

**DOES SEARCH MATTER? : USING ONLINE CLICK STREAM DATA TO EXAMINE  
THE RELATIONSHIP BETWEEN ONLINE SEARCH AND PURCHASE BEHAVIOR**

Neveen F. Awad, Joni L. Jones, Jian Zhang

**Introduction**

The ability to track consumer behavior online is considered a prime advantage that online retailing has over brick and mortar retailers. However, firms are still grappling with how to understand their consumers through click stream data. Initially, firms were concerned with visit volume, conjecturing that by increasing the number of visitors, firm revenue would also increase. In fact, visit measures continue to be among the most widely used and cited online metric (Demers and Lev 2001). However, visit measures clearly do not tell the full story (Moe and Fader 2004). Thus practitioners and researchers alike are presented with an important issue, namely firms are collecting myriads of click stream data online, but they do not know what the data is telling them about the relationship between consumer search and purchase behavior. Thus, in this study, we explore the relationship between search behavior and purchase behavior at two different levels of analysis; specifically, we examine search breadth and search depth in association with purchase behavior in the extended and session-level contexts. Prior literature has examined inter-firm search analytically (Bakos 1997), in relation to posted price (Brynjolfsson and Smith 2000), in an experimental setting (Lynch and Ariely 2000), and in relation to individual tendency to search over time (Johnson, Moe et al. 2004). In this study, we extend this prior research by examining the relationship of online search behavior to online purchase behavior across sites. Our goal is assess whether a relationship between search and purchase behavior exists, and how that relationship appears different at different levels of analysis.

**Hypotheses**

Howard and Sheth (Howard and Sheth 1969) began a stream of research that has examined consumer consideration set formation. When choosing to make a purchase, consumers classically begin by first screening the set of all possible online stores to a relevant set called a store consideration set (Fotheringham 1988), then they make purchase/consumption decisions from the retailers in the search set.

Because consumers have limited information-processing abilities and limited information acquisition abilities (Manrai and Andrews 1998) they limit the number of websites that they visit. We refer to this limited set of firms as the consumer's search set. There has been a great deal of research on the details of this process; for this paper, we are interested in how the depth and breath of the search set is associated with the consumer purchase amounts.

Prior literature has assessed breadth of search as a sum of the sources of search (Karson and Fisher 2005). In addition, consumer search has been theoretically modeled as a process in which the consumer seeks out additional information as a function of the expected benefit of that added information (Diamond 1987). As a consumer obtains more information from visiting additional online domains, the expected benefit of the new information decreases. As a result, the probability of increasing the search breadth, by visiting more sites, is likely to decrease as a function of the number of sites that a consumer visits. Aligned with prior literature (Johnson, Moe et al. 2004), we model search breadth as the probability that an individual consumer ( $i$ ) searches an  $x^{th}$  site as a function of visiting the  $(x-1)^{th}$  site:

$$\Pr[X_i = x_i] = \frac{(x_i - 1)\theta_i}{x_i} \Pr[X_i = x_i - 1], \quad x_i = 2, 3, \dots,$$

The parameter  $\theta_i$  is an individual-specific search propensity parameter ( $0 < \theta_i < 1$ ). Because the maximum value of the search parameter  $\theta_i$  is 1, the probability of searching an additional site is a decreasing function of  $x_i$ . Since the probability of search equation is a recursive relationship, we can work backwards through it to obtain a logarithmic distribution (Johnson, Kotz et al. 1993):

$$\Pr[X_i = x_i] = \frac{a_i \theta_i^{x_i}}{x_i}, \quad x_i = 1, 2, \dots, \quad \text{where } a = -[\ln(1 - \theta_i)]^{-1}$$

Based on this theory, we utilize the logarithmic model in our empirical estimate. Our theoretical construct of search breadth, therefore is the logarithm of the number of unique online retailer domains that a consumer visits. Our theoretical construct of search depth is developed in a similar manner, whereby the probability that a consumer visits another page within a specific domain is likely to decrease as a function of the number of pages within that site that a consumer visits. As shown above, our theoretical construct of search *depth*, therefore is the logarithm of the number of unique pages within an online retailer domain that a consumer visits. Prior literature suggests that consumers that purchase online do so from either an implicit favorite firm (Fader and Hardie 2001), or from a small set of a few sites (Johnson, Moe et al. 2004). Based on this prior search literature, along with our theory of diminishing returns of the number of domains, we expect that in a single session, an increased search breadth will be associated with a *decreased* amount of purchases. Hypothesis 1 follows:

*Hypothesis 1:* In the session-level search analysis, search *breadth* will be negatively associated with an increased amount of purchase.

Our theoretical construct of search *depth* was developed in a similar way as the theoretical construct of search *breadth*. We expect that increased search depth is actually a signal of increased likelihood to purchase. It is likely that consumers with a higher session-level search depth have gone into areas of the site including the shopping cart, or the order fulfillment page, and thus are more likely to complete a transaction. As such, we expect search depth to be positively associated with purchase amount at the session level. Hypothesis 2 follows:

*Hypothesis 2:* In the session-level search analysis, search *depth* will be positively associated with an increased amount of purchase.

We examine the theoretical constructs of search breadth and search depth on two different levels, on a session-level basis, and over an extended-two month search period. Prior literature has generally observed a greater level of search activity from more frequent shoppers (Johnson, Moe et al. 2004) over time, which suggests that over an extended period of time, some households are heavier searchers, as well as heavier shoppers, than other households. Consequently, we expect to see a positive relationship between search *breadth* and purchase amount over an extended search analysis period of two months. Hypothesis 3 follows:

*Hypothesis 3:* In the extended search analysis, search *breadth* will be positively associated with an increased amount of purchase.

In addition, we expect to see a positive relationship between search *depth* and purchase amount over the extended search period, as we believe that both metrics are measures of search activity, and therefore shopping activity, over an extended period of time. Hypothesis 4 follows:

*Hypothesis 4:* In the extended search analysis, search *depth* will be positively associated with an increased amount of purchase.

## Methodology

We study the relationship between search and purchase activity directly, through the examination of a large panel of Internet users over time. The data used in this study is Internet click stream data collected by comScore/Media Metric, Inc.. We expect to find evidence of extended search in the online setting, and thus we define two levels of analysis in order to examine whether results differ when you explore extended search versus session-specific search. We examine extended-search over a two-month period of analysis at the household level; whereas we examine session-level search at the individual session level. Below we explain the methodology used at each level.

### Extended-Search Analysis

For the extended-search analysis, purchase behavior is examined in relation to search behavior as well as household demographic information. In the full paper, we plan to examine four different dependent variables in the extended-search analysis; these variables include: 1) Number of Purchases, calculated as the total number purchases per household; 2) Purchase spending, calculated as the total amount of dollars spent per household; 3) Session density, calculated as the number of household sessions that resulted in a purchase divided by the total number of household sessions; and 4) Domain density, calculated as the number of domains from which the household made a purchase divided by the total number of domains. However, due to space limitations, the results presented in the extended-search results section will be for just the first one dependent variable: purchase amount.

The independent variables in the extended-search analysis include search breadth, search depth, product price, and demographics information including: 1) household eldest age; 2) household size; 3) household highest education; and 4) Census region. Equation (1) is as follows:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \beta_6 x_{i6} + \varepsilon_i \quad (1)$$

For the machine  $i$ ,  $y_i$  is the purchase behavior,  $\beta_0$  is the intercept,  $x_{i1}$  is the number of search sessions of the person,  $x_{i2}$  is average product price per session per person,  $x_{i3}$ ,  $x_{i4}$ ,  $x_{i5}$ ,  $x_{i6}$  are household eldest age, household size, household highest education, and location region respectively. And  $\varepsilon_i$  is the error term. The parameter  $\beta_1$  is assumed to be affected by two dimensions of search: search breadth (average domains per session for the person), and search depth of the sessions (average views per domain per session for the person). Therefore,

$$\beta_1 = \alpha_0 + \alpha_1 z_{i1} + \alpha_2 z_{i2} \quad (2)$$

where  $\alpha_0$  is the intercept, and  $z_1$  and  $z_2$  are search breadth and search depth, respectively.

Consequently, combining equation (1) and (2), we obtain the equation to be estimated:

$$y_i = \beta_0 + \alpha_0 x_{i1} + \alpha_1 (x_{i1} z_{i1}) + \alpha_2 (x_{i1} z_{i2}) + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \beta_6 x_{i6} + \varepsilon_i \quad (3)$$

Lastly, aligned with our theory, we transform the linear model into a log linear exponential regression model:

$$\ln y_i = \beta_0 + \alpha_0 \ln x_{i1} + \alpha_1 \ln(x_{i1} z_{i1}) + \alpha_2 \ln(x_{i1} z_{i2}) + \beta_2 \ln x_{i2} + \beta_3 \ln x_{i3} + \beta_4 \ln x_{i4} + \beta_5 \ln x_{i5} + \beta_6 \ln x_{i6} + \varepsilon_i \quad (4)$$

### Session Level Analysis

The session-level purchase behavior is modeled as a function of search breadth, search depth, demographic information, and product price. The function can be expressed as follows:

$$y_{it} = \beta_0 + X'_{it} \beta_1 + Y'_i \beta_2 + Z'_{it} \beta_3 + \varepsilon_{it} \quad (5)$$

Where  $y_{it}$  is the purchase behavior of household  $i$  during search session  $t$ . The purchase behavior is measured as two separate dependent variables: 1) the number of purchases per session; and 2) the total spending per session. Due to space limitations, we only report the results for the second dependent variable, total spending per session.  $\beta_0$  is the intercept of the function, representing some features that cannot be observed (It is assumed that these features are person-specific features. However, relaxing of this assumption does not change the results).  $X'_{it}$  is the transpose of the column vector which provides the information of the search session  $t$  of the person  $i$ . It consists of four items: 1) search breadth (domains per session); 2) search depth (views per domain per session); 3) search time; and 4) A dummy variable to separate out heavy searchers (household whose average search activity fell into a range one standard deviation about the mean were classified as heavy searchers).  $Y'_i$  is the transpose of the column vector which contains the demographic information of household  $i$ . In all the search sessions,  $Y'_i$  is therefore composed of the demographic variables: 1) household eldest age; 2) household size; 3) highest education level in the household; 4) census region; and 5) unobservable factors, such as the household's preference set. Finally,  $Z'_{it}$  represents the vector of the product price of the search session. In this study, we use average purchase price of the product for  $Z'_{it}$ .  $\beta_1$ ,  $\beta_2$ ,  $\beta_3$  are the column vectors of coefficients for the variables  $X'_{it}$ ,  $Y'_i$ , and  $Z'_{it}$ , respectively.  $\varepsilon_{it}$  is the error term, and is assumed to be composed of independent identically-distributed random variables (iid)<sup>1</sup>.

For each household  $i$ , we replicate the above search session equation for another search session in the following equation (6). In this function, the household  $i$  is conducting the search session  $m$ :

$$y_{im} = \beta_0 + X'_{im} \beta_1 + Y'_i \beta_2 + Z'_{im} \beta_3 + \varepsilon_{im} \quad (6)$$

Subtracting equation (5) from the equation (6), we obtain the equation (7) for estimation:

$$y_{it} - y_{im} = (X'_{im} - X'_{it}) \beta_1 + (Z'_{it} - Z'_{im}) \beta_3 + (\varepsilon_{it} - \varepsilon_{im}) \quad (7)$$

The equation (3) measures the effects of the change of the search behaviors and the target products on the customer's purchase behavior.  $\varepsilon_{it} - \varepsilon_{im}$  is still iid due to the fact that both  $\varepsilon_{it}$  and  $\varepsilon_{im}$  are assumed to be iid (again, relaxing this result does not change the results).

In accordance with out theory of search breadth and search depth, we transform the linear model into a log linear exponential regression model for estimation, expressed as:

$$\ln(y_{it} - y_{im}) = (\ln(X'_{im} - X'_{it})) \beta_1 + (\ln(Z'_{it} - Z'_{im})) \beta_3 + e \quad (7)$$

A fixed effects technique assumes that differences across categories can be captured in differences in the constant term. One methodological option is to use this technique to control for differences across families in the session level analysis. However, the technique potentially requires a data set with small number of categories and large number of observations for each category. Otherwise, the degree of freedom would be greatly reduced and the estimate would be highly inefficient, or even inconsistent. In particular, in our data, there are totally 135,748 different families (represented by 135,748 different machine\_id's) and 1,870,300 different search sessions. On average, there is only 13.78 search sessions for each family. Empirical analysis using a fixed effects model on this data structure is therefore not ideal.

## Results

### Extended-Search Results

| Table 1. Extended-Search level analysis of Purchase Amount |                         |               |                   |
|--|-------------------------|---------------|-------------------|
| Number of Observations:                                    | 104,305                 |               |                   |
| R-square   | 0.9951                  |               |                   |
| F-value  | 1972691                 |               |                   |
|  | Estimates (StdDev)      | t Value       | Pr >  t           |
| <b>Intercept</b>   | <b>0.1053 (0.008)**</b> | <b>13.56</b>  | <b>&lt;0.0001</b> |
| <b>Sessions per Person</b>                                 | <b>0.1478 (0.001)**</b> | <b>120.78</b> | <b>&lt;0.0001</b> |
| <b>Average Search Breadth</b>                              | <b>0.1950 (0.004)**</b> | <b>49.21</b>  | <b>&lt;0.0001</b> |
| <b>Average Search Depth</b>                                | <b>0.0066 (0.001)**</b> | <b>5.47</b>   | <b>&lt;0.0001</b> |

<sup>1</sup> In fact, we find evidence of heteroscedasticity through the use of both White's test and the Breusch-Pagan test. We therefore also used Generalized Method of Moments (GMM) to control for the heteroscedasticity, and the results remained consistent the OLS results. In addition, even with the presence heteroscedasticity, we can still expect to obtain consistent results with some loss of efficiency.

|                                   |                          |                |                   |
|-----------------------------------|--------------------------|----------------|-------------------|
| <b>Average Product Price</b>      | <b>1.0340 (0.0002)**</b> | <b>4285.18</b> | <b>&lt;0.0001</b> |
| <b>Household Eldest Age</b>       | <b>0.0143 (0.002)**</b>  | <b>7.38</b>    | <b>&lt;0.0001</b> |
| <b>Household Size</b>             | <b>-0.0098 (0.002)**</b> | <b>-4.47</b>   | <b>&lt;0.0001</b> |
| <b>Census Region 2</b>            | <b>-0.0077 (0.003)*</b>  | <b>-2.27</b>   | <b>0.0231</b>     |
| <b>Census Region 3</b>            | <b>-0.0114 (0.003)**</b> | <b>-3.66</b>   | <b>0.0003</b>     |
| Census Region 4                   | -0.0022 (0.004)          | -0.61          | 0.5401            |
| Household Most Education 1        | -0.0032 (0.006)          | -0.53          | 0.5949            |
| <b>Household Most Education 2</b> | <b>0.0124 (0.006)**</b>  | <b>2.11</b>    | <b>0.0350</b>     |
| Household Most Education 3        | 0.00424 (0.006)          | 0.66           | 0.5096            |
| <b>Household Most Education 4</b> | <b>0.0331 (0.006)**</b>  | <b>5.38</b>    | <b>&lt;0.0001</b> |
| <b>Household Most Education 5</b> | <b>0.0379 (0.006)**</b>  | <b>5.78</b>    | <b>&lt;0.0001</b> |

Table 1 lists the estimates of the extended-search model; There are 104, 305 observations in the estimation. The adjusted R-square is 0.9951 and the F-value of 1972691.00 indicates the overall estimates are significant with more than 99% confidence. We see from the results of the extended-search model that both search depth (0.0066) and search breadth (0.1950) are significant and positively correlated with number of online transactions, confirming Hypotheses 3 and 4. In addition, when search breadth and depth are kept fixed, the number of sessions (0.1478) is significant and positively correlated with dollar amount of online purchases; therefore just going online often is a significant signal that a household is a heavy consumer household. The comparative impact of the search breadth (0.1950) is greater than the impact of the search depth (0.0066) by more than a factor of ten, suggesting that at the aggregate, a better signal of greater purchase amount is a greater number of sites searched. The relationship between search breadth, search depth, and purchase amount remains consistent when the dependent variable is number of transactions rather than total dollar amount, lending further credibility to the extended-search level results. However, due to space limitations, we are unable to present the results of the number of transactions dependent variable.

### Session-level Results

| <b>Table 2. Session level analysis of Purchase Amount</b> |                           |                |                    |
|---|---------------------------|----------------|--------------------|
| Number of Observations:                                   | 881,881                   |                |                    |
| R-square  | 0.9987                    |                |                    |
| F-value   | 8.21E+07                  | <u>t Value</u> | <u>Pr &gt;  t </u> |
| <u>Explanatory Variables</u>                              | <u>Estimates (StdDev)</u> |                |                    |
| <b>Search Breadth</b>                                     | <b>-0.0129 (0.0001)*</b>  | <b>29.2</b>    | <b>&lt;0.0001</b>  |
| <b>Search Depth</b>                                       | <b>0.0070 (0.0002)*</b>   | <b>28.92</b>   | <b>&lt;0.0001</b>  |
| <b>(Search Depth)<sup>2</sup></b>                         | <b>0.0009 (0.0000)*</b>   | <b>77.34</b>   | <b>&lt;0.0001</b>  |
| <b>Search Duration Time</b>                               | <b>0.0158 (0.0002)*</b>   | <b>-11.07</b>  | <b>&lt;0.0001</b>  |
| <b>(Search Duration Time)<sup>2</sup></b>                 | <b>-0.0002 (0.0000)*</b>  | <b>67.26</b>   | <b>&lt;0.0001</b>  |
| <b>SBreadth × Duration Time</b>                           | <b>0.0017 (0.0000)*</b>   | <b>8.25</b>    | <b>&lt;0.0001</b>  |
| <b>Heavy search</b>                                       | <b>0.0135 (0.0016)*</b>   | <b>11216</b>   | <b>&lt;0.0001</b>  |
| <b>Average Price</b>                                      | <b>1.0131 (0.0000)*</b>   | <b>23720.9</b> | <b>&lt;0.0001</b>  |

Table 2 lists the estimates of the session-level model. There are 881,881 observations in the estimation. The R-square is 0.9987 and the F value of 8.21E+07 indicates the overall estimates are significant with more than 99% confidence. The model is set up as a logarithmic model, which provides estimates of the people's search behaviors' elasticity for their online purchase spending, and is aligned with out theory on search breadth and search depth. We control for average price of purchased goods in the session-level model to assess the spending amount in relation to the price of the good being searched for and purchase; clearly the more expensive the product, the greater the amount of purchase spending, as the results significantly confirm.

The results based on the session level analysis show that search depth is significant and positively associated with purchase amount, confirming hypotheses 2. Specifically, a 1% increase in search depth increases is associated with a 0.0079% increase in spending. Search breadth, however, is significant and negatively associated with session-level purchase amount, suggesting that the greater number of sites that a consumer searches across in given session, the lower the purchase amount that he will spend during that session, confirming hypotheses 1. Specifically we find that, for the individual session-level analysis, a 1% increase in search breadth is associated with a 0.0129% reduction in the amount of money spent during a session, keeping all other factors fixed. Another interesting result is the quadratic relationship between session duration time and amount spent. Our results show a significant positive relationship between time and purchase amount, but a significant negative relationship between duration time-squared and purchase amount, suggesting that consumers spend more money up to a certain period of time in the session, after which an increase in the time is actually associated with a decrease in the amount spent. Specifically, 1% increase in duration time

is associated with a 0.0158% in purchase amount. However, a 1% increase in duration time-squares is associated with a 0.0002% decrease in purchase amount. The inclusion of the square terms of both search depth and search depth is aligned with prior literature (Moe, 2003), which suggests that search is valuable up to a certain threshold. Our results confirm such theory in the case of search duration time, where after too much search duration time, the value is negative. The results also show that heavy searchers are significantly associated with greater purchase amounts.

## Discussion

One of the most cited advantages of online retailing for firms is that they can capture data regarding all of their consumers' activities. However, to date, minimal research has examined the relationship between online search behavior and online purchase behavior using such click stream data. This study explores the relationship between search behavior and purchase behavior at two different levels of analysis; specifically, we examine search breadth and search depth in association with purchase behavior in the extended-search level and session-level contexts. In the extended search context, search depth was found to be positive and significantly associated with number of purchase transactions, and total amount spent; this result suggests that, in the aggregate, search breadth is a good signal of consumers that shop often and that spend a lot of money.

On the other hand, in the session-level analysis, search breadth was found to be significant and negatively associated with amount spent in a session, suggesting that the more sites that a consumer visits in a session, the less money the consumer spends. This result may suggest that higher search breadth in a session indicates either: 1) price-sensitive consumers (Kim, Srinivasan et al. 1999); or 2) non-directed shoppers (Moe 2003). Future research is needed to explore this distinction in greater detail. Search depth is positive and significant in the session-level analysis, indicating that an increase in the number of pages viewed within a given site is associated with an increase in the total amount of dollars spent in a given session. This result, in accordance with the quadratic relationship between duration time and total amount spent, corresponds with previous literature that suggests that online consumers tend to be more time constrained than the average consumer (Bellman et al. 1999). Thus, the online consumers that searches deeper into a domain in a given session are likely to be focused on purchasing during that session, and thus are correlated with a greater session-level spending amount.

## Conclusion

This paper establishes a direct relationship between search behavior and purchase behavior online. In addition, we illustrates that the nature of the relationship between online search and online purchase behavior is dependent upon the level of analysis. From the perspective of the firm, these results can inform customer segmentation and strategies: consumers that search across more sites over an extended period of time may be more frequent purchasers, and may spend more money. On the other hand, consumers that search across more firms within a given session may spend less money. Lastly, consumers that search more *within* a firm's site within a given session are likely to spend more money. Thus, firms may want to consider up-selling to consumers with a longer search path in a given session, or may want to offer various purchase incentives during the session. This research underscores the value of looking at different levels of analysis, depending upon the information that a firm is seeking to find. We see this research as an initial step towards exploring the relationship between search behavior and purchase behavior online *directly* through click stream and transaction data. Our goal for future iterations of this paper is to further develop theory of the types of searches that consumers perform online, in a way that can be tested using efficient empirical methods. There is a wealth of theory to be developed and tested in the area of online search and purchase behavior; we believe this paper takes in an important first step in that direction.

## References

- Bakos, J. Y. (1997). "Reducing Buyer Search Costs: Implication for Electronic Marketplaces." Management Science **43**(12).
- Brynjolfsson, E. and M. Smith (2000). "Frictionless Commerce? A Comparison of Internet and Conventional Retailers." Management Science **46**(4).
- Demers, E. and B. Lev (2001). "A Rude Awakening: Internet Shakeout in 2000." Review of Accounting Studies **6**(2 - 3): 331 - 359.
- Diamond, P. (1987). Search Theory in The New Palgrave. A Dictionary of Economics. J. E., M. M. & P. K. N. London, Stockton Press.
- Fader, P. S. and B. G. S. Hardie (2001). "Forecasting Repeat Sales at CDNOW: A Case Study." Interfaces **31**(3): 94-107.
- Fotheringham, A. (1988). "Consumer Store Choice and Choice set Definition." Marketing Science **7**(3): 299-311.
- Howard, J. A. and J. N. Sheth (1969). The Theory of Buyer Behavior. New York, John Wiley.
- Johnson, E. J., W. Moe, et al. (2004). "On the Depth and Dynamics of Online Search Behavior." Management Science **50**(3): 299-308.
- Johnson, N. L., S. Kotz, et al. (1993). Univariate Discrete Distributions. New York, John Wiley & Sons, Inc.
- Karson, E. J. and R. J. Fisher (2005). "Predicting intentions to return to the Web site: Extending the dual mediation hypothesis." Journal of Interactive Marketing **19**(3): 2-14.
- Kim, B.-D., K. Srinivasan, et al. (1999). "Identifying price sensitive consumers: the relative merits of demographic vs. purchase pattern information - A Meta-Analysis of Econometric Models of Sales." Journal of Retailing **75**(2): 173-193.
- Lynch, J. G. and D. Ariely (2000). "Wine Online: Search Costs Affect Competition on Price, Quality, and Distribution." Marketing Science **19**(1): 83-103.
- Moe, W. W. (2003). "Buying, Searching, or Browsing: Differentiating between Online Shoppers Using In-Store Navigational Clickstream." Journal of Consumer Psychology **13**(3): 29-39.