

Методика выбора конфигурируемых гиперпараметров интеллектуального классификатора неструктурируемых текстовых данных по степени конфиденциальности на основе метода анализа иерархий

В.В. Шадский

Краснодарское высшее военное училище

Аннотация: В данной статье приведена структурная модель интеллектуального классификатора неструктурируемых текстовых данных по степени конфиденциальности, представляющего собой двухуровневый каскадный ансамбль моделей-классификаторов. Выделена наибольшим образом влияющая на эффективность классификации мета-модель полносвязной нейросетевой архитектуры. Произведена декомпозиция многокритериальной задачи конфигурирования интеллектуального классификатора на задачу выбора конфигурируемых гиперпараметров мета-модели и задачу подбора их значений. С учетом выделенных гиперпараметров нейросетевой мета-модели, многокритериальная задача выбора подлежащих конфигурированию гиперпараметров представлена в виде иерархии, включающей в себя цель, критерии и альтернативы. Разработана методика выбора конфигурируемых гиперпараметров интеллектуального классификатора неструктурируемых текстовых данных по степени конфиденциальности на основе метода анализа иерархий.

Ключевые слова: DLP-система, неструктурируемые текстовые данные, интеллектуальный классификатор, гиперпараметры, метод анализа иерархий.

Стремительный рост количества утечек защищаемой информации [1,2] наряду с результатами проведенных исследований [3-5] подчеркивает необходимость внедрения в реализуемые DLP-системами процесс анализа контента интеллектуальных методов. Способствуя расширению перечня детектируемых в составе анализируемых электронных документов конфиденциальных классификационных признаков, данные методы, в конечном итоге, позволяют повысить защищенность обрабатываемой информации.

Ввиду того, что заключительным этапом процесса контентного анализа является решение задачи бинарной классификации электронных документов по степени конфиденциальности, в целях повышения ее эффективности

разрабатываются интеллектуальные модели классификации, основанные на использовании алгоритмов машинного обучения.

В соответствии с [6-8], для решения задач бинарной классификации извлеченных из анализируемых электронных документов неструктурируемых текстовых данных по степени конфиденциальности разработана модель интеллектуального классификатора, состав и структура которой представлены на рис.1.

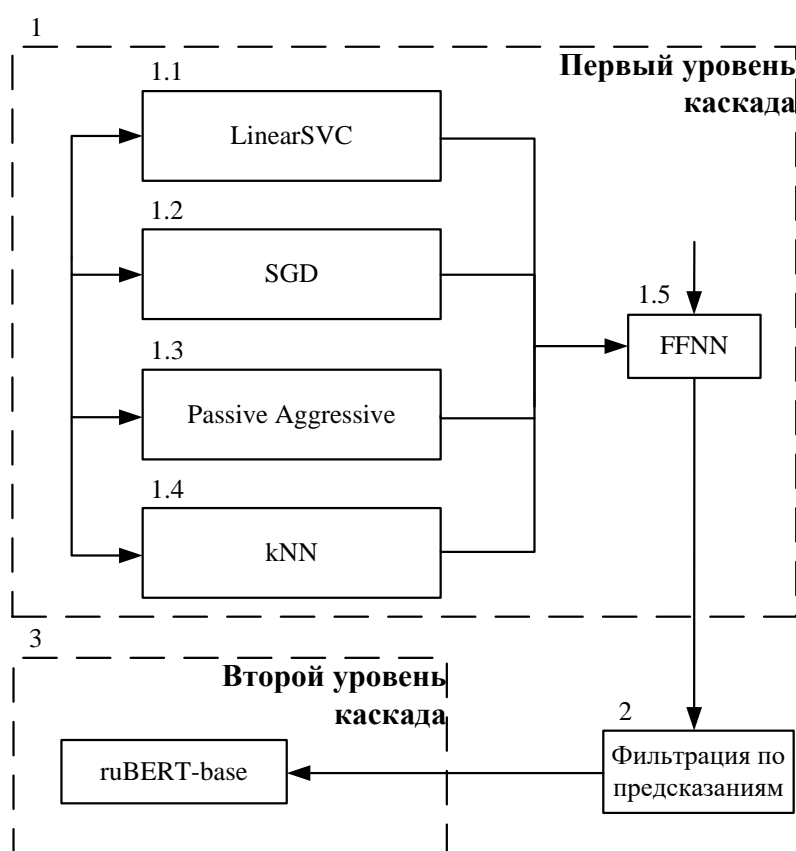


Рис. 1. – Структурная модель интеллектуального классификатора неструктурируемых текстовых данных по степени конфиденциальности

Данная модель представляет собой двухуровневый каскадный ансамбль. Первый уровень каскада (бл. 1 на рис.1) представлен стекингвым ансамблем моделей слабых учеников, в качестве которых выступают

классические модели на основе метода опорных векторов LinearSVC (бл. 1.1), параметрическая модель SGD (бл. 1.2), модель последовательного обучения Passive Aggressive (бл. 1.3), метрическая модель kNN (бл. 1.4), а также мета-модель полносвязной нейросетевой архитектуры FFNN (бл. 1.5). В качестве модели-классификатора второго уровня каскада используется предобученная базовая модель ruBERT-base (бл. 3).

Стоит отметить, что значительная роль в повышении эффективности интеллектуального классификатора отводится наилучшим образом подобранным значениям гиперпараметров, входящих в состав ансамбля моделей-классификаторов. Ввиду относительной немногочисленности данных гиперпараметров моделей слабых учеников LinearSVC, SGD, Passive Aggressive и kNN, а также детерминированности архитектуры модели-классификатора ruBERT-base, наибольшее влияние на эффективность классификации оказывает мета-модель FFNN [9]. Следовательно, процесс конфигурирования интеллектуального классификатора неструктурируемых текстовых данных по степени конфиденциальности сводится к выбору наибольшим образом влияющих на эффективность классификации гиперпараметров мета-модели FFNN и последующему подбору их наилучших значений применительно к конкретной решаемой задаче классификации. В соответствии с работами [10,11], а также исследованиями отечественных ученых [12,13], данный процесс не имеет строгой методологической основы, носит интуитивный характер и базируется на эмпирическом опыте лица, осуществляющего конфигурирование. Кроме того, ввиду многочисленного количества гиперпараметров мета-модели FFNN, данный процесс является слабоинтерпретируемым.

С учетом вышеизложенного возникает необходимость в разработке методики выбора подлежащих конфигурированию гиперпараметров интеллектуального классификатора неструктурируемых текстовых данных

по степени конфиденциальности. После выделения основных гиперпараметров мета-модели FFNN [14], многокритериальная задача выбора наибольшим образом влияющих на эффективность интеллектуального классификатора и подлежащих конфигурированию гиперпараметров представлена в виде иерархии, изображенной на рис.2.



Рис. 2. – Иерархическая декомпозиция многокритериальной задачи выбора конфигурируемых гиперпараметров интеллектуального классификатора

Исходя из рис.2, гиперпараметры мета-модели FFNN представлены в виде двухуровневой иерархии альтернатив, первый уровень которой включает в себя группы гиперпараметров, а второй – соответствующие данным группам гиперпараметры.

Методика выбора конфигурируемых гиперпараметров. Методика основывается на использовании предложенного американским ученым

Томасом Саати метода анализа иерархий [15], являющегося математическим инструментом системного подхода к сложным проблемам принятия решений. Данная методика состоит из двух этапов: подготовительной части и выбора подлежащих конфигурированию гиперпараметров.

I. Подготовительная часть.

а) определение значимости критериев.

Составляется квадратная матрица парных сравнений критериев Cr и вычисляется ее собственный вектор (вектор-столбец) X_{Cr} , значения элементов которого укажут на приоритет того или иного критерия.

б) определение значимости групп гиперпараметров.

Составляются матрицы парных сравнений альтернатив первого уровня иерархии, представляющих собой группы гиперпараметров, по каждому отдельному критерию Al_{Res} , Al_{Rez} .

Для каждой из них вычисляются собственные векторы-столбцы, определяющие приоритет альтернатив первого уровня иерархии по каждому из критериев:

$$X_{AlRes} = \begin{pmatrix} x_{AlRes1} \\ x_{AlRes2} \\ x_{AlRes3} \end{pmatrix},$$

$$X_{AlRez} = \begin{pmatrix} x_{AlRez1} \\ x_{AlRez2} \\ x_{AlRez3} \end{pmatrix}.$$

в) определение значимости гиперпараметров.

составляются матрицы парных сравнений альтернатив второго уровня иерархии, представляющих собой гиперпараметры соответствующих вышеприведенных групп, по каждому отдельному критерию:

- для группы «Архитектурные гиперпараметры» Ar_{Res} , Ar_{Rez} ;

- для группы «Гиперпараметры обучения и оптимизации» L_{Res} , L_{Rez} .

Для каждой из них вычисляются собственные векторы-столбцы, определяющие приоритет альтернатив второго уровня иерархии:

$$X_{ArRes} = \begin{pmatrix} x_{ArRes1} \\ x_{ArRes2} \end{pmatrix},$$

$$X_{ArRez} = \begin{pmatrix} x_{ArRez1} \\ x_{ArRez2} \end{pmatrix},$$

$$X_{LRes} = \begin{pmatrix} x_{LRes1} \\ x_{LRes2} \\ x_{LRes3} \end{pmatrix},$$

$$X_{LRez} = \begin{pmatrix} x_{LRez1} \\ x_{LRez2} \\ x_{LRez3} \end{pmatrix}.$$

II. Оценка гиперпараметров.

Ввиду самостоятельной оценки групп гиперпараметров осуществляется взвешивание полученных векторов приоритетов X_{ArRes} X_{ArRez} X_{LRes} X_{LRez} значениями приоритетов соответствующей группы по каждому из критериев:

$$V_{ArRes} = x_{ArRes1} \begin{pmatrix} x_{ArRes1} \\ x_{ArRes2} \end{pmatrix},$$

$$V_{ArRez} = x_{ArRez1} \begin{pmatrix} x_{ArRez1} \\ x_{ArRez2} \end{pmatrix},$$

$$V_{LRes} = x_{LRes2} \begin{pmatrix} x_{LRes1} \\ x_{LRes2} \\ x_{LRes3} \end{pmatrix},$$

$$V_{LRez} = x_{ALRez2} \begin{pmatrix} x_{LRez1} \\ x_{LRez2} \\ x_{LRez3} \end{pmatrix}.$$

Так как группа «Гиперпараметры регуляризации» представлена одним единственным гиперпараметром, значение приоритета данного гиперпараметра будет определяться значением приоритета группы по каждому их критериев x_{ALRes3} , x_{ALRez3} .

Получение общих оценок гиперпараметров осуществляется посредством перемножения матриц, составленных из взвешенных вектор-столбцов приоритетов гиперпараметров по каждому отдельному критерию, на вектор-столбец приоритетов критериев:

$$V_{Ar} = (V_{ArRes} V_{ArRez})X_{Cr} = \begin{pmatrix} v_{Ar1} \\ v_{Ar2} \end{pmatrix},$$

$$V_L = (V_{LRes} V_{LRez})X_{Cr} = \begin{pmatrix} v_{L1} \\ v_{L2} \\ v_{L3} \end{pmatrix}.$$

Значения элементов полученных вектор-столбцов укажут на приоритеты гиперпараметров данных групп. В свою очередь, получение общего значения приоритета гиперпараметра «Значение коэффициента L2-регуляризации» осуществляется посредством перемножения вектор-столбца, составленного из соответствующих группе «Гиперпараметры регуляризации» значений приоритетов по каждому из критериев на транспонированный вектор-столбец приоритетов критериев:

$$V_{Reg} = \begin{pmatrix} x_{ALRes3} \\ x_{ALRez3} \end{pmatrix} X_{Cr}^T = v_{Reg}.$$

Таким образом, полученные количественные оценки приоритетов укажут степень влияния того или иного гиперпараметра мета-модели FFNN на эффективность интеллектуального классификатора, позволяя тем самым осуществить последующий подбор их значений под условия конкретной решаемой задачи классификации неструктурируемых текстовых данных по степени конфиденциальности.

Литература

1. Отчёт об утечках данных за 1 полугодие 2022 года // InfoWatch URL: infowatch.ru/sites/default/files/analytics/files/otchyot-ob-utechkakh-dannykh-za-1-polugodie-2022-goda_1.pdf (дата обращения: 04.04.2023).
2. 2021 Data Breach Investigations Report // Verizon: [сайт]. – URL: verizon.com/business/resources/reports/2021-data-breach-investigations-report.pdf (дата обращения: 04.04.2023).
3. Машечкин И.В., Петровский М.И., Царев Д.В. Применение методов интеллектуального анализа текстовой информации для предотвращения утечек данных // Программирование. –2015. – № 1. – С. 32-43.
4. Дятлов А.В., Коннова Н.С. Применение методов семантического анализа текста в системах предотвращения утечек информации // Безопасные информационные технологии. – 2017. – № 1. – С. 187-192.
5. Hart Michael, Manadhata Pratyusa, Johnson Rob. Text Classification for Data Loss Prevention // HP Laboratories. – 2011. – № 3. – С. 1-21.
6. Шадский В.В., Сизоненко А.Б., Козленко С.Л., Шишков А.В., Шестаков А.К., Ильин Н.А. Подсистема интеллектуальной адаптивной классификации электронных текстовых документов систем предотвращения утечек информации // Материалы II Межведомственной научно-практической конференции «Кибербезопасность: угрозы, тенденции, технологии защиты». – Краснодарское высшее военное училище. – 2022. – С. 76-79.

7. Шадский В.В., Сизоненко А.Б., Чекмарев М.А., Шишков А.В., Исакин Д.А. Исследование способов векторизации неструктурируемых текстовых документов на естественном языке по степени их влияния на качество работы различных классификаторов // Научно-технический вестник информационных технологий, механики и оптики. – 2022. – № 1. – С. 114-119. doi: 10.17586/2226-1494-2022-22-1-114-119.

8. Шадский В.В., Сизоненко А.Б., Чекмарев М.А., Дудко А.Л. Методика определения архитектуры нейронных сетей для решения задач семантического поиска с использованием метода анализа иерархий // Материалы Всероссийской конференции «Информационные технологии в деятельности органов внутренних дел». – Московский университет Министерства внутренних дел Российской Федерации им. В.Я. Кикотя. – 2021. – С. 129-133.

9. Шадский В.В., Сизоненко А.Б., Шишков А.В., Середин Д.В., Посохов Д.А., Козленко С.Л. Математическая модель оценки эффективности подсистем контентного анализа документов систем предотвращения утечек информации, реализующих дополнительный функционал анализа контента // Электронный сетевой политематический журнал «Научные труды КубГТУ». – 2023. – № 1. – С. 36-47.

10. Greeshma K.V., Sreekumar K. Hyperparameter Optimization and Regularization on Fashion-MNIST Classification // International Journal of Recent Technology and Engineering. – 2019. – № 8. – С. 3713-3719.

11. Kun Cheng Ke, Ming-Shyan Huang. Enhancement of multilayer perceptron model training accuracy through the optimization of hyperparameters: a case study of the quality prediction of injection-molded parts // The International Journal of Advanced Manufacturing Technology. – 2021. – № 118. – С. 2247–2263.

12. Тимофеев А.В. Метод выбора гиперпараметров в задачах машинного обучения для классификации стохастических объектов // Научно-

технический вестник информационных технологий, механики и оптики. – 2020. – № 5. – С. 667-676.

13. Косых Н.Е. Оценка гиперпараметров при анализе тональности русскоязычного корпуса текстов // Интеллектуальные технологии на транспорте. – 2020. – № 3. – С. 41-44

14. Pramoditha Rukshan. Classification of Neural Network Hyperparameters // towards data science: [сайт]. – URL: towardsdatascience.com/classification-of-neural-network-hyper-parameters-c7991b6937c3 (дата обращения: 04.04.2023).

15. Саати Т. Принятие решений: Метод анализа иерархий: учебное пособие / Пер. с англ. Р.Г. Вачнадзе. – Москва: Радио и связь, 1993. – 315 с.

References

1. Otchjot ob utechkah dannyh za 1 polugodie 2022 goda. InfoWatch. URL: infowatch.ru/sites/default/files/analytics/files/otchyot-ob-utechkakh-dannykh-za-1-polugodie-2022-goda_1.pdf (data obrashhenija: 04.04.2023).

2. 2021 Data Breach Investigations Report. Verizon. URL: [verizon.com/business/resources/reports/2021-data-breach-investigations-report.pdf](https://www.verizon.com/business/resources/reports/2021-data-breach-investigations-report.pdf) (data obrashhenija: 04.04.2023).

3. Mashechkin I.V., Petrovskij M.I., Carev D.V. Programmirovaniye. 2015. № 1. pp. 32-43.

4. Djatlov A.V., Konnova N.S. Bezopasnye informacionnye tehnologii. 2017. № 1. pp. 187-192.

5. Hart Michael, Manadhata Pratyusa, Johnson Rob. Text Classification for Data Loss Prevention. HP Laboratories. 2011. № 3. pp. 1-21.

6. Shadskij V.V., Sizonenko A.B., Kozlenko S.L., Shishkov A.V., Shestakov A.K., Il'in N.A. Podsystema intellektual'noj adaptivnoj klassifikacii jelektronnyh tekstovyh dokumentov sistem predotvrashhenija utechek informacii. Materialy II Mezhdomstvennoj nauchno-prakticheskoy konferencii «Kiberbezopasnost':

ugrozy, tendencii, tehnologii zashhity». Krasnodarskoe vyssee voennoe uchilishhe. 2022. pp. 76-79.

7. Shadskij V.V., Sizonenko A.B., Chekmarev M.A., Shishkov A.V., Isakin D.A. Nauchno-tehnicheskij vestnik informacionnyh tehnologij, mehaniki i optiki. 2022. № 1. pp. 114-119. doi: 10.17586/2226-1494-2022-22-1-114-119.

8. Shadskij V.V., Sizonenko A.B., Chekmarev M.A., Dudko A.L. Metodika opredelenija arhitektury nejronnyh setej dlja reshenija zadach semanticheskogo poiska s ispol'zovaniem metoda analiza ierarhij. Materialy Vserossijskoj konferencii «Informacionnye tehnologii v dejatel'nosti organov vnutrennih del». Moskovskij universitet Ministerstva vnutrennih del Rossijskoj Federacii im. V.Ja. Kikotja. 2021. pp. 129-133.

9. Shadskij V.V., Sizonenko A.B., Shishkov A.V., Seredin D.V., Posohov D.A., Kozlenko S.L. Jelektronnyj setevoj politematicheskij zhurnal «Nauchnye trudy KubGTU». 2023. № 1. pp. 36-47.

10. Greeshma K.V., Sreekumar K. Hyperparameter Optimization and Regularization on Fashion-MNIST Classification. International Journal of Recent Technology and Engineering. 2019. № 8. pp. 3713-3719.

11. Kun Cheng Ke, Ming-Shyan Huang. The International Journal of Advanced Manufacturing Technology. 2021. № 118. pp. 2247–2263.

12. Timofeev A.V. Nauchno-tehnicheskij vestnik informacionnyh tehnologij, mehaniki i optiki. 2020. № 5. pp. 667-676.

13. Kosyh N.E. Intellektual'nye tehnologii na transporte. 2020. № 3. pp. 41-44

14. Pramoditha Rukshan. Classification of Neural Network Hyperparameters. Towards data science: [sajt]. URL: towardsdatascience.com/classification-of-neural-network-hyper-parameters-c7991b6937c3 (data obrashhenija: 04.04.2023).

15. Saati T. Prinjatie reshenij: Metod analiza ierarhij [Decision-making: Hierarchy Analysis Method]: uchebnoe posobie. Per. s angl. R.G. Vachnadze. Moskva: Radio i svjaz', 1993. 315 p.
