



Opinion Fraud Detection via Neural Autoencoder Decision Forest

Manqing Dong^{a,**}, Lina Yao^a, Xianzhi Wang^a, Boualem Benatallah^a, Chaoran Huang^a, Xiaodong Ning^a

^a*School of Computer Science and Engineering, University of New South Wales, Sydney 2052, Australia*

ABSTRACT

Online reviews play an important role in influencing buyers' daily purchase decisions. However, fake and meaningless reviews, which cannot reflect users' genuine purchase experience and opinions, widely exist on the Web and pose great challenges for users to make right choices. Therefore, it is desirable to build a fair model that evaluates the quality of products by distinguishing spamming reviews. We present an end-to-end trainable unified model to leverage the appealing properties from Autoencoder and random forest. A stochastic decision tree model is implemented to guide the global parameter learning process. Extensive experiments were conducted on a large Amazon review dataset. The proposed model consistently outperforms a series of compared methods.

© 2018 Elsevier Ltd. All rights reserved.

1. Introduction

As a result of the popularity of the e-commerce and experience sharing websites, online reviews influence increasingly on people's purchase decisions. More and more people are referring to others' opinions before buying something. In view of the significance of online reviews to the success of e-commerce, many sellers intentionally post fake reviews, either positive to promote the sales of their own products or negative ones to pull down the sales of other sellers' products. In fact, nowadays, online fraud, including the posting of fake reviews, has become a common phenomenon driven by profits (Wu et al., 2015). In the worst cases, some even hire a group of people for fraud.

Under such circumstances, it is paramountly important to evaluate products based on those credible reviews rather than the suspicious (or fake) ones, to help online users make wise purchase decisions. Until now, there have been a large body of research efforts on detecting fake or spam reviews (Sandulescu and Ester, 2015; Mukherjee et al., 2012), as well as spammers (i.e., the person who provides spam reviews) and spammer groups (Mukherjee et al., 2012). Most existing approaches achieve these purposes by extracting features from the review text, ratings, product meta-data (category and price), or review feedback (helpful/unhelpful votes etc.) (Jindal and Liu, 2008). Some researchers improve performance by considering

user behavior patterns (e.g. user's interactions with products or other users) (Wu et al., 2015) or the probabilistic distribution of users' behaviors (spamming and non-spamming) (Li et al., 2017). Rayana and Akoglu (2015) design a collective unified framework to exploit linguistic clues of deception, behavioral footprints, or relational ties between agents in a review system.

As for the learning techniques, several previous works mainly use traditional classification methods such as Support Vector Machine (SVM) (Ott et al., 2011), Naive Bayes Classifiers, and logistic regression (Jindal and Liu, 2008). Some researchers consider to use text analysis, for example, learning the written patterns (Ren et al., 2017). Some used probabilistic methods such as unsupervised Bayesian approach (Mukherjee et al., 2013), and Hidden Markov models (Li et al., 2017). And limited works considered using deep learning based methods (LeCun et al., 2015) for detecting the fake reviews (Wang et al., 2016).

Deep learning methods, such as Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN) (LeCun et al., 2015), have achieved a great amount of success in works such as Image Recognition (He et al., 2016). They show advantages in dealing with high dimensional data and utilize the non-linear combination of input features – where for detecting the fake news, we are also dealing with various features and try to find their link to the signal of fake. Some related works includes Ma's attempt (Ma et al., 2016), they used RNN for detecting rumors in microblogs; and in Baohua's work (Wang et al., 2016), they designed a semi-supervised recursive autoencoders for de-

^{**}Corresponding author:
e-mail: dongmanqing@gmail.com (Manqing Dong)

detecting review spam in microblogs. However, is there a way that could combine the deep learning methods and traditional classification methods, to fully extend the advantages of different methods?

In this work, we present an end-to-end trainable unified model to leverage the appealing properties of Autoencoder(Lange and Riedmiller, 2010) and random forest(Liaw et al., 2002). And the reason for combining these two methods is shown in the following. Autoencoder has been proved to be a robust algorithm which can produce unsupervised representations in feature patterns, and such representations are often more concise. Random forest is the ensemble of several decision trees, which can help prevent over fitting in each tree and shows good performance in practice. In our unified model, we use autoencoder to generate the hidden representations of the features, and take them as the input of random forest. The entire model is trained jointly via the stochastic and differentiable decision tree model, and the decision forest generates the final prediction. To summarize, we make the following contributions to the field of opinion fraud detection.

- We employ statistical analysis to define a list of quality measures, which utilize multiple metrics such as reviewer behavioral patterns, and review content analysis. These measures will be used as discriminative indicators of spamicity.
- We propose to detect spam reviews with the learned quality features to infer the salient correlations between reviews and heterogeneous cross-domain historical user information. We propose a joint model by fusing autoencoder and random forest to harvest the merits in an end-to-end trainable way. The model is learned by back propagation using stochastic and differentiable decision tree model. As such, the global optimization of model parameters like splits, leafs and nodes can be obtained.
- Experimental results on the real dataset demonstrate that the proposed model achieves the best prediction accuracy compared to several baselines and state-of-art methods. We have made the related code and dataset open-sourced for reuse.

2. Related Work

Opinion spam detection has been an active field of research in recent years, which covers broad topics such as detecting singleton fake reviewers, fake reviews, or even spam groups. For example, Jindal et al. (Jindal and Liu, 2008) built a logistic regression classifier with review feedback features, content characteristics, and rating related features to detect fake reviewers; (Feng et al., 2012) investigated syntactic stylometry for deceptive reviews detection by using SVR classifiers; (Mukherjee et al., 2012) propose a model for group opinion spam spotting with assistance of human experts.

Most detection methods are learning discriminate features from the linguistic (Ott et al., 2011; Feng et al., 2012), relational (Akoglu et al., 2013), and behavioral aspects (Feng et al.,

2012; Mukherjee et al., 2012) based on review text, ratings, product meta-data (category and price), and review feedback (helpful votes), reviewer behavioral patterns and linkage structures. Original fake review detection methods mainly consider only the review text information, for example, Ott et al.(Ott et al., 2011) consider text categorization and the sentiment of the text. And later on, a lot of works combine those different aspects of features for improvements. In Rayana’s work (Rayana and Akoglu, 2015), they combine features such as review texts, time-stamps, user behavioral information and the review network. However, limited work found have conducted quality feature analysis (e.g. the importance of a feature) for evaluating the selected features.

The techniques used in detecting the opinion fraud vary a lot. Some researchers use classical classification methods such as Support Vector Machine (SVM) Ott et al. (2011), Naive Bayes Classifiers, and logistic regression(Jindal and Liu, 2008). Such as Vlad et al.(Sandulescu and Ester, 2015) employ bag-of-words and bag-of-opinions Latent Dirichlet Allocation (LDA) model to detect singleton review spammers based on semantic similarity, in which they use cosine similarity as a measure, and topics and extracted from the reviews text, with parts-of-speech (POS) as patterns. And others may prefer statistical methods: Mukherjee et al. (Mukherjee et al., 2013) use unsupervised Bayesian approach and considered user’s behavior features and bigrams; A most recent work (Li et al., 2017) models the distribution of reviewers’ posting rates and utilizes a coupled hidden Markov model to capture both reviewers’ posting behaviors and co-bursting signals. There are also some trials on deep learning, which achieves great success in several areas in recent years. Several attempts include Wang’s work(Wang, 2017), which proposes using a hybrid convolutional neural network to detect fake news. Baohua(Wang et al., 2016) et al. designed a semi-supervised recursive autoencoders for detecting review spam in Weibo, a Chinese alternative to Twitter. Still, to the best of our knowledge, very limited work on hybrid deep learning models are proposed for fake opinion detection.

In this work, we propose to integrate neural networks with neural random forest, which is a differentiable model that could be fused into deep learning methods. And there have been already several attempts on leveraging the benefits of deep learning methods and random forest. One of the most important work in this area is the studies that cast the neural random forest as the structure of neural networks. The first few attempts started in the 1990s when SethiSethi (1990) proposed entropy net, which encoded decision trees into neural networks. Welbl et. al Welbl (2014) further proposed that any decision trees can be represented as two-layer Convolutional Neural Networks (CNN). Since the classical random forest cannot conduct back-propagation, other researchers tried to design differentiable decision tree methods, which can be traced from Montillo’s work Montillo et al. (2013). They investigated the use of sigmoidal functions for the task of differentiable information gain maximization. In Peter et al.’s workKontschieder et al. (2015), they introduce a stochastic, differentiable, back propagation compatible version of decision trees, guiding the representation learning in lower layers of deep convolutional networks. Different

from their work, we develop a hybrid model by combining neural decision forest with lightweight autoencoder to detect fraud opinion.

3. Quality Feature Representation

For detecting the fake reviews, the first thing is defining the features. In this paper, we mainly considered two scopes of features, one is user’s behavioral features and one is user’s review content features. The selected numerical features and the description of those features are listed in Table 1. And the categorical features will be shown later.

We extensively studied the reviewer behavior patterns in this section for the good understanding of the difference between spamicity and non-spamicity.

First, we investigate the difference between spammers and non-spammers. We used Wilcoxon signed-rank test(Kerby, 2014) to evaluate this difference over the continuous data such as entropy of ratings and review time entropy. This is a non-parametric statistical hypothesis test used to compare two related samples, and determines whether their population means ranks differ. Also there is no need to make an assumption on data distribution while using the *Wilcoxon Signed-Rank Test*.

Table 2 shows the p-value of the Wilcoxon Signed-Rank Test of three different alternative hypotheses: not equal to 0, less than 0, and greater than 0. A small p-value, especially a p-value is less than 0.05, stands for the high possibility to accept the alternative hypothesis, which also indicates a significant difference between spammer and non-spammer. For example, from the first line from table 2, we could know that the entropy of ratings for spammers is significantly less than non-spammers, which indicates the spammers are tending to give similar rates. The features in bold typeface are the features that show low significance in the distinguishing spammers and non-spammers.

Figure 1 shows a set of distinguishable patterns of review meta-data regarding density distribution. The red column stands for non-spammers and the green column stands for spammers. Figure 1(a) shows the entropy of ratings for the spammers is lower than non-spammers, which is in consistency with our p-value showed above. Figure 1(b) reflects that the spammers are holding the higher value than non-spammers in both low and high positive rate ratio, but lower than non-spammers in the middle. This can be explained that the spammers would likely to keep his/her rating habit, that is, a user who always rates very low scores or very high scores may likely be a spammer. Figure 1(c) describes the number of words of the summary, we could say that spammers would like to use fewer words than non-spammers; Figure 1(d) indicates that the spammers tend to comment early on the products; Figure 1(e) shows that spammers would be more likely appeared in products which is released for a long time; and from Figure 1(f) we could see spammers comment more randomly than non-spammers.

For categorical features, we apply Pearson χ^2 test(Corder and Foreman, 2014) to evaluate the correlation between the features and the labels (whether a reviewer is a spammer or not). Pearson χ^2 test is a statistical test applied to sets of categorical data

Table 1: Feature Explanation

User’s behavioral data
Number of reviewed products
Length of name
Reviewed product structure: the products include 54 different categories, such as book and music, here indicates the ratio of user’s reviewed product’s categories.
Minimum score and maximum score
Ratio of each score: 1 to 5
Number of each score
Ratio of positive and negative ratings: the proportion of high rates (score 4 and 5) and low rates (score 1 and 2) of a user
Entropy of ratings: as a measure of skewness in ratings. The entropy of the user’s rating can be explained by $-\sum_{i=1}^5 p_i \log p_i$, where p_i is the proportion that the user rated score i .
Total number of helpful and unhelpful votes
Average number of helpful and unhelpful votes
The ratio of the helpful and unhelpful votes
Median, min, max number of helpful and unhelpful votes.
Day gap: the time gap between the user’s first review and last review by days.
Review time entropy: as a measure of skewness in user’s review time by years, which could be denoted by formula $-\sum_{j=1}^m t_j \log t_j$, where m means the year gap between the user’s first review and last review, and t_j represent the proportion of the number of reviews a user generated for year j .
Same date indicator: whether the user’s first comment and the last comment are on the same date, 1 indicate the same date.
Active ratio: the proportion that a user was active during m years.
User’s review data
The Average review’s ratings for the product
Number of reviews
Entropy of the scores: as a measure of the skewness for the product’s review ratings. Can be denoted by $-\sum_{i=1}^5 s_i \log s_i$, where s_i denotes the proportion for score i in the whole reviews.
Comment time gap: the duration of the first review and the last review by days.
Entropy of product’s comment time: as a measure of the skewness for the product’s review times, <i>month</i> stands for the duration of the reviews for the product by month, then the entropy of product’s comment time can be denoted by $-\sum_{i=1}^{month} r_i \log r_i$, where r_i means the ratio of the reviews created in month i .
User’s rate for the product.
User’s helpful and unhelpful votes from others.
User’s comment time gap ratio: the days gap between user’s comment time and the first comment divided by whole comment time gap for the product.
User’s comment time rank: rank equal to 1 means this user is the first person giving the product review ratio.
The ratio of the user’s comment time rank and the total number of comments.
Review summary length: the number of words of the summary text.
Review text semantic: the semantic value for the review content, 1 stands for positive, 0 stands for neutral, -1 stands for negative.

to evaluate how likely it is that any observed difference between the sets arose by chance. The null hypothesis of the test is stating that the frequency distribution of certain events observed in a sample is consistent with a particular theoretical distribu-

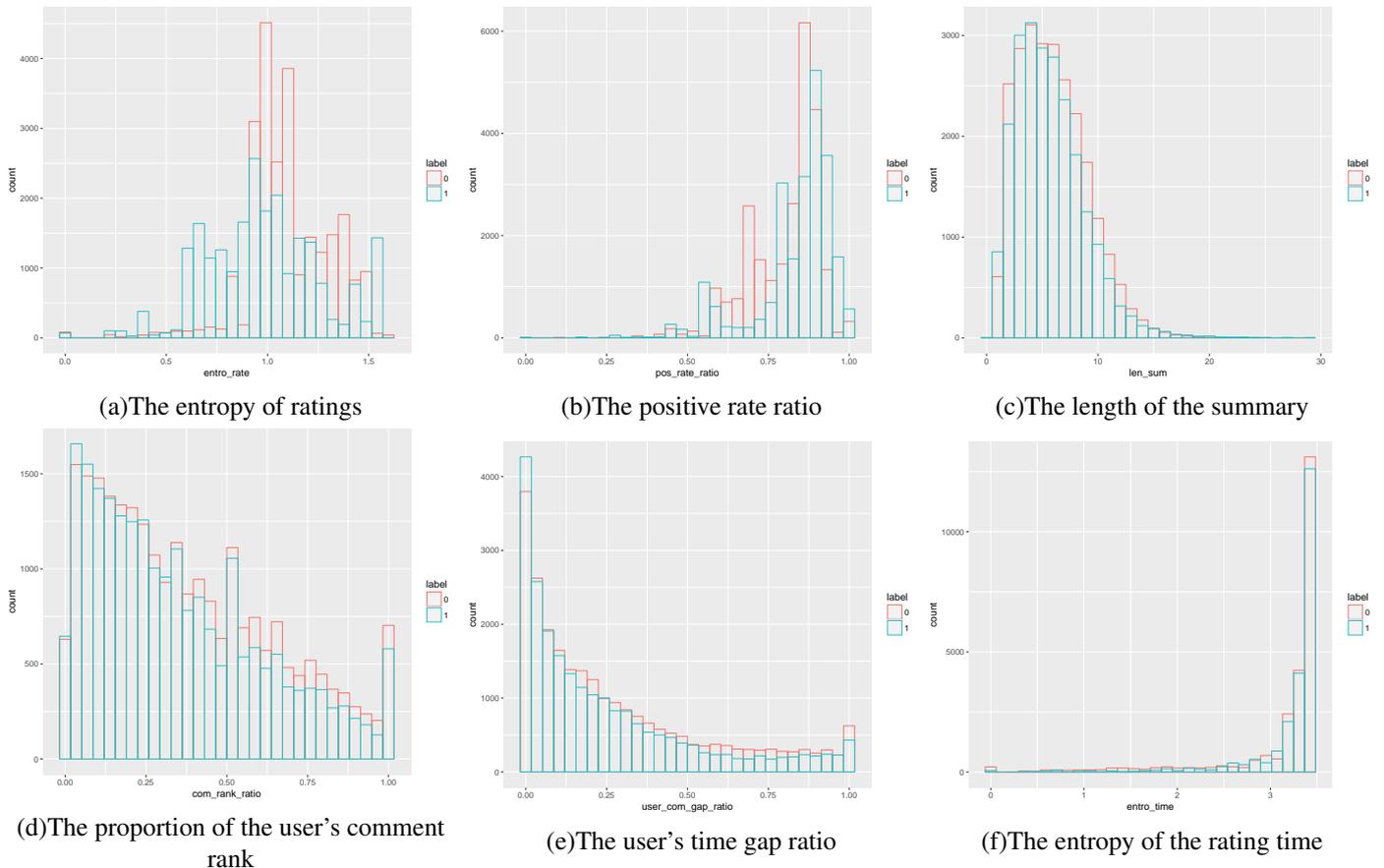


Fig. 1: The histogram of the spammer and non-spammer in (a) the entropy of the ratings, (b) the positive rate ratio, (c) the length of the summary, (d) the proportion of the user's comment rank around all reviews, (e) the user's time gap ratio for the comment time, (f) the entropy of the rating time.

tion. So if p -value is less than 0.05 stands the feature has a high correlation with the labels.

In Table 3, min/max rate is the minimal/maximal score a user had through his/her review history; common username is a binary variable – we extracted the normal English name list from world-english.org, the weirdness shows whether user's name showed in the name list, if so is valued as 0 otherwise is 1; having memo or not stands for whether a user has self-description; semantic of the summary stands for the semantic value for the review summary content; Semantic of the summary is the text analysis over review comments. The features in bold typeface don't demonstrate the significant difference between spammers and non-spammers.

4. Proposed Methodology

In this section, we introduce our proposed methodology, which can be referred as Figure 2. The whole idea is firstly to initialize all the parameters, and feed the features into the input layer X_I of autoencoder; also the cleaned features (the reconstruction of input layer X_I) are regarded as X_C . Then we get the hidden layer H , whose nodes will be followed by fully connected layers and reaches the nodes in the tree. After that, the nodes in the tree layer are delivered to the nodes in k decision trees. Decision nodes will make a decision and give a prediction by averaging the decision results. Then we update all the

parameters by minimizing the cost between the prediction and the real value, and iteratively optimize the whole process till to reach the optimal results.

Auto-encoder is a neural network (LeCun et al., 2015) that is trained to attempt to copy its input to its output. And it always consists of two parts, namely the encoder and the decoder. Internally, it has a hidden layer H describes a code used to represent the input.

An auto-encoder takes an input $X_I \in \mathbb{R}^n$, where X is the input feature vectors, and n is the number of the features.

First we map it (with an encoder) to a hidden representation $H \in \mathbb{R}^m$,

$$H = f_E(W_E X_I + b_E) \quad (1)$$

where $f_E(*)$ is a non-linearity such as the sigmoid function (Hecht-Nielsen, 1992). And W_E and b_E are weights and bias variables used in encoder part

The latent representation H or code is then mapped back (with a decoder) into a reconstruction X_C ,

$$X_C = f_D(W_D H + b_D) \quad (2)$$

where X_C is seen as a prediction of X_I , given the code H . And W_D and b_D are weights and bias variables used in decoder part, they are optimized such that the average reconstruction error is minimized.

The reconstruction error can be measured in many ways, depending on the appropriate distributional assumptions on the

Table 2: p -value of Wilcoxon Signed-Rank Test for Continuous Features

Features	$(p\text{-value})$ with <i>Alternative Hypothesis</i>		
	Not equal to 0	Less than 0	Greater than 0
User's history features			
Entropy of ratings	2.8e-05	1	1.4e-05
Number of reviewed products	0.0492	0.0246	0.9754
Ratio of positive ratings	0.0003	0.0001	0.9998
Ratio of negative ratings	0.1007	0.9497	0.0503
Number of helpful votes	0.0001	5.3e-05	0.9999
Sum unhelp	4.2e-06	2.1e-06	1
Average help	8.1e-08	4.1e-08	1
Average unhelp	1.5e-11	7.3e-12	1
Time gap	0.0024	0.0012	0.9988
Entropy of rating time	0.0006	0.0003	0.9997
Active ratio	0.0085	0.9958	0.0042
Memo length	0.8915	0.5544	0.4457
Mean rate	0.0004	0.0002	0.9998
User's review data			
Overall product score	5.8e-09	2.9e-09	1
Number of comments	2.2e-16	1	2.2e-16
Number of first day's comments	0.4217	0.7892	0.2108
Product's comment time gap	2.2e-16	2.2e-16	1
Comment rank	2.2e-16	1	2.2e-16
Comment rank ratio	2.2e-16	1	2.2e-16
User help	2.2e-16	2.2e-16	1
User unhelp	2.2e-16	2.2e-16	1
Length of the summary	2.2e-16	1	2.2e-16
Length of the review text	0.5556	0.2778	0.7222
User comment time gap	5.4e-05	2.7e-05	1
User comment time gap ratio	2.2e-16	1	2.2e-16

Table 3: p -value of Pearson χ^2 test for categorical features

Features	p-value
Min rate	1.4e-06
Max rate	0.1158
Common name	0.1721
Having memo or not	1
User rate	2.2e-16
Semantic of the summary	0.1554
Semantic of the review text	2.2e-16

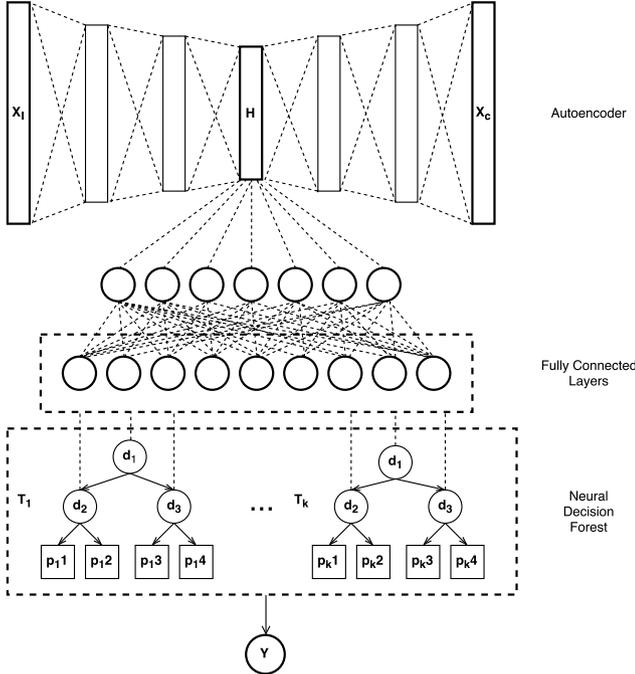


Fig. 2: Proposed Neural Autoencoder Decision Forest

input given the code. For example, the traditional squared error $L_{AE}(X_I, X_C) = \|X_I - X_C\|^2$ can be used.

Then let the hidden layer H go through several fully connected layers, we get X_T as the input nodes for the neural decision forest

$$X_T = g(W_F H + b_F) \quad (3)$$

where W_F and b_F are the weights used in the fully connection layers, and similarly, the $g(*)$ is the activation function. And in the next we will introduce how input those tree nodes X_T into the tree layer.

For an end-to-end neural network, each part should be differentiable, where traditional random forest is not suitable. Thus here in our work we utilized a differentiable version of random forest, which is called neural decision forest.

The neural decision forest (Kontschieder et al., 2015) is designed by Kontschieder et al., each node in the decision tree holds probabilistic distribution thus make the decision tree differentiable. Below explain the details.

Suppose we have number of K trees, for each tree, it is a structured classifier consisting of decision (or split) nodes \mathcal{D} and prediction (or leaf) nodes \mathcal{L} . Where the leaf nodes are the terminal nodes of the tree, and each prediction node $\ell \in \mathcal{L}$ holds a probability distribution p_ℓ over Y . Each decision node $d \in \mathcal{D}$ holds a decision function $D_d(X_T; \Theta) \in [0, 1]$, which stands for the probability that a sample reached decision node d and then be sent to the left sub-tree, for each decision nodes, suppose it holds the representation:

$$D_d(X_T; \Theta) = \sigma(f_d(X_T)) \quad (4)$$

where $\sigma(x) = (1 + e^{-x})^{-1}$ is the sigmoid function, and $f_d(X_T)$ is the transfer function for X_T ,

$$f_d(X_T) = W_T X_T \quad (5)$$

Θ stands for the set of previous parameters $\{W_E, W_D, W_F, W_T, b_E, b_D, b_F\}$.

About constructing the tree, here we suppose all the trees are following a classical binary tree structure, which means each node has two subtrees, then after defining the depth n_depth of a tree, suppose the depth is 1, just like what shows in figure 1, then the decision nodes in depth 1 is 2 and the number of leaf nodes is 4. More generally, once define the depth n_depth , the number of decision nodes would be 2^{n_depth} , and the number of leaf nodes would be $2^{n_depth+1}$.

And then the probability of a sample reach tree k to be predicted as class y would be:

$$\mathbb{P}_{T_k}[y|X_T, \Theta, P] = \sum_{\ell \in \mathcal{L}} P_\ell (\prod_{d \in \mathcal{D}} D_d(X_T; \Theta)^{\mathbb{1}_{left}} \overline{D}_d(X_T; \Theta)^{\mathbb{1}_{right}}) \quad (6)$$

where $\overline{D}_d(h; \Theta) = 1 - D_d(h; \Theta)$, and $\mathbb{1}_{left}$ indicates the indicator function for the nodes go left.

For example, in tree T_1 in figure 1, to get to the leaf node p_1 , here the $\prod_{d \in \mathcal{D}} D_d(h; \Theta)^{\mathbb{1}_{left}} \overline{D}_d(h; \Theta)^{\mathbb{1}_{right}}$ equals to $d_1 d_2$, which means the probability that a sample reach leaf node p_1 , and for p_2 , this probability will be $d_1 \overline{d_2}$, and P_ℓ means the probability

for the nodes in leaf ℓ predicted to be label y . Here, in our problem, the P_{ℓ_y} will have the form $P_{\ell_y} = (p_0, p_1)$, where p_0 stands for the probability of the leaf node to be 0 non-spammer.

Then for the forest of decision trees, it is an ensemble of decision trees $\mathcal{F} = \{T_1, \dots, T_K\}$, which delivers a prediction for sample x by averaging the output of each tree, which can be showed from:

$$\mathbb{P}_{\mathcal{F}}[y|x] = \frac{1}{K} \sum_{k=1}^K \mathbb{P}_{T_k}[y|x] \quad (7)$$

Then the prediction for the label y would be:

$$\hat{y} = \underset{y}{\operatorname{argmax}} \mathbb{P}_{\mathcal{F}}[y|x] \quad (8)$$

The loss for the neural decision forest is the average of the loss of each tree. And in the training process, the total loss is the average of all the training samples.

So the whole loss functions can be defined as:

$$L_F(x, y, \Theta, P) = E_{x \in X} (E_{T \in \mathcal{F}} (L_T(x, y, \Theta, P))) \quad (9)$$

And here, for we are dealing with a classification problem, the loss function we chose here is defined as:

$$\begin{aligned} L_T(x, y, \Theta, P) &= -(y \times \log(\mathbb{P}_T[y|x, \Theta, P]) \\ &\quad + (1 - y) \log(\mathbb{P}_T[1-y|x, \Theta, P])) \\ &= -\log(\mathbb{P}_T[y|x, \Theta, P]) \end{aligned} \quad (10)$$

where $y=0$ or 1 .

Then our final optimal function is minimizing the whole loss function which is finding:

$$\mathcal{L}(X_I, y, \Theta, P) = \arg \min_{\Theta} L_{AE}(X_I, X_C) + L_F(x, y, \Theta, P) \quad (11)$$

For learning the parameters, the optimization function we chose here is the RMSProp (Ruder, 2016), which is a mini-batch version of RProp, and RProp is equivalent to using the gradient but also dividing by the size of the gradient, the parameters are updating by:

$$g_t = \nabla \mathcal{L}(\Theta_{t-1}; \mathcal{B}, P) \quad (12)$$

$$G_t = G_t + g_t \odot g_t \quad (13)$$

$$\Theta_t = \Theta_{t-1} - \frac{\eta}{\sqrt{G_t} + \epsilon} \odot g_t \quad (14)$$

where g_t is the gradient, the η and ϵ are learning rates, and \mathcal{B} is a random subset of training data set.

For learning the leaf node distribution probability P , we still follow a RProp process but add a softmax function to the updated parameters.

$$g_t^P = \nabla \mathcal{L}(P_{t-1}; \mathcal{B}, \Theta) = \nabla L_F(P_{t-1}; x, y, \Theta) \quad (15)$$

$$G_t^P = G_t^P + g_t^P \odot g_t^P \quad (16)$$

Algorithm 1 Training Process

Require: Feature vector of a review X_I , Label for the review (fake or not) y , number of epoch n_epoch , number of trees n_tree , number of depth for a tree n_depth

- 1: Generalizing the structure for the forest with n_tree and n_depth
- 2: Initialize parameters Θ, P randomly
- 3: Randomly shuffle the data sets
- 4: **for** i **in** $1:n_epoch$ **do**
- 5: break data sets with n_batch of $batch_size$ piece of data sets
- 6: **for** j **in** $1:n_batch$ **do**
- 7: Update Θ with RProp optimization.
- 8: **end for**
- 9: Update P by (15) to (18)
- 10: **end for**

Algorithm 2 Prediction with proposed approach

Require: Feature vector of a review X

Ensure: Label for the review (fake or not) y

- 1: Get hidden representation for X by an encoder:
- 2: $H = f_E(W_E X_I + b_E)$
- 3: Get tree input layer from fully connected layers:
- 4: $X_T = g(W_F H + b_F)$
- 5: **for** i **in** $1:I$ **do**
- 6: $p(y=i) = 0$
- 7: **for** k **in** $1:K$ **do**
- 8: $p = P_k[y = i | X_T, \Theta, P]$
- 9: $p(y = i) = p(y = i) + p$
- 10: **end for**
- 11: $p(y=i) = \frac{1}{K} p(y = i)$
- 12: **end for**
- 13: $y = \max_i p(y = i)$
- 14: **return** y

$$P_t = P_{t-1} - \frac{\eta^P}{\sqrt{G_t^P} + \epsilon} \odot g_t^P \quad (17)$$

$$P_t = \operatorname{softmax}(P_t) \quad (18)$$

The details for how to calculate the derivative of loss function for each parameters could refer from paper Kotschieder et al. (2015). And below is the pseudo-code for our algorithm for both training process and testing process.

4.1. Computational Complexity

Our proposed process mainly contains three parts, the autoencoder part, the fully connected layers, and the neural random forest part.

For each epoch, for the autoencoder part, the time complexity is $O(\sum_l^{2*n_{AE}} n_{l-1}^{AE} n_l^{AE})$, where n_{AE} is the number of layers from input layer to hidden layer, l is the layer index, n_l^{AE} is the number of nodes in an autoencoder layer l .

For the fully connected layers part, the time complexity is $O(\sum_l^{n_{FC}} n_{l-1}^{FC} n_l^{FC})$, where n_{FC} is the number of layers for fully

connected layers, and similarly, n_l^{FC} is the number of nodes in fully connected layer l .

As for the neural random forest part, the time complexity is $O(n_{n_FC} \times n_tree \times n_leafnodes)$, where n_{n_FC} is the number of output nodes from fully connected layers, n_tree is the number of trees in a forest, and $n_leafnodes$ is the number of leafnodes in a tree.

For the whole training process, the whole time complexity should multiply n_epoch and n_batch , which are the number of epochs and number of batch of the dataset.

5. Experiments

We evaluate proposed approach using a real-world dataset. In the following section we first describe our dataset, and then compared existing methods with the presentation of results and analysis. Also, experiment settings and parameter tuning are explained, as well as the impact of features.

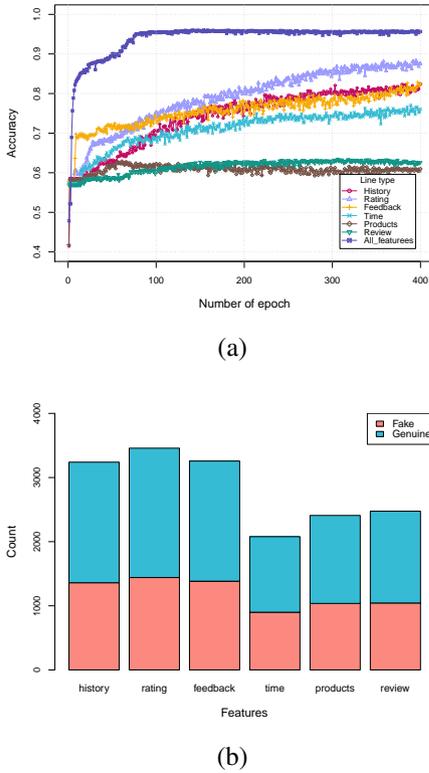


Fig. 3: (a) accuracy with different feature subset, the blue line shows the accuracy using full feature set; (b) predicted results using different features.

5.1. Dataset Description

The dataset used for our work is based on the public collected Amazon review data set from He and McAuley (2016), and ground-truth spamming review dataset from Mukherjee et al. (2012). We preprocessed the dataset as follows. We obtained a manual labeled subset for 815 unique users in more than 2000 potential groups. After averaging the manually labeled spamming score (ranging from 0 and 1, more the score approaches 1, more likely this user is a spammer.), we labeled users with

the average score less than 0.5 as a non-spammer and labeled them 0, otherwise, the spammers are labeled 1. Filtering is also carried out to handle missing values. To control the high appearance of some users, who wrote even higher than 500 reviews, we set a threshold(20 reviews) for selecting user’s reviews. Then we got 7950 reviews as our dataset, which is used for spamming activity detection. We regard reviews written by spammers as fake reviews and whose number counts to 3363.

5.2. Impact of Features

In this section, we evaluated the impact of different types of feature subsets on the performance. All features are stated in Table 1, and we divide the features into 6 scopes:

- History records: number of reviewed products; length of name; reviewed product structure.
- Rating signal: minimum score and maximum score; Ratio of each score (1 to 5) ; number of each score; Ratio of positive and negative ratings; entropy of ratings.
- Feedback signal: total number of helpful and unhelpful votes; average number of helpful and unhelpful votes; the ratio of the helpful and unhelpful votes; median, min, max number of helpful and unhelpful votes.
- Time signal: day gap; review time entropy; same date indicator; active ratio.
- Product’s comment information: the average reviews ratings for the product; number of reviews; entropy of the scores; comment time gap; entropy of products comment time.
- User’s review information: users rate for the product; users helpful and unhelpful votes from others; users comment time gap ratio; users comment time rank; the ratio of the users comment time rank and the total number of comments; review summary length; review text semantic.

We run the proposed neural autoencoder decision forest using these 6 types of feature, as well as full feature set. And we draw the accuracy line chart and the proportion of the predicted results for each type of features, which is shown in figure 3.

From Figure 3, we can observe that the performance with full feature set consistently outperforms the prediction results with single type of feature set; among the six type of features, user’s rating signal is most helpful for predicting fake reviews. Which demonstrate the effectiveness of designated various types of features in this work, and the effectiveness of our proposed model for combining different scope of features to final prediction.

5.3. Parameter Tuning

Last section we showed the effectiveness of our model for combining the features, and in this part, we are going to tuning the parameters. We mainly considered the parameters: the number of depth of each tree in the neural decision forest; with

or without fully connected layers between autoencoder and neural decision forest; number of layers in autoencoder; the normalization methods (Ioffe and Szegedy, 2015); the batch size; and the number of trees in a forest.

In the training process, among the 7950 reviews, we randomly selected 4000 samples, and initialized the weights and biases of autoencoder with normal distributed random numbers. The default set for our model is 2 layers of autoencoder, with one fully connected layer after the autoencoder, and 5 trees with depth 2 in the neural random forest. The batch size is 100 and we take z-score as our normalization method. Then we iterated the model with 400 epochs.

We compared the accuracy and time cost among different parameters, which is shown in figure 4.

From Figure 4 (a)-(f) we can see that the the random forest with only 1 depth outputs a bad prediction in performance, and the accuracy becomes better with 2 and onwards in depths; fully connected layers can help with the prediction accuracy; a model with 2 layers in autoencoder can get great prediction accuracy with first 100 epochs, and more layers for each autoencoder results in better prediction yet longer to converge and give the result; the normalized dataset shows better performance than the unnormalized one, and z_score has a better performance than min-max in terms of normalization; batch size 50 produces the optimal prediction, with a comparable short computation time at the highest prediction accuracy; and finally, the prediction improves with increased number of trees in a forest.

Figure 4 (g)-(l) show the time consumption under different parameter settings. The axis y is taking the *log* for the time (seconds). We could see deeper the depth of the tree, more time will be cost; add fully connected layers will increase the time cost; the time cost will increase with bigger number of autoencoder layers; unnormalized dataset cost more time; less batch size takes more training time; and higher tree depths incurs more training time.

Thus, considering the time cost and the prediction performance, we choose a batch size with 50, 5 trees in a random forest with depth of 3, z-score as normalization method, one fully connected layer, and 2 layers in autoencoder. In the test we get correctly predicted fake reviews and 2192 correctly predicted genuine reviews, 65 falsely predicted fake reviews and 99 falsely predicted genuine reviews. This means an accuracy of approximately 95.85%.

5.4. Overall Comparison

We compared our model with a set of existing works, which also use Amazon dataset but are different in extracted features and solutions. A brief introduction for these papers is listed below:

- Jindal and Liu (Jindal and Liu, 2008) first manually labeled easy identify spam reviews and then considered multi-features, they used logistic regression to do the classification.
- Lai et al. (Lai et al., 2010): considered the content of the reviews, they used an unsupervised probabilistic language model and a supervised SVM model.

- Mukherjee et al. (Mukherjee et al., 2013): This work used unsupervised Bayesian approach and considered user's behavior features and bigrams.
- Rout et al. (Rout et al., 2017): considered the sentiment analysis and they extracted sentiment score, linguistic features and unigram as feature set.
- Shreesh et al. (Bhat and Culotta, 2017): mainly focused on the content based information and they combined positive unlabeled learning with domain adaptation to train a classifier.
- Heydari et al. (Heydari et al., 2016): Heydari et al. considered the rating deviation, content based factors and activeness of reviewers, and they captured suspicious time intervals from time series by a pattern recognition technique.
- Xu and Zhang (Xu and Zhang, 2015): considered target based, rating based, temporal based and activity based signals. A Latent Collusion Model (LCM) is proposed to work unsupervisedly and supervisedly under different conditions.
- Zhang et al. (Zhang et al., 2016): considered un verbal features and used four classical supervised classification methods: SVM, Naive Bayes, decision tree and random forest.

To evaluate the performance of our fake news detection algorithms, we used the following criteria (Powers, 2011):

$$\begin{aligned}
 Accuracy &= \frac{TP + TN}{TP + FP + TN + FN} \\
 Precision &= \frac{TP}{TP + FP} \\
 Recall &= \frac{TP}{TP + FN} \\
 F1 &= \frac{2TP}{2TP + FP + FN}
 \end{aligned} \tag{19}$$

where TP, FP, TN, FN represent the counts for true positive, false positive, true negative and false negative, respectively.

Table 3 shows the comparison results. And the results for each paper are extracted from the published papers. We could see that our proposed model gives the best prediction accuracy, precision and F1 score values.

Comparing with these methods, we could conclude that our approach demonstrates readily good performance in predicting the fake reviews. It also indicates that the proposed neural autoencoder random forest model is at least a appreciable classification model in predicting the review spam.

6. Conclusion

In this paper, we focus on detecting fake reviews, which is an attractive topic and has drawn attention of certain researchers. While previous work mainly focused on using classical classification methods, here we proposed an end to end trainable

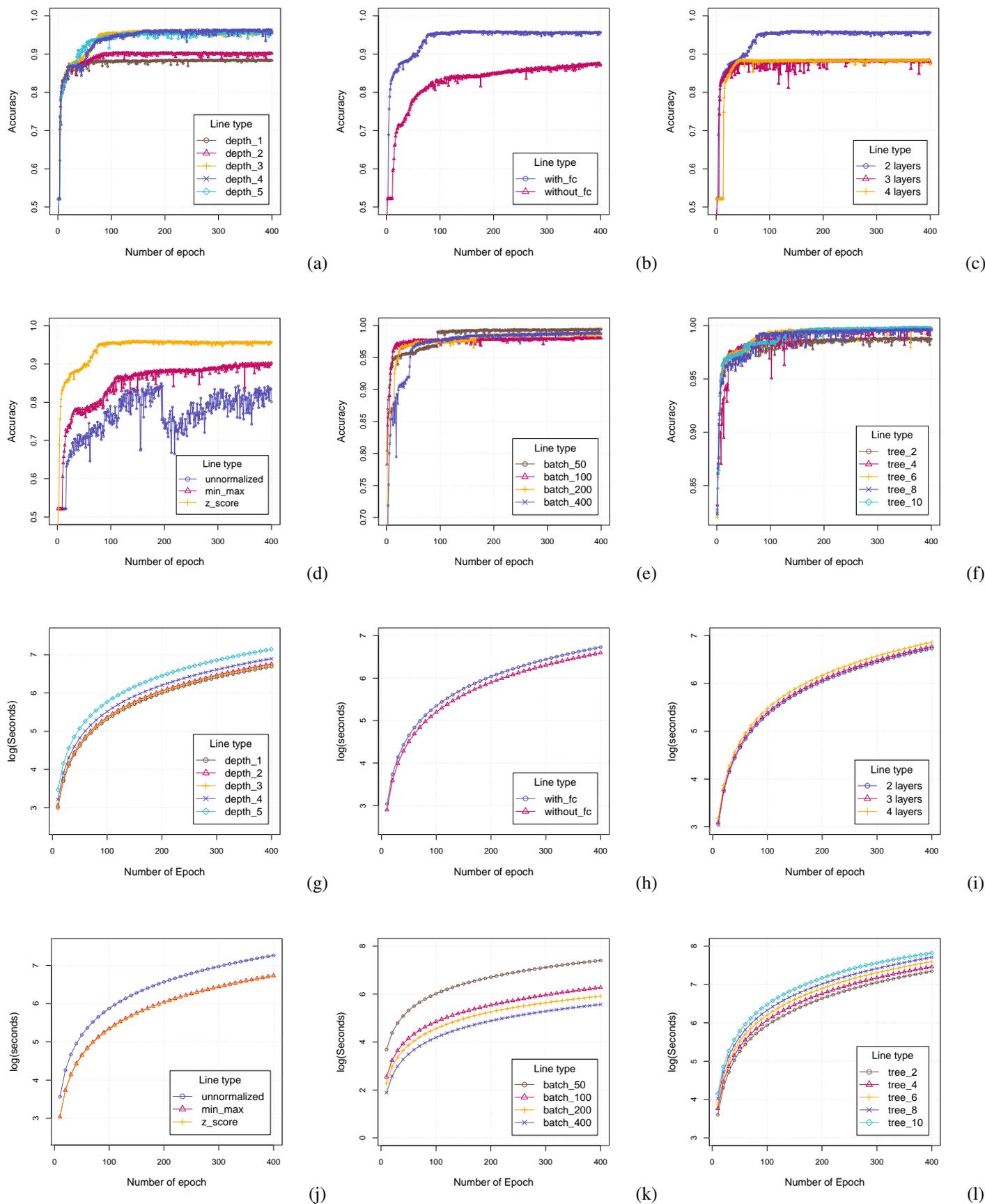


Fig. 4: The accuracy for 400 epochs of the parameter: (a)the number of depths of the tree, (b)with or without fully connected layers between autoencoder and neural decision forest, (c)the number of autoencoder layers, (d)the normalization methods, (e) the batch size, and (f) the number of trees in a forest. And the time for 400 epochs of the according parameters: (g)-(l).

Table 4: Comparison with the State-of-art Methods

Methods	Precision	Accuracy	F1	Recall
Jindal and Liu	-	78%	-	-
Lai et al.	-	-	-	96.38%
Mukherjee et al.	79.6%	86.1%	77.3%	75.1%
Rout et al.	-	92.11%	-	-
Shreesh et al.	87.6%	-	84.3%	82.8%
Heydari et al.	82%	-	86%	88%
Xu and Zhang	86.4%	85.2%	-	-
Zhang et al.	87.12%	87.81%	88.31%	89.63%
Ours	96.08%	95.85%	95.11%	94.15%

joint model combining autoencoder with neural random forest, to detect the fake reviews. Firstly, we did quality feature analysis to mining the feature efficiency in detecting opinion fraud, and those efficient features are employed as input in our model. And then we declared the design of our model. The model starts with an autoencoder, and then fully connected layers and ends in a differentiable random forest, where the differentiable random forest is trained back propagation. The forest averages the prediction result of each tree and outputs our final prediction. We also analyzed the impact of features and tuned various parameters in our model. The extensive experimental results demonstrate that our method beats a series of state-of-the-art methods yielding 96% of accuracy.

References

- Akoglu, L., Chandy, R., Faloutsos, C., 2013. Opinion fraud detection in online reviews by network effects. *ICWSM 13*, 2–11.
- Bhat, S.K., Culotta, A., 2017. Identifying leading indicators of product recalls from online reviews using positive unlabeled learning and domain adaptation. *arXiv preprint arXiv:1703.00518*.
- Corder, G.W., Foreman, D.I., 2014. *Nonparametric statistics: A step-by-step approach*. John Wiley & Sons.
- Feng, S., Banerjee, R., Choi, Y., 2012. Syntactic stylometry for deception detection, in: *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers-Volume 2*, Association for Computational Linguistics. pp. 171–175.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
- He, R., McAuley, J., 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering, in: *proceedings of the 25th international conference on world wide web*, International World Wide Web Conferences Steering Committee. pp. 507–517.
- Hecht-Nielsen, R., 1992. Theory of the backpropagation neural network, in: *Neural networks for perception*. Elsevier, pp. 65–93.
- Heydari, A., Tavakoli, M., Salim, N., 2016. Detection of fake opinions using time series. *Expert Systems with Applications* 58, 83–92.
- Ioffe, S., Szegedy, C., 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift, in: *International conference on machine learning*, pp. 448–456.
- Jindal, N., Liu, B., 2008. Opinion spam and analysis, in: *Proceedings of the 2008 International Conference on Web Search and Data Mining*, ACM. pp. 219–230.
- Kerby, D.S., 2014. The simple difference formula: An approach to teaching nonparametric correlation. *Comprehensive Psychology* 3, 11–17.
- Kontschieder, P., Fiterau, M., Criminisi, A., Buló, S.R., 2015. Deep neural decision forests, in: *Computer Vision (ICCV), 2015 IEEE International Conference on*, IEEE. pp. 1467–1475.
- Lai, C., Xu, K., Lau, R.Y., Li, Y., Jing, L., 2010. Toward a language modeling approach for consumer review spam detection, in: *e-Business Engineering (ICEBE), 2010 IEEE 7th International Conference on*, IEEE. pp. 1–8.
- Lange, S., Riedmiller, M., 2010. Deep auto-encoder neural networks in reinforcement learning, in: *Neural Networks (IJCNN), The 2010 International Joint Conference on*, IEEE. pp. 1–8.
- LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *nature* 521, 436.
- Li, H., Fei, G., Wang, S., Liu, B., Shao, W., Mukherjee, A., Shao, J., 2017. Bimodal distribution and co-bursting in review spam detection, in: *Proceedings of the 26th International Conference on World Wide Web*, International World Wide Web Conferences Steering Committee. pp. 1063–1072.
- Liaw, A., Wiener, M., et al., 2002. Classification and regression by randomforest. *R news* 2, 18–22.
- Ma, J., Gao, W., Mitra, P., Kwon, S., Jansen, B.J., Wong, K.F., Cha, M., 2016. Detecting rumors from microblogs with recurrent neural networks., in: *IJ-CAI*, pp. 3818–3824.
- Montillo, A., Tu, J., Shotton, J., Winn, J., Iglesias, J.E., Metaxas, D.N., Criminisi, A., 2013. Entanglement and differentiable information gain maximization, in: *Decision Forests for Computer Vision and Medical Image Analysis*. Springer, pp. 273–293.
- Mukherjee, A., Kumar, A., Liu, B., Wang, J., Hsu, M., Castellanos, M., Ghosh, R., 2013. Spotting opinion spammers using behavioral footprints, in: *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM. pp. 632–640.
- Mukherjee, A., Liu, B., Glance, N., 2012. Spotting fake reviewer groups in consumer reviews, in: *Proceedings of the 21st international conference on World Wide Web*, ACM. pp. 191–200.
- Ott, M., Choi, Y., Cardie, C., Hancock, J.T., 2011. Finding deceptive opinion spam by any stretch of the imagination, in: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, Association for Computational Linguistics. pp. 309–319.
- Powers, D.M., 2011. Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation.
- Rayana, S., Akoglu, L., 2015. Collective opinion spam detection: Bridging review networks and metadata, in: *Proceedings of the 21th acm sigkdd international conference on knowledge discovery and data mining*, ACM. pp. 985–994.
- Ren, P., Chen, Z., Ren, Z., Wei, F., Ma, J., de Rijke, M., 2017. Leveraging contextual sentence relations for extractive summarization using a neural attention model, in: *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, ACM. pp. 95–104.
- Rout, J.K., Singh, S., Jena, S.K., Bakshi, S., 2017. Deceptive review detection using labeled and unlabeled data. *Multimedia Tools and Applications* 76, 3187–3211.
- Ruder, S., 2016. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*.
- Sandulescu, V., Ester, M., 2015. Detecting singleton review spammers using semantic similarity, in: *Proceedings of the 24th international conference on World Wide Web*, ACM. pp. 971–976.
- Sethi, I.K., 1990. Entropy nets: from decision trees to neural networks. *Proceedings of the IEEE* 78, 1605–1613.
- Wang, B., Huang, J., Zheng, H., Wu, H., 2016. Semi-supervised recursive autoencoders for social review spam detection, in: *Computational Intelligence and Security (CIS), 2016 12th International Conference on*, IEEE. pp. 116–119.
- Wang, W.Y., 2017. "liar, liar pants on fire": A new benchmark dataset for fake news detection. *arXiv preprint arXiv:1705.00648*.
- Welbl, J., 2014. Casting random forests as artificial neural networks (and profiting from it), in: *German Conference on Pattern Recognition*, Springer. pp. 765–771.
- Wu, Z., Wang, Y., Wang, Y., Wu, J., Cao, J., Zhang, L., 2015. Spammers detection from product reviews: a hybrid model, in: *Data Mining (ICDM), 2015 IEEE International Conference on*, IEEE. pp. 1039–1044.
- Xu, C., Zhang, J., 2015. Towards collusive fraud detection in online reviews, in: *Data Mining (ICDM), 2015 IEEE International Conference on*, IEEE. pp. 1051–1056.
- Zhang, D., Zhou, L., Kehoe, J.L., Kilic, I.Y., 2016. What online reviewer behaviors really matter? effects of verbal and nonverbal behaviors on detection of fake online reviews. *Journal of Management Information Systems* 33, 456–481.