

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 11

Методы повышения эффективности дискретных процедур распознавания

- Одним из подходов к повышению качества распознавания и снижению вычислительных затрат является проведение предварительного анализа обучающей информации. Целью такого анализа является оценка основных информативных характеристик обучающей выборки, в частности, оценка информативности признаков и выделение «шумящих» признаков, а также выделение наиболее представительных (типичных) для своего класса обучающих объектов.

1. Методы оценки информативности признаков

- Традиционно в качестве меры важности признака x_j , $j \in \{1, 2, \dots, n\}$, рассматривалась величина

$$I_j = |C_j^A|/|C^A|,$$

где C_j^A - подмножество таких элементарных классификаторов из C^A , в которых содержится признак x_j . При этом в качестве A использовался тестовый алгоритм или алгоритм голосования по представительным наборам.

- Во многих случаях такой способ оценки информативности признаков даёт не очень хорошие результаты. Например, если признак является тупиковым тестом, тогда он входит лишь в один тупиковый тест и согласно формуле для I_j имеет очень маленький вес. Однако, очевидно, что указанный признак является важным и обладает большой информативностью.

- С другой стороны, многозначные признаки, являющиеся малоинформативными «шумящими» признаками порождают большое число тестов (представительных наборов) и, следовательно, имеют большой вес I_j .
- Пусть A - алгоритм голосования по представительным наборам. Частично проблему шумящих признаков можно решить, введя порог p_{min} минимальной встречаемости представительного набора в обучающей выборке. То есть рассматривать только те представительные наборы, которые в обучающей выборке встречаются не менее p_{min} раз. Используя это дополнительное условие, шумящие многозначные признаки не получают самой большой оценки по информативности. Однако не всегда можно взять достаточно большой порог p_{min} из-за того, что структура класса такова, что каждый представительный набор встречается не очень большое число раз. В данном случае предлагается следующее.

- Положим $\bar{K} = \{K_1, \dots, K_l\} \setminus K$. Зададим целые числа p и q такие, что
- $1 \leq p \leq \min_{1 \leq i \leq l} |K_i|$ $0 \leq q \leq \min_{1 \leq i \leq l} |\bar{K}_i|$, $p > q$.
- Элементарный классификатор (σ, H) назовём (p, q) - представительным набором для класса K_i , если не менее чем для p объектов S' из класса $|K_i|$ справедливо $B(\sigma, S', H) = 1$ и не более чем для q объектов S'' из \bar{K}_i справедливо $B(\sigma, S'', H) = 1$. В частности, представительный набор является $(1, 0)$ -представительным набором.
- Пусть (σ, H) является (p, q) -представительным набором класса K , p_σ - его встречаемость в обучающей выборке в классе K , q_σ - встречаемость в остальных классах ($p_\sigma \geq p, q_\sigma \leq q$). Тогда информативный вес представительного набора (σ, H) будем вычислять по следующей формуле

$$v(\sigma, H) = (p_\sigma - p) + (q - q_\sigma).$$

- За счет варьирования p и q можно существенно снизить влияние шумящих признаков. Дело в том, что разрешив представительным наборам встречаться в других классах (увеличив q), можно увеличить минимальную встречаемость в своём классе (увеличить p). В практической реализации для снижения вычислительных затрат используется следующий приём. Множество (p, q) -представительных наборов строится по случайной подвыборке, которая составляет **60 - 70%** обучающей выборки. После этого осуществляется проверка, являются ли построенные представительные наборы (p, q) -представительными наборами для всей обучающей выборки, и в случае положительного результата вычисляются их информативные веса. Далее вычисляется информативный вес каждого признака как отношение суммы весов (p, q) -представительных наборов, которые он порождает, к сумме весов всех (p, q) -представительных наборов для данного класса.

- Для решения проблемы шумящих признаков не обязательно добиваться того, чтобы оценка их информативности была низкой на любой подвыборке. Достаточно, чтобы оценка информативности шумящих признаков при различных разбиениях была неустойчива по разбиению информации на две подвыборки, при том, что по остальным признакам наблюдается некоторая устойчивость.

- Достаточно эффективным и не требующим больших вычислительных затрат является предлагаемый ниже метод оценки информативности признаков и отдельных значений признаков, основанный на вычислении близости между парами объектов по отдельным признакам.

- Пусть $S' \in K_i$, $i \in \{1, 2, \dots, l\}$, $j \in \{1, 2, \dots, n\}$. Положим

$$\mu_{ij}^{(1)}(S') = \frac{1}{|K_i|} \sum_{S'' \in K_i} B(S', S'', \{x_j\}),$$

$$\mu_{ij}^{(2)}(S') = \frac{1}{|\bar{K}_i|} \sum_{S'' \in \bar{K}_i} B(S', S'', \{x_j\}),$$

- Величины $\mu_{ij}^{(1)}(S')$ и $\mu_{ij}^{(2)}(S')$ характеризуют близость объекта S' соответственно к своему классу и к другим классам по признаку x_j . Величину

$$\mu_{ij}(S') = \mu_{ij}^{(1)}(S') - \mu_{ij}^{(2)}(S')$$

назовём весом значения признака x_j , для объекта S' . Будем говорить, что значение признака x_j является типичным для объекта S' , если $\mu_{ij}(S') > \mu$, где μ - порог информативности значения признака, $-1 < \mu < 1$.

- Например, можно положить $\mu = 0$. Тогда значение признака будет являться типичным для класса, если в этом классе оно встречается чаще, чем в остальных.

- Множество типичных значений признаков в таблице обучения образует так называемую информативную зону. Далее при построении множества элементарных классификаторов имеет смысл анализировать только те значения признаков, которые попадают в информативную зону, тем самым уменьшается перебор при построении распознающего алгоритма. Кроме того, в информативную зону не попадают значения признаков, которые очень редко встречаются во всех классах (шумящие). Поэтому использование информативной зоны позволяет также снизить влияние шумящих признаков.
- Исследование типичности значений признаков позволяет также оценить сложность (в смысле возможности построения качественного алгоритма распознавания) решаемой задачи.

2. Выделение типичных объектов в классе

- При решении прикладной задачи распознавания интересно попытаться оценить эффективность построенного алгоритма при распознавании объектов, не входящих в обучающую выборку. Например, воспользоваться хорошо известным методом скользящего контроля. К сожалению, в ряде прикладных задач алгоритмы не всегда показывают достаточно высокую эффективность. Такая ситуация возникает, когда классы плохо отделяются друг от друга (в каждом классе есть много объектов, описания которых похожи на объекты, не принадлежащие данному классу). В этом случае построенный алгоритм зачастую, хотя и хорошо распознает "известные" ему объекты (объекты, которые участвовали в построении алгоритма), но плохо распознает "новые" объекты. В данном разделе предлагается подход, позволяющий повысить качество распознающих алгоритмов за счёт выделения "типичных" для своих классов объектов. Этот подход продемонстрирован ниже на примере модели голосования по представительным наборам.

- Объекты, лежащие на границе между классами, плохо распознаются и, по-видимому, они не позволяют строить короткие представительные наборы. Пусть описание обучающего объекта, не принадлежащего классу **K**, похоже на описания некоторых объектов из **K**. Тогда данный объект "лишает" класс **K** некоторого множества коротких представительных наборов, и это существенно снижает эффективность алгоритма. Для решения указанной проблемы предлагается разбить обучающую выборку на две подвыборки, по первой (базовой) построить множество представительных наборов, по второй (контрольной) вычислить их веса. Причем разбить нужно таким образом, чтобы объекты, находящиеся на границе между классами, попали в контрольную подвыборку, а все остальные (типичные) объекты - в базовую подвыборку. Практические эксперименты на прикладных задачах подтверждают гипотезу о том, что такое разбиение увеличивает число коротких представительных наборов и тем самым позволяет повысить качество алгоритма распознавания.

- Описанный метод позволяет очень быстро оценить типичность обучающих объектов и отдельных их фрагментов по отношению к своим классам. Его недостатком является то, что информативность (типичность) значения признака вычисляется независимо от других признаков. Не учитывается тот факт, что, вообще говоря, фрагмент описания некоторого объекта может быть типичен для одного из классов (в этом классе данный фрагмент встречается значительно чаще, чем в остальных классах), но при этом значения признаков, которые составляют этот фрагмент, встречаются в разных классах примерно одинаковое число раз и не являются типичными ни для одного из классов. Кроме того, необходимо настраивать параметры μ и p , что не очень удобно.
- Ниже предлагается метод выделения типичных объектов на основе проведения процедуры скользящего контроля, который заключается в следующем.

- Из обучающей выборки исключается один из объектов, например объект S_i , $i \in \{1, 2, \dots, m\}$. По подвыборке $\{S_1, \dots, S_m\} \setminus S_i$ строится распознающий алгоритм (например, используется модель алгоритма голосования по представительным наборам в классическом варианте). Далее этот алгоритм применяется для распознавания объекта S_i . Объект S_i считается типичным для своего класса, если построенный алгоритм распознал его правильно, и нетипичным для своего класса, если алгоритм отнёс его к другому классу или отказался от распознавания. Описанная процедура повторяется для всех объектов обучающей выборки.
- Очевидно, что для больших задач процедура скользящего контроля в алгоритме голосования по представительным наборам требует существенных вычислительных затрат. Существует быстрый способ вычисления требуемых оценок, позволяющий сократить время счета примерно в m раз (см. Дюкова Е.В., Песков Н.В. Построение распознающих процедур на базе элементарных классификаторов // Математические вопросы кибернетики. 2005. № 14. С. 57-92, <http://library.keldysh.ru/mvk.asp?id=2005-57>).

3. Вычисление информативности эл.кл.

- Способ вычисления информативности эл.кл. проиллюстрируем на примере представительного набора. Разделим обучающую выборку на две подвыборки: базовую и контрольную. По базовой подвыборке построим множество представительных наборов. Сопоставим каждому построенному представительному набору некий вес, который вычислим по контрольной подвыборке.
- Пусть (σ, H) - представительный набор класса K , $K \in \{K_1, \dots, K_l\}$, порождаемый объектом из базовой выборки, и пусть $N(K, \sigma, H)$ - число объектов в контрольной выборке, за которых данный представительный набор голосует "правильно", $\bar{N}(K, \sigma, H)$ - число объектов в контрольной выборке, за которых он голосует "неправильно". Тогда в качестве веса представительного набора (σ, H) можно рассматривать, например, следующие величины:

$$v_1(\sigma, H) = N(K, \sigma, H) ;$$

$$v_2(\delta, H) = \frac{1 + N(K, \sigma, H)}{1 + \bar{N}(K, \sigma, H)} .$$

- Принадлежность объекта S классу K будем оценивать величиной

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

- В качестве информативного веса признака x_j будем рассматривать величину

$$I_j = \frac{\sum_{(\sigma, H) \in C^A(K), x_j \in H} v(\sigma, H)}{\sum_{(\sigma, H) \in C^A(K)} v(\sigma, H)}$$

- В общем случае функция, по которой вычисляется вес представительного набора, должна обладать следующим свойством. Эта функция должна монотонно возрастать по числу объектов из контрольной выборки, за которые представительный набор голосует правильно, и монотонно убывать по числу объектов, за которые он голосует не правильно.

- УПРАЖНЕНИЯ

- 1. Пусть один класс один представлен прецедентами $\{(1, 0, 0, 1), (1, 0, 1, 1), (1, 1, 0, 1)\}$, а другой $\{(1, 1, 0, 0), (1, 0, 0, 0), (1, 0, 1, 0)\}$. Найдите все $(2, 1)$ представительные наборы.
- 2. Будет ли для (p, q) элементарного классификатора применима для оценки информативности $\nu_1(\delta, H) ? \nu_2(\delta, H) ?$