

Нейросетевое обобщение классических статистических критериев для обработки малых выборок биометрических данных

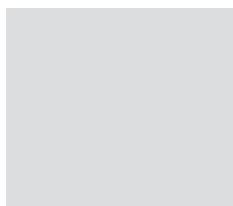
Александр И. Иванов, АО «Пензенский научно-исследовательский электротехнический институт», Российская Федерация, Пенза

Евгений Н. Куприянов, ФБГОУ ВПО «Пензенский государственный университет», Российская Федерация, Пенза

Сергей В. Туреев, НИИ систем связи и управления, Российская Федерация, Москва



Александр И. Иванов



Евгений Н. Куприянов



Сергей В. Туреев

Резюме. Целью работы является повышение мощности статистических критериев за счет их совместного применения для снижения требований к объему тестовой выборки. **Методы.** Классические статистические критерии: хи-квадрат, Крамера фон-Мизеса и Шапиро-Уилка предложено объединять через использование эквивалентных им искусственных нейронов. Каждый нейрон выполняет сравнение входных статистик с заранее вычисленным порогом и имеет два выходных состояния. Это позволяет получать три разряда бинарного выходного кода сети из трех искусственных нейронов. **Результаты.** Показано, что каждый из этих критериев на малых выборках биометрических данных дает большие значения ошибок первого и второго рода при проверке гипотезы нормальности. Нейросетевое объединение трех рассматриваемых критериев позволяет существенно снизить вероятности ошибок первого и второго рода. Приведены результаты парных нейросетевых обобщений, а также неросетевое обобщение тройки рассматриваемых статистических критериев. **Выводы.** Дается прогноз ожидаемых вероятностей ошибок первого и второго рода для нейросетевых обобщений 10 и 30 классических статистических критериев для малых выборок, содержащих 21 опыт. Важным элементом технологии прогнозирования является симметризация задачи, когда вероятности ошибок первого и второго рода делаются одинаковыми и усредняются. Усредняются также модули коэффициентов парной корреляции выходных состояний сумматоров искусственных нейронов. Только в этом случае связь числа обобщаемых критериев с ожидаемыми вероятностями ошибок первого и второго рода становится линейной в логарифмических координатах.

Ключевые слова: статистические критерии: хи-квадрат, Крамера фон-Мизеса, Шапиро-Уилка; искусственные нейронные сети, малые выборки, проверка гипотезы нормального закона распределения данных.

Формат цитирования: Иванов А.И., Куприянов Е.Н., Туреев С.В. Нейросетевое обобщение классических статистических критериев для обработки малых выборок биометрических данных // Надежность. 2019. №2. С. 22-27. DOI: 10.21683/1729-2646-2019-19-2-22-27

Проблема контроля закона распределения данных малых выборок

Проблемы обеспечения надежности уникальных ответственных систем [1, 2] многогранны и могут быть решены только комплексом организационно-технических мероприятий. Особо ярко эти проблемы проявляются в нейросетевой биометрии. Каждый из нас имеет уникальный биометрический образ, который необходимо однозначно преобразовать криптографический ключ или длинный пароль доступа из случайных символов. Однозначность преобразования обеспечивается обучением нейронной сети, при этом обучающая выборка имеет близкий к нормальному многомерный закон распределения данных. Проблема состоит в том, что обучающие выборки малы. В частности стандартный алгоритм [3] обучения хорошо справляется со своей задачей на выборке в 20 примеров, если эта выборка корректно получена и не имеет выбросов (грубых ошибок).

В случае, когда выборка биометрических данных велика (200 опытов и более) проверить гипотезу нормальности их распределения не сложно. Можно воспользоваться хи-квадрат критерием или иным другим статистическим критерием [4]. Одной из проблем биометрии [5] является то, что ее пользователи не желают предъявлять автомату обучения искусственной нейронной сети [3] 200 и более примеров своего биометрического образа. Пользователи комфортно себя чувствуют, предъявляя обучающую выборку объемом от 10 до 20 примеров своего уникального биометрического образа, например, рукописного пароля или голосовой парольной фразы. Требования предъявить больше примеров для обучения, воспринимаются пользователями негативно.

Аналогичная ситуация возникает в ботанике, биологии, медицине. Селекционер ботаник или биолог в короткие сроки не имеет возможности получить выборку в 200 особей животных (экземпляров растений), обладающих нужными редкими свойствами. Достаточную для корректных статистических оценок выборку удастся

получить через достаточно длительный интервал времени, фактически выделив и закрепив желаемые редкие свойства в нескольких поколениях.

В медицине наблюдается аналогичная ситуация: для проверки статистических гипотез требуются большие выборки. Тематика статистической обработки малых выборок популярна, однако известные рекомендации [6, 7] не дают кардинальных улучшений ситуации. Как правило, добиться улучшения удастся, если применять несколько статистических критериев [8].

Можно попытаться усилить известные статистические критерии [9], однако этот путь не дает значительных улучшений. Как правило, новые статистические критерии или модификация старых по отдельности плохо работают.

Основной идеей данной работы является нейросетевое объединение классических статистических критериев [4, 10, 11]. Успехи, достигнутые российской нейросетевой биометрией, весьма и весьма значительны. Силами регуляторов отечественного рынка информационной безопасности создан пакет российских национальных стандартов ГОСТ Р 52633.хх, регламентирующих ряд очень жестких требований к нейросетевой биометрии. В рамках данной статьи мы фактически предпримем попытку переложить хорошо отработанные математические приемы нейросетевой биометрии в новые предметные области. При этом мы попытаемся показать, что рекордно жесткие требования регуляторов отечественного рынка информационной безопасности к вероятностям ошибок первого и второго рода будут выполнимы и в других предметных областях, если придерживаться основных рекомендаций пакета стандартов нейросетевой биометрии ГОСТ Р 52633.хх.

Синтез и настройка хи-квадрат нейрона с 5 входами

При проверке гипотезы нормальности на практике наиболее часто используется хи-квадрат критерий Пирсона. Для малой выборки объемом в 21 опыт формула для вычисления значения хи-квадрат критерия имеет следующий вид:

$$\chi^2 = 21 \cdot \sum_{i=1}^5 \frac{\left(\frac{n_i}{21} - \Delta \tilde{P}_i\right)^2}{\Delta \tilde{P}_i}, \quad (1),$$

где n_i – число опытов, попавших в i -тый интервал гистограммы, $\Delta \tilde{P}_i$ – ожидаемая вероятность попадания опытов в i -тый интервал гистограммы при нормальном законе распределения данных проверяемой выборки.

Заметим, что в соответствии с отечественными стандартизованными рекомендациями [10] среднее число опытов, попавших в каждый из интервалов гистограмм должно быть близко к 5. Именно по этой причине в формуле (1) используется суммирование по 5 интервалам гистограммы для малой выборки в 21 опыт.

При создании своей формулы в 1900 году Пирсон не мог воспользоваться вычислительной техникой. По этой причине он вынужден был искать асимптотические соотношения для бесконечно больших выборок. Сегодня ситуация изменилась, любой студент способен написать программу, воспроизводящую миллионы выборок объемом в 21 опыт. На рисунке 1 приведено распределение плотности вероятности значений хи-квадрат критерия для выборок с нормальным и равномерным законом распределения значений.

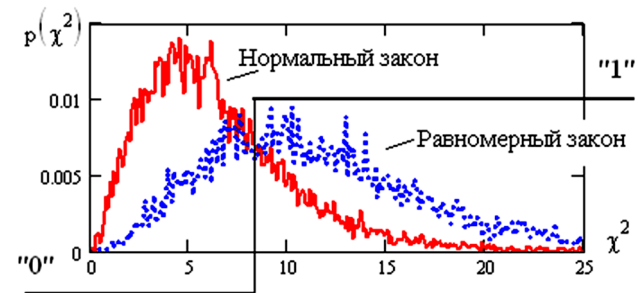


Рисунок 1 – Распределения значений хи-квадрат критерия для выборок объемом в 21 опыт с равномерным и нормальным законами распределения значений

Следует отметить, что искусственные нейроны настраиваются таким образом, чтобы эффективно разделять на два класса «нормальные» и «равномерные» входные данные [12]. На рисунке 1 пороговый элемент хи-квадрат нейрона делит континуум выходных на две области: «0» – «нормальные» данные и «1» – «равномерные» данные. Настройка выходного квантователя хи-квадрат нейрона выполняется исходя из условия равновероятного значения ошибок первого и второго рода $P_1 = P_2 = P_{EE} = 0,292$.

Для того, чтобы получить пять входных параметров хи-квадрат нейрона, необходимо выполнить сортировку данных проверяемой выборки по их значениям:

$$x = \text{sort}(x). \quad (2)$$

Кроме того, необходимо вычислить ширину интервалов гистограммы:

$$\Delta x = \frac{x_{20} - x_0}{5}. \quad (3)$$

Далее выполняют вычисление положения краев интервалов:

$$X_i = x_0 + \Delta x \cdot i \text{ при } i = 0, 1, \dots, 5. \quad (4)$$

Только после этого удастся подсчитать число попаданий в каждый из интервалов гистограммы и сформировать вектор входных параметров $\{n_1, n_2, \dots, n_5\}$ для нейрона (1). Итоговый результат подвергается квантованию:

$$\begin{cases} z(\chi^2) \leftarrow "0" & \text{если } \chi^2 \leq 7,72; \\ z(\chi^2) \leftarrow "1" & \text{если } \chi^2 > 7,72. \end{cases} \quad (5)$$

В итоге мы имеем полное формальное описание реализации хи-квадрат нейрона для выборки, состоящей из 21 опыта.

Синтез и настройка нейрона Шапиро-Уилка с 10 входами

Очевидно, что к той же выборке из 21 реализации может быть применен критерий Шапиро-Уилка [4, 11]. Его значение вычисляется по следующей формуле:

$$v^2 = \frac{1}{(\sigma(x))^2} \cdot \left\{ \sum_{i=0}^9 a_i \cdot (x_{20-i} - x_i) \right\}^2, \quad (6)$$

где x_i – упорядоченные значения, проверяемой выборки, $\sigma(x)$ – стандартное отклонение, a_i – табличные значения коэффициентов Шапиро-Уилка.

Распределения значений этого критерия для равномерного и нормального законов приведены на рисунке 2.

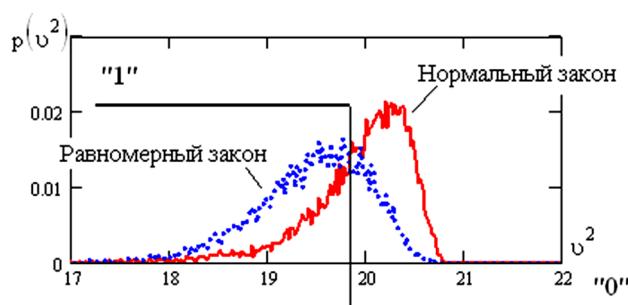


Рисунок 2 – Распределения значений критерия Шапиро-Уилка для выборки объемом в 21 опыт с равномерным и нормальным законами распределения значений

Если рассматривать функционал (6) как некоторый искусственный нейрон, то его выходами будут 10 разностей данных, исследуемой выборки, а выходной квантователь будет описываться следующими соотношениями:

$$\begin{cases} z(v^2) \leftarrow "0" & \text{если } v^2 \geq 19,8; \\ z(v^2) \leftarrow "1" & \text{если } v^2 < 19,8. \end{cases} \quad (7)$$

Такая настройка порога срабатывания квантователя обеспечивает одинаковые значения вероятности ошибок первого и второго рода $P_1 = P_2 = P_{EE} = 0,303$.

Синтез и настройка нейрона Крамера фон-Мизеса с 20 входами

Если сравнивать хи-квадрат нейрон (1) и нейрон Шапиро-Уилка (6) мы наблюдаем рост размерности их входной размерности (числа входов их сумматоров). Еще большей входной размерностью обладает нейрон Крамера фон-Мизеса:

$$\omega^2 = \sum_{i=0}^{19} \left(\frac{i+1}{21} - \tilde{P}(x_i) \right)^2 \cdot \frac{x_{i+1} - x_i}{x_{20} - x_0}. \quad (8)$$

Распределения значений на выходе сумматора нейрона Крамера фон-Мизеса приведены на рисунке 3.

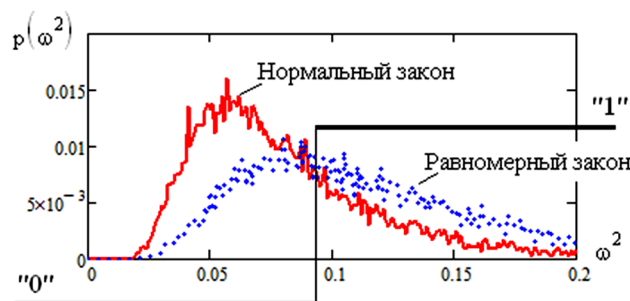


Рисунок 3 – Распределения значений критерия Крамера фон-Мизеса для выборок объемом в 21 опыт с равномерным и нормальным законами распределения значений

Настроенный выходной квантователь сумматора нейрона описывается следующим образом:

$$\begin{cases} z(\omega^2) \leftarrow "0" & \text{если } \omega^2 \leq 0,087; \\ z(\omega^2) \leftarrow "1" & \text{если } \omega^2 > 0,087. \end{cases} \quad (9)$$

Такая настройка порога срабатывания квантователя обеспечивает одинаковые значения вероятности ошибок первого и второго рода $P_1 = P_2 = P_{EE} = 0,342$.

Совместное использование трех статистических критериев

Описанные выше статистические критерии линейно независимы (имеют модули коэффициентов корреляции менее единицы):

$$\begin{cases} \text{corr}(\chi^2, v^2) \approx +0,559; \\ \text{corr}(\chi^2, \omega^2) \approx -0,708; \\ \text{corr}(\omega^2, v^2) \approx -0,667. \end{cases} \quad (10)$$

Отсутствие полной линейной зависимости (10) выходных состояний трех критериев позволяет объединить их для совместного использования. В этом случае выходной код трех нейронов «000» будет соответствовать трехкратному подтверждению гипотезы нормальности данных, исследуемой выборки. Инверсное состояние этого кода «111» будет соответствовать трехкратному подтверждению гипотезы равномерного закона распределения данных малой выборки.

По аналогии с практикой применения нейросетевых преобразователей биометрия-код будем принимать решение о принятии одной из двух гипотез по большинству состояний «0» или «1» в выходном коде сети трех нейронов. В этой ситуации каждому из четырех кодовых состояний «нормальное» распределение будет соответствовать своя вероятность ошибок, эти данные сведены в таблицу 1.

Таблица 1. Вероятности появления ошибок для кодовых состояний «нормальное» распределение

Код	«000»	«001»	«010»	«100»
P_1	0,0404	0,0423	0,0441	0,0621

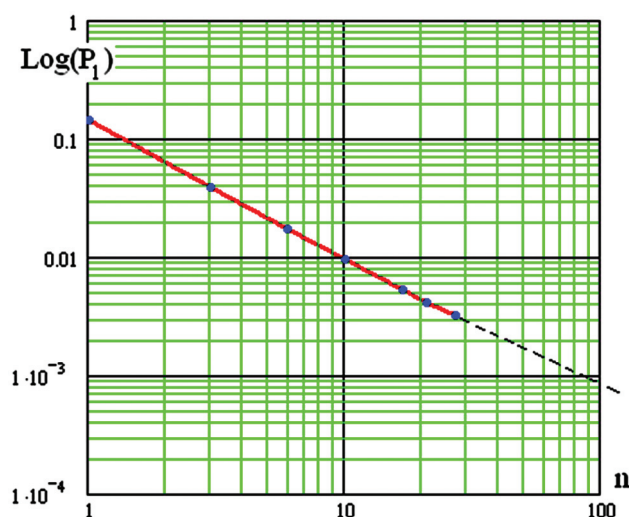


Рисунок 5 – Линия снижения вероятности ошибок первого рода из-за совместного применения нескольких статистических критериев с коэффициентами корреляции 0,645

Линия построена по 7 группам, состоящим из 1, 3, 6, 10, 16, 21, 27 нейронов. При проведении численного эксперимента использовалась выборка из 10 000 000 реализаций, время вычислений составляет примерно 9 минут на обычной вычислительной машине. Следует отметить, что применяя эту же вычислительную машину трудно провести численный эксперимент для группы из 100 нейронов, так как придется ждать несколько месяцев. Сократить время удастся путем экстраполяции (пунктирная линия на рисунке 5).

В конечном итоге прогнозируемое значение вероятности ошибок для нейросетевого обобщения 10 критериев должно составить $P_1 = 0,01$, а при обобщении 100 критериев вероятность ошибок должна снизиться до величины 0,0009. Столь существенное снижение вероятности ошибок является хорошим стимулом для организации работ по синтезу новых статистических критериев [13-17].

Заключение

Пирсон, создавший в 1900 году хи-квадрат критерий, по сути дела, начал революцию в статистической обработке. Путь развития, обнаруженный Пирсоном, оказался очень плодотворным и за прошедшие 119 лет, его последователями были созданы десятки разных статистических критериев.

Нейронные сети начали активно изучаться с середины XX века, однако только в начале XXI века эта технология обработки была доведена до промышленного применения и стандартизована [3].

Основным утверждением данной статьи является возможность объединения двух, казалось бы, разных ветвей математики. Для объединения вполне достаточно использовать стандартизованные в России технологии нейросетевой обработки биометрических данных, применив их к трем или более классическим

статистическим критериям. Для рассмотренной тройки статистических критериев этот подход дает снижение вероятности ошибок более чем в 7 раз. При этом становится очевидным тезис о целесообразности расширения номенклатуры существующих статистических критериев. Чем больше размер группы, обобщаемых нейронами статистических критериев, тем лучше должен быть конечный результат.

В этом контексте кардинально меняется подход к синтезу новых статистических критериев. После Пирсона математики старались найти новый критерий, мощность которого выше чем у его предшественников. Огромное число критериев, которые были исследованы, но имели относительно низкую мощность, не публиковалось. При нейросетевом объединении множества статистических критериев мощность каждого из них перестает играть основную роль. Крайне важным оказывается еще и то, каковы корреляционные связи добавленного критерия с группой других критериев. В нашем случае два объединяемых критерия имеют примерно одинаковую мощность, однако в этой группе присутствует особый критерий Шапиро-Уилка, который имеет низкую коррелированность с основными критериями хи-квадрат и Крамера-фон Мизеса.

Как следствие, необходимо повторить работы по исследованию возможного многообразия статистических критериев, принимая во внимание не только их относительную мощность, но и значения их коэффициентов корреляции в группах с другими, наиболее востребованными статистическими критериями. Новые статистические критерии с относительно низкой мощностью разделения гипотез ранее отбраковывались и не публиковались, теперь ситуация коренным образом изменилась. Куда важнее становится то, как новый критерий дополняет уже исследованные статистические критерии. Скорее всего, в ближайшее время потребуются создавать некоторую таблицу уровня родственности (коррелированности) уже известных и перспективных статистических критериев. Объединять в группы и создавать для них нейросетевые обобщения выгоднее всего линейно независимые (слабо коррелированные) статистические критерии.

Библиографический список

1. Похабов Ю.П. Проблемы надежности и пути их решения при создании уникальных высокоответственных систем [Текст] / Ю.П. Похабов // Надежность. – Том 19. – № 1. – С. 10-17.
2. Похабов Ю.П. Обеспечение надежности уникальных высокоответственных систем [Текст] / Ю.П. Похабов // Надежность. – Том 17. – № 3. – С. 17-23.
3. ГОСТ Р 52633.5-2011. Защита информации. Техника защиты информации. Автоматическое обучение нейросетевых преобразователей биометрия-код доступа [Текст]. – Введ. 2012-04-01. – М.: Стандартинформ, 2012. – IV, 15 с.

4. **Кобзарь А.И.** Прикладная математическая статистика: для инженеров и научных работников [Текст] / А.И. Кобзарь. – М.: ФИЗМАТЛИТ, 2006 г. – 816 с.
5. **Язов Ю.К.** и др. Нейросетевая защита персональных биометрических данных [Текст] / Ю.К. Язов, В.И. Волчихин, А.И. Иванов, В.А. Фунтиков, И.Г. Назаров / Под ред. Ю.К. Язова. – М.: Радиотехника, 2012 г. – 157 с.
6. **Сухорученков Б.И.** Анализ малой выборки. Прикладные статистические методы [Текст] / Б.И. Сухорученков. – М.: Вузовская книга, 2010. – 384 с: ил.
7. **Дерффель К.** Статистика в аналитической химии [Текст] / К. Дерффель. – М.: Мир, 1994. – 258 с.
8. **Даев Ж.А., Нурушев Е.Т.** Применение статистических критериев для улучшения эффективности методов оценки рисков. // Надежность. – Том 18. – № 2. – С. 42-45.
9. **Ахметов Б.Б., Иванов А.И.** Оценка качества малой выборки биометрических данных с использованием более экономичной формы хи-квадрат критерия. // Надежность. – 2016. – № 2(57). – С. 43-48.
10. **Р 50.1.037-2002.** Рекомендации по стандартизации. Прикладная статистика. Правила проверки согласия опытного распределения с теоретическим. Часть I. Критерии типа χ^2 [Текст]. – М.: Госстандарт России, 2001. – 140 с.
11. **Р 50.1.037-2002.** Прикладная статистика. Правила проверки согласия опытного распределения с теоретическим. Часть II. Непараметрические критерии [Текст]. – М.: Госстандарт России, 2002. – 123 с.
12. **Хайкин Саймон.** Нейронные сети: полный курс [Текст] / С. Хайкин. – М.: «Вильямс», 2006. – 1104 с.
13. **Серикова Н.И.** Оценка правдоподобия гипотезы о нормальном распределении по критерию Джини для числа степеней свободы, кратного числу опытов [Текст] / Н.И. Серикова, А.И. Иванов, Ю.И. Серикова. // Вопросы радиоэлектроники. – 2015. – № 1(1). – С. 85-94.
14. **Перфилов К.А.** Критерий среднего геометрического, используемый для проверки достоверности статистических гипотез распределения биометрических данных [Текст] / К.А. Перфилов. / Труды научно-технической конференции кластера пензенских предприятий, обеспечивающих безопасность информационных технологий. – Пенза, 2014. – Том 9. – С. 92-93. – URL: <http://www.pniei.penza.ru/RV-conf/T9/C92>.
15. **Иванов А.И.** Оценка соотношения мощностей семейства статистических критериев «среднего геометрического» на малых выборках биометрических данных [Текст] / А.И. Иванов, К.А. Перфилов / XI Всероссийская научно-практическая конференция «Современные охраняемые технологии и средства обеспечения комплексной безопасности объектов». Пенза-Заречный. 20 октября 2016 г. – 2016. – С. 223-229.
16. **Иванов А.И.** Многомерный статистический анализ качества биометрических данных на предельно малых выборках с использованием критериев среднего геометрического, вычисленного для анализируемых функций вероятности [Текст] / А.И. Иванов, К.А. Перфилов, Е.А. Малыгина // Измерение. Мониторинг. Управление. Контроль. – 2016. – № 2(16). – С. 64-72.
17. **Иванов А.И.** Оценка качества малых выборок биометрических данных с использованием дифференциального варианта статистического критерия среднего геометрического [Текст] / А.И. Иванов, К.А. Перфилов, Е.А. Малыгина // Вестник СИБГАУ. – 2016. – №4(17). – С. 864-871.
18. **Малыгин А.Ю.** Быстрые алгоритмы тестирования нейросетевых механизмов биометрико-криптографической защиты информации [Текст] / А.Ю. Малыгин, В.И. Волчихин, А.И. Иванов, В.А. Фунтиков. – Пенза, Издательство Пензенского государственного университета, 2006. – 161 с.
19. **Ахметов Б.С.** Алгоритмы тестирования биометрико-нейросетевых механизмов защиты информации [Текст] / Б.С. Ахметов, В.И. Волчихин, А.И. Иванов, А.Ю. Малыгин. – Алматы, КазНТУ им. Сатпаева, 2013. – 152 с. – URL: <http://portal.kazntu.kz/files/publicate/2014-01-04-11940.pdf>

Сведения об авторах

Александр И. Иванов, доктор технических наук, доцент, ведущий научный сотрудник лаборатории биометрических и нейросетевых технологий АО «Пензенский научно-исследовательский электротехнический институт», Российская Федерация, Пенза, e-mail: ivan@pniei.penza.ru

Евгений Н. Куприянов, аспирант кафедры «Технические средства информационной безопасности» ФГБОУ ВО «Пензенский государственный университет», Российская Федерация, Пенза, e-mail: ibst@pnzgu.ru

Сергей В. Туреев, начальник научно-технического центра, НИИ систем связи и управления, Российская Федерация, Москва, e-mail: niissu@niissu.ru

Поступила: 22.01.2019