

«Оптимизация и математические методы принятия решений»

ст. преп. каф. СС и ПД
Владимиров Сергей Александрович

Лекция 4

Методы математической статистики в задачах принятия решений

СОДЕРЖАНИЕ

Введение

Учебные вопросы :

1. Постановка задачи и общий алгоритм анализа случайных последовательностей при принятии решений с использованием методов математической статистики.
2. Алгоритмы получения эмпирических оценок числовых характеристик, вероятностей и законов распределения случайных последовательностей и анализ их качества.

Заключение

Литература:

1. Гмурман В. Е. Теория вероятностей и математическая статистика . Учебное пособие для вузов – Изд.. 7-е, стер. – М.: Высш. шк. 2001.
2. Вентцель Е.С., Овчаров Л.А. Теория случайных процессов и ее инженерные приложения. Учебное пособие для вузов – Изд.. 2-е, стер. – М.: Высш. шк. 2000.
3. Вентцель Е.С., Овчаров Л.А. Теория вероятностей и ее инженерные приложения. Учебное пособие для вузов – Изд.. 2-е, стер. – М.: Высш. шк. 2000.

Введение

На прошлой лекции были рассмотрены случайные факторы, их распределения и параметры оценок величин, то есть что можно получить, имея такие данные. На этой лекции мы рассмотрим вопрос — как это делать и каким математическим аппаратом производить анализ, оценку и оптимизацию данных в задачах принятия решений.

При принятии решений проблемы с оценкой получаемых данных подразделяются на три класса:

1) хорошо структурированные или количественно сформулированные данные и соответственно анализ проблемы, в которых получают численные оценки;

2) неструктурированные или качественно выраженные данные и отсюда проблемы, в которых количественные зависимости между признаками и характеристиками совершенно неизвестны;

3) слабо структурированные или смешанные данные и следовательно проблемы, содержащие как количественные, так и качественные элементы, причем последние имеют тенденцию к доминированию.

Постановка задачи и общий алгоритм анализа случайных последовательностей с использованием методов математической статистики.

Наиболее полным описанием случайной последовательности является функция распределения вероятностей ее значений и задача анализа в общем случае сводится к получению эмпирических вероятностных характеристик по доступным выборочным данным и проверке гипотез о их соответствии некоторым стандартным характеристикам, определяющим различные классы случайных последовательностей и отдельные их свойства. Часто в качестве стандартной случайной последовательности (СП) X выступает последовательность, например, с нормальным распределением $N(M_X; D_X)$ и вычисляемыми числовыми характеристиками:

- M_X - математическое ожидание;
- D_X - дисперсия случайной последовательности.

Общий алгоритм анализа случайной последовательности с учетом вводимой стандартной случайной последовательности может включать следующие этапы.

1. Определение эмпирических вероятностных характеристик анализируемой случайной последовательности (математического ожидания, дисперсии, корреляционного момента, вероятностей событий и функции распределения вероятностей). Важно, чтобы качество полученных эмпирических оценок соответствовало выдвигаемым априорно требованиям к допустимому отклонению от истинных значений характеристик (доверительному интервалу и доверительной вероятности), а также определялось требуемым для этого размером выборки. На основе полученных характеристик могут быть установлены свойства симметрии распределения (совпадение значений среднего, моды и медианы, либо равенство значений вероятностей превышения и

не превышения среднего значения) и близости его формы к некоторому стандартному (см. лекцию 3), например, к нормальному.

2. Построение гистограммы вероятностей и восстановление эмпирического распределения случайной последовательности на основе полученных вероятностных характеристик и выдвижение гипотезы о виде распределения СП.

3. Проверка верности выдвинутой гипотезы по критериям соответствия эмпирических и аналитических вероятностных характеристик, а также определение класса и основных свойств случайной последовательности с оценкой показателей качества полученных оценок и решений.

Основные этапы анализа случайных последовательностей в предположении выполнения условия стационарности выборочных данных.

Вероятностной характеристикой θ случайной величины X , определяемой непосредственно путем эксперимента, является некоторое число - математическое ожидание, дисперсия, вероятность события $\alpha < X < \beta$. Символ θ означает истинное значение характеристики. Путем обработки результатов экспериментального исследования X получают экспериментальное значение характеристики $\bar{\theta}$, как статистическую характеристику или как оценку $\tilde{\theta}$ характеристики θ .

Экспериментальное исследование случайной величины X с целью определения $\tilde{\theta}$ - оценки (приближенного значения) θ , заключается в проведении N **опытов** (испытаний, наблюдений) и получении (путем соответствующих измерений) ряда значений $x_1, \dots, x_i, \dots, x_N$ — реализаций X . В результате обработки экспериментальных данных определяется $\tilde{\theta} = \Psi_{\theta}(x_1, \dots, x_i, \dots, x_N)$ как функция эксперимента.

Если провести *еще одну серию из N опытов*, то будет получен следующий ряд других реализаций $x'_1, \dots, x'_i, \dots, x'_N$ случайной величины X и другое значение $\tilde{\theta}'$ оценки искомой характеристики θ . Значение x_i случайной величины X , полученное в результате i -ого опыта в серии, можно рассматривать как значение случайной величины X_i а оценку $\tilde{\theta}$ - как реализацию более общей случайной величины

$$\tilde{\Theta} = \Psi_{\theta}(X_1, \dots, X_i, \dots, X_N), \quad (1)$$

Вероятностными характеристиками системы двух случайных величин (X, Y) , определяемыми непосредственно на основании эксперимента, являются математические ожидания, дисперсии, корреляционный момент, вероятность события $\alpha_x < X < \beta_x, \alpha_y < Y < \beta_y$. Эксперимент заключается в проведении N опытов и получении ряда значений $(x_1, y_1), \dots, (x_i, y_i), \dots, (x_N, y_N)$ реализаций случайных величин X, Y . В результате обработки экспериментальных данных получается оценка

$$\tilde{\theta} = \Psi_{\theta}((x_1, y_1), \dots, (x_i, y_i), \dots, (x_N, y_N)),$$

как реализация случайной функции аналогичной (1).

Полученная в результате аналитическая или гипотетическая функция

$$\tilde{\Theta} = \Psi_{\theta}(X_1, Y_1; \dots, X_i, Y_i; \dots, X_N, Y_N), \quad (2)$$

Погрешность приближения оценки $\tilde{\Theta}$ к θ равная

$$\Delta\tilde{\Theta} = \tilde{\Theta} - \theta, \quad (3)$$

является, как и $\tilde{\Theta}$, случайной величиной.

Функцию Ψ_{θ} желательно выбирать так, чтобы выполнялось три условия

1. Математическое ожидание $\Delta\tilde{\Theta}$ равно нулю:

$$M[\Delta\tilde{\Theta}] = 0. \quad (4)$$

2. Дисперсия $\Delta\tilde{\Theta}$ стремится к нулю с увеличением N

$$\lim_{N \rightarrow \infty} D(\Delta\tilde{\Theta}) = 0. \quad (5)$$

3. Дисперсия $D(\Delta\tilde{\Theta})$ при данной Ψ_{θ} должна быть наименьшей.

$$D(\Delta\tilde{\Theta}) \Rightarrow \min$$

Определения

При выполнении условия (4) **оценка** $\tilde{\theta}$ называется **несмещенной**,
условий (4), (5) - **состоятельной**,
всех трех условий - **эффективной**.

Вследствие случайного характера погрешности (3) для характеристики точности приближенного равенства $\tilde{\theta} = \theta$ необходимо располагать вероятностью p_Δ того, что абсолютное значение погрешности не превзойдет некоторого предела

$$P(|\Delta\tilde{\theta}| \leq \Delta_\theta) = p_\Delta. \quad (6)$$

Интервал от $\tilde{\theta} - \Delta_\theta$ до $\tilde{\theta} + \Delta_\theta$, в котором с вероятностью p_Δ находится истинное значение θ , называется **доверительным интервалом**, его границы - **доверительными границами**, а вероятность p_Δ - **доверительной вероятностью**.

Если число экспериментальных данных N достаточно велико, то погрешность (3) состоятельной оценки $\tilde{\theta}$ можно практически считать распределенной нормально с математическим ожиданием (4), дисперсией $D(\Delta\tilde{\theta})=D(\tilde{\theta})=D_{\theta}$ и средним квадратичным отклонением $\sigma(\Delta\tilde{\theta})=\sigma(\tilde{\theta})=\sigma_{\theta}=\sqrt{D_{\theta}}$. При этом выражение (6) имеет

вид:

$$p_{\partial} = \Phi\left(\frac{\Delta_{\tilde{\theta}}}{\sigma_{\tilde{\theta}}}\right) - \Phi\left(-\frac{\Delta_{\tilde{\theta}}}{\sigma_{\tilde{\theta}}}\right) = 2\Phi\left(\frac{\Delta_{\tilde{\theta}}}{\sigma_{\tilde{\theta}}}\right) = 2\Phi(t), \quad (7)$$

где $\Phi(t)$ - функция Лапласа, $t = \frac{\Delta_{\tilde{\theta}}}{\sigma_{\tilde{\theta}}}$.

С помощью этой формулы решается задача определения доверительной вероятности p_{∂} по известным данным $\Delta_{\tilde{\theta}}, \sigma_{\tilde{\theta}}$.

Функция Лапласа $\tau = \Phi(t)$ выражает зависимость τ от t . Обратная $t = \Phi^{-1}(\tau)$ выражает зависимость t от τ . При $t = \frac{\Delta_{\tilde{\theta}}}{\sigma_{\tilde{\theta}}}$, $\tau = \frac{p_{\partial}}{2}$ имеем

$$\frac{\Delta_{\tilde{\theta}}}{\sigma_{\tilde{\theta}}} = \Phi^{-1}(\tau). \quad (8)$$

С помощью формулы (8) и обратной функции Лапласа решается задача определения доверительного интервала $\Delta_{\tilde{\theta}}$ по известным p_{∂} и $\sigma_{\tilde{\theta}}$ и необходимого числа испытаний по известным p_{∂} и $\Delta_{\tilde{\theta}}$.

При решении первой задачи согласно (8) определяется $\Delta_{\tilde{\theta}}$. При решении второй задачи согласно (8) определяется $\sigma_{\tilde{\theta}}$, а затем N .

Для проведения анализа СП обычно приводят к стандартному виду. Для случая двоичной нуль - единичной последовательности это достигается перекодировкой исходной последовательности в симметричную $-1,1$ -ю последовательность в соответствии с правилом

$$x_i(k) = 2x_i'(k) - 1.$$

Здесь $x_i(k), x_i'(k)$ - элементы стандартной и исходной последовательностей соответственно.

Алгоритмы получения эмпирических оценок числовых характеристик, вероятностей и законов распределения случайных последовательностей и анализ их качества.

Определение математического ожидания

$$\hat{M}_x = \frac{1}{N} \sum_{i=1}^N x_i, \quad \tilde{M}_x = \frac{1}{N} \sum_{i=1}^N X_i, \quad (9)$$

где X_i - независимые случайные величины с одинаковыми $M_{x_i} = M_x$ и $D_{x_i} = D_x$.

Математическое ожидание погрешности оценки среднего равно

$$M[\Delta \tilde{M}_x] = M[\tilde{M}_x - M_x] = M\left[\frac{1}{N} \sum_i^N X_i\right] - M_x = \frac{1}{N} NM_x - M_x = 0. \quad (10)$$

Дисперсия погрешности оценки среднего равна

$$D(\Delta \tilde{M}_x) = D(\tilde{M}_x) = D_{\tilde{M}_x} = D\left(\frac{1}{N} \sum_i^N X_i\right) = \frac{1}{N^2} D\left(\sum_i^N X_i\right) = \frac{1}{N^2} ND_x = \frac{D_x}{N} = \frac{\sigma_x^2}{N}. \quad (11)$$

Среднеквадратическое отклонение оценки математического ожидания (9)

$$\sigma_{\tilde{M}_x} = \sigma_x \sqrt{N} \approx \tilde{\sigma}_x \sqrt{N}. \quad (12)$$

Оценка (9) – несмещенная, состоятельная и эффективная.

Определение оценки дисперсии и ее среднеквадратического отклонения

$$\text{Оценка дисперсии } D_x \quad D_x = \frac{1}{N} \sum_i^N (x_i - M_x)^2 .$$

Так как значение M_x априори неизвестно, то принимают $M_x \approx \tilde{M}_x$ и тогда

$$\tilde{D}_x = \frac{1}{N} \sum_i^N X_i^2 - \left(\frac{1}{N} \sum_i^N X_i \right)^2 . \quad (13)$$

Математическое ожидание погрешности оценки равно

$$M[\Delta \tilde{D}_x] = M[\tilde{D}_x - D_x] = -\frac{D_x}{N} , \quad (14)$$

что означает, что оценка (14) является *смещенной*.

Смещение пропорционально D_x и обратно пропорционально N . Это означает, что оценка D_x , полученная согласно (14), - *состоятельная*.

Смещение устраняется с переходом к $\tilde{D}'_x = \frac{N}{N-1} \tilde{D}_x$.

$$\text{При этом вместо (13) имеем } \tilde{D}'_x = \frac{1}{N-1} \sum_i^N x_i^2 - \frac{N}{N-1} \tilde{M}_x^2 . \quad (15)$$

При больших значениях N результаты расчета по формулам (13) и (15) практически будут одинаковыми.

Зависимость среднего квадратического отклонения $\sigma_{\tilde{D}_x}$ от его точного значения σ_x определяется выражением

$$\sigma_{\tilde{D}_x} \approx \sqrt{\frac{2}{N-1}} \sigma_x .$$

Определение корреляционного момента и коэффициента корреляции

Экспериментальное значение корреляционного момента R_{xy} как оценка смешанного центрального момента m_{11} системы двух случайных величин равно

$$R_{xy} = \frac{1}{N} \sum_1^N (x_i - M_x)(y_i - M_y). \quad (16)$$

Так как значения M_x , M_y неизвестны, то принимают $M_x \approx \tilde{M}_x$, $M_y \approx \tilde{M}_y$ и тогда

$$\tilde{R}_{xy} = \frac{1}{N} \sum_1^N (x_i - \tilde{M}_x)(y_i - \tilde{M}_y) = \frac{1}{N} \sum_1^N x_i y_i - \left(\frac{1}{N} \sum_1^N \bar{x}_i \right) \left(\frac{1}{N} \sum_1^N \bar{y}_i \right)$$

или

$$\tilde{R}_{xy} = \frac{1}{N} \sum_1^N X_i Y_i - \left(\frac{1}{N} \sum_1^N \bar{X}_i \right) \left(\frac{1}{N} \sum_1^N \bar{Y}_i \right). \quad (17)$$

Погрешность оценки \tilde{R}_{xy}

$$\Delta \tilde{R}_{xy} = \tilde{R}_{xy} - R_{xy} \quad (18)$$

Математическое ожидание погрешности (18)

$$M[\Delta \tilde{R}_{xy}] = M\left[\frac{1}{N} \sum_1^N X_i Y_i\right] - M\left[\left(\frac{1}{N} \sum_1^N \bar{X}_i\right)\left(\frac{1}{N} \sum_1^N \bar{Y}_i\right)\right] = -\frac{R_{xy}}{N}$$

Это означает, что оценка (17) - смещена и равна

$$\tilde{R}_{xy} = \frac{N-1}{N} R_{xy} . \quad (19)$$

Можно показать, что она является и состоятельной.

Смещение устраняется с переходом от \tilde{R}_{xy} к $\tilde{R}'_{xy} = \frac{N-1}{N} \tilde{R}_{xy}$.

При этом вместо (17) имеем
$$\tilde{R}_{xy} = \frac{1}{N-1} \sum_1^N x_i y_i - \frac{N}{N-1} \tilde{M}_x \tilde{M}_y . \quad (20)$$

Среднеквадратическое значение погрешности (18) равно среднему квадратическому отклонению оценки (20):

$$\sigma_{\tilde{R}_{xy}} \approx \sqrt{(\tilde{R}_{xy}^2 + D_x D_y) / (N-1)} . \quad (23)$$

Оценка коэффициента корреляции определяется согласно

$$\tilde{r}_{xy} = \tilde{R}_{xy} / (\tilde{\sigma}_x \tilde{\sigma}_y) . \quad (24)$$

Определение вероятности события через его повторяемость

Экспериментальное значение вероятности P некоторого события - это его повторяемость [1-3]

$$W = \frac{n}{N} = \tilde{P} = \frac{1}{N} \sum_1^N x_i, \quad (26)$$

причем число n появлений события в серии из N испытаний можно рассматривать как сумму N независимых случайных слагаемых:

$$W = \frac{1}{N} \sum_1^N X_i, \quad (27)$$

каждое из которых может принимать только два значения 1 и 0 с вероятностями P и $1 - P$.

Математическое ожидание и дисперсия случайной величины X_i :

$$M_{x_i} = P; \quad D_{x_i} = P(1 - P). \quad (28)$$

Погрешность оценки (26) равна

$$\Delta W = W - P. \quad (29)$$

Математическое ожидание погрешности и ее дисперсия:

$$M[\Delta W] = 0; \quad D(\Delta W) = P(1 - P)/N = D(W). \quad (30)$$

Таким образом, оценка (26) - несмещенная и состоятельная. Среднее квадратическое отклонение оценки (26)

$$\sigma_w = \sqrt{P(1 - P)/N}.$$

На практике принимают

$$\sigma_w \approx \sqrt{W(1 - W)/N}. \quad (31)$$

Определение законов распределения случайной величины (строим гистограмму).

Если случайная величина X - дискретная, то определяются \tilde{M}_x , \tilde{D}_x и оценки \tilde{P}_{x_i} значений функции вероятности $P(x_i)$ или оценки $\tilde{F}(x_i)$ значений функции распределения $F(x_i)$.

Если случайная величина X - непрерывная, то определяются M_x , D_x и оценки плотности вероятности $f_x(x)$ и функции распределения $F_x(x)$.

При оценивании законов распределения непрерывной случайной величины процесс обработки экспериментальных данных - реализаций x_1, \dots, x_N , начинается с выбора границ a и $b > a$ интервала, заключающего возможные значения X , и деления этого интервала на k равных элементарных промежутков $c = (b - a) / k$.

При расчете c значения a и b следует для удобства округлять, принимая, например, вместо $b = 3,341$, $a = -2,63$ значения $3,4$ и $-2,7$. Во всех случаях округление производится в сторону увеличения разности $b-a$. Значение k выбирается в пределах от 8 до 20. Удобно принять $k=10$.

После этого определяют границы x'_j всех элементарных промежутков и составляют таблицу (табл.1), в которой $x'_0 = a$, $x'_k = b$. Значение \tilde{n}_j - это число реализаций X , оказавшихся в пределах j -ого интервала от x'_{j-1} , до x'_j . Значения \tilde{P}_j и \tilde{F}_j :

$$\tilde{P}_j = \tilde{n}_j / N \approx P(x'_{j-1} < X < x'_j); \quad (32)$$

$$\tilde{F}_j = \tilde{P}_1 + \dots + \tilde{P}_{j-1} = P(X < \tilde{x}_j). \quad (33)$$

При группировке реализаций X по отдельным интервалам может оказаться что некоторые из них придутся точно на границу двух смежных промежутков. В этих случаях необходимо прибавить к числам \tilde{n}_j и \tilde{n}_{j+2} смежных интервалов по 1/2.

Таблица 1

| x_i | x'_0 | x'_1 | x'_2 | ... | x'_{k-1} | x'_k |
|---------------|--------|---------------|---------------|-----|---------------|--------|
| n_j | | \tilde{n}_1 | \tilde{n}_2 | ... | \tilde{n}_k | |
| \tilde{P}_j | | \tilde{P}_1 | \tilde{P}_2 | ... | \tilde{P}_k | |
| \tilde{F}_j | | \tilde{F}_1 | \tilde{F}_2 | ... | \tilde{F}_k | |

По данным таблицы могут быть построены эмпирические гистограмма и график функции распределения.

Затем возникает весьма сложная задача подбора аналитического закона распределения, достаточно хорошо согласующегося с результатами эксперимента. Основанием для выбора аналитического выражения плотности вероятности $f_x(x)$ могут служить соображения о том, чтобы простейшие числовые характеристики теоретической случайной величины были равны экспериментальным значениям этих характеристик. Если, например, теоретический закон определяется двумя параметрами, то их выбирают так, чтобы совпали два момента (m_1, m_2).

Для оценки существенности или несущественности расхождения между теоретическим и эмпирическим распределениями используют различные критерии

Критерий интервальных оценок

Располагая результатами эксперимента согласно (31) рассчитывают средние квадратические отклонения: $\sigma_{\tilde{P}_j} = \sqrt{\tilde{P}_j(1-\tilde{P}_j)N}$; $\sigma_{\tilde{F}_j} = \sqrt{\tilde{F}_j(1-\tilde{F}_j)N}$. (34)

Согласно (8) рассчитываются доверительные интервалы $\Delta_{\tilde{P}_j} = 3\sigma_{\tilde{P}_j}$; $\Delta_{\tilde{F}_j} = 3\sigma_{\tilde{F}_j}$ и границы изменения ВВХ $\tilde{P}_j \pm \Delta_{\tilde{P}_j}$; $\tilde{F}_j \pm \Delta_{\tilde{F}_j}$, (35)

соответствующие доверительной вероятности $p_0 = 0,9972$ и $\Phi^{-1}(0,4986) = 3$.

Располагая выбранным аналитическим выражением плотности вероятности

$f_x(x)$, рассчитываются теоретические значения: $P_j = P(x'_{j=1} < X < x'_j) = \int_{x'_{j=1}}^{x'_j} f_x(x) dx$; (36)
 $F_j = P(x'_{j=1} < X < x'_j) = \int_{x'_{j=1}}^{x'_j} f_x(x) dx$

Критерием согласия теоретического и экспериментального распределения

является соблюдение неравенств:

$$\begin{aligned} \tilde{P}_j - \Delta_{\tilde{P}_j} < P_j < \tilde{P}_j + \Delta_{\tilde{P}_j}; \\ \tilde{F}_j - \Delta_{\tilde{F}_j} < F_j < \tilde{F}_j + \Delta_{\tilde{F}_j}. \end{aligned} \quad (37)$$

Критерий χ^2

Рассчитав P_j согласно (35), находят значения $n_j = NP_j$ (38)

и рассчитывают

$$\chi^2 = \sum_1^k (\tilde{n}_j - n_j)^2 / n_j. \quad (39)$$

Если расхождение между экспериментальным и теоретическим распределением несущественно, то распределение случайной величины (39) близко к нормальному с математическим ожиданием $M_{\chi^2} = s$ и средним квадратическим отклонением $\sigma_{\chi^2} = \sqrt{2s}$, где s - так называемое *число степеней свободы* и согласно (8) с доверительной вероятностью $p_0 = 0,997$ справедливо неравенство

$$(\chi^2 - s) / \sqrt{2s} < 3. \quad (40)$$

Число степеней свободы $s = k - u$ - это разность между числом интервалов k , выбираемых произвольно, и числом условий u , которым должно удовлетворять эмпирическое распределение случайной величины. Этих условий обычно три: сумма всех \tilde{P}_j равна единице, математическое ожидание равно \tilde{M}_j , дисперсия равна \tilde{D}_x .

Заключение.

В заключение отметим, что все используемые в настоящее время методы анализа случайных последовательностей не выходят за рамки представленного общего подхода, однако в некоторых случаях позволяют существенно упростить процедуры их классификации если учитывают специфику анализируемых последовательностей.