



SHILLING ATTACK DETECTION IN RECOMMENDER SYSTEMS USING CLASSIFICATION TECHNIQUES

Parneet Kaur

Department of Computer Science & Engg.
Thapar University, Patiala, Punjab, India

Shivani Goel

Department of Computer Science & Engg.
Thapar University, Patiala, Punjab, India

Abstract— Collaborative filtering based recommender system is prone to shilling attacks because of its open nature. Shillers inject pseudonomous profiles in the system's database with the intent of manipulating the recommendations to their benefits. Prior study has shown that the system's behavior can be easily influenced by even a less number of shilling profiles. In this paper, we simulated various attack models on Movie-Lens¹ dataset and used machine learning techniques to detect the attacks. We compared five classification algorithms and proposed a new model by integrating two models with high performances. In our experiments, we investigated and proved that the combination of random forest and adaptive boosting algorithm is more accurate than simple random forest model.

Keywords— Collaborative filtering, recommender system, shilling attack, prediction shift, precision, recall, f-measure, classification

I. INTRODUCTION

Recommender systems (RSs) have become popular in e-commerce, which provide recommendations of items to a customer that might be of his interest by predicting the ratings that would be assigned to an item by him. Collaborative filtering recommender system (CFRS) is vulnerable to profile injection attack [1]. Malicious user, who is indistinguishable from genuine users, inserts fake profiles in the system's database by using an automated tool to manipulate the system's output.

CF algorithms collect the profiles of the users, represent the preferences of individuals and provide them recommendations and predictions based on the taste of other like-minded profiles. The attacker injects his/her fake profile into the database. There is a possibility that his profile becomes similar to the actual users and produce recommendations in his favor. Some efforts are necessary to mount an attack. First is the knowledge needed about the system. In low knowledge attack,

details about the system are not required whereas in high knowledge attack, an attacker must have knowledge regarding the rating distribution in a system. The second effort required is to put the fake profiles and ratings in the system. Injection of fake profiles is difficult to prevent. Therefore, to maintain the trust of the system, attacker's profiles should be detected flawlessly. To solve this problem, we compared the performance of classification algorithms and proposed an integrated model by using voting approach.

Our remaining paper is arranged as follows. Section II describes the research done in this area. In Section III, we explained CF algorithm and various attack models. Attack profile classification, several detection attributes and metrics for evaluating the classifiers are introduced in section IV. Section V introduces extensive experiments executed and their results. In section VI, conclusion is given with potential direction for future work.

II. RELATED WORK

The term "shilling" was given by Riedl and Lam who introduced two models for mounting the attacks : Random and Average Bot and demonstrated that item based algorithm is more advantageous than user based algorithm [1]. However, item based approach also suffers from shilling attacks [3]. It has been analyzed that even with the less knowledge of the system attacks can be implemented successfully [4-6]. Segment based attacks against CFRSs has been introduced in [7]. Chirita introduced an algorithm of evaluation metrics to detect and remove the attack profiles from the system [8]. Shilling attack can be detected by using supervised, unsupervised or semi-supervised techniques.

Burke studied various detection attributes and used three classifiers i.e. SVM, kNN, C4.5 to classify the profiles and to improve the strength of the system. Among them, SVM is the best performer [9-10]. In [11], six supervised models are compared and they observed that neural network, random forest and SVM have higher performance. They ensemble these models and built a new model which outperforms in most of the cases. Zhang and Zhou proposed a detection model by including ensemble technique and back propagation



neural network and ensemble technique that improves the low precision of existing supervised approaches [12].

Unsupervised techniques require some prior knowledge. Lee and Zhu developed an attack detection technique by using clustering algorithms and multidimensional scaling but it is not useful for recognizing attack with small filler sizes [13]. In order to improve the power of the item based CF algorithm, a new collaborative filtering technique has been proposed in [14] by building various user models and DBSCAN clustering technique is used to detect the malicious users. The authors in [15] presented a new unsupervised attack detection approach i.e. RD-TIA. They also introduced a new detection attribute, DegSim' which succeeded in detecting segment and other group attacks. Dhimmar and Chauhan used ECLARANS and PAM clustering algorithms to detect the spam users and proved that former algorithm has higher accuracy than later [16].

In addition to supervised and unsupervised techniques, there are semi-supervised approaches that can be used to detect the attackers. In [17] semi-SAD detector has been introduced by applying semi-supervised techniques. It uses unlabeled profiles for improving the performance of detection.

We have used supervised learning approach in our paper to distinguish the malicious profiles and proposed a novel model.

III. BACKGROUND

A. Recommendation Algorithms –

User based CF algorithm obtains the k most similar users and predict the ratings of the target user by using k nearest neighbor algorithm [2]. Similarity $S_{m,n}$ between user m and n can be calculated by using Pearson's correlation measure as follows:

$$S_{m,n} = \frac{\sum_{x \in I} (r_{m,x} - \bar{r}_m) \cdot (r_{n,x} - \bar{r}_n)}{\sqrt{\sum_{x \in I} (r_{m,x} - \bar{r}_m)^2} \cdot \sqrt{\sum_{x \in I} (r_{n,x} - \bar{r}_n)^2}} \quad (1)$$

where, I denotes a set of all items, $r_{m,x}$ and $r_{n,x}$ are the ratings given to an item x by user 'm' and its neighbor 'n' respectively. We took 30 as a neighborhood size in our experiments. We only considered those neighbors who have similarity greater than zero to prevent the negative correlations. Rating is predicted for item x by using (2):

$$P_{m,x} = \bar{r}_m + \frac{\sum_{n \in Q} S_{m,n} (r_{n,x} - \bar{r}_n)}{\sum_{n \in Q} |S_{m,n}|} \quad (2)$$

where, Q represents k similar users, $r_{n,x}$ denotes rating by the user n to item x , \bar{r}_n represents overall mean of ratings.

B. Shilling Attack Models –

The attacks contain attack profiles which influence the output of the system, biased data and target items. Profile injection attacks can be classified as nuke attack and push attack. In

nuke attack, fictitious users provide lowest rating to the target items in order to downgrade them whereas, in the push attack, maximum score is given to target items with the aim of promoting them. In this paper, we concentrated on push attacks. The attacker creates unscrupulous profiles using attack models that require less knowledge and have high impact on the system. Four attack models: random, average, bandwagon and segment attack are used to inject attack profiles into the database [3].

A fictitious profile contains vector of t -dimensional ratings where t denotes the count of items in system. The t -dimensional vector is partitioned into four sets: I_N, I_T, I_F, I_S .

- I_N : A set of items which are unrated.
- I_T : A set of target items. Rating, r_{max} is given to these items in case of push attack and r_{min} in nuke attack.
- I_F : A set of randomly selected filler items.
- I_S : A set of items those are chosen randomly.

The characteristics of different attack models are summarized in table 1.

Table -1. Attack Model's Characteristics

Attacks	Random	Average	Bandwagon	Segment
I_F	overall mean	item mean	overall mean	r_{min}/r_{max}
I_S	empty	empty	r_{max}/r_{min}	r_{max}/r_{min}
I_T (push/nuke)	r_{max}/r_{min}	r_{max}/r_{min}	r_{max}/r_{min}	r_{max}/r_{min}
I_N	empty	empty	empty	empty

1) Random Attack

To mount an attack, less knowledge about the system is necessary. This attack is mounted by selecting filler items (I_F) randomly. These items are rated using standard deviation and average rating of the system distributed normally. Selected item set is empty in this case. Depending on the type of attack (nuke or push), lowest or highest ratings are assigned to targeted items. This attack has very less impact on the system [1].

2) Average Attack

This attack is difficult to implement as it requires details of the system. Attackers randomly choose filler items and provide them rating in the same manner as in random attack. But the only difference is that it uses average ratings of each item instead of global average of the system. The pattern for rating the targeted items is also similar to the previous attack. Average attack has maximum impact on user based algorithm [3].

3) Bandwagon Attack

In this model, attacker creates the skewed profiles which consist of those items that are rated by many users. Hence, the probability of attackers being similar to the genuine users is high. Maximum ratings are provided to the set of targeted items, I_T and to frequently rated items, I_S . The filler item set is



selected in a same way as in previous attacks. This attack model is considered as low knowledge attack because less effort is required to obtain the popular items.

4) *Segment Attack*

It requires less knowledge to mount an attack. The main intention of this attack is to publicize the target items among a group of target users [7]. For example, a producer of an action movie wants his/her movie to be suggested to the viewers who are the fans of “The Dark Knight”, not to the ones who like romantic movies. Maximum ratings are assigned to segmented items, I_S . Segment items are the popular items among a particular segment. Filler items set, I_F are given lowest ratings so that the attack would have maximum impact.

IV. ATTACK PROFILE DETECTION

A. Classification of attack profiles –

We have used training set to train a classifier in order to differentiate malicious profiles from the authentic profiles. Then two types of detection attributes have been created i.e. generic and model specific attributes. Generic attributes are generated by considering a profile as a whole whereas model specific attributes are created to find the features of a particular attack model. We compared the performance of five classification algorithms: Random forest (RF), Naive bayes (NB), J48, ZeroR, Radial basis function network (RBF n/w). k-fold cross validation technique is used to evaluate these predictive models.

B. Evaluation Metrics –

Recall, precision and f-measure are the metrics that have been used to measure the performance of classifiers. These metrics can be calculated as follows:

$$recall = \frac{TP}{TP+FN} \tag{3}$$

$$precision = \frac{TP}{TP+FP} \tag{4}$$

$$f - measure = \frac{2 \cdot precision \cdot recall}{precision + recall} \tag{5}$$

where false negatives (FN) is the fake profiles that are not correctly identified, false positives (FP) means the count of genuine profiles that are incorrectly classified, true positive (TP) means the fake profiles identified accurately.

C. Detection Attributes –

Our objective is to classify the profile as a genuine user or a fake using detection attributes. These attributes are categorized as generic or model specific attributes. Generic attributes are basic metrics for all profiles. Type specific attributes are those attributes which are used to distinguish fake profiles based on the features of particular attack model.

Model specific attributes are more effective than generic attributes.

1) *Generic Attributes*

Various general attributes have been introduced by Chirita for detecting anonymous profiles [8]. Suppose U is a universal set of users in the system, p_a denotes a profile of user a , M_a specifies the count of items given by user a , $r_{a,x}$ is the rating that user a gave to some item x . R_x is total count of ratings given to item x . \bar{r}_x denotes the mean rating of an item x .

- *Length Variance (LengthVar):*

$$LengthVar_a = \frac{|\#r_a - \bar{\#r}|}{\sum_{a \in U} (\#r_a - \bar{\#r})^2} \tag{6}$$

where $\#r_a$ is number of ratings for user a .

- *Weighted degree of agreement (WDA):*

$$WDA_a = \sum_{x=0}^{M_a} \frac{|r_{a,x} - \bar{r}_x|}{R_x} \tag{7}$$

- *Rating deviation from mean agreement (RDMA):*

$$RDMA_a = \frac{\sum_{x=0}^{M_a} \frac{|r_{a,x} - \bar{r}_x|}{R_x}}{M_a} \tag{8}$$

- *Degree of similarity with top neighbors (DegSim):*

$$DegSim_b = \frac{\sum_{b \in neighbors(a)} S_{a,b}}{l} \tag{9}$$

- *Weighted deviation from mean agreement (WDMA):*

$$WDMA_a = \frac{\sum_{x=0}^{M_a} \frac{|r_{a,x} - \bar{r}_x|}{R_x^2}}{M_a} \tag{10}$$

2) *Model Specific Attributes*

Previous studies have shown that only generic attributes are not sufficient in differentiating the fake profiles and original users. So, there is a need of augmenting generic attributes with the attributes of particular attack type. These attributes are:

- *Mean Variance (MeanVar):*

It is used to detect average attacks in the system. We computed this metric as follows. $p_{a,t}$ denotes the set of ratings of target i.e. $r_{a,i} = r_{max}$ and $p_{a,f} = p_a - p_{a,t}$.

$$MeanVar_a = \frac{\sum_{i \in p_{a,f}} (r_{a,i} - \bar{r}_a)^2}{|p_{a,f}|} \tag{11}$$

- *Filler Mean Target Difference (FMTD):*

This metric detects the bandwagon attack profiles.

$$FMTD_a = \left| \left(\frac{\sum_{i \in p_{a,t}} r_{a,i}}{|p_{a,t}|} \right) - \left(\frac{\sum_{j \in p_{a,f}} r_{a,j}}{|p_{a,f}|} \right) \right| \tag{12}$$



V. EXPERIMENTAL RESULTS

A. Dataset Description –

We have used MovieLens-100k dataset in our experiments which consists of 943 viewers who have given 1,00,000 ratings to 1682 movies. It contains only those users who rated atleast 20 movies. Integers 1 to 5 are given to ratings where 1 is assigned to unfavorable items and 5 to most favorable items. 50 movies and 63 users are selected randomly such that distribution of ratings is close to the overall rating distribution. Attack is performed on each movie individually.

B. Experimental Setup –

The training set is created by selecting the set of profiles from system’s database that doesn’t contains malicious profiles and labeled them as *genuine*. Then a mixture of attacker’s data at several attack sizes and filler sizes are pushed into this training set and labeled as *fake*. Detection attributes for each profile in the training set are generated. Table 2 shows the structure of training dataset.

Table -2. Structure Of Training Dataset

Attribute ₁	Attribute ₂	...	Attribute _n	Label
------------------------	------------------------	-----	------------------------	-------

Classifiers i.e. RF, NB, J48, ZeroR, RBF n/w are trained using training dataset in WEKA. We then performed k fold cross validation technique to estimate the predictive models. Their performances are analyzed on the basis of recall, precision (PR), and f-measure (F-m) and we found NB and RF are the best performers. We then combined these models using vote method and built our new integrated model. Performance results of the models on attack sizes i.e. 5%, 10%, 15%, 20% and filler size of 50% for bandwagon and average attack are shown in table 3 and 4 respectively. Now, performance is analyzed on 25% attack size and 20%, 30%, 40%, 50% filler sizes for bandwagon and average attacks. The results are presented in table 5 and 6. Fig. 1 shows result of k-fold cross validation for average attack at different values of k.

We also ensemble random forest with three models using ensemble techniques. Firstly, it is combined with adaboostM1 model using boosting method, then with bagging and last with

```

Tester:      weka.experiment.PairedCorrectedTTester
Analysing:   Percent_correct
Datasets:    1
Resultsets:  4
Confidence:  0.05 (two tailed)
Sorted by:   -
Date:        14/5/16 5:49 PM
    
```

Dataset	(1) trees.RandomFo	(2) meta.AdaBoo	(3) meta.Bagging	(4) meta.Stackin
'all attributes'	(100) 99.71(2.01)	99.86(1.43) *	100.00(0.00)	100.00(0.00)
	(v/ /*)	(0/0/1)	(0/1/0)	(0/1/0)

Fig. 2. Improved accuracy of ensemble Random Forest model

stacking method. The accuracy of these models are shown in fig. 2, and we observed that boosted RF has 99.86% accuracy whereas simple RF has 99.71% accuracy. We can also see that there is a * in front of the accuracy of boosted RF which means that this algorithm is meaningful. Accuracy without * is meaningless. Hence boosted RF is more accurate than simple RF.

Cross Validation for average attack

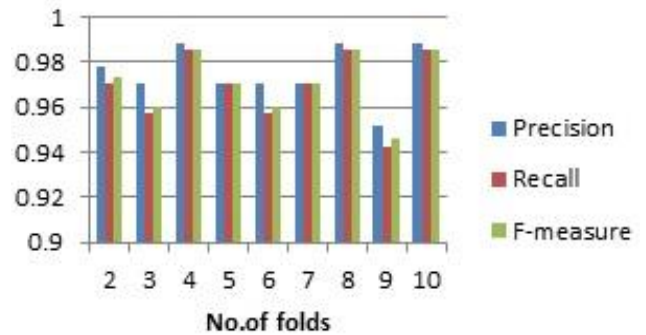


Fig. 1. k fold cross validation at 25% attack size and 50% filler size

VI. CONCLUSION AND FUTURE SCOPE

Detection of shillers is a key component for robust recommender system. Therefore, training sets are created by generating generic and model specific attributes, various classification models are demonstrated to distinguish the attack profiles. Their performance is analyzed using precision, recall and f-measure and we identified that NB and RF models are best performers. When integrated, k-fold cross validation proved that the new model outperforms in most of the cases. We also determined that integrated RF model is more accurate than simple one. More detection attributes can be used to detect the malicious profiles and this procedure can be implemented on other classification models and accuracy can also be improved.



Table -3. Performance Analysis Of Models For Bandwagon Attack At 50% Filler Size.

Attack Size	5%			10%			15%			20%		
	PR	Recall	F-m	PR	Recall	F-m	PR	Recall	F-m	PR	Recall	F-m
RBF n/w	0.972	0.972	0.963	0.952	0.967	0.962	0.948	0.949	0.948	0.891	0.893	0.896
J48	0.949	0.946	0.945	0.949	0.948	0.957	0.956	0.955	0.956	0.952	0.953	0.956
ZeroR	0.821	0.875	0.862	0.726	0.731	0.759	0.801	0.813	0.816	0.823	0.821	0.827
NB	0.976	0.979	0.98	0.981	0.98	0.976	0.98	0.979	0.976	0.970	0.971	0.971
RF	0.952	0.961	0.954	0.972	0.973	0.968	0.979	0.978	0.978	0.988	0.985	0.986
Integrated (NB+RF)	0.972	0.963	0.965	0.981	0.98	0.975	0.981	0.984	0.982	0.987	0.986	0.986

Table -4. Performance Analysis Of Models For Average Attack At 50% Filler Size.

Attack Size	5%			10%			15%			20%		
	PR	Recall	F-m	PR	Recall	F-m	PR	Recall	F-m	PR	Recall	F-m
RBF n/w	0.986	0.986	0.985	0.978	0.986	0.984	0.96	0.954	0.923	0.892	0.813	0.924
J48	0.952	0.942	0.946	0.941	0.943	0.95	0.942	0.946	0.957	0.947	0.951	0.92
ZeroR	0.834	0.913	0.872	0.83	0.826	0.862	0.766	0.875	0.817	0.792	0.791	0.793
NB	0.988	0.986	0.986	0.975	0.92	0.941	0.971	0.973	0.968	0.974	0.968	0.971
RF	0.96	0.957	0.958	0.986	0.984	0.974	0.987	0.986	0.986	0.981	0.982	0.983
Integrated (NB+RF)	0.978	0.971	0.973	0.98	0.986	0.985	0.988	0.981	0.973	0.971	0.975	0.976

Table -5. Performance Analysis Of Models For Bandwagon Attack At 25% Attack Size.

Filler Size	20%			30%			40%			50%		
	PR	Recall	F-m	PR	Recall	F-m	PR	Recall	F-m	PR	Recall	F-m
RBF n/w	0.89	0.892	0.893	0.902	0.906	0.904	0.921	0.924	0.923	0.935	0.933	0.934
J48	0.965	0.962	0.962	0.959	0.953	0.956	0.949	0.946	0.947	0.937	0.931	0.935
ZeroR	0.846	0.85	0.852	0.861	0.864	0.861	0.842	0.845	0.847	0.832	0.835	0.836
NB	0.978	0.976	0.975	0.983	0.984	0.984	0.976	0.975	0.978	0.98	0.982	0.981
RF	0.968	0.967	0.968	0.965	0.964	0.964	0.97	0.978	0.978	0.981	0.983	0.983
Integrated (NB+RF)	0.978	0.976	0.97	0.970	0.973	0.975	0.981	0.981	0.979	0.973	0.976	0.976

Table -6. Performance Analysis Of Models For Average Attack At 25% Attack Size.

Filler Size	20%			30%			40%			50%		
	PR	Recall	F-m	PR	Recall	F-m	PR	Recall	F-m	PR	Recall	F-m
RBF n/w	0.885	0.886	0.883	0.87	0.872	0.873	0.865	0.867	0.865	0.875	0.873	0.874
J48	0.95	0.952	0.954	0.961	0.965	0.95	0.956	0.955	0.953	0.946	0.948	0.949
ZeroR	0.81	0.813	0.816	0.847	0.846	0.842	0.802	0.803	0.803	0.791	0.79	0.786
NB	0.962	0.962	0.965	0.97	0.972	0.971	0.98	0.979	0.978	0.972	0.975	0.973
RF	0.972	0.971	0.974	0.964	0.964	0.965	0.974	0.976	0.976	0.98	0.982	0.983
Integrated (NB+RF)	0.974	0.973	0.973	0.973	0.976	0.975	0.978	0.979	0.978	0.981	0.983	0.984

VII. REFERENCES

- [1] S. Lam, J. Riedl, "Shilling Recommender Systems for Fun and Profit", *In Proceedings of the 13th international conference on World Wide Web*, ACM, pp. 393-402, 2004.
- [2] J. Herlocker, J. Konstan, A. Borchers, J. Riedl, "An algorithmic framework for performing collaborative filtering", *In Proceedings of the 22nd ACM Conference on Research and Development in Information Retrieval (SIGIR '99)*, pp. 230-237, 1999.
- [3] B. Mobasher, R. Burke, R. Bhaumik, C. Williams, "Effective Attack Models for Shilling Item-Based Collaborative Filtering Systems", *In Proc. of the 2005 WebKDD Workshop, Chicago, Illinois*, 2005.
- [4] M. Mahony, N. Hurley, G. Silvestre, "Recommender Systems: Attack Types and Strategies", *American Association for Artificial Intelligence*, pp. 334-339, 2005.
- [5] B. Mobasher, R. Burke, R. Bhaumik, J. Sandvig, "Attacks and Remedies in Collaborative Recommendation", *IEEE Intell. Syst.*, Vol. 22, No. 3, pp. 56-63, 2007.
- [6] M. O'Mahony, N. Hurley, N. Kushmerick, G. Silvestre, "Collaborative recommendation: A robustness analysis", *ACM Transactions on Internet Technology*, pp. 344-377, 2004.
- [7] R. Burke, B. Mobasher, R. Bhaumik, C. Williams, "Segment-based injection attacks against collaborative filtering recommender systems", *In Data Mining, Fifth IEEE International Conference*, 2005.



- [8] P. Chirita, W. Nejdl, C. Zamfir, "Preventing shilling attacks in online recommender systems", *In WIDM '05: Proc. of the 7th annual ACM Int'l workshop on Web information and data management*, pp. 67-74, 2005.
- [9] R. Burke, B. Mobasher, C. Williams, R. Bhaumik, "Classification Features for Attack Detection in Collaborative Recommender Systems", *In Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 542-547, 2006.
- [10] C. Williams, B. Mobasher, R. Burke, "Defending recommender systems: detection of profile injection attacks", *SOCA*, Vol. 1, No. 3, pp. 157-170, 2007.
- [11] A. Kumar, D. Garg, P. Rana, "Ensemble Approach to Detect Profile Injection Attack in Recommender System", *Fourth International Conference on Advances in Computing, Communications and Informatics (ICACCI-2015)*, IEEE, pp. 1734-1740, 2015.
- [12] F. Zhang, Q. Zhou, "Ensemble detection model for profile injection attacks in collaborative recommender systems based on BP neural network", *IET Information Security*, Vol. 9, No. 1, pp. 24-31, 2015.
- [13] J. Lee, D. Zhu, "Shilling Attack Detection—A New Approach for a Trustworthy Recommender System", *INFORMS Journal on Computing*, Vol. 24, No. 1, pp. 117-131, 2012.
- [14] M. Gao, B. Ling, Q. Yuan, Q. Xiong, L. Yang, "A Robust Collaborative Filtering Approach Based on User Relationships for Recommendation Systems", *Mathematical Problems in Engineering*, Vol. 2014, pp. 1-8, 2014.
- [15] W. Zhou, J. Wen, Y. Koh, Q. Xiong, M. Gao, G. Dobbie, S. Alam, "Shilling Attacks Detection in Recommender Systems Based on Target Item Analysis", *PLOS ONE*, Vol. 10, No. 7, pp. e0130968, 2015.
- [16] J. Dhimmar, R. Chauhan, "An accuracy Improvement of Detection of Profile-Injection Attacks in Recommender Systems using Outlier Analysis", *International Journal of Computer Applications*, Vol. 122, No. 10, pp. 22-27, 2015.
- [17] Z. Wu, J. Wu, J. Cao, D. Tao, "HySAD: A Semi-Supervised Hybrid Shilling Attack Detector for Trustworthy Product Recommendation", *Proc. of the 18th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, pp. 985-993, 2011.