

Лекция 1

Случайная выборка

Точечные оценки

Кое-какие организационные вопросы

Зачем нужна статистика?

Типичное статистическое исследование:

Для изучения некоторой генеральной совокупности объектов была сделана случайная выборка из этой совокупности, состоящая из n объектов. У каждого обследованного объекта был измерен некий признак, так что имеется ряд результатов измерения:

$$x_1, x_2, \dots, x_n.$$

Задачи:

- > представление данных в удобном для восприятия виде (сведение к осмысленным характеристикам: среднее, стандартное отклонение и т. п.; таблицы, графики);
- > распространение результатов выборочного обследования на генеральную совокупность (оценивание параметров генеральной совокупности, проверка гипотез об этих параметрах).

статистические признаки

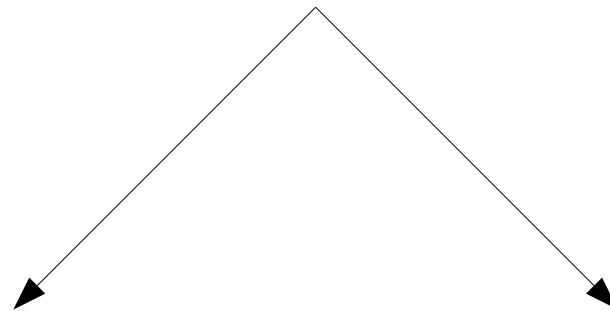


КОЛИЧЕСТВЕННЫЕ

*определены
арифметические
операции*



качественные



порядковые
(ординальные)

*определено отношение
порядка: можно
сравнивать «больше-
равно-меньше»*

НОМИНАЛЬНЫЕ

*можно сравнивать
«равно-не равно»*

Основной приём математической статистики

Мы считаем, что располагаемые данные x_1, x_2, \dots, x_n суть реализации (значения) случайных величин X_1, X_2, \dots, X_n .

Величины X_1, X_2, \dots, X_n будем называть *случайной выборкой*.

Числа x_1, x_2, \dots, x_n — *реализация* случайной выборки.

Что даёт нам право так считать?

Простейший случай: величины X_1, X_2, \dots, X_n независимы и одинаково распределены (обозначается $X_i \sim i.i.d.$ - independent identically distributed).

Когда эта предпосылка выполняется? Когда нарушается?

Замечание. На практике и числа x_1, x_2, \dots, x_n , и случайные величины X_1, X_2, \dots, X_n называют выборкой.

Точечные оценки

Пусть распределение признака в генеральной совокупности известно нам с точностью до параметра θ (возможно — с точностью до нескольких параметров $\theta_1, \dots, \theta_p$).

Например, мы знаем, что $X_i \sim N(\mu, \sigma^2)$, но значения параметров μ и σ^2 нам неизвестны.

Оценкой для параметра θ по выборке X_1, \dots, X_n называется случайная величина, определённая на том же пространстве элементарных исходов, что и случайные величины X_1, \dots, X_n , и являющаяся функцией от X_1, \dots, X_n .

*«Данное определение удивляет своей широтой, переходящей в бессмысленность»
(А.С. Шведов, «Теория вероятностей и математическая статистика», с. 114)*

Итак, $\hat{\theta} = f(X_1, \dots, X_n)$ - оценка для θ .

Помните: оценка — случайная величина, оцениваемый параметр — нет!

А почему?

Свойства оценок

Пусть Θ — множество допустимых значений параметра θ .

I. Несмещённость

Оценка $\hat{\theta}$ для параметра θ называется *несмещённой*, если

$$E(\hat{\theta}) = \theta \quad \forall \theta \in \Theta.$$

Что это значит?

Смещение оценки $\hat{\theta}$ для параметра θ :

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta} - \theta) = E(\hat{\theta}) - \theta.$$

несмещённая оценка — unbiased estimator

Пример

Имеется случайная выборка

$$X_1, X_2, \dots, X_n,$$

такая что

$$E(X_i) = \mu.$$

Проверьте на несмещённость оценки для μ :

$$\hat{\mu}_1 = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n};$$

$$\hat{\mu}_2 = X_1;$$

$$\hat{\mu}_3 = X_1 - 3X_2.$$

Пример

Имеется случайная выборка

$$X_1, X_2, \dots, X_n,$$

такая что

$$E(X_i) = \mu.$$

Проверьте на несмещённость оценки для μ :

$$\hat{\mu}_1 = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n};$$

$$\hat{\mu}_2 = X_1;$$

$$\hat{\mu}_3 = X_1 - 3X_2.$$

Решение.

$$E(\hat{\mu}_1) = E\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{E(X_1) + \dots + E(X_n)}{n} = \frac{n\mu}{n} = \mu \quad \Rightarrow \text{несмещённая.}$$

Пример

Имеется случайная выборка

$$X_1, X_2, \dots, X_n,$$

такая что

$$E(X_i) = \mu.$$

Проверьте на несмещённость оценки для μ :

$$\hat{\mu}_1 = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n};$$

$$\hat{\mu}_2 = X_1;$$

$$\hat{\mu}_3 = X_1 - 3X_2.$$

Решение.

$$E(\hat{\mu}_1) = E\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{E(X_1) + \dots + E(X_n)}{n} = \frac{n\mu}{n} = \mu \quad \Rightarrow \text{несмещённая.}$$

$$E(\hat{\mu}_2) = E(X_1) = \mu \quad \Rightarrow \text{несмещённая.}$$

Пример

Имеется случайная выборка

$$X_1, X_2, \dots, X_n,$$

такая что

$$E(X_i) = \mu.$$

Проверьте на несмещённость оценки для μ :

$$\hat{\mu}_1 = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n};$$

$$\hat{\mu}_2 = X_1;$$

$$\hat{\mu}_3 = X_1 - 3X_2.$$

Решение.

$$E(\hat{\mu}_1) = E\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{E(X_1) + \dots + E(X_n)}{n} = \frac{n\mu}{n} = \mu \quad \Rightarrow \text{несмещённая.}$$

$$E(\hat{\mu}_2) = E(X_1) = \mu \quad \Rightarrow \text{несмещённая.}$$

$$E(\hat{\mu}_3) = E(X_1) - 3E(X_2) = \mu - 3\mu = -2\mu \quad \Rightarrow \text{смещённая.}$$

Пример

Имеется случайная выборка

$$X_1, X_2, \dots, X_n,$$

такая что

$$E(X_i) = \mu.$$

Проверьте на несмещённость оценки для μ :

$$\hat{\mu}_1 = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n};$$

$$\hat{\mu}_2 = X_1;$$

$$\hat{\mu}_3 = X_1 - 3X_2.$$

Решение.

$$E(\hat{\mu}_1) = E\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{E(X_1) + \dots + E(X_n)}{n} = \frac{n\mu}{n} = \mu \quad \Rightarrow \text{несмещённая.}$$

$$E(\hat{\mu}_2) = E(X_1) = \mu \quad \Rightarrow \text{несмещённая.}$$

$$E(\hat{\mu}_3) = E(X_1) - 3E(X_2) = \mu - 3\mu = -2\mu \quad \Rightarrow \text{смещённая.}$$

$$\text{Bias}(\hat{\mu}_3) = E(\hat{\mu}_3) - \mu = -2\mu - \mu = -3\mu.$$

Ещё пример

Разведчик Антон обследовал 25 вражеских танков и передаёт в штаб информацию об их состоянии: 1 — хорошее состояние, 0 — плохое. Из-за помех при передаче возникают ошибки, так что с вероятностью 0.8 штаб получает то, что передаёт Антон, а иначе — либо 0, либо 1 равновероятно.

Пусть p — доля вражеских танков в хорошем состоянии. Для оценивания этого параметра связист Зоя использует величину

$$\tilde{p} = \frac{X_1 + \dots + X_{25}}{20} + c,$$

где X_i — принятое штабом число, характеризующее состояние танка i (или зашумление), c — придуманная Зоей поправка на шум.

Каким должно быть число c , чтобы \tilde{p} была несмещённой оценкой для p ?

Ещё пример

Разведчик Антон обследовал 25 вражеских танков и передаёт в штаб информацию об их состоянии: 1 — хорошее состояние, 0 — плохое. Из-за помех при передаче возникают ошибки, так что с вероятностью 0.8 штаб получает то, что передаёт Антон, а иначе — либо 0, либо 1 равновероятно.

Пусть p — доля вражеских танков в хорошем состоянии. Для оценивания этого параметра связист Зоя использует величину

$$\tilde{p} = \frac{X_1 + \dots + X_{25}}{20} + c,$$

где X_i — принятое штабом число, характеризующее состояние танка i (или зашумление), c — придуманная Зоей поправка на шум.

Каким должно быть число c , чтобы \tilde{p} была несмещённой оценкой для p ?

Решение. Для начала поймём, как распределение X_i связано с параметром p .

$$P(X_i=1) = 0.8 p + 0.2 \cdot 0.5 = 0.8 p + 0.1.$$

$$P(X_i=0) = 0.8(1-p) + 0.2 \cdot 0.5 = 0.9 - 0.8 p.$$

$$X_i \sim \begin{pmatrix} 0 & 1 \\ 0.9 - 0.8 p & 0.1 + 0.8 p \end{pmatrix}$$

$$E(X_i) = 0 \times (0.9 - 0.8 p) + 1 \times (0.1 + 0.8 p) = 0.1 + 0.8 p.$$

Ещё пример (2)

Разведчик Антон обследовал 25 вражеских танков и передаёт в штаб информацию об их состоянии: 1 — хорошее состояние, 0 — плохое. Из-за помех при передаче возникают ошибки, так что с вероятностью 0.8 штаб получает то, что передаёт Антон, а иначе — либо 0, либо 1 равновероятно.

Пусть p — доля вражеских танков в хорошем состоянии. Для оценивания этого параметра связист Зоя использует величину

$$\tilde{p} = \frac{X_1 + \dots + X_{25}}{20} + c,$$

где X_i — принятое штабом число, характеризующее состояние танка i (или зашумление), c — придуманная Зоей поправка на шум.

Каким должно быть число c , чтобы \tilde{p} была несмещённой оценкой для p ?

Решение. Итак, $E(X_i) = 0.1 + 0.8 p$.

Теперь найдём математическое ожидание оценки:

$$E(\tilde{p}) = \frac{E(X_1) + \dots + E(X_{25})}{20} + c = \frac{25 \times (0.1 + 0.8 p)}{20} + c = \frac{2.5 + 20 p}{20} + c = 0.125 + p + c.$$

Чтобы оценка была несмещённой, нужно чтобы $E(\tilde{p}) = p$, так что $c = -0.125$.

Ответ. $c = -0.125$.

Свойства оценок

II. Состоятельность

Последовательность оценок $\{ \hat{\theta}_n \}$ для параметра θ называется *состоятельной*, если

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| > \epsilon) = 0 \quad \forall \epsilon > 0, \quad \forall \theta \in \Theta.$$

Индекс n означает объём выборки:

$$\hat{\theta}_1 = f(X_1), \quad \hat{\theta}_2 = f(X_1, X_2), \dots$$

Слово «последовательность» в дальнейшем будем опускать.

Что это значит?

состоятельность - consistency

Достаточное условие состоятельности

Пусть для оценки $\hat{\theta}_n$ при любом $\theta \in \Theta$ выполняются равенства:

$$1. \lim_{n \rightarrow \infty} E(\hat{\theta}_n) = \theta; \quad (\text{асимптотическая несмещённость})$$

$$2. \lim_{n \rightarrow \infty} D(\hat{\theta}_n) = 0.$$

Тогда $\hat{\theta}_n$ - состоятельная оценка для параметра θ .

А зачем нам всё это нужно?

Что мы хотим от оценки?

Пример

Пусть $X_i \sim \text{i.i.d.}$, $E(X_i) = \mu$, $D(X_i) = \sigma^2 < \infty$.

Рассмотрим оценку для параметра μ :

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n} \quad (\text{выборочное среднее}).$$

Оценка несмещённая: $E(\bar{X}) = \frac{E(X_1) + \dots + E(X_n)}{n} = \frac{n\mu}{n} = \mu.$

Найдём дисперсию:

$$\begin{aligned} D(\bar{X}) &= \frac{1}{n^2} D(X_1 + \dots + X_n) = [\text{в силу независимости } X_i] = \\ &= \frac{D(X_1) + \dots + D(X_n)}{n^2} = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n} \rightarrow 0 \quad \text{при } n \rightarrow \infty. \end{aligned}$$

Значит, \bar{X} - состоятельная оценка для μ .

(Закон Больших Чисел в форме Чебышёва)

Свойства оценок

III. Эффективность.

Несмещённая оценка $\hat{\theta}$ для параметра θ называется *эффективной*, если для любой другой несмещённой оценки $\tilde{\theta}$ параметра θ по той же выборке выполняется неравенство:

$$D(\hat{\theta}) \leq D(\tilde{\theta}) \quad \forall \theta \in \Theta.$$

Почему нам не нравится дисперсия?

Относительная эффективность несмещённых оценок $\hat{\theta}_1$ и $\hat{\theta}_2$:

$$RE(\hat{\theta}_1, \hat{\theta}_2) = \frac{D(\hat{\theta}_2)}{D(\hat{\theta}_1)}.$$

эффективность — efficiency

относительная эффективность — relative efficiency

Свойства оценок

III. Эффективность.

Несмещённая оценка $\hat{\theta}$ для параметра θ называется *эффективной*, если для любой другой несмещённой оценки $\tilde{\theta}$ параметра θ по той же выборке выполняется неравенство:

$$D(\hat{\theta}) \leq D(\tilde{\theta}) \quad \forall \theta \in \Theta.$$

Почему нам не нравится дисперсия?

Как измерить точность оценки?

MSE (Mean Squared Error, средний квадрат ошибки):

$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2].$$

У несмещённых оценок $E(\hat{\theta}) = \theta$, так что

$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = E[(\hat{\theta} - E(\hat{\theta}))^2] = D(\hat{\theta}).$$

Пример

Имеются независимые случайные величины

$$X_1, \dots, X_n,$$

такие что

$$E(X_i) = \mu, \quad D(X_i) = \sigma^2 < \infty$$

Сравните эффективность двух несмещённых оценок для μ :

$$\hat{\mu}_1 = \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n};$$
$$\hat{\mu}_2 = X_1.$$

Пример

Имеются независимые случайные величины

$$X_1, \dots, X_n,$$

такие что

$$E(X_i) = \mu, \quad D(X_i) = \sigma^2 < \infty.$$

Сравните эффективность двух несмещённых оценок для μ :

$$\begin{aligned}\hat{\mu}_1 &= \bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}; \\ \hat{\mu}_2 &= X_1.\end{aligned}$$

Решение. Дисперсии оценок:

$$D(\hat{\mu}_1) = \frac{\sigma^2}{n}; \quad D(\hat{\mu}_2) = D(X_1) = \sigma^2.$$

$$D(\hat{\mu}_1) < D(\hat{\mu}_2) \quad \text{при } n > 0 \Rightarrow \hat{\mu}_1 \text{ эффективнее } \hat{\mu}_2.$$

Относительная эффективность $\hat{\mu}_1$ по сравнению с $\hat{\mu}_2$:

$$RE(\hat{\mu}_1, \hat{\mu}_2) = \frac{D(\hat{\mu}_2)}{D(\hat{\mu}_1)} = n.$$

Упражнение. Докажите полезную формулу:

$$\text{MSE}(\hat{\theta}) = \text{Bias}^2(\hat{\theta}) + D(\hat{\theta}).$$

Ещё упражнение. Попробуйте доказать, что оценка \bar{X} - эффективная в классе линейных несмещённых оценок для μ .

Иначе говоря, покажите, что оценка \bar{X} имеет наименьшую дисперсию

среди всех оценок вида $\tilde{X} = \alpha_1 X_1 + \alpha_2 X_2 + \dots + \alpha_n X_n$,

удовлетворяющих условию $E(\tilde{X}) = \mu$.

Предполагается, что $X_i \sim \text{i.i.d}$, $E(X_i) = \mu$, $D(X_i) = \sigma^2 < \infty$.

Пример: сравнение смещённой и несмещённой оценки

Случайные величины X_1, \dots, X_9 независимы, $E(X_i) = \theta$, $D(X_i) = \theta^2$.

Сравните по MSE оценки для θ :

$$(1) \quad \bar{X} = \frac{X_1 + \dots + X_9}{9},$$

$$(2) \quad \hat{\theta} = \frac{X_1 + \dots + X_9}{10}.$$

Пример: сравнение смещённой и несмещённой оценки

Случайные величины X_1, \dots, X_9 независимы, $E(X_i) = \theta$, $D(X_i) = \theta^2$.

Сравните по MSE оценки для θ :

$$(1) \quad \bar{X} = \frac{X_1 + \dots + X_9}{9},$$

$$(2) \quad \hat{\theta} = \frac{X_1 + \dots + X_9}{10}.$$

Решение. Воспользуемся формулой $MSE(\hat{\theta}) = \text{Bias}^2(\hat{\theta}) + D(\hat{\theta})$.

Найдём смещения оценок:

$$\text{Bias}(\bar{X}) = E(\bar{X}) - \theta = 0, \quad \text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta = \frac{9\theta}{10} - \theta = -\frac{\theta}{10}.$$

Пример: сравнение смещённой и несмещённой оценки

Случайные величины X_1, \dots, X_9 независимы, $E(X_i) = \theta$, $D(X_i) = \theta^2$.

Сравните по MSE оценки для θ :

$$(1) \quad \bar{X} = \frac{X_1 + \dots + X_9}{9},$$

$$(2) \quad \hat{\theta} = \frac{X_1 + \dots + X_9}{10}.$$

Решение. Воспользуемся формулой $MSE(\hat{\theta}) = \text{Bias}^2(\hat{\theta}) + D(\hat{\theta})$.

Найдём смещения оценок:

$$\text{Bias}(\bar{X}) = E(\bar{X}) - \theta = 0,$$

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta = \frac{9\theta}{10} - \theta = -\frac{\theta}{10}.$$

Теперь дисперсии:

$$D(\bar{X}) = \frac{9\theta^2}{9^2} = \frac{\theta^2}{9},$$

$$D(\hat{\theta}) = \frac{9\theta^2}{10^2} = \frac{9\theta^2}{100}.$$

Пример: сравнение смещённой и несмещённой оценки

Случайные величины X_1, \dots, X_9 независимы, $E(X_i) = \theta$, $D(X_i) = \theta^2$.

Сравните по MSE оценки для θ :

$$(1) \quad \bar{X} = \frac{X_1 + \dots + X_9}{9}, \quad (2) \quad \hat{\theta} = \frac{X_1 + \dots + X_9}{10}.$$

Решение. Воспользуемся формулой $MSE(\hat{\theta}) = Bias^2(\hat{\theta}) + D(\hat{\theta})$.

Найдём смещения оценок:

$$Bias(\bar{X}) = E(\bar{X}) - \theta = 0, \quad Bias(\hat{\theta}) = E(\hat{\theta}) - \theta = \frac{9\theta}{10} - \theta = -\frac{\theta}{10}.$$

Теперь дисперсии:

$$D(\bar{X}) = \frac{9\theta^2}{9^2} = \frac{\theta^2}{9}, \quad D(\hat{\theta}) = \frac{9\theta^2}{10^2} = \frac{9\theta^2}{100}.$$

Наконец, MSE:

$$MSE(\bar{X}) = D(\bar{X}) = \frac{\theta^2}{9}, \quad MSE(\hat{\theta}) = \left(-\frac{\theta}{10}\right)^2 + \frac{9\theta^2}{100} = \frac{\theta^2 + 9\theta^2}{100} = \frac{\theta^2}{10}.$$

$$MSE(\bar{X}) \geq MSE(\hat{\theta}) \quad \Rightarrow \text{смещённая оценка оказалась точнее.}$$

Пример: сравнение смещённой и несмещённой оценки

Случайные величины X_1, \dots, X_9 независимы, $E(X_i) = \theta$, $D(X_i) = \theta^2$.

Сравните по MSE оценки для θ :

$$(1) \quad \bar{X} = \frac{X_1 + \dots + X_9}{9},$$

$$(2) \quad \hat{\theta} = \frac{X_1 + \dots + X_9}{10}.$$

Решение. Воспользуемся формулой $MSE(\hat{\theta}) = \text{Bias}^2(\hat{\theta}) + D(\hat{\theta})$.

Найдём смещения оценок:

$$\text{Bias}(\bar{X}) = E(\bar{X}) - \theta = 0,$$

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta = \frac{9\theta}{10} - \theta = -\frac{\theta}{10}.$$

Теперь дисперсии:

$$D(\bar{X}) = \frac{9\theta^2}{9^2} = \frac{\theta^2}{9},$$

$$D(\hat{\theta}) = \frac{9\theta^2}{10^2} = \frac{9\theta^2}{100}.$$

Наконец, MSE:

$$MSE(\bar{X}) = D(\bar{X}) = \frac{\theta^2}{9},$$

$$MSE(\hat{\theta}) = \left(-\frac{\theta}{10}\right)^2 + \frac{9\theta^2}{100} = \frac{\theta^2 + 9\theta^2}{100} = \frac{\theta^2}{10}.$$

$$MSE(\bar{X}) \geq MSE(\hat{\theta}) \quad \Rightarrow \text{смещённая оценка оказалась точнее.}$$

На сегодня хватит.

Теперь организационные вопросы.

Расчёт итоговой оценки

$$O_{\text{итог}} = 0.4 \cdot O_{\text{д/з}} + 0.2 \cdot O_{\text{онлайн-курс}} + 0.4 \cdot O_{\text{экзамен}}.$$

Домашнее задание — от 5 до 8 задач, которые сдаются постепенно в течение семестра с устной защитой.

Онлайн-курс — дам ссылку и объяснения чуть позже.

Экзамен — письменный, не блокирующий.

Автоматы?

Учебники

Пока что рекомендую две основные книги:

J.L. Devore, K.N. Berk «Modern mathematical statistics with applications».

P. Newbold, «Statistics for business and economics».

При надобности буду давать другие источники или присылать подготовленные материалы.

Все лекции и семинары мы стараемся записывать, а потом давать на них ссылки.
Но иногда забываем.

Если видите, что занятие идёт без записи, напоминайте, пожалуйста.

В следующий раз:

некоторые часто используемые точечные оценки

оценки для распределения в целом

описательная статистика и наглядное представление данных