

# Метрики в А/В- тестировании

Курс «Продуктовые метрики»

# Как выбрать метрику?

## **Требования:**

- Измеримая (measurable)
- Устойчивая (stable)
- Объяснимая (attributable)
- Чувствительная (sensitive)

# Как выбрать метрику?

**ОЕС** (Overall Evaluation Criterion) — количественная мера, отражающая основную цель эксперимента.

## Типы метрик в эксперименте:

- Таргет (goal) — метрика успеха организации
- Драйверы (driver / proxy) — влияющие на успех
- Guardrail / контр-метрики — метрики-ограничители, которые сигнализируют о нарушении

+

- Качество данных (data quality) — метрики валидности эксперимента
- Диагностики (diagnostic) / информативные / debug — дополнительная гранулярность, контекст

# Как выбрать метрику?

## Правила выбора:

- Таргет (goal):
  - Понятность
  - Устойчивость
- Драйверы:
  - Сонаправленность с таргет-метрикой
  - Управляемость
  - Чувствительность
  - Устойчивость к манипуляциям

# Как выбрать метрику?

## Как на практике?

### Метрик-сет:

- Таргет-метрика (P0)
- Метрики-драйверы (P1) — при дизайне эксперимента мы считаем «таргет»-метрикой именно драйвер, то есть H0 чаще всего формулируется для них
- Информативные гранулярные метрики, дающие контекст (P1-P2-P3)
- Guardrails (P3-P4; P1-P2 других команд, если есть риск каннибализации)
- Debug (P3-P4)
- Качество данных (метрики свойств эксперимента)

Не забываем, что при увеличении числа метрик в эксперименте появляется проблема множественного тестирования. Самый простой вариант решения — снижать  $\alpha$ .

# Как выбрать метрику?

**Пример:** Тестирование push-уведомления о брошенной корзине в e-commerce приложении.

## **Метрик-сет:**

- Таргет-метрики: GMV per User, Orders per User
- Метрики-драйверы: CR from Cart to Purchase
- Информативные: AOV, CTR, Orders per User by Category
- Guardrails: Unsubscribe Rate, Uninstall Rate, Errors, User Complaints
- Debug: Delivery Rate
- Качество данных: Balance Test / Control

# Как выбрать метрику?

## Как принять решение о раскатке?

Среди ключевых метрик эксперимента (таргет + драйверы + guardrails):

- все метрики значимы положительно / отрицательно — катим / не катим;
- стат значима лишь одна и эффект положительный — катим;
- стат значима лишь одна и эффект отрицательный — не катим;
- все не стат значимы — не катим, думаем о перезапуске (повышение чувствительности, изменение дизайна);
- часть значима положительно, часть отрицательно — принимаем решение на основе трейд оффа.

# Как выбрать метрику?

## Как принять решение о раскатке?

Для целей A/B-тестирования можно разработать взвешенную метрику, чтобы проще принимать решения в последнем кейсе.

Пример ОЕС для Email-рассылок в Amazon:

$$OEC = \frac{\Delta \text{Revenue} - (s \times \text{Unsubscribe\_Lifetime\_Loss})}{\text{Number\_of\_Users}}$$

- $\Delta \text{Revenue}$  — инкрементальная выручка от email рассылки в результате эксперимента
- $s$  — число пользователей, отписавшихся от получения писем
- $\text{Unsubscribe\_lifetime\_loss}$  — оценка потери выручки от отсутствия коммуникаций с юзером в email
- $\text{Number\_of\_Users}$  — число пользователей в тестовой группе



Поиск прокси-метрики

# Поиск прокси-метрики

**Прокси-метрика** — краткосрочный индикатор, предсказывающий долгосрочный эффект таргет-метрики.

## Требования:

- Чувствительность — позволяет выявить стат значимый эффект за более короткий период, чем требуется для таргета.
- Сонаправленность (Label Agreement) — изменения прокси сонаправлены с изменениями таргета, поэтому можно принимать решение о раскатке, опираясь на неё.

# Поиск прокси-метрики

**Единого алгоритма поиска прокси-метрики не существует.**

Но есть набор проверок, на которых можно построить алгоритм внутри компании:

- Binary Sensitivity (бинарная чувствительность)
- MSE с таргет-метрикой
- Корреляция с таргет-метрикой
- Proxy Score (или Label Agreement)

Для поиска прокси-метрик необходимо иметь накопленный корпус экспериментов. Мы оцениваем способность прокси предсказывать таргет именно по эффектам в экспериментах, а не по динамике самих временных рядов.

# Поиск прокси-метрики

**Binary Sensitivity** — доля экспериментов, в которых метрика показала статистически значимый эффект.

Рассчитывается для каждой метрики (и для прокси, и для таргет) по формуле:

$$\text{Binary Sensitivity}(\text{metric}) = \frac{|\{\text{experiments where metric is significant}\}|}{|\{\text{all experiments}\}|}$$

Обратите внимание: Binary Sensitivity ничего не говорит о сонаправленности, она лишь помогает оценить, получаем ли мы для кандидатов в прокси стат значимые результаты чаще, чем для таргета.

# Поиск прокси-метрики

**MSE (Mean Squared Error)** показывает, насколько величина эффекта прокси отклоняется от эффекта таргета в разных экспериментах.

Рассчитывается попарно для каждой прокси и таргета по формуле:

$$MSE = \frac{1}{N} \sum_{i=1}^N (\Delta_{\text{proxy},i}^{rel} - \Delta_{\text{target},i}^{rel})^2, \text{ где } \Delta^{rel} = \frac{\bar{x}_{test} - \bar{x}_{control}}{\bar{x}_{control}}$$

Меньшее значение MSE означает, что прокси ближе к таргету по оценке эффекта. Применяется только вместе с другими метриками для комплексной оценки.

# Поиск прокси-метрики

**Корреляция** измеряет силу и направление связи между эффектами прокси и таргета в исторических экспериментах.

Рассчитывается попарно для каждой прокси и таргета по формуле:

$$\rho = \frac{\sum_{i=1}^N \left( \Delta_{\text{proxy},i}^{\text{rel}} - \overline{\Delta_{\text{proxy}}^{\text{rel}}} \right) \left( \Delta_{\text{target},i}^{\text{rel}} - \overline{\Delta_{\text{target}}^{\text{rel}}} \right)}{\sqrt{\sum_{i=1}^N \left( \Delta_{\text{proxy},i}^{\text{rel}} - \overline{\Delta_{\text{proxy}}^{\text{rel}}} \right)^2} \cdot \sqrt{\sum_{i=1}^N \left( \Delta_{\text{target},i}^{\text{rel}} - \overline{\Delta_{\text{target}}^{\text{rel}}} \right)^2}}, \text{ где } \Delta^{\text{rel}} = \frac{\bar{x}_{\text{test}} - \bar{x}_{\text{control}}}{\bar{x}_{\text{control}}}$$

В примере корреляция Пирсона, но можно использовать любую другую.

# Поиск прокси-метрики

**Proxy Score** показывает, насколько прокси-метрика согласуется с таргет-метрикой в исторических экспериментах.

Рассчитывается попарно для каждой прокси и таргета по формуле:

$$\text{Proxy Score} = \frac{\text{Detections} - \text{Mistakes}}{\#\{\text{experiments where target is significant}\}}$$

- Detections — число экспериментов, где прокси и таргет значимы и сонаправлены по направлению эффекта.
- Mistakes — число экспериментов, где прокси значима, но не совпадает с таргетом по направлению.

Значение Proxy Score лежит в диапазоне от -1 до 1. Чем ближе к 1, тем надежнее прокси отражает таргет.

# Поиск прокси-метрики

**Матрица согласия-несогласия знаков** — это вспомогательный инструмент. Она визуализирует совпадения / расхождения по знаку эффекта в исторических экспериментах.

		North Star Metric			
		Long-Term Effect			
Proxy Metric	Short-Term Effect		negative	neutral	positive
		negative	34	203	3
		neutral	15	571	5
		positive	0	93	76



# Поиск прокси-метрики

## Принятие решения

- Оценка идет сразу по всем метрикам: Binary Sensitivity, MSE и/или корреляция, Proxy Score.
- Если результаты согласуются, прокси можно рассматривать как кандидата.
- Если есть противоречия — прокси не надежна.

## Предостережения

- Убедитесь, что прокси действительно нужны.
- Сначала попробуйте улучшить дизайн эксперимента.
- Выбирайте прокси со здравым смыслом.
- Валидируйте и мониторьте прокси.