

НЕЙРОСЕТЕВЫЕ ТЕХНОЛОГИИ ОБРАБОТКИ ДАННЫХ: АНАЛИЗ И ПРИМЕНЕНИЕ

Головко В. А.

Кафедра Интеллектуальных информационных технологий, Брестский государственный технический университет

Брест, Республика Беларусь

E-mail: gva@bstu.by

В данной статье рассматривается применение нейронных сетей для обработки хаотических процессов с целью диагностики эпилепсии, а также для обнаружения и распознавания атак на компьютерные сети.

ВВЕДЕНИЕ

В настоящее время все больше возрастает тенденция проектирования систем искусственного интеллекта на основе биологически инспирированных подходов, таких как нейронные и иммунные сети, эволюционное программирование и т.д. [1]. Способность нейронных сетей к обучению, обобщению результатов и решению трудно формализуемых задач создает предпосылки для проектирования на их базе различного рода интеллектуальных систем. Такие нейросетевые системы применяются для решения задач управления, аппроксимации функций, классификации и распознавания образов, прогнозирования, диагностики и т.д. Переход на нейросетевой базис позволяет улучшить качество функционирования различных технических систем по сравнению с традиционными подходами. В настоящее время исследования в области искусственных нейронных сетей ориентированы в основном на создание специализированных систем для решения конкретных задач. На этом пути существует ряд проблем, препятствующих эффективному развитию нейронных сетей. Основными задачами здесь является выбор моделей нейронных сетей, формирование обучающей выборки, отображение задачи на нейронную сеть, интеграция нейронных сетей в интеллектуальную систему. Данные проблемы являются актуальными при проектировании нейросетевых моделей для решения различных задач, например, обработка хаотических процессов, управление мобильными роботами, прогнозирование и обработка данных.

I. ОБНАРУЖЕНИЕ ЭПИЛЕПТОФОРМНОЙ АКТИВНОСТИ

В данном разделе описывается нейросетевая диагностическая система для детектирова-

ния сегментов с эпилептической активностью в сигналах EEG. В качестве диагностического критерия используется значение старшего показателя Ляпунова, которое снижается при наступлении эпилептических припадков. Старший показатель Ляпунова характеризует среднюю скорость экспоненциального расхождения двух близко лежащих траекторий. Наличие у системы положительной экспоненты Ляпунова свидетельствует о том, что любые две близкие траектории быстро расходятся с течением времени, то есть имеет место чувствительность к значениям начальных условий. На рисунке 1 представлена нейросетевая система обнаружения эпилептической активности в сигналах EEG. На вход системы поступает набор сигналов EEG одной регистрации. Эти сигналы описывают динамику нелинейной хаотической системы, которая характеризует электрическую активность нейронов головного мозга. Каждый EEG сигнал снимается с определенного участка головного мозга, характеризует электрическую активность множества нейронов соответствующего участка головного мозга и содержит различные артефакты (помехи, появляющиеся на ЭЭГ в результате моргания, движения подбородком и т.п.). Поэтому на первом этапе необходимо произвести предобработку EEG сигналов, чтобы отфильтровать их от различного рода артефактов и получить максимально независимые сигналы. Для такой обработки EEG сигналов используется метод независимых компонент (ICA - Independent Component Analysis), который позволяет выделить независимые сигналы из их смеси. Результатом предобработки являются чистые сигналы EEG, содержащие электрическую активность нейронов головного мозга.

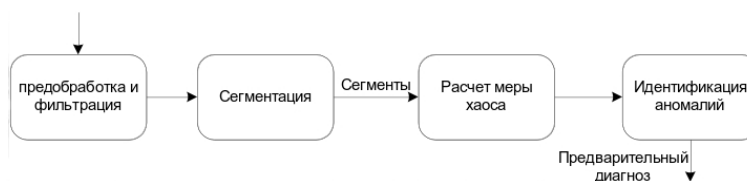


Рис. 1 – Нейросетевая система обнаружения эпилептиформной активности; на вход системы подается набор сигналов электроэнцефалограмм; $\lambda(t)$ -ряд значений старшего показателя Ляпунова

Каждый сигнал, полученный после ICA фильтрации, подвергается адаптивной сегментации при помощи многослойного перцептрона (MLP). В результате сегментации происходит разбиение каждого сигнала на квазистационарные участки, где поведение сигнала не изменяется. Затем для каждого выделенного сегмента производится вычисление оценки старшего показателя Ляпунова. В результате получается детерминированный ряд показателей Ляпунова для каждого чистого сигнала EEG.

$$\lambda(t) = (\lambda_1, \lambda_2, \dots, \lambda_p)$$

, где p – количество временных выделенных сегментов в сигнале EEG. Если различные сегменты имеют одинаковые значения старшего показателя Ляпунова, то они объединяются в один сегмент. На последнем этапе происходит идентификация эпилептической активности в соответствии со следующим критерием:

$$\begin{cases} \lambda > 0, \text{ нормальная активность;} \\ \lambda \leq 0, \text{ эпилептиформная активность.} \end{cases}$$

В результате выполнения данной процедуры для каждого сигнала EEG выделяются временные сегменты с эпилептической и нормальной активностью. Рассмотрим результаты экспериментов по тестированию разработанной диагностической системы. В качестве исходных данных использовались как стандартные EEG данные, взятые из департамента эпилептологии боннского университета и данные, полученные из 5-ой клинической больницы г. Минска. Первый набор данных представляет собой очищенные от артефактов данные, которые состоят из множеств (А-Е) EEG сигналов. Каждое множество состоит из 100 сигналов, и каждый сигнал содержит 4096 отсчетов продолжительностью 23,6 секунды. Множества А и В состоят из EEG сигналов, полученных от здоровых пациентов с открытыми глазами (множество А) и закрытыми глазами (множество В) соответственно. Множества С and D включают EEG фрагменты больных эпилепсией во время отсутствия эпилептического состояния. Множество С состоит из EEG фрагментов, полученных из эпилептической зоны, а множество D состоит из фрагментов, полученных из противоположного полушария головного мозга. Множество Е содержит EEG фрагменты с эпилептической активностью. На рис. 2 представлены результаты экспериментов.

Как видно из рисунка, система не имеет ложных обнаружений на множестве А. Такие же результаты были показаны на множестве В. Таким образом, для здоровых пациентов диагностическая система не имеет ложных обнаружений. В таблице 1 представлены общие резуль-

таты экспериментов для всех множеств. Интересно, что в множестве С, которое состоит из EEG сигналов, полученных с эпилептической зоны каждый раз выделялся один сегмент с эпилептической активностью (таблица 1, 6% обнаружения эпилептической активности). Следует отметить, что для множеств С и D, система обнаруживает эпилептические события. Возможно, это происходит из-за скрытых эпилептических процессов. Если не брать в расчет множества С и D, то общая точность классификации составляет 97,7%, что соответствует лучшим результатам, полученных при анализе данных сигналов EEG.

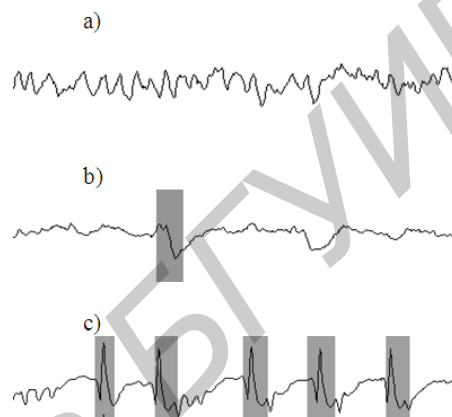


Рис. 2 – Анализ EEG фрагментов множества А, С and Е. а) Во фрагменте множества А не выявлено эпилептической активности. б) Фрагмент множества С имеет один сегмент с эпилептической активностью (выделено серым цветом). в) Пять сегментов с эпилептической активностью выделено (серый цвет) в фрагменте Е.

Таблица 1 – Результаты классификации

Set	Класс 1: эпилептическая активность	Класс 2: нормальная активность
А	0%	100%
В	0%	100%
С	6%	94%
Д	32%	68%
Е	92%	8%

Преимуществом данной системы является то, что она способна не только определять эпилептическую активность в сигналах EEG, но также выделять временные сегменты, где эта активность имеется. Следующие эксперименты были осуществлены с использованием EEG сигналов, полученных из 5-ой клинической больницы г. Минска. EEG сигналы представляют набор данных, которые содержат 21 регистрацию для 8 взрослых пациентов. Каждая регистрация состоит из 16 сигналов EEG. То есть всего было получено 336 сигналов EEG. Каждый EEG сигнал представляет собой временной ряд, состоящий из 2000 отсчетов. Значения статистических параметров, которые характеризуют качество классификации, показаны в таблице 2.

Таблица 2 – Значения статистических параметров

Статистические параметры	Значения
Специфичность	99.7%
Чувствительность	90.6%
Общая точность классификации	99.6%

Как следует из таблицы, ошибка второго рода, что характеризует очень малое значение ложных срабатываний.

II. ОБНАРУЖЕНИЕ И РАСПОЗНАВАНИЕ АТАК НА КОМПЬЮТЕРНЫЕ СЕТИ

Одной из форм глобализации мирового пространства является информационная глобализация, которая связана с повсеместным распространением сети Интернет. В результате этого значительно возросло количество атак и злоупотреблений в сфере высоких технологий. Поэтому вопросу безопасности компьютерных систем уделяется все больше и больше внимания. Задачей Систем Обнаружения Атак (Intrusion Detection Systems - IDS) является защита компьютерных сетей. В последнее время системы IDS активно изучаются. Они должны выполнять свои функции в режиме реального времени. Основным недостатком таких систем является их неспособность обнаруживать новые или неизвестные ата-

ки, т.е. записи о которых в системе отсутствуют. В настоящее время разрабатывается большое количество различных технологий защиты компьютерных сетей, которые базируются на применении нейронных сетей (neural networks), на технологиях извлечения данных (data mining), статистическом анализе и т.п. К недостаткам существующих моделей IDS, в первую очередь, можно отнести уязвимость к новым атакам, низкая точность и скорость работы. Современные системы обнаружения вторжений плохо приспособлены к работе в реальном режиме времени, в то время как возможность обрабатывать большой объем данных в реальном режиме времени – это определяющий фактор практического использования систем IDS. В данном разделе рассматриваются различные варианты архитектур систем IDS, которые базируются на применении рециркуляционной нейронной сети (Recirculation Neural Network - RNN) и многослойного перцепторона. Задачей RNN является сжатие входного пространства образов с целью получения главных компонент. Многослойный перцептрон производит основные вычисления, связанные с распознаванием входного вектора, используя информацию, предоставленную рециркуляционной нейронной сетью. Общая схема обнаружения атаки в сетевом трафике приведена на рис.3. Она включает три этапа.

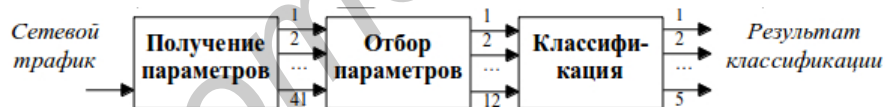


Рис. 3 – Процесс обнаружения

На первом этапе осуществляется захват трафика сети. Сбор необходимых данных выполняет специальное программное средство (sniffer). Для обучения нейронных сетей используется база данных KDD-99. Она содержит около 5 000 000 записей о соединениях. Каждая запись в этой базе представляет собой образ сетевого соединения. Соединение – последовательность TCP пакетов за некоторое конечное время, моменты начала и завершения которого четко определены, в течение которого данные передаются от IP-адреса источника на IP-адрес приемника (и в обратном направлении) используя некоторый определенный протокол. Каждая запись о соединении включает 41 параметр сетевого трафика и промаркирована как “атака” или “не атака”. Второй этап связан с уменьшением размерности

входного вектора данных. Для этого используется рециркуляционная нейронная сеть (RNN), которая сжимает входной сигнал в сигнал главных компонент. В результате экспериментов было определено оптимальное число главных компонент – 12. Третий этап состоит в обнаружении и распознавании атак. Для этих целей используется многослойный перцептрон. В базе KDD-99 представлено 22 типа атак. При этом атаки делятся на четыре основные класса: DoS, U2R, R2L и Probe. Каждый класс в свою очередь состоит из отдельных типов атак. В качестве выходных данных используется 5-мерный вектор, где 5 - это количество классов атак плюс нормальное состояние. Комбинируя RNN и MLP нейронные сети можно получить различные архитектуры систем обнаружения атак.

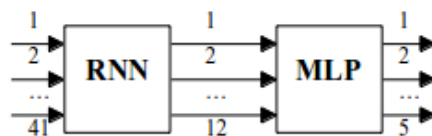


Рис. 4 – Первый вариант IDS (модель 1)

Так на рис. 4 приведена система обнаружения атак, которая состоит из рекуррентной нейронной сети (RNN) и многослойного персептрона (MLP), которые соединены последовательно. Задачей RNN является сжатие входного

41-размерного вектора в 12-размерный выходной вектор. Многослойный персептрон осуществляет обработку сжатого пространства входных образцов (главных компонент) с целью распознавания класса атаки.

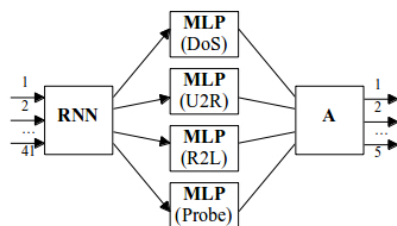


Рис. 5 – Второй вариант IDS (модель 2)

На рис. 5 приведена вторая схема системы обнаружения атак. Она характеризуется тем, что главные компоненты с выходов RNN одновременно поступают на 4 отдельных многослойных персептрона, каждый из которых соответствует определенному классу атаки: DoS, U2R, R2L и Probe. С выходов MLP данные поступают на арбитр, который и принимает окончательное решение о состоянии системы. В качестве арбитра может использоваться линейный или многослойный персептрон. Тогда обучение его будет производиться после обучения RNN и MLP. Такая схема может осуществлять иерархическую классификацию атак. В этом случае арбитр определяет один из 5 классов атаки, а соответствующий многослойный персептрон – тип атаки.

накопленные экспертами, в общее решение, которое имеет приоритет над каждым решением отдельного эксперта. Каждый эксперт представляет собой отдельную систему классификации. В качестве эксперта используется модель 1. Обучение каждого эксперта происходит на отдельном множестве данных, т.е. данные для обучения каждого последующего эксперта формировались с учетом результатов обучения предыдущих экспертов. Алгоритм, используемый для такого обучения, называют алгоритмом усиления за счет фильтрации (boosting by filtering). После обучения нейронные сети способны обнаруживать атаки. В режиме тестирования на вход каждого эксперта подается исходный 41-размерный вектор. Арбитр принимает окончательное решение. Чтобы оценить эффективность предложенных подходов обнаружения вторжений, был проведен ряд экспериментов. Алгоритм усиления за счет фильтрации, который использовался в случае модели 3, предполагает наличие большого (в идеале – бесконечного) множества примеров. Поэтому использовалась 10% выборка из базы KDD (почти 500 000 записей!). Для обучения нейронных сетей были отобраны 6186 примеров. Далее вся 10% выборка применялась для тестирования. Те же наборы данных использовались для обучения и тестирования модели 1 и модели 2, что позволяет сравнивать производительности рассматриваемых в статье систем обнаружения атак друг с другом. Предлагаемые системы обнаружения вторжений осуществляют распознавание 5 классов атак, встречающихся в базе KDD, а именно: DoS, U2R, R2L, Probe и Normal. Для исследования характеристик предложенных архи-

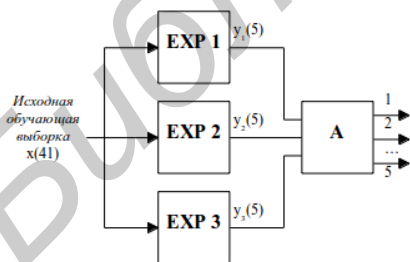


Рис. 6 – Третий вариант IDS (модель 3 - режим тестирования)

Сложные вычислительные задачи могут решаться при помощи их разбиения на множество небольших и простых задач с последующим объединением полученных решений. Вычислительная простота достигается за счет распределения задачи обучения среди множества экспертов (EXP) (рис.6). Такая схема интегрирует знания,

текстур использовались три основных параметра: доля обнаруженных, доля распознанных атак по каждому классу и число ложных срабатываний системы. Доля обнаруженных атак определяется как число образов атак отдельного класса, обнаруженных системой, деленное на общее количество записей об атаках этого класса в базе данных. Подобным образом определяется и доля распознанных атак. Ложные срабатывания указывают общее число нормальных данных сети, классифицированные как атаки. Сводные данные по каждому из рассмотренных вариантов построения системы обнаружения атак приведены в табл. 3:

Таблица 3 – Сводные данные по результатам тестирования каждой модели

Модель	Обнаруж. атаки	Распозн. атаки	Ложн. срабатывания	Общая доля распознанных %
Модель 1	396696 (99.98%)	375522 (94.65%)	46446 (47.75%)	86.30%
Модель 2	395949 (99.80%)	375391 (94.61%)	13398 (13.77%)	92.97%
Модель 3	396549 (99.95%)	375730 (94.70%)	12549 (12.90%)	93.21%

Как видно из таблицы, модель 3 характеризуется высокой точностью (93,21%) и наименьшим числом ложных срабатываний. При использовании модели 1 были распознаны 86,3% вход-

ных образов, а модели 2 – 92,97%. Модели 2 и 3 могут успешно применяться для работы с большими наборами сложных по структуре данных. Таким образом, путем комбинирования двух различных нейронных сетей, а именно RNN и MLP, можно идентифицировать и распознавать атаки на компьютерные сети с достаточно высокой степенью точности. Основными преимуществами использования подходов, основанных на нейронных сетях, является способность адаптироваться к динамическим условиям и быстрота функционирования, что особенно важно при работе системы в режиме реального времени.

III. ЗАКЛЮЧЕНИЕ

В данной статье рассмотрены основные принципы построения нейросетевых систем для диагностики эпилепсии и обнаружения атак на компьютерные сети. Детектирование эпилепсии базируется на вычислении старшего показателя Ляпунова, для определения которого используется многослойный персептрон. Для обнаружения атак на компьютерные сети применяются различные комбинации рециркуляционной нейронной сети и многослойного персептрона. Эксперименты показали эффективность предложенных подходов.

IV. СПИСОК ЛИТЕРАТУРЫ

1. Головки, В. А. 1. Нейронные сети: обучение, организация и применение: // Кн. 4: учеб. пособие для вузов/ Общая ред. А.И. Галушкина.–М.: ИПРЖР, 2001. – С. 256