

МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РОССИЙСКОЙ
ФЕДЕРАЦИИ

Федеральное государственное бюджетное образовательное учреждение
высшего образования
«МОСКОВСКИЙ АВИАЦИОННЫЙ ИНСТИТУТ
(НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ)»

Кафедра 319 «Системы интеллектуального мониторинга»

КУРСОВАЯ РАБОТА

по дисциплине «Технология разработки программного
обеспечения»

**«Проектирование и разработка веб-приложения
классификации новостей с применением методов
машинного обучения»**

Студент _____ Пожидаев Е.В.

Группа _____ МЗО – 221М – 19

Руководитель _____ Полицына Е.В.

Оценка _____ Дата защиты «_____» декабря 2020 г.

Москва 2020

МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РОССИЙСКОЙ
ФЕДЕРАЦИИ

Федеральное государственное бюджетное образовательное учреждение
высшего образования
«МОСКОВСКИЙ АВИАЦИОННЫЙ ИНСТИТУТ
(НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ)»

Кафедра 319 «Системы интеллектуального мониторинга»

З А Д А Н И Е

на курсовую работу по дисциплине

Технология разработки программного обеспечения

Студент МЗО – 221М – 19 Пожидаев Егор Витальевич
(№ группы, Ф. И. О.)

Тема Проектирование и разработка веб-приложения классификации
новостей с применением методов машинного обучения

Перечень вопросов, подлежащих разработке в курсовой работе:

1. Проектирование архитектуры системы и выбор средств разработки
 2. Проектирование и реализация ядра системы – классификатора
 3. Проектирование и реализация фронтенда
 4. Проектирование и реализация бекенда
-

Рекомендуемая литература

1. Django [Электронный ресурс]. – Режим доступа: <https://docs.djangoproject.com/>. – Заглавие с экрана. – (Дата обращения 25.12.20).
 2. Scikit-learn [Электронный ресурс]. – Режим доступа: <https://scikit-learn.org/>. – Заглавие с экрана. – (Дата обращения 25.12.20).
-

Задание выдано « 12 » сентября 2020 г.

Руководитель Полицына Екатерина Валерьевна, к.т.н., доцент
кафедры 319 МАИ

(Ф. И. О., должность, подпись)

Студент Пожидаев Е.В.

(подпись)

СОДЕРЖАНИЕ

1	ОПИСАНИЕ ВОЗМОЖНОСТЕЙ ПРИЛОЖЕНИЯ.....	4
1.1	Общие сведения	4
1.2	Назначение и цель	4
2	ТРЕБОВАНИЯ К ПРИЛОЖЕНИЮ.....	5
2.1	Функциональные требования.....	5
2.2	Нефункциональные требования	5
3	АРХИТЕКТУРА СИСТЕМЫ	6
3.1	Схема архитектуры.....	6
3.2	Протоколы взаимодействия.....	6
3.3	Используемые технологии.....	13
3.4	Паттерны проектирования	14
4	ОПИСАНИЕ КЛАССИФИКАТОРА.....	15
4.1	Классы	15
4.2	Вектор признаков	15
4.3	Модель машинного обучения классификатора	15
4.4	Обучающая и тестовая коллекция документов.....	15
4.5	Оценка точности классификации	16
5	ОПИСАНИЕ ИНФРАСТРУКТУРЫ РАЗРАБОТКИ.....	17
5.1	Система контроля версий.....	17
5.2	Сборка и развертывание приложения	17
5.3	Обновление модели классификатора	17
6	РЕЗУЛЬТАТЫ ТЕСТИРОВАНИЯ.....	18
7	ИНТЕРФЕЙС И ВОЗМОЖНОСТИ СИСТЕМЫ.....	19
8	АНАЛИЗ РЕЗУЛЬТАТОВ	23
	СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	24

1 ОПИСАНИЕ ВОЗМОЖНОСТЕЙ ПРИЛОЖЕНИЯ

1.1 Общие сведения

Веб-приложение классификации новостных статей представляет веб-сервис, который позволяет классифицировать новости и управлять коллекцией собранных новостей с информационного портала «interfax.ru».

1.2 Назначение и цель

С помощью веб-приложения возможно определить, к какой категории относится новостная статья, введенная пользователем. Сервис предоставляет возможность просмотреть собранные новостные статьи, которые представляли собой входные данные для обучения классификатора с возможностью их фильтрации. Также веб-приложение позволяет удалять, изменять собранные новостные статьи. У веб-приложения есть возможность добавить новостную статью в коллекцию.

2 ТРЕБОВАНИЯ К ПРИЛОЖЕНИЮ

2.1 Функциональные требования

Функциональные требования к системе:

- добавление новостной статьи;
- удаление новостной статьи;
- изменение новостной статьи;
- вывод новостных статей по страницам;
- фильтрация новостных статей по категориям и диапазону даты публикации;
- классификация текста по категории новостной статьи с помощью методов машинного обучения;
- обновление модели классификатора.

2.2 Нефункциональные требования

Система должна быть реализована на базе клиент-серверной архитектуры и представлять из себя веб-сервис для классификации новостных статей с графическим пользовательским интерфейсом.

Веб-приложение должно быть разработано с помощью следующих технических средств:

- backend должен быть реализован на языке программирования Python с помощью фреймворка Django;
- СУБД SQLite;
- frontend должен быть реализован с помощью HTML и языка программирования Javascript;
- проект и вся необходимая документация и данные по нему должны храниться в системе контроля версий git на сайте сервиса GitHub.

Система должна быть отказоустойчивой при работе пользователя с системой.

3 АРХИТЕКТУРА СИСТЕМЫ

3.1 Схема архитектуры

Архитектура разработанного приложения изображена на рисунке 1.

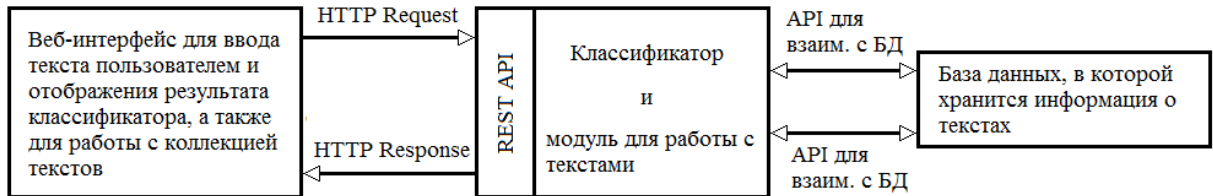


Рисунок 1 – Архитектура приложения

Frontend часть приложения, представляющая собой веб-интерфейс для ввода текста новости и работы с коллекцией текстов взаимодействует с backend частью приложения посредством HTTP запросов и ответов. Backend часть приложения состоит из базы данных, в которой хранится информация о новостных статьях, классификатора и модуля для работы с коллекцией новостей. Классификатор и модуль для работы с коллекцией новостей взаимодействует с базой данных с помощью API для взаимодействия с БД.

3.2 Протоколы взаимодействия

Разработанная система является RESTful веб-сервисом. Используемые в системе API:

- Страница классификатора.

GET /

Возвращает html-документ страницы классификатора.

- Страница для работы с коллекцией текстов.

GET /collection

Возвращает html-документ страницы для работы с коллекцией текстов.

- Страница новостной статьи.

GET /article/{articleId}

Возвращает html-документ страницы новостной статьи.

- Фильтрация новостных статей.

POST /api/articles/filter

Данная точка входа фильтрует новостные статьи по категориям и по диапазону даты публикации, при этом позволяет отсортировать полученную выборку по данным полям, а также провести выборку по страницам, с возможностью задания количества элементов на них. Возвращает в теле ответа json, который содержит информацию о новостных статьях, а также их количестве. Json-схема json ответа фильтрации новостных статей представлена в листинге 1. Точка входа получает в теле запроса json, который содержит информацию о параметрах для фильтрации. Json-схема json запроса представлена в листинге 2.

Листинг 1 – Json-схема при ответе точки входа для фильтрации новостных статей.

```
{
  "$schema": "http://json-schema.org/draft-04/schema#",
  "title": "filterResult",
  "description": "information about filtered articles",

  "properties":{
    "len":{
      "type": "integer"
    },
    "articles":{
      "type": "array",
      "items":{
        "type": "object",
        "properties":{
          "id":{
            "type": "integer"
          },
          "title":{
            "type": "string"
          },
          "category":{
            "type": "string"
          },
          "datetime":{
            "type": "string"
          },
          "description":{
            "type": "string"
          }
        }
      }
    }
  }
}
```


информацию о добавляемой статье. Json-схема json запроса представлена в листинге 3.

Листинг 3 – Json-схема при запросе точки входа для добавления новостной статьи.

```
{
  "$schema": "http://json-schema.org/draft-04/schema#",
  "title": "articleInformationForAdd",
  "description": "information about article for add to
collection",
  "properties":{
    "source":{
      "type": "string"
    },
    "category":{
      "type": "string"
    },
    "datetime":{
      "type": "string"
    },
    "title":{
      "type": "string"
    },
    "description":{
      "type": ["string", "null"]
    },
    "text":{
      "type": "string"
    },
    "tags":{
      "type": ["array", "null"],
      "items":{
        "type": "string"
      }
    }
  }
}
```

- Вывод новостной статьи.

GET /api/articles/{articleId}

Возвращает в теле ответа json, который содержит информацию о новостной статье. Json-схема json овета представлена в листинге 4.

Листинг 4 – Json-схема при ответе точки входа для вывода новостной статьи.

```
{
  "$schema": "http://json-schema.org/draft-04/schema#",
```

```

    "title": "articleInformationForGet",
    "description": "information about output article",

    "properties":{
        "source":{
            "type": "string"
        },
        "category":{
            "type": "string"
        },
        "datetime":{
            "type": "string"
        },
        "title":{
            "type": "string"
        },
        "description":{
            "type": "string"
        },
        "text":{
            "type": "string"
        },
        "tags":{
            "type": "array"
            "items":{
                "type": "string"
            }
        }
    }
}

```

- **Изменение новостной статьи.**

PUT /api/articles/{articleId}

Возвращает html ответ со статусом 200 при успешном изменении новостной статьи. Точка входа в теле запроса получает json, который содержит информацию об измененной новостной статье. Json-схема json запроса представлена в листинге 5.

Листинг 5 – Json-схема при ответе точки входа для изменения новостной статьи.

```

{
    "$schema": "http://json-schema.org/draft-04/schema#",
    "title": "articleInformationForUpdate",
    "description": "information about article for update",

    "properties":{
        "source":{
            "type": "string"
        }
    }
}

```

```

    },
    "category":{
        "type": "string"
    },
    "datetime":{
        "type": "string"
    },
    "title":{
        "type": "string"
    },
    "description":{
        "type": ["string", "null"]
    },
    "text":{
        "type": "string"
    },
    "tags":{
        "type": ["array", "null"],
        "items":{
            "type": "string"
        }
    }
}
}

```

- Удаление новостной статьи.

DELETE /api/articles/{articleId}

Возвращает html ответ со статусом 200 при успешном удалении новостной статьи.

- Получение информации о классификаторе.

GET /api/classifier

Возвращает в теле ответа json, который содержит информацию о классификаторе. Json-схема json ответа представлена в листинге 6.

Листинг 6 – Json-схема при ответе точки входа для получения информации о классификаторе.

```

{
    "$schema": "http://json-schema.org/draft-04/schema#",
    "title": "informationClassifier",
    "description": "information about accuracy of classifier",

    "properties":{
        "precision":{
            "type": "array",
            "items":{

```

```

        "type": "number"
    },
    },
    "precisionMacro":{
        "type": "number"
    },
    "precisionMicro":{
        "type": "number"
    },
    "recall":{
        "type": "array",
        "items":{
            "type": "number"
        }
    },
    "recallMacro":{
        "type": "number"
    },
    "recallMicro":{
        "type": "number"
    },
    "f1Score":{
        "type": "array",
        "items":{
            "type": "number"
        }
    },
    "f1ScoreMacro":{
        "type": "number"
    },
    "f1ScoreMicro":{
        "type": "number"
    },
    "accuracy":{
        "type": "number"
    },
    "categoryList":{
        "type": "array",
        "items":{
            "type": "string"
        }
    }
}
}

```

- Классификация новостной статьи.

POST /api/classifier

Возвращает в теле ответа json, который содержит категорию – результат работы классификатора. Данная точка входа в теле запроса получает json, который содержит текст классифицируемой новостной статьи.

Json-схема json ответа представлена в листинге 7. Json-схема json запроса представлена в листинге 8.

Листинг 7 – Json-схема при ответе точки входа для классификации новостной статьи.

```
{
  "$schema": "http://json-schema.org/draft-04/schema#",
  "title": "classifierResult",
  "description": "Result of classifier",

  "properties": {
    "category": {
      "type": "string"
    }
  }
}
```

Листинг 8 – Json-схема при запросе точки входа для классификации новостной статьи.

```
{
  "$schema": "http://json-schema.org/draft-04/schema#",
  "title": "TextForclassifier",

  "properties": {
    "text": {
      "type": "string"
    }
  }
}
```

- Обновление модели классификатора.

PUT /api/classifier/fit

Возвращает html ответ со статусом 200 при успешном обновлении модели классификатора.

Для взаимодействия с базой данных используется ORM, которая является частью фреймворка Django.

3.3 Используемые технологии

Для реализации серверной части использовался язык программирования Python, а также фреймворк Django.

В качестве СУБД используется SQLite.

Для создания модели классификатора используется модуль языка Python scikit-learn.

В качестве средства реализации пользовательского интерфейса используется HTML и Javascript.

В качестве системы контроля версий использовался сервис GitHub, представляющий из себя графический интерфейс для технологии git – распределенной системы управления версиями.

3.4 Паттерны проектирования

При разработке использовался паттерн MVC, который позволяет разрабатывать бизнес-логику и визуальное представление отдельно.

4 ОПИСАНИЕ КЛАССИФИКАТОРА

4.1 Классы

Классификатор определяет следующие классы:

- в мире;
- в России;
- Москва;
- культура;
- спорт;
- экономика.

4.2 Вектор признаков

В качестве модели представления текста новостных статей использовался метод bag-of-words. При этом проводилась начальная подготовка текста: удалялись лишние неинформативные символы, а также различные стоп-слова. Вектором признаков является масштабированное отношение количества вхождений слова среди всех новостей на количество слов в тексте по методу tf-idf.

4.3 Модель машинного обучения классификатора

В качестве модели машинного обучения был выбран метод опорных векторов.

4.4 Обучающая и тестовая коллекция документов

В качестве обучающей коллекции были взяты 700 новостных статей по каждой из категории из общей коллекции текстов. Таким образом, классификатор обучался на 4200 новостных статьях.

В качестве тестовой коллекции были взяты 300 новостных статей по каждой из категории из общей коллекции текстов. Таким образом классификатор тестировался на 1800 новостных статьях.

4.5 Оценка точности классификации

Оценка точности разработанного классификатора оценивалась по нескольким критериям: точности (precision), полноте (recall), F-мере и точности (accuracy).

Точность (accuracy) классификации равна 83%. Оценки точности по категориям представлены в таблице 1.

Таблица 1 – Оценка точности обученного классификатора по категориям

Категория	Точность (precision)	Полнота (recall)	F-мера
В мире	0.79	0.8	0.8
В России	0.79	0.50	0.61
Москва	0.83	0.83	0.83
Культура	0.86	0.94	0.9
Спорт	0.96	0.97	0.97
Экономика	0.75	0.92	0.82
Взвешенное среднее	0.83	0.83	0.83

5 ОПИСАНИЕ ИНФРАСТРУКТУРЫ РАЗРАБОТКИ

5.1 Система контроля версий

Проектирование и реализация приложения велась локально на компьютере разработчика при помощи Microsoft Visual Studio 2017. При разработке системы использовалась система контроля версий git. Исходный код хранится локально на компьютере разработчика, и удаленно, на сайте сервиса GitHub: <https://github.com/XaZdarova/Classifier>.

5.2 Сборка и развертывание приложения

Для успешного запуска веб-приложения необходимо наличие интерпретатора Python 3.0 (64 бит) или выше, а также операционной системы Windows 7 и выше. Необходимо скачать архив с проектом, разархивировать папку ClassifierProject. Далее через командную строку активировать виртуальное окружение python с помощью команд:

- `cd {pathTo}\ClassifierProject;`
- `env\Scripts\activate.bat.`

После активации окружения, достаточно запустить web-приложение по данной команде: `manage.py runserver`.

5.3 Обновление модели классификатора

Для обновления модели классификатора, достаточно послать POST запрос по ссылке `/api/classifier/fit`.

6 РЕЗУЛЬТАТЫ ТЕСТИРОВАНИЯ

Для тестирования системы применялось ручное тестирование. Тестирование проводилось методом черного ящика.

В процессе тестирования использовался веб-браузер Google Chrome версии. Основные тестовые сценарии и результаты по ним представлены в таблице 2.

Таблица 2 – Основные тестовые сценарии и результаты тестирования

Сценарий тестирования	Результат тестирования
Ввод текста новостной статьи для классификации и последующая его классификация	Был выведен корректный результат классификации текста
Ввод пустого текста новостной статьи для классификации и последующая его классификация	Было выведено сообщение о том, что был введен пустой текст
Добавление новой новостной статьи	Было выведено сообщение о том, что статья успешно добавлена.
Добавление новой новостной статьи с некорректными параметрами	Было выведено сообщение о том, что параметры некорректны
Изменение новостной статьи	Было выведено сообщение о том, что статья успешно изменена
Изменение новостной статьи с некорректными параметрами	Было выведено сообщение о том, что параметры некорректны
Удаление новостной статьи	Было выведено сообщение о том, что статья удалена
Вывод страницы с новостными статьями	При загрузке страницы была выведена первая страница с коллекцией новостей
Вывод новостной статьи	При загрузке страницы с новостной статье, информация по ней была успешно отображена
Вывод несуществующей новостной статьи	Было выведено сообщение о том, что такой статьи не существует

7 ИНТЕРФЕЙС И ВОЗМОЖНОСТИ СИСТЕМЫ

Система предоставляет пользователю классифицировать текст и просматривать коллекцию новостных статей, а также производить операции с отдельными новостями.

Классификация текста происходит при нажатии кнопки «Классифицировать». Интерфейс страницы классификации текстов представлен на рисунке 2.

Классификация текстов

- [Классификатор](#)
- [Коллекция текстов](#)

Ввод текста



Рисунок 2 – Интерфейс классификации новостной статьи

Для фильтрации коллекции текстов необходимо задать желаемые параметры и нажать кнопку «Найти». Интерфейс для работы с коллекцией новостных статей представлен на рисунке 3.

Коллекция текстов

- [Классификатор](#)
- [Коллекция текстов](#)

Фильтр

Категория

Левая граница даты публикации

Правая граница даты публикации

Сортировка

Тип сортировки

Кол-во элементов на странице

Рисунок 3 – Интерфейс для работы с коллекцией новостей

Найденные новостные статьи можно просматривать постранично, переключаясь между страницами с помощью кнопок «Предыдущая» и «Следующая», а также переходить по предоставленным ссылкам для получения более детальной информации, изменения или удаления статьи. Для изменения используется кнопка «Изменить новостную статью», а для удаления «Удалить новостную статью». Интерфейс для отображения коллекции новостей представлен на рисунке 4. Интерфейс для работы с отдельной новостью представлен на рисунке 5.

[Новостная статья №647941](#)

Заголовок: Третьяковскую галерею закрыли для посетителей после кражи картины Куинджи

Дата публикации: 2019-01-27 20:49:00

Теги: третьяковская галерея

[Новостная статья №647943](#)

Заголовок: Минкультуры экстренно внесло картину Куинджи "Ай-Петри. Крым" в реестр похищенных

Дата публикации: 2019-01-27 21:10:00

Теги: третьяковская галерея

[Новостная статья №647945](#)

Заголовок: Составлен фоторобот похитителя картины Куинджи

Дата публикации: 2019-01-27 21:30:00

Теги: третьяковская галерея

1 [Предыдущая](#) [Следующая](#) 1732

Рисунок 4 – Интерфейс для отображения коллекции новостей

Новостная статья

- [Классификатор](#)
- [Коллекция текстов](#)

№ 647941

Ссылка на источник:

<https://www.interfax.ru/cu>

Категория:

Культура ▼

Дата публикации:

2019-01-27 20:49

Заголовок:

Третьяковскую галерею

Описание:

Текст:

Москва. 27 января. INTERFAX.RU - Главный вход в Третьяковскую галерею закрыт для посетителей после того, как из музея была похищена картина художника Архипа Куинджи "Ай-Петри. Крым". Как передает корреспондент "Интерфакса", у здания Третьяковской галереи в Лаврушинском переулке в центре Москвы стоят несколько полицейских автомобилей. Также здесь собрались журналисты различных изданий. Вход в музей закрыт, за дверью выставлено небольшое ограждение. У входа стоят люди в штатском, не пропускающие людей внутрь. В холле здания видны все еще оставшиеся посетители галереи, которые фотографируются на фоне пресс-волла, посвященного выставке Куинджи. В здании дежурят несколько полицейских. По словам вышедших из здания посетителей, перед тем, как выпустить из Третьяковской галереи, сотрудники правоохранительных органов просили их показать содержимое сумок. Со своей стороны, источник "Интерфакса" сообщил, что в данный момент на месте работает следственно-оперативная группа, опрашивают очевидцев похищения картины, а также изымаются записи видеонаблюдения. Ведется работа по установлению личностей подозреваемых. Ранее информированный источник сообщил "Интерфаксу", что картина Архипа Куинджи исчезла из зала Третьяковской галереи в Лаврушинском переулке, где проходит выставка художника. Она пропала на глазах у многочисленных посетителей экспозиции. Мужчина, которого приняли за работника музея, спокойно подошел к полотну, достал его из рамы и ушел.

Теги:

третьяковская галерея

[Удалить новостную статью](#)

[Изменить новостную статью](#)

Рисунок 5 – Интерфейс для работы с отдельной новостью

Для добавления новой статьи необходимо ввести данные о ней в форму, затем нажать на кнопку «Добавить статью». Интерфейс для добавления новой статьи представлен на рисунке 6.

Добавить новостную статью

Ссылка на статью

Категория

Дата публикации

Заголовок

Описание

Текст статьи

Теги

Рисунок 6 – Интерфейс для добавления новой статьи

8 АНАЛИЗ РЕЗУЛЬТАТОВ

Таким образом, был спроектировано и реализовано web-приложение для классификации новостных статей и работы с их коллекцией.

Разработанное web-приложение не требует больших затрат на поддержку, сопровождение и обновление.

К приложению легко добавить новые web-страницы и функционал, однако при изменении уже разработанного функционала требуется приложить много усилий. Также при изменении структуры новостной статьи требуется изменить большое количество программных единиц.

В плане отказоустойчивости, web-приложение защищено от ввода некорректных данных введенных пользователем. Однако при изменении инфраструктуры приложения, а также различных метафайлов возможно нарушение работы приложения.

Классификатор имеют неплохую точность распознавания категорий, а именно 83% процента. Однако надо в дальнейшем повышать данный результат, за счет лучшей подготовки текста. Также у классификатора имеется ещё одна проблема, которая заключается в низкой точности определения категории «В России», что также надо исправлять.

Также frontend разрабатывался в большей степени для браузера Chrome, поэтому могут возникнуть проблемы при использовании приложения на других браузерах.

Необходимо значительно переработать дизайн интерфейса для повышения удобства эксплуатации.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Interfax [Электронный ресурс]. – Режим доступа: <https://www.interfax.ru/>. – Заглавие с экрана. – (Дата обращения: 25.12.20).
2. Python [Электронный ресурс]. – Режим доступа: <https://www.python.org/>. – Заглавие с экрана. – (Дата обращения: 25.12.20).
3. Django [Электронный ресурс]. – Режим доступа: <https://docs.djangoproject.com/>. – Заглавие с экрана. – (Дата обращения 25.12.20).
4. Scikit-learn [Электронный ресурс]. – Режим доступа: <https://scikit-learn.org/>. – Заглавие с экрана. – (Дата обращения 25.12.20).
5. Htmlbook [Электронный ресурс]. – Режим доступа: <http://htmlbook.ru/>. – Заглавие с экрана. – (Дата обращения 25.12.20).
6. Современный учебник javascript [Электронный ресурс]. – Режим доступа: <https://learn.javascript.ru/>. – Заглавие с экрана. – (Дата обращения 25.12.20).