

Does Betrayal Hide in Linguistics?

Machine Learning Betrayal Detection in Diplomacy Game

Liu Yehao, Lin Yuanhui, Li Hangqian

Ecole Polytechnique Fédérale de Lausanne

1 Abstract

Friendship is fickle, and there may be clues hidden in the words from friendship to betrayal[2]. In our work, we mainly based on a previous research on linguistic harbingers of betrayal occurs in one online strategy game *Diplomacy*. We use machine learning methods including random forest, support vector machine, logistic regression, and unsupervised learning to build a classifier that can judge whether a game will end up betrayal.

2 Introduction

Diplomacy is a strategy board game, which is shown in Fig.1. In this game, each player plays a role in one of the seven countries in World War I. The goal of this game is to control everywhere on the map. There are many seasons in one game. In each season, players can chat with each other to build friendship or move their armies to attack somewhere else on the map. There are thousands of different strategies in this game, and one of the most important decisions is to decide when to betray your allies so that you can get the most benefit.



Figure 1: Diplomacy game

Authors of the original paper analyzed many linguistic features from 500 games. They used different techniques to extract features from players' conversation and tried to find whether there are some signals that can foretell betrayal. These features including politeness level of the message, number of sentences exchanged and so on. They found out that the imbalances of features can be used to predict betrayal.

Can we predict whether a game will end up in betrayal with a higher accuracy using linguistic features? We no-

tice that in part four of the paper, the authors mentioned there are imbalances of features between the two players. Hence, it is worth trying to use these imbalances to predict if a game will end up betrayal. Also, in part five, the authors mentioned that there are imbalance of features before the seasons leading up to betrayal. Therefore, we put the features of four seasons before betrayal together and try to train model to predict whether a betrayal will happen using different machine learning techniques.

3 Data set

In our project, we use the data set from the original paper[2], which contains a collection of features of exchanged messages between players from 500 games in the Diplomacy games. A game consists of consecutive game seasons in which the two players can chat, support or attack each other. The reason to use this data set is that we could compare our result with the original paper directly. If we use a different data set, we might have a high probability to have different results due to different pre-processing methods and different data set. It is more convincing to say our model is better if our method increases the accuracy when we use the same data set.

When processing the data set, we first get 15 different features based on description in the original paper. After we examined the imbalance level and the correlation between the features, we keep 11 features as the features we consider for future prediction tasks. We get our training features in two ways. The first is to keep the following 11 variables before the last act of friendship in each game. The variables and descriptions are shown in Table 1.

Feature	Description
sent pos	Positive sentiment
sent neu	Neutral sentiment
sent neg	Negative sentiment
discourse comp	Discourse complexity
plan	Planning level
argu claim	Argumentation level of claim
argu premise	Argumentation level of premise
n request	Number of requests
politeness	Politeness
n words	Number of words
n sentences	Number of sentences

Table 1: Linguistic features

In this way, we get 719 imbalances of betrayers and victims without betrayal and 663 with betrayal. The second way is to obtain the 11 variables of the seasons before the last support season(included), which is a 4-season series of variables. If one of the players in the game does not speak at all for 4 seasons, then the game will be eliminated because there will be too many 0s in the data set.

4 Related work

In the previous work of Vlad, they explored linguistic features which are related to betrayal. They found that there are subtle but consistent patterns in how people communicate when they are going to betray. Friendships which will end up betrayal are imbalanced. For example, betrayers show more positiveness, plan less and be more polite than victims before the betrayal. Friendships which are doomed to break up are more likely to have an unbalanced talkativeness. Linguistic cues are used to predict whether the betrayal will happen in the end using logistic regression. The cross-validation accuracy is 0.57, which is higher than the human-predict result 0.52.

They also found out that changes in balance could mark imminent betrayal. As the breakdown approaches, the betrayer becomes more positive but less polite. While the victims request more and show more politeness. They also trained a classifier using the the features from the older seasons preceding the last friendly support to predict eventually betrayal. The best model achieved an F1 score of 0.31.

5 Methods

In this section, we are trying to predict whether betrayal will happen given the conversation between the two players. In the original work, the authors mentioned that there are imbalances of linguistic features between eventual betrayer and victim for those games end up betrayal. Hence, we try to use these features to predict whether the game will end up betrayal in the next turn. They also mentioned that there are imbalances of features for the seasons leading up to betrayal. Hence, we put the features of 4 seasons (which are 3 seasons before the last friendly support, and the last support season) together and try to train different classifiers to predict whether the game will end in betrayal or not.

5.1 Overall features prediction

5.1.1 Distribution change of feature imbalance

We first plot the imbalances of the 11 features between the betrayers and victims before the last support into a histogram, which is shown in Fig.2. After comparing, we find that imbalances peak moves significantly to the right in betrayal games, which means imbalance between the betrayer and victims become larger in requests, words, and sentences number in betrayal games than in non-betrayal games.

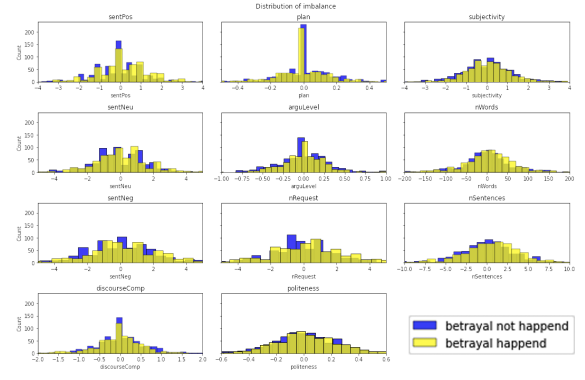


Figure 2: Distribution of imbalance in betrayal games and non-betrayal games

5.1.2 Pearson’s correlation of features

For each feature in the data set, we compare its distribution in the treated group with its distribution in the control group and compute the Pearson correlation, which is shown in Fig.3. We find that planing is highly related to discourse complexity, sentiment negative has a high correlation with words number, sentence number, and request number. Sentiment positive shows slightly lower correlation with words, sentences, and requests number. Comparison is highly correlated with expansive. Also words number, sentences number, and request number are highly correlated with each other.

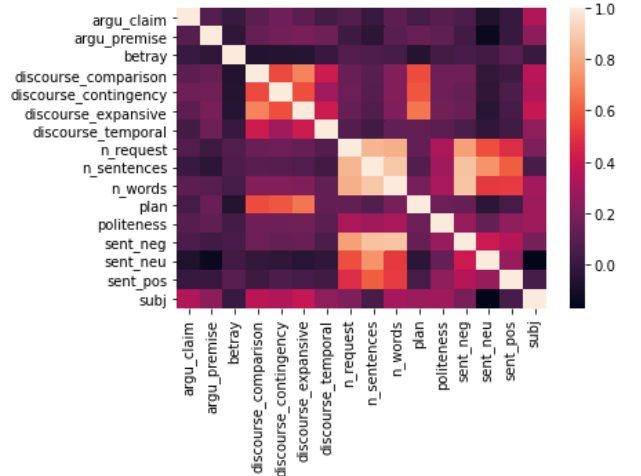


Figure 3: Heatmap of feature Pearson correlation

In our next work, we will use the highly correlated features for logistic regression analysis; meanwhile we will use cross-validation method.

5.1.3 Predictive power

We add the label betrayal to the data set, and mark it as 1 if betrayal occurs, and mark it as 0 if no betrayal occurs. we randomly split the data set into 0.80 training data set and 0.20 test set, and use logistic regression to train the classifier. The best accuracy is obtained by only

keeping request number, words number, sentences number, and the interaction between Negative sentiment and Sentences number features, which is 0.57. The AUC value is 0.56.

Then, we use 5-fold cross-validation to compute the accuracy. However, when we apply it on logistic regression, random forest and SVM, the highest accuracy we can get is only 0.53. So we are going to use the 4 seasons' time series to predict eventually betrayal.

5.2 Time series features prediction

Here we are going to present the second method we use. In this section, we combine the 11 variables of four seasons before betrayal together to get a 44 variables data features. Instead of treating each season individually, we treat each game individually and predict whether betrayal will happen in a certain season after given four seasons of conversation. Here we use some non-linear machine learning method as well. For example, we use support vector machine and random forest to fit the data set and predict whether they will betray or not.

5.2.1 Support vector machine

At first we use support vector machine. we split 80 percent of data as training set and the rest 20 percent as testing set. We try many different feature engineering but cannot find a good way to increase the accuracy. So, we do standardization for the training data set. We test a 5-fold cross-validation on our training set to find the best hyperparameters, and get a 0.565 accuracy on the cross validation. For the test set, we get a 0.622 accuracy. Compared with random guessing, which follows distribution $\text{Bin}(82, 0.5)$, the accuracy is much better. Using the command $1 - \text{pbinom}(51, 82, 0.5)$ in R language, we can get 0.018, which is the p-value. Our hypothesis is that our model is the same as random guessing. The p-value we just calculated is 0.018 lower than 0.05. Thus, we reject the null hypothesis that our model is the same as random guessing and take the alternative hypothesis that our model is better than random guessing.

5.2.2 Random Forest

For the random forest, we split 0.8 of data as training set and the rest 0.2 as testing set. We use a 5-fold cross-validation on our training set to find the best hyperparameters and get a 0.577 of accuracy on the cross validation. For the testing set, we get a 0.597 of accuracy on classifying them. Compared with random guessing, which follows distribution $\text{Bin}(82, 0.5)$. Use the command $1 - \text{pbinom}(49, 82, 0.5)$ in R language we get a p-value of 0.03. Our null hypothesis is that our model is the same as random guessing. Since the p-value we get is 0.03, we reject the null hypothesis that our model is the same as random guessing and take the alternative hypothesis that our model is better than random guessing.

Therefore, we can indeed find a strategy that performs better than random guessing on deciding if a game will end up betrayal.

5.2.3 K-Means

We also use unsupervised learning to classify the data set. But when we use TSNE and PCA to reduce the 11-dimension data set to two dimensions, we find the betrayal and non-betrayal data highly overlap, which means we cannot effectively classify it into two sets. The result of k-means also shows that the data set cannot be effectively classified in this way, which is shown in Fig.4.

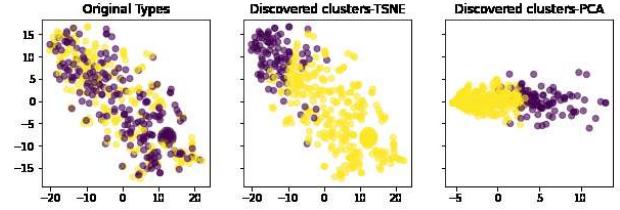


Figure 4: K-means result

6 Conclusions

We use features' imbalance based on all seasons before last friendly support and features data based on the four-seasons time series, and find it can predict the betrayal significantly better than human guessing which accuracy is 0.52. Hence, we can develop a software to detect whether the relationship will end up in betrayal or not in the next turn by inputting all the sentences.

In linguistic emotion detection, there may be some emotions that are not easy to predict, such as negative or Sarcasm Detection[1], which may have a negative impact on the true prediction of betrayal.. Vriji(2004) reported that experienced professionals such as police officers have an average accuracy of 65% when asked to detect lies [3]. Compared with it, our classifier do not have any information about players' facial expression or body language but can perform a 0.60 accuracy.

One of the challenging future work is to increase the accuracy of predicting. Maybe there are different kinds of way to process these verbal data so that the accuracy of lying detection can be increased.

References

- [1] T. Jain, N. Agrawal, G. Goyal, and N. Aggrawal. Sarcasm detection of tweets: A comparative study. In *2017 Tenth International Conference on Contemporary Computing (IC3)*, pages 1–6, Los Alamitos, CA, USA, aug 2017. IEEE Computer Society.
- [2] Vlad Niculae, Srijan Kumar, Jordan Boyd-Graber, and Cristian Danescu-Niculescu-Mizil. Linguistic harbingers of betrayal: A case study on an online strategy game. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1650–1659, Beijing, China, July 2015. Association for Computational Linguistics.

- [3] Aldert Vrij. Why professionals fail to catch liars and how they can improve. *Legal and criminological psychology*, 9(2):159–181, 2004.