

## How to fail: Optimize learnings from every test that you run





Ronnie Cheung Senior Strategy Consultant Optimizely



Laura Shelton Optimisation Manager Channel 4







Charlotte Golding Head of Optimisation Virgin Media

# Let's talk about failures in testing



# Can you guess what is the average win rate across all experiments ran by Optimizely? 10%





### **Common "Failures"**



### Statistically Significant, but Negative



### f Inconclusive



### **Flawed Test Design**



### Why does it suck so much?



### The feeling of going back to square one



Wasted time and resources



You weren't able to prove someone wrong





### How people read statistically significant negative test results

How people SHOULD read statistically significant negative test results



### Every negative result is still a win

#### Validated a feature / an experience

Instead of going straight build, you were able to identify that this variation will not generate an positive improvement

#### Learned about your customers / users

An underperforming variation tells you something about your users' preferences, which you didn't know about

#### **Concluded results with confidence**

Reaching statistical significance is typically a sign of a well designed, sufficiently powered test.



### What were the interesting things you learned about your users through experiments with stat. negative results?



#### SUCCESSFUL



- Exclusive rewards and competitions
- And other delicious All 4 Exclusives

#### Opt in and register

#### Register

By continuing, you confirm your date of birth is correct, you confirm you have read our Privacy Policy, and you accept the Terms and Conditions and Channel 4's use of cookies as set out in our Cookies Policy.



By continuing, you confirm your date of birth is correct, you confirm you have read our Privacy Policy, and you accept the Terms and Conditions and Channel 4's use of cookies as set out in our Cookies Policy.

#### UNSUCCESSFUL



Don't miss new entertainment...



Opt in to hear about:

- Our top entertainment
- From MAFS and MiC to Gogglebox and Bake Off
- Brand-new and exclusive shows
- Competitions and rewards

#### Opt in and register

#### Register

By continuing, you confirm your date of birth is correct, you confirm you have read our Privacy Policy, and you accept the Terms and Conditions and Channel 4's use of cookies as set out in our Cookies Policy.

#### UNSUCCESSFUL

#### Register

Don't miss new entertainment...



- Our highest-rating entertainment
- Reality hits from the UK, US and Australia
- · Brand-new and exclusive shows
- Competitions and rewards

#### Opt in and register

#### Register

By continuing, you confirm your date of birth is correct, you confirm you have read our Privacy Policy, and you accept the Terms and Conditions and Channel 4's use of cookies as set out in our Cookies Policy.

# 4

#### . . . . . . . . . . . . . . .

What advice would you give if someone receives a statistically negative results, in terms of the actions to take?





### Not all inconclusive are the same

Statistical Significance 🕜	Confidence Interval 💿	Statistical Significance ⑦	Confidence Interval 🕜
 Baseline		 Baseline	
32% 43,034 visitors remaining	0	<1% >100,000 visitors remaining	0

Statistical Significance ⑦	Confidence Interval ⑦
Baseline	
88% 5,785 visitors remaining	0



### What does it mean?

0	Confidence Interval 🕐	Statistical Significance ⑦
iseline		Baseline
<1%	0	32%
aining		43,034 visitors remaining

#### Some impact being seen in treatment, low confidence

dence Interval () Statistical Significance (	Ð
Bas	eline
<u>,</u> ,	<1%
>100,000 visitors rema	iining

#### No impact being seen in treatment

0

Confidence Interval ⑦	Statistical Significance 🕜
	Baseline
0	88%
	5,785 visitors remaining

#### Impact being seen in treatment, but not statistically significant



### To run or to pause?



Continue running if: Haven't reached sample size and Visitors remaining is low

Statistical Significance 🕜
 Baseline
88%
5,785 visitors remaining

#### CX-2 Click (PRIMARY METRIC)

Day 9

**Day 20** 

Unique conversions per visitor for CX-2 Click event

	Unique Conversions ⑦ Visitors	Conversion Rate 🕐	Improvement ⑦	Confidence Interval	Statistical Significance
Original	<b>34</b> 9,368	0.36%			 Baseline
Variation #1	<b>297</b> 9,444	3.14%	+766.5%	0	29% 328 visitors remaining

Edit

Variation #1	<b>1,056</b> 19,130	5.52%	+1,333%	0	81% 83 visitors remaining
Original	<b>74</b> 19,207	0.39%			 Baseline
	Unique Conversions ⑦ Visitors	Conversion Rate	Improvement ?	Confidence ?	Statistical ⑦
Unique conversions p	er visitor for CX-2 Click e	event			Edit
CX-2 Click PRIM	ARY METRIC				

#### "Variation is showing +1,333%, why is it still not reaching stat sig?"

Your baseline conversion matters. 74 conversions mean there is not enough sample data from the control to validate your results.

CX-2 CIICK PRIMARY	METRIC				Edit
Unique conversions per v	isitor for CX-2 Click e	event			
	Unique Conversions ⑦ Visitors	Conversion Rate ?	Improvement ?	Confidence ⑦ Interval	Statistical ⑦
Original	74	0.39%			
e lightai	19,207			٦	Baseline
Variation #1	<b>1,056</b> 19,130	5.52%	+1,333%	0	81% 83 visitors remaining

### That being said, statistical significance is about risk tolerance.

If you are confident that the variation performs better, then there's no reason why you shouldn't implement this.

CX-2 Click PRIMARY	METRIC				Edit	
Unique conversions per	visitor for CX-2 Click e	event				
	Unique Conversions ? Visitors	Conversion Rate ?	Improvement ?	Confidence ⑦	Statistical ⑦ Significance	
Original	74	0.39%				
onginar	19,207				Baseline	
Variation #1	<b>1,056</b> 19,130	5.52%	+1,333%	<u> </u>	81% 83 visitors remaining	
					and the second se	

### Think contextually about your user journey.

You may have a high **risk tolerance** towards a test that's running on the final checkout page, versus a test that's running on a confirmation page.



### Do you have any examples where you saw the opportunities from an inconclusive experiment?



#### INCONCLUSIVE



A Place in the Sun

SUCCESSFUL

Hey TEST!

It's A Sin

Gogglebox

The Inbetweeners

24 Hours in A&E









Friday Night Dinner

Submit shows

Married at First Sight Australia











The Great British Bake Off

Please pick 3 or more shows you like for better recommendations

Derry Girls



Hollyoaks

# What advice would you give if a team keeps getting inconclusive results?





### Mistakes we see again and again...



#### No qual/quant. backing

This prevents you from being able to identify and address your user problems

#### Testing that home buyers in consideration phase prefer "softer copy" and removed the home loan interest rate.

#### ✓ Enquire

Unique conversions per visitor for Enquire event					
	Unique Conversions Visitors	Conversion Rate 🕐	Improvement 🕜	Confidence Interval 🕐	Statistical ③ Significance
Original	<b>75</b> 2,521	2.98%			 Baseline
Variation #1	<b>28</b> 2,554	1.10%	-63.15%	0	>99%

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0

#### Not assigning the right metrics

Your primary metric should also be the direct action that a user can take.

### This experiment only included one metric in the set up.

The variation was unsuccessful, but we couldn't understand why without the ability to evaluate other behavioural metrics.

	219 406	0 5 2 2			
	Total Conversions Visitors	Conversions per ⑦ Visitor	Improvement ?	Confidence Interval 🕐	Statistical Significance
Total conversions per visit	or for home_sign_in_click ev	rent			
$\checkmark$ home_sign_in_click					
Primary Metric					
nonio_olgi_in_toxtiloid			48.14%		0.423
home sign in textfield			381,421		-20.44%
home_sign_in_button			410,832 51.86%		0.532
Variations			Visitors		home_sign_in_click
Summary					

0.423

-20.44%

Baseline

>99%

home\_sign\_in\_button

home\_sign\_in\_textfield

410,832

161,389

381,421

#### Not doing the pre-test calculations

You start launching experiment(s), only to realise none of them are going to reach stat sig.

### This CTA colour change test has been running for almost 10 weeks.

#### Visit Page: JP - Contact Sales PRIMARY METRIC

Unique conversions per visitor for Visit Page: JP - Contact Sales event

	Unique Conversions Visitors	Conversion Rate	Improvement	Confidence Interval	Statistical Significance	
Control	<b>52</b> 33,510	0.16%			 Baseline	
Variation #1	<b>82</b> 33,469	0.25%	+57.89%	0	19% 29,283 visitors remaining	
View Graph 🗸						

Edit

# What were some of the painful lessons you learned in the past about test designs?





### How do you make sure all teams follow best practices when setting up experiments?





### What have we learned?

### **Every (stat. sig) negative result is still a win.**

### Not all inconclusive are the same.

## A proper test design can improve your win rate and conclusive rate.





Laura Shelton Optimisation Manager Channel 4





Charlotte Golding Head of Optimisation Virgin Media