

## Increasing Statistical Power of A/B Tests: High Likelihood of Detecting True Affects

Sneha Dingre\*

Data Analyst/ Modeler, Miami, FL, USA

**Citation:** Dingre S. Increasing Statistical Power of A/B Tests: High Likelihood of Detecting True Affects. *J Artif Intell Mach Learn & Data Sci* 2022, 1(1), 68-70. DOI: doi.org/10.51219/JAIMLD/sneha-dingre/33

**Received:** August 3, 2022; **Accepted:** August 11, 2022; **Published:** August 13, 2022

\*Corresponding author: Sneha Dingre, Data Analyst/ Modeler, Miami, FL, USA

**Copyright:** © 2022 Dingre S., This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

### ABSTRACT

This paper provides a comprehensive overview of A/B testing in experimental design. It emphasizes the role of statistical analysis in deriving meaningful insights. The paper explores strategies to enhance the statistical power of A/B tests, including estimating required sample sizes, increasing sample sizes, selecting appropriate significance levels, optimizing test durations, employing pre-post design strategies, and choosing relevant outcome metrics.

**Keywords:** A/B testing, statistical significance, experimental, sample size

### 1. Introduction

A/B testing, or split testing, is a powerful method used in experimental design to assess the impact of changes or interventions on a particular outcome, often applied in fields such as marketing, product development, and user experience optimization. This technique involves dividing a sample population into two groups, A and B, exposing each group to different variations (such as different webpage layouts, marketing messages, or product features), and then comparing their responses. The success of A/B testing relies heavily on statistical analysis to derive meaningful insights from the collected data. Statistics play a crucial role in evaluating the likelihood of success by providing a systematic framework to interpret results, distinguish between random variations and genuine effects, and quantify the level of confidence in observed outcomes.

Through statistical methods, A/B testing allows practitioners to make data-driven decisions with a level of certainty, mitigating the influence of chance. Techniques like hypothesis testing and confidence intervals help assess whether observed differences between groups are statistically significant, indicating a reliable impact rather than mere randomness. In this way, statistics provides a rigorous foundation for understanding the efficacy of interventions, enabling organizations to optimize strategies

based on evidence rather than intuition in the dynamic landscape of experimentation and decision-making.

#### 2.1. Improving the statistical power of A/B testing

Improving the statistical power of A/B tests is crucial for increasing the likelihood of detecting true effects and avoiding false negatives. Statistical power is the probability that a test will correctly reject a false null hypothesis. Below are several techniques to enhance the statistical power of A/B tests.

##### 2.1.1. Estimating the required sample size using power analysis:

Before conducting the A/B test, perform a power analysis to estimate the required sample size for a given effect size, significance level, and desired power. This ensures that the test has adequate power to detect meaningful effects. Conducting an A/B test with a power analysis involves a systematic process to estimate the required sample size, ensuring the experiment has sufficient power to detect meaningful effects. To select an appropriate statistical test based on the experimental design and estimate the variability in the data, typically, standard deviation is used. statistical softwares can also be employed to perform a power analysis, utilizing formulas relevant to the chosen test. It is crucial to evaluate the results of the power analysis, ensuring the estimated sample size is practical within the experiment's constraints. If the calculated size is impractical, the effect size, significance level, or power need to be adjusted to find a

reasonable balance. Subsequently, A/B test is conducted with the determined sample size, ensuring random assignment of participants to control and treatment groups. Throughout the test, data is collected according to the specified experimental design, and then statistical analysis is performed on the collected data using the chosen test. Subsequently, effect sizes, confidence intervals, and p-values are calculated to assess the significance of observed differences. Interpret the findings in the context of the initial hypotheses, considering whether the effect size is practically significant. This systematic approach ensures that researchers not only estimate an appropriate sample size but also implement a well-designed A/B test that has sufficient power to detect meaningful effects. Adjusting parameters and carefully interpreting results contribute to the overall reliability and validity of the experiment, supporting informed decision-making based on the observed outcomes<sup>1</sup>. discusses sample size determination and emphasizes the importance of continuing data collection until the desired statistical information is achieved to ensure more accurate estimation of treatment effects in group sequential trials.

### 2.1.2. Increasing the sample size to boost statistical power:

One of the most effective ways to boost statistical power is to increase the sample size. A larger sample size reduces variability and provides a more accurate representation of the population.

Increasing the sample size in A/B testing is crucial for enhancing statistical power and the reliability of results. Several methods can be employed to achieve this. Firstly, extending the test duration allows for a larger dataset by capturing more user interactions over time. Secondly, utilizing historical data supplements current data, but relevance and representativeness are crucial considerations. Thirdly, boosting traffic or user participation through marketing efforts or promotions can increase the overall sample size. Segmenting and analyzing subgroups provide detailed insights, while leveraging multiple platforms or channels taps into diverse user groups. Employing a multi-arm bandit approach dynamically allocates more traffic to winning variations, optimizing resource use. Proxy metrics, which are more sensitive to changes, can be used when the primary metric requires a large sample size. Sequential testing allows periodic assessments and potential early test termination. Combining A/B test results from related tests, using a pre-post design, collaborating with external platforms, adjusting traffic allocation, and exploring partner collaborations are additional strategies<sup>2</sup>. explores adopting various sample size methodologies and how they can be leveraged to improve the statistical significance of the experiment. However, maintaining test integrity and considering potential biases are essential when implementing these methods.

**2.1.3. Choosing the right significance level:** While the standard significance level is often set at 0.05, consider adjusting it based on the goals and constraints of the experiment. Lowering the significance level (e.g., 0.01) increases the power but may lead to a higher chance of Type II errors. Choosing an appropriate significance level ( $\alpha$ ) in A/B testing is a critical decision that influences the trade-off between Type I and Type II errors. The commonly used significance level is 0.05, but several considerations can guide the choice: It is important to consider industry standards and common practices within the field, but also be aware of specific practices within your organization. Additionally, it is crucial to consider the consequences of both Type I and Type II errors. For instance, in fields like medical research, avoiding false positives might be more crucial than

avoiding false negatives<sup>3</sup>. As suggests, how the legal system leverages different significance levels based on the type of crime to make a decision. Scientific rigor is essential, particularly in research or publications, where adherence to standard practices is common. Take into account the preferences of stakeholders involved in decision-making, as their risk aversion can influence the chosen significance level. Multiple testing correction methods can be applied when conducting multiple tests simultaneously to adjust the significance level. Practical significance of observed effects should be considered, ensuring that statistical significance aligns with meaningful practical outcomes. Documenting the rationale behind the chosen significance level is crucial as it ensures transparency and reproducibility in the research or decision-making process.

**2.1.4. Optimizing the test duration:** It is crucial to ensure that the A/B test runs for a sufficient duration to collect a representative sample of data as short-duration tests may result in lower power, especially if the effect takes time to manifest. Ensuring an A/B test runs for a sufficient duration is critical for obtaining reliable and representative results, particularly when short-duration tests may compromise statistical power. A pre-test analysis can be conducted to estimate the expected duration and elements like effect size, baseline conversion rates, and user behavior variability can be factored in. Seasonality and trends in user behavior should be accounted for, adjusting the test duration to cover complete cycles<sup>4</sup>. discusses that as per Google analytical tool, it is advised to run at least for 2 weeks to gather the right amount of data. However, the duration time for test varies per tool. Monitoring cumulative effects and recognizing potential lagged effects is crucial, as some changes may only manifest after a certain period. Running pilot tests or phased rollouts allows observation of initial user reactions, aiding in estimating the time required for the effect to stabilize. Implementing sequential testing methodologies allows for periodic assessments and potential early stopping if significant results are achieved. Historical data can be utilized to understand typical durations of similar interventions and set a minimum run time for the A/B test based on various considerations. Statistical stability should be regularly monitored, ensuring the collected data reflects consistent patterns. By strategically considering these factors, researchers can ensure that A/B tests are conducted for an adequate duration, resulting in robust and reliable insights that support informed decision-making based on the outcomes of the test.

**2.1.5. Evaluate effects of intervention:** Employing a pre-post design strategy when conducting an A/B test is instrumental in mitigating individual variations and refining the accuracy of effect estimations. Executing a pre-post design within an A/B test involves adopting a methodical approach to amplify the precision of effect estimates by capturing baseline metrics before implementing any treatment or modification. Firstly, defining the A/B test objectives and selecting metrics pertinent to the desired goals is done. Baseline measurements can be gathered from the target population to establish the initial state, thereby accounting for individual differences and variability in the absence of the intervention. Once the baseline is established, introduce the treatment to the selected group (Treatment Group) while maintaining existing conditions for another group (Control Group). Post-treatment metrics should be monitored after the intervention has been in effect for a sufficient duration, focusing on user interactions, conversion rates, or other pertinent indicators. The analysis entails a comparative examination of

baseline and post-treatment metrics for both groups, enabling the identification of individual differences and providing insights into the intervention's impact. Statistical tests should be used to assess the significance of observed changes, controlling for confounding variables if applicable. When interpreting the results, consider both statistical and practical significance, and meticulously document the findings for future reference and reporting purposes. The pre-post design not only augments the precision of effect estimates but also facilitates a more holistic understanding of experimental outcomes. If necessary, refine the design and implementation based on insights gained to optimize the treatment or inform future experiments. Overall, a pre-post design fortifies the robustness of A/B tests by systematically addressing individual differences and offering a structured approach to evaluate the effects of intervention.

**2.1.6. Selecting relevant metrics:** Choosing relevant outcome metrics are more likely to capture the true impact of the intervention. It is important to use sensitive and relevant metrics increases the chances of detecting meaningful effects. Choosing outcome metrics for A/B testing involves selecting user-centric, sensitive, and relevant measures aligned with intervention objectives. Metrics should be prioritized with quick response to changes, avoiding vanity metrics. Precision can be opted for, considering lag effects, and ensuring alignment with business goals. User experience metrics and segmentation possibilities enhance insights. A holistic approach, combining leading and lagging indicators, provides a comprehensive view. Pilot testing helps assess metric sensitivity and relevance before the full-scale A/B test, increasing the chances of capturing meaningful intervention effects. Some of the metrics can be click through rate or number of sessions. The study in<sup>5</sup> aims to solve the problem of needing large sample sizes in online search A/B testing by using machine learning to identify sensitive metric combinations, resulting in significant improvements in sensitivity.

**2.1.7. Replicate experiments to ensure consistent results:** Replicating experiments under similar conditions helps in validating findings and obtaining consistent results across multiple experiments. This increases confidence in the observed effects. It is also important to consider segmenting users based on relevant characteristics and analyze results for each segment separately. This helps in identifying effects that may be specific to certain user groups.

### 3. Conclusion

The paper discusses the significance of robust statistical methods in A/B testing, providing a structured guide for practitioners across diverse domains. The emphasis on replicating experiments and segmenting user groups contributes to the overall credibility and applicability of A/B testing outcomes. By emphasizing techniques to enhance statistical power, optimize test parameters, and evaluate intervention effects, the paper empowers decision-makers to derive reliable insights from experimental data. Lastly, this paper discusses continual adaptation and learning from experimental outcomes of A/B tests. This is crucial for organizations to gain a competitive advantage, make crucial decisions based on concrete data, and optimize their strategies by identifying what resonates best with their audience.

### 4. References

1. C R. Mehta, A A. Tsiatis. Flexible sample size considerations using Information-Based interim monitoring. *Drug Information Journal*, 2001; 35: 1095-1112.
2. H. Kang. Sample size determination and power analysis using the G\*Power software. *Journal of Educational Evaluation for Health Professions*, 2021; 18: 17.
3. Kim, Jae. How to choose the level of significance: A pedagogical Note. *Munich Personal REPEC Archive*, 2015.
4. C. Pala. How long should you run an A/B test?. *NorthStar Analytics*, 2022.
5. E. Kharitonov, A. Druitsa, and P. Serdyukov. Learning Sensitive Combinations of A/B Test Metrics. *ACM Digital Library*, 2017; 40: 651-659.