

# Online Versus Offline Experimentation

Roger Longbotham<sup>1</sup>

Microsoft Corporation, One Microsoft Way, Redmond, WA 98052

## Abstract

Most statisticians are taught Design of Experiments (DOE) as part of their education. This field was developed before the internet and deals almost exclusively with offline experiments. Having a career with extensive experience in both online and offline experimentation, my objective with this article is to point out how online experimentation is different from offline experimentation for those who have had some training in DOE. Since many statisticians may be considering conducting online experiments, I end with advice on what support one would need to do that.

**Key Words:** experimentation, online, offline, design of experiments, DOE, AB testing

## Introduction

First, an overview of online experimentation. Almost all large online companies run experiments - Google, Amazon, Facebook, LinkedIn, Microsoft, Etsy, Booking.com, Expedia, Yahoo!, etc. These companies run many experiments every year and most sites have multiple experiments running at any one time. These experiments are not in the usability lab but include most or all of their daily customers. The earliest online experiments I am aware of were run by Amazon in the late 90's, although other sites could have been doing experiments earlier. The practice of online experimentation goes by many names: A/B Testing, Multivariate tests, parallel tests, bucket tests, split tests, online field experiments, and more.

What can be tested on a website? Almost any change someone wants to make to the site can, and should, be tested. This could be simple changes like the color or font of some of the text, which image to show, size of images, etc. Here is a more complete list of potential changes a site could test:

## Changes that can be tested

### User Interface (UI)

- Any change to text or content
- Design of page
- Colors, fonts, images and image size, etc.

---

<sup>1</sup> Roger Longbotham was instrumental in conducting hundreds of experiments offline while an employee at Rockwell International and as Technical Director for QualPro, Inc. He has also been involved in more than a thousand online experiments at Amazon.com, Microsoft (Bing, MSN, Office online, etc.) and other companies.

### User Support Options

- Where to place user support options
- Which user support options (phone number, chat, link to support FAQs, etc.)

### Algorithms

- Search algorithms
- Personalization and recommendation algorithms
- Content importance algorithms

### Apps

#### Device specific changes

- Changes to the site that only affect tablets, phones, etc.

#### Underlying codebase

- Changes to content management system, etc.

Therefore, online experiments can and should include any visual change that a visitor may notice as well as “behind-the-scenes” changes that a visitor would not notice such as search results, recommendations etc. Even bug fixes can be tested if it is uncertain what the full extent of the effects of the bug fixes will be. In online experiments, the website without the change is to be compared to the website with the change. Common terminology for the website without the change is the control or version A. You may consider this the default, or baseline. The website with the change is the treatment or version B. The simplest online experiment is just one treatment versus control, but many experiments will have multiple treatments or versions that are to be compared to control.

### Responses from Online Experiments

The data that is collected and metrics that are formed to judge the effect of a treatment vary according to the objectives of the site and the sophistication of the site instrumentation. For example, an online retailer is keenly interested in whether a visitor makes a purchase or not and how large a purchase they make. A content site that primarily has articles (or pictures or videos) for visitors is interested in how many articles the visitor reads, how long they stay on the site and how often they return to the site. All sites should be measuring some basic metrics

#### Engagement

- Number of visitors
- Activity of each visitor
  - Number of pages viewed
  - Number of sessions<sup>2</sup>

---

<sup>2</sup> A session is defined as a set of activities/events by a visitor if there is no gap in the events of 30 minutes or more. If there is a 30 minute gap, a new session is started. The length of a session is the last time stamp minus the first time stamp.

- Average length of session
- Goals achieved, depends on the site. Some common ones:
  - Number of downloads
  - Sign up for newsletter

Loyalty

- Number of days visitor returns per month
- Number of sessions per month

In addition, metrics should be broken down by some easily collected information from the user agent such as country, language, browser, operating system and more. The response metrics an online site collects for an experiment can easily number in the hundreds.

The most common unit of randomization is the visitor, or user.

**Major Differences Between Online and Offline Experimentation**

<i>Dimension</i>	Online	Offline
<i>Sample Size</i>	<ul style="list-style-type: none"> <li>• Thousands of experimental units (EUs) to many millions</li> <li>• In most cases EU is a visitor to a website</li> </ul>	<ul style="list-style-type: none"> <li>• Often quite small (less than 20 to dozens to thousands) but some exceptions (marketing experiments)</li> <li>• Fewer EUs available</li> </ul>
<i>Signal-to-Noise Ratio</i>	<ul style="list-style-type: none"> <li>• Low S-N ratio</li> <li>• High noise and want power to see a relatively small (1-2%) change</li> </ul>	<ul style="list-style-type: none"> <li>• Usually higher S-N</li> <li>• Often can get fairly low variability metrics and precise measurements</li> <li>• Commonly looking for larger changes (10-50%)</li> </ul>
<i>Standardization and Costs</i>	<ul style="list-style-type: none"> <li>• Standardization:</li> <li>• Software platforms automate many of the elements of experimentation:                             <ul style="list-style-type: none"> <li>– randomization</li> <li>– collection of data</li> <li>– turn on/off</li> <li>– analysis, etc.</li> </ul> </li> <li>• EU cost is free or minimal (standardized to be visitor)</li> </ul>	<ul style="list-style-type: none"> <li>• Standardization:</li> <li>• Not automated                             <ul style="list-style-type: none"> <li>– design</li> <li>– experimental unit</li> <li>– randomization</li> <li>– data collection</li> <li>– analysis</li> <li>– monitoring (diagnostics)</li> <li>– control during experiment</li> </ul> </li> <li>• EU cost is anywhere from free to quite expensive (e.g. in manufacturing)</li> </ul>
<i>Nature of Variation</i>	<ul style="list-style-type: none"> <li>• Large common cause variation – visitor to visitor.</li> <li>• Robots and other outliers.</li> <li>• Very large temporal instability.</li> <li>• Most metrics are highly skewed. (Normally need &gt;1000 EUs for</li> </ul>	<ul style="list-style-type: none"> <li>• Small to large common cause (within) variation.</li> <li>• Few outliers which can often be attributed to special causes and corrected.</li> <li>• Often assumed to be stable environment.</li> </ul>

	CLT to hold. Many more needed for some metrics)	<ul style="list-style-type: none"> <li>Often metrics are not far from Normal (i.e. need &lt;100 EUs for CLT to hold)</li> </ul>
<i>Design of Experiments</i>	<ul style="list-style-type: none"> <li>Mantra: keep it simple. Most practitioners test one factor (or idea) at a time.</li> <li>May have large number of people in a company running experiments and attempts to coordinate/combine tests of multiple factors is seen as more trouble than worth (in most cases).</li> <li>Must run all treatment combinations concurrently due to temporal instability.</li> </ul>	<ul style="list-style-type: none"> <li>Experiments are usually much more manual and costly to set up and sometimes costly for experimental units so need to get the most information from each experiment.</li> <li>The incremental cost from adding more factors to test is small relative to the cost of running the experiment.</li> <li>Often can run treatments combinations sequentially.</li> </ul>
<i>Multi-factor Experiments</i>	<ul style="list-style-type: none"> <li>Tests of more than one factor in an experiment is easily carried out as a full factorial with EUs being independently randomized within each factor.</li> </ul>	<ul style="list-style-type: none"> <li>Often need to run experiments with fewer treatment combinations than a full factorial due to <ul style="list-style-type: none"> <li>sample size restrictions,</li> <li>to get better estimate of within standard deviation,</li> <li>complexity of running experiment with many TCs is logistically challenging.</li> </ul> </li> </ul>
<i>Analysis</i>	<ul style="list-style-type: none"> <li>Almost always automated. Can be simple statistical t tests or more complex tests such as permutation tests, etc.</li> <li>May have hundreds of metrics per experiment.</li> </ul>	<ul style="list-style-type: none"> <li>Ideally simple. Usually ad hoc (i.e. set up just for each experiment.)</li> <li>Relative few metrics/experiment (usually &lt;10)</li> </ul>
<i>What Can Go Wrong</i>	<ul style="list-style-type: none"> <li>Robots</li> <li>Triggering (logic that says who is in the experiment)</li> <li>Other tests at same time with same EUs (interactions)</li> </ul>	<ul style="list-style-type: none"> <li>Many sources of potential problems <ul style="list-style-type: none"> <li>Test design</li> <li>Test execution</li> <li>People not carrying out their recipe</li> <li>Measurements</li> </ul> </li> </ul>
<i>Who Runs Experiments</i>	<ul style="list-style-type: none"> <li>Many people in organization. Anyone who wants to make a change to the website</li> <li>No one else needed to carry it out.</li> </ul>	<ul style="list-style-type: none"> <li>Experiment set up and analysis by expert.</li> <li>May need help from many people to carry out.</li> </ul>
<i>Number of Experiments</i>	<ul style="list-style-type: none"> <li>Many</li> <li>Some websites have many thousands/year. (Some &gt;10,000/yr)</li> </ul>	<ul style="list-style-type: none"> <li>Relatively few for most organizations.</li> </ul>
<i>Focus of Statistical Development</i>	<ul style="list-style-type: none"> <li>Increase power.</li> </ul>	<ul style="list-style-type: none"> <li>Creative designs to</li> </ul>

- |   |   |
|---|---|
| <ul style="list-style-type: none"> <li>• Able to run more experiments with current EUs without interfering with each other.</li> <li>• Control/minimize FDR while maintaining power.</li> </ul> | <ul style="list-style-type: none"> <li>– Test more factors with fewer EUs,</li> <li>– Get exact information needed (e.g. regarding interactions), etc.</li> </ul> |
|---|---|

### What Do You Need To Run Online Experiments?

If you are comfortable running offline experiments, what else is needed to run online experiments?

- **Need a software platform.** You may have a long development cycle unless you get a start from a vendor. Ad hoc experiments are much more costly than with a platform, so most sites will develop a platform for continuous experimentation.
- **Can't do it alone.** Need support from programmers, computer scientists to set up and maintain experimentation platform. (Some of this could be provided by a vendor, but you'll need in-house expertise as well.)
- **Software skills** to do off-line/ad hoc analysis and investigate special causes.
- **Conduct training.** Large websites may have many (hundreds) of people who occasionally or constantly run experiments so need to conduct training so they are competent to run, monitor and interpret results correctly.

### Bibliography

Deng, Alex, Ya Xu, Ron Kohavi, and Toby Walker. 2013. "Improving the Sensitivity of Online Controlled Experiments by Utilizing Pre-Experiment Data." WSDM 2013.

Deng, Shaojie, Roger Longbotham, Toby Walker, and Ya Xu. 2011. "Choice of Randomization Unit in Online Controlled Experiment." Joint Statistical Meetings Proceedings. 4866-4877.

Hochberg, Yosef, and Yoav Benjamini. 1995. "Controlling the false discovery rate: a practical and powerful approach to multiple testing Series B." Journal of the Royal Statistical Society 57 (1): 289-300.

Kaushik, Avinash. 2006. "Experimentation and Testing: A Primer." Occam's Razor. May 22. Accessed 2008. <http://www.kaushik.net/avinash/2006/05/experimentation-and-testing-a-primer.html>. 10

Kohavi, Ron, Alex Deng, Brian Frasca, Roger Longbotham, Toby Walker, and Ya Xu. 2012. "Trustworthy online controlled experiments: Five puzzling outcomes explained." Proceedings of the 18th Conference on Knowledge Discovery and Data Mining. <http://bit.ly/expPuzzling>.

Kohavi, Ron, Alex Deng, Brian Frasca, Toby Walker, Ya Xu, and Nils Pohlmann. 2013. "Online Controlled Experiments at Large Scale." KDD 2013: Proceedings of the 19th

ACM SIGKDD international conference on Knowledge discovery and data mining.  
<http://bit.ly/ExpScale>.

Kohavi, Ron, Alex Deng, Roger Longbotham, and Ya Xu. 2014. "Seven Rules of Thumb for Web Site." Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '14). <http://bit.ly/expRulesOfThumb>.

Kohavi, Ron, Roger Longbotham, and Toby Walker. 2010. "Online Experiments: Practical Lessons." IEEE Computer, September: 82-85. <http://bit.ly/expPracticalLessons>.

Kohavi, Ron, Roger Longbotham, Dan Sommerfield, and Randal M. Henne. 2009. "Controlled experiments on the web: survey and practical guide." Data Mining and Knowledge Discovery 18: 140-181. <http://bit.ly/expSurvey>.

McFarland, Colin. 2012. Experiment!: Website conversion rate optimization with A/B and multivariate. New Riders.

Schrage, Michael. 2014. The Innovator's Hypothesis: How Cheap Experiments Are Worth More than Good Ideas. MIT Press.

Tang, Diane, Ashish Agarwal, Deirdre O'Brien, and Mike Meyer. 2010. "Overlapping Experiment Infrastructure: More, Better, Faster Experimentation." KDD.