

Improving Duckworth-Lewis: Statistical Methods for Revising Score Targets in Limited-Overs Cricket

Matthew Knowles



UNIVERSITY
of York

Department of Mathematics,
University of York,
United Kingdom

A dissertation submitted in partial fulfillment of the requirements for the degree of Master of Mathematics

Abstract

Write this last

Acknowledgements

The author would like to thank Dr. Jessica Hargreaves for their invaluable supervision and advice throughout this project. Secondly, my eternal thanks and admiration for my friends in York, who have made my time here so wonderful from the very beginning, and for all the support you have given me every step of the way. Finally, to the UYMCC, who have enriched my university experience beyond belief.

Notes

All code written in support of this project can be found on GitHub at:
<https://github.com/mattnowles314/mastersThesis>.

Contents

1	Background	4
2	Data	6
3	The DLS Method in Detail	7
4	Overview of Pattern Recognition Techniques	11
5	Exploratory Data Analysis	12
6	Implementing SPR Methods on Match Data	14
7	Analysis and Results	15
8	Conclusion	16
A	Python Code for Chapter 2	17

List of Figures

3.1	Graph showing how the rate at which runs are scored decays as a game progresses.	7
3.2	Plot of average runs scored at the fall of each wicket. Averages taken over 1437 innings. .	10
5.1	Density of all teams for first wicket falling	12
5.2	Overall density plot for FOW 1	12
5.3	Density of all teams for first wicket falling	12
5.4	Overall density plot for FOW 5	12
5.5	Density of all teams for first wicket falling	13
5.6	Overall density plot for FOW 9	13

Chapter 1

Background

1.1 Limited-Overs Cricket and the need for Duckworth-Lewis

Cricket is a game played by two teams, each with 11 players. The objective for the batting team is to score as many “runs” as possible without losing 10 of their batsmen¹, who can be “out” in a variety of ways. Cricket is played in “overs”, each over lasting 6 deliveries. Traditionally, the game lasted for 4 days and there was no limit to the number of overs the bowling team could bowl.

In 1963, the cricketing calendar in the UK had for the first time a different format of the game amended to it. The “Gillette Cup” introduced a form of cricket wherein each team has 1 innings, lasting 50 overs. The idea was that this competition, coming to a conclusion within the space of a day, would increase spectator numbers and by extension, ticket sales.

However, given the tournament-based nature of this new competition, we have the natural need for a definitive result. This is something that is not always guaranteed in “first class” cricket, where draws are common.² As such, in order to allow for a result to be determined when a game is cut short, the idea of target-revision was introduced. This meant that a team could be set a roughly equivalent run target off of a reduced number of overs that would still classify them as the winners.

1.2 Problems raised with DLS

In [?], the authors found that DLS has a bias towards not only the team batting first, but whoever won the coin toss at the start of the game. The winner of the toss chooses whether their side will bat or bowl in the first innings. They go on to propose a simple extension for reducing these biases, but we won’t cover that in this project.

1.3 Project Aims

bar

1.4 Review of current Literature

We begin to look at literature on this topic with the paper published by F. Duckworth himself in 1998 [?]. In this paper, Duckworth proposes the “D/L” method based on 5 principles. However, the issue with the sentence that appears after defining the function $Z(u, w)$. “Commercial confidentiality prevents the disclosure of the mathematical definitions of these functions”. The claim is that these functions have been defined via experimentation. This however gives rise to the first issue: the first T20 game of cricket was not played until 2003. So clearly, D/L was not designed with T20 in mind, and so the functions

¹There have to be two batters on the pitch, so if 10 wickets are lost, it leaves one batter stranded.

²Note that “draw” and “tie” are not interchangeable terminology in cricketing terms.

derived from experimentation and research will not be accurate for T20 cricket.

In 2004, a year after the first T20 matches were played, Duckworth and Lewis published another paper [?], in which they report that the table used for calculating the method can be employed by

Chapter 2

Data

2.1 Data Origin

The primary source of data for carrying out this project was downloaded from “cricsheet”¹ and stored locally on a private server. In total there are 2167 individual matches of data. Each in a JSON format. These cover matches ranging from the 3rd of January, 2004. Up to the 20th of July, 2021.

2.2 Attributes

Each JSON file contains a considerable amount of metadata surrounding the match in question. Along with ball-by-ball data for the entire match. We have access to attributes such as the date, where the match was played, the entire teamsheet for both teams, who the officials were, who won- and by what margin, who won the toss; and many others.

We also have the ball-by-ball data. So for every ball bowled, it gives who were the striking and non-striking batsmen, how many runs were scored and how. It also details if a wicket was taken that ball, and how.

2.3 Pre-Processing

Many Python scripts were written to get data in an appropriate form for analysis.

¹<https://cricsheet.org/>

Chapter 3

The DLS Method in Detail

We now look at the mathematics behind the D/L, and DLS methods. D/L being the original method, and DLS the method that Stern helped to revise. In the original paper, the authors state “Commercial confidentiality prevents the disclosure of the mathematical definitions of these functions.” [?]. Which is naturally a slight problem for this, but what we can do is instead use sample values based on data we have to look at how these functions behave, so all is certainly not lost.

3.1 Origins: Duckworth and Lewis

We begin by looking at the original paper. The first thing to establish is how many runs are scored, on average, in a given number of overs. This is given by the equation:

$$Z(u) = Z_0[1 - \exp(-bu)] \quad (3.1)$$

Where u is the number of overs, b is the exponential decay constant, and Z_0 is the average total score in first class cricket, but with one-day rules imposed.

Now because we don't have access to actual values for Z_0 or b , a plot to see what equation 3.1 looks like was created by using 3 sample values for Z_0 . For b , it was a process of trial and error to find a value that resulted in the graph having a similar shape to original figures in the D/L paper.

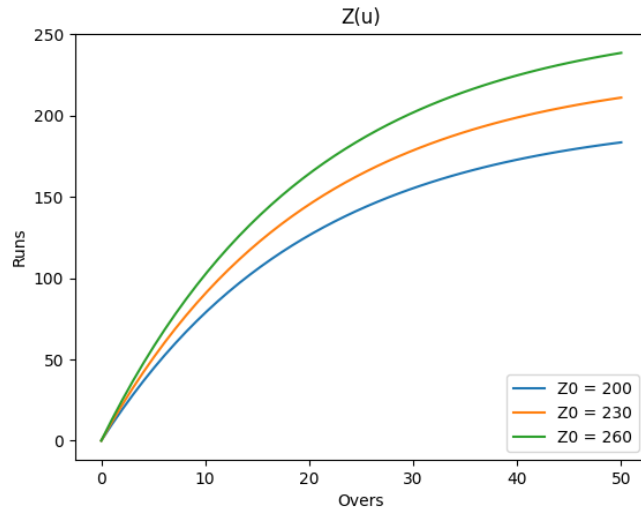


Figure 3.1: Graph showing how the rate at which runs are scored decays as a game progresses.

However, we have not yet looked at what happens when wickets are lost. To introduce this metric, equation 3.1 is revised to incorporate the scenario that w wickets have been lost, and that there are u overs remaining. The revised equation is given as follows:

$$Z(u, w) = Z_0(w)[1 - \exp(-b(w)u)] \quad (3.2)$$

Where now, we have $Z_0(w)$ giving the average total score from the last $10 - w$ wickets in first class cricket. We now also have $b(w)$ as the exponential decay constant, which now changes depending on wickets lost.

With this in mind, we now look at the specific case of equation 3.2 with $u = N$ and $w = 0$, namely, the conditions at the start of an N-over innings. We have:

$$Z(N, 0) = Z_0[1 - \exp(-bN)]. \quad (3.3)$$

Which we then incorporate into the ratio

$$P(u, w) = \frac{Z(u, w)}{Z(N, 0)}. \quad (3.4)$$

The ratio 3.4 gives, keeping in mind there are u overs still to be bowled, with w wickets lost, the average proportion of the runs that still need to be scored in the innings. It is this ratio that is where the revised scores come from. Let us now look a bit more at how that works practically.

Example 3.1.1. Assume there is a break in the second innings (due to rain or similar), which results in the second team missing some overs. Let u_1, u_2 be the number of overs played before the break, and available after it respectively. We impose the condition that $u_2 < u_1$. At the time of the break, the second team had lost w wickets. The aim is to adjust the required score to account for the $u_1 - u_2$ overs they have lost. The winning “resources” available are given by

$$R_2 = [1 - P(u_1, w) + P(u_2, w)].$$

Which means, if the first team batting scored S runs, then the new target is given by

$$T = \lceil SR_2 \rceil$$

3.2 Improvements by Stern

We now look at the DLS method, introduced by Stern in [?] to extend the original D/L method. The motivation for DLS was that very high scoring cricket matches presented a pattern of straightening towards average run rate which is not uniform across the innings considered (around 500). To account for this, Stern introduces an additional damping factor to the equation for D/L. The updated equation is now given by:

$$Z_{dls}(u, w, \lambda) = Z_0 F_w \lambda^{n_w + 1} \left\{ 1 - e^{-ubg(u, \lambda)/F_w \lambda_w^n} \right\}. \quad (3.5)$$

This equation assumes u overs remaining, w wickets down in an M-over match with a final total of S . Here, we have:

$$g(u, \lambda) = \left(\frac{u}{50} \right)^{-(1 + \alpha_\lambda + \beta_\lambda u)} \quad (3.6)$$

In the above, we define $\alpha_\lambda = -1/\{1 + c_1(\lambda - 1)e^{-c_2(\lambda - 1)}\}$ and $\beta_\lambda = -c_3(\lambda - 1)e^{-c_4(\lambda - 1)}$. The constants c_1, \dots, c_4 are based on what Stern describes as a “detailed examination” of high scoring cricket matches.

3.3 Extension: Bayesian Modelling

This section will look at the work of paper [?], which used Bayesian modelling to try and improve the score predictions of the D/L method. The authors used the same dataset as we are using for the dissertation, although over a slightly smaller range of games. Only games between 2005-2017 are used, whereas our dataset goes up to 2021. Note that to keep consistent with the fact we are talking about a different model now, we switch from using $Z(u, w)$, to using $R(u, w)$, as to keep consistency with the different papers being looked at. They start by introducing the following nonlinear regression model:

$$\bar{R}(u, w) \sim N(m(u, w; \theta), \frac{\sigma^2}{n_{uw}}) \quad (3.7)$$

Where $\bar{R}(u, w)$ is the sample average of runs scored by a team from the total number of matches in the data set. We have $m(u, w; \theta)$ as the corresponding modeled population average of runs scored by a team when a considerable amount of games are taken into account. θ denotes a vector of parameters that will be specified later on. Since $R(u, w)$ is not calculated in all games, the average score is taken over all matches where scores are present, this is given by n_{uw} . The sample mean $\bar{R}(u, w)$ is the calculated over this quantity. This is actually the point which motivates using Bayesian statistics to extend the D/L method. The authors report that approximately 26.8% of values for $R(u, w)$ were missing in the dataset, so by using a Bayesian inference model, we can use the posterior predictive distribution to account for missing values. Which in turn should increase the predictive accuracy of this model.

We return now to look at the $m(u, w; \theta)$ parameter that appears in 3.7. This mean function, based on the plots produced in the original D/L paper, (see 3.1 for an example), the authors adopted an exponential decay. We can see why this choice makes sense from figure 3.1. The resources considered by this method, namely runs and wickets, do decrease exponentially as a game progresses. $m(u, w; \theta)$ depends on the parameters $a_w = Z_0(w)$ and $b_w = b(w)$, each one depending on the number of wickets fallen at that given time, with u overs remaining. We therefore have:

$$m(u, w; \theta) = a_w(1 - \exp(-b_w u)) \quad (3.8)$$

$$\theta = \{(a_w, b_w); w \in \mathcal{W} = \{0, 1, \dots, 9\}\} \quad (3.9)$$

Before proceeding with defining the prior specifications on the parameters for this model, we need the following distribution:

Definition 3.3.1. $X \sim Ga(a, b)$ has a **Gamma Distribution** with mean $\frac{a}{b}$ and variance $\frac{a}{b^2}$ if its probability density function is

$$\frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}.$$

We can now define the prior specifications on the parameters. First, we fix A_0 and B_0 large enough such that $R(50, 0) < A_0$. We initialise the model with $a_0 \sim U(0, A_0)$ and $b_0 \sim U(0, B_0)$. Then given any pair (a_0, b_0) and for $w = 0, 1, \dots, 8$, we generate:

$$a_{w+1} | \sigma^2, a_w, b_w \sim U(0, a_w) \quad (3.10)$$

$$b_{w+1} | \sigma^2, a_{w+1}, b_w, a_w \sim U\left(0, \frac{a_w b_w}{a_{w+1}}\right) \quad (3.11)$$

$$\frac{1}{\sigma^2} \sim Ga(a, b). \quad (3.12)$$

Above, $U(a, b)$ represents a uniform distribution over the open interval (a, b) . Now that the prior distribution is obtained, the likelihood function is then given by:

$$L(\theta, \sigma^2) = \left(\frac{1}{\sigma / \sqrt{n_{uw}}} \right)^{500} \exp \left\{ -\frac{1}{2(\sigma^2 / n_{uw})} \sum_{u=1}^{50} \sum_{w=0}^9 (R(u, w) - m(u, a_w, b_w))^2 \right\} \quad (3.13)$$

The authors chose their parameters A_0 , B_0 , a and b such that the prior is not sensitive to the posterior inference. The values chosen were therefore $A_0 = 200$, $B_0 = 100$ and $a = b = 0.1$.

One thing briefly mentioned in chapter 3 of this paper is the average score at the fall of each wicket. This isn't discussed further, but it's interesting to see how it behaves. So we take a brief detour to look at it. This can be seen in figure 3.2.

Initially, one may think that this is normally distributed. It certainly holds a bell curve-like shape. However, a Kolmogorov-Smirnov test [?] was performed to see if these values were normally distributed,

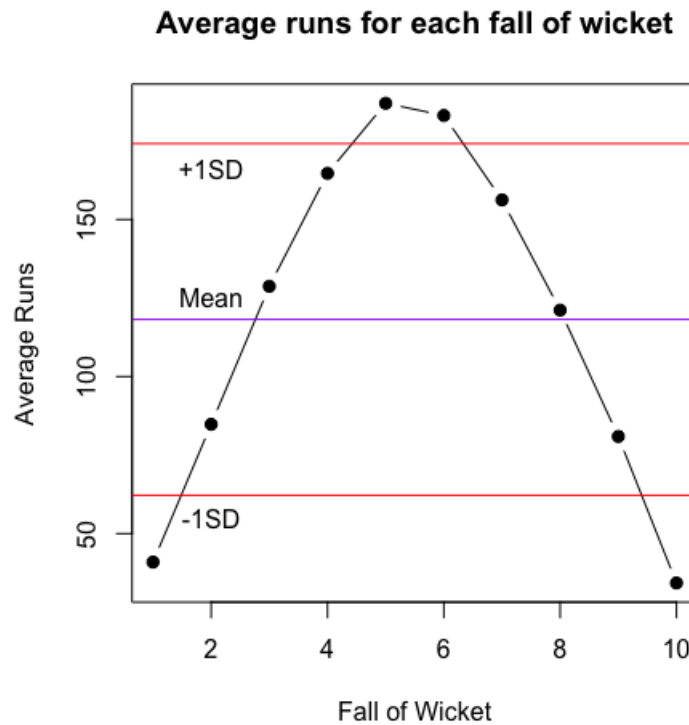


Figure 3.2: Plot of average runs scored at the fall of each wicket. Averages taken over 1437 innings.

and it turns out they aren't. This test was performed using the R function *ks.test()*, which returned a p-value of 2.2×10^{-16} for the two-sided alternative hypothesis parameter.

Does the shape of this curve make sense? Well yes, because usually batters 3/4 are actually the best in the side, so the fact they put on more runs is not a surprise. But what this does do is reinforce the idea behind using an exponential decay for modelling resources. This is because cumulatively, the tail end of batters (the last 4) will not put on as many runs as the higher/middle order.

Most of the work from this paper went into creating a better resource table for calculating revised scores. Which isn't particularly relevant to this project, but what is is the section on score prediction. The authors split every match at the 35th over to try and predict the final score. They found the bayesian model gives a better prediction of final match scores. Figure 3 of their paper outlines these results.

Chapter 4

Overview of Pattern Recognition Techniques

4.1 Model 1

Bar

Chapter 5

Exploratory Data Analysis

In this chapter, we are looking only at the first innings of the games, and only those games in which the full 50 overs were played. The reason for this is the models we will build are going to try and predict a score as if a full innings has been played.

We begin our exploration of the data with a look at how the density of the runs scored per fall of wicket changes. This has been done for each individual team in the dataset, and in figures 5.2-5.6, we can see how this evolves.

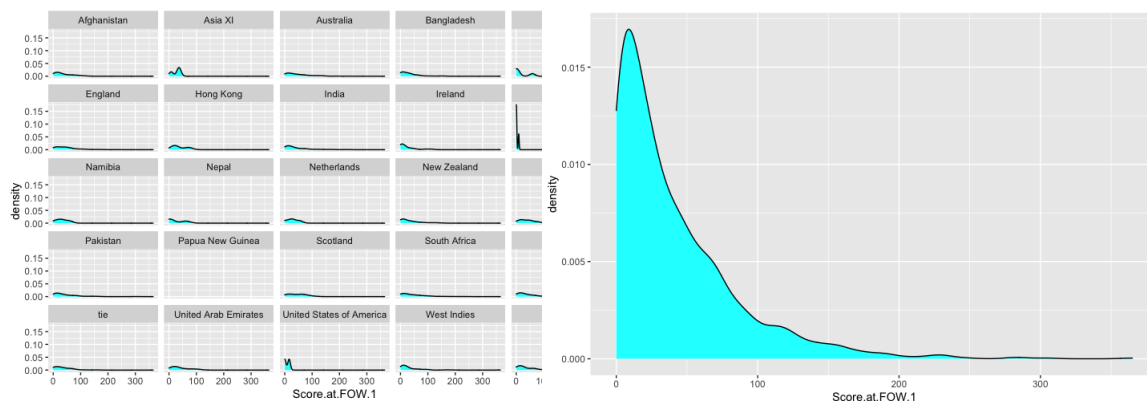


Figure 5.1: Density of all teams for first wicket falling

1

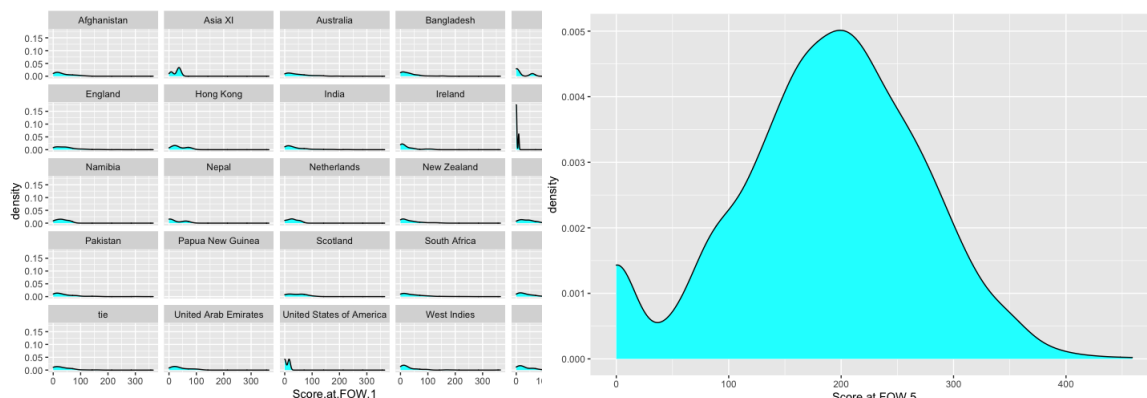


Figure 5.3: Density of all teams for first wicket falling

5

In 5.2, we see the density is heavily skewed to the left. This makes sense, as the bowling team will presumably be starting their innings by using their best bowlers, who will be hunting to get wickets early on. In 5.4, we see a much more normally distributed density function. But in actual fact, we see

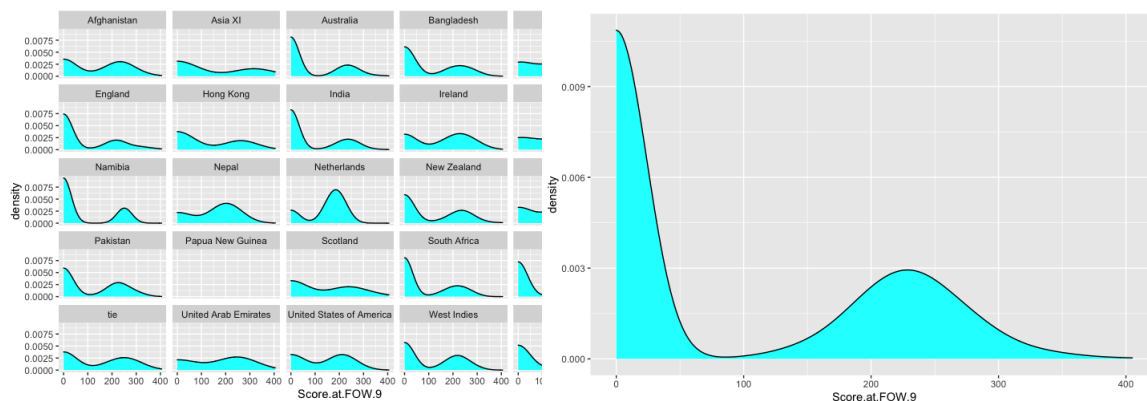


Figure 5.5: Density of all teams for first wicket falling
 Figure 5.6: Overall density plot for FOW 9

this interesting second, smaller peak appearing lower down in the score. Does this make sense? It's certainly not surprising. What these two peaks exemplify is the fact games can go heavily in favour of the bowling team, which can be seen in the first small peak, wherein they have taken a lot of wickets in quick succession, meaning the later order batters are coming in earlier than usual. Secondly, it shows when the batting team is having a good day, because we have this much larger peak around the 200 runs mark.

Finally, in 5.6, we can see that the earlier bowling advantage peak is much higher, because the lower order batters are traditionally less skilled at batting, and so the bowling team have a distinct advantage in taking wickets against these players. But we also see the second, batting-favoured peak is no much lower. This corresponds to the scenario in which the earlier batters have laid a good foundation of the game, and the lower-order batters have not had to contribute much to the score.

Chapter 6

Implementing SPR Methods on Match Data

6.1 Method 1

Foo

Chapter 7

Analysis and Results

7.1 Method 1

foo

Chapter 8

Conclusion

foo

Appendix A

Python Code for Chapter 2

Listing A.1: Data Wrangling Code

```
import json
import os
import csv

#Moves all the games from a set of files that were decided by DLS
def get_dls_results(inputPath, outputPath):
    os.chdir(inputPath)
    for f in os.listdir():
        with open(inputPath+f, "r") as file:
            data = json.load(file)
            if("method" in data["info"]["outcome"]):
                if(data["info"]["outcome"]["method"] == "D/L"):
                    os.system("mv_" + f + "_" + outputPath)

def get_match_data(matchPath, csvPath):
    os.chdir(matchPath)
    for f in os.listdir():
        with open(matchPath+f, "r") as file:
            data = json.load(file)

            homeTeam = data["info"]["teams"][0]
            awayTeam = data["info"]["teams"][1]
            ground = data["info"]["venue"]

            if "winner" in data["info"]["outcome"]:
                winner = data["info"]["outcome"]["winner"]
                if "runs" in data["info"]["outcome"]["by"]:
                    runsWonBy = data["info"]["outcome"]["by"]["runs"]
                    wicksWonBy = 0
                else:
                    wicksWonBy = data["info"]["outcome"]["by"]["wickets"]
                    runsWonBy = 0
            else:
                winner = "tie"

            targetOvers = data["innings"][1]["target"]["overs"]
            targetRuns = data["innings"][1]["target"]["runs"]
```

```
row = homeTeam,awayTeam,ground , winner , str (runsWonBy) , str (wicksWonBy  
    ) , str (targetRuns) , str (targetOvers)  
csvF = open(csvPath , "a")  
writer = csv.writer(csvF)  
writer.writerow(row)  
csvF.close()
```

```
get_match_data("/Users/matthew/Documents/Uni/Dissertation/Data/dls/" , "/"  
    Users/matthew/Documents/Uni/Dissertation/Data/dlsMaster.csv")
```