

## Определение классов

Классы:

1. Автомобили.
2. Технологии.
3. Политика.
4. Трейлер.
5. Игры.
6. Экономика.
7. Общество
8. Аниме.
9. Манга.
10. Интернет.
11. Наука.
12. Происшествия.
13. Спорт.
14. Культура.
15. Прочее.

### Выбранные методы векторизации

Для обучения и сравнения были выбраны следующие 2 метода векторизации:

- Bag of Words
- Word2Vec

### Вектор признаков для bag of words

Пусть:  $A$  – множество уникальных слов на русском языке (кроме служебных частей речи, слов на латинице и неизвестных слов) в рамках обучаемой выборки. Тогда:

$$A = \{a_1, a_2, \dots, a_n\},$$

где  $a_i$  –  $i$ -ое уникальное слово в нормальной форме;  
 $n$  – размер множества  $A$ .

Пусть:  $V$  – вектор признаков. Тогда:

$$V = \{v_1, v_2, \dots, v_n\},$$

где  $v_i$  – частота встречаемости слова  $a_i$  для входного текста.

## Способ векторизации текстов для bag of words

Шаги:

1. Вектор  $V$  для заданного текста автоматически заполняется нулями.
2. Из текста извлекаются нормальные формы слов  $B = \{b_1, b_2, \dots, b_m\}$ , где  $m$  – количество слов, которые не принадлежат служебным частям речи, написаны не на латинице и найдены в словаре морфологического анализатора.
3. Проверяется принадлежность каждого  $b_j$  множеству  $A$ .
4. Если  $b_j \in A$  и  $b_j = a_i, i \in 1..n$ , то  $v_i = v_i + 1$ .

## Выбранные модели для обучения и их сравнения

Для обучения были использованы следующие модели:

- LogisticRegression
- LinearSVC
- GaussianNB
- MultinomialNB (только Bag of Words)
- RandomForestClassifier
- LogisticRegressionCV

Оценки результатов изображены на таблицах 1-6 для Bag of Words и 7-11 для Word2Vec.

Для обучения был использован алгоритм MultiOutputClassifier, который позволяет обученной модели иметь несколько выходов. В случае нашей обучаемой модели – это массив  $C^{pred} = \{c_1^{pred}, c_2^{pred}, \dots, c_m^{pred}\}$ , где  $m$  – кол-во классов (категорий),  $c_j^{pred}, j = 1..m$  – значение 1 или 0 – принадлежит или не принадлежит статья к данной категории.

Таблица 1 – Оценка Bag of Words, модель LogisticRegression

<b>Accuracy</b>	0.559
<b>Precision</b>	0.883
<b>Recall</b>	0.636
<b>F1 score</b>	0.740
<b>Hamming loss</b>	0.035

Таблица 2 – Оценка Bag of Words, модель LinearSVC

<b>Accuracy</b>	0.601
<b>Precision</b>	0.891
<b>Recall</b>	0.669
<b>F1 score</b>	0.764
<b>Hamming loss</b>	0.032

Таблица 3 – Оценка Bag of Words, модель GaussianNB

<b>Accuracy</b>	0.286
<b>Precision</b>	0.898
<b>Recall</b>	0.321
<b>F1 score</b>	0.473
<b>Hamming loss</b>	0.056

Таблица 4 – Оценка Bag of Words, модель MultinomialNB

<b>Accuracy</b>	0.334
<b>Precision</b>	0.916
<b>Recall</b>	0.378
<b>F1 score</b>	0.535
<b>Hamming loss</b>	0.051

Таблица 5 – Оценка Bag of Words, модель RandomForestClassifier

<b>Accuracy</b>	0.296
<b>Precision</b>	0.944
<b>Recall</b>	0.346
<b>F1 score</b>	0.507
<b>Hamming loss</b>	0.052

Таблица 6 – Оценка Bag of Words, модель LogisticRegressionCV

<b>Accuracy</b>	0.569
<b>Precision</b>	0.883
<b>Recall</b>	0.646
<b>F1 score</b>	0.746
<b>Hamming loss</b>	0.034

Таблица 7 – Оценка Word2Vec, модель LogisticRegression

<b>Accuracy</b>	0.525
<b>Precision</b>	0.847
<b>Recall</b>	0.625
<b>F1 score</b>	0.719
<b>Hamming loss</b>	0.038

Таблица 8– Оценка Word2Vec, модель LinearSVC

<b>Accuracy</b>	0.548
<b>Precision</b>	0.861
<b>Recall</b>	0.644
<b>F1 score</b>	0.737
<b>Hamming loss</b>	0.036

Таблица 9 – Оценка Word2Vec, модель GaussianNB

<b>Accuracy</b>	0.043
<b>Precision</b>	0.290
<b>Recall</b>	0.931
<b>F1 score</b>	0.443
<b>Hamming loss</b>	0.182

Таблица 10 – Оценка Word2Vec, модель RandomForestClassifier

<b>Accuracy</b>	0.601
<b>Precision</b>	0.862
<b>Recall</b>	0.675
<b>F1 score</b>	0.757
<b>Hamming loss</b>	0.033

Таблица 11 – Оценка Word2Vec, модель LogisticRegressionCV

<b>Accuracy</b>	0.665
<b>Precision</b>	0.871
<b>Recall</b>	0.757
<b>F1 score</b>	0.810
<b>Hamming loss</b>	0.027

На основании полученных данных был выбран метод Word2Vec с моделью LogisticRegressionCV.