

marketingsherpa marketingexperiments

optimization summit 2012

MEASURE.
TEST.
CONVERT.

June 11-14 · Denver

Validity for Decision Makers

A guide to making optimization
decisions with confidence

Sergio Balegno

Director of Research

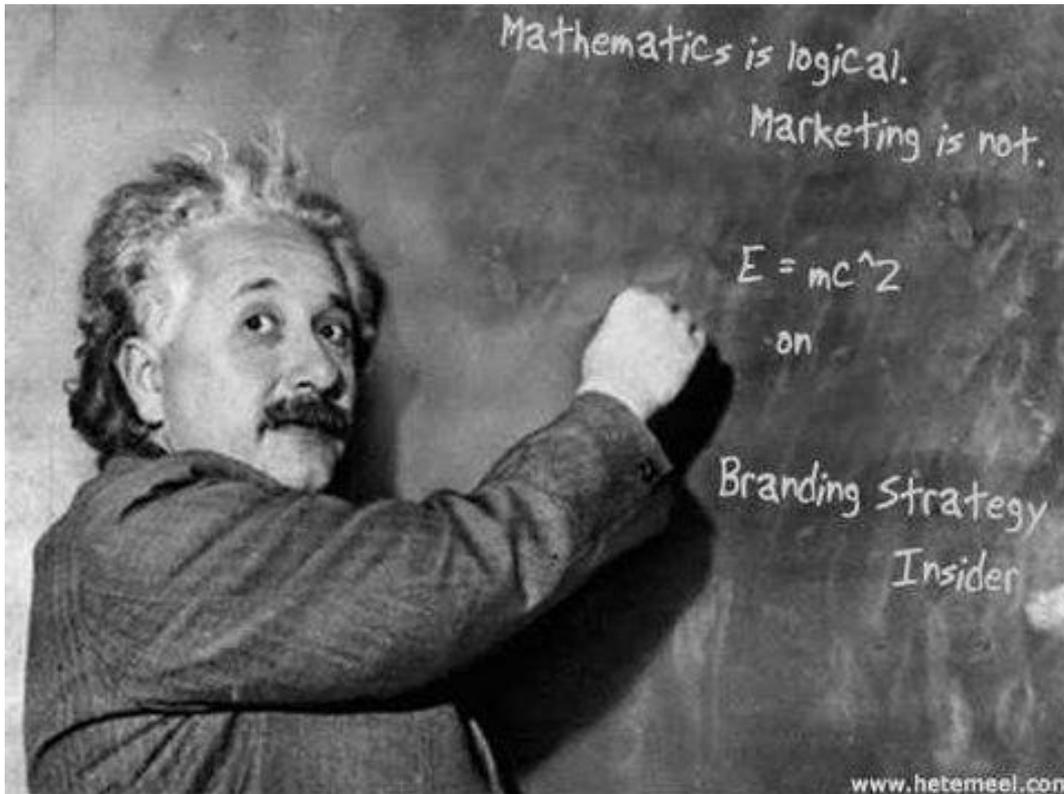
MECLABS

Bob Kemper

Sr. Director of Sciences

MECLABS

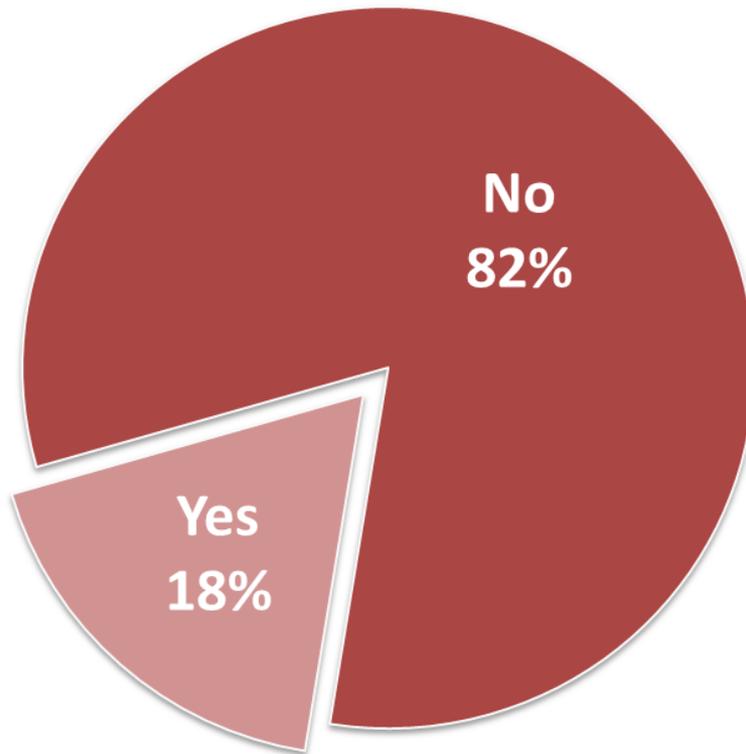
What we don't know about validity can hurt us



- Do you monitor the validity of your tests?
- How well do you understand the math used to determine statistical validity?

82% of marketers surveyed do not monitor optimization test validity

Regularly monitor optimization test validity



- If you don't, then you're in good company...

... or at least you're surely not alone!

 marketing **sherpa** Source: ©2012 MarketingSherpa Website Optimization Benchmark Survey
Methodology: Fielded April 2012, N=2677

Validity for Decision Makers

1 An invalid test...

... is worse than no test at all

Principles and Taxonomy of Optimization Testing

2 Test validity:

There's more to it than 'Math'

3 Validity Threats

History Effect

Instrumentation Effect

Sample Selection Effect

Sample Distortion Effect

4 Math and practical examples



Experiment ID: *Protected*

Location: MarketingExperiments Research Library

Test Protocol Number: TP2047

Research Notes:

Background: An ecommerce site focusing on special occasion gifts

Goal: To increase email clickthrough and conversion-to-sale

Primary research question: Which email design will yield the highest conversion rate?

Approach: Series of A/B variable cluster split tests

Validity for Decision Makers – CS 1

- In a series of daily email tests lasting 5 weeks, we tested 7 different email templates designed for their most loyal customer segment. Below are examples of three of those email templates tested.

Control Template



Treatment Template #1



Treatment Template #2



...

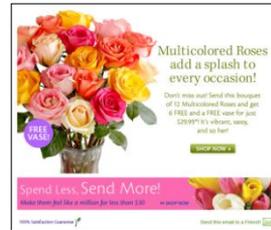
Validity for Decision Makers – CS 1

Week 1

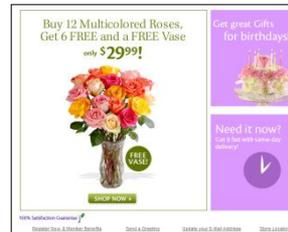
Week 2

Week 3

Control Template →



Treatment Template #1 →



Treatment Template #2 →





74% Increase in Conversion

Simple side-by-side layout outperformed the Control

Week 1 Results

Template Version	CR	Rel. Diff.
Control Template	14.01%	-
Treatment Template #1	17.06%	25.68%
Treatment Template #2	24.38%	74.05%

- ✓ **What you need to understand:** After a week of testing, Treatment 2 converted at a rate 74.05% higher than the control.

However, as subsequent samples were collected, there was a noticeable shift in results.

Validity for Decision Makers – CS 1

- In the subsequent 2 weeks, the relative conversion rates of the Experimental Treatment templates declined to as low as -6% for Treatment 1, and as low as 3% for Treatment 2 compared to Control.

Week 2 Results

Template Version	CR	Rel. Diff.
Control Template	24.08%	-
Treatment Template #1	22.59%	-6.17%
Treatment Template #2	24.89%	3.38%

Week 3 Results

Template Version	CR	Rel. Diff.
Control Template	19.04%	-
Treatment Template #1	19.09%	0.26%
Treatment Template #2	20.74%	8.93%

Validity for Decision Makers – CS 1

- For the remainder of the 5 week test, the relative CR never again exceeded +9%, but stabilized indicating that something had potentially invalidated the first week of tests.

	Week 1	Week 2	Week 3	Week 4	Week 5
Control	14.01%	24.08%	19.04%	19.77%	20.05%
Treatment #1	17.06%	22.59%	19.09%	19.42%	19.52%
Treatment #2	24.38%	24.89%	20.74%	17.93%	20.50%
Rel. Diff. (T2)	74%	3%	9%	-9%	2%

- As we drilled down into the numbers, we learned that there was a problem during the first week related to incoming traffic to the Control.

We'll come back to this test a little later, when we talk about specific validity threats...



74% Increase in Conversion

Treatment 2 converted 74.05% more recipients than the Control

Week 1 Results		
Template Version	CR	Rel. Diff.
Control Template	14.01%	-
Treatment Template #1	17.06%	25.68%
Treatment Template #2	24.38%	74.05%



What you need to understand: The apparent significant performance boost indicated after Week 1 was in fact a ‘phantom’ gain.

There actually was no statistically significant difference among the treatments.

An invalid test...

... is worse than no test at all.

Those who neglect to test know the risk they're taking, and make their changes cautiously and with a healthy trepidation.

Those who conduct invalid tests are blind to the risk they take, and make their changes boldly and with an unhealthy sense of confidence.

Validity for Decision Makers

The purpose and taxonomy of
optimization testing

 **Key Principle**

In online testing, the ultimate measure of a test is its **Utility**:

In marketing optimization...

A **useful** test is one that helps you predict future customer behavior.

Five elements compose the Utility of an experimental test. For us, they are expressed as...

Online Testing Heuristic:

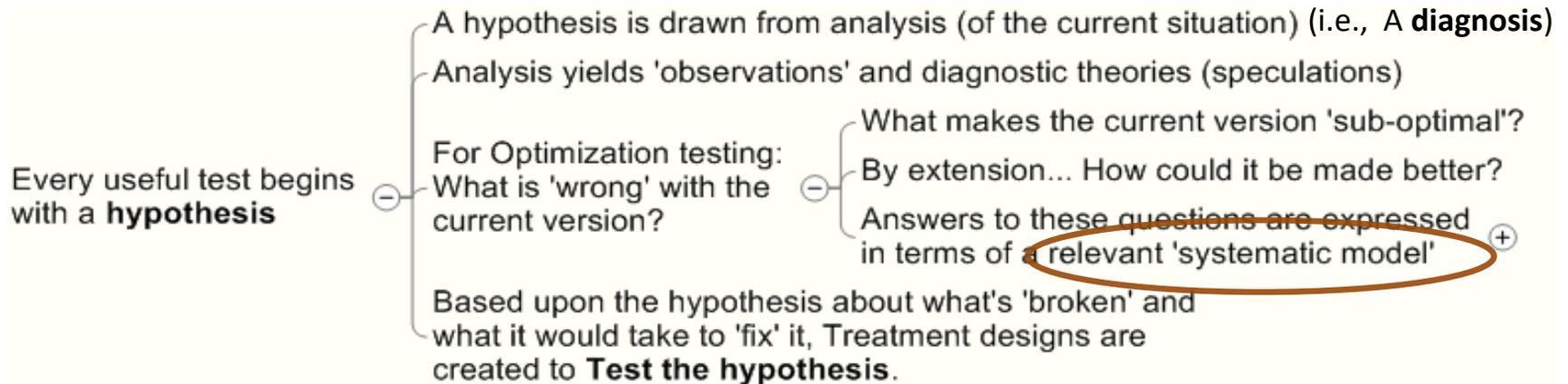
$$u = 2q + t + m + 2v + i \text{ ©}$$

- u = Utility
- q = Research Question
- t = Treatment
- m = Metric System
- v = Validity Factor
- i = Interpretation

 **Key Principle**

In online testing, the ultimate measure of a test is its **Utility**:

In marketing optimization...



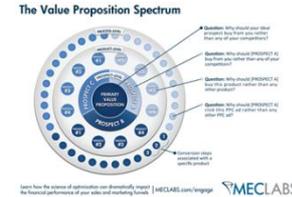
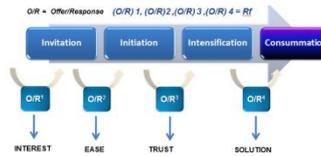
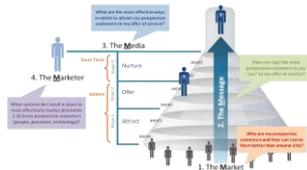
The hypothesis is one that proposes a cause or a **reason that customers are currently behaving** as they do...

... and describes **how that behavior would differ** if something specific were altered about the experience they encounter (when they receive your email, when they arrive at your site, ...)

Process: **Diagnosis** → **Hypothesis** → **Treatment Design & Development**

Systematic Models

Overarching Meta-theory



Landing Pages:

100% Satisfaction Guarantee

$C = 4m + 3v + 2(i-f) - 2a$ ©

The Most Accurate Mailing Lists Available!

Opr > Oprn > Ocn ©

Setup your FREE access to

- Search our business and consumer database
- Build a list 24 hours a day, 7 days a week
- Preview and download mailing lists
- Get expert advice on how to most effectively turn leads into sales
- Access our exclusive Resource Center which includes FREE white papers

First Name:

Last Name:

Company:

What People Are Saying

It's a powerful tool for small businesses to market like the big guys

Haydens
Sparta, New Jersey

I would recommend **_____** to anyone looking for speedy service, accurate listings and great customer service. It is refreshing to receive a follow-up phone call within a week, just to make sure that I am satisfied.

Ruhler Auction & Realty Inc.
Hastings, Nebraska

It's great, one-stop shopping.

Email effectiveness...

$$eme = rv(of+ i) - (f + a) \text{ ©}$$

$$ec < op < ct < lp \text{ ©}$$

 **Key Principles**

- **Ultimately, a test is designed for one purpose...**

To help us “**Accept**” or “**Reject**” a hypothesis.

- **The extent to which it does this is the measure of its Utility:**

Validity for Decision Makers

Validity Threats



Key Principles

In online testing, there are four common validity threats that must be considered:

- History Effect
- Instrumentation Effect
- Selection Effect
- Sampling Distortion Effect



Essential Definitions

History Effect – when a test variable is affected by the passage of time



Experiment ID: *Protected*

Location: MarketingExperiments Research Library

Research Notes:

Background: Online sex offender registry service for parents concerned about their surrounding areas

Goal: To increase the click-through rate of a PPC advertisement

Primary research question: Which ad headline will produce the most click-through?

Test Design: A/B/C/D split test focusing on the headlines of a PPC advertisement

Validity for Decision Makers – History Effect – Experiment

We prepared a headline test using Google AdWords as the split-testing platform. The headlines were chosen by the participants of the certification course from a pool which they created. The test was conducted for seven days and received 55,000 impressions.

Child Predator Registry

Identify sex offenders living in your area. Protect your kids today.

www.XXXXXXXXXXXXXXXXXX.com

Predators in Your Area

Identify sex offenders living in your area. Protect your kids today.

www.XXXXXXXXXXXXXXXXXX.com

Is Your Child Safe?

Identify sex offenders living in your area. Protect your kids today.

www.XXXXXXXXXXXXXXXXXX.com

Find Child Predators

Identify sex offenders living in your area. Protect your kids today.

www.XXXXXXXXXXXXXXXXXX.com

Validity for Decision Makers – History Effect – Experiment



The screenshot shows the MSNBC website interface. At the top, there is a search bar, a timestamp 'Updated: 6:34 p.m. ET May 4, 2006', and navigation links for 'MSNBC Home', 'Hotmail', and 'Sign In'. The main header features the NBC peacock logo and the text 'MSNBC Home >> Dateline NBC >> 'To Catch a Predator''. Below this, the 'DATELINE TO CATCH A PREDATOR' title is displayed in large, stylized letters. A navigation menu on the left lists various categories like 'Dateline NBC', 'Control', 'Predator', etc. The main content area features a large photo of a man in a brown jacket sitting at a bar, talking to a man in a suit. Below the photo is the headline 'Admitted child abuser caught in Dateline sting' and a sub-headline 'To Catch a Predator: Cameras follow the potential predators from chatting about sex online, to shuffling through the criminal justice system. Plus, the most frightening man caught in this investigation. Chris Hansen reports. • FULL STORY'. To the right of the main article is a video player titled 'VIDEOS The next investigation' with a 'Launch' button. Below the video player is a section titled 'WEDNESDAYS ON NBC' with text about the show's schedule. At the bottom of the page, there is an 'advertisement' placeholder.

- During the test, Dateline aired a special called “*To Catch a Predator*” viewed by approximately 10 million individuals
- Throughout this program sex offenders are referred to as “predators”

Validity for Decision Makers – History Effect – Experiment

Headline	Impressions	Clicks	CTR
<i>Predators</i> in Your Area	21,096	1,423	6.74%
Child <i>Predator</i> Registry	14,712	652	4.43%
Find Child <i>Predators</i>	18,459	817	4.42%
Is Your Child Safe?	15,128	437	2.89%

-  **What you need to understand:** In the two days following the Dateline special, there was a considerable spike in overall click-through, but a relative difference between those ads with "predator" in the headline and those ads without "predator" of up to 133%. **So, in effect, an extraneous variable (the Dateline special) associated with the passage of time jeopardized the validity of our experiment.**



Key Principles

1. In online testing, there are four common validity threats that must be addressed: History Effect, Instrumentation Effect, Selection Effect, and Sampling Distortion Effect.
2. To minimize the **History Effect**, monitor external events, compare current with historical data, and look for anomalies in the testing data.

Examples of external events that commonly impact test results:

- *Media events (news and entertainment)*
- *Competitor marketing initiatives*
- *Internal marketing initiatives*
- *Economic changes (esp. 'shocks')*
- *Political changes*



Essential Definitions

History Effect – when a test variable is affected by the passage of time

Instrumentation Effect – when a test variable is affected by a change in the measurement instrument



Experiment ID: *(Protected)*

Location: MarketingExperiments Research Library

Research Notes:

Background: Online ‘people search’ service offering background checks about criminal records and other public information about individuals, limited to the U.S.

Goal: To increase the rate of Clickthrough to the order-start page, and ultimately to final conversion to sale.

Primary research question: Which landing page will produce the highest Order-start page click-through?

Test Design: Multi-factor (MV) test with 5 factors at 2 levels each.

Validity for Decision Makers – Instrumentation Effect – Experiment

Control Values

The screenshot shows the NetDetective website interface. The 'Member Login' section at the top right is highlighted with an orange box. The main navigation bar includes 'HOME', 'RESOURCES', 'USERS', 'FACTS & QUESTIONS', and 'SIGN UP NOW'. Below the navigation, a banner reads 'Search Over 211 Million Records, Find People, Run Background Checks & Criminal Records'. The main content area features a 'WELCOME TO YOUR PERSONAL PRIVATE INVESTIGATOR' message and a 'Customer Comments' section. A blue box highlights a search-related image on the left. A purple box highlights the 'CURRENT DATABASE STATISTICS' section, which lists search capabilities like finding U.S. residents, happy users, and instant access. A green box highlights the 'ENDORSED BY THE National Association of Independent Private Investigators' logo. A red box highlights the sign-up process, which includes a 'Become A NetDetective In Just 2 Steps' button, a 'Step 1: Enter Email Address' input field, and a 'Continue To Step 2' button.

Treatment Values

The screenshot shows the NetDetective website interface, similar to the control version. The 'Member Login' section at the top right is highlighted with an orange box. The main navigation bar includes 'HOME', 'RESOURCES', 'USERS', 'FACTS & QUESTIONS', and 'SIGN UP NOW'. Below the navigation, a banner reads 'Search Over 211 Million Records, Find People, Run Background Checks & Criminal Records'. The main content area features a 'WELCOME TO YOUR PERSONAL PRIVATE INVESTIGATOR' message and a 'Customer Comments' section. A blue box highlights a search-related image on the left. A purple box highlights the 'Current Database Statistics' section, which lists search capabilities like finding U.S. residents, happy users, and instant access. A green box highlights the 'ENDORSED BY THE National Association of Independent Private Investigators' logo. A red box highlights the sign-up process, which includes a 'Start Your Search in Just 2 Steps!' button, a 'Step 1: Enter Email Address' input field, and a 'Become a Net Detective' button.

- In this test, we were comparing five variables (Factors) with two differing values (Levels) each (highlighted above).

Experiment Notes:

- We discovered that in the testing software, a “fail safe” feature was enabled specifically that, if for any reason the treatment page was not running correctly, the page would default back to the control page.
 - This was enabled by loading hidden Control versions of experimental variables whenever the Treatment was loaded.
 - This caused the Treatment pages to have substantially longer load times than the Control when they were rendered on the visitor’s browser.

Validity for Decision Makers – Instrumentation Effect – Experiment

Page Load Time Reference Chart (in seconds)							
	Page Size(kb)	Connection Rate					
		14.4	28.8	33.6	56	128 (ISDN)	1440 (T1)
Control Page							
	50	35.69	17.85	15.30	9.18	4.02	0.36
	75	53.54	26.77	22.95	13.77	6.02	0.54
	84.9	60.61	30.30	25.98	15.59	6.82	0.61
Treatment Page							
	100	71.39	35.69	30.60	18.36	8.03	0.71
	125	89.24	44.62	38.24	22.95	10.04	0.89
	137	97.80	48.90	41.92	25.15	11.00	0.98
	150	107.08	53.54	45.89	27.54	12.05	1.07
	175	124.93	62.47	53.54	32.13	14.05	1.25
	Add. Load Time (s)	37.19	18.60	15.94	9.56	4.18	0.37

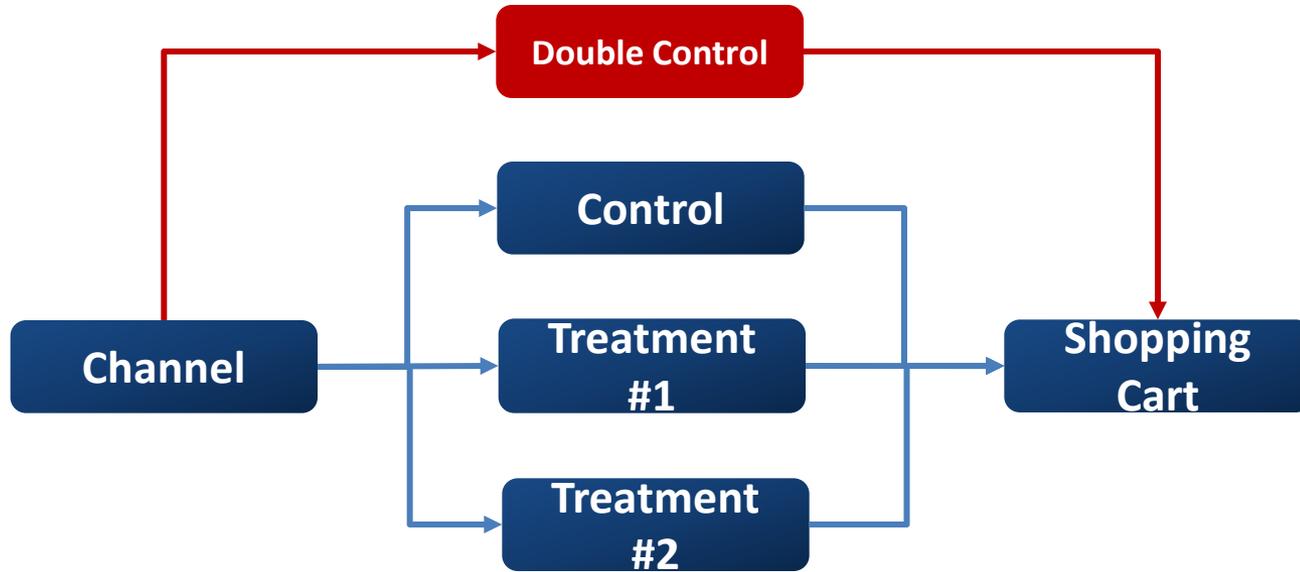
- ✓ **What you need to understand:** The Treatment page took 60% longer to render than Control. At speeds slower than T1, this translated to several seconds. The difference in load times caused the “user experience” to be asymmetric among the page versions, thereby threatening test validity.



Key Principles

1. In online testing, there are four common validity threats that must be addressed: History Effect, Instrumentation Effect, Selection Effect, and Sampling Distortion Effect.
2. To minimize the **History Effect**, monitor external events, compare current with historical data, and look for anomalies in the testing data.
3. To address the **Instrumentation Effect**, never presume your analytics system is accurate; carefully monitor and test for anomalies.

Validity for Decision Makers – Instrumentation Effect – Double Control



- An additional safeguard against instrumentation anomalies is to include a 'double control' treatment that is identical to the Control for comparison.
- A high variance between Control and Double Control indicates potential problems with instrumentation.
- Note: This approach innately increases test duration



Essential Definitions

History Effect – when a test variable is affected by an extraneous variable associated with the passage of time

Instrumentation Effect – when a test variable is affected by a change in the measurement instrument

Selection Effect – when a test variable is affected by different types of subjects not being evenly distributed among experimental treatments



Experiment ID: *Protected*

Location: MarketingExperiments Research Library

Test Protocol Number: TP2047

Research Notes:

Background: An ecommerce site focusing on special occasion gifts

Goal: To increase email clickthrough and conversion-to-sale

Primary research question: Which email design will yield the highest conversion rate?

Approach: Series of A/B variable cluster split tests

Validity for Decision Makers – Selection Effect – Experiment

- In a series of tests lasting 5 weeks, we tested 7 different email templates designed for their most loyal customer segment. Below are examples of three of those email templates tested.

Control Template



Treatment Template #1



Treatment Template #2



Validity for Decision Makers – Selection Effect – Experiment

Control Template



Week 1



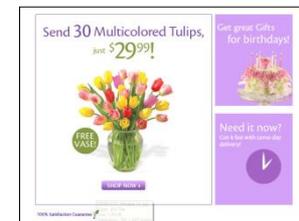
Week 2



Week 3



Treatment Template #1



Treatment Template #2





74% Increase in Conversion

Simple side-by-side layout outperformed the Control

Week 1 Results

Template Version	CR	Rel. Diff.
Control Template	14.01%	-
Treatment Template #1	17.06%	25.68%
Treatment Template #2	24.38%	74.05%

- ✓ **What you need to understand:** After a week of testing, Treatment 2 converted at a rate 74.05% higher than the control.

However, as the test samples continued during the 2nd and subsequent weeks, there was a noticeable shift in results.

Validity for Decision Makers – Selection Effect – Experiment

- In the subsequent 2 weeks, the relative conversion rates of the Experimental Treatment templates declined to as low as -6% for Treatment 1, and as low as 3% for Treatment 2 compared to Control.

Week 2 Results

Template Version	CR	Rel. Diff.
Control Template	24.08%	-
Treatment Template #1	22.59%	-6.17%
Treatment Template #2	24.89%	3.38%

Week 3 Results

Template Version	CR	Rel. Diff.
Control Template	19.04%	-
Treatment Template #1	19.09%	0.26%
Treatment Template #2	20.74%	8.93%

Validity for Decision Makers – Selection Effect – Experiment

- For the remaining test duration, the relative CR never exceeded +9%, but stabilized indicating that something has potentially invalidated the first week of tests.

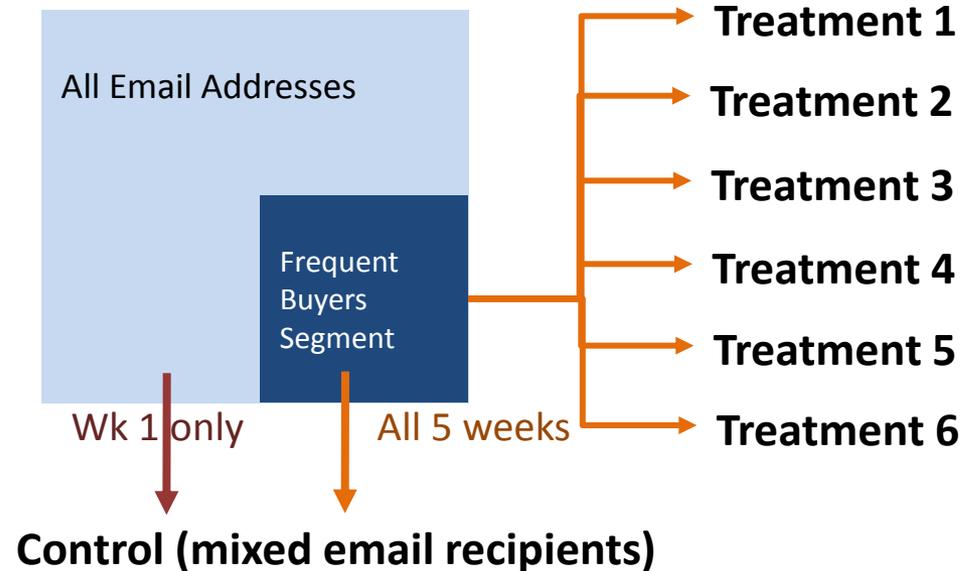
	Week 1	Week 2	Week 3	Week 4	Week 5
Control	14.01%	24.08%	19.04%	19.77%	20.05%
Treatment #1	17.06%	22.59%	19.09%	19.42%	19.52%
Treatment #2	24.38%	24.89%	20.74%	17.93%	20.50%
Rel. Diff. (T2)	74%	3%	9%	-9%	2%

The diagram features a blue bracket spanning from the start of Week 2 to the end of Week 5. Two red arrows point upwards: one from below Week 1 and another from below Week 3. An orange oval highlights the Control row for Weeks 1 through 5, and another orange oval highlights the Rel. Diff. (T2) row for Weeks 2 through 5. An orange arrow points from the 74% value in Week 1 of the Rel. Diff. (T2) row to the 17.06% value in Week 1 of the Treatment #1 row.

- As we drilled down into the numbers, we learned that there was a problem during the first week related to incoming traffic to the Control...

Validity for Decision Makers – Selection Effect – Experiment

- During the entire period, the Treatments and the Control both received evenly distributed traffic coming from a specific segment consisting of their most frequent buyers.
- However, during the first week only, the Control also received traffic from a mixture of the rest of their email list.



Validity for Decision Makers – Selection Effect – Experiment Results



74% Increase in Conversion

Treatment 2 converted 74.05% more recipients than the Control

Week 1 Results		
Template Version	CR	Rel. Diff.
Control Template	14.01%	-
Treatment Template #1	17.06%	25.68%
Treatment Template #2	24.38%	74.05%



What you need to understand: The systematic difference in distribution of arriving traffic among the treatments caused a Selection Effect validity threat which invalidated the first week's test results. So, the apparent significant performance boost indicated after Week 1 was in fact a 'phantom' gain, and there was no 'real' performance difference among the Control and the Experimental treatments.



Key Principles

1. In online testing, there are four common validity threats that must be addressed: History Effect, Instrumentation Effect, Selection Effect, and Sampling Distortion Effect.
2. To minimize the **History Effect**, monitor external events, compare current with historical data, and look for anomalies in the testing data.
3. To address the **Instrumentation Effect**, never assume your analytics system is accurate and carefully monitor and test for any anomalies.
4. To address the **Selection Effect**, ensure the profile of your test subjects matches as closely as possible the profile of your visitors, and that the assignment of treatments to test subjects is random.



Essential Definitions

History Effect – when a test variable is affected by an extraneous variable associated with the passage of time

Instrumentation Effect – when a test variable is affected by a change in the measurement instrument

Selection Effect – when a test variable is affected by different types of subjects not being evenly distributed among experimental treatments

Sample Distortion Effect – the effect on a test outcome caused by failing to collect a sufficient number of observations

Validity for Decision Makers – Sample Distortion Effect – Case Study



Experiment ID: *(Protected)*

Location: MarketingExperiments Research Library

Research Notes:

Background: Consumer company that offers online brokerage services

Goal: To increase the number of accounts created online

Primary research question: Which page design will generate the highest rate of conversion?

Test Design: A/B/C/D multi-factor test

Validity for Decision Makers – Sample Distortion Effect – Case Study

Control

CONTACT US | CHAT | SEARCH

- OPEN AN ACCOUNT
- DEMO SIGN UP
- COMMUNITY ACCESS
- EDUCATIONAL COURSES
- AND MORE

GET STARTED NOW

ROTATING BANNER

Forex

... offers something unique in the forex marketplace. We took the power of our ... routing technology and connected it to banks and pools of liquidity to create a true non-dealing desk platform. Traders have the ability to route orders via fifteen [order types](#) to qualified destinations through our ... platform. They also have the ability to cross orders anonymously in our internal order book via our The technology facilitates "best execution" for you and generates price competition, which means traders can "skip the middle man" and save money. We even let you trade between the spread and see your order reflect to the entire marketplace. [Open a forex account](#) with ... and put yourself in control.

CHAT WITH A LIVE AGENT
CLICK HERE

TESTIMONIALS

"You guys do a great job and I love the navigator. IMHO I do not believe there is a better platform out there, and you just keep making it better."

[Read More](#)

**The testimonial is not indicative of future success

NFA
NATIONAL FUTURE ASSOCIATION
CLICK TO SEE OUR REGISTRATION

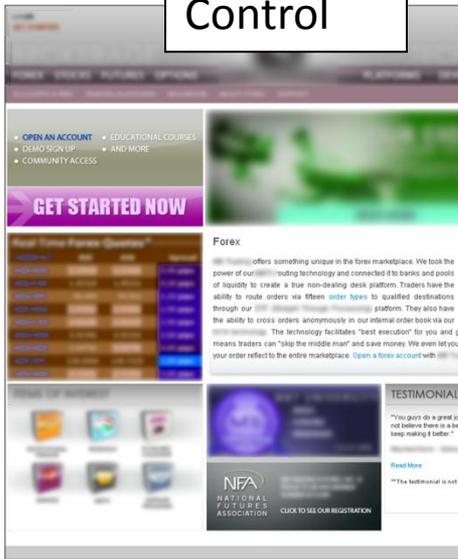
- Heavily competing imagery and messages
- Multiple competing calls-to-action

Validity for Decision Makers – Sample Distortion Effect – Case Study

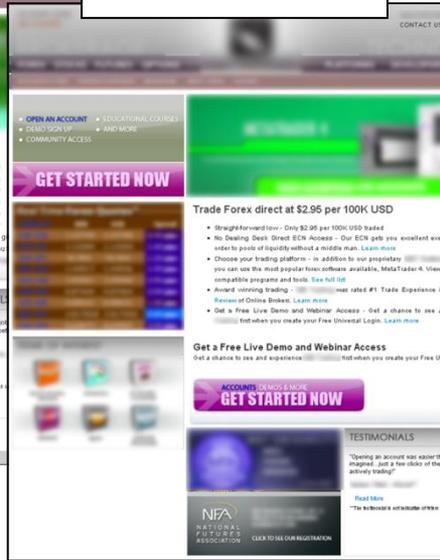
Hypotheses:

- Difficulty-oriented Friction from multiple competing CTA's is causing drop-out
- Clarity of Value proposition mitigated by lack of visual cues and large block of unbroken text.

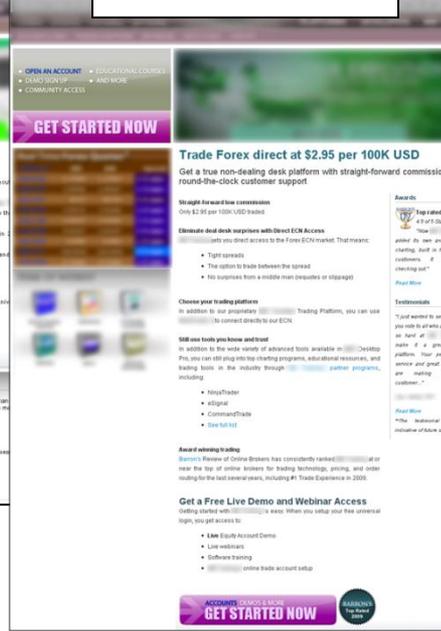
Control



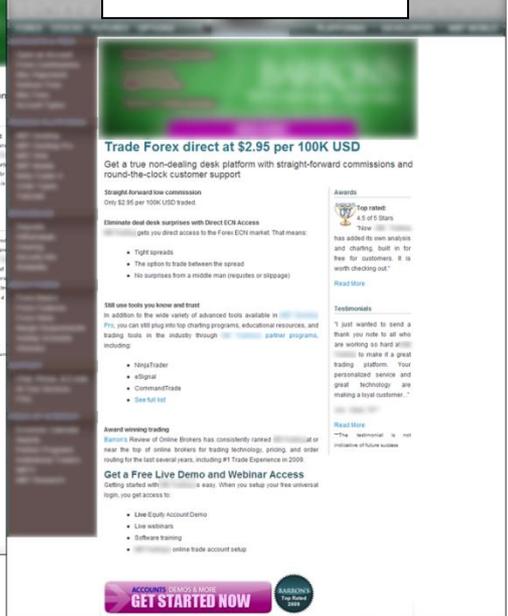
Treatment 1



Treatment 2



Treatment 3



Observations (from CIA):

- Heavily competing imagery & messages
- Multiple competing CTA's

- Most elements on the page remain the same
- Headline added
- Bulleted copy highlights VP items
- Large, clear CTA added

- Left column is the same, but footer items removed
- Long copy, vertical flow
- Added awards & testimonials in right column
- Large, clear CTA – Like T1

- Similar to Treatment 2, except left-column reduced even further
- Still a long copy, vertical flow, single call-to-action design

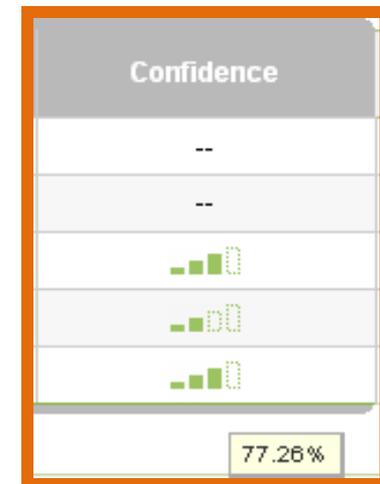
Validity for Decision Makers – Sample Distortion Effect – Case Study



No Significant Variation

None of the Experimental treatments conclusively outperformed the Control

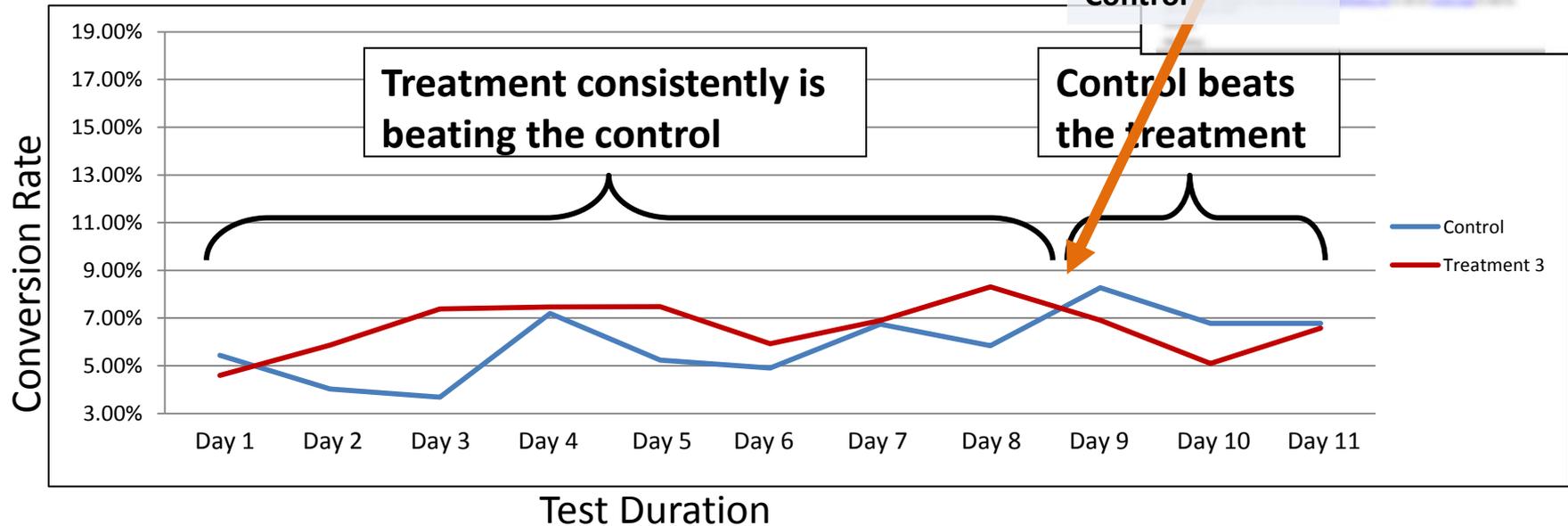
Test Designs	Conversion Rate	Relative Diff.
Control	5.95%	-
Treatment 1	6.99%	17.42%
Treatment 2	6.51%	9.38%
Treatment 3	6.77%	13.70%



What you need to understand: After eleven days, the test management platform reports remained inconclusive. None of the treatments significantly outperformed the control.

Validity for Decision Makers – Sample Distortion Effect

- However, a review of the daily sample data showed a distinct shift in relative performance beginning on Day 8 and lasting 2 days.
- Subsequent investigation revealed that during the test, an email was sent that drove traffic to Control, skewing the sampling distribution.

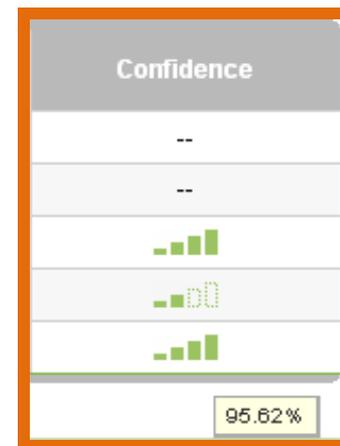




31% Increase in Conversions

The highest performing treatment outperformed the Control by 31%

Test Designs	Treatment	Variation
Control	5.35%	-
Treatment 1	6.66%	25%
Treatment 2	6.13%	15%
Treatment 3	7.03%	31%



- ✓ **What you need to understand:** Inclusion of only the data collected prior to the renegade email send established that Treatment 3 significantly outperformed Control. Measured Level of Confidence: 95.6%.

This underscores the need for a simple way to address the statistical validity threats associated with the testing sample .

Validity for Decision Makers

Validity Threats: Sample Distortion Effect

Validity for Decision Makers – Sample Distortion – The Challenge

QUESTION: At what point could we confidently assume that one web page will outperform another?

The Weekend Starts Here.
Put yourself in the center of the conversation by getting [redacted] delivered directly to your home. And be the first to know which movies to see, which destinations to visit, which investments to pursue, and much, much more. Order now to take advantage of our introductory rate — 50% off the regular subscription price.*

Start a New Subscription.
[redacted] (Friday-Sunday) for just **\$5.20** per week [Details](#)

Delivery Information.

First name
Last name
Address

City, State, ZIP **JACKSONVILLE, FL 32250**
Phone () -
E-mail address

Billing Information.

Billing address, if different from delivery address

Pay by credit card and get your first **12 weeks** at 50% off
 Pay by mail and get your first **8 weeks** at 50% off

Credit Card
Card Number
Exp. Date MM YY YY
Card ID Number [What's this?](#)

ABOUT OUR CERTIFICATION

Friday. Saturday. Sunday.

Vs.

The Weekend Starts Here.
Put yourself in the center of the conversation by getting [redacted] delivered directly to your home. And be the first to know which movies to see, which destinations to visit, which investments to pursue, and much, much more. Order now to take advantage of our introductory rate — 50% off the regular subscription price.*

Start a New Subscription.

[redacted] (Friday-Sunday) for just **\$5.20** per week [Details](#)
 Daily Delivery (7 Days) for just **\$7.40** per week [Details](#)
 Weekday (Monday-Friday) for just **\$3.70** per week [Details](#)
 Sunday Only for just **\$3.75** per week [Details](#)

Delivery Information.

First name
Last name
Address

City, State, ZIP **JACKSONVILLE, FL 32250**
Phone () -
E-mail address

Billing Information.

Billing address, if different from delivery address

Pay by credit card and get your first **12 weeks** at 50% off
 Pay by mail and get your first **8 weeks** at 50% off

Credit Card
Card Number
Exp. Date MM YY YY
Card ID Number [What's this?](#)

ABOUT OUR CERTIFICATION

Friday. Saturday. Sunday.

After Three days?
Five days?
Seven days?

Fifty seven days?

Validity for Decision Makers

Statistical Confidence (in Theory)



Essential Definitions

Sample Distortion Effect – the effect on a test outcome caused by failing to collect a sufficient number of observations

Statistical Level of Confidence – the statistical probability that there really is a performance difference between the control and experimental treatments based upon the samples collected to date

H_0 is that there is NO difference...



Key Principles

1. There are essentially three factors that determine an outcome's statistical certainty:
 - **The Amount of Difference** – The magnitude of the difference in the primary metric between the Control and an Experimental treatment. (usually expressed as a percentage value)
 - **The Success Rate** – The number of experimental trials in which the outcome was favorable with respect to the primary success metric divided by the total samples received. (usually expressed as a percentage value)
 - **The Samples Received** – The total number of experimental trials used to compute the measured sample success rate. (e.g., number of visitors to a landing page, number of email messages sent).
2. An outcome's true statistical certainty is then compared with the desired level of certainty (chosen prior to experimentation).

Validity for Decision Makers – Sample Distortion – Underlying Statistical Principles

THE MATH: *Statistical principles underlying Significance and Level of Confidence as it relates to online optimization testing*

Principle: Distinguish between Random Variation and a ‘Real’ systematic difference. I.e., what are the chances that what I am observing could happen even though the two treatments were really IDENTICAL?

Mathematical Method:

Determine:

What is the best estimate of the **magnitude of difference** between the observed relative performance of the treatments?

Using: Observed absolute ‘Difference in Sample Means’

$$(\text{Avg \%success}_T) - (\text{Avg \%success}_C)$$

$$\bar{s} = \frac{\sum s_T}{n_T} - \frac{\sum s_C}{n_C}$$

Relative %Difference in Avg success rates

$$\text{Rel \%Diff} = \frac{\%S_T - \%S_C}{\%S_C}$$

Validity for Decision Makers – Sample Distortion – Underlying Statistical Principles

THE MATH: *Statistical principles underlying Significance and Level of Confidence as it relates to online optimization testing*

Principle: Distinguish between Random Variation and a ‘Real’ systematic difference. I.e., what are the chances that what I am observing could happen even though the two treatments were really IDENTICAL?

Mathematical Method:

Determine:

What is the best measure of the **Dispersion (amt. of variation)** within the treatment performance observations?

Using: Standard Deviation (sd) of ‘Difference in Sample Means’

and Presumption of Normal distribution of sample data
(this is not an unreasonable presumption for a well-designed optimization test, and allows us to use the Central Limit theorem)

Margin of Error (2*sd)

$$sd = \sqrt{\frac{p_a(1-p_a)}{n_a} + \frac{p_b(1-p_b)}{n_b}}$$

Probability Density: Normal distribution

$$P(x) = \frac{e^{-x^2/2}}{\sqrt{2\pi}}$$

$$d = 2 \sqrt{\frac{p_a(1-p_a)}{n_a} + \frac{p_b(1-p_b)}{n_b}}$$

Validity for Decision Makers – Sample Distortion – Underlying Statistical Principles

THE MATH: *Statistical principles underlying Significance and Level of Confidence as it relates to online optimization testing*

Principle: Distinguish between Random Variation and a ‘Real’ systematic difference. I.e., what are the chances that what I am observing could happen even though the two treatments were really IDENTICAL?

Mathematical Method:

Determine:

What is the likelihood that the observed difference could just be the result of random variation, even if the two treatments are identical?

Using: Statistical Level of Confidence

Confidence Interval

$$\hat{p} \pm z_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

\hat{p} is the proportion of successes

α is the error percentile

n is the sample size

Validity for Decision Makers – Sample Distortion – Statistical Confidence Level

Additional Notes:

- Though perhaps slightly counterintuitive, statistical confidence level is really aimed at answering the question: *“What are the chances that the apparent ‘winning’ Treatment is, in fact, **not really better** (...is perhaps even worse) than the Control”*

Example:

A 95% level of confidence would mean that there is a 5% chance that the Treatment’s performance is not better or worse than (i.e., is not ‘different from’) the Control.

- The complement of the statistical level of confidence is referred to as the **statistical significance level**. This is what is used in the MECLABS Test Protocol.

Validity for Decision Makers – Sample Distortion – Underlying Statistical Principles

THE PURPOSE: *The purpose of Sample Distortion related validity calculations in Optimization is to enable you to ...*

- **Plan Tests**
- **Assess the outcome of tests**
- **Facilitate interpretation of the outcome to improve your customer theory**

Objective: Distinguish between Random Variation and a ‘Real’ systematic difference. I.e., what are the chances that what I am observing could happen even though the two treatments were really IDENTICAL?

Decide: What is the desired level of certainty?
How big a difference do I need to make it worthwhile to make a change?

Measure: Observed ‘difference in sample means’
How many samples (trials, visits, ...)?
How many ‘successes’ (conversions, clicks, ...)?

Know: Your traffic rates (# visits per day, # in email list, ...)
Number of treatments you intend to test

Validity for Decision Makers – Sample Distortion – Underlying Statistical Principles

THE OBJECTIVE: *The purpose of Sample Distortion related validity calculations in Optimization is to enable you to : Plan, Assess, Facilitate interpretation ...*

What you'll use the Validity math for:

Plan Tests:

How long should I expect the test to take if I achieve my desired min. difference amount?
What are my options if the estimated test duration is longer than I can afford to wait?

Assess Outcome of Tests:

Have I reached a sufficient sample size to achieve test validity?
Are the results of my test Conclusive?
(i.e., can I “Accept” or “Reject” my hypothesis based on the data I have collected to-date?)

Examples: Using MECLABS Test Protocol – Validity Tool



Key Principles

1. The four basic steps to establish sample size validity

STEP 1: Estimate the success rate, sample size, and amount of difference

STEP 2: Estimate the duration of the test

STEP 3: Gather and record the actual data

STEP 4: Determine if the data collected is sufficient to end the test

MECLABS Test Protocol – Validity Tool

Example Test Scenario

Example Scenario

- You believe your primary offer page could convert better.
- You perform a diagnostic analysis on the existing offer page and arrive at a Hypothesis that form length and layout are contributing to excessive Friction.
- You design and develop **3 treatments** that will serve to test your hypothesis.
- Your recent **conversion rate** has been around **10%**.
- You want to be at least **95%** confident that the results of your test will **detect a difference** in conversion rate of **1% or more (10% relative)**.
- The average **traffic level** for this page has been about **500 visitors per day**.

How long will it take to achieve a valid sample?

Example Test Scenario

Minimum sample size estimation

Metric	Qty
Minimum size difference you wish to detect:	1.00%
Expected average success rate in the test: (Use 50% to be most conservative)	10.0%
Statistical significance level: (Remember, this is the probability of accidentally concluding there is a difference, due to sampling error, when really there is no difference.)	5%
Number of standard deviations for this significance level:	1.95996
Required standard deviation:	0.005102
Expected 2p(1-p)	0.18
Number of observations per treatment:	6,915

Test duration projection

Metric	Qty	
Arrival rate (arrivals/day)	500	arrivals/day
Number of treatments	4	
#Sample success observations / day	50	Successes/day
Estimated minimum acceptable sample size (#)	27,659	Arrivals
Estimated minimum test duration (days)	55	days

Inputs

- **1.0%** is the minimum difference you wish to detect
- **10%** is your expected average conversion rate
- Since you decided on a 95% confidence level, you enter its complement, a **5%** statistical significance level
- Based on historical data, traffic level for the page is about **500** visitors per day
- Since you are testing three experimental treatment against the Control page, the total number of treatments is **4**

Output

- The estimated minimum test duration is calculated to be **55** days

Lets say you can only afford to wait 21 days. What are your options?

Example Test Scenario (continued)

What are my options?

- Traffic? (X)
- # Experimental treatments? Try 2... Try 1.
- Confidence Level (1- α)? Try 90%. ($\alpha=10\%$)
- Bigger size difference? Try 20% instead of 10%.

Now how long is it expected to take?

Example Test Scenario (continued)

Test gets underway.

15 days later...

- Control: Visits: 3,728 Conversions: 371
 - Treatment: Visits: 3,781 Conversions: 403
- Enough data? Conclusive?

21 days later...

- Control: Visits: 5,971 Conversions: 587
 - Treatment: Visits: 5,948 Conversions: 638
- Enough data? Conclusive?

25 days later...

- Control: Visits: 6,488 Conversions: 657
 - Treatment: Visits: 6,496 Conversions: 711
- Enough data? Conclusive?

Validity for Decision Makers

Key Principles



Essential Definitions

Sample Distortion Effect – the effect on a test outcome caused by failing to collect a sufficient number of observations

Statistical Level of Confidence – the statistical probability that there really is a performance difference between the control and experimental treatments based upon the samples collected to date



Key Principles

1. There are essentially three factors that determine an outcome's statistical certainty:
 - **The Amount of Difference** – The magnitude of the difference in the primary metric between the Control and an Experimental treatment. (usually expressed as a percentage value)
 - **The Success Rate** – The number of experimental trials in which the outcome was favorable with respect to the primary success metric divided by the total samples received. (usually expressed as a percentage value)
 - **The Samples Received** – The total number of experimental trials used to compute the measured sample success rate. (e.g., number of visitors to a landing page, number of email messages sent).
2. An outcome's true statistical certainty is then compared with the desired level of certainty (chosen prior to experimentation).



Key Principles

1. The four basic steps to establish sample size validity

STEP 1: Estimate the success rate, sample size, and amount of difference

STEP 2: Estimate the duration of the test

STEP 3: Gather and record the actual data

STEP 4: Determine if the data collected is sufficient to end the test

Validity for Decision Makers

Thank You

Practical Example: MECLABS Test Protocol – Detail 1

ID: TP2047_Follow Up Radical Template

1. Question

$$u = 2q + t + m + 2v + i \text{ ©}$$

1.1) Primary Research Question	Which email design yields the highest click through rate?
1.2) Secondary Research Question	Which email design yields the highest conversion rate?

2. Treatments (These Variables / Values will determine the test treatments)

2.1) What is the variable?	Removal of top nav and improvements to layout, simplified approach
2.2) What are the values?	Control
	Treatment 1- w/2 side offers
	Treatment 2 - w/color background
	Treatment 3 - w/white background
	Treatment 4 - w/white background and side nav
	Treatment 5 - w/white background and colored nav
	Treatment 6- w/ white background & more products links

3. Metrics

3.1) What must we measure in order to determine the best performing value?

Metrics	Notes
Click-through Rate	
Conversion Rate	

4. Validation (Refer to separate Validation worksheet)

Note: Do not enter data in sections 4.2-4.4. They are linked to the separate Validation worksheet.

Practical Example: MECLABS Test Protocol – Detail 2

4. Validation (Refer to separate Validation worksheet)

Note: Do not enter data in sections 4.2-4.4. They are linked to the separate Validation worksheet.

4.1) Have you considered the impact of these validation threats?

Validity Threat Types	Considered (Y/N)	Describe Possible Threats For This Test
History effect	Yes	Product category is highly seasonal; especially with regard to holidays, et al. Normalizing using priors.
Instrumentation effect	Yes	3 text CTA's in each treatment share a common CTID -- Aggregating for CTR sub-analysis.
Selection effect	Yes	Using an asymmetrical send profile. Validation methodology is already profile-aware --> not a factor.
Sample Distortion effect	Yes	Splitter has been verified in prior tests using Control-Control method. Passed.
Other Validity effect	Yes	No known.

4.2) What is the sample size needed to ensure that the results will be predictive?

Projected minimum acceptable sample size:	38223	From validation worksheet - Primary Research Metric.
---	-------	--

4.3) How long should it take to obtain the minimum sufficient sample size?

Estimated Test Duration	0	Days	From validation worksheet - Primary Research Metric.
-------------------------	---	------	--

4.4) How will you verify during the test that you have reached a sufficient number of sample observations?

Sample Size: (Validation worksheet)	Notes
Events	976517
Successes	4955
Sample Size Valid?:	NO (based upon the corresponding Validation worksheet - Primary Research Metric.)

5. Metrics and Results

5.1) What are the test results?

Treatments	CTR Total	Relative Difference %	Validation Status	CR Total	Relative Difference %	Validation Status	Notes
Control	0.51%		YES	0.11%		YES	For the test cluster to-date, the Control is outperforming all treatments at the aggregate level. A quick subsetting dive shows signif. Perf. Asymmetries. Now looking for core-factors. (BK)
Treatment 1- w/2 side offers	0.44%	-14.22%	YES	0.09%	-17.96%	YES	
Treatment 2 - w/color background	0.42%	-16.80%	YES	0.09%	-15.60%	YES	
Treatment 3 - w/white background	0.46%	-10.55%	YES	0.11%	-3.81%	NO	
Treatment 4 - w/white background and side nav	0.50%	-1.31%	NO	0.10%	-7.14%	NO	
Treatment 5 - w/white background and colored nav	0.48%	-5.08%	NO	0.10%	-6.10%	NO	
Treatment 6- w/ white background & more	0.49%	-3.75%	NO	0.10%	-9.68%	NO	

6. Interpretation

Practical Example: MECLABS Test Protocol – Detail 3

5. Metrics and Results

5.1) What are the test results?

Treatments	CTR Total	Relative Difference %	Validation Status	CR Total	Relative Difference %	Validation Status	Notes
Control	0.51%		YES	0.11%		YES	For the test cluster to-date, the Control is outperforming all treatments at the aggregate level. A quick subsetting dive shows signif. Perf. Asymmetries. Now looking for corel-factors. (BK)
Treatment 1- w/2 side offers	0.44%	-14.22%	YES	0.09%	-17.96%	YES	
Treatment 2 - w/color background	0.42%	-16.80%	YES	0.09%	-15.60%	YES	
Treatment 3 - w/white background	0.46%	-10.55%	YES	0.11%	-3.81%	NO	
Treatment 4 - w/white background and side nav	0.50%	-1.31%	NO	0.10%	-7.14%	NO	
Treatment 5 - w/white background and colored nav	0.48%	-5.08%	NO	0.10%	-6.10%	NO	
Treatment 6- w/ white background & more	0.49%	-3.75%	NO	0.10%	-9.68%	NO	

6. Interpretation

What insights can we gain from this test?

6.1) Objective Interpretation based on Test Data	For the test cluster to-date, the Control is outperforming all treatments at the aggregate level. There is significant variance depending upon whether Product or Collection page is the landing page. Treatments 3 and 4 outperformed Control by 12% - 22% (+3% MoE) in 3 of 4 sends when Category page is the landing page. [ref. pivot & correlation tabs]. (BK)
6.2) Expert Speculation on Possible causes of the Test Outcome	A quick subsetting dive shows significant performance asymmetries, with periods during which treatments (esp. Treatment 3) substantially outperforms Control. Now looking for potential correlative factors. First set includes whether landing page is Product or Collection page (link study). (BK)
6.3) Suggested Follow-Up Tests	