

# Престратификация

Материал написан на основе статьи [Xie, H., & Aurisset, J. \(2016, August\). Improving the sensitivity of online controlled experiments: Case studies at netflix. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining \(pp. 645-654\).](#)

## 1) Стратификация (stratification sampling)

Мы помним, что снижение дисперсии снижает необходимое число наблюдений.

Предположим, что мы выбрали переменную, по которой мы можем разбить нашу выборку на группы (страты). Пусть таких групп  $K$  штук, а  $n_k$  – численность каждой из них, то есть  $\sum_{k=1}^K n_k = N$ , где  $N$  величина генеральной совокупности. Тогда вероятность попасть в одну из групп равняется  $p_k = \frac{n_k}{N}$  – доля людей из  $k$ -й страты в генеральной совокупности.

Пусть наша зависимая переменная равна  $Y$ , причем  $\mathbb{E}(Y) = \mu$  и  $\text{var}(Y) = \sigma^2$ . Тогда среднее значение зависимой переменной равно:

- Обычное среднее по всей выборке:

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

- Средневзвешенное при стратификации равно сумме средних в каждой страте, взвешенных по вероятностям попасть в каждую из этих страт:

$$\hat{Y}_{\text{strat}} = \sum_{k=1}^K \underbrace{\frac{n_k}{N}}_{p_k} \bar{Y}_k = \sum_{k=1}^K p_k \bar{Y}_k, \text{ где } \bar{Y}_k = \frac{1}{n_k} \sum_{j=1}^{n_k} Y_{kj} \text{ среднее значение метрики внутри } k\text{-й страты}$$

Распишем среднее при стратификации подробнее

$$\underbrace{\sum_{k=1}^K p_k \bar{Y}_k}_{\hat{Y}_{\text{strat}}} = \sum_{k=1}^K \frac{n_k}{N} \frac{1}{n_k} \sum_{j=1}^{n_k} Y_{kj} = \underbrace{\frac{1}{N} \sum_{k=1}^K \sum_{j=1}^{n_k} Y_{kj}}_{\bar{Y}} \text{ обе средние равны}$$

Найдем среднее и дисперсию зависимой переменной при стратификации.

- Среднее

$$\mathbb{E}(\hat{Y}_{\text{strat}}) = \mathbb{E}\left[\sum_{k=1}^K p_k \bar{Y}_k\right] = \sum_{k=1}^K p_k \mathbb{E}(\bar{Y}_k) = \sum_{k=1}^K p_k \mu_k = \boxed{\mu}$$

- Дисперсия

$$\text{var}(\hat{Y}_{\text{strat}}) = \text{var}\left[\sum_{k=1}^K p_k \bar{Y}_k\right] = \sum_{k=1}^K p_k^2 \text{var}(\bar{Y}_k) = \sum_{k=1}^K \frac{n_k^2}{N^2} \frac{1}{n_k} n_k \sigma_k^2 = \frac{1}{N} \sum_{k=1}^K \frac{n_k}{N} \sigma_k^2 = \boxed{\frac{1}{N} \sum_{k=1}^K p_k \sigma_k^2} \quad (1)$$

## 2) Простая рандомизация (simple random sampling)

Аналогично найдем среднее и дисперсию зависимой переменной при классическом эксперименте.

- Среднее

$$\mathbb{E}(\bar{Y}) = \mathbb{E}\left[\frac{1}{N} \sum_{k=1}^K \sum_{j=1}^{n_k} Y_{kj}\right] = \frac{1}{N} \sum_{k=1}^K \sum_{j=1}^{n_k} \mathbb{E}(Y_{kj}) = \frac{1}{N} \sum_{k=1}^K \sum_{j=1}^{n_k} \mu = \frac{1}{N} \sum_{k=1}^K n_k \mu = \frac{1}{N} N \mu = \boxed{\mu}$$

- Дисперсия

$$\text{var}(\bar{Y}) = \text{var}\left[\frac{1}{N} \sum_{k=1}^K \sum_{j=1}^{n_k} Y_{kj}\right] = \frac{1}{N^2} \sum_{k=1}^K \sum_{j=1}^{n_k} \text{var}(Y_{kj}) = \frac{1}{N^2} \sum_{k=1}^K n_k \sigma^2 = \frac{N \sigma^2}{N^2} = \boxed{\frac{\sigma^2}{N}} \quad (2)$$

Дисперсия зависимой переменной при рандомизированном эксперименте может быть представлена в виде суммы внутригрупповой и межгрупповой дисперсии.

Для этого нам понадобится total variance law (см. доказательство на последней странице):

$$var(Y) = var[\mathbb{E}(Y | X)] + \mathbb{E}[var(Y | X)]$$

Пусть  $Z$  – номер страты от 1 до  $K$ , тогда:

$$var(Y) = \mathbb{E}(var(Y | Z)) + var(\mathbb{E}(Y | Z)) =$$

$$= \{I(Z = k) \text{ индикаторная переменная, равная 1, если } Z = k, \text{ и равная нулю иначе}\} =$$

$$= \mathbb{E} \left[ \sum_{k=1}^K \sigma_k^2 I(Z = k) \right] + var \left[ \sum_{k=1}^K \mu_k I(Z = k) \right] =$$

$$= \sum_{k=1}^K \sigma_k^2 \mathbb{E}[I(Z = k)] + \mathbb{E} \left[ \sum_{k=1}^K \mu_k I(Z = k) \right]^2 - \left[ \mathbb{E} \left[ \sum_{k=1}^K \mu_k I(Z = k) \right] \right]^2 =$$

$$= \sum_{k=1}^K \sigma_k^2 p_k + \sum_{k=1}^K \mu_k^2 p_k - \mu^2 = \{*\} = \sum_{k=1}^K \sigma_k^2 p_k + \sum_{k=1}^K p_k (\mu_k - \mu)^2 \quad (3)$$

$$(*) : \sum_{k=1}^K p_k (\mu_k - \mu)^2 = \sum_{k=1}^K p_k (\mu_k^2 - 2\mu\mu_k + \mu^2) =$$

$$= \sum_{k=1}^K p_k \mu_k^2 - 2\mu \sum_{k=1}^K \mu_k p_k + \mu^2 \sum_{k=1}^K p_k = \sum_{k=1}^K \mu_k^2 p_k - 2\mu^2 + \mu^2 = \sum_{k=1}^K \mu_k^2 p_k - \mu^2$$

Из (2) и (3) следует:

$$var(\bar{Y}) = \frac{\sigma^2}{N}$$

$$var(Y) = \sum_{k=1}^K \sigma_k^2 p_k + \sum_{k=1}^K p_k (\mu_k - \mu)^2$$

$$var(\bar{Y}) = \underbrace{\frac{1}{N} \sum_{k=1}^K p_k \sigma_k^2}_{\text{внутригрупповая дисперсия}} + \underbrace{\frac{1}{N} \sum_{k=1}^K p_k (\mu_k - \mu)^2}_{\text{межгрупповая дисперсия}} \quad (4)$$

Из (1) и (4) следует:

$$var(\hat{Y}_{\text{strat}}) = \frac{1}{N} \sum_{k=1}^K p_k \sigma_k^2$$

$$var(\bar{Y}) = \underbrace{\frac{1}{N} \sum_{k=1}^K p_k \sigma_k^2}_{\text{внутригрупповая дисперсия}} + \underbrace{\frac{1}{N} \sum_{k=1}^K p_k (\mu_k - \mu)^2}_{\text{межгрупповая дисперсия}}$$

Дисперсия среднего значения зависимой переменной всегда больше, чем дисперсия среднего значения зависимой переменной в случае стратификации, поскольку межгрупповая дисперсия всегда неотрицательная. Таким образом,  $var(\bar{Y}) \geq var(\hat{Y}_{\text{strat}})$

## Приложение

### Условные обозначения

- $Y$  – зависимая переменная, причем
  - $\mathbb{E}(Y) = \mu$
  - $\text{var}(Y) = \sigma^2$
- Выборка разбита на  $k$  страт по показателю  $X$ , то есть
  - $\mu_k$  – среднее в страте
  - $\sigma_k^2$  – дисперсия в страте
  - $n_k$  – численность в страты, то есть  $\sum_{k=1}^K n_k = N$
  - $p_k = \frac{n_k}{N}$  – доля людей из  $k$ -й страты в генеральной совокупности
  - $Y_{kj}$  – метрика  $j$ -го человека из  $k$ -й страты
- Тогда среднее значение зависимой переменной
  - $\bar{Y} = \frac{1}{N} \sum_{k=1}^K \sum_{j=1}^{n_k} Y_{kj}$  обычное среднее
  - $\hat{Y}_{\text{strat}} = \sum_{k=1}^K p_k \bar{Y}_k$  среднее при стратификации
    - \*  $\bar{Y}_k = \frac{1}{n_k} \sum_{j=1}^{n_k} Y_{kj}$  среднее значение метрики внутри  $k$ -й страты
  - $\underbrace{\sum_{k=1}^K p_k \bar{Y}_k}_{\hat{Y}_{\text{strat}}} = \sum_{k=1}^K \frac{n_k}{N} \frac{1}{n_k} \sum_{j=1}^{n_k} Y_{kj} = \underbrace{\frac{1}{N} \sum_{k=1}^K \sum_{j=1}^{n_k} Y_{kj}}_{\bar{Y}}$  обе средние равны

### Total variance law

**Доказать:**

$$\text{var}(Y) = \text{var}[\mathbb{E}(Y | X)] + \mathbb{E}[\text{var}(Y | X)]$$

**Доказательство:**

$$\begin{aligned} \text{var}(Y) &= \mathbb{E}(Y^2) - (\mathbb{E}(Y))^2 = \\ &= \{ \text{закон повторного мат. ожидания } \mathbb{E}(Y) = \mathbb{E}(\mathbb{E}(Y | X)) \} = \\ &= \mathbb{E}(\mathbb{E}(Y^2 | X)) - [\mathbb{E}(\mathbb{E}(Y | X))]^2 = \\ &= \left\{ \text{добавим и вычтем } \mathbb{E}[(\mathbb{E}(Y | X))^2 | X] \right\} = \\ &= \underbrace{\mathbb{E}[\mathbb{E}(Y^2 | X) - (\mathbb{E}(Y | X))^2 | X]}_{\mathbb{E}[\text{var}(Y|X)]} + \underbrace{\mathbb{E}[(\mathbb{E}(Y | X))^2 | X] - \mathbb{E}[\mathbb{E}(Y | X)]\mathbb{E}[\mathbb{E}(Y | X)]}_{\text{var}[\mathbb{E}(Y|X)]} = \\ &= \mathbb{E}[\text{var}(Y | X)] + \text{var}[\mathbb{E}(Y | X)] \end{aligned}$$