

NANYANG TECHNOLOGICAL UNIVERSITY

March Machine Learning Mania 2016
Predict the 2016 NCAA Basketball Tournament

Submitted in Fulfilment of the Coursework Requirements
for the CE/CZ 4041 Machine Learning
by

| | |
|-------------------|-----------|
| Joe Tan Chin Yong | U1521434C |
| Liu Zeyan | U1421784C |
| Wei Yumou | U1320554F |
| Xie Dai | U1340229K |

School of Computer Science and Engineering

THIS PAGE INTENTIONALLY LEFT BLANK

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 5 |
| 1.1 | Background | 5 |
| 1.2 | Problem & Evaluation Method | 6 |
| 1.3 | Data Sets | 8 |
| 1.4 | Role Assignment | 8 |
| 2 | Preliminary Analysis | 9 |
| 2.1 | Challenges | 9 |
| 2.1.1 | Feature Selection | 9 |
| 2.1.2 | Interrelationship | 11 |
| 2.1.3 | The Curse of Model Popularity | 11 |
| 3 | Landmark Graph Construction | 13 |
| 4 | Time-Dependent Edge Cost Estimation | 15 |

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 1

Introduction

The National Collegiate Athletic Association (NCAA) Division I Men’s Basketball Tournament is one of the most popular annual sport festivities in the United States. Every year, the Tournament attracts a sizeable pool of audience, with the national champion becoming one of the hottest topics throughout the whole year. In the meanwhile, interests in predicting the winning team of a particular tournament match are escalating [3] and online machine-learning communities like Kaggle are organising annual competitions to encourage creative solutions. The purpose of this project is to present a feasible machine learning-based solution to Kaggle’s competition “March Machine Learning Mania 2016” [1] that accurately predicts a team’s *probability* of winning a particular tournament match based on historical match data provided by Kaggle.

1.1 Background

According to Wikipedia [5], the Tournament is played during every March and April based on the rule of single-elimination. Out of the 68 participating college basketball teams, 32 *conference* match champions are automatically qualified for the Tournament, while the other 36 teams are admitted at the discretion of a NCAA selection committee based on a criterion known as Rating Percentage Index [6].

The total 68 teams are then ranked by the selection committee from 1 to 68, and distributed amongst the four regions nominally known as East, West, South and Midwest. The top four teams receiving a rank from 1 to 4 are distributed to and given a *seed* of 1 in each of the four regions, followed by the next four teams with a rank of 5-8 that receive a seed of 2 in each region. The process continues until only the last eight teams are left whereby they have to fight with one of the other teams for the 16th seed position for each region, which marks the commencement of

the Tournament and is officially known as the *First Four* round.

At this point, the number of contesting teams are reduced to 64. During the next *First Round* in each region, a team with a higher seed position plays against a team with a lower seed position. For example, there are matchups between teams with the 1st seed and teams with the 16th seed, between the 2nd seed teams and the 15th seed teams and so on. Subsequently, the 32 winning teams advance to the *Second Round* where they play against one of the other teams, after which the 16 remaining teams are known as the *Sweet Sixteen*.

The number of contesting teams continues to halve until the four regional champions are determined. During the *National Semi-final*, the regional champion with the 1st seed position plays against the regional champion with the 4th seed position, while the other two teams play against each other. *National Final* is the last round and conducted between the two winners of the National Semi-final to determine the National Champion. However, there is no consolation game for the third place in the Tournament.

1.2 Problem & Evaluation Method

The problem of this Kaggle competition is to predict a team's probability of winning a particular tournament match. In a real tournament, there should be 68 teams and 67 matches in total because of the rule of single elimination (only the champion is not *eliminated* after the 67 matches). But how teams are paired up in these matches are not known in advance, except for the first 32 matches whose contesting pairs can be derived from the seed results. Therefore, Kaggle requires participants submit predictions for all possible matchups between any two of the 68 teams, which amount to $68 \times (68 - 1)/2 = 2278$ matches according to the Handshaking lemma.

According to Kaggle [1], the method for evaluating predictive models and compiling the leaderboard is based on cross entropy or log loss:

$$LogLoss = -\frac{1}{n} \sum_{i=1}^n [y_i \ln \hat{y}_i + (1 - y_i) \ln(1 - \hat{y}_i)] \quad (1.1)$$

where

- n is the number of games played
- \hat{y}_i is the predicted probability of team 1 beating team 2
- y_i is 1 if team 1 wins, 0 if team 2 wins

The goal of any predictive models is to *minimise* the log loss. Based on the LogLoss evaluation method, a completely incorrect prediction will lead to a score of ∞ , for example, a prediction of 1 when the actual outcome is 0. To avoid such an unpleasant score, Kaggle uses a threshold function to scale the submitted probability into a reasonable range as follows:

$$p' = \max(\min(p, 1 - 10^{-15}), 10^{-15}) \quad (1.2)$$

where p is the submitted probability and p' is the adjusted probability.

When using the LogLoss function to evaluate submissions, Kaggle assumes the ground truth variable y_i is discrete, namely $y_i \in \{0, 1\}$. But in general, y_i should be a continuous variable representing the *true* probability that Team 1 will beat Team 2, and \hat{y}_i still denotes the *predicted* probability. In this case, the cross entropy measures how close these two probabilities are. If there exists a *perfect knowledge predictor* that *always* gives correct predictions, namely $\hat{y}_i = y_i$ always holds, then its prediction performance can be described in Figure 1-1. Moreover, if the true probability follows an uniform distribution $y_i \sim U(0, 1)$, then in the long run the perfect knowledge predictor will have an expected score of

$$-\int_0^1 [y_i \ln y_i + (1 - y_i) \ln(1 - y_i)] dy_i = 0.5 \quad (1.3)$$

which is graphically equivalent to the height of a rectangle that shares the same base of length 1 as the performance curve's.

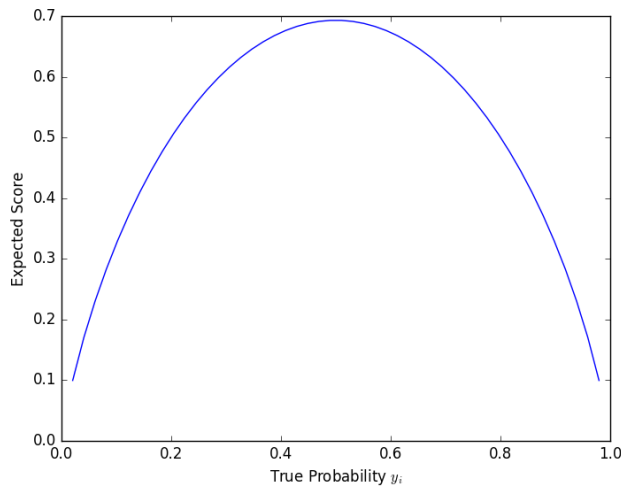


Figure 1-1: Performance of a Perfect Knowledge Predictor

1.3 Data Sets

Kaggle provides almost three decades' data about NCAA basketball matches, including matches during the regular seasons as well as the tournaments. The data sets are summarised in Table 1.1.

| Data Set | Description |
|------------------------------|---|
| RegularSeasonCompactResults | Game-by-game results during regular seasons from 1985-2015 |
| RegularSeasonDetailedResults | More detailed results during regular seasons from 2003-2015 |
| Seasons | Different seasons present in the dataset |
| Teams | Different college teams present in the dataset |
| TourneyCompactResults | Game-by-game tournament results from 1985-2015 |
| TourneyDetailedResults | More detailed tournament results from 2003-2015. |
| TourneySeeds | The seeds for all teams in a tournament |
| TourneySlots | The pair-ups between two teams based on their seeds |

Table 1.1: A Summary of the Data Sets

However, not all data sets are useful. Only two of them are used in this project, namely **RegularSeasonCompactResults** and **Teams**.

1.4 Role Assignment

Table 1.2 shows the assignment of roles for each team member in this project.

| Team member | Role | Responsibility |
|-------------------|----------------|-----------------------------------|
| Joe Tan Chin Yong | Data Scientist | To build predictive models |
| Liu Zeyan | Data Analyst | To optimise selected models |
| Wei Yumou | Data Scientist | To build predictive models |
| Xie Dai | Statistician | To provide mathematical knowledge |

Table 1.2: Role Assignment

Chapter 2

Preliminary Analysis

2.1 Challenges

In the attempt to tackle the problem, various challenges arise that are worth mentioning.

2.1.1 Feature Selection

Features are pieces of information that may be useful for predictions. It is generous of Kaggle to provide a comprehensive collection of data. However, not all data is relevant to making predictions. Before a set of features can be selected, careful choices must be made on which data sets to use and which portion of that data set is relevant to solving the problem.

Regardless of the machine-learning technique used, the most obvious choice of data sets to use is the historical match records for both regular and tournament seasons. However, some data sets contain match records that date back to as early as 1985, which are no longer relevant in today's context. After all, the NCAA tournament teams consist of *college students* who can only stay with a team for a maximum of four years before graduation. As players constantly leave and join the team, it is difficult to quantify the effect brought by changes in a team's composition on the strength of that team, given only the team-level data. Moreover, there were some substantial updates on the Tournament's rules at the beginning of the 2008-2009 season [4], which also affected the strategies teams used in the subsequent tournaments. Therefore, only the most recent match records are useful for prediction. For the purpose of this project, a four-year window from 2013 to 2016 is selected and all data used in this project fall within this window. In addition, only the match records from the regular seasons are actually used, since there are only a few records from the tournament seasons which add little value to making predictions.

There are a number of potential features that can be selected from the data set **RegularSeasonDetailedResult**. The potential features are further categorised as basic and additional features. Table 2.1 gives a summary of the basic features.

| Feature | Description |
|---------|---|
| Wteam | The id number of the team that won the game |
| Wscore | The number of points scored by the winning team |
| Lteam | The id number of the team that lost the game |
| Lscore | The number of points scored by the losing team |
| Wloc | The location of the winning team |
| Numot | The number of overtime periods in the game |

Table 2.1: A Summary of the Basic Features

The additional features describe both *offensive* and *defensive* strengths of both winning team and losing team. Table 2.2 gives a summary of features describing offensive strengths, while Table 2.3 focuses on defensive strengths. Although the two tables list features from the winning team’s perspective, another duplicated set of features also exists for the losing team.

| Feature | Description |
|---------|---|
| Wfgm | The number of field goals made by the winning team |
| Wfgm3 | The number of three pointers made by the winning team |
| Wftm | The number of free throws made by the winning team |
| Wor | The number of offensive rebounds by the winning team |
| Wast | The number of assists by the winning team |
| Wstl | The number of steals by the winning team |
| Wblk | The number of blocks by the winning team |

Table 2.2: A Summary of the Offensive Features

| Feature | Description |
|---------|--|
| Wfga | The number of field goals attempted by the winning team |
| Wfga3 | The number of three pointers attempted by the winning team |
| Wfta | The number of free throws attempted by the winning team |
| Wdr | The number of defensive rebounds by the winning team |

Table 2.3: A Summary of the Defensive Features

The potentially many features make model selection challenging, because they create an exponential number of possible models. If a feature has m possible values, to select the best models based on n features, at least $\Theta(n^m)$ models have to be examined, which is impractical given that only one submission is allowed at one time. Therefore, this project uses a greedy strategy where all other features are fixed at their optimal values when one feature varies to select a suboptimal model.

2.1.2 Interrelationship

No team is in isolation. The tournament matches are interactive processes whereby complex interrelationships exist amongst all the contesting teams, which adds another layer of complication in the attempt to predict match results. For example, given the history match records that Team A beaten Team B and Team C lost to Team B, one should intuitively conclude that Team A should have a higher probability of beating Team C. But what if another record shows that Team A once lost to Team C? In that case, the relative strength levels of the three teams will be hard to determine based purely on that intuition. Moreover, the match whose result is to be predicted may be the very *first* match ever between two teams. In other words, there are no historical records that give a direct assessment on the two teams' relative strengths. Such lack of knowledge must be complemented by some forms of inference based on the interrelationships amongst the teams. So a good predictive model should not only be able to take into consideration the current game record, but also explore the interrelationships amongst all the game records and generalise on unseen matches.

2.1.3 The Curse of Model Popularity

The problem to solve can be categorised as a classification problem under supervised learning. But since the final answers to be submitted are in fact *probabilities*, non-probabilistic models like support vector machines will hardly work. Therefore, the first a few models attempted in this project include the most popular ones: logistic regression and multilayer perceptrons.

The inputs to these popular models are the seed positions of each team, since the seeds represent an official view on the relative strength of each team. However, the performances of these models are not as good as expected: the logistic model only gives a 283rd position on the leaderboard. Moreover, only teams in a tournament are assigned a seed but the number of tournament matches is not large enough to support accurate predictions. Although [2] suggests a combination of logistic regression and Markov chain to predict match outcomes, these popular models alone based on simply seeds are unlikely to give a good result.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 3

Landmark Graph Construction

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 4

Time-Dependent Edge Cost Estimation

THIS PAGE INTENTIONALLY LEFT BLANK

Bibliography

- [1] Kaggle Inc., *March machine learning mania 2016*, March 2016.
- [2] Paul Kvam and Joel S. Sokol, *A logistic regression/markov chain model for ncaa basketball*, Naval Research Logistics **53** (2006).
- [3] Net Prophet, *Exploring algorithms for predicting ncaa basketball games*, April 2017.
- [4] Net Prophet, *Five mistakes kaggle competitors should avoid*, February 2015.
- [5] Wikipedia, *Ncaa division I men's basketball tournament*, April 2017.
- [6] Wikipedia, *Rating percentage index*, March 2017.