



DS-поток, 3 курс, осень 2022

Статистика



1. Введение

1.1. Основная задача математической статистики

1.2. Вероятностно-статистическая модель

1.3. Виды подходов к статистике

Пример



Введение

Теория вероятностей

Зная природу случайного явления,
посчитать характеристику этого явления.

Математическая статистика

По результатам экспериментальных данных
высказать суждение о том, какова была природа этого явления.



Классический пример

На третьем курсе N студентов; из них M выбирает DS-поток.

Задача в теории вероятностей

$P(\text{среди } n \text{ чел. ровно } m \text{ слушателей DS-потока})=?$

Предполагается, что M известно.

Задача в математической статистике

Среди случайных n чел. есть m слушателей DS-потока.

Оценить M .

Предполагается, что M не известно.



Еще пример

$\xi \sim \mathcal{N}(a, \sigma^2)$ — случайная величина

Задача в теории вероятностей

Известно, что $a = 2.3, \sigma = 7.1$

$$P(\xi \in [0, 1]) - ?$$

$$E\xi - ?$$

Задача в математической статистике

x_1, \dots, x_n — независимые реализации случайной величины ξ .

Оценить a и σ .



Задача математической статистики

Пусть x_1, \dots, x_n — численные характеристики n -кратного повторения некоторого явления.

Будем их воспринимать как независимые реализации $\xi \sim P$.

Задача: по значениям x_1, \dots, x_n высказать некоторое суждение о распределении P .

Решение: *статистический вывод или обучение.*



1. Введение

1.1. Основная задача математической статистики

1.2. Вероятностно-статистическая модель

1.3. Виды подходов к статистике

Пример



Однократный эксперимент

\mathcal{X} — *выборочное пространство* = множество всех возможных значений эксперимента;

$\mathcal{B}_{\mathcal{X}}$ — некоторая σ -алгебра на \mathcal{X} ;

P — некоторое неизвестное распределение на $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$;

Предполагается $P \in \mathcal{P}$ — некоторое семейство распределений.

Вероятностно-статистическая модель

$$(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, \mathcal{P}).$$

Для оперирования с результатами эксперимента как со случайными величинами, определим случайную величину $X : \mathcal{X} \rightarrow \mathcal{X}$ по правилу $X(x) = x \ \forall x \in \mathcal{X}$, которую будем называть *наблюдением*.



Многократный эксперимент

Вероятностно-статистическая модель

$$(\mathcal{X}^n, \mathcal{B}_{\mathcal{X}^n}, \mathcal{P}^n),$$

- ▶ $\mathcal{X}^n = \mathcal{X} \times \dots \times \mathcal{X}$;
- ▶ $\mathcal{B}_{\mathcal{X}^n} = \sigma(B_1 \times \dots \times B_n, B_i \in \mathcal{B}_{\mathcal{X}})$;
- ▶ $\mathcal{P}^n = \{P^n, P \in \mathcal{P}\}$, причем
$$P^n(B_1 \times \dots \times B_n) = P(B_1) \dots P(B_n) \quad \forall B_i \in \mathcal{B}_{\mathcal{X}}.$$

Для любой $P \in \mathcal{P}$ определенная таким образом P^n существует и единственна по теореме о продолжении вероятностной меры.



Многократный эксперимент

Наблюдение, соответствующее i -му эксперименту:

Сл. вел. $X_i : \mathcal{X}^n \rightarrow \mathcal{X}$, т.ч. $X_i(x) = x_i \ \forall x \in \mathcal{X}^n$.

Случайный вектор $X = (X_1, \dots, X_n)$ — *выборка* размера n .

Выборка является вектором независимых
одинаково распределенных случайных величин,
каждая компонента которого имеет распределение P .



Бесконечный эксперимент

Вероятностно-статистическая модель

$$(\mathcal{X}^\infty, \mathcal{B}_{\mathcal{X}^\infty}, \mathcal{P}^\infty),$$

- ▶ $\mathcal{X}^\infty = \mathcal{X} \times \mathcal{X} \times \dots;$
- ▶ $\mathcal{B}_{\mathcal{X}^\infty} = \sigma(B_1 \times \dots \times B_n \times \mathcal{X}^\infty, B_i \in \mathcal{B}_{\mathcal{X}}, n \in \mathbb{N});$
- ▶ $\mathcal{P}^\infty = \{P^\infty, P \in \mathcal{P}\},$ причем

$$P^\infty(B_1 \times \dots \times B_n \times \mathcal{X} \times \dots) = P^n(B_1 \times \dots \times B_n) \forall B_i \in \mathcal{B}_{\mathcal{X}}.$$

Для любой $P \in \mathcal{P}$ определенная таким образом P^n существует и единственна по теореме о продолжении вероятностной меры.



Бесконечный эксперимент

Наблюдение, соответствующее i -му эксперименту:

Сл. вел. $X_i : \mathcal{X}^\infty \rightarrow \mathcal{X}$, т.ч. $X_i(x) = x_i \forall x \in \mathcal{X}^\infty$.

Случайная последовательность $X = (X_1, X_2, \dots)$ —
выборка неограниченного размера.

Выборка неограниченного размера является последовательностью независимых одинаково распределенных случайных величин, каждая компонента которого имеет распределение P .



Далее

- ▶ Для простоты будем опускать индексы n и ∞ ;
- ▶ Будем считать, что в качестве сигма-алгебры используется борелевская, если не сказано обратное;
- ▶ Да и вообще забудем про ее существование :)



1. Введение

1.1. Основная задача математической статистики

1.2. Вероятностно-статистическая модель

1.3. Виды подходов к статистике

Пример



Классификация по типу методов вывода

1. Параметрический

Предполагается, что истинное распределение P принадлежит некоторому классу распределений \mathcal{P} ,
которое параметризовано параметром $\theta \in \Theta$.

$$P \in \{P_\theta \mid \theta \in \Theta\}$$

2. Непараметрический

Предполагается, что истинное распределение P принадлежит некоторому классу распределений \mathcal{P} ,
на котором не введен параметр.



Параметрический подход

1. $X_1, \dots, X_n \sim \text{Exp}(\theta)$, где $\theta > 0$, подразумевает:

$$\mathcal{X} = (0, +\infty), \mathcal{P} = \{\text{Exp}(\theta) \mid \theta > 0\}, \Theta = (0, +\infty).$$

Статистический вывод: указание числа из множества Θ .

2. Схема испытаний Бернулли X_1, \dots, X_n подразумевает:

$$\mathcal{X} = \{0, 1\}, \mathcal{P} = \{\text{Bern}(\theta) \mid 0 \leq \theta \leq 1\}, \Theta = [0, 1].$$

Статистический вывод: указание числа из множества Θ .

3. $X_1, \dots, X_n \sim \mathcal{N}(a, \sigma^2)$, где оба параметра неизвестны:

$$\mathcal{X} = \mathbb{R}, \mathcal{P} = \{\mathcal{N}(a, \sigma^2) \mid a \in \mathbb{R}, \sigma > 0\},$$

$$\theta = (a, \sigma), \Theta = \mathbb{R} \times (0, +\infty).$$

Статистический вывод: указание пары чисел из множества Θ .



Пример

$$\mathcal{X} = \mathbb{R};$$

$$\mathcal{P} = \{U(0, \theta) \mid \theta > 0\};$$

Дана выборка $(X_1, X_2, X_3) = (1, 2, 3)$.

Может ли истинное значение θ быть равным 100 , 3 , 1.5 , -1?

- ▶ 100 и 3 — да;
- ▶ -1 — нет, поскольку $\theta > 0$;
- ▶ 1.5 — да, поскольку $\mathcal{X} = \mathbb{R}$
⇒ возможны любые вещественные числа, правда вероятность получения хотя бы одного числа вне отрезка $[0, \theta]$ равна нулю.



Непараметрический подход

1. В отсутствии предположений:

$$\mathcal{X} = \mathbb{R};$$

\mathcal{P} — все распределения на \mathbb{R} .

В качестве статистического вывода можно некоторым образом оценить функцию распределения;

2. Предполагается, что выборка взята из непрерывного распредел.:

$$\mathcal{X} = \mathbb{R};$$

\mathcal{P} — все непрерывные распределения на \mathbb{R} .

В качестве статистического вывода можно оценить плотность.



- ▶ Любое семейство распределений может рассматриваться как в параметрическом подходе, так и в непараметрическом;
- ▶ Смысл деления на два типа — принципиально разные методы к оценке неизвестного распределения:
 - ▶ Методы параметрического подхода как-либо оценивают параметр, соответствующий неизвестному распределению. На практике обычно $\Theta \subset \mathbb{R}^d$, размерность d фиксирована;
 - ▶ Непараметрические методы пытаются некоторым способом напрямую оценить неизвестное распределение. На практике обычно содержат нефиксированное количество параметров.



Классификация по способу вывода

Два доминирующих подхода:

1. Частотный

Построение суждения о распределении P происходит только на основе выборки X_1, \dots, X_n .

Все распределения из класса \mathcal{P} равноправны.

2. Байесовский

На распределениях из \mathcal{P} задано некоторое распределение ("априорное знание"), которое учитывается при построении суждения, как правило, с помощью формулы Байеса.

Пример: $\mathcal{P} = \{\mathcal{N}(\theta, 1) \mid \theta \in \mathbb{R}\}$,

и предполагается, что истинное значение θ_0 выбрано из $\mathcal{N}(0, 1)$.



1. Введение

1.1. Основная задача математической статистики

1.2. Вероятностно-статистическая модель

1.3. Виды подходов к статистике

Пример



Пример

$X_1, \dots, X_n \sim \mathcal{N}(a, \sigma^2)$, где оба параметра неизвестны.

Строго супер формально: $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathcal{P}^n)$,

где $\mathcal{P} = \{\mathcal{N}(a, \sigma^2) \mid a \in \mathbb{R}, \sigma > 0\}$, $\theta = (a, \sigma)$, $\Theta = \mathbb{R} \times (0, +\infty)$.

Это мы упрощаем до $(\mathbb{R}, \mathcal{B}(\mathbb{R}), \mathcal{P})$, работаем только с \mathbb{R} и \mathcal{P} .

Найдем оценку по методу моментов. Решаем систему

$$\begin{cases} E_{\theta} X_1 = \bar{X}; \\ E_{\theta} X_1^2 = \overline{X^2}. \end{cases} \quad \begin{cases} a = \bar{X}; \\ a^2 + \sigma^2 = \overline{X^2}. \end{cases} \quad \begin{cases} a = \bar{X}; \\ \sigma = \sqrt{\overline{X^2} - a^2}. \end{cases}$$

Получаем оценки

Оценка \hat{a} сильно состоятельна: $\hat{a} = \bar{X} \xrightarrow{P_{\theta-\text{п.н.}}} E_{\theta} \bar{X} = a$.

Строго супер формально:

1. беск. выборка $X_1, X_2, \dots \sim \mathcal{N}(a, \sigma^2)$ на пр-ве $(\mathbb{R}^{\infty}, \mathcal{B}(\mathbb{R}^{\infty}), \mathcal{P}^{\infty})$;
2. посл-ть оценок $\hat{a}_n = \frac{1}{n} \sum_{i=1}^n X_i$ сильно сост., т.к. $\hat{a}_n \xrightarrow{P_{\theta-\text{п.н.}}} a$.



ВСЁ!