
**ПРИБОРОСТРОЕНИЕ ДЛЯ БИОЛОГИИ
И МЕДИЦИНЫ**

УДК 543.51+ 681.2–5

© Л. В. Новиков, В. В. Манойлов, А. Г. Кузьмин, Ю. А. Титов,
И. В. Заруцкий, А. О. Нефедов, А. В. Нефедова, А. И. Арсеньев, 2020

**ЭКСПРЕСС-ДИАГНОСТИКА ЗАБОЛЕВАНИЙ
ПО ВЫДЫХАЕМОМУ ВОЗДУХУ
НА ОСНОВЕ КВАДРУПОЛЬНОГО МАСС-СПЕКТРОМЕТРА**

Данная работа является развитием исследований, опубликованных в журнале "Научное приборостроение" № 1, 2019. В сравнении с указанной статьей в данной работе расширен литературный обзор, впервые разработан и апробирован алгоритм расчета вероятности заболеваний. Приводятся результаты обработки данных больных, находящихся на лечении в двух онкологических клиниках. Расчет вероятности заболевания по данным масс-спектрометрического анализа выдыхаемого воздуха основан на определении принадлежности масс-спектра тестируемого больного соответствующей контрольной группе. Каждая контрольная группа формируется путем набора массива спектров не менее десяти пациентов с одним и тем же заболеванием. Диагностика выполняется путем преобразования матрицы спектров контрольной группы и спектра тестируемого больного в пространство главных компонент. Вероятность заболевания определяется по евклидову расстоянию координат больного от центроида контрольной группы в многомерном пространстве главных компонент.

Кл. сл.: экспресс-диагностика заболевания, метод главных компонент, многомерная плотность вероятности, обработка многомерных данных

ВВЕДЕНИЕ

Диагностика заболеваний, основанная на анализе дыхания, в последние годы широко используется в клинической практике [1, 2] с применением самых разнообразных методик и аналитических приборов [3]. Такая диагностика предполагает выявление заболеваний на ранних стадиях, что является перспективным. Установлено, что в выдыхаемом воздухе содержатся важные биомаркеры для ранней диагностики целого ряда заболеваний: сердечно-сосудистых, онкологических, нейродегенеративных, органов дыхания, диабета и др.

В работах [4, 5] исследуются биомаркеры дыхания при онкологических заболеваниях легких, в том числе у курильщиков. Летучие вещества в образцах предварительно концентрировали путем твердофазной микроэкстракции (SPME) с последующим анализом в тандеме газовый хроматограф – масс-спектрометр (GC-MS). Полученные концентрации анализировались в пространстве главных компонент (PCA) с использованием первых трех координат (PCA1, PCA2, PCA3).

В работе [6] приведен подробный обзор результатов исследований методов диагностики путем анализа количественного и качественного составов выдыхаемого воздуха (ВВ), рассмотрены

требования к аналитическим методам и приборам исследования многокомпонентных газовых смесей. Отмечено, что анализ на основе ВВ имеет ряд преимуществ перед диагностикой, основанной на лабораторных методах, — он является относительно дешевым, занимает немного времени и позволяет обнаруживать детектируемые вещества в минимальных концентрациях. Он не предполагает инвазивных вмешательств и может проводиться с любой кратностью, предоставляя возможность тщательного изучения динамики физиологических процессов. Приведен ряд маркеров самых различных заболеваний: органов дыхания, желудочно-кишечного тракта, онкологии, центральной нервной системы и др. По исследованиям, проведенным в США, ранняя диагностика заболеваний, в частности органов дыхания, позволила снизить смертность на 20 %. Определены требования к средствам обнаружения молекул-маркеров с учетом разработанных методик и пределов обнаружения: возможность одновременного измерения и идентификации нескольких биомаркеров; высокая чувствительность и высокая точность определения концентраций (на уровне 10 ppb – 10 ppt); селективность (регистрация и идентификация веществ на фоне высоких концентраций азота, кислорода, воды и углекислого газа); быстрое действие (время проведения анализа с усредне-

нием по нескольким выдохам около 5–10 с); простота использования в клинических условиях; цена "одного замера" и стоимость прибора. В качестве приборов для анализа ВВ рассматривают газовую хроматографию (ГХ), масс-спектрометрию (МС), масс-спектрометрию с газохроматографическим разделением (МСГХ) и ИК-спектроскопию (ИКС).

Авторы работы [7] отмечают необходимость разработки удобного, простого в использовании, неинвазивного метода для скрининга хронического заболевания почек и качества диализа в режиме реального времени. Для этого предложено измерять концентрацию аммиака в выдыхаемом воздухе, однако использование масс-спектрометрии, спектрометрии подвижности ионов, оптических методов лазерного поглощения и др. оборудования дорого, сложно и неудобно для клинического использования. Предложено использовать недавно разработанные газовые сенсоры — полупроводниковый сенсорный чип на основе полимера с нанопорами.

Для сравнения результатов обследования нескольких групп пациентов использовались различные статистические подходы: парный t-тест Стьюдента, критерий Вилкоксона, тест Краскела – Уоллиса, тест Дана, тест Фридмана. Коэффициент корреляции Пирсона и простая линейная регрессия использовались для оценки корреляции между концентрациями выдыхаемого аммиака и др. диагностическими измерениями.

В работе [8] исследуется возможность диагностики эпилепсии путем сравнения состава выдыхаемого воздуха у больных и здоровых пациентов с использованием "электронного носа". Данные прибора обрабатывались нейронной сетью, которая была обучена данными предварительно отобранных пациентов с заболеванием эпилепсией и ее отсутствием. Исследование эффективности диагностики эпилепсии по ВВ показали ее несомненные достоинства по сравнению с традиционными тестами с использованием электроэнцефалограмм.

В работе [9] анализируются различные методы и приборы для сбора пробы в полевых условиях с последующим анализом в хроматографе с масс-спектрометром в качестве детектора, снабженным приставкой для термодесорбции (TD-GC-MS анализ). Статистическая обработка результатов экспериментов по измерению, в частности концентрации изопрена, проводилась с использованием t-теста Стьюдента и смешанного линейного моделирования.

В работе [10] показано, что метод анализа выдыхаемого воздуха с использованием газового хроматографа – спектрометра ионной подвижности (GC-IMS) пригоден для ранней неинвазивной диагностики нейродегенеративных заболеваний,

в частности болезни Альцгеймера, характерной для пожилых людей. Комплекс является высокопроизводительным диагностическим инструментом для диагностики в реальном времени в клинических условиях.

Существует острая необходимость в разработке эффективных диагностических инструментов, особенно таких, которые позволяют надежно выявлять заболевания на ранних стадиях (перед проведением сложных лабораторных исследований) с использованием неинвазивных подходов. Ключевой проблемой в анализе ВВ является надлежащий статистический анализ и интерпретация больших и разнородных наборов данных, полученных из исследований ВВ. В настоящей работе использованы квадрупольный масс-спектрометр и метод обработки результатов масс-спектрометрического анализа ВВ, основанный на определении вероятности заболевания путем определения принадлежности масс-спектра тестируемого пациента соответствующей контрольной группе.

КВАДРУПОЛЬНЫЙ МАСС-СПЕКТРОМЕТР

Характеристики

Квадрупольный масс-спектрометр, применяемый в работе для анализа выдыхаемых газов, описан в работе [11]. При малых габаритах и весе он обеспечивает обнаружение компонентов ВВ в диапазоне масс — от 1 до 200 а.е.м., разрешение по массовым числам — 0.5, скорость регистрации — до 1 масс-спектра в 10 с с чувствительностью по концентрациям примесей — от 0.1 ppm. Он прост в использовании, обеспечивает высокую скорость анализа, начиная с забора воздуха до получения результата. Время получения результатов анализа составляет не более минуты.

Краткое описание работы прибора

Анализируемый газ при давлении атмосферы через капиллярный ввод подается в камеру ионизации источника ионов с электронным ударом. Образовавшиеся ионы вводятся в масс-анализатор квадрупольного типа. Получившиеся в процессе регистрации масс-спектрометрические сигналы обрабатываются с помощью специализированного программного обеспечения и сравниваются со спектрами в библиотеке стандартных масс-спектров, затем проводится идентификация отдельных компонент спектра и определение их концентрации. Прогреваемая капиллярная система ввода пробы в масс-спектрометр позволяет проводить анализ на расстоянии до 5 м от прибора. Расход пробы составляет около 5 мкл в 1 с. В вакуумной

системе используется турбомолекулярный и мембранный насосы.

Важным преимуществом прибора являются его малые габариты (см. [11]), что позволяет его оперативное перемещение из одного медицинского учреждения в другое для проведения скрининговых обследований больших групп населения.

ОБРАБОТКА ДАННЫХ

При массовом обследовании состояния здоровья населения, как правило, отсутствует предварительная информация о характере заболевания пациента. При проведении диагностики состояния здоровья пациентов методом анализа выдыхаемого воздуха целесообразно регистрировать весь спектр выдыхаемых газов (маркеров заболеваний) и затем сопоставлять его с данными одной из контрольных групп с известным заболеванием. Данные ВВ каждого пациента этой группы можно представить в виде точки в многомерном пространстве. Объединенные данные всех пациентов группы образуют "облако" исходных параметров. Чем ближе точка, определяющая спектр тестируемого больного, находится к центру (центроиде) этого "облака", тем больше вероятность того, что пациент либо страдает заболеванием контрольной группы, либо здоров, если это контрольная группа здоровых пациентов.

Пусть $x_{i,j}$ — j -й параметр (в данном случае — интенсивность спектральной компоненты) i -го пациента контрольной группы, причем $i = 1, 2, \dots, I$ и $j = 1, 2, \dots, J$, и набор из значений для I пациентов по J регистрируемых компонент ВВ образуют (I, J) обучающую матрицу \mathbf{X} , столбцы которой обозначим как \mathbf{X}_j : $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_j, \dots, \mathbf{X}_J]$ [12].

Вектор \mathbf{X}_j является случайной величиной. Для простоты предположим, что компоненты этого вектора имеют нормальное распределение со средним значением \bar{X}_j и дисперсией σ_j^2 . Тогда элементы матрицы \mathbf{X} образуют в J -мерном пространстве "облако" (обозначим его как \mathbf{G}) из $I \times J$ точек с центроидой в точке $\bar{\mathbf{X}}$, координаты которой равны $\bar{\mathbf{X}} = [\bar{X}_1, \bar{X}_2, \dots, \bar{X}_J]$. Тестируемый пациент страдает тем же заболеванием (или здоров), что и контрольная группа, если измеренный вектор параметров ВВ $\mathbf{X}_d = [x_{d,1}, x_{d,2}, \dots, x_{d,J}]$, где $x_{i,j}$ — значение параметра, принадлежит J -мерному пространству \mathbf{G} , т.е. $\mathbf{X}_d \in \mathbf{G}$. Это условие выполняется, если вероятность P отклонения точки \mathbf{X}_d от центроиды $\bar{\mathbf{X}}$ не превышает некото-

рый порог α . Для определения этой вероятности построим многомерное распределение вероятности принадлежности события пространству \mathbf{G} , предполагая, что это распределение подчиняется нормальному закону [13]:

$$P(\mathbf{X}_d) = W \cdot \exp \left\{ -\frac{1}{2} (\mathbf{X}_d - \bar{\mathbf{X}})' \mathbf{K}_X^{-1} (\mathbf{X}_d - \bar{\mathbf{X}}) \right\}, \quad (1)$$

$$\mathbf{X}_d \in \mathbf{G},$$

где \mathbf{K}_X — ковариационная матрица $\mathbf{K}_X = E \left[(\mathbf{X} - \bar{\mathbf{X}}) \cdot (\mathbf{X} - \bar{\mathbf{X}})' \right]$, E — символ математического ожидания, $(\cdot)'$ — символ транспонирования матрицы, W — нормирующий множитель. Очевидно, что вероятность P этого события равна единице, если для очередного пациента окажется, что $\mathbf{X}_d \equiv \bar{\mathbf{X}}$. Это условие выполняется, если в формуле (1) положить $W = 1$. Тогда условие принадлежности пациента контрольной группе имеет вид: $P(\mathbf{X}_d) < \alpha$, где величина α выбирается методом экспертной оценки.

Однако непосредственное использование формулы (1) для расчета величины P сопряжено со значительными вычислительными трудностями, связанными с наличием большого числа параметров J и корреляционных связей между столбцами матрицы \mathbf{X} . Для сжатия данных, сокращения размерности пространства \mathbf{G} используют ортогональное преобразование данных в пространство главных компонент — метод главных компонент (МГК) [14].

Для перехода в пространство ГК формируется новая матрица, состоящая из всех строк матрицы \mathbf{X} и строки \mathbf{X}_d . Обозначим эту матрицу как $\mathbf{X}\mathbf{I}$. Тогда в новой системе координат:

$$\mathbf{X}\mathbf{I} = \mathbf{T} \cdot \mathbf{P}' + \mathbf{e} = \sum_{j=1}^A \mathbf{t}_j \mathbf{p}_j' + \mathbf{e},$$

где \mathbf{p}_j — собственные функции ковариационной матрицы \mathbf{K}_X . Матрицу \mathbf{T} называют матрицей *счетов* $\mathbf{T} = [\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_A]$, ее размерность — $(I \times A)$; матрицу \mathbf{P} называют матрицей *нагрузок*, ее размерность — $(I \times A)$; \mathbf{e} — это матрица *остатков* (шумов) размерности $(I \times J)$; векторы-столбцы \mathbf{T}_j ($j = (1, 2, \dots, A)$) называют *главными компонентами* (ГК), A — число главных компонент. Величина A значительно меньше числа переменных J . Это означает, что основная информация о состоянии пациента сосредоточена в нескольких первых ГК.

Последняя строка этой матрицы, вектор \mathbf{T}_d — координаты параметров тестируемого пациента

в пространстве ГК: $\mathbf{T}_d = [t_{d,1}, t_{d,2}, \dots, t_{d,A}]$.

Средние значения столбцов матрицы \mathbf{T} равны нулю, а дисперсия — вектор σ^2 с элементами $\sigma_j^2 = \lambda_j$, т.е. собственными числами ковариационной матрицы. Свойство разложения по ГК таково, что дисперсия быстро уменьшается уже к четвертой ГК, а столбцы матрицы \mathbf{T} некоррелированы, т.е.

$$\mathbf{T}_m \mathbf{T}'_n = \begin{cases} 0 & \text{при } n \neq m, \\ \lambda_j & \text{при } n = m. \end{cases}$$

Учитывая это обстоятельство, в новой системе координат формула (1) имеет вид:

$$P(\mathbf{T}_d) = \exp\left\{-\frac{1}{2} \mathbf{T}_d \sigma^{-2} \mathbf{T}'_d\right\} = \exp\left\{-\frac{1}{2} \sum_{j=1}^A \frac{t_{d,j}^2}{\sigma_j^2}\right\}, \quad (2)$$

и евклидово расстояние от пациента с индексом d до центроиды обучающей группы равно

$$D(d) = \sqrt{\sum_{j=1}^A t_{d,j}^2}. \quad (3)$$

ОПИСАНИЕ АЛГОРИТМА

Обработка данных состоит из двух этапов: *обучение и диагностика*.

На этапе *обучения* формируется K матриц \mathbf{X}^k , $k = 1, 2, \dots, K$ с данными по интенсивности спектра ВВ для каждого вида заболеваний. При этом число строк I^k каждой матрицы должно быть в несколько раз больше числа столбцов J^k .

На этапе *диагностики* выполняется следующая последовательность операций:

1. Измеряются параметры ВВ диагностируемого пациента и формируется вектор-строка $\mathbf{X}_d = \{x_{d,1}, x_{d,2}, \dots, x_{d,J}\}$.

2. Строится новая матрица $\mathbf{X}\mathbf{I}^k = [\mathbf{X}^k; \mathbf{X}_d]$ размерности $((I^k + 1) \times J^k)$, последней строкой которой является вектор X_d .

3. Выполняется нормировка этой матрицы, например, на максимальное значение ее элементов.

4. Используя алгоритм метода ГК, вычисляется матрица счетов \mathbf{T}^k и вектор собственных чисел λ_j^k [15].

5. Определяется число первых столбцов A матрицы \mathbf{T}^k из условия $\sigma_A^k = \sqrt{\lambda_A^k} < \varepsilon$. При этом без потери достоверности можно положить

$\varepsilon = 0.001-0.01$. Последняя строка матрицы $\mathbf{T}_d^k = \{t_{d,1}^k, t_{d,2}^k, \dots, t_{d,A}^k\}$ — главные компоненты параметров ВВ диагностируемого пациента.

6. Определяется вероятность $P(\mathbf{T}_d^k)$ по формуле (2).

7. Переход к п. 2 ($k = k + 1$) для определения вероятности принадлежности тестируемого пациента другой группе заболеваний.

8. Анализ результатов вычисления вероятности $P(\mathbf{T}_d^k)$, $k = 1, 2, \dots, K$ для всех K заболеваний с целью определения наиболее вероятного.

ПРОВЕРКА АЛГОРИТМА

Проиллюстрируем изложенную выше теорию на одной группе из 43 больных, с диагностированной онкологией одного типа¹⁾. Эта группа разделена на две: обучающую (36 пациентов) и контрольную (7 пациентов). При этом концентрация компонентов ВВ первого пациента из контрольной группы для примера совпадает с центроидой обучающей. Привлечена к проверке эффективности алгоритма еще одна контрольная группа²⁾ с другими онкологическими заболеваниями из 10 больных. Пациентам обеих контрольных групп присвоены порядковые номера соответственно $d = (1, 2, 3, \dots, 7)$ и $d = (1, 2, 3, \dots, 10)$.

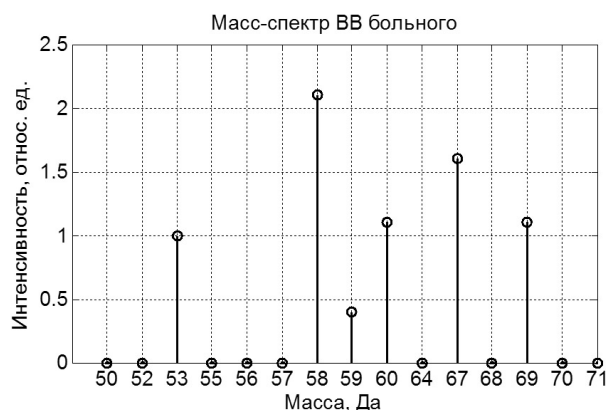


Рис. 1. Масс-спектр ВВ обследуемого пациента. Масса 53 — изопрен, 58 — ацетон, 60 — уксус, 67 — изопрен (?), 69 — пентин; 59 — некий нераспознанный пока биомаркер, выявленный у некоторых пациентов с онкологией

¹⁾ Данные о ВВ этой группы пациентов предоставлены ГБУЗ "Санкт-Петербургский клинический научно-практический центр специализированных видов медицинской помощи (онкологический), (СПб КНпЦСВМП(о)).

²⁾ Данные о ВВ этой группы пациентов предоставлены Ленинградским областным клиническим онкологическим диспансером.

На рис. 1 показан пример масс-спектра ВВ одного из пациентов первой (обучающей) группы. Из каждой контрольной группы поочередно выбираем данные по ВВ и выполняем п. 1–6 алгоритма. При этом оказывается, что достаточно положить $A=6$, т.к. $\sigma_7^k \cong 0$. По формулам (3) и (2) вычисляются евклидово расстояние от центроиды обучающей выборки и вероятность диагностики заболевания онкологией контрольного пациента. Для каждой контрольной группы эти результаты заносятся в табл. 1 и 2. На рис. 2 в пространстве первых двух ГК (T1 и T2) изображены координаты пациентов обучающей (o) и первой контрольной (+) группы. Из данных табл. 1 видно, что по мере удаления координат ВВ от центроиды для этой контрольной группы вероятность заболевания онкологией уменьшается, но остается большой. Следует обратить внимание, что эта вероятность зависит также от места расположения координат пациента в шестимерном (в данном случае) про-

Табл. 1. Измеренная вероятность болезни для пациентов первой контрольной группы

Пациент	Евклидово расстояние	Вероятность заболевания
1	0.0116	0.9854
2	0.1076	0.1888
3	0.1499	0.2055
4	0.2044	0.1164
5	0.2554	0.2864
6	0.2892	0.3458
7	0.3247	0.1197

Табл. 2. Измеренная вероятность болезни для пациентов второй контрольной группы

Пациент	Евклидово расстояние	Вероятность заболевания
1	2.8791	0.0159
4	4.5739	0.0000
3	4.7337	0.0000
2	5.5927	0.0000
5	5.6321	0.0000
8	5.7129	0.0000
9	5.7322	0.0000
6	5.8520	0.0000
10	5.8635	0.0000
7	5.9168	0.0000

странстве ГК: чем больше проекция на первые ГК, тем больше эта вероятность.

На рис. 3 в пространстве первых двух ГК (T1 и T2) изображены координаты пациентов обучающей (o) и второй контрольной группы (+).

Оказалось, что все пациенты этой группы (кроме первого) находятся вне "облака" обучающей выборки. Из данных табл. 2 видно, что вероятность наличия заболевания онкологией пациентов этой группы ничтожно мала.

ЗАКЛЮЧЕНИЕ

Система диагностики заболеваний по ВВ, состоящая из квадрупольного масс-спектрометра

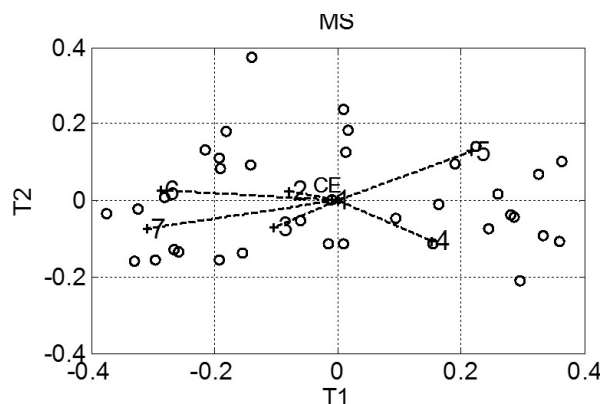


Рис. 2. Пациенты обучающей (o) и первой контрольной (+) групп пациентов в пространстве первых двух ГК. CE — положение центроиды, 1–7 — номера пациентов

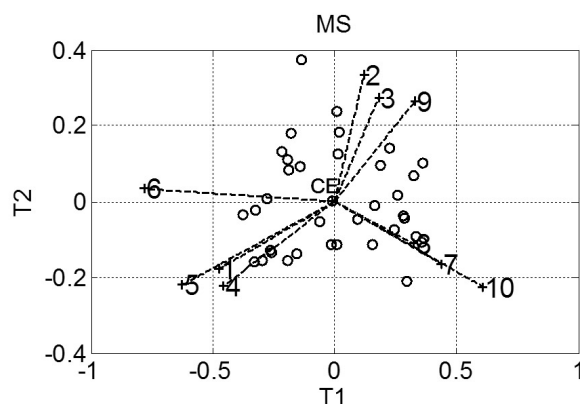


Рис. 3. Пациенты обучающей (o) и второй контрольной (+) групп пациентов в пространстве первых двух ГК. CE — положение центроиды, 1–10 — номера пациентов

и комплекта обработки данных с возможностью обучения и определения вероятности заболевания, обеспечивает неинвазивный экспресс-анализ состояния больного. Диагностика с определением вероятности принадлежности больного одной из обучающих групп позволяет за несколько десятков секунд составить общую картину вероятных заболеваний пациента и квалифицированно назначить дальнейшее обследование.

СПИСОК ЛИТЕРАТУРЫ

1. Cao W., Duan Y. Breath analysis: potential for clinical diagnosis and exposure assessment // *Clinical Chemistry*. 2006. Vol. 52, no. 5. P. 800–811. DOI: 10.1373/clinchem.2005.063545
2. Pereira J., Porto-Figueira P., Cavaco C., Taunk K., Rapole S., Dhakne R., Nagarajaram H., Cãmara J.S. Breath analysis as a potential and non-invasive frontier in disease diagnosis: an overview // *Metabolites*. No. 5. 2014. P. 3–55. DOI: 10.3390/metabo5010003
3. Kononov A., Korotetsky B., Jahatspanian I., Gubal A. et al. Online breath analysis using metal oxide semiconductor sensors (electronic nose) for diagnosis of lung cancer // *Journal of Breath Research*. 2020. Vol. 14, no. 1. DOI: 10.1088/1752-7163/ab433d
4. Kischkel S., Miekisch W., Sawacki A., Straker E.M., et al. Breath biomarkers for lung cancer detection and assessment of smoking related effects — confounding variables, influence of normalization and statistical algorithms // *Clinica Chimica Acta*. 2010. Vol. 411, is. 21–22. P. 1637–1644.
- 5.
6. Вакс В.Л., Домрачева Е.Г., Собакинская Е.А., Черняева М.Б. Анализ выдыхаемого воздуха: физические методы, приборы и медицинская диагностика // *УФН*. 2014. Т. 184, № 7. С. 739–758. DOI: 10.3367/UFNr.0184.201407d.0739
7. Chen I Ch.-Ch., Hsieh Ju-Ch., Chao Ch.-H., Yang W.-Sh., et al. Correlation between breath ammonia and blood urea nitrogen levels in chronic kidney disease and dialysis patients // *Journal of Breath Research*. 2020. Vol. 14, no. 3. DOI: 10.1088/1752-7163/ab728b
8. van Dartel D., Schelhaas H.J., Colon A.J., Kho K.H., de Vos C.C. Breath analysis in detecting epilepsy // *Journal of Breath Research*. 2020. Vol. 14, no. 3. DOI: 10.1088/1752-7163/ab6f14
9. Harshman S.W., Pitsch R.L., Davidson Ch.N., Scott A.M., et al. Characterization of standardized breath sampling for off-line field use // *Journal of Breath Research*. 2019. Vol. 14, no. 1. DOI: 10.1088/1752-7163/ab55c5
10. Tielel A., Wicaksono A., Daulton E., Ifeachor E., Eyre V., et al. Breath-based non-invasive diagnosis of Alzheimer's disease: A pilot study // *Journal of Breath Research*. 2020. Vol. 14, no. 2. DOI: 10.1088/1752-7163/ab6016
11. Манойлов В. В., Кузьмин А. Г., Заруцкий И.В., Титов Ю.А., Самсонова Н.С. Методы обработки и исследование возможностей классификации масс-спектров выдыхаемых газов // *Научное приборостроение*. 2019. Т. 29, № 1. С. 106–110. URL: <http://iairas.ru/mag/2019/abst1.php#abst16>
12. Большаков А.А. Методы обработки многомерных данных и временных рядов / А.А. Большаков, Р.Н. Каримов. М.: Горячая линия-Телеком, 2007. 522 с.
13. Ширяев А.Н. Вероятность. Том 1. М.: Изд-во МЦНМО, 2007. 552 с.
14. Айвазян С.А., Бухштабер В.М., Енюков И.С., Мешалкин Л.Д. Прикладная статистика. Классификация и снижение размерности. М.: Финансы и статистика, 1989. 607 с.
15. Эсбенсен К. Анализ многомерных данных / пер. с англ. под ред. О.Е. Родионовой. Черноголовка: Изд-во ИПХФ РАН, 2005. 161 с.

Институт аналитического приборостроения РАН, Санкт-Петербург (Новиков Л.В., Манойлов В.В., Кузьмин А.Г., Титов Ю.А., Заруцкий И.В.)

Санкт-Петербургский клинический научно-практический центр специализированных видов медицинской помощи (онкологический) (Нефедов А.О., Нефедова А.В.)

Научно-исследовательский центр онкологии им. Н.Н. Петрова Минздрава России, Санкт-Петербург (Арсеньев А.И.)

Контакты: Манойлов Владимир Владимирович, manoilov_vv@mail.ru

Материал поступил в редакцию 15.10.2020

EXPRESS DIAGNOSTICS OF DISEASES BASED ON A QUADRUPOLE MASS SPECTROMETER ANALYSIS OF EXHALED AIR

L. V. Novikov¹, V. V. Manoilov¹, A. G. Kuzmin¹,
Yu. A. Titov¹, I. V. Zarutsky¹, A. O. Nefedov²,
A. V. Nefedova², A. I. Arseniev³

¹*Institute for Analytical Instrumentation of RAS, Saint Petersburg, Russia*

²*St. Petersburg Clinical Scientific and Practical Center for Specialized Types
of Medical Care (Oncological), Russia*

³*Scientific Research Center of Oncology named after N.N. Petrov
of Ministry of Health of Russia, Saint Petersburg*

The method is proposed for express diagnostics of diseases according to the data of mass spectrometric analysis of exhaled air. An algorithm for calculating the probability of diseases has been developed and tested. The results of data processing of patients treated in two oncological clinics are presented. The calculation of the probability of disease according to the data of mass spectrometric analysis of exhaled air is based on attributing the mass spectrum of the tested patient to the mass-spectra of the corresponding control group. Each control group is formed by collecting an array of spectra from at least ten patients with the same disease. Diagnostics is performed by transforming the matrix of spectra of the control group and the spectrum of a patient being tested into the space of the principal components. The probability of a disease is determined by the Euclidean distance of the patient's coordinates from the centroid of the control group in the multidimensional space of these principal components.

Keywords: express diagnostics of the disease, principal component analysis, multivariate probability density, multivariate data processing

INTRODUCTION

Diagnosis of diseases based on the analysis of respiration has been widely used in clinical practice [1, 2] using a variety of techniques and analytical instruments [3] in recent years. Such a diagnosis involves the detection of diseases at an early stage, which is promising. It has been established that the exhaled air contains important biomarkers for the early diagnosis of a number of diseases: cardiovascular, oncological, neurodegenerative, diabetes, respiratory, etc.

In works [4, 5], biomarkers of respiration are studied in concern to lung cancer, including cases of smokers. Volatiles in the samples were pre-concentrated by solid phase microextraction (SPME) followed by analysis in tandem of gas chromatograph and mass spectrometer (GC-MS). The obtained concentrations were analyzed in the space of the principal components (PCA) using the first three coordinates (PCA1, PCA2, PCA3).

The work [6] provides a detailed review of the results of research of diagnostic methods for quantitative and qualitative analysis of exhaled air (EA) composition, considers the requirements for analytical methods and instruments for studying multicomponent gas mixtures. It is noted that the analysis based on EA has a number of advantages over the diagnos-

tics based on laboratory methods — it is relatively cheap, takes little time and allows the detection of searched substances in minimal concentrations. It does not imply invasive interventions and can be performed with any frequency, providing an opportunity for a thorough study of the dynamics of physiological processes. A number of markers of a wide variety of diseases are given: for respiratory organs, gastrointestinal tract, oncology, central nervous system, etc. According to studies conducted in the USA, early diagnosis of diseases, in particular of the respiratory system, has reduced mortality by 20 %. Requirements are defined for means of detection of marker molecules taking into account the developed methods and detection limits: the ability to simultaneously measure and identify several biomarkers; high sensitivity and high accuracy of concentration determination (at the level of 10 ppb–10 ppm); selectivity (registration and identification of substances against the background of high concentrations of nitrogen, oxygen, water and carbon dioxide); speed (time of analysis with averaging over several expirations is about 5–10 s); ease of use in a clinical setting; the price of one test and the cost of the device. Gas chromatography (GC), mass spectrometry (MS), mass spectrometry with gas chromatographic separation (MSGC) and IR spectroscopy are considered as in-

struments for the analysis of EA.

The authors of [7] note the need to develop a convenient, easy-to-use, non-invasive method for screening chronic kidney disease and the quality of dialysis in real time. For this, it is proposed to measure the concentration of ammonia in exhaled air, however, the use of mass spectrometry, ion mobility spectrometry, optical methods of laser absorption, and other equipment is expensive, difficult and inconvenient for clinical use. It is proposed to use the recently developed gas sensors — a semiconductor sensor chip based on a polymer with nanopores.

To compare the results of the examination of several groups of patients, various statistical approaches were used: paired Student's t-test, Wilcoxon's test, Kruskal-Wallis test, Dan's test, Friedman's test. Pearson's correlation coefficient and simple linear regression were used to assess the correlation between expired ammonia concentrations and other diagnostic measurements.

In work [8], the possibility of diagnosing epilepsy by comparing the composition of exhaled air of sick and healthy patients using the "electronic nose" is investigated. The device data was processed by a neural network that was trained with data from pre-selected patients with epilepsy and its absence. The study of the effectiveness of EA diagnostics of epilepsy showed its undoubted advantages in comparison with traditional tests using electroencephalograms.

The work [9] analyzes various methods and devices for collecting samples in the field followed by analysis in a chromatograph with a mass spectrometer as a detector, equipped with an attachment for thermal desorption (TD-GC-MS analysis). Statistical processing of the results of measurement experiments, in particular the concentration of isoprene, was carried out using the Student's t-test and mixed linear modeling.

In work [10], it was shown that the method of analyzing exhaled air using a gas chromatograph — ion mobility spectrometer (GC-IMS) is suitable for early non-invasive diagnosis of neurodegenerative diseases, in particular, Alzheimer's disease, which is characteristic of the elderly. The complex is a high-performance diagnostic tool for real-time diagnostics in a clinical setting.

There is an urgent need to develop effective diagnostic tools, especially those that can reliably detect diseases in the early stages (before complex laboratory studies) using non-invasive approaches. A key issue in EA analysis is proper statistical analysis and interpretation of large and heterogeneous datasets derived from exploration studies. In this work, we used a quadrupole mass spectrometer and a method for processing the results of mass spectrometric analysis of EA, based on determining the likelihood of a disease by determining the assignment of the mass spectrum of the tested patient to the mass spectra of

the corresponding control group.

QUADRUPOLE MASS SPECTROMETER

Specifications

A quadrupole mass spectrometer used in this work for analyzing exhaled gases is described in [11]. With its small dimensions and weight, it provides the detection of EA components in the mass range — from 1 to 200 u, resolution in mass numbers — 0.5, registration rate — up to 1 mass spectrum in 10 s with sensitivity for impurity concentrations — from 0.1 ppm it is easy to use and provides high speed analysis from air sampling to result. Time for obtaining the analysis results is no more than a minute.

Brief description of device operation

The analyzed gas under atmospheric pressure is fed into electron impact ionization chamber through a capillary inlet. The resulting ions are introduced into a quadrupole mass analyzer. The mass spectrometric signals obtained during the registration process are processed using specialized software and compared with the spectra in the library of standard mass spectra, then the individual components of the spectrum are identified and their concentration is determined. The heated capillary system for introducing a sample into the mass spectrometer allows analysis at a distance of up to 5 m from the instrument. The sample flow rate is about 5 $\mu\text{L/s}$. The vacuum system uses a turbomolecular and a diaphragm pump.

An important advantage of the device is its small dimensions (see [11]), which allows its rapid movement from one medical institution to another for screening examinations of large groups of the population.

DATA PROCESSING

During a mass examination of the state of health of the population, as a rule, there is no preliminary information about the nature of a patient's illness. When diagnosing the state of health of patients by analyzing the exhaled air, it is advisable to register the entire spectrum of exhaled gases (markers of diseases) and then compare it with the data of one of the control groups with a known disease. The EA data of each patient in this group can be represented as a point in a multidimensional space. The combined data of all patients in a group forms a "cloud" of initial parameters. The closer the point defining the spectrum of the tested patient is to the center (centroid) of this "cloud", the more likely it is that the patient either suffers from a disease of the control group, or is healthy if it is a control group of healthy patients.

Let $x_{i,j}$ be the j -th parameter (in this case, the intensity of the spectral component) of the i -th patient in the control group, $i=1,2,\dots,I$ and $j=1,2,\dots,J$, and a set of values for patients I for registered EA components J form (I,J) the training matrix \mathbf{X} , the columns of which will be denoted as \mathbf{X}_j : $\mathbf{X}=[\mathbf{X}_1,\dots,\mathbf{X}_j,\dots,\mathbf{X}_J]$ [12].

The vector \mathbf{X}_j is a random variable. For simplicity, we will assume that the components of this vector are normally distributed with average value \bar{X}_j and variance σ_j^2 . Then the elements of the matrix \mathbf{X} form in the J -dimensional space a "cloud" (let us denote it as \mathbf{G}) of $I \times J$ points with a centroid at a point $\bar{\mathbf{X}}$ which coordinates are $\bar{\mathbf{X}}=[\bar{X}_1,\bar{X}_2,\dots,\bar{X}_J]$. The tested patient suffers from the same disease (or is healthy) as the control group, if the measured vector of the parameters of the EA is $\mathbf{X}_d=[x_{d,1},x_{d,2},\dots,x_{d,J}]$, where $x_{i,j}$ is the parameter value, that refers to the J -dimensional space \mathbf{G} , i.e. $\mathbf{X}_d \in \mathbf{G}$. This condition is satisfied if the probability P of the deviation of the point \mathbf{X}_d from the centroid $\bar{\mathbf{X}}$ does not exceed a certain threshold α . To determine this probability, we construct a multivariate probability distribution of the event referring to the space \mathbf{G} , assuming that this distribution obeys the normal law [13]:

$$P(\mathbf{X}_d) = W \cdot \exp\left\{-\frac{1}{2}(\mathbf{X}_d - \bar{\mathbf{X}})' \mathbf{K}_X^{-1}(\mathbf{X}_d - \bar{\mathbf{X}})\right\}, \quad (1)$$

$$\mathbf{X}_d \in \mathbf{G},$$

where \mathbf{K}_X is the covariance matrix $\mathbf{K}_X = E\left[(\mathbf{X} - \bar{\mathbf{X}}) \cdot (\mathbf{X} - \bar{\mathbf{X}})'\right]$, E — the expected value symbol, $(\cdot)'$ — the symbol of the matrix transposition, W — the normalizing factor. Obviously, the probability P of this event is equal to one if for the next patient it turns out that $\mathbf{X}_d \equiv \bar{\mathbf{X}}$. This condition is satisfied taking $W=1$ in formula (1). Then the condition of the patient's referring to the control group has the form: $P(\mathbf{X}_d) < \alpha$, where the value α is selected by the method of expert evaluation.

However, the direct use of formula (1) to calculate the value of P is associated with significant computational difficulties caused by the presence of a large number of parameters J and correlations between the columns of the matrix \mathbf{X} . To compress data, reduce the dimensionality of space \mathbf{G} , orthogonal data transforma-

tion into the space of principal components is used — the Principal Component Analysis (PCA) [14].

To go to the principal component (PC) space, a new matrix is formed, consisting of all rows of the matrix \mathbf{X} and a row \mathbf{X}_d . Let's designate this matrix as $\mathbf{X}\mathbf{I}$. Then in the new coordinate system:

$$\mathbf{X}\mathbf{I} = \mathbf{T} \cdot \mathbf{P}' + \mathbf{e} = \sum_{j=1}^A \mathbf{t}_j \mathbf{p}'_j + \mathbf{e},$$

where \mathbf{p}_j are the eigenfunctions of the covariance matrix \mathbf{K}_X . The matrix \mathbf{T} is called the *T-scores* matrix $\mathbf{T}=[\mathbf{T}_1,\mathbf{T}_2,\dots,\mathbf{T}_A]$, its dimensionality is $(I \times A)$; the matrix \mathbf{P} is called the matrix of *loads*, its dimension is $(I \times A)$; \mathbf{e} is a matrix of *residuals* (noises) of dimension $(I \times J)$; column vectors \mathbf{T}_j ($j=(1,2,\dots,A)$) are called *principal components* (PC), A is the number of principal components. The value A is significantly less than the number of variables J . This means that the main information about the patient's condition is concentrated in the first few PCs. The last row of this matrix, the vector \mathbf{T}_d is the coordinates of the parameters of the tested patient in the PC space: $\mathbf{T}_d=[t_{d,1},t_{d,2},\dots,t_{d,A}]$.

The average values of the \mathbf{T} matrix columns are zero, and the variance is a vector $\boldsymbol{\sigma}^2$ with elements $\sigma_j^2 = \lambda_j$, i.e. eigenvalues of the covariance matrix. The property of the PC expansion is such that the variance rapidly decreases already to the fourth PC, and the columns of the matrix are uncorrelated, i.e.

$$\mathbf{T}_m \mathbf{T}'_n = \begin{cases} 0 & \text{при } n \neq m, \\ \lambda_j & \text{при } n = m. \end{cases}$$

Taking this circumstance into account in the new coordinate system, formula (1) has the form:

$$P(\mathbf{T}_d) = \exp\left\{-\frac{1}{2} \mathbf{T}_d \boldsymbol{\sigma}^{-2} \mathbf{T}'_d\right\} = \exp\left\{-\frac{1}{2} \sum_{j=1}^A \frac{t_{d,j}^2}{\sigma_j^2}\right\}, \quad (2)$$

and the Euclidean distance from the patient with index d to the centroid of the training group is

$$D(d) = \sqrt{\sum_{j=1}^A t_{d,j}^2}. \quad (3)$$

DESCRIPTION OF THE ALGORITHM

Data processing consists of two stages: *training and diagnostics*.

At the **training stage** K of matrices is formed \mathbf{X}^k , $k = 1, 2, \dots, K$ with data on the intensity of the EA spectrum for each type of diseases. The number of rows I^k of each matrix must be several times greater than the number of columns J^k .

At the **diagnostic stage** the following sequence of operations is performed:

1. The parameters of the EA of the diagnosed patient are measured and a row vector $\mathbf{X}_d = \{x_{d,1}, x_{d,2}, \dots, x_{d,J}\}$ is formed.

2. A new matrix $\mathbf{X}I^k = [\mathbf{X}^k; \mathbf{X}_d]$ of dimensions $((I^k + 1) \times J^k)$ is constructed, the last row of which is a vector X_d .

3. This matrix is normalized, for example, to the maximum value of its elements.

4. Using the algorithm of the PC method, the score matrix \mathbf{T}^k and the vector of eigenvalues λ_j^k are calculated [15].

5. The number of the first columns A of the matrix

\mathbf{T}^k is determined from the condition $\sigma_A^k = \sqrt{\lambda_A^k} < \varepsilon$.

In this case, without loss of reliability, we can set $\varepsilon = 0.001-0.01$. The last row of the matrix

$\mathbf{T}_d^k = \{t_{d,1}^k, t_{d,2}^k, \dots, t_{d,A}^k\}$ is the principal components of the EA parameters of the diagnosed patient.

6. The probability $P(\mathbf{T}_d^k)$ is determined by the formula (2).

7. Go to item 2 ($k = k + 1$) to determine the probability of attributing tested patient to other group of diseases.

8. Analysis of the results of calculating the probability $P(\mathbf{T}_d^k)$, $k = 1, 2, \dots, K$ for all diseases K in order to determine the most likely.

ALGORITHM CHECK

Let us illustrate the above theory on one group of 43 patients diagnosed with the same type of oncolo-

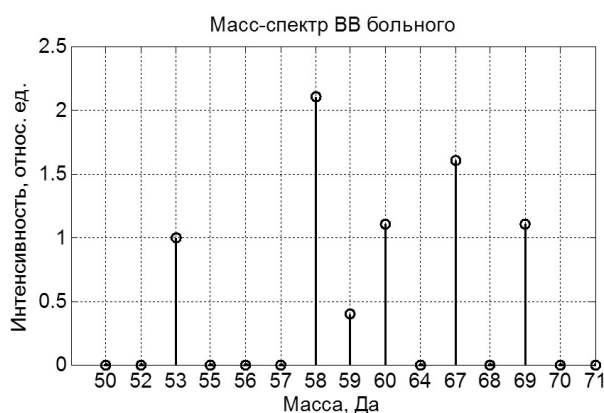


Fig. 1. Mass spectrum of EA of the examined patient. Weight 53 — isoprene, 58 — acetone, 60 — vinegar, 67 — isoprene (?), 69 — pentin; 59 — a still unrecognized biomarker identified in some oncological patients

Tab. 1. Measured probability of illness for patients of the first control group

Patient	Euclidean distance	Probability of disease
1	0.0116	0.9854
2	0.1076	0.1888
3	0.1499	0.2055
4	0.2044	0.1164
5	0.2554	0.2864
6	0.2892	0.3458
7	0.3247	0.1197

Tab. 2. Measured probability of illness for patients in the second control group

Patient	Euclidean distance	Probability of disease
1	2.8791	0.0159
4	4.5739	0.0000
3	4.7337	0.0000
2	5.5927	0.0000
5	5.6321	0.0000
8	5.7129	0.0000
9	5.7322	0.0000
6	5.8520	0.0000
10	5.8635	0.0000
7	5.9168	0.0000

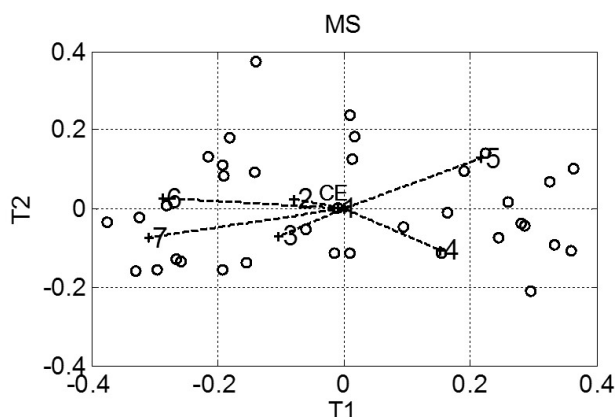


Fig. 2. Patients of the training (○) and the first control (+) groups of patients in the space of the first two PCs. CE — centroid position, 1–7 — patient numbers

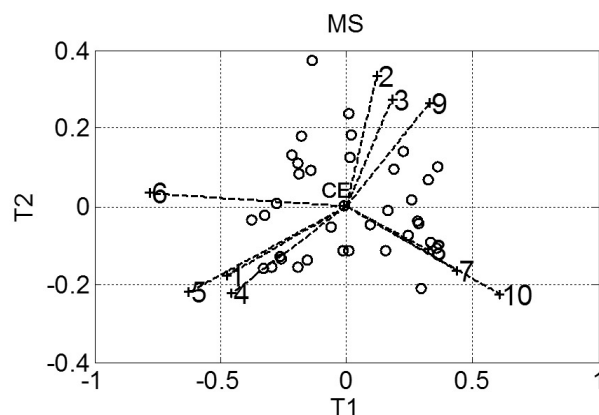


Fig. 3. Patients of the teaching (○) and second control (+) groups of patients in the space of the first two PCs. CE — centroid position, 1–10 — patient numbers

gy¹⁾. This group is divided into two: training (36 patients) and control (7 patients). In this case, the concentration of EA components of the first patient from the control group, for example, coincides with the centroid of the training group. Another control group of 10 patients with other oncological diseases was recruited to test the effectiveness of the algorithm either.²⁾ Patients of both control groups were assigned serial numbers, respectively, $d = (1, 2, 3, \dots, 7)$ and $d = (1, 2, 3, \dots, 10)$.

Fig. 1 shows an example of a mass spectrum of EA of one of the patients in the first (training) group. Out of each control group, we select data on EA in turn and perform pp. 1–6 of algorithm. It turns out that it is enough to assign $A = 6$, since $\sigma_7^k \cong 0$.

Using formulas (3) and (2), the Euclidean distance from the centroid of the training sample and the probability of diagnosing the oncological disease in the control patient are calculated. For each control group, these results are recorded in Tab. 1 and 2. In Fig. 2 in the space of the first two PCs (T1 and T2), the coordinates of the patients of the training (○) and the first control (+) group are shown. According to the data of Tab. 1 it can be seen that as the EA coordinates move away from the centroid for this control group, the probability of oncology decreases, but remains high. It should be noted that this probability

also depends on the location of the patient's coordinates in the six-dimensional (in this case) PC space: the larger the projection onto the first PCs, the greater this probability.

In Fig. 3 in the space of the first two PCs (T1 and T2), the coordinates of the patients of the training (○) and the second control group (+) are shown.

It turned out that all patients in this group (except for the first one) are outside the "cloud" of the training sample. Tab. 2 shows that the likelihood of a presence of oncological diseases for this group of patients is negligible.

CONCLUSION

The system for the diagnosis of diseases by EA, consisting of a quadrupole mass spectrometer and a data processing set with the possibility of training and determining the probability of the disease provides a non-invasive rapid analysis of the patient's condition. Diagnostics with the determination of the likelihood of a patient attributing to one of the training groups allows for a few tens of seconds to compile an overall picture of the patient's probable diseases and to prescribe a further examination in a qualified manner.

REFERENCES

1. Cao W., Duan Y. Breath analysis: potential for clinical diagnosis and exposure assessment. *Clinical Chemistry*, 2006, vol. 52, no. 5, pp. 800–811. DOI: 10.1373/clinchem.2005.063545
2. Pereira J., Porto-Figueira P., Cavaco C., Taunk K., Rapole S., Dhakne R., Nagarajaram H., Câmara J.S. Breath analysis as a potential and non-invasive frontier in disease diagno-

¹⁾ Data on EA of this group of patients were provided by the state budgetary healthcare institution "St. Petersburg Clinical Scientific and Practical Center for Specialized Types of Medical Care (Oncological).

²⁾ EA data for this group of patients are provided by the Leningrad Regional Clinical Oncology dispensary.

- sis: an overview. *Metabolites*, 2014, no. 5, pp. 3–55. DOI: 10.3390/metabo5010003
3. Kononov A., Korotetsky B., Jahatspanian I., Gubal A. et al. Online breath analysis using metal oxide semiconductor sensors (electronic nose) for diagnosis of lung cancer. *Journal of Breath Research*, 2020, vol. 14, no. 1. DOI: 10.1088/1752-7163/ab433d
 4. Kischkel S., Miekisch W., Sawacki A., Straker E.M., et al. Breath biomarkers for lung cancer detection and assessment of smoking related effects — confounding variables, influence of normalization and statistical algorithms. *Clinica Chimica Acta*, 2010, vol. 411, is. 21–22, pp. 1637–1644.
 5. Fuchs P., Loeseke Ch., Schubert J.K., Miekisch W. Breath gas aldehydes as biomarkers of lung cancer. *International Journal of Cancer*, 2010, vol. 126, no. 11, pp. 2663–2670.
 6. Vaks V.L., Domracheva E.G., Sobakinskaya E.A., Chernyaeva M.B. [Exhaled breath analysis: physical methods, instruments and medical diagnostics]. *Uspekhi Fizicheskikh Nauk* [Advances in Physical Sciences], 2014, vol. 184, no. 7, pp. 739–758. DOI: 10.3367/UFNr.0184.201407d.0739 (In Russ.).
 7. Chen1 Ch.-Ch., Hsieh Ju-Ch., Chao Ch.-H., Yang W.-Sh., et al. Correlation between breath ammonia and blood urea nitrogen levels in chronic kidney disease and dialysis patients. *Journal of Breath Research*, 2020, vol. 14, no. 3. DOI: 10.1088/1752-7163/ab728b
 8. van Dartel D., Schelhaas H.J., Colon A.J., Kho K.H., de Vos C.C. Breath analysis in detecting epilepsy. *Journal of Breath Research*, 2020, vol. 14, no. 3. DOI: 10.1088/1752-7163/ab6f14
 9. Harshman S.W., Pitsch R.L., Davidson Ch.N., Scott A.M., et al. Characterization of standardized breath sampling for off-line field use. *Journal of Breath Research*, 2019, vol. 14, no. 1. DOI: 10.1088/1752-7163/ab55c5
 10. Tiele1 A., Wicaksono A., Daulton E., Ifeachor E., Eyre V., et al. Breath-based non-invasive diagnosis of Alzheimer’s disease: A pilot study. *Journal of Breath Research*, 2020, vol. 14, no. 2. DOI: 10.1088/1752-7163/ab6016
 11. Manoilov V.V., Kuzmin A.G., Zarutskiy I.V., Titov Yu.A., Samsonova N.S. [Methods of processing and investigation of the possibilities of classification of mass spectra of exhaled gases]. *Nauchnoe Priborostroenie* [Scientific Instrumentation], vol. 29, no. 1, pp. 106–110. DOI: 10.18358/np-29-1-i106110 (In Russ.). URL: <http://iairas.ru/mag/2019/abst1.php#abst16>
 12. Bolshakov A.A., Karimov R.N. ed. *Metody obrabotki mnogomernykh dannyh i vremennykh ryadov* [Methods of processing multidimensional data and time series]. Moscow, Goryachaya liniya-Telekom, 2007. 522 p. (In Russ.).
 13. Shiryaev A.N. *Veroyatnost. Tom 1* [Probability. Vol. 1]. Moscow, MCNMO Publ., 2007. 552 p. (In Russ.).
 14. Ajvazyan S.A., Buhstaber V.M., Enyukov I.S., Meshalkin L.D. *Prikladnaya statistika. Klassifikaciya i snizhenie razmernosti* [Application statistics. Classification and dimensioning]. Moscow, Finansy i statistika, 1989. 607 p. (In Russ.).
 15. Esbensen K.H. *Multivariate Data Analysis*. Norway, CAMO Software AS, 480 p. (Russ. ed.: Esbensen K. *Analiz mnogomernykh dannyh*. Translate O.E. Rodionova. Chernogolovka, IPCP OF RAS Publ., 2005, 161 p.). URL: <https://www.amazon.com/Multivariate-Data-Analysis-introduction-Analytical/dp/826911040X> (In Russ.).

Contacts: *Manoilov Vladimir Vladimirovich*,
manoilov_vv@mail.ru

Article received by the editorial office on 15.10.2020