

Duplicate Question Detection
{dougc, shgu}@stanford.edu

Applications: Quora/Slack/Twitter/email

Substantial amount of communication contains Questions.

How to automate answering of questions?

<https://youtu.be/fSpvfPnMXuQ>

Data Set Kaggle Unbalanced dataset

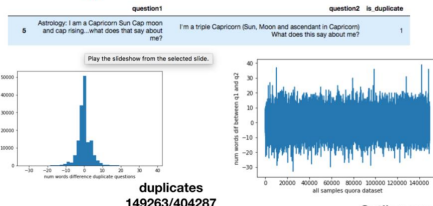


Table with columns: question1, question2, duplicate, and a list of question pairs with their duplicate status.

duplicates 149263/404287 37%
Outliers waste memory
nonduplicates 255024/404287 63%

biLSTM Architectures
Start with basic biLSTM + Linear
Good: no overfitting
Bad: poor test accuracy

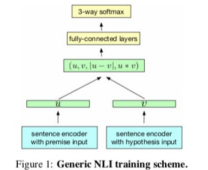


Figure 1: Generic NLI training scheme.
SNLI RTE 2017
What has changed?
Bigger GPUs, V100 16GB

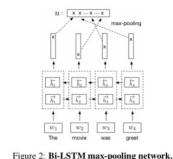


Figure 2: Bi-LSTM max-pooling network.
Duplicate/Not Duplicate
Modified to 2 classes, num_layers > 1

biLSTM Architectures
 $S = W_1 \cdot W_2 \cdot W_3 \dots W_n$
 $\sum_{k=1}^K e^{s_k}$

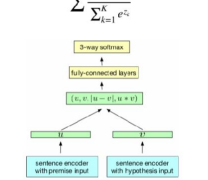


Figure 1: Generic NLI training scheme.

Modified to 2 classes, num_layers > 1
Increase accuracy 16%

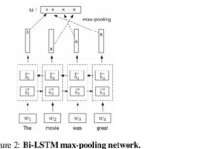


Figure 2: Bi-LSTM max-pooling network.

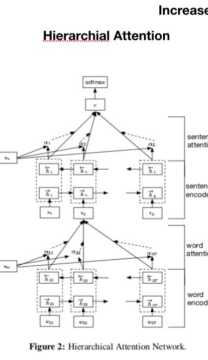


Figure 2: Hierarchical Attention Network.

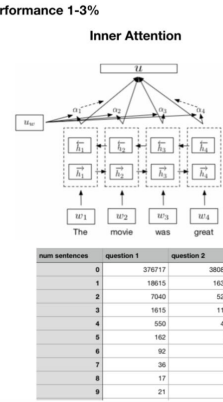
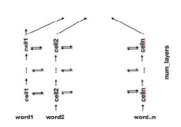


Table with columns: num sentences, question 1, question 2. Shows a distribution of sentence counts for two questions.

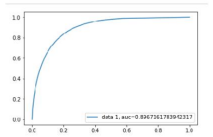
Regularization to Reduce Overfitting

Best performance is with Bi-LSTM @ 83%



Best model used q1,q2, IA + dropout + BN+Tanh 16GB limit

Table with columns: Model, train, dev, test, dev. Lists various models and their performance metrics.

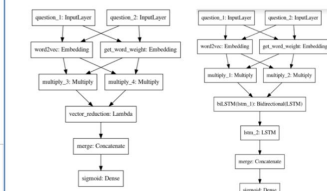


Approach II: Siamese Networks

Siamese Networks using pre-trained word2vec (vocab_size=400k)

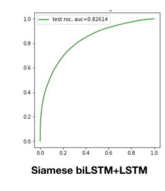
Table with columns: Model, dev_train, dev_dev, accuracy_train, accuracy_dev. Compares different Siamese network architectures.

Table 1: Experiment Results



(a) siamese_v1

(b) siamese biLSTM+LSTM



Siamese biLSTM+LSTM

Further Work

LSTMs are hard to regularize, try Merity's AWD-LSTM/QRNN, IF LSTMs are production worthy for an api call

Different dataset, directed intent from Twitter thread postings

LSTMs don't scale, Transformer architecture scales to multiple GPUs but slow inference times, Bert