

Экспериментальное исследование свойств специализированных методов индуктивного формирования знаний для онтологий медицинской диагностики

Клещев А.С., Смагин С.В.

*Институт автоматики и процессов управления ДВО РАН,
ул. Радио, д.5, г. Владивосток, 690041, Россия.*

kleshev@iacp.dvo.ru, smagin@iacp.dvo.ru

Аннотация. *Индуктивное формирование знаний (ИФЗ) лежит в основе многих направлений исследований, таких как машинное обучение, распознавание образов и т.д., каждое из которых характеризуется собственным подходом к проблеме ИФЗ, собственными постановками задач и методами их решения. В рамках данного исследования разработан подход к организации компьютерных экспериментов для изучения свойств методов ИФЗ на выборках модельных данных разного объема. Предложенный подход реализован на примере экспериментального исследования свойств специализированного метода ИФЗ (метода случайной расстановки границ периодов динамики) для онтологии медицинской диагностики. Полученные результаты показывают, что если рассматривать метод не с точки зрения практического применения, а с точки зрения выяснения того, как на самом деле устроена природа, то этот метод решает задачу плохо. В общем случае проведенное экспериментальное исследование показывает, что хорошие внешние оценки метода ИФЗ могут сочетаться с его плохими внутренними оценками.*

Ключевые слова: экспериментальное исследование свойств, специализированный метод, индуктивное формирование знаний, модельные данные, онтология медицинской диагностики, внешняя оценка, внутренняя оценка.

1. Введение

В работе [1] представлен обзор существующих экспериментальных исследований в области изучения свойств методов индуктивного формирования знаний (ИФЗ), а также приведена постановка задачи таких исследований. В работах [1,2] предложен общий подход к проведению компьютерных экспериментов по изучению свойств методов ИФЗ, основанный на использовании модельных данных. В [1,3] в качестве одного из примеров применения общего подхода рассмотрена предметная область (ПО) медицинской диагностики. В [3] приведена ее непримитивная упрощенная онтология, определяющие соотношения, а также алгоритмы генерации случайной базы знаний (СБЗ) и случайной выборки данных на основе СБЗ. Настоящая работа посвящена реализации предложенного общего подхода на примере экспериментального исследования свойств специализированного метода ИФЗ (метода случайной расстановки границ периодов динамики) для упрощенной онтологии медицинской диагностики.

2. Постановка задачи индуктивного формирования знаний для упрощенной онтологии медицинской диагностики и метод ее решения

В работе [4] приводятся постановки задач ИФЗ в терминах непримитивных онтологий ПО, а также предлагается специализированный метод их решения для непримитивной упрощенной онтологии медицинской диагностики. В этой онтологии рассматривается один вид причинно-

следственных отношений – клинические проявления заболеваний, учитываются многократные наблюдения пациента, результаты которых зависят от времени наблюдения.

Упрощенная онтология медицинской диагностики и ее определяющие соотношения представлены в [3,4]. Базой знаний (БЗ) здесь является набор значений терминов для описания знаний – неинтересных параметров (*признаки, заболевания, возможные значения и клиническая картина*) и интересных параметров (*нормальные значения, число периодов динамики, значения для периода, верхняя граница и нижняя граница*). Выборкой данных здесь является совокупность историй болезней (ИБ), диагнозами которых являются заболевания из БЗ. Каждая ИБ представляет собой информацию, описывающую течение одного заболевания у пациента – набор значений терминов для описания действительности – наблюдаемых неизвестных (*диагноз, признаки, моменты наблюдения, значения признаков в моменты наблюдения*). Значения ненаблюдаемых неизвестных (*число периодов динамики признака, границы периодов динамики признака*) в выборке отсутствуют.

Задача ИФЗ для упрощенной онтологии медицинской диагностики заключается в поиске индуктивно формируемой базы знаний (ИФБЗ) по обучающей выборке, состоящей из выборки примеров и выборки контрпримеров. В работе [4] показано, что если все неизвестные онтологии являются наблюдаемыми, то для такого класса онтологий можно построить определяющие соотношения, по которым ИФБЗ может быть сформирована всего за один проход по обучающей выборке примеров, т.е. без перебора. Такая задача названа задачей ИФЗ с полной информацией. Задача, в которой часть неизвестных является ненаблюдаемыми, названа задачей ИФЗ с неполной информацией. В этом случае суть метода ИФЗ заключается в определении значений ненаблюдаемых неизвестных по обучающей выборке примеров и контрпримеров, что сводит задачу с неполной информацией к задаче с полной информацией. Ее точным решением является полный перебор наборов всех возможных значений ненаблюдаемых неизвестных. Однако, такой способ решения связан с существенной вычислительной работой, время которой очевидно превысит любые допустимые пределы.

В [5] предложен метод случайной расстановки границ периодов динамики для упрощенной онтологии медицинской диагностики, который является одним из вариантов методов Монте-Карло [6]. Идея метода заключается в многократной случайной генерации (в соответствии с рядом ограничений) значений ненаблюдаемых неизвестных онтологии, используя обучающую выборку примеров (см. Рис.1).

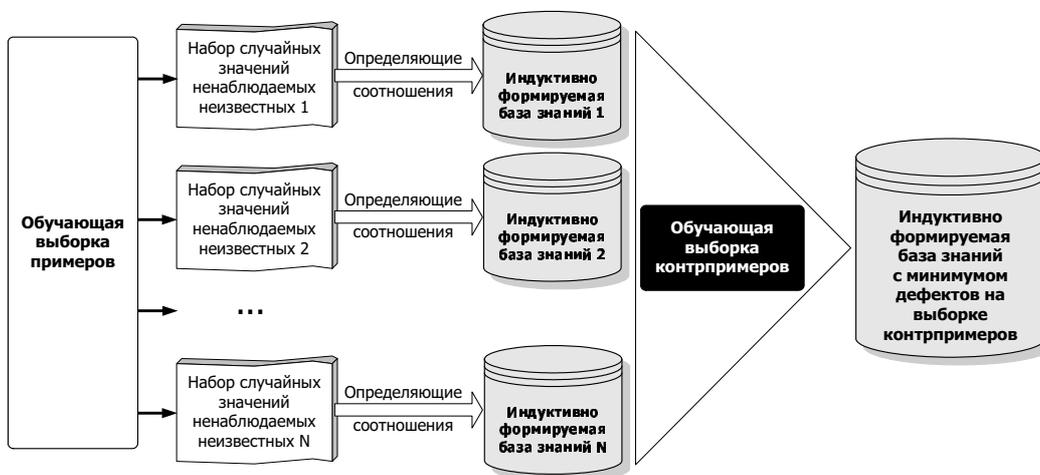


Рис.1. Метод (случайной расстановки границ периодов динамики) решения задачи индуктивного формирования знаний с неполной информацией

Количество случайных генераций задается в качестве параметра метода. Для каждого сгенерированного варианта этих значений по определяющим соотношениям находится ИФБЗ. Далее среди найденных ИФБЗ по выборке контрпримеров осуществляется выбор такой из них, которая имеет минимум дефектов (см. Глоссарий в [1,3]) на обучающей выборке контрпримеров.

Применительно к упрощенной онтологии медицинской диагностики решение задачи ИФЗ осуществляется отдельно для каждого заболевания онтологии. Т.о. все ИБ обучающей выборки, диагнозом которых является заболевание, для которого решается задача ИФЗ, считаются его примерами, а все остальные ИБ – контрпримерами. Каждое заболевание обладает клинической картиной (КК) – набором таких признаков, значения которых зависят от заболевания и проявляются в его клинических проявлениях.

Значения ненаблюдаемых неизвестных онтологии случайным образом генерируются для каждого признака, входящего в клиническую картину заболевания, для которого решается задача ИФЗ, на основе его примеров. Вначале определяется *число периодов динамики признака*, после чего генерируются *границы периодов динамики признака* на основе указанных ограничений. Далее для сгенерированного варианта значений ненаблюдаемых неизвестных для признака по определяющим соотношениям вычисляются значения интересных параметров *значения для периода, верхняя граница и нижняя граница*. Т.о. формируется описание клинического проявления (ОКП) заболевания, для которого решается задача ИФЗ, по признаку, входящему в его клиническую картину. Для признаков, не входящих в клинические картины заболеваний, формируется множество значений интересного параметра *нормальные значения признака*, как результат решения задачи с полной информацией. Совокупность ОКП каждого заболевания по каждому признаку образует ИФЗ.

Т.к. для одного признака предполагается генерация множества вариантов значений ненаблюдаемых неизвестных, то в результате для него будет получено множество ОКП заболевания. Выбор наилучших среди них осуществляется на основе значений признака в ИБ контрпримеров заболевания, для которого решается задача ИФЗ. Наилучшим среди ОКП заболевания по признаку считается то, которое отвергает максимум ИБ контрпримеров этого заболевания. Если таких ОКП получается несколько, они становятся конкурирующими. После получения одного или нескольких конкурирующих ОКП заболевания по каждому признаку, осуществляется комбинирование ОКП заболевания по разным признакам. Целью этого является нахождения наилучшей комбинации ОКП по разным признакам, отвергающей максимальное число ИБ контрпримеров заболевания, для которого решается задача ИФЗ.

Алгоритм, реализующий метод случайной расстановки границ периодов динамики, представлен в работе [7], а схема его распараллеливания – в работе [5].

3. Экспериментальное исследование свойств метода случайной расстановки границ периодов динамики

В работе применен общий подход к экспериментальному исследованию свойств методов индуктивного формирования знаний [1-3], основанный на использовании модельных данных. Идея этого подхода состоит в том, чтобы для получения модельных данных, необходимых для экспериментальных исследований свойств метода ИФЗ, явным образом представить описание класса баз знаний (которое лежит в основе любого метода ИФЗ). В нашем случае таким описанием является онтология медицинской диагностики. На основе этой онтологии разработаны алгоритмы случайной генерации баз знаний из данного класса, а также алгоритмы генерации случайных выборок модельных данных (обучающих и контрольных) различного объема на основе этих баз знаний. Эти алгоритмы представлены в работе [3]. Применяя исследуемый метод ИФЗ к обучающим выборкам модельных данных разного объема, для каждой из них находилась ИФЗ. После этого, используя внешние оценки СБЗ и ИФЗ на контрольных выборках, а также сравнивая ИФЗ с СБЗ, были получены оценки свойств метода ИФЗ: времени его работы, внешнего и внутреннего качества результатов, а также оценки устойчивости этих свойств.

Экспериментальное исследование метода ИФЗ на основе модельных данных организовано в виде серий. Отдельная серия включает в себя ряд экспериментов, каждый из которых проводился на основе одной обучающей выборки заданного объема, в соответствии с параметрами, едиными для всей серии. После завершения всех экспериментов серии, их результаты оценивались и визуализировались совместно. Для серии экспериментов было сгенерировано 10 различных СБЗ на основе следующих значений параметров:

Таблица 1. Параметры генерации случайной базы знаний и их значения в серии экспериментов

| Параметры генерации случайной базы знаний | Значения |
|---|----------|
| Количество признаков | 100 |
| Количество заболеваний | 2 |
| Размер области пересечения признаков в клинических картинах разных заболеваний | 80% |
| Размер области возможного пересечения признаков в клинических картинах разных заболеваний | 20% |
| Количество возможных значений признака | 10 |
| Ограничение на число периодов динамики | 5 |
| Ограничение на верхнюю границу | 24 часа |

На основе каждой СБЗ и приведенных ниже параметров было сгенерировано 8 различных объемов обучающих выборок, включающих от 10 до 1280 (8 степеней двойки, умноженные на 10) ИБ каждого заболевания, а также максимальный объем контрольной выборки. Для исследования устойчивости свойств метода каждая обучающая выборка случайным образом была сгенерирована в 10 различных экземплярах и разделена на выборки примеров и контрпримеров для каждого заболевания своей СБЗ. Т.о. в серии было проведено 800 (10 СБЗ * 8 объемов * 10 экземпляров) экспериментов и получено 800 ИФБЗ.

Таблица 2. Параметры генерации случайных выборок и их значения в серии экспериментов

| Параметры генерации случайных выборок | Значения |
|--|----------|
| Количество экземпляров каждой случайной выборки одного и того же объема в наборе для СБЗ | 10 |
| Количество примеров (историй болезни) каждого заболевания в случайной выборке | 10..1280 |
| Ограничение на количество моментов наблюдения признака, входящего в клиническую картину заболевания | 10 |
| Ограничение на количество моментов наблюдения признака, не входящего в клиническую картину заболевания | 5 |
| Ограничение на длительность наблюдения признака, не входящего в клиническую картину заболевания | 30 часов |

Серия экспериментов проводилась на основе значений параметров:

Таблица 3. Параметры проведения эксперимента и их ограничения

| Параметры проведения эксперимента | Ограничения |
|---|----------------|
| Ограничение на длительность проведения серии экспериментов | 900 минут |
| Максимальное число периодов наблюдения признака в эксперименте | 5 |
| Количество попыток расстановки границ для числа периодов динамики от 2 до 5 | 10, 20, 40, 80 |
| Количество стартов эксперимента для каждой случайной выборки набора | 1 |

Согласно общему подходу, проведение серии экспериментов было разбито на последовательно выполняемые этапы, описание которых представлено ниже.

3.1. Подготовка модельных данных

На данном этапе была разработана подсистема генерации модельных данных DataGen, реализующая соответствующие алгоритмы из работы [3]. С ее помощью была осуществлена генерация СБЗ в количестве, определяемом описанием серии экспериментов, а также генерация наборов обучающих и контрольных выборок различных объемов для каждой СБЗ (количество выборок в наборе и диапазон их объемов также определяются описанием серии).

3.2. Применение исследуемого метода ИФЗ

На данном этапе была разработана подсистема индуктивного формирования знаний MethMonte, реализующая алгоритм из работы [7]. В этой подсистеме для набора обучающих выборок каждой СБЗ проводятся эксперименты. Эксперимент – это решение задачи ИФЗ, входными данными которой является обучающая выборка конкретного объема, сгенерированная на основе одной СБЗ, а результатом – ИФБЗ, состоящая из ОКП всех заболеваний по каждому признаку. При проведении эксперимента возможно получение конкурирующих ОКП некоторого заболевания по одному и тому же признаку (или их комбинации). По завершении эксперимента из конкурирующих ОКП заболевания выбирается одно, а остальные удаляются.

3.3. Оценка и визуализация результатов

На данном этапе была разработана подсистема оценки и визуализации ShowRes результатов экспериментального исследования свойств метода случайной расстановки границ периодов динамики. С ее помощью была проведена внешняя и внутренняя оценка свойств исследуемого

метода, а затем визуализация их значений. Была оценена зависимость времени решения задачи ИФЗ от объема обучающих выборок. На контрольных выборках максимального объема было оценено качество сгенерированных СБЗ и ИФБЗ (была проведена диагностика по всем признакам). Также была установлена степень сходства СБЗ и соответствующих им ИФБЗ, полученных на основе обучающих выборок разного объема. Для каждого свойства была исследована его устойчивость относительно разных обучающих выборок.

При диагностике по всем признакам на основе контрольных выборок ошибка 1-ого и 2-ого рода (см. Глоссарий в [1,3]) вычисляются для ситуаций, когда в базе знаний находятся всего два взаимоисключающих заболевания, а результатом решения задачи диагностики для конкретной ИБ становится только одно из них. В случаях, когда последнее условие не выполняется, возможен один из следующих результатов диагностики:

1. точно чужое (результатом становится неправильное заболевание),
2. отказ от диагностики (ни одно из заболеваний не становится результатом),
3. оба заболевания (результатом становятся оба заболевания),
4. точно свое (результатом становится правильное заболевание).

В случаях, когда в базе знаний находится больше двух заболеваний, результат 3 может распасться на ряд вариантов, смысл которых состоит в установлении того, насколько полученный результат неточен (сколько в него войдет неправильных заболеваний).

4. Результаты экспериментального исследования метода случайной расстановки границ периодов динамики

В экспериментальном исследовании были получены внешние (1-2) и внутренние (3-6) оценки свойств метода случайной расстановки границ периодов динамики. Для каждой из них указаны ее название, способ вычисления, общая интерпретация, график зависимости (полученный с помощью подсистемы ShowRes), а также его описание.

4.1. Зависимость времени решения задачи ИФЗ от объема обучающих выборок

Способ вычисления: среди длительностей решения задачи ИФЗ (длительностей экспериментов) на каждом из 8 объемов обучающих выборок выбираются соответствующие минимумы и максимумы. На Рис.2 приведен график временной зависимости в логарифмическом масштабе.

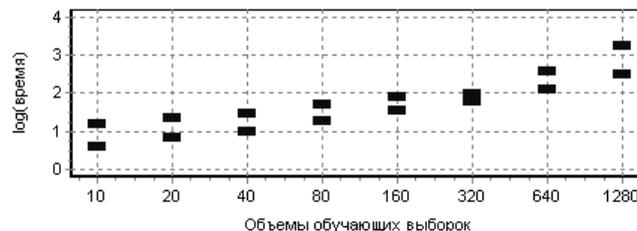


Рис.2. Вид временной зависимости при различных объемах обучающих выборок

Описание: полученный вид временной зависимости показывает, что на малых и средних объемах обучающих выборок время ИФЗ растет почти линейно. Это связано с тем, что на таких объемах основную часть времени занимает процесс формирования ОКП заболевания по признаку, зависящий только от объема выборки примеров заболевания. Выбор наилучших ОКП зависит от объема выборки контрпримеров этого заболевания. На больших же объемах время растет более чем линейно за счет того, что увеличивается доля работы по комбинированию ОКП заболевания по разным признакам.

4.2. Зависимость результатов диагностики для ИФБЗ от объема обучающих выборок

Способ вычисления:

- для каждой ИФБЗ на основе соответствующей ей контрольной выборки максимального объема решается задача диагностики по всем признакам и вычисляются проценты каждого из четырех возможных результатов диагностики, описанных выше (см. 3.3),

- для каждой из 100 ИФБЗ, сформированной на основе обучающей выборки одного из 8 объемов, среди оценок каждого результата диагностики выбираются их минимумы и максимумы, которые отображаются в виде столбцов на графиках, что демонстрирует разброс качества ИФБЗ для данного объема обучающей выборки,
- также для каждой СБЗ на основе соответствующей ей контрольной выборки максимального объема, аналогично предыдущему пункту, решается задача диагностики по всем признакам, и для каждого результата вычисляются минимумы и максимумы процентов случаев, которые графически образуют коридоры, демонстрирующий разброс качества СБЗ.

Описание: при исследовании разброса качества 10 СБЗ, сгенерированных для серии экспериментов, возможным было получение только двух результатов диагностики: 3 (оба заболевания), либо 4 (точно свое), причем сумма процентов случаев их возникновения должна равняться 100%. В исследовании выяснилось, что результат 4 (точно свое) возникает в 100% случаев. Это связано с тем, что для эксперимента было сгенерировано достаточно признаков, чтобы среди них естественным образом нашлись такие, которые сделали эти заболевания разделимыми, т.е. непохожими друг на друга.

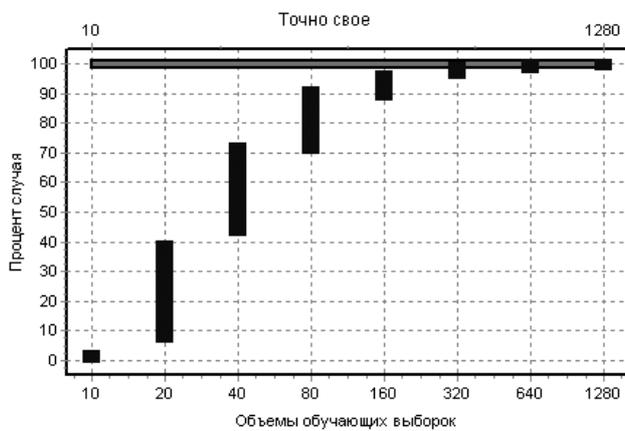


Рис.3. Результат 4 (точно свое) для ИФБЗ при различных объемах обучающих выборок

При исследовании разброса качества 800 ИФБЗ, было возможно получение любого из четырех результатов диагностики. В серии экспериментов получилось, что результаты 1 (точно чужое) и 3 (оба заболевания) не возникают, что связано с хорошей делимостью заболеваний СБЗ, по которым генерировались обучающие и контрольные выборки. Распределение результата 4 (точно свое) представлено на Рис.3. Очевидно, что с ростом обучающих выборок качество ИФБЗ растет: проценты случая 4 – увеличиваются (а проценты случая 2 в свою очередь уменьшаются). В то же время с ростом обучающих выборок разброс процентов обоих результатов уменьшается, что свидетельствует об устойчивости метода, а проценты случая 4 стремятся к своему максимуму – коридору СБЗ, что свидетельствует о том, что на больших объемах обучающих выборок качество ИФБЗ становится сопоставимым с качеством соответствующих СБЗ.

4.3. Зависимость отношений между множествами значений признаков, не входящих в клинические картины заболеваний в ИФБЗ и СБЗ, от объема обучающих выборок

Способ вычисления:

цикл по всем ИФБЗ:

цикл по всем заболеваниям:

- для ОКП заболевания по признакам, не входящим в клиническую картину этого заболевания, проводится сравнение значений признаков в ИФБЗ со значениями этих же признаков в соответствующей СБЗ – возможны 2 отношения между множествами их значений:
 1. {значения признака в ИФБЗ} \equiv {значения признака в СБЗ},
 2. {значения признака в ИФБЗ} \subset {значения признака в СБЗ},

- если x – число признаков в отношении 1-2, а y – число признаков, не входящих в клинические картины заболеваний, то процент каждого отношения между множествами значений таких признаков в ИФБЗ и СБЗ вычисляется по формуле $(x/y)*100$,
- для каждой из 100 ИФБЗ, сформированной на основе обучающей выборки одного из 8 объемов, среди процентов каждого отношения выбираются их минимумы и максимумы, которые отображаются в виде столбцов на графиках, что демонстрирует разброс качества ИФБЗ для данного объема обучающей выборки.

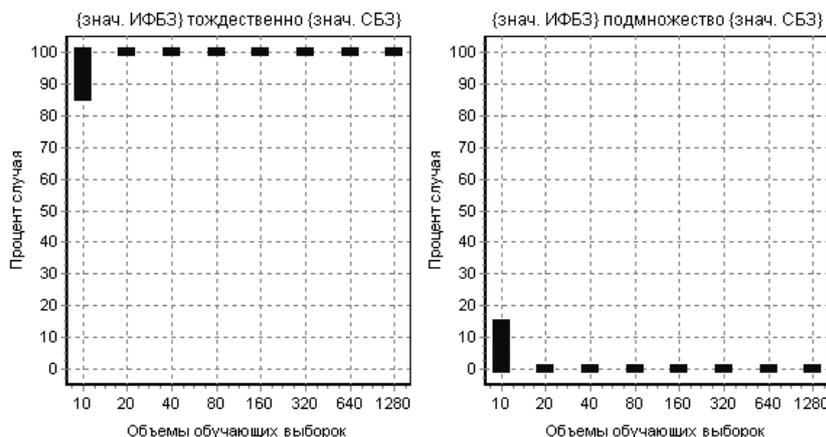


Рис.4. Отношения между множествами значений признаков, не входящих в клинические картины заболеваний в ИФБЗ и СБЗ, при различных объемах обучающих выборок

Описание: результаты сравнения, представленные на Рис.4, свидетельствуют о полном совпадении множеств нормальных значений (т.к. признаки не входят в клинические картины заболеваний) признаков в ИФБЗ и СБЗ при объемах обучающих выборок больше 10 ИБ каждого заболевания. Это объясняется тем, что при формировании ОКП заболеваний по признакам, не входящим в клинические картины этих заболеваний, решается задача ИФЗ с полной информацией, в которой нет ненаблюдаемых неизвестных, а вычисления ведутся по определяющим соотношениям.

4.4. Зависимость процента совпадения ЧПД одних и тех же признаков, входящих в клинические картины заболеваний в ИФБЗ и СБЗ, от объема обучающих выборок

Способ вычисления:

цикл по всем ИФБЗ:

цикл по всем заболеваниям:

- для ОКП заболевания по каждому признаку, входящему в клиническую картину этого заболевания, проводится сравнение числа периодов динамики (ЧПД) признака с ЧПД этого же признака в соответствующей СБЗ – между ними возможны 3 отношения:
 1. ЧПД признака из ИФБЗ < ЧПД признака из СБЗ,
 2. ЧПД признака из ИФБЗ = ЧПД признака из СБЗ,
 3. ЧПД признака из ИФБЗ > ЧПД признака из СБЗ,
- если x – число признаков, у которых ЧПД находятся в отношении 1-3, а y – число признаков, входящих в клинические картины заболеваний, то процент каждого отношения вычисляется по формуле $(x/y)*100$,
- для каждой из 100 ИФБЗ, сформированной на основе обучающей выборки одного из 8 объемов, среди процентов каждого отношения выбираются их минимумы и максимумы, которые отображаются в виде столбцов на графиках, что демонстрирует разброс качества ИФБЗ для данного объема обучающей выборки.

Описание: результаты сравнения, представленные на Рис.5, свидетельствуют о том, что с ростом обучающей выборки процент совпадения ЧПД не растет, а исследуемый метод стремится увеличивать ЧПД – так ему проще найти лучшую расстановку границ периодов динамики.

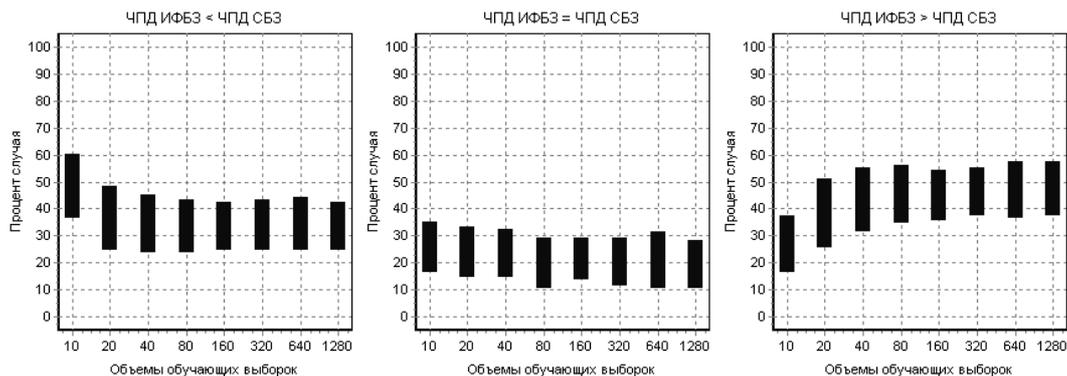


Рис.5. Отношения между ЧПД признаков, входящих в клинические картины заболеваний в ИФБЗ и СБЗ, при различных объемах обучающих выборок

4.5. Зависимость отношений между множествами значений признаков, входящих в клинические картины заболеваний, ЧПД которых в ИФБЗ и СБЗ совпали, от объема обучающих выборок

Способ вычисления:

цикл по всем ИФБЗ:

цикл по всем заболеваниям:

- для ОКП заболевания по каждому признаку, входящего в клиническую картину этого заболевания, ЧПД которого в ИФБЗ и СБЗ совпали, проводится сравнение значений признака со значениями этого же признака в соответствующей СБЗ – возможны 3 отношения между множествами их значений:
 1. {значения признака в ИФБЗ} \equiv {значения признака в СБЗ},
 2. {значения признака в ИФБЗ} \subset {значения признака в СБЗ},
 3. остальные,
- если x – число признаков в отношении 1-3, а y – число признаков, ЧПД которых в ИФБЗ и СБЗ совпали, то процент каждого отношения между множествами значений таких признаков в ИФБЗ и СБЗ вычисляется по формуле $(x/y)*100$,
- для каждой из 100 ИФБЗ, сформированной на основе обучающей выборки одного из 8 объемов, среди процентов каждого отношения выбираются их минимумы и максимумы, которые отображаются в виде столбцов на графиках, что демонстрирует разброс качества ИФБЗ для данного объема обучающей выборки.

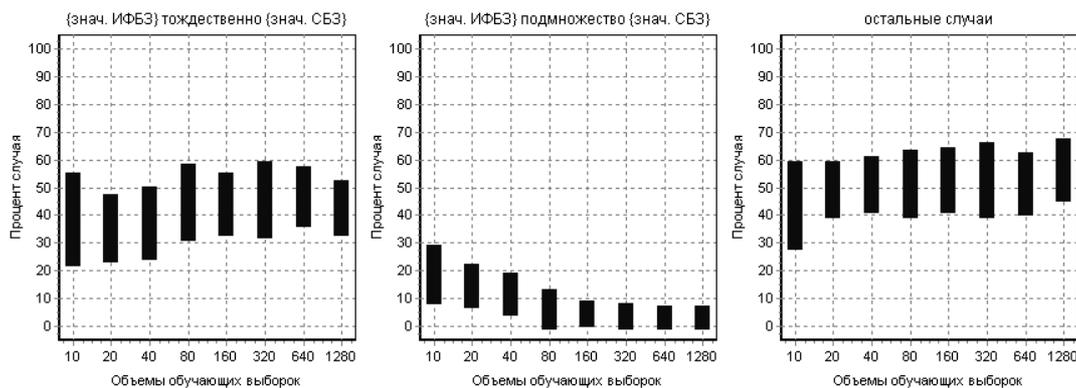


Рис.6. Отношения между множествами значений признаков, входящих в клинические картины заболеваний в ИФБЗ и СБЗ, с совпавшими ЧПД при различных объемах обучающих выборок

Описание: результаты сравнения, представленные на Рис.6, свидетельствуют о том, что с ростом обучающей выборки процент совпадения областей значений не улучшается.

4.6. Зависимость отношений между границами периодов динамики признаков, входящих в клинические картины заболеваний, ЧПД которых в ИФБЗ и СБЗ совпали, от объема обучающих выборок

Способ вычисления:

цикл по всем ИФБЗ:

цикл по всем заболеваниям:

- для ОКП заболевания по признаку, входящего в клиническую картину этого заболевания, ЧПД которых в ИФБЗ и СБЗ совпали, вычисляются разности верхних (соответственно нижних) границ периодов динамики признака в СБЗ и верхних (нижних) границ периодов динамики этого признака в соответствующей ИФБЗ,
- для каждого значения разности, если x – количество периодов динамики с этим значением разности, y – суммарное количество периодов динамики, то процент каждого случая разности вычисляется по формуле $(x/y)*100$,
- для каждой из 100 ИФБЗ, сформированной на основе обучающей выборки одного из 8 объемов, среди процентов каждого отношения выбираются их минимумы и максимумы, которые отображаются в виде столбцов на графиках, что демонстрирует разброс качества ИФБЗ для данного объема обучающей выборки.

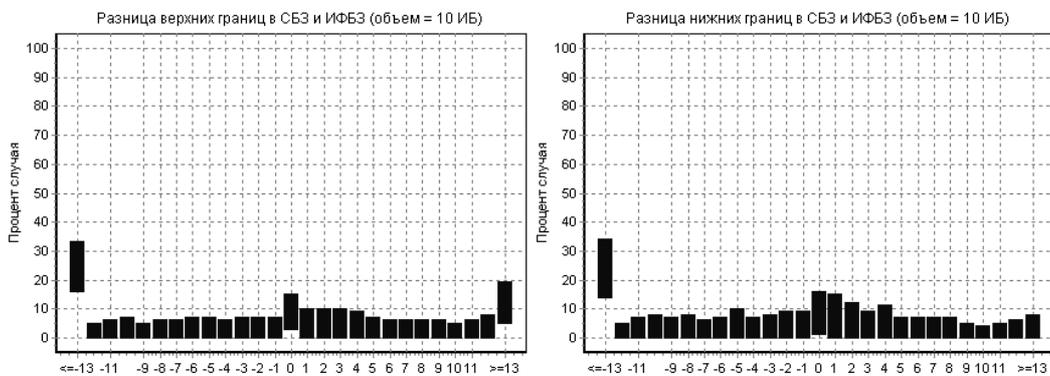


Рис.7. Отношения между границами периодов признаков, входящих в КК заболеваний, ЧПД которых в ИФБЗ и СБЗ совпали, при объеме обучающих выборок равном 10 ИБ каждого заболевания

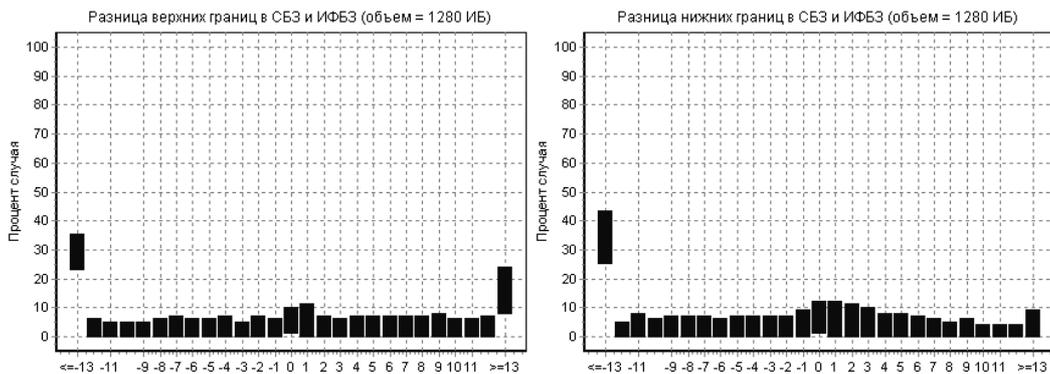


Рис.8. Отношения между границами периодов признаков, входящих в КК заболеваний, ЧПД которых в ИФБЗ и СБЗ совпали, при объеме обучающих выборок равном 1280 ИБ каждого заболевания

Описание: результаты сравнения, представленные на Рис.7 и Рис.8, свидетельствуют о том, что с ростом обучающих выборок значения разности верхних (нижних) границ в СБЗ и ИФБЗ распределены практически равномерно.

5. Заключение

1. Исследованный метод случайной расстановки границ периодов динамики с точки зрения временных затрат является очень быстрым и время работы при увеличении объемов обучающих выборок растет почти линейно в достаточно широком диапазоне. *Этот результат был ожидаем, т.к. методы Монте-Карло не предполагают больших переборных вариантов.*

2. Внешняя оценка метода получилась очень хорошей и к тому же с увеличением объемов обучающих выборок очень устойчивой. Используя только внешнюю оценку, можно сказать, что в этих условиях применения метод дает очень хорошие результаты. *В силу сложности баз знаний по структуре и объему, внешняя оценка ожидалась не столь хорошей.*
3. Когда решается задача ИФЗ с полной информацией (при сравнении областей значений признаков, не входящих в клинические картины заболеваний), внутренняя оценка оказывается очень хорошей. *Этот результат был ожидаем, т.к. такие области вычисляются по определяющим соотношениям.*
4. Когда решается задача с неполной информацией (при сравнении ЧПД признаков, входящих в клинические картины заболеваний, их областей значений, а также верхних и нижних границ периодов), внутренняя оценка оказывается очень плохой. Т.е. ИФЗ очень мало походит на СБЗ. Если рассматривать метод не с точки зрения практического применения, а с точки зрения выяснения того, как на самом деле устроена природа, то этот метод решает задачу плохо. *Предполагалось, что внутренняя сходимость исследуемого метода (сходимость по внутренней оценке) будет не высока.*
5. В общем случае проведенное экспериментальное исследование показывает, что хорошие внешние оценки метода ИФЗ могут сочетаться с его плохими внутренними оценками. *Этот результат оказался неожиданным.*

Опыт экспериментального исследования метода случайной расстановки границ периодов динамики показал, что дальнейшее экспериментальное изучение свойств методов ИФЗ для онтологий медицинской диагностики целесообразно проводить с учетом ограничений:

- в СБЗ должны генерироваться плохо разделяемые заболевания, за счет уменьшения количества признаков при значении интересного параметра *число периодов динамики*, равного единице, и введением условия на то, что области значений одних и тех же признаков у различных заболеваний должны существенно пересекаться,
- в онтологию должно быть добавлено соглашение о том, что в соседних периодах признака области его значений в этих периодах не пересекаются,
- наряду с полной обследованностью в экспериментах должна быть смоделирована неполная обследованность признаков в контрольных выборках.

6. Благодарности

Работа выполнена при финансовой поддержке ДВО РАН в рамках Программы №2 фундаментальных исследований Президиума РАН “Интеллектуальные информационные технологии, математическое моделирование, системный анализ и автоматизация”, проект “Развитие систем управления базами знаний с коллективным доступом”. Авторы выражают признательность некоммерческой организации “Благотворительный фонд культурных инициатив (Фонд Михаила Прохорова)” за финансирование поездки на конференцию.

Литература

- [1] Клещев А.С., Смагин С.В. Организация компьютерных экспериментов по индуктивному формированию знаний // НТИ. Сер. 2. – 2008. – №1. – С. 16-24.
- [2] Клещев А.С., Смагин С.В. Общий подход к проведению компьютерных экспериментов по индуктивному формированию знаний // Программные продукты и системы. – 2008. – №1. – С. 56-58. [<http://www.swsys.ru/index.php?page=article&id=101>]
- [3] Клещев А.С., Смагин С.В. Организация компьютерных экспериментов по индуктивному формированию знаний. Владивосток: ИАПУ ДВО РАН, 2007. 36 с. [<http://iacp.dvo.ru/is/publications/2007-Kleshev,Smagin-Organizing.pdf>]
- [4] Клещев А.С. Задачи индуктивного формирования знаний в терминах непримитивных онтологий предметных областей. // НТИ. – 2003. – Сер. 2. – № 8. – С. 8-18.
- [5] Клещев А.С., Смагин С.В. Распараллеливание вычислений при решении задачи индуктивного формирования баз знаний. // Искусственный интеллект. – 2006. – №3. – С. 421-428.
- [6] Соболев И.М. Метод Монте-Карло. – М.: Наука, 1968. – 64 с.
- [7] Клещев А.С., Смагин С.В. Компьютерный эксперимент по исследованию свойств метода случайной расстановки границ периодов динамики. Владивосток: ИАПУ ДВО РАН, 2009. 44 с. [<http://iacp.dvo.ru/is/publications/2009-Kleshev,Smagin-ExperOne.pdf>]