# Estimating Bidders' Valuation Distributions in Online Auctions

Albert Xin Jiang, Kevin Leyton-Brown
Department of Computer Science,
University of British Columbia
{jiang, kevinlb}@cs.ubc.ca

## Abstract

There is much active research into the design of automated bidding agents, particularly for environments that involve multiple decoupled auctions. These settings are complex partly because an agent's strategy depends on information about other bidders' interests. When bidders' valuation distributions are not known *ex ante*, machine learning techniques can be used to approximate them from historical data. It is a characteristic feature of auctions, however, that information about some bidders' valuations is systematically concealed. This occurs in the sense that some bidders may fail to bid at all because the asking price exceeds their valuations, and also in the sense that a high bidder may not be compelled to reveal his valuation. Ignoring these "hidden bids" can introduce bias into the estimation of valuation distributions. To overcome this problem, we propose an EM-based algorithm. We validate the algorithm experimentally using both synthetic and real-world (eBay) datasets, and show that our approach estimates bidders' valuation distributions and the distribution over the true number of bidders significantly more accurately than more straightforward density estimation techniques.

## 1 Introduction

There has been much research on the study of automated bidding agents for auctions and other market-based environments. The Trading Agent Competitions (TAC) and the TAC Supply Chain Management competitions (TAC-SCM) have attracted much interest [14]. There have also been research efforts on bidding agents and bidding strategies in other auction environments [5, 4, 7, 3, 6, 2]. Although this body of work considers many different auction environments, bidding agents always face a similar task: given a valuation function, the bidding agent needs to compute an optimal bidding strategy that maximizes expected surplus. (Some environments such as TAC-SCM also require agents to solve additional, e.g., scheduling tasks.)

The "Wilson Doctrine" in mechanism design argues that mechanisms should be constructed so that they are "detail-free"—that is, so that agents can behave rationally in these mechanisms even without information about the distribution

of other agents' valuations. For example, under the VCG mechanism it is a weakly dominant strategy to bid exactly one's valuation, regardless of other agents' beliefs, valuations or actions. Under common assumptions (in particular, independent private values) single-item English auctions are similar: an agent should remain in the auction until the bidding reaches the amount of his valuation.

While detail-free mechanisms are desirable, however, they are not ubiquitous. Very often, agents are faced with the problem of deciding how to behave in games that do not have dominant strategies and where other agents' payoffs are strategically relevant. For example, we may want to participate in a series of auctions run by different sellers at different times.

## 1.1 Game-Theoretic and Decision-Theoretic Approaches

Depending on the assumptions we choose to make about other bidders, two approaches to computing bidding strategies suggest themselves: a game theoretic approach and a decision theoretic approach. The game theoretic approach assumes that all agents are perfectly rational and that this rationality is common knowledge; the auction is modeled as a Bayesian game (see, e.g., the survey in [8]). Under this approach, a bidding agent would compute a Bayes-Nash equilibrium of the auction game, and play the equilibrium bidding strategy. For example, for environments with multiple, sequential auctions for identical items and in which each bidder wants only a single item, Milgrom and Weber [9, 13] identified Bayes-Nash equilibria. Such equilibria very often depend on the distribution of agents' valuation functions and the number of bidders. Although this information is rarely available in practice, it *is* usually possible to estimate these distributions from the bidding history of previous auctions of similar items. Note that this involves making the assumption that past and future bidders will share the same valuation distribution.

The game-theoretic approach has received a great deal of study, an is perhaps the dominant paradigm in microeconomics. In particular, there are very good reasons for seeking strategy profiles that are resistant to unilateral deviation. However, this approach is not always useful to agents needing to decide what to do in a particular setting, especially when the rationality of other bidders is in doubt, when the computation of equilibria is intractable, or when the game has multiple equilibria. In such settings, it may be more appropriate to rely on decision theory. A decision theoretic approach treats other bidders as part of the environment, and ignores the possibility that they may change their behavior in response to the agent's actions. As above, we again make the assumption that the other bidders come from a population that exhibit stationary bidding behavior; however, this time we model agents' bid amounts directly, rather than modeling their valuations and then applying an equilibrium strategy. We then solve the resulting single-agent decision problem to find a bidding strategy that maximizes expected payoff. We could also use a reinforcement-learning approach, where we continue to learn the bidding behavior of other bidders while participating in the auctions.

This paper does not attempt to choose between these two approaches; it is

our opinion that each has domains for which it is the most appropriate. The important point is that regardless of which approach we elect to take, we are faced with the subproblem of estimating two distributions from the bidding history of past auctions: the distribution on the number of bidders, and the distribution of bid amounts (for decision theoretic approaches) or of valuations (for game theoretic approaches).

## 1.2 Hidden Bids

It might seem that there is very little left to say on this topic: we learn the distributions of interest from historical data and then compute a bidding strategy based on that information for the current auction. However, bidding histories often systematically omit relevant information. For example, in sealed bid auctions, the auctioneer may choose not to reveal the bid amounts except the price the winner pays. An English auction is stopped when there is only one active bidder left (i.e., when the second-highest bidder drops out), meaning that the valuation of the highest bidder is not revealed. How can we learn valuation distributions when the data available to us is biased in this way?

For concreteness, in this paper we focus on a single domain; however, our techniques are broadly applicable. Here we consider sequential English auctions in which a full bidding history is revealed, such as the online auctions run by eBay. We are thus concerned with two kinds of missing information. First, some bidders may come to the auction when it is already in progress, find that the current price already exceeds their valuation, and leave without placing a bid. Second, the amount the winner was willing to pay is never revealed.

Ignoring these sources of bias would lead to poor estimates of the underlying valuation distribution. We propose a novel learning approach based on the Expectation Maximization (EM) algorithm, which iteratively generates hidden bids consistent with the observed bids, and then computes maximum-likelihood estimations of the valuation distribution based on the completed set of bids. Considering both synthetic data (in which true valuation distributions are known) and real-world data from eBay, we show that our approach outperforms more straightforward distribution estimation techniques which do not attempt to account for this missing data.

The rest of the paper is organized as follows. Section 2 introduces our auction setting and describes our generative probabilistic model for the bidding process. Section 3 focuses on the estimation problem, and describes our EM learning approach. Section 4 discusses the computation of the optimal strategy under the decision theoretic approach. In Section 5 we present experimental results on synthetic data sets as well as on data collected from eBay, which show that our EM learning approach makes better estimates of the distributions, and gets more payoff under the decision theoretic model, as compared to the straightforward approach which ignores hidden bids.

## 2 A Model of the Bidding Process of Online Auctions

Online auctions present unique challenges to agents trying to estimate the underlying valuation distributions, because we do not get to observe the true

number of potential bidders, instead we only see the bidders that has decided to make a bid; whereas in an auction house environment, we could estimate the number of potential bidders by the number of people present at the auction. Most online auctions nowadays are English auctions, which are ascending-price auctions during which bidders raise the current price level until only one bidder is willing to pay the current price for the item.

In this paper we analyze online English auctions as implemented by eBay. Other auction sites' rules are similar to eBay's. One common feature of online auction sites including eBay is the proxy bidding system, which allows bidders to enter their maximum willingness to pay as proxy bids. The proxy bidding system will then make the bids automatically for the bidder, up to the specified amount of the proxy bid[1]. Each eBay auction has a fixed closing time; when the auction closes, the highest bidder is declared the winner of the auction.

We now present a generative model for the bidding process of an eBay auction, which describes how the potential bids are generated, which of them become visible and which of them are hidden. There are $m$ potential bidders interested in a certain eBay auction of a single item. We assume that bidders have independent private values(IPV). We assume that $m$ is drawn from a discrete distribution $g(m)$ with support $[2, \infty)$. Bidders' potential bids are independently drawn from a continuous distribution $f(x)$.

A decision theoretic agent's task is to estimate $f(x)$ and $g(m)$ and use the estimated distributions to compute an optimal bidding strategy. If we instead are using the game theoretic approach, we are interested in the bidders' valuations. In that case we would use a slightly different model, where bidders' valuations are independently drawn from a distribution $f(v)$, and each bidder bids according to a known Bayes-Nash equilibrium. Our game theoretic agent's task is then to estimate $f(v)$ and $g(m)$.

The $m$ potential bidders submit their potential bids in a sequential order. When a proxy bid is submitted, eBay compares it to the current price level, which is the second-highest proxy bid so far plus a small bid increment. For simplicity, in this paper we ignore this small increment and assume that the current price level is the second-highest proxy bid so far. If the submitted bid is no greater than the current price level, the bid is dropped and nothing is observed. If the submitted bid is higher than the current price level but lower than the highest proxy bid so far, then the price level is increased to equal the submitted bid. If the submitted bid is higher than the previous highest bid, then the price level is increased to equal the previous highest bid. At the end of the auction, the item is awarded to the bidder who placed the highest bid, and the final price level will be equal to the second highest bid.

Our model of the bidding process is quite general. Notice that when a bidder observes that the price level is higher than her potential bid, she may decide not to bid in this auction. This is equivalent to our model in which she always submits the bid, because dropped bids do not appear in the bidding history. Also our model covers the case of last-minute bidding, which happens quite

---

[1] http://pages.ebay.com/help/buy/proxy-bidding.html

often in eBay auctions [10]: even though last-minute bids may be submitted almost simultaneously, eBay processes the bids in sequence.

With the proxy bidding system, and when agents have IPV, there is no strong motivation to bid more than once in an auction. However, in practice eBay bidders quite often make multiple bids in one auction. One possible motivation of these bids is to reveal more information about the proxy bid of the current high bidder [11]. However, only the last bid of the bidder represents her willingness to pay. Given a bidder's last bid, her earlier bids carry no extra information. Therefore, we will be interested in only the last bid from each bidder[2]. We can preprocess the bidding histories by removing all bids except the last bids from each bidder, without losing much information.

## 3  Estimating the Distributions

Given the model of the bidding process, the first task of our bidding agent is to estimate the distributions $f(x)$ and $g(m)$ from the bidding history. Suppose we have access to the bidding history of $K$ auctions of the same item.

### 3.1  The Simple Approach

One simple approach is to ignore the missing bids, and try to directly estimate $f(x)$ and $g(m)$ from observed data. The observed number of bidders, $n$, is used to estimate $g(m)$. To estimate $f(x)$ we use the observed bids $x_v$, which consists of $(n-1)$ bids for each auction, since the bids of the highest bidders are missing. This approach and its variations have been used in e.g. [3, 5, 7].

### 3.2  EM Learning Approach

We would like to have an estimation strategy that accounts for the missing data and any bias introduced by its absence. Let us denote the hidden bids by $x_h$. According to our model, there are two types of missing bids:

1. The highest bid of each auction $x_{hi}$.

2. The dropped bids $x_d$ that are not observed due to being lower than the current price. Since there are $m$ potential bidders in total, $(n-1)$ bids are visible, and one bid is the highest bid $x_{hi}$, there are $(m-n)$ dropped bids in $x_d$.

Suppose $f(x)$ belongs to a class of distributions parameterized by $\theta$: $f(x|\theta)$, and $g(m)$ belongs to a class of distributions parameterized by $\lambda$: $g(m|\lambda)$. We want to find the maximum likelihood estimates of $\theta$ and $\lambda$, given the observed data $x_v$.

Suppose that we could actually observe the hidden bids $x_h$ as well as $x_v$. Then estimating $\theta$ and $\lambda$ from the completed data set $(x_v, x_h)$ would be a much easier task. Unfortunately we do not have $x_h$. Given $x_v$, and with the knowledge of the bidding process, we can generate $x_h$ if we know $\theta$ and $\lambda$. Unfortunately we do not know $\theta$ and $\lambda$.

---

[2]In a common value model, the earlier bids does carry some information, and we would not be able to simply ignore those bids.

A popular strategy for learning this kind of model with missing data is the Expectation Maximization (EM) algorithm. EM is an iterative procedure that alternates between E steps which generate the missing data given current estimates for the parameters and M steps which compute the maximum likelihood (or maximum a posteriori) estimates for the parameters based on the completed data, which consists of the observed data and current estimates for the missing data.

Formally, the E step computes

$$Q(\theta) = \int \log(p(x_h, x_v | \theta)) p(x_h | x_v, \theta^{(old)}, \lambda^{(old)}) dx_h \tag{1}$$

The M step does the following optimization:

$$\theta^{(new)} = \arg\max_\theta (Q(\theta)) \tag{2}$$

Similar computations are done to estimate $\lambda$, the parameter for $g(m|\lambda)$. The integral in Equation (1) is generally intractable for this complex bidding process. Instead, we can compute a Monte Carlo approximation of the integral: we draw $N$ samples from the distribution $p(x_h | x_v, \theta^{(old)}, \lambda^{(old)})$, and approximate the integral by a small sum over the samples (see e.g. [1]).

Applied to our model, in each E step our task is to generate samples from the distribution $p(x_h | x_v, \theta^{(old)}, \lambda^{(old)})$. Recall that $x_h$ consists of the highest bid $x_{hi}$ and the dropped bids $x_d$.

Given $\theta^{(old)}$ and the second highest bid (which is observed), the highest bid $x_{hi}$ can easily be sampled from the distribution $f(x|\theta^{(old)})$ truncated at the second highest bid. Sampling the dropped bids $x_d$ is a more difficult task. We use the following procedure, which is based on simulating the bidding process:

1. Sample $m$ from $g(m|\lambda^{(old)})$.

2. If $m < n$, reject the sample and go back to step 1.

3. Simulate the bidding process using $x_v$ and $m - n$ dropped bids:

   (a) Repeatedly draw a sample bid from $f(x|\theta^{(old)})$, and compare it to the current price level. If it is lower than the price level, add the bid to the set of dropped bids $x_d$. Otherwise, the current price level is increased to the next visible bid from $x_v$.

   (b) If the number of bids in $x_d$ exceeds $m - n$, or if we used up all the bids in $x_v$ before we have $m - n$ dropped bids in $x_d$, we reject this sample and go back to step 1. Only when we used up all bids in $x_v$ and we have $m - n$ bids in $x_d$, do we accept the sample of $x_d$.

4. Repeat until we have generated N samples of $x_d$.

Our task at each M step is to compute the maximum likelihood (ML) estimates of $\lambda$ and $\theta$ from $x_v$ and the generated samples of $x_h$. For many standard parametric family of distributions, There are analytical solutions for the ML

estimates. If analytical solutions does not exist we can use numerical optimization methods such as simulated annealing. The EM algorithm terminates when $\lambda$ and $\theta$ converges.

## 4 Computing an Optimal Bidding Strategy

In Section 2 and 3 we presented a model of the bidding process for a single auction, and proposed methods to estimate the distributions of bids and number of bidders in an auction. But our work is not done yet: how do we make use of the estimated distributions to compute a bidding strategy? If we only participate in one English auction, under the IPV model it is a dominant strategy to bid up to our valuation of the item, and we do not even need to estimate the distributions. However if we are interested in buying from multiple auctions, good estimates of the distributions $f(x)$ and $g(m)$ is essential in computing a good bidding strategy.

In this section we develop a decision theoretic bidding agent for finitely repeated auctions. We choose this setting because it is a fairly accurate model of the decision theoretic problem we would face when we want to buy one item from an online auction site. Our estimation algorithm can easily be applied to more complex decision theoretic models such as infinite horizon models with discount factors, and combinatorial valuation models, as well as game-theoretic bidding models.

### 4.1 Repeated Auctions

Suppose we only want to buy one item, say a Playstation from eBay, where multiple auctions of similar Playstation systems are held regularly. If we successfully win one item, our utility will be equal to our valuation for the item minus the price we pay. So our bidding agent's task is to compute a bidding strategy that will maximize this utility. We are only interested in the next $k$ auctions after we arrived at the auction site. One motivation for such a restriction is that usually we prefer to have the item soon, i.e. we do not want the bidding agent to spend too much time in order to get the best deal. If we fail to win an item from the $k$ auctions, we lose interest of the item and leave the auction site, and our utility is 0. An alternative model would be that we could go and buy an item from a store after we leave the auction site, in which case we would get some other constant utility. Some of the $k$ auctions may overlap in time, but since eBay auctions have strict closing times, this can be modeled as a sequential decision problem, where our agent makes bidding decisions right before each auction closes.

Number the auctions $1 \ldots k$ according to their closing times. Let $v_j$ denote our valuation for the item from auction $j$. Note that this allows the items in the auctions to be non-identical. Let $b_j$ denote our agent's bid for auction $j$. Let $U_j$ denote our agent's expected payoff from participating in auctions $j \ldots k$, assuming we did not win before auction $j$. Let $U_{k+1}$ be our payoff if we fail to win any of the auctions. For simplicity, in this paper we define $U_{k+1} = 0$. Suppose for each auction $j$, the number of other bidders is drawn from $g_j(m)$ and each bidder's bid is drawn from $f_j(x)$. Since each auction $j$ is an English

auction, only the highest bid from other bidders affects our payoff. Let $f_j^1(x)$ and $F_j^1(x)$ respectively denote the probability density function and cumulative density function (CDF) of the highest bid from other bidders in the $j$-th auction. Then $F_j^1(x) = \sum_{m=2}^{\infty} g_j(m)(F_j(x))^m$, where $F_j(x)$ is the CDF of $f_j(x)$. Now $U_j$ can be expressed as the following function of the future bids $b_{j:k} = (b_j, \ldots, b_k)$ and valuations $v_{j:k} = (v_j, \ldots, v_k)$:

$$U_j(b_{j:k}, v_{j:k}) = \int_{-\infty}^{b_j} (v_j - x) f_j^1(x) dx + (1 - F_j^1(b_j)) U_{j+1}(b_{j+1:k}, v_{j+1:k}) \quad (3)$$

The first term is the expected payoff from the $j$-th auction; the second term is the expected payoff from the later auctions. Let $U_j^*(v_{j:k})$ denote the expected payoff under optimal strategy $b_{j:k}^*$. We can optimize $U_j$ from the $k$-th auction to the first one, in a manner similar to backward induction. By solving the first-order conditions of $U_j$, we obtain the optimal bidding strategy:

$$b_j^* = v_j - U_{j+1}^*(v_{j+1:k}) \quad (4)$$

In other words, our agent should shade her bids by the "option value", i.e. the expected payoff of participating in future auctions, except for the $k$-th auction. In the latter case there are no future auctions and the optimal bid is $b_k^* = v_k$.

The computation of the optimal bidding strategies requires the computation of the expected payoffs $U_j^*$, which involves an integral over the distribution $f_j^1(x)$. In general this cannot be done analytically, but we can compute its Monte Carlo approximation if we can sample from $f_j^1(x)$. If we can sample from $f_j(x)$ and $g_j(m)$, we can use the following straightforward procedure to generate a sample from $f_j^1(x)$: first draw $m$ from $g_j(m)$, then draw $m$ samples from $f_j(x)$ and take the maximum.

The bidding strategy $b_{1:k}^*$ computed using Equations 4 and 3 is optimal, provided that the distributions $f_j(x)$ and $g_j(m)$ are the correct distributions of bids and number of bidders for all $j \in 1 \ldots k$. Of course in general we do not know the true $f_j(x)$ and $g_j(m)$ and the focus of this paper is to estimate the distributions from the bidding history and use the estimated distributions to compute the bidding strategy. As a result, the computed bidding strategy should be expected to achieve less than the optimal expected payoff. However, it is reasonable to think that better estimates of $f(x)$ and $g(m)$ should give bidding strategies with higher expected payoffs. This is confirmed in our experiments across a wide range of data sets, which we discuss in Section 5.

### 4.2 Auctions with Partial Bidding Activity

So far we have been estimating the distribution of the highest bid $f^1(x)$ using $f(x)$ and $g(m)$. In practice, we often observe some early bidding activity by other bidders in auctions $j + 1, \ldots, k$ before we have to make a bid on auction $j$. This allows us to make even more informed estimates on $f^1(x)$, based on $f(x)$, $g(m)$ and the observed bids. Suppose we have observed $n - 1$ early bids, denoted by $x_v$; the current highest bid $x_{hi}$ is not revealed (but can be sampled from $f(x)$ truncated at the current price). Since the auction is not over, there

could be some set of future bids $x_{future}$. When the auction closes, the highest bid from the other bidders will be $\max\{x_{hi}, x_{future}\}$. We can generate $x_{future}$ if we know the number of future bids. We know the total number of bids $m$ is drawn from $g(m)$, and the number of bids made so far is $n + |x_d|$, where $x_d$ are the dropped bids so far, so the number of future bids is $m - n - |x_d|$. Now we have a procedure that samples from $f^1(x)$:

1. Simulate the auction using our model in Section 2 to generate $x_d$, the dropped bids so far.

2. Sample the total number of bids $m$ from $g(m)$.

3. Compute the number of future bids, $m - n - |x_d|$. If this quantity is negative, reject the sample. Otherwise generate $x_{future}$.

4. Generate $x_{hi}$ and take the maximum of $x_{future}$ and $x_{hi}$.

## 5  Experiments

We evaluated both our EM learning approach and the simple approach on several synthetic data sets and on real world data collected from eBay. For each data set, both approaches were used to estimate the distributions of number of bidders and bid amounts, and these estimates were then used to compute bidding strategies and expected payoffs under the decision theoretic model of Section 4. We compared the approaches in two ways:

1. Which approach gives better estimates of the distributions $f(x)$, $g(m)$ and $f^1(x)$? This is important because better estimation of these distributions should give better results, regardless of whether agents take a decision theoretic approach or a game theoretic approach to bidding. We measure the closeness of an estimated distribution to the true distribution using the Kullback-Leibler (KL) Divergence from the true distribution to the estimated distribution. The smaller the KL Divergence, the closer the estimated distribution to the true one.

2. Which approach gives better expected payoff under the decision theoretic bidding model[3] as described in Section 4?

Our experiments show that the EM learning approach outperforms the simple approach on both aspects, across a wide range of data sets. In this section we present results on four representative data sets:

- Data Set 1 has auctions of identical items, and we know the family of distributions that f(x) and g(m) belong to.

- Data Set 2 has auctions of non-identical items, but we know the bid distribution $f(x)$ is influenced linearly by an attribute $a$.

---

[3]We recognize that game theoretic models are also important, and game theoretic bidding agents would benefit from better estimation of the valuation distributions. However in a game theoretic model we cannot simply use expected payoff as a measure of performance, since our payoff depends very much on the other agents' strategies.

- Data Set 3 has auctions of identical items, but we do not know what kind of distributions $f(x)$ and $g(m)$ are. We use nonparametric estimation techniques to estimate the distributions.

- Data Set 4 is real-world auction data on identical items, collected from eBay.

## 5.1 Synthetic Data Set 1: Identical Items

In this data set, the items on sale in all auctions are identical, so the number of bidders and bid amounts come from stationary distributions $g(m)$ and $f(x)$. $f(x)$ is a Normal distribution $N(4, 3.5)$. $g(m)$ is a Poisson distribution shifted to the right: $g(m - 2) = P(40)$, i.e. the number of bidders is always at least 2. The bidding history is generated using our model of the bidding process as described in Section 2. Each instance of the data set consists of bidding history from 40 auctions. We generated 15 instances of the data set.

Both estimation approaches are informed of the parametric families from which $f(x)$ and $g(m)$ are drawn; their task is to estimate the parameters of the distributions, $(\mu, \sigma)$ for $f(x)$ and $\lambda$ for $g(m)$. At the M step of the EM algorithm, standard ML estimates for $\mu$, $\sigma$, and $\lambda$ are computed, i.e. sample mean of the bid amounts for $\mu$, standard deviation of the bid amounts for $\sigma$, and the mean of the number of bidders minus 2 (due to the shifting) for the Poisson parameter $\lambda$.

Our results show that the EM approach outperforms the simple approach in the quality of its estimates for the distributions $f(x)$, $g(m)$ and $f^1(x)$. Figure 1 shows typical estimated distributions[4] and the true distributions. We observe that the plot of the estimated $f(x)$ by the simple approach is significantly shifted to the right of the true distribution, i.e. the simple approach overestimated $f(x)$. We have also calculated KL Divergences from the true distributions to the estimated distributions, and the EM estimations have consistently lower KL Divergences. This difference was verified to be significant, using the non-parametric Wilcoxon sign-rank test.

Then, estimates from both approaches are used to compute bidding strategies for an auction environment with 8 sequentially held auctions of the same kind of items, using the decision theoretic model. The agent's "actual" expected payoffs $U_1(b, v)$ under these bidding strategies are then computed, using the true distributions. The optimal bidding strategy and its expected payoff are also computed.

Our results show that the EM approach gives bidding strategies closer to the optimal strategy, and achieves higher expected payoffs, compared to the simple approach. Figure 1 has a plot of the bidding strategies in the first auction, and a box plot of the regrets, which is the differences between optimal expected payoffs and actual expected payoffs. From the box plot we observe that the mean regret of the EM approach is much smaller than that of the simple approach.

---

[4]The distributions shown are randomly chosen from the 15 instances of the data set. We have verified that the plots of the other distributions are qualitatively similar.
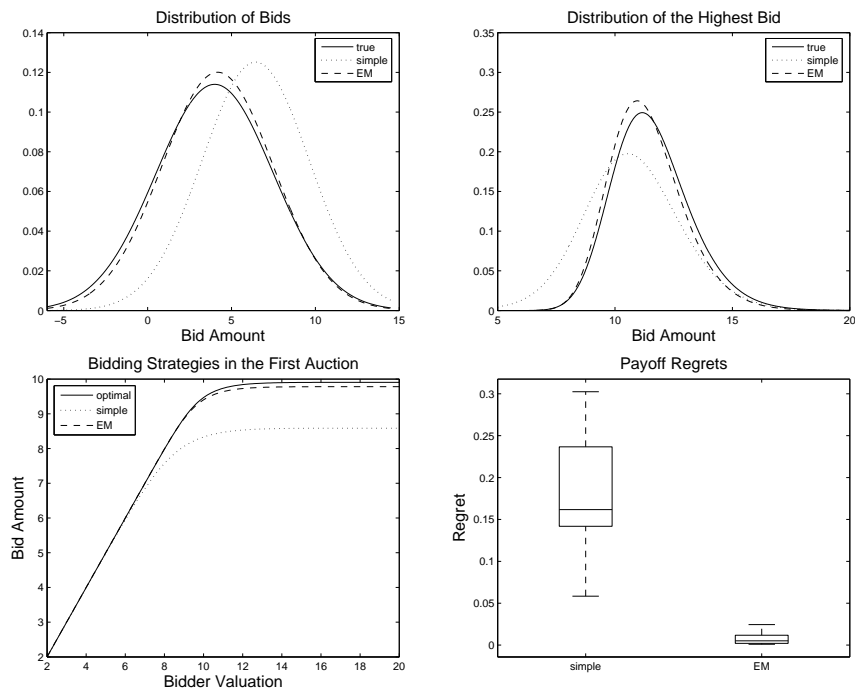
Figure 1: Results for Data Set 1: Distribution of bids $f(x)$ (top-left). Distribution of highest bids $f^1(x)$ (top-right). Bidding strategies in the first auction (bottom-left). Box plot of payoff regrets of the two approaches (bottom-right).

We also used the estimated distributions on the decision-theoretic model with partial bidding activity, as described in Section 4.2. Again our results show that the EM approach achieves higher expected payoffs compared to the simple approach.

## 5.2 Synthetic Data Set 2: Non-identical Items

In our second data set, the items on sale are not identical; instead the distribution of valuations are influenced by an observable attribute $a$. In this data set the dependence is linear: $f(x|a) = N(1.1a + 1.0, 3.5)$. $g(m)$ is a Poisson distribution as before: $g(m - 2) = P(35)$. For each auction, $a$ is sampled uniformly from the interval $[3, 9]$. In other words, this data set is similar to data set 1, except that the bid distribution $f(x)$ is drawn from a different parametric family. Both approaches now use linear regression to estimate the linear coefficients.

Again, our results show that the EM approach outperforms the simple approach for this data set, in terms of its estimates for $f(x)$ and $g(m)$. Figure 2 shows the estimated linear relation between the mean of $f(x|a)$ and $a$. From the figure we can see that the EM approach gives a much better estimate to the linear function. The simple approach again significantly overestimates the bid amounts. In fact the simple approach has consistently overestimated $f(x)$ for
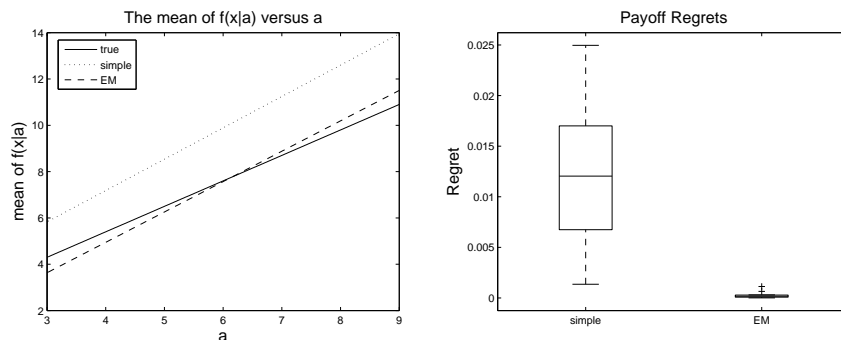
11

Figure 2: Results for Data Set 2: Linear relationship between the mean of $f(x|a)$ and $a$ (left). Box plot of payoff regrets (right).

all the synthetic data sets we tested. This shows that given our model of the bidding process, the estimated $f(x)$ of the simple approach is biased.

We then used the estimated distributions to compute a decision-theoretic agent's bidding strategies and expected payoffs of an auction environment with 8 sequential auctions, where the attribute $a$ of each item is observed. The EM approach also gives better expected payoff, the statistical significance of which is confirmed by Wilcoxon's sign-rank test. Figure 2 has a box plot of regrets from different instances of data sets, which shows that the EM approach is consistently getting higher payoffs.

### 5.3 Synthetic Data Set 3: Unknown Distributions

We go back to the identical items model with stationary distributions $f(x)$ and $g(m)$. For this data set, $f(x)$ is a Gamma distribution with shape parameter 2 and scale parameter 3. $g(m)$ is a mixture of two Poisson distributions: $P(4)$ with probability 0.6 and $P(60)$ with probability 0.4. But now the estimation approaches does not know the types of the true distributions. Rather than guessing the types of distributions, we use kernel density estimation (kernel smoothing), a nonparametric estimation strategy. Essentially, given N samples from a distribution $p(x)$, we estimate $p(x)$ by a mixture of N kernel functions centered at the N samples.

A Gaussian kernel is used for estimating $f(x)$ and a uniform kernel is used for estimating $g(m)$. At each M step of the EM algorithm, the bandwidth parameters of the two kernel estimations need to be selected. We use the simple "rule of thumb" strategy [12] for bandwidth selection. The same type of kernel estimation and bandwidth selection technique is implemented for the simple approach.

Our results show that the EM approach gives better estimation than the simple approach does. Figure 3 shows typical estimated distributions and true distributions. From the figure we can observe that the EM estimates of $f(x)$, $g(m)$ and $f^1(x)$ are much closer to the true distributions that the simple estimates. The EM estimates have significantly smaller KL Divergences compared
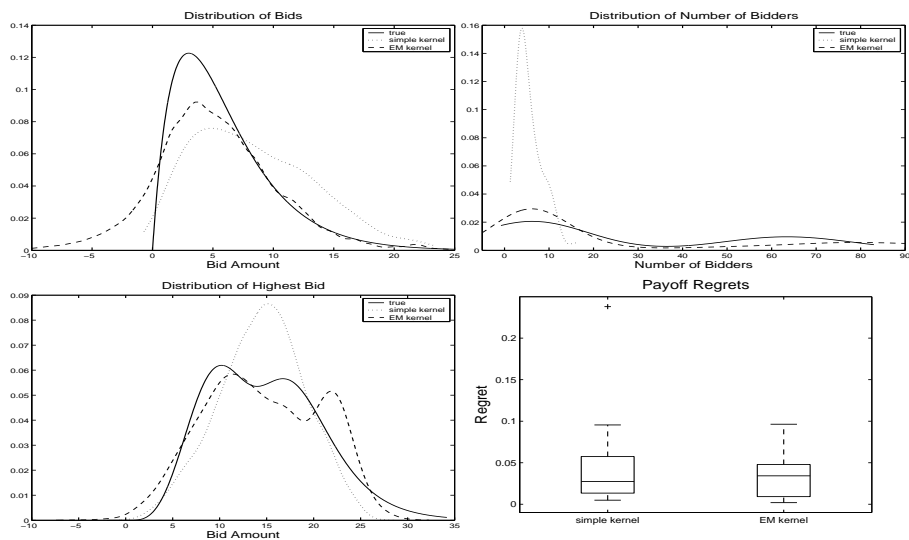
Figure 3: Results for Data Set 3: Distribution $f(x)$ (top-left). Distribution $g(m)$ (top-right). Distribution $f^1(x)$ (bottom-left). Box plot of payoff regrets (bottom-right).

to the simple estimates, verified by Wilcoxon's sign-rank test.

We then computed the expected payoffs under the decision theoretic model with 8 sequential auctions. The expected payoffs of the EM approach were not significantly better than that of the simple approach, as shown by the box plot in Figure 3. One possible explanation is that although the simple estimates were further from the correct distributions than the EM estimates, under this particular decision-theoretic model the bidding strategy computed using the simple estimates happened to achieve high payoffs. In other words, it is not because the EM approach did badly, rather it is because the simple approach happened to get high payoffs in this setting.

### 5.4 eBay Data on Sony Playstation-2 Systems

Our experiments on synthetic data sets showed that our EM approach gave good estimates of the true distributions in several different settings. However, the synthetic data sets are generated using our model for the bidding process. Thus, the above experiments do not tell us whether our model for the bidding process is an accurate description of what happens in real world online auctions. To answer this question, we need to test our approach on real world bid data. On eBay, the bidding histories of completed auctions are available for the most recent 30 days. Unfortunately, information on the hidden bids, especially the proxy bids of the winners of the auctions, is not publicly available. So unlike in the synthetic data experiments, we cannot compare our estimated distributions with the "true" distributions.

To get around this problem, we used the following approach: first we collect

13

bidding histories from a set of eBay auctions. Now we pretend that those highest bids were not placed, and the previously second highest bids are the highest bids. We can now "hide" these new highest bids of each auction, and use our estimation approaches to try to predict the distribution of the highest bid, $f^1(x)$. We can now compare our estimated distributions to the bids we have hidden, and also compute expected payoffs under the decision theoretic model. While it is true that this approach of hiding bids will introduce bias into our estimations, we are not trying to learn the true distribution of bids of those eBay auctions. Instead we are trying to learn the distribution of this "shifted" data set, which is nonetheless collected from the real world and thus should have similar characteristics to the actual bidding history data. If our model of the bidding process is correct, then our EM approach should be able to correctly account for the hidden bids in this data set and produce good estimates to $f^1(x)$.

We have collected bidding history of eBay auctions on brand new Sony Playstation 2 (Slim Model) consoles, over a period of a month in March 2005. Shah et al. [11] analyzed eBay auction data on an earlier version of Sony Playstation, where they argued that bidders' valuations on Playstations tend to be close to the private value model. We considered only auctions that lasted one day, and had at least 3 bidders. We observed 60 auctions that satisfied these requirements. The data was then randomly divided into a training set and a testing set. We tested four learning approaches: the EM and simple approaches that estimates a Normal distribution for $f(x)$ and a Poisson distribution for $g(m)$, and the EM and simple approaches that use kernel density estimation to estimate $f(x)$ and $g(m)$. Each approach tried to estimate $f^1(x)$ from the training set, and the estimates were compared against the highest bids from the test set. We did 8 runs of this experiment with different random partitions of training set and testing set, and aggregated the results. The KL Divergences of $f^1(x)$ of the approaches are similar, and no one approach is significantly better than the others.

We then computed the expected payoffs under the decision theoretic model. The EM approaches achieved significantly higher payoffs than the simple approaches, as shown in Figure 4. The approaches using parametric models achieved similar payoffs to the corresponding approaches with kernels. The good performance of the parametric estimation EM approach for the eBay data set indicates that the Normal and Poisson models for $f(x)$ and $g(m)$ may be adequate models for modeling bidding on eBay.

# 6 Conclusion and Future Work

In this paper we have focused on an important problem faced by bidding agents, that of estimating the distributions of the number of bidders and bid amounts from incomplete auction data. We proposed a learning approach based on the EM algorithm that takes into account the missing bids by iteratively generating missing bids and doing maximum likelihood estimates on the completed set of bids. We conducted experiments on both synthetic data as well as on eBay data, and compared our approach against the straightforward approach of ignoring
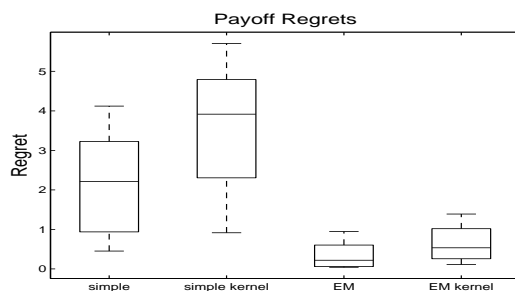
Figure 4: Box plot of payoff regrets on the eBay Data Set

the missing data. Our results show that our approach never did worse and often did much better than the simple approach, both in terms of the quality of the estimates and in terms of expected payoffs under a decision theoretic bidding model. Our EM learning approach is not limited to decision theoretic auction models; it can also be used by game-theoretic bidding agents. We are currently investigating such applications of our EM learning approach in a game-theoretic setting.

# References

[1] C. Andrieu, N. de Freitas, A. Doucet, and M.I. Jordan. An introduction to MCMC for machine learning. *Machine Learning*, 2003.

[2] P. Anthony, W. Hall, V.D. Dang, and N. Jennings. Autonomous agents for participating in multiple online auctions. In *Proc. of the IJCAI Workshop on EBusiness and the Intelligent Web*, 2001.

[3] A. Arora, H. Xu, R. Padman, and W. Vogt. Optimal bidding in sequential online auctions. Working Paper.

[4] C. Boutilier, M. Goldszmidt, and B. Sabata. Sequential auctions for the allocation of resources with complementarities. In *Proceedings of 16th IJCAI*, 1999.

[5] A Byde. A comparison among bidding algorithms for multiple auctions. In *Agent-Mediated Electronic Commerce IV*, 2002.

[6] G. Cai and P.R. Wurman. Monte Carlo approximation in incomplete-information, sequential-auction games. Technical report, North Carolina State University, 2003.

[7] A. Greenwald and J. Boyan. Bidding under uncertainty: Theory and experiments. In *Proc. UAI-04*, 2004.

[8] P. Klemperer. Auction theory: A guide to the literature. In P. Klemperer, editor, *The Economic Theory of Auctions*. Edward Elgar, 2000.

[9] P. Milgrom and R. Weber. A theory of auctions and competitive bidding, II. In P. Klemperer, editor, *The Economic Theory of Auctions*. Edward Elgar, 2000.

[10] A.E. Roth and A. Ockenfels. Last-minute bidding and the rules for ending second-price auctions: Evidence from ebay and amazon auctions on the internet. *American Economic Review*, 2002.

[11] H.S. Shah, N.R. Joshi, A. Sureka, and P.R. Wurman. Mining for bidding strategies on ebay. *Lecture Notes on Artificial Intelligence*, 2003.

[12] B.W. Silverman. *Density Estimation*. Chapman and Hall, London, 1986.

[13] R. Weber. Multi-object auctions. In R. Engelbercht-Wiggans, M. Shubik, and R. Stark, editors, *Auctions, Bidding, and Contracting: Uses and Theory*, pages 165–191. New York University Press, 1983.

[14] M.P. Wellman, A. Greenwald, P. Stone, and P.R. Wurman. The 2001 trading agent competition. In *Proceddings of IAAI*, 2002.