

Теория вероятностей и математическая статистика
Доверительные интервалы III.

Глеб Карпов

ФКН ВШЭ

Доверительные интервалы для разности истинных долей

Мотивация

- В реальном мире часто возникает необходимость **сравнивать** доли между двумя группами:
 - Какая версия сайта лучше конвертирует посетителей в покупателей?
 - Какая маркетинговая кампания эффективнее привлекает клиентов?
 - На какой производственной линии меньше брака?
 - Какое лекарство эффективнее лечит заболевание?
- Пусть p_1 — истинная доля в первой генеральной совокупности, p_2 — истинная доля во второй генеральной совокупности. Нас интересует **разность** $\theta = p_1 - p_2$.
- Для бизнеса, маркетинга, медицины и социологии критически важно знать, есть ли **статистически значимая разница** между долями в двух группах. Это позволяет принимать обоснованные решения о выборе стратегии, продукта или лечения.
- Новый вопрос в статистике: как построить доверительный интервал для разности $p_1 - p_2$?
- **Важная идея:** если доверительный интервал для $p_1 - p_2$ не содержит нуль, это означает статистически значимую разницу между долями в двух группах.

Точечная оценка для разности истинных долей

- Предположим, у нас есть случайная выборка $X = \{X_1, \dots, X_n\}$ из распределения Бернулли с $P(X_i = 1) = p_1$, и случайная выборка $Y = \{Y_1, \dots, Y_m\}$ из распределения Бернулли с $P(Y_i = 1) = p_2$. Тогда обе случайные выборки - процессы Бернулли длины n и m соответственно.
- Нас интересует разность истинных долей (или, то же самое, разность вероятностей успеха):

$$\theta = p_1 - p_2$$

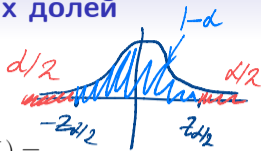
- Если $n, m > 30$, то по ИТМЛ:

$$\hat{p}_1 \sim \mathcal{N}\left(p_1, \frac{p_1(1-p_1)}{n}\right), \quad \hat{p}_2 \sim \mathcal{N}\left(p_2, \frac{p_2(1-p_2)}{m}\right)$$

- Введём точечную оценку $\hat{\theta} = \hat{p}_1 - \hat{p}_2$ — разность двух выборочных долей.
- Свойства точечной оценки: $E[\hat{\theta}] = p_1 - p_2$, $Var[\hat{\theta}] = \frac{p_1(1-p_1)}{n} + \frac{p_2(1-p_2)}{m}$.
- Так как сумма двух нормальных случайных величин — нормальная случайная величина:

$$\hat{\theta} \sim \mathcal{N}\left(p_1 - p_2, \frac{p_1(1-p_1)}{n} + \frac{p_2(1-p_2)}{m}\right) \quad \frac{\hat{\theta} - (p_1 - p_2)}{std[\hat{\theta}]} \sim \mathcal{N}(0, 1)$$

Доверительные интервалы для разности истинных долей



- Если выполнены условия ИТМЛ, то действуем знакомым способом:

$$1 - \alpha = P(L < p_1 - p_2 < U) = P(-U < -(p_1 - p_2) < -L) =$$
$$P\left(\frac{\hat{\theta} - U}{\sqrt{\text{Var}(\hat{\theta})}} < \frac{\hat{\theta} - \theta}{\sqrt{\text{Var}(\hat{\theta})}} < \frac{\hat{\theta} - L}{\sqrt{\text{Var}(\hat{\theta})}}\right) = P\left(-z_{\alpha/2} < Z < z_{\alpha/2}\right) = 1 - \alpha$$

- Находим точку $z_{\alpha/2}$, такую, что $P(Z > z_{\alpha/2}) = \frac{\alpha}{2}$, выполняем обратное преобразование, и находим теоретические границы такого доверительного интервала:

$$L = \hat{p}_1 - \hat{p}_2 - z_{\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n} + \frac{p_2(1-p_2)}{m}}, \quad U = \hat{p}_1 - \hat{p}_2 + z_{\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n} + \frac{p_2(1-p_2)}{m}}$$

где p_1 и p_2 — истинные параметры (константы), которые заменяются на их точечные оценки на практике, так как истинные параметры неизвестны.

Доверительные интервалы для разности истинных долей

На практике $(1 - \alpha)100\%$ доверительный интервал для разности долей генеральных совокупностей:

$$p_1 - p_2 \in \left(\tilde{p}_1 - \tilde{p}_2 - z_{\alpha/2} \sqrt{\frac{\tilde{p}_1(1 - \tilde{p}_1)}{n} + \frac{\tilde{p}_2(1 - \tilde{p}_2)}{m}}, \tilde{p}_1 - \tilde{p}_2 + z_{\alpha/2} \sqrt{\frac{\tilde{p}_1(1 - \tilde{p}_1)}{n} + \frac{\tilde{p}_2(1 - \tilde{p}_2)}{m}} \right),$$

где \tilde{p}_1, \tilde{p}_2 — реализации выборочных долей.

Иллюстративные задачи $1-\alpha = 0.95$

$$\alpha = 0.05$$

$$Z_{\alpha/2}: P(Z > Z_{\alpha/2}) = 0.025 = 1 - F_Z(Z_{\alpha/2})$$

Пример 1: сравнение версий сайта

Компания проводит сравнение двух версий главной страницы сайта. Версия А показана 500 пользователям, из них 85 совершили покупку. Версия В показана 480 пользователям, из них 112 совершили покупку. Постройте 95% доверительный интервал для разности долей конверсии между версиями.

X — Bern, покупка на версии А

$$X \sim \text{Bern}(p_1)$$

Y — Bern, — // — В

$$Y \sim \text{Bern}(p_2)$$

$$L, U: 1-\alpha = P(L < p_1 - p_2 < U)$$

$$\theta = p_1 - p_2$$

Иллюстративные задачи

Решение

$$\tilde{p}_1 = \frac{85}{500} = 0.17, \tilde{p}_2 = \frac{112}{480} = 0.233, z_{0.025} = 1.96:$$

$$p_1 - p_2 \in \left(0.17 - 0.233 - 1.96 \sqrt{\frac{0.17 \cdot 0.83}{500} + \frac{0.233 \cdot 0.767}{480}}, 0.17 - 0.233 + 1.96 \sqrt{\frac{0.17 \cdot 0.83}{500} + \frac{0.233 \cdot 0.767}{480}} \right)$$
$$= (-0.103, -0.023)$$

Интервал не содержит ноль, можем сделать вывод, что версия сайта В имеет статистически значимо более высокую конверсию.



Доверительные интервалы для разности математических ожиданий

Мотивация

- В коммерческих и научных исследованиях бывает необходимость сравнить средние значения между двумя группами:
 - Какая производственная линия более производительна?
 - Какой метод обучения даёт лучшие результаты?
 - В каком регионе выше средний доход населения?
- Пусть μ_X — истинное математическое ожидание в первой генеральной совокупности, μ_Y — истинное математическое ожидание во второй генеральной совокупности. Нас интересует **разность** $\theta = \mu_X - \mu_Y$.
- Важно знать, есть ли **статистически значимая разница** между математическими ожиданиями в двух группах. Это позволяет принимать дальнейшие обоснованные решения.
- Новый вопрос: как построить доверительный интервал для разности $\mu_X - \mu_Y$?
- **Идея:** если доверительный интервал для $\mu_X - \mu_Y$ не содержит нуль, это означает статистически значимую разницу между средними значениями в двух группах.

Точечная оценка разности математических ожиданий

Если дисперсии известны

- Предположим, у нас есть две независимые выборки: $\mathcal{X} = \{X_1, \dots, X_n\}$, $\mathcal{Y} = \{Y_1, \dots, Y_m\}$. Характеристики называем $\mu_X \equiv E[X_i]$, $\sigma_X^2 \equiv Var[X_i]$, и соответственно $\mu_Y \equiv E[Y_i]$, $\sigma_Y^2 \equiv Var[Y_i]$
- Дисперсии σ_X^2 и σ_Y^2 предполагаем **известными**.
- Нас интересует разность истинных математических ожиданий:

$$\theta = \mu_X - \mu_Y$$

- Введём точечную оценку $\hat{\theta} = \bar{X} - \bar{Y}$ — разность двух выборочных средних.
- Свойства точечной оценки: $E[\hat{\theta}] = \mu_X - \mu_Y$, $Var[\hat{\theta}] = \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}$.
- При $n, m > 30$ работает ЦПТ и распределение точечной оценки для θ :

$$\hat{\theta} \sim \mathcal{N}\left(\mu_X - \mu_Y, \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}\right)$$
$$\frac{\hat{\theta} - (\mu_X - \mu_Y)}{std[\hat{\theta}]} \sim \mathcal{N}(0, 1)$$

Доверительные интервалы для разности математических ожиданий

Если дисперсии известны

- Если выполнены условия (большие выборки или нормальное распределение), то действуем знакомым способом:

$$1 - \alpha = P(L < \mu_X - \mu_Y < U) = P(-U < -(\mu_X - \mu_Y) < -L) =$$
$$P\left(\frac{\bar{X} - \bar{Y} - U}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}} < \frac{\bar{X} - \bar{Y} - (\mu_X - \mu_Y)}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}} < \frac{\bar{X} - \bar{Y} - L}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}}\right) = P\left(-z_{\alpha/2} < Z < z_{\alpha/2}\right)$$

- Находим точку $z_{\alpha/2}$, такую, что $P(Z > z_{\alpha/2}) = \frac{\alpha}{2}$, выполняем обратное преобразование, и находим теоретические границы такого доверительного интервала:

$$L = \bar{X} - \bar{Y} - z_{\alpha/2} \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}, \quad U = \bar{X} - \bar{Y} + z_{\alpha/2} \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}$$

где σ_X^2 и σ_Y^2 — известные дисперсии (константы).

Доверительные интервалы для разности математических ожиданий

Если дисперсии известны

На практике $(1 - \alpha)100\%$ доверительный интервал для разности математических ожиданий:

$$\mu_1 - \mu_2 \in \left(\bar{x} - \bar{y} - z_{\alpha/2} \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}, \bar{x} - \bar{y} + z_{\alpha/2} \sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}} \right),$$

где \bar{x} , \bar{y} — реализации выборочных средних, а σ_X^2 и σ_Y^2 — известные дисперсии исследуемых случайных процессов.

Иллюстративные задачи

Пример 2: Сравнение производительности двух производственных линий

Производственная компания хочет сравнить среднюю производительность двух производственных линий. Известно, что стандартное отклонение производительности для первой линии составляет $\sigma_1 = 12$ единиц продукции в час, а для второй линии — $\sigma_2 = 15$ единиц продукции в час.

Было проведено тестирование: - Первая линия: выборка из $n = 40$ часов работы, средняя производительность $\bar{x} = 145$ единиц/час - Вторая линия: выборка из $m = 35$ часов работы, средняя производительность $\bar{y} = 138$ единиц/час

76%

1. Постройте 95% доверительный интервал для разности средних производительностей двух линий.
2. Можете ли вы сделать вывод о том, какая линия более производительна?

X — произв. линии 1, $E[X] = \mu_x$, $\text{Var}[X] = \sigma_x^2 = \sigma_1^2$

Y — — — 2, $E[Y] = \mu_y$, $\text{Var}[Y] = \sigma_y^2 = \sigma_2^2$

$$3. 1 - \alpha = 0.76$$

$$\alpha = 0.24$$

$$\alpha/2 = 0.12 \Rightarrow$$

$$z_{\alpha/2}: P(Z > z_{\alpha/2}) = 0.12 = 1 - F_Z(z_{\alpha/2})$$

$$z_{\alpha/2} = 1.18$$

$$(l, u) = (7 - 1.18 \cdot 3.17, 7 + 1.18 \cdot 3.17) = (3.26, 10.74)$$

Иллюстративные задачи

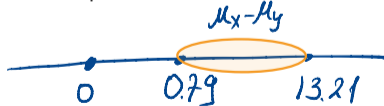
Решение

1. $n = 40, m = 35, \bar{x} = 145, \bar{y} = 138, \sigma_1 = 12, \sigma_2 = 15, z_{0.025} = 1.96$:

$$\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}} = \sqrt{\frac{12^2}{40} + \frac{15^2}{35}} = \sqrt{3.6 + 6.43} = \sqrt{10.03} \approx 3.17$$

$$\mu_1 - \mu_2 \in (145 - 138 - 1.96 \cdot 3.17, 145 - 138 + 1.96 \cdot 3.17) = (0.79, 13.21)$$

2. Интервал $(0.79, 13.21)$ не содержит нуль и полностью находится в положительной области. Это означает, что мы можем с определенной степенью полагать, что первая линия имеет статистически значимо более высокую производительность, чем вторая.



X - с. вел., кепр.,

Иллюстративные задачи

$L, U:$

$$E[X] \equiv \mu_X \text{ - ?}$$

$$1-d = P(L < \mu_X < U)$$

$$\text{Var}[X] \equiv \sigma_X^2 \quad \checkmark$$

Пример 3: средний чек в ресторане

Менеджер ресторана хочет оценить среднюю сумму m , которую посетитель тратит на обед. По выборке из $n = 36$ посетителей получено выборочное среднее $\bar{x} = 3.60$ (в долларах). Известно, что стандартное отклонение расходов одного посетителя равно $\sigma = 0.72$.

Найдите уровень доверия, которому соответствует доверительный интервал $(3.5; 3.7)$ для m .

$$n = 36$$

$$\bar{X} \sim N\left(\mu_X, \frac{\sigma_X^2}{n}\right), \text{ CLT}$$

$$\bar{x} - z_{d/2} \frac{\sigma}{\sqrt{n}}, \quad \bar{x} + z_{d/2} \frac{\sigma}{\sqrt{n}} = 3.7$$

3.6

$$z_{d/2} \cdot \frac{\sigma}{\sqrt{n}} = 0.1$$

$$z_{d/2} = \frac{0.1 \cdot \sqrt{n}}{\sigma} = \frac{0.1 \cdot 6}{0.72} = 0.83$$

$$d/2 = 1 - 0.7967 \approx 0.203$$

$$1-d = 0.594$$

$$\tilde{p}_x = \frac{13}{30} \approx 0.43$$

$$\tilde{p}_y = \frac{15}{45} \approx 0.33$$

Иллюстративные задачи

$$0.27 \text{ (U)} = \tilde{p}_x - \tilde{p}_y + z_{d/2} \cdot \sqrt{\frac{\tilde{p}_x(1-\tilde{p}_x)}{n} + \frac{\tilde{p}_y(1-\tilde{p}_y)}{m}}$$

$$z_{d/2} = \frac{0.17}{\dots} = \frac{0.17}{0.1144} \approx 1.49$$

$$d/2 = 1 - 0.9319 \Rightarrow d = 0.136, t-d = 0.864$$

Пример 4: разность долей конверсии

Сравнивают две рекламные кампании. В первой опросили $n = 30$ клиентов, из них 13 совершили покупку. Во второй — $m = 1.5n = 45$ клиентов, из них 15 совершили покупку. Построенный доверительный интервал для разности долей конверсии $p_1 - p_2$ равен $(-0.07; 0.27)$.

Найдите уровень доверия этого интервала.

$$X \sim \text{Ber}(p_1)$$

$$Y \sim \text{Ber}(p_2)$$

$$\Theta = p_1 - p_2 \quad \checkmark$$

$$\hat{p}_x = \frac{\sum X_i}{n} ; \hat{p}_x \sim \mathcal{N}$$

$$\hat{p}_y = \frac{\sum Y_i}{m} ; \hat{p}_y \sim \mathcal{N}$$