

Aggregation Bias in Sponsored Search Data:

The Curse and The Cure

Vibhanshu Abhishek<sup>†</sup>, Kartik Hosanagar<sup>‡</sup>, and Peter S. Fader<sup>‡</sup>

{vibs@cmu.edu, faderp@wharton.upenn.edu, kartikh@wharton.upenn.edu}

<sup>†</sup>Heinz College, Carnegie Mellon University

<sup>‡</sup>Wharton School, University of Pennsylvania

## Aggregation Bias in Sponsored Search Data: The Curse and The Cure

### Abstract

There has been significant recent interest in studying consumer behavior in sponsored search advertising (SSA). Researchers have typically used daily data from search engines containing measures such as average bid, average ad position, total impressions, clicks and cost for each keyword in the advertiser’s campaign. A variety of random utility models have been estimated using such data and the results have helped researchers explore the factors that drive consumer click and conversion propensities. However, virtually every analysis of this kind has ignored the intra-day variation in ad position. We show that estimating random utility models on aggregated (daily) data without accounting for this variation will lead to systematically biased estimates – specifically, the impact of ad position on click-through rate (CTR) is attenuated and the predicted CTR is higher than the actual CTR. We demonstrate the existence of the bias analytically and show the effect of the bias on the equilibrium of the SSA auction. Using a large dataset from a major search engine, we measure the magnitude of bias and quantify the losses suffered by the search engine and an advertiser using aggregate data. The search engine revenue loss can be as high as 11% due to aggregation bias. We also present a few data summarization techniques that can be used by search engines to reduce or eliminate the bias.

*Key words:* Sponsored search, generalized second-price auctions, consumer choice models, Hierarchical Bayesian estimation, latent instrumental variables, aggregation bias.

# 1 Introduction

Sponsored search advertising (SSA) has not only transformed the way companies conduct their marketing activities, but it has also been a tremendous resource to academic researchers who seek to better understand how consumers respond to such ads. A myriad of researchers have turned to SSA data to uncover new insights about consumer search (Ghose and Yang, 2009; Rutz and Bucklin, 2011), choice and related purchasing behaviors (Jeziorski and Segal, 2009; Yang and Ghose, 2010; Agarwal et al., 2011) and advertiser/search engine strategies (Animesh et al., 2009; Yao and Mela, 2011; Rutz et al., 2012). Many of these papers have used random utility models to study the effect of ad position, keyword length, presence or absence of brand name, etc. on the click-through and conversion rates of the ads.

Sponsored search refer to ads that are displayed alongside organic search results when a user issues a query at a search engine. The advertisers submit bids for keywords that are relevant to them, along with these ads.<sup>1</sup> When a user enters a query, the search engine identifies the advertisers bidding on keywords closely related to the query and uses data on bids and ad quality/performance to rank order the ads in a list of sponsored results. The most widely used pricing model is the *pay-per-click* model, in which the advertiser pays only when a user clicks on his ad. The advertiser's cost per click or CPC is determined using a generalized second price auction (GSP), i.e. whenever a user clicks on an ad at a particular position, the advertiser pays an amount equal to the minimum bid needed to secure that position.

Although SSA is a relatively new practice, it already has fairly well-established data standards associated with it. Most researchers who have modeled SSA-related issues have worked with a data structure such as the one illustrated in Table 1. Advertisers also obtain similar datasets from search engines and analyze them to design their bidding policies. In almost all cases, these data are aggregated to the daily level and contain summary statistics for the day, such as the number of ad impressions, average position of the ad, number of clicks received and the average CPC. It should be clear how this kind of dataset lends itself to the types of models mentioned above, as well as analysis of a variety of other customer behaviors (and related firm actions). But despite the creativity and methodological prowess that has been demonstrated in this growing body of

---

<sup>1</sup>The term *keyword* refers to term or phrase on which an advertiser bids. *Query* (or *query term*) is the search phrase entered by the consumer when conducting the search

**Table 1.** Sample dataset for a particular keyword

date	impressions	clicks	avg. pos	avg. bid	avg. CPC
01/11/09	180	1	19.33	1.00	0.30
01/12/09	202	0	18.42	1.00	0.00
01/13/09	223	5	8.19	2.00	1.24
01/14/09	198	3	7.94	2.00	0.89
01/15/09	197	5	8.08	2.00	1.21
01/16/09	321	21	2.00	3.00	2.12

literature, we believe that these modeling efforts are plagued by a major problem: an aggregation bias due to the way that the raw (search-by-search) data are “rolled” up into Table 1.

In practice, the position of an ad can vary substantially within a day and aggregated data fail to capture this variation. It is also widely known that the impact of position on CTR is non-linear. For example, an ad at the topmost position tends to receive a disproportionately large number of clicks as compared to ads on other positions. The convexity in the CTR, coupled with the intra-day variation in position suggests that the daily aggregation might lead to estimation bias. The goal of this paper is to provide a thorough evaluation of the nature of this bias, to demonstrate its effects and provide recommendations to eliminate the bias.

The paper makes the following contributions. Firstly, we show that applying a logistic model to aggregated SSA data can lead to biased estimation of the parameters of a random utility model. Due to the bias, the effect of position on CTR is attenuated and the predicted CTR is higher than the actual CTR. Surprisingly, we find that if all the advertisers use aggregate data, the consequence of the bias is entirely borne by the search engine. Secondly, we quantify the magnitude of the bias and measure its economic impact. While some researchers studying the use of aggregated consumer data in other marketing settings suggest that aggregation bias is not important (Gupta et al., 1996; Russell and Kamakura, 1994), others emphasize the need to address aggregation bias explicitly (Narayanan and Nair, 2012; Neslin and Shoemaker, 1989). Using a unique and large disaggregate dataset from a major search engine, we show that aggregation bias is very significant in this context. We find that the search engine losses can be as high as 11% on average due to aggregation bias. Our findings raise serious concerns for SSA researchers and practitioners and also question the adequacy of the data standards that have become common in SSA. Finally, we present alternative

data summarization techniques and modeling approaches that can reduce or eliminate the bias. We find that sharing harmonic means instead of arithmetic means or complementing arithmetic means by other summary information such as the variance can help reduce the bias significantly. These results are robust in our analysis of a large real-world dataset as well as simulations that attempt to vary several aspects of the sponsored search market.

The rest of the paper is organized as follows. Section 2 discusses related work and positions our work in the literature. In Section 3, we analytically prove the existence of the bias and build a game-theoretic model to study the economic impact of the bias. In Section 4, we analyze a large disaggregate dataset from a search engine using the HB-LIV model. We present the managerial implications of the bias in Section 5. In Section 6, we present some cures for the problem of aggregation – data summarization techniques and modeling approaches that can reduce the bias. Finally, we discuss the implications of the bias on research and practice and conclude the study in Section 7.

## **2 Related work**

There has been a considerable amount of work on auction design and consumer choice models in SSA (Weber and Zheng, 2007; Liu et al., 2007; Hao et al., 2009; Goldfarb and Tucker, 2007). More specifically, there are two streams of work that are closely related to our study, namely empirical research on consumer click and conversion behavior in SSA and work related to aggregation biases in choice models.

### **Empirical Research in Sponsored Search**

There has been a lot of recent interest in trying to understand the factors driving keyword performance in SSA. Craswell et al. (2008) and Ali and Scarr (2007) propose individual keyword-level models to study how consumers navigate sponsored links. Other researchers have used logit models to measure the influence of factors like ad position and keyword characteristics on consumer behavior in SSA (Rutz et al., 2012; Rutz and Bucklin, 2011; Ghose and Yang, 2009; Agarwal et al., 2011). Rutz et al. (2012) compare the performance of several logit models in predicting the conversions for various keywords. Their results show that keywords are heterogeneous in their conversion rate and

a significant portion of this variation can be explained by the presence of brand or geographical information in the keyword. In another paper, Rutz and Bucklin (2011) measure the spill-over effect of generic keywords on branded keywords. Ghose and Yang (2009) use a random effect logit model to understand the relationship between different metrics such as CTR, conversion rates, bid prices and ad position using the advertiser's aggregate data. They show that keywords containing retailer information have a higher CTR whereas keywords that are more specific or contain brand information have a lower CTR. Recent work by Agarwal et al. (2011) uses a logit model to show that although the CTR decreases with position, the conversion rate is non-monotonic in position. They point out that the topmost position is not necessarily the revenue maximizing position.

Most of this stream of research uses aggregate data to estimate the parameters of the model. The aggregate data obfuscate the variation in ad position and research in this area has overlooked this fact. Ignoring this variation can lead to potential biases in the estimation of parameters and ultimately affect the conclusions from these studies.

## **Aggregation Bias**

Though researchers have grappled with the issue of data aggregation for many years, there is no clear consensus on this issue. The problems associated with aggregation have been commonly encountered in spatial and demographic studies and are referred to as the Yule-Simpson effect (Good and Mittal, 1987). The drawbacks of aggregation have also been pointed out in various studies in the economics and marketing literature. Neslin and Shoemaker (1989) point out the limitations of aggregate data by refuting the claim that sales promotions undermine the consumer's repeat-purchase propensity. They show that even if the individual purchase propensities do not change before and after promotions, statistical aggregation would lead to lower average repeat probabilities for post promotional purchases. Yatchew and Griliches (1985) discuss the implications of aggregation in the context of probit models. Issues related to data aggregation in the case of logit models have been presented by Kelejian (1995). He discusses why aggregation bias might occur when logit models are estimated on aggregate data and proposes a test for the existence of this bias.

On the other hand, several researchers believe that the effect of aggregation is negligible or absent when the disaggregate model can be approximated by the aggregate model (Gupta et al.,

1996; Russell and Kamakura, 1994). Using household-level panel data and store-level purchase data, Gupta et al. (1996) show that the price elasticity estimated from the two models differ by a very small amount (4.7%). Allenby and Rossi (1991) present an analytical proof for the non-existence of aggregation bias in nested logit models of consumer choice when the products are close substitutes of each other, though they assume that the micro-level consumer behavior is approximately linear in the product attributes. More specifically, in the context of sponsored search, Rutz and Trusov (2011) use a latent instrumental variable approach that addresses endogeneity in search auctions. Their model has an added benefit that it addresses some aspects of aggregation bias, such as that due to aggregation of data from heterogeneous consumers with different preferences and that due to strategic bidding. However, as we show later in Section 5, the approach does not fully address aggregation the bias resulting from aggregated ad performance data when ads are shown in multiple positions.

The discussion reveals two themes. First, a number of recent studies have applied the logit model on aggregated SSA data to study consumer choice behavior. Second, although aggregation bias has been shown to exist in a number of environments, its non-existence has also been demonstrated in several other environments. It is not clear which of these arguments is most applicable in the SSA context, and it is thus not clear whether and to what extent aggregation bias affects SSA research. This paper uses a theoretical model to show why data aggregation might lead to biased estimates in the SSA context and how this bias affects the outcome in a search engine auction. An extensive disaggregate search engine dataset is used to empirically measure the extent of aggregation bias in SSA research, which we believe has been done for the first time. Finally, we suggest ways in which the bias can be reduced or eliminated.

### 3 Aggregation Bias

In this section, we explore the estimation bias due to the aggregation of SSA data. We begin by pointing out the distinction between the complete (disaggregate) and the summary (aggregate) data that have been referred to in the paper. Table 2 is a stylized example that presents impression level data for an ad, reflecting every search query for a particular term. Each observation contains the date on which the impression occurred, position of the ad, bid placed by the advertiser, whether

the consumer clicked on the ad and finally the CPC. Search engines usually do not provide such granular data to advertisers or researchers. They provide aggregated data at the daily level as shown earlier in Table 1 that mask the intra-day variation in position. Discussions with several search engines reveal that they do not provide impression level data because it is expensive for advertisers to store, manage and analyze such huge amounts of data. In addition, concerns about user privacy and click-fraud further reduce the incentive to provide disaggregate data.

The intra-day variation arises due to two major factors – Firstly, SSA auctions are extremely dynamic with advertisers entering and exiting the auction or changing their bids continuously. Changes in competitors’ behavior lead to changes in the ad position. Secondly, most of the ads, specifically *broad match* and *phrase match* ads are shown for a number of different queries.<sup>2</sup> As the set of competitors can be different for different queries, the position of the ad also varies across queries.

**Table 2.** Complete dataset for a particular keyword

impression	date	click	pos	bid	CPC
1	01/11/09	0	16	1	0.00
2	01/11/09	0	20	1	0.00
...	...	...	...	...	...
8190	02/15/09	0	6	2	0.00
...	...	...	...	...	...
9145	02/23/09	1	1	3	2.31

### 3.1 Analytical Proof of Aggregation Bias

#### Logit utility model

The logit utility model has been extensively used in economics and marketing to explain consumer choice behavior. Researchers have primarily focused on keyword-level models to analyze the effect that factors like ad position, specificity of the keyword and presence of brand name have on the

<sup>2</sup>An *exact match* occurs when the user’s query term exactly matches the advertiser’s keyword. A *phrase match* occurs when the advertiser’s keyword appears anywhere within the user’s query. Finally, a *broad match* occurs when user query is determined to be broadly similar to advertiser’s keyword. Broad match is commonly used by advertisers as it maximizes the number of ad impressions.



consumer’s propensity to click on the ad. The consumer’s utility has been modeled as

$$U = X'\boldsymbol{\beta} + \epsilon, \tag{1}$$

where  $X$  is a vector of covariates,  $\boldsymbol{\beta}$  is the consumer’s sensitivity to these attributes and  $\epsilon \sim \text{Logistic}(0, 1)$ . In binary choice models, this utility is not observed but constitutes a latent variable. The consumer clicks on an ad when  $U > 0$ . Variable  $Y \in \{0, 1\}$  denotes whether a click was made or not. In conjugation with prior research, we build a keyword-level model that ignores other ad characteristics and focuses our attention on the impact of ad position on click through rate, i.e.  $CTR = f(\text{position})'$ . The simple model allows us to clearly identify the existence and direction of the bias. However, this assumption does not impose any restrictions on the model as all keyword characteristics (which typically do not change during the day) are subsumed in the intercept term and the we focus our attention on position which varies intra-day. Although we focus on a keyword-level model in this paper, our finding are applicable for different levels of analysis.

### Estimation using complete dataset

We now discuss the estimation of  $\boldsymbol{\beta}$  when the model is estimated using the complete dataset. Let  $V_i$  be a random variable denoting the ad position on the  $i^{\text{th}}$  impression. We assume that  $V_i$  is independent and identically distributed and has a cumulative distribution function (c.d.f.) given by  $F_V(\cdot)$ , which is assumed to be constant during the period of observation.<sup>3</sup> The consumer’s utility is given by the following expression

$$U_i = \beta_0 + \beta_1 V_i + \epsilon_i. \tag{2}$$

As  $\epsilon_i$  is extreme value distributed, the probability of clicking on an ad (CTR) is  $p_i = 1/(1 + \exp -\{\beta_0 + \beta_1 v_i\})$ . Note that  $p_i$  might vary across impressions due to variation in the ad position. Let  $\hat{\boldsymbol{\beta}}_{\mathbf{c}}$  denote the maximum likelihood estimate from the complete dataset. It can be easily shown that  $\hat{\boldsymbol{\beta}}_{\mathbf{c}}$  is a consistent and unbiased estimator of  $\boldsymbol{\beta}$  (Hayashi, 2000, Proposition 7.6).

---

<sup>3</sup>Let  $\mu_V = \mathbb{E}[V]$  and  $\sigma_V^2 = \text{Var}(V)$ .

## Estimation using aggregate dataset

Researchers do not observe  $V_i$  when aggregate data are used. They only observe the mean daily position  $W$  which is given by

$$W = \frac{V_1 + V_2 + \dots + V_N}{N}, \quad (3)$$

where  $N$  is the random number of ad impressions on a particular day and  $V_1, \dots, V_N$  is the ad position during each of those impressions. The distribution of  $W$  is  $F_W(\cdot)$  and it depends on the distribution of  $V$  and  $N$ . If the effect of position is estimated from aggregate data, the consumer utility from clicking the ad is effectively modeled as

$$U_d = \beta_0 + \beta_1 W_d + \epsilon_d, \quad (4)$$

where  $W_d$  is the average position for the day and  $\epsilon_d$  is the logistically distributed error term. This formulation causes a mis-specification as the consumers do not observe the ad at a position  $W_d$  but at position  $V$ . As the variable  $Z = V - W_d$  ( $\mathbb{E}[Z|W] = 0$ ), which affects the consumer's click behavior, is not accounted for in the regression, the mis-specification is similar to omitted variables bias pointed out by Yatchew and Griliches (1985) and Wooldridge (2001). However, this issue arises primarily due to data aggregation and our approach is closely related to prior work in Marketing by Christen et al. (1997), Steenkamp et al. (2005) and Gupta et al. (1996). Further, as  $Var(Z)$  is not constant (it depends on the number of impressions in a days), the findings of Yatchew and Griliches (1985) and Wooldridge (2001) are not directly applicable in this context. Hence, we derive an important relationship between  $W$  and  $V$  in order to prove the aggregation bias.

**Lemma 1**  $W \leq_{cx} V$ , i.e.  $W$  is less than  $V$  in convex order. <sup>4</sup>,

This relationship between  $W$  and  $V$  is very general and holds for any distribution,  $F_V(\cdot)$ . Using Lemma 1, we prove an important result of this paper.

**Proposition 1**  $\hat{\beta}_s$  is a biased estimator of  $\beta$ .

---

<sup>4</sup> $X$  is less than  $Y$  in *convex order* if  $E[f(X)] \leq E[f(Y)]$  for all real convex functions  $f$  such that the expectation exists. All proofs appear in the appendix.

If  $\hat{\beta}_c$  is equal to  $\hat{\beta}_s$  then the convex order between  $W$  and  $V$  implies that

$$\mathbb{E} \left[ \frac{\exp\{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V\}}{1 + \exp\{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V\}} \right] > \mathbb{E} \left[ \frac{\exp\{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W\}}{1 + \exp\{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W\}} \right].$$

Since both the L.H.S. and the R.H.S. equal the overall *CTR* during the observation period (as shown in the appendix), this inequality is incorrect and hence  $\hat{\beta}_c$  cannot equal  $\hat{\beta}_s$ . Since  $\hat{\beta}_c$  is a consistent and unbiased estimator of  $\beta$ ,  $\hat{\beta}_s$  is biased. This finding is contrary to earlier work by Allenby and Rossi (1991), Gupta et al. (1996) and Russell and Kamakura (1994), which prove that aggregation bias in market or store-level scanner data is negligible. Aggregation bias is significantly reduced in their context as products are very close substitutes of each other and the consumers (or house-holds) are exposed to very similar marketing activities. However, position has a very strong effect in sponsored search (Craswell et al., 2008) and ads in different positions may be perceived very differently by consumers. Coupled with variation in ad position, aggregation bias can be quite substantial, which is formalized in the following proposition.

**Proposition 2** *The direction of aggregation bias is such that (i) the CTR estimated from the summary data is greater than or equal to the actual CTR at any position, (ii)  $\hat{\beta}_{1,s} > \hat{\beta}_{1,c}$ , and under certain conditions, (iii)  $\hat{\beta}_{1,s} \xrightarrow{p} \beta_1 / \sqrt{\beta_1^2 \{ \sigma_V^2 + (3\sigma_V^2 + \mu_V^2)\phi \} + 1}$ , where  $\phi = \mathbb{E}[1/N]$ .<sup>5</sup>*

$\hat{\beta}_s$  is biased and it predicts a CTR that is higher than the actual CTR at any position. Incorrect estimation of CTR might lead advertisers to make suboptimal choices in sponsored search auctions. The second part of Proposition 2 is consistent with Wooldridge (2001) which shows that the estimates are scaled towards zero. Furthermore, the omitted variable  $Z$  increases the disturbance term in the regression. The estimated error can be computed as a convolution of  $\epsilon$  and  $Z$ , but this is analytically intractable as  $\epsilon$  is logistically distributed. However, Proposition 2(iii) holds when  $\beta_1 Z + \epsilon$  can be approximated closely by a logistic distribution. Proposition 2 shows that, if the variation in the intra-day position is known, the actual  $\beta_1$  can be approximated by multiplying  $\hat{\beta}_{1,s}$  by the scaling factor. However, this approach suffers from several simplifying assumptions that are required for analytical tractability. We provide more general and robust empirical methods to remove the bias in Section 5. Although we prove that the bias exists for a logistic model, it is

<sup>5</sup>The expectation of the inverse of the number of daily impression,  $\mathbb{E}[1/N]$ , is bounded by  $(1 - (1 - e^{-\lambda}))/\lambda \leq \phi \leq (1 - [1 - F_\lambda(1) + 3(1 - F_\lambda(2))/m])$  where  $\lambda = \mathbb{E}[N]$  and  $F_\lambda(\cdot)$  is the CDF of  $N$ . When there are a large number of daily impressions, i.e.  $\lambda \rightarrow \infty$ ,  $\phi \rightarrow 1$ .

easy to show that these results would be applicable in any binary choice model where the choice behavior is convex in position.

### 3.2 Effect on Equilibrium Behavior

As advertisers use estimates from historical data to bid in SSA auctions, incorrect estimation of the CTR might have a negative impact on their revenues. In this section, we build a game-theoretic model to analyze the impact of aggregation on the advertisers' and search engine's revenues.

For our analysis, we make the following assumptions – (i) there are  $K$  advertising slots and  $K+1$  advertisers, (ii) the advertisers' valuation for a click,  $s_k$  are drawn from a continuous distribution with support on  $[0, \infty)$ , (iii) advertisers know their own valuation and the distribution of competing bids and finally, (iv) the click-through rate,  $\alpha_i$ , decreases with the position  $i$ . In addition, we also assume that advertisers estimate  $\alpha_i$  from historical data and are unaware of the aggregation bias. The advertisers are indexed in decreasing order of their valuations, i.e.  $s_1 > s_2 > \dots > s_{K+1}$  and their bids are  $b_1, \dots, b_{K+1}$ , respectively.<sup>6</sup> In addition, let  $h = (b_i, \dots, b_{K+1})$  refers to the history of bids prior to assignment of position  $i$ .<sup>7</sup> The case where complete data are used is analyzed first. Here the advertisers correctly estimate  $\alpha_i$ . Edelman et al. (2007) show that under the above assumptions, there exists a unique *envy-free* perfect Bayesian equilibrium and the optimal strategy for advertiser  $k$  under this equilibrium is to bid as follows,

$$b_k(s_k, i, h) = s_k - \frac{\alpha_i}{\alpha_{i-1}}(s_k - b_{i+1}). \quad (5)$$

This is the maximum CPC that the advertiser is willing to pay to move to position  $i-1$  and receive more clicks. At this point, Advertiser  $k$  is indifferent between getting position  $i-1$  at a CPC of  $b_k(s_k, i, h)$  and position  $i$  at  $b_{i+1}$ . This is an ex-post equilibrium, i.e., it is optimal for Advertiser  $k$  to follow the equilibrium strategy for any realization of other advertisers' valuations. The search-engine revenue is  $\Pi_S^C = \sum_{i=1}^K \alpha_i b_{i+1}$  and the payoffs for advertiser  $i$  is  $\Pi_i^C = \alpha_i(s_i - b_{i+1})$ . It should be noted that this equilibrium ensures an assortive match, i.e., if  $s_i > s_j$  then advertiser  $i$  bids higher than advertiser  $j$  and occupies a slot above advertiser  $j$  in equilibrium. This case serves as a reference for the ensuing discussion.

---

<sup>6</sup> $b_{K+1}$  is set to 0.

<sup>7</sup>Positions  $K$  through  $i+1$  are assignment before bidding for position  $i$  starts.

Next, we consider the case in which aggregate data are used. Let the CTR estimated from aggregate data be denoted as  $\alpha'_i$ .

**Remark 1** *When aggregate data are used to estimate CTR, the ratio  $\alpha_i/\alpha_{i-1}$  is overestimated due to the presence of aggregation bias.*<sup>8</sup>

Due to this overestimation, advertisers might bid incorrectly. As the equilibrium considered here is ex-post, the advertisers' bidding strategies depend neither on their beliefs about each others valuations nor on the fact that some advertisers might be using aggregate data. The bidding strategies continue to be similar to the one outlined in Equation (5), but the bids in this case,  $b'_1, \dots, b'_K$ , might be different.<sup>9</sup> We consider two extreme cases to study the impact of aggregation bias – (i) all advertisers except one use complete data and, (ii) all advertisers use aggregate data. Let the search-engine revenue in Case I(II) be denoted by  $\Pi_S^{AI(II)}$  and advertisers' profit by  $\Pi_i^{AI(II)}$ .

**Case I:** Suppose advertisers other than Advertiser  $j$  have access to complete data and can compute  $\alpha_i$  correctly. Only Advertiser  $j$  uses aggregate data and overestimates  $\beta_1$ . This leads him to overestimate the ratio  $\alpha_i/\alpha_{i-1}$  and he bids in the following manner.

$$b'_j(s_j, i, h) = s_j - \frac{\alpha'_i}{\alpha'_{i-1}}(s_j - b'_{i+1}). \quad (6)$$

As Advertiser  $j$  bids lower in equilibrium, he occupies a position  $j' \geq j$ . The following proposition characterizes the equilibrium in this case (detailed analysis and proofs are provided in the appendix).

**Proposition 3** (i) *If only Advertiser  $j$  uses aggregate data, the top advertisers ( $i \leq j$ ) bid lower, advertisers in between ( $j < i \leq j'$ ) bid higher and the remaining advertisers ( $i > j'$ ) bid the same as they would have when everyone had complete data.* (ii) *The payoffs of the search engine and Advertiser  $j$  decrease ( $\Pi_S^{AI} < \Pi_S^C, \Pi_j^{AI} < \Pi_j^C$ ) while all other advertisers receive payoffs that are either the same or higher than payoffs they would have received if all advertisers were using complete data ( $\Pi_i^{AI} \geq \Pi_i^C, i \neq j$ ).*

---

<sup>8</sup>Writing the CTR in terms of the logit model we get,  $\frac{\alpha_i}{\alpha_{i-1}} = \frac{\exp\{\beta_0 + \beta_1 i\}}{1 + \exp\{\beta_0 + \beta_1 i\}} \times \frac{1 + \exp\{\beta_0 + \beta_1(i-1)\}}{\exp\{\beta_0 + \beta_1(i-1)\}} \approx \exp\{\beta_1\}$  when CTRs are small. Since  $\hat{\beta}_{1,s} > \hat{\beta}_{1,c} \Rightarrow \alpha_i/\alpha_{i-1} < \alpha'_i/\alpha'_{i-1}$ .

<sup>9</sup>We continue to assume that Advertiser  $K + 1$  still bids 0.

Advertiser  $j$  underestimates the impact of position and incorrectly bids less, which might move him to a lower position. In turn, some advertisers who were below him move up one position. The ordering of these advertisers does not change, which is a consequence of the bidding policy (Edelman et al., 2007). As  $b'_j < b_{j'}$ , bids required to acquire all positions above  $j'$  decrease. As bids for all position are (weakly) lower, the search engine loses revenue. Clearly, Advertiser  $j$ 's payoff is lower, because he deviates from the optimal policy. However, the loss in revenue for the search engine is substantially higher than the loss in revenue for the advertiser using aggregate data. Interestingly, all these losses are transferred to the other advertisers ( $\neq j$ ) as excess surplus since GSP is a zero-sum game. Hence, the search engine suffers the most due to aggregation bias and all advertisers apart from  $j$  are better off due to aggregation. In the subsequent case, we observe that the search engine internalizes all the negative impact of aggregation.

**Case II:** When all advertisers use aggregate data, their estimate of the CTR,  $\alpha'_i$  are greater than the actual CTR as shown in Proposition 3. For simplicity, we assume that all advertisers arrive at the same estimates for  $\alpha'_i$ .<sup>10</sup> As a result, Advertiser  $k$  adopts the following bidding strategy,

$$b'_k(s_j, i, h) = s_k - \frac{\alpha'_i}{\alpha'_{i-1}}(s_k - b'_{i+1}).$$

It is easy to see that the bid placed by advertiser  $K$  (the last advertiser) is less than the bid he would have placed had he estimated CTR from complete data. Proceeding in an iterative fashion we can show that all advertisers place a lower bid. The equilibrium in this case is specified in the following proposition.

**Proposition 4** (i) *When all advertisers use aggregate data, the advertisers are arranged in as-sortive order. The resulting bids are lower than the bids when complete data are used ( $b'_i < b_i, i = 1, \dots, K$ ).* (ii) *Search-engine revenue is lower ( $\Pi_S^{A2} < \Pi_S^C$ ) and advertisers' payoff are higher ( $\Pi_i^{A2} > \Pi_i^C$ ) as compared to the complete case.*

In the appendix, we show that the advertisers bid less than what they would have had they known the actual CTR. As all the advertisers use the same incorrect estimate of the CTR, the

---

<sup>10</sup>This result continues to hold even if the advertisers arrive as different estimates of  $\alpha'_i$  as long as  $\alpha_i/\alpha_{i-1} < \alpha'_i/\alpha'_{i-1}$ , which always hold true due to aggregation bias as shown earlier.

eventual ranking remains the same as in the complete case. They receive the same number of clicks but at a lower CPC and hence their payoffs are higher. Surprisingly, the search-engine revenue suffers the most when all advertisers use aggregate data even though the advertisers make the wrong decisions. These results question the data standards that have become common in SSA and underscore the need to provide better data to advertisers. We also show that it is incentive compatible for the search engine to provide richer/better data to advertisers.

Note that an advertiser always receives a higher payoff when he uses complete data as compared to aggregate data, irrespective of the fraction of advertisers using aggregate data. This intuition is formalized in the following proposition.

**Proposition 5** *An advertiser can always increase his payoff from SSA by unilaterally using complete data instead of aggregate data.*

The difference in payoff between the two cases (complete v/s aggregate) can be considered as the value of complete data or alternately the disutility from aggregate data. Although advertisers cannot get impression level data, they can periodically crawl the search engine and estimate the empirical distribution of the ad position. Moreover, with recent improvements in ad tagging and user tracking techniques, advertisers might be able to collect impression level data for a few customers.

## 4 Empirical Analysis

In the previous section, we analytically show that aggregation at a daily level leads to a bias. However, as pointed earlier, several papers show that aggregation bias is negligible in various marketing data (Allenby and Rossi, 1991; Gupta et al., 1996; Russell and Kamakura, 1994). They argue that aggregation bias is significantly reduced because the products considered in their analysis are very close substitutes. These findings might not hold in SSA as Craswell et al. (2008), Ghose and Yang (2009) and Abhishek and Hosanagar (2013) show that ad position has a very strong effect in SSA. We perform the following empirical analysis on large representative search engine data to examine and conclusively prove that ad position has a strong influence on consumer click behavior, which leads to significant aggregation bias. Furthermore, we want to measure the extent of the bias and observe its economical significance.

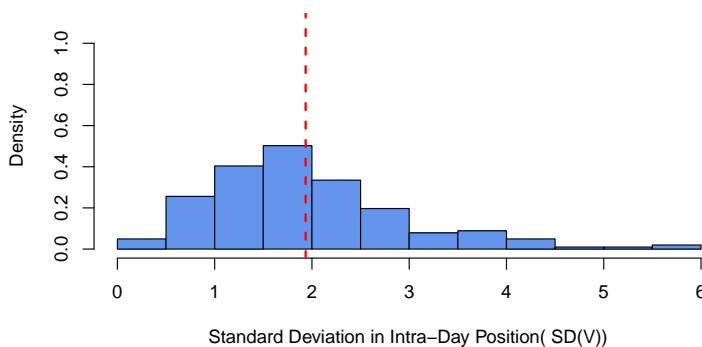
## 4.1 Data Description

We analyze a large disaggregate data from a major search engine, which is extremely representative of consumer behavior in SSA. The data contain around 8 million unique impressions chosen randomly from all user queries between August 10, 2007 and September 25, 2007. These are very unique data as search engines rarely provide impression level data to advertisers or researchers. For every impression, the dataset contains the user query, ads shown on the page and number of ads on the preceding pages. Each ad is identified by a unique ad identifier, though the dataset does not contain any ad-specific information. The dataset also contains information about clicks during this period of observation. We construct an ad-level dataset, that contains information about the keyword the advertiser was bidding on and all the impressions of the ad associated with the keyword, which is similar to the one presented in Table 2.

**Table 3.** Summary Statistics

Total impressions	8,142,210
Unique queries	24,235
Unique ads	229,960
Ads with more than one impression	184,481
Mean impressions for every ad	64.4
Median impressions for every ad	7.0

There is evidence that there is substantial variation in position and reporting average position alone results in the loss of information on actual position as shown in Figure 1. We next investigate the impact of data aggregation.



**Figure 1.** Mean Intra-day Variation in Position.

The ad level data described earlier are summarized at a daily level to create aggregate data.



The data thus generated are similar to the campaign summaries that search engines make available to the advertisers (as presented in Table 1).

## 4.2 Hierarchical Bayesian Model

We estimate a random-effect logit model using Hierarchical Bayesian (HB) techniques that are commonly used in SSA. As our data do not contain any ad-specific attributes, the only covariate included in our models is position. The effect of ad characteristics is captured in the ad-specific intercept term. We demonstrate the aggregation bias for HB model.

We extend the binary choice logit model proposed earlier in Section 3 to account for multiple keywords. Under this specification, the consumer’s utility from clicking on ad  $k$  during impression  $i$  is given by

$$U_{ik} = \beta_{0k} + \beta_{1k}V_{ik} + \epsilon_{ik} \quad (7)$$

where  $\epsilon_{ik}$  is the idiosyncratic, logistically distributed error term.  $\beta_{\mathbf{k}} = (\beta_{0k}, \beta_{1k})'$  are keyword specific parameters which are assumed to be random and heterogeneous across ads. They drawn from a multivariate normal distribution in the following manner:

$$\beta_{\mathbf{k}} \sim N_2(\boldsymbol{\mu}_{\beta}, V_{\beta}) \quad \text{where} \quad V_{\beta} = \begin{pmatrix} \sigma_{\beta_0} & \sigma_{\beta_0\beta_1} \\ \sigma_{\beta_0\beta_1} & \sigma_{\beta_1} \end{pmatrix}.$$

Similar models have been extensively used in prior research in SSA (Ghose and Yang, 2009; Yang and Ghose, 2010). Note that although this random coefficient model captures heterogeneity across ads, it still fails to accounts for the intra-day variation in  $V_{ik}$  if the model is estimated on aggregate data. As a result, we expect the aggregation bias to extend to the random-coefficient model as well. To test this hypothesis, we take a random sample of ads from our data and apply the model in Equation (7) to both the disaggregate and aggregate datasets and compare the estimates.

The log-likelihood function for the complete data is as follows

$$LL(\boldsymbol{\beta}|\text{complete data}) \propto \sum_{k=1}^K \sum_{i=1}^{I_K} Y_{ik} \log p_{ik} + (1 - Y_{ik}) \log(1 - p_{ik}), \quad (8)$$

where  $Y_{ik}$  is the indicator variable that denotes whether the  $i^{th}$  impression of keyword  $k$  received a

click or not and  $p_{ik}$ , the click-through probability is given by

$$p_{ik} = \frac{\exp\{\beta_{0k} + \beta_{1k}v_{ik}\}}{1 + \exp\{\beta_{0k} + \beta_{1k}v_{ik}\}}. \quad (9)$$

The log-likelihood function for the aggregate data is as follows

$$LL(\boldsymbol{\beta}|\text{aggregate data}) \propto \sum_{k=1}^K \sum_{d=1}^D c_{dk} \log p_{dk} + (n_{dk} - c_{dk}) \log(1 - p_{dk}), \quad (10)$$

where  $n_{dk}$  and  $c_{dk}$  denote the number of impressions and clicks on day  $d$  respectively and  $p_{dk}$ , the click-through probability is given by

$$p_{dk} = \frac{\exp\{\beta_{0k} + \beta_{1k}w_{dk}\}}{1 + \exp\{\beta_{0k} + \beta_{1k}w_{dk}\}}. \quad (11)$$

As the data on clicks are often sparse for most keywords in sponsored search, the SSA literature primarily uses Hierarchical Bayesian models. We use a similar approach and assume that the mean and variance-covariance matrix for  $\boldsymbol{\beta}_k$  have the following priors

$$\boldsymbol{\mu}_\beta \sim N_2(\boldsymbol{\mu}, \Sigma), \quad (12)$$

$$V_\beta^{-1} \sim \text{Wishart}(\nu, \Delta). \quad (13)$$

The parameters  $\boldsymbol{\mu}, \Sigma, \nu, \Delta, \boldsymbol{\mu}_\beta$  and  $V_\beta^{-1}$  are estimated separately from the complete and the aggregate datasets using a Markov Chain Monte Carlo (MCMC) approach. Before discussing the details of the MCMC estimation procedure, we discuss some identification issues associated with the model presented here.

### 4.3 Identification

The ad position in the previous exposition has been assumed to be exogenous. However, the position is decided by the bids placed by the advertiser. In addition, we know that past performance affects the quality score of the ad which in turn affects the position. The auction process and historical performance jointly determine the position which is one of the most important strategic variables that advertisers focus on in SSA. This indicates that position is endogenous (i.e.  $\mathbb{E}[\text{pos}_t \epsilon_t] \neq 0$ )

and the endogeneity should be explicitly incorporated in the HB model presented earlier.<sup>11</sup>

Endogeneity has been a major concern in the SSA literature and researchers have proposed several techniques to address this issue. Ghose and Yang (2009) and Yang and Ghose (2010) use a simultaneous equation model to address this problem. Their simultaneous model forms a triangular system of equations which can be identified without any further identification constraints. Agarwal et al. (2011) use a series of random bids to address the endogenous nature of position. In their specification, position is completely determined by the random bids and quality score which are exogenous. Recent econometric advances have led to the development of the latent instrument variable (LIV) framework (Ebbes et al., 2005), which has been used by Rutz and Trusov (2011) and Rutz et al. (2012) to account for position endogeneity. The LIV framework uses a likelihood based approach, which can be easily integrated with the HB model proposed earlier, and can be estimated using the MCMC estimator.

In a LIV formulation, the endogenous covariate is decomposed into a stochastic term that is uncorrelated with the error and another one that is possibly correlated with the error, i.e.  $X = \theta + \eta$  where  $\theta$  is the uncorrelated part of  $X$  such that  $\mathbb{E}[\theta\epsilon] = 0$  and  $\mathbb{E}[\eta\epsilon] = \sigma_{\eta\epsilon}$ . As  $\theta$  varies in the dataset, it is possible to identify the correlation between  $\eta$  and  $\epsilon$ , denoted by  $\sigma_{\eta\epsilon}$ . For the sake of simplicity, we modify the model presented in Equation (7) such that,  $\epsilon_{ik} = \zeta_{ik} + \vartheta_{ik}$ , where  $\zeta_{ik}$  is correlated with ad position and  $\vartheta_{ik}$  is orthogonal to position and logistically distributed. Clearly,  $\mathbb{E}[\eta\zeta] = \mathbb{E}[\eta\epsilon] = \sigma_{\eta\epsilon}$ , which we denote by  $\sigma_{\eta\zeta}$  for the sake of exposition. The LIV approach is extended to binary choice models by introducing a latent categorical variable with  $M$  categories (Rutz et al., 2012). Position can assume any of these  $M$  categorical values with a probability  $\Pi = \{\pi_1, \pi_2, \dots, \pi_M\}$ ,  $\sum_{m=1}^M \pi_m = 1$ . More specifically,

$$pos_{tk} = \Theta'_{kt}\gamma + \mathbf{Z}'_{kt}\delta + \eta_{kt}, \quad (14)$$

where  $\Theta \sim Multinomial_M(\Pi)$  is the stochastic part of  $pos_{kt}$  and is exogenous whereas  $\eta_{kt}$  is

---

<sup>11</sup>We thank the anonymous reviewer for this suggestion.

endogenous.<sup>12</sup> The errors  $(\eta_{kt}, \epsilon_{kt})'$  are MVN distributed in the following manner

$$\begin{pmatrix} \eta_{tk} \\ \zeta_{tk} \end{pmatrix} = MVN \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_\eta^2 & \zeta_{\eta\zeta}^2 \\ \sigma_{\eta\zeta}^2 & \sigma_\zeta^2 \end{pmatrix} \right]. \quad (15)$$

$\mathbf{Z}_{kt}$  represents the observed instruments, and in our analysis we use lagged position as IVs which is similar to the approach adopted by Rutz et al. (2012). Similar to Rutz et al. (2012), we estimate the parameters  $\Pi, \Theta$  and  $V_\xi$  jointly for the entire observation period when the estimation is performed on aggregate data. In the case of complete data, we exploit the richness of the data to estimate the parameters,  $\Pi, \Theta$  and  $V_\xi$ , on a daily basis. Such an approach can account for variations over time, due to changes in advertisers' bids or other auction related factors, and might lead to better performance. The exact difference between estimation on the complete versus aggregate data is explained in detail in the Online Appendix.

#### 4.4 Estimation Results

We estimate both the HB model and HB model with LIV (HB-LIV) to draw comparisons between the two methods. A sample size of 200 ads is chosen for estimating the parameters. We make this choice primarily for computational convenience as estimating the model on disaggregate data takes a long time. The disaggregate dataset contains a large number of observations, hence the estimation on the disaggregate dataset is really slow.<sup>13</sup>

We start off with diffused priors ( $\mu = 0, \Sigma = 100I, \nu = 5, \Delta = \nu I$ ) and refine them as the estimation proceeds. The exact estimation procedure is outlined in the Online Appendix. We run the MCMC simulation for 100,000 draws and the first 50,000 sample are discarded. The MCMC chains are stationary after the burn-in period. The MCMC chains are thinned to remove autocorrelation between draws and every 10th draw in the stationary period is used for the subsequent analysis.

The estimation results are presented in Table 4. We observe that there are significant differences in the  $\mu_\beta$  estimated on the complete and aggregate data, for both the HB and HB-LIV models.  $\mu_1$  is overestimated by 10% when the estimation is performed on aggregate data indicating that aggregation bias exists when a HB model is used. The bias in  $\mu_1$  increases to 12.9% when a HB-LIV

<sup>12</sup> $pos_{tk} = v_{ik}$  or  $w_{dk}$  depending on the context, where  $t$  represents the unit of time.

<sup>13</sup>We estimate the HB model on 25 different samples and the qualitative findings remain the same.

**Table 4.** Estimates of Parameters

Parameters	HB		HB-LIV	
	Complete	Aggregate	Complete	Aggregate
<b><math>\mu_\beta</math>:</b>				
$\mu_{0\beta}$	-1.495(0.000)	-1.459(0.001)	-1.672(0.301)	-1.612(0.189)
$\mu_{1\beta}$	-0.793(0.085)	-0.727(0.098)	-0.642(0.120)	-0.558(0.074)
<b><math>V_\beta</math>:</b>				
$\sigma_{\beta_0}$	0.654(0.128)	0.753(0.170)	0.678(0.107)	0.689(0.192)
$\sigma_{\beta_1}$	0.153(0.037)	0.155(0.038)	0.147(0.033)	0.162(0.042)
$\sigma_{\beta_1\beta_2}$	0.025(0.067)	-0.085(0.076)	0.019(0.052)	-0.090(0.076)
<b><math>V_\xi</math>:</b>				
$\sigma_\zeta$	0.263(0.067)	0.389(0.082)	0.142(0.047)	0.196(0.065)
$\sigma_\eta$			1.733(0.238)	2.113(0.414)
$\sigma_{\eta\zeta}$			0.136(0.039)	0.176(0.052)
<b><math>\Theta</math>:</b>				
$\theta_1$			0.962 (0.083)	0.748 (0.091)
$\theta_2$			2.838 (0.029)	2.364 (0.037)
$\theta_3$			4.552 (0.012)	5.927 (0.045)
<b><math>\Pi</math>:</b>				
$\pi_1$			0.523 (0.255)	0.322 (0.181)
$\pi_2$			0.238 (0.173)	0.193 (0.127)
<b>Instrument Variable:</b>				
$pos_{t-1}$			0.885(0.096)	0.766(0.127)

Note that the values reported for parameters  $V_\xi$ ,  $\Theta$  and  $\Pi$  in the complete data case are means of the daily estimates.

model is used. This results strongly indicates that aggregation bias is significant in the context of SSA. This finding is of concern as the extant literature on SSA (Agarwal et al., 2011; Ghose and Yang, 2009; Yang and Ghose, 2010) is silent about this issue and only very recently, Rutz and Trusov (2011) acknowledge its existence and account for it in their model. Rutz and Trusov (2011) significantly advance the SSA literature by proposing a novel methodology to address endogeneity which also addresses aggregation to some extent. However, as demonstrated in Table 4, the aggregation bias is not completely eliminated, which points to the value of complete information. Aggregate data is prone to two primary disadvantages – (i) it does not capture intra-day variations, and (ii) it is difficult to identify temporal patterns due to the limited amount of data (one observation per day). When the estimation is performed on complete data, these variations can be explicitly captured and accurately estimated underscoring the importance of richer data in SSA.

Our findings demonstrate that this bias is non-trivial in practice and future research in SSA should be informed about the pitfalls of aggregate data. In conjunction with prior literature, we also observe a statistically significant difference in the estimates from the HB and the HB-LIV model, confirming the need to control for endogeneity of ad position in sponsored search. Aggregation bias acts in addition to the bias due to endogeneity, and it is just as important to correct.

We use the differences between the estimates from summary and complete data for a random sample of 5000 exact-match keywords to compute the empirical distribution of the error ( $\varepsilon = \hat{\beta}_s - \hat{\beta}_c$ ) due to aggregation. This empirical distribution is used in Section 5 to quantify the impact of aggregation bias on search engine and advertiser revenues.

## 5 Managerial Implications

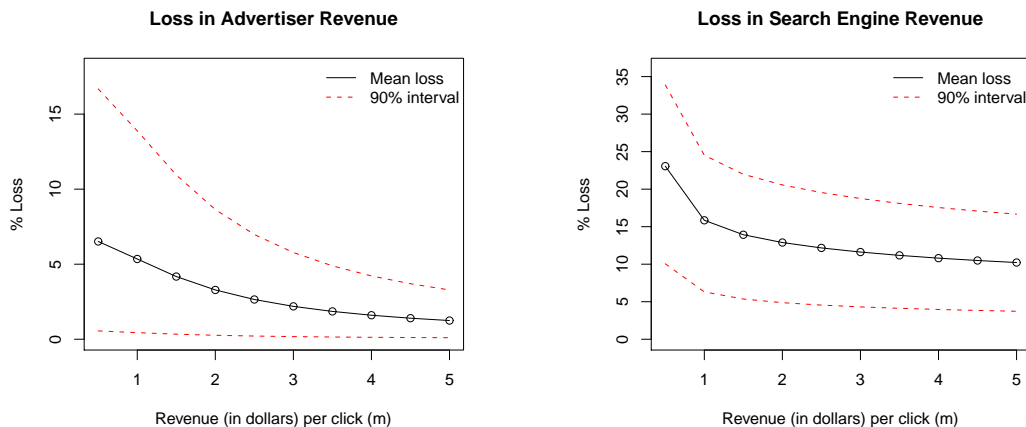
In the previous section, we show the existence of aggregation bias using the HB and HB-LIV models. In this section, we discuss the managerial implications of the bias. The analysis presented in Section 3.2 provides evidence that aggregation bias can lead to a decrease in the search engine and possibly an advertiser’s revenue. Here, we extend the analysis to characterize the loss, exploiting the large dataset available to us. Suppose  $m$  is the expected revenue *per-click*, then the advertiser’s profit per-impression is given by

$$\Pi = CTR(CPC) \times \{m - CPC\}. \quad (16)$$

It is easy to see the trade-off between bidding high to get more clicks ( $CTR$  increases with  $CPC$ ) and bidding low to earn greater profit per click. The optimal bids can be computed by substituting the expression for  $CTR$  (as a function of advertisers  $CPC$ ) into the profit function. Unfortunately, our data do not contain bidding information ( $CPC$ ), and as a result, we use estimates commonly found in the extant literature to illustrate the magnitude of the loss. We assume that the relationship between position and  $CPC$  is given by  $pos = e^{2(1-CPC)}$  (Ghose and Yang, 2009) and the  $CTR$  for the average keyword is a logit with  $\hat{\beta}_c = (-1.672, -0.642)'$  (Table 4). The optimal bid can be derived by substituting these relationships in Equation (16). When the advertiser uses aggregate data, his estimate of  $\beta$  is given by  $\hat{\beta}_s = \hat{\beta}_c + \varepsilon$ , where  $\varepsilon$  is the estimation error computed in the previous section. We sample  $\hat{\beta}_s$  from this empirical distribution, and for each  $\hat{\beta}_s$ , compute the optimal bid that maximizes the advertiser’s profit in Equation (16). The bids computed using  $\hat{\beta}_c$

are greater than almost all the bids computed using samples of  $\hat{\beta}_s$ , which supports our earlier claim that aggregation bias results in an advertiser placing a lower-than-optimal bid.

Next, we use the computed *CPCs* to estimate the effect of aggregation bias on the advertiser’s and search-engine’s revenues. Figure 2 shows the % loss suffered by the advertiser due to aggregation bias. There is a great deal of variation in this loss as the valuation changes. The impact of aggregation bias is more pronounced when the advertiser’s valuation for the click is low. In this situation, he bids lower and even small deviations from the optimal bid can lead to significant changes in the position leading to significantly lower payoffs. On the other hand, when the valuation is high, the bid is correspondingly higher and deviations from the optimal bid have little impact on position, and hence the loss is relatively smaller.



**Figure 2.** Loss in revenues caused due to Aggregation bias.

The effect of aggregation is considerably higher for a search engine. On average, it loses more than 11% of its payments from the advertiser as a result of aggregation bias. Lower bids, as a consequence of aggregation bias, negatively impact search engine revenues in two ways – Firstly, lower bids imply that the search engine generates lesser revenue *per-click*. Secondly, the ad appears at a lower position due to the lower bid, which in turn leads to fewer clicks. The advertiser pays the search engine for fewer clicks and pays less for each click. In this example, we estimate that the search engine loses 1.4¢ for every impression of the representative ad.

In the preceding, analysis we consider how payments from *one advertiser* to the search engine are affected due to aggregation bias. Naturally, the dynamics are more complicated when we try to quantify the effect that the bias has on payments by all advertisers and the overall profitability

of the search engine. For example, when an advertiser bids sub-optimally and moves to a lower position, another advertiser moves up to occupy the vacant position. Though the search engine loses revenue from the advertiser that moves down, it earns more from the advertiser that moves up, reducing the overall loss suffered by the search engine. It is difficult to accurately estimate this loss as we do not have data from multiple advertisers. Therefore, we leave this as direction for future research.

## 6 Suggested Cures

In the previous discussion, we outlined the problem of aggregation bias and some of its implications. As mentioned earlier, aggregation bias arises due to inadequate data. It might be infeasible for search engines to store and report impression level data due to the size of such datasets and potential privacy concerns. However, a search engine can provide different data to reduce the effect of aggregation. Kendall and Stuart (1977) show that information about a few moments of a distribution can be used to create a good approximation of the distribution. Accordingly, we explore various summary statistics and measure the improvements they offer over the standard aggregate data provided by search engines. We also consider a few modeling approaches, which explicitly account for the variation in position and estimate the improvements offered by them.<sup>14</sup> We first outline these approaches and subsequently draw comparisons between them. These approaches are compared in two ways – (i) using the large representative search engine data available to us, and (ii) an extensive evaluation using simulated data.

### 6.1 Proposed Summarization Techniques

Here, we present a few approaches that might be able to address the problem of aggregation bias.

#### 6.1.1 Sample Mean

##### Different Ways of Aggregation

An important reason for bias is the non-linearity of the position-CTR curve. As a result, linear aggregation of the position does not yield the correct underlying response parameters. Christen

---

<sup>14</sup>We thank the Associate Editor and the anonymous reviewers for recommending this extension.



et al. (1997) and Danaher et al. (2008) show that when the response is multiplicative, i.e. of the form  $\alpha x_1^{\beta_1} x_2^{\beta_2} \dots$ , where  $x_i$  are marketing mix variables, an aggregate model should use geometric means to correctly estimate the coefficients. Unfortunately, there is no analytical analog of this result when the underlying model is logit. We use both the geometric and harmonic means and empirically compare which method of aggregation works better in SSA. As the position-CTR curve is convex in nature, both these aggregation methods might perform better than linear aggregation.

### Modeling Position Variation using a Poisson

We also consider a model where  $V_i$  is drawn from a Poisson distribution with mean equal to the daily (arithmetic) mean  $u$ . The log-likelihood of observing the data in this case is given by

$$LL(\beta|\text{data}, \lambda) \propto \sum_{k=1}^K \sum_{d=1}^D c_{dk} \log \left\{ \sum_{i=0}^{\infty} P(V_{idk} = v) p_{idk} \right\} + (n_{dk} - c_{dk}) \log \left\{ 1 - \sum_{i=0}^{\infty} P(V_{idk} = v) p_{idk} \right\}, \quad (17)$$

where  $\lambda$  is a  $K \times D$  matrix, and every column of  $\lambda$  contains the scale parameter of the Poisson distribution for every day. The position of every impression  $V_{idk}$  for keyword  $k$  on the  $d^{th}$  day is drawn from a Poisson distribution with  $\lambda_{kd} = u_{kd}$  and the probability  $P(V_{idk} = v) = \lambda_{dk}^v e^{-\lambda_{dk}} / v!$ .

### Pooling Positions Across Days

In the previous approach, we try to model the data generating process (DGP) for position. Here, we extend this analysis to model the data generating process by pooling data across multiple days. Let's assume that the  $\mathbb{E}[V_i] = \mu$  and  $Var(V_i) = \sigma_V^2$ .<sup>15</sup> Using law of large numbers, it is easy to show that the mean daily position,  $U \sim N(\mu, \sigma_V^2/n)$ , where  $n$  is the number of impressions on day  $d$ . To estimate  $\beta$ , we follow a two-step process. In the first step, we estimate  $\mu$  and  $\sigma_V^2$  for each keyword by maximizing the following likelihood,

$$LL(\mu, \sigma_V^2|\text{data}) \propto \prod_{d=1}^D \phi \left( \frac{u_d - \mu}{\sigma_V^2/n_d} \right),$$

where  $\phi(\cdot)$  represents the p.d.f. of the standard normal distribution. In the second step, we approximate  $V_i$  by a normal distribution such that  $P(V_{idk} = v) = \Phi \left( \frac{v - \mu_k}{\sigma_{V_k}^2} \right) - \Phi \left( \frac{v - 1 - \mu_k}{\sigma_{V_k}^2} \right)$ , where

<sup>15</sup>We drop the day and keyword subscripts for simplicity.

$\Phi(\cdot)$  represents the c.d.f. of a standard normal distribution. Substituting this p.d.f. in Equation (17) gives us the overall log-likelihood. Note that this approach relies on the implicit assumption that  $F_V(\cdot)$  does not change over time.

### 6.1.2 Higher Order Statistics

If a search engine provides higher order moments in addition to the mean, the aggregation bias may be reduced significantly. We discuss three approaches with increasing data requirements.

#### Mean and Variance

When the mean ( $\mu_{dk}$ ) and variance ( $\sigma_{dk}^2$ ) of position are provided, we assume that the position has a negative binomial distribution (NBD) with probability of success  $p_{dk}$  and (real) number of trails  $r_{dk}$ . The NBD makes intuitive sense as the success of probability  $p_{dk}$  of an NBD, can be thought of as the probability of a competing advertiser placing a higher bid. The log-likelihood function is similar to Equation (17), but the distribution of ad position in this model is given by

$$P(V_{idk} = v) = \binom{v + r_{dk} - 1}{v} (1 - p_{dk})^{r_{dk}} p_{dk}^v,$$

where  $p_{dk} = 1 - \frac{\mu_{dk}}{\sigma_{dk}^2}$  and  $r_{dk} = \frac{\mu_{dk}^2}{\sigma_{dk}^2 - \mu_{dk}}$ .

#### Empirical Distribution

In the preceding approaches, the variation in position is modeled in a parametric manner due to the limitation in the numbers of moments reported. A search engine can provide further moments, e.g. skewness and kurtosis, which might help the research model the randomness in position more accurately. For the sake of brevity, we adopt the extreme case and assume that the search engine provides the empirical distribution of the ad position as shown in Table 5, in addition to the daily summary. In this case, the variation in position can be modeled non-parametrically. The log-likelihood is given by Equation (17), where  $P(V_{idk} = v)$  is provided by the empirical distribution, e.g.  $P(V_{idk} = 5) = 0.29$ . As pointed out earlier, this data can also be independently collected by an advertiser by periodically crawling ads from a search engine.

**Table 5.** Empirical distribution of ad position

position	frequency	$P(V = v)$
2	20	0.13
3	10	0.07
4	50	0.32
5	45	0.29
6	30	0.19

### Position-Level Summary

The complete dataset entirely eliminates aggregation bias, but it might be difficult for search engines to provide this data due to privacy or technical concerns. Instead, the search engine can provide a position-level summary, which is *sufficient statistics* for the logit model as we show below. The position-level summary reports the keyword performance measures for every position (where the ad appeared) separately. It is easy to show that Equation (8) can be simplified to

$$LL(\boldsymbol{\beta}|\text{complete data}) \propto \sum_{k=1}^K \sum_{d=1}^D \sum_{v=1}^{\infty} c_{vdk} \log p_{vdk} + (n_{vdk} - c_{vdk}) \log(1 - p_{vdk}), \quad (18)$$

where  $n_{vdk}$  and  $c_{vdk}$  are the number of impressions and clicks on the  $k^{\text{th}}$  ad at position  $v$  on the  $d^{\text{th}}$  day, respectively. Since Equation (18) is identical to the likelihood function of the position-level data, a position-level summary is sufficient statistic to correctly estimate the parameters of the model.

## 6.2 Application to Search Engine Data

We apply the summarization and modeling techniques presented earlier to the search engine data available to us. Since the data used in Section 4 are quite extensive and representative of typical data in SSA, the results presented here aim to provide real world validation for the suggested summarization techniques. The results demonstrating the performance of these techniques vis-a-vis the search engine data are presented in Table 6. The comparisons between the different summarization techniques are performed using the mean average percentage error (MAPE) in the estimates of  $\beta_0$  and  $\beta_1$ , measured as  $|\hat{\beta}_{0c} - \hat{\beta}_{0s}|/\hat{\beta}_{0c}$  and  $|\hat{\beta}_{1c} - \hat{\beta}_{1s}|/\hat{\beta}_{1c}$ , respectively. Note that the HB-LIV model is used for this analysis.

**Table 6.** Comparative performance of various data summarization approaches on search engine data.

Method	Data Requirement	MAPE( $\varepsilon_0$ )	MAPE( $\varepsilon_1$ )
Different Ways of Aggregation			
Arithmetic Mean	O(X)	53.4%	28.9%
Geometric Mean	O(X)	42.1%	17.5%
Harmonic Mean	O(X)	40.3%	10.2%
Poisson Model	O(X)	43.3%	22.8%
Pooling Data Across Days	O(X)	21.8%	10.8%
Mean and Variance	O(2X)	15.2%	7.8%
Empirical Distribution	O(NX)	9.8%	4.2%
Position-level Summary	o(NX)	0%	0%

Firstly, we observe that both geometric and harmonic means perform better than the arithmetic mean. This reduction in bias is due to the convexity of both these aggregation techniques, which match the convexity of the *position-CTR* curve. We also observe that the harmonic mean performs better than the geometric mean. This implies that if a search engine wants to provide only the mean position in the campaign reports, it should provide the harmonic mean of the position. This result also suggests that researchers who aggregate sponsored search data at a weekly or a monthly level for lack of sufficient data or computational reasons should use the harmonic mean for aggregation. Secondly, we observe that modeling techniques have mixed performance. Modeling the variation in position as a Poisson random variable does not work well as shown in Table 6. This approach does not perform very well because the ad position does not follow a Poisson distribution, a hypothesis we confirm using the Neyman-Scott test. However, pooling data across days leads to significant decrease in the bias, but the effect is heterogeneous across keywords. This approach works best when the distribution of the ad position does not change across days, e.g. when the bids are held constant in the observation period. If the bids change, then the performance of this technique drops considerably. Thirdly, we observe that richer data that provide higher order moments lead to significant improvement in the parameter estimation. From Table 6, we observe that using

both the mean and the variance to model the variation in the ad position significantly improves the estimates. Not only do we observe a reduction in the aggregation bias but there is also a considerable decrease in the error in  $\hat{\beta}_{0s}$ . On average, there is more than 70% reduction in the estimation error associated with  $\beta_0$  and  $\beta_1$ . Using the empirical distribution marginally improves the estimation performance. This summarization technique performs better than all the preceding techniques and the MAPE is considerably lower for  $\beta_0$  and  $\beta_1$ . When the position-level summary is used, the aggregation bias is completely eliminated. However, this technique requires substantially more data as compared to the other techniques.

### 6.3 Application to Simulated Data

In order to perform better characterization of these techniques and ascertain their performance under different conditions, we turn to simulated data.<sup>16</sup>

#### 6.3.1 A Model of Sponsored Search Auctions

We model a generalized second price auction to determine the position of the ad across different impressions. We adopt the approach presented by Abhishek and Hosanagar (2013) and assume that the competing advertisers draw their bids from a Weibull distribution which is given by

$$F(x; \psi, c) = 1 - \exp \left\{ - \left( \frac{x}{c} \right)^\psi \right\},$$

where  $\psi$  is the shape parameter and  $c$  is the scale parameter. A Weibull distribution is quite flexible and has been shown to capture the bid distribution better than other commonly used distributions (Abhishek and Hosanagar, 2013). Variation in  $\psi$  can give rise to different types of competing bid distributions, whereas  $c$  changes the magnitude of the bids. Low values of  $\psi$  lead to heterogeneous bids, whereas higher values of  $\psi$  lead to relatively homogeneous bids. These competing bids are stochastically changed during the course of the day, leading to intra-day variation in position.<sup>17</sup> The probability that the bids are changed before an impression is denoted by  $\kappa$ . Intuitively, the intra-day variation in position increases in  $\kappa$ . In Section 3, we had assumed that the distribution of the ad position remains constant during the period of observation for analytical tractability.

---

<sup>16</sup>We thank the Associate Editor and the anonymous reviewers for suggesting this extension.

<sup>17</sup>The findings are qualitatively similar for different values of  $\psi$ .

Here, we relax this assumption such that the advertiser can change his bid several times during the observation period (but not within a day), changing  $F_V(\cdot)$ . To model consumer choice, a logit data generating process is used and the inputs to the simulation are the coefficients of the logit model  $(\beta_0, \beta_1)$  and  $\kappa$ . The consumer’s decision to click on the ad (of the focal advertiser) conditional on its position is simulated for each query using the random utility model specified in Equation (2). The complete data record the position and the binary click decision for every impression. In addition, the daily total number of impressions, clicks and the mean position for the ad are computed and recorded in the aggregate data. 100 different runs are generated for every tuple  $\beta_0, \beta_1$  and  $\kappa$ , where  $\beta_0 \in [-2, -1], \beta_1 \in [-1, -0.25]$  and  $\kappa \in [0, 1]$ .

For the sake of brevity, the exact characterization of the bias is presented in the Online Appendix. There are three main insights from this analysis – (i) the bias (and the intra-day variation) increases with  $\kappa$  as shown earlier in Proposition 2, (ii) when  $F_V(\cdot)$  is not constant during the simulation period, the bias increases super-linearly in  $|\beta_1|$ , and (iii) the error in the estimate of intercept term increases in  $|\beta_1|$  and  $\kappa$ .<sup>18</sup> We now apply the techniques proposed in Section 6.1 on the simulated data to obtain a greater understanding of their effectiveness.

### 6.3.2 Analysis on Simulated Data

Similar to Section 6.2, we use MAPE to compare the various techniques, but instead of measuring the error with respect to  $\hat{\beta}_c$ , we use the simulation parameter  $\beta$ .<sup>19</sup> We begin by presenting the overall performance of the aforementioned techniques and subsequently discuss how their performance varies across different values of  $\beta$ .

The comparisons between these methodologies are presented in Table 7. We continue to observe that the harmonic mean performs considerably well as compared to the arithmetic and geometric means. The magnitude of the errors reported in Table 7 is lower in comparison to the errors reported in Table 6. This is due to the fact that we include a larger range of  $\beta_0$  and  $\beta_1$  in our simulation analysis than what is observed in practice. Pooling data across days, to derive the underlying distribution of position, reduces the bias significantly. In the absence of adequate data, this approach can be used by researchers and academics to reduce aggregation bias. This approach

---

<sup>18</sup>These results are presented in Tables 1 and 2 in the Online Appendix.

<sup>19</sup>There is no statistical difference between  $\beta$  and  $\hat{\beta}_c$ .

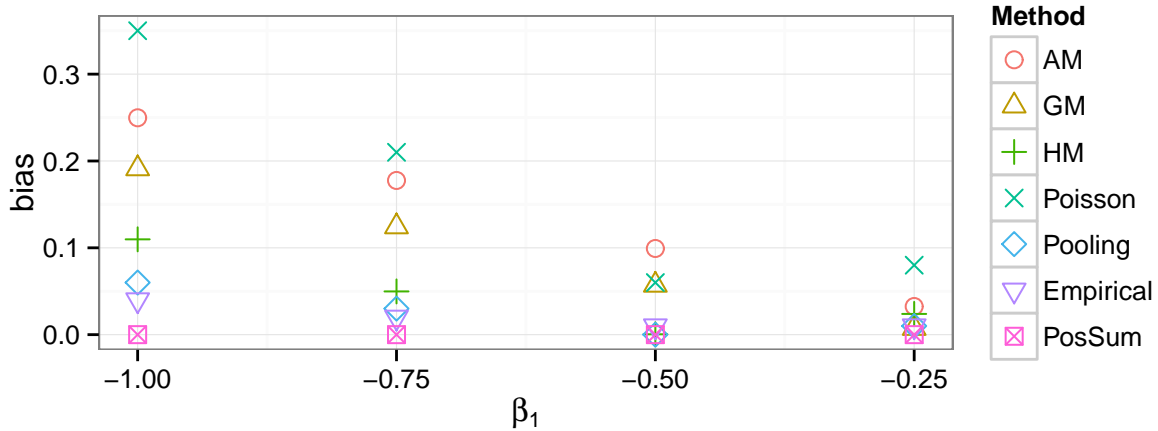
**Table 7.** Comparative performance of various data summarization approaches on simulated data.

Method	MAPE( $\varepsilon_0$ )	MAPE( $\varepsilon_1$ )
Different ways of aggregation		
Arithmetic mean	10.3%	20.1%
Geometric mean	8.1%	10.6%
Harmonic mean	4.4%	5.8%
Poisson model	8.1%	9.7%
Pooling data across days		
Not accounting for bid changes	6.2%	7.7%
Accounting for bid changes	3.4%	4.2%
Mean and Variance	2.8%	3.6%
Empirical Distribution	2.1%	2.5%
Position-level Summary	0%	0%

is more effective when changes in the position distribution (e.g. due to bid changes) are explicitly accounted for in the estimation. If there are significant differences in the position distribution across days, then this approach might further increase the bias.

Daily summaries that report more data, e.g. the daily mean and variance, decrease the estimation error to a significant extent in the simulations. The bias when both the daily mean and variance are reported is quite similar to the bias when data is pooled across days. However, it should be kept in mind that pooling data across days is effective only when changes in  $F_V(\cdot)$  are explicitly accounted for, which can be challenging in a real-world setting. Providing the daily empirical distribution of position further decreases the estimation error, but the incremental benefit is small. Not surprisingly, the position-level summary has no bias.

Figure 3 demonstrates the bias for different techniques. Clearly, we observe that the magnitude of the bias increases non-linearly for all techniques as  $\beta_1$  increases. Furthermore, we observe that using the empirical distribution out performs all the other techniques, while using the harmonic mean leads to the least amount of bias when only the daily means are reported.



**Figure 3.** Magnitude of the bias for different values of  $\beta_1$  when  $\beta_0 = 2.0$  and  $\kappa = 0.5$ .

## 6.4 Summary of Results

The empirical analysis presented here has two distinct themes. The first theme suggests that provisioning of better data by a search engine can lead to significant reduction in the aggregation bias. As we show in the preceding sections, it is incentive compatible for the search engine to provide better data to advertisers. The appropriate dataset should be determined as a trade-off between the loss due to aggregation and the costs associated with providing richer data to advertisers. If a search engine can provide only daily means, it should report the harmonic mean. On the hand, if a search engine is at a liberty to provide any data, it should provide the position-level summary. The second theme points to steps that can be taken by an advertiser or a researcher to explicitly account for the variation in position using modeling techniques (e.g. pooling data across days) or collecting addition information (e.g. by crawling the search engine to generate the empirical distribution). A unilateral reduction in the aggregation bias can prove to be profitable for an advertiser as shown in Corollary 1.

## 7 Conclusions

Search engine advertising is fast emerging as an important and popular medium of advertising for several firms. The medium offers rich data for advertisers on consumer click and conversion behavior. As a result, there has been considerable interest in analyzing SSA data among practitioners and researchers. Several models have been proposed to study consumer behavior and inform advertiser



strategies.

This paper makes three main contributions. First, we demonstrate the existence of aggregation bias and its effect on the equilibrium of the SSA auction. We show that equilibrium bids are lower when advertisers use aggregate data. As a result the search engine’s revenues are always lower due to the bias. Second, we use a large search engine dataset and quantify the magnitude of the bias and measure its economic impact. Third, we present various summarization techniques that can be used by search engines to provide better datasets to advertisers.

These findings have important managerial and economic implications. Advertisers commonly use aggregate data provided by search engines to guide their bidding strategies. Our results suggest that advertisers might not be bidding optimally in these auctions because they overestimate the clicks obtainable at a given position. This not only impacts the advertisers negatively, but also leads to a reduction in the revenue of the advertiser. Given the size of the SSA industry, these losses can translate into several million dollars of lost revenues for the search engines. Our study points out that the current format of the data provided to advertisers is not adequate, and search engines should take steps to address this problem. We recognize that it might be infeasible for search engines to store and report impression level data due to the size of such datasets and potential privacy concerns. However, these constraints do not imply that it is infeasible to provide adequate data to advertisers. We provide guidelines to search engines about the nature of datasets that can be provided to researchers and quantify the reduction in the bias that each of these techniques offer.

We also find that, as a result of aggregation bias, consumer response to other ad attributes, such as ad text or branding, may also have been incorrectly estimated. Thus advertisers must be cautious in applying the biased estimates to guide key managerial decisions such as ad design and keyword selection. In the absence of adequate data from search engines, advertisers and researchers must take into account the variation in ad position within a day. This can be determined by examining if multiple queries are matched to a single keyword (match type is broad) and if bids of competitors change considerably within a given day. If the ad position for a keyword is somewhat stable across impressions within a day, the bias is likely to be low and existing random utility models can be applied on aggregate data.

Our study focuses mainly on demonstrating the existence and direction of aggregation bias in

the coefficient of position and identifying some economic consequences of this bias. An interesting and related question is how aggregation affects ad attributes like wordographics, presence of brand information, ad creative etc. and whether their coefficients also suffer from aggregation bias. In this paper, the effect of ad attributes is subsumed in the intercept term as we do not have data on ad attributes. A richer dataset that contains ad characteristics might help in a more extensive analysis of this issue. Another direction for future research is building models that endogenize the variation in position. The variation in position can be modeled using probabilistic models or structural methods. In ongoing work, we are developing a model that explicitly accounts for the intra-day variation in position and would therefore not suffer from aggregation bias.

SSA presents an exciting opportunity to understand consumer behavior and drivers of firms' advertising strategy. Through this paper we hope to inform the practitioners about the inadequacies of the data standards commonly used in SSA so that they can take steps to address these problems. We also identify issues with some common modeling techniques in SSA so that subsequent research in this emerging area is informed about these issues.

## Appendix: Proofs of Propositions

### Proof of Lemma 1

We use the following result from Muller and Stoyan (2002, page 27) to prove this result. Let  $V_1, \dots, V_n$  be iid random variables and  $f_1, \dots, f_n$  measurable real functions. Define the function  $\bar{f}$  by

$$\bar{f}(v) = \frac{1}{n} \sum_{i=1}^n f_i(v).$$

$$\text{Then, } \sum_{i=1}^n \bar{f}(V_i) \leq_{\text{cx}} \sum_{i=1}^n f_i(V_i).$$

Using this result we now prove that  $W_n \leq_{\text{cx}} V$ , where  $W_n$  is the average position when there are exactly  $n$  impressions on the day. Let  $f_i(v) = v/(n-1)$  for all  $i = 1, \dots, n-1$  and  $f_n(v) \equiv 0$ .

$$\bar{f}(v) = \frac{1}{n} \sum_{i=1}^{n-1} v/(n-1) = \frac{v}{n}.$$

$$\text{Since } \sum_{i=1}^n \bar{f}(V_i) = \frac{1}{n} \sum_{i=1}^n V_i \quad \text{and} \quad \sum_{i=1}^n f_i(V_i) = \frac{1}{n-1} \sum_{i=1}^{n-1} V_i \Rightarrow \frac{1}{n} \sum_{i=1}^n V_i \leq_{\text{cx}} \frac{1}{n-1} \sum_{i=1}^{n-1} V_i.$$

Proceeding in a recursive manner

$$\begin{aligned} \frac{1}{n-1} \sum_{i=1}^{n-1} V_i &\leq_{\text{cx}} \frac{1}{n-2} \sum_{i=1}^{n-2} V_i, \\ &\dots \\ \frac{V_1 + V_2}{2} &\leq_{\text{cx}} V. \end{aligned}$$

$$\Rightarrow W_n = \frac{1}{n} \sum_{i=1}^n V_i \leq_{\text{cx}} \frac{1}{n-1} \sum_{i=1}^{n-1} V_i \leq_{\text{cx}} \dots \leq_{\text{cx}} \frac{V_1 + V_2}{2} \leq_{\text{cx}} V \quad \blacksquare$$

Let  $g$  be any convex function

$$\begin{aligned} \text{As } \mathbb{E}[g(W)] &= \sum_{i=1}^{\infty} g(W_n) P(n), \Rightarrow \mathbb{E}[\mathbb{E}[g(W)]] = \mathbb{E} \left[ \sum_{i=1}^{\infty} g(W_n) P(n) \right], \\ \Rightarrow \mathbb{E}[\mathbb{E}[g(W)]] &= \sum_{i=1}^{\infty} \mathbb{E}[g(W_n)] P(n) \leq \sum_{i=1}^{\infty} \mathbb{E}[g(V)] P(n) = \mathbb{E}[g(V)] \sum_{i=1}^{\infty} P(n) = \mathbb{E}[g(V)], \end{aligned}$$

as  $P(n)$  is a probability measure. Therefore  $W \leq_{\text{cx}} V$ .  $\blacksquare$

## Proof of Proposition 1

Let  $Y_i$  denote an indicator variable that equals 1 if the  $i^{\text{th}}$  impression resulted in a click and zero otherwise. We assume that the clicks are independent of each other and hence  $Y_i$ s are independent. The log likelihood of observing the dataset with a total of  $I$  ad impressions is given by

$$LL(\boldsymbol{\beta} | \text{complete data}) = \sum_{i=1}^I y_i \log p_i + (1 - y_i) \log(1 - p_i). \quad (19)$$

The first order condition (F.O.C) for Equation (19) is as follows

$$\begin{aligned} \frac{\partial LL}{\partial \boldsymbol{\beta}} &= \sum_{i=1}^I \{y_i(1 - p_i) - (1 - y_i)p_i\} \mathbf{x}'_i = 0, \\ &= \sum_{i=1}^I \{y_i - p_i\} \mathbf{x}'_i = 0. \end{aligned}$$

where  $\mathbf{x}_i = (1 \ v_i)'$ . Since we know that  $LL(\beta|data)$  is a convex function in  $\beta$  (Hayashi, 2000) this *F.O.C* gives us the following two equations

$$C = \sum_{i=1}^I p_i, \quad (20)$$

$$\sum_{i=1}^I y_i v_i = \sum_{i=1}^I v_i p_i. \quad (21)$$

Dividing Equation (20) by  $I$  we get

$$obsctr = \frac{C}{I} = \frac{1}{I} \sum_{i=1}^I p_i. \quad (22)$$

If the number of impressions on day  $d$  is  $n_d$  and the number of clicks is  $c_d$ . The log-likelihood of observing the aggregate data for  $D$  days is given by

$$LL(\beta|\text{aggregate data}) = \sum_{d=1}^D c_d \log p_d + (n_d - c_d) \log(1 - p_d) \quad (23)$$

Evaluating the first order condition for Equation (23)

$$\begin{aligned} \frac{\partial LL}{\partial \beta} &= \sum_{d=1}^D \{c_d(1 - p_d) - (n_d - c_d)p_d\} \bar{\mathbf{x}}_d' = 0, \\ &= \sum_{d=1}^D \{c_d - n_d p_d\} \bar{\mathbf{x}}_d' = 0, \end{aligned}$$

where  $\bar{\mathbf{x}}_d = (1 \ w_d)'$ , which in turn gives us

$$\sum_{d=1}^D c_d = C = \sum_{d=1}^D n_d p_d, \quad (24)$$

$$\sum_{d=1}^D w_d c_d = \sum_{d=1}^D n_d w_d p_d. \quad (25)$$

Dividing Equation (24) by  $I$  we get

$$obsctr = \frac{C}{I} = \sum_{d=1}^D \frac{n_d p_d}{I}. \quad (26)$$

Note that the *obsctr* is the same in both cases.

Assuming  $I$  is large we can apply Chebychev's law of large numbers to rewrite Equation (22) as

$$obsctr = \mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V}}{1 + e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V}} \right]. \quad (27)$$

If we have a large enough observation period, Equation (26) can be simplified as

$$obsctr = \sum_{d=1}^D \frac{n_d e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}w_n}}{I(1 + e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}w_n})} = \mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W}}{1 + e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W}} \right].$$

As the observed *ctr*, *obsctr* is same in both the cases,

$$\mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V}}{1 + e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V}} \right] = \mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W}}{1 + e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W}} \right]. \quad (28)$$

As the convex ordering in Lemma 1 holds and logit is a convex in position for  $\beta_0 < 0$  (which is a reasonable assumption in SSA as the *CTR* on the topmost position is less than 0.2 in all cases), it follows from the definition of convex ordering that

$$\mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V}}{1 + e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V}} \right] \geq \mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W}}{1 + e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W}} \right], \quad (29)$$

if  $\hat{\beta}_c = \hat{\beta}_s$ . The equality holds only when  $F_V(\cdot) = F_W(\cdot)$  i.e. which hold only under few special cases (e.g. when there is exactly one impression every day or there is no intra-day variation in position). Since Equation (28) and Equation (29) cannot simultaneously be true and Equation (28) always holds we prove by contradiction that  $\hat{\beta}_c \neq \hat{\beta}_s$ . ■

## Proof of Propositions 2 (i) and (ii)

i) Since we know that

$$\mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V}}{1 + e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c}V}} \right] = \mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W}}{1 + e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s}W}} \right],$$

and by definition of convex order

$$\mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c} V}}{1 + e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c} V}} \right] \geq \mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c} W}}{1 + e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c} W}} \right],$$

we can say that

$$\mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s} W}}{1 + e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s} W}} \right] \geq \mathbb{E} \left[ \frac{e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c} W}}{1 + e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c} W}} \right].$$

As this result hold for any distribution of  $W$ , this relation should hold pointwise for the two function.

$$\Rightarrow \frac{e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s} x}}{1 + e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s} x}} \geq \frac{e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c} x}}{1 + e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c} x}} \quad \blacksquare$$

ii) The preceding relationship implies that

$$\begin{aligned} e^{\hat{\beta}_{0,s} + \hat{\beta}_{1,s} x} &\geq e^{\hat{\beta}_{0,c} + \hat{\beta}_{1,c} x} \quad \forall x \geq 0 \\ \Rightarrow \hat{\beta}_{1,s} &> \hat{\beta}_{1,c} \quad \blacksquare \end{aligned}$$

### Proof of Propositions 2 (iii)

If the true model is  $U = \beta_{0k} + \beta_{1k} V_i + \epsilon$ , but we use the aggregate data for estimation, then we have

$$\begin{aligned} U &= \beta_0 + \beta_1 (W_d + Z_i) + \epsilon, \\ &= \beta_0 + \beta_1 W_d + \beta_1 Z_i + \epsilon. \end{aligned}$$

The disturbance term in this case is  $\beta_1 Z_i + \epsilon$ . Assuming that  $\beta_1 Z_i + \epsilon$  is approximately logistically distributed,  $\hat{\beta}_{1,s} \rightarrow \beta_1 / \sqrt{\text{Var}(\beta_1 Z_i) + 1}$ , where  $\sigma_\epsilon = 1$  for the purpose of identification. The variance  $\text{Var}(\beta_1 Z_i)$  can be computed in the following manner.

$$\begin{aligned} \text{Var}(Z_i) &= \text{Var}(V_i - W_d), \\ &= \text{Var}(V_i) + \text{Var}(W_d) - 2\text{Cov}(V_i, -W_d), \end{aligned}$$

$$\begin{aligned}
&= \sigma_V^2 + \{\sigma_V^2 + \mu_V^2\} \mathbb{E} \left[ \frac{1}{N} \right] + 2\sigma_V^2 \mathbb{E} \left[ \frac{1}{N} \right], \\
&= \sigma_V^2 + \{3\sigma_V^2 + \mu_V^2\} \mathbb{E} \left[ \frac{1}{N} \right].
\end{aligned}$$

Hence,

$$\hat{\beta}_{1,s} \rightarrow \frac{\beta_1}{\sqrt{\beta_1^2 \{\sigma_V^2 + (3\sigma_V^2 + \mu^2)\phi\} + 1}},$$

where  $\phi = \mathbb{E}[1/N]$ .

### Proof of Proposition 3

(i) We assume that advertiser  $j$  moves to position  $j'$  if he uses aggregate data. Since  $j'$  ends up higher than  $j$  in equilibrium,  $b'_j(s_j, j', h) < b_{j'}(s_{j'}, j', h)$  or

$$\frac{1}{\alpha'_{j'-1}} ((\alpha'_{j'-1} - \alpha_{j'})s_{j'} + \alpha_{j'}b_{j'+1}) > \frac{1}{\alpha'_{j'-1}} ((\alpha'_{j'-1} - \alpha'_{j'})s_j + \alpha'_{j'}b_{j'+1}), \quad (30)$$

which also implies that his bid in the equilibrium decreases, i.e.  $b'_j < b_{j'}$ . In addition, as  $b'_j$  is lower than  $b_j$ ,  $j' \geq j$ . The bids for all the advertisers in this case are as follows:

$$\begin{aligned}
b'_i &= \frac{1}{\alpha_{i-1}} \left( \sum_{k=i}^K (\alpha_k - \alpha_{k+1})s_k \right) \quad \text{for } i > j', \\
b'_j &= \frac{1}{\alpha'_{j'-1}} \left( (\alpha'_{j'-1} - \alpha'_{j'})s_j + \frac{\alpha'_{j'}}{\alpha'_{j'}} \sum_{k=i}^{j'+1} (\alpha_k - \alpha_{k+1})s_k \right), \\
b'_i &= \frac{1}{\alpha'_{i-2}} \left( \sum_{k=i}^{j'} (\alpha_{i-2} - \alpha_{i-1})s_i + \alpha_{j'+1}b'_j \right) \quad \text{for } j' \geq i > j, \\
b'_i &= \frac{1}{\alpha'_{i-1}} \left( \sum_{k=i}^{j-1} (\alpha_{i-1} - \alpha_i)s_i + \alpha_{j-1}b'_{j+1} \right) \quad \text{for } i < j.
\end{aligned}$$

As  $h$  do not change for advertisers below  $j'$ , therefore their bids remain the same. It is easy to see that advertisers  $j+1$  to  $j'$  end up bidding higher, i.e.  $b'_i > b_i$  for  $j < i \leq j'$  though they moves up by 1 position.

We now show that the bid associated with every position  $\geq j'$  is lower than when complete

data is used. Lets consider the bid  $b'_{j'}$ , placed by advertiser  $j'$  who occupies position  $j' - 1$ . We start off by showing that  $b'_{j'} < b_{j'-1}$ .

$$b'_{j'} - b_{j'-1} = \frac{1}{\alpha_{j'-2}} ((\alpha_{j'-2} - \alpha_{j'-1}) \underbrace{(s_{j'} - s_{j'-1})}_{<0 \text{ by construction}} + \alpha_{j'-1} \underbrace{(b'_j - b_{j'})}_{<0 \text{ by assumption}}) < 0.$$

So the bid for position  $j' - 1$  is lower than the bid in the complete case. Proceeding in a similar manner it is easy to show that bids for all positions above  $j' - 1$  will also be lower. This implies that  $b'_i < b_i$  for  $i < j$  (these ads do not change position). To summarize  $b'_j < b_j$ ,  $b'_i < b_i$  for  $i < j$ ,  $b'_i > b_i$  for  $j < i \leq j'$  and  $b'_i = b_i$  for  $i > j'$ .

(ii) As all the bids are either the same or lower in this case, search-engine revenue is lower ( $\Pi_S^{A2} < \Pi_S^C$ ). The payoff of advertiser  $j$  is lower as any deviation from the optimal bidding policy results in a strictly lower payoff. Advertisers  $j' + 1$  onwards receive the same payoff and all other advertisers are better off due to suboptimal bid by advertiser  $j$ . Starting off with advertiser  $j'$ ,

$$\begin{aligned} \Pi_{j'}^{A2} - \Pi_{j'}^C &= (\alpha_{j'-1} - \alpha_{j'})s_{j'} - \alpha_{j'-1}b'_j + \alpha_{j'}b'_{j'+1} \\ &= (\alpha_{j'-1} - \alpha_{j'})s_{j'} + \alpha_{j'}b'_{j'+1} - \frac{\alpha_{j'-1}}{\alpha_{j'-1}}((\alpha'_{j'-1} - \alpha'_{j'})s_j + \alpha'_{j'}b_{j'+1}) \\ &> 0 \quad (\text{by Equation 30}) \end{aligned}$$

Similarly, using induction we can show that  $\Pi_i^{A2} > \Pi_i^C$  for advertiser  $i$ , s.t.  $j < i \leq j'$ . For  $i < j$ , the revenues remain the same but the payment to the search engines is lower, hence their payoff are higher for these advertisers too ( $\Pi_i^{A2} > \Pi_i^C, i \leq j$ ).

## Proof of Proposition 4

(i) If all advertisers use the same incorrect estimate of  $\alpha_i$ , the optimal bidding policy is the one proposed by Edelman et al. (2007). They just use  $\alpha'_i$  instead of  $\alpha_i$  to compute the optimal bids. We can show by induction that  $b'_j \leq b_j$  for all advertisers.

Step 0: Let  $b'_{K+1} = b_{K+1} = 0$ .

Step 1:  $b'_K = s_K \left(1 - \frac{\alpha'_K}{\alpha_{K-1}}\right) < s_K \left(1 - \frac{\alpha_K}{\alpha_{K-1}}\right) = b_K$ .



Assuming  $b'_{j+1} < b_{j+1}$ ,

Step  $j$ :

$$\begin{aligned}
b'_j &= s_j \left( 1 - \frac{\alpha'_j}{\alpha'_{j-1}} \right) + \frac{\alpha'_j}{\alpha'_{j-1}} b'_{j+1} \\
&< s_j \left( 1 - \frac{\alpha'_j}{\alpha'_{j-1}} \right) + \frac{\alpha'_j}{\alpha'_{j-1}} b_{j+1} \\
&< s_j \left( 1 - \frac{\alpha_j}{\alpha_{j-1}} \right) + \frac{\alpha_j}{\alpha_{j-1}} b_{j+1} \\
&< b_j.
\end{aligned}$$

Hence,  $b'_j < b_j \forall j \leq K$ .

(ii) Since all advertisers occupy the same position as they did earlier and pay less, search engine profits are lower ( $\Pi_S^{A1} < \Pi_S^C$ ). The advertisers' payoff in the case are higher:  $\Pi_i^{A1} = \alpha_i(s_i - b'_{i+1}) > \alpha_i(s_i - b_{i+1}) = \Pi_i^C$  from Proposition 5 (i).

## Proof of Proposition 5

Lets assume that Advertiser  $j$  uses aggregate data and appears at a position  $j'$ . Let the equilibrium bids be denoted by  $b'_1, \dots, b'_K, 0$ . In the equilibrium,  $\alpha_{j'}(s_j - b_{j'}) > \alpha_i(s_j - b_{i+1}) \forall i \neq j'$ . Now suppose that Advertiser  $j$  does not have access to aggregate data and overestimates  $\alpha_i/\alpha_{i-1}$ . As a result, he bids lower and moves to position  $j'' \geq j'$ . Following the argument in the Proof of Proposition 4, the bids for all positions  $i, i \leq j''$  decrease and all *other* advertisers are better off. As the bids are (weakly) lower, the search-engine revenue is lower. The payoff to Advertiser  $j$  is  $\alpha_{j''}(s_j - b_{j''}) < \alpha_{j'}(s_j - b_{j'})$  as he found it optimal to bid for position  $j'$  when he could correctly estimate the CTR. This implies that he is worse off using aggregate data.

## References

- V. Abhishek and K. Hosanagar. Optimal bidding in multi-item multislot sponsored search auctions. *Operations Research*, 61(4):855–873, 2013.
- A. Agarwal, K. Hosanagar, and M. D. Smith. Location, location, location: An analysis of profitabil-

- ity and position in online advertising markets. *Journal of Marketing Research*, 48(6):1057–1073, 2011.
- K. Ali and M. Scarr. Robust methodologies for modeling web click distributions. In *Second Workshop on Sponsored Search Auctions*, 2007.
- G. M. Allenby and P. E. Rossi. There is no aggregation bias: Why macro logit models work. *Journal of Business and Economic Statistics*, pages 1–14, 1991.
- A. Animesh, V. Ramachandran, and S. Viswanathan. Research Note—Quality Uncertainty and the Performance of Online Sponsored Search Markets: An Empirical Investigation. *Information Systems Research*, 2009.
- M. Christen, S. Gupta, J. C. Porter, R. Staelin, and D. R. Wittink. Using market-level data to understand promotion effects in a nonlinear model. *Journal of Marketing Research*, 34(3):pp. 322–334, 1997.
- N. Craswell, O. Zoeter, M. Taylor, and B. Ramsey. An experimental comparison of click position-bias models. In *ACM International Conference on Web Search and Data Mining*, 2008.
- P. J. Danaher, A. Bonfrer, and S. Dhar. The effect of competitive advertising interference on sales for packaged goods. *Journal of Marketing Research*, 45(2):211–225, 2008.
- P. Ebbes, M. Wedel, T. Steerneman, and U. Boeckenholt. Solving and testing for regressor-error (in)dependence when no instrumental variables are available: With new evidence for the effect of education on income. *Quantitative Marketing and Economics*, 3:365–392, 2005.
- B. Edelman, M. Ostrovsky, and M. Schwarz. Internet advertising and the generalized second-price auction: Selling billions of dollars worth of keywords. *American Economic Review*, 97(1):242–259, 2007.
- A. Ghose and S. Yang. An empirical analysis of search engine advertising: Sponsored search in electronic markets. *Management Science*, 55(10):1605–1622, 2009.
- A. Goldfarb and C. Tucker. Search Engine Advertising: Pricing Ads to Context. *SSRN eLibrary*, 2007.

- I. J. Good and Y. Mittal. The amalgamation and geometry of two-by-two contingency tables. *The Annals of Statistics*, 15(2):694–711, 1987.
- S. Gupta, P. Chintagunta, A. Kaul, and D. R. Wittink. Do household scanner data provide representative inferences from brand choices: A comparison with store data. *Journal of Marketing Research*, 33(4):383–398, 1996.
- J. Hao, S. Menon, S. Raghunathan, and S. Sarkar. A comparison of rbb and rbr ranking mechanisms in sponsored search. In *Proc Conference on Information Systems and Technology*, 2009.
- F. Hayashi. *Econometrics*. Princeton University Press, 2000.
- P. Jeziorski and I. Segal. What Makes Them Click: Empirical Analysis of Consumer Demand for Search Advertising. *SSRN eLibrary*, 2009.
- H. H. Kelejian. Aggregated heterogeneous dependent data and the logit model: A suggested approach. *Economic Letters*, pages 243–248, 1995.
- M. Kendall and A. Stuart. *The Advanced Theory of Statistics, Volume 1: Distribution Theory*. Macmillan Publishing Co., 1977.
- D. Liu, J. Chen, and A. B. Whinston. Competing keyword auctions. In *Proc Workshop on Information Systems and Economics*, 2007.
- A. Muller and D. Stoyan. *Comparison Methods for Stochastic Models and Risks*. Wiley, 2002.
- S. Narayanan and H. S. Nair. Estimating causal installed-base effects: A bias-correction approach. Working papers, NET Institute, 2012.
- S. A. Neslin and R. W. Shoemaker. An alternative explanation for lower repeat rates after promotion purchases. *Journal of Marketing Research*, 26(2):205–213, 1989.
- G. J. Russell and W. A. Kamakura. Understanding brand competition using micro and macro scanner data. *Journal of Marketing Research*, 31(2):289–303, 1994.
- O. J. Rutz and R. E. Bucklin. From generic to branded: A model of spillover dynamics in paid search advertising. *Journal of Marketing Research*, 48(1), 2011.

- O. J. Rutz and M. Trusov. Zooming in on paid search ads—a consumer-level model calibrated on aggregated data. *Marketing Science*, 30(5):789–800, 2011.
- O. J. Rutz, R. E. Bucklin, and G. P. Sonnier. A latent instrumental variables approach to modeling keyword conversion in paid search advertising. *Journal of Marketing Research*, 49(3), June 2012.
- J.-B. E. M. Steenkamp, V. R. Nijs, D. M. Hanssens, and M. G. Dekimpe. Competitive reactions to advertising and promotion attacks. *Marketing Science*, 24(1):35–54, 2005.
- T. A. Weber and Z. E. Zheng. A model of search intermediaries and paid referrals. *Information Systems Research*, 18(4):414–436, 2007.
- J. M. Wooldridge. *Econometric Analysis of Cross Section and Panel Data*. The MIT Press, October 2001.
- S. Yang and A. Ghose. Analyzing the relationship between organic and sponsored search advertising: Positive, negative, or zero interdependence? *Marketing Science*, page mksc.1100.0552, 2010.
- S. Yao and C. F. Mela. A dynamic model of sponsored search advertising. *Marketing Science*, 30: 447–468, May 2011.
- A. Yatchew and Z. Griliches. Specification error in probit models. *The Review of Economics and Statistics*, 67(1):134–139, February 1985.