

# Outcome Prediction of DOTA2 Based on Naive Bayes Classifier

Kaixiang Wang

School of Computer Science  
Communication University of China  
Beijing, China  
wangkaxiang1996@163.com

Wenqian Shang

School of Computer Science  
Communication University of China  
Beijing, China  
shangwenqian@163.com

**Abstract**—Although DOTA2 is a popular game around the world, no clear algorithm or software are designed to forecast the winning probability by analyzing the lineups. However, the author finds that Naive Bayes classifier, one of the most common classification algorithm, can analyze the lineups and predict the outcome according to the lineups and gives an improved Naive Bayes classifier. Using the DOTA2 data set published in the UCI Machine Learning Repository, we test Naive Bayes classifier's prediction of respective winning probability of both sides in the game. The results show that Naive Bayes classifier is a practical tool to analyze the lineups and predict the outcome based on players' choices.

**Keywords**—Naive Bayes; Classification Algorithm; DOTA2

## I. INTRODUCTION

Defense of the Ancients 2 (DOTA2) is an independent multiplayer online game, which develops from a map in the WarCraft3. DOTA2 swept the world with more than 50 million players. There have been over 50 million players of DOTA2, and hundreds of clubs and players participate in its global professional leagues every year, whereas no one has predicted the outcome of a game by analyzing the lineups in electronic sports field by now.

There are two teams, known as Radiant and Dire, that respectively consist of five players in each game. Players control one of 113 heroes (by December, 2016), each of which can be chosen only once. According to the Principle of Permutation and Combination, there are

$$C_{113}^5 \times C_{108}^5 \times \frac{1}{2} \approx 4.69 \times 10^{17} \quad (1)$$

different lineups in the game. Different lineups have different features, as each hero has its own strengths and weaknesses. Thus, not only professional players but also new players need scientific analysis of lineups eagerly.

Bayes classifier, which is based on Bayes theorem, has a solid mathematical foundation and stable efficiency of classification. Since the algorithm was invented over 250 years ago, it has been in an unparalleled position in both mathematics and machine learning fields. Comparing Naive Bayes classifier with other machine learning algorithms including decision tree and neural network algorithms, Michie (1994) et al found that Naive Bayes classifier is competitive in learning algorithms in

some specific cases. Bayes classification algorithm has two categories: one is Bayes network classifier, which is still developing, when considered the possible links between the various attributes; the other is Naive Bayes classifier, which assumes that the attributes are mutually independent.

In DOTA2, heroes are independent of one another, so there is no dependency and relevance among heroes. What's more, the game will not end in a draw, but absolutely decides victory and defeat. Therefore, we choose Naive Bayes classifier in our experiment. To examine the feasibility of it, we use the DOTA2 data set published in the UCI Machine Learning Repository to test Naive Bayes classifier's prediction of winning probability.

## II. BAYES CLASSIFIER AND IMPROVEMENT

### A. Bayes Classifier

A naive Bayes Classifier (NBC) is a probabilistic classifier which applies the Bayes' theorem with a strong (naive) assumption: it is assumed that the features describing the objects to be classified are statistically independent each other [8].

Suppose  $x(c, a_1, a_2, \dots, a_n)$  is an item to be classified, which includes  $n$  feature variables  $a_1, a_2, \dots, a_n$ . Bayes Classifier is based on Bayes theorem, so we know:

$$P(C|a_1, a_2, \dots, a_n) = \frac{P(a_1, a_2, \dots, a_n|C)}{P(a_1, a_2, \dots, a_n)} \quad (2)$$

Naive Bayes classifier is more advanced under the condition that all attributes are mutually independent. Equation (2) can be simplified to

$$P(C|a_1 a_2 \dots a_n) = \prod_{i=1}^n P(a_i|C) \quad (3)$$

Compared with other classifiers, Naive Bayes classifier has the following features:

- When the attributes are mutually independent, this classification is accurate.
- Not only a small amount of parameters need to be estimated in Naive Bayes Classification model, but it is less sensitive to missing data than Bayes network classifier, so the algorithm of Naive Bayes Classification is relatively simple.

- Naive Bayes algorithm is based on the conditional independent assumption, and the effect of the attribute value on a given class is independent of other attribute values.

### B. Improvement

If the data item to be classified by Naive Bayes classifier satisfies the following conditions:

- The data item has many attributes, but the number of attributes that determine the classification are significantly less than the number of attributes in total.
- Attributes are mutually independent.
- The attribute takes only two values or similar.

We can ignore the data that does not affect the classification so as to reduce the complexity of computations. Assuming  $V$  is a value of a dependent attribute, equation 3 can be rewritten as follow:

$$P(C|x) = P(C) \times \prod_{i=1}^n P(a_i = V | C). \quad (4)$$

## III. EXPERIMENT AND ANALYSIS

### A. Preprocess Data

The training set has 92,649 data in total, each of which has 117 attributes, shown as table I. The first attribute represents the result of a game with the value of 1 or -1. 1 means victory for the team Radiant and -1 represents victory for the team Dire. The second attribute to the fourth are the game ID information, having no impact on the result. Attribute 5 to the 117 represent the selection of heroes with the value of 1, 0 and -1. If the hero is selected by Radiant, the value will be 1, and -1 for team Dire. In addition, 0 means the hero has not been selected by any team. In this way, we retain game result attribute and hero selection attributes after processing through `DataFrame(GameResult, a1, a2, ..., a113)`.

TABLE I. DATA ITEM

No.	ID1	ID2	ID3	H1	H2	H3	...	H112	H113
1	223	8	2	1	0	-1	...	0	-1

Figure 1 illustrates the frequencies of heroes in the training set, and figure 2 shows the teams' percentage of wins.

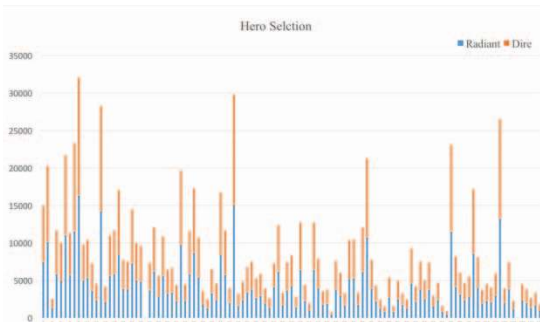


Fig. 1. Hero Selection in Training Set

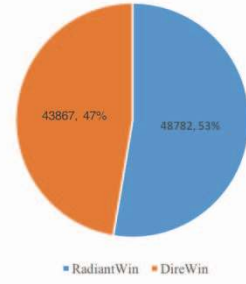


Fig. 2. Total Number of Wins for Each Team

### B. Analyze Goal

To get forecast result, we calculate

$$P(\text{Result}|a_1 a_2 \dots a_{113}) = \frac{P(a_1 a_2 \dots a_{113} | \text{Result})}{P(a_1 a_2 \dots a_{113})}. \quad (5)$$

As mentioned before, it can be simplified to

$$P(\text{Result}|a_1 a_2 \dots a_{113}) = \prod_{i=1}^{113} P(a_i | \text{Result}). \quad (6)$$

### C. Calculate Prior Probability

If the selection of heroes is  $a_i$  in a piece of data to be classified, its prior probability will be

$$P(a_i | \text{Radiant}) = \frac{\text{Count}(a_i)}{\text{Count}(\text{RadiantWin})}, \quad (7)$$

$$P(a_i | \text{Dire}) = \frac{\text{Count}(a_i)}{\text{Count}(\text{DireWin})}. \quad (8)$$

### D. Calculate Posterior Probability

Assuming the result is  $x$ , we calculate probability of  $x$  is 1 or -1:

Winning probability for Radiant,

$$P(\text{RadiantWin}|x) = \quad (9)$$

$$P(\text{RadiantWin}) \times \prod_{i=1}^{113} P(a_i | \text{RadiantWin}),$$

Winning probability for Dire,

$$P(\text{DireWin}|x) = \quad (10)$$

$$P(\text{DireWin}) \times \prod_{i=1}^{113} P(a_i | \text{DireWin}).$$

### E. Calculate Relative Winning Probability

Relative winning probability for Radiant:

$$P(\text{Radiant}) = \frac{P(\text{Radiant}|x)}{P(\text{Radiant}|x) + P(\text{Dire}|x)}, \quad (11)$$

Relative winning probability for Dire:

$$P(\text{Dire}) = \frac{P(\text{Dire}|x)}{P(\text{Radiant}|x) + P(\text{Dire}|x)}. \quad (12)$$

The Python code in Figure 3 calculates prior probabilities, posterior probabilities and relative winning probabilities.

```

1 #Prior Probability
2 def Prior(data,result):
3     count = 0
4     for item in data:
5         #first column is result
6         if item[ResultColumn] == result:
7             count += 1
8     return count/len(data)
9
10 #Posterior Probability
11 def Posterior(data,ai,i,result):
12     count1 = 0
13     count2 = 0
14     for item in data:
15         if item[ResultColumn] == result:
16             count1 += 1
17             if item[i] == ai:
18                 count2 += 1
19     return count2/count1
20
21
22 def CalculateProbability(dataTrain,dataTest):
23     count = 0
24     for item in dataTest:
25
26         PRadiant = Prior(dataTrain,1)
27         PDire = Prior(dataTrain,-1)
28
29         for i in range(2,len(item)):
30             PRadiant *= Posterior(dataTrain,item[i],i,1)
31             PDire *= Posterior(dataTrain,item[i],i,-1)
32
33     Radiant = PRadiant/(PRadiant+PDire)
34     Dire = PDire/(PRadiant+PDire)

```

Fig. 3. Python Code 1

#### F. Improve Classifier

As we mentioned in part II, ignoring unselected heroes enables us to reduce the quantities of prior probabilities to be calculated from 113 to 10, which greatly reduce computational complexity. Figure 4 shows the optimized code.

```

21 #Optimizing
22 def CalculateProbability(dataTrain,dataTest):
23     count = 0
24     for item in dataTest:
25
26         PRadiant = Prior(dataTrain,1)
27         PDire = Prior(dataTrain,-1)
28
29         for i in range(1,len(item)):
30             if (item[i]!=0):
31                 PRadiant *= Posterior(dataTrain,item[i],i,1)
32                 PDire *= Posterior(dataTrain,item[i],i,-1)
33
34     Radiant = PRadiant/(PRadiant+PDire)
35     Dire = PDire/(PRadiant+PDire)

```

Fig. 4. Python Code 2

#### G. Analyze Result

After testing 92649 pieces of data in the training set, we find that the accuracy of training set on training set is 85.33%.

Test set has 10294 pieces of data. The accuracy of test set on training set is lower than that of training set on training set. However, with the quantity of training set increasing, the correct rate tends to be stable at 58.99%.

TABLE II. ACCURACY OF 2 DATASET

	Training Set on Training Set	Test Set on Training set
Accuracy	85.33%	58.99%

It can be verified that the difference of the accuracy between original classifier and optimized classifier is at most 1%, which means the optimized classifier is feasible. This is shown in figure 3.

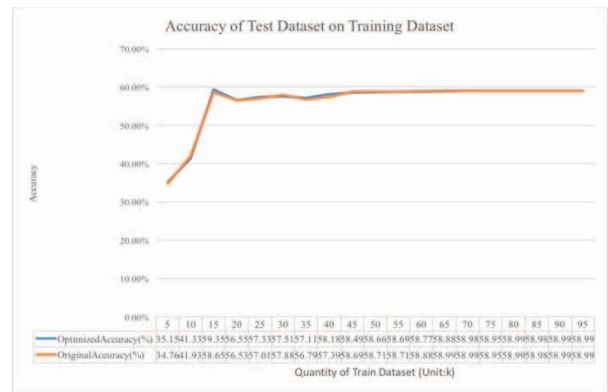


Fig. 5. Accuracy of Test Daset on Training Dataset

#### IV. CONCLUSION

Although the independent assumptions are often inaccurate in fact, some of the properties of the Naive Bayes classifier make it surprisingly effective in practice. The author provides a way to analyze lineups and forecast winning probability in Dota2 with Naive Bayes classifier, introduces the basic idea of how to analyze the game the Naive Bayes classifier and verifies the feasibility of analyze the game with quantitative data in Naive Bayes classifier model.

#### ACKNOWLEDGMENT

This paper is supported by National Training Programs of Innovation and Entrepreneurship for Undergraduates “ Mobile Service System of Communication University of China for Strudents”.

I would also like to express my sincere gratitude to Jiaqi Sun, my friend in Sun Yat-sen University, for revising the English version of this paper.

#### REFERENCES

- [1] J.M. Li, L.H. Sun, Q.R. Zhang and C.S. Zhang, Application of naive Bayes classifier to text Classification, Journal of Harbin Engineering University, Vol. 24 No.1, Feb. 2003.
- [2] Nino Prasetyo Hamal Pratama, Eko Mulyanto Yuniarno and Supeno Mardi Susiki Nugroho, Fuzzy Controller based AI for Dynamic Difficulty Adjustment for Defense of the Ancient 2 (DotA2), 2016 International Seminar on Intelligent Technology and Its Application (ISITIA), pp. 95-100.
- [3] Mahit Mertiya and Ashima Singh, Naive Bayes and Adjective Analysis for Sentiment Detection on Twitter, 978-1-5090-1286-2, ISBN 2016.
- [4] Z. Fu, Q. Chen, H. Zhang and S. Liu, Application of Naive Bayes Classifier in Stampede Risk Early-warning of Large-scale Activities, 2016 International Conference on Industrail Information - Computing Technology, Intelligent Technology, Industrial Information Integration (ICHCII), pp.174-180.
- [5] Peter Harrington, Machine Learning in Action, Posts & Telecom Press.
- [6] Tom M. Mitchell, Machine Learning, China Machine Press.
- [7] Zhu Jun and Hu Wenbo, Recent Advances in Bayesian Machine Learning, Journal of Computer Research and Development, 2015, pp.16-26.
- [8] C. De Stefano, F. Fontanella, A. Scotto di Freca, A Novel Naive Bayes Voting Strategy for Combining Classifiers, 2012 International Conference on Frontieres in Handwriting Recognition, pp.467-472.