

Решение пробного итогового теста по дисциплине
«Интеллектуальный анализ данных»
для направлений НФИ, НПИ, НБИ

7 октября 2024 г.

1. Подготовка данных

Дан набор данных с числовыми признаками. Вычислить:

- математическое ожидание, медиану, дисперсию заданного признака

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i, \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

$$m = \begin{cases} x_{k+1}, & n = 2k + 1 \\ \frac{1}{2}(x_k + x_{k+1}), & n = 2k \end{cases}, x_1 \leq x_2 \leq \dots \leq x_n.$$

- значение эмпирической CDF признака в заданной точке

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x)$$

- ковариацию (корреляцию) между заданными признаками

$$\hat{\sigma}_{12} = \frac{1}{n} \sum_{i=1}^n (x_{i1} - \hat{\mu}_1)(x_{i2} - \hat{\mu}_2)$$
$$\hat{\rho}_{12} = \frac{\hat{\sigma}_{12}}{\hat{\sigma}_1 \hat{\sigma}_2} = \frac{\sum_{i=1}^n (x_{i1} - \hat{\mu}_1)(x_{i2} - \hat{\mu}_2)}{\sqrt{\sum_{i=1}^n (x_{i1} - \hat{\mu}_1)^2 \sum_{i=1}^n (x_{i2} - \hat{\mu}_2)^2}}$$

Задание 1. Дан следующий набор данных с двумя числовыми признаками:

$$D = \{(1, 1), (3, 4), (2, 3), (0, 4W)\},$$

где W – некоторое числовое значение.

Вычислить медиану первого признака и эмпирическую CDF второго признака в точке $x = 3.5$ при $W = \frac{1}{2}$.

Решение: Упорядочим значения первого признака по возрастанию и получим последовательность 0, 1, 2, 3, поэтому медиана равна среднему арифметическому двух значений в середине последовательности:

$$\frac{1+2}{2} = 1.5$$

Упорядочим значения второго признака по возрастанию и получим последовательность 1, 2, 3, 4. Для точки $x = 3.5$ три элемента последовательности из четырех меньше или равны x , поэтому эмпирическая CDF равна

$$\frac{3}{4} = 0.75$$

2. Наивный байесовский классификатор

Задание 2. Дан размеченный набор данных с одним признаком (первая компонента) и метками классов (вторая компонента):

$$\mathbf{D} = \{(0, 1), (1, 1), (2, 1), (4, 2), (6, 2)\}$$

Используя (наивный) байесовский подход, найти прогнозируемый класс для заданной точки $x = 3$

Решение: Определяем априорные вероятности классов c_1 и c_2

$$n = |\mathbf{D}| = 5,$$

$$n_1 = |\mathbf{D}_1| = |\{(0, 1), (1, 1), (2, 1)\}| = 3,$$

$$n_2 = |\mathbf{D}_2| = |\{(4, 2), (6, 2)\}| = 2,$$

$$\mathbb{P}[c_1] = \frac{n_1}{n} = \frac{3}{5} = 0.6, \mathbb{P}[c_2] = \frac{n_2}{n} = \frac{2}{5} = 0.4$$

Вычисляем математические ожидания признака для классов c_1 и c_2 :

$$\mu_1 = \mathbb{E}[X | c_1] = \frac{1}{3}(0 + 1 + 2) = 1, \mu_2 = \mathbb{E}[X | c_2] = \frac{1}{2}(4 + 6) = 5$$

и дисперсии признака для классов c_1 и c_2 :

$$\sigma_1^2 = \frac{1}{3} \left((0 - 1)^2 + (1 - 1)^2 + (2 - 1)^2 \right) = \frac{2}{3} = 0.667,$$

$$\sigma_2^2 = \frac{1}{2} \left((4 - 5)^2 + (6 - 5)^2 \right) = \frac{2}{2} = 1$$

Для точки $x = 3$ апостериорные вероятности классов c_1 и c_2 будут пропорциональны значениям:

$$p_1 = \frac{1}{\sigma_1} \exp \left\{ -\frac{(x - \mu_1)^2}{2\sigma_1^2} \right\} \mathbb{P}[c_1] = \frac{1}{\sqrt{\frac{2}{3}}} \exp \left\{ -\frac{(3 - 1)^2}{2 \cdot \frac{2}{3}} \right\} \frac{3}{5} = \frac{\sqrt{3}}{5\sqrt{2}} \exp \{-3\}$$

$$p_2 = \frac{1}{\sigma_2} \exp \left\{ -\frac{(x - \mu_2)^2}{2\sigma_2^2} \right\} \mathbb{P}[c_2] = \frac{1}{\sqrt{1}} \exp \left\{ -\frac{(3 - 5)^2}{2} \right\} \frac{2}{5} = \frac{2}{5} \exp \{-2\}$$

Так как

$$\frac{p_1^2}{p_2^2} = \frac{\frac{3}{25 \cdot 2} \exp \{-6\}}{\frac{2}{25} \exp \{-4\}} = \frac{3}{4} \exp \{-2\} < 1,$$

для заданной точки $x = 3$ будет спрогнозирован класс 2.

3. Поиск ассоциативных правил

Дана база транзакций \mathbf{D} . Для заданного набора предметов X найти:

- поддержку (support)

$$\sup(X, \mathbf{D}) = |\{t \mid \langle t, \mathbf{i}(t) \rangle \in \mathbf{D}, X \subseteq \mathbf{i}(t)\}| = |\mathbf{t}(X)|$$

- относительную поддержку (relative support)

$$\text{rsup}(X, \mathbf{D}) = \frac{\sup(X, \mathbf{D})}{|\mathbf{D}|}$$

Для заданного ассоциативного правила $X \rightarrow Y$ найти:

- поддержку (support)

$$s = \sup(X \rightarrow Y) = |\mathbf{t}(XY)| = \sup(XY)$$

- относительную поддержку (relative support)

$$\text{rsup}(X \rightarrow Y) = \frac{\sup(XY)}{|\mathbf{D}|} = \mathbb{P}[X \wedge Y]$$

- достоверность (confidence)

$$c = \text{conf}(X \rightarrow Y) = \mathbb{P}[Y \mid X] = \frac{\mathbb{P}[X \wedge Y]}{\mathbb{P}[X]} = \frac{\sup(XY)}{\sup(X)}$$

- лифт (lift)

$$\text{lift}(X \rightarrow Y) = \frac{\mathbb{P}[XY]}{\mathbb{P}[X]\mathbb{P}[Y]} = \frac{\text{rsup}(XY)}{\text{rsup}(X)\text{rsup}(Y)} = \frac{\text{conf}(X \rightarrow Y)}{\text{rsup}(Y)}$$

- рычаг (leverage)

$$\text{leverage}(X \rightarrow Y) = \mathbb{P}[XY] - \mathbb{P}[X]\mathbb{P}[Y] = \text{rsup}(XY) - \text{rsup}(X)\text{rsup}(Y)$$

- убежденность (conviction)

$$\text{conv}(X \rightarrow Y) = \frac{\mathbb{P}[X]\mathbb{P}[\neg Y]}{\mathbb{P}[X \neg Y]} = \frac{1}{\text{lift}(X \rightarrow \neg Y)} = \frac{1 - \text{rsup}(Y)}{1 - \text{conf}(X \rightarrow Y)}$$

Задание 3. Дана база транзакций

$$\mathbf{D} = \{(1, ABDE), (2, BCE), (3, ABDE), (4, ABCE), (5, ABCDE), (6, BCD)\}$$

Для набора предметов $\{BC\}$ вычислить поддержку (support).

Решение: По определению поддержки набора предметов

$$\sup(X, \mathbf{D}) = |\{t \mid \langle t, \mathbf{i}(t) \rangle \in \mathbf{D}, X \subseteq \mathbf{i}(t)\}| = |\mathbf{t}(X)|$$

Поэтому

$$\sup(\{BC\}, \mathbf{D}) = |\{(2, BCE), (4, ABCE), (5, ABCDE), (6, BCD)\}| = 4$$

4. Кластеризация данных

Дано разбиение набора данных с числовыми признаками на k кластеров. Найти расстояние между кластерами: • методом одиночной связи (single link)

$$\delta(C_i, C_j) = \min_{\bar{\mathbf{x}}, \bar{\mathbf{y}}} \{ \delta(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mid \bar{\mathbf{x}} \in C_i, \bar{\mathbf{y}} \in C_j \}$$

- методом полной связи (complete link)

$$\delta(C_i, C_j) = \max_{\bar{\mathbf{x}}, \bar{\mathbf{y}}} \{ \delta(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mid \bar{\mathbf{x}} \in C_i, \bar{\mathbf{y}} \in C_j \}$$

- методом средней связи (group average)

$$\delta(C_i, C_j) = \frac{1}{n_i n_j} \sum_{\bar{\mathbf{x}} \in C_i} \sum_{\bar{\mathbf{y}} \in C_j} \delta(\bar{\mathbf{x}}, \bar{\mathbf{y}})$$

- методом Уарда (Ward's measure)

$$\delta(C_i, C_j) = \Delta SSE_{ij} = SSE_{ij} - SSE_i - SSE_j, \quad SSE_i = \sum_{\bar{\mathbf{x}} \in C_i} \|\bar{\mathbf{x}} - \bar{\mu}_i\|^2$$

$$\delta(C_i, C_j) = \frac{n_i n_j}{n_i + n_j} \|\bar{\mu}_i - \bar{\mu}_j\|^2$$

Задание 4. Дано разбиение набора данных с двумя числовыми признаками на два кластера:

- кластер 1 (C_1): $\{(0, 0), (0, 1), (1, 0)\}$
- кластер 2 (C_2): $\{(2, 2), (3, 3)\}$

Найти евклидово расстояние между кластерами методом одиночной связи (single link).

Решение: Расстояние между кластерами методом одиночной связи (single link) вычисляется по формуле

$$\delta(C_i, C_j) = \min_{\bar{\mathbf{x}}, \bar{\mathbf{y}}} \{ \delta(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \mid \bar{\mathbf{x}} \in C_i, \bar{\mathbf{y}} \in C_j \}$$

Евклидово расстояние вычисляется по формуле

$$\rho(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}, \quad \bar{\mathbf{x}} = (x_1, \dots, x_d), \quad \bar{\mathbf{y}} = (y_1, \dots, y_d)$$

Наиболее близкими точками двух кластеров являются точка $(0, 1)$ (или точка $(1, 0)$) и точка $(2, 2)$, поэтому

$$\delta(C_1, C_2) = \rho((0, 1), (2, 2)) = \sqrt{2^2 + 1^2} = \sqrt{5} = 2.236$$

5. Кластеризация данных

Дана таблица (матрица) сопряженности кластеризации. Вычислить:

- чистоту (purity) или точность (precision) кластеров и чистоту кластеризации

$$p_j = prec_j = \frac{1}{m_j} \max_{i=1, \overline{k}} n_{ij} = \frac{n_{ij}}{m_j}, i_j = \max_{i=1, \overline{k}} n_{ij}, j = \overline{1, r}$$

$$p = \sum_{j=1}^r \frac{m_j}{n} p_j = \frac{1}{n} \sum_{j=1}^r \max_{i=1, \overline{k}} n_{ij},$$

- полнота (recall) кластеров

$$recall_j = \frac{n_{ij}}{|T_{ij}|} = \frac{n_{ij}}{n_{ij}}, n_{ij} = |T_{ij}|, j = \overline{1, r},$$

- F-меру кластеров и F-меру кластеризации

$$F_j = \frac{2}{\frac{1}{prec_j} + \frac{1}{recall_j}} = \frac{2n_{ij}}{n_{ij} + m_j}, j = \overline{1, r}, F = \frac{1}{r} \sum_{j=1}^r F_j$$

- энтропию кластеризации (разбиения на кластеры)

$$H(\mathcal{C}) = - \sum_{j=1}^r p_{C_j} \log_2 p_{C_j}, p_{C_j} = \frac{m_j}{n}$$

- энтропию разбиения на классы

$$H(\mathcal{T}) = - \sum_{i=1}^k p_{T_i} \log_2 p_{T_i}, p_{T_i} = \frac{n_i}{n}$$

- условную энтропию относительно кластера

$$H(\mathcal{T} | \mathcal{C}_j) = - \sum_{i=1}^k \left(\frac{n_{ij}}{n_i} \right) \log_2 \left(\frac{n_{ij}}{n_i} \right)$$

- условную энтропию относительно кластеризации

$$H(\mathcal{T} | \mathcal{C}) = \sum_{j=1}^r \frac{m_j}{n} H(\mathcal{T} | \mathcal{C}_j) = - \sum_{j=1}^r \sum_{i=1}^k p_{ij} \log_2 \left(\frac{p_{ij}}{p_{C_j}} \right), p_{ij} = \frac{n_{ij}}{n}$$

- парные меры FN, FP, TP, TN

$$TP = |\{(\mathbf{x}_i, \mathbf{x}_j) : y_i = y_j, \hat{y}_i = \hat{y}_j\}| = \sum_{i=1}^k \sum_{j=1}^r \binom{n_{ij}}{2} = \frac{1}{2} \left(\left(\sum_{i=1}^k \sum_{j=1}^r n_{ij}^2 \right) - n \right)$$

$$FN = |\{(\mathbf{x}_i, \mathbf{x}_j) : y_i = y_j, \hat{y}_i \neq \hat{y}_j\}| = \sum_{i=1}^k \binom{n_i}{2} - TP = \frac{1}{2} \left(\sum_{i=1}^k n_i^2 - \sum_{i=1}^k \sum_{j=1}^r n_{ij}^2 \right)$$

$$FP = |\{(\mathbf{x}_i, \mathbf{x}_j) : y_i \neq y_j, \hat{y}_i = \hat{y}_j\}| = \sum_{j=1}^r \binom{m_j}{2} - TP = \frac{1}{2} \left(\sum_{j=1}^r m_j^2 - \sum_{i=1}^k \sum_{j=1}^r n_{ij}^2 \right)$$

$$TN = |\{(\mathbf{x}_i, \mathbf{x}_j) : y_i \neq y_j, \hat{y}_i \neq \hat{y}_j\}| = N - (TP + FN + FP) =$$

$$= \frac{1}{2} \left(n^2 - \sum_{i=1}^k n_i^2 - \sum_{j=1}^r m_j^2 + \sum_{i=1}^k \sum_{j=1}^r n_{ij}^2 \right)$$

- индекс Rand

$$Rand = \frac{TP + TN}{TP + TN + FP + FN}$$

- индекс Жаккара (Jaccard Index)

$$Jaccard = \frac{TP}{TP + FP + FN}$$

- индекс Фоулкса – Мэллоуса (Fowlkes-Mallows Index)

$$FM = \sqrt{\frac{TP}{TP + TN} \frac{TP}{TP + FP}}$$

Задание 5. Дана таблица сопряженности кластеризации с двумя классами и тремя кластерами:

$$\begin{pmatrix} 5 & 5 & 5 \\ 3 & 10 & K \end{pmatrix},$$

где K – параметр, принимающий целое положительное значение, причем $K > 5$.

Вычислить чистоту (purity) каждого кластера.

Решение: Если разбиение на классы задано в виде $\mathcal{T} = \{T_1, T_2, \dots, T_k\}$, а разбиение на кластеры имеет вид $\mathcal{C} = \{C_1, \dots, C_r\}$, то (расширенная) таблица (матрица) сопряженности имеет вид

$\mathcal{T} \setminus \mathcal{C}$	C_1	C_2	\dots	C_r	\sum
T_1	n_{11}	n_{12}	\dots	n_{1r}	n_1
T_2	n_{21}	n_{22}	\dots	n_{2r}	n_2
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
T_k	n_{k1}	n_{k2}	\dots	n_{kr}	n_k
\sum	m_1	m_2	\dots	m_r	n

Чистота (purity) кластеров задается формулами

$$p_j = \frac{1}{m_j} \max_{i=1, \dots, k} n_{ij} = \frac{n_{i_j j}}{m_j}, \quad i_j = \arg \max_{i=1, \dots, k} n_{ij}, \quad j = \overline{1, r}$$

Строим расширенную таблицу сопряженности:

$$\left(\begin{array}{ccc|c} 5 & 5 & 5 & 15 \\ 3 & 10 & K & K+13 \\ 8 & 15 & K+5 & K+28 \end{array} \right)$$

и определяем чистоту кластеров:

$$p_1 = \frac{5}{8} = 0.625, \quad p_2 = \frac{10}{15} = 0.667, \quad p_3 = \frac{K}{K+5}$$