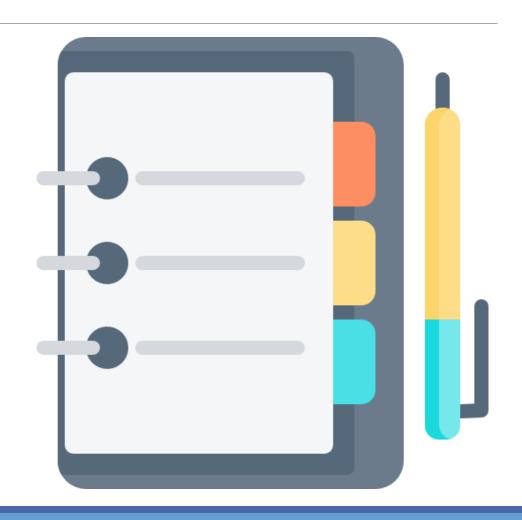


GOOGLE PLAY REVIEWS SENTIMENT ANALYSIS

Agenda

- INTRODUCTION
- PIPELINE
- DATA COLLECTION
- DATA ANALYSIS
- DATA PREPROCESSING
- TEXT REPRESENTATION
- MODELS
- API
- DEPLOYMENT
- FUTURE WORK
- RESOURCES



Introduction



Introduction



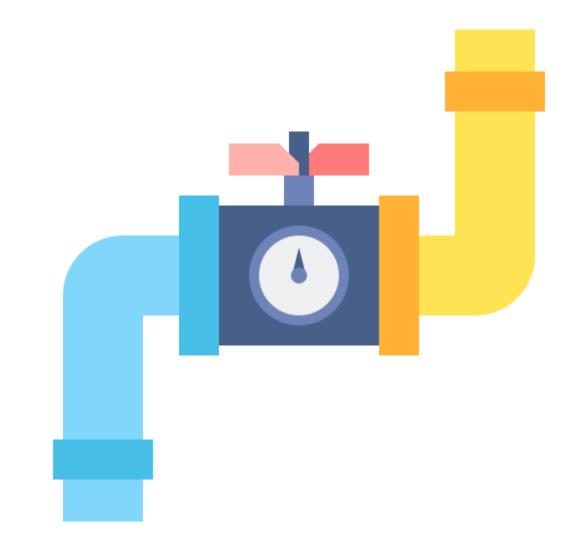
- What is Sentiment Analysis?
- Business Domain





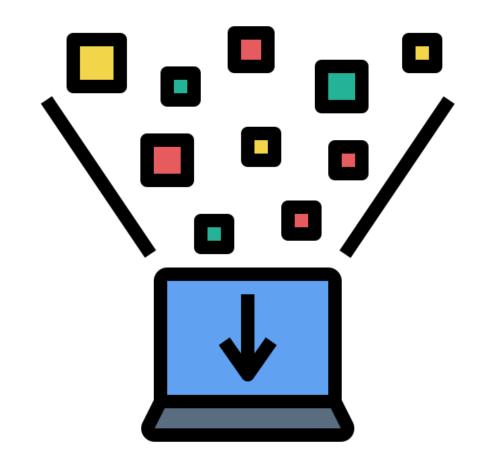
Sentiment analysis

Pipeline



Pipeline







1. Scrap Tool Definition:

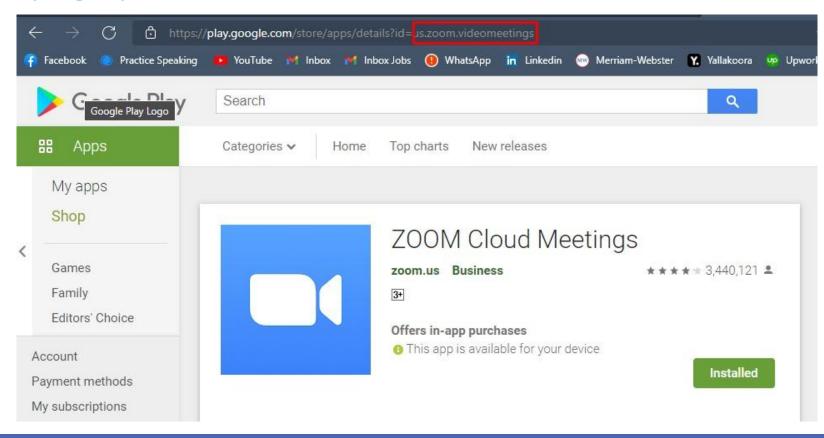
google-play-scraper 1.0.3

pip install google-play-scraper 🕒





2. Scraping Input:





3. Scraping Output

```
"title": "Pokémon GO",
"appId": "com.nianticlabs.pokemongo",
"url": "https://play.google.com/store/apps/details?id=com.nianticlabs.p
"description": "NEW! It's time to interact with your Pokémon like never
"descriptionHTML": "NEW! It's time to interact with your Pokémon like n
"summary": "Step outside and catch Pokémon in the real world! Collect &
"summaryHTML": "Step outside and catch Pokémon in the real world! Colle
"installs": "100,000,000+",
"minInstalls": 100000000,
"score": 4.126332,
"ratings": 12001843,
"reviews": 4772986.
"originalPrice": None,
"sale": False,
"saleText": None,
"saleTime": None,
```

```
"userName": "Alyssa Williams",
"userImage": "https://lh3.googleusercontent.com/-cVEHKr7mzv8/AAAAAA
"content": "This is literally the best idle game I have ever played
"score": 5,
"thumbsUpCount": 54,
"reviewCreatedVersion": "1.16",
"at": datetime.datetime(2020, 2, 24, 17, 19, 34),
"replyContent": "Hello, We will gradually improve the various syste
"repliedAt": datetime.datetime(2020, 2, 24, 18, 30, 42),
"reviewId": "gp:AOqpTOEOIy5S9Je1F8W1BgCl6l_TCFP_QN4qGtRATX3PeB5VV9a
```

Apps Reviews



4. Applications We Scraped:





















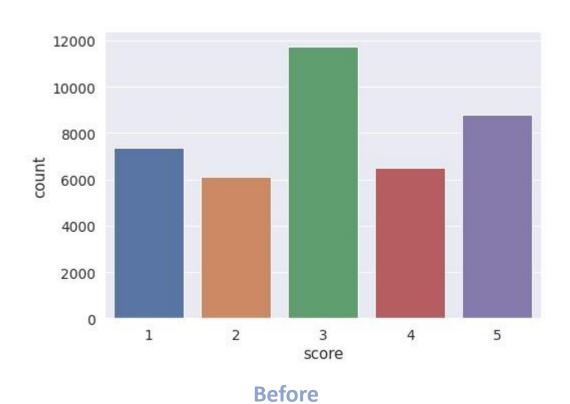


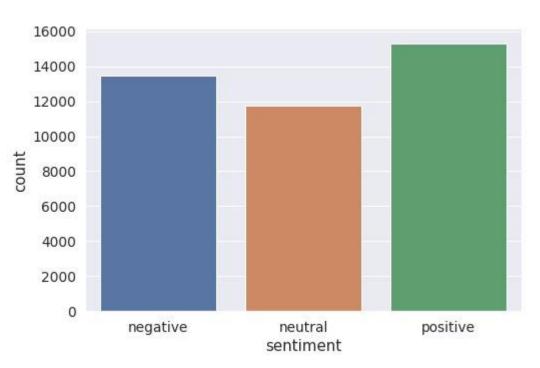


#	columns (total 12 col Column	Non-Null Count	Dtype
	202222		
0	reviewId	40546 non-null	object
1	userName	40546 non-null	object
2	userImage	40546 non-null	object
3	content	40518 non-null	object
4	score	40546 non-null	int64
5	thumbsUpCount	40546 non-null	int64
6	reviewCreatedVersion	33389 non-null	object
7	at	40546 non-null	object
8	replyContent	7943 non-null	object
9	repliedAt	7943 non-null	object
10	sortOrder	40546 non-null	
11	appId	40546 non-null	object

•	reviewId	0
	userName	0
	userImage	0
	content	28
	score	0
	thumbsUpCount	0
		7457



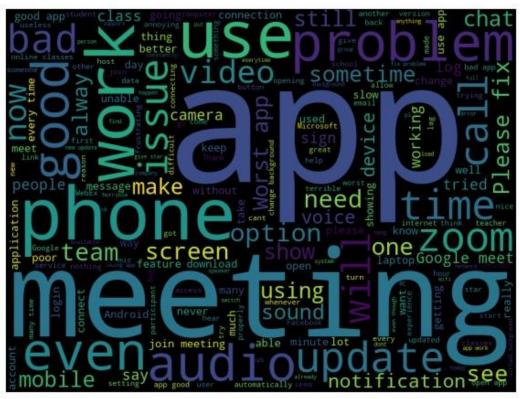




After







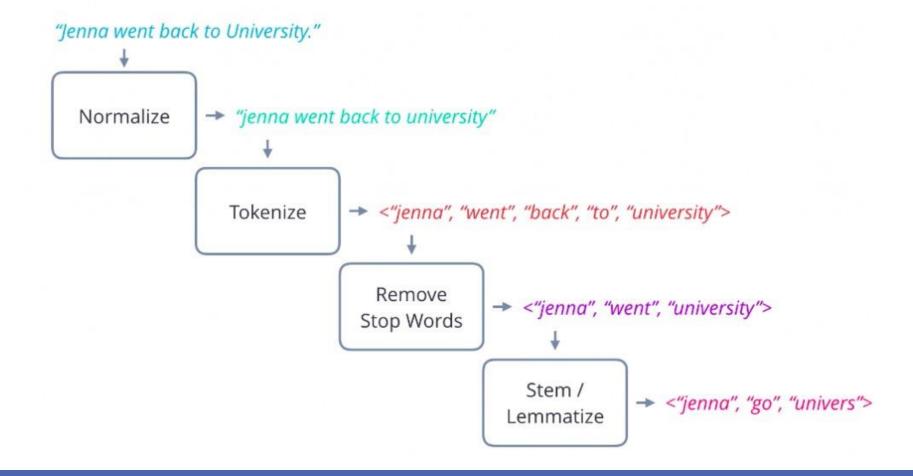
Positive Negative

Data Preprocessing



Data Preprocessing





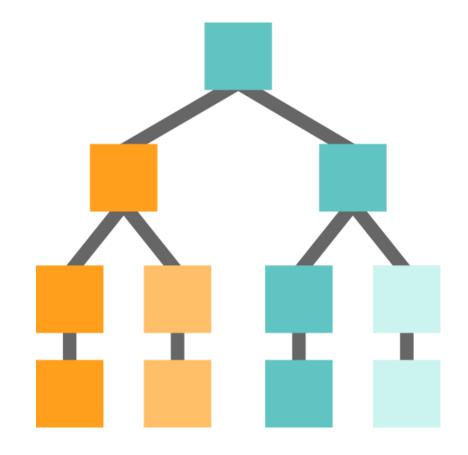
Data Preprocessing



```
1 print(f"Orignal Text : {data.review[13]}")
2 print()
3 print(f"Preprocessed Text : {preprocess_text(data.review[13])}")

Orignal Text : I like the mobility the app provides. however, it does tend to be a bit glitchy

Preprocessed Text : ['i', 'like', 'mobil', 'app', 'provid', 'howev', 'tend', 'bit', 'glitchi']
```





Count Vectorizer:

- Count the occurrences of each word in the corpus
- Doesn't care about order
- Easy to generate and understand

Color	Red	Yellow	Green
Red			
Red	1	0	0
Yellow	1	0	0
Green	0	1	0
Yellow	0	0	1



TF-IDF:

- Focus on unique words in the corpus
- Give words/n-gram a weight based on it's frequency
- very useful in many domains

$$w_{x,y} = tf_{x,y} \times log(\frac{N}{df_x})$$

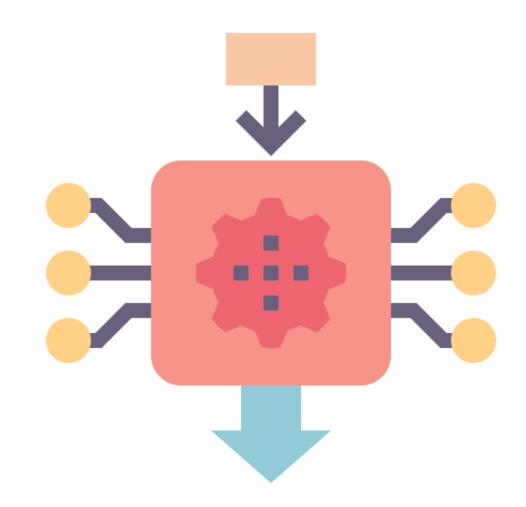


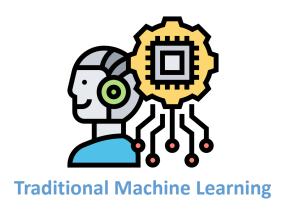
tf_{x,y} = frequency of x in y
df_x = number of documents containing x
N = total number of documents



Word2vec(pretrained):



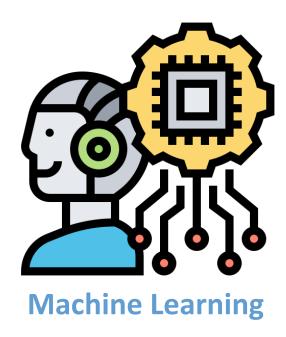




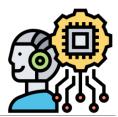




- logistic regression
- SVM
- Naïve Bayes

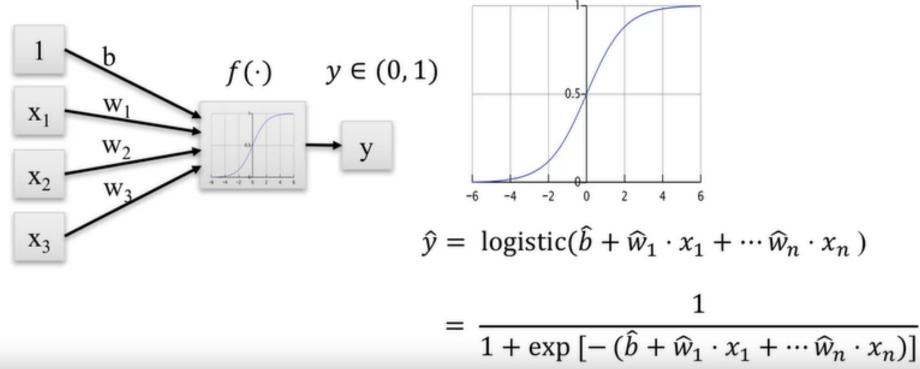




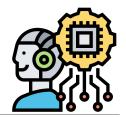


logistic regression(Architecture):

Input features



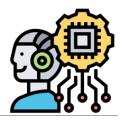




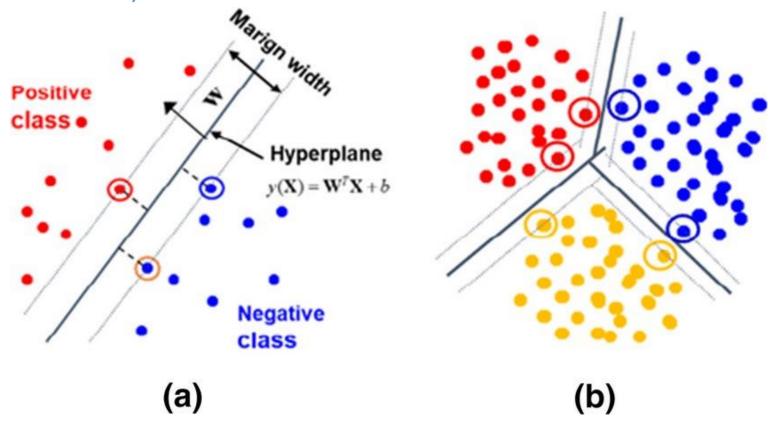
logistic regression (Evaluation):

	precision	recall	f1-score	support	
0	0.67	0.66	0.66	2693	
1	0.52	0.38	0.44	2350	
2	0.66	0.80	0.72	3061	
accuracy			0.63	8104	
macro avg	0.62	0.61	0.61	8104	
weighted avg	0.62	0.63	0.62	8104	
accuracy: 0.6	532				

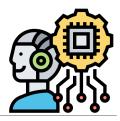
Models Machine Learning:



SVM(Architecture):







SVM(Evaluation):

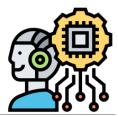
	precision	recall	f1-score	support	
0	0.59 0.47	0.72 0.25	0.65 0.32	2693 2350	
2	0.64	0.76	0.70	3061	
accuracy macro avg weighted avg	0.57 0.57	0.57 0.60	0.60 0.56 0.57	8104 8104 8104	

SVM with word2vec

р	recision	recall	f1-score	support	
0	0.70	0.72	0.71	2693	
1	0.59	0.48	0.53	2350	
2	0.71	0.79	0.75	3061	
accuracy			0.68	8104	
macro avg	0.67	0.67	0.66	8104	
weighted avg	0.67	0.68	0.67	8104	
accuracy: 0.680					

SVM with TF-IDF

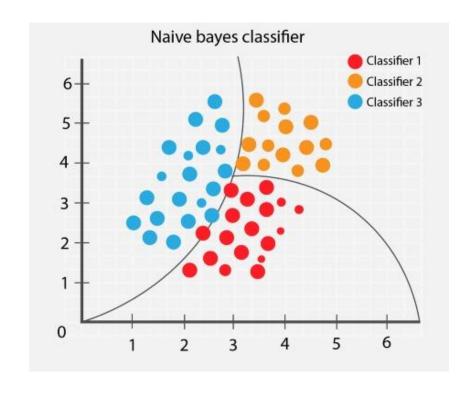




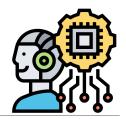
Naïve Bayes (Architecture):

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

using Bayesian probability terminology, the above equation can be written as



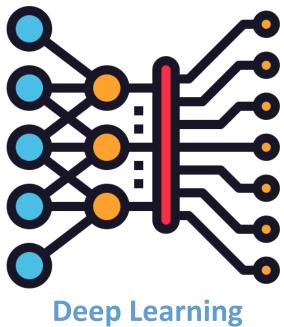




Naïve Bayes (Evaluation):

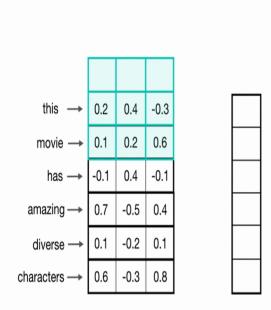
	Naive Bayes REPORT.				
	precision	recall	f1-score	support	
ø	0.73	0.61	0.66	4790	
1	0.30	0.56	0.39	1943	
2	0.78	0.66	0.71	5423	
accuracy			0.62	12156	
macro avg	0.60	0.61	0.59	12156	
weighted avg	0.68	0.62	0.64	12156	

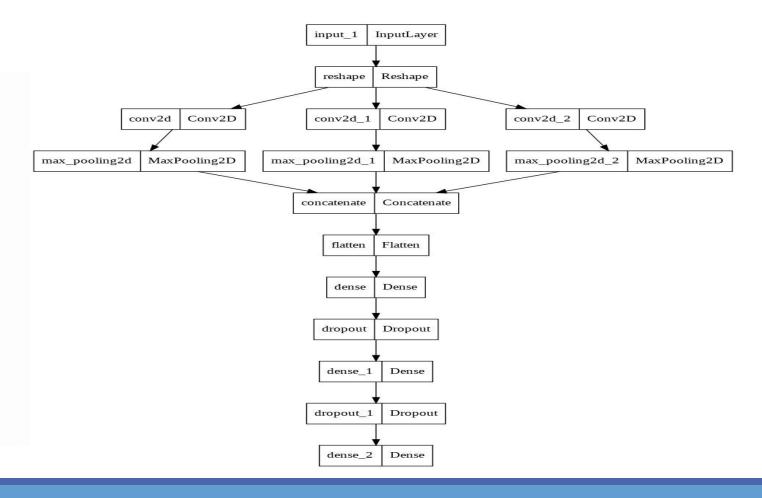
- CNN
- GRU
- LSTM

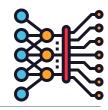






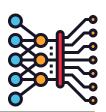






CNN (Evaluation):

	precision	recall	f1-score	support	
0 1 2	0.56 0.41 0.74	0.67 0.41 0.65	0.61 0.41 0.69	2219 2381 3504	
accuracy macro avg weighted avg	0.57 0.59	0.58 0.58	0.58 0.57 0.59	8104 8104 8104	



GRU (Architecture):

GRU

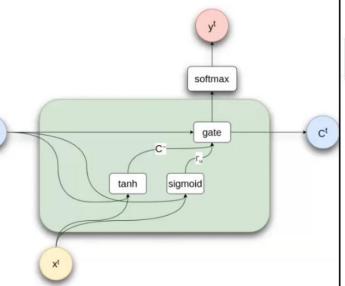
Sequential Modeling

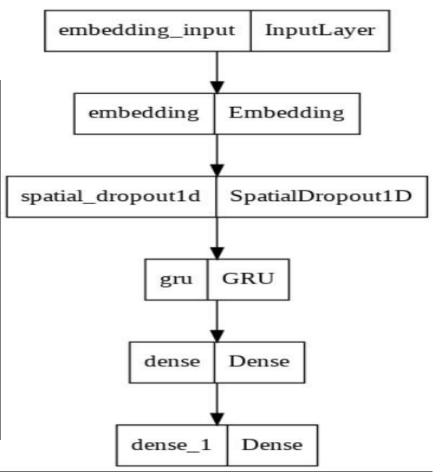
$$\tilde{C} = tanh(W_c[C^{t-1}, X^t] + b_c)$$

$$\Gamma_u = \sigma(W_u[C^{t-1}, X^t] + b_u)$$

$$C^t = \Gamma_u * \tilde{C} + (1 - \Gamma_u) * C^{t-1}$$

.



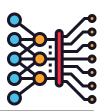




GRU (Evaluation):

	precision	recall	f1-score	support
0	0.65	0.75	0.70	2210
1	0.53	0.38	0.44	1873
2	0.70	0.76	0.73	2521
accuracy			0.65	6604
macro avg	0.63	0.63	0.62	6604
weighted avg	0.64	0.65	0.64	6604

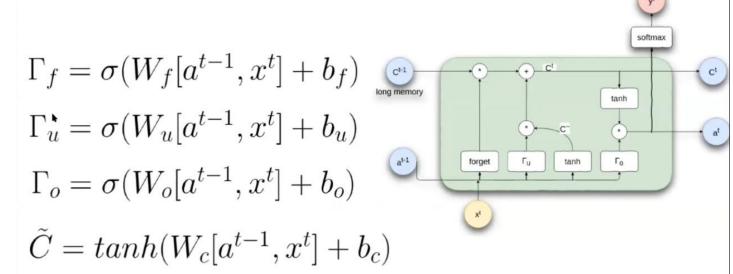
Models Deep Learning:

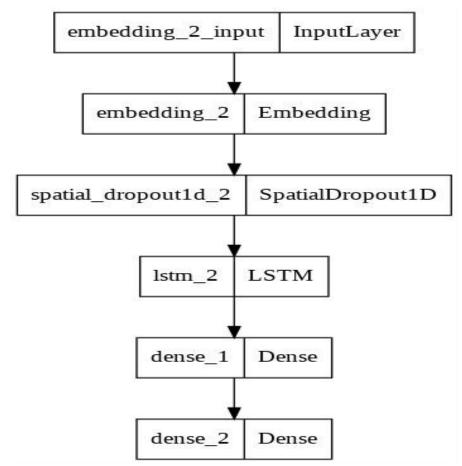


LSTM (Architecture):

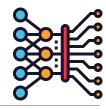
LSTM

Sequential Modeling





Models Deep Learning:



LSTM (Evaluation):

	precision	recall	f1-score	support	
0	0.70	0.68	0.69	960	
1	0.55	0.42	0.48	830	
2	0.65	0.79	0.72	1004	
accuracy			0.65	2794	
macro avg	0.63	0.63	0.63	2794	
weighted avg	0.64	0.65	0.64	2794	

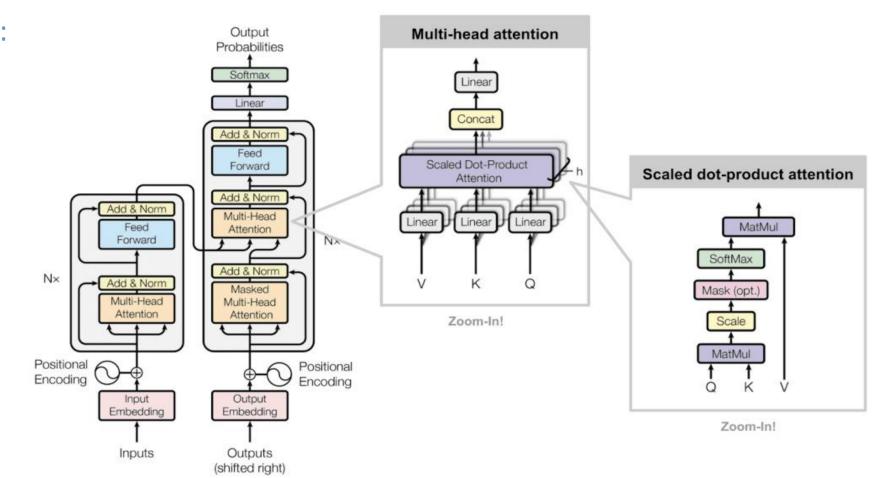
Transformer (pretrained):
BERT







Transformers (Architecture):

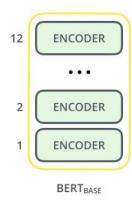


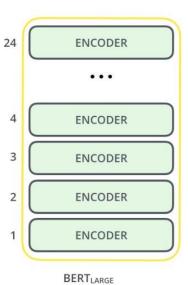


Transformers: BERT (Architecture)

- BERT-Base, Uncased: 12-layer, 768-hidden, 12-heads, 110M parameters.
 - BERT-Large, Uncased: 24-layer, 1024-hidden, 16-heads, 340M parameters.
 - BERT-Base, Cased: 12-layer, 768-hidden, 12-heads, 110M parameters.
 - BERT-Large, Cased: 24-layer, 1024-hidden, 16-heads, 340M parameters.

Bidirectional Encoder Representations from Transformers

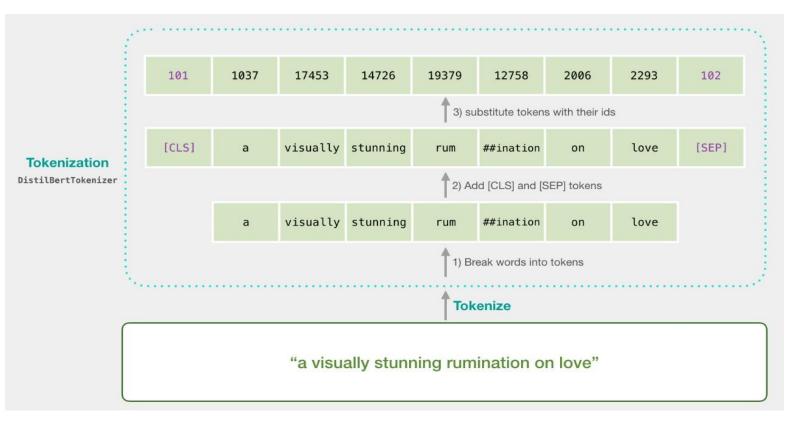






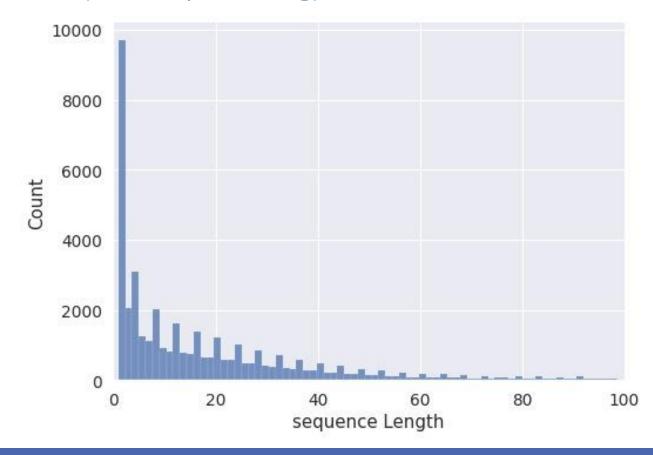
Transformers: BERT (Text Preprocessing)

```
print(len(encoding['input_ids'][0]))
encoding['input ids'][0]
tensor([ 101, 2296, 7292, 1045, 8833, 1999, 1012, 1012, 1012, 1996,
     10439, 5466, 2033, 2041, 1998, 2360, 1000, 2115, 5219, 2038,
     3092, 3531, 8833, 1999, 2153, 1000, 1045, 26618, 1999,
     1019, 2335, 1999, 5216, 1012, 1045, 2109, 2023, 10439,
     1016, 2086, 1998, 2023, 1996, 2034, 2051, 2000, 2227,
     3291, 1012, 102,
print(len(encoding['attention_mask'][0]))
encoding['attention mask']
0, 0, 0, 0, 0, 0, 0, 0]])
```





Transformers: BERT (Text Preprocessing)



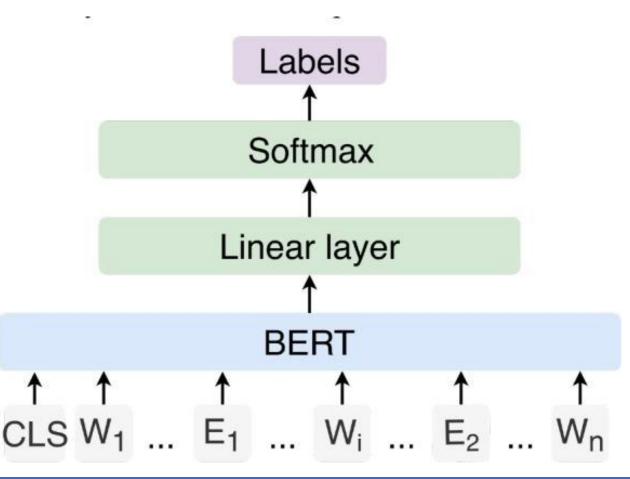


Transformers: BERT (Fine Tuning)

• Batch size: 16, 32

• Learning rate (Adam): 5e-5, 3e-5, 2e-5

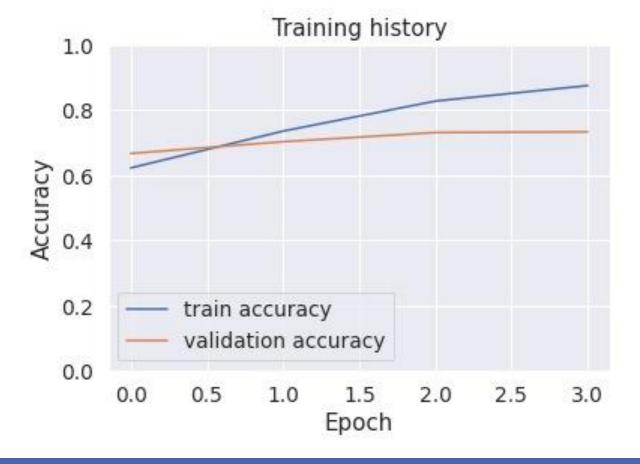
• Number of epochs: 2, 3, 4





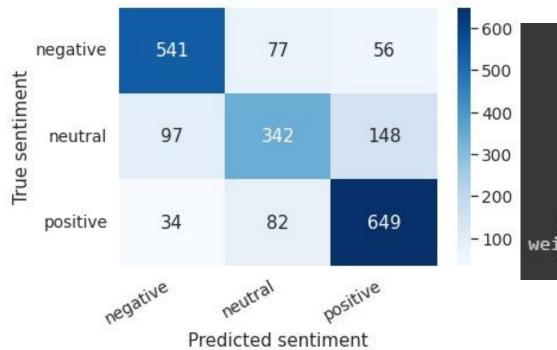
Transformers: BERT (Training)

```
Epoch 1/4
Train loss 0.8296178524170006 accuracy 0.6224976690615917
    loss 0.7464140923693776 accuracy 0.6668311944718657
Epoch 2/4
Train loss 0.6353758487915784 accuracy 0.7353973564416169
    loss 0.7077576396986842 accuracy 0.702862783810464
Epoch 3/4
Train loss 0.4508616898993128 accuracy 0.8273734437558273
     loss 0.7428961358964443 accuracy 0.7309970384995064
Epoch 4/4
Train loss 0.3412153007626011 accuracy 0.8748423188723743
    loss 0.8411004249937832 accuracy 0.7329713721618953
CPU times: user 48min 18s, sys: 18min 19s, total: 1h 6min 38s
Wall time: 1h 6min 39s
```





Transformers: BERT (Evaluation)



	precision	recall	f1-score	support
negative	0.81	0.80	0.80	674
neutral	0.68	0.58	0.63	587
positive	0.76	0.85	0.80	765
accuracy			0.76	2026
macro avg	0.75	0.74	0.74	2026
weighted avg	0.75	0.76	0.75	2026

What is the Best Model?



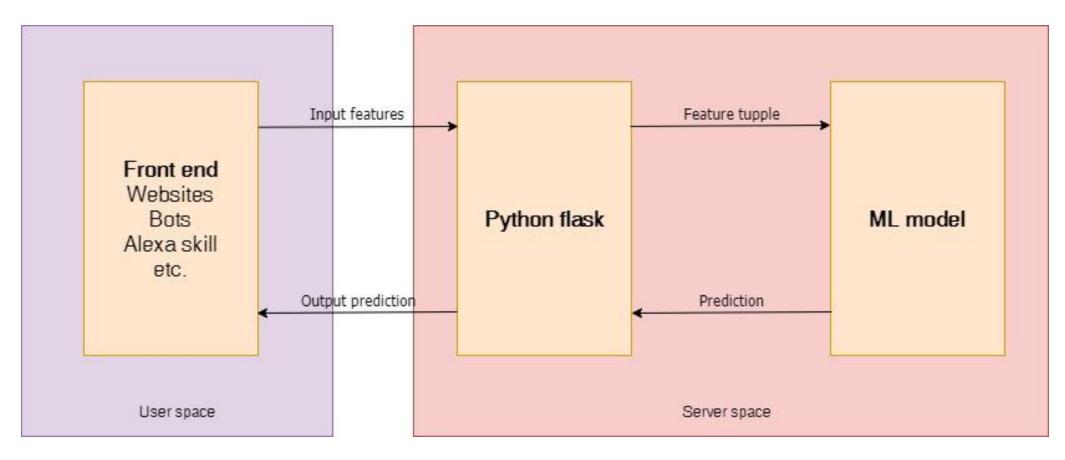
Model	Туре	Accuracy	F1 Score				
				C	£ 1	Pros	Cons
SVM	Traditional ML	0.68	0.75	0.53	0.71	FastLow Resources	 Low performance with large data
LSTM	DL from scratch	0.65	0.72	0.48	0.69	 Represent sequential data Saving short and long history 	 No parallelization Slow Vanishing gradient
BERT T	Transfer Learning	0.76	0.80	0.63	0.80	 Parallelization Good result 	 large data high resource

API

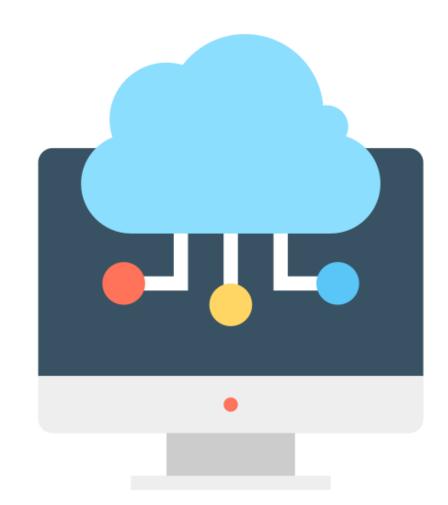


API



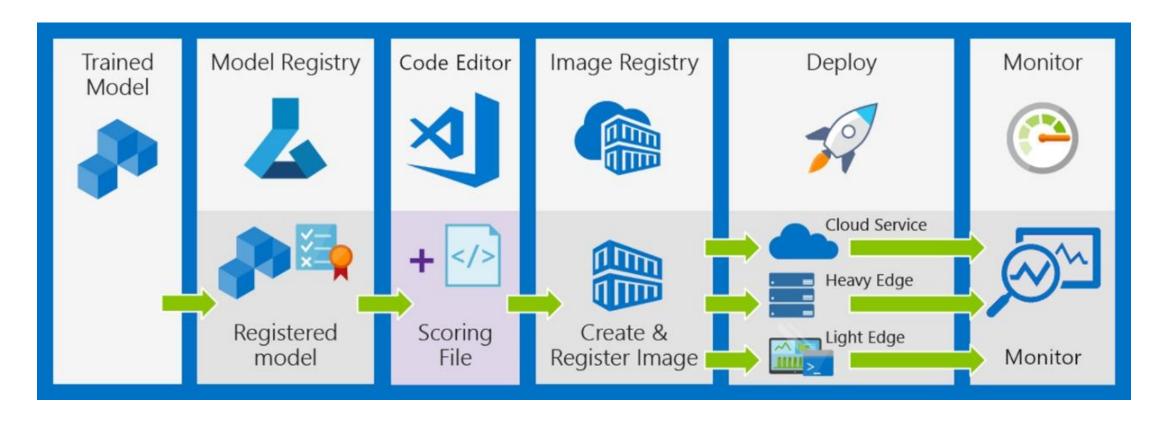


Deployment



Deployment





Future Work



- Scrapping data automatically and store it at database
- Chatbot



https://google-play-reviews-sentiment-analysis.azurewebsites.net/

Papers:

- Pre-training of Deep Bidirectional Transformers for Language Understanding
- Deep Learning for Sentiment Analysis on Google Play Consumer Review
- Study of Sentiment on Google Play Store Applications
- Attention Is All You Need



