimate the amount by around 6-8% (we do not consider
1041 the lower but older - thus less significant - errors for the year 1993). For
1042 our purpose, it could really be the best compromise, including a reasonable
1043 prudence margin.

1044 Once the best ensembles are selected, the last step is the estimation of the
1045 reserve on the test dataset, that is, the allocation for the claims reported in
1046 1997 and 1998 (see Subsection 5.1). A further assumption is immediately
1047 necessary: given that the reporting year is a categorical variable, the new
1048 reporting years are not included in our models, so they will be replaced by
1049 the last year available, that is, 1996 - the most significant one in terms of
1050 timing. The results are stored in Figure 40. As specified in the column fields,
1051 remember that these two reporting years contain information for a limited
1052 number of closing delay categories: 0, 1 and 2 for 1997, while 0 and 1 for
1053 1998. This is because the dataset was extracted at year end 1999. Whatever
1054 the ensemble is, the overall claim amounts are materially overestimated for
1055 both of the years. A total 46M claim amount in 1997 is predicted to be about

		<i>Reporting year 1997</i> (closing delay 0-2 years)	<i>Reporting year 1998</i> (closing delay 0-1 years)
	Actual claim amount	46.006.511	9.683.722
<i>k-nearest neighbors & regression tree</i>	Predicted amount	63.174.500	22.272.094
	Delta	17.167.988	12.588.372
<i>classification tree & regression tree</i>	Predicted amount	64.860.982	24.113.611
	Delta	18.854.471	14.429.888
<i>k-nearest neighbors & neural network</i>	Predicted amount	68.002.405	22.774.882
	Delta	21.995.894	13.091.160
<i>classification tree & neural network</i>	Predicted amount	68.079.346	23.632.134
	Delta	22.072.834	13.948.412

Fig. 40: Summary results on the test dataset (reporting years 1997 and 1998)

1056 50% greater, between 63M and 68M. In 1998, this delta reach about 150%!
1057 Actually, this is not surprising, if we take into account the performance of
1058 both k-nearest neighbors and classification tree in predicting closing delay
1059 (see Subsection 5.2). Basically, the ensembles correctly predict amounts,
1060 but they tend to allocate them in wrong closing delay categories. More
1061 specifically, some claims are allocated by lower delays - 0 and 1 - whereas
1062 they should be allocated by higher delays - 2 and 3. Additionally, we do not
1063 have got many actual claims by higher delay categories yet, so this bias is not
1064 balanced out by them, as opposed to the reporting years 1993-1996. However,
1065 this happens in those years as well, if we separate claims by reporting year. In
1066 Figure 41-44), this effect is represented by the gap between the dashed lines
1067 (actual cumulative amount payed) and the four colored lines, representing
1068 the four ensembles.
1069 On the one hand, the low-closing-delay payments tend to be overestimated,
1070 but the overall amount per reporting year still converges to the correct one.
1071 As it's said, keeping in mind that a prudence margin is always required,
1072 such ensembles may be still used for reserving purposes. On the other hand,
1073 it would be legitimate to let the company reduce such a prudence margin.
1074 Ideally, it should own enough data to get better predictions for the closing
1075 delay, in order to reduce the gaps shown in Figure 41-44.

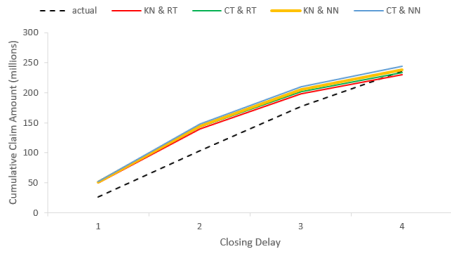


Fig. 41: Cumulative amount in 1993

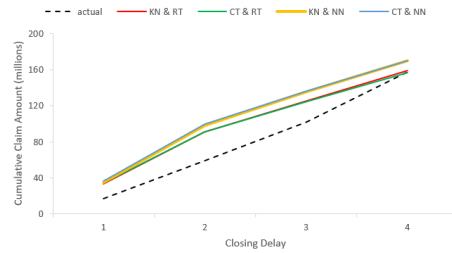


Fig. 42: Cumulative amount in 1994

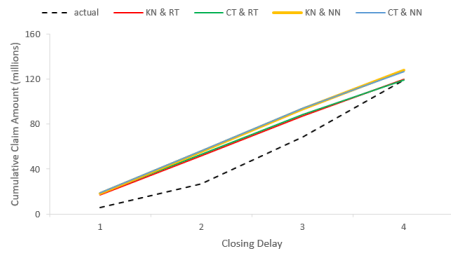


Fig. 43: Cumulative amount in 1995

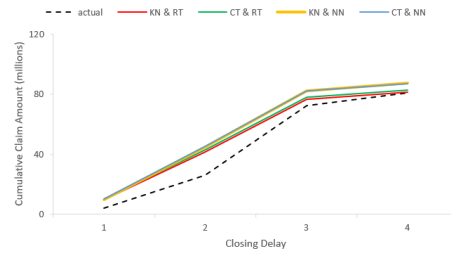


Fig. 44: Cumulative amount in 1996

1076 **6. Limitations, extensions, and conclusions**

1077 In this paper, we demonstrated the great potential of machine learning
 1078 in solving a traditional actuarial problem such as claim reserving in non-life.
 1079 Nonetheless, we focused on a specific application to automobile bodily injury
 1080 claim data. As a consequence, no general conclusion may be drawn (which is
 1081 typical when it comes with machine learning), but of course the ensembles we
 1082 used can be easily adapted to different datasets. And this is possible thank
 1083 to the unique flexibility of these tools.

1084 On the other hand, we should point out some major drawbacks we faced
 1085 during the analysis:

- 1086 • data availability is the crucial constraint, and it may happen that a
 1087 dataset is extremely useful for some target variable, while being very
 1088 poor for others (for instance, compare closing delay performance in
 1089 Subsection 5.2 and payment amount performance in Subsection 5.3);
- 1090 • poor predictions for the closing delay is not a problem if claim reserves
 1091 are evaluated as at REPORTING year, but it would if they were eval-

1092 uated as at PAYMENT year (this is actually another way to state the
1093 problem we faced in Subsection 5.4);

1094 • a complete analysis should include greater closing delays, whose claims
1095 tend to be as rare as severe and expensive, as well as zero-amount
1096 claims, which would reduce the total reserve, but we couldn't because
1097 of data availability issues (see Subsection 5.1);

1098 • IBNYR reserve might be a relevant component of the total claim re-
1099 serve, but it is not considered in this paper (see Section 2 for further
1100 details).

1101 As regards the last bullet point, we should probably build a complete different
1102 model for IBNYR prediction. And it would not be strictly "individual", be-
1103 cause the company does not get any individual claim data before the report-
1104 ing date of the claim itself. If we want to use machine learning techniques for
1105 this purpose, we should rather rely on different data. More specifically, cross-
1106 sectional data related to the policyholder (age, family, address, habits, etc.),
1107 and external information like economic environment (unemployment rate, in-
1108 flation rate, financial distress, etc.), weather conditions, natural catastrophes
1109 (storm, flood, earthquake, etc.), and so on. Of course, such additional data
1110 may be useful for better prediction of the RBNYS reserve as well.

1111 A last remark is important to conclude the paper. Even if data is materially
1112 informative, traditional parametric methods could still outperform nonpara-
1113 metric tools, or return comparable results. For instance, see the payment
1114 amount predictions for closing delays 0 and 1 in Figure 31-34. Actually,
1115 gamma regression predictions seem more accurate than those of regression
1116 tree and neural network. That's just an example of the fact that machine
1117 learning is NOT generally superior as compared to traditional methods. The
1118 strength of machine learning relates to the ability to catch intricate dependen-
1119 cies among data, which is however not always necessary. More importantly,
1120 a greater effort in predicting such dependencies could paradoxically lead to
1121 worse performance on regular data. This is also clear in Figure 31-34: bet-
1122 ter predictions for rare, severe claims (closing delays 3 and 4), but slightly
1123 worse predictions for numerous, regular claims (closing delays 0 and 1). In
1124 fact, if data is actually regular enough to fulfill the regression assumptions
1125 about residuals and multicollinearity, there is no reason to use other tools:
1126 regression would return the best performance by definition.

1127 In spite of it, it is worth noting how these new methods can be convenient for

1128 non-life companies. Although we only used basic machine learning tools and
1129 combined them together, we have still got very accurate predictions per re-
1130 porting year. The process may be improved and automatized, little by little,
1131 until traditional triangle-based models will be completely dismissed. This
1132 evolution is going to outcome two major results: instantaneous, automatic,
1133 accurate reserve estimation, and a brand new field for non-life reserving ac-
1134 tuaries.

1135 **Acknowledgements.** The author wishes to thank his promoters and supervisors
1136 Prof. Susanna Levantesi (University “La Sapienza”, Rome) and Dr. Joseph Lo
1137 (Institute and Faculty of Actuaries, London) for their continuous support and
1138 precious feedbacks since the very first draft of this paper.

1139 **References**

- 1140 [1] K. Antonio, R. Plat, *Micro-level stochastic loss reserving for general insur-*
1141 *ance*, Scandinavian Actuarial Journal 2014/7, 649-669, 2014.
- 1142 [2] R. L. Bornhuetter, R. E. Ferguson, *The actuary and IBNR*, Proceedings of
1143 the Casualty Actuarial Society 59, 181-195, 1972.
- 1144 [3] L. Breiman, J. H. Friedman, R. A. Olshen, C. J. Stone, *Classification and*
1145 *Regression Trees*, Wadsworth Statistics/Probability Series, 1984.
- 1146 [4] P. De Jong, G. Z. Heller, *Generalized linear models for insurance data*, Cam-
1147 bridge University Press, 2008.
- 1148 [5] C. Dutang, A. Charpentier, *Package ‘CASdatasets’*, 4-5, 2016.
- 1149 [6] T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning.*
1150 *Data Mining, Inference, and Prediction*, Springer Series in Statistics, 2009.
- 1151 [7] M. Hiabu, C. Margraff, M. D. Martinez-Miranda, J. P. Nielsen, *The link*
1152 *between classical reserving and granular reserving through double chain ladder*
1153 *and its extensions*, British Actuarial Journal 21/1, 97-116, 2016.
- 1154 [8] G. James, D. Witten, T. Hastie, R. Tibshirani, *An Introduction to Statistical*
1155 *Learning. With Applications in R*, Springer Texts in Statistics, 2015.
- 1156 [9] A. H. Jessen, T. Mikosch, G. Samorodnitsky, *Prediction of outstanding pay-*
1157 *ments in a Poisson cluster model*, Scandinavian Actuarial Journal 2011/3,
1158 214-237, 2011.
- 1159 [10] K. Larsen, *Generalized naive Bayes classifiers*, SIGKDD Explorations, 7(1),
1160 76-81, 2005.
- 1161 [11] T. Mack, *Distribution-free calculation of the standard error of chain ladder*
1162 *reserve estimates*, ASTIN Bulletin 23/2, 213-225, 1993.
- 1163 [12] M. D. Martinez-Miranda, J. P. Nielsen, R. J. Verrall, M. V. Wüthrich, *The*
1164 *link between classical reserving and granular reserving through double chain*
1165 *ladder and its extensions*, Scandinavian Actuarial Journal 2015/5, 383-405,
1166 2015.
- 1167 [13] M. Pigeon, K. Antonio, M. Denuit, *Individual loss reserving with the multi-*
1168 *variate skew normal framework*, ASTIN Bulletin 43/3, 399-428, 2013.

- 1169 [14] G. Shmueli, N. R. Patel, P. C. Bruce, *Data Mining for Business Intelligence*,
1170 Wiley, 160-240, 2010.
- 1171 [15] G. Taylor, G. McGuire, J. Sullivan, *Individual claim loss reserving conditioned*
1172 *by case estimates*, *Annals of Actuarial Science* 3/1-2, 215-256, 2008.
- 1173 [16] R. Trippi, E. Turban, *Neural networks in finance and investing*, McGraw-Hill,
1174 1996.
- 1175 [17] R. J. Verrall, M. V. Wüthrich, *Understanding reporting delay in general in-*
1176 *surance*, *Risks* 4/3, 25, 2016.
- 1177 [18] M. V. Wüthrich, *Machine learning in individual claims reserving*, Swiss Fi-
1178 nance Institute, Research Paper Series 16-67, 2016.
- 1179 [19] M. V. Wüthrich, *Neural networks applied to Chain-Ladder reserving*, SSRN
1180 Manuscript ID 2966126, 2017.
- 1181 [20] M. V. Wüthrich, C. Buser, *Data analytics for non-life insurance pricing*, ETH
1182 Zurich, Lecture Notes, 2016.
- 1183 [21] M. V. Wüthrich, M. Merz, *Stochastic claims reserving manual: advances in*
1184 *dynamic modeling*, SSRN Manuscript ID 2649057, 2015.