

УДК 621.391

Оценка эффективности методов классификации состояний телекоммуникационной сети

О. С. Высочина, Салман Амер Мухсин, С. И. Шматков

Харьковский национальный университет имени В.Н. Каразина, Украина

Харьковский национальный университет радиоэлектроники, Украина

В статье оценивается эффективность различных алгоритмов классификации при распознавании состояния телекоммуникационной сети. Анализ алгоритмов классификации проводился при помощи системы анализа данных Weka. Исходным набором данных являлась динамика изменения оценок показателей качества телекоммуникационной сети. Анализ алгоритмов классификации показал, что наиболее эффективным методом классификации являются нейронные сети.

Ключевые слова: телекоммуникационная сеть, нейронные сети, система Weka.

В статті оцінюється ефективність різних алгоритмів класифікації при розпізнаванні стану телекомунікаційної мережі. Аналіз алгоритмів класифікації проводився за допомогою системи аналізу даних Weka. Вхідним набором даних була динаміка зміни оцінок показників якості телекомунікаційної мережі. Аналіз алгоритмів класифікації показав, що найбільш ефективним методом класифікації є нейронні мережі.

Ключові слова: телекомунікаційна мережа, нейронні мережі, система Weka.

In the present paper efficiency of different classification algorithms is estimated at recognition of the telecommunications network state. The analysis of classification algorithms was carried out by means of the Weka data analysis system. The initial data set was a dynamics of quality metrics estimations change of the telecommunications network. The analysis of classification algorithms showed that one of the most effective methods of classification is neural network.

Key words: telecommunications network, neural network, Weka data analysis system.

1. Общая постановка задачи и её актуальность

Возрастающие темпы использования новых технологий информационного обеспечения приводят к увеличению количества сервисов, предоставляемых телекоммуникационной сетью. В связи с этим выдвигаются новые требования к качеству обслуживания. Для обеспечения требуемого качества обслуживания необходимо не только иметь информацию о текущем состоянии сети, но и уметь его прогнозировать. Таким образом, возникает необходимость в разработке новых методов управления телекоммуникационной сетью. Одной из компонент подобной системы управления является система мониторинга.

Системы мониторинга телекоммуникационных сетей присутствуют на рынке уже около восьми лет. Проведенный анализ показал, что подобные системы способны выдать пользователю статистику по ограниченному набору параметров сети, без учета их взаимосвязи. Поэтому для более полного анализа состояния сети в такие системы необходимо включать дополнительные модули обработки статистической информации, работа которых формально сводится к решению задачи классификации состояния телекоммуникационной сети.

Использование подобных модулей позволит оценить, а главное даст возможность спрогнозировать изменение показателей качества сети, с учетом

взаимовлияния и доминирования информационных потоков.

Состояние телекоммуникационной сети характеризуется определенным набором показателей качества. Анализ динамики изменения основных показателей качества работы сети (среднепутевая задержка, джиттер, количество потерянных пакетов, среднее время простоя в очереди и т.д.) позволяет выделить наиболее типовые тренды этих величин. Конкретный вид зависимости определяется текущим состоянием сети, взаимным влиянием различных показателей качества между собой, внешними воздействиями и т.д. Таким образом, возможно по статистически полученной динамике изменения показателей качества распознавать и прогнозировать состояние сети.

Из вышесказанного можно сделать вывод, что задача выбора алгоритма классификации для распознавания состояния телекоммуникационной сети является актуальной.

2. Анализ публикаций

Для решения задачи классификации состояния телекоммуникационной сети существует множество разнообразных подходов и алгоритмов. Классические решения данной задачи рассматриваются в теории распознавания образов [1].

Работы [2–4] посвящены байесовской теории принятия решений, применением разделяющих функций и решением вопросов проверки гипотез. В работе [5] подробно рассмотрен метод потенциальных функций. В работе [6] особое внимание уделено статистической теории распознавания и методу "обобщенный портрет". Метод комитетов описан работе [7]. В работе [8] в качестве метода классификации предложен метод группового учета аргументов. Алгоритмы таксономии и анализа знаний представлены в работе [9]. В работе [10] предложены логические методы распознавания и поиска зависимостей. Многие работы в области теории распознавания и классификации связаны с применением искусственных нейронных сетей [11].

Проведенный анализ показал, что спектр методов достаточно широк, от классического статистического анализа до аппарата искусственных нейронных сетей. Решения, найденные различными алгоритмами, могут существенно отличаться друг от друга. Поиск наилучшего решения затруднен отсутствием общепризнанных универсальных критериев качества решений, поэтому для выбора наиболее эффективного метода использовалась система анализа данных Weka. Система анализа данных Weka представляет собой библиотеку программ, реализующих линейные, комбинаторно-логические, статистические, нейросетевые, гибридные методы прогноза, классификации и извлечения знаний из прецедентов, а также коллективные методы прогноза и классификации [12].

3. Цель статьи

Целью настоящей работы является оценка эффективности различных алгоритмов классификации состояний телекоммуникационной сети при помощи системы анализа данных Weka.

4. Основная часть

Для сбора статистической информации о состоянии сети проведен натурный

експеримент на базі обладнання Cisco. Побудовано сегмент телекомунікаційної мережі, що складається з 5 маршрутизаторів і підтримує роботу 100 абонентів. Мережа працює в штатному режимі (підтримка роботи серверної бази даних, IP-телефонії, електронного документооборота). В процесі експерименту за допомогою SNMP-клієнта з кожного маршрутизатора знімалися значення показувачів бази даних MIB [13], після чого виконувалося усереднення, і визначалася кореляція між станом маршрутизатора і об'ємом трафіку в мережі.

На основі накопленої статистичної інформації виділено 6 типових закономірностей зміни показувачів якості телекомунікаційної мережі. Динаміка зміни основних показувачів якості при різних станах мережі носить повністю визначений характер. Основні залежності інтерполюються відомими формульними співвідношеннями:

$$y = \frac{1}{(1 + e^{-x})} \dots D(f) = R, \dots E(f) = [0;1]; \quad (4.1)$$

$$y = ax + b, \dots a, b \in R, \dots D(f) = R, \dots E(f) = R \dots (a \neq 0), \dots E(f) = \{b\} (a = 0); \quad (4.2)$$

$$y = e^x, \dots D(f) = R, \dots E(f) =]0;+\infty[; \quad (4.3)$$

$$y = \ln x, \dots D(f) =]0;+\infty[, \dots E(f) = R; \quad (4.4)$$

$$y = \log_a x, \dots a > 0, \dots a \neq 1, \dots D(f) =]0;+\infty[, \dots E(f) = R; \quad (4.5)$$

де, $D(f)$ – область визначення функції $f : X \rightarrow R : X = D(f)$;

$E(f)$ – множина значень функції $f : E(f) = \{f(x) | x \in X\} = f(X)$;

В процесі роботи телекомунікаційної мережі показувачі якості приймають значення, що характеризують поточний стан мережі. При зміні стану мережі найбільш чутливі до причини зміни стану показувачі якості відходять від усталених значень згідно своєї динаміки. Задача алгоритму класифікації полягає в розпізнаванні подібних ситуацій.

На вхід алгоритму класифікації надходять значення всіх показувачів якості, а результатом його роботи є розпізнавання, який з показувачів якості динамічно змінюється.

Розглянуто більше 60 алгоритмів класифікації. Оцінка ефективності роботи алгоритмів проводилася за ймовірності правильного розпізнавання стану телекомунікаційної мережі.

Результати проведеного чисельного експерименту і аналіз літератури показали, що найбільш ефективною при розв'язанні задачі класифікації стану телекомунікаційної мережі є наступні алгоритми (Рис. 1):

1. Многослойный перцептрон;
2. Нейросеть с радиальными базисными функциями;

3. Метод опорных векторов;
4. Метод "ближайшего соседа";
5. Наивно-байесовский подход;
6. RandomForest;
7. Алгоритм C4.8;
8. Метод байесовских сетей;
9. Алгоритм OneR;
10. AdaBoost.

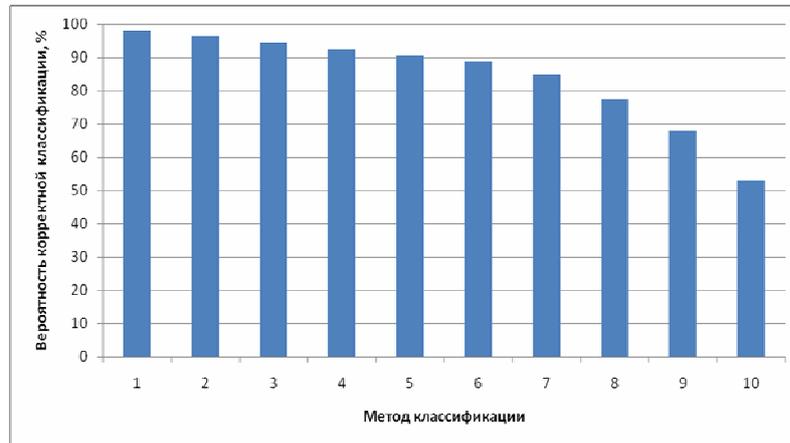


Рис.1. Результаты работы наиболее эффективных алгоритмов классификации при распознавании состояния телекоммуникационной сети

Однако следует отметить у каждого из методов свои особенности применения и недостатки.

AdaBoost [14]

- Склонность к переобучению при наличии значительного уровня шума в данных.
- Требование достаточно длинных обучающих выборок.

Алгоритм OneR [15]

- Невозможность прямой обработки непрерывных переменных, данные необходимо преобразовывать в дискретный вид, что приводит к потере закономерностей.

Метод байесовских сетей [16]

- Вычислительная сложность.
- При попытке учесть большое количество зависимостей между переменными, оценки условных вероятностей приобретают большую дисперсию, таким образом, оценки параметров становятся недостоверными, что в итоге приводит к ухудшению качества классификации.
- Ориентированность на обучающие данные из-за большого количества параметров, что приводит к хорошим результатам классификации на обучающих данных и неудовлетворительным результатам на тестовых данных, т.е. модель описывает не общие закономерности в структуре данных, а скорее набор частных случаев в обучающей выборке.

Алгоритм C4.8 [17]

- Медленная работа на сверхбольших и зашумленных наборах данных.

RandomForest [18]

- Построенная модель занимает большое количество памяти.
- Склонность к переобучению.
- Неспособность к экстраполяции.

Наивно-байесовский подход [19]

- Невозможность непосредственной обработки непрерывных переменных – необходимо их преобразование к интервальной шкале, чтобы атрибуты были дискретными, однако подобные преобразования приводят к потере значимых закономерностей.
- Влияние на результат классификации только индивидуальных значений входных переменных, комбинированное влияние пар или троек значений разных атрибутов не учитывается.

Метод "ближайшего соседа" [20]

- Сложность выбора меры "близости" (метрики). От данной меры главным образом зависит объем множества записей, которые нужно хранить в памяти для достижения удовлетворительной классификации или прогноза.
- Необходимость полного перебора обучающей выборки при распознавании, следствие этого – вычислительная трудоемкость.

Метод опорных векторов [21]

- Для классификации используется не все множество образцов, а лишь их небольшая часть, которая находится на границах.

5. Выводы по результатам и направления дальнейших исследований

При помощи системы анализа данных Weka произведен анализ алгоритмов классификации при распознавании состояния телекоммуникационной сети. Оценка эффективности работы производилась по вероятности правильного распознавания. Рассмотрено более 60 алгоритмов. Исходным набором данных являлась динамика изменения показателей качества телекоммуникационной сети. Анализ полученных результатов и литературы показали, что наилучшим образом при распознавании состояния телекоммуникационной сети себя проявил многослойный персептрон и нейронная сеть с радиальными базисными функциями [11], которые относятся к классу нейронных сетей. Поэтому можно сделать вывод о том, использование нейронных сетей при решении задачи классификации состояния телекоммуникационной сети является одним из наиболее эффективных методов. Поэтому перспективы дальнейших исследований в этом направлении связаны, в первую очередь, с выбором типа нейронной сети, и на основе этого с проведением более детальных модельных экспериментов.

ЛИТЕРАТУРА

1. Фукунага К. Введение в статистическую теорию распознавания образов.– М.: Наука, 1979. – 368с.

2. Neyman, J., Pearson, E.S. On the Problem of the Most Efficient Tests of Statistical Hypotheses // *Phil. Trans. R. Soc.* – 1933. – Series A, №231. – pp. 289–337.
3. Fisher R.A. The use of multiple measurements in taxonomic problems, *Ann. Eugenics.* – 1936. – Part II, №7. – pp. 179–188.
4. Wald A. Contributions to the theory of statistical estimation and testing of hypotheses, *Ann.Math.Stat.* – 1939. – №10. – pp. 299–326.
5. Айзерман М.А., Браверманн Э.М., Розоноэр Л.И. Метод потенциальных функций в теории обучения машин. – М.: Наука, 1970.–384 с.
6. Вапник В.Н., Червоненкис А.Я. Теория распознавания образов (статистические проблемы обучения). – М.: Наука, 1974.–415 с.
7. Мазуров В.Д., Хачай М.Ю. Комитеты систем линейных неравенств// *Автоматика и телемеханика.* – 2004. – №.2. – С. 43–54.
8. Ивахненко А.Г. Системы эвристической самоорганизации в технической кибернетике. – Киев: Техніка. – 1971.–372 с.
9. Загоруйко Н.Г. Прикладные методы анализа данных и знаний. Новосибирск: Изд-во Института математики, 1999.
10. Лбов Г.С. Методы обработки разнотипных экспериментальных данных. – Новосибирск.: Наука, 1981. – 160 с.
11. Хайкин Саймон. Нейронные сети. – М.: Вильямс, 2006. – 1104 с.
12. <http://www.cs.waikato.ac.nz/~ml/weka/>
13. Семенов Ю.А. Протоколы Internet. Энциклопедия. – М.: Горячая линия – Телеком, 2001. – 1100 с.
14. Schapire. Robert E. The boosting approach to machine learning: An overview / Robert.E. Schapire // In MSRI Workshop on Nonlinear Estimation and Classification. – 2002.
15. Мурыгин К.В. Обнаружение объектов на изображении на основе каскада классификаторов / К.В. Мурыгин // Искусственный интеллект. – 2007. – № 2. – С. 104–108.
16. Тулупьев А. Л., Николенко С. И., Сироткин А. В. Байесовские сети: логико-вероятностный подход. СПб.: Наука, 2006. – 607 с.
17. Lopez D., Espana S. Error-correcting tree language inference. *Pattern Recognition Letters* – 2002. – №23. – pp. 1–12.
18. Hastie, T., Tibshirani R., Friedman J. Chapter 15. Random Forests // The Elements of Statistical Learning: Data Mining, Inference, and Prediction. – 2nd ed. – Springer-Verlag, 2009. – 746 p.
19. Brand E., Gerritsen R Naive-Bayes and Nearest Neighbor BMS, 1998. – № 7.
20. Дуда Р., Харт П., Распознавание образов и анализ сцен. – М.: Мир, 1976. – 511 с.
21. Christopher J.C. Burges. A Tutorial on Support Vector Machines for Pattern Recognition, Appeared in: *Data Mining and Knowledge Discovery 2*, 1998 – pp. 121–167.