

СПЕЦКУРС

Логический анализ данных в распознавании (Logical data analysis in recognition)

лектор д.ф.-м.н. Елена Всеволодовна Дюкова

Спецкурс посвящён вопросам применения аппарата дискретной математики в задачах интеллектуального анализа данных. Излагаются общие принципы, лежащие в основе логического подхода к задачам машинного обучения. Описываются методы конструирования процедур классификации по прецедентам с использованием понятий теории булевых функций и теории покрытий булевых матриц. Рассматриваются основные модели логических процедур классификации, вопросы сложности их реализации и качества решения прикладных задач.

Спецкурс для бакалавров 2-4 курсов ВМК МГУ им. М.В. Ломоносова.

По спецкурсу издано учебное пособие:

<http://www.ccas.ru/frc/papers/djukova03mp.pdf>

Лекция 3

Общая схема конструирования дискретных (логических) процедур распознавания с использованием понятия элементарного классификатора. Модели алгоритмов голосования по антипредставительным наборам и по покрытиям класса.

- В лекции 2 рассматривались модели классических логических алгоритмов распознавания. В этих алгоритмах используется процедура голосования по фрагментам описаний обучающих объектов, т.е. по таким наборам значений признаков, которые встречаются в описаниях прецедентов. В последнее время построены новые модели, в которых при вычислении оценки за класс K используется процедура голосования по наборам значений признаков, не встречающихся в описаниях прецедентов из K .
- Для того, чтобы описать общую схему конструирования логических процедур распознавания нам понадобится понятие элементарного классификатора.

- Пусть \mathbf{H} – набор из r различных признаков вида $\mathbf{H} = \{x_{j_1}, \dots, x_{j_r}\}$, $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_r)$, σ_i – допустимое значение признака x_{j_i} , $i = 1, 2, \dots, r$. Пару $(\boldsymbol{\sigma}, \mathbf{H})$ назовем элементарным классификатором (эл.кл.).
- Близость объекта $\mathbf{S} = (a_1, \dots, a_n)$ из M и эл.кл. $(\boldsymbol{\sigma}, \mathbf{H})$, $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_r)$, $\mathbf{H} = (x_{j_1}, \dots, x_{j_r})$, будем оценивать величиной $\mathbf{B}(\boldsymbol{\sigma}, \mathbf{S}, \mathbf{H})$ равной 1, если $a_{j_t} = \sigma_t$ при $t = 1, 2, \dots, r$, и равной 0 в противном случае.
- Если $\mathbf{B}(\boldsymbol{\sigma}, \mathbf{S}, \mathbf{H}) = 1$, то будем говорить, что объект \mathbf{S} содержит эл.кл. $(\boldsymbol{\sigma}, \mathbf{H})$.
- Нетрудно видеть, что фрагмент описания обучающего объекта вида $(\mathbf{S}', \mathbf{H})$, где $\mathbf{S}' = (a'_{j_1}, \dots, a'_{j_r})$, $\mathbf{H} = \{x_{j_1}, \dots, x_{j_r}\}$, порождает эл.кл. $(\boldsymbol{\sigma}, \mathbf{H})$, где $\boldsymbol{\sigma} = (a'_{j_1}, \dots, a'_{j_r})$.
- Множество всех эл.кл. обозначим через \mathbf{C} . Итак, $\mathbf{C} = \{(\boldsymbol{\sigma}, \mathbf{H})\}$, где $\mathbf{H} \subseteq (x_1, \dots, x_n)$, $\mathbf{H} = (x_{j_1}, \dots, x_{j_r})$, $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_r)$, σ_i – допустимое значение признака x_{j_i} , $i = 1, 2, \dots, r$.

- Будем считать, что каждый распознающий алгоритм A для каждого класса K , $K \in \{K_1, \dots, K_l\}$, строит некоторое подмножество $C^A(K)$ множества C . Обозначим

$$C^A = \bigcup_{j=1}^l C^A(K_j).$$

- Распознавание объекта S осуществляется на основе вычисления величины $B(\sigma, S, H)$ для каждого эл.кл. (σ, H) множества C^A , т.е. по каждому элементу множества C^A осуществляется процедура голосования. В результате для каждого класса K , $K \in \{K_1, \dots, K_l\}$, вычисляется оценка $\Gamma(S, K)$ принадлежности объекта S классу K . Таким образом, алгоритм A из рассматриваемого семейства распознающих алгоритмов определяется множеством эл.кл. C^A . Далее мы увидим, что алгоритмы будут отличаться и способом вычисления оценки $\Gamma(S, K)$, которая получается на основе голосования по эл.кл. из $C^A(K)$.

- В рассмотренных в лекции 2 моделях число голосов, поданных парами из $\mathcal{C}^A(K)$ за принадлежность объекта S к классу K , в простейших модификациях вычисляется по формуле

$$\Gamma_1(S, K) = \frac{1}{|\mathcal{C}^A(K)|} \sum_{(\sigma, H) \in \mathcal{C}^A(K)} P_{(\sigma, H)} B(\sigma, S, H),$$

где $P_{(\sigma, H)}$ – вес эл.кл. (σ, H) . В качестве $P_{(\sigma, H)}$ обычно берётся число обучающих объектов из K , содержащих (σ, H) .

- Например, если A – алгоритм вычисления оценок, то множество \mathcal{C}^A состоит из таких эл.кл., которые порождаются опорными множествами алгоритма A . Если же A – алгоритм голосования по представительным наборам, то $\mathcal{C}^A(K)$ порождается некоторым подмножеством множества представительных наборов класса K . В обоих случаях оценка $\Gamma(S, K)$ получается на основе суммирования величин $P_{(\sigma, H)} B(\sigma, S, H)$, где $(\sigma, H) \in \mathcal{C}^A(K)$. Далее объект S относится к тому классу, для которого оценка принадлежности наибольшая (если таких классов несколько, то происходит отказ от распознавания).

- Будем говорить, что элементарный классификатор (σ, H) является **корректным для класса K** , если нельзя указать пару обучающих объектов S' и S'' таких, что $S' \in K$, $S'' \notin K$ и $B(\sigma, S'', H) = B(\sigma, S', H) = 1$. Нетрудно видеть, что если A – тестовый алгоритм или алгоритм голосования по представительным наборам, то множество \mathcal{C}^A состоит только из корректных элементарных классификаторов.
- В общем случае элементарный классификатор (σ, H) , по отношению к классу K может обладать одним из следующих трех свойств:
 - 1) каждый обучающий объект S' из класса K содержит (σ, H) ;
 - 2) не все, а лишь некоторые обучающие объекты S' из класса K содержат (σ, H) ;
 - 3) ни один обучающий объект S' из класса K не содержит (σ, H) .

- Первая ситуация встречается крайне редко, поэтому работать с эл.кл., для которых выполняется свойство 1, не представляется возможным.
- Существенное различие в информативности следующих двух свойств заключается в том, что свойство 2 характеризует лишь некоторое подмножество обучающих объектов из класса K , а свойство 3 все объекты из K .
- Следовательно, в случае, когда важно рассматривать класс K изолированно от других классов, напрашивается вывод о большей информативности таких эл.кл., для которых выполнено свойство 3. В указанном случае аргументом за отнесение распознаваемого объекта S в класс K более естественно считать ситуацию, когда набор значений признаков не присутствует у всех объектов из класса K и не присутствует у объекта S .

- Классические логические процедуры классификации (тестовые алгоритмы, алгоритмы голосования по представительным наборам) основаны на построении корректных эл.кл., обладающих свойством 2. В этой лекции мы рассмотрим корректные процедуры классификации, основанные на построении корректных эл.кл., обладающих свойством 3. Этими моделями являются **модель голосования по покрытиям класса** и **модель голосования по антипредставительным наборам**. В ряде случаев указанные модели позволяют повысить качество распознавания и требуют меньших вычислительных затрат.
- Описываемые ниже модели (модель голосования по покрытиям класса и модель голосования по антипредставительным наборам класса) основаны на построении для каждого класса только таких эл.кл., которые не содержатся в описаниях ни одного объекта класса. Если распознаваемый объект S также не содержит подобный эл.кл., то считается, что этот объект близок к рассматриваемому классу. Таким образом, ищутся закономерности, присущие всем обучающим объектам рассматриваемого класса, т.е. каждый эл.кл. характеризует весь класс целиком. Использование данных моделей позволяет несколько снизить вычислительные затраты в случае большого числа классов.

- Эл.кл. (σ, H) называется *покрытием класса K* , если любой обучающий объект из K не содержит (σ, H) . Покрытие класса K называется *тупиковым*, если любое его собственное подмножество не является покрытием класса K .
- В моделях голосования по (тупиковым) покрытиям класса множество $C^A(K)$ состоит из таких эл.кл., которые являются (тупиковыми) покрытиями класса K . В отличие от классического алгоритма здесь эл.кл. из $C^A(K)$ голосует за принадлежность распознаваемого объекта классу K , если этот эл.кл. не встречается в описании рассматриваемого объекта. Принадлежность объекта S классу K (в простейшей модификации) оценивается величиной

$$\Gamma_2(S, K) = \frac{1}{|C^A(K)|} \sum_{(\sigma, H) \in C^A(K)} (1 - B(\sigma, S, H)).$$

- Переформулируем понятие представительного набора класса K , используя понятие эл.кл. Пусть $\bar{K} = \{K_1, \dots, K_l\} \setminus K$. Будем рассматривать \bar{K} как отдельный класс, т.е. будем считать, что у нас всего два класса K и \bar{K} .
- Эл.кл. (σ, H) называется *представительным набором для класса K* , если ни один обучающий объект из \bar{K} не содержит (σ, H) и хотя бы один обучающий объект из K содержит (σ, H) . Таким образом, (тупиковый) представительный набор (σ, H) является (тупиковым) покрытием для \bar{K} и не является покрытием для K .

- Рассмотрим теперь модель с антипредставительными наборами.
- Эл.кл. (σ, H) называется **(тупиковым) антипредставительным набором**, если (σ, H) является (тупиковым) покрытием класса K и хотя бы один обучающий объект из \bar{K} содержит (σ, H) , т.е. (σ, H) не является покрытием для \bar{K} . Принадлежность объекта S классу K (в простейшей модификации) оценивается величиной

$$\Gamma_3(S, K) = \frac{1}{|C^A(K)|} \sum_{(\sigma, H) \in C^A(K)} P_{(\sigma, H)} (1 - B(\sigma, S, H)),$$

где $P_{(\sigma, H)}$ – число обучающих объектов и \bar{K} , содержащих (σ, H) .

- Заметим, что представительный набор для класса K является антипредставительным для \bar{K} . Поэтому нетрудно показать, в случае двух классов при использовании обеих моделей объект S будет отнесен к одному и тому же классу.

• В самом деле, пусть A_1 – алгоритм голосования по всем представительным наборам, A_2 – алгоритм голосования по всем антипредставительным наборам. Как уже было сказано, $C^{A_1}(K_1) = C^{A_2}(K_2) = C_1$, $C^{A_1}(K_2) = C^{A_2}(K_1) = C_2$. Пусть в распознаваемом объекте S с представительными наборами класса K_1 совпадают q_1 фрагментов, а с представительными наборами класса K_2 – q_2 фрагментов. Пусть далее $\Gamma_{A_i}(S, K_j)$, $i \in \{1, 2\}$, $j \in \{1, 2\}$ – результат голосования для алгоритма A_i за класс K_j . Тогда имеем следующие оценки

$$\Gamma_{A_1}(S, K_1) = \frac{q_1}{|C^{A_1}(K_1)|} = \frac{q_1}{C_1}, \quad \Gamma_{A_1}(S, K_2) = \frac{q_2}{|C^{A_1}(K_2)|} = \frac{q_2}{C_2},$$

$$\Gamma_{A_2}(S, K_1) = \frac{|C^{A_2}(K_1)| - q_2}{C^{A_2}(K_1)} = \frac{C_2 - q_2}{C_2} = 1 - \frac{q_2}{C_2},$$

$$\Gamma_{A_2}(S, K_2) = \frac{|C^{A_2}(K_2)| - q_1}{C^{A_2}(K_2)} = 1 - \frac{q_1}{C_1}$$

Очевидно, если $\Gamma_{A_1}(S, K_1) > \Gamma_{A_1}(S, K_2)$, то и $\Gamma_{A_2}(S, K_1) > \Gamma_{A_2}(S, K_2)$.

- Рассмотрим теперь случай, когда классов больше двух. Очевидно, что представительный набор для класса K_i , $i \in \{1, 2, \dots, l\}$, является антипредставительным для каждого из остальных классов. Однако антипредставительный набор для K_i может не являться представительным набором для любого из остальных классов. В самом деле, если некий фрагмент описания встречается в описаниях объектов из классов K_{i_1} и K_{i_2} , но не встречается в K_i , тогда он является антипредставительным набором для K_i , но не является представительным ни для K_{i_1} , ни для K_{i_2} . Таким образом, если A_1 - алгоритм голосования по представительным наборам, A_2 - алгоритм голосования по антипредставительным наборам, то справедливо

$$C^{A_1} \subseteq C^{A_2}$$

- Наиболее трудоемким этапом при построении рассматриваемых алгоритмов является нахождение требуемого множества покрытий. В классических моделях покрытия строятся для \bar{K} , в новых моделях покрытия строятся для K . Таким образом, при большом числе классов новые модели требуют меньших вычислительных затрат.

УПРАЖНЕНИЯ

- 1. В обучающей выборке

$\{S_1 = (0, 1, 1), S_2 = (1, 2, 0)\}$ – класс K_1

$\{S_3 = (1, 2, 1), S_4 = (1, 1, 0)\}$ – класс K_2

найти все корректные эл.кл. для каждого из классов K_1 и K_2 , порождаемые одним признаком.

- 2. Пусть A_1 – алгоритм голосования по представительным наборам, A_3 – алгоритм голосования по покрытиям классов. Верно ли хотя бы одно из включений: 1) $C^{A_1} \subseteq C^{A_3}$ 2) $C^{A_3} \subseteq C^{A_1}$