

Cross Channel Optimized Marketing by Reinforcement Learning

Naoki Abe, Naval Verma and Chid Apte
Mathematical Sciences Dept.
IBM T. J. Watson Res. Ctr.
Yorktown Heights, NY 10598
nabe, nverma, apte@us.ibm.com

Robert Schroko
Database Marketing
Saks Fifth Avenue
12 E. 49th Street, New York, NY 10017
Robert.Schroko@s5a.com

ABSTRACT

The issues of cross channel integration and customer life time value modeling are two of the most important topics surrounding customer relationship management (CRM) today. In the present paper, we describe and evaluate a novel solution that treats these two important issues in a unified framework of Markov Decision Processes (MDP). In particular, we report on the results of a joint project between IBM Research and Saks Fifth Avenue to investigate the applicability of this technology to real world problems. The business problem we use as a testbed for our evaluation is that of optimizing direct mail campaign mailings for maximization of profits in the store channel. We identify a problem common to cross-channel CRM, which we call *the Cross-Channel Challenge*, due to the lack of explicit linking between the marketing actions taken in one channel and the customer responses obtained in another. We provide a solution for this problem based on old and new techniques in reinforcement learning. Our in-laboratory experimental evaluation using actual customer interaction data show that as much as 7 to 8 per cent increase in the store profits can be expected, by employing a mailing policy automatically generated by our methodology. These results confirm that our approach is valid in dealing with the cross channel CRM scenarios in the real world.

Categories and Subject Descriptors

I.2.6 [Artificial Intelligence]: Learning

General Terms

Verification

Keywords

customer life time value, CRM, cost sensitive learning, reinforcement learning, targeted marketing

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

1. INTRODUCTION

The issues of cross channel integration and customer life time value modeling are undoubtedly two of the most important topics surrounding customer relationship management (CRM) today. Despite the wide-spread acknowledgement of their importance, there has not been a satisfactory solution to these issues in the market. In many cases, vendors provide infrastructure that enables cross channel integration of customer data, but provide no special-purpose analytics. Applying existing data analytics tools will not fully leverage the integrated data. This is because existing tools do not help optimize decision making in the respective channels for maximization of future profits across different channels. In the present paper, we propose a novel solution that treats these two important issues in a unified framework.

The proposed solution is based on our earlier work on the application of reinforcement learning to sequential cost-sensitive decision making [8]. In that paper, it was demonstrated that by combining reinforcement learning and scalable data mining technologies, decision rules that are optimized with respect to long term benefits can be automatically generated solely from data analytics. We use variants of this basic technology to perform customer life time value modeling in a cross-channel setting, and optimize marketing actions with respect to long term, multi-channel profits.

We have conducted a joint study between IBM Research and Saks Fifth Avenue to investigate the applicability of this technology on a practical problem, which we will report on in the current paper. The business problem we elected as a testbed for our investigation is that of optimizing interactions between the direct mail channel and the store channel. Of the various channels of customer interactions that Saks Fifth Avenue owns and operates, such as the web, telemarketing, direct mailing and the store, we chose to focus on the interactions between the latter two channels, based on the relative ease of evaluation and readiness of the relevant data.

The largest obstacle that we faced in our effort is one that is characteristic of a cross-channel scenario, and is of general interest to most applications involving cross-channel interactions. The problem, which we call *the Cross Channel Challenge*, is the lack of explicit linking between the marketing actions taken in one channel and the responses (profits or rewards) obtained in another. This translates, in practice, to very low correlations observed between marketing actions and their effects across channels. Therefore, applying off-the-shelf regression methods to model the rewards as

a function of various variables, is likely to produce a model that is independent of the marketing actions, thus leading to useless marketing rules.

We resolve this challenge by invoking the reinforcement learning technology, and devising a number of modifications to it. One of these modifications, based on an existing method of reinforcement learning called *advantage updating* [3], is particularly notable. It manages to solve the cross channel challenge by focusing on learning the *difference* in the effects on the rewards of competing actions, thereby bypassing the accurate estimation of the noisy reward function.

We conducted experimental evaluation of the proposed methods using actual customer interaction data from Saks Fifth Avenue. The results of our in-laboratory evaluation experiments suggest that we can expect as much as 7 to 8 per cent increase in the store profits by employing targeting rules automatically produced by our methodology, as compared to the current mailing policy used at Saks. These results seem to confirm that our approach is valid in dealing with the cross channel CRM scenarios in the real world.

The rest of the paper is organized as follows. We begin by describing the business problem that we address, in Section 2. Section 3 presents the methodology including detailed descriptions of newly devised methods. Section 4 will describe the experiments we conducted and the results we obtained. Section 5 concludes with a summary.

2. PROBLEM DESCRIPTION

2.1 The business problem

The business problem we address is that of optimizing direct mail catalogue mailings over multiple campaigns, to maximize the effect on the profits/revenue obtained in the store channel. At Saks Fifth Avenue over 60 major direct mail campaigns are conducted each year. These campaigns vary from mailings of general store catalogues to those specific to particular product groups, such as women's apparel and cosmetics. Some are seasonal campaigns, such as Christmas season campaigns. Some may involve store coupons, while others may provide information on upcoming sales and their contents and durations. Many of these campaign features are available for use in data analysis.

Currently, the generation of mailing list for each of these campaigns is based on a number of criteria and constraints, and is not fully automated. There are intricate issues surrounding the process of generating these mailing lists. Our goal is not necessarily a full and immediate automation of this process, but rather in demonstrating the potential use of sophisticated data analytics in assisting and improving it. As a step in this endeavor, our *technical task* is to analyze the past data and automatically generate targeting rules that can be used to construct mailing lists for future campaigns.

2.2 The cross channel challenge

The problem just described is a challenging one, and turns out not to be solvable by straightforward applications of existing modeling techniques. For example, the simplest solution would be to model the short term profits (or revenues) in the store generated by a particular customer, say in a window of one month, as a function of various features of that customer, including the control variable of whether a given catalogue is mailed to that customer. As we elaborate

in the section on experiments, this will result in a model that is independent of the control variable, thus giving no interesting information on the effect of the mailing action.

At the heart of the problem is the *credit assignment* problem. That is, there is no explicit information in the data, linking the actions taken in the direct mail channel to the responses observed in the store channel. This would be possible, if a good part of the transactions were associated with coupons issued in the outbound channel of interest, and this fact were recorded. This is rarely the case in practice, and in particular is not the situation we face in our current problem.

To place further burden on the modeling task, our problem setting involves events with variable length time intervals, i.e. the intervals between the decision points (campaign mailings) are variable in length. This adds considerable amount of noise in the data, and makes the task of modeling the responses of marketing actions even more difficult.

3. METHODOLOGY

3.1 Cross channel life time value maximization and MDP

Common practice in database marketing and CRM today is to organize customer data into a table consisting of fields representing various attributes of customers and response fields, model a response field of interest as a function of those attributes, and then optimize marketing actions against the obtained model. Here we go a step further: We use time stamped sequences of such data to represent time-varying sequences of customer's attributes (state) and marketing actions. We then model the process of customer interactions as a dynamic process, and optimize marketing actions with respect to such a model. The technical framework we employ is the so-called Markov Decision Process (MDP) model, popular in dynamic programming and reinforcement learning. Here we refer the reader to the literature for detailed descriptions of theory and methods in MDP, e.g. [5, 9].

In the current application, we can use the attribute vector corresponding to a customer at a given point in time to represent the state for that customer at that point in time. Cross-channel integration of data allows us to represent the entire history of interactions between a given customer and the enterprise, across all channels, thereby providing a unified view of the customer. The maximization of life time value of a given customer, across all channels, can then be naturally formulated as the maximization of the discounted cumulative rewards in the standard MDP terminology. It follows that life time value maximization in the cross-channel setting is reduced to solving for the optimum policy of the MDP with the cross-channel state representation. Since the obtained policy is a mapping from customer attribute vectors to actions, it can be readily translated to generic if-then style rules for use in any rules engine.

3.2 Q-learning with variable time intervals

As we mentioned in Introduction, for our problem it is necessary to extend the MDP framework to a formulation involving variable time intervals. The *variable time interval MDP* we consider here is identical to the standard discrete time MDP, except that every event is timed. Here we assume that the time at the initial state is 0, and then all subsequent events will have positive time associated with

them. The process starts in some initial state s_1 at $t_1 = 0$, and then the learner repeatedly takes actions, resulting in a sequence of action, state, reward and time quadruples, $\{(s_i, a_i, r_i, t_i)\}_{i=1}^{\infty}$. The goal of the learner is then defined as the maximization of the total discounted rewards, with the discounted factors determined as a function of the time durations. That is, it is to maximize the cumulative reward R ,

$$R = \sum_{i=1}^{\infty} \gamma^{t_i} r_i \quad (1)$$

We note that the model we introduce here is different from and simpler than the extension of MDP to the continuous time setting e.g. SMDP proposed by [4], in that we still assume discrete time steps, though having variable interval length.

The challenge in devising a learning method in the variable time interval MDP is in determining how the rewards in various time intervals should be normalized, in order to lessen the effect of noise introduced by the varying interval length. For estimating the immediate rewards, it is clear that the reward received in a time interval can be normalized by dividing it by the time interval. For estimating the *Q-value function* (the expected cumulative rewards), which is the objective function in an MDP, the situation is significantly more complicated. In particular, we must account for the fact that, at any learning iteration, the interval over which the discounted rewards is summed (for approximating the value function) is incremented. Furthermore, the effective interval is affected by the learning rate, since a larger learning rate corresponds to assigning greater weights to the rewards received far into the future. Here we propose a normalization scheme in which both of these factors are taken into account in determining the effective time interval to be used for normalization. More specifically, we use an analogous update rule for the normalization factor, as that for the Q-value estimate. The resulting variant of the batch Q-learning method for the variable time interval set-up, Vari-RL(Q), is presented in Figure 1. (The update rules for both Q-values, $v_{i,j}^{(k)}$, and normalization factors, $Z_{i,j}^{(k)}$, are in block 4.2 of the pseudo-code.)

3.3 Batch Advantage Updating

A number of past papers have addressed the problem of extending Q-learning and other related learning methods to variable time intervals and more generally to the continuous time setting, e.g. [3, 4]. Of these, the work by Baird [3] is of particular interest to us, since his solution appears to be addressing closely related problems to the one we currently face: function approximation in Q-learning involving variable time intervals is often difficult due to the noise introduced by the varying intervals. Baird proposes a novel method, called *advantage updating*, which tries to learn the relative advantage of competing actions in any given state. This procedure avoids having to explicitly estimate the Q-value function, thereby bypassing the noisy estimation problem. We briefly describe this method and some modifications we made to make it work for our problem.

Advantage updating is based on the notion of *advantage* of an action a relative to the optimal action at a given state s , written $A(s, a)$. The following is one of alternative defi-

Procedure Vari-RL(Q)

Premise:

A base learning module, Base, for regression is given.

Input data: $D = \{e_i | i = 1, \dots, N\}$ where

$$e_i = \{(s_{i,j}, a_{i,j}, r_{i,j}, t_{i,j}) | j = 1, \dots, l_i\}$$

(e_i is the i -th episode, and l_i is the length of e_i .)

1. For all $e_i \in D$

1.1 For $j = 1$ to l_i , $\Delta t_{i,j} = t_{i,j+1} - t_{i,j}$

2. For all $e_i \in D$

2.1 For $j = 1$ to $l_i - 1$

$$Z_{i,j}^{(0)} = \Delta t_{i,j}$$

$$v_{i,j}^{(0)} = r_{i,j}$$

$$D_i^{(0)} = \{(s_{i,j}, a_{i,j}), \frac{v_{i,j}^{(0)}}{Z_{i,j}^{(0)}} | j = 1, \dots, l_i\}$$

3. $\tilde{Q}^{(0)} = \text{Base}(\bigcup_{i=1, \dots, N} D_i^{(0)})$

4. For $k = 1$ to K

4.1 Set α_k , e.g. $\alpha_k = \frac{1}{k}$

4.2 For all $e_i \in D$

For $j = 1$ to $l_i - 1$

$$Z_{i,j}^{(k)} = (1 - \alpha_k) Z_{i,j}^{(k-1)} + \alpha_k (\Delta t_{i,j} + \gamma^{\Delta t_{i,j}} \cdot Z_{i,j+1}^{(k-1)})$$

$$v_{i,j}^{(k)} = (1 - \alpha_k) \tilde{Q}^{(k-1)}(s_{i,j}, a_{i,j}) \cdot Z_{i,j}^{(k-1)}$$

$$+ \alpha_k (r_{i,j} + \gamma^{\Delta t_{i,j}} \max_a \tilde{Q}^{(k-1)}(s_{i,j+1}, a) \cdot Z_{i,j+1}^{(k-1)})$$

$$D_i^{(k)} = \{(s_{i,j}, a_{i,j}), \frac{v_{i,j}^{(k)}}{Z_{i,j}^{(k)}} | j = 1, \dots, l_i - 1\}$$

4.3 $\tilde{Q}^{(k)} = \text{Base}(\bigcup_{i=1, \dots, N} D_i^{(k)})$

5. Output the final unnormalized model, i.e. $\tilde{Q}^{(K)} \cdot Z^{(K)}$.

Figure 1: Variable time interval batch Q-learning.

nitions of this quantity.

$$A^*(s, a) = \frac{1}{\Delta t_s} (r + \gamma^{\Delta t_s} E_{s'} [V^*(s')] - V^*(s)) \quad (2)$$

In the above, note that V^* is defined as $V^*(s) = \max_a Q^*(s, a)$, where $Q^*(s, a)$ is the Q-value of an optimal policy. We used Δt_s to denote the time interval immediately following state s , and s' to denote the state reached after s . Alternatively, the advantage can be written in terms of the Q-value by:

$$A^*(s, a) = \frac{1}{\Delta t_s} (Q^*(s, a) - \max_{a'} Q^*(s, a')) \quad (3)$$

This quantity is extremely interesting for us for two reasons: It factors out the dependence of the value function on the time interval, *and* on the state. Given this notion of advantage, *advantage updating* is an on-line learning method that learns this function iteratively, by a coupled set of update rules for the estimates for A and V , and a normalization step for $A^*(s, a)$ which drives $\max_{a'} A^*(s, a')$ towards zero. We exhibit a batch version of this method in Figure 2.

A couple of other modifications were necessary, before the two methods just presented, Vari-RL(Q) and batch-AU, could be made to work satisfactorily. One modification has to do with the initialization of the quantities being estimated in the two methods, the Q-value and A-value, using the empirical cumulative rewards observed in the data, rather than the immediate rewards as in the original on-line methods. The other modification is allowing optional applications of function approximation, e.g. applying function approximation in every n -th iteration. Due to space limitations, we will omit the details of these modifications, and refer the reader to our technical report [1].

Procedure Batch-AU

Premise:

A base learning module, Base, for regression is given.

Input data: $D = \{e_i | i = 1, \dots, N\}$ where

$$e_i = \{(s_{i,j}, a_{i,j}, r_{i,j}) | j = 1, \dots, l_i\}$$

(e_i is the i -th episode, and l_i is the length of e_i .)

1. For all $e_i \in D$

1.1 For $j = 1$ to l_i , $\Delta t_{i,j} = t_{i,j+1} - t_{i,j}$

2. For all $e_i \in D$

$$D_i^{(0)} = \{(s_{i,j}, a_{i,j}, \frac{r_{i,j}}{\Delta t_{i,j}}) | j = 1, \dots, l_i\}$$

3. $A^{(0)} = \text{Base}(\bigcup_{i=1, \dots, N} D_i^{(0)})$

4. For all $e_i \in D$ and for $j = 1$ to $l_i - 1$, initialize

4.1 $A_{i,j}^{(0)} = A^{(0)}(s_{i,j}, a_{i,j})$

4.2 $Amax_{i,j}^{(0)} = \max_{a'} A^{(0)}(s_{i,j}, a')$

4.3 $V_{i,j}^{(0)} = Amax_{i,j}^{(0)}$

5. For $k = 1$ to K

5.1 Set α_k, β_k and ω_k , e.g. $\alpha_k = \beta_k = \omega_k = \frac{1}{k}$

5.2 For all $e_i \in D$

For $j = 1$ to $l_i - 1$

$$A_{i,j}^{(k)} = (1 - \alpha_k) A_{i,j}^{(k-1)} + \alpha_k (Amax_{i,j}^{(k-1)} + \frac{r_{i,j} + \gamma \Delta t_{i,j} V_{i,j+1}^{(k-1)} - V_{i,j}^{(k-1)}}{\Delta t_{i,j}})$$

$$D_i^{(k)} = \{(s_{i,j}, a_{i,j}, A_{i,j}^{(k)}) | j = 1, \dots, l_i - 1\}$$

5.3 $A^{(k)} = \text{Base}(\bigcup_{i=1, \dots, N} D_i^{(k)})$

5.4 For all $e_i \in D$ and for $j = 1$ to $l_i - 1$, update

$$A_{i,j}^{(k)} = A^{(k)}(s_{i,j}, a_{i,j})$$

$$Amax_{i,j}^{(k)} = \max_{a'} A^{(k)}(s_{i,j}, a')$$

$$V_{i,j}^{(k)} = (1 - \beta_k) V_{i,j}^{(k-1)} + \beta_k (\frac{Amax_{i,j}^{(k)} - Amax_{i,j}^{(k-1)}}{\alpha_k} + V_{i,j}^{(k-1)})$$

5.5 For all $e_i \in D$ and for $j = 1$ to $l_i - 1$, normalize

$$A_{i,j}^{(k)} = (1 - \omega_k) A_{i,j}^{(k)} + \omega_k (A_{i,j}^{(k)} - Amax_{i,j}^{(k)})$$

6. Output the final advantage model, $A^{(K)}$.

Figure 2: Batch reinforcement learning based on advantage updating.

4. EXPERIMENTS

We evaluated the proposed methods using actual customer interaction data from Saks Fifth Avenue from the past years. Below we will describe the data we used for this experimentation in some detail. Then we will discuss the challenge we face in evaluating the methods using past data. We then describe the experimental results.

4.1 Data

As briefly explained in Introduction, the data we used for our data analysis can be categorized into the following four types.

1. *Customer data* for 1.6 million customers, whose recent annual spendings exceeded a certain threshold. These contain demographic and other types of information on the customers. We note that privacy sensitive information was stripped off before they were used for data analytics.
2. *Transaction data* for the said 1.6 million customers for the past three years. The transaction data are time stamped and contain the entire point of sales data,

including the categories of purchased items and sales price, among other things.

3. *Campaign data* for the major campaigns for the year 2002. There were 69 campaigns. These data contain timing of mailing, duration of sales, types of catalogues sent, and product groups (divisions) being targeted by the campaigns.
4. *Product data* for the purchased items in the transactions data. These contain taxonomy information on them, ranging in granularity from product groups down to the SKU level.

These data were then used to generate time stamped sequences of feature vectors containing summarized information on the history of cross-channel interactions. The features we generated and elected to use, representing the state of a customer at any given point in time, are summarized in Table 1. (Note that the features also fall into four types according to the data sources.)

Note that as the third column in Table 1 the correlation coefficients between each of the features and the response variable (reward) are listed. Here, the response variable was calculated simply by summing the observed profits in the data, over a fixed period of time since the time of the mailing in question, and it does not exactly correspond to the *cumulative discounted profits* that we wish to maximize. It is worth noting, nonetheless, that a very low correlation is observed between the control variable (mailing action) and the response variable, as compared to some of the other features. The control variable has the third lowest correlation coefficient (at 0.008), and is magnitudes lower than those of typical (transaction and campaign) features. This explains the nature of the *cross-channel challenge*, and in particular why we tend to get a model that is independent of the control variable, when we run a standard regression engine to model the response variable as a function of these features.

4.2 Evaluation

A common problem in performance evaluation of reinforcement learning methods is that often it is difficult or impossible to conduct real life experiment in which the learning methods have access to on-line interactions with the MDP. Our current application domain of CRM and database marketing is no exception. To conduct such performance evaluation reliably using only static historical data is itself a big challenge. The problem is that we need to evaluate the policy generated by the learning procedure using only past data, which presumably were collected using some policy that is different from it.

Here we propose a solution to this problem, and use it to conduct performance evaluation for our methodology. Our solution is based on a notion recently proposed by Kakade and Langford [6] called *policy advantage*, and a bias correction technique based on importance sampling, c.f. [10]. We elaborate on this below.

First, we define a discrete time version of the notion of *advantage* introduced earlier (in Eq. 3), with respect to any policy π .

$$A_\pi(s, a) = Q_\pi(s, a) - \max_{a'} Q_\pi(s, a') \quad (4)$$

Then the policy advantage of a new policy π' with respect to an old (or sampling) policy π and initial state distribution

| Features | Descriptions | Cor. |
|------------------------------|--|--------|
| full_line_store_of_residence | if a full-line store exists in the area | 0.004 |
| off_fifth_store_of_residence | if an off-fifth store exists in the area | -0.004 |
| loyalty_program_level | loyalty program level | 0.074 |
| fvrt_store_channel | favorite store channel (web/store) | 0.046 |
| purchase_amt_lm | amount of purchase in last month | 0.085 |
| purchase_amt_2_3m | amount of purchase in last 2-3 month period | 0.098 |
| purchase_amt_6m | amount of purchase in last 4-6 month period | 0.102 |
| purchase_amt_ly | amount of purchase in last year | 0.139 |
| purchase_amt_tot | total amount of purchase (in 3 years) | 0.155 |
| promo_purchase_ratio | ratio of purchases in promotion periods | 0.047 |
| cur_div_purchase_amt_lm | purchase amount last month in current division | 0.090 |
| cur_div_purchase_amt_2_3m | purchase amount last 2-3 month in current division | 0.080 |
| cur_div_purchase_amt_6m | purchase amount last 4-6 month in current division | 0.091 |
| cur_div_purchase_amt_ly | purchase amount last 1 year in current division | 0.128 |
| cur_div_purchase_amt_tot | total purchase amount in current division | 0.147 |
| div_purchase_amt_tot_j | total purchase amount in division j | 0.093* |
| n_cat_lm | number of catalogues sent in last month | 0.021 |
| n_cat_2_3m | number of catalogues sent in last 2-3 month | 0.028 |
| n_cat_4_6m | number of catalogues sent in last 4-6 month | 0.063 |
| n_cat_tot | total number of catalogues sent | 0.086 |
| cur_div_n_cat_lm | number of catalogues sent in last month targeting current division | 0.028 |
| cur_div_n_cat_2_3m | number of catalogues sent in last 2-3 month targeting current division | 0.025 |
| cur_div_n_cat_4_6m | number of catalogues sent in last 4-6 month targeting current division | 0.062 |
| cur_div_n_cat_tot | total number of catalogues sent targeting current division | 0.062 |
| action | to mail or not to mail | 0.008 |

Table 1: Features used in our experiments: Features, their descriptions, and their correlations with the reward field. The features are listed in 4 groups: (1) customer features; (2) transaction features; (3) campaign features; and (4) product group specific campaign features. (* An example correlation value is exhibited for several features of this type.)

μ , written $A_{\pi, \mu}(\pi')$, is defined as follows.

$$A_{\pi, \mu}(\pi') = E_{s \sim \pi, \mu} [E_{a \sim \pi'(a|s)} [A_{\pi}(s, a)]] \quad (5)$$

Intuitively, the policy advantage measures how much advantage can result by replacing the action of the old policy by that of the new policy at a random state selected by the sampling policy, while all other actions remain unchanged (specified by the sampling policy). In some sense, this measure quantifies how much local improvement is attained by changing the old policy by the new policy, assuming that the overall state distribution is not significantly affected by that change.¹

The policy advantage can be estimated using only data collected by the sampling policy π , using a bias correction technique based on importance sampling as follows.

$$A_{\pi, \mu}(\pi') = E_{s, a \sim \pi, \mu} \left[\frac{\pi'(a|s)}{\pi(a|s)} [A_{\pi}(s, a)] \right] \quad (6)$$

Note that $\pi'(a|s)$ is known since it is the (possibly stochastic) policy generated by reinforcement learning, but $\pi(a|s)$ need be estimated from the data, since we do not know the sampling policy explicitly. It should be pointed out that this quantity becomes impossible to estimate for a deterministic sampling policy, since the bias correction factor $\frac{\pi'(a|s)}{\pi(a|s)}$ diverges for actions a never chosen by the sampling policy π . In real world settings, however, the state information is often not sufficient to determine the chosen action deterministically, as is the case in our setting.

¹Kakade and Langford [6] have established a theoretical result which implies that a new policy with a positive policy advantage can be used to define a new policy, which provably has better performance than the old policy.

4.3 Experimental Results

We used the evaluation method just introduced, namely that of bias corrected estimation of policy advantage in the data, to validate the performance of the proposed methods. We used both of the proposed methods, Vari-RL(Q) and Batch-AU. In both cases, we used IBM's scalable regression engine, ProbETM [7, 2], as the basic regression module. For both methods, we used the feature to initialize the Q-value and A-value estimates using the empirical life time values. Also, for both methods, we used the option of applying function approximation at every fourth learning iteration. In our evaluation, we randomly sampled approximately 1.0 percent of the individual customers from the entire data (approximately 16 thousand customers from 1.6 million). The episodic data were then generated, for use in training, by randomly selecting a sub-episode of length 10 (consisting of 10 events) corresponding to each of the sampled individuals. A separate test data set consisting of 5,000 randomly selected individuals was also sampled for calculating the policy advantage. The policy advantage was calculated using these individual data over the entire 68 campaigns, and for each learning iteration, the results were averaged over 10 random runs.

The results of this evaluation are exhibited in Figures 3 and 4. Figure 3 plots how the policy advantage of the policy output by the Vari-RL(Q) method changes as the learning iteration progresses in a typical run. Figure 4 shows the analogous graph for the batch-AU method. In each of the graphs, the y-axis is the policy advantage shown as the percentage over the value of the old policy. Strictly speaking, what is shown on the x-axis is not the number of iterations, but rather is the number of times function approximation is performed. In both cases, we chose to run function approximation at every fourth learning iteration.

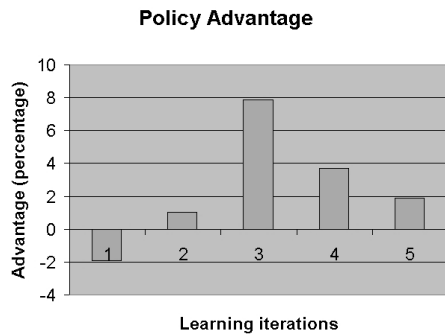


Figure 3: The policy advantage for the variable time interval Q-learning method (in a typical run.)

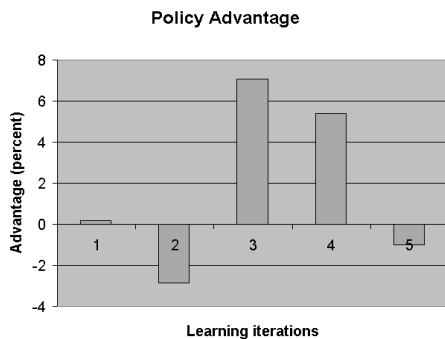


Figure 4: The policy advantage for the Batch Advantage Updating method (in a typical run.)

A definite trend can be read off from these graphs. For both of the methods, a typical run starts with a policy that is relatively uninformative which does not show any advantage over the sampling policy. It is worth noting that, since both methods were initialized with the empirical life time value observed in the data, this shows that direct modeling with empirical LTV does not lead to any advantage over the existing mailing policy. At the third function approximation, or after 9 learning iterations, the policy advantage peaks and then it starts declining again. This behavior is thought to be attributable, in part, to the nature and limitation of the evaluation method. Policy advantage measures the advantage of a new policy with respect to an old policy, using the old policy as the sampling policy. It is therefore more effective when the two policies are relatively similar. As the learning progresses and the two policies start diverging, the measure becomes less and less reliable.

Even with the limitation mentioned above, the obtained results are quite encouraging. With the assumption that the new policy does not significantly change the state distribution, the results imply that as much as 7 to 8 percent increase in the store profits can be expected, by employing the policy output by our methodology, over the current mailing policy used at Saks.

5. CONCLUSIONS

We validated our reinforcement learning based approach to life time value modeling and cross-channel optimized marketing on a real world problem. In the course of our investi-

gation, we identified a general problem common to modeling cross-channel interactions, and proposed a solution based on old and new techniques of reinforcement learning. We also provided a solution to the sampling bias problem in evaluation of learned policies, and used it to evaluate the proposed approach. Some issues for future investigation include the following: (1) Easing deployment by reducing the need to customize; (2) Handling various channel constraints including budget constraints;

6. ACKNOWLEDGMENTS

We wish to thank Bill Franks and Sheri Wilson-Gray of Saks Fifth Avenue for their executive leadership in making the joint project possible. We also wish to thank Edwin Pednault and Bianca Zadrozny of IBM Research and John Langford of TTI at Chicago for helpful discussions and assistance.

7. REFERENCES

- [1] N. Abe, N. Verma, C. Apte, and R. Schroko. Cross channel optimized marketing by reinforcement learning. Technical Report RC23132(W0403-021), IBM Research, March 2004.
- [2] C. Apte, E. Bibelnicks, R. Natarajan, E. Pednault, F. Tipu, D. Campbell, and B. Nelson. Segmentation-based modeling for advanced targeted marketing. In *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, pages 408–413. ACM, 2001.
- [3] L. C. Baird. Reinforcement learning in continuous time: Advantage updating. In *Proceedings of the International Conference on Neural Networks*, June 1994.
- [4] S. Bradtke and M. Duff. Reinforcement learning methods for continuous-time Markov decision problems. In *Advances in Neural Information Processing Systems*, volume 7, pages 393–400. The MIT Press, nov 1995.
- [5] L. P. Kaelbling, M. L. Littman, and A. W. Moore. Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, 4, 1996.
- [6] S. Kakade and J. Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the 19th International Conference on Machine Learning*, July 2002.
- [7] R. Natarajan and E. Pednault. Segmented regression estimators for massive data sets. In *Second SIAM International Conference on Data Mining*, Arlington, Virginia, 2002. to appear.
- [8] E. Pednault, N. Abe, B. Zadrozny, H. Wang, W. Fan, and C. Apte. Sequential cost-sensitive decision making with reinforcement learning. In *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM Press, 2002. To appear.
- [9] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- [10] B. Zadrozny. *Policy mining: Learning decision policies from fixed sets of data*. PhD thesis, University of California, San Diego, 2003.