

# Senty

Балакший Андрей    Сухочев Александр  
Куратор: Юрий Курочкин

19 мая 2015

## О проекте

**Сентиментальная разметка коротких текстов с использованием "чистой" разметки обучающего множества.**

## Вдохновение

Прочитали статью исследователей из Стэнфордского университета о сентиментальной разметке твитов. Они производили "грязную" разметку твитов(разметка по смайликам) и обучались по ней.

[http://cs.stanford.edu/people/alecmgo/papers/  
TwitterDistantSupervision09.pdf](http://cs.stanford.edu/people/alecmgo/papers/TwitterDistantSupervision09.pdf)

## Вдохновение

Прочитали статью исследователей из Стэнфордского университета о сентиментальной разметке твитов. Они производили "грязную" разметку твитов(разметка по смайликам) и обучались по ней.

<http://cs.stanford.edu/people/alecmgo/papers/TwitterDistantSupervision09.pdf>

Захотелось чего-то аналогичного, но с русским языком и с "чистой" разметкой (разметкой вручную).

## Цель

**Хотим по короткому тексту понимать настроение его автора.**

## Этапы выполнения

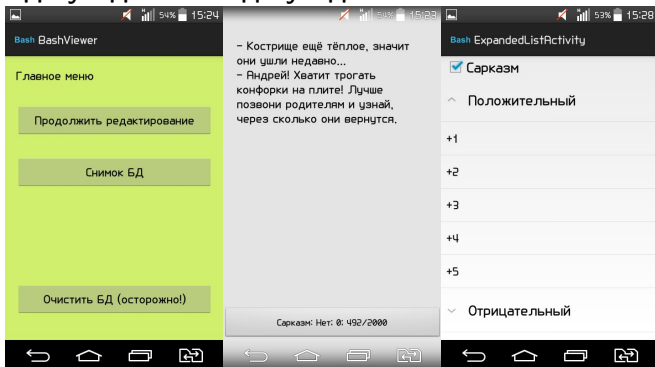
- Найти источник материалов для обучения, содержащий высказывания людей на различные темы.
- Разметить эти материалы вручную на положительные и отрицательные.
- Придумать различные признаки, которые помогут определять настроения авторов текстов (далее для упрощения будем называть эти признаки фичами).
- Реализовать машинное обучение по этим признакам.

## Получение материалов для обучения

Материалы для обучения спарсили с цитатника рунета **bash.im**. Для парсинга воспользовались питоновской библиотекой **Beautiful Soup**.

## "Чистая" разметка обучающего множества

Для удобства "чистой" разметки текстов (далее башей) написали приложение под Android, чтобы размечать баши где угодно и когда угодно.





## Фичи

Исходя из логических соображений, имеющихся знаний и результатов работы классификации для различных башей, было создано множество фич, являющихся надстройками над стандартным экстрактором термов. В стандартном экстракторе мы просто убрали знаки препинания и нормализовывали слова с помощью питоновской библиотеки `pymystem3`.

## Примеры работы некоторых фич

- **Исходный текст:**

"УЖАААСНО не люблю вставать по понедельникам и идти на работу ((((... и по вторникам тоже )))"

## Примеры работы некоторых фич

- **Исходный текст:**

"УЖАААСНО не люблю вставать по понедельникам и идти на работу ((((... и по вторникам тоже )))"

- **Работа стандартного экстрактора с pystem3:**

"ужааасно не любить вставать по понедельник и идти на работа ((( ( и по вторник тоже )))"

## Примеры работы некоторых фиц

- **Исходный текст:**

"УЖАААСНО не люблю вставать по понедельникам и идти на работу ((((... и по вторникам тоже )))"

- **Работа стандартного экстрактора с pystem3:**

"ужааасно не любить вставать по понедельник и идти на работа ((( ( и по вторник тоже )))"

- **Добавление фици, убирающей предлоги, союзы и местоимения:**

"ужааасно не любить вставать понедельник идти работа ((( ( вторник тоже )))"

## Примеры работы некоторых фиц

- **Исходный текст:**

"УЖАААСНО не люблю вставать по понедельникам и идти на работу ((((... и по вторникам тоже )))"

- **Работа стандартного экстрактора с pystem3:**

"ужааасно не любить вставать по понедельник и идти на работа ((( ( и по вторник тоже )))"

- **Добавление фици, убирающей предлоги, союзы и местоимения:**

"ужааасно не любить вставать понедельник идти работа ((( ( вторник тоже )))"

- **Добавление фици, убирающей повторные буквы:**

"ужасно не любить вставать понедельник идти работа ((( ( вторник тоже )))"

## Примеры работы некоторых фич

- **Исходный текст:**

"УЖАААСНО не люблю вставать по понедельникам и идти на работу ((((... и по вторникам тоже )))"

- **Добавление фичи, склеивающей "не" со словом, которому эта частица принадлежит:**

"ужасно нелюбить вставать понедельник идти работа ((((вторник тоже)"

## Примеры работы некоторых фич

- **Исходный текст:**

"УЖАААСНО не люблю вставать по понедельникам и идти на работу ((((... и по вторникам тоже )))"

- **Добавление фичи, склеивающей "не" со словом, которому эта частица принадлежит:**

"ужасно нелюбить вставать понедельник идти работа (((  
вторник тоже"

- **Добавление фичи, выделяющей смайлики:**

"ужасно нелюбить вставать понедельник идти работа  
**плохосмайл** вторник тоже **хоросмайл**"

## Примеры работы некоторых фич

- **Исходный текст:**

"УЖАААСНО не люблю вставать по понедельникам и идти на работу ((((... и по вторникам тоже )))"

- **Добавление фичи, склеивающей "не" со словом, которому эта частица принадлежит:**

"ужасно нелюбить вставать понедельник идти работа (((  
вторник тоже"

- **Добавление фичи, выделяющей смайлики:**

"ужасно нелюбить вставать понедельник идти работа  
**плохосмайл** вторник тоже **хоросмайл**"

- **Также были созданы фичи, выделяющие n-граммы, убирающие иностранные слова и многие другие.**



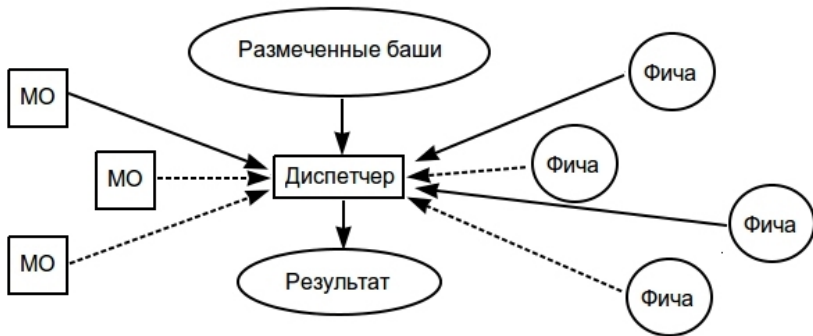
# Выбор классификаторов

- Naive Bayes
- SVM
- MaxEntropy

Для их реализации использовали питоновскую библиотеку `scikit-learn`.

## Структура проекта

Для использования различных комбинаций фич с различными классификаторами была использована следующая структура.



# Проблемы

- Малый объём выборки(примерно каждый четвёртый баш является эмоционально окрашенным, много сарказма, "чистая" разметка - дело не быстрое). Всего 9811 термов было выделено стандартным экстрактором.
- Сильный "перекос" в сторону отрицательных башей в обучающей выборке: из 1115 башей, имеющих эмоциональную окраску, 732 являются отрицательными (65.65%).

## Дальнейшие перспективы

- Двойная классификация: сначала классифицировать на нейтральное/с эмоциональной окраской, потом уже эмоциональные на положительные/отрицательные.
- Прикрутить поиск по упоминаниям по базе башей для определения в скольких упоминаниях о предмете поиска отзывались положительно/отрицательно/нейтрально, с возможностью получения детальной информации по башам (по какому-то фильтру например)
- Увеличение материала для обучения

# Итоги практики

- Работа в команде
- Изучение Python(Mystem, Scikit-learn, Beautiful Soup)
- Работа с GitHub
- Работа с  $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ (собственно создание самой презентации)



# Спасибо за внимание

<https://github.com/cscenter/senty>