

# Контрольная работа по теме "Классификация текста"

Выполнил: Пакало Александр Сергеевич ИУ5-22М

## Задание

Необходимо решить задачу классификации текстов на основе любого выбранного Вами датасета (кроме примера, который рассматривался в лекции). Классификация может быть бинарной или многоклассовой. Целевой признак из выбранного Вами датасета может иметь любой физический смысл, примером является задача анализа тональности текста.

Необходимо сформировать два варианта векторизации признаков - на основе CountVectorizer и на основе TfidfVectorizer.

В качестве классификаторов необходимо использовать два классификатора по варианту для Вашей группы: RandomForestClassifier, LogisticRegression

```
In [1]: # Loading extension for reloading editable packages (pip install -e .)
%load_ext autoreload
```

```
In [2]: # Reloading editable packages.
%autoreload
from charts.main import get_metrics_grouped_bar_chart
```

```
In [3]: RANDOM_SEED = 13
```

## Набор данных

Проведём классификацию текста используя набор данных [BBC News Archive \(https://www.kaggle.com/datasets/hgultekin/bbcnewsarchive\)](https://www.kaggle.com/datasets/hgultekin/bbcnewsarchive)

Подготовка переменных для работы с данными

```
In [4]: from pathlib import Path

data_path = Path("../..data")
external_data_path = data_path / "external"
raw_data_path = data_path / "raw"

dataset_filename = "bbc-news-data.zip"
```

Разархивирование набора данных

```
In [5]: import os
import shutil

raw_data_path.mkdir(exist_ok=True)

file_path = external_data_path / dataset_filename
raw_data_path = external_data_path / dataset_filename

if not os.path.isfile(raw_data_path):
    shutil.unpack_archive(file_path, extract_dir=raw_data_path)
    # file_path.unlink() # Remove archive after extracting it.
```

### Загрузка данных из csv

```
In [6]: import pandas as pd

df = pd.read_csv(raw_data_path, sep="\t")
```

## Разведочный анализ данных

Ознакомимся немного с данными, с которыми собираемся работать

Основные характеристики датасета

```
In [7]: df.head()
```

Out [7]:

	category	filename	title	content
0	business	001.txt	Ad sales boost Time Warner profit	Quarterly profits at US media giant TimeWarne...
1	business	002.txt	Dollar gains on Greenspan speech	The dollar has hit its highest level against ...
2	business	003.txt	Yukos unit buyer faces loan claim	The owners of embattled Russian oil giant Yuk...
3	business	004.txt	High fuel prices hit BA's profits	British Airways has blamed high fuel prices f...
4	business	005.txt	Pernod takeover talk lifts Domecq	Shares in UK drinks and food firm Allied Dome...

```
In [8]: df.tail()
```

Out [8]:

	category	filename	title	content
2220	tech	397.txt	BT program to beat dialler scams	BT is introducing two initiatives to help bea...
2221	tech	398.txt	Spam e-mails tempt net shoppers	Computer users across the world continue to i...
2222	tech	399.txt	Be careful how you code	A new European directive could put software w...
2223	tech	400.txt	US cyber security chief resigns	The man making sure US computer networks are ...
2224	tech	401.txt	Losing yourself in online gaming	Online role playing games are time- consuming,...

### Размер датасета

```
In [9]: num_of_rows, num_of_columns = df.shape  
print(f'Размер датасета: {num_of_rows} строка, {num_of_columns} колонок')
```

Размер датасета: 2225 строка, 4 колонок

### Определение типов

```
In [10]: df.dtypes
```

Out [10]: category      object  
filename      object  
title      object  
content      object  
dtype: object

### Проверка на наличие пустых значений

```
In [11]: df.isnull().sum()
```

Out [11]: category      0  
filename      0  
title      0  
content      0  
dtype: int64

### Обработки пустых значений не требуется

### Проверка на уникальные значения

```
In [12]: pd.Series(df["category"].unique())
```

```
Out[12]: 0      business
1  entertainment
2      politics
3         sport
4          tech
dtype: object
```

## Подготовка корпуса

Некоторые колонки имеют неверные типы данных, их следует преобразовать.

Строки вместо `object` сделаем `string`, а колонку "category" сделаем типа `category`.

```
In [13]: df = df.astype({
          "category": "category",
          "filename": "string",
          "title": "string",
          "content": "string",
        })
df.dtypes
```

```
Out[13]: category      category
filename  string[python]
title     string[python]
content   string[python]
dtype: object
```

## Токенизация

Загрузка модели spacy.

```
In [14]: import spacy

spacy_prefer_gpu = spacy.prefer_gpu()
nlp = spacy.load("en_core_web_sm")
```

Токенизация текстовых значений набора данных (кроме названия файла)

```
In [15]: from spacy.tokens.doc import Doc

corpus: list[Doc] = []

for text in (df["title"] + df["content"]).values:
    corpus.append(nlp(text))

corpus[:3]
```

```
Out[15]: [Ad sales boost Time Warner profit Quarterly profits at US media giant TimeWarner jumped 76% to $1.13bn (£600m) for the three months to December, from $639m year-earlier. The firm, which is now one of the biggest investors in Google, benefited from sales of high-speed internet connections and higher advert sales. TimeWarner said fourth quarter sales rose 2% to $11.1bn from $10.9bn. Its profits were buoyed by one-off gains which offset a profit dip at Warner Bros, and less users for AOL. Time Warner said on Friday that it now owns 8% of search-engine Google. But its own internet business, AOL, has mixed fortunes. It lost 464,000 subscribers in the fourth quarter profits were lower than in the preceding three quarters. However, the company said AOL's underlying profit before exceptional items rose 8% on the back of stronger internet advertising revenues. It hopes to increase subscribers by offering the online service free to TimeWarner internet customers and will try to sign up AOL's existing customers for high-speed broadband. TimeWarner also has to restate 2000 and 2003 results following a probe by the US Securities Exchange Commission (SEC), which is close to concluding. Time Warner's fourth quarter profits were slightly better than analysts' expectations. But its film division saw profits slump 27% to $284m, helped by box-office flops Alexander and Catwoman, a sharp contrast to year-earlier, when the third and final film in the Lord of the Rings trilogy boosted results. For the full-year, TimeWarner posted a profit of $3.36bn, up 27% from its 2003 performance, while revenues grew 6.4% to $42.09bn. "Our financial performance was strong, meeting or exceeding all of our full-year objectives and greatly enhancing our flexibility," chairman and chief executive Richard Parsons said. For 2005, TimeWarner is projecting operating earnings growth of around 5%, and also expects higher revenue and wider profit margins. TimeWarner is to restate its accounts as part of efforts to resolve an inquiry into AOL by US market regulators. It has already offered to pay $300m to settle charges, in a deal that is under review by the SEC. The company said it was unable to estimate the amount it needed to set aside for legal reserves, which it previously set at $500m. It intends to adjust the way it accounts for a deal with German music publisher Bertelsmann's purchase of a stake in AOL Europe, which it had reported as advertising revenue. It will now book the sale of its stake in AOL Europe as a loss on the value of that stake. ,
```

Dollar gains on Greenspan speech The dollar has hit its highest level against the euro in almost three months after the Federal Reserve head said the US trade deficit is set to stabilise. And Alan Greenspan highlighted the US government's willingness to curb spending and rising household savings as factors which may help to reduce it. In late trading in New York, the dollar reached \$1.2871 against the euro, from \$1.2974 on Thursday. Market concerns about the deficit has hit the greenback in recent months. On Friday, Federal Reserve chairman Mr Greenspan's speech in London ahead of the meeting of G7 finance ministers sent the dollar higher after it had earlier tumbled on the back of worse-than-expected US jobs data. "I think the chairman's taking a much more sanguine view on the current account deficit than

he's taken for some time," said Robert Sinche, head of currency strategy at Bank of America in New York. "He's taking a longer-term view, laying out a set of conditions under which the current account deficit can improve this year and next." Worries about the deficit concerns about China do, however, remain. China's currency remains pegged to the dollar and the US currency's sharp falls in recent months have therefore made Chinese export prices highly competitive. But calls for a shift in Beijing's policy have fallen on deaf ears, despite recent comments in a major Chinese newspaper that the "time is ripe" for a loosening of the peg. The G7 meeting is thought unlikely to produce any meaningful movement in Chinese policy. In the meantime, the US Federal Reserve's decision on 2 February to boost interest rates by a quarter of a point - the sixth such move in as many months - has opened up a differential with European rates. The half-point window, some believe, could be enough to keep US assets looking more attractive, and could help prop up the dollar. The recent falls have partly been the result of big budget deficits, as well as the US's yawning current account gap, both of which need to be funded by the buying of US bonds and assets by foreign firms and governments. The White House will announce its budget on Monday, and many commentators believe the deficit will remain at close to half a trillion dollars.

,  
Yukos unit buyer faces loan claim The owners of embattled Russian oil giant Yukos are to ask the buyer of its former production unit to pay back a \$900m (£479m) loan. State-owned Rosneft bought the Yugansk unit for \$9.3bn in a sale forced by Russia to partly settle a \$27.5bn tax claim against Yukos. Yukos' owner Menatep Group says it will ask Rosneft to repay a loan that Yugansk had secured on its assets. Rosneft already faces a similar \$540m repayment demand from foreign banks. Legal experts said Rosneft's purchase of Yugansk would include such obligations. "The pledged assets are with Rosneft, so it will have to pay real money to the creditors to avoid seizure of Yugansk assets," said Moscow-based US lawyer Jamie Firestone, who is not connected to the case. Menatep Group's managing director Tim Osborne told the Reuters news agency: "If they default, we will fight them where the rule of law exists under the international arbitration clauses of the credit." Rosneft officials were unavailable for comment. But the company has said it intends to take action against Menatep to recover some of the tax claims and debts owed by Yugansk. Yukos had filed for bankruptcy protection in a US court in an attempt to prevent the forced sale of its main production arm. The sale went ahead in December and Yugansk was sold to a little-known shell company which in turn was bought by Rosneft. Yukos claims its downfall was punishment for the political ambitions of its founder Mikhail Khodorkovsky and has vowed to sue any participant in the sale. ]

```
In [16]: assert len(corpus) == num_of_rows
```

## Классификация

### Подготовка данных для классификации

Заметим, что хоть `спрасу` при печати и выводит текст, на самом деле это объект. Наши модели ожидают увидеть строки.

```
In [17]: spacy_text = nlp('training: nlp!')
         spacy_text, type(spacy_text), type(spacy_text[0])
```

```
Out[17]: (training: nlp!, spacy.tokens.doc.Doc, spacy.tokens.token.Token)
```

```
In [18]: [token.text for token in spacy_text]
```

```
Out[18]: ['training', ':', 'nlp', '!']
```

Поэтому преобразуем наш corpus в упрощённый формат, совместимый с word2vec.

```
In [19]: str_corpus = [spacy_text.text for spacy_text in corpus]
```

Выберем X и y среди нашего набора данных

```
In [20]: X = str_corpus
         y = df["category"].values
```

Составим выборки для обучения

```
In [21]: from sklearn.model_selection import train_test_split

         X_train, X_test, y_train, y_test = train_test_split(
             X, y,
             test_size=0.33,
             random_state=RANDOM_SEED
         )
```

## Составим pipeline

```

In [22]: import numpy as np
from IPython.display import display
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, balanced_accuracy_score

def accuracy_score_for_classes(
    y_true: np.ndarray,
    y_pred: np.ndarray) -> dict[int, float]:
    """
    Вычисление метрики accuracy для каждого класса
    y_true - истинные значения классов
    y_pred - предсказанные значения классов
    Возвращает словарь: ключ - метка класса,
    значение - Accuracy для данного класса
    """
    # Для удобства фильтрации сформируем Pandas DataFrame
    d = {'t': y_true, 'p': y_pred}
    df = pd.DataFrame(data=d)
    # Метки классов
    classes = np.unique(y_true)
    # Результирующий словарь
    res = dict()
    # Перебор меток классов
    for c in classes:
        # отфильтруем данные, которые соответствуют
        # текущей метке класса в истинных значениях
        temp_data_flt = df[df['t']==c]
        # расчет accuracy для заданной метки класса
        temp_acc = accuracy_score(
            temp_data_flt['t'].values,
            temp_data_flt['p'].values)
        # сохранение результата в словарь
        res[c] = temp_acc
    return res

def print_accuracy_score_for_classes(
    y_true: np.ndarray,
    y_pred: np.ndarray):
    """
    Вывод метрики accuracy для каждого класса
    """
    accs = accuracy_score_for_classes(y_true, y_pred)
    results = pd.DataFrame(data={ "Категория": accs.keys(), "Точность":
    accs.values() })

    display(results)

    return results

```



```
In [23]: class EmbeddingVectorizer(object):
'''
    Для текста усредним вектора входящих в него слов
'''
    def __init__(self, model):
        self.model = model
        self.size = model.vector_size

    def fit(self, X, y):
        return self

    def transform(self, X):
        return np.array([np.mean(
            [self.model[w] for w in words if w in self.model]
            or [np.zeros(self.size)], axis=0)
            for words in X])
```

```
In [24]: from sklearn.pipeline import Pipeline

def classifier_pipeline(v, c, scaler=None):
    pipeline_steps = [
        ("vectorizer", v),
    ]

    if scaler:
        pipeline_steps.append(("scaler", scaler))

    pipeline_steps.append(("classifier", c))

    pipeline = Pipeline(pipeline_steps)

    classifier_X_train = X_train
    classifier_y_train = y_train
    classifier_X_test = X_test
    classifier_y_test = y_test

    pipeline.fit(classifier_X_train, classifier_y_train)
    y_pred = pipeline.predict(classifier_X_test)

    return print_accuracy_score_for_classes(classifier_y_test, y_pred)
```

## Проверка результатов

```
In [25]: ClassifierName = str
ModelName = str
CategoryName = str

# For seeing results of each category for each classifier.
MetricsDataPerCategory = dict[ClassifierName, dict[CategoryName, float]]
metrics_data_per_category: MetricsDataPerCategory = {}

# For seeing general results of each classifier for each model (category
# metrics are generalized then).
MetricsDataPerModel = dict[ClassifierName, dict[ModelName, float]]
metrics_data_per_model: MetricsDataPerModel = {}
```

```
In [26]: def add_metrics_data(classier_name: ClassifierName, model_name: ModelName):
    if not metrics_data_per_model.get(classier_name):
        metrics_data_per_model[classier_name] = {}

    metrics_data_per_model[classier_name][model_name] = np.mean(results_per_category[classier_name][model_name])

    results_per_category = dict(zip(results["Категория"], results["Точность"]))
    metrics_data_per_category[classier_name] = results_per_category

    return metrics_data_per_model, metrics_data_per_category
```

## CountVectorizer

Scaler не нужен (и его даже невозможно применить, ведь CountVectorizer возвращает разреженную матрицу)

```
In [27]: from sklearn.feature_extraction.text import CountVectorizer

count_vectorizer = CountVectorizer(ngram_range=(1, 3))
```

## LogisticRegression

```
In [28]: add_metrics_data("LogisticRegression", "count vectorizer", classifier_p
```

	Категория	Точность
0	business	0.968553
1	entertainment	0.912409
2	politics	0.942446
3	sport	0.993711
4	tech	0.921986

```
Out[28]: ({'LogisticRegression': {'count vectorizer': 0.9478209537671576}},
 {'LogisticRegression': {'business': 0.9685534591194969,
 'entertainment': 0.9124087591240876,
 'politics': 0.9424460431654677,
 'sport': 0.9937106918238994,
 'tech': 0.9219858156028369}})
```

## RandomForestClassifier

```
In [29]: from sklearn.ensemble import RandomForestClassifier

add_metrics_data("RandomForestClassifier", "count vectorizer", classifi
```

	Категория	Точность
0	business	0.968553
1	entertainment	0.824818
2	politics	0.856115
3	sport	0.993711
4	tech	0.836879

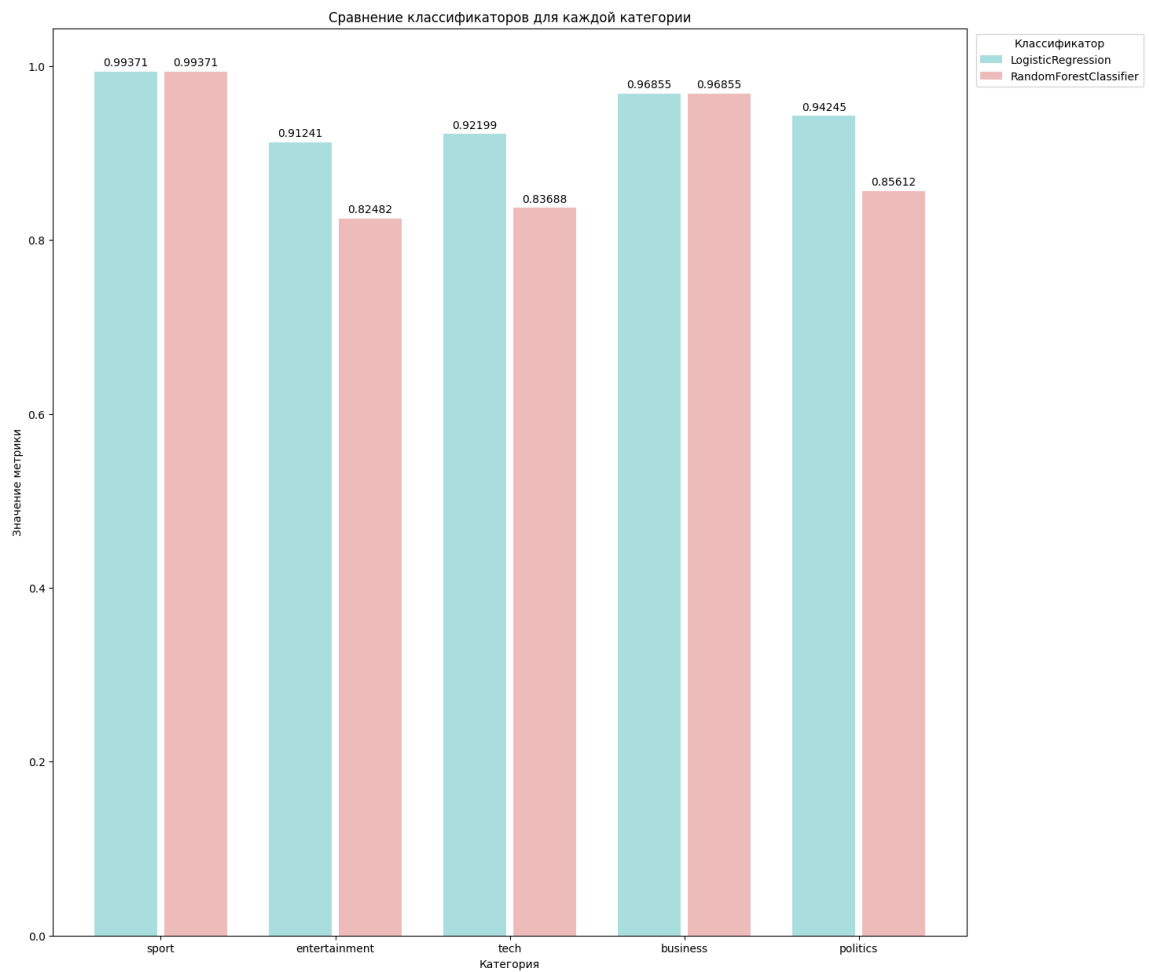
```
Out[29]: ({'LogisticRegression': {'count vectorizer': 0.9478209537671576},
  'RandomForestClassifier': {'count vectorizer': 0.896015241945870
7}},
{'LogisticRegression': {'business': 0.9685534591194969,
  'entertainment': 0.9124087591240876,
  'politics': 0.9424460431654677,
  'sport': 0.9937106918238994,
  'tech': 0.9219858156028369},
'RandomForestClassifier': {'business': 0.9685534591194969,
  'entertainment': 0.8248175182481752,
  'politics': 0.8561151079136691,
  'sport': 0.9937106918238994,
  'tech': 0.8368794326241135}})
```

### Сравнение результатов классификаторов для каждой категории

```
In [30]: def show_classifiers_bar_chart(metrics_data_per_category: MetricsDataPe
classifiers_bar_chart = get_metrics_grouped_bar_chart(metrics_data_
classifiers_bar_chart["plt"].title('Сравнение классификаторов для к
classifiers_bar_chart["plt"].xlabel('Категория')
classifiers_bar_chart["plt"].ylabel('Значение метрики')

classifiers_bar_chart["ax"].legend(title='Классификатор', bbox_to_a
classifiers_bar_chart["plt"].show()
```

```
In [31]: show_classifiers_bar_chart(metrics_data_per_category)
```



## TFIDF

Scaler не нужен (и его даже невозможно применить, ведь tfidf возвращает разреженную матрицу)

```
In [32]: from sklearn.feature_extraction.text import TfidfVectorizer  
tfidf = TfidfVectorizer(ngram_range=(1, 3))
```

## LogisticRegression

```
In [33]: add_metrics_data("LogisticRegression", "tfidf", classifier_pipeline(tfi
```

	Категория	Точность
0	business	0.955975
1	entertainment	0.919708
2	politics	0.971223
3	sport	0.993711
4	tech	0.992908

```
Out[33]: ({'LogisticRegression': {'count vectorizer': 0.9478209537671576,
  'tfidf': 0.9667048773578898},
  'RandomForestClassifier': {'count vectorizer': 0.896015241945870
7}},
{'LogisticRegression': {'business': 0.9559748427672956,
  'entertainment': 0.9197080291970803,
  'politics': 0.9712230215827338,
  'sport': 0.9937106918238994,
  'tech': 0.9929078014184397},
  'RandomForestClassifier': {'business': 0.9685534591194969,
  'entertainment': 0.8248175182481752,
  'politics': 0.8561151079136691,
  'sport': 0.9937106918238994,
  'tech': 0.8368794326241135}})
```

## RandomForestClassifier

```
In [34]: from sklearn.ensemble import RandomForestClassifier

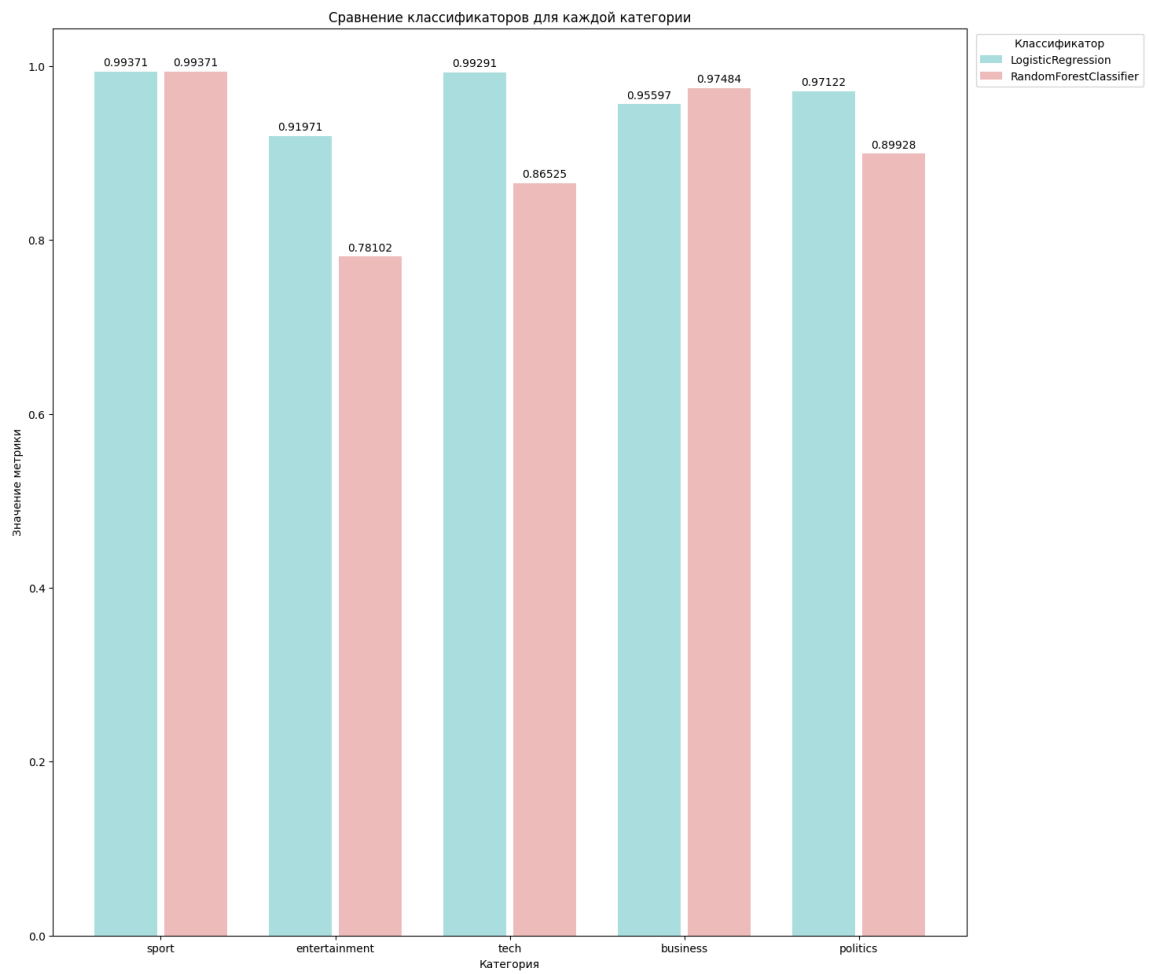
add_metrics_data("RandomForestClassifier", "tfidf", classifier_pipeline
```

	Категория	Точность
0	business	0.974843
1	entertainment	0.781022
2	politics	0.899281
3	sport	0.993711
4	tech	0.865248

```
Out[34]: ({'LogisticRegression': {'count vectorizer': 0.9478209537671576,
  'tfidf': 0.9667048773578898},
  'RandomForestClassifier': {'count vectorizer': 0.8960152419458707,
  'tfidf': 0.9028208318839278}},
 {'LogisticRegression': {'business': 0.9559748427672956,
  'entertainment': 0.9197080291970803,
  'politics': 0.9712230215827338,
  'sport': 0.9937106918238994,
  'tech': 0.9929078014184397},
  'RandomForestClassifier': {'business': 0.9748427672955975,
  'entertainment': 0.781021897810219,
  'politics': 0.8992805755395683,
  'sport': 0.9937106918238994,
  'tech': 0.8652482269503546}})
```

**Сравнение результатов классификаторов для каждой категории**

```
In [35]: show_classifiers_bar_chart(metrics_data_per_category)
```

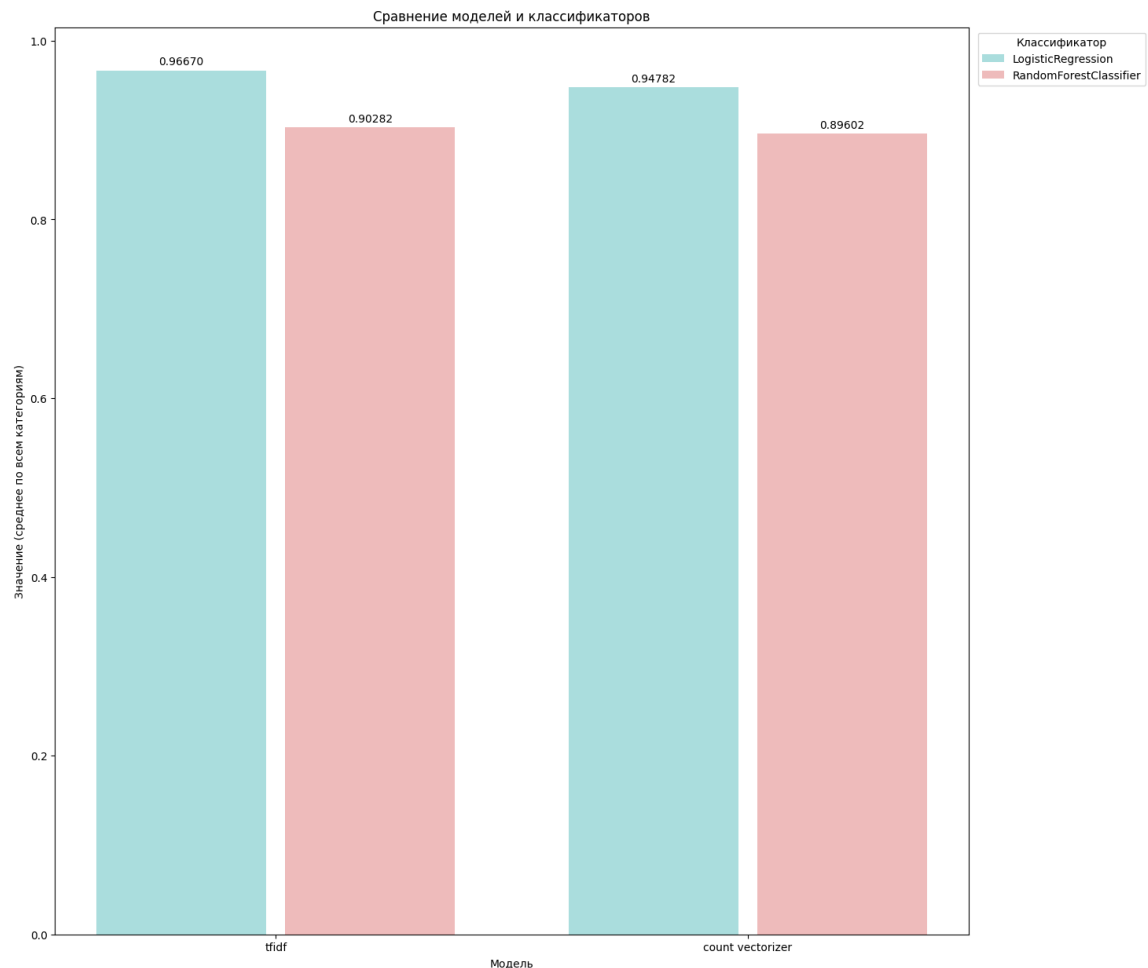


**Сравнение результатов моделей и классификаторов**

```
In [36]: models_bar_chart = get_metrics_grouped_bar_chart(metrics_data_per_model)
models_bar_chart["plt"].title('Сравнение моделей и классификаторов')
models_bar_chart["plt"].xlabel('Модель')
models_bar_chart["plt"].ylabel('Значение (среднее по всем категориям)')

models_bar_chart["ax"].legend(title='Классификатор', bbox_to_anchor=(1.

models_bar_chart["plt"].show()
```



## Вывод

Как видно по графику сравнения моделей и классификаторов, наиболее успешной оказалась связка TFIDF и LogisticRegression.

В общем, значения моделей с классификатором логистической регрессии выше на 5%, чем у соответствующих моделей с RandomForestClassifier.

При этом значения моделей при одинаковых классификаторах почти не отличаются: и TFIDF, и CountVectorizer демонстрируют хорошие результаты.