

# Contents

---

<i>Preface – How to Read This Book</i>	<i>page xv</i>
<i>Acknowledgments</i>	<i>xvii</i>
PART I INTRODUCTORY TOPICS FOR EVERYONE	1
<b>1 Introduction and Motivation</b>	<b>3</b>
Online Controlled Experiments Terminology	5
Why Experiment? Correlations, Causality, and Trustworthiness	8
Necessary Ingredients for Running Useful Controlled Experiments	10
Tenets	11
Improvements over Time	14
Examples of Interesting Online Controlled Experiments	16
Strategy, Tactics, and Their Relationship to Experiments	20
Additional Reading	24
<b>2 Running and Analyzing Experiments: An End-to-End Example</b>	<b>26</b>
Setting up the Example	26
Hypothesis Testing: Establishing Statistical Significance	29
Designing the Experiment	32
Running the Experiment and Getting Data	34
Interpreting the Results	34
From Results to Decisions	36
<b>3 Twyman’s Law and Experimentation Trustworthiness</b>	<b>39</b>
Misinterpretation of the Statistical Results	40
Confidence Intervals	43
Threats to Internal Validity	43
Threats to External Validity	48

	Segment Differences	52
	Simpson's Paradox	55
	Encourage Healthy Skepticism	57
<b>4</b>	<b>Experimentation Platform and Culture</b>	58
	Experimentation Maturity Models	58
	Infrastructure and Tools	66
	<b>PART II SELECTED TOPICS FOR EVERYONE</b>	79
<b>5</b>	<b>Speed Matters: An End-to-End Case Study</b>	81
	Key Assumption: Local Linear Approximation	83
	How to Measure Website Performance	84
	The Slowdown Experiment Design	86
	Impact of Different Page Elements Differs	87
	Extreme Results	89
<b>6</b>	<b>Organizational Metrics</b>	90
	Metrics Taxonomy	90
	Formulating Metrics: Principles and Techniques	94
	Evaluating Metrics	96
	Evolving Metrics	97
	Additional Resources	98
	SIDEBAR: Guardrail Metrics	98
	SIDEBAR: Gameability	100
<b>7</b>	<b>Metrics for Experimentation and the Overall Evaluation Criterion</b>	102
	From Business Metrics to Metrics Appropriate for Experimentation	102
	Combining Key Metrics into an OEC	104
	Example: OEC for E-mail at Amazon	106
	Example: OEC for Bing's Search Engine	108
	Goodhart's Law, Campbell's Law, and the Lucas Critique	109
<b>8</b>	<b>Institutional Memory and Meta-Analysis</b>	111
	What Is Institutional Memory?	111
	Why Is Institutional Memory Useful?	112
<b>9</b>	<b>Ethics in Controlled Experiments</b>	116
	Background	116
	Data Collection	121
	Culture and Processes	122
	SIDEBAR: User Identifiers	123

PART III COMPLEMENTARY AND ALTERNATIVE TECHNIQUES TO CONTROLLED EXPERIMENTS	125
<b>10 Complementary Techniques</b>	127
The Space of Complementary Techniques	127
Logs-based Analysis	128
Human Evaluation	130
User Experience Research (UER)	131
Focus Groups	132
Surveys	132
External Data	133
Putting It All Together	135
<b>11 Observational Causal Studies</b>	137
When Controlled Experiments Are Not Possible	137
Designs for Observational Causal Studies	139
Pitfalls	144
SIDEBAR: Refuted Observational Causal Studies	147
PART IV ADVANCED TOPICS FOR BUILDING AN EXPERIMENTATION PLATFORM	151
<b>12 Client-Side Experiments</b>	153
Differences between Server and Client Side	153
Implications for Experiments	156
Conclusions	161
<b>13 Instrumentation</b>	162
Client-Side vs. Server-Side Instrumentation	162
Processing Logs from Multiple Sources	164
Culture of Instrumentation	165
<b>14 Choosing a Randomization Unit</b>	166
Randomization Unit and Analysis Unit	168
User-level Randomization	169
<b>15 Ramping Experiment Exposure: Trading Off Speed, Quality, and Risk</b>	171
What Is Ramping?	171
SQR Ramping Framework	172
Four Ramp Phases	173
Post Final Ramp	176

<b>16</b>	<b>Scaling Experiment Analyses</b>	177
	Data Processing	177
	Data Computation	178
	Results Summary and Visualization	180
	<b>PART V ADVANCED TOPICS FOR ANALYZING EXPERIMENTS</b>	183
<b>17</b>	<b>The Statistics behind Online Controlled Experiments</b>	185
	Two-Sample t-Test	185
	p-Value and Confidence Interval	186
	Normality Assumption	187
	Type I/II Errors and Power	189
	Bias	191
	Multiple Testing	191
	Fisher's Meta-analysis	192
<b>18</b>	<b>Variance Estimation and Improved Sensitivity: Pitfalls and Solutions</b>	193
	Common Pitfalls	193
	Improving Sensitivity	196
	Variance of Other Statistics	198
<b>19</b>	<b>The A/A Test</b>	200
	Why A/A Tests?	200
	How to Run A/A Tests	205
	When the A/A Test Fails	207
<b>20</b>	<b>Triggering for Improved Sensitivity</b>	209
	Examples of Triggering	209
	A Numerical Example (Kohavi, Longbotham et al. 2009)	212
	Optimal and Conservative Triggering	213
	Overall Treatment Effect	214
	Trustworthy Triggering	215
	Common Pitfalls	216
	Open Questions	217
<b>21</b>	<b>Sample Ratio Mismatch and Other Trust-Related Guardrail Metrics</b>	219
	Sample Ratio Mismatch	219
	Debugging SRMs	222

<b>22</b>	<b>Leakage and Interference between Variants</b>	226
	Examples	227
	Some Practical Solutions	230
	Detecting and Monitoring Interference	234
<b>23</b>	<b>Measuring Long-Term Treatment Effects</b>	235
	What Are Long-Term Effects?	235
	Reasons the Treatment Effect May Differ between Short-Term and Long-Term	236
	Why Measure Long-Term Effects?	238
	Long-Running Experiments	239
	Alternative Methods for Long-Running Experiments	241
	<i>References</i>	246
	<i>Index</i>	266