

В.С. Казанцев

О ВОЗМОЖНОСТЯХ ИСПОЛЬЗОВАНИЯ ПАКЕТА КВАЗАР ПРИКЛАДНЫХ ПРОГРАММ РАСПОЗНАВАНИЯ ОБРАЗОВ В РЕШЕНИИ ЗАДАЧ МЕДИЦИНЫ И ЗДРАВООХРАНЕНИЯ

ГАУДПО «Уральский институт управления здравоохранением имени А. Б. Блохина»,
г. Екатеринбург, Российская Федерация

Резюме. Введение. Обсуждается актуальность применения современных методов обработки и анализа медицинских данных. Акцентируется внимание на необходимости системного подхода к анализу данных и использовании многофакторного анализа на основе методов распознавания образов (РО). **Цель работы** — ознакомить врачей и специалистов по анализу медицинских данных с возможностями применения методов распознавания образов и пакета программ КВАЗАР для решения практических задач диагностики, прогнозирования и поиска закономерностей в материалах исследований и данных статистической отчетности. **Материалы и методы.** Методы распознавания образов являются эффективным средством решения сложных, плохо формализуемых задач и успешно используются при решении задач классификации, диагностики, прогнозирования, управления. Одним из программных средств решения задач распознавания образов является пакет КВАЗАР, разработанный в Институте математики и механики УрО РАН. Подробно рассматриваются возможности пакета по решению основных задач РО: обучения по прецедентам и таксономии, а также вопросы контроля качества обрабатываемых данных и заполнения пропусков в них. Материалами для решения задач РО с помощью пакета КВАЗАР служили данные государственной статистики, а также массивы наблюдений, полученные при выполнении специальных исследований. В двух приведенных примерах использовались данные о состоянии здоровья населения, демографии и ресурсном обеспечении систем здравоохранения муниципальных образований (МО) Свердловской области за 2016–2019 годы. **Результаты.** Перечислены наиболее интересные исследования, выполненные с помощью пакета КВАЗАР. Кроме того, в статье рассматривается задача построения и практического использования модели классификации МО области по смертности трудоспособного населения от заболеваний острыми нарушениями мозгового кровообращения (ОНМК). Приводятся результаты «проигрывания» на модели управленческих сценариев, направленных на снижение смертности от ОНМК. Также приводится пример использования пакета КВАЗАР для кластерного анализа муниципальных образований области по девяти статистическим показателям. **Выводы:**

1. Опыт исследований, выполненных с использованием пакета КВАЗАР, показывает, что методы РО являются эффективным средством решения задач классификации, диагностики, прогнозирования и могут успешно применяться при анализе данных в медицине и здравоохранении.
2. Результаты математического моделирования свидетельствуют о том, что увеличение доли госпитализаций по СМП больных с артериальной гипертензией и хронической ишемической болезнью сердца способствует снижению смертности трудоспособного населения от заболеваний ОНМК.
3. Многомерный таксономический анализ МО области по девяти показателям показал наличие двух крупных кластеров, различающихся уровнями общей смертности населения и показателями оснащенности системы оказания медицинской помощи.

Ключевые слова: анализ медицинских данных, распознавание образов, обучение по прецедентам, моделирование, таксономия, заполнение пропусков

Конфликт интересов отсутствует.

Контактная информация автора, ответственного за переписку:

Казанцев Владимир Сергеевич

kvs222@yandex.ru

Дата поступления 16.11.2022 г.

Образец цитирования:

Казанцев В.С. О возможностях использования пакета КВАЗАР прикладных программ распознавания образов в решении задач медицины и здравоохранения. [Электронный ресурс] Вестник уральской медицинской академической науки. 2022, Том 19, №5, с. 533–546, DOI: 10.22138/2500-0918-2022-19-5-533-546

Введение

Широкое внедрение информатизации и компьютерных технологий в медицину и здравоохранение способствуют накоплению больших объемов данных в электронном виде, грамотное использование и анализ которых могут оказать большую помощь в повышении качества диагностики заболеваний, лечении больных, решении различных задач управления здравоохранением в целом.

В последние два десятилетия в англоязычной литературе по обработке и анализу данных получил распространение термин «Data Mining» [1]. Обычно его переводят на русский язык как «извлечение информации» или «добыча данных». Основная суть подхода состоит в получении новой полезной информации из большого количества «сырых», необработанных данных. Существующие трактовки в определении «Data Mining» имеют свои различия, но в целом они сходятся в следующем:

- «Data Mining» — это обнаружение нетривиальных, практически полезных и интерпретируемых знаний, информации;
- эти знания нужны для принятия решений;
- источником этих знаний являются большие (иногда очень большие) массивы многомерных данных;
- алгоритмы извлечения этих знаний не сводятся к традиционным статистическим методам анализа;
- использование этих алгоритмов в силу их сложности и больших объемов обрабатываемых данных практически невозможно без применения современной вычислительной техники.

Методологическую базу этого направления составляют различные методы анализа многомерных данных, например, такие как множественный регрессионный анализ, факторный анализ, нейронные сети, временные ряды, деревья решений, кластерный анализ, распознавание образов и др. Эти методы находят успешное применение при решении сложных, плохо формализуемых задач медицинской диагностики, прогнозирования, анализа влияния различных факторов на заболеваемость и др. В большинстве случаев наиболее адекватными поставленным задачам являются методы, удовлетворяющие принципам системного подхода.

Сущность системного подхода состоит в том, что явления и процессы, протекающие в природе и обществе, различные вещи и предметы рассматриваются как целостные системы, состоящие из подсистем и определенного набора элементов. Множество элементов, находящихся в отношениях и связях между собой, образуют определенную целостность, единство. Поэтому представление об изучаемом объекте как о целостной системе является исходным пунктом системного подхода, всякого системного изучения [2].

Все это в полной мере относится к объектам и ситуациям, с которыми приходится иметь дело при решении многих сложных задач. Так, в задачах медицинской диагностики, прогнозирования или выбора метода лечения объектом исследования является человек – чрезвычайно сложная система, содержащая большое количество взаимосвязанных подсистем и элементов. При решении задач эпидемиологического анализа работают с комплексом факторов, обуславливающих заболеваемость тех или иных групп населения, и т. д. Полноценное исследование таких систем невозможно без привлечения соответствующих многофакторных методов анализа. К числу таких методов, в частности, относится распознавание образов.

Материалы и методы

1. Распознавание образов и пакет КВАЗАР

Распознавание образов можно охарактеризовать как одно из направлений искусственного интеллекта, связанное с моделированием способности человека узнавать (классифицировать) различные объекты и явления окружающего мира, а также обучаться этому на основе наблюдения соответствующих объектов и явлений. Методы распознавания образов являются эффективным средством решения

сложных, плохо формализуемых задач и успешно используются при решении задач классификации, диагностики, прогнозирования, управления.

В последние годы много говорят о нейронных сетях. Нейронные сети — это технология, программное или программно-аппаратное моделирование нейронной сети человеческого мозга и его функции классифицировать объекты окружающего мира, а также обучаться этому на примерах. Можно встретить много информации о различных применениях нейронных сетей, но, по сути дела, все эти применения сводятся к способностям сетей решать задачи классификации.

По отношению к распознаванию образов нейронная сеть — это один из подходов, один из алгоритмов решения задач обучения. В распознавании образов для этого используются также другие методы и алгоритмы. Это наиболее близкие по своей идеологии к нейронным сетям методы комитетного распознавания [3, 4], а также алгоритмы на основе метода потенциальных функций [5], методы логического распознавания образов [6], а также вероятностные, алгебраические и некоторые другие.

Проблематика распознавания образов включает две основные задачи:

1) задача «обучения с учителем», которую также называют задачей обучения по прецедентам или задачей дискриминантного анализа;

2) задача «обучения без учителя», называемая также задачей автоматической классификации, таксономии, кластерного анализа.

В настоящее время существует большое количество программных средств, предназначенных для решения задач распознавания образов. Одним из них является пакет КВАЗАР, созданный в Институте математики и механики УрО РАН [7]. Пакет предназначен для решения задач классификации, диагностики и прогнозирования на основе использования методов таксономии и дискриминантного анализа. Для построения правил классификации на основе принципов обучения с учителем в пакете используются пять различных алгоритмов:

- 1) алгоритм обучения с использованием однородных комитетов большинства [3, 8];
- 2) алгоритм обучения с использованием комитетов старшинства [9];
- 3) рекуррентный алгоритм линейного разделения выпуклых оболочек двух множеств [10];
- 4) алгоритм обучения на основе метода потенциальных функций [5];
- 5) алгоритм классификации методом k ближайших соседей [6].

С помощью пакета можно также оценивать качество и достоверность материала наблюдений, заполнять пропуски в данных, решать задачи регрессионного анализа, анализа межгрупповых различий средних арифметических значений показателей с помощью T -критерия Стьюдента, рассчитывать отдельные статистические характеристики материала наблюдений.

Рассмотрим здесь более подробно содержательные постановки основных задач распознавания образов, а также некоторые вопросы их практического использования в прикладных исследованиях.

2. Задача обучения по прецедентам

Это одна из наиболее важных задач в распознавании образов, позволяющая «обучить» компьютер (или правильнее сказать, «заложенную» в него специальную программу) узнаванию каких-либо объектов, явлений или ситуаций на примерах. Узнать объект — это значит классифицировать его, то есть отнести к некоторому классу объектов (или образу). Близкими здесь являются понятия диагностики и прогнозирования, поскольку и то и другое из них — это тоже узнавание, классификация, хотя и есть некоторая специфика в их использовании. Так, термин «диагностика» употребляется обычно в тех случаях, когда говорят о классификации состояния, в котором находится в настоящее время тот или иной объект. При прогнозировании же ставится целью указать некоторые будущие состояния объекта.

Существуют, по крайней мере, два способа обучения классификации объектов. Один из них состоит в том, что обучаемому (человеку или компьютеру) сообщаются готовые правила (алгоритмы) классификации, которыми он в последующем сможет воспользоваться. Однако такой способ обучения иногда оказывается неэффективным или вообще не применимым из-за отсутствия готовых, известных правил классификации.

Второй способ — обучение классификации на примерах (прецедентах) — состоит в том, что обучаемому предъявляется некоторое количество примеров объектов, которые ему предстоит научиться классифицировать, и обучаемый должен сам выработать правило классификации. В ряде случаев второй способ оказывается более эффективным. Для обучения компьютера классификации объектов существуют специальные алгоритмы и программы распознавания образов. Обучение компьютера со-

стоит в «показе» объектов, которые необходимо научиться распознавать с указанием при этом их классификации, то есть указанием того, к какому классу принадлежит данный объект. Слово «показ» здесь взято в кавычки, поскольку, как правило, компьютеру показывается не сам объект (хотя это и возможно при решении задач обучения распознавания зрительных образов). Обычно при обучении компьютеру «показывается» некая числовая модель объекта. Как правило, это числовой вектор (строка чисел), координатами которого являются представленные в виде чисел некоторые характеристики (признаки) реального объекта.

Примером наиболее эффективного применения РО в медицине является решение задач дифференциальной диагностики в ситуациях, когда врач не располагает точным алгоритмом постановки диагноза, потому что на современном уровне, возможно, такого алгоритма просто нет. Еще относительно недавно, до начала применения компьютерной томографии, к таким задачам можно было отнести задачу диагностики характера мозгового инсульта, когда при различной этиологии заболевания нередко внешние его проявления были схожи настолько, что без специальных методов обследования (ангиография, пунктирование) точный диагноз поставить было сложно. Задачи такого рода встречаются и сейчас.

На практике для решения таких задач формируется обучающая выборка, включающая в себя примеры диагностируемых заболеваний. Каждый конкретный случай заболевания описывается числовым вектором, представляющим собой последовательность значений некоторого заранее выбранного набора признаков. Векторы вводятся в память компьютера и в определенном порядке предъявляются специальной программе обучения классификации. При этом программе сообщается, к какому заболеванию каждый из них относится. Программа должна построить диагностическое (решающее) правило, которое в последующем, может быть использовано для дифференциальной диагностики новых случаев заболеваний. Решающее правило может быть записано математически в виде дискриминантной (разделяющей) функции или оформлено в виде распознающего алгоритма. Для оценки качества полученного решающего правила предварительно организуется распознавание проверочной, или экзаменующей, выборки, включающей векторы с известной классификацией, но не участвовавшие в процедуре обучения. По результатам экзамена принимается решение о возможности практического использования решающего правила.

Геометрически суть задачи дискриминантного анализа можно пояснить следующим образом. Предположим, решается задача построения решающего правила для дифференциальной диагностики двух схожих по своим проявлениям заболеваний, представленных для обучения двумя наборами числовых векторов. Каждый n -мерный вектор может быть представлен в виде точки в многомерном пространстве влияющих признаков (факторов). В результате в этом пространстве будет образовано «облако» точек, соответствующее одному из диагностируемых заболеваний, и второе «облако» точек — другому. Задача дискриминантного анализа состоит в том, чтобы в n -мерном пространстве построить и математически описать границу, разделяющую два множества. При этом множества точек в пространстве признаков могут быть хорошо разделимы, и в этом случае граница будет простой, может быть, даже линейной, а могут и взаимно пересекаться; в этом случае граница будет сложной. Математическое описание разделяющей множества точек границы и есть искомая математическая модель, которую в последующем можно использовать для моделирования различных управляющих воздействий.

Рассмотренный метод может успешно применяться не только для диагностики состояния больных, но и при анализе здоровья населения на популяционном уровне. Как есть больные люди, так могут быть, например, и «больные территории» в том смысле, что наряду с благополучными территориями (муниципальными образованиями) в плане того или иного вида заболеваемости или смертности, как правило, есть и неблагополучные. Каждая территория также может быть представлена моделью в виде числового вектора, координатами которого могут быть значения различных показателей, характеризующих уровень медицинского обслуживания населения, демографическую и социально-экономическую ситуацию и т. д. Так, признаками в таком векторе могут быть, например, значения показателей обеспеченности населения врачами, койками, мощности поликлиник и другие.

В случае построения качественных решающих правил (обеспечивающих не менее 80–90% правильного распознавания векторов экзаменующей выборки) возможно использование их для моделирования результатов тех или иных действий, направленных на улучшение здоровья населения. Так, после проведения соответствующих процедур обучения с помощью построенной математической модели

можно получить ответ на вопрос «как изменится смертность от болезней системы кровообращения, если обеспеченность врачами-кардиологами будет такой-то?» или, например, «как повлияет на уровень заболеваемости раком легкого снижение доли курильщиков в той или иной популяции?» и т. д. Воспользовавшись рассмотренной выше геометрической моделью разделения двух множеств в n -мерном пространстве, суть «проигрывания» можно интерпретировать следующим образом.

Положение каждой точки в пространстве определяется числовым вектором, которым представлен объект (например, территория). И если изменить какую-либо координату (то есть значение соответствующего влияющего фактора) этого вектора, положение точки в n -мерном пространстве также изменится. Поэтому, если во всех «неблагополучных» векторах увеличить, например, численное значение показателя обеспеченности врачами-кардиологами, то и точки «неблагополучного» множества изменят свое положение в пространстве признаков и, возможно, часть из них окажется по другую сторону границы, то есть на стороне «благополучного» множества. Установить это как раз и позволяет математическая модель разделения множеств. Остаётся только подсчитать процент «неблагополучных» точек, перешедших на сторону «благополучных», чтобы оценить эффект от изменения показателя обеспеченности врачами.

3. Задача таксономии

Определенный научный и практический интерес при анализе данных может представлять многомерная классификация наблюдений с помощью алгоритмов таксономии. Суть задачи таксономии состоит в том, чтобы разбить имеющееся множество объектов, описываемых заданным набором признаков, на некоторое число групп (кластеров), содержащих близкие между собой объекты. Содержательный смысл задачи популярно объясняет в [6] Н. Г. Загоруйко, приводя слова Демокрита, написанные еще во II в. до нашей эры в «Письме ученому соседу»: «Если тебе, дорогой друг, нужно разобраться в сложном нагромождении фактов или вещей, ты сначала разложи их на небольшое число куч по схожести. Картина прояснится, и ты поймешь природу этих вещей».

Чтобы лучше понять смысл задачи автоматической классификации многомерных наблюдений, рассмотрим алгоритм «корреляционных плеяд», предложенный П. В. Терентьевым [11] для кластеризации параметров (столбцов таблицы) на основе использования коэффициентов парной корреляции в качестве меры близости между ними. В последующем метод стал применяться и для таксономии числовых векторов (объектов). Рассмотрим его кратко применительно к таксономии векторов.

1. Задается некоторое пороговое расстояние R . Если координатами векторов являются вещественные числа, для вычисления расстояния между ними обычно используют метрику Евклида. При работе с векторами, координаты которых представлены в двоичном виде, применяется метрика Хемминга.

2. Выбирается один из векторов таксономируемого множества, например, первый из них, и включается в таксон.

3. В таксономируемом множестве ищутся векторы, удаленные от выбранного не больше, чем на R . Если таких векторов не обнаружено, формирование этого таксона заканчивается. Если такие векторы обнаружены, они включаются в таксон.

4. Затем для каждого из вновь включенных в таксон векторов таким же образом ищутся ближайшие и тоже включаются в таксон и т. д. до тех пор, пока близкие векторы перестанут находиться. На этом формирование таксона заканчивается.

5. Если остались векторы, не включенные в таксон, выбирается один из таких векторов, с которого начнется формирование нового таксона в соответствии с пунктами 3 и 4 и так до тех пор, пока все векторы не будут включены в таксоны. При этом некоторые таксоны могут состоять из одного вектора.

Количество таксонов, получаемых в результате работы алгоритма, зависит от величины R . При очень малых значениях R каждый таксон будет состоять из одного объекта, при слишком больших R все объекты могут объединиться в один таксон. Важной особенностью алгоритма является то, что результат разбиения множества векторов на таксоны не зависит от выбора начального вектора. Рис. 1 иллюстрирует результат работы алгоритма при заданном пороговом расстоянии R .

Данный алгоритм таксономии реализован в пакете КВАЗАР в виде циклической процедуры, где таксономия последовательно выполняется при различных значениях порогового расстояния R . При этом в качестве начального значения R выбирается минимальное расстояние между векторами в таксономируемом множестве. Обычно в результате выполнения первого шага работы алгоритма образуется один таксон, содержащий два ближайших вектора. Остальные векторы множества при этом остаются

не сгруппированными и представляют собой «единичные» таксоны. Следует заметить, что начальное расстояние R может равняться нулю, если в множестве присутствуют одинаковые векторы. В этом случае в процессе выполнения первого шага алгоритма могут быть получены один или несколько таксонов с одинаковыми между собой векторами. Остальные векторы также при этом образуют «единичные» таксоны. После завершения работы алгоритма при начальном значении R пороговое расстояние увеличивается на некоторую заданную или автоматически выбираемую величину ΔR , после чего выполняется следующий шаг таксономии и так далее. По мере увеличения порогового расстояния таксоны сливаются, укрупняются, пока, наконец, все множество векторов не объединится в один таксон. При решении практических задач таксономии представляется интересным анализ связей в получаемых таксонах, а также средних значений отдельных признаков в таксонах. В пакете КВАЗАР предусмотрена возможность проведения такого анализа. Пример практической задачи таксономии с анализом полученных таксонов приведен ниже.

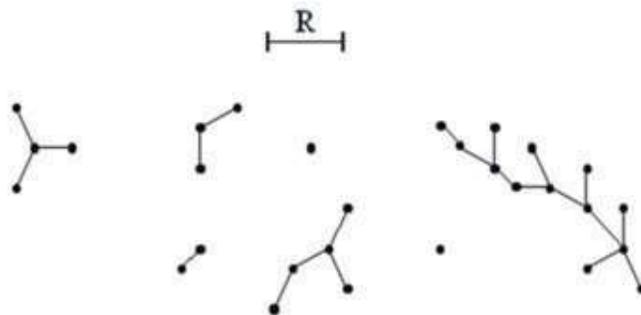


Рис. 1. Вариант таксономии объектов при заданном значении R .
Fig. 1. A variant of object taxonomy for a given value of R .

4. Заполнение пропусков и анализ качества данных

Результаты решения рассмотренных выше задач распознавания образов в значительной мере зависят от качества используемых данных, поэтому предварительному анализу и подготовке материалов наблюдений следует уделять самое серьезное внимание. Частой проблемой при сборе значительных объемов информации является отсутствие отдельных значений в таблицах обрабатываемых данных, то есть наличие так называемых пропусков, которые иногда встречаются в экспериментальных данных, в результатах анкетирования и т. д. Нередко также таблицы наблюдений содержат недостоверные, ошибочные данные. В настоящее время существуют различные методы заполнения пропусков и выявления недостоверных данных; некоторые из них рассмотрены ниже.

Природа пропусков может быть различной. Рассмотрим этот вопрос на примере данных о работе коечного фонда кардиологической службы муниципальных образований в некотором регионе, представленных в виде таблицы. Предположим, каждая строка таблицы представляет набор среднегодовых значений показателей работы коечного фонда в том или ином муниципальном образовании, в числе которых есть, например, показатель обеспеченности кардиологическими койками в расчете на 10 тысяч жителей данной территории, показатель работы койки (дней в году), показатель летальности на кардиологической койке и другие. Возможное отсутствие отдельных значений данных показателей в таблице может быть обусловлено разными причинами. Так, пропуск значения показателя обеспеченности койками для некоторой территории может возникнуть из-за того, что соответствующие данные по каким-то причинам не были предоставлены. Другой причиной пропуска значения этого показателя может быть фактическое отсутствие кардиологических коек в медицинском учреждении того или иного населенного пункта, и в этом случае на месте пропуска следовало бы поставить значение «0». Можно говорить и еще об одной причине присутствия пропусков в таблице: любое значение показателя работы койки для муниципального образования, не имеющего кардиологических коек, было бы бессмысленно, а поэтому соответствующая клетка таблицы в таком случае может быть оставлена пустой. То же самое относится и к показателю летальности. Для обозначения пропусков такого рода обычно используют различные условные обозначения, например, символ “-“. Другие же пропуски могут быть заполнены с помощью специальных методов и алгоритмов.

В некоторых случаях наличие пропусков не является препятствием для обработки данных. Так,

средние значения показателей могут быть рассчитаны и при наличии пропусков. Иногда проблему решают путем удаления строк и столбцов, содержащих пропуски. Нередко используют простейшие методы заполнения пропусков средними значениями показателей, рассчитываемых по всему множеству строк или значениями, взятыми из ближайшей строки. При оценке расстояния между строками обычно используется евклидова метрика. Более сложные методы заполнения пропусков основываются на применении кластерного, регрессионного, корреляционного анализа. Все методы заполнения пропусков данных подразделяют на глобальные и локальные. В основе глобальных методов лежит поиск и использование общих закономерностей в таблице наблюдений. Локальные же методы предполагают использование предсказывающих матриц, содержащих близлежащие элементы.

Широкую известность получило семейство локальных алгоритмов ZET, разработанных в Институте математики Сибирского отделения РАН [6]. В основе этих алгоритмов лежит предположение об информационной избыточности числовых массивов, организованных в виде таблиц. При этом имеется в виду, что между столбцами таблицы часто существуют корреляционные зависимости, а строки в той или иной мере также могут быть похожими. Чем сильнее взаимосвязи в таблице, тем более точно могут быть восстановлены в ней пропущенные значения. Идеальным примером числового массива с высокой информационной избыточностью может служить таблица умножения. Пропуски в ней могут быть заполнены безошибочно. Обратным примером может служить таблица, полученная с использованием генератора случайных чисел. Заполнить пропуски в такой таблице не сможет даже самый совершенный алгоритм.

Для заполнения пропуска алгоритм ZET строит предсказывающую матрицу из элементов строк и столбцов, наиболее близких к строке и столбцу, на пересечении которых имеется пропуск. В качестве меры близости между столбцами таблицы рассчитываются коэффициенты парной корреляции, а поиск ближайших строк производится по величине евклидова расстояния между ними. Далее строятся линейные регрессионные зависимости, связывающие известные элементы строки и столбца матрицы, содержащих пропущенное значение, с другими строками и столбцами предсказывающей матрицы. В случае построения зависимостей, обеспечивающих заданную точность предсказания известных элементов предсказывающей матрицы, регрессионная модель может быть использована для прогнозирования пропущенного значения.

Алгоритм ZET включен в пакет ОТЭКС [6], разработанный под руководством известного специалиста в области распознавания образов профессора Н. Г. Загоруйко. Модификация этого алгоритма также используется и в пакете КВАЗАР в виде программного модуля SPACE.

Следует отметить, что алгоритм SPACE может использоваться как для заполнения пропусков, так и для оценки качества обрабатываемых данных. С этой целью в пакете предусмотрен специальный режим работы, при котором элементы таблицы по очереди «объявляются» пропусками, и алгоритм предсказывает их значения. Затем производится сравнение имеющихся в таблице значений с предсказанными и рассчитываются оценки степени достоверности отдельных элементов, а также соответствующие усредненные оценки для строк, столбцов и всей таблицы в целом.

5. Материалы исследований

В разделе «Результаты» приводятся примеры многочисленных исследований, выполненных с использованием пакета КВАЗАР. Материалами для решения соответствующих задач служили данные государственной статистики, а также массивы наблюдений, полученные при проведении специальных исследований.

В разделе также рассмотрены результаты двух исследований, выполненных в рамках данной работы. Материалами для их выполнения явились статистические данные о состоянии здоровья населения муниципальных образований Свердловской области за 2016-2019 годы, а также данные о демографии и ресурсном обеспечении системы здравоохранения этих территорий. Более подробно эти материалы описаны при изложении соответствующих результатов.

Результаты

1. Примеры исследований, выполненных с помощью пакета КВАЗАР

Пакет КВАЗАР был разработан в 1982 году и с тех пор активно используется при решении сложных, плохо формализуемых задач классификации, диагностики, прогнозирования в самых различных областях применения. В качестве примеров успешного применения пакета в медицине и здравоохранении

можно назвать некоторые выполненные в разные годы с его помощью исследования:

- дифференциальная диагностика характера инсульта головного мозга;
- прогнозирование риска послеоперационных осложнений при желчнокаменной болезни;
- прогнозирование патологических состояний у новорожденных;
- прогнозирование результатов лечения рака гортани;
- оценка отдаленных последствий влияния Восточно-Уральского радиоактивного следа на состояние организма жителей г. Каменска-Уральского;
- комплексная оценка влияния различных факторов на развитие рака легких, желудка, молочной железы;
- многофакторный анализ канцерогенной опасности радона;
- комплексный анализ факторов, определяющих психическое здоровье населения Курганской области;
- прогнозирование результатов лечения отслойки сетчатки с использованием силиконовой тампоны;
- оценка генетической и индивидуальной предрасположенности к различным заболеваниям.

В качестве примеров исследований, выполненных с использованием пакета в течение последних нескольких лет, можно привести такие, как:

- таксономический анализ медико-социальных характеристик пациентов с хроническими облитерирующими заболеваниями артерий нижних конечностей [12];
- многофакторный анализ причин приобщения школьников к алкоголю и табаку [13];
- оценка взаимосвязи здоровья и образа жизни родителей со здоровьем и образом жизни их детей [14];
- оценка влияния кадровых ресурсов муниципального здравоохранения на уровень смертности населения [15];
- оценка влияния деятельности межмуниципальных медицинских центров на смертность населения от болезней системы кровообращения [16].

Рассмотрим здесь также результаты решения двух задач многофакторного анализа статистических данных о здоровье и медицинском обслуживании населения муниципальных образований Свердловской области.

2. Моделирование результатов управленческих сценариев на основе решающих правил

С помощью пакета КВАЗАР решалась задача построения и использования модели классификации муниципальных образований Свердловской области по смертности трудоспособного населения от заболеваний острыми нарушениями мозгового кровообращения (ОНМК). При анализе использовались данные за три года, предшествующих пандемии covid19 (2016–2019 гг.). Для решения задачи необходимо было сформировать группы благополучных и неблагополучных в плане смертности трудоспособного населения от ОНМК территорий. В разные годы одна и та же территория может быть как относительно благополучной, так и неблагополучной. Поэтому если имеются данные по смертности от ОНМК в муниципальных образованиях (МО) области за несколько лет, то изучаемой единицей и объектом включения в выборку правильно будет считать не «МО», а «МО в таком-то году». Таким образом, если мы, к примеру, располагаем данными по 59 МО за 3 года, то в нашем распоряжении будет $59 \times 3 = 177$ объектов. В данном случае в благополучную группу были включены 78 территорий с уровнем смертности от 8,4 до 30,0 случаев на 100 тыс. человек в год, в неблагополучную — такое же количество территорий с уровнем смертности от 40,0 до 123,6 случаев на 100 тыс. человек. Территории со значениями показателя смертности от ОНМК от 30,0 до 39,99 на 100 тыс. человек при анализе не использовались с целью повышения различимости классов при решении задачи обучения. Поскольку выделенным группам приписывается смысл классов территорий «условно благополучных» и «условно неблагополучных», то целесообразность исключения МО, имеющих средние значения показателя смертности, становится очевидной.

Каждая территория была представлена 24-мерным числовым вектором своих характеристик (признаков). В их число входили 12 показателей медицинского обслуживания населения, включая показатели ресурсного обеспечения кардиологической и неврологической служб, 4 показателя госпитализации по СМП больных с болезнями системы кровообращения, а также 3 демографических показате-

ля и 5 показателей социально-экономического благополучия территорий. С помощью специального алгоритма все множество векторов было разбито на обучающую и экзаменующую выборки. Объем экзаменующей выборки составил 20 векторов. Остальные векторы вошли в обучающую выборку. Задача решалась с помощью пакета КВАЗАР. По результатам анализа информативности признаков для решения задачи обучения были отобраны 10 наиболее информативных признаков и, соответственно, вся дальнейшая работа проводилась с 10-мерными векторами. Для построения решающего правила использовался алгоритм обучения на основе метода потенциальных функций. В результате было построено решающее правило, безошибочно распознающее как материал обучения, так и экзаменующую выборку. Высокое качество решающего правила позволило использовать его для оценки эффективности некоторых управленческих сценариев. В частности, одна из поставленных задач состояла в оценке влияния уровня госпитализации по СМП больных с артериальной гипертензией (АГ) и хронической ишемической болезнью сердца (ХИБС) на смертность от ОНМК лиц трудоспособного возраста. С этой целью было выполнено моделирование влияния на смертность соответствующих факторов, представленных следующими показателями:

- доля госпитализированных по СМП при АГ, всего, в %;
- доля госпитализированных по СМП при ХИБС, всего, в %.

Диапазон значений первого из этих показателей варьировал в неблагоприятной группе от 11,8% до 100% при среднем арифметическом значении 41,2% и значении медианы — 38,9%. При этом в 8-ми муниципальных образованиях значение показателя не превышало 10%, а в 24-х МО оно было меньше 30%.

По второму показателю диапазон значений составлял 0–85,6% при среднем арифметическом 23,9% и значении медианы 20,5%. При этом в 16-ти МО значение показателя было меньше 10%, а в 39-ти — меньше 20%.

Таблица 1

Потенциальная результативность некоторых управленческих сценариев, направленных на снижение смертности трудоспособного населения от ОНМК населения муниципальных образований Свердловской области

Table 1

Potential effectiveness of some management scenarios aimed at reducing the mortality of the working-age population from CVAs in the municipalities of the Sverdlovsk region

Сценарий/ Scenario	Сокращение группы МО с высоким уровнем смертности трудоспособного населения от ОНМК после коррекции факторов, в %/ Reduction in the group of municipalities with a high CVA mortality level in the working-age population after a correction of factors, %
Увеличение минимальной доли госпитализированных по СМП при артериальной гипертензии/ Reduction in the group of municipalities with a high CVA mortality level in the working-age population after a correction of factors, %	
- до/ up to 30%	3,85%
- до/ up to 40%	10,26 %
- до/ up to 50%	14,10%
Увеличение минимальной доли госпитализированных по СМП при хронической ишемической болезни сердца/ Increase in the minimum proportion of emergency hospitalizations for chronic coronary heart disease	
- до/ up to 25%	3,85 %
- до/ up to 30%	6,41 %

Представлялось интересным посмотреть, как увеличение значений данных показателей для территорий с низкими их значениями повлияло бы на показатель смертности трудоспособного населения от ОНМК.

Моделирование производилось с помощью специальной программы, функция которой состояла в заданной корректировке векторов неблагоприятной группы и последующей их классификации с помощью ранее полученного решающего правила. Корректировка векторов состояла в замене значений показателей госпитализации более высокими, задаваемыми в режиме диалога с программой. Так, например, при задании значения «25», этим числом заменялись соответствующие координаты векторов

с меньшими значениями, то есть, если значение показателя госпитализации по СМП при АГ в некотором МО составляло 14,2%, это значение в соответствующем векторе заменялось на 25. После выполнения всех необходимых корректировок векторы неблагополучной группы классифицировались с помощью решающего правила и подсчитывался процент «перехода» в благополучную группу.

Результаты моделирования управленческих сценариев, направленных на снижение численности территорий с высоким уровнем смертности трудоспособного населения от ОНМК, приведенные в таблице 1, показывают, что рост госпитализации по СМП больных с сердечно-сосудистыми заболеваниями способствовал бы снижению уровня данного вида смертности.

3. Кластерный анализ территорий Свердловской области по девяти показателям

Задача решалась пакетом КВАЗАР с помощью метода таксономии на основе алгоритма корреляционных плеяд. Таксономируемое множество состояло из 59 векторов, содержащих данные об общей смертности населения в муниципальных образованиях (МО) Свердловской области в 2019 году, а также некоторые демографические показатели и показатели медицинского обслуживания населения МО области. В частности, при таксономии использовались следующие показатели:

1. Общая смертность населения, на 1000 человек
2. Доля мужского населения, в %
3. Доля лиц старше трудоспособного возраста, в %
4. Доля городского населения, в %
5. Обеспеченность врачами-физлицами, всего, на 10 тыс. человек
6. Укомплектованность врачами-физлицами, в %
7. Удельный вес врачей с высшей квалификационной категорией, в %
8. Удельный вес врачей-организаторов здравоохранения с высшей квалификационной категорией, в %
9. Обеспеченность больничными койками, на 10 тыс. человек

Таксономия проводилась в соответствии с изложенным выше алгоритмом. В одном из вариантов разбиения было получено два крупных таксона, включающих 21 и 13 векторов. Также в этом варианте был получен таксон из трех векторов и два таксона, содержащих по два вектора каждый, а остальные 18 векторов представляли собой единичные таксоны. Интересно было проанализировать, в чем состоят наиболее существенные различия между МО, входящими в два самых крупных таксона. С этой целью был выполнен расчет средних арифметических значений рассматриваемых показателей в этих таксонах (табл. 2). Соответственно количеству векторов, вошедших в каждый таксон, в таблице они имеют названия «Таксон 21» и «Таксон 13».

В «Таксон 13» входят, в основном, крупные городские округа. В частности, в него вошли города Нижний Тагил, Каменск-Уральский, Первоуральск, Верхняя Пышма, Березовский и другие. Средняя доля городского населения территорий данного таксона составляет 84,3%. Более многочисленный «Таксон 21» со средней долей городского населения 57,4% включает в себя менее крупные города и районы. Для этого таксона характерно незначительное превышение средних значений показателей доли мужского населения (46,8 против 45,7% в «Таксоне 13») и лиц старше трудоспособного возраста (28,7 и 27,1% соответственно). Значительно более существенные различия между территориями, входящими в эти два таксона, отмечаются по показателям медицинского обслуживания населения. Так, средняя обеспеченность врачами-физическими лицами по территориям, входящим в «Таксон 13», составляет 21,4 на 10 тысяч населения, в то время как по территориям «Таксона 21» она значительно ниже: 16,0 на 10 тыс. Средние значения укомплектованности штатов врачами-физическими лицами также различается: 63,1 и 53,0% соответственно. Различен и уровень квалификации врачебного персонала. Так, доля врачей с высшей квалификационной категорией для территорий «Таксона 13» в среднем составляет 30,8%, а для территорий «Таксона 21» — 21,6%. Еще больше данные таксоны различаются по среднему уровню квалификации врачей-организаторов здравоохранения: 31,9 и 1,2% соответственно. Что касается обеспеченности населения больничными койками, то и здесь среднее значение величины показателя в «Таксоне 13» составляет 50,2 на 10 тысяч населения, в то время как для «Таксона 21» эта величина равна 37,4.

Таблица 2
Средние значения показателей в таксонах
Table 2
Average values of indices in taxa

№	Показатель/ Index	Среднее значение показателя в таксоне/ Mean value of the index in taxa	
		Таксон «13»/ Taxon «13»	Таксон «21»/ Taxon «21»
1	Общая смертность населения, на 1000 человек/ Total mortality of the population, per 1000 people	14,1	16,0
2	Доля мужского населения, в %/ Proportion of the male population, %	45,7	46,8
3	Доля лиц старше трудоспособного возраста, в %/ Proportion of people older than working age, %	27,1	28,7
4	Доля городского населения, в %/ Proportion of the urban population, %	84,3	57,4
5	Обеспеченность врачами-физлицами, всего, на 10 тыс. человек/ Medical service density, doctors per 10 000 people	21,4	16,0
6	Укомплектованность врачами-физлицами, в %/ Staffing level for doctors, %	63,1	53,0
7	Удельный вес врачей с высшей категорией, в %/ Proportion of doctors with highest category, %	30,8	21,6
8	Удельный вес врачей-организаторов здравоохранения с высшей категорией, в %/ Proportion of health care organizers with highest category, %	31,9	1,2
9	Обеспеченность больничными койками, на 10 тыс. человек/ Provision of hospital beds, per 10000 people	50,2	37,4

Выводы

1. Опыт многочисленных исследований, выполненных с использованием пакета КВАЗАР, показывает, что методы распознавания образов являются эффективным средством решения задач классификации, диагностики, прогнозирования, поиска зависимостей и могут успешно применяться при анализе данных в медицине и здравоохранении.

2. Результаты математического моделирования, выполненного в данной работе, свидетельствуют о том, что увеличение доли госпитализаций по СМП больных с артериальной гипертензией и хронической ишемической болезнью сердца (ХИБС) способствует снижению смертности трудоспособного населения от заболеваний ОНМК. Так, увеличение до 50% минимальной доли госпитализированных по СМП при артериальной гипертензии в муниципальных образованиях с высокой смертностью (текущее среднее значение показателя равно 41,2%) позволило бы снизить последнюю примерно на 14%, а эффект от соответствующего увеличения минимальной доли госпитализаций при ХИБС до 30% (текущее среднее — 23,9%) был бы равен примерно 6,4%.

3. Многомерный таксономический анализ муниципальных образований Свердловской области по девяти показателям показал наличие двух крупных кластеров, различающихся уровнями общей смертности населения и оснащенностью системы оказания медицинской помощи. Один такой кластер включает в себя 13 муниципальных образований. В основном, это крупные города с относительно низкими показателями общей смертности и достаточно высоким уровнем обеспеченности населения высококвалифицированным врачебным персоналом и больничными койками. Второй кластер включает 21 территорию со средней долей городского населения 57,4%, более высокой смертностью и более низкими показателями медицинского обслуживания.

ЛИТЕРАТУРА

1. Тучкова А.С., Кондрашева П.П. Термин «Data mining». Задачи, решаемые методами data mining // Тенденции развития науки и образования. – 2019. – №55-2. – С.27-30. DOI: 10.18411/lj10-2019-26
2. Гасников В.К. Основы научного управления и информатизации в здравоохранении (учебное посо-

бие). Ижевск: изд-во «Вектор», 1997. - 169с.

3. Мазуров Вл.Д. Метод комитетов в задачах оптимизации и классификации. - М.: Наука. - 1990. - 245 с.
4. Mazurov V.D., Polyakova E.Y. Committees: history and applications in machine learning. Communications in Computer and Information Science. 2019. V. 1090. pp. 3-16. DOI: 10.1007/978-3-030-33394-2_1
5. Аркадьев А.Г., Браверман Э.М. Обучение машин классификации объектов. - М.: Наука. - 1971. - 192 с.
6. Загоруйко Н.Г. Прикладные методы анализа данных и знаний. Новосибирск, Изд-во Ин-та математики, 1999. - 270 с.
7. Казанцев В.С. Задачи классификации и их программное обеспечение. - М.: Наука. - 1990. - 135 с.
8. Тягунов Л. И. Алгоритмы комитетного распознавания и их применение для решения практических задач классификации: Автореф. дис. ... канд. техн. наук. Свердловск, 1973. 24 с.
9. Osborne W.L. The seniority logic – a logic for committee machine // IEEE Trans. Comput.- 1977.- V. C-26, N 12.- P.1302 – 1306.
10. Козинец В. Н. Рекуррентный алгоритм разделения выпуклых оболочек двух множеств // Алгоритмы обучения распознаванию образов/Под ред. В. Н. Вапника. М.: Сов. радио, 1972. С. 43-50.
11. Терентьев П. В. Метод корреляционных плеяд // Вестн. ЛГУ. Биология. 1959. № 9. С. 137-141.
12. Казанцев В.С., Михайлова Д.О., Погосян В.А. Использование метода таксономии для анализа медико-социальных характеристик пациентов с хроническими облитерирующими заболеваниями артерий нижних конечностей в Свердловской области // Современные проблемы здравоохранения и медицинской статистики. 2019. № 5. С. 73-74.
13. Липанова Л.Л., Насыбуллина Г.М., Казанцев В.С. Распространенность потребления школьниками психоактивных веществ и многофакторный анализ причин приобщения к алкоголю и табаку // Профилактическая и клиническая медицина. 2019. № 1 (70). С. 4-9.
14. Липанова Л.Л., Насыбуллина Г.М., Бабилова А.С., Казанцев В.С. Взаимосвязь здоровья и образа жизни родителей со здоровьем и образом жизни их детей // В книге: Актуальные проблемы образования и здоровья обучающихся. Монография. Под редакцией В.И. Стародубова, В.А. Тутельяна. Москва, 2020. С. 429-446.
15. Алленов А.М., Соловьев И.Р., Казанцев В.С. Оценка влияния кадровых ресурсов муниципального здравоохранения на уровень смертности населения // сб. ст. по материалам XIV Международной научно-практической конференции «Современная медицина: новые подходы и актуальные исследования». – № 8(12). – М., Изд. «Интернаука», 2018. – С. 32-38.
16. Мальков Н.А., Казанцев В.С., Столбиков С.А. Оценка влияния деятельности межмуниципальных медицинских центров на смертности населения от болезней системы кровообращения // Естественные науки и медицина: теория и практика: сб. ст. по матер. VII-VIII междунар. науч.-практ. конф. № 2-3(5). – Новосибирск: СибАК, 2019. – С. 36-41.

Автор:

Казанцев Владимир Сергеевич

ГАУДПО «Уральский институт управления здравоохранением имени А.Б. Блохина»

Кандидат технических наук, старший научный сотрудник, ведущий научный сотрудник

Российская Федерация, 620075, г. Екатеринбург, ул. Карла Либкнехта, 86

kvs222@yandex.ru

V.S. Kazantsev

ON THE POSSIBILITIES OF USING THE KVAZAR PATTERN RECOGNITION SOFTWARE PACKAGE IN SOLVING MEDICINE AND HEALTHCARE PROBLEMS

A.B. Blokhin Ural Institute of Health Management,
Yekaterinburg, Russian Federation

Abstract. Introduction. The relevance of the application of modern methods of processing and analysis of medical data is discussed. Attention is focused on the need for a systematic approach to data analysis and the

use of multivariate analysis based on pattern recognition methods. **The aim** of the paper is to inform medical professionals and specialists in the analysis of medical data about the possibilities of using pattern recognition methods and the KVAZAR software package in solving practical problems of diagnostics, forecasting and searching for patterns in research materials and statistical reporting data. **Materials and methods.** Pattern recognition methods are an effective tool for dealing with complex, poorly formalized problems and are successfully used in solving problems of classification, diagnostics, forecasting, and control. The KVAZAR software package, which was developed at the Institute of Mathematics and Mechanics, Ural Branch of the Russian Academy of Sciences, is a tool for solving pattern recognition problems. The possibilities of the package for solving the basic pattern recognition problems, in particular, learning from examples, taxonomy, and issues of controlling the quality of the processed data and filling gaps in them, are considered in detail. The materials for solving pattern recognition problems using KVAZAR were government statistics data as well as arrays of observations data obtained during special studies. In the two examples given in the paper, data on the state of health of the population, demography, and resources of health care systems of municipalities of the Sverdlovsk region for 2016–2019 were used. Results. The most interesting studies carried out with the help of KVAZAR are reported. The paper also deals with the problem of the construction and practical use of a model that classifies the municipalities of the region according to the mortality of the working-age population from acute cerebrovascular accidents (CVAs). The results modeling some management scenarios aimed at reducing mortality from CVAs are presented. An example of using KVAZAR for cluster analysis of municipalities of the region according to nine statistical indicators is presented. **Conclusions:** 1. The analysis performed with the use of the KVAZAR package shows that pattern recognition methods are an effective tool for solving problems of classification, diagnostics, and forecasting and can be successfully applied in data analysis in medicine and public health. 2. The results of mathematical modeling prove that an increase in the proportion of emergency hospitalizations of patients with arterial hypertension and chronic coronary heart disease helps to reduce the mortality of the working-age population from CVAs. 3. A multivariate taxonomic analysis of the municipalities based on nine indices has shown the presence of two large clusters that differ in the general mortality of the population and the level of equipment of the medical care system.

Keywords: medical data analysis, pattern recognition, learning from examples, modeling, taxonomy, gap filling

There is no conflict of interest.

Contact information of the author responsible for correspondence:

Vladimir S. Kazantsev

kvs222@yandex.ru

Received 16.11.2022

For citation:

Kazantsev V.S. On the possibilities of using the KVAZAR pattern recognition software package in solving medicine and healthcare problems. [Online] Vestn. Ural. Med. Akad. Nauki. = Journal of Ural Medical Academic Science. 2022, Vol. 19, no. 5, pp. 533–546. DOI: 10.22138/2500-0918-2022-19-5-533-546 (In Russ)

REFERENCES

1. Tuchkova A.S., Kondrasheva P.P. The term «data mining». Problems solved by data mining methods. Trends in the development of science and education, 2019, No. 55-2, P. 27-30. (in Russ) DOI: 10.18411/lj10-2019-262
2. Gasnikov V.K. Fundamentals of scientific management and informatization in health care (textbook). Vektor, Izhevsk, 1997. (in Russ)
3. Mazurov V.I.D. The method of committees in optimization and classification problems. Nauka, Moscow, 1990. 248 p. (in Russ)
4. Mazurov V.D., Polyakova E.Y. Committees: history and applications in machine learning. Communications in Computer and Information Science. 2019. Vol. 1090, pp. 3-16. DOI: 10.1007/978-3-030-33394-2_1
5. Arkadiev A.G., Braverman E.M. Machine learning for object classification. Nauka, Moscow, 1971. (in Russ)

6. Zagoruiko N.G. Applied methods of data and knowledge analysis. Novosibirsk, Izd. Inst. Mat., 1999. (in Russ)
7. Kazantsev V.S. Problems of classification and software for them. Nauka, Moscow, 1990, 135 p. (in Russ)
8. Tyagunov L.I. Algorithms of committee recognition and their application for solving practical problems of classification: Abstract of candidate's dissertation in technical sciences. Sverdlovsk, 1973, 24 p. (in Russ)
9. Osborne W.L. The seniority logic – a logic for committee machine. IEEE Trans. Comput., 1977, Vol. C-26, No. 12, pp. 1302-1306.
10. V. N. Kozinets, A recursive algorithm for separating convex hulls of two sets. Algorithms for pattern recognition, Ed. by V.N. Vapnik. Sov. radio, Moscow, 1972, pp. 43-50. (in Russ)
11. Terentiev P.V. Method of correlation pleiades. Vestn. LGU. Biol. 1959. No. 9, pp. 137–141. (in Russ)
12. Kazantsev V.S., Mikhailova D.O., Pogosyan V.A. Using the taxonomy method to analyze the medical and social characteristics of patients with chronic obliterating diseases of the arteries of the lower extremities in the Sverdlovsk region. Modern problems of healthcare and medical statistics. 2019, No. 5, pp. 73-74. (in Russ)
13. Lipanova L.L., Nasybullina G.M., Kazantsev V.S. Prevalence of consumption of psychoactive substances by schoolchildren and multivariate analysis of the reasons for initiation to alcohol and tobacco. Preventive and Clinical Medicine, 2019, No. 1 (70), pp. 4-9. (in Russ)
14. Lipanova L.L., Nasybullina G.M., Babikova A.S., Kazantsev V.S. The relationship of health and lifestyle of parents with the health and lifestyle of their children. In the book: Actual problems of education and health of students. Monograph. Edited by V.I. Starodubov, V.A. Tutelyan, Moscow, 2020, pp. 429-446. (in Russ)
15. Allenov A.M., Soloviev I.R., Kazantsev V.S. Evaluation of the impact of human resources of municipal health care on the level of mortality of the population. Collection of papers based on the materials of the XIV International Scientific and Practical Conference «Modern Medicine: New Approaches and Current Research». No. 8(12), Moscow, Internauka, 2018, pp. 32-38. (in Russ)
16. Malkov N.A., Kazantsev V.S., Stolbikov S.A. Evaluation of the impact of the activities of intermunicipal medical centers on the mortality of the population from diseases of the circulatory system. Natural Sciences and Medicine: Theory and Practice. Collection of papers based on the materials of VII-VIII International Scientific and Practical Conference. No. 2-3(5), Novosibirsk, SibAK, 2019, pp. 36-41. (in Russ)

Author

Vladimir S. Kazantsev

A.B. Blokhin Ural Institute of Health Management

Candidate of Technical Sciences, Senior Researcher, Leading Researcher

8b st. Karl Liebknecht, Yekaterinburg, Russian Federation, 620075

kvs222@yandex.ru