

Predicting sentiment on Amazon review data

https://github.com/spitko/amazon-sentiment-analysis





About

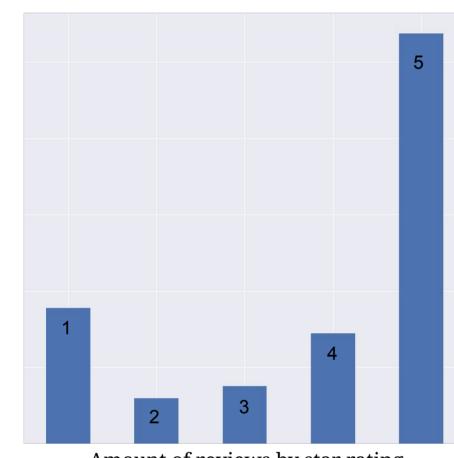
Sentiment analysis is contextual mining of text which identifies and extracts subjective information in source material. One common use case for this is classifying opinions as either positive or negative. Our goal with this project was to predict how a reviewer felt about their purchase based on their Amazon review by developing a model capable of classifying given text's polarity.

Data

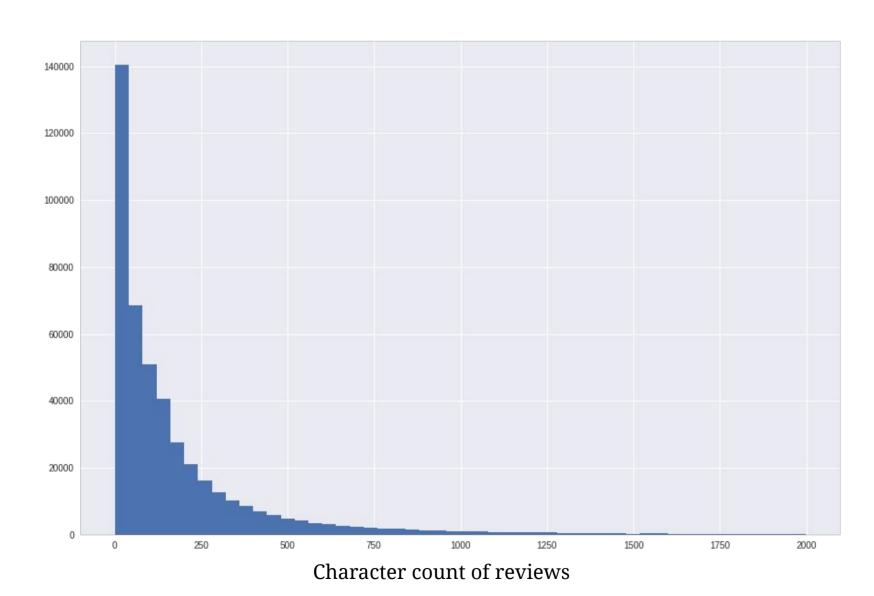
Our data consisted of two separate datasets found on Kaggle, which combined contained data about over 480k different reviews left on different Amazon mobile phone product listings.

After filtering out reviews with over 2000 characters and neutral (3 star) reviews, we were left with 454k rows of data of which 74% were positive (4 or 5 stars).

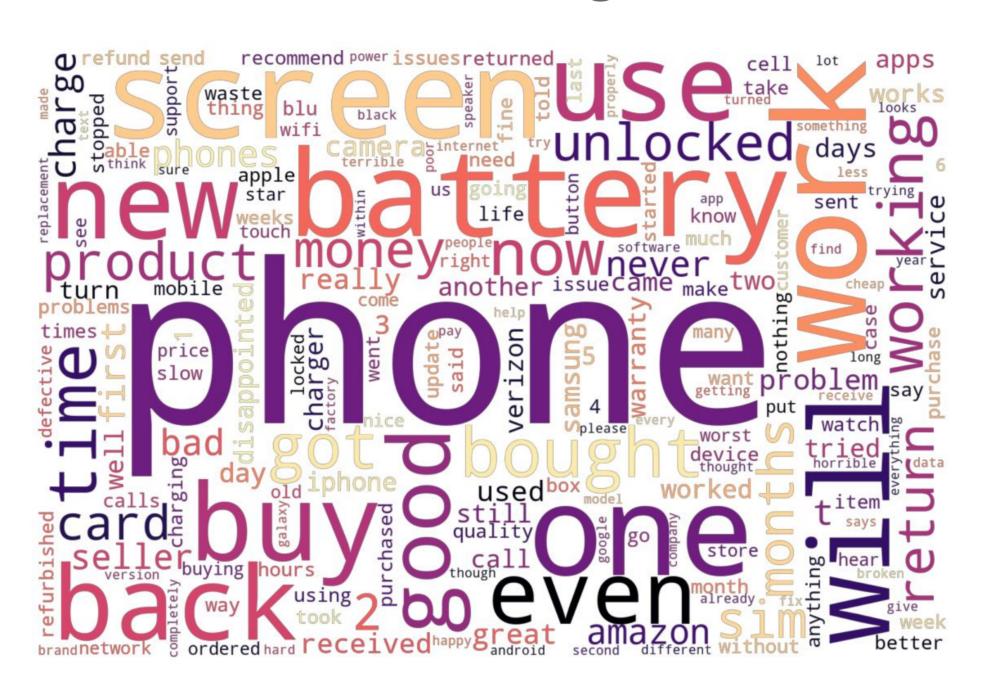
Median number of words in a review was 18.



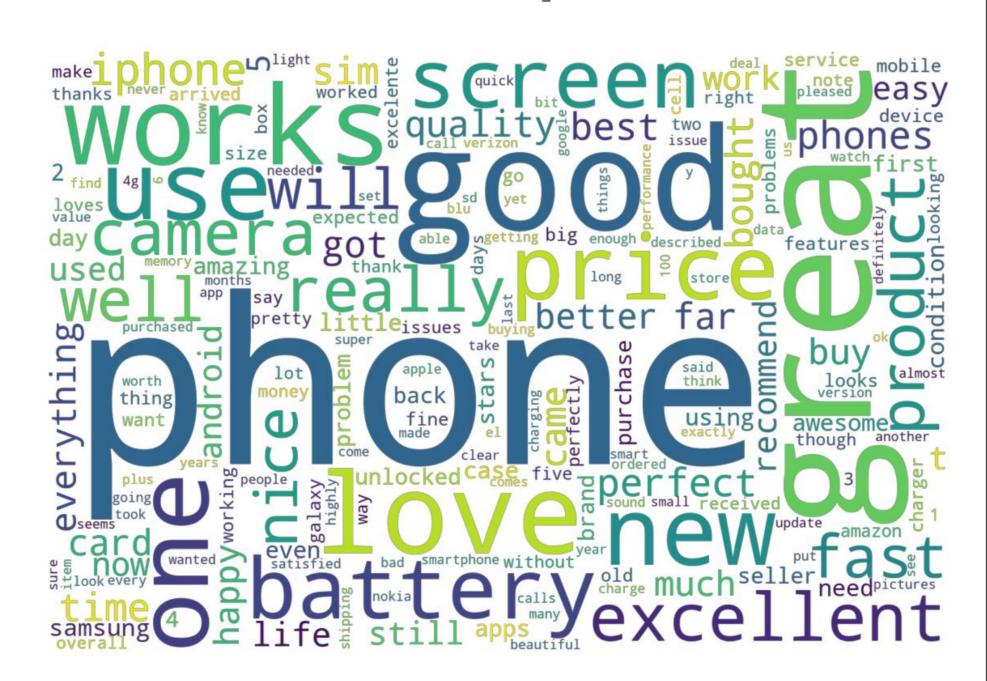
Amount of reviews by star rating



Most common words in negative reviews



Most common words in positive reviews



Preprocessing for ML

As the remaining data was fairly unbalanced with only 26% negative reviews, the first step was to downsample the positive reviews, leaving us with an equal amount of positive and negative samples.

From the remaining data a vocabulary was created containing 20000 most common words and assigning each a numeric index. The review texts were then replaced with sequences (lists) of word indexes.

Model

For machine learning, a convolutional neural network was utilised in TensorFlow using Keras as the frontend.

Given the large amount of data available, no pre-trained word embeddings were used and all embeddings were learned from scratch.

Model was trained on a Kaggle kernel, that included an Nvidia P100 GPU.

Conclusions

Our model achieved 96.8% accuracy on validation data. Considering our set goal to reach 85% accuracy, we can safely consider this project a

			е				
embedding (Embedding)	(None,	200,	250)		:		
dropout (Dropout)			250)	0			
conv1d (Conv1D)							
global_max_pooling1d (Global	(None,	64)		0			
dense (Dense)	(None,	1)		65			
Total params: 5,080,129							
Trainable params: 5,080,129							
Non-trainable params: 0							
Train on 176061 samples, val:	idate o						
Epoch 1/3							
176061/176061 - 118s - loss:	0.1794	- ac	c: 0.9319 - va	l_loss: 0.1	225	- val_acc:	0.957
Epoch 2/3			D 12222			2	
176061/176061 - 113s - loss:	0.0959	- ac	c: 0.9685 - va	l_loss: 0.1	037 ·	- val_acc:	0.966°
Epoch 3/3	0.0625		. 0 0010	1 1000 0 1	000		0.060
176061/176061 - 108s - loss:	0.0035	- ac	c. 0.9810 - va	11_10SS: 0.1	002	- val_acc:	0.968

Model summary and training

Contact information

Samuel Pitko Helen Pung

https://github.com/spitko https://github.com/helenpung