



# Debiasing Irrelevant Words in Natural Language Processing

Shi Feng, Qihang Chen

Institute for Interdisciplinary Information Sciences, Tsinghua University



交叉信息研究院  
Institute for Interdisciplinary  
Information Sciences

## Introduction

### Motivation:

- There are some leakage features caused by correlation of data when evaluating a model.
- Some people use them to get extraordinary performances in some competitions.
- How to eliminate the impact of leakage features in our model to get generalization?

### Challenge:

- Decouple the leakage features and useful features.
- Compute weights to debias the datasets.

### State of the art:

- Writing style debiasing
- Data frequency debiasing

### Contributions:

- Visualizing the bias of neutral words.
- Training an unbiased sentiment analysis model.
- Give better performance in general cases.

This movie is great and I will watch it again.  
It's a waste of money to go see this movie.  
Seeing this movie made me have deep doubts about the director's level.  
I was watching this movie and wondering if it was going to be a box office winner.  
I've watched this movie three times and I can't get enough of it every time.  
This movie almost made me vomit, and I never want to see it again.

Train

Machine Learning Model (BiLSTM)

Predict

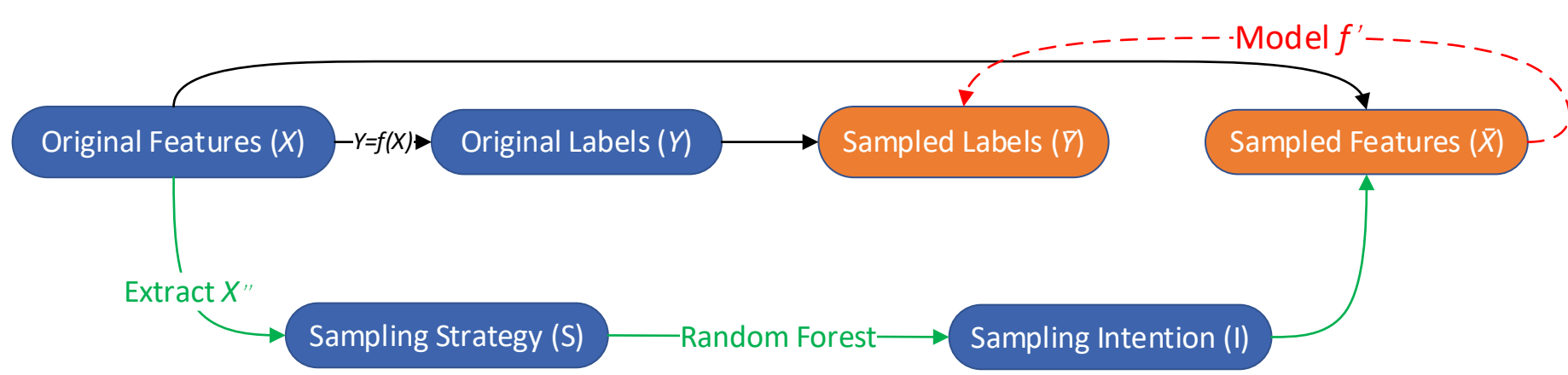
watch = good ,  
see = bad ?  
Wrong !

Scan for more information:



## Method

### Causal explanation of the data generation process:



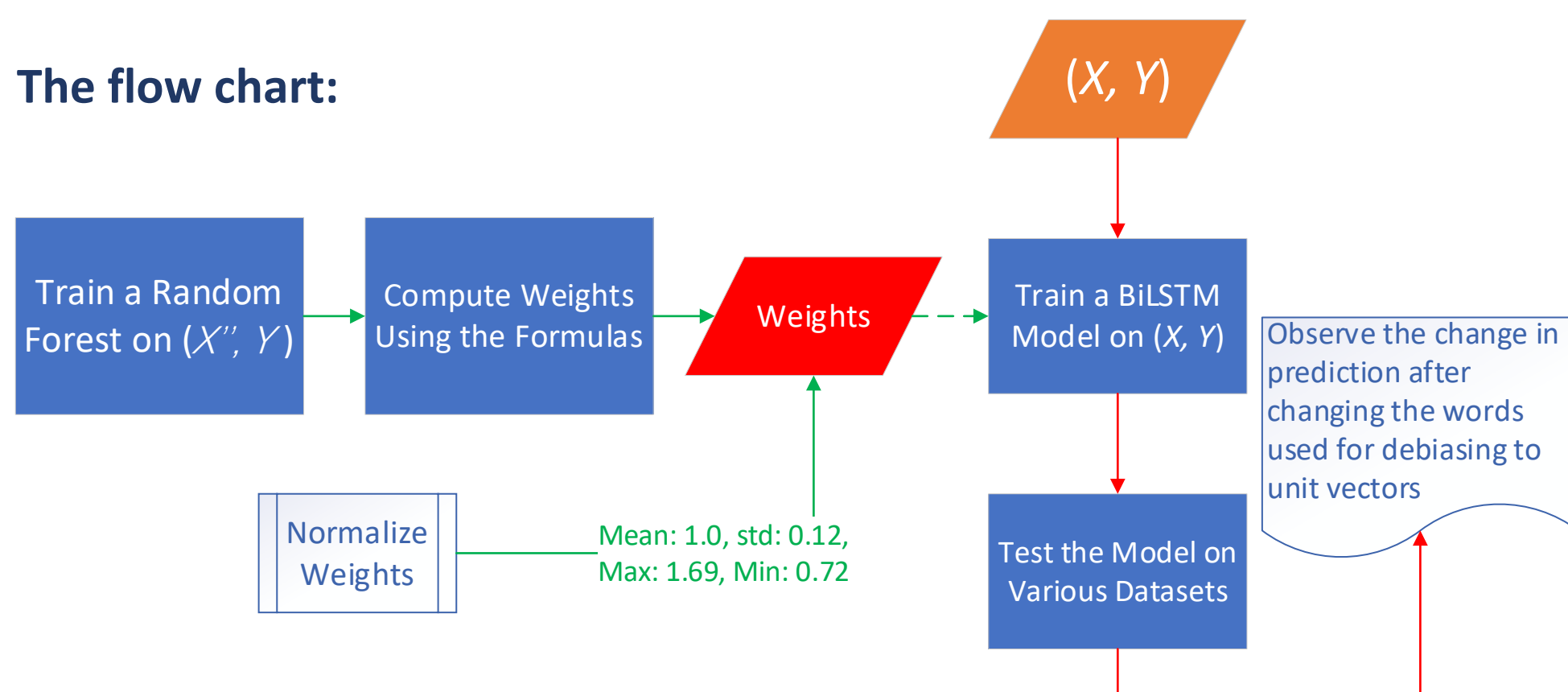
### Estimate sampling intention $I$ for some sampling strategy:

- Predict sentiment labels by random forest, using  $X'' = [\text{word} \in \text{sentence}]$  as features.
- The model is used for estimating the conditional probability relationship:  $P(I = y|s)$ .

### Calculate the weights of samples, using $P(I = y|s)$ :

- $$P(I = y|S) = \frac{P(Y=y)P_{\bar{A}}(Y = 1 - y|S)}{P(Y=0)P_{\bar{A}}(Y = 1|S) + P(Y=1)P_{\bar{A}}(Y = 0|S)}$$
- Normalize the mean of the weights to 1.

### The flow chart:



## Theoretical Analysis

**Theorem:** Suppose the classifier is  $f'$  and the training data  $X$  is divided into two parts,  $X'$  and  $X''$ . Here,  $X''$  is the biased features (leakage features) and  $X'$  is the useful features that fit the practical problem. We denote the sampling intention and strategy as  $I, S$ , then  $w = \frac{P(I=Y)}{P(I=y|s)}$  is an unbiased weight. Or equivalently,

$$E_{x,y,s \sim \bar{A}}[wL(f'(x', g^{-1}(s)), y)] = E_{x,y,s \sim A}[L(f'(x', g^{-1}(s)), y)]$$

Here,  $\bar{A}$  is the biased joint distribution of  $X, Y, S, I$  and  $A$  is the original distribution. Moreover,  $L$  is the loss function and  $g$  is the functional relation between sampling strategy  $S$  and biased features  $X''$ .

## Experimental Settings

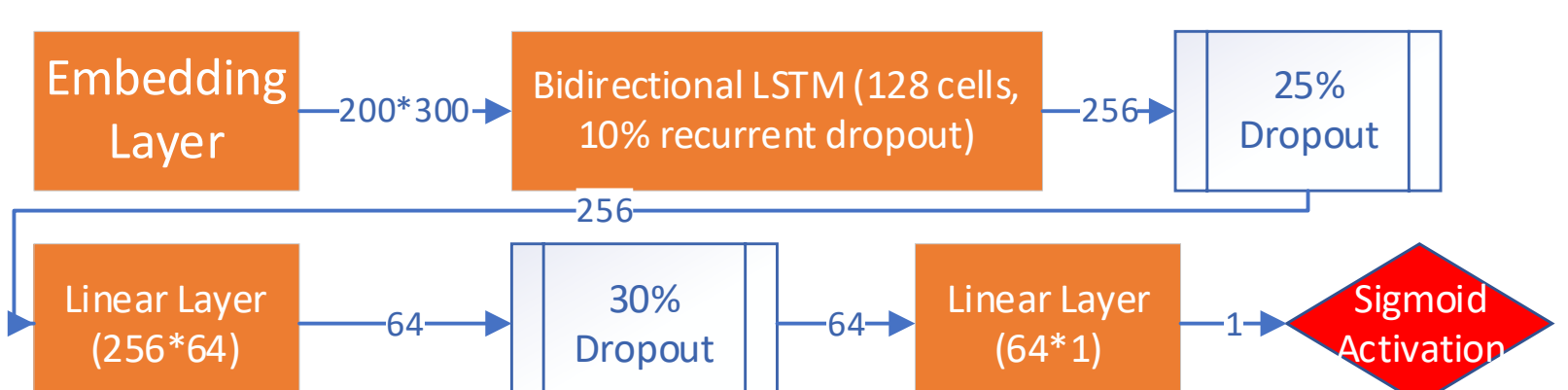
**Neural Words:** which, would, could, car, tree, etc.

### Datasets:

- Training set: training set of *IMDB Review* (IMDB).
- Test sets: test set of *IMDB review*; training set and test set of *Bag of Words Meets Bags of Popcorn* (BWMBP); training set and test set of *515K Hotel Reviews Data in Europe*.

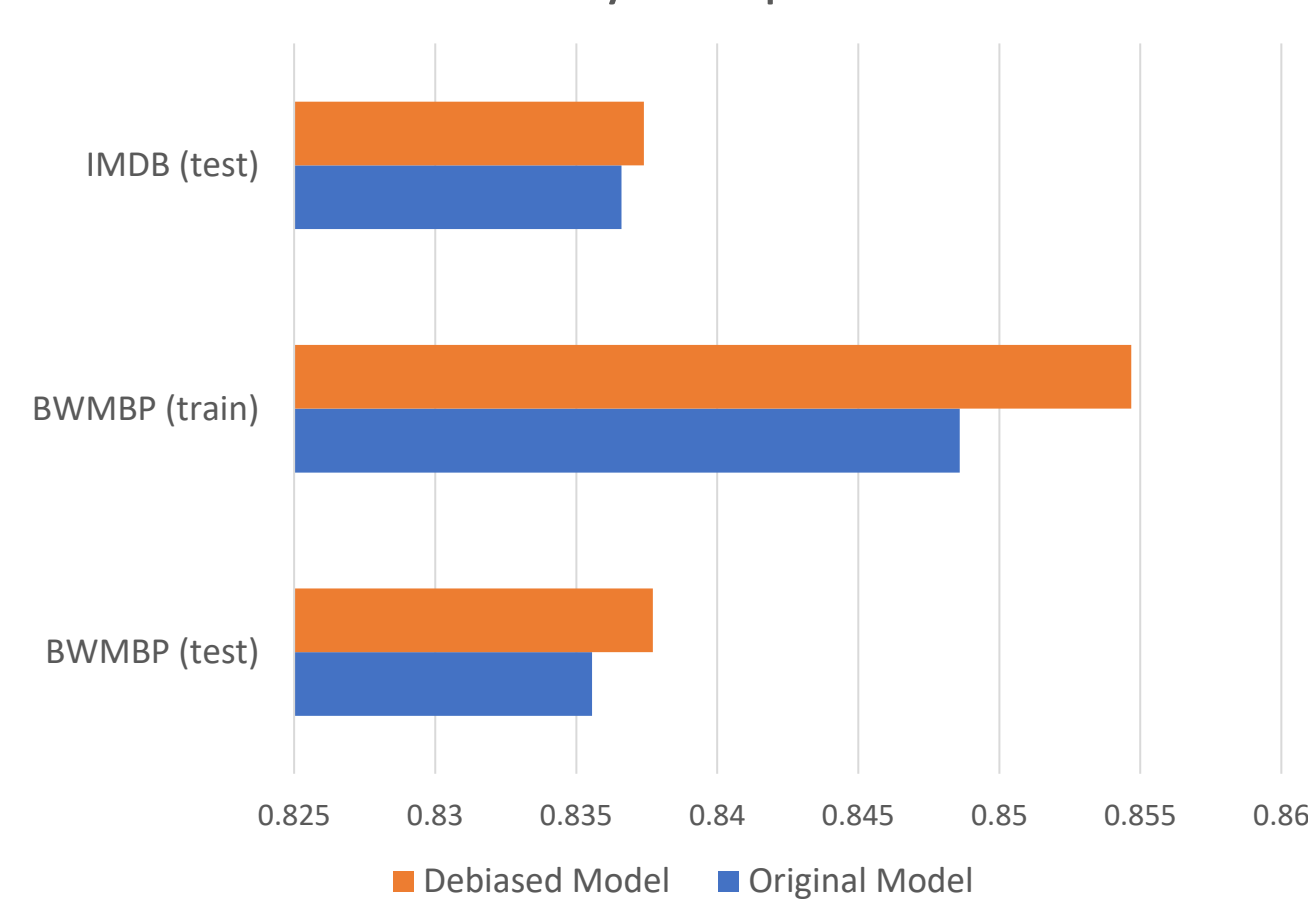
### Models:

- Model used to calculate the weights: random forest with maximum depth of 6.
- Main Model: BiLSTM with *GoogleNews-vectors-negative300* as the pre-trained word vectors.



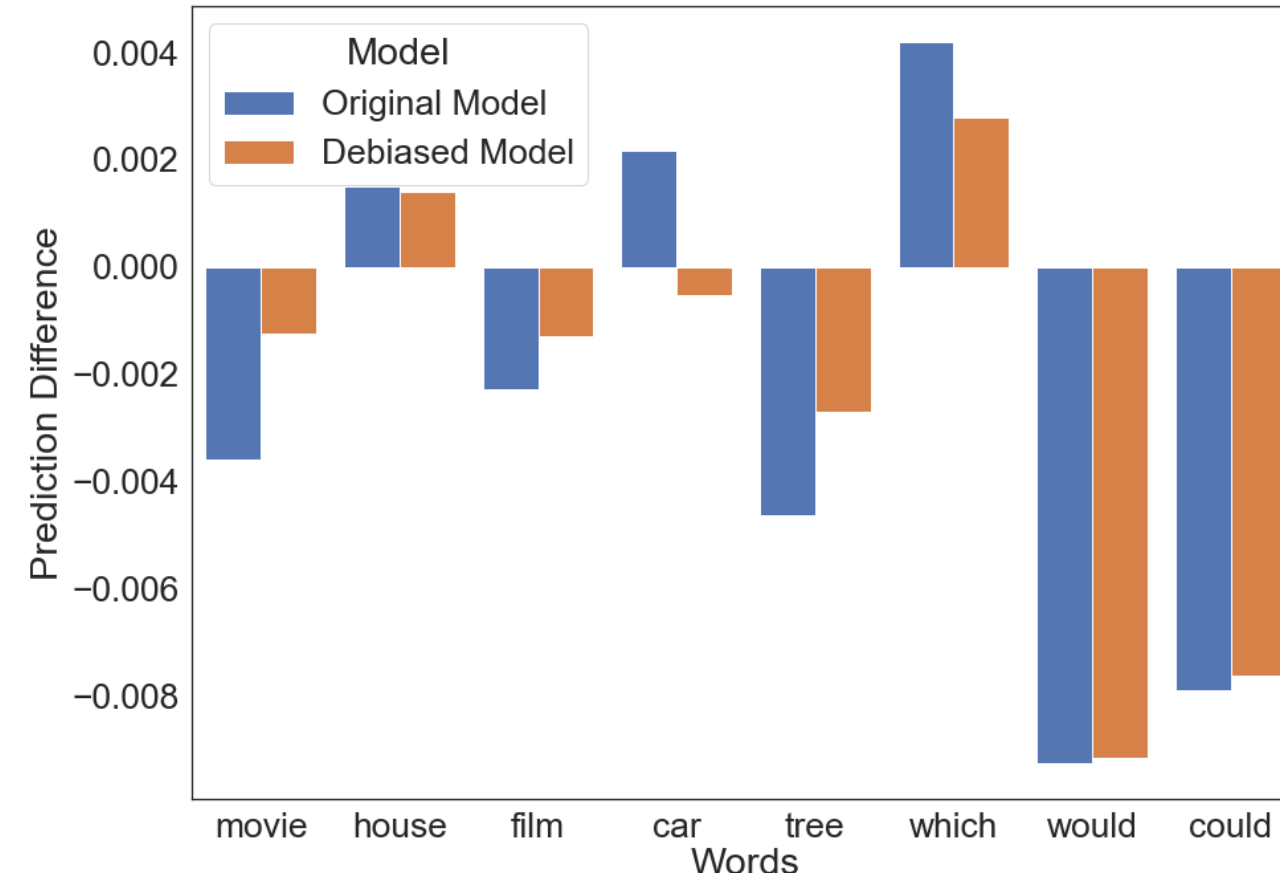
## Empirical Results

Accuracy Comparison



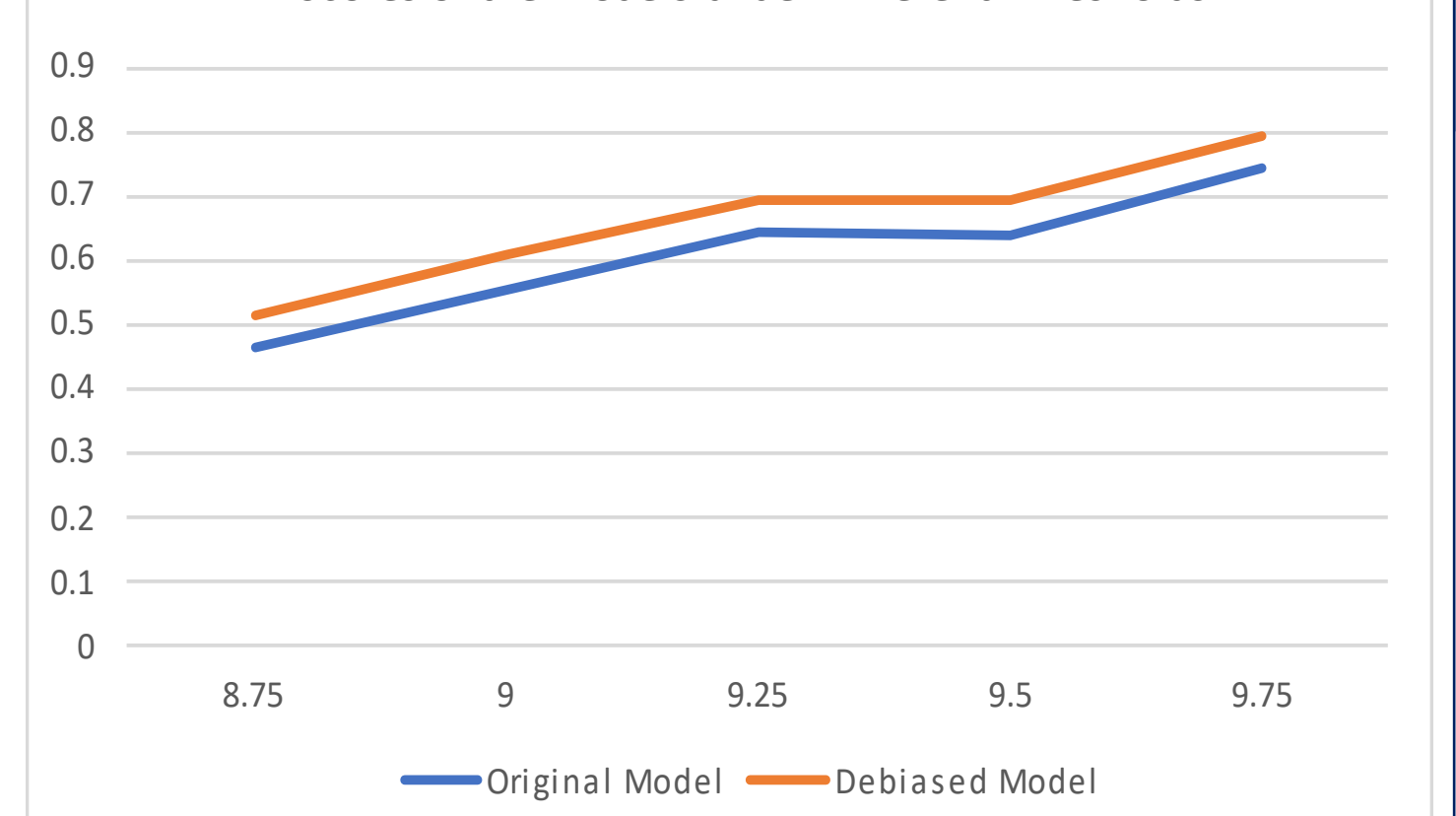
1. The debaised model performs better and generalizes better than the original model on data sets other than the training set.

The Effect of Words on Prediction



2. After the model being debaised, the effect of irrelevant neutral words in the prediction has been significantly reduced. The new model can be used for a wider range of tasks.

F1-scores of the Models under Different Thresholds



3. The model performance we obtain in hotel evaluation dataset, using different thresholds (used to discriminate between positive and negative). The generalization of the model after debiasing is better.