

Duplicate Question Detection

Neda Zolaktaf and Vaishnavi Malhotra

1. Problem Statement

Duplicate question detection in forums

- Duplicates pairs are labeled 1, non-duplicates are labeled 0.

Motivation

- Users often ask same question
- Discovering duplicates is valuable for users to get an answer immediately

3. Traditional Models

Two stage approach

- Feature extraction** needs vocabulary, use the frequency of words, or use TFIDF term weighting [1]
- Classifier** Logistic Regression, Random Forest, etc

Example: TFIDF vector + Logistic Regression

Limitations

- Extracted features may not be good, word order can be lost, used different n-grams sizes, linear classifier, no natural notion of similarity between words

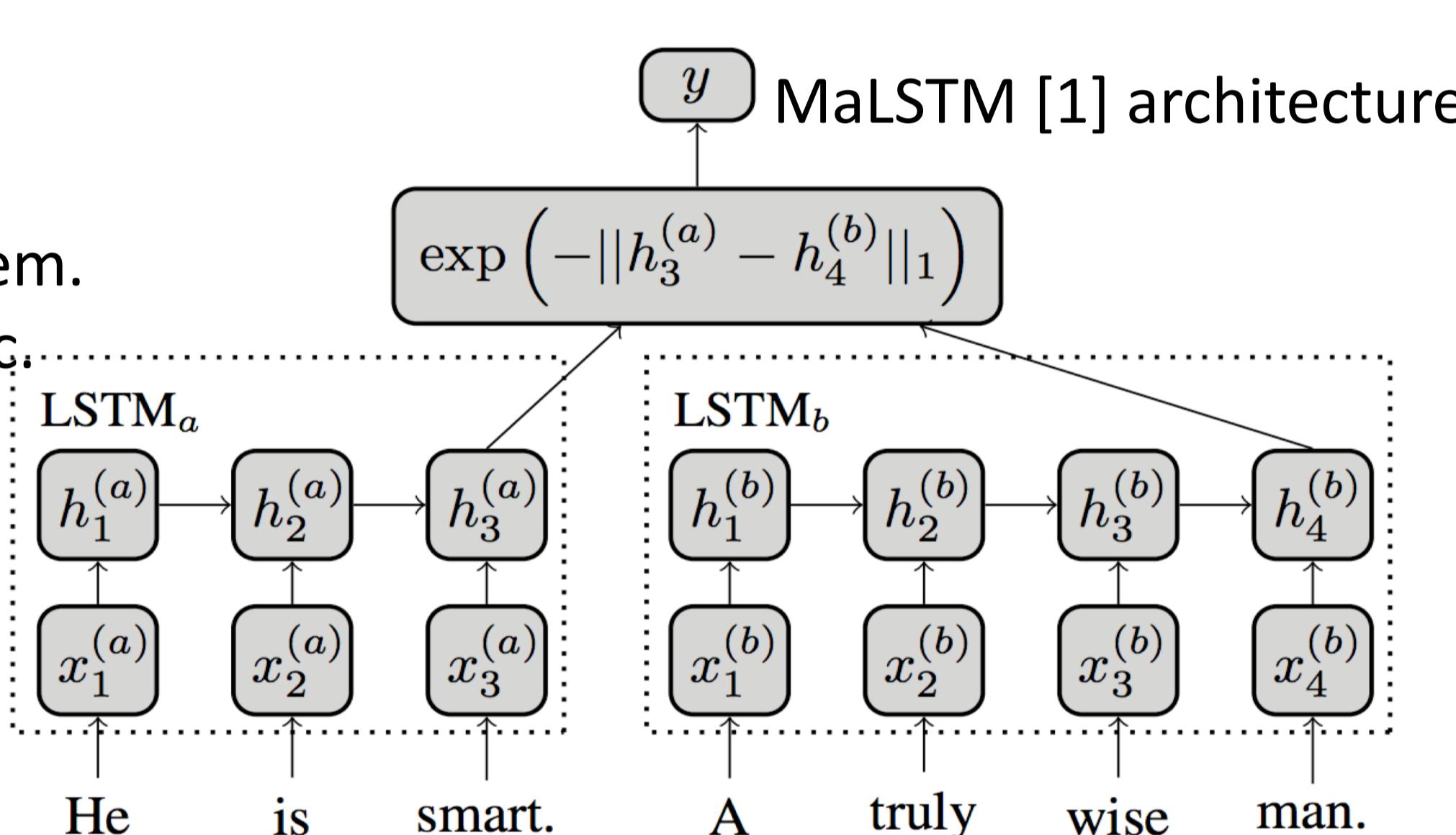
4. Siamese Neural Network

Siamese Network

- networks that have two or more identical sub-networks in them.
- used for sentence similarity, recognizing forged signatures, etc.

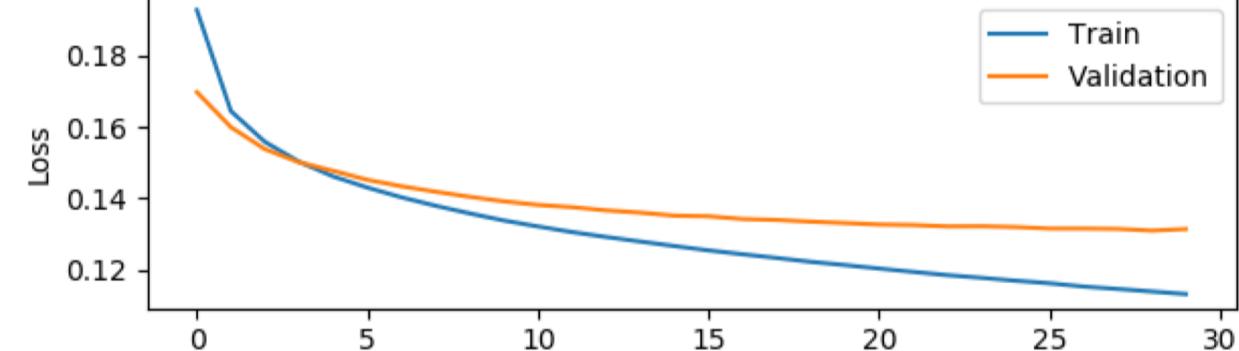
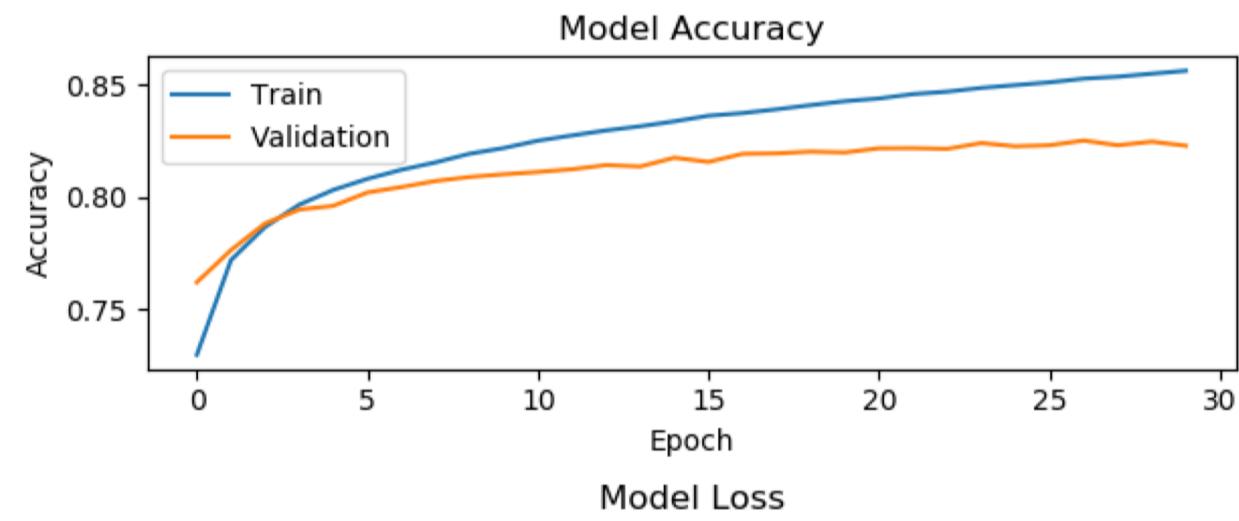
Manhattan LSTM (MaLSTM)[1]

- Long-Term Short-Term Memory (LSTM)
- Use output as semantic representation of question
- Computes sentence similarity using Manhattan distance
- Implemented in Keras
 - Input layer uses word2vec (word embeddings dimension = 300)
 - One-layer LSTM (hidden dimension =50 and 100 for quora)
 - epochs = 5 (30 for quora)
 - batchsize = 1024
 - Adam Optimizer



Objective function is Mean Squared Error.
We tried changing the architecture for cross entropy loss but did not obtain improvements.

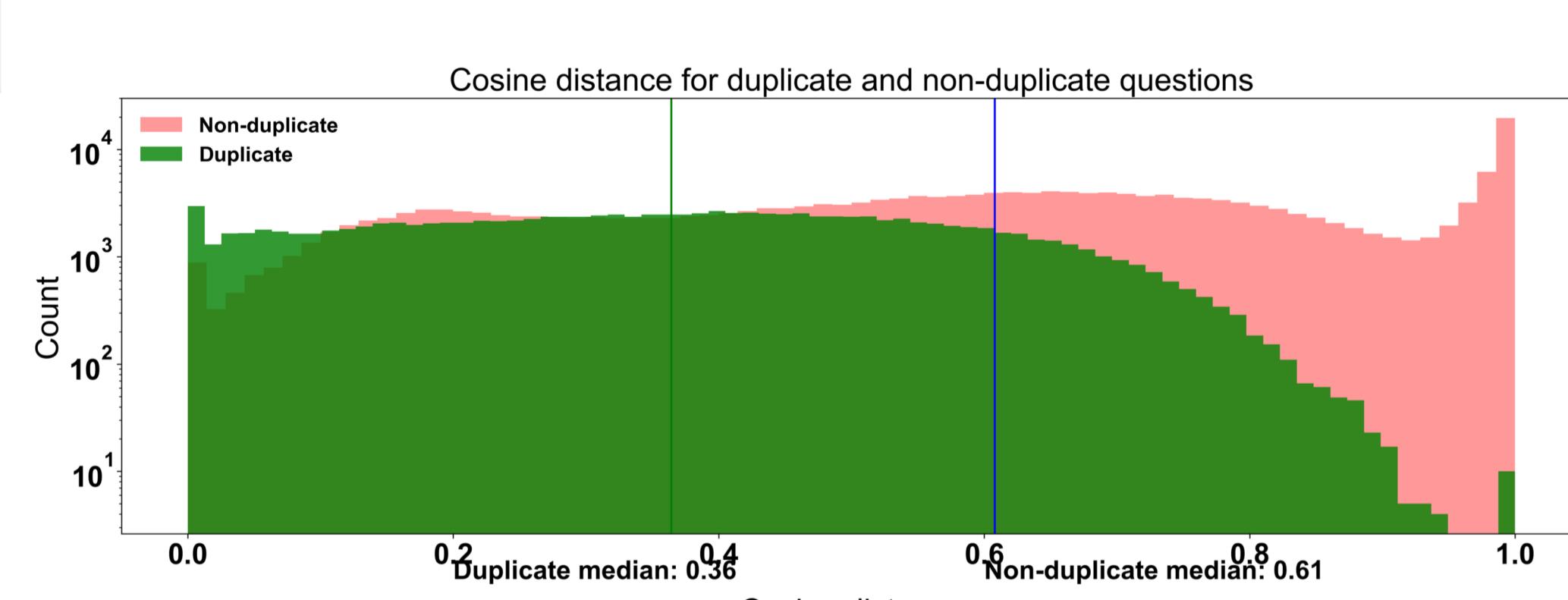
6. Experiments



Evaluating MaLSTM on Quora

Model	Quora		Apple		Android		Sprint		Superuser		AskUbuntu	
	Accuracy	Fscore	Accuracy	Fscore	Accuracy	Fscore	Accuracy	Fscore	Accuracy	Fscore	Accuracy	Fscore
TFIDF + Logistic Regression	0.7554	0.6336	0.9906	0.1161	0.9901	0.1960	0.9901	0.01	0.9902	0.0319	0.9899	0.0091
TFIDF + Random Forest	0.7783	0.6248	0.9911	0.1948	0.9922	0.4335	0.9900	0.0010	0.9909	0.1639	0.9917	0.3370
MaLSTM	0.6756	0.2814	0.991	0.002	0.9899	0.001	0.9917	0.05	0.9904	0.006	0.9900	0.003

Evaluation on individual datasets.



Model	Precision	Recall	Fscore
BOW + Xgboost	0.82	0.61	0.70
Word level TFIDF + Xgboost	0.82	0.62	0.71
n-gram tfidf + Xgboost	0.81	0.41	0.54
Character-level TFIDF + Xgboost	0.84	0.72	0.78

Integrated Quora and AskUbuntu and evaluated on integrated dataset.

7. Interface

- Exploratory Data analysis
- Demo
- Project information

8. Conclusion & Future Work

- Evaluated different methods on 5 datasets
 - On quora TFIDF + Logistic obtains best Fscore on class 1
- Integrated Quora and AskUbuntu datasets into one dataset
 - Evaluated TFIDF at different levels + Xgboost.
 - Character-level TFIDF + Xgboost obtains best Fscore on class 1
- The datasets are highly imbalanced
 - Consider other objective functions
 - Consider alternative neural network models.

9. References

- Jonas Mueller and Aditya Thyagarajan, "Siamese Recurrent Architectures for Learning Sentence Similarity", AAAI, 2016.
- <https://medium.com/mlreview/implementing-malstm-on-kaggles-quora-question-pairs-competition-8b31b0b16a07>
- train.py for [2] based on code from:
<https://github.com/likejazz/Siamese-LSTM>

