

Elsevier Editorial System(tm) for European Journal of Operational Research
Manuscript Draft

Manuscript Number: EJOR-D-06-01563R2

Title: Modeling Churn Using Customer Lifetime Value

Article Type: Regular Paper

Section/Category: Marketing

Keywords: Data Mining, Decision Support Systems, Marketing, Churn Prediction

Corresponding Author: Nicolas Glady, Ph.D. Student

Corresponding Author's Institution: K.U. Leuven

First Author: Nicolas Glady

Order of Authors: Nicolas Glady; Bart Baesens, Ph.D.; Christophe Croux, Ph. D.

Manuscript Region of Origin:

Abstract: The definition and modeling of customer loyalty have been central issues in customer relationship management since many years. Recent papers propose solutions to detect customers that are becoming less loyal, also called churners. The cherner status is then defined as a function of the volume of commercial transactions. In the context of a Belgian retail financial service company, our first contribution is to redefine the notion of customer loyalty by considering it from a customer-centric viewpoint instead of a product-centric one. We hereby use the customer lifetime value (CLV) defined as the discounted value of future marginal earnings, based on the customer's activity. Hence, a cherner is defined as someone whose CLV, thus the related marginal profit, is decreasing. As a second contribution, the loss incurred by the CLV decrease is used to appraise the cost to misclassify a customer by introducing a new loss function. In the empirical study, we compare the accuracy of various classification techniques commonly used in the domain of churn prediction, including two cost-sensitive classifiers. Our final conclusion is that since profit is what really matters in a commercial environment, standard statistical accuracy measures for prediction need to be revised and a more profit oriented focus may be desirable.

EJOR-D-06-01563R2: "Modeling Churn Using Customer Lifetime Value," by Glady, Baesens and Croux

We would like to thank the editor and the reviewers for pointing out the remaining errors. All these have been corrected. We also added two more recent references from the OR literature.

Thanks again for the consideration given to our paper.

Modeling Churn Using Customer Lifetime Value

Nicolas Glady^a Bart Baesens^{a,b,c} Christophe Croux^a *

^a Faculty of Business and Economics, K.U.Leuven, Naamsestraat 69, B-3000 Leuven, Belgium

^b School of Management, University of Southampton, SO17 1BJ, UK

^c Vlerick Leuven Ghent Management School, Reep 1, B-9000, Ghent, Belgium.

Abstract

The definition and modeling of customer loyalty have been central issues in customer relationship management since many years. Recent papers propose solutions to detect customers that are becoming less loyal, also called churners. The churner status is then defined as a function of the volume of commercial transactions. In the context of a Belgian retail financial service company, our first contribution is to redefine the notion of customer loyalty by considering it from a customer-centric viewpoint instead of a product-centric one. We hereby use the customer lifetime value (CLV) defined as the discounted value of future marginal earnings, based on the customer's activity. Hence, a churner is defined as someone whose CLV, thus the related marginal profit, is decreasing. As a second contribution, the loss incurred by the CLV decrease is used to appraise the cost to misclassify a customer by introducing a new loss function. In the empirical study, we compare the accuracy of various classification techniques commonly used in the domain of churn prediction, including two cost-sensitive classifiers. Our final conclusion is that since profit is what really matters in a commercial environment, standard statistical accuracy measures for prediction need to be revised and a more profit oriented focus may be desirable.

Keywords: Churn Prediction, Classification, Customer Lifetime Value, Prediction Models.

*We thank ING Belgium for their support and useful information, especially Martine George head of the customer intelligence department. We also thank three anonymous reviewers for their precious comments. All correspondence should be sent to the first author Nicolas Glady: Naamsestraat 69, B-3000 Leuven; Nicolas.Glady@econ.kuleuven.ac.be

1 Introduction

In a time of cost-cutting and intensive competitive pressure, it becomes of crucial importance for companies to fully exploit their existing customer base. Consequently, customer retention campaigns are implemented. Therefore, when the future duration of the relationship between customers and the company is not known, it is of crucial importance to detect the customers decreasing their loyalty to the company, also called churners. This paper proposes a new framework for the churning detection process, using the earnings a customer brings to the company.

A churning is often defined with respect to the longevity of his/her historical monetary value. However, Reinartz and Kumar (2000) criticize this method and demonstrate that profit and life-cycle are not necessarily related. Rust et al. (2004) emphasize that marketing strategies should focus on projected future financial return, and that customer equity, defined as the total value of the customer base, is of central interest. In order to predict this value, Dwyer (1997) and Berger and Nasr (1998) provide a framework using the lifetime value of a customer. Following this idea, Gupta et al. (2004) show that the profit, and hence the firm's value, is a function of the total *Customer Lifetime Value* (CLV). Venkatesan and Kumar (2004) demonstrate the usefulness of CLV as a metric for customer selection, since "customers who are selected on the basis of their lifetime value provide higher profits in future periods than do customers selected on the basis of several other customer-based metrics". Finally, in a recent paper, Neslin et al. (2006) compare several churn classifiers with regard to the CLV change they incur.

This paper contributes to the existing literature by using the customer lifetime value as a basis concept for the modeling and prediction of churn in a non-contractual setting. That is, when the future duration and the modalities of the relationship between the customers and the focal company is not known. First, in order to define the value of a customer, we define the CLV as the present value of future cash flows yielded by the customer's product usage, without taking into account previously spent costs. Subsequently, to detect churning behavior, we consider Baesens et al. (2003) who proposed solutions to estimate the slope of the customer life-cycle, giving an insight on future spending evolutions. Combining these two ideas, we predict churn on the basis of the slope of the customer lifetime value in time, thereby moving from a product-centric viewpoint to a customer centric one. A churning is then defined as someone with a customer lifetime value decreasing over time.

Consequently, we are able to compute the actual loss caused by a bad prediction (with

no or inefficient action) by defining a new type of profit-sensitive loss function. Our key point is that in any business activity, to lose only a few profitable customers is worse than to lose many non-profitable ones. That is why usual statistical accuracy measures may not be most ideal in this context.

Next, we use and contrast several classifiers for churn prediction. A decision tree and a neural network are compared to a baseline logistic regression model. A cost-sensitive design is provided by Turney (1995) and Fan et al. (1999). These papers provide tools to optimize classifiers using boosting with regard to a cost function. Such algorithms are called meta-classifiers, since they only optimize other “base” classifiers, see Lemmens and Croux (2006) for an example in the churn prediction context. Applying this idea, we implement a decision tree generated on a cost-sensitive training sample, and the classifier AdaCost, a variant proposed by Fan et al. (1999) of the well-known AdaBoost algorithm, which has been described in Freund and Schapire (1997). For the sake of simplicity, the only predictor variables in these models are the RFM (recency, frequency and monetary) type: Buckinx and Van den Poel (2005) and Fader et al. (2005) show that RFM variables are good predictors for the CLV.

In our empirical study, using data provided by a retail banker, the loss function presented is applied to assess various common classification techniques for the detection of churn. The purpose of this paper is not to provide a new way to model the CLV, or a new classification technique, but instead, under some assumptions defined later, to construct a framework using a profit-sensitive loss function for the selection of the best classification techniques with respect to the estimated profit.

Our paper is organized as follows: in Section 2, we discuss the general definition of churn in order to propose a new one using the CLV. Likewise, in Section 3, we discuss the usual loss functions for churn prediction and we provide a new one using the CLV. In Section 4.1, we describe the data set used in Section 4.2 in order to compare in Section 5 usual classification techniques used in churn prediction. In the last section, we discuss the assumptions made and the results obtained. Finally, we propose issues for further research.

2 Definitions of Churn

Churn is a marketing-related term characterizing a consumer who is going from one company to another. As a customer, he still has a relationship with the focal company, but will go to the competitor in the near future. If the company wants to prevent him from leaving, a

1 retention action is required. Modeling churn is only interesting from a retention perspective.
2 The population of interest is therefore the customers that have already been acquired.

3 First, we have to define the condition under which a customer has to be considered as
4 decreasing his/her loyalty, and hence as churning. The issue in a competitive environment
5 is that most people have more than one supplier. For instance, in retail banking, a customer
6 could have a current account in a first bank and a mortgage loan in another. Most people
7 have several current accounts even if they do not use them (so-called “sleeping” accounts).
8 We need to find a definition of a churner applicable to non-contractual products, as opposed
9 to contractual products. Contractual products are for instance insurance, mortgage, cellular
10 phone (if high entry or exit barriers and fixed price), in other words all products with
11 “contractual” cash-flows. On the other hand, non-contractual products could be catalog
12 sales, cellular phones (if low entry and exit barriers and marginal price), etc. In the empirical
13 study, we will focus on the private person checking accounts of a Belgian financial institution.
14 A checking account corresponds to non-contractual products because even if the general
15 relationship is long and contractual, the price for the customer to stop using it is low and
16 the product usage is at the customer’s discretion.
17
18
19
20
21
22
23
24
25
26
27

28 **2.1 Previous Definitions of Churners**

29 Most definitions of churn use the product activity of a customer and a threshold fixed by
30 a business rule. If the activity of the customer has fallen below the threshold, (or equal to
31 zero), this customer is considered as a churner. Van den Poel and Larivière (2004), define
32 a churner as someone who closed all his accounts, i.e. with no activity. Buckinx and Van
33 den Poel (2005) define a partial defector as someone with the frequency of purchases below
34 the average and the ratio of the standard deviation of the interpurchase time to the mean
35 interpurchase time above the average. The retail banker of our retail application defines a
36 churner as a customer with less than 2500 Euros of assets (savings, securities or other kinds
37 of products) at the bank. We claim that this threshold approach is not always relevant and
38 that one should observe the evolution in the customer activity instead.
39
40
41
42
43
44
45
46

47 As an example, consider a business rule labeling all customers with a product activity
48 below 5 transactions per year as churners. If a customer has made 4 transactions in the
49 current year, he/she will be considered as a churner, even though during past years 5 trans-
50 actions were made annually. On the other hand, if another customer had an activity of 100
51 transactions per year for 10 years, but has made 6 transactions only this year, he/she will
52
53
54
55
56
57

not be considered as a churner. This is problematic since it is not sure that the loyalty of the first customer has decreased, whereas the product usage of the second customer has obviously changed. A churner status definition based on a major change in the activity would be more appropriate.

Furthermore, if one has to wait until the customer has ended his/her relationship with the company, it is too late to take any preemptive action. The ultimate purpose is to increase the earnings yielded by the customers, by detecting churning behavior at the very beginning. Moreover, the idea to define a churner for a non-contractual product based on life-cycle duration only, has been challenged by Reinartz and Kumar (2000). Consequently, as noted by Rust et al. (2004), only future earnings (that is what we will later define as the CLV) are relevant to take any potential preemptive action, even though assumptions for the future are obviously made considering the past.

2.2 Churner Status Indicator Based on the Slope of the Product Usage

In a more dynamic approach, Baensens et al. (2003) describe methods to estimate the slope of future spending for long-life customers, hereby providing qualitative information for marketers. Our contribution is to propose a framework to resolve the heterogeneity in the customer population by identifying the more profitable customers such that they can be carefully approached using future actions. Instead of looking at the past to observe whether the customer has churned, we will focus on the future in order to estimate whether the relationship will remain profitable.

Consequently, as a first definition for the churner status, we could consider that if the slope of the product usage in time is below a certain value (let us say 1, when the product usage is decreasing), then the customer should be considered as churning. With $x_{i,j,t}$ being the product j usage, during period t , of customer i , then we define $\alpha_{i,j,t}$ as the slope of the product usage:

$$x_{i,j,t+1} = \alpha_{i,j,t}x_{i,j,t}. \quad (1)$$

The slope of the product usage $\alpha_{i,j,t}$ could then be interpreted as a growth rate for $\alpha_{i,j,t} \gg 1$, a retention rate for $\alpha_{i,j,t} \simeq 1$ and a churning rate for $\alpha_{i,j,t} \ll 1$. The purpose of this paper is to focus on the third case, when the customer is churning. Baensens et al. (2003) defined the indicator function of the churner status $y_{i,j,t}$ for the customer i during period t for product

1 *j* as,

$$2 \quad y_{i,j,t}^{(1)} = I(\alpha_{i,j,t} < 1). \quad (2)$$

3
4 In other words, a customer *i* is then considered as a churner for product *j* during period *t* if
5 his/her product usage will be decreasing in the near future (*t* + 1).
6

7 Although the definition of Baesens et al. (2003) is simple and easy to understand, it is
8 product-centric. The products are considered separately, whereas a customer could have
9 several products. The same customer could then be considered as a churner for one product
10 but loyal for another. On the opposite, according to many authors such as Dwyer (1997),
11 Rust et al. (2004) and Gupta et al. (2004), all marketing campaigns should be customer-
12 centric. The churner status should ideally be defined based on the entire customer activity.
13 That is the issue we will try to address in the next section.
14
15
16
17
18
19

20 **2.3 A New Definition of Churner Using the Customer Lifetime** 21 **Value**

22
23
24
25 Our first goal is to detect the customers with a decreasing loyalty, now defined as those
26 decreasing their future customer lifetime value. Secondly, we need to identify those for
27 which a retention action will be profitable.
28
29
30

31 **2.3.1 Definition of Customer Lifetime Value**

32
33
34 Customer valuation is a major topic since many years and has been discussed by several
35 papers in the customer relationship management literature, see Dwyer (1997), Berger and
36 Nasr (1998), Rust et al. (2004) and Malthouse and Blattberg (2005). Nowadays, one can
37 see a proliferation of valuation methods using concepts such as “Customer Lifetime Value”
38 or “Customer Equity”, for an overview, see Pfeifer et al. (2005). This paper follows Gupta
39 et al. (2004), defining the value of a customer as “the expected sum of discounted future
40 earnings [...] where a customer generates a margin [...] for each period [...].”
41
42
43
44

45 The CLV is a function of all the transactions a customer will make, for the *q* products
46 the company is selling, but it does not take into account cross-individual (word of mouth)
47 effects. Consequently, the customer lifetime value of the customer *i*, for the horizon *h* from
48 the period *t* is the sum of the net cash flows $CF_{i,j,t+k}$, yielded by the transaction on product
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

j , discounted at the rate r (assumed constant)¹ and defined as

$$CLV_{i,t} = \sum_{k=1}^h \sum_{j=1}^q \frac{1}{(1+r)^k} CF_{i,j,t+k}. \quad (3)$$

Since we are focussing on retention and not acquisition, all customers were acquired in the past and only marginal earnings are to be accounted, disregarding acquisition cost and any sunk or fixed costs². Hence, if we denote the marginal profit by unit of product usage for product j as π_j , assumed fixed by product³, we can define the net cash flow, $CF_{i,j,t}$, generated by a product j sold to a customer i during period t as a function of the product usage $x_{i,j,t}$,

$$CF_{i,j,t} = \pi_j x_{i,j,t}. \quad (4)$$

Using (3), this gives the CLV for the customer i at t for the q products,

$$CLV_{i,t} = \sum_{k=1}^h \sum_{j=1}^q \frac{1}{(1+r)^k} \pi_j x_{i,j,t+k}. \quad (5)$$

Note that, if π_j is low for the set of products considered, the company should work on forcing churn or letting it happen by natural attrition. As observed in Reinartz and Kumar (2000), the CLV could be high not only if the product usage remains positive for longer horizons, but also if the product usage $x_{i,j,t}$ itself is high as well. That is our main argument to say that one should focus on profitability instead of longevity only.

2.3.2 Churner Status Indicator Based on Marginal Action Profit

Improving the churner status definition, we could use the decrease of the CLV instead of the slope of the product usage $x_{i,j,t}$ to identify the churners. First, using (1) and (4), we could re-state the product profit (net cash flow) as follows:

$$CF_{i,j,t+1} = \pi_j \alpha_{i,j,t} x_{i,j,t}. \quad (6)$$

Next, we reformulate the present value of future earnings for the customer i during period t for the product j (that is the part of the CLV of the customer i due to the product j),

$$CLV_{i,j,t} = \sum_{k=1}^h \frac{\prod_{v=0}^{k-1} \alpha_{i,j,t+v}}{(1+r)^k} \pi_j x_{i,j,t}, \quad (7)$$

¹For simplicity purposes, we will consider the discount as if all cash flows were obtained end-of-month.

²In our empirical application, the marginal profit considered is nearly equal to the transaction price paid by the customer, since the marginal costs of the transactions are negligible.

³It may depend on the type of customer, thus on i . Customers may have preferential conditions according to their status. For simplicity reasons, we will consider an average product margin.

where v is an index accounting for the time. The gain in CLV due to a retention action is an opportunity gain. It is the difference between the CLV, after the retention action⁴, and the CLV without action. We will call it the marginal action profit ($MAP_{i,j,t}$) and it will be denoted as

$$\begin{aligned}
MAP_{i,j,t} &= \Delta CLV_{i,j,t} \\
&= CLV_{i,j,t}(\text{with action}) - CLV_{i,j,t}(\text{without action}) \\
&= \sum_{k=1}^h \frac{1}{(1+r)^k} \pi_j x_{i,j,t} - \sum_{k=1}^h \frac{\prod_{v=0}^{k-1} \alpha_{i,j,t+v}}{(1+r)^k} \pi_j x_{i,j,t}.
\end{aligned} \tag{8}$$

However, equation (8) is not easy to implement. Indeed, we would need to know all the information for h periods in advance in order to have all the $\alpha_{i,j,t+v}$ values, before being able to compute the CLV and knowing whether a customer is a churner or not. Instead, we will consider that $\alpha_{i,j,t}$ is constant during h periods without action⁵. This number of periods h will obviously be finite and constant for convenience purpose. The equation (8) becomes⁶

$$\begin{aligned}
MAP_{i,j,t} &= \sum_{k=1}^h \frac{1}{(1+r)^k} \pi_j x_{i,j,t} - \sum_{k=1}^h \frac{\alpha_{i,j,t}^k}{(1+r)^k} \pi_j x_{i,j,t} \\
&= \pi_j x_{i,j,t} \left(\frac{1}{r} \left(1 - \frac{1}{(1+r)^h} \right) - \frac{\alpha_{i,j}}{1+r-\alpha_{i,j}} \left(1 - \left(\frac{\alpha_{i,j}}{1+r} \right)^h \right) \right).
\end{aligned} \tag{9}$$

We will use this value as a lower bound of the profit for a customer who has in mind to churn but has been stopped from doing so by a retention action. When the customer was

⁴That formula could be modified with any other value than $\alpha = 1$, with the assumption that a customer retention campaign should at least not decrease the *CLV* or even, increase it.

⁵A constant retention rate for customer valuation was also accepted by Gupta et al. (2004). Therefore, for simplification purposes, since we consider a small horizon and under the smoothing conditions described below, we will assume the constant character of $\alpha_{i,j,t}$ in order to have a minimum delay before to be able to assess the model.

⁶In order to have the total present value of the possible future loss for the churning behavior of customer i during period t for product j , one could use the convergence of (9) in h ,

$$\lim_{h \rightarrow \infty} MAP_{i,j,t} = \pi_j x_{i,j,t} \left(\frac{1}{r} - \frac{\alpha_{i,j,t}}{1+r-\alpha_{i,j,t}} \right),$$

and passing from a single product view to a customer view (all products), we have,

$$\lim_{h \rightarrow \infty} MAP_{i,t} = \sum_{j=1}^q \pi_j x_{i,j,t} \left(\frac{1}{r} - \frac{\alpha_{i,j,t}}{1+r-\alpha_{i,j,t}} \right).$$

But since it may be unlikely that α remains constant, this value should be used as an informal indication only.

not intending to churn, the action does not have any effect. Then the lower bound of the marginal action profit is the action effect on the customer cash flows for all the products q ,

$$MAP_{i,t} = \sum_{j=1}^q MAP_{i,j,t}, \quad (10)$$

$$\text{with } \begin{cases} MAP_{i,j,t} = 0 & \text{for } \alpha_{i,j,t} \geq 1 \\ MAP_{i,j,t} = MAP_{i,j,t} & \text{for } \alpha_{i,j,t} < 1. \end{cases} \quad (11)$$

Finally, if our purpose is to have an efficient action and if the marginal action cost (MAC) is assumed fixed but not negligible, we arrive at the following customer-centric churner definition:

$$y_{i,t} = I(MAP_{i,t} > MAC). \quad (12)$$

In other words, a churner is defined as someone for whom a retention action is profitable.

This new indicator function offers three major advantages compared with (2), where a customer is labeled as a churner if his product usage is decreasing. First, churners not worthy of dealing with will be neglected. The second advantage is that it is a cross-product, customer-centric definition of a churner instead of a product-oriented definition. As discussed in Shah et al. (2006), the CLV is a customer-centric concept that should drive a firm's strategy. Even though one could argue that our empirical application defines CLV as a function of the earnings of a single product and is therefore product-centric, we nevertheless claim that our approach is customer-centric since our definition is based on the CLV and consequently easily extendable to many products. Finally, once the parameters (action cost, product profit, etc.) have been defined, this definition is applicable to every type of business.

In reality, it may be difficult to find the exact unitary action marginal cost (MAC), the exact marginal product revenue (π_j) and the exact effect of the action on the product usage (the value of $\alpha_{i,j,t}$ if the action is taken). However, if the scale of these parameters is approximately correct, this valuation gives an insight about the profit of a retention action. Moreover, that will enable us to compare the financial value of various churner detection techniques.

3 Definition of the Loss Function

During the empirical study, several classifiers will be compared. In order to assess the accuracy of each classifier, the loss incurred by wrong predictions needs to be quantified;

1 a loss function needs to be defined. The most common measure of loss (or gain), is the
2 *Percentage of Correctly Classified* (PCC) observations. This measure implicitly assumes
3 equal misclassification costs, which is most often not the case. Moreover, this measure is
4 very sensitive to the class distribution and the choice of the cut-off value used to map the
5 classifier output to classes, as we will see below.
6
7

8 Another well-known classification performance metric is the *Receiver Operating Char-*
9 *acteristic* curve (ROC), described in Egan (1975). A ROC curve is a graphical plot of the
10 sensitivity (percentage of true positive) versus 1-specificity (percentage of false positive), let-
11 ting the classification cut-off vary between its extremes. The AUROC, the *Area Under the*
12 *Receiver Operating Characteristic* curve, is then a summary measure of classification perfor-
13 mance. This second measure provides a better evaluation criterion, since it is independent
14 of any cut-off.
15
16
17
18
19

20 This paper also implements a bias analysis as defined in Kohavi and Wolpert (1996). We
21 will measure the bias, the variance and the noise of the classifiers. The bias can be regarded
22 as a measure of the difference between the actual and predicted distributions of churners and
23 non-churners. The variance expresses the variability of the classifier's predictions regardless
24 its accuracy. The noise measures the variability of the actual classes. For precise definitions
25 we refer to Kohavi and Wolpert (1996).
26
27
28
29

30 Other papers study the classification performances according to a certain cost function.
31 For instance, Drummond and Holte (2006) introduce cost curves for visualizing the perfor-
32 mance of 2-class classifiers over possible misclassification costs. Nevertheless, misclassifica-
33 tions are not always causing the same loss for different individuals. In a business context, a
34 very profitable customer (with a high misclassification cost) has to be monitored very closely,
35 whereas churners that are not yielding any profit (with low misclassification costs) may be
36 less interesting to consider. In the next subsection, we will use the CLV in order to define a
37 new loss function proportional to the decrease in earnings generated by a bad classification.
38
39
40
41
42

43 In what follows, two kinds of errors are distinguished. The first one is the false positive
44 type, when a customer is classified as a churning customer whereas he/she is not decreasing loyalty. In
45 this case, an action is taken that was not necessary. The loss is the action cost, which is
46 assumed to be the same for every customer. The second one is the false negative type, when
47 a churning customer is not detected by the classifier. Here, the loss function is the difference between
48 the earnings generated without action, and the earnings that would have been generated if
49 the customer would have been stopped from churning (i.e. with $\alpha_{i,j,t} = 1$).
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

We define the loss function for a customer i during period t , using (10), as follows⁷

$$L(x_{i,j,t}, \alpha_{i,j,t}, y_{i,t}, \hat{y}_{i,t}) = \begin{cases} 0 & \text{for } y_{i,t} = \hat{y}_{i,t} \\ MAC & \text{for } y_{i,t} = 0 \text{ and } \hat{y}_{i,t} = 1 \\ MAP_{i,t}(x_{i,j,t}, \alpha_{i,j,t}) - MAC & \text{for } y_{i,t} = 1 \text{ and } \hat{y}_{i,t} = 0. \end{cases} \quad (13)$$

Here, the churning status $y_{i,t}$ is defined in (12), and $\hat{y}_{i,t}$ is its prediction using a certain classification method and threshold (see Section 4.3). More profitable customers that are churning will cause a bigger loss (if misclassified) than those who are less profitable.⁸

In order to be able to compare our loss function with the PCC, we first compute the ratio between the losses incurred by the classification model, and the worst case scenario, yielding a number between 0 and 1. The worst case scenario assumes that every customer is misclassified. We denote this ratio as the cumulative loss percentage,

$$L_{tot} = \frac{\sum L(x_{i,j,t}, \alpha_{i,j,t}, y_{i,t}, \hat{y}_{i,t})}{\sum L(x_{i,j,t}, \alpha_{i,j,t}, y_{i,t}, 1 - y_{i,t})}, \quad (14)$$

where the sum is over all indices i , t and j . Finally, we define the cumulative profit percentage as the opposite of the cumulative loss percentage

$$\bar{L}_{tot} = 1 - L_{tot}. \quad (15)$$

Most classifiers are giving a probability to belong to one of the two classes instead of a binary outcome. We need a threshold (or cut-off value, denoted by τ) to distinguish one class from another. Let $p_{i,t}$ be the churning probability estimated by the classifier for customer i during period t . The cut-off value τ is the value between 0 and 1, such that, if $p_{i,t} \geq \tau$, then the customer is classified as a churner. Accordingly, the profit curve (PROC) $f(\tau)$ becomes:

$$f(\tau) = \bar{L}_{tot}(\tau), \quad (16)$$

for $0 < \tau < 1$. We can then define the area under the profit curve (AUPROC) as a profit based measure of classification performance which is independent of the cut-off. This curve may then also be used to set the cut-off in a profit optimal way.

⁷The reader has to keep in mind that we are doing an incremental analysis: what are the incremental consequences on the CLV of a retention action? Similarly, we are assessing a classifier with regard to the change in the CLV it will yield. Given this opportunity cost or opportunity gain approach, we can state that the cost incurred by a good classification is zero.

⁸The reader should not forget that $MAP_{i,t}$, thus the loss function defined in (13), is only a lower bound of the opportunity cost of a misclassification, since it is most likely that the action effect will be more than only preventing the customer from churning, but may also increase product consumption, hence the product profit.

1 To compute the AUPROC, one could use a discrete integration under the curve with an
 2 arbitrary precision parameter pr . Consider the set of $\lfloor \frac{1}{pr} \rfloor$ cut-off values $pr, 2pr, \dots, 1$, then
 3 the approximation of the AUPROC is computed as follows
 4

$$5 \quad 6 \quad 7 \quad 8 \quad 9 \quad AUPROC = pr \sum_{l=1}^{\lfloor \frac{1}{pr} \rfloor} \bar{L}_{tot}(l \cdot pr). \quad (17)$$

10 Defined as in (17), it is obvious that AUPROC does not depend on any cut-off value, similar
 11 to AUROC. On the other hand, whereas the AUROC is only sensitive to the ranking of the
 12 predictions, the AUPROC will also depend on their numerical values.
 13
 14
 15

16 17 **4 The Empirical Study**

18 19 **4.1 Description of the Data Set**

20 We study the current account transactions (number of invoices last month, amount invoiced
 21 last month, number of withdrawals, etc.) provided by a Belgian financial service company
 22 for a sample of $n = 10,000$ customers and $s = 9$ months (from January 2004 till September
 23 2004). The population consists of new, old and sleeping (without any activities since many
 24 months) customers. All transactions are aggregated at the customer level. We consider two
 25 different product usages, the total number of debit transactions and the total amount debited
 26 by month. Credit transactions, for simplification purposes, are not taken into account.
 27
 28
 29
 30
 31
 32
 33
 34

35 Before estimating and assessing the classification models, we separate the sample into
 36 a training set (66% of the observations) to design the classifiers and a test set (33% of the
 37 observations) for the performance assessment. The training set is composed of the product
 38 transactions from January 2004 till June 2004 (6 months). The test set contains the product
 39 transactions for the same customers, but from July 2004 till September 2004 (3 months).
 40
 41
 42
 43
 44

45 46 **4.2 Implementation Details**

47 Since the action profit (9) is very sensitive to the value of $\alpha_{i,j,t}$, we first smooth the values
 48 of both $x_{i,j,t}$ and $\alpha_{i,j,t}$ in order to remove the noise, seasonality, and other instability in the
 49 churning status. Indeed, it could happen that the slope of the product usage goes slightly up
 50 and down from one month to another. Since we are studying the trend of the product usage,
 51 we need to have a smoothed value of this slope. Rearranging (1), we applied a Holt-Winters
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

exponential smoothing scheme as described in Brockwell and Davis (2002). If we denote $\tilde{x}_{i,j,t}$ the smoothed value of $x_{i,j,t}$, and $\tilde{\alpha}_{i,j,t}$ the smoothed value of $\alpha_{i,j,t}$, then

$$\tilde{x}_{i,j,t} = ax_{i,j,t} + (1 - a)(\tilde{x}_{i,j,t-1} + T_{i,j,t-1}), \quad (18)$$

where

$$T_{i,j,t} = b(\tilde{x}_{i,j,t} - \tilde{x}_{i,j,t-1}) + (1 - b)T_{i,j,t-1}, \quad (19)$$

giving

$$\tilde{\alpha}_{i,j,t} = \frac{\tilde{x}_{i,j,t+1}}{\tilde{x}_{i,j,t}}. \quad (20)$$

The smoothing parameters a and b are set at 0.8, as determined using experimental evaluation. Next, each observation is rearranged as follows

$$\mathbf{x}_{i,t} = [\tilde{x}_{i,1,t}, \dots, \tilde{x}_{i,1,t-m}, \dots, \tilde{x}_{i,q,t}, \dots, \tilde{x}_{i,q,t-m}], \quad (21)$$

whereby $\tilde{x}_{i,j,t}$ represents the smoothed value of explanatory variable j for customer i observed during time period t . The maximum number of lags considered is $m = 3$. The vector $\mathbf{x}_{i,t}$ contains then the values of the predictor variables to be used in the classification procedures (to be discussed in Section 4.3). Note that the variables $x_{i,1,t}$ and $x_{i,2,t}$, i.e. the number of debit transactions and the total amount debited in month t for customer i , are function of the frequency and monetary value of the customer. A customer with no activity during a certain period will have a product usage of 0 for the related months. Therefore, this provides a recency value. However, this information is only partial, since not taking into account the full transaction history. The vector $\mathbf{x}_{i,t}$ is completely observed for the training and test sample for $i = 1 \dots n$ and $t = 4, 5, 6$. The corresponding $y_{i,t}$ is then computed according to (12). Note that for the models assessment on the test set, the smoothed values of the $\mathbf{x}_{i,t}$ are in-sample, and the values of the $y_{i,t}$ are out-of-sample. In the following, we denote an observation i as a couple (\mathbf{x}_i, y_i) , with $i = 1 \dots N$ for the training set, dropping the dependency on time. Note that $N = 3n = 30,000$, yielding a very huge training sample size. During the parameter estimation process and the models comparison, we discard the most extreme percentile of the customer base, i.e. customers with the 1% largest value of $x_{i,j,t}$. The results are therefore more robust. Moreover these ‘‘high spending’’ customers are closely followed by branch agents and a global model may be not appropriate in this matter.

For the computation of the CLV, the product yield considered is directly proportional to the transaction volume (product usage 1), $\pi_1 = 0.1\%$. There is no fixed contribution by transaction (product usage 2), $\pi_2 = 0\%$. These values are fixed by a business rule and

correspond to the real profit per transaction on average. These values can slightly differ from the real ones at the individual level; what really matters is the relative CLV changes and not the numerical values. The discount rate applied is the weighted average cost of capital disclosed in the 2004 financial statement of the financial service provider, $r = 8.92\%$ yearly, giving a monthly discount rate of 0.7146% .

In order to compare short-term and long-term CLV, the study is made for two distinct values of the time horizon (h). The first measures are made by quarter, $h = 3$. The longer-term view is computed for a semester, $h = 6$. Finally, the churning status is defined using (12) with marginal action cost (MAC) fixed at 2 EUR, which is our best guess for an upper bound of the marginal average cost of a mailing retention campaign.

We denote the AUPROC computed in (17) as $AUPROC_3$ for the quarterly view and $AUPROC_6$ for the semester view. These values have to be compared with the non cost-sensitive AUROC values. We denote $\bar{L}_3 = 1 - L_3$ the cumulative profit percentage for the quarterly view and $\bar{L}_6 = 1 - L_6$ the cumulative profit percentage for the semester view, see (14) and (15). Both measures are compared with the non cost-sensitive percentage of correctly classified observations (PCC). These performance measures are computed over the test set, where the indices in (14) range from $i = 1, \dots, n$, $j = 1, 2$ and $t = 7, 8$. Note that we cannot include the last month, $t = 9$, in the test set since $\alpha_{i,j,t}$ is not computable for it. This yields $2n = 20,000$ observations (\mathbf{x}_i, y_i) in the test sample. Such a large testing sample size guarantees precise estimation of the performance measures.

4.3 Description of the Classifiers

4.3.1 Logistic Regression, Decision Trees and Neural Networks

The first classifiers applied are a selection of well-known data mining algorithms: a logistic regression, a decision tree and a neural network. Note that Smith et al. (2000) used the same three classifiers in a customer retention problem in the insurance industry.

The first classifier, a logistic regression, is a standard statistical binary regression model, a reference is Agresti (2002). Decision trees are recursive partitioning algorithms, which are estimated using e.g. information theoretic concepts so as to arrive at a comprehensible tree-based decision model, that is evaluated in a top-down way as discussed in Quinlan (1992). A *Multi-Layer Perceptron* (MLP) neural network is a non-linear predictive model whereby inputs are transformed to outputs by using weights, bias terms, and activation functions. These last two models have been included in our study, because non-linear relationships were

found in Fader et al. (2005) between CLV and RFM explanatory variables. The software used for the implementation was Matlab 7.4 using the PRtools toolbox of Duin et al. (2007).

4.3.2 Description of the Cost-Sensitive Classifiers

AdaCost

This paper implements a version of AdaCost algorithm as proposed by Fan et al. (1999). Other cost-sensitive approaches could also be applied, as discussed in Viaene and Dedene (2005). AdaCost is basically an extension of AdaBoost, giving better performance with regard to the cumulative loss percentage (14). It selects repeatedly a random sample (bootstrap) of the original training set, each time estimating a classifier, $h(\mathbf{x}_i)$. Whereas for AdaBoost the probability of selection is higher for misclassified observations, see Freund and Schapire (1997), in AdaCost the probability for an observation i to be selected in the bootstrap is proportional to its misclassification cost, c_i , here defined as

$$c_i = \frac{L(\mathbf{x}_i, y_i)}{\sum_{i=1}^N L(\mathbf{x}_i, y_i)}, \quad (22)$$

where \mathbf{x}_i has been defined in (21) and, $L(\mathbf{x}_i, y_i) = L(\tilde{x}_{i,j,t}, \tilde{\alpha}_{i,j,t}, y_{i,t}, 1 - y_{i,t})$ as defined in (13). The difference between AdaCost and AdaBoost then lies in the probability of selection of an observation in each iteration. For AdaCost, this probability is a function of the misclassification cost, and for AdaBoost, it is a function of the binary classification status. The algorithm is outlined in Figure 1. We used decision trees as base classifiers $h(\mathbf{x}_i)$, because the aggregation of decision trees has been reported in Neslin et al. (2006) and Lemmens and Croux (2006) to be an efficient approach to consider for defection detection.

The choices for w_l , r_l and $\beta_l(i)$ in step 4 are the same as in Fan et al. (1999). The number of iterations in the AdaCost algorithm was the usual number of iterations in the AdaBoost-like algorithm, $L = 50$.

Cost-Sensitive Decision Tree

The last classifier we will study is a special version of AdaCost. If there is only one iteration (without re-weighting), the classifier becomes a decision tree trained on a cost-weighted bootstrap. Such a technique is very fast, straightforward, and more readable, it may be an interesting alternative to consider.

For all classification models, we study the performances by comparing the $AUPROC_h$, AUROC, L_h , the PCC and the percentage of true positives at horizons $h = 3$ and $h = 6$.

- Given the training sample $S = \{(\mathbf{x}_1, y_1, c_1), \dots, (\mathbf{x}_N, y_N, c_N)\}$, with $\mathbf{x}_i \in \mathbb{R}^{m \times q}$, y_i recorded such that $y_i \in \{-1, 1\}$ and $c_i > 0$
- Initialize $c_1(i) = c_i$, according to (22) for $1 \leq i \leq n$
- For $l = 1 \dots L$
 1. Create bootstrap sample B_l using bootstrap weights $c_l(i)$.
 2. Train base learner h_l for bootstrap sample B_l .
 3. Compute the classifier $h_l: \mathbb{R}^{m \times q} \rightarrow [-1, 1]$ on the set S .
 4. Compute $w_l = \frac{1}{2} \times \ln\left(\frac{1+r}{1-r}\right)$ where $r_l = \sum_{i=1}^N c_l(i) h_l(x_i) \beta_l(i) y_i$, and $\beta_l(i) = 0.5 + 0.5 \times c_l(i)$ for misclassified observations and $\beta_l(i) = 0.5 - 0.5 \times c_l(i)$ for correctly classified observations.
 5. Update the costs according to $c_{l+1}(i) = c_l(i) \exp(-w_l h_l(\mathbf{x}_i) \beta_l(i) y_i)$ and rescale them such that they sum to one.
- Output the final AdaCost classifier $\hat{f}(x_i) = \sum_l^L w_l \times h_l(\mathbf{x}_i)$

Figure 1: General AdaCost algorithm.

In order to assess the sensitivity to the cut-off, we consider three different cut-off values for the horizon $h = 3$. The first value is the naive one, $\tau = 0.5$. The second cut-off considered maximizes the PCC metric on the training set. The last one maximizes the cumulative profit percentage, \bar{L}_3 , on the training set.

5 Empirical Results

In this section, we describe our empirical results. First, some descriptive statistics are presented, showing that churners are substantially more expensive to misclassify than non-churners. Next, the accuracy of the classifiers previously described is discussed. Two points are made. First, the new loss function provides different results than the standard measures of accuracy. Secondly, cost-sensitive classifiers are presented as an interesting alternative to the usual techniques.

5.1 Frequency of Churners

The churners and non-churners, defined according to (12), are distributed as indicated in Tables 1 and 2. The first line contains statistics for the total data set (training set and test set) and the second line only for the test set. In the first two columns, one can see the relative frequencies of non-churners and churners, assuming each observation has the same weight. The next two columns contain relative frequencies expressed in a cost-weighted way. For non-churners this is

$$\frac{\sum_i I(y_i = 0)c_i}{\sum_i c_i}, \quad (23)$$

and for churners

$$\frac{\sum_i I(y_i = 1)c_i}{\sum_i c_i}. \quad (24)$$

Obviously, to misclassify a churning is, on average, far more expensive than to misclassify a non-churner. For a longer horizon ($h = 6$, see Table 2), we have evidently more churners. For a longer period of CLV computation, the retention action profit increases and thus, is more likely to be greater than the action cost.

The reader has to keep in mind that the reported frequencies depend on the product yield π_j and the marginal action cost. First, all other parameters being equal, the larger the marginal action cost, the less it is cost-effective to target the customers with only moderate churning behavior ($\alpha_{i,j,t}$ close to 1). On the contrary, if the product yield was greater, these customers would be considered as worthy to start an action.

Table 1: Frequency of churners and non-churners, for $h = 3$

Data Set	Relative frequency		Cost-adjusted frequency		Total Number
	Non-Churners	Churners	Non-Churners	Churners	
Total	87.45%	12.55%	43.90%	56.10%	49500
Test Set	86.72%	13.28%	35.06%	64.94%	19800

Table 2: Frequency of churners and non-churners, for $h = 6$

Data Set	Relative frequency		Cost-adjusted frequency		Total Number
	Non-Churners	Churners	Non-Churners	Churners	
Total	78.36%	21.64%	26.73%	73.27%	49500
Test Set	77.10%	22.90%	25.63%	74.37%	19800

From Tables 1 and 2, one could observe that there are proportionally less churners in the total data set than in the test set. This is due to the way the data sets have been constructed. In the long-run, everybody dies, or, in our case, churns. Since the test set consisted of customers sampled during the first month and observed six months later, churning behavior is of course going to increase when customers are observed in later time periods.

5.2 Comparison of Classifiers

The classification results on the test set of the various techniques are depicted in Tables 3 and 4, for $h = 3$ and 6, respectively. Five classifiers are compared: a logistic regression, a multi-layer perceptron neural network, a decision tree, a cost-sensitive decision tree and the AdaCost boosting method previously described. Their performance is measured by the area under the profit curve, AUPROC, defined in (17) and the area under the receiver operating curve (AUROC). We also assess the classifiers by computing, for a cut-off value $\tau = 0.5$, the cumulative profit percentage \bar{L}_h , the percentage of correct classifications (PCC) and the percentage of churners predicted as churners, also called the true positives. Table 5 reports these three measures of performance (\bar{L}_h , the PCC and the percentage of true positives) for $h = 3$ and two different values of the cut-off. The first cut-off value considered maximizes the PCC on the training set, the second one maximizes \bar{L}_3 , also on the training set. During our investigation, we tested the significance of the differences between these

Table 3: Performance of classifiers with $h = 3$, as measured by the cumulative profit percentage \bar{L}_3 , and the area under the profit curve $AUPROC_3$, together with the percentage of correctly classified observations (PCC), the AUROC, and the percentage of true positives.

Models	$AUPROC_3$	AUROC	\bar{L}_3	PCC	True Pos.
Logistic Regression	64.94%	95.31%	64.62%	90.84%	41.52%
Neural Network	75.82%	96.39%	96.12%	92.53%	57.26%
Decision Tree	82.11%	94.94%	95.35%	91.35%	61.86%
AdaCost	95.81%	95.39%	96.42%	91.09%	84.30%
Cost-Sensitive Tree	94.56%	95.19%	95.71%	91.28%	68.21%

Table 4: As in Table 3, but now for $h = 6$.

Models	$AUPROC_6$	AUROC	\bar{L}_6	PCC	True Pos.
Logistic Regression	85.04%	91.91%	87.91%	83.38%	38.91%
Neural Network	78.60%	93.52%	84.33%	86.01%	51.70%
Decision Tree	82.72%	91.70%	89.60%	84.80%	59.64%
AdaCost	91.94%	90.01%	93.89%	81.22%	94.88%
Cost-Sensitive Tree	93.42%	91.66%	93.76%	82.25%	88.07%

results. For the traditional measures of accuracy (AUROC and PCC), all the differences were significant. Unfortunately, we did not have a rigorous test for the significance of the cost-sensitive measures of accuracy, $AUPROC_h$ and L_h , at our disposal. We therefore were not able to test their significance, but we claim that the differences and the number of observations are large enough for the results to be considered significant.

The profit curves, being defined in (16), are plotted in Figure 2, for the logistic regression, the neural network, the decision tree, the AdaCost classifier and the cost-sensitive tree. The profit curve plots the cumulative profit percentage as a function of the cut-off value being used for classifying the observations as being churners or not. The plots for $h = 3$ are presented, the results for $h = 6$ being similar.

These profit curves are useful in deciding on the optimal cut-off value τ . The cut-off can be set at the maximum of the profit curve, hereby correcting for the asymmetry in the misclassification costs and the class distributions. In practice, however, such an optimal cut-off needs to be determined from the training set. As we see from Table 5, this cut-off leads

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 5: Performance of classifiers with $h = 3$, for two different cut-offs. The first four columns contains the cumulative profit percentage, the PCC and the percentage of true positives for a cut-off τ selected as the best one for the PCC on the training set. The second four columns are for the best cut-off for the cumulative profit percentage on the training set.

Models	Optimal PCC				Optimal \bar{L}_3			
	\bar{L}_3	PCC	True Pos.	τ	\bar{L}_3	PCC	True Pos.	τ
Logistic Regression	65.29%	91.28%	48.71%	0.39	66.12%	91.67%	63.95%	0.23
Neural Network	96.30%	92.58%	58.75%	0.50	96.98%	92.62%	68.67%	0.43
Decision Tree	95.35%	91.35%	61.86%	0.49	95.72%	91.41%	68.40%	0.42
AdaCost	94.24%	85.73%	96.46%	0.13	93.83%	84.77%	97.30%	0.11
Cost-Sensitive Tree	95.67%	91.56%	62.78%	0.66	95.73%	91.57%	65.25%	0.63

Table 6: Bias, Variance and Noise of the classifiers for $h = 3$.

Models	Bias	Variance	Noise
Logistic Regression	0.00	0.60	0.62
Neural Network	0.00	0.60	0.62
Decision Tree	0.00	0.61	0.62
AdaCost	0.00	0.66	0.62
Cost-Sensitive Tree	0.00	0.62	0.62

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

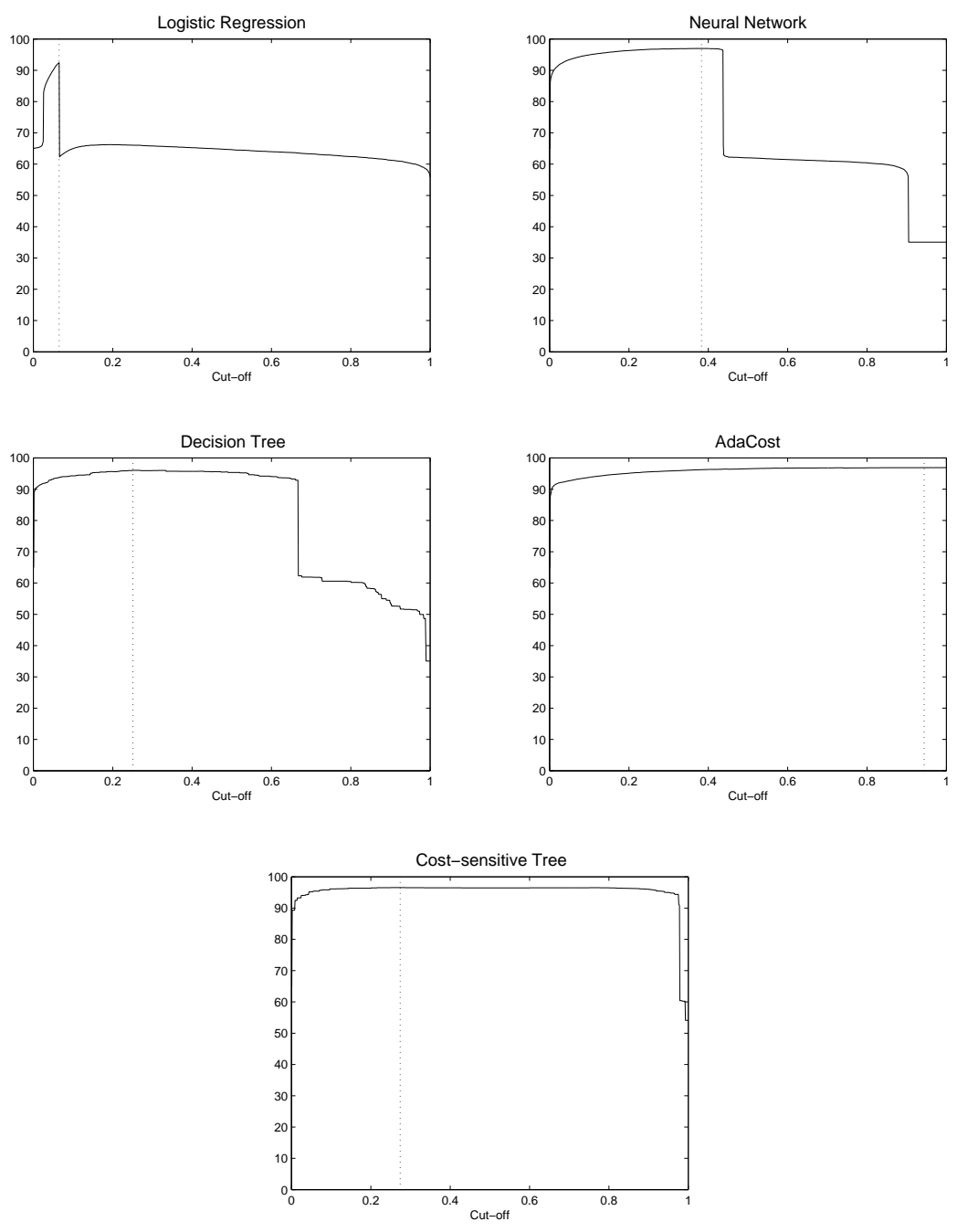


Figure 2: Profit Curves for $h = 3$, for the logistic regression, the neural network, the decision tree, the AdaCost classifier, and the cost-sensitive tree. The dashed line indicates the maximum of the profit curve.

1 to a suboptimal choice on the test set. Note that for AdaCost and the cost-sensitive decision
2 tree, the induced asymmetry is taken into account in the construction of the classifier, hence
3 for these methods we use the standard cut-off value $\tau = 0.5$. For the non cost-sensitive
4 classifiers, since false negatives are more expensive than false positives, all maxima are
5 situated in the left half of the plots. Hence, when setting the cut-off using the profit curve,
6 more customers are classified as churners. Table 5 shows that a sensitive choice of the cut-off
7 can improve the results for the neural network and the decision tree.

8
9
10
11 The area under the profit curve, summarizing the profit curve in a single number, provides
12 an insight regarding the performance of the classifier predictions. The closer the predicted
13 probabilities are to the extremes (0 for assumed perfect non-churners or 1 for assumed perfect
14 churners), the higher will be the value of the area under the profit curve (AUPROC).
15
16
17

18 Finally, we carried out a bias analysis, reported in Table 6. The higher values of the
19 variance for AdaCost and the cost-sensitive tree express that the predicted values are more
20 “extreme” than for the other classifiers, explaining the flat AUPROC curves for these two
21 techniques. If the customers with a high cost of misclassification have a \hat{y}_i close to 0 or 1,
22 the variation of the cut-off will less affect the cumulative loss percentage, resulting in a flat
23 profit curve.
24
25
26
27

28 29 30 **5.3 Discussion**

31
32 From Tables 3 and 4, it follows that the classifiers achieving the best results in our empirical
33 application are the AdaCost classifier and the cost-sensitive tree. They attain the highest
34 values for the AUPROC at both horizons. Since these classifiers directly include cost infor-
35 mation in designing the classification models, it comes as no surprise that both give the best
36 results in terms of profit. The other three classifiers are yielding a lower profit in general. It
37 is interesting to note that the neural network and the decision tree are sensitive to the choice
38 of the cut-off value and that, by selecting this cut-off sensibly (on the left side as shown in
39 Figure 2), these classifiers can achieve very good results.
40
41
42
43
44

45 One can see that it is well possible that two classification methods have similar values
46 for the PCC (or the AUROC), but perform very differently according to the profit-sensitive
47 measures. As a matter of fact, if one would select a classifier on the basis of a standard mea-
48 sure of accuracy (e.g. AUROC), one would choose the neural network. The neural network
49 has however a poor AUPROC value. This difference is mainly explained by the fact that
50 the misclassification cost is, on average, greater for churners than for non-churners. Conse-
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

quently, the total profit for the classifiers that manage to correctly classify the churners (e.g. the cost-sensitive classifiers) is better than those that do not (e.g. the logistic regression). Nevertheless, even though the empirical study shows that, for a selected cut-off value, the proportion of true positives is crucial with respect to the profit generated, one cannot only consider the true positive accuracy. For example, as one can see from Table 5, the AdaCost classifier identifies the highest percentage of churners, whereas the cost-sensitive decision tree still has a higher cumulative profit percentage.

Overall, the cost-sensitive decision tree achieved very good empirical results, in a computationally efficient way. It provides a good trade-off between classifier construction simplicity and profit maximization.

6 Conclusion

In this paper, we provide a framework for evaluating churner classification techniques based on a financial measure of accuracy, i.e. the profit loss incurred by a misclassification, considered from a customer lifetime value perspective. Note that the concept of customer lifetime value, originating from marketing, did recently received attention in the OR literature as well, see Crowder et al. (2007) and Ma et al. (2008). First, using a customer-centric approach, we define a churner as someone whose CLV is decreasing in time. Second, we emphasize the fact that not all customers are equal, neither are all misclassifications. Therefore, we propose a CLV-sensitive loss function and area based measure to evaluate the classifiers. In our empirical application, we use both traditional as well as cost-sensitive classifiers. We show that the cost-sensitive approaches achieve very good results in terms of the defined profit measure, emphasizing the point that, besides achieving a good overall classification, it is important to correctly classify potentially profitable churners.

We can identify different topics for further research. As we have seen, the product usage growth rate α has a large impact on the CLV. In this paper, we assumed α to be constant. It would be interesting to allow varying α and investigate the impact on our findings. Further developments could focus on a more accurate prediction of this value or a more accurate prediction of the CLV. Some of the parameters used in its computation could surely be different in other empirical applications. Also, the model we used to define the CLV has some limitations: we study only non-contractual product types, without taking into account either cross-product effects (cross-selling), or cross-individual effects (word-to-mouth). The new cost-sensitive measures of performance provide a new way to appraise the classifiers, but

1 one could wonder if the differences in performances between the classifiers are statistically
2 significant. A rigorous method of inference to test the significance of these differences would
3 be interesting in this matter.
4
5

6 7 **References** 8

- 9
10 Agresti, A., 2002. Categorical data analysis. Wiley, Hoboken, New Jersey.
11
12 Baesens, B., Verstraeten, G., Van den Poel, D., Egmont-Petersen, M., Van Kenhove, P.,
13 Vanthienen, J., 2003. Bayesian network classifiers for identifying the slope of the customer
14 lifecycle of long-life customers. *European Journal of Operational Research* 156 (2), 508–
15 523.
16
17
18
19
20 Berger, P. D., Nasr, N. I., 1998. Customer lifetime value: Marketing models and applications.
21 *Journal of Interactive Marketing* 12 (1), 17–30.
22
23
24 Brockwell, P. J., Davis, R. A., 2002. Introduction to Time Series and Forecasting. Springer,
25 New York.
26
27
28 Buckinx, W., Van den Poel, D., 2005. Customer base analysis: Partial defection of
29 behaviorally-loyal clients in a non-contractual fmcg retail setting. *European Journal of*
30 *Operational Research* 164 (1), 252–268.
31
32
33
34 Crowder, M., Hand, D. J., Krzanowski, W., 2007. On optimal intervention for customer
35 lifetime value. *European Journal of Operational Research* 183 (3), 1550–1559.
36
37
38 Drummond, C., Holte, R. C., 2006. Cost curves: An improved method for visualizing clas-
39 sifier performance. *Machine Learning* 65 (1), 95–130.
40
41
42
43 Duin, R. P. W., Juszczak, P., Paclik, P., Pekalska, E., de Ridder, D., Tax, D. M. J., 2007.
44 PRTools Version 4.0.23, A Matlab Toolbox for Pattern Recognition. Delft University of
45 Technology.
46
47
48 Dwyer, F. R., 1997. Customer lifetime valuation to support marketing decision making.
49 *Journal of Direct Marketing* 11 (4), 6–13.
50
51
52
53 Egan, J. P., 1975. Signal detection theory and roc analysis. In: *Series in Cognition and*
54 *Perception*. Academic Press, New York.
55
56

- 1
2 Fader, P. S., Hardie, B. G. S., Ka Lok Lee, 2005. RFM and CLV: Using iso-value curves for
3 customer base analysis. *Journal of Marketing Research* 42 (4), 415–430.
- 4 Fan, W., Stolfo, S. J., Zhang, J., Chan, P. K., 1999. Adacost: misclassification cost-sensitive
5 boosting. In: *In Proc. 16th International Conf. on Machine Learning*. Morgan Kaufmann,
6 San Francisco, CA, pp. 97–105.
- 7
8
9 Freund, Y., Schapire, R. E., 1997. A decision-theoretic generalization of on-line learning and
10 an application to boosting. *Journal of Computer and System Sciences* 55 (1), 119–139.
- 11
12
13 Gupta, S., Lehmann, D. R., Stuart, J. A., 2004. Valuing customers. *Journal of Marketing*
14 *Research* 41 (1), 7–18.
- 15
16
17
18 Kohavi, R., Wolpert, D. H., 1996. Bias plus variance decomposition for zero-one loss func-
19 tions. In: Saitta, L. (Ed.), *Machine Learning: Proceedings of the Thirteenth International*
20 *Conference*. Morgan Kaufmann, pp. 275–283.
- 21
22
23 Lemmens, A., Croux, C., 2006. Bagging and boosting classification trees to predict churn.
24 *Journal of Marketing Research* 43 (2), 276–286.
- 25
26
27
28 Ma, M., Li, Z., Chen, J., 2008. Phase-type distribution of customer relationship with marko-
29 vian response and marketing expenditure decision on the customer lifetime value. *Euro-*
30 *pean Journal of Operational Research* 187 (1), 313–326.
- 31
32
33
34 Malthouse, E. C., Blattberg, R. C., 2005. Can we predict customer lifetime value? *Journal*
35 *of Interactive Marketing* 19 (1), 2–16.
- 36
37
38
39 Neslin, S. A., Gupta, S., Kamakura, W., Junxiang, L., Manson, C. H., 2006. Defection de-
40 tection: Measuring and understanding the predictive accuracy of customer churn models.
41 *Journal of Marketing Research* 43 (2), 204–211.
- 42
43
44
45 Pfeifer, P. E., Haskins, M. R., Conroy, R. M., 2005. Customer lifetime value, customer prof-
46 itability, and the treatment of acquisition spending. *Journal of Managerial Issues* 17 (1),
47 11–25.
- 48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1 Reinartz, W. J., Kumar, V., 2000. On the profitability of long-life customers in a non con-
2 tractual setting: An empirical investigation and implications for marketing. *Journal of*
3 *Marketing* 64 (4), 17–35.
4
5
6 Rust, R. T., Lemon, K. N., Zeithaml, V. A., 2004. Return on marketing: Using customer
7 equity to focus marketing strategy. *Journal of Marketing* 68 (1), 109–127.
8
9
10 Shah, D., Rust, R. T., Parasuraman, A., Staelin, R., Day, G. S., 2006. The path to customer
11 centrality. *Journal of Service Research* 9 (2), 113–124.
12
13
14 Smith, K. A., Willis, R. J., Brooks, M., 2000. An analysis of customer retention and insurance
15 claim patterns using data mining: A case study. *Journal of the Operational Research*
16 *Society* 51 (1), 532–541.
17
18
19
20 Turney, P. D., 1995. Cost-sensitive classification: Empirical evaluation of a hybrid genetic
21 decision tree induction algorithm. *Journal of Artificial Intelligence Research* 2, 369–409.
22
23
24 Van den Poel, D., Larivière, B., 2004. Customer attrition analysis for financial services using
25 proportional hazard models. *European Journal of Operational Research* 157 (1), 196–217.
26
27
28 Venkatesan, R., Kumar, V., 2004. A customer lifetime value framework for customer selection
29 and resource allocation strategy. *Journal of Marketing* 68 (4), 106–125.
30
31
32
33 Viaene, S., Dedene, G., 2005. Cost-sensitive learning and decision making revisited. *Euro-*
34 *pean journal of operational research* 166 (1), 212–220.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65