



Marketing Incrementality Measurement: An In-Depth Analysis

By sig.ai Contact: info@sig.ai

02-05-2025

Introduction

Marketing incrementality measurement is a data-driven approach to determine the true causal impact of advertising and marketing efforts on outcomes like conversions, sales, or revenue. Unlike attribution models that often rely on correlation, incrementality testing uses controlled experiments to isolate the “lift” generated by marketing – in other words, the additional value that would **not** have occurred without the marketing intervention. By comparing a group exposed to ads (treatment group) against a comparable group that isn’t exposed (control group), marketers can measure the difference in outcomes and attribute that difference to the marketing activity. This report provides a structured, academic-style overview of incrementality measurement, covering methodologies, special techniques (ghost ads and PSA tests), statistical significance considerations, real-world case studies in e-commerce, and a practical implementation checklist. All discussions are supported by empirical research and real examples, with citations in APA style.

Incrementality Testing Methodologies

Incrementality testing is fundamentally rooted in experimental design. The goal is to set up a **randomized controlled trial (RCT)** in a marketing context: one group of users is subjected to the marketing treatment (ads or promotions), and another group is held out as a control. Below, we breakdown several common methodologies for incrementality testing, including their principles, advantages, and challenges.

A/B Testing (User-Level Randomized Controlled Trials)

A/B testing is the classic approach to incrementality measurement and is often considered the gold standard. In a marketing A/B test, eligible users are randomly split into two groups: Group A (treatment) and Group B (control). The treatment group is exposed to the marketing campaign (for example, seeing a particular ad or promotion), while the control group is not exposed. Because the assignment is random at the individual level, and the groups are drawn from the same audience pool, any systematic difference in outcomes (like conversion rate or sales)

between the groups can be attributed to the marketing intervention with high confidence. This approach mirrors clinical trials in medicine and ensures an “apples-to-apples” comparison between users who saw the ads and those who did not.

Advantages: Properly executed user-level A/B tests directly establish causality. Randomization ensures the two groups are statistically equivalent on observed and unobserved traits, so differences in outcomes reflect the effect of the ads, not underlying biases. A/B tests are versatile and can be applied to almost any digital marketing channel where individual users can be randomly assigned to see or not see ads.

Limitations: In some cases, A/B testing at the user level may be infeasible or even invalid. If users can influence each other (for example, one user’s treatment status affects another user’s exposure or behavior), it breaks the independence of the groups. Similarly, some channels (like certain offline media or walled gardens) do not support user-level randomization. Additionally, “holdout” control groups in A/B tests mean deliberately not marketing to a portion of potential customers, which can have an opportunity cost. Despite these challenges, A/B testing remains a foundational incrementality method whenever it can be implemented properly.

Holdout Testing (Intent-to-Treat Approach)

Holdout testing is a specific implementation of A/B testing widely used in advertising: a percentage of the target audience is “**held out**” from a campaign to serve as a control group. For example, an advertiser might show ads to 90% of the eligible audience, while 10% are withheld (suppressed) and receive no ads. After the campaign, the advertiser compares conversion or revenue metrics between the exposed group and the holdout group. Because the holdout users saw no ads, any **incremental** lift observed in the exposed group’s metrics (beyond the holdout’s baseline) can be attributed to the advertising.

Holdout tests follow an intent-to-treat philosophy: everyone designated in the treatment group is analyzed as part of that group, even if some of those users never actually saw an ad (due to under-delivery or frequency capping, etc.). Likewise, the control group is analyzed in full, including all users who would have been eligible to see ads had they not been held out. This simplicity can introduce some “noise” – for instance, users in the treatment group who were never reached by the ads dilute the measured lift. However, the intent-to-treat holdout method is easy to implement and interpret.

Advantages: Holdout experiments are straightforward to set up on most ad platforms (many platforms have built-in support to designate a control group that is not targeted). They are low-cost since the control group simply receives no ads (unlike some methods that show alternative ads to control). This method provides an unbiased estimate of incremental effect as long as randomization is done correctly. It’s particularly useful for ongoing **lift measurement** – e.g. keeping a always-on 10% holdout to continually gauge incremental ROI of campaigns.

Limitations: The main challenge is statistical noise. Not all users in the treatment group will necessarily see the ad (due to delivery constraints), and those unexposed users dilute the

measurable difference between treatment and control. This can make it harder to detect a statistically significant lift, especially if the true effect is small. In practice, the measured lift in a holdout test may understate the per-exposed-user effect because the analysis includes many “unexposed” users in the treatment group. Another limitation is **exposure logging** – because the control group sees nothing, we may not know which control users *would* have seen the ad. This complicates deeper analysis like per-exposure lift or understanding treatment on the treated (as opposed to intent-to-treat). More advanced methodologies (PSA and ghost ads, below) address this by ensuring control group “exposure” can be logged in some form.

Geo Experimentation (Geographical Holdouts or Matched Markets)

Geo experimentation involves using geographic regions as the unit of experiment rather than individual users. In a geo experiment, different geographic areas (cities, regions, or markets) are assigned to treatment or control groups. For example, a retailer might turn off a marketing channel in a set of test cities while continuing it as usual in a set of control cities, then compare sales between these geos. Wayfair, for instance, employs geo-experiments by dividing comparable regions into control vs. treatment to measure marketing impact when user-level experiments are not possible. Geo tests have been used to measure the incrementality of channels like TV, direct mail, or platform-wide campaigns where individual targeting is impractical.

To implement a geo experiment, one must define the geo units (e.g., ZIP codes, cities, states) and ensure they are comparable. Often geos are **matched** or paired based on historical performance and demographics, then one of each pair is randomly assigned to treatment and the other to control. This matching controls for external differences across regions. The experiment then runs by applying the marketing intervention in treatment geos and withholding it in control geos for a set period. Outcomes (e.g. total sales, new customers) are measured and the difference in trends between treatment and control geos indicates the lift, usually analyzed via difference-in-differences or time-series models.

Advantages: Geo experiments can measure incrementality at a macro level, capturing **halo effects** and offline impact that user-level digital experiments might miss. They are often the only viable method for channels like offline media (billboards, TV) or situations where individual randomization is not feasible. When designed well (using matched areas and proper statistical models), geo tests can provide robust causal measurement. They also avoid the issue of user-level cross-contamination (one user seeing both treatment and control, which can happen if users move or share devices – geo boundaries can minimize this).

Limitations: Geographic tests typically require larger scale and longer duration to get sufficient statistical power, since the number of experimental units (geos) is smaller than in user-level tests. Variability in local factors (weather, local events, competition) can introduce noise. Care must be taken to select similar geos for fair comparison; otherwise, results may be confounded by regional differences. Another challenge is **spillover** – if media from one geo inadvertently affects another (e.g., a TV signal bleeding into a neighboring region, or people traveling across geo boundaries), it can contaminate the control group. Geo experiments also measure