

Bag of N-Gram Document Classification

Jyotirmoy Mohapatra

10/10/2018

[GitHub Repo](#)

1 Comparing Performance by tuning N-Gram size, Optimizer and Learning Rate

In this section, an experiment was conducted with several combinations of N-gram size, Optimizer and learning rate in order to find the combination that gave the best validation accuracy. The tokenization scheme used here removes the punctuation and converts every token to lowercase. Other hyperparameters that were kept constant were : Batch Size= 32, Vocabulary Size = 10000, Embedding Size = 100.

Hyperparameter Comparison				
N-Grams	Optimizer	Learning Rate	Train Acc.	Validation Acc.
1	Adam	0.01	99.445	78.62
1	Adam	0.001	96.155	83.78
1	SGD	0.01	64.785	65.06
1	SGD	0.1	71.095	70.34
1	SGD	$0.1(\gamma = 0.03/\text{epoch})$	66.77	67.54
1	Adam	$0.001(\gamma = 0.03/\text{epoch})$	90.135	83.76
1	Adam	$0.001(\gamma = 0.17/\text{epoch})$	92.75	85.18
2	Adam	0.001	94.225	85.02
2	Adam	$0.001(\gamma = 0.03/\text{epoch})$	88.435	84.6
2	Adam	$0.001(\gamma = 0.17/\text{epoch})$	91.145	85.66
3	Adam	0.001	93.29	85.12
3	Adam	$0.001(\gamma = 0.03/\text{epoch})$	87.995	83.96
3	Adam	$0.001(\gamma = 0.17/\text{epoch})$	92.56	85.24
4	Adam	0.001	92.225	84.76
4	Adam	$0.001(\gamma = 0.03/\text{epoch})$	87.543	84.54

From the table above, it was observed that Adam optimizer performs considerably better than the Stochastic Gradient Descent (SGD). SGD with a lower learning rate tends to overshoot and does not converge well. Hence, a high learning rate is better for SGD optimizer in this case. Adam optimizer performs better with a lower learning rate. Thus, in order to further improve the performance with Adam, learning rate annealing was performed with different decay rates. It was observed that a decay rate of 0.5 every 3 epochs performs better than a decay rate of 0.1 every 3 epochs. Secondly, different n-gram sizes were experimented upon. In this experiment, a unigram model consisted of single word tokens, a bi-gram model consisted of both single word tokens and 2 word tokens and a similar approach was used for tri-gram and 4-gram models. The unigram, bi-gram and tri-gram models gave comparable accuracy but there was a slight dip in accuracy for the 4-gram model. The Bi-gram model with adam optimizer and learning rate of 0.001 with decay rate of 0.5 every 3 epochs was used for conducting further experiments.

2 Exploring different Tokenization Schemes

To further improve the performance of the model, three tokenization schemes were explored. The first tokenization scheme involved no pre-processing of tokens. The second tokenization scheme removed the punctuations and converted every token to lowercase. The third tokenization scheme replaced each token with its lemma and also removed few HTML tags.

Tokenization Scheme Comparison			
Tokenization Scheme	N-Grams	Train Acc.	Validation Acc.
No Pre-Processing	1	92.16	83.84
No Pre-Processing	2	90.005	84
No Pre-Processing	3	89.5	83.84
Lower Case and no punctuation	1	92.75	85.18
Lower Case and no punctuation	2	91.145	85.66
Lower Case and no punctuation	3	92.56	85.24
Stemming	1	91.41	83.4
Stemming	2	89.675	83.74
Stemming	3	88.445	83.02

The tokenization scheme with punctuation removal gave the best results. The stemming tokenization scheme reduced the dictionary size and gave more meaningful vocabulary, however, it did not perform as well as expected. One possible reason for this might be because the punctuation were not removed in this scheme which might have diluted the quality of the vocabulary.

3 Exploring Other Hyperparameters

Further experiments were performed using the Bi-gram model with Adam optimizer, a learning rate of 0.001 and decay factor of 0.5 every 3 epochs along with tokenization scheme which removed punctuation and converted every token to lowercase.

3.1 Vocabulary Size and Embedding Size

Vocabulary and Embedding size Comparison			
Vocab Size	Embedding Size	Train Acc.	Validation Acc.
10000	100	91.145	85.66
10000	300	92.87	85.28
20000	100	94.32	85.14
20000	300	95.595	85.38
40000	100	95.42	85.18
40000	300	95.445	85.41

Varying the embedding size and the vocabulary size did not really have a huge impact on the model's performance. Increasing the vocabulary size and embedding size increased the number of computation and thus increased the overall training time.

3.2 Batch Size

Batch size Comparison		
Batch Size	Train Acc.	Validation Acc.
16	92.48	85.06
32	91.145	85.66
64	89.155	84.72

Batch sizes of 16 and 32 had comparable performance. However, increasing the batch size to 64 resulted in dip in accuracy.

4 Comparison Graphs

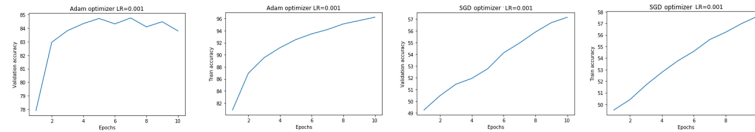


Figure 1: Adam vs SGD Optimizer.

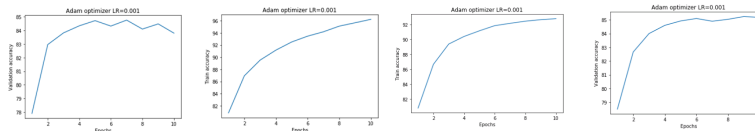


Figure 2: Learning Rate annealing

5 Correct and Incorrect Review

5.1 Correct

vincent junk plays the part of paul an ex con junk to an office job where he meets junk a secretary who is junk junk when she has her hearing junk in junk junk when not played by junk junk together they help each other to develop as junk /i/br /iwhat was particularly interesting about this film was the junk of the characters not fitting into obvious stereotypes paul appears junk in the office environment is it that he 's just not cut out for work this belief is junk when he gets a job in a bar and junk /i/br /ithe film has a certain junk which i find refreshing and showed how easy is to act junk even if we think it is junk or junk /i/br /ifinally it is a film full of great junk both touching and humorous one is when junk is junk and is trying to junk a screaming baby she continues to junk it but takes her hearing junk out for her own junk junk junk plays the the part part of junk junk junk junk junk to an junk junk junk where he he meets junk junk junk junk who is junk junk junk when she she has -positive

i can agree with other comments that there was n't an enormous amount of history junk in the movie but it was n't a documentary it was meant to entertain and i think it did a very good job at it./i/br /i/i agree with the black family the scenes with them seemed out of place like all of a sudden it would be thrown in but i did catch on to the story and the connection between the families later on and found it pretty good./i/br /i/despite it was n't a re junk of the 60s it did bring into the light very big and important junk junk of the decade i found it very entertaining and worth my while to watch i can junk agree with junk junk junk that there there was was n't junk junk junk amount of junk junk junk in the the movie movie but but it it was was n't n't a a documentary junk it was was meant meant to junk junk and i i think think it it did did a a very very good good job junk junk it./i/br /i/br /i/i agree with with the the black junk junk the scenes scenes with with them junk junk out of of place junk like all all of of a junk junk it would would be junk thrown in junk but i i did junk junk on to to the the story story and and the junk -positive

for those of you that do n't that reference junkl was 4 junkl hitting one body ... jbr /l jbr junkl onto the junkl /l jbr /l i miss junkl saturday night some of my favorite wrestling moments took place on this stage i remember watching stunning steve junkl rick junkl brian junkl junkl jack junkl junkl johnny b. junkl junkl in his junkl days lord steven junkl junkl heat junkl junkl junkl ... junkl be here a while junkl everyone point is junkl had an awesome junkl in the pre junkl days and they were producing entertaining television junkl junkl on commentary in it 's later years gave me a whole new reason to watch when i started smoking junkl as a teenager ... i really wish junkl would put him on the junkl for a show or two maybe at the next great american junkl they junkl here comes junkl junkl he was junkl /l jbr junkl for those those of of you you that that do do n't n't that junkl junkl junkl junkl junkl junkl junkl junkl junkl junkl ... jbr jbr /l jbr junkl junkl onto the junkl junkl /l jbr /l i junkl junkl junkl saturday night junkl some of of my my favorite junkl junkl junkl junkl junkl on this junkl junkl i remember junkl junkl junkl junkl -positive

5.2 Incorrect

sitting down to watch the junkl season of the junkl on the junkl of love i knew i would be in for an interesting time i had watched some of the previous seasons of the junkl in passing watching an episode or two and missing the next three or so i find that the junkl is often appealing and intriguing though its quality and morality are often junkl /l jbr junkl the junkl of love details the journey taken by jake a junkl year old commercial pilot from junkl texas to find true love as true a love as one can find in a season long reality drama junkl show jake meets 25 beautiful girls from all over the country he begins to get to know them a bit but it is mostly superficial how well can you get to know someone in a few 5 minute conversations jake tries to make his true intentions known from the very beginning at least to the audience he noted that he does n't just want love or a good time but he wants a junkl or wife we can only assume that he has made this clear to the women in the competition -negative

anyone who has a remote interest in science fiction should start at the junkl everyone says star wars and star trek are the best science fiction films to begin at which is fine but the truth is the junkl and this movie junkl green are far better choices than those series junkl is probably science fiction 's best kept secret it remains one of the biggest yet most forgotten films but the impact of its setting is becoming more a reality with each passing day junkl heston junkl his role yet it works edward junkl junkl in his final role makes the most out of it in junkl green more than anyone else and his final scenes are junkl /l jbr /l it is manhattan in junkl the world is junkl and food is an unbelievable fortune a small junkl of junkl junkl costs junkl a big executive for the junkl company is murdered and police detective junkl is on the junkl /l jbr /l the secret of junkl green is not a mystery if you do research on the movie junkl is enjoyable to watch but the whole screenplay is a joke it is just as cheap as the entire production the screenplay -negative

an junkl although repetitive and rather junkl junkl of exploitation and junkl in a situation where there is no way forward or up where the attempts to make yourself feel better by junkl and putting down whoever is below you seems to be the only junkl but even here in this junkl junkl of lost dreams and no future that does not work and junkl out to something or someone to junkl and share with a simple act of junkl gives some junkl even if it just makes the present junkl by junkl memories of the junkl /l jbr /l although there is little actual on screen violence this is a harsh and brutal film about the small junkl of junkl junkl and personally that does not make for easy entertainment clearly based on a play with a small cast a junkl more junkl junkl to the general social and political environment would possibly have helped the film to reach a junkl audience junkl junkl junkl junkl junkl junkl junkl junkl junkl in a junkl junkl junkl there is is no no way junkl junkl junkl junkl where the junkl attempts to to make junkl junkl junkl junkl junkl junkl junkl junkl junkl junkl junkl junkl -positive

6 Conclusions

Multiple experiments were conducted to find out a combination of hyper-parameters that gave the best accuracy on the validation set. The following hyper-parameter combination gave the best result: n-grams:2, optimizer: adam, learning rate: 0.001 with 0.5 decay every 3 epochs, vocab size: 10000, embedding size: 100, batch size: 32, tokenization scehme: no punctuation. The above model achieved an accuracy of 85.18 on the Test dataset.