

АҚПАРАТТЫҚ ҚАУІПСІЗДІК
ИНФОРМАЦИОННАЯ БЕЗОПАСНОСТЬ
INFORMATION SECURITY

DOI 10.51885/1561-4212_2024_2_210
MPHTI 81.93.29

А.Е. Хохлова¹, О.Е. Бакланова², Ш.Т. Тезекпаева³

Восточно-Казахстанский технический университет имени Д. Серикбаева,
г. Усть-Каменогорск, Казахстан

¹E-mail: nsox2021@mail.ru

²E-mail: obaklanova@ektu.kz*

³E-mail: shtezekpaeva@ektu.kz

ПРИМЕНЕНИЕ МЕТОДА НАИВНОГО БАЙЕСА ПРИ РЕШЕНИИ ЗАДАЧИ ФИЛЬТРАЦИИ СПАМА

СПАМДЫ СҮЗУ МӘСЕЛЕСІН ШЕШУДЕ АҢҒАЛ БАЙЕС ӘДІСІН ҚОЛДАНУ

APPLICATION OF THE NAIVE BAYES METHOD IN SOLVING THE SPAM FILTERING PROBLEM

Аннотация. Актуальность исследования темы фильтрации спама заключается в том, что в настоящее время спам остается одной из самых крупных проблем, с которыми сталкивается интернет-сообщество. В настоящий момент эта проблема не решена до конца, так как появляются всё новые способы организации вредоносной или просто неприятной рассылки.

Один из локальных методов фильтрации спама – байесовская фильтрация. Это широкий класс алгоритмов классификации, основанный на принципе максимума апостериорной вероятности. Цель исследования заключается в оценке эффективности применения метода наивного Байеса в качестве метода фильтрации спама и выявить его преимущества и недостатки по сравнению с другими методами. Целесообразность исследования заключается в том, что наивный Байес является одним из самых популярных и доступных методов фильтрации спама. Этот метод основан на простых принципах и позволяет эффективно решать проблему фильтрации спама, используя минимум ресурсов. Практическая значимость исследования данного метода заключается в том, что по сравнению с другими алгоритмами классификации, метод наивного Байеса имеет высокую скорость обучения и может обрабатывать большое количество функций.

Ключевые слова: спам; фильтрация; условная вероятность; конфиденциальность; вредоносные рассылки; электронная почта, байесовская фильтрация

Аңдатпа. Спамды сүзу тақырыбын зерттеудің өзектілігі-қазіргі уақытта спам Интернет-қауымдастықтың алдында тұрған ең үлкен мәселелердің бірі болып қала береді. Қазіргі уақытта бұл мәселе толығымен шешілмеген, өйткені зиянды немесе жай ғана жавымсыз ақпараттық бюллетеньді ұйымдастырудың барлық жаңа тәсілдері пайда болады

Байес фильтрациясы – максималды артқы ықтималдық принципіне негізделген жіктеу алгоритмдерінің кең класы. Зерттеудің мақсаты – аңғал Байесті спамды сүзу әдісі ретінде қолданудың тиімділігін бағалау және оның басқа әдістермен салыстырғанда артықшылықтары мен кемшіліктерін анықтау. Зерттеудің орындылығы-аңғал Байес спамды сүзудің ең танымал және қол жетімді әдістерінің бірі болып табылады. Бұл әдіс қарапайым принциптерге негізделген және ең аз ресурстарды пайдалана отырып, спамды сүзу мәселесін тиімді шешуге мүмкіндік береді. Бұл әдісті зерттеудің практикалық маңыздылығы басқа жіктеу алгоритмдерімен салыстырғанда, аңғал Байес әдісі жоғары оқу жылдамдығына ие және көптеген функцияларды өңдей алады.

Түйін сөздер: спам; сүзу; шартты ықтималдық; құпиялылық; зиянды хабарлар; электрондық пошта, Байес сүзгісі

Abstract. The relevance of the study on spam filtering lies in the fact that spam remains one of the biggest problems faced by the internet community. Currently, this problem has not been fully solved as new ways of organizing malicious or just unpleasant mailing are appearing.

In this work, one of the local methods of filtering, Bayesian filtering, has been considered. Bayesian filtering is a wide class of classification algorithms based on the principle of maximum a posteriori probability. The goal of the research is to assess the effectiveness of using Naive Bayes as a method of spam filtering. The expediency of the study lies in the fact that Naive Bayes is one of the most popular and affordable methods of spam filtering. This method is based on simple principles and allows you to effectively solve the problem of spam filtering using a minimum of resources. The practical significance of the study of this method lies in the fact that, compared with other classification algorithms, the Naive Bayes method has a high learning rate and can process a large number of functions.

Keywords: spam; filtering; conditional probability; confidentiality; malicious mailings; email, Bayesian filtering.

Введение. Спам – это массово рассылаемые рекламные сообщения конкретных людей или организаций для тех, кто не желает получать таковые сообщения. С появлением в современном мире сети Интернет и связанной с ней электронной почтой резко возникла потребность в фильтрации спама в связи с тем, что распространение спам-сообщений стало очень частым явлением для владельцев виртуальных ящиков. В настоящий момент эта проблема не решена до конца, так как появляются всё новые способы организации вредоносной или просто неприятной рассылки [1,2].

Фильтрация спама является важным аспектом информационной безопасности в электронной почте. Существует множество методов и алгоритмов, которые используются для фильтрации спама, одним из них является метод наивного Байеса [3]. Наивный Байес – это простой и эффективный метод классификации, который основывается на предположении о независимости признаков. Цель данного исследования – эффективность и применимость метода наивного Байеса в процессе фильтрации спама, а также выявить возможные проблемы и пути их решения. Основными задачами исследования являются:

1. Исследовать существующие методы фильтрации спама и оценить их эффективность.
2. Разработать и описать метод наивного Байеса для фильтрации спама.
3. Подтвердить или опровергнуть гипотезу о том, что метод наивного Байеса является эффективным инструментом для фильтрации спама.
4. Рассмотреть перспективы дальнейшего развития и применения метода наивного Байеса для фильтрации спама.

Литературный обзор. Фильтрация спама является одной из самых важных проблем в электронной почте. Она необходима, чтобы отделить ценную информацию от мусора. Решение этой проблемы стало важным направлением в компьютерных науках, и в последние годы метод наивного Байеса стал одним из наиболее изучаемых и используемых методов фильтрации спама. Множество исследований было посвящено эффективности этого метода, а также его сравнению с другими методами фильтрации спама. Например, в статье «An overview of recent trends in email spam filtering techniques» [4] были рассмотрены различные методы фильтрации спама, включая методы машинного обучения, байесовскую фильтрацию и использование базы данных черных списков. Был сделан вывод, что наилучшие результаты показали методы машинного обучения. При рассмотрении различных методов машинного обучения в статье «Spam filtering using machine learning algorithms: a review» [5] был проведен анализ сильных и слабых сторон различных методов и выработаны рекомендации для выбора наилучшего метода в зависимости от конкретных потребностей. По результатам сравнения одним из лучших алгоритмов машинного обучения для фильтрации спама является байесовская фильтрация. Для повышения эффективности фильтрации был предложен метод, который объединяет в себе несколько байесовских классификаторов [6]. Также для улучшения точности фильтрации спама применяется использование метода наивного Байеса в сочетании с другими методами машинного обучения, такими как дерево решений [7].

Материалы и методы исследования. На сегодняшний день для борьбы со спамом существует множество методов по его устранению. Все методы борьбы можно разделить по способу их организации на две категории: распределенные и локальные [8] (рис. 1).

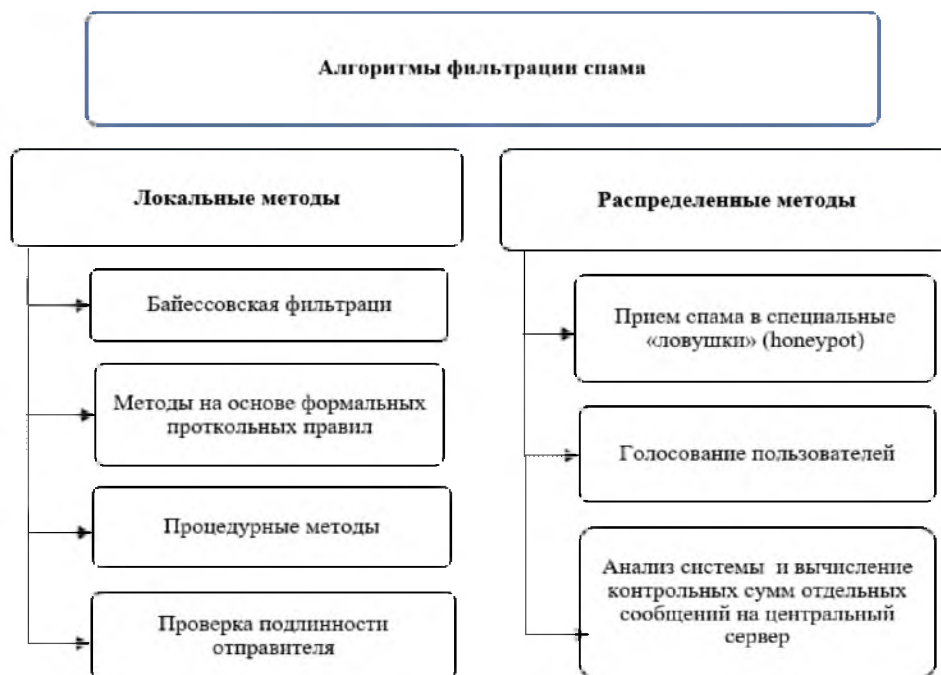


Рисунок 1. Алгоритмы фильтрации спама

Распределенные методы предполагают участие в сборе информации о спаме от большого числа независимых почтовых систем, которые обмениваются данными между собой. Каждая из систем-участниц предоставляет необходимую (специфичную для каждого метода) информацию о проходящем через нее почтовом трафике, тем самым пополняя базу данных информации о спаме. Качество фильтрации достигается привлечением как можно большего числа участников и совершенствованием механизмов сбора и анализа информации о спаме. Чем больше точек сети предоставляют информацию о спаме и чем качественнее эта информация, тем полнее становится картина действий спамеров и тем эффективнее можно с ними бороться. Однако в рамках распределенных методов фильтрации спама отсутствует возможность тонкой настройки фильтра в отдельно взятой почтовой системе.

Локальные методы работают в рамках одной почтовой системы и не используют для работы внешних ресурсов. Так как эти методы не предполагают получения информации о спаме из внешних источников, то каждый раз при изменении вида входящих писем или тактики спамеров, приводящих к большому числу ошибок фильтра, настройка фильтра под характер почтового трафика и работа по повышению качества фильтрации полностью ложится на администратора. Но, в отличие от распределенных методов фильтрации, локальные методы изначально имеют возможность тонкой адаптации под конкретную почтовую систему.

Наивные байесовские классификаторы [9, 10] – популярный локальный метод фильтрации электронной почты. Обычно они используют функции набора слов для идентификации спама по электронной почте – подход, используемый при классификации

текстов. Для классифицируемого объекта вычисляются функции правдоподобия каждого из классов, по ним вычисляются апостериорные вероятности классов. Объект относится к тому классу, для которого апостериорная вероятность максимальна.

Формула Байеса позволяет «переставить причину и следствие» [11]: по известному факту события вычислить вероятность того, что оно было вызвано данной причиной. События, отражающие действие «причин», в данном случае называют гипотезами, так как они – предполагаемые события, повлекшие данное. Безусловную вероятность справедливости гипотезы называют априорной (насколько вероятна причина вообще), а условную вероятность с учётом факта произошедшего события – апостериорной (насколько вероятна причина оказалась с учётом данных о событии). Основная формула:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}, \quad (1)$$

где: $P(A|B)$ – вероятность гипотезы A при наступлении события, условная (апостериорная) вероятность; $P(B|A)$ – вероятность наступления события B при истинности гипотезы A ; $P(A)$ – априорная вероятность гипотезы A ; $P(B)$ – полная вероятность наступления события B .

Применим теорему Байеса для описания вероятности того, что электронное письмо является спам-сообщением [12].

Электронное письмо является спам-сообщением, если встречается слово w , и определяется вероятностью того, что это слово w находится в спам-сообщении s , умноженной на общую вероятность того, что электронное письмо является спам-сообщением s . Это делится на вероятность появления этого слова в электронном письме (спам и не спам вместе взятые):

$$P(s|w) = \frac{P(W|S)P(s)}{P(w|s \cup h)}. \quad (2)$$

Чтобы решить, является ли электронное письмо спамом, нужно получить единую вероятность P для всего электронного письма, а не только для отдельных слов ($p_1 \dots p_n$). Поскольку спам-фильтр использует байесовский подход, достичь этого можно, перемножив вероятности для каждого слова и разделив их произведение на сумму совокупной вероятности того, что каждое слово окажется в спам-сообщении, и совокупной вероятности того, что каждое слово не окажется в спам-сообщении. Формула для достижения этой цели выглядит следующим образом:

$$P = \frac{p_1 p_2 \dots p_n}{p_1 p_2 \dots p_n + (1-p_1)(1-p_2) \dots (1-p_n)}. \quad (3)$$

Теперь, когда у нас есть вероятность того, что электронное письмо является спамом, нужно решить, является ли это электронное письмо спамом или нет. Отнесение письма к «спаму» производится при превышении его «веса» некой планки, заданной пользователем (обычно берут 60-80 % [13]). Если количество спама в электронном письме выше порогового значения, оно классифицируется как спам и получает индекс 1 (спам), если ниже – это безвредное электронное письмо и оно получает индекс 0 (не спам).

Для начала работы потребуется набор данных с помеченными спам-письмами и безвредными электронными письмами. Используемый набор данных содержит 5728 электронных писем и два столбца:

- Столбец «text», строка, содержащая сообщение, полученное из электронной почты;
- Столбец «спам», в двоичном формате, который классифицирует электронное письмо как спам (1) или не спам (0). Пример спам-сообщения приведен на рис. 2.

```
'Subject: unbelievable new homes made easy im wanting to show you this h
omeowner you have been pre - approved for a $ 454 , 169 home loan at a 3
. 72 fixed rate . this offer is being extended to you unconditionally and
your credit is in no way a factor . to take advantage of this limited tim
e opportunity all we ask is that you visit our website and complete the
1 minute post approval form look foward to hearing from you , dorcas pit
tman'
```

Рисунок 2. Пример сообщения, содержащего спам

Чтобы создать функционирующий спам-фильтр, необходимо выполнить следующие шаги:

1. Удаление стоп-слов и специальных символов [14]. Многие слова не несут никакого значения для спам-фильтра. Эти слова рассматриваются как стоп-слова и могут быть полезны, только если мы хотим рассмотреть контекст электронного письма. Но поскольку теорема Байеса рассматривает каждое слово независимо (наивно), мы не будем использовать контекст. После очищения текста в нашем наборе мы понижаем регистр всех слов и избавляемся от заголовка «Subject: (Тема)» в каждом электронном письме.

2. Разделение набора данных на тренировочный/тестовый набор [15]. Набор данных содержит в общей сложности 5728 электронных писем с индексом 0 для не спама и 1 для спама. Более четверти электронных писем являются спамом (1368 из 5728). Следующий шаг – разделить набор данных на тестовый и обучающий наборы. Установим размер тестового набора на 0,3 или 30 %, чтобы обеспечить последовательность устанавливаем random_state функции train_test_split равным 42. Итог: набор данных теперь разделен на тренировочный (4009 электронных писем) и тестовый (1719 электронных писем) наборы.

3. Вычисление общей вероятности того, что электронное письмо является спамом или не спамом. После разделения на обучающий и тестовый наборы, мы можем рассчитать вероятность того, является электронное письмо спамом или нет. Для этого количество писем со спамом (n_{spam}) и не спамом (n_{ham}) делим на общее количество писем (n_{emails}):

$$P(s) = \frac{n_{spam}}{n_{emails}} = \frac{1368}{5728} = 0,23, \quad (4)$$

$$P(\neg s) = \frac{n_{ham}}{n_{emails}} = \frac{4360}{5728} = 0,77. \quad (5)$$

4. Вычисление условной вероятности появления слова в спаме ($P(w|s)$). Вероятность того, что слово будет найдено в спам-письме ($P(w|s)$), определяется количеством повторений этого слова в спам-письмах деленным на общее количество спам-писем. То же самое касается вероятности того, что слово находится в электронном письме, которое не является спамом. Этот подход не учитывает, встречается или нет одно и то же слово в электронном письме несколько раз и в каком отношении эти слова относятся друг к другу. Это означает, что контекст ни на что не влияет. Предложения вроде «БЕСПЛАТНАЯ! БЕСПЛАТНАЯ! Вещь», которые указывают на спам, не интерпретируются как таковые, а вместо этого рассматриваются как уникальное количество слов со словами «бесплатно» и «вещь». В нашем случае ($P(w|s)$) – это вероятность того, что эти слова появятся в спам-сообщении.

$$P(w|s) = \frac{\text{number of spam email constaining } w}{\text{total number of spam emails}} \quad (6)$$

$$P(w|\neg s) = \frac{\text{number of ham email constaining } w}{\text{total number of spam emails}} \quad (7)$$

Для большей эффективности мы создаем словарь, содержащий все слова (кроме стоп-слов) электронного письма в обучающем наборе, и предварительно рассчитываем вероятность того, что слово окажется в спаме или не спаме.

5. Расчет коэффициента правдоподобия (LR). После вычисления вероятности для каждого слова можно рассчитать коэффициент правдоподобия. Это соотношение дает информацию о том, насколько полезно слово в качестве индикатора спама. Коэффициент правдоподобия показывает, во сколько раз больше шансов получить определенное слово в категории не спам вместо спама. Таким образом, чем меньше число, тем больше вероятность того, что это слово появляется только в спам-письмах. Коэффициент правдоподобия для слова рассчитывается по следующей формуле:

$$L = \frac{P(w|s)}{P(w|h)}, \quad (8)$$

где $P(w|s)$ – вероятность попадания слова в спам; $P(w|h)$ – общая вероятность попадания слова в не спам.

6. Расчет общей вероятности P (спама) для каждого электронного письма. Для эффективного расчета необходимо объединить формулы (2) и (3) в одну:

$$P = \frac{P(s) \prod_{i=1}^n P(w_i|s)}{\prod_{i=1}^n P(w_i|s \cup h)}. \quad (9)$$

Вероятность того, что слово будет найдено в спам-письме ($P(w|s)$), определяется количеством появлений этого слова в спам-письмах, деленным на общее количество спам-писем. То же самое касается вероятности того, что слово находится в электронном письме, не содержащем спам ($P(w|\neg s)$). Визуализируем наш результат, построив гистограмму (рис. 3). Как видим, байесовский подход имеет тенденцию доводить значения до крайних точек (0 и 1).

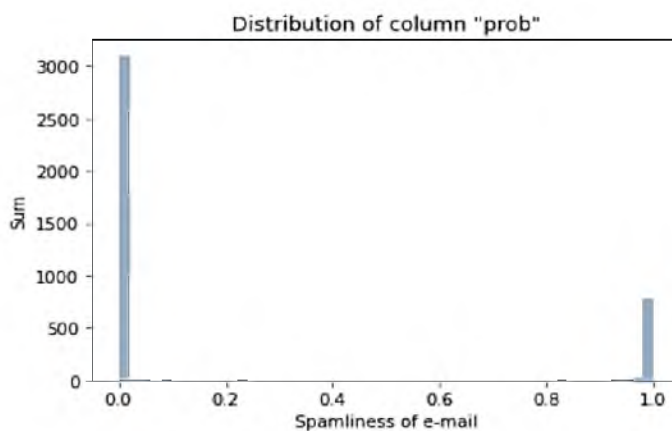


Рисунок 3. Результат расчета спама с помощью байесовского подхода для каждого письма

7. Выбор порогового значения для определения спама. Теперь, когда есть функция, которая вычисляет вероятности для наших электронных писем, можно запустить ее в нашем тестовом наборе, предварительно протестировав несколько различных пороговых значений в нашей функции классификации. Необходимо свести к минимуму количество ложноположительных результатов, т.к. мы не хотим, чтобы электронные письма, не являющиеся спамом, попадали в ящик спама. Для начала необходимо рассчитать показатели качества классификации [16]. Введем следующие переменные:

1. TR (true positive) – верно классифицированные не спам письма;
2. TN (true negative) – верно классифицированные спам-письма;
3. FR (false positive) – не верно классифицированные не спам-письма;
4. FN (false positive) – не верно классифицированные спам-письма;
5. P (actual positive) – фактическое количество не спам-писем;
6. N (actual negative) – фактическое количество спам-писем.

Рассчитаем TPR – долю верно классифицированных положительных не спам-писем по отношению к общему количеству не спам-писем в тренировочном наборе. TPR называется чувствительностью классификации и рассчитывается по следующей формуле:

$$TPR = \frac{TP}{P} = \frac{TP}{TP+FN}. \quad (10)$$

Рассчитаем TNR – долю верно классифицированных спам-писем по отношению к общему количеству спам-писем в тренировочном наборе. TNR называется специфичностью классификации и рассчитывается по следующей формуле:

$$TNR = \frac{TN}{N} = \frac{TN}{TN+FP}. \quad (11)$$

Рассчитаем OCR – долю верно классифицированных спам-писем от общего количества всех писем электронной почте. OCR – называется общей точностью классификации и рассчитывается по следующей формуле:

$$OCR = \frac{TP+TN}{P+N} = \frac{TP+TN}{TP+TN+FP+FN}. \quad (12)$$

Нам необходимо свести к минимуму количество ложных срабатываний. Для этого тестируем несколько различных пороговых значений в нашей функции классификации и выбираем пороговое значение, при котором минимальное значение у параметра FN и максимальная точность OCR . Построим график зависимости порогового значения и показателей качества классификации (рис. 4).

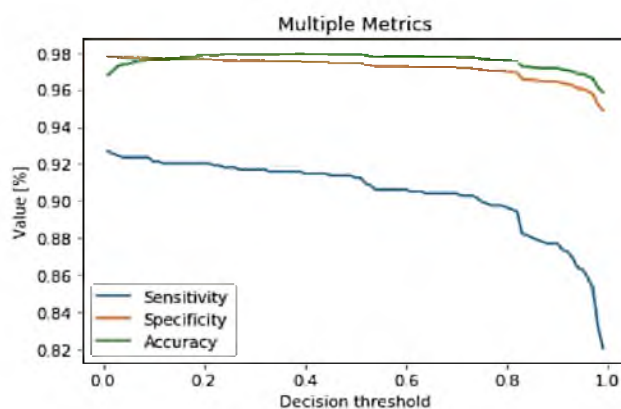


Рисунок 4. Зависимость показателей качества классификации от порогового значения

Результаты и их обсуждение. Используя пороговое значение, полученное на предыдущем этапе, мы можем протестировать наш классификатор на тестовом наборе. Для оценки качества результатов обучения построим матрицу ошибок (confusion matrix). Confusion matrix – это матрица, используемая в машинном обучении для оценки качества классификации [17]. Она показывает количество правильных и неправильных классификаций модели, распределенных по каждой категории. Это помогает выявить слабые и сильные стороны модели и принять решение о ее дальнейшей оптимизации.

Как видно в приведенной матрице ошибок (рис. 5), в результате обучения с помощью байесовского подхода мы получаем хороший результат: только два электронных письма в нашем тестовом наборе ошибочно помечены как спам (ложноположительный результат), при этом 90 спам-писем классифицированы как безвредные электронные письма (ложноотрицательный результат). Это следствие того, что в обучающей выборке было недостаточное количество спам-писем и спам-слов, чтобы классифицировать данные верно. Точность этого классификатора составляет более 94 % на тестовом наборе.

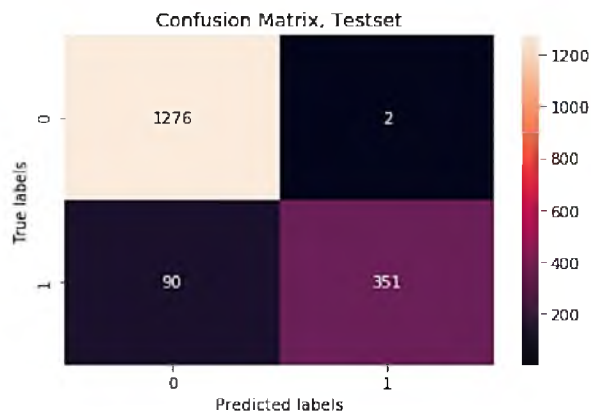


Рисунок 5. Матрица ошибок для тестового набора данных

Заключение. Метод наивного Байеса все еще остается одним из наиболее популярных и надежных методов для фильтрации спама. Дальнейшее развитие этого метода зависит от постоянно растущего объема информации и новых технологий, таких как искусственный интеллект и глубокое обучение. Одной из перспектив развития является совершенствование алгоритмов обучения, что позволит улучшить точность фильтрации спама. Уже разрабатываются сложные модели машинного обучения, которые способны выявлять спам на основе множества различных факторов, включая текстовые, графические и структурные элементы.

По сравнению с другими алгоритмами классификации, главное преимущество метода наивного Байеса состоит в том, что он может обрабатывать большое количество функций. В нашем примере тысячи разных слов и каждое слово рассматривается как функция. Кроме того, даже если есть нерелевантные функции, это дает хороший эффект. Следующее важное преимущество – относительная простота: метод можно использовать напрямую, и параметры редко нужно корректировать, если только данные распределения не известны и не нуждаются в корректировке. И последним преимуществом метода является быстрое обучение и скорость прогнозирования относительно объема данных, которые он может обработать.

Однако метод наивного Байеса имеет и ряд недостатков: во-первых, он не может обрабатывать взаимодействия между функциями, что может привести к низкой точности предсказания в некоторых случаях; во-вторых, может быть неустойчив к выбросам или выбросам данных, так как метод основывается на доверительных предположениях о распределениях.

Список литературы

1. Мурашов А. В., Сафонов В. В. ПРОБЛЕМА СПАМА И ЕЕ РЕШЕНИЕ //Наука и Образование. – 2020. – Т. 3. – №. 2.
2. Отчет «Лаборатории Касперского» по спаму и фишингу за 2021 год – Электронный ресурс.

3. M.G. Hossain, M.Z. Islam Naive Bayes Algorithm: A Comprehensive Study. – 2021.
4. Ehsan. M. and Ali. S. An overview of recent trends in email spam filtering techniques//Journal of Ambient Intelligence and Humanized Computing. – 2021. 1-22.
5. Tiwari. A. and Sharma. P. Spam filtering using machine learning algorithms: a review//Journal of Ambient Intelligence and Humanized Computing – 2020. – 11(4), 1-22.
6. J.M. Kim, H.S. Kim Naive Bayes Algorithm for Multi-class Text Classification: A Survey. – 2021.
7. Лютова Е.И, Коломойцева И.А. Анализ алгоритмов фильтрации спама. – 2020. – С. 116-120.
8. Е.В. Шарапова, Р.В. Шарапов Обнаружение почтового спама на основе сигнатур электронных писем // V Междунар. конф. и молодёжная школа «Информационные технологии и нанотехнологии» (ИТНТ-2019). Владимирский гос. ун-т, Муромский ин-т. – 2019. – С. 924-930.
9. Буртолик Д.О. Байесовские методы классификации // Прикладная математика: современные проблемы математики, информатики и моделирования. – 2020. – С. 320-324.
10. Ломкина Л.С., Субботин А. Н. Классификация потоковых данных на основе байесовского критерия // Моделирование, оптимизация и информационные технологии. – 2020. – Т. 8. – № 1 (28). – С. 18.
11. Цицина А.С., Хоменко Т.В. формирование принятия решений на основе теоремы Байеса для горнодобывающих предприятий региона. – 2021.
12. Мальцева Д.Н., Лукин Д.В. Применение теоремы Байеса для фильтрации спама // Общество-наука-инновации. – 2021. – С. 6-8.
13. Gao H., Zeng X., Yao C. Application of improved distributed naive Bayesian algorithms in text classification // The Journal of Supercomputing. – 2019. – Т. 75. – С. 5831-5847.
14. Sarica S., Luo J. Stopwords in technical language processing // Plos one. – 2021. – Т. 16. – № 8.
15. Vabalas A. et al. Machine learning algorithm validation with a limited sample size // PloS one. – 2019. – Т. 14. – № 11.
16. Старовойтов, В.В. Сравнительный анализ оценок качества бинарной классификации / В.В. Старовойтов, Ю.И. Голуб // Информатика. – 2020. – Т. 17. – № 1. – С. 87-101.
17. Liang J. Confusion Matrix: Machine Learning // POGIL Activity Clearinghouse. – 2022. – Т. 3. – № 4.

References

1. Murashov A.V., Safonov V.V. THE PROBLEM OF SPAM AND ITS SOLUTION //Science and Education. – 2020. – Т. 3. – No. 2.
2. Отчет «Лаборатории Касперского» по спаму и фишингу за 2021 год – Электронный ресурс
3. M.G. Hossain, M.Z. Islam Naive Bayes Algorithm: A Comprehensive Study. – 2021.
4. Ehsan. M. and Ali. S. An overview of recent trends in email spam filtering techniques//Journal of Ambient Intelligence and Humanized Computing. – 2021. 1-22.
5. Tiwari. A. and Sharma. P. Spam filtering using machine learning algorithms: a review//Journal of Ambient Intelligence and Humanized Computing – 2020. – 11(4), 1-22
6. J.M. Kim, H.S. Kim Naive Bayes Algorithm for Multi-class Text Classification: A Survey. – 2021.
7. Lyutova E. I, Kolomojceva I.A Analiz algoritmov fil'tracii spama. – 2020. – S.116-120
8. E.V. SHarapova, R.V. SHarapov Obnaruzhenie pochtovogo spama na osnove signatur elektronnyh pisem // V Mezhdunar. konf. i molodyozhnaya shkola «Informacionnye tekhnologii i nanotekhnologii» (ITNT-2019). Vladimirsij gosudarstvennyj un-t, Muromskij in-t. – 2019. – S.924-930
9. Burtolik D. O. Bajesovskie metody klassifikacii //Prikladnaya matematika: sovremennye problemy matematiki, informatiki i modelirovaniya. – 2020. – S. 320-324.
10. Lomakina L. S., Subbotin A. N. Klassifikaciya potokovyh dannyh na osnove bajesovskogo kriteriya // Modelirovanie, optimizaciya i informacionnye tekhnologii. – 2020. – Т. 8. – № 1 (28). – S. 18.
11. Cicina A. S., Homenko T. V. Formirovanie prinyatiya reshenij na osnove teoremy Bajesa dlya gornodobyvayushchih predpriyatij regiona. – 2021.
12. Mal'ceva D.N., Lukin D.V. Primenenie teoremy Bajesa dlya fil'tracii spama // Obshchestvo-nauka-innovacii. – 2021. – S. 6-8.
13. Gao H., Zeng X., Yao C. Application of improved distributed naive Bayesian algorithms in text classification // The Journal of Supercomputing. – 2019. – Т. 75. – С. 5831-5847.
14. Sarica S., Luo J. Stopwords in technical language processing //Plos one. – 2021. – Т. 16. – № 8.
15. Vabalas A. et al. Machine learning algorithm validation with a limited sample size // PloS one. – 2019. – Т. 14. – № 11.
16. Starovojtov, V.V. Sravnitel'nyj analiz ocenok kachestva binarnoj klassifikacii / V.V. Starovojtov, Yu.I. Golub // Informatika. – 2020. – Т.17, №1. – S. 87-101.
17. Liang J. Confusion Matrix: Machine Learning // POGIL Activity Clearinghouse. – 2022. – Т. 3. – № 4.