

Gubela, Robin; Bequé, Artem; Gebert, Fabian; Lessmann, Stefan

Working Paper

Conversion uplift in e-commerce: A systematic benchmark of modeling strategies

IRTG 1792 Discussion Paper, No. 2018-062

Provided in Cooperation with:

Humboldt University Berlin, International Research Training Group 1792 "High Dimensional Nonstationary Time Series"

Suggested Citation: Gubela, Robin; Bequé, Artem; Gebert, Fabian; Lessmann, Stefan (2018) : Conversion uplift in e-commerce: A systematic benchmark of modeling strategies, IRTG 1792 Discussion Paper, No. 2018-062, Humboldt-Universität zu Berlin, International Research Training Group 1792 "High Dimensional Nonstationary Time Series", Berlin

This Version is available at:

<https://hdl.handle.net/10419/230773>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



Conversion uplift in e-commerce: A systematic benchmark of modeling strategies

Robin Gubela*
Artem Bequé*
Fabian Gebert*²
Stefan Lessmann*



* Humboldt-Universität zu Berlin, Germany

*² Akanoo GmbH, Germany

This research was supported by the Deutsche
Forschungsgemeinschaft through the
International Research Training Group 1792
"High Dimensional Nonstationary Time Series".

<http://irtg1792.hu-berlin.de>
ISSN 2568-5619

International Research Training Group 1792

Conversion uplift in e-commerce: A systematic benchmark of modeling strategies

Robin Gubela¹, Artem Bequé¹, Fabian Gebert², Stefan Lessmann¹

¹ School of Business and Economics, Humboldt-University of Berlin, Unter-den-Linden 6,
10099 Berlin, Germany

² Technology & Data Science Department, Akanoo GmbH, Mittelweg 121, 20148 Hamburg

Abstract

Uplift modeling combines machine learning and experimental strategies to estimate the differential effect of a treatment on individuals' behavior. The paper considers uplift models in the scope of marketing campaign targeting. Literature on uplift modeling strategies is fragmented across academic disciplines and lacks an overarching empirical comparison. Using data from online retailers, we fill this gap and contribute to literature through consolidating prior work on uplift modeling and systematically comparing the predictive performance and utility of available uplift modeling strategies. Our empirical study includes three experiments in which we examine the interaction between an uplift modeling strategy and the underlying machine learning algorithm to implement the strategy, quantify model performance in terms of business value and demonstrate the advantages of uplift models over response models, which are widely used in marketing. The results facilitate making specific recommendations how to deploy uplift models in e-commerce applications.

Keywords: e-commerce analytics, machine learning, uplift modeling, real-time targeting

1 Introduction

The meteoric rise of electronic commerce and continuous growth in internet adoption leads many business organizations to adopt digital channels for expanding their market presence (Bagchi & Mukhopadhyay, 2006). For example, the Digital Commerce 360¹ report predicts that electronic commerce will be 17 percent of the U.S. retail sales by 2022. However, online marketplaces also increase competition and pose several challenges. For example, lower search cost increase price competition and diminish seller profits (Bakos, 1998). Similarly, negative consumer reviews adversely affect company reputation and may cause financial losses (Lee et al., 2008). To cope with these challenges, companies execute digital marketing strategies to achieve their business objectives and to survive in a challenging business environment.

Digital marketing uses analytic methods to extract relevant insights from massive amounts of data and to drive the company toward growth and profitability. Empirical marketing decision models support all stages of a customer lifecycle including acquisition, development, and retention management (Ascarza et al., 2017; Rhouma & Zaccour, 2018). Digital marketing is most successful when it is personalized and well-targeted (Huang & Tsui, 2016). To target marketing communication, marketers use response models that predict customer behavior (Baecke & Van den Poel, 2010; Chen, 2006) and, in particular, the likelihood of customers to respond to a marketing offer (Coussement et al., 2015). There are many examples of response modeling in digital marketing. Use cases include targeting customers with email-based digital coupons (Sahni et al., 2016), a dynamic adaptation of websites to infer user intentions (Ding et al., 2015), or prediction of the success of social media initiatives (Ballings & Van den Poel, 2015).

In the context of marketing campaign planning, response models suffer a limitation. They fail to distinguish different customer segments (Kondareddy et al., 2016; Michel et al., 2017). To illustrate this, consider a marketing campaign aimed at soliciting digital coupons (Ieva et al., 2018). To efficiently allocate the marketing budget available for the campaign, a marketer wants to offer a coupon only to those customers who do not buy without such price reduction (Zhao & Zhu, 2010). Response models ignore the causal link between the marketing action and customer response (Rzepakowski & Jaroszewicz, 2012b). Instead, they recommend targeting customers with highest likelihood to buy. Such targeting inevitably leads to soliciting customers who would also buy without an incentive and thus wastes marketing resources (Radcliffe, 2007). Uplift models add the element of causality that response models miss. They identify customers who buy because

¹ Data from the official website retrieved on 22.09.2017 (<https://www.digitalcommerce360.com/2017/08/09/e-commerce-grow-17-us-retail-sales-2022/>)

of a marketing action and enable better campaign targeting (Hansotia & Rukstales, 2002b; Jaroszewicz & Rzepakowski, 2014). More specifically, uplift models identify customers who are likely to change their behavior in response to a marketing message (Kane et al., 2014). This is equivalent to modeling the differential (i.e., causal) effect of a marketing incentive on customer behavior.

Several approaches for uplift modeling have appeared in the literature (see De Vriendt et al., 2018 for a recent survey). This paper focuses on uplift modeling strategies that work together with existing algorithms from supervised machine learning. The important advantage of corresponding strategies for corporate practice is that they facilitate predicting uplift and overcome the limitations of response models, while avoiding the need to invest in new technology. Leveraging supervised machine learning, a technology widely available and used in corporate environments (Melli et al., 2012), corresponding uplift modeling strategies are relatively easy to adopt.

As we detail in the review of related literature, previous work on uplift models does not emphasize the advantage of uplift modeling strategies to avoid large upfront investments (for example into new software) through reusing supervised machine learning. Moreover, little attempt has been made to systematically explore their potential for conversion modeling. The need for a comprehensive benchmark emerges because available uplift modeling strategies come from different academic disciplines. Furthermore, the few papers that employ uplift modeling strategies consider only a small set of learning algorithms – typically only one – and do not examine interactions between different classification algorithms and uplift modeling strategies. Consequently, guidance which classifiers work well with which uplift modeling strategy is missing. The goal of the paper is to fill these research gaps. To achieve this, we integrate previous literature on uplift modeling, evaluate the effectiveness of alternative uplift modeling strategies for campaign planning, and examine the degree to which this effectiveness depends on the learning algorithm to implement the modeling strategy.

In pursuing its objective, the paper makes the following contributions. First, we consolidate the state-of-the-art in uplift modeling. The comprehensive literature examination helps us understand and clarify the conceptual differences between different approaches and, through studying different streams of research, provides an update on modern uplift modeling research. Second, we empirically evaluate the performance of uplift modeling strategies for conversion modeling through large-scale experimentation. In particular, we benchmark several strategies across numerous data sets of different product lines from multiple geographies. The benchmark experiment provides a reference point for other academics and practitioners in campaign planning and uplift modeling. In total, our empirical study includes three experiments. For the first experiment, we

consider multiple machine learning algorithms for the experimentation that we pair with each uplift modeling strategy in a full-factorial setup. Thus, we shed light on the interactions between uplift modeling strategies and underlying learning algorithms, and provide specific recommendations on their relative suitability. Based on the benchmark results, in the second experiment, we quantify the degree to which targeting marketing campaigns using uplift modeling increases business value. That is, we explain which strategy contributes most (least) to business value. To clarify differences in performance between response and uplift modeling, we compare the former with the latter in a third and final experiment.

2 Conversion modeling using uplift vs. response models

The term conversion modeling encompasses a set of marketing decision support models that estimate the probability of customers to react toward a marketing action in a way intended by the marketer. The goal of developing a conversion model is to allocate marketing resources efficiently. For example, marketers use conversion models to identify the most suitable channel to contact a customer in multi-channel advertising (Zantedeschi et al., 2016), to select responsive customers for email surveys (Michaelidou & Dibb, 2006), or, more generally, to inform targeting decisions (Daskalova et al., 2017; Ieva et al., 2018). These examples illustrate how conversion modeling finds broad application in marketing and e-commerce to anticipate customer behavior and to increase conversion rates.

We distinguish conversion models into response and uplift models. Response models rely on supervised classification algorithms (hereafter, base learners), which estimate a functional relationship between a binary class label (i.e., response vs. no response) and a set of explanatory variables that characterize customers. Such variables often include demographic, behavioral, and attitudinal information or, more generally, any piece of information an analyst believes to be possibly linked to customers' response behavior.

To target campaigns using a response model, candidate recipients are ordered according to model-estimated conversion probabilities and a fraction of the top-ranked recipients is contacted, whereby the size of the target group depends on the available budget and/or other business considerations. Uplift models do not predict conversion probabilities. Their objective is to predict how much the campaign changes the conversion probabilities of individual customers (Lessmann et al., 2018; Rzepakowski & Jaroszewicz, 2012b). The important implication for campaign planning is that an uplift model will recommend a target group of *persuadable* customers, whose conversion probability raises because of the campaign. The target group recommended by a response model,

on the other hand, will consist of *responsive* customers who may be influenced by the campaign or not. More formally, an uplift model estimates a conditional average treatment effect (Chernozhukov et al., 2018) and establishes a causal link between the marketing action and how it alters customer behavior. Causality is crucial in campaign targeting to maximize the efficiency of resource utilization. Marketing budget should be spent on those customers where it increases conversion probabilities the most. Table 1 further elaborates on the connection between the action and the behavioral change it induces by distinguishing four groups of customers. Without loss of generality, we assume in the following that a marketing campaign aims at direct selling. Hence, successful conversion implies that a customer purchases the offered item. We denote corresponding customers as buyers. Instead of campaign, we use the term treatment, which is more general than campaign and used in the econometrics literature on causal inference (Athey & Imbens, 2017). A customer with a treatment is one who received the marketing action (e.g., email-based digital coupon).

Table 1 Customer types as per uplift modeling

| | | Buyer without treatment | |
|----------------------|-----|-------------------------|--------------|
| | | Yes | No |
| Buyer with treatment | Yes | Sure Things | Persuadables |
| | No | Do-Not-Disturbs | Lost Causes |

According to Table 1, customers classified as *sure things* buy regardless of the treatment while *lost causes* never buy. Clearly, contacting these groups with a marketing message is a waste of resources. Even worse, the effect of a treatment is detrimental for customers in the *do-not-disturbs* group. Their conversion probability decreases when being treated. Last, the *persuadables* buy if being treated and refrain from buying otherwise. This means they buy because of the treatment and are thus the only group worth considering in targeted marketing actions (Rzepakowski & Jaroszewicz, 2012b). Targeting *persuadables* allows marketing managers to maximize the incremental number of purchases which implies an efficient use of marketing resources.

Table 1 also reveals a conceptual difference between response and uplift models. Response models require a labeled data set of customers the actual buying behavior of whom is known from a past campaign. This is the standard setting in supervised learning. In addition to a target label, developing an uplift model also requires data from two groups of customers, a treatment group who received the marketing action and a control group who did not. Random trials, pilot campaigns, or A/B tests are common instruments to obtain corresponding data. Subsequently, a

straightforward way to develop an uplift model, called the two model uplift method, involves estimating two classification models from the treatment and control group data, respectively. To estimate conversion uplift as the difference in the predicted purchase probabilities with and without treatment, predictions for new customers (e.g., potential recipients of an upcoming campaign) are subtracted (Radcliffe & Surry, 1999). A formal derivation of this approach follows in the fourth section of this paper.

3 Prior work in uplift modeling

In general, uplift models estimate the entry probability of an event of interest through relating the event to a set of explanatory variables. A crucial difference to ordinary regression models or supervised machine learning is that uplift models aim at estimating how the probability of the event changes with specific actions. Estimating how treatment with a certain medication changes the survival probability of a patient exemplifies this approach in a medical application context. Here, the action is to apply the medication or refrain from doing so (Jaskowski & Jaroszewicz, 2012). Another example arises in marketing where a marketer is interested to estimate how actions in the form of marketing messages (newsletters, telephone calls, etc.) alter the purchase behavior of customers (Dost et al., 2014). These examples hint at the variety of applications in which an uplift model may be useful. The methodology underlying such models, however, is the same, which explains why prior work on uplift models spreads across different academic disciplines.

In terms of methodology, previous work on uplift models splits into three streams. The first stream comprises studies that develop uplift models using machine learning algorithms (Jaskowski & Jaroszewicz, 2012; Lo, 2002; Tian et al., 2014). We use the term *uplift modeling strategy* to refer to corresponding approaches because they embed a conventional learning algorithm into an overall modeling framework that facilitates predicting uplift. The second stream of literature develops new learning algorithms to predict uplift (Guelman et al., 2015; Rzepakowski & Jaroszewicz, 2012a; Zaniewicz & Jaroszewicz, 2013). We summarize corresponding approaches using the umbrella term *uplift algorithm*. Finally, the development of an uplift model is only one step in an overall modeling process. The third stream of research includes studies that focus on process steps preceding uplift model development such as feature selection or variable importance assessment, as employed by Hua (2016) for instance. Similarly, some studies concentrate on tasks that follow uplift model building. Nassif et al. (2013b) exemplify this approach through proposing new evaluation measures for uplift model assessment. Table 2 summarizes previous work on uplift

models along four dimensions: central focus of the study, uplift literature stream, experiments, industry/science, and data origin.

Table 2 reveals that much previous research is directed toward developing uplift algorithms. Corresponding works often draw inspiration from decision trees and modify algorithms for tree induction so as to predict uplift. Early work of Radcliffe and Surry (1999) introduces tree-based uplift algorithms. Aiming at classifying recipients of a direct marketing campaign into buyers and non-buyers, their idea was to alter the splitting criterion, which governs tree growing, in such a way that it maximizes the difference between the response rate of customers in the treatment and control group. An explicit consideration of the treatment and control group alongside the class variable is the main difference to ordinary classification trees, which consider only the class variable when inducing splits. Several later studies employ a similar approach and propose improved ways to induce uplift trees. Examples include Hansotia and Rukstales (2002b), who extend the X^2 criterion of the CHAID algorithm to accommodate uplift, Chickering and Heckerman (2000), who propose an approach to grow uplift trees so as to maximize expected profits, or Rzepakowski and Jaroszewicz (2012a) who further elaborate on tree induction through maximization of treatment and control group class distributions, and introduce novel splitting criteria based on conditional divergence. The tree-based uplift algorithm of Radcliffe and Surry (2011) assesses the statistical significance of the differences among class probabilities between treatment and control group observations. Guelman et al. (2014) and Guelman et al. (2015) propose uplift random forest, which they derive from embedding conditional inference trees and other uplift trees in an ensemble framework. Specifically, they mimic the original random forest classifier and combine bagging with random subspace to ensemble member (uplift) models (Breiman, 2001). Sołtys et al. (2015) systematize existing and contribute new uplift ensemble methods and evaluate them in marketing and medical applications.

Table 2 also illustrates that relatively few studies concentrate on uplift modeling strategies. Lo (2002) as well as Tian et al. (2014) introduce modeling strategies based on transformed data input spaces to facilitate uplift predictions. Pursuing the same goal, Jaskowski and Jaroszewicz (2012) propose a methodology to modify the data output space (i.e., response variable). Shaar et al. (2016) refer to disturbance effects of uplift models that limit prediction reliability. To cope with these effects, authors combine diverse uplift modeling strategies, including the uplift model of Lai et al. (2006) and reflective uplift modeling in a weighted procedure to derive a pessimistic uplift score. Building on the ideas by Lai et al. (2006), Kane et al. (2014) introduce a generalized weighting procedure of class probabilities.

In a benchmarking experiment, Kane et al. (2014) empirically compare some of the above strategies. De Vriendt et al. (2018) also contrast alternative uplift models amongst which they consider uplift modeling strategies. However, as both studies exemplify, prior literature on uplift modeling strategies considers a relatively small set of supervised learning algorithms. Irrespective of the development of an uplift modeling strategy or uplift algorithm, studies generally employ a single base learner without further empirical testing. As described above, tree-based algorithms are especially popular and used, amongst others, by Radcliffe and Surry (1999); Hansotia and Rukstales (2002b); Chickering and Heckerman (2000). Other studies consider base learners such as logistic regression (Lo, 2002), neural networks (Manahan, 2005), and k-nearest-neighbors (Larsen, 2010). We also observe some authors to use support vector machines for uplift modeling (Jaroszewicz & Zaniewicz, 2016; Kuusisto et al., 2014; Zaniewicz & Jaroszewicz, 2013). Due to the focus of previous research to consider a specific learning algorithm, empirical evidence related to interactions between uplift modeling strategies and learning algorithms is lacking. Therefore, one objective of this paper is to implement uplift modeling strategies using a set of alternative classification algorithms, which we believe to offer original insights related to the relative suitability of different learners to implement specific uplift modeling strategies.

The overarching conclusion emerging from Table 2 for e-commerce in general and marketing campaign planning is that available approaches in uplift modeling come from diverse strands of literature. This motivates a systematic comparison of the performance of such approaches, which, according to Table 2, is missing. Given that marketers typically use response models (Coussement et al., 2015), modification of such models to account for uplift effects would mean additional efforts and, eventually, sacrifice of well-timed performance. In contrast to the individual development of uplift algorithms, we therefore regard uplift modeling strategies as more beneficial for e-commerce since they make it possible to apply several supervised learning algorithms for uplift modeling without the need of modification. From current literature we observe that there are only few papers that focus on such strategies and an empirical comparison of available strategies with several supervised learning algorithms is lacking. Instead, studies in uplift modeling focus on single models which is why specific recommendations which models are comparably most valuable to apply are missing. To close this research gap, we benchmark available modeling strategies for conversion uplift that we pair with multiple base learners.

Table 2 Prior work in uplift modeling

| Study | Main topic | Research stream | Experiment | Industry / science | Data origin |
|------------------------------------|---|--------------------------|--|-------------------------------|--|
| Cai et al. (2011) | Two-stage estimation procedure for treatment differences for HIV-infected patients | Uplift algorithm | Treatment 1: Therapy based on drug combination (zidovudine, lamivudine) Treatment 2: Therapy based on drug combination (zidovudine, lamivudine, indinavir) | Clinical trials | Licensed open-source real-world data (AIDS Clinical Trials Group; see study ACTG 320 (Cole & Stuart, 2010)) |
| Chickering and Heckerman (2000) | Greedy decision-tree learning algorithms (FORCE vs. NORMAL) | Uplift algorithm | Mail advertisement for MSN subscription | Software | Private real-world data (anonymized authority) |
| De Vriendt et al. (2018) | Literature survey and empirical analysis of uplift models for marketing decision support | Specific task | Treatment 1: Marketing of insurances Treatment 2: Email marketing Treatment 3: Catalog mailing Treatment 4: Retention marketing | Various | R-package information ² Open-source real-world data (Hillstrom, 2008) Data from an Udemy online course ³ Private real-world data from a retention program |
| Dost et al. (2014) | Willingness-to-pay (WTP) range-based targeting approach | Uplift algorithm | Experiment 1: Discount offer Experiment 2: WOM (T1), visual (T2), information (T3) Experiment 3: Discount (T1), guarantee (T2) Experiment 4: Participation | Various | Surveys in different settings Participants from Amazon Mechanical Turk Participants from Amazon Mechanical Turk Students from a German university Consumers from an agency panel |
| Guelman (2014) | Personalized treatment learning problem, uplift random forest and uplift causal conditional inference forest | Uplift algorithm | Email promotion to buy a certain product at a bank | Financial services | Private real-world data (anonymized authority) |
| Guelman et al. (2012) | Uplift random forests | Uplift algorithm | Treatment 1: Letter (retention) Treatment 2: Letter plus outbound courtesy call (retention) | Insurance | Private real-world data (anonymized authority) |
| Guelman et al. (2014) | Causal conditional inference trees in personalized treatment learning | Uplift algorithm | Direct mail campaign (cross-selling) | Insurance | Private real-world data (anonymized authority) |
| Guelman et al. (2015) | Uplift random forests | Uplift algorithm | Information letter plus courtesy call (as one treatment) | Insurance | Private real-world data (anonymized authority) |
| Hansen and Bowers (2008) | Stratification to balance the distributions of pre-treatment variables | Specific task | Especially: GOTV field experiment (GOTV messages: personal visit, phone call, mailing) and simulation studies | Social and political sciences | Get-Out-The-Vote (GOTV) field experiment (Green & Gerber, 2015) |
| Hansotia and Rukstales (2002a) | Concept of uplift tree-based approaches | Uplift algorithm | - | - | - |
| Hansotia and Rukstales (2002b) | CHAID decision tree with ΔP split criterion | Uplift algorithm | Mail promotion (\$10 off a purchase of at least \$100 basket value) | Holiday retail | Private real-world data (anonymized authority) |
| Hua (2016) | Uplift random forests in capital market research with focus on results of embedded variable selection procedure | Specific task | - | Financial services | Licensed open-source real-world data (different data sources) |
| Imai and Ratkovic (2013) | Estimation of heterogeneous treatment effects as a variable selection problem with modified support vector machines | Specific task | GOTV field experiment (GOTV messages: personal visit, phone call, mailing) and simulation studies | Social and political sciences | Get-Out-The-Vote (GOTV) field experiment (Green & Gerber, 2015) |
| Jaroszewicz and Rzepakowski (2014) | Uplift modeling for survival analysis | Uplift algorithm | Chemotherapy against colon cancer Treatment 1: Therapy with Levamisole Treatment 2: Therapy with Levamisole plus 5-Fluorouracil | Clinical trials | Open-source real-world data (Dheeru & Karra Taniskidou, 2017) |
| Jaroszewicz and Zaniewicz (2016) | Uplift support vector machines with Székely regularization | Uplift algorithm | Therapy with right heart catheterization procedure (RCH) | Clinical trials | Open-source real-world data (Connors et al., 1996) |
| Jaskowski and Jaroszewicz (2012) | Response variable transformation | Uplift modeling strategy | Experiment 1: Therapy with peripheral blood transplant Experiment 2: Therapy with tamoxifen plus radio therapy against breast cancer Experiment 3: Therapy with steroids against hepatitis | Clinical trials | Open-source real-world data (Pintilie, 2006) Open-source real-world data (Pintilie, 2006) Open-source real-world data (Dheeru & Karra Taniskidou, 2017) |

² <https://cran.r-project.org/web/packages/Information/index.html>³ <https://www.udemy.com/uplift-modeling>

| | | | | | |
|------------------------------------|--|--------------------------|--|--|---|
| Kane et al. (2014) | Generalized weighting procedure of class probabilities, comparison of uplift approaches; signal-to-noise (S/N) ratio | Uplift modeling strategy | Experiment 1: Direct mail (paper) Experiment 2: Email Experiment 3: Direct mail (paper) | Financial services, online merchandise, retail office supplies | Private real-world data (anonymized authority) |
| Kuusisto et al. (2014) | Uplift support vector machines | Uplift algorithm | Simulated marketing activity | - | Simulation data |
| Lai et al. (2006) | Transformation scheme with weighted class probabilities | Uplift modeling strategy | Loan product promotion | Financial services | Canadian Imperial Bank of Commerce (CIBC) |
| Larsen (2010) | Uplift k-nearest neighbor and variable selection | Uplift algorithm | - | - | - |
| Lo (2002) | Interaction term approach | Uplift modeling strategy | - | - | Simulation data |
| Lo and Pachamanova (2015) | Multiple treatment optimization approach for prescriptive uplift analytics | Uplift algorithm | Email campaign (men and women separately targeted) | Online retail | Open-source real-world data (Hillstrom, 2008) |
| Manahan (2005) | Uplift neural network implementation with SAS | Uplift algorithm | Contract renewal campaign | Telecommunication | Private real-world data (Cingular) |
| (Nassif et al., 2013a) | Multi-relational uplift modeling system for medical research (SAYL algorithm) | Uplift algorithm | Therapy against breast cancer | Clinical trials | Open-source real-world data (Nassif et al., 2012) |
| (Nassif et al., 2013b) | Alternative uplift evaluation measures (ROC) | Specific task | Therapy against breast cancer | Clinical trials | Open-source real-world data (Nassif et al., 2010) |
| Radcliffe (2007) | Uplift evaluation measures | Specific task | Experiment 1: Catalogue mailing Experiment 2: Retention marketing Experiment 3: Cross-selling | Retail, telecommunication, financial services | Private real-world data (anonymized authority) |
| Radcliffe and Surry (1999) | Fundamental idea of uplift modeling with reference to differential response analysis | Uplift algorithm | - | - | - |
| Radcliffe and Surry (2011) | Significance-based uplift decision trees with several key features, uplift evaluation measures | Uplift algorithm | - | - | - |
| (Rzepakowski & Jaroszewicz, 2012a) | Uplift modeling for multiple treatments | Uplift algorithm | No campaign conducted (artificial allocation of observations to either treatment or control group in 16 datasets) | - | Open-source real-world data (Dheeru & Karra Taniskidou, 2017) |
| (Rzepakowski & Jaroszewicz, 2012b) | Uplift decision trees with different split criteria | Uplift algorithm | Email campaign (men and women separately targeted) | Online retail | Open-source real-world data (Hillstrom, 2008) |
| Shaar et al. (2016) | Pessimistic uplift modeling approach to minimize disturbance effects | Uplift modeling strategy | Simulated campaigns/treatments in marketing and medicine | - | Open-source real-world data (Dheeru & Karra Taniskidou, 2017; Hillstrom, 2008; Pintilie, 2006) |
| Sołtys et al. (2015) | Ensemble methods for uplift modeling (bagging, random forest) | Uplift algorithm | Simulated campaigns/treatments in marketing and medicine | - | Open-source real-world data (Dheeru & Karra Taniskidou, 2017; Hillstrom, 2008; Pintilie, 2006) |
| Su et al. (2012) | Causal inference trees and uplift k-nearest neighbor approach in assessing treatment effects | Uplift algorithm | Synthetic data creation (uniform distribution) | Machine learning research | Simulation data |
| Tian et al. (2014) | Investigation of the effects of a transformation of input space on a certain outcome of interest in medical research | Uplift modeling strategy | 1. Study of the implications of ACE inhibitors on lowering cardiovascular risk for patients with stable coronary artery disease and normal or reduced left ventricular function 2. Study of interactions between gene expression levels and Tamoxifen treatment in breast cancer patients | Clinical trials | 1. Preventive of Events with Angiotension Converting Enzyme Inhibition (PEACE) study (Braunwald et al., 2004) 2. Breast cancer dataset consisting of 414 patients in the cohort GSE6532 (Loi et al., 2007) |
| Yong (2015) | Prediction inference procedure with stratification to obtain generalizable predictions for medical examinations | Specific task | Several | Clinical trials | Several; among them the Mayo liver study data |
| Zaniewicz and Jaroszewicz (2013) | Uplift support vector machines (USVM) | Uplift algorithm | Simulated campaigns/treatments in marketing and medicine | - | Open-source real-world data (Dheeru & Karra Taniskidou, 2017; Hillstrom, 2008; Pintilie, 2006) |

4 Uplift modeling strategies

In this study, we empirically benchmark eight uplift modeling strategies. We depict these strategies in Table 3. The strategies have been proposed in previous work and used in diverse settings. Evidence on their relative effectiveness in a given context is lacking and, thus, originally provided here. We reintroduce the modeling strategies in subsequent sections and distinguish between basic, advanced and special strategies for conversion uplift. The latter exhibits a comparable level of complexity as advanced strategies but does not necessarily focus on data transformation schemes. Rather, related strategies have their own distinct characteristics and are based on most recent research.

With the choice of strategies, we are confident to provide a wide portfolio of state-of-the-art uplift modeling strategies. Recall that the strategies enhance execution of standard classification procedures for uplift modeling. As a result, the strategies can be practiced directly in e-commerce initiatives such as customer acquisition, customer development (Kane et al., 2014), or customer retention (Guelman et al., 2015) without a need to modify base learners.

Table 3 Uplift modeling strategies

| Category | Uplift modeling strategy | Acronym | Source |
|----------|--|-------------|----------------------------------|
| Basic | Two Model Uplift Method | TWO_MODEL | <i>Various</i> |
| Advanced | Interaction Term Method | ITM | Lo (2002) |
| | Treatment-Covariates Interactions Approach | TCIA | Tian et al. (2014) |
| | Class Variable Transformation | CVT | Jaskowski and Jaroszewicz (2012) |
| | Lai's Weighted Uplift Method | LWUM | Lai et al. (2006) |
| | Lai's Generalized Weighted Uplift Method | LGWUM | Kane et al. (2014) |
| Special | Reflective Uplift Modeling | REFLECTIVE | Shaar et al. (2016) |
| | Pessimistic Uplift Modeling | PESSIMISTIC | Shaar et al. (2016) |

Consider a training set $TRAIN_m = \{(x_i, y_i)\}_{i=1}^m$ of m customers gathered, for example, by means of a pilot campaign. Every customer is characterized by a set of explanatory variables x_i and a binary variable $y_i \in \{0, 1\}$ that indicates whether a conversion has been observed. We refer

to y_i as the target variable that we seek to explain. Let T_i and C_i indicate the membership of customer i to the treatment or control group, with prior probability distributions $P(T_i)$ and $P(C_i)$. Then, $P(Y_i = 1|T_i, X_i)$ and $P(Y_i = 1|C_i, X_i)$ denote the conditional probabilities of conversion for treatment and control group customers, respectively. For notational convenience, we refer to these conditional probabilities as $P(Y_i|T_i)$ and $P(Y_i|C_i)$ in the following. Furthermore, we define the four unconditional probabilities as follows: $P(T_i \cap Y_i)$ treated and response, $P(T_i \cap \bar{Y}_i)$ treated and non-response, $P(C_i \cap Y_i)$ non-treated and response, and $P(C_i \cap \bar{Y}_i)$ non-treated and non-response.

4.1 Basic uplift modeling strategy

The *two model* uplift method (e.g., Radcliffe, 2007; Radcliffe & Surry, 1999) captures the difference in class probabilities by providing a mechanism to differentiate between structures of customers' motivation:

$$Uplift_i^{TWO_MODEL} = P(Y_i|T_i) - P(Y_i|C_i) \quad (1)$$

Building and predicting with two equal learning algorithms given these two samples constitutes the methodology of the two model uplift method. In contrast, response models predict $P(Y_i|T_i)$.

4.2 Advanced uplift modeling strategies

Lo (2002) proposes a modification of the explanatory variables. He introduces a dummy variable $D_i \in \{0, 1\}$ for control and treatment group, respectively. D_i is multiplied with the entire input space X_i to gain an interaction term that is used in model prediction:

$$Uplift_i^{ITM} = P(Y_i|X_i, D_i, X_i \cdot D_i) \quad (2)$$

More specifically, *ITM* of Lo (2002) first develops an uplift model from the training data where D_i is known for all customers. Unlike the two model uplift method, which develops individual classification models for treatment and control group customers, ITM estimates only one model. Then, to estimate uplift for a novel customer with characteristics X_i^{new} , this single model is evaluated twice; setting $D_i = 1$ and $D_i = 0$ in the first and second evaluation, respectively. As is evident from (2), the resulting probability predictions will differ because of D_i . The former represents

the customer's conversion probability if treated while the prediction resulting from setting $D_i = 0$ approximates the conversion probability without treatment. Similar to the two model uplift method, the estimate of conversion uplift is given by the difference between the two predictions.

Independently from Lo (2002), Tian et al. (2014) propose an uplift modeling strategy, called **TCIA** in the following, which is conceptually similar to ITM. Differences to ITM are minute and limited to the coding and scaling of the interaction terms. In particular, Tian et al. (2014) obtain a set of interaction terms, D_i^* , as $D_i^* = \frac{X_i^* \cdot D_i}{2}$, whereby X_i^* denotes the original covariates, X_i , after mean centering. Another difference relates to D_i , which in the case of ITM, represents a zero-one dummy variable for control and treatment group, respectively. Thus, ITM captures differences in the treatment effect via movements of the intercept. This differs in TCIA where Tian et al. (2014) set $D_i \in \{-1, 1\}$. As a result, TCIA captures the treatment effect by subtracting treatment and control group probabilities within one functional form. With these modifications, TCIA predicts uplift as:

$$Uplift_i^{TCIA} = P(Y_i | X_i, D_i^*) \quad (3)$$

Tian et al. (2014) have applied their uplift modeling strategy to study interactions between gene expression levels and drug substances regarding breast cancer patients. Guelman et al. (2014) further validated this modeling strategy by means of a simulation.

Jaskowski and Jaroszewicz (2012) present a transformation procedure - **CVT** - that develops a novel target variable based on the original target (i.e., binary conversion response) and the membership of the respective customer to either the treatment or control group. Let Z_i denote the binary transformed target variable corresponding to customer i . Then, $Z_i = 1$ if $(T_i \cap Y_i) \cup (C_i \cap \bar{Y}_i)$ is given; otherwise $Z_i = 0$. Thus, $Z_i = 1$ captures treated customers with response as well as non-treated customers without response. On the contrary, for treated customers without response as well as non-treated customers with response, $Z_i = 0$. The definition of Z_i is based on the link between the desired behavior and a marketing action. More specifically, $Z_i = 1$ reflects customers that convert due to an incentive, but do not convert if not being solicited. The focus of this modeling strategy is to target these customers because they are likely to be persuaded. Hence, as opposed

to previous uplift modeling strategies, the uplift effect is based on the distribution of the transformed conversion variable and defined as:

$$Uplift_i^{CVT} = 2 \cdot P(Z_i = 1) - 1 \quad (4)$$

Lai et al. (2006) presents an extension of CVT - **LWUM** - that weights probabilities of positive and negative classes. LWUM assumes that the positive uplift lies in correctly identified *persuadables* (here, treatment-group responders and control-group non-responders), whilst the negative uplift can be found in the *do-not-disturbs* group (here, treatment-group non-responders and control-group responders). Therefore, let W be the number of positive observations divided by the total population. The uplift effect is then defined as:

$$Uplift_i^{LWUM} = P(Z_i = 1) \cdot W - P(Z_i = 0) \cdot (1 - W) \quad (5)$$

LWUM, thus, seeks to maximize the positive uplift while decreasing negative uplift in the first decile.

Kane et al. (2014) present **LGWUM** as the generalized version of LWUM with weighted probability scores that realize the influence of the fraction of treatment and control group customers on the lift measure and is defined as:

$$Uplift_i^{LGWUM} = P(Y_i|T_i) + P(\bar{Y}_i|C_i) - P(\bar{Y}_i|T_i) - P(Y_i|C_i) \quad (6)$$

4.3 Special uplift modeling strategies

Shaar et al. (2016) present the **reflective** uplift modeling strategy by two separate models that are built to learn the treatment effect in the conversion and non-conversion groups. The authors recognize disturbance effects when applying uplift models. The first one is a response effect that takes place due to correlation between explanatory variables and a binary class label, and the second effect – a partitioning effect – appears when the treatment indicator depends on the covariates. To overcome these negative effects, reflective uplift modeling has been introduced. The uplift effect is then calculated, whereas the groups are treated as positive and negative as in CVT:

$$Uplift_i^{REFLECTIVE} = P(T_i \cap Y_i) \cup P(C_i \cap \bar{Y}_i) - P(T_i \cap \bar{Y}_i) \cup P(C_i \cap Y_i) \quad (7)$$

Thus, the probabilities for positive and negative groups are obtained from two different models. To determine a score in terms of pessimistic uplift modeling, LWUM is again considered. The final *pessimistic* uplift modeling strategy is defined as:

$$Uplift_i^{PESSIMISTIC} = 0.5 \cdot (Uplift_i^{LWUM} + Uplift_i^{REFLECTIVE}) \quad (8)$$

4.4 Conceptual evaluation

In this section, we examine the relative merits of the modeling strategies for conversion uplift from a conceptual perspective. First, we consider the two model uplift method that is presented by the difference between the class probabilities (i.e., treatment vs. non-treatment). This modeling strategy suffers from poor approximation, since both probability estimates originate from two separate samples (e.g., Guelman et al., 2012; Jaroszewicz & Rzepakowski, 2014). ITM (Lo, 2002) and TCIA (Tian et al., 2014) manipulate the data input space through interaction terms with dummy variables indicating the treatment effect. Incorporating interaction effects for all variables, these uplift modeling strategies increase dimensionality. Therefore, ITM and TCIA appear less suitable for data sets where the number of original variables is large. CVT as in Jaskowski and Jaroszewicz (2012), on the contrary, changes the response variable to facilitate focusing on *persuadables* and *do-not-disturbs* and improves targeting decisions. However, CVT does not regard the difference between the relative sizes of positive and negative observations. This is why Lai et al. (2006) introduced the weights as per proportion of positive and negative observations and developed LWUM. This uplift modeling strategy should address differences in these proportions. However, we expect the accuracy of LWUM to suffer when the ratio of treatment and non-treatment assignments is not approximately equal. LGWUM (Kane et al., 2014) overcomes this and is designed to combat disturbance such as multicollinearity. Shaar et al. (2016) presents the reflective uplift modeling strategy that estimates the uplift effect from the conversion and non-conversion groups. The authors further extend it through pessimistic uplift modeling that combines LWUM and the reflective uplift modeling strategy into one model. The combination is claimed to overcome the disadvantage of the two model uplift method, where the separated estimation of response

probabilities among treatment and control group customers deteriorates the accuracy of uplift predictions (e.g., Guelman et al., 2012; Jaroszewicz & Rzepakowski, 2014; Shaar et al., 2016).

5 Experimental setup

We involve numerous data sets that belong to the field of e-commerce, indicating the goal to categorize the customer base into two classes: buyer and non-buyer. In the following, we elaborate the campaign process and underlying data, base learners to be paired with the aforementioned uplift modeling strategies, and finally the performance metrics.

5.1 Campaign process and data

The experimental setup involves 27 data sets from several digital marketing campaigns. These campaigns were executed by Akanoo⁴, a company specializing in analytics-as-a-service solutions for online shops. Akanoo provided us with a fully anonymized version of real-world campaign data in the scope of a research collaboration. The data is sensitive and can, therefore, not be disclosed to the public. It includes multiple campaigns that were carried out in different electronic marketplaces and designed so that customers who show specific behavioral patterns during their shop visit, as identified by an uplift model, are targeted with a digital coupon. Customers that leave the respective shop by having activated this coupon obtain a discount of 10% off their final basket value. A real-time targeting process has been applied to identify customers to receive the coupon. Every customer has been assigned either to the treatment or control group by chance or by a model. In the latter case, the individual online behavior of new customers is considered after five pageviews and that of returning customers after three pageviews. The derived predictive scores determine whether the customer is likely to be persuadable (i.e., customer with high probability to respond if being treated with coupon). As a result, the model qualifies the customers to the treatment group. The systematic component of the targeting process creates a selection bias that leads to a quasi-experiment.

⁴ <https://akanoo.com/>

Table 4 summarizes the available data sets in terms of product line, geographical location, the number of observations and responses in the treatment and control group, respectively, and the uplift.

Table 4 Summary of e-retail data sets

| Shop | Product line | Geographical location | No. of cases: treatment/control | No. of responses: treatment/control | Uplift |
|------|----------------------------|-----------------------|---------------------------------|-------------------------------------|--------|
| 1 | Apparel | Poland | 206,148 / 69,177 | 6909 / 2289 | 0,04% |
| 2 | Apparel | Germany | 128,469 / 43,467 | 8277 / 2523 | 0,64% |
| 3 | Apparel | Germany | 36,288 / 12,327 | 3054 / 879 | 1,29% |
| 4 | DIY products | United Kingdom | 216,534 / 72,978 | 5160 / 1560 | 0,25% |
| 5 | Apparel | Czech Republic | 46,983 / 16,284 | 2733 / 1005 | -0,35% |
| 6 | Apparel | Germany | 8733 / 2877 | 786 / 234 | 0,87% |
| 7 | Books and multimedia | Germany | 9003 / 3030 | 360 / 114 | 0,24% |
| 8 | Toys | Germany | 898,734 / 300,847 | 96,318 / 31,874 | 0,12% |
| 9 | DIY products | Germany | 92,961 / 31,125 | 1800 / 525 | 0,25% |
| 10 | DIY products | France | 9471 / 3309 | 501 / 129 | 1,39% |
| 11 | Pharmaceuticals | Germany | 5319 / 1680 | 2436 / 807 | -2,24% |
| 12 | Special apparel (hats) | Germany | 16,734 / 5580 | 1911 / 603 | 0,61% |
| 13 | Apparel/household items | France | 47,964 / 15,900 | 135 / 24 | 0,13% |
| 14 | Fan articles and toys | Germany | 9534 / 3303 | 777 / 168 | 3,06% |
| 15 | Apparel | Germany | 18,417 / 6033 | 2472 / 708 | 1,69% |
| 16 | Apparel | The Netherlands | 5520 / 1806 | 348 / 75 | 2,15% |
| 17 | Alcoholic beverages | Germany | 6996 / 2400 | 1803 / 624 | -0,23% |
| 18 | Pharmaceuticals | Germany | 6699 / 1998 | 3411 / 990 | 1,37% |
| 19 | Sports apparel/accessories | Germany | 83,865 / 27,765 | 13,428 / 4599 | -0,55% |
| 20 | Pet food | Germany | 16,881 / 5601 | 3456 / 1143 | 0,07% |
| 21 | Apparel | Germany | 89,424 / 30,141 | 6060 / 1926 | 0,39% |
| 22 | Shoes and accessories | Germany | 244,506 / 81,726 | 14,643 / 5031 | -0,17% |
| 23 | Pharmaceuticals | Germany | 4104 / 1239 | 2304 / 651 | 3,60% |
| 24 | Apparel | Austria | 20,913 / 6855 | 684 / 207 | 0,25% |
| 25 | Shoes | Germany | 2403 / 801 | 99 / 39 | -0,75% |
| 26 | Special apparel (hats) | The Netherlands | 7863 / 2589 | 396 / 114 | 0,63% |
| 27 | Outdoor apparel | Germany | 45,210 / 14,928 | 2469 / 846 | -0,21% |

Table 4 indicates that the consumer goods relate to different sorts of apparel, toys, garden articles, books and multimedia, pet food, and many other. In addition, sports and outdoor articles are also sold in a few shops. Businesses operate in Austria, the Czech Republic, France, Germany, the Netherlands, Poland, and the United Kingdom. The total number of cases across all data sets

is roughly three million. On average, we observe a treatment to control group ratio of 3:1 meaning that three out of four customers have received a digital coupon. Based on the number of responses in the treatment and control group, we capture the impact of the shop-wise marketing campaigns. The last column in Table 4 reports campaign uplift per shop, which we calculate as the differences between the relative response rate in the treatment and control group, respectively. Table 4 reveals low (positive) uplift for almost every shop. For some shops, we even observe negative uplift resulting from the response rate in the control group exceeding the response rate in the treatment group. The average uplift across the 27 shops is 0.54%.

The data contains 60 features that profile the customers' behavior. Every observation relates to the shop-based journey performed during a certain time span (i.e., from entering to leaving the shop). Cookie technology allows to differentiate between new and returning customers. Most features are numeric while some are factors. These features provide information on numerous customer activities during the shop visit, for instance, how much time the customer spends on certain page types, whether the customer has purchased a specific product in the same shop in the past and how much time has passed since the customer added an item to the shopping cart. Further examples relate to how many views the customer has made on a sale-related page and how many products lie in the customer's shopping basket during the current session. Inspiration on data collection has been gained from Van den Poel and Buckinx (2005). Furthermore, meta dimensions that are crucial for the use of uplift models have been collected, particularly, an indicator of treatment or control group assignment, and a variable that captures the purchase event. Being structured in terms of current session, previous session(s), and identifiables of the respective customer, Table 5 lists and describes all features used for the empirical study.

Table 5 Clickstream features used for empirical study

| Setting | # | Name of feature | Description | Based on Van den Poel and Buckinx (2005) |
|----------------------------|----|-----------------------------------|---|--|
| Current session | 45 | InitBasketNonEmpty | State of the initial basket (empty vs. non-empty) | |
| | | HadBasketAdd | Whether the visitor has added at least one product to the basket | |
| | | TimeToBasketAdd | Amount of time since a product has been added to the basket | |
| | | BasketQuantity | Number of products in current basket | |
| | | NormalizedBasketSum | Normalized value of customer basket (for comparisons across shops) | |
| | | TimeToFirst (pagetype) | Time span from shop arrival to first click on page type 'cart' / 'overview' / 'product' / 'sale' / 'search' | |
| | | TimeSinceFirst (pagetype) | Amount of time since first click on page type 'cart' / 'overview' / 'product' / 'sale' / 'search' | |
| | | TimeSinceOn (pagetype) | Duration on page type 'cart' / 'overview' / 'product' / 'sale' / 'search' until leave of online shop | |
| | | TimeOn (pagetype) | Duration on page type 'cart' / 'overview' / 'product' / 'sale' / 'search' until leave of page type | |
| | | HourOfDay | Hour of the day (1 – 24) when the visitor has entered the online shop | |
| | | SessionTime | Duration of current visitor session | X |
| | | ScrollHeight (overview) | Scroll height for pages of type 'overview' | |
| | | ScreenWidth | Screen width of customer device | |
| | | TabSwitch (product) | Number of total tab switches for pages of type 'product' during session | |
| | | Clicks (product) | Number of clicks for pages of type 'product' | |
| | | TimeSinceClick | Time span from first click to shop leave | |
| | | TimeSinceTabSwitch | Time span from first switch of tabs | |
| | | ViewCount | Number of views in the current session | X* |
| | | ViewedBefore (cart) | Whether visitor has already viewed a specific page from page type 'cart' | |
| | | ViewsOn (pagetype) | Number of views on page type 'cart' / 'overview' / 'product' / 'sale' / 'search' | X* |
| | | InitPageWas (overview) | Whether initial page had page type 'overview' | |
| | | InitPageWas (product) | Whether initial page had page type 'product' | |
| | | InitPageWas (sale) | Whether initial page had page type 'sale' | |
| Previous session(s) | 8 | NumberOfDifferentPages (overview) | Number of views on different pages from page type 'overview' | X* |
| | | NumberOfDifferentPages (product) | Number of views on different pages from page type 'product' | X* |
| | | TimeSinceLastConversion | Amount of time since last product purchase | X |
| | | VisitCountLastWeek | Number of shop visits within the previous week | |
| | | VisitCountToday | Number of shop visits during day of session-of-interest | |
| | | PreviousVisitCount | Number of previous shop visits | X |
| | | TimeSinceFirstVisit | Amount of time since first shop visit | |
| | | TimeSinceLastVisit | Amount of time since last shop visit | X |
| | | DurationLastVisit | Time span of previous shop visit | X* |
| | | ViewCountLastVisit | Number of views during last shop visit | |
| Identifiables | 7 | VisitorKnown | Whether the visitor has already visited the shop in the past | X* |
| | | WasConvertedBefore | Whether the visitor has already purchased a product in a previous session | X |
| | | Conversion | Whether the visitor has purchased a product in session-of-interest | |
| | | Normalized revenue | Normalized amount of revenue (for comparisons across shops) | |
| | | Treatment/control group | Whether the visitor has been shown the e-coupon | |
| | | Shop-ID | Unique shop identifier | |
| | | Timestamp | Point in time when visitor has entered the online shop | |

* Based on Buckinx and Van den Poel (2005) but slightly adapted

Another important concern relates to data partitioning. We have created three partitions from the available data: 40% training partition that we use to train the strategies, 30% for a parameter-tuning partition that we use to validate the meta-parameter tuning, and another 30% for a test partition. To guarantee a reliable evaluation, we apply a 10-fold cross validation scheme “through time” to reflect the situation in marketing practice and increase the size of observations by resampling. For all uplift modeling strategies, the stated models first predict on the training and parameter-tuning partitions together. Strategy-wise models with the best candidate settings are then validated on the validation sample to assure a reliable benchmark.

5.2 Base learners

The experimental design includes six base learners to ensure a vast benchmark study. Recall that we benchmark modeling strategies for conversion uplift that can be paired with any base learner. Thus, we secure every possible combination between uplift modeling strategy and base learners. The experiment is performed in Python and builds upon libraries for data manipulation, statistics, visualization and data science; namely NumPy, Pandas, Matplotlib and Scikit-learn (Pedregosa et al., 2011). We consider a wide range of meta-parameters for every base learner (see Table 6). Every model is tuned automatically and transmitted to the cross-validation technique discussed previously. In total, we involve 245 models.

We pair base learners and modeling strategies for conversion uplift in a full-factorial experimental setup. Recall that we involve eight uplift modeling strategies and, additionally, response modeling. This, thus, results in 2,205 models in total. We choose the base learners due to their popularity in response and uplift modeling. In response modeling, for example, they are often questioned in pivotal benchmark studies (Baesens et al., 2003; Lessmann et al., 2015). SGDC and RFC demonstrate excellent performance in real-world experiments (Guelman et al., 2015). Due to the fact that RFC is less sensitive to meta-parameter adaptations than SGDC (Ogutlu et al., 2011), we consider for RFC a smaller number of models. In uplift modeling, LogR (Lo, 2002), KNN (Larsen, 2010), and SVC (Kuusisto et al., 2014; Zaniwicz & Jaroszewicz, 2013) have gained a strong research interest. As a standard base learner without meta-parameters, we add a NB algorithm to the library of base learners.

Table 6 Meta-parameters of the base learners

| Base learner | Acronym | No. of models | Meta-parameter | Candidate setting |
|--|---------|---------------|--|--|
| Logistic regression | LogR | 34 | Regularization term Regularization factor | [L1, L2] [1e-8, 1e-7, ..., 1e8] |
| Support vector machines with linear kernel | SVC | 42 | Regularization factor Calibration method | [1e-10, 1e-9, ..., 1e10] [Sigmoid, Isotonic] |
| k-Nearest-Neighbor | KNN | 20 | Number of nearest neighbors | [1, 5, 10, 20, ..., 100, 200, ..., 500, 1000, 2000, ..., 4000] |
| Naïve Bayes | NB | 1 | - | - |
| Stochastic gradient descent for classification | SGDC | 144 | Loss function Regularization term Alpha Learning rate | [Log, Mod. Huber, Hidge, Percep.] [L1, L2, Elastic Net] [1e-6, 1e-5, ..., 1e-1] [Optimal, Invscaling] |
| Random forest for classification | RFC | 4 | Max. no. of covariates Min. no. of samples | [8, 9] [1000, 2000] |

5.3 Validation measures

Typically, the performance of predictive models grounds on a comparison of actual versus predicted outcomes. In uplift modeling, however, this is not reasonable since a customer cannot be part of both the treatment and control group. This phenomenon is known as the fundamental problem of causal inference (Holland, 1986). Consequently, today’s best practice is a decile-based evaluation approach to identify uplift. Hence, model performance is captured in terms of Qini coefficient Q and visualized in uplift gains charts by means of Qini curves (Radcliffe, 2007). This includes the assumption that similarly scored cases behave likewise, i.e., the k percent highest scores on treatment out-of-sample test data are compared to the k percent highest scores on control out-of-sample test data and with the subtraction of the top gains from both groups a meaningful estimate of uplift can be derived (Jaskowski & Jaroszewicz, 2012). Q is, thus, defined as the area between a model’s Qini curve and a random targeting line (Radcliffe & Surry, 2011). Because typically uplift gains charts display Qini curves that relate to a cumulative measure, we further consider uplift bar charts that mask the effect of cumulativeness to provide a decile-isolated analysis of model performance.

6 Empirical results

The experimental results consist of the performance estimates for every combination of 6 levels of base learners, 9 levels of modeling strategies (response modeling included), and 27 levels of data sets. The performance measures capture the degree to which the marketing campaign strategy improves via application of uplift modeling strategies in terms of Qini coefficient and cumulative (non-cumulative) number of incremental purchases.

6.1 Examination of the interaction between uplift modeling strategies and base learners

To identify synergy effects between the modeling strategies for conversion uplift and base learners, we now examine their interaction. Table 7 summarizes the corresponding results. To obtain them, we pair every base learner with all uplift modeling strategies and capture the predictive performance on the out-of-sample test set in terms of Qini coefficient. These values are averaged over the data sets. We express the Qini coefficient in percentage terms, i.e., Q_{pct} , by subtracting the control group response rate from the treatment group response rate for every decile. In contrast to the general Q coefficient (Radcliffe & Surry, 2011), Q_{pct} makes comparisons across the data sets with different number of observations possible and, thus, requires no normalization procedure. To increase the readability of Q_{pct} , we multiply its values with a factor of 1,000. We use bold face for every best combination (i.e., uplift modeling strategy and base learner). For example, the value in the last column for CVT is marked in bold face indicating that CVT interacts best with RFC.

Table 7 Qini coefficient of uplift modeling strategies

| Uplift modeling strategy | Base learner | | | | | |
|--------------------------|--------------|--------------|--------|--------|--------------|--------------|
| | KNN | LogR | NB | SGDC | SVC | RFC |
| CVT | 3.171 | 3.348 | -0.951 | -1.041 | 2.017 | 6.145 |
| ITM | 3.991 | 2.901 | 3.770 | 0.979 | 8.017 | 3.216 |
| LGWUM | -0.230 | 3.767 | -4.459 | 1.831 | -0.932 | 5.593 |
| LWUM | 3.171 | 4.258 | -0.945 | 0.203 | 2.049 | 6.130 |
| PESSIMISTIC | 1.418 | 4.269 | -1.626 | 0.720 | 2.010 | 6.606 |
| REFLECTIVE | -1.526 | 3.310 | -2.914 | 0.868 | -0.727 | 2.303 |
| TCIA | 1.043 | -1.950 | -2.821 | 1.222 | 3.893 | 0.403 |
| TWO_MODEL | 7.267 | 4.305 | 3.297 | 0.688 | 2.806 | 5.401 |

Table 7 reveals multiple important findings. First, the best possible interaction is between ITM and SVC with Q_{pct} of 8.017. This is followed by the two model uplift method coupled with KNN with Q_{pct} of 7.267 and CVT with RFC of 6.145. This strongly signals in favor of ITM as a modeling strategy for conversion uplift and of SVC as a base learner. This view is only strengthened when we look at the pair of TCIA and SVC, where SVC is the best performer. However, we recommend RFC as a base learner for uplift modeling since it collects the biggest number of wins. More specifically, RFC is the best performer when coupled with CVT, LGWUM, LWUM, and pessimistic uplift modeling. We observe that KNN performs best when paired with the two model uplift method and the differences in the performance compared to other uplift modeling strategies are substantial. For example, the pair of the two model uplift method and KNN achieves Q_{pct} of 7.267 compared to the second-best performer pair of ITM and KNN with Q_{pct} of 3.991 and the worst performer pair of the reflective uplift modeling strategy and KNN with Q_{pct} of -1.526. As a result, we can only recommend considering KNN when coupled with the two model uplift method. We also observe that the reflective uplift modeling strategy performs best coupled with LogR. However, LogR shows also high and better potential when interacting with other strategies. For example, Q_{pct} of couples of pessimistic uplift modeling, LWUM, and the two model uplift method with LogR is higher than that of reflective uplift modeling with LogR. Thus, LogR seems to be more flexible than KNN for uplift modeling. On the contrary, due to the weak performance compared to other base learners, NB and SGDC have no wins. Thus, we cannot recommend executing them for uplift modeling. This recommendation is supported by the fact that for many uplift modeling strategies, NB collects negative Q_{pct} values. The same applies to the pair of CVT and SGDC. We also would like to stress that the best pessimistic uplift model outperforms all base learners related to LWUM and reflective uplift modeling. This is interesting since LWUM and reflective uplift modeling hold equal shares in creation of the pessimistic modeling strategy. LGWUM does not add more value than LWUM. With SGDC being the only exception, all base learners paired with LWUM outperform their equivalents for LGWUM. Analogous picture we see for covariate transformations. All base learners but SGDC and SVC paired with CVT obtain higher Q_{pct} values when compared to respective TCIA counterparties.

To support findings of Table 7, we examine the robustness of the uplift modeling strategies. To achieve this, we capture the performance of the modeling strategies coupled with base learners in a 10-fold cross validation (see Figure 1). Every boxplot portrays base learners on the x -axis and the performance measured in Q_{pct} on the y -axis. We scale the Q_{pct} values to ease comparability.

Figure 1 highlights the performance of RFC. RFC is the best performer when coupled with, e.g., CVT, LGWUM, or pessimistic uplift modeling. Furthermore, RFC shows relatively small variance. This can be especially emphasized on the combination of RFC with the two model uplift method. Thus, Figure 1 further supports the view that RFC is a very suitable base learner for uplift modeling. Figure 1 also stresses weak performance of NB and SGDC. We observe that the mean values of NB are negative for the pessimistic and the reflective uplift modeling strategies as well as for CVT. The same we see on the couple of CVT and SGDC, whereby SGDC also exhibits higher variance than NB. These findings caution from execution of these base learners for uplift modeling. Third, we see a comparably high variance of SVC when coupled with ITM. This finding injects doubt on the previous insight where the couple ITM and SVC is the best performer (Q_{pct} of 8.017). Hence, we conclude that ITM paired with SVC does not provide a reliable estimate. In contrast, we observe that KNN paired with the two model uplift method shows low variance that makes this couple more promising than ITM with SVC (given findings from Table 7). More precisely, ITM-based SVC shows a standard deviation of 8.3 compared to two model-based KNN with 6.0. As a result, KNN has a 28% lower standard deviation than SVC. Note that the standard deviation values are percentages derived from taking the mean of all decile-wise values. At the same time, we conclude that KNN and SVC (ITM being exception) show stable results in terms of variance when coupled with other uplift modeling strategies. Same conclusion can be drawn for LogR which moreover enjoys comparably high stable results across the uplift modeling strategies. In general, we would like to conclude that ITM and the two model uplift method show the most promising results when interacting with all base learners (SGDC being exception). These uplift modeling strategies do not show negative Q_{pct} values, relatively low variance, and comparable results among the base learners. Reflective uplift modeling and TCIA demonstrate opposite performance and, thus, can be regarded as least beneficial modeling strategies involved in this study.

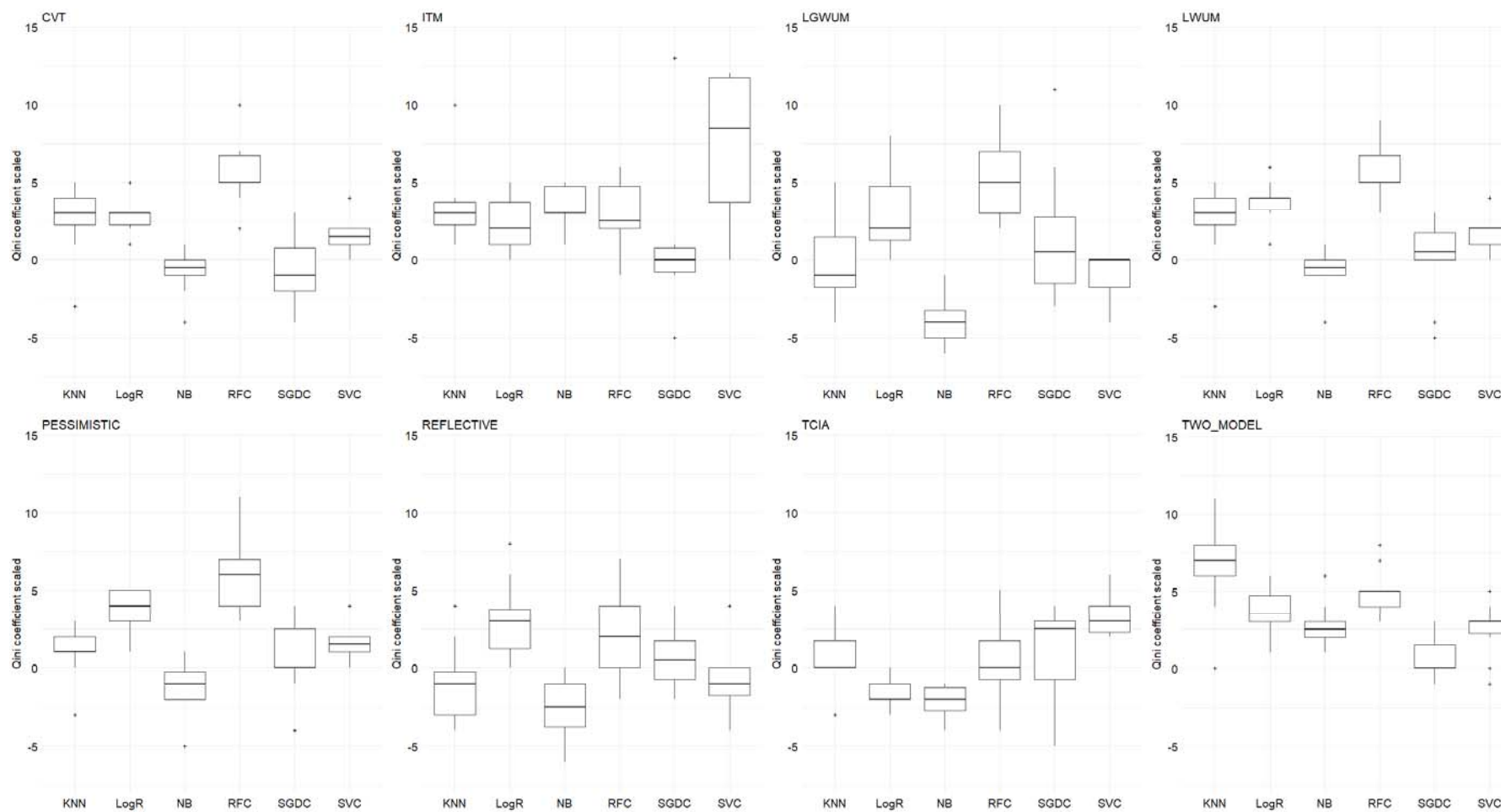


Figure 1 Scaled Qini coefficient across uplift modeling strategies

6.2 Examination of the impact of uplift modeling strategies on business value

We now examine the potential of the modeling strategies for conversion uplift to increase business value. To do so, we analyze the weighted model performance for every targeting decile in terms of the cumulative (and later non-cumulative) number of incremental purchases. One can think of a marketing campaign similar or identical to that we describe in this paper that targets a certain fraction of customers from the customer base. The purpose of this targeting is the product purchase that customers perform. That is, we capture the degree uplift modeling strategies contribute to the increase of those purchases. Again, we describe the effect of every uplift modeling strategy coupled with all base learners. Since the increased number of the incremental purchases results in increased revenue, we argue that uplift modeling might contribute to the increase of business value. To quantify the impact of the modeling strategies on business value, we first provide a tabular view of the decile-wise model performance. Table 8 presents the results obtained on the out-of-sample test set, across the uplift modeling strategies and base learners. We highlight in italic face the winner among the base learners within the uplift modeling strategy and in bold face a global winner (i.e., across all uplift modeling strategies) in every decile. Consider the very left (upper) column. We contact a 10% fraction of the customer base via marketing campaign. CVT enhances RFC to achieve 883 purchases. We mark this estimate in italic face indicating that RFC is the winner within the 10% fraction across the classifiers paired with CVT. Another example (same column) is the pair of pessimistic uplift modeling and LogR. This pair achieves 935 purchases within the first decile and is marked in both italic and bold face. The former indicates that LogR is the winner regarding the classifiers paired with the pessimistic modeling strategy for the first decile while the latter highlights that the pair of pessimistic uplift modeling and LogR presents the global winner in the first decile across all uplift modeling strategies.

Table 8 Summary of cumulative number of incremental purchases

| Uplift modeling strategy / base learner | Cumulative number of incremental purchases per decile | | | | | | | | | |
|---|---|------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|------|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| CVT | | | | | | | | | | |
| KNN | 301 | 655 | 831 | <i>1148</i> | 1213 | 1303 | 1332 | 1423 | 1486 | 1671 |
| LogR | 819 | 611 | 712 | 795 | 1133 | 1352 | 1401 | 1444 | 1547 | 1671 |
| NB | -234 | 278 | 421 | 213 | 874 | 1098 | 1281 | 1405 | 1533 | 1671 |
| SGDC | 66 | 225 | 422 | 602 | 754 | 918 | 1079 | 1300 | 1440 | 1671 |
| SVC | -209 | 239 | 484 | 823 | 1158 | <i>1601</i> | <i>1603</i> | 1631 | 1571 | 1671 |
| RFC | <i>883</i> | <i>983</i> | <i>1066</i> | 1110 | <i>1297</i> | 1456 | 1597 | <i>1641</i> | <i>1698</i> | 1671 |
| ITM | | | | | | | | | | |
| KNN | 418 | 673 | 868 | 901 | 1161 | 1381 | 1395 | 1722 | 1735 | 1671 |

| | | | | | | | | | | |
|-------------|------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|------|
| LogR | 487 | 765 | 1195 | 1332 | 1409 | 1755 | 1819 | 1690 | 1711 | 1671 |
| NB | 457 | 687 | 958 | 1564 | 1080 | 1058 | 1248 | 1445 | 1607 | 1671 |
| SGDC | 96 | 215 | 610 | 741 | 880 | 894 | 783 | 1751 | 2223 | 1671 |
| SVC | 395 | 1108 | 1292 | 1383 | 1602 | 1679 | 1885 | 1834 | 1838 | 1671 |
| RFC | 482 | 807 | 994 | 1043 | 1198 | 1205 | 1384 | 1341 | 1270 | 1671 |
| LGWUM | | | | | | | | | | |
| KNN | 347 | 431 | 256 | 499 | 667 | 968 | 1146 | 1392 | 1655 | 1671 |
| LogR | 591 | 867 | 878 | 855 | 1127 | 1207 | 1289 | 1564 | 1725 | 1671 |
| NB | 364 | 373 | 285 | 11 | 16 | 728 | 901 | 746 | 1039 | 1671 |
| SGDC | 434 | 642 | 871 | 1002 | 1109 | 1139 | 1226 | 1407 | 1425 | 1671 |
| SVC | 271 | 229 | 301 | 609 | 720 | 1010 | 1030 | 1280 | 1432 | 1671 |
| RFC | 302 | 689 | 842 | 1148 | 1501 | 1672 | 1701 | 1787 | 1713 | 1671 |
| LWUM | | | | | | | | | | |
| KNN | 301 | 655 | 831 | 1148 | 1213 | 1303 | 1332 | 1423 | 1486 | 1671 |
| LogR | 855 | 951 | 909 | 998 | 1143 | 1180 | 1422 | 1436 | 1543 | 1671 |
| NB | -243 | 286 | 408 | 232 | 871 | 1097 | 1282 | 1407 | 1533 | 1671 |
| SGDC | 172 | 385 | 557 | 734 | 869 | 936 | 1150 | 1346 | 1511 | 1671 |
| SVC | -208 | 245 | 474 | 808 | 1188 | 1602 | 1617 | 1625 | 1571 | 1671 |
| RFC | 884 | 991 | 1071 | 1103 | 1294 | 1448 | 1598 | 1636 | 1694 | 1671 |
| PESSIMISTIC | | | | | | | | | | |
| KNN | 224 | 461 | 719 | 960 | 1065 | 1127 | 1261 | 1357 | 1319 | 1671 |
| LogR | 935 | 932 | 1023 | 934 | 1015 | 1121 | 1350 | 1547 | 1587 | 1671 |
| NB | -148 | 226 | 325 | 674 | 867 | 1005 | 1110 | 1219 | 1127 | 1671 |
| SGDC | 216 | 425 | 655 | 788 | 967 | 1070 | 1124 | 1330 | 1438 | 1671 |
| SVC | -149 | 226 | 486 | 787 | 1165 | 1586 | 1626 | 1577 | 1590 | 1671 |
| RFC | 876 | 1017 | 1122 | 1271 | 1418 | 1495 | 1581 | 1581 | 1685 | 1671 |
| REFLECTIVE | | | | | | | | | | |
| KNN | -123 | -6 | 493 | 550 | 813 | 932 | 1140 | 1274 | 1399 | 1671 |
| LogR | 403 | 821 | 970 | 1022 | 1055 | 1236 | 1300 | 1398 | 1584 | 1671 |
| NB | 50 | 7 | 195 | 192 | 680 | 997 | 1273 | 998 | 1131 | 1671 |
| SGDC | 222 | 378 | 605 | 772 | 975 | 1076 | 1201 | 1374 | 1511 | 1671 |
| SVC | 170 | 257 | 368 | 444 | 490 | 816 | 1179 | 1462 | 1837 | 1671 |
| RFC | -55 | 276 | 667 | 897 | 1150 | 1464 | 1544 | 1500 | 1658 | 1671 |
| TCIA | | | | | | | | | | |
| KNN | 133 | 441 | 654 | 795 | 994 | 1220 | 1322 | 1305 | 1371 | 1671 |
| LogR | 103 | 60 | 88 | -96 | 399 | 962 | 1171 | 1573 | 1922 | 1671 |
| NB | 11 | 229 | 84 | 64 | 463 | 710 | 965 | 1221 | 1838 | 1671 |
| SGDC | 309 | 470 | 711 | 843 | 980 | 1065 | 1235 | 1388 | 1482 | 1671 |
| SVC | -17 | 261 | 1033 | 1190 | 1281 | 1677 | 1578 | 1498 | 1685 | 1671 |
| RFC | 249 | 423 | 642 | 802 | 961 | 1026 | 1085 | 1157 | 1454 | 1671 |
| TWO_MODEL | | | | | | | | | | |
| KNN | 321 | 732 | 1202 | 1576 | 1730 | 1775 | 1810 | 1760 | 1594 | 1671 |
| LogR | 864 | 1126 | 796 | 1002 | 1257 | 1059 | 1319 | 1503 | 1542 | 1671 |
| NB | 82 | 583 | 976 | 1183 | 1345 | 1364 | 1410 | 1347 | 1488 | 1671 |
| SGDC | 162 | 421 | 609 | 727 | 929 | 1010 | 1223 | 1378 | 1536 | 1671 |
| SVC | -49 | 62 | 250 | 900 | 1285 | 1655 | 1785 | 1881 | 1675 | 1671 |
| RFC | 877 | 1111 | 1097 | 1117 | 1165 | 1273 | 1437 | 1502 | 1641 | 1671 |

Multiple important findings can be derived from Table 8. First, we would like to emphasize the performance of RFC another time. In particular, we observe that RFC performs well with multiple uplift modeling strategies. For example, within CVT, RFC gets the largest number of wins across the deciles in terms of the cumulative number of purchases compared to the remaining base learners. The same conclusion can be drawn, e.g., for LGWUM and LWUM. RFC is especially successful in the first deciles. Given this, we recommend RFC for the suggestion of Lo (2002) to limit the targeting to the top 10% most valuable customers. However, the success of RFC can be interrupted in the middle deciles. For example, for the 4th decile, the pair CVT and KNN compared

to CVT and RFC gets 1,148 and 1,110 cumulative number of purchases, respectively. The pair CVT and SVC outperforms CVT-based RFC in the 6th and 7th deciles. Identical picture can be seen in terms of the pessimistic uplift modeling strategy, whereby SVC gets 1,586 and 1,626 cumulative number of purchases compared to 1,495 and 1,581 of RFC in the 6th and 7th deciles. Thus, we conclude that there are differences in the impact on business value depending on the size of the targeted fraction of the customer base. In general, we see the larger cumulative numbers of purchases in the middle deciles than in the first ones. To give an example, see a steady increase of cumulative purchases for the pair LWUM and SGDC from the first to the last decile. However, this does not indicate that targeting a larger fraction results in a higher cumulative number of purchases. See, for example, the pair of the two model uplift method and KNN in the 7th and 8th deciles (1,810 and 1,760 purchases, respectively). Therefore, our results clearly show that targeting the whole population of the customers – a mail-to-all strategy according to Chickering and Heckerman (2000) – is not the best choice. Most importantly, we now are confident to identify the best combination of base learner and uplift modeling strategy in terms of business value. These pairs are CVT and RFC, ITM and SVC, LGWUM and RFC, LWUM and RFC, pessimistic uplift modeling and RFC, reflective uplift modeling and LogR, TCIA and SVC, and finally the two model uplift method and KNN. They demonstrate the largest numbers of wins on the deciles. This finding is also supported in terms of Qini coefficient (see 6.1). In the following, therefore, we concentrate on these pairs.

To provide specific recommendations which pair works best, we now present uplift gain charts in Figure 2. These charts much resemble common gain charts. However, while the performance of models in gain charts in customer acquisition campaigns is typically illustrated by the number of purchases on the *y-axis*, uplift gain charts draft Qini curves that are by nature capable to signal incrementality. This implies that the number of purchases is replaced by the incremental number of purchases in uplift gains charts. The incremental number of purchases is a helpful indicator to support decision making in marketing practice and can be derived by comparing the purchase rate in the treatment group with the purchase rate in the control group. In both the traditional and uplift case, the purchase indicator is a function of the fraction of people targeted from the campaign’s total population, being mapped on the *x-axis* (Radcliffe, 2007). Qini curves summarize the decile-wise performance of their underlying uplift models. A diagonal line reflects random targeting and therefore presents a baseline for all strategy-based combinations. Recall that we present the uplift gain charts only for the winner pairs identified before. We also draw the average performance line – AVG – across the winner pairs to better judge on the performance.

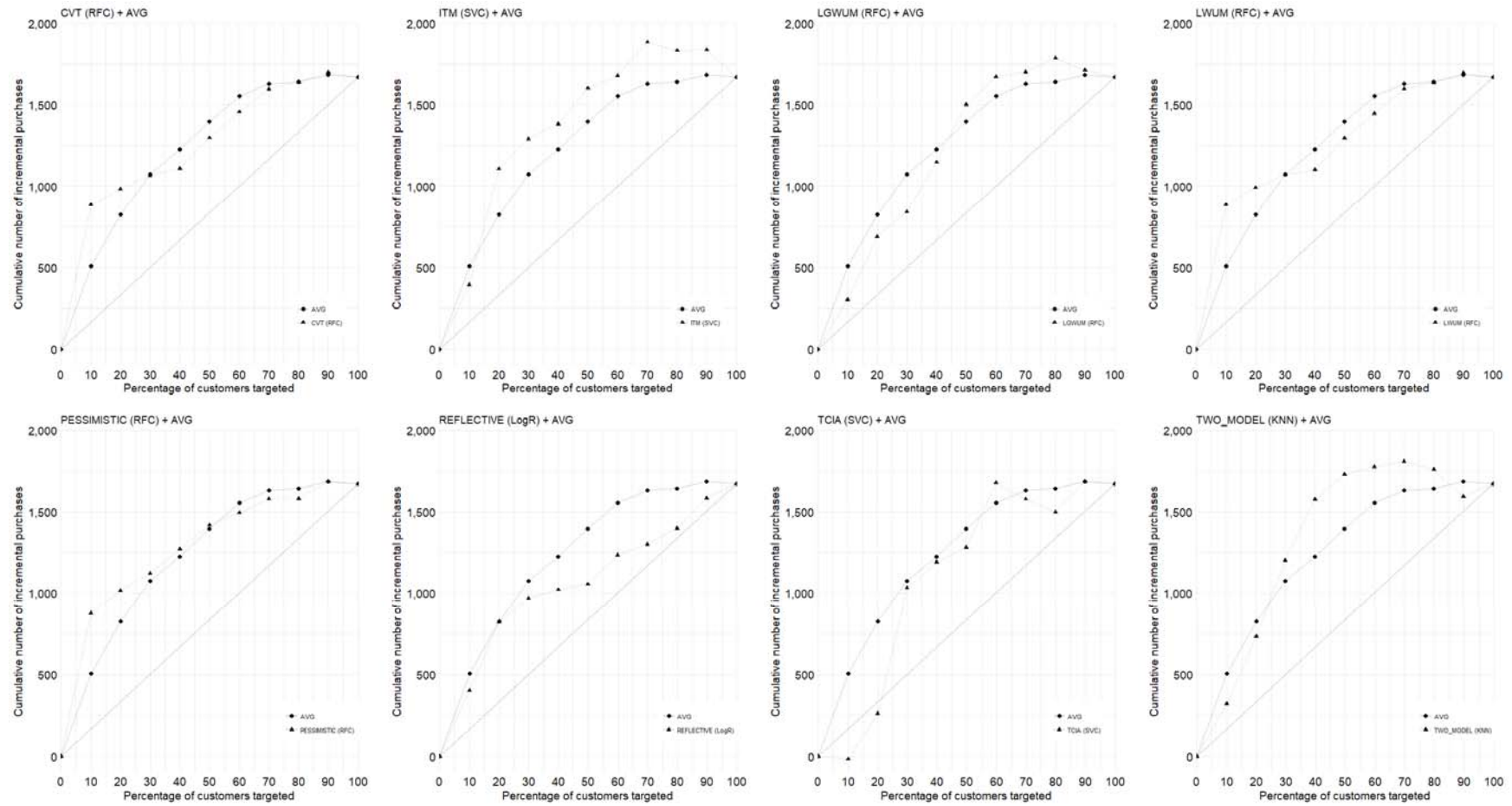


Figure 2 Uplift gain charts across uplift modeling strategies

Figure 2 provides new insights into the performance of the winner modeling pairs identified before. First, we observe that all pairs, even though unequally, contribute to higher cumulative number of purchases than the baseline. We see that the higher the fraction of customers targeted, the higher is the cumulative number of purchases. Every pair is capable to increase that number right from the beginning. Only TCIA coupled with SVC fails to achieve this. Second, we now clearly see that ITM coupled with SVC and the two model uplift method coupled with KNN outperform all other uplift modeling strategies. See, for example, that both couples perform better than the average performance starting from the 3rd decile. We also note that the performance of the two model uplift method paired with KNN deteriorates starting from the 7th decile and becomes even lower than the average rate in the 9th decile. This is not valid for the ITM-based SVC. However, we kindly remind that ITM-based SVC has shown extreme variance in the previous analysis (see 6.1). That is, we conclude that there are more signals in favor of the couple of the two model uplift method and KNN. This couple outperforms all other pairs (including ITM-based SVC) starting from 4th and ending with the 6th deciles. Third, we regard CVT, pessimistic uplift modeling, LGWUM, and LWUM as second-best choice since these strategies perform similar to the average. For example, pessimistic uplift modeling paired with RFC performs slightly better than the average in the first deciles, similar to average in the middle deciles, and underperforms in the last deciles. On the contrary, LGWUM coupled with RFC underperforms the average until the 5th decile and thereafter slightly outperforms the average. Fourth, we observe that combinations of reflective uplift modeling and LogR as well as TCIA and SVC show the weakest performance. Both are clearly inferior to the average. This is especially relevant for the pair of reflective uplift modeling and LogR, since we observe the underperformance in every decile. Thus, we cannot recommend adopting these modeling strategies for similar marketing campaigns. Given that RFC is the best choice in terms of base learners, Figure 2 suggests that it best performs coupled with the pessimistic uplift modeling strategy since it demonstrates until the 5th decile better or identical performance as average does; this is not given by other combinations.

To get more confidence in the findings obtained so far, we present the *non-cumulative* numbers of incremental purchases in the subsequent experiment. As before, the results are based on the out-of-sample test set. Figure 3 summarizes the respective results for the winner pairs on a decile-level.

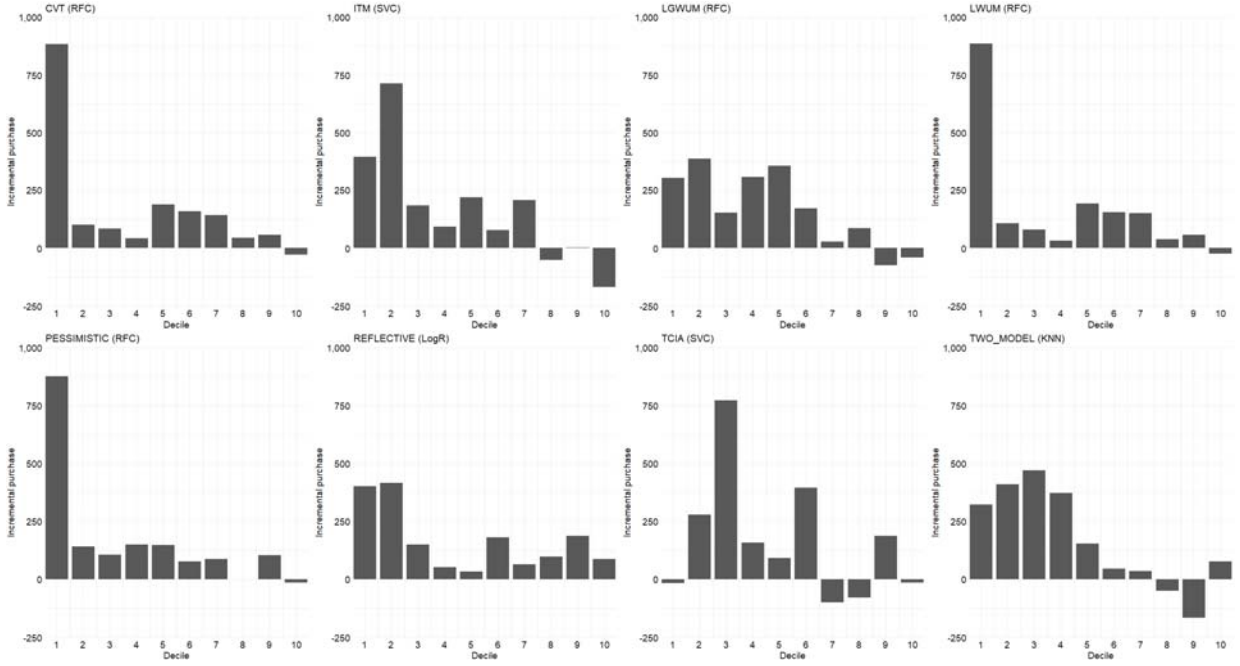


Figure 3 Non-cumulative numbers of incremental purchases

Figure 3 provides further findings. Given that truly valuable uplift models are capable to sort customers with high uplift to the first deciles and customers with comparably lower uplift or even negative to latter deciles (Kane et al., 2014), we first conclude that CVT, LWUM, and pessimistic uplift modeling perform quite well in the first decile. Second, comparing the winner pairs as per uplift gain charts – ITM coupled with SVC and the two model uplift method coupled with KNN – we now are more confident that there are more signals in favor of the latter pair. This is because the two model uplift method paired with KNN is able to assign customers who are likely to induce positive uplift to the first deciles and negative to the latter gradually. Although ITM-based SVC presents a powerful alternative achieving similar results, we observe that it assigns more customers in the latter deciles than two model uplift method-based KNN. See, for example, the 5th, 6th, and 7th deciles. Beyond this, we observe that ITM-based SVC assigns less customers in the 4th decile than in the 5th, less in the 6th than in the 7th, indicating unstable results. Third, Figure 3 provides more confidence in fact that the combinations TCIA and SVC as well as reflective uplift modeling and LogR present the least valuable alternatives. This is because the former allocates customers with negative uplift to the first decile and the latter presents a modeling strategy with no negative uplift in any decile. Given the shortcoming of the pair of TCIA and SVC in the first decile and the fact that it exhibits more variance in the latter deciles, we conclude that this pair presents the worst uplift modeling strategy considered in this study. However, we caution from execution of both methods.

6.3 Performance comparison between response and uplift modeling

Our final experiment is devoted to the examination of the performance of response modeling, a conventional method in marketing applications, vis-à-vis *the best* – two model uplift method paired with KNN – and *the worst* – TCIA paired with SVC – uplift modeling strategies. To provide a holistic picture on the performance of response modeling, we re-iterate all previous experiments, re-present the performance of the best and the worst uplift modeling strategies, and extend these experiments by the estimates obtained from response modeling. To secure fair empirical comparisons, we execute response modeling to the same out-of-sample test set for all experiments. We examine the interaction between the modeling strategies and involve Q_{pct} . Table 9 mimics the same setup for the interaction examination and adds response modeling to the modeling strategies for conversion uplift (see last row of the table).

Table 9 Qini coefficient of selected modeling strategies

| Modeling strategy | Base learner | | | | | |
|-------------------|--------------|--------|--------|-------|--------------|--------------|
| | KNN | LogR | NB | SGDC | SVC | RFC |
| TCIA | 1.043 | -1.950 | -2.821 | 1.222 | 3.893 | 0.403 |
| TWO MODEL | 7.267 | 4.305 | 3.297 | 0.688 | 2.806 | 5.401 |
| RESPONSE | 4.752 | 4.263 | 0.432 | 0.546 | 1.893 | 5.679 |

Table 9 shows that response modeling outperforms TCIA. That is because it achieves higher Q_{pct} values for two thirds of all base learners (i.e., KNN, LogR, NB, and RFC). Furthermore, we observe that the highest Q_{pct} value of response modeling coupled with RFC is higher than that of TCIA coupled with SVC, 5.679 and 3.893, respectively. This indicates that response modeling might be more beneficial than modern uplift modeling strategies. However, we also see that response modeling fails to outperform the two model uplift method. Apart from RFC, the two model uplift method is superior compared to response modeling in every combination. We observe that the interaction of the two model uplift method with KNN contributes to higher Q_{pct} value than the best combination of response modeling, 7.267 and 5.679, respectively. We understand that response modeling interacts best with RFC. This generalizes our finding that RFC is the winning base learner in terms of interaction with uplift modeling strategies. On the contrary, NB and SGDC show worst results when interacting with response modeling; finding that alerts to not execute these base learners for neither uplift nor response modeling.

Figure 4 presents the robustness procedure, aggregation of the results across the 10-fold cross validation, to judge about the variance in the results.

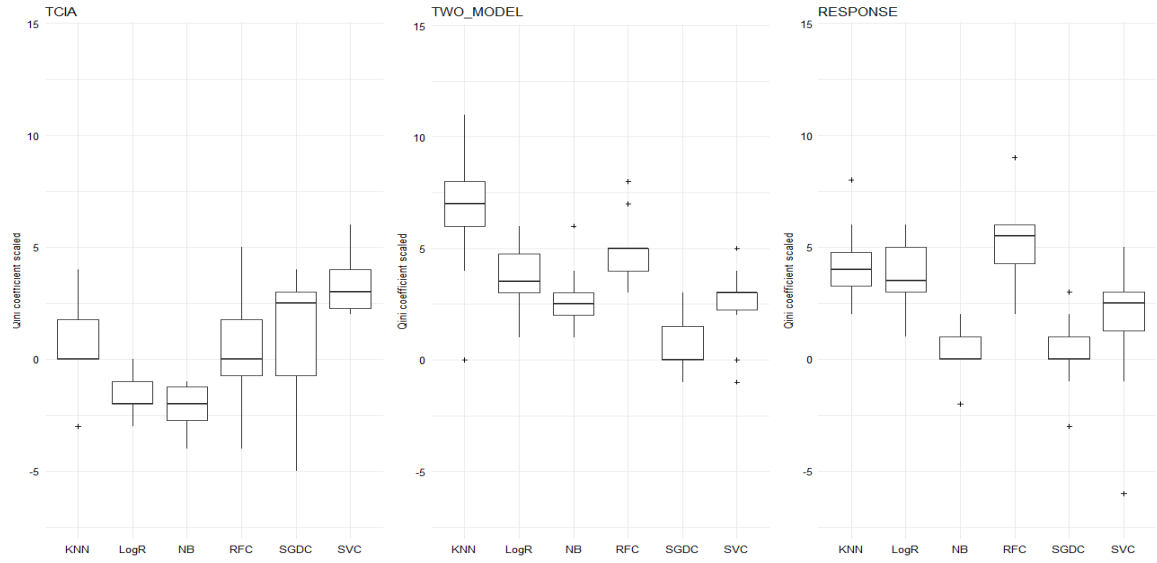


Figure 4 Scaled Qini of selected modeling strategies

Figure 4 illustrates that response modeling is superior to TCIA since it exhibits smaller variance in the estimates (see, for example, RFC or SGDC) and better interacts with NB and SGDC than TCIA does. We now also see that response modeling interacts with KNN and LogR comparably well to RFC and conclude that the former two base learners are promising when being paired with response modeling. Figure 4 also confirms that response modeling is inferior to the two model uplift method. We understand that the big share of NB, SGDC, and SVC estimates negative scaled values for Qini, while this is only the case for SGDC when paired with the two model uplift method (outliers not considered). Moreover, we observe that response modeling interacting with SVC and RFC exhibits higher variance than the two model uplift method with the same base learners.

We now examine the potential of response modeling to contribute to business value in terms of cumulative and non-cumulative incremental purchases. We echo the same experiments from 6.2 and extend them by the estimates of response modeling. First, we examine the tabular view of the cumulative number of incremental purchases. Recall that figures marked in *italic* and **bold face** indicate the same logic as in 6.2.

Table 10 Summary of cumulative number of incremental purchases

| Modeling strategy / base learner | Cumulative number of incremental purchases per decile | | | | | | | | | |
|--|---|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|------|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| TCIA | | | | | | | | | | |
| KNN | 133 | 441 | 654 | 795 | 994 | 1220 | 1322 | 1305 | 1371 | 1671 |
| LogR | 103 | 60 | 88 | -96 | 399 | 962 | 1171 | <i>1573</i> | 1668 | 1671 |
| NB | 11 | 229 | 84 | 64 | 463 | 710 | 965 | 1221 | <i>1709</i> | 1671 |
| SGDC | <i>309</i> | <i>470</i> | 711 | 843 | 980 | 1065 | 1235 | 1388 | 1582 | 1671 |
| SVC | -17 | 261 | <i>1033</i> | <i>1190</i> | <i>1281</i> | <i>1677</i> | <i>1578</i> | 1498 | 1685 | 1671 |
| RFC | 249 | 423 | 642 | 802 | 961 | 1026 | 1085 | 1157 | 1454 | 1671 |
| TWO MODEL | | | | | | | | | | |
| KNN | 321 | 732 | <i>1202</i> | <i>1576</i> | <i>1730</i> | <i>1775</i> | <i>1810</i> | 1760 | 1594 | 1671 |
| LogR | 864 | <i>1126</i> | 796 | 1002 | 1257 | 1059 | 1319 | 1503 | 1542 | 1671 |
| NB | 82 | 583 | 976 | 1183 | 1345 | 1364 | 1410 | 1347 | 1364 | 1671 |
| SGDC | 162 | 421 | 609 | 727 | 929 | 1010 | 1223 | 1378 | 1536 | 1671 |
| SVC | -49 | 62 | 250 | 900 | 1285 | 1655 | 1785 | <i>1881</i> | <i>1675</i> | 1671 |
| RFC | 877 | 1111 | 1097 | 1117 | 1165 | 1273 | 1437 | 1502 | 1641 | 1671 |
| RESPONSE | | | | | | | | | | |
| KNN | 295 | 626 | 933 | 1179 | <i>1417</i> | <i>1509</i> | <i>1612</i> | 1562 | 1641 | 1671 |
| LogR | <i>917</i> | <i>1014</i> | 733 | 1202 | 1154 | 1368 | 1287 | 1318 | 1445 | 1671 |
| NB | -206 | 214 | 442 | 623 | 1285 | 1197 | 1349 | 1358 | 1555 | 1671 |
| SGDC | 154 | 388 | 601 | 724 | 946 | 1004 | 1195 | 1363 | 1521 | 1671 |
| SVC | -87 | -34 | 406 | 784 | 1158 | 1586 | 1507 | <i>1684</i> | <i>1813</i> | 1671 |
| RFC | 897 | 853 | <i>1199</i> | <i>1307</i> | 1340 | 1393 | 1457 | 1475 | 1486 | 1671 |

Table 10 confirms the superiority of response modeling over TCIA in terms of business value. We observe that response modeling holds two global wins, i.e., in the first and in the 9th deciles (i.e., 917 and 1,813 cumulative incremental purchases, respectively), while TCIA none. However, response modeling is inferior to the two model uplift method, since the latter holds global wins starting from the 2nd and ending with the 8th deciles. Table 10 also reveals that response modeling interacts successfully with LogR, KNN, and SVC apart from RFC (see number of wins; marked in italic face). Although the pair of response modeling and RFC holds only two wins compared to three wins of the pair of response modeling and KNN, we conclude that the former is the best choice, since this finding is previously supported by the examination of Qini coefficient and robustness procedure. Therefore, we now examine the performance of this best pair compared to the other two best pairs. Recall that TCIA performs best with SVC and the two model uplift method with KNN. Figure 5 presents related uplift gain charts.

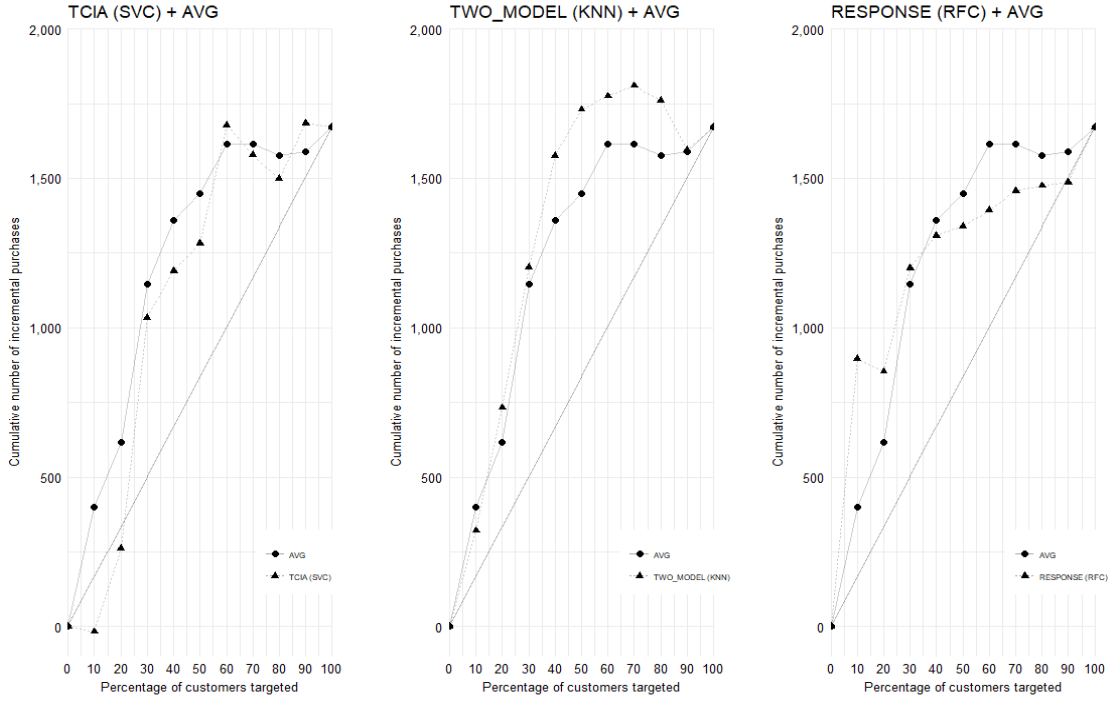


Figure 5 Uplift gain chart for response modeling

Figure 5 provides new insights. First, we see that response modeling is more successful in the first three deciles compared to the average. Recall that we now average the performance of only these three winner pairs. The performance of response modeling coupled with RFC deteriorates from the 4th decile. The pair TCIA and SVC outperforms response paired with RFC in the latter deciles. See, for example, the 7th, the 8th, and the 9th decile. Figure 5, thus, indicates that TCIA-based SVC might be more beneficial when contacting a larger fraction of customers than response-based RFC. Figure 5 also confirms that the two model uplift method coupled with KNN is superior over response modeling and RFC combination in every decile.

We now further examine the performance of the winning pairs as per non-cumulative number of incremental purchases. Figure 6 presents the corresponding results in bar charts. We mimic again the same experimental setup as in 6.2.

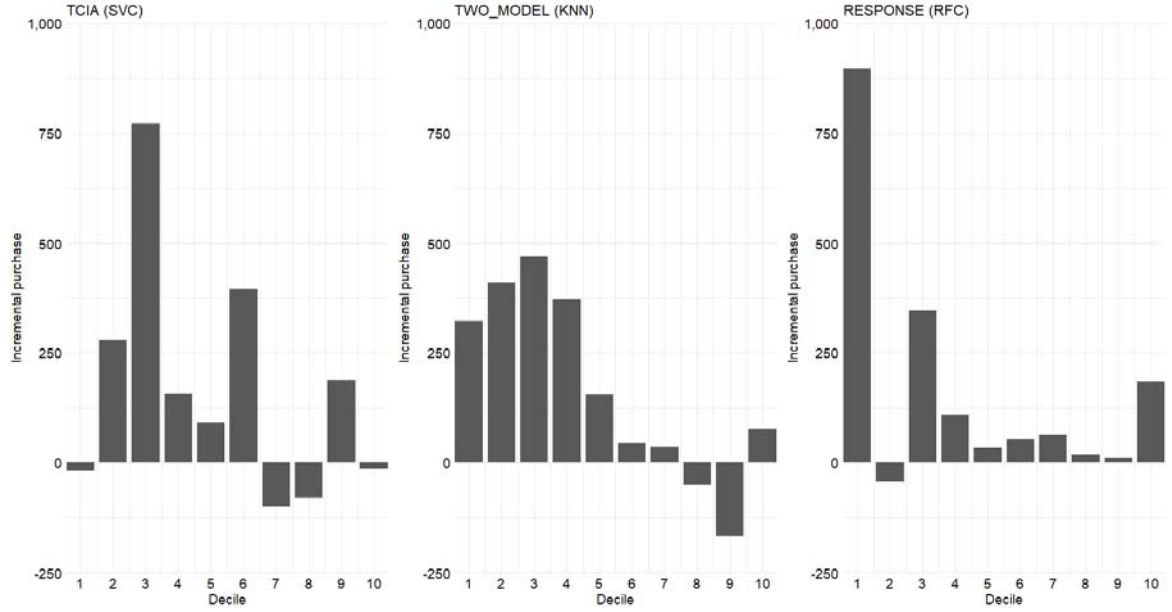


Figure 6 Non-cumulative number of incremental purchases of selected modeling strategies

Figure 6 provides the following insights. First, it becomes apparent that response-based RFC performs better than TCIA-based SVC in the first decile. However, we also observe that the former fails to perform in the 2nd decile. Recall that a good strategy aggregates a high number of non-cumulative purchases in the first deciles and small (or even negative) in the latter. Furthermore, response-based RFC fails to assign negative uplift in the latter deciles. We, therefore, conclude that response modeling paired with RFC represents a weak alternative for uplift modeling compared to the two model uplift method coupled with KNN.

After all, we would like to highlight two fundamental findings. First, response modeling which is usually practiced in marketing applications (Coussement et al., 2015) represents a powerful strategy that leads to success in such marketing campaigns that we describe in this study. We clearly see that it might outperform uplift modeling strategies that have been developed with the purpose to explain the causal relationship between marketing campaigns and an event of interest. And, second, most importantly, that response modeling might be inferior to selected uplift modeling strategies in many experimental dimensions. We, thus, conclude that our study makes it clear that marketers should be aware of the differences among the uplift modeling strategies and apply the best choice in real-world practice.

7 Conclusion

We set out to examine how different modeling strategies for conversion uplift contribute towards increasing the fit of marketing strategies for real-world applications. Uplift modeling can be seen as a technique that patterns causal effect of a marketing incentive on customer behavior.

Empirical examination goes alongside multiple dimensions and involves numerous data sets that stem from different geographies and represent distinct product lines. Given that uplift modeling strategies have been proposed in different strands of literature and no attempt has been made to systematically compare predictive performance of them, specific recommendations which strategies achieve highest relative performance have been missing. This study aims to close this research gap through multi-faceted experimentation.

Our study consolidates previous work in conversion uplift and provides a holistic picture of the state-of-the-art in predictive modeling for retail electronic commerce; more specifically, personalized marketing targeting through couponing. From an academic viewpoint, an important question is whether efforts invested to the development of novel uplift algorithms are worthwhile. Our study raises some critical concerns. We find the proposed method to generalize LWUM with weighted probability scores to account for the fraction of treatment and control group customers by Kane et al. (2014) fails to outperform the original LWUM developed by Lai et al. (2006) in terms of Qini coefficient. Similar picture is obtained in the field of covariates manipulation. We find that TCIA proposed by Tian et al. (2014), which to a large extent mimics the procedure of ITM, is inferior to original ITM developed by Lo (2002). On the contrary, we find that ITM as well as the straightforward two model uplift method (Radcliffe & Surry, 1999) that captures differences in class probabilities of customers' motivation represent modeling strategies of first choice for conversion uplift. Our study, therefore, implies that the progress has stalled, and efforts invested to the methodological advancement must be accompanied by a rigorous assessment of new uplift modeling strategies vis-à-vis challenging benchmark. We identify the two model uplift method and ITM as best performers according to our experiments and advise to compare novel modeling strategies in the field of uplift modeling against them.

An important question to answer in future research concerns the origins of the interaction between uplift modeling strategies and the underlying base learner. We have identified base learners that work specifically well for conversion uplift in digital marketing. However, our study does not seek to explain their success. We believe this is a fruitful avenue for future research; while our study may be regarded as a first move toward gaining insights to this question. For example, we find RFC to interact best with the majority of strategies. This is not given by other base learners. Moreover, RFC performs quite well in the first deciles of targeting and, therefore, can be strongly recommended for campaigns with little budget so that only the 10% most valuable customers are subject to treatment. We find SVC as a valid alternative, although it exhibits high variance in estimates as per robustness procedure presented in this paper. Surprisingly, KNN, usually seen as weak in predictive modeling, has shown appealing results, especially interacting with the two

model uplift method. On the contrary, SGDC and NB have shown poor results in every experiment. We, therefore, forewarn from considering these base learners for uplift modeling.

From a practitioner’s viewpoint, it is important to reason whether the observed results can be generalized to real world applications. On the one hand, we believe that numerous data sets from online shops, several cross-validation repetitions, and performance examination from different perspectives make our results relevant for the task of real-time targeting digital coupons in e-commerce. We also believe our main performance criterion, cumulative number of incremental purchases, to approximate the business value of an uplift model, which also raises the relevance of results from a practical point of view. However, we acknowledge that all data sets come from the same provider and exhibit similar features. Uplift models and base learners may behave differently when processing different feature sets. Consequently, we strongly encourage future research to study the behavior of uplift modeling strategies in other marketing and non-marketing applications using different feature sets. Without claiming external validity, our study may aid corresponding initiatives in pre-selecting promising and less promising modeling strategies. For example, we caution from using TCIA, the poor performance of which discourages its considerations in future experiments.

References

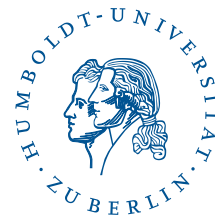
- E. Ascarza, P. S. Fader, and B. G. Hardie, *Marketing models for the customer-centric firm*, Handbook of Marketing Decision Models (Springer, New York, 2017), pp. 297-329.
- S. Athey, and G. W. Imbens, The State of Applied Econometrics: Causality and Policy Evaluation, *Journal of Economic Perspectives* **31**(2) (2017) 3-32.
- P. Baecke, and D. Van den Poel, Improving purchasing behavior predictions by data augmentation with situational variables, *International Journal of Information Technology & Decision Making* **9**(6) (2010) 853-872.
- B. Baesens, T. Van Gestel, S. Viaene, M. Stepanova, J. Suykens, and J. Vanthienen, Benchmarking state-of-the-art classification algorithms for credit scoring, *Journal of the Operational Research Society* **54**(6) (2003) 627-635.
- K. Bagchi, and S. Mukhopadhyay, Predicting global internet growth using augmented diffusion, fuzzy regression and neural network models, *International Journal of Information Technology & Decision Making* **5**(1) (2006) 155-171.
- Y. Bakos, The emerging role of electronic marketplaces on the Internet, *Communications of the ACM* **41**(8) (1998) 35-42.
- M. Ballings, and D. Van den Poel, CRM in social media: Predicting increases in Facebook usage frequency, *European Journal of Operational Research* **244**(1) (2015) 248-260.
- E. Braunwald, M. Domanski, S. Fowler, N. Geller, B. Gersh, J. Hsia, M. Pfeffer, M. Rice, Y. Rosenberg, and J. Rouleau, Angiotension-Converting-Enzyme Inhibition in Stable Coronary Artery Disease, *New England Journal of Medicine* **351** (2004) 2058-2068.
- L. Breiman, Random forests, *Machine Learning* **45**(1) (2001) 5-32.
- T. Cai, L. Tian, P. H. Wong, and L. J. Wei, Analysis of randomized comparative clinical trial data for personalized treatment selections, *Biostatistics* **12**(2) (2011) 270-282.
- Z. Chen, From data to behavior mining, *International Journal of Information Technology & Decision Making* **5**(4) (2006) 703-711.
- V. Chernozhukov, M. Demirer, E. Duflo, and I. Fernandez-Val, Generic machine learning inference on heterogenous treatment effects in randomized experiments, *Corr arXiv:1712.04802v3* (2018).
- D. M. Chickering, and D. Heckerman, A decision theoretic approach to targeted advertising, Proc. of the 16th Conf. on Uncertainty in Artificial Intelligence, Stanford, CA, USA, 2000 (Morgan Kaufmann Publishers Inc., pp. 82-88.
- S. R. Cole, and E. A. Stuart, Generalizing Evidence From Randomized Clinical Trials to Target Populations: The ACTG 320 Trial, *American Journal of Epidemiology* **172**(1) (2010) 107-115.
- A. F. J. Connors, T. Speroff, N. V. Dawson, C. Thomas, F. E. H. Jr, D. Wagner, N. Desbiens, L. Goldman, A. W. Wu, R. M. Califf, W. J. F. Jr, H. Vidaillet, S. Broste, P. Bellamy, J. Lynn, and W. A. Knaus, The effectiveness of right heart catheterization in the initial care of critically ill patients, *The Journal of the American Medical Association* **276**(11) (1996) 889-897.
- K. Coussement, P. Harrigan, and D. F. Benoit, Improving direct mail targeting through customer response modeling, *Expert Systems With Applications* **42**(22) (2015) 8403-8412.
- N. Daskalova, F. R. Bentley, and N. Andalibi, It's All About Coupons: Exploring Coupon Use Behaviors in Email, Proc. of the 2017 CHI Conf. on Human Factors in Computing Systems, 2017 (ACM, New York), pp. 1152-1160.
- F. De Vriendt, D. Moldovan, and W. Verbeke, A Literature Survey and Experimental Evaluation of the State-of-the-Art in Uplift Modeling: A Stepping Stone Toward the Development of Prescriptive Analytics, *Big Data* **6**(1) (2018) 13-41.
- D. Dheeru, and E. Karra Taniskidou, UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml/>] (2017)
- A. W. Ding, S. Li, and P. Chatterjee, Learning user real-time intent for optimal dynamic web page transformation, *Information Systems Research* **26**(2) (2015) 339-359.
- F. Dost, R. Wilken, M. Eisenbeiss, and B. Skiera, On the edge of buying: A targeting approach for indecisive buyers based on willingness-to-pay ranges, *Journal of Retailing* **90**(3) (2014) 393-407.
- D. P. Green, and A. S. Gerber, *Get Out the Vote: How to Increase Voter Turnout* (Brookings Institution Press, 2015).
- L. Guelman. (2014). *Optimal Personalized Treatment Learning Models with Insurance Applications*. PhD in Economics, Universitat de Barcelona.
- L. Guelman, M. Guillén, and A. M. Pérez-Marín, *Random Forests for Uplift Modeling: An Insurance Customer Retention Case*, (eds.) K. J. Engemann, A. M. Gil-Lafuente & J. M. Merigó, Modeling and Simulation in Engineering, Economics and Management: International Conference, MS 2012, New Rochelle, NY, USA, May 30 - June 1, 2012. Proceedings (Springer Berlin Heidelberg, Berlin, Heidelberg, 2012), pp. 123-133.
- L. Guelman, M. Guillén, and A. M. Pérez-Marín, Uplift Random Forests, *Cybernetics and Systems* **46**(3-4) (2015) 230-248.
- L. Guelman, M. Guillén, and A. M. Pérez Marín, Optimal personalized treatment rules for marketing interventions: A review of methods, a new proposal, and an insurance case study, *UB Riskcenter Working Paper Series, 2014/06* (2014).
- B. B. Hansen, and J. Bowers, Covariate balance in simple, stratified and clustered comparative studies, *Statistical Science* **23**(2) (2008) 219-236.
- B. Hansotia, and B. Rukstales, Direct marketing for multichannel retailers: Issues, challenges and solutions, *Journal of Database Marketing & Customer Strategy Management* **9**(3) (2002a) 259-266.
- B. Hansotia, and B. Rukstales, Incremental value modeling, *Journal of Interactive Marketing* **16**(3) (2002b) 35-46.
- K. Hillstrom, MineThatData (2008), <https://blog.minethatdata.com/2008/03/minethatdata-e-mail-analytics-and-data.html>
- P. W. Holland, Statistics and causal inference, *Journal of the American Statistical Association* **81**(396) (1986) 945-960.
- S. Hua, What makes underwriting and non-underwriting clients of brokerage firms receive different recommendations? An application of uplift random forest model, *International Journal of Finance & Banking Studies* **5**(3) (2016) 42-56.
- E. Y. Huang, and C.-j. Tsui, Assessing customer retention in B2C electronic commerce: An empirical study, *Journal of Marketing Analytics* **4**(4) (2016) 172-185.

- M. Ieva, F. De Canio, and C. Ziliani, Daily deal shoppers: What drives social couponing?, *Journal of Retailing and Consumer Services* **40** (2018) 299-303.
- K. Imai, and M. Ratkovic, Estimating treatment effect heterogeneity in randomized program evaluation, *The Annals of Applied Statistics* **7**(1) (2013) 443-470.
- S. Jaroszewicz, and P. Rzepakowski, Uplift Modeling With Survival Data, ACM SIGKDD Workshop on Health Informatics (HI KDD'14), New York, USA, (2014) (ACM).
- S. Jaroszewicz, and Ł. Zaniewicz, *Székel Regularization for Uplift Modeling*, (eds.) S. Matwin & J. Mielniczuk, Challenges in Computational Statistics and Data Mining (Springer International Publishing, Switzerland, 2016), pp. 135-154.
- M. Jaskowski, and S. Jaroszewicz, Uplift Modeling for Clinical Trial Data, ICML 2012 Workshop on Clinical Data Analysis, Edinburgh, Scotland, 2012, pp.
- K. Kane, S. Y. V. Lo, and J. Zheng, Mining for the truly responsive customers and prospects using true-lift modeling: Comparison of new and existing methods, *Journal of Marketing Analytics* **2**(4) (2014) 218-238.
- S. P. Kondareddy, S. Agrawal, and S. Shekhar, Incremental Response Modeling Based on Segmentation Approach Using Uplift Decision Trees, Industrial Conference on Data Mining, 2016 (Springer, Berlin), pp. 54-63.
- F. Kuusisto, V. S. Costa, H. Nassif, E. Burnside, D. Page, and J. Shavlik, *Support Vector Machines for Differential Prediction*, (eds.) T. Calders, F. Esposito, E. Hüllermeier & R. Meo, Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2014, Nancy, France, September 15-19, 2014. Proceedings, Part II (Springer Berlin Heidelberg, Berlin, Heidelberg, 2014), pp. 50-65.
- Y.-T. Lai, K. Wang, D. Ling, H. Shi, and J. Zhang, Direct Marketing When There Are Voluntary Buyers, Proceedings of the 6th International Conference on Data Mining (ICDM), Hong Kong, China, 2006 (IEEE Computer Society, Washington, DC, USA), pp. 922-927.
- K. Larsen. (2010). Net Lift Models. Slides of a talk given at the A2010 - Analytics Conference, September 2-3, Copenhagen, Denmark.
- J. Lee, D.-H. Park, and I. Han, The effect of negative online consumer reviews on product attitude: An information processing view, *Electronic Commerce Research and Applications* **7**(3) (2008) 341-352.
- S. Lessmann, B. Baesens, H.-V. Seow, and L. C. Thomas, Benchmarking state-of-the-art classification algorithms for credit scoring: A ten-year update, *European Journal of Operational Research* **247**(1) (2015) 124-136.
- S. Lessmann, K. Coussement, K. W. D. Bock, and J. Haupt, Targeting customers for profit: An ensemble learning framework to support marketing decision making (2018), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3130661
- V. S. Lo, The true lift model: a novel data mining approach to response modeling in database marketing, *ACM SIGKDD Explorations Newsletter* **4**(2) (2002) 78-86.
- V. S. Lo, and D. A. Pachamanova, From predictive uplift modeling to prescriptive uplift analytics: A practical approach to treatment optimization while accounting for estimation risk, *Journal of Marketing Analytics* **3**(2) (2015) 79-95.
- S. Loi, B. Haibe-Kains, C. Desmedt, F. Lallemand, A. M. Tutt, C. Gillet, P. Ellis, A. Harris, J. Bergh, J. A. Foekens, J. G. Klijn, D. Larsimont, M. Buyse, G. Bontempi, M. Delorenzi, M. J. Piccart, and C. Sotiriou, Definition of Clinically Distinct Molecular Subtypes in Estrogen Receptor-Positive Breast Carcinomas Through Genomic Grade, *Journal of Clinical Oncology* **25** (2007) 1239-1246.
- C. Manahan, A proportional hazards approach to campaign list selection, *Proc. of the SAS User Group Intern. Meeting* (2005).
- G. Melli, X. Wu, P. Beinat, F. Bonchi, L. Cao, R. Duan, C. Faloutsos, R. Ghani, B. Kitts, B. Goethals, G. McLachlan, J. Pei, A. Srivastava, and O. Zaïane, Top-10 data mining case studies, *International Journal of Information Technology & Decision Making* **11**(2) (2012) 389-400.
- N. Michaelidou, and S. Dibb, Using email questionnaires for research: Good practice in tackling non-response, *Journal of Targeting, Measurement and Analysis for Marketing* **14**(4) (2006) 289-296.
- R. Michel, I. Schnakenburg, and T. von Martens, Effective customer selection for marketing campaigns based on net scores, *Journal of Research in Interactive Marketing* **11**(1) (2017) 2-15.
- H. Nassif, F. Kuusisto, E. S. Burnside, D. Page, J. W. Shavlik, and V. S. Costa, Score As You Lift (SAYL): A Statistical Relational Learning Approach to Uplift Modeling, Machine Learning and Knowledge Discovery in Databases - European Conference, ECML/PKDD, Prague, Czech Republic, 2013a (Springer, pp. 595-611.
- H. Nassif, F. Kuusisto, E. S. Burnside, and J. W. Shavlik, Uplift Modeling with ROC: An SRL Case Study, Late Breaking Papers of the 23rd International Conference on Inductive Logic Programming (ILP'13), Rio de Janeiro, Brazil, 2013b, pp. 40-45.
- H. Nassif, D. Page, M. Ayvaci, J. Shavlik, and E. S. Burnside, Uncovering age-specific invasive and DCIS breast cancer rules using inductive logic programming, Proceedings of the 1st ACM International Health Informatics Symposium, Arlington, Virginia, USA, (2010) (ACM, 1883005).
- H. Nassif, Y. Wu, D. Page, and E. S. Burnside, Logical Differential Prediction Bayes Net. Improving Breast Cancer Diagnosis for Older Women, American Medical Informatics Association Symposium (AMIA), Chicago, 2012, pp. 1330-1339.
- J. O. Ogutu, H.-P. Piepho, and T. Schulz-Streeck, A comparison of random forests, boosting and support vector machines for genomic selection, BMC proceedings, 2011 (BioMed Central, pp. S11.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and É. Duchesnay, Scikit-learn: Machine Learning in Python, *Journal of Machine Learning Research* **12** (2011) 2825-2830.
- M. Pintilie, *Competing Risks: A Practical Perspective* (John Wiley & Sons, Ltd Chichester, 2006).
- N. J. Radcliffe, Using control groups to target on predicted lift: Building and assessing uplift models, *Direct Marketing Analytics Journal* (2007) 14-21.
- N. J. Radcliffe, and P. D. Surry, Differential Response Analysis: Modeling True Responses by Isolating the Effect of a Single Action, Proceedings of Credit Scoring and Credit Control IV, Edinburgh, Scotland, (1999) (Credit Research Centre, University of Edinburgh Management School.

- N. J. Radcliffe, and P. D. Surry. (2011). Real-World Uplift Modelling with Significance-Based Uplift Trees. Portrait Technical Report, TR-2011-1.
- T. B. Rhouma, and G. Zaccour, Optimal marketing strategies for the acquisition and retention of service subscribers, *Management Science* **64**(6) (2018) 2609-2627.
- P. Rzepakowski, and S. Jaroszewicz, Decision trees for uplift modeling with single and multiple treatments, *Knowledge and Information Systems* **32**(2) (2012a) 303-327.
- P. Rzepakowski, and S. Jaroszewicz, Uplift modeling in direct marketing, *Journal of Telecommunications and Information Technology* **2** (2012b) 43-50.
- N. S. Sahni, D. Zou, and P. K. Chintagunta, Do targeted discount offers serve as advertising? Evidence from 70 field experiments, *Management Science* **63**(8) (2016) 2688-2705.
- A. Shaar, T. Abdessalem, and O. Segard, Pessimistic Uplift Modeling, *CoRR* **abs/1603.09738** (2016).
- M. Softys, S. Jaroszewicz, and P. Rzepakowski, Ensemble methods for uplift modeling, *Data Mining and Knowledge Discovery* **29**(6) (2015) 1531-1559.
- X. Su, J. Kang, J. Fan, R. A. Levine, and X. Yan, Facilitating score and causal inference trees for large observational studies, *Journal of Machine Learning Research* **13** (2012) 2955-2994.
- L. Tian, A. A. Alizadeh, A. J. Gentles, and R. Tibshirani, A simple method for estimating interactions between a treatment and a large number of covariates, *Journal of the American Statistical Association* **109**(508) (2014) 1517-1532.
- D. Van den Poel, and W. Buckinx, Predicting online-purchasing behaviour, *European Journal of Operational Research* **166**(2) (2005) 557-575.
- F. H.-L. Yong. (2015). *Quantitative methods for stratified medicine*. Doctoral dissertation, Harvard University.
- L. Zaniewicz, and S. Jaroszewicz, Support vector machines for uplift modeling, Proc. of the 13th IEEE Intern. Cong. on Data Mining Workshops (ICDMW), 2013 (IEEE, Piscataway), pp. 131-138.
- D. Zantedeschi, E. M. Feit, and E. T. Bradlow, Measuring multichannel advertising response, *Management Science* **63**(8) (2016) 2706-2728.
- L. Zhao, and J. Zhu, Internet marketing budget allocation from practitioner's perspective, *International Journal of Information Technology & Decision Making* **9**(5) (2010) 779-797.

IRTG 1792 Discussion Paper Series 2018

For a complete list of Discussion Papers published, please visit
irtg1792.hu-berlin.de.



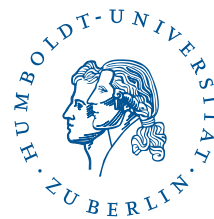
- 001 "Data Driven Value-at-Risk Forecasting using a SVR-GARCH-KDE Hybrid" by Marius Lux, Wolfgang Karl Härdle and Stefan Lessmann, January 2018.
- 002 "Nonparametric Variable Selection and Its Application to Additive Models" by Zheng-Hui Feng, Lu Lin, Ruo-Qing Zhu and Li-Xing Zhu, January 2018.
- 003 "Systemic Risk in Global Volatility Spillover Networks: Evidence from Option-implied Volatility Indices " by Zihui Yang and Yinggang Zhou, January 2018.
- 004 "Pricing Cryptocurrency options: the case of CRIX and Bitcoin" by Cathy YH Chen, Wolfgang Karl Härdle, Ai Jun Hou and Weining Wang, January 2018.
- 005 "Testing for bubbles in cryptocurrencies with time-varying volatility" by Christian M. Hafner, January 2018.
- 006 "A Note on Cryptocurrencies and Currency Competition" by Anna Almosova, January 2018.
- 007 "Knowing me, knowing you: inventor mobility and the formation of technology-oriented alliances" by Stefan Wagner and Martin C. Goossen, February 2018.
- 008 "A Monetary Model of Blockchain" by Anna Almosova, February 2018.
- 009 "Deregulated day-ahead electricity markets in Southeast Europe: Price forecasting and comparative structural analysis" by Antanina Hryshchuk, Stefan Lessmann, February 2018.
- 010 "How Sensitive are Tail-related Risk Measures in a Contamination Neighbourhood?" by Wolfgang Karl Härdle, Chengxiu Ling, February 2018.
- 011 "How to Measure a Performance of a Collaborative Research Centre" by Alona Zharova, Janine Tellingner-Rice, Wolfgang Karl Härdle, February 2018.
- 012 "Targeting customers for profit: An ensemble learning framework to support marketing decision making" by Stefan Lessmann, Kristof Coussement, Koen W. De Bock, Johannes Haupt, February 2018.
- 013 "Improving Crime Count Forecasts Using Twitter and Taxi Data" by Lara Vomfell, Wolfgang Karl Härdle, Stefan Lessmann, February 2018.
- 014 "Price Discovery on Bitcoin Markets" by Paolo Pagnottoni, Dirk G. Baur, Thomas Dimpfl, March 2018.
- 015 "Bitcoin is not the New Gold - A Comparison of Volatility, Correlation, and Portfolio Performance" by Tony Klein, Hien Pham Thu, Thomas Walther, March 2018.
- 016 "Time-varying Limit Order Book Networks" by Wolfgang Karl Härdle, Shi Chen, Chong Liang, Melanie Schienle, April 2018.
- 017 "Regularization Approach for Network Modeling of German EnergyMarket" by Shi Chen, Wolfgang Karl Härdle, Brenda López Cabrera, May 2018.
- 018 "Adaptive Nonparametric Clustering" by Kirill Efimov, Larisa Adamyan, Vladimir Spokoiny, May 2018.
- 019 "Lasso, knockoff and Gaussian covariates: a comparison" by Laurie Davies, May 2018.

IRTG 1792, Spandauer Straße 1, D-10178 Berlin
<http://irtg1792.hu-berlin.de>

This research was supported by the Deutsche
Forschungsgemeinschaft through the IRTG 1792.

IRTG 1792 Discussion Paper Series 2018

For a complete list of Discussion Papers published, please visit
irtg1792.hu-berlin.de.



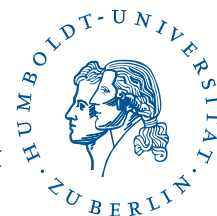
- 020 "A Regime Shift Model with Nonparametric Switching Mechanism" by Haiqiang Chen, Yingxing Li, Ming Lin and Yanli Zhu, May 2018.
- 021 "LASSO-Driven Inference in Time and Space" by Victor Chernozhukov, Wolfgang K. Härdle, Chen Huang, Weining Wang, June 2018.
- 022 "Learning from Errors: The case of monetary and fiscal policy regimes" by Andreas Tryphonides, June 2018.
- 023 "Textual Sentiment, Option Characteristics, and Stock Return Predictability" by Cathy Yi-Hsuan Chen, Matthias R. Fengler, Wolfgang Karl Härdle, Yanchu Liu, June 2018.
- 024 "Bootstrap Confidence Sets For Spectral Projectors Of Sample Covariance" by A. Naumov, V. Spokoiny, V. Ulyanov, June 2018.
- 025 "Construction of Non-asymptotic Confidence Sets in 2 -Wasserstein Space" by Johannes Ebert, Vladimir Spokoiny, Alexandra Suvorikova, June 2018.
- 026 "Large ball probabilities, Gaussian comparison and anti-concentration" by Friedrich Götze, Alexey Naumov, Vladimir Spokoiny, Vladimir Ulyanov, June 2018.
- 027 "Bayesian inference for spectral projectors of covariance matrix" by Igor Silin, Vladimir Spokoiny, June 2018.
- 028 "Toolbox: Gaussian comparison on Euclidian balls" by Andzhey Koziuk, Vladimir Spokoiny, June 2018.
- 029 "Pointwise adaptation via stagewise aggregation of local estimates for multiclass classification" by Nikita Puchkin, Vladimir Spokoiny, June 2018.
- 030 "Gaussian Process Forecast with multidimensional distributional entries" by Francois Bachoc, Alexandra Suvorikova, Jean-Michel Loubes, Vladimir Spokoiny, June 2018.
- 031 "Instrumental variables regression" by Andzhey Koziuk, Vladimir Spokoiny, June 2018.
- 032 "Understanding Latent Group Structure of Cryptocurrencies Market: A Dynamic Network Perspective" by Li Guo, Yubo Tao and Wolfgang Karl Härdle, July 2018.
- 033 "Optimal contracts under competition when uncertainty from adverse selection and moral hazard are present" by Natalie Packham, August 2018.
- 034 "A factor-model approach for correlation scenarios and correlation stress-testing" by Natalie Packham and Fabian Woebbecking, August 2018.
- 035 "Correlation Under Stress In Normal Variance Mixture Models" by Michael Kalkbrener and Natalie Packham, August 2018.
- 036 "Model risk of contingent claims" by Nils Detering and Natalie Packham, August 2018.
- 037 "Default probabilities and default correlations under stress" by Natalie Packham, Michael Kalkbrener and Ludger Overbeck, August 2018.
- 038 "Tail-Risk Protection Trading Strategies" by Natalie Packham, Jochen Papenbrock, Peter Schwendner and Fabian Woebbecking, August 2018.

IRTG 1792, Spandauer Straße 1, D-10178 Berlin
<http://irtg1792.hu-berlin.de>

This research was supported by the Deutsche
Forschungsgemeinschaft through the IRTG 1792.

IRTG 1792 Discussion Paper Series 2018

For a complete list of Discussion Papers published, please visit irtg1792.hu-berlin.de.



- 039 "Penalized Adaptive Forecasting with Large Information Sets and Structural Changes" by Lenka Zbonakova, Xinjue Li and Wolfgang Karl Härdle, August 2018.
- 040 "Complete Convergence and Complete Moment Convergence for Maximal Weighted Sums of Extended Negatively Dependent Random Variables" by Ji Gao YAN, August 2018.
- 041 "On complete convergence in Marcinkiewicz-Zygmund type SLLN for random variables" by Anna Kuczmazewska and Ji Gao YAN, August 2018.
- 042 "On Complete Convergence in Marcinkiewicz-Zygmund Type SLLN for END Random Variables and its Applications" by Ji Gao YAN, August 2018.
- 043 "Textual Sentiment and Sector specific reaction" by Elisabeth Bommers, Cathy Yi-Hsuan Chen and Wolfgang Karl Härdle, September 2018.
- 044 "Understanding Cryptocurrencies" by Wolfgang Karl Härdle, Campbell R. Harvey, Raphael C. G. Reule, September 2018.
- 045 "Predicative Ability of Similarity-based Futures Trading Strategies" by Hsin-Yu Chiu, Mi-Hsiu Chiang, Wei-Yu Kuo, September 2018.
- 046 "Forecasting the Term Structure of Option Implied Volatility: The Power of an Adaptive Method" by Ying Chen, Qian Han, Linlin Niu, September 2018.
- 047 "Inferences for a Partially Varying Coefficient Model With Endogenous Regressors" by Zongwu Cai, Ying Fang, Ming Lin, Jia Su, October 2018.
- 048 "A Regime Shift Model with Nonparametric Switching Mechanism" by Haiqiang Chen, Yingxing Li, Ming Lin, Yanli Zhu, October 2018.
- 049 "Strict Stationarity Testing and GLAD Estimation of Double Autoregressive Models" by Shaojun Guo, Dong Li, Mui Li, October 2018.
- 050 "Variable selection and direction estimation for single-index models via DC-TGDR method" by Wei Zhong, Xi Liu, Shuangge Ma, October 2018.
- 051 "Property Investment and Rental Rate under Housing Price Uncertainty: A Real Options Approach" by Honglin Wang, Fan Yu, Yinggang Zhou, October 2018.
- 052 "Nonparametric Additive Instrumental Variable Estimator: A Group Shrinkage Estimation Perspective" by Qingliang Fan, Wei Zhong, October 2018.
- 053 "The impact of temperature on gaming productivity: evidence from online games" by Xiaojia Bao, Qingliang Fan, October 2018.
- 054 "Topic Modeling for Analyzing Open-Ended Survey Responses" by Andra-Selina Pietsch, Stefan Lessmann, October 2018.
- 055 "Estimation of the discontinuous leverage effect: Evidence from the NASDAQ order book" by Markus Bibinger, Christopher Neely, Lars Winkelmann, October 2018.
- 056 "Cryptocurrencies, Metcalfe's law and LPPL models" by Daniel Traian Pele, Miruna Mazurencu-Marinescu-Pele, October 2018.
- 057 "Trending Mixture Copula Models with Copula Selection" by Bingduo Yang, Zongwu Cai, Christian M. Hafner, Guannan Liu, October 2018.

IRTG 1792, Spandauer Straße 1, D-10178 Berlin
<http://irtg1792.hu-berlin.de>

This research was supported by the Deutsche
Forschungsgemeinschaft through the IRTG 1792.

IRTG 1792 Discussion Paper Series 2018

For a complete list of Discussion Papers published, please visit irtg1792.hu-berlin.de.



- 058 "Investing with cryptocurrencies – evaluating the potential of portfolio allocation strategies" by Alla Petukhina, Simon Trimborn, Wolfgang Karl Härdle, Hermann Elendner, October 2018.
- 059 "Towards the interpretation of time-varying regularization parameters in streaming penalized regression models" by Lenka Zbonakova, Ricardo Pio Monti, Wolfgang Karl Härdle, October 2018.
- 060 "Residual's Influence Index (R_{infin}), Bad Leverage And Unmasking In High Dimensional L2-Regression" by Yannis G. Yatracos, October 2018.
- 061 "Plug-In L2-Upper Error Bounds In Deconvolution, For A Mixing Density Estimate In \mathbb{R}^d And For Its Derivatives" by Yannis G. Yatracos, October 2018.
- 062 "Conversion uplift in e-commerce: A systematic benchmark of modeling strategies" by Robin Gubela, Artem Bequé, Fabian Gebert and Stefan Lessmann, November 2018.

IRTG 1792, Spandauer Straße 1, D-10178 Berlin
<http://irtg1792.hu-berlin.de>

This research was supported by the Deutsche
Forschungsgemeinschaft through the IRTG 1792.