

Формула размера выборки

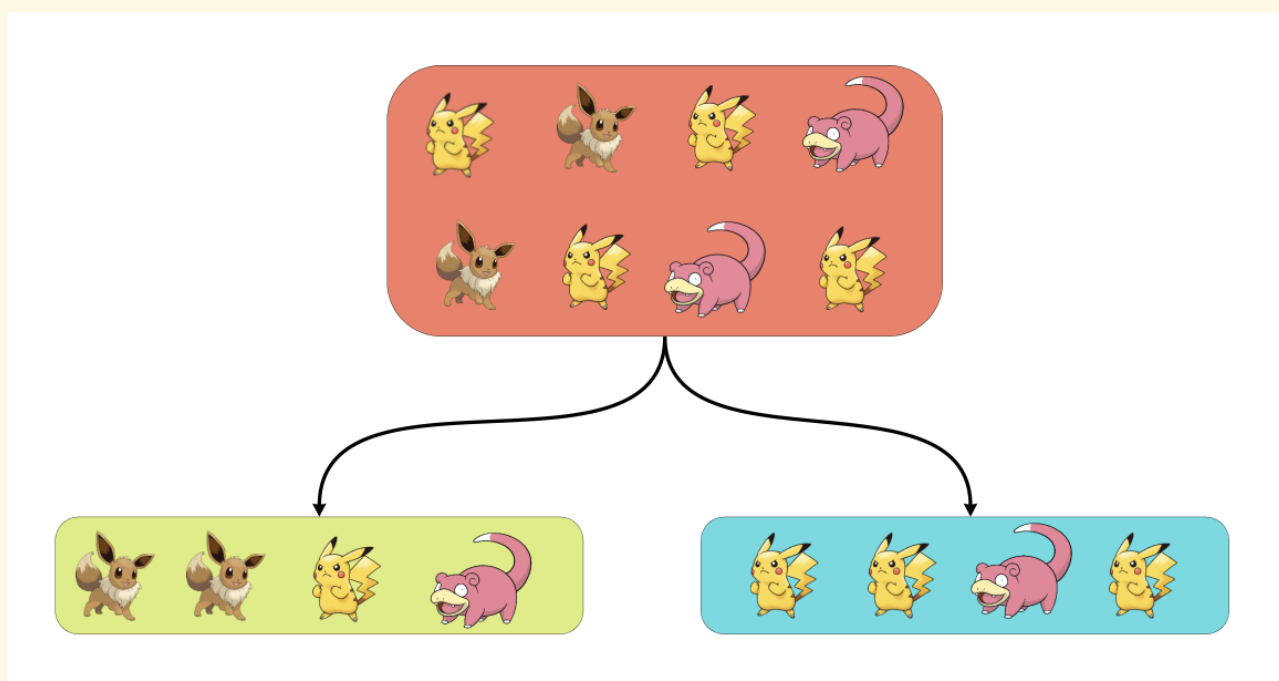
$$n = \frac{[\varphi^{-1}(1-\alpha) + \varphi^{-1}(1-\beta)]^2 (\sigma_x^2 + \sigma_y^2)}{MDE^2}$$

Стратификация

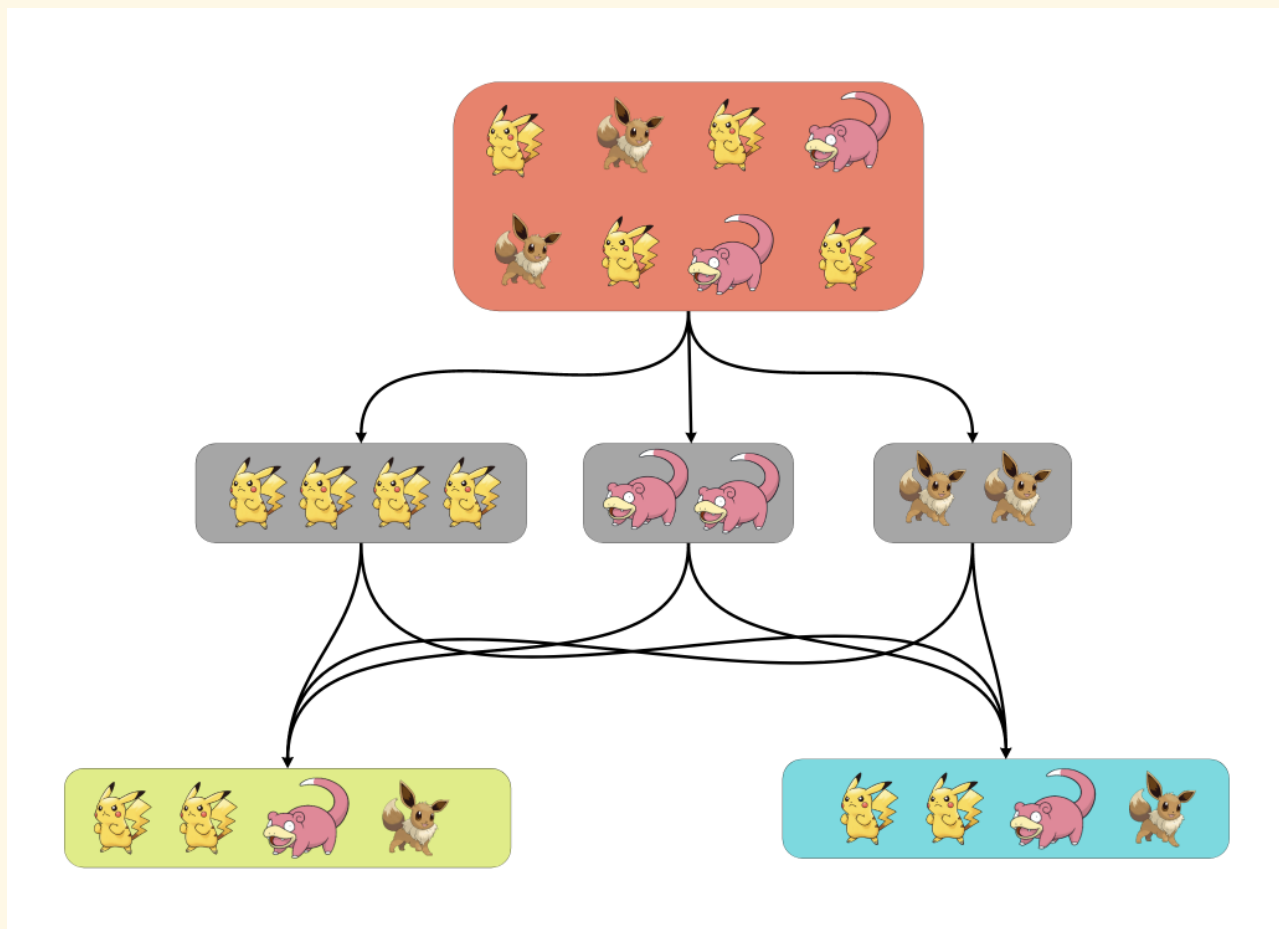
Стратификация
Страт. семплирование
+ стратифиц. оценка

Пост-стратификация
случайная выборка +
стратифицированные оценки
на основе заранее известных
вероятностей попадания в
группу

Случайное семплирование



Стратифицированное семплирование



Обозначения

- $\mu = E(Y)$ - популяционное среднее
- $\sigma^2 = D(Y)$
- μ_k, σ_k^2 - среднее и дисперсия для k -й страты
- w_k - доля k -й страты в популяции
- n_k - число пользователей k -й страты
- $n = \sum_{k=1}^K n_k$ - общий размер группы
- $Y_{11}, \dots, Y_{1n_1}, \dots, Y_{K1}, \dots, Y_{Kn_K}$ - выборка из Y , где Y_{kj} - метрика для j -го пользователя k -й страты

Простое среднее

$$\bar{Y} = \frac{1}{n} \sum_{k=1}^K \sum_{j=1}^{n_k} Y_{kj}$$

Стратифицированное среднее

$$\bar{Y}_{str} = \sum_{k=1}^K w_k \bar{Y}_k; \quad \bar{Y}_k = \frac{1}{n_k} \sum_{j=1}^{n_k} Y_{kj}$$

	Обычное среднее $\bar{Y} = \frac{1}{n} \sum_{k=1}^K \sum_{j=1}^{n_k} Y_{kj}$	Стратифицированное среднее $\bar{Y}_{str} = \sum_{k=1}^K w_k \bar{Y}_k$
Случайное семплирование	Классический подход Без стратификации	Постстратификация
Стратифицированное семплирование	Не контролирует вероятность ошибки I рода	Стратификация

Свойства

- $\bar{Y}_{strat} = \bar{Y}$ (оценки равны)
- $E(\bar{Y}_{strat}) = E(\bar{Y}) = E(\bar{Y}_{post}) = \mu$
→ пост-стратифицир.

Дисперсия

$$1. D(\bar{Y}) = \frac{1}{n} \sum_{k=1}^K w_k \sigma_k^2 + \underbrace{\left(\frac{1}{n} \sum_{k=1}^K w_k (\mu_k - \mu)^2 \right)}_{\text{межгрупповая дисперсия}}$$

$$2. D(\bar{Y}_{strat}) = \frac{1}{n} \sum_{k=1}^K w_k \sigma_k^2$$

$$D(\bar{Y}_{post-strat}) = \frac{1}{n} \sum_{k=1}^K w_k \sigma_k^2 + \underbrace{\left(\frac{1}{n^2} \sum_{k=1}^K (1-w_k) \sigma_k^2 \right)}_{\text{дисперсия}} + O\left(\frac{1}{n^2}\right)$$

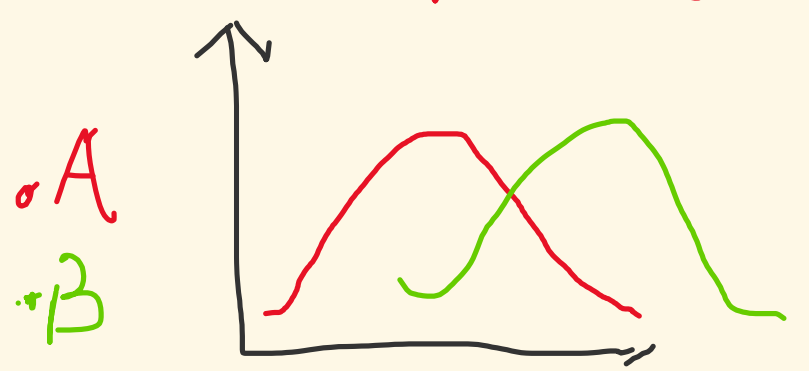
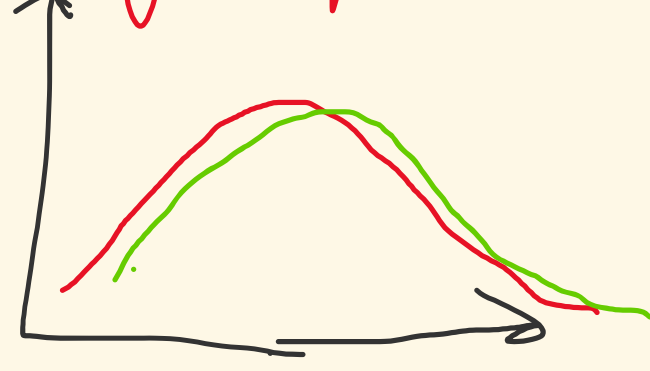
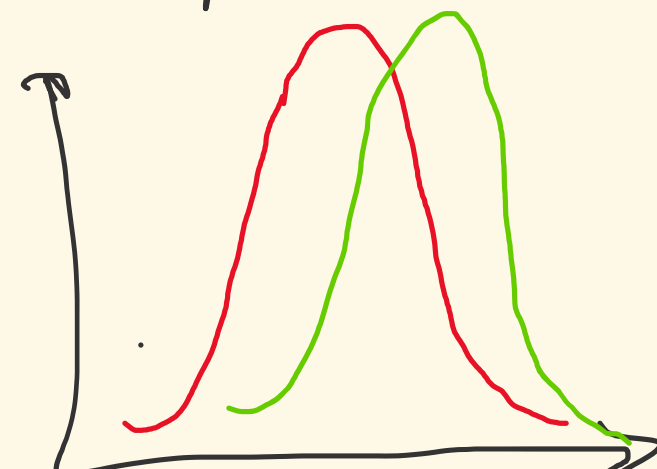
$$D(\bar{Y}_{strat}) = D(\bar{Y}_{post-strat}) + O\left(\frac{1}{n^2}\right) = D(\bar{Y}) + O\left(\frac{1}{n^2}\right)$$

$$D(\bar{Y}_{strat}) \leq D(\bar{Y}_{post-strat}) \leq D(\bar{Y})$$

CUPED

Controlled - experiment Using Pre-Experimental data

Зачем сокращать дисперсию у метрики?

Большой эффект
высокая дисперсияМаленький эффект
высокая дисперсияМаленький эффект
маленькая дисперсия

Как сокращать дисперсию?

Можно использовать дополнительную информацию (историю)

1) Стратификация

2) Использование ковариат

 $\hat{\theta}$ - ненамеченные оценки на θ

$$E(\hat{\theta}) = \theta$$

$$\hat{\theta}_{cuped} = \hat{\theta} - \alpha \hat{\theta}_a, \quad E(\hat{\theta}_a) = 0 \text{ - для ненамеченности}$$

ковариата

$$D(\hat{\theta}_{cuped}) = D(\hat{\theta}) + \alpha^2 D(\hat{\theta}_a) - 2\alpha \underbrace{\text{cov}(\hat{\theta}, \hat{\theta}_a)}_{\text{оценки сильно коррелируют}}$$

Оптимизируем по α , чтобы $D(\hat{\theta}_{cuped}) \rightarrow \min$

$$\frac{\partial}{\partial \alpha} = 2\alpha D(\hat{\theta}_a) - 2 \text{cov}(\hat{\theta}, \hat{\theta}_a) = 0$$

$$\alpha = \frac{\text{cov}(\hat{\theta}, \hat{\theta}_a)}{D(\hat{\theta}_a)}$$

$$\hat{\theta}_{cuped} = \hat{\theta} - \frac{\text{cov}(\hat{\theta}, \hat{\theta}_a)}{D(\hat{\theta}_a)} \cdot \hat{\theta}_a$$

$$D(\hat{\theta}_{cuped}) \leq D(\hat{\theta})$$

На практике, чтобы $E(\hat{\theta}_a) = 0$ делают так:

$$\hat{\theta}_{cuped} = \hat{\theta} - \alpha \cdot (\hat{\theta}_a - \hat{\theta}_a^{\text{среднее}})$$

! Среднее по ковариате ($\hat{\theta}_a$) и α считаются по всем данным (тестовая группа + контрольная + те, которые вы не взяли в группу)