



ODS QA

Финальный проект курса:
DL in NLP Spring 2020

Докладчик: Юрий Зеленский
Куратор: Алексей Сорокин

DeepPavlov.ai

Команда проекта ODS QA

Илья Сироткин (@Satel1ite) - основной вклад

Екатерина Карпова (@vengodelsur) - автор изначальной идеи

Юрий Зеленский (@yzelensky) - организация, идеи, доклад

Вадим Нарейко (@vnareyko) - консультации

Алексей Сорокин - куратор

Вехи

“Автоответчик ODS” в slack #ods_pet_projects

- Open domain QA != SQuAD
- Open domain QA ещё не с нами

Кластеризация вопросов

- “Родная” задача для USE с опубликованным baseline (Quora, AskUbuntu)
- ODS Community - сообщество высокой культуры

Классификация вопросов по каналам в Slack ODS:

- сравнительный анализ современных моделей
- робот recommending тематический канал



ekatkar 🤖 9:55 PM

Wednesday, March 18th ▾

Привет!

В ожидании уточнений про то, как будет проходить наш онлайн-формат, приглашаем и вас тоже рассказывать про идеи проектов. Пока — ещё один проект во имя сообщества под уже мелькавшим кодом "ODS QA".

Вас должна воодушевить идея, если вы грустите, когда:

- пишут по много раз одни и те же вопросы в канал, ну можно же по слаку поискать! 🦎
- задаёшь вопрос, а в ответ только «было уже, поищите», да пытал я уже этот поиск по слаку, ну тяжко ему с русской морфологией!
- хочется попилить что-нибудь NLPшное для души, а на кагле одни картиночки!

Мы хотим сделать «автоответчик» для ODS. Это значит, что можно:

- получить плюшечки от сообщества, как это всегда бывает с проектами во славу ODS 🧑
- покрутить модные модельки (yorko не прочь поделиться опытом успешнейшего гоняния альбертов на TPU)
- да что угодно в целом, хотите попробовать написать модель на ассемблере или сделать демку с анимацией на HTML5 — мы чужие вкусы не осуждаем...
- и отдельная интересная часть, с которой всё начнётся: поразвлекаться с краудсорсингом. Нам под силу сделать из нашего чудесного дампа слака датасет, достойный стать Новым Важным Бенчмарком в Question Answering.

Пока неизвестно, что может встать на нашем пути: саркастические ответы? Грязная ложь в комментариях, противоречащая документации rutorch? Расползание обсуждения по нескольким тредам? Так что попутно помимо вопросов-ответов может понадобиться ещё много какая разметка, идущая рука об руку с новыми задачами, так что не волнуйтесь, тем для курсовых проектов на 1Pшных курсах хватит всем.

Ставьте +, кому интересно поучаствовать. Тогда я не позволю вам пропустить ссылочку на репозиторий:) (edited)

+ 16 🔥 1 😊

10 replies Last reply 4 months ago

Вехи

“Автоответчик ODS” в slack #ods_pet_projects

- Open domain QA != SQuAD
- Open domain QA ещё не с нами

Кластеризация вопросов

- “Родная” задача для USE с опубликованным baseline (Quora, AskUbuntu)
- ODS Community - сообщество высокой культуры

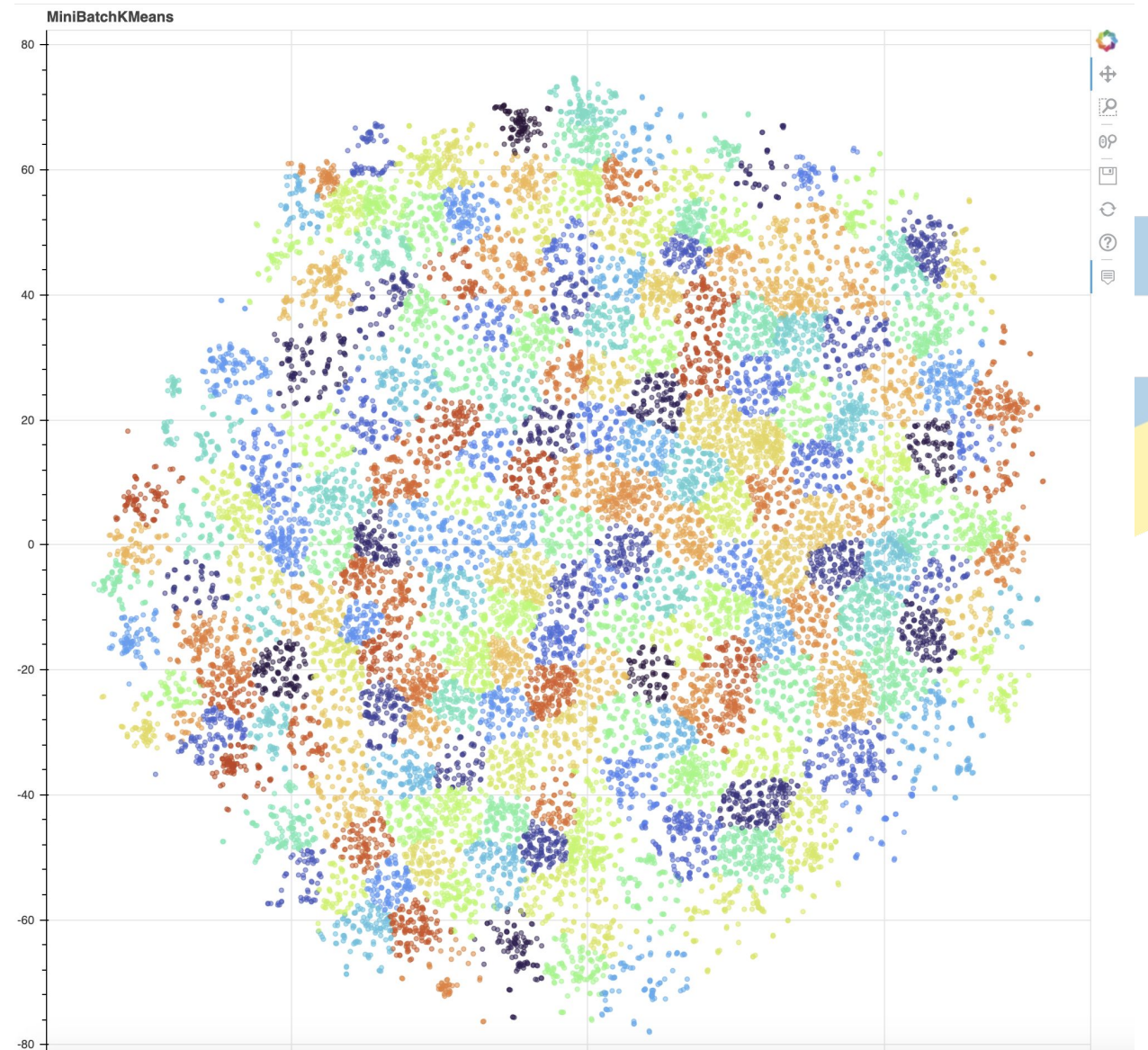
Классификация вопросов по каналам в Slack ODS:

- сравнительный анализ современных моделей
- робот recommending тематический канал



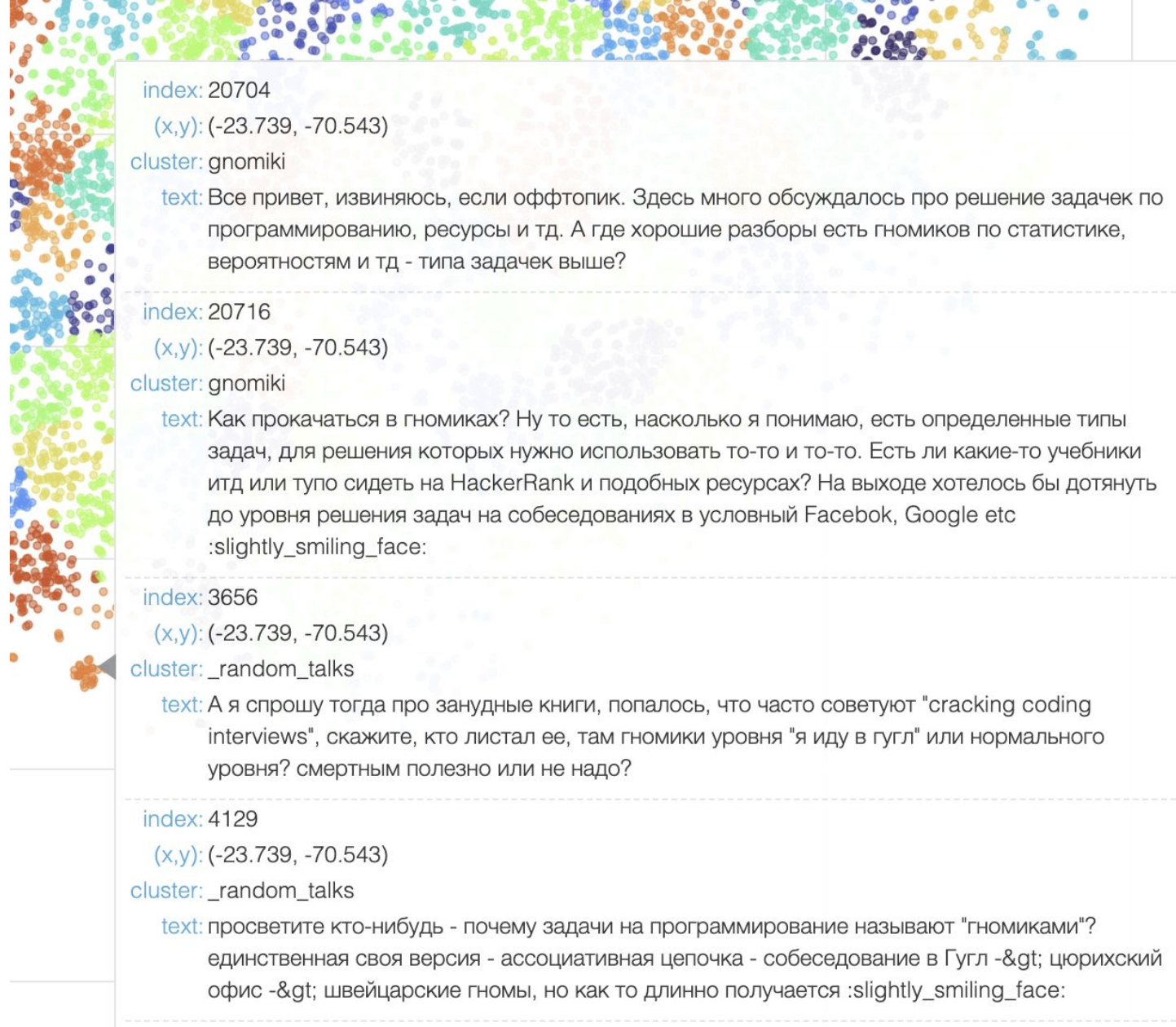
EDA

bokeh
tsna, k-means
embeddings
jittering



EDA

Редкий пример вопросов дубликатов



Вехи

“Автоответчик ODS” в slack #ods_pet_projects

- Open domain QA != SQuAD
- Open domain QA ещё не с нами

Кластеризация вопросов

- “Родная” задача для USE с опубликованным baseline (Quora, AskUbuntu)
- ODS Community - сообщество высокой культуры

Классификация вопросов по каналам в Slack ODS:

- **сравнительный анализ современных моделей**
- **робот рекомендующий тематический канал**



Итоговая задача

Мета-идея:

- сравнительный анализ современных моделей на примере данных дампа ODS с адекватной метрикой выраженной одним числом.

Постановка:

- классификации принадлежности **вопроса** к slack каналу ODS

Метрика:

- Pairwise macro average multiclass ROC AUC score (pwROC-AUC).
- A Simple Generalisation of the Area Under the ROC Curve for Multiple Class Classification Problems <https://link.springer.com/article/10.1023/A:1010920819824>

Данные

- Дамп slack ODS с 2015-03-12 по 2020-01-14
- Выбраны треда
- Начало треда - вопрос, сообщения - ответы
- Фильтр по наличию знака вопроса
- Фильтр по каналам - 100+ вопросов

	Train	Val	Test
Вопросы	215,054	71,685	71,685
Слова	10,765,047	3,587,705	3,622,738
Классы	79	79	79

EDA словоформы (regex, r morphology)

normal_form	forms_count	word_count_sum	forms
использовать	38	6684	использовать использую использует использовал используете используя использую...
хороший	37	4139	хороший хорошие лучший хорошая хорошего лучшие лучших хорошее хороших хорошо...
написать	31	2880	написать написал написано написали напишите написала напишу написаны написан...
реализовать	31	671	реализовать реализовал реализовали реализован реализовано реализована реализ...
найти	29	4902	найти нашел нашёл нашла найду нашли найдет найди найдёт найдены найдешь найд...
сделать	29	6982	сделать сделал сделали сделано сделает сделаю сделаны сделала сделан сделает...
получить	29	3428	получить получил получили полученных получим получу получит полученный получ...
выбрать	26	1402	выбрать выбрали выбрал выбрала выбранных выбрано выбраны выбран выбери выбра...
построить	26	892	построить построил построили постройте построено построен построишь построен...
обучать	25	1470	обучать обучаю обучающей обучаем обучающую обучал обучали обучающая обучающи...
сохранить	25	421	сохранить сохранил сохранив сохранена сохранено сохраненные сохраненную сохр...
дать	25	10481	данных данные данными дать данным данной дали дайте дал даст даны дало дадут...

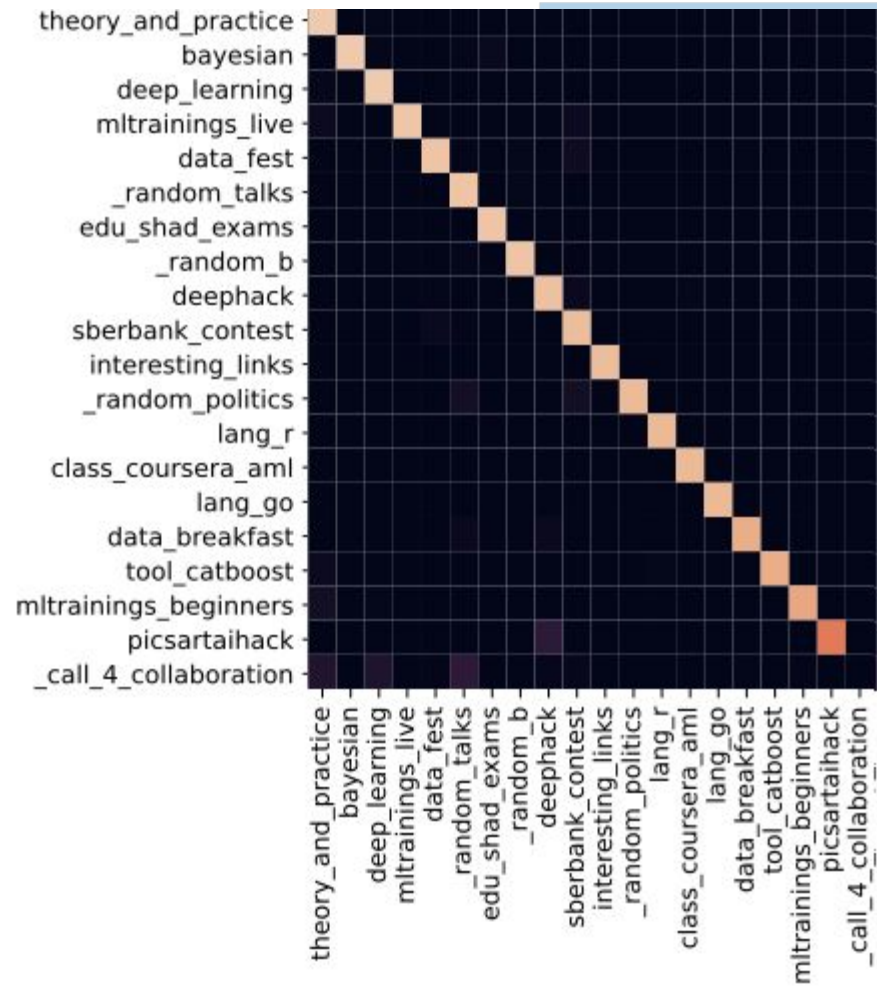
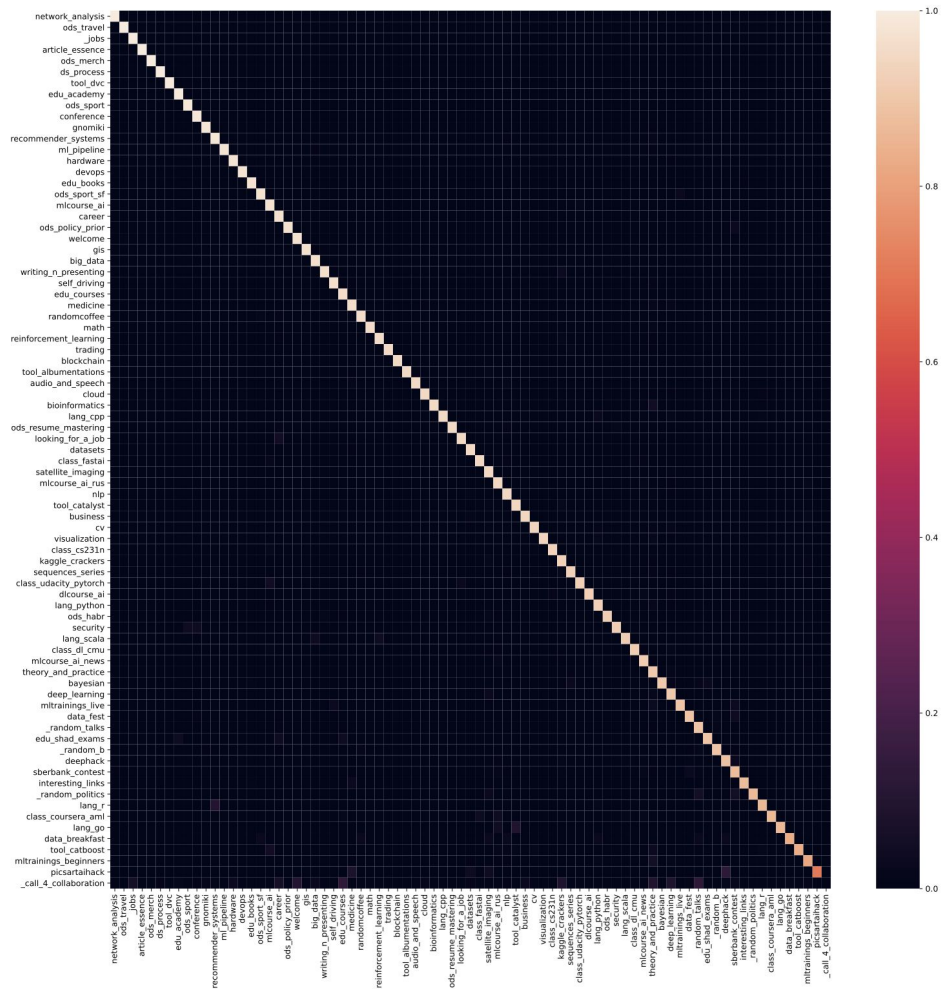
Результаты

Model	pwROC-AUC score
USE_baseline	0.61
USE_finetuned	0.985
SBERT_finetuned	0.991
DeepPavlov ru-SBERT	0.995
DeepPavlov mul-SBERT	0.993

<https://tfhub.dev/google/universal-sentence-encoder-multilingual/3>

<https://github.com/UKPLab/sentence-transformers>

Результаты



Выводы

- Дообучение проанализированных моделей, сильно повысило их результативность на задаче классификации
- Современным моделям не нужна предварительная токенизация
- Даже дообучение требовательно к объему доступной GPU RAM
- Высокие результаты (и USE и SBERT) открывают путь к практическому использованию
- Высокие абсолютные показатели, свидетельствуют о высоком качестве разделения (и разделимости) исходных данных (и ODS - сообщество высокой культуры и эвристика построения датасета сработала)

Дальнейшие исследования и разработки

- Добавить к сравнительному анализу “классические” подходы к topic modelling (LDA, SVM on TF-IDF, etc)
- Добавить к сравнительному анализу другие современные модели, в частности справившись с их ограничениями на объем GPU RAM (ULM-Fit, ELMO)
- Провести hard sample mining анализ
- Понять причину сложностей с классификацией сообщений канала _call_4_colaboration
- Провести слепой тест, на свежих данных, исключить “утечку”
- Разработать Slack-бота, ненавязчиво рекомендующего “правильный” канал
- Вернуться к исходным, более “тяжелым” постановка задачи

Спасибо за внимание



Отчет <https://gist.github.com/yzelensky/72148d7b4e8ad62d2f480a2e3ccb6228>

