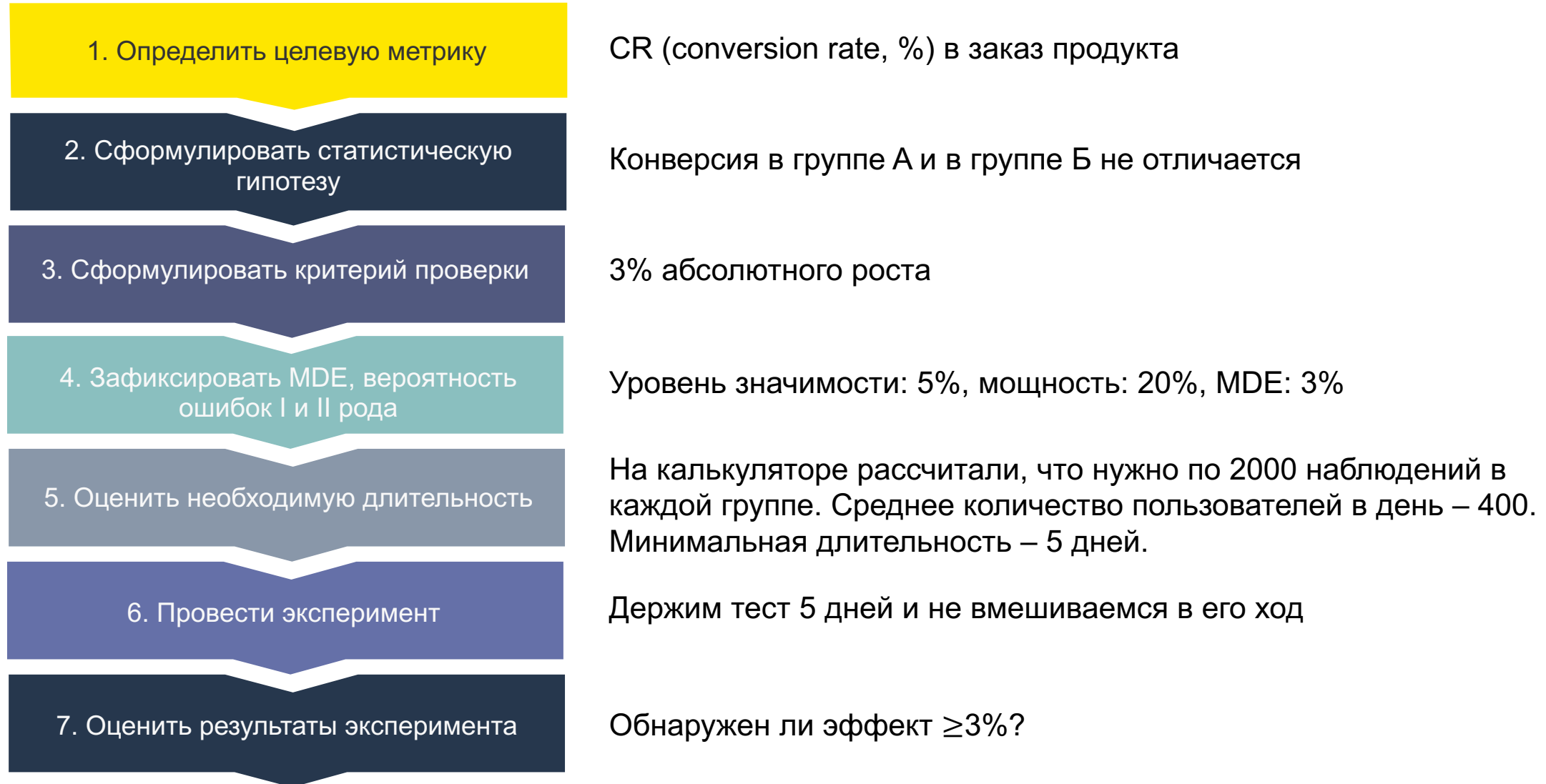


Дизайн экспериментов

Алгоритм проведения АБ-теста



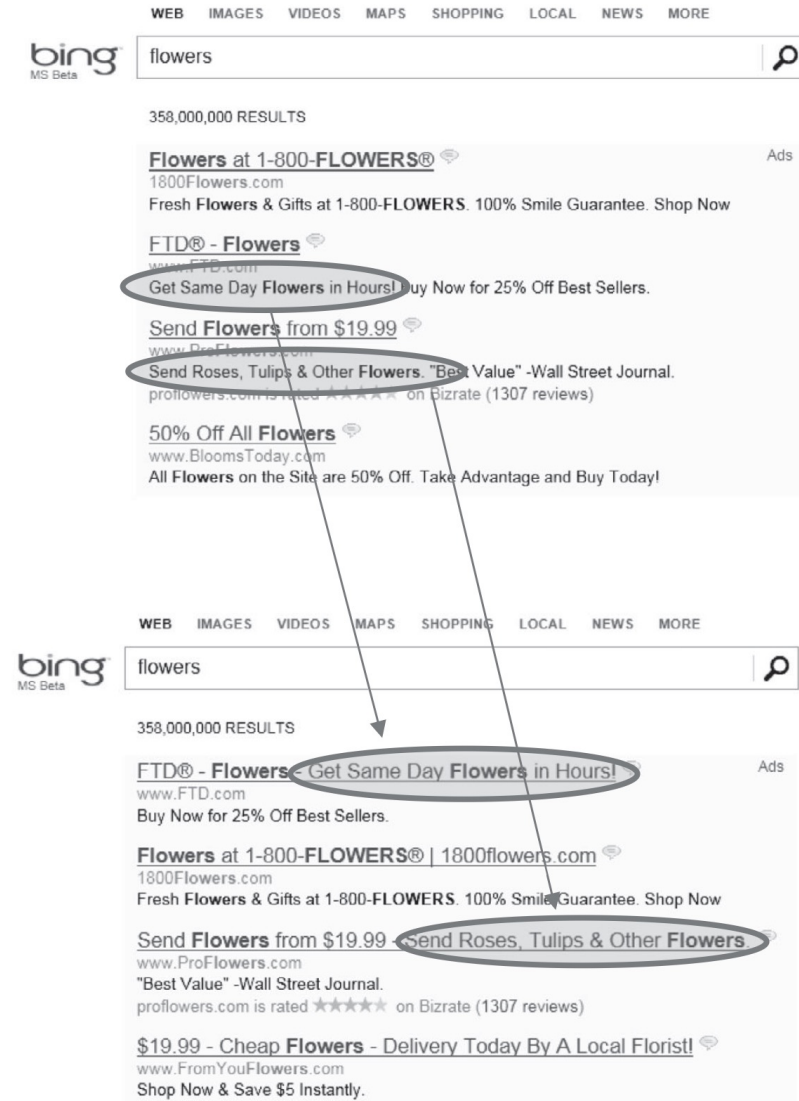
Примеры А/В тестов

Кейс Bing

В 2012 сотрудник поисковой системы Microsoft Bing предложил изменить способ отображения заголовков: дописывать первую строку под заголовком в сам заголовок

Через несколько часов сработала система алерта об очень высокой выручке.

Выручка выросла на 12%, что принесло более \$100М в год только на рынке США.



Примеры A/B тестов

Кейс Google:

41 shades of blue



В 2009 году Google протестировал 41 оттенок синего цвета на странице поиска.

Гугл не раскрыл точные цифры, но отметили, что это привело к существенному увеличению вовлеченности пользователей по их метрикам.

Microsoft Bing также показали, что цветовые настройки их поисковой системы увеличили годовую выручку на более чем \$10M в США ежегодно, а также улучшилась метрика Time-to-success

Примеры А/В тестов

Кейс Amazon:

Making an offer at Right time

В 2004 году Amazon разместила на главной странице предложение кредитной карты. Это было очень прибыльно, но имело очень низкий CTR

Команда переместила это предложение в корзину после добавления товара, а также начала показывать экономию с картой от каждой покупки.

В таком случае предложение показывалось в нужное время, и данный тест увеличил годовую прибыль Amazon на десятки миллионов долларов

You could save \$30 today with the Amazon Visa® Card:



Your current subtotal: \$32.20
Amazon Visa discount: - \$30.00
Your new subtotal: \$2.20

[Find out how](#)

Save \$30 off your first purchase, earn **3% rewards**, get a **0% APR***, and pay **no annual fee**.

Примеры A/B тестов

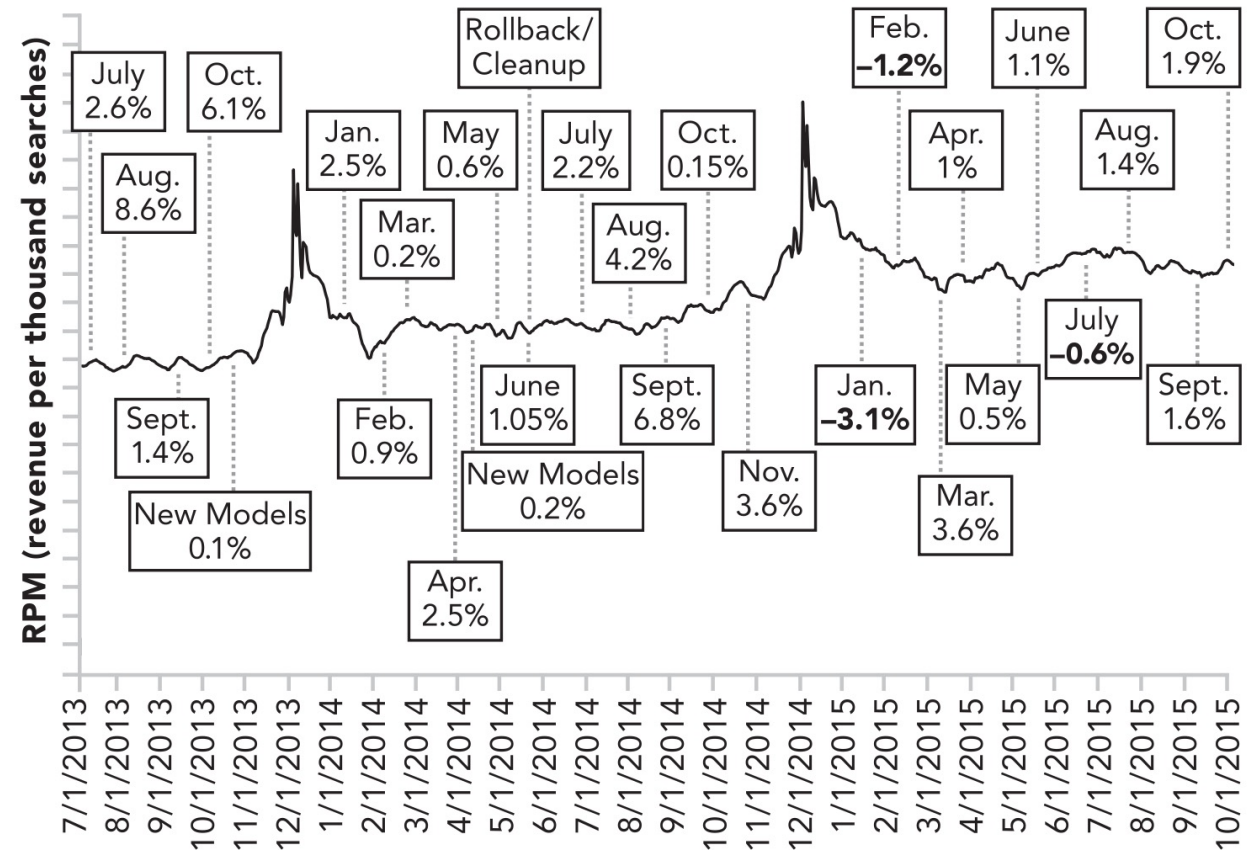
Кейс Bing:

Ads team

Команда по рекламе в Bing последовательно увеличивала доход на 15–25 % с 2013-2025, как показано на рисунке справа.

Часто в больших компаниях проводится куча тестов, и улучшения от них небольшие (а иногда даже отрицательные)

Но именно это и есть правильная культура тестирования гипотез в компании:
маленькие изменения приводят в итоге к большим результатам



(*) Numbers have been perturbed for obvious reasons

Как подобрать метрику?

Кейс Amazon: E-mail at Amazon

В Amazon была система рекомендаций книг на основе следующих правил:

- 1) Рекомендуем новую книгу автора, если пользователь уже покупал книги данного автора
- 2) Программа рекомендаций на основе историй покупок
- 3) «Перекрестное опыление»: очень специфичные точечные рекомендации, отправляемые после покупок специальных категорий

Какую метрику подобрать под данную систему рекомендаций?

Как подобрать метрику?

Кейс Amazon: E-mail at Amazon

В Амазон была система рекомендаций книг на основе следующих правил:

- 1) Рекомендуем новую книгу автора, если пользователь уже покупал книги данного автора
- 2) Программа рекомендаций на основе историй покупок
- 3) «Перекрестное опыление»: очень специфичные точечные рекомендации, отправляемые после покупок специальных категорий

Какую метрику подобрать под данную систему рекомендаций?

Amazon выбрал метрику: выручка с кликов по эмейлам

В чем проблема такой метрики?

Как подобрать метрику?

Кейс Amazon: E-mail at Amazon

В Амазон была система рекомендаций книг на основе следующих правил:

- 1) Рекомендуем новую книгу автора, если пользователь уже покупал книги данного автора
- 2) Программа рекомендаций на основе историй покупок
- 3) «Перекрестное опыление»: очень специфичные точечные рекомендации, отправляемые после покупок специальных категорий

Какую метрику подобрать под данную систему рекомендаций?

Amazon выбрал метрику: выручка с кликов по эмейлам

В чем проблема такой метрики?

Выручка монотонно увеличивается с увеличением количества писем и рекламных кампаний, что в конечном итоге приведет к спаму.

Что случится в долгосрочном итоге?

Как подобрать метрику?

Кейс Amazon: E-mail at Amazon

В Амазон была система рекомендаций книг на основе следующих правил:

- 1) Рекомендуем новую книгу автора, если пользователь уже покупал книги данного автора
- 2) Программа рекомендаций на основе историй покупок
- 3) «Перекрестное опыление»: очень специфичные точечные рекомендации, отправляемые после покупок специальных категорий

Какую метрику подобрать под данную систему рекомендаций?

Amazon выбрал метрику: выручка с кликов по эмейлам

В чем проблема такой метрики?

Выручка монотонно увеличивается с увеличением количества писем и рекламных кампаний, что в конечном итоге приведет к спаму.

Что случится в долгосрочном итоге?

Произойдет отток пользователей

Как подобрать метрику?

Кейс Amazon: E-mail at Amazon

В Амазон была система рекомендаций книг

Какую метрику подобрать под данную систему рекомендаций?

Amazon выбрал метрику: выручка с кликов по эмейлам

В чем проблема такой метрики?

Выручка монотонно увеличивается с увеличением количества писем и рекламных кампаний, что в конечном итоге приведет к спаму.

Что случится в долгосрочном итоге?

Произойдет отток пользователей

Первое решение:

Amazon сначала пытались решить эту проблему ограничениями на количество писем пользователю.

Но появилась проблема оптимизации:

какое письмо отправить пользователю каждые X дней, когда несколько программ рекомендаций хотят порекомендовать этому пользователю продукт?

Как подобрать метрику?

Кейс Amazon: E-mail at Amazon

В Амазон была система рекомендаций книг

Какую метрику подобрать под данную систему рекомендаций?

Amazon выбрал метрику: выручка с кликов по эмейлам

В чем проблема такой метрики?

Выручка монотонно увеличивается с увеличением количества писем и рекламных кампаний, что в конечном итоге приведет к спаму.

Оптимальная метрика

$$\text{Metric} = (\sum_i Rev_i - s * \text{unsubscribe_lifetime_loss}) / n$$

- i – счетчик для получателей e-mail
- s – количество отписавшихся пользователей
- $\text{unsubscribe_lifetime_loss}$ – ожидаемое потеря дохода из-за невозможности отправить ни одно письмо пользователю за все время
- n – количество пользователей

Как подобрать метрику?

Кейс Amazon: E-mail at Amazon

В Амазон была система рекомендаций книг

Какую метрику подобрать под данную систему рекомендаций?

Amazon выбрал метрику: выручка с кликов по эмейлам

Оптимальная метрика

$$\text{Metric} = (\sum_i Rev_i - s * \text{unsubscribe_lifetime_loss}) / n$$

- i – счетчик для получателей e-mail
- s – количество отписавшихся пользователей
- $\text{unsubscribe_lifetime_loss}$ – ожидаемое потеря дохода из-за невозможности отправить ни одно письмо пользователю за все время
- n – количество пользователей

После внедрения такой метрики более половина кампаний показали отрицательное значение метрики.

Что еще более интересно, осознание того, что отказ от рассылок несет такие большие потери, привело к созданию другой страницы отказа от рассылок, где по умолчанию была отписка только от данной «кампании», а не от всех электронных писем Amazon, что резко снизило количество отписок от всех писем

Метрики

Нечувствительные

- Цена акции компании
- Доля пользователей, продливших подписки в годовых сервисах
- Total выручка

Чувствительные

Плохие

- CTR новой кнопки на главной странице
- Время, проведенное на сайте

Хорошие

- CTR в таргетное действие (покупка)
- Time-to-success целевого действия

Как собрать одну метрику?

Вообще очень сложно.

Бизнес постоянно должен оптимизировать разные вещи (вовлеченность пользователя, деньги и тд)

Но есть такая темка:

1. Берем несколько метрик, которые хотим оптимизировать
2. Нормализуем их в заданный диапазон (например, 0-1)
3. Даем вес каждой метрике
4. Итоговая метрика – взвешенная сумма отнормированных показателей

$\text{Metric} = 0.5 * \text{CTR} + 0.5 * (1 / \text{Time-to-success}) \rightarrow \max$

Важно: если не можете сделать 1 метрику, то делайте их меньше 5!

Если нулевая гипотеза верна (нет изменений), то $P(\text{p-value} < 0.05) = 0.05$, т.е. O1P = 5% для 1 метрики

Если у вас k независимых метрик, то вероятность того, что хотя бы одна метрика будет иметь $\text{p-value} < 0.05$:

$$P(\text{p-value} < 0.05) = 1 - (1 - 0.05)^k$$

Для k = 5 вероятность того, что одна из метрик будет статистически значима = 23%

Для k=10 будет 40%

Типичные ошибки в A/B тестировании



Множественное тестирование



Преждевременная остановка теста (подглядывание)



Неверный выбор чувствительности критериев (MDE)



Оценка эффекта по зависимым выборкам



Оценка эффекта по части выборки



Незафиксированные уровень значимости и мощность

Множественное тестирование

Описание теста

Тест на увеличение конверсии в заказ продукта.

А группа: контрольная

В группа: чекбоксы с дополнительными продуктами

С группа: картинки с дополнительными продуктами

Как тестировали?

Сравнили конверсии

A vs B

A vs C

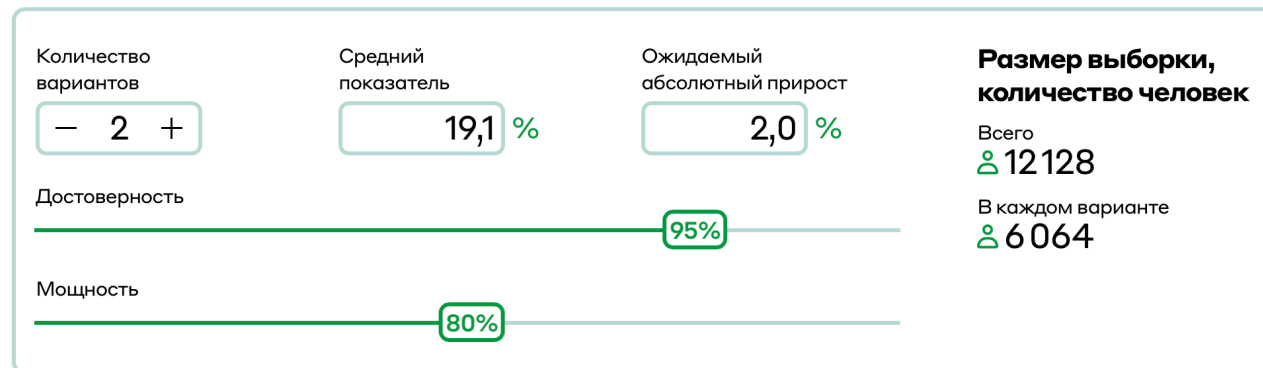
Важно: нельзя использовать одну и ту же контрольную выборку дважды. Снижается достоверность теста.

Решение:

1) в этом тесте нужно удвоить контрольную выборку. В подобных тестах можем делить трафик в соотношении:

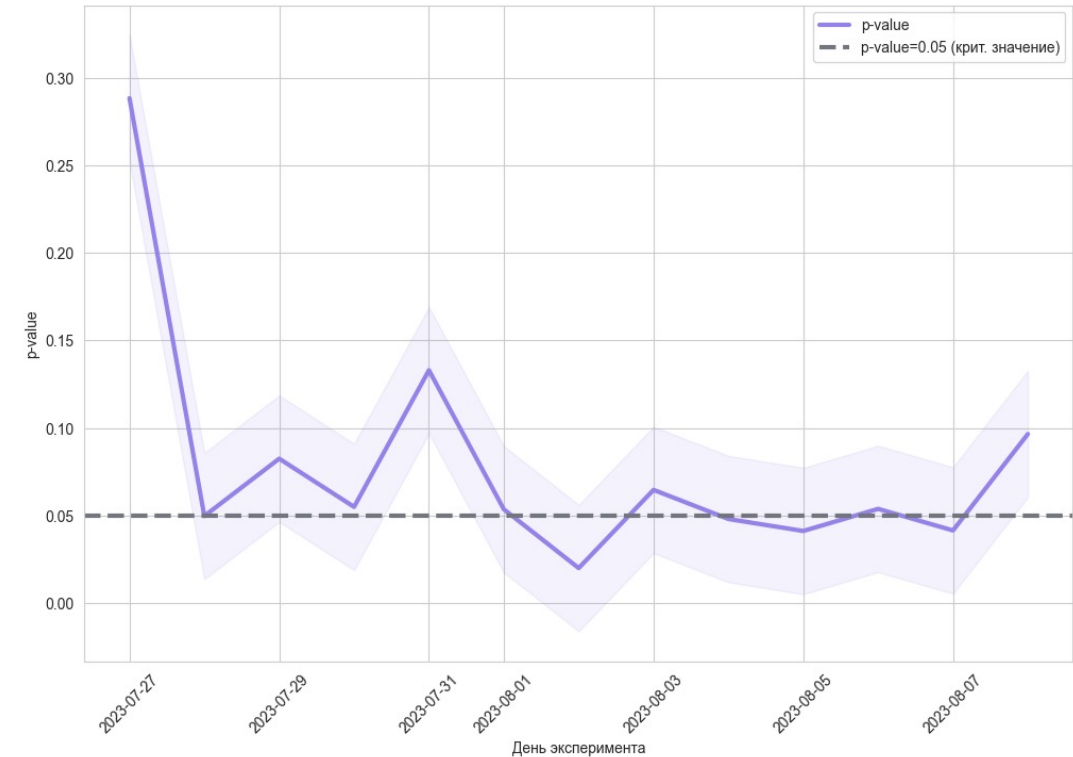
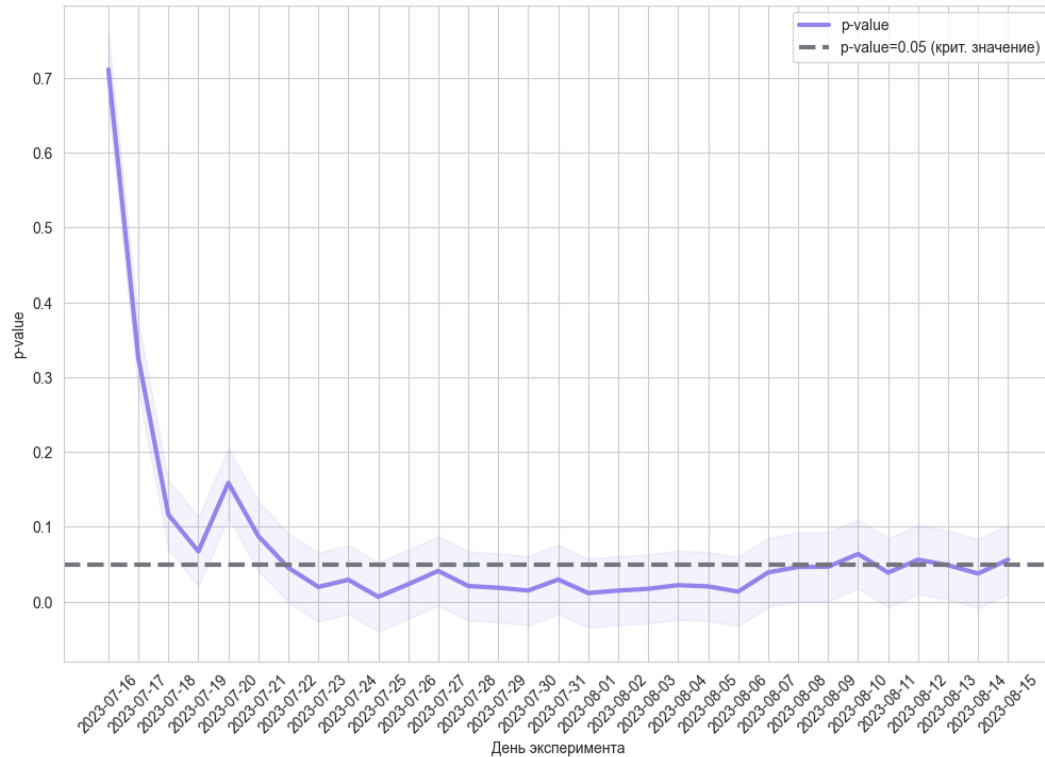
- 50% - контрольная группа
- 25% - чекбоксы
- 25% - баннеры

2) Множественное тестирование



Преждевременная остановка теста (подглядывание)

Если остановить тест раньше рассчитанного времени, можно получить ложную прокраску:

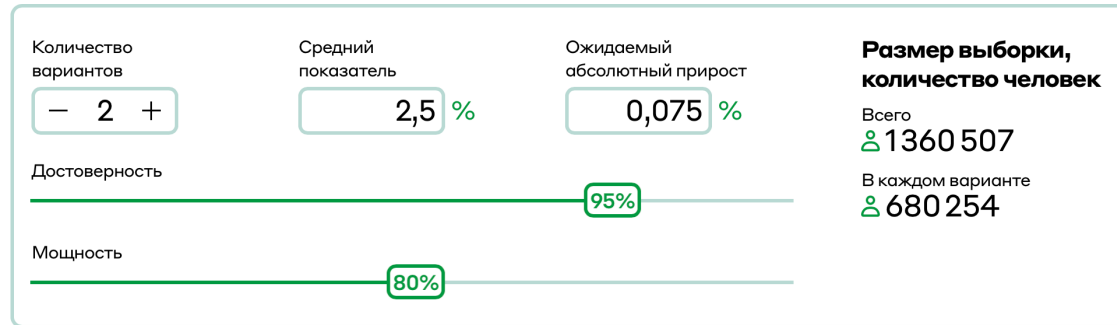


Решение:

- 1) Использовать последовательное тестирование, если есть вероятность, что тест придется остановить раньше или если рассчитанная длительность слишком большая
- 2) Ждать столько времени, сколько нужно для fixed-horizon тестирования

Выбор чувствительности критериев (MDE)

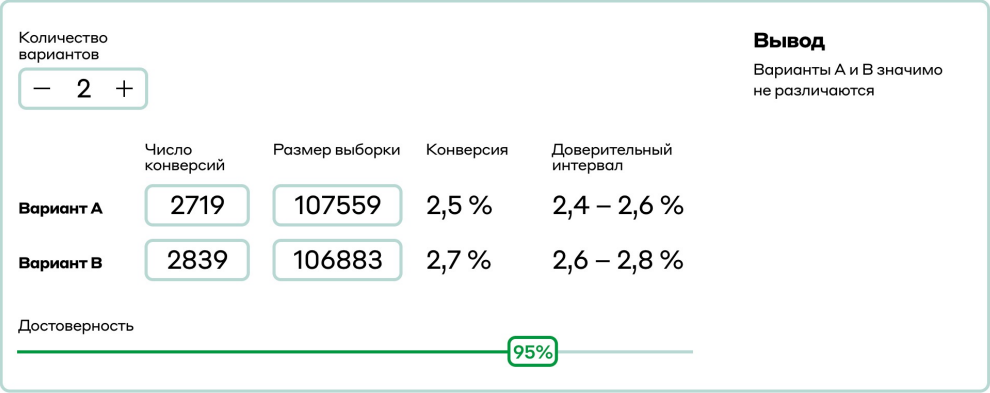
В тесте хотим детектить **относительный** рост конверсии на 3%: то есть при средней конверсии 2,5% нужно заметить **абсолютный** прирост 0,075%. Тогда необходимый размер выборки:



Это очень много. Как решить эту проблему:

- Выбирать больший MDE для тестов с низкой базовой конверсией. В этом случае мы не сможем уловить совсем небольшие изменения, но зато тест будет корректен.
- Ориентироваться на абсолютный прирост, чтобы понимать масштаб детектируемого эффекта.

Оценка эффекта по зависимым выборкам



Группировки	Метрики	Выберите цель
<input checked="" type="checkbox"/> Тип устройства	<input checked="" type="radio"/> Визиты	<input checked="" type="radio"/> Посетители
	<div>%</div>	<div>%</div>
<input type="checkbox"/> Итого и средние	107 559	≈ 106 883
<input checked="" type="checkbox"/> ПК	55,1 %	≈ 55,1 %
<input checked="" type="checkbox"/> Смартфоны	43,8 %	≈ 43,8 %
<input type="checkbox"/> Планшеты	1,1 %	+1.07x 1,18 %

Одному пользователю может принадлежать несколько визитов. Такие визиты будут зависимыми друг от друга. Тогда предпосылки нашего теста не соблюдаются, и результат нельзя интерпретировать. Самое простое решение – всегда использовать **поюзерные** метрики.

Оценка эффекта по сегменту выборки

- Проверка значимости разницы между контролем и экспериментом может проводиться только по достижении заранее рассчитанного размера выборки.
- Если тестировать на сегменте выборки – например, только на Мобайле – может возникнуть проблема подглядывания.

Решение:

- 1) Делать выводы только на основе всех данных (все устройства)
- 2) делать стратификацию по сегментам – тогда размер выборки не сократится, а сегменты с большим весом привнесут больший вклад в среднее по выборке.

Мобайл:

Количество вариантов
– 2 +

Вывод
Вариант В лучше варианта А

	Число конверсий	Размер выборки	Конверсия	Доверительный интервал
Вариант А	1906	47094	4,0 %	3,9 – 4,2 %
Вариант В	2014	46772	4,3 %	4,1 – 4,5 %

Достоверность

95%

Декстоп:

Количество вариантов
– 2 +

Вывод
Варианты А и В значимо не различаются

	Число конверсий	Размер выборки	Конверсия	Доверительный интервал
Вариант А	800	59275	1,3 %	1,3 – 1,4 %
Вариант В	815	58844	1,4 %	1,3 – 1,5 %

Достоверность

95%

Последствия ошибок

Ошибки в дизайне АБ-тестов приводят к некорректным результатам, которые нельзя интерпретировать: нам могло повезти и тест покрасился, хотя на самом деле эффекта нет; и наоборот, тест мог не покраситься, хотя эффект есть.