

Определение эмоциональной окраски постов пользователей в социальных сетях.

Под постом подразумевается нечто, что состоит из картинки и текста, в объединении обладающими смысловой нагрузкой. То есть видео и аудио данные рассматриваться не будут.

Поскольку предполагается работать с текстовыми данными и изображениями, то предлагается использовать две модели: одна - отвечает за текст, вторая - за изображения.

Также, для простоты, ограничимся английским языком. Данные.

Поскольку предполагается решать задачу классификации, причём двумя различными способами, то придётся использовать два датасета: в первом тексты и метки им соответствующие (эмоциональные окраски), аналогично с картинками. Причем метки классов должны быть одинаковыми в обоих наборах.

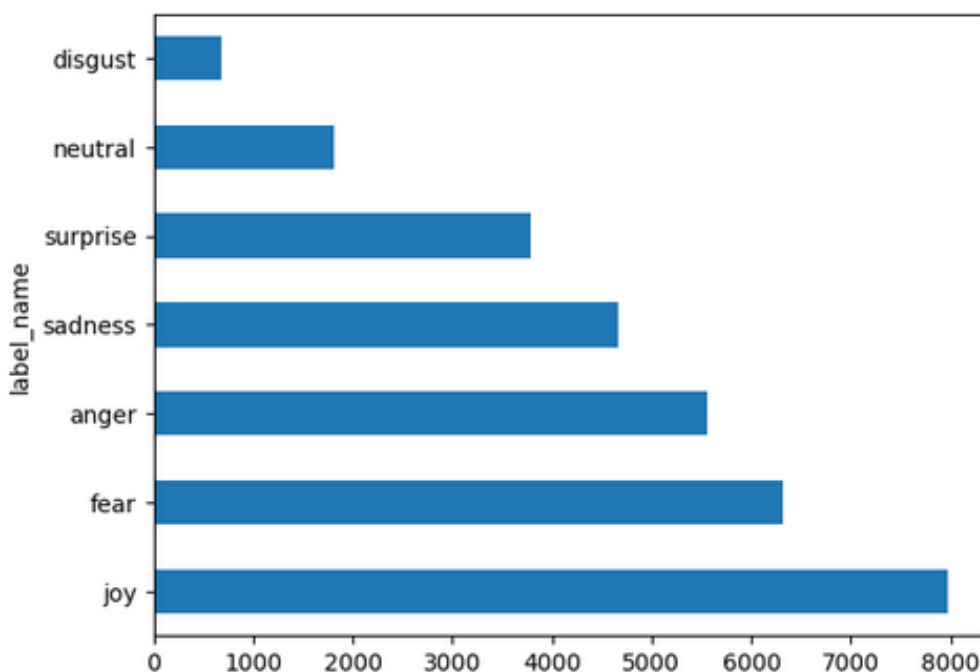


Figure 1: Количество примеров для каждого лэйбла

Топ 7 слов в каждом из классов:

Как этот топ 7 выглядит в терминах облаков слов.

В первую очередь подбирались данные с картинками, потому что такие наборы искать было труднее, чем текстовые. Поэтому в качестве датасета с картинками был выбран [*https://huggingface.co/*](https://huggingface.co/)

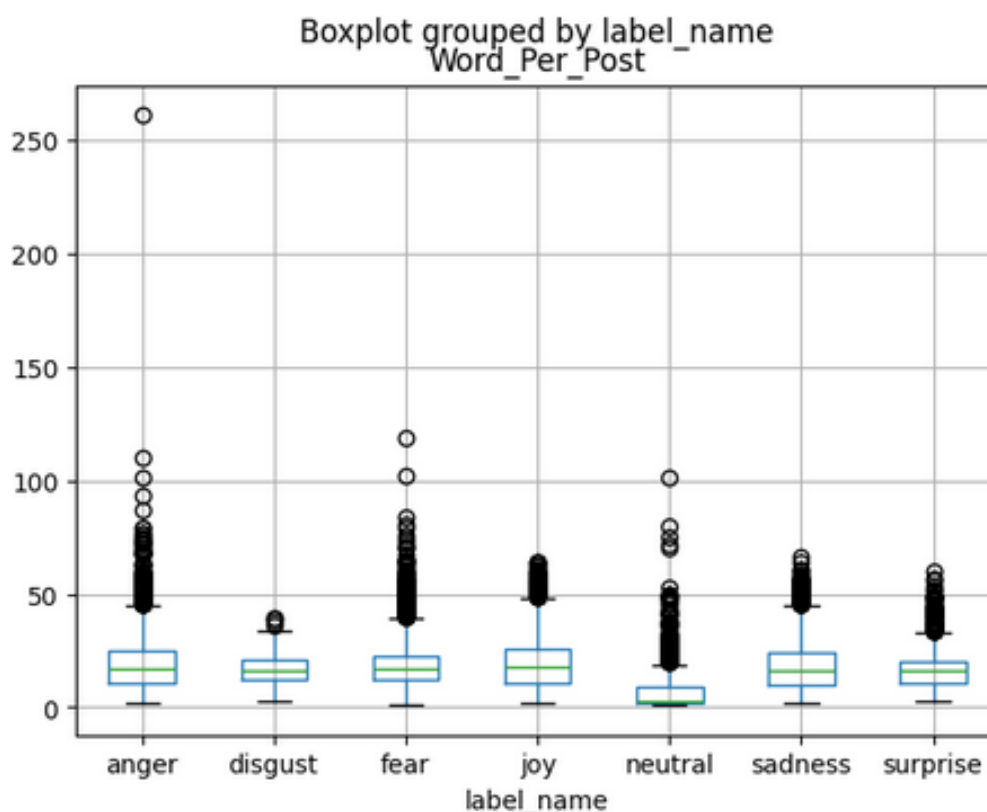


Figure 2: Длина текстов в каждом из лэйблов

содержащий в себе 9400 изображений лиц людей, разбитых по 8-ми категориям. В дальнейшем класс "content" пришлось убрать, т.к. соответствующих текстовых данных подобрать не удалось.

В качестве текстового датасета были выбраны два набора, метки которых объединении содержали все классы, присущие картинкам (лишние, разумеется пришлось устранить). Данные распределены равномерно по лэйблам (судя по Huggingface)

По итогу, получилось два датасета с семью классами.

Модели.

Как говорилось выше - будет задействовано две модели. Для анализа текста была выбрана модель ($7 * 10^7$ параметров) <https://huggingface.co/distilbert/distilbert-base-uncased>

Эта модель в первую очередь предназначена для тонкой настройки на задачах, которые используют целое предложение (потенциально замаскированное) для принятия решений, таких как классификация последовательностей, классификация токенов или ответ на вопросы.

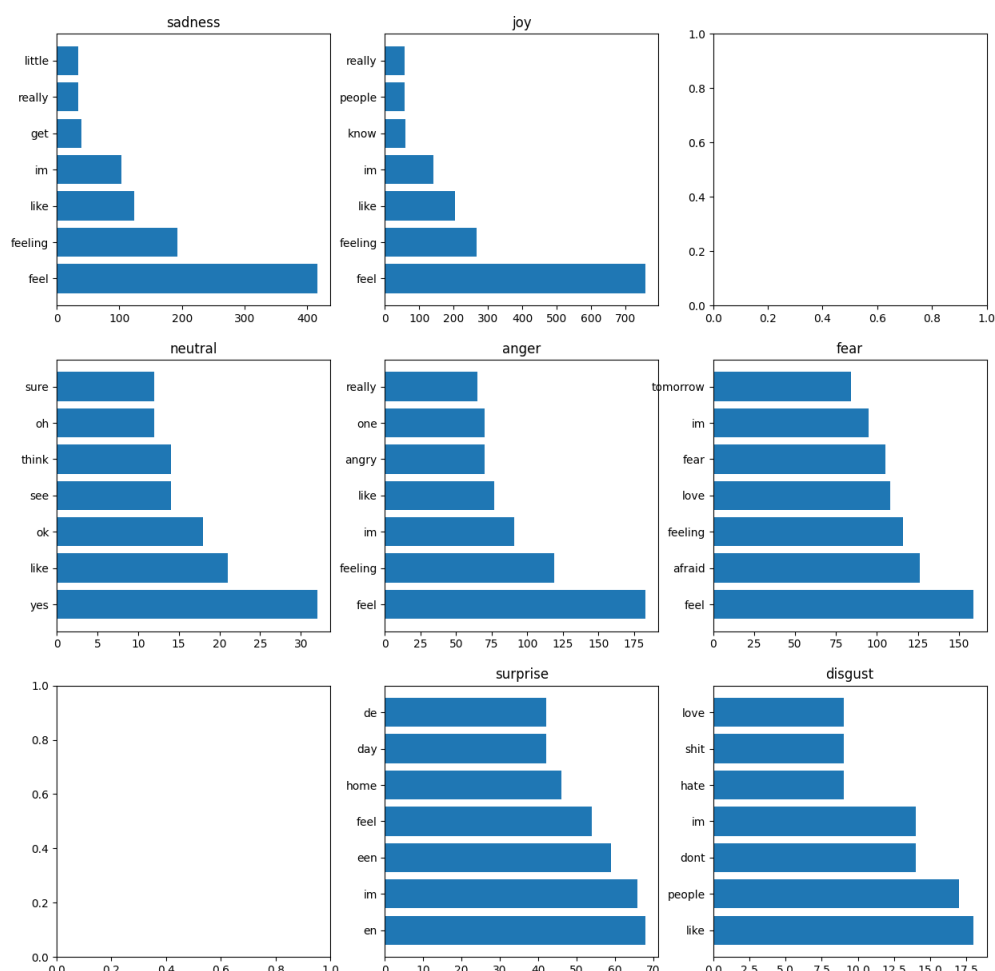


Figure 3: Топ 7 слов в каждом из классов

Для анализа изображений - задействована модель сопоставимых масштабов. <https://huggingface.co/google/vit-base-patch16-224-in21k> Vision Transformer (ViT) — это модель кодировщик-трансформер (похожая на BERT), предварительно обученная на большой коллекции изображений контролируемым образом, а именно ImageNet-21k, с разрешением 224x224 пикселей.

Изображения представляются модели как последовательность фрагментов фиксированного размера (разрешение 16x16). Также добавляется токен [CLS] в начало последовательности, чтобы использовать ее для задач классификации. Также добавляются абсолютные позиционные эмбединги перед подачей последовательности на слой модели.

Обучение.

Distilbert обучалась в 6 эпох на датасете из 30803 примеров,

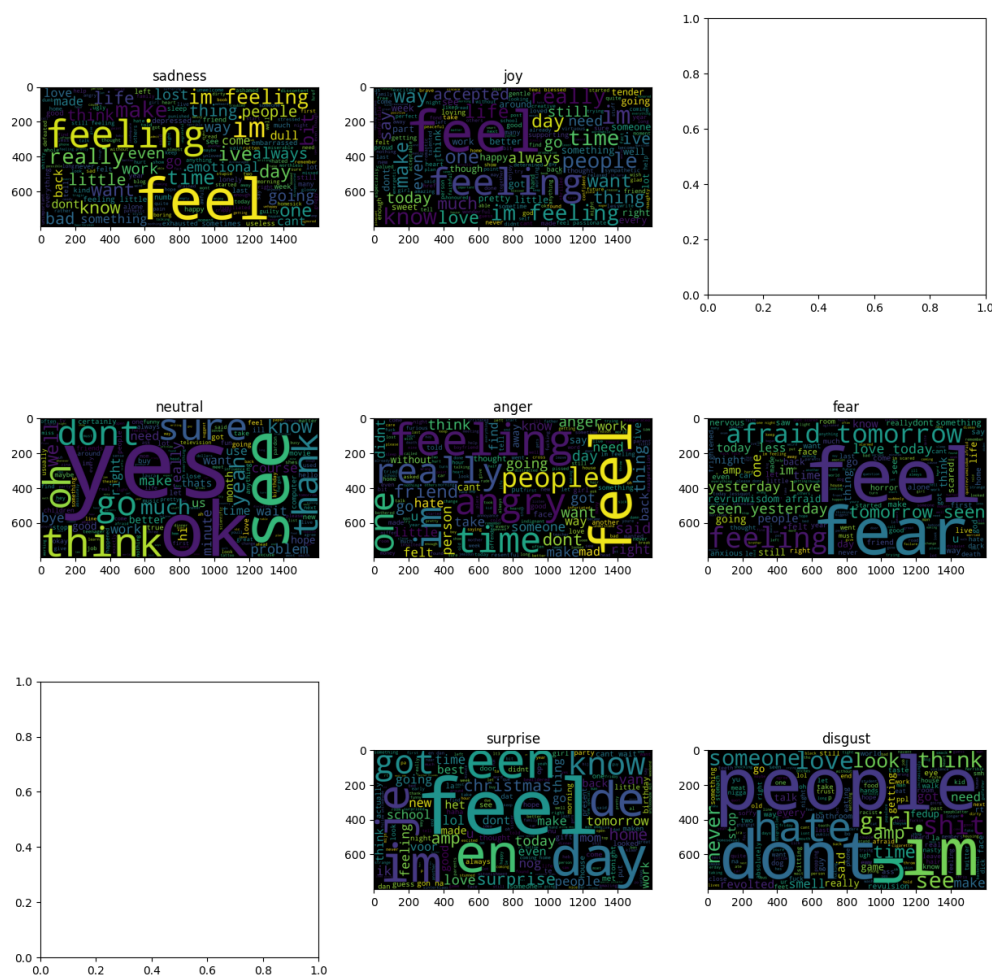


Figure 4: топ 7 выглядит в терминах облаков слов

при помощи стандартных методов Trainer.

ViT обучалась в 6 эпох на датасете из 7254 примеров, при помощи стандартных методов Trainer.

Эксперименты.

Для текстовой модели посмотрим, с какой точностью ей удаётся определять каждый из классов. Поскольку датасет разбалансирован (класс "disgust" в дефиците), то следует ожидать, что самый малочисленный класс будет определяться хуже остальных. То есть на инференсе представитель самого бедного класса может часто падать в другой, но в малочисленный класс редко будут попадать по ошибке. См. рис.

Посмотрим, какие слова чаще всего попадают в предложения, отнесённых к определённому классу.

То же самое, но на облаках слов.

Epoch	Training Loss	Validation Loss	Accuracy	F1
1	No log	0.541501	0.805689	0.788816
2	0.711000	0.391791	0.864689	0.861598
3	0.322300	0.264887	0.916619	0.916233
4	0.236000	0.185536	0.946445	0.946267
5	0.174200	0.141388	0.960585	0.960343
6	0.128900	0.125051	0.965461	0.965403

Figure 5: Обучение Distilbert

Step	Training Loss	Validation Loss	Accuracy
100	1.186600	1.256971	0.533713
200	0.751900	1.039499	0.624881
300	0.714200	0.946472	0.654321
400	0.373200	0.920718	0.685660
500	0.254000	0.942241	0.685660

Figure 6: Обучение ViT

Видно, что почти везде встречаются слова "feel", "feeling", "like", "de", "en". Попробуем выкинуть их, чтобы большую роль играли уникальные слова, что, возможно, приведет к лучшему качеству классификации. Проверим гипотезу, рассмотрев аналогичные картинки. Из таблиц видно, что столь агрессивный подход привёл к заметной потере качества. Что говорит о важности не конкретных слов, а их комбинаций.

Рассмотрим качество классификации в зависимости от длины исходного предложения рис 11.

Для ViT посмотреть какие элементы картинки оказывают наибольшее влияние на классификацию - более сложная процедура, поэтому рассмотрим с какой точностью ей удаётся определять каждый из классов. Исходя из более скромных показателей ViT

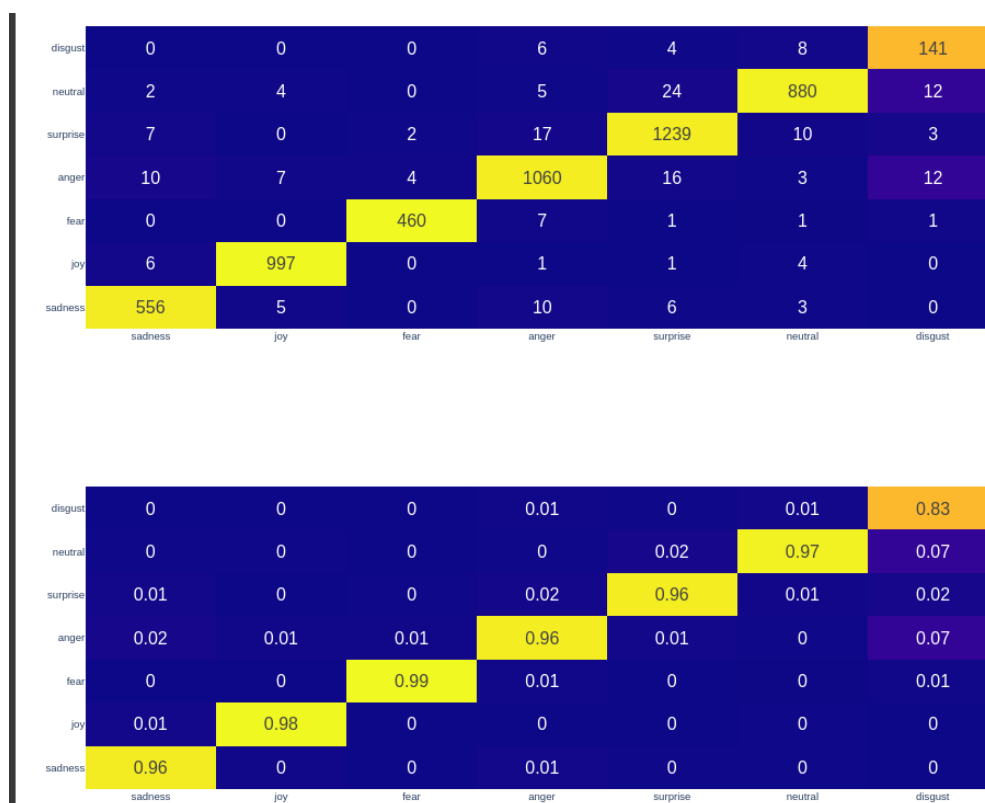


Figure 7: зависимость предсказаний модели от истинного класса. Сверху - количественные оценки на тестовом наборе, снизу - частотные

по сравнению с Distilbert, не трудно догадаться, что большее количество тестовых объектов попадёт вне побочной диагонали. Здесь скорее больший интерес представляет вопрос, а какие классы труднее всего отличать друг от друга. рассмотрим рисунок.

12

Предсказание в комплексе.

Изначально планировалось проводить классификацию с опорой на две модели. Однако, анализ текстов оказался более простой задачей. Поэтому можно ожидать, что во время инференса на объекте, состоящем и из картинки, и из текста, стоит серьёзнее относиться к предсказаниям Distilbert.

Web-интерфейс. Наверное, один из самых разумных способов моделировать пост - отправлять телеграм боту картинку с подписью. Однако реализовать его я не успел.

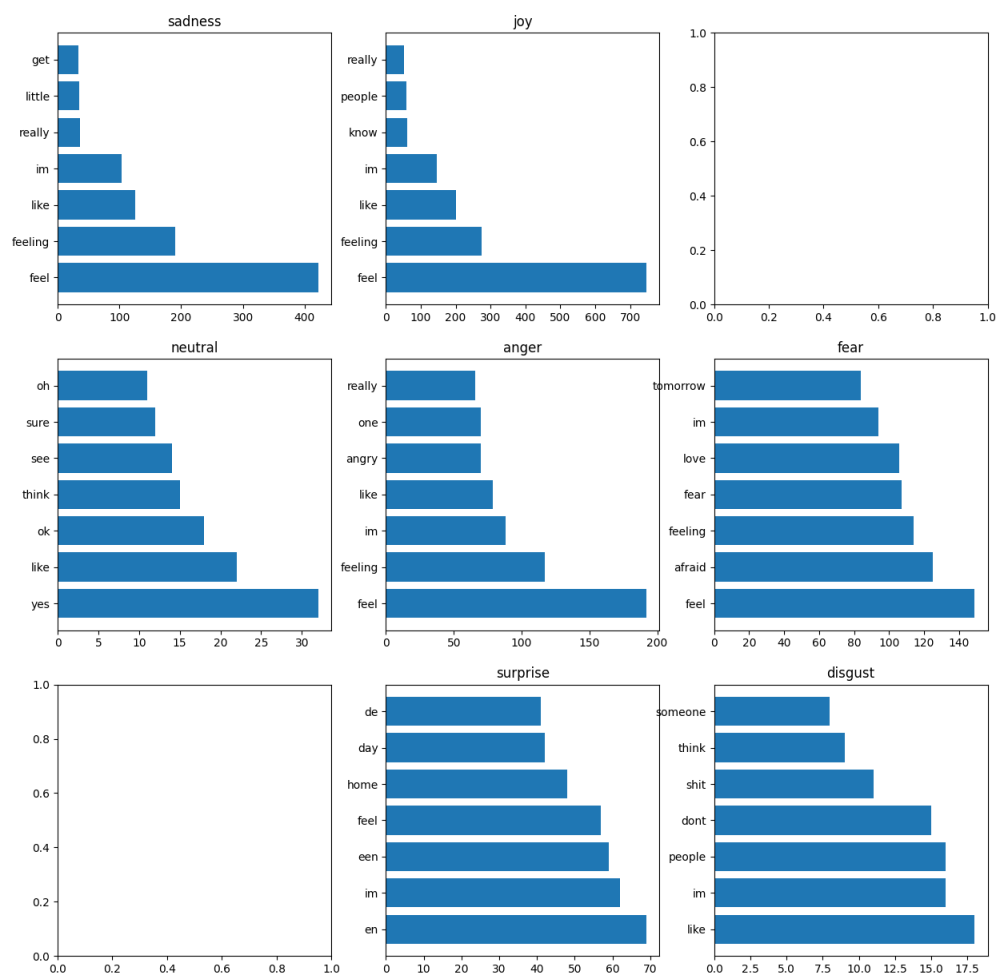


Figure 8: Топ 7 слов в каждом из классов с точки зрения модели

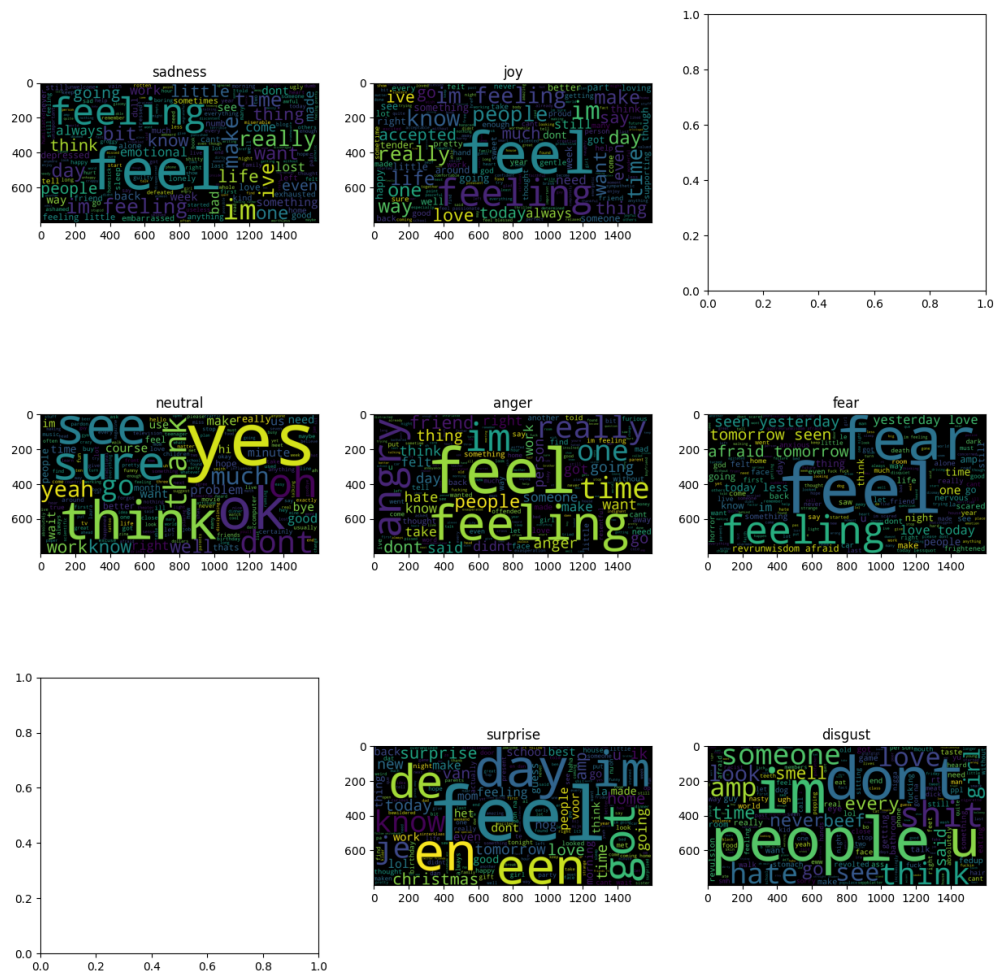


Figure 9: Облака слов для оп 7 слов в каждом из классов с точки зрения модели

disgust	1	7	2	18	12	5	21
neutral	401	657	127	797	949	850	132
surprise	24	44	15	68	237	26	10
anger	5	7	12	116	16	5	3
fear	150	298	310	107	77	23	3
joy	0	0	0	0	0	0	0
sadness	0	0	0	0	0	0	0
	sadness	joy	fear	anger	surprise	neutral	disgust

disgust	0	0.01	0	0.02	0.01	0.01	0.12
neutral	0.69	0.65	0.27	0.72	0.74	0.94	0.78
surprise	0.04	0.04	0.03	0.06	0.18	0.03	0.06
anger	0.01	0.01	0.03	0.1	0.01	0.01	0.02
fear	0.26	0.29	0.67	0.1	0.06	0.03	0.02
joy	0	0	0	0	0	0	0
sadness	0	0	0	0	0	0	0
	sadness	joy	fear	anger	surprise	neutral	disgust

Figure 10: зависимость предсказаний модели после выкидывания некоторых слов от истинного класса. Сверху - количественные оценки на тестовом наборе, снизу - частотные

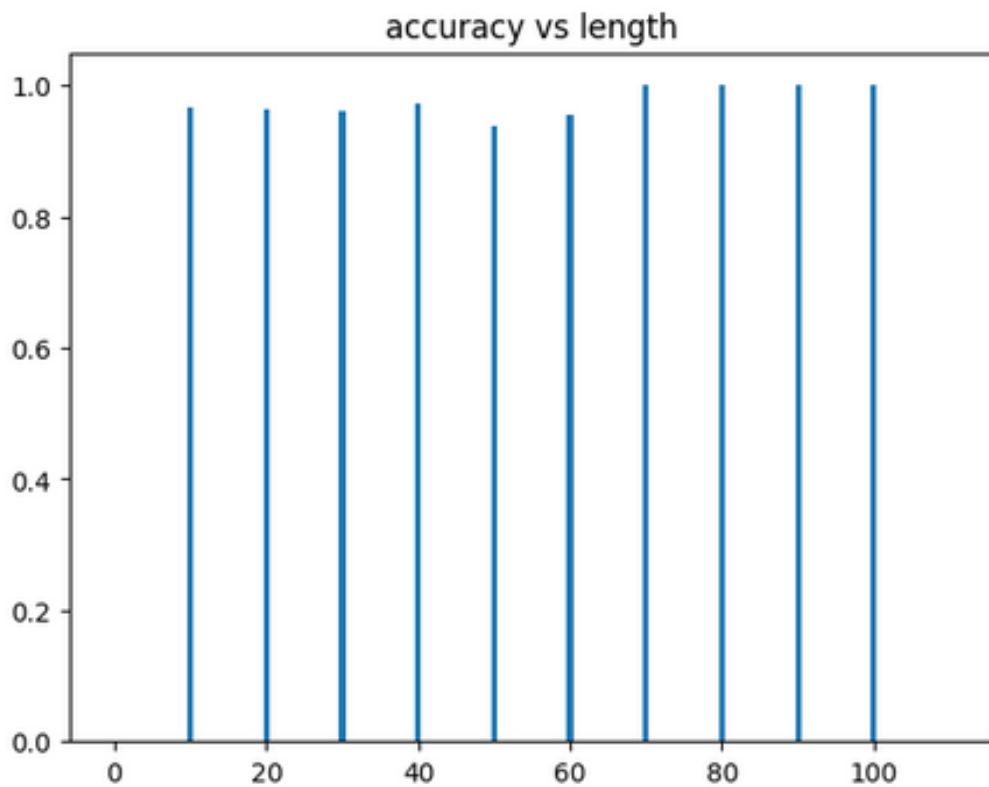


Figure 11: качество от длины промпта

surprise	9	4	12	17	6	7	88
disgust	14	26	6	2	6	106	4
happy	2	0	2	0	120	6	2
fear	7	7	1	116	0	11	35
neutral	18	14	115	3	1	6	9
anger	5	80	6	3	0	10	4
sad	97	19	11	8	7	14	7
	sad	anger	neutral	fear	happy	disgust	surprise

surprise	0.06	0.03	0.08	0.11	0.04	0.04	0.59
disgust	0.09	0.17	0.04	0.01	0.04	0.66	0.03
happy	0.01	0	0.01	0	0.86	0.04	0.01
fear	0.05	0.05	0.01	0.78	0	0.07	0.23
neutral	0.12	0.09	0.75	0.02	0.01	0.04	0.06
anger	0.03	0.53	0.04	0.02	0	0.06	0.03
sad	0.64	0.13	0.07	0.05	0.05	0.09	0.05
	sad	anger	neutral	fear	happy	disgust	surprise

Figure 12: зависимость предсказаний модели по картинкам от истинного класса