### Understanding On-line Archival Use through Web Analytics

Chris Prom, Assistant University Archivist, University of Illinois at Urbana-Champaign
prom@uiuc.edu

Prepared for ICA-SUV Seminar, Dundee, Scotland, August 14, 2007.

For many years, archivists have conducted user surveys, gathered reference statistics, and consulted informally with the users of archival records and manuscript collections. Our information about users' needs and preferences, as well as their information-seeking behaviors, has been generated by direct contact with users via conversations, interviews, surveys, and other research. However, online archival databases, image repositories, and other electronic sources have opened our archives to new audiences and have provided traditional users a way to access materials without an archivist's mediation. It is very difficult to get accurate information about how such users interact with our resources, since many of them never contact the archives. When we have better information about our virtual users, we can design more effective web sites and, hopefully, increase our users' satisfaction with our services. In this paper, I will introduce the concept of web analytics and demonstrate how tested its potential to improve archival services in a pilot project at the University of Illinois Archives.

According to Wikipedia, Web analytics is "the study of the behaviour of website visitors. In a commercial context, web analytics especially refers to the use of data collected from a web site to determine which aspects of the website work towards the business objectives; for example, which landing pages encourage people to make a purchase."[1] Web analytics has also been defined as "the measurement, collection, analysis and reporting of Internet data for the purposes of understanding and optimizing Web usage."[2] There is nothing specifically archival about web analytics. It is a business tool and arose from a commercial need. Online businesses maximize profit when they design websites that make it easy for people to buy goods or services. Web analytics software helps online businesses improve their websites so that they can make more money. Although as archivists, we are typically not profit-oriented (at least

---

1. http://en.wikipedia.org/wiki/Web_analytics, accessed November 19, 2007.

in a commercial sense) we can use this tool to increase our users' satisfaction with our online services.

The tools and concepts that help businesses understand (and exploit) user needs and behaviors to

maximize profit can also be used to help archives understand (and exploit) user needs and behaviors to

maximize archival use.

Web analytics includes two discrete elements: the installation of a software tool to collect,

measure, and report user data; and methods to interpret the reported information so that it can be used to

improve a website's features. The reports that the software generates may facilitate effective analysis and

decision making since the types of information reported are more discrete and detailed than others to

which archivists might have access, such as server logs.  For example, the software may help an archivist

understand how a particular resource is being used.

Website usability improvements have a basis in useful and accurate data, but reporting

mechanisms such as server logs provide rudimentary and irrelevant data.  For example, our internal logs

show that the University of Illinois Archives received 803,838 'hits' in July 2007 alone, but we estimate

that only a small percentage represent human use (although we really can't say precisely how many). [3]

Aside from the fact that internal log data is misleading, it cannot serve as the basis for a decision or action

because it is so trivial.

Web analytics software, on the other hand, provides very accurate and deep information about the

actual use of your online resources.  For example, it reports which pages referred people to your site, how

long users stayed on there, and how many pages they viewed.  The report excludes traffic generated by

non-human agents (such as web crawlers), so only actual use is shown.

---

2. Neil Mason, "The Four Parts of Web Analytics,"  http://www.google.com/analytics/cu/ac_the_four_parts.html, accessed November 9, 2007.
3. Server logs typically tracks 'hits.'  However, a hit is simply a request from another computer to a web server.  In response to the request, the web server sends a response—typically the information needed to construct a particular view a particular page, including embedded images or other content (such as headers or footers).  As a data point, the concept of the 'hit' is inherently problematic for several reasons.  Hits originate from both humans and from non-human agents, such as a search engine's web crawler.  The number of hits does not accurately represent human interaction with the website because server logs rarely segregate actual page views from the activity of non-human agents.  In addition, a

Web analytics software has another advantage: it provides an unobtrusive means to study how users interact with on-line resources while they are seeking information relevant to an actual research need. Most other methods of studying user behavior involve some intervention, user contact, or observation, which inevitably affects the results of the study. Other studies have users run canned searches. I do not mean to suggest that web analytics can be substituted for talking to your users or conducting these studies, but when it is used in conjunction with other, more traditional methods, web analytics forces us to ask new questions about users and their information-seeking behaviors. For example, Google Analytics, the application we implemented at the University of Illinois, allows us to study what keyword searches lead users to our site, to trace a typical user's browsing path through our site, and to measure the "click-thru" rates to certain types of resources.[4] This information may help us understand what types of content archival users want and show us how they interact with it, so we can redesign our site to make their interactions more productive and satisfying.

**Project Preparation**

In July of 2007, we implemented a pilot project to test the utility of Google Analytics (www.google.com/analytics) for measuring and analyzing use of our website. We hoped that the project would provide preliminary information about on-line use and demonstrate whether the concept of web analytics might have broader applicability in the archival community, both in the US and internationally.

Before beginning the project, we analyzed existing data about the use of our archives, developed a set of questions to examine during the pilot project, and developed a privacy policy. While these steps were time consuming, we believe they helped ensure that the information we gathered and analyzed was more accurate and useful.

page view is often recorded as several 'hits' since each embedded image, header, and footer is recorded separately.
4. While most of the web analytics literature targets commercial web sites, the principles of measuring web usage can be applied to any web resource. The specific types of information that can be measured are covered in detail by E. T. Peterson, *Web Site Measurement Hacks: Tips and Tools to Help Optimize Your Online Business* (O'Reilly: 2005).

It seemed reasonable to mine our existing use statistics before beginning our study of on-line use. The University of Illinois Archives has gathered reference statistics since its founding in 1963, and William Maher has described our use trends through 1985.[5] While this paper is not the appropriate forum to update that analysis, it is clear that use of our materials has been increasing and migrating from on-site to remote services. In our fiscal year 2005/06, 782 (28%) of our 2,782 users contacts were initiated via email. An additional 367 (13%) were initiated via phone. Few of these users subsequently visited the Archives. We provided them a variety of remote ranging from photocopies and scanning to complimentary and fee-based research. Using the reference statistics system that produced these figures, it is impossible to gather any information regarding those who did **not** contact us but used our website in some fashion.

Assessment of website use must be conducted against measurable goals. Our website serves several complementary purposes. Specifically, our web presence provides basic information regarding the Archives, promotes our programs and services, encourages users to find relevant descriptive information, and facilitates user contact with archival staff. Although it would be tempting to gather information regarding all of these goals, we felt that based on our use trends it was most important to find out whether our website's design encouraged or discouraged archival use and contacts with archival staff.

The web measurement literature notes that a successful analytics program will collect a limited amount of data and seek to answer a few discrete and measurable questions and at any given time.[6] While it is tempting to collect large amounts of data and prepare numerous reports, web analytics software can be configured to provide specific types of information if you have precise questions in mind. We developed four such questions:

- How do people get to our site?

- What are the most popular pages/groups of content?

- What are most popular searches?

5. William J. Maher, "The Use of User Studies," *Midwestern Archivist* 11:1 (1986) 15-26.

- How do users navigate through our archival descriptions?

Asking these questions is easy but answering them requires the right tools and careful planning. Using the analytics software we developed reports to answer the first three questions. To gather data regarding the last question, we defined certain user behaviors that we deemed desirable, such as viewing a descriptive record or emailing the archives, then configured the software to report how and when users exhibited the wanted behaviors.

Before attempting to install the software and configure it so that it could report some data regarding these questions, we developed a privacy policy regarding the installation of the analytics program. Not only was this required by state law, but it was good ethical and professional practice. It forced us to think clearly about the information we were gathering and how it would be used. It is important to note that reports generated by Google Analytics do not include data regarding individual users or IP addresses. Aggregate data is reported, so users stay anonymous. Nevertheless, archivists and others using web analytics software should prepare a privacy policy prior to implementing an analytics program, so that users are informed of the tracking activities and, ideally, provided a way to opt out.[7] It is the right thing to do.

**Installing Web Analytics Software**

Those interested in web analytics can select from many vendors or services. For archivists, Google Analytics is probably the most feasible option not only because it is it free and easy to install, but because it reports very useful data. To get started, we registered an account at www.google.com/analytics and loaded a piece of tracking code on every page on our website. For most pages, this was a simple process because our site includes a common footer on every page. If your site does not include a common footer, you can use a directory-level "find and replace" to insert the code. You do not need to install an application on your desktop since Google Analytics is a web application.

---

6. Peterson, *Measurement Hacks*, 54-55, 72.

After completing basic installation, we configured the software to ensure that the data reported in each area was as useful and accurate as possible. For example, we defined four "goals." Goals allow you to track progress toward a desired user behavior. In a commercial setting, a goal might be viewing a thank-you page, which is only shown after a user makes a purchase. We defined goals that we believed would help us answer the last of our four research questions: "How do users navigate through our archival descriptions?" (Our other three questions could be answered by the standard reports provided by the software.) For example, we defined a goal to track how users navigated to the page in our 'Archon' database which shows series-level descriptions. Similarly, we defined a goal to track how many users clicked a link to email the archives or download a PDF finding aid. As the analytics literature recommends, we attempted to track user progress through the site in a linear fashion. As users progress through the site, one can imagine them being funneled toward a goal, with some users dropping off at each step in the process.[8]

After we verified that data was being accurately received and reported, we viewed and downloaded reports using the customizable Google Analytics 'dashboard,' shown in figure one. The dashboard provides basic data, such as the number of visits and visitors, the top 'referrers,' (i.e. the websites that users visited immediately prior to viewing your page) and the most viewed pages. All of the reports can be deeply drilled to provide specific information regarding a single page or group of pages. For example, we mined the reports in the 'content' area to determine which Google searches referred the most users to our site and to see which individual pages caused users to leave our site in the greatest relative numbers.

---

7. A copy of the UIUC privacy policy is available at http://www.library.uiuc.edu/ahx/about/privacy.php (accessed November 9, 2007.)

8 Google Analytics concisely defines a funnel as "a series of pages through which a visitor must pass before reaching the conversion goal. The name comes from a graph of visitors who reach each page - the first page counts the most visitors, and each successive page shows less visitors as they drop off before reaching the final goal." The concept of the funnel is useful in helping to analyze how efficiently the website direct visitors toward a website goal. Pages that are not user-friendly will see higher drop-offs.
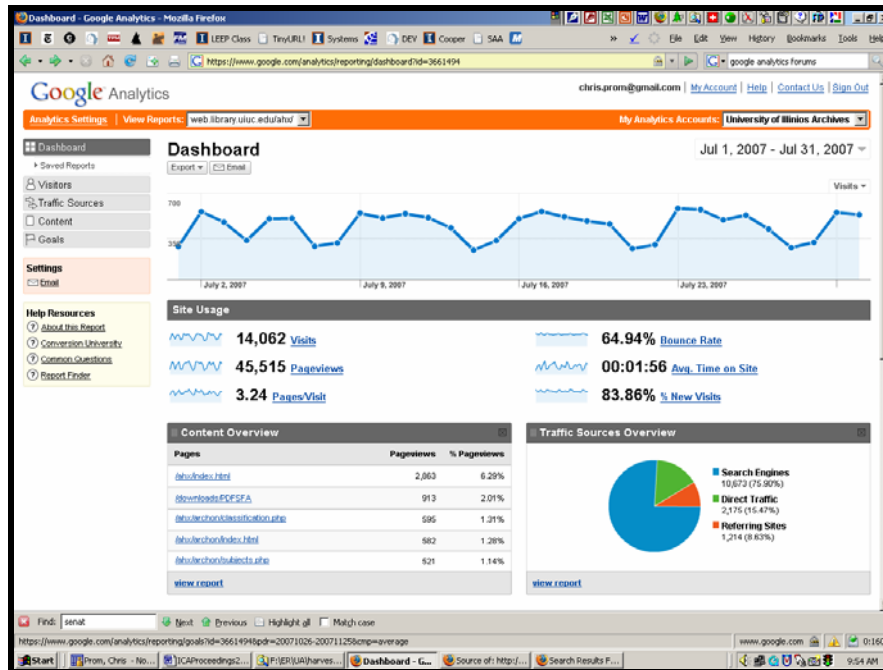
**Figure One: Google Analytics "Dashboard"**

**Interpreting Analytics Reports**

One must, of course, be very careful in making changes to a website, much less your archival program, based on one source of data. After carefully analyzing of the results returned by the Google Analytics report tools, we noticed several trends regarding how users interact with our site and made change to our site to rectify problems.

Since there is not enough time during this short talk to review results for each of the four questions we defined above, let me focus my comments on one in particular: How do users navigate through our site? W attempted to answer this question by tracking user progress toward four goals. Specifically, we tracked how often users completed searches in our holdings database, viewed our series-level descriptions (known informally as 'control cards'), downloaded or viewed a full finding aid, and clicked a link to send us an email message.

My hypothesis was that typical users would enter our site through our home page, complete a search, look at the holdings records, view the full finding aid and—if they felt the records described

online might meet their information need —contact us for more information.  The progression I envisaged

matched the concept of the 'goal conversion funnel' that is used in analytics software.  Our site was

designed in a way that encouraged users to view increasingly detailed descriptive records then contact us

so that an information need could be directly fulfilled by the reference staff.  Figure two shows our home

page design, with the search box placed in the 'sweet spot' to which users' eyes gravitate.  Figure three

shows the result of a search, our 'control-card' pages describing a record series.  We include an 'email us'

link in the navigation bar.



**Figure Two: UIUC Archives Homepage**

**Figure Three: Record Series view with "E-mail Us" link**

The basic goal report is shown in figure four. Few if any of the people using our website conformed to the linear 'funnel' which I had envisaged for the typical user. During July 2007, users viewed our 'control card' descriptive records 8,035 times but emailed the archives only 49 times. The descriptive records were viewed many more times than the database was searched since relatively few users enter our site through the home page and then complete a search. A "landing pages" report showed that 6,733 of the 'control-card' views were by visitors who entered our site directly from a search engine (usually Google.) Direct hits on our database after a Google search accounted for nearly 50% of the visits to our site.[9] Only 12% of visitors entered through our home page.[10]

---

9 . Over 75% of visitors to our site as a whole are referred to us by search engines.
10. The latter figure probably overestimates that number of actual users who enter through the home page since staff use was measured during the trial period, and all staff computers have the Archives' website set as 'home.'

**Figure Four: Goal Conversions in Google Analytics**

This conclusion was in one sense heartening, because it seemed to verify that our descriptive records were being correctly indexed by Google and that search engines were driving many users to our site. On the other hand, few if any of these users seemed to be finding the information they needed or contacting us for follow-up information. Nearly 71% of them bounced directly out of our site after viewing only the page on which they landed. This figure is high when compared to non-archival 'content sites.[11] Nevertheless, 1,977 users did look at additional pages. While this is a substantial number, few of them contacted us via email.

Many of the users visitors who landed on our site ran into an informational dead-end. For example, 37 users entered our site on the page describing the scrapbook of alumus Clara Hamilton. However, 36 of them left immediately. Since this page is the first one linked if you search Google for "Clara Hamilton," we can conclude that users did not find the information they were seeking. Perhaps it was the wrong Clara Hamilton, or perhaps they simply wanted actual information about her, not a description of a scrapbook. In any case, the landing page, shown in figure five, would probably be disorienting to a user, since it does not provide much visual interest or enough information to decide whether the information is relevant. The same could be said for many of our pages, since we typically do

---

11. One well-regarded analytics blog notes that "**Content websites** with high search visibility (often for irrelevant terms) can bounce at **40-60%**." http://blackbeak.conversionchronicles.com/2006/04/12/bounce-rate-or-single-page-access-industry-averages/. Accessed November 26, 2007.

not provide archival context information (information about the creator) in a prominent location on the page.



**Figure Five: Typical 'control card' series description.**

We also recognized that while some elements of our web design work well, many pages could be optimized to make them more useful to users. For example, we discovered that Google is indexing our controlled subject terms and that those terms lead visitors to our site. However, the design of the 'landing page' liked from these Google search pages was far from optimal. Figure six shows the Google result set for the search "strip mine illiniois" and figure seven the landing page the user saw after clicking the link. By making a few simple changes to our database software, we will put the text "Strip Mines" instead of "Search Results" into the Google link and will load a more relevant landing page for the user, one that shows actual links to archival descriptions or digital content.
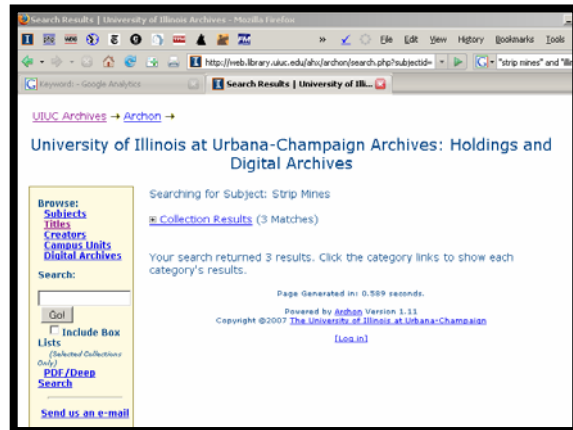
**Figure Six: Link to Search Results**



**Figure Seven: Landing Page from Link**

Time limits prevent me from discussing other examples in detail, but we plan to implement further improvements based on the results of this pilot project. For example, we will analyze users' keyword searches to identify highly-sought materials for potential digitization, we will integrate better information into page <title> elements to improve our Google page ranks, and we will move more important information to the top of common landing pages, so that users are able to see the full context of their 'hit.'

**Analysis and Conclusion**

In July of 2007, 12,008 "absolute unique visitors" landed on somewhere on the University of Illinois Archives website. Of these, only 49 emailed the archives. Quite a few of them may have phoned us or visited in person (our reference statistics show 734 users for the month). While it will always be a goal to provide good reference services to on-site our users, we must make it as easy as possible for remote users to find relevant information without an archivist's mediation—and for them to contact us

when necessary.   Remote use will drive our future growth.


Based on our pilot project, we plan to continue monitoring online use in order to improve our web presence.  Perhaps the most useful part of our study has been the fact that we now know more about our online users and can begin to tailor services and information around their needs.  It would be a mistake to intuit users' intentions based on the results of this study.

However, people use our site in very differently ways than we as staff assumed that they would.  They do not walk a prescribed path through our site to fulfill the goal—contacting the archives—that we had in mind for them.  Users referred to us by Google or another search engine thus represent a large and relatively untapped user base.   What more can be done to interest these users in our collections and services?

At the time of the study there was relatively little digital content available through our site.  It seems probable that fewer users would leave our site if they could directly access information from our holdings.  We certainly want to retain more of these potential users and thus expand our user base. While it is unlikely that we will digitize more than a fraction of our holdings in the near future, the fact that we see many users leaving our archival descriptive records without contacting us means we must design a system that seamlessly integrates archival description and digital content.  With such a system in place, we could place more emphasis on providing direct access to archival information. The system would also make it possible for users to browse descriptions of online holdings (and, when necessary, contact the archives) but these goals should be subsidiary to the main goal of providing direct access to archival records and information.

We have begun to integrate such features into our website and into the 'Archon' software that powers it.  Over time, we hope to move our website from its current focus on archival description toward direct provision of archival records and information.  The Google analytics reports that we run over the next several years will help us determine whether we are providing these records and information in a

way that matches user needs and information-seeking habits.