

BRLAB

Анализ тональности: как выделить главное из тысяч отзывов и мнений



BRLAB

Виталий Горбачев

Руководитель направления R&D

✉ vitaliy.g@brl.ru

✈ @imbarus

f [vitaly.gorbachev.5](https://www.facebook.com/vitaly.gorbachev.5)



РБ Рейтинг Букмекеров

ODDS.RU

ΜΣΤΑ RATINGS

План выступления

Что такое анализ тональности?

- Какой он бывает
- В чем сложности

Методы анализа

- Экспертные
- Обучение без учителя
- Обучение с учителем
- Инструменты

Наши результаты

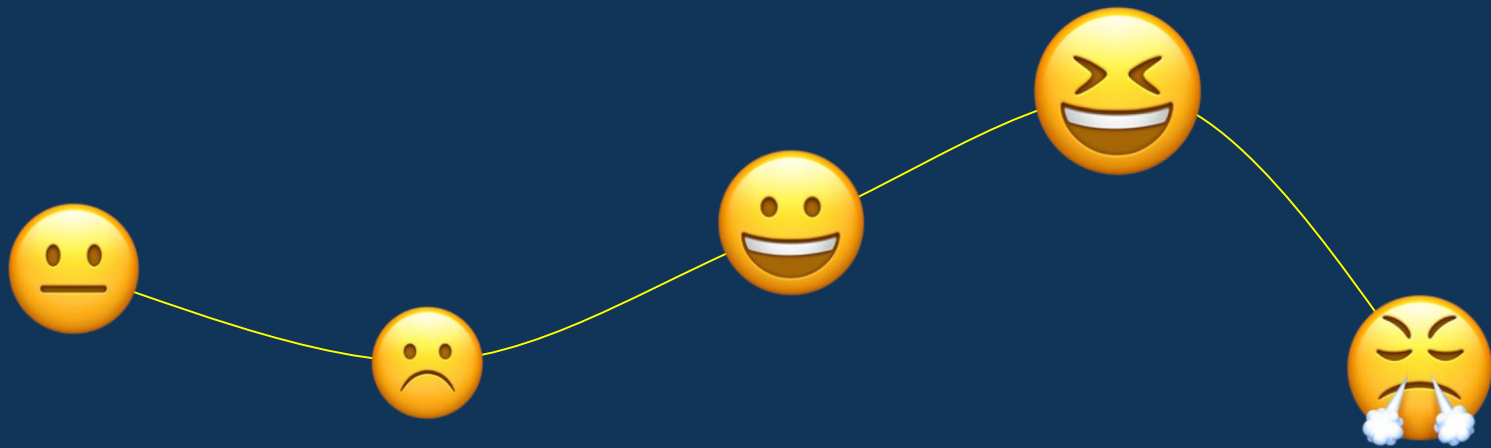
Как проанализировать тысячи отзывов и не свихнуться



BRLAB



Анализ тональности — это определение оценки, которую несёт текст или фраза.



Оценка текста или фраз

Бинарная

- Позитивн/негатив
- Для простой статистики

Многоклассовая

- Позитив/негатив/нейтрально/спам/etc

Числовая

- По любой шкале (0-1,0-100)
- Полезно в выявлении конфликтов

Тональность можно оценивать не по всему отзыву, а по характеристикам (аспектам)

Нормальный букмекер, все нормально и вводом денег и выводом. Тех поддержка ну 3,5, кэфы правда маленькие. А так все хорошо. Меньше года пользуюсь и не было проблем ещё.

Нормальный букмекер, все нормально и вводом денег и выводом. Тех поддержка ну 3,5, кэфы правда маленькие. А так все хорошо. Меньше года пользуюсь и не было проблем ещё.

Плюсы: **вывод и пополнение**
Минусы: **коэффициенты**
Нейтрально: **поддержка**

...все нормально?

... 3,5/5?

Почему?

Потому что русский язык:

- Требуется предобработка из-за сложности
- Очень метафоричен
- Мало готовых инструментов

Monkeylearn, aylien, etc...

The image displays three overlapping screenshots of the Monkeylearn web interface, illustrating the machine learning workflow for text classification.

Top Screenshot (Training Phase): Shows the 'Tag Data' screen under the 'TRAINING' tab. The instruction reads: "Choose one or more tags that apply and click confirm. As new texts appear, the model will learn from your criteria." The text being tagged is "Bad wif". The 'TAGS' panel on the right has the 'negative' tag selected. A progress bar at the bottom indicates "Tagged texts: 215 out of 305". A green 'Confirm (1)' button is visible.

Middle Screenshot (Testing Phase): Shows the 'Test with your own text' screen under the 'Run' tab. The text input field contains "Location, staff, good service, clean big rooms". The 'Results' table shows the following classification:

TAG	CONFIDENCE
staff	83.3%
location	82.2%
rooms	88.5%

A green '82.2%' is displayed above the results table. A blue 'Classify Text' button is at the bottom left. A note at the bottom right says: "Not the result you expected? Build more accuracy by training the model."

Bottom Screenshot (Training Phase): Shows another 'Tag Data' screen. The text is "We haven't been able to access our accounts since this morning, any updates on status?". The 'TAGS' panel has 'Accounts', 'Performance', and 'Updates' listed. A progress bar at the bottom indicates "Tagged texts: 0 out of 12". A green 'Confirm (0)' button is visible.

Почему?

Потому есть анафора:

“Из недостатков отнесу то, что с недавних пор часть трансляций перестала работать. Это не критично, но доставляет неудобства определенные. Надеюсь, это поправят скоро. А так приятная контора”

Почему?

Потому есть анафора:

“Первая интернет контора, где я начал играть. Раньше у них была еще классическая версия сайта, которую уже закрыли месяца 2 назад, но вот новая — просто бомба.”

Почему?

Потому что нет универсальности:

“...идите читать книгу!”:

Отзыв на книгу

Отзыв на фильм по книге

“... очень высокие”

Отзыв на букмекера

(коэффициенты)

Отзыв на ресторан (цены)

Почему?

Потому что есть контекст.

“Казалось бы, должна быть идеальная букмекерская контора. Хорошие коэффициенты, широкая линия, качественная поддержка -- как в рекламе. Но на деле всё рушится.”

Почему?

Потому что отзывы пишут люди

“играла команда хоккей : ТОРОНТО МАРЛИС -БЕЛЬВИЛЬ ПОСТАВИЛ НА П1 НА ВТОРОМ ТАЙМЕ НАЧИНАЕТСЯ ТРЕТИЙ ТАЙМ ДО КОНЦА ИГРЫ ОСТАЕТСЯ 8 МИНУТ ИГРЫ ПОКАЗЫВАЕТ ЧТО ТОРОНТО ЗАБИВАЕТ Я ЗНАЧИТ АГА МОЛОДЦЫ СДЕЛАЛИ НА ЧТО ШЕЛ РАНО РАДОВАТЬСЯ МАТЧ ВООБЩЕ ПРОПАЛ С ЛАЙВА НУ Я ЗНАЧИТ ОБНОВЛЯЮ ЕГО ОН КАК ПРОПАЛ БОЛЬШЕ НЕ ПОКАЗАЛСЯ СМОТРЮ ИСТОРИЮ ПОКАЗЫВАЕТ ПРОЙГРАЛ МАЙКОРЕ ПОСМОТРЕЛ ИГРУ ТАМ ПОКАЗАЛО ЧТО ВТОРАЯ КОМАНДА ЗАБИВАЕТ .А ОЛИМП НЕ ДАЛА ТАКУЮ ВОЗМОЖНОСТЬ МНЕ .ЕСЛИ БЫ Я ВИДЕЛ ЧТО КОГДА ЕЩЕ ШЛА ИГРА И ЗАБИЛА ВТОРАЯ КОМАНДА Я БЫ СМОГ ПРОДАТЬ СТАВКУ ЧТОБ ИГРАТЬ ДАЛЬШЕ НО ОНИ ЭТОГО НЕ ХОТЕЛИ СДЕЛАТЬ ПОТОМУ ЧТО МНОГО ЛЮДЕЙ СТАВИЛО НА ТОРОНТО ПОТОМУ ЧТО ОНИ ХОРОШО ИГРАЮТ И БЫЛО МНОГО НА НИХ ПОСТАВЛЕНО ЕСЛИ БЫ ЭТОМ МАТЧ ВИСЕЛ В ЛАЙВЕ ЛЮДИ КОТОРЫЕ ВИДЕЛИ ЧТО ИХ СТАВКА ПРОИГРЫВАЕТ ОНИ БЫ СДЕЛАЛИ КАК И Я ЖЕ ПРОДАЛИ ЕЕ ,НО ДЛЯ ОЛИМПА ЭТО НЕ ВЫГОДНО ВОТ ОНИ И УБРАЛИ МАТЧ С ГРУППЫ ИГРЫ .ПИСАЛ ИМ ОНИ МНЕ ГОВОРЯТ ВАША СТАВКА СЫГРАНО ТАК КАК ВЫ ПОСТАВИЛИ НА КОМАНДУ ПЕРВУЮ ВОЗВРАТА НЕ БУДЕТ У НАС БЫЛА ТЕХНИЧЕСКАЯ НЕ ПОЛАДКА НУ КАК ВСЕГДА ОНИ ЭТО ГОВОРЯТ Я ТОЖЕ НЕ СОВЕТУЮ ИГРАТЬ В ЭТОЙ КОНТОРЕ ЕСТЬ ХОРОШИЕ КОНТОРЫ ПО СРАВНЕНИЮ ОЛИМП И ФОНБЕТА”1

Почему?

Потому что отзывы пишут люди

“БК ЛЕОН БК”

**“в крымю нельзя - конч енная бк
апвап”**

**“После этого перестал играть в
этой конторе. А так все очень
даже неплохо у них.”**

И что делать?

Свихнуться, но сделать всё красиво
и правильно.

Методы: экспертный

Разметка:

1. **Многоклассовая** - размечаем по классам
2. **По шкале** - размечаем по шкале либо размечаем слова внутри отзыва и высчитываем тональность для отзыва
3. **По характеристикам:**
 - a. Все фразы/слова, относящиеся к характеристике + Все фразы/слова, дающие оценку этой характеристике
 - b. Lite -- просто ставим тэг характеристика-тональность отзыву

Методы: rule-based

Смотрим на данные, придумываем правила, которые покрывают максимальное количество случаев.

В любом случае нужна экспертная разметка для:

1. Словарей
2. Оценки результата

“если сказуемое ("нравится") входит в положительный набор глаголов ("нравится", "люблю", "одобряю" ...) и в предложении не имеется отрицаний, то классифицировать тональность как "positive"

Методы: обучение без учителя

Автоматически выделяем из текстов слова/наборы слов по метрикам:

- **TF/IDF**
- **PMI**
- **Logreg**
- **Частоте**

В любом случае нужна экспертная разметка для:

1. **Улучшения словарей**
2. **Оценки результата**

Методы: обучение с учителем

1. **Берем** данные
2. **Предобрабатываем**, если необходимо
3. **Размечаем** часть как нам нужно
4. **Обучаем** классификатор.

Методы: обучение с учителем | предобработка

Прежде чем засовывать текст в классификаторы, нам нужно как-то преобразовать его в цифры.

Как?

1. TF-IDF
2. Вектора (`CountVectorizer`, `FastText`)
3. К-к-комбо

Методы: обучение с учителем | предобработка

Но перед этим пайплайн:

1. Убираем **стоп-слова** (руками/nltk)
2. Убираем **мусор (ссылки)** (регулярками)
3. **Лемматизируем** (переводим в начальную форму) (pymorphy, mystem)
4. **Ловеркейс (lower)**
5. Делаем **биграммы-триграммы**, если нужно (руками или через nltk)
6. **Векторизируем** (превращаем в цифры)
7. **Подаем** в классификатор. **Вы прекрасны!**

Методы: обучение с учителем| классификаторы

- **Алгоритмы классификации:** SVM, LinearRegression, RandomForest, GradientBoosting, LogisticRegression
- Чем **больше** классов - тем **сложнее** классифицировать
- **Маленькая** выборка - **большие** проблемы
- **Несбалансированные** классы - тоже **проблемы**
- Для **бинарной классификации** на общую тональность работают **хорошо**, для **аспектов** - **плохо**

Методы: обучение с учителем| результаты

- Разбиваем на **train** и **test**, смотрим:
 - a. **Accuracy**
 - b. **Precision**
 - c. **Recall**
 - d. **F1**
- Сравниваем с **бенчмарками**.

Инструменты, которые вам понадобятся:



NLTK



pandas
 $y_i = \beta' x_i + \mu_i + \epsilon_i$



TensorFlow

Инструменты, которые вам понадобятся:

Словари тональности:

- <http://linis-crowd.org/about/>
- Четверкина-Лукашевич (по запросу)
- http://web-corpora.net/wsgi/senti_game.wsgi/about

Инструменты, которые вам понадобятся:

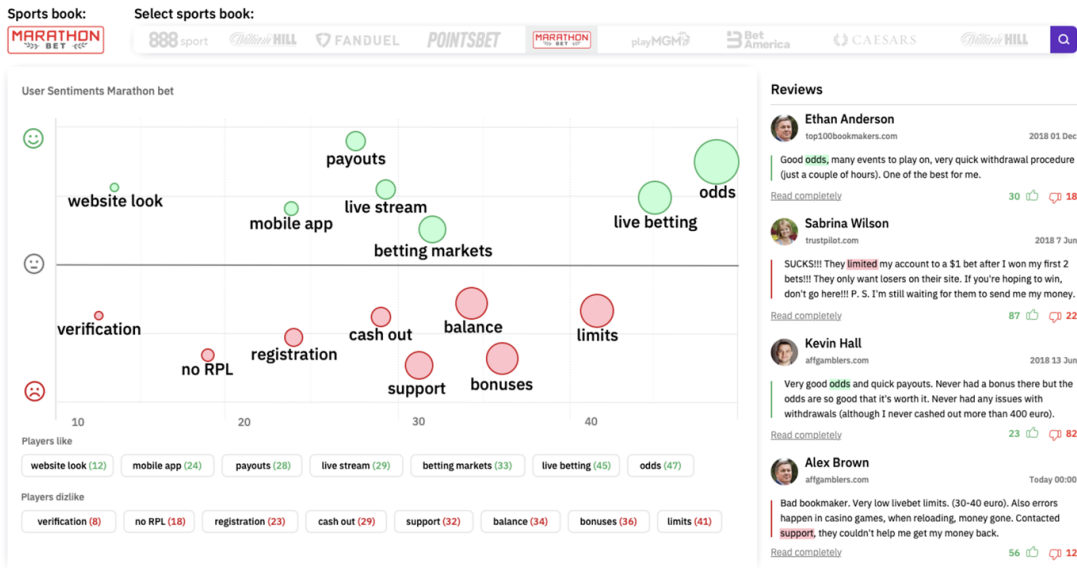
Открытые корпуса для тренировок:

- https://github.com/oldaandozerskaya/auto_reviews
- <https://study.mokoron.com/>
- http://alt.qcri.org/semEval2016/task5/data/uploads/trial-data/restaurants_trial_russian.xml

Или достаем сами с помощью парсеров

Для чего это
нам

User Review Sentiments: Layer 1



Что получилось у нас:

Бинарная классификация для общей тональности:

```
tf-idf + LogisticRegression = precision/recall/f-measure ~ 0.92-0.93
```

Что получилось у нас:

Классификация по аспектам
(15 аспектов, 2 тональности: позитив, негатив):

```
deep memory network with attention model: test-acc=0.70-0.71
```


СПАСИБО!

**... и не забудьте оставить
отзыв!**

BR LAB



BRLAB

Виталий Горбачев

Руководитель направления R&D

✉ vitaliy.g@brl.ru

✈ @imbarus

f [vitaly.gorbachev.5](https://www.facebook.com/vitaly.gorbachev.5)

