

Copyright is owned by the Author of the thesis. Permission is given for a copy to be downloaded by an individual for the purpose of research and private study only. The thesis may not be reproduced elsewhere without the permission of the Author.

Statistical Methods for Cricket Team Selection

A THESIS PRESENTED IN PARTIAL FULFILMENT
OF THE REQUIREMENT OF THE DEGREE OF
MASTER OF APPLIED STATISTICS
AT MASSEY UNIVERSITY, ALBANY
NEW ZEALAND

Paul J. Bracewell

1999

Abstract

Cricket generates a large amount of data for both batsmen and bowlers. Methods for using this data to select a cricket team are examined. Utilising the assumption that an individual's natural ability is expressed via performance outputs, this thesis seeks to describe and understand the underlying statistical processes of player performance. Randomness is tested for and then the distributional properties of the data are sought.

This information is then used to monitor the estimate of natural ability via widely accepted control methods, such as Shewhart control charts, CUSUM, EWMA and multivariate versions of these procedures. To accommodate the distribution presented by batting scores, a new control chart based on quartiles is also studied.

Further, ranking and selection procedures employ the estimates of individual ability to select the best individuals and note the probability of correct selection.

Major contributions of this study include:

- a) Development of performance measures for cricket
- b) 2 – Dimensional runs test, with further applicability outside cricket.
- c) Statistical interpretation specific to cricket
 - Outliers are very important
 - Form is autocorrelation
 - Zone rules for cricket needed to detect good/poor performance
 - Relatively short nominal ARL's
- d) Control Chart based on quantiles to preserve outlier influences in a non-parametric procedure.
- e) The recommendation of appropriate tools for monitoring batting, bowling and all-rounder performance and also choosing man of the match.
- f) Discriminates between different types of bowlers using the consistency of their performance measures.
- g) Evaluates the members of a team relative to potential contenders.

Acknowledgements

I am fortunate to have had two exceptional people help me through this thesis that I need to thank. The first, my supervisor, Dr Denny Meyer with whom I'm very grateful to have been taught, guided and inspired by in the pursuit of knowledge. I know that I'm not the only student to be appreciative of Denny's assistance.

Secondly, Jenita Higgs who's ongoing support enabled the continuation of full-time study.

I would also like to thank Dr Howard Edwards for his help in the fourth Chapter, on ranking and selection.

Contents

Acknowledgements.....	iii
Table of Contents.....	iv
List of Figures.....	vi
List of Tables.....	vii
Chapter 1 An Overview.....	1
1.1 Introduction.....	1
1.2 General Overview of Cricket Statistics.....	2
1.3 Statistics and Team Selection.....	5
1.4 The Application of Quality Control to Cricket.....	6
1.5 Major Contributions of this Study.....	8
Chapter 2 Analysing Individual Data Characteristics.....	10
2.0 Introduction.....	10
2.1 Literature Review.....	12
2.1.0 Performance Output Measures.....	12
2.1.1 Investigation of Bowling Measures.....	16
2.1.2 Randomness.....	16
2.1.3 Distribution of Performance Measures.....	18
2.2 Data.....	20
2.2.1 Calculation of Batting Measures.....	20
2.2.2 Calculation of Bowling Measures.....	22
2.3 Methods.....	29
2.3.0 Introduction.....	29
2.3.1 Tests for Randomness.....	30
2.3.2 Distribution Fitting.....	34
2.4 Results.....	38
2.4.0 Introduction.....	38
2.4.1 Batting Results.....	38
2.4.2 Bowling Results.....	40

2.5 Discussion.....	41
Chapter 3 Monitoring Player Performance.....	43
3.1 Introduction.....	44
3.2 Previous Studies.....	46
3.2.1 Applying the Methodology of Previous Studies.....	48
3.3 Univariate Quality Control Methodology.....	54
3.3.1 Shewhart Control Charts.....	55
3.3.2 CUSUM.....	63
3.3.3 EWMA.....	66
3.4 Proposed Control Chart Based on Quartiles.....	68
3.5 Comparison of Univariate Control Chart Performance.....	75
3.6 Multivariate Charts.....	82
3.6.1 Hotelling T^2 Chart.....	83
3.6.2 MEWMA.....	90
3.7 Summary.....	92
Chapter 4 Ranking and Selection as Applied to Cricket.....	94
4.1 Ranking the Estimate of an Individual's Ability.....	96
4.2 Homogeneity of Variance.....	97
4.3 Establishing the Probability of Correct Selection.....	101
4.4 New Zealand Statistical XI 94-98.....	103
Chapter 5 Implementation of Statistics in Selection Methodology...	107
Appendix A Data Description.....	110
Appendix B Not Out: The Eleventh Form of Dismissal?.....	112
Appendix C Bowling Results.....	117
Appendix D Batting Results.....	118
Appendix E Ducks and Runs Distribution Theoretical Quartiles.....	120
Appendix F Eligible Players for NZ Statistical XI.....	126
References.....	127

List of Figures

Figure 1. Relative Effectiveness of the Attack Index	27
Figure 2. Relative Effectiveness of the Economy Index	28
Figure 3. Distributional Comparison of Contribution and Score	42
Figure 4. Non-parametric EWMA for B.R. Hartland	50
Figure 5. Non-parametric EWMA for M.J. Horne	50
Figure 6. Distribution-free CUSUM for B.R. Hartland	52
Figure 7. Distribution-free CUSUM for M.J. Horne	52
Figure 8. Control Chart with Warning Lines	57
Figure 9. Control Chart with Warning Lines for Cricket	58
Figure 10. Shewhart Control Chart of M.J. Horne's Transformed Batting Contribution With Zone Run Rules	61
Figure 11. Shewhart Control Chart of B.R. Hartland's Transformed Batting Contribution With Zone Run Rules	62
Figure 12. \bar{C} USUM Control Chart of M.J. Horne's Transformed Batting Contribution . . .	64
Figure 13. CUSUM Control Chart of B.R. Hartland's Transformed Batting Contribution . .	65
Figure 14. EWMA Control Chart of M.J. Horne's Transformed Batting Contribution With Zone Run Rules	67
Figure 15. EWMA Control Chart of B.R. Hartland's Transformed Batting Contribution With Zone Run Rules	67
Figure 16. Fitted Line Plot of Mean Contribution Vs Mean Score	69
Figure 17. Quartile Control Chart	71
Figure 18. Establishing Number of Consecutive Increasing Points for Alarm	73
Figure 19. Histograms of simulated Score data with and without Transformation	78
Figure 20. Histograms of simulated Contribution data with and without Transformation . .	79
Figure 21. Quartile Control Chart for M. J. Horne	80
Figure 22. Quartile Control Chart for B. R. Hartland	81
Figure 23. T2 Control Chart of C.M. Brown for Bowling Indices	85
Figure 24. T2 Control Chart of P.J. Wiseman for Bowling Indices	86
Figure 25. Bivariate Control Chart of C.M. Brown for Bowling Indices	87
Figure 26. Bivariate Control Chart of P.J. Wiseman for Bowling Indices	87
Figure 27. T2 Control Chart of A.C. Barnes	89
Figure 28. MEWMA Control Chart of Bowling Indices for C.M. Brown	90
Figure 29. MEWMA Control Chart of Bowling Indices for P.J. Wiseman.	91
Figure 30. Plot Showing Ranked Nature of Population Batting order.	113

List of Tables

Table 1: Tests for randomness in Batting Performance Measures	38
Table 2: Distribution Fitting for Individual Batting Scores	39
Table 3: Distribution Fitting for Individual Batting Contribution	39
Table 4: Tests for randomness in Batting Performance Measures	40
Table 5: Normality test for Individual Bowling Indices	40
Table 6: Signal Length in Consecutive Points	74
Table 7: Comparative ARL's	76
Table 8: Ratio's Comparing Run Length From In-Control to Out-of-Control.	77
Table 9: Theoretical Quartile Limits for Horne and Hartland	79
Table 10: Comparison of the Number of Points to First Alarm for Horne and Hartland for Different Methods	92
Table 11: P-Values from Test for Homogeneity of Variance for Bowlers	97
Table 12: Probability of Correct Selection for NZ Statistical XI	105

Chapter 1.

An Overview

1.1 Introduction

Cricket is a game of numbers. The very core of the sport is entwined with numerical values that translate ultimately to a match result. These sport statistics are a natural by-product of competitive sport and have been around along as contested sport has existed. Currently sport reporters and commentators bombard observers with a vast array of numerical values designed to describe an individual's performance at a particular skill. These added extras contribute to the entertainment value of professional sport. However, is this information of use to coaches and selectors of cricket teams?

This thesis is aimed primarily at the selectors of top-level cricket teams. An attempt is made to keep the statistics involved as simple as possible so that all levels of selectors may apply this methodology to their teams.

There are several key reasons for measuring and evaluating performance in team sport. Organisational Behaviour Theory proves particularly useful in drawing together sport statistics and selection. According to Greenberg and Baron (1997) to build high performance teams appropriate performance measures are required. Tests and measurements are tools that can be used for evaluation of an individual's performance (Franks, B. & Deutsch, H., 1973).

Having found suitable measures of performance these indicators can be used in the selection process. For a high performance team, the right team members need to be selected (Greenberg et al, 1973). This means combining all available evidence, quantitative and qualitative, to make correct selection decisions.

While few really important decisions are made purely on the basis of objective evidence (Franks et al, 1973), selection decisions cannot be based upon subjective evidence alone. The correct balance needs to be implemented. In order to use sport statistics successfully, a deeper understanding of the numerical values involved is necessary. Firstly the nature of cricket statistics is discussed.

1.2 General Overview of Cricket Statistics

Cricket statistics are meticulously collated ball by ball. The vocabulary of the game continually refers to abstract statistical concepts such as average, aggregate and form – without divulging the secrets of what these mystical values contain. For the sake of simplicity, all values involved are reduced to one dimension. However, this leaves the cricket observer to assume and speculate as to the base values involved. The written media has recently taken to describing bowling performance by listing the number of wickets taken followed by the bowler's average. This form is limited; the basis behind this statement is discussed later in the evolution of the bowling indices.

In recent times, an increasing number of studies have been undertaken to understand the statistical processes at work in the game of cricket. G.H. Wood and W.P. Elderton started the ball rolling in 1945, analysing individual batsmen in an attempt to find a general model that would describe individual scores. This is in accordance with the general trend, where most work to date has revolved around batsmen. This seems like an apparent contradiction, as the first skill taught to junior cricketers is how to bowl, for without bowlers the game cannot be played. However, with advent of one-day cricket and now Cricket Max, both geared towards entertainment, the game is becoming increasingly batsmen orientated. "Batsmen have always received the highest accolades. Most histories of cricket are written around them, with the bowlers regarded merely as a necessary evil." (Nigel Smith, 1994, p.177). The reasoning behind the domination of batsmen in statistical papers may be due to the perceived ease of evaluation.

This leads to the definition of the statistics utilised in analysis of player performance, enabling a better understanding of the statistics involved.

Sport Statistics can be separated into two broad categories; Performance Indicators and Performance Outputs.

A Performance Indicator is a quantitative measure that indicates individual performance in a particular facet of the game. These values are collated during the game in progress. Effectively, the game is dissected into small manageable slices, such that a numerical value can be assigned as a descriptive measure. An example of a Performance Indicator, associated with fielding performance, is Ground Ball Efficiency, defined as the number of times the ball is fielded cleanly divided by the total number of times fielded. These values do not have a direct impact on the match figures.

In contrast a Performance Output is a numerical expression detailing the direct result of participation in an event. For cricket these are summary measures detailed in a score book at the completion of an innings, such as score, wickets taken, overs bowled and so forth. As a consequence these values have a direct impact on the match figures.

It stands to reason that these two categories are related in some manner. However, only performance outputs will be examined in this study, due to the ease of data collection and availability. Investigating the possible relationship between performance indicators and performance outputs will be analysed in future research.

Statistically, assessing the performance of a batsman is relatively simple, as this can be given by a single variable, either runs scored in an innings, aggregate, average, or average contribution.

'Aggregate' refers to the total number of runs scored by the individual over a specified period of time. A player's 'Average' is then calculated by dividing the aggregate by the number of times the individual was dismissed during the specified time. 'Contribution' is the percentage of runs the individual provides the team total in an innings. Each value on its own can effectively describe performance.

Describing performance by bowlers is more complicated. A typical bowling analysis gives four values; Runs, maidens, overs and wickets. 'Runs' corresponds to the number of runs conceded by the individual. The number of maidens bowled, refers to the number of completed overs where no runs are penalised against the bowler (leg byes and byes are not added to a bowlers total). 'Overs' refers to the number of six ball sets a bowler has delivered. Finally, 'wickets' are the number of dismissals credited to the bowler. Alone, these individual factors give little insight into how well a bowler performed. Together, they are more meaningful, but not until compared to a full score card can the value of the performance be evaluated. The use of the bowling average attempts to describe performance in one dimension. This is found by dividing the number of runs conceded by the number of wickets taken. However, no time frame is suggested by this value. Essentially it is assumed that a bowler will concede 3-4 runs per over. Over a long period of time this assumption becomes more valid, but is not suitable for a game by game situation.

The Deliotte Ratings create a one-dimensional measure of performance in Test Cricket. This involves an algorithm that takes into consideration several factors and weightings. This is currently the method of determining the best players in the world. Whilst the formulae involved are extremely thorough, the histories of all players need to be known and equally thorough. An attempt to create a one-dimensional index, using both factor analysis and principal component analysis, failed to provide meaningful results (Bracewell (1), 1998). Intuitively this is obvious as two basic concepts are involved.

Ideally, two dimensions need to be considered, one involving the players attacking ability, the other involving the ability to restrict runs. Kimber (1993) gives a graphical method for comparing bowlers. This utilises two dimensions; the attacking ability (strike rate) and the ability to restrict runs (economy rate). Bracewell (2)(1998) proposed two independent normally distributed indices, based upon strike rate and economy rate, to describe performance. The first index deals with a bowler's ability to take wickets, the second with the ability to restrict runs.

Both indices are evaluated using simple variations of formula that are already used, taken relative to the team performance. The section dealing with assessing bowlers relies heavily on these indices. Having defined the performance outputs to be assessed it is necessary to discuss the relevance in a selection situation.

1.3 Statistics and Team Selection

With the wealth and quality of data available in cricket, it makes sense to utilise this quantitative information in the selection of individuals to maximise the formation of a collective unit (the team). The main assumption underpinning the work in this thesis is that a player's natural ability is expressed by individual performance outputs collated following the completion of a match.

Statistics are not the only factors considered when selecting a team. However, Former New Zealand Coach Glenn Turner (1998) discusses the importance of statistics in choosing players in his book *Lifting the Covers*. In particular the second chapter reveals the emphasis placed on statistics in comparing and selecting individuals. In this instance it is used particularly to justify the non-selection of players, (Andrew Jones and Ken Rutherford) then to defend the selection of Lee Germon.

"Late in 1995 Francis Payne, cricket author and statistician, provided me with statistics which mostly confirmed what we had known before we picked our first test team." (p42).

Glenn Turner (1998), Former New Zealand Coach

Since statistics are used to make and confirm selection decisions it is necessary to attempt to understand the nature of the data being generated by participation in sport. A greater understanding leads directly to better implementation and hopefully a competitive edge, for the selected team.

Former Australian captain, Richie Benaud, remarked on the simplistic nature of selection and the use of statistics stating, "All a captain needs is the confidence that his bowlers are each capable of taking five wickets in an innings, his batsmen are capable of scoring a century and that everyone can field like Viv Richards (Benaud, 1995, p169)." Obviously the captain deals with the player s on the field and is not responsible for those selected to take the field, this lies in the hands of the selectors. The captain must believe that he has been given the best men to compete. It then becomes the job of the selectors to ensure that the best combination of players available takes the field. If statistics are to be used in the selection process they must be meaningful, and secondly they must be used in an appropriate manner. This means that a relevant application of statistical methodology is that of monitoring individual ability.

1.4 The Application Of Quality Control to Cricket

The idea of monitoring performance is as useful to the selector and the player as it is to the arm chair critic. An ideal method for monitoring an individual's performance is with control charts. The control chart is a useful tool in statistical process control. First developed by W.A. Shewhart, the shewhart charts are widely accepted as standard tools for monitoring process of univariate independent and nearly normal measurements (Liu & Tang, 1996). Control charts have found frequent applications in both manufacturing and non-manufacturing settings (Montgomery, 1997). With slight adjustments shewhart charts can be applied to cricket.

Provided the measurements of the individual's performance are reflective of quality, function, or performance then the nature of the 'thing' being measured has no bearing on the general applicability of control charts (Montgomery, 1997).

Montgomery (1997) discloses several reasons for the popularity of control charts. At least 3 draw direct parallels to cricket. Possibly most important is that control charts provide diagnostic information. This can identify flaws in technique, or the tendency for a player to struggle under certain conditions. Also control charts are proven at improving productivity, which translates to pushing a player and not allowing complacency.

In Cricket we are interested in selecting individuals that will maximise team performance and ensure the best chance of victory. Whilst Cricket is a team sport, the nature of the game allows for individual aspects to stand out. Indeed, when we look at the possible selection of an individual, it is the performance outputs of the individual that is of primary concern. Therefore to ensure the right selections are made, it is important the right statistics are used.

Due to the awkward nature of bowling performance outputs, this leads to the evolution of the bowling indices. These two independent, random, standard normal indices are a simple and effective way of allowing bowling performance to be measured from the post match statistics. They are more useful than the current convention used in the written media of quoting the number of wickets taken and the bowling average of an individual.

Utilising the assumption that an individual's worth is expressed via performance outputs, this thesis seeks to describe and understand the underlying statistical processes that shape our impression of player performance in the second chapter. Randomness is tested for and then distributional properties of the data are sought.

Armed with information generated in the second chapter, the third chapter assesses methods for monitoring the estimate of natural ability.

Widely accepted control methods, such as Shewhart control charts, CUSUM, EWMA and multivariate versions of these procedures are implemented and the performance for both batting and bowling is discussed. To accommodate the distribution presented by batting scores a new control chart based on quartiles is also studied.

Further, ranking and selection procedures utilise the estimates of individual ability to select the best individuals and note the probability of correct selection in chapter four.

Chapter Five then details how this information can be drawn together and applied in selecting a side with the assistance of statistics based upon performance outputs.

1.5 Major Contributions of this Study

A number of new and novel approaches are presented in this thesis, these include:

- a) the further development of individual performance measures for the main disciplines of batting and bowling for cricket.
- b) A 2 – Dimensional runs test, utilising the T^2 statistic, with further applicability outside cricket.
- c) Statistical interpretation of assumptions and results specific to cricket namely:
 - Outliers are very important in determining the estimate of ability for an individual.
 - Form is autocorrelation.
 - Zone rules for cricket are needed to detect good/poor performance.
 - Relatively short nominal ARL's to accommodate the restricted number of sampling opportunities presented in a season.

-
- d) A new Control Chart based on quantiles to preserve outlier influences in a non-parametric procedure.
 - e) The recommendation of appropriate tools for monitoring batting, bowling and all-rounder performance and also choosing man of the match.
 - f) a selection procedure for bowlers that discriminates between different types of bowlers using the consistency of their performance measures.
 - g) Following selection, an evaluation of the probability of correct selection of individuals to a team, relative to potential contenders.
-

Chapter 2

Analysing the Characteristics of Individual Cricketer Performance Data

Viz Testing for Randomness Over Time and Fitting Appropriate Distributions in Order to Settle the Assumptions of Statistical Inference

2.0 Introduction

This chapter looks at measures of individual performance in terms of batting and bowling. Having established appropriate performance measures, each attribute is tested for randomness. Then the distributional properties of the sampling distributions, for each discipline, are investigated. The properties of these distributions are used to estimate the player's true 'ability'.

Randomness is associated with performance through the concept of form. From a statistical viewpoint 'form' is the occurrence of trends over time in an individual's performance. Form is similar to consistency, conforming to a regular pattern or style (Oxford Senior Dictionary). At first class level it is assumed that individuals participating have the necessary skills to adequately compete. If this is so the previous performance relative to the team should have no impact on the current performance relative to the team. This assumption may be challenged when the player is too good, but everyone fails at some time or another; or when the player is not good enough, in which case selection may be discontinued if that is the sole domain of selection. This is based on the assumption that when the skill level of competing individuals is approximately equal, then it is expected the performance of an individual is dictated more by chance than skill. However, when the skill level is unmatched between individual's, superior skill will offset chance to a degree, as will inferior skill be more dependent on chance. If performance is random, the job for the selectors is simplified.

Previous work, involving the analysis of both batsman and bowler performance outputs has assumed random performance. This assumption needs to be clarified before further progress can be made. The initial thrust of the thesis is the identification of what constitutes form. Form can be likened to autocorrelation, in that an individual displays patterns or trends in performance over time.

It is expected that two extremes may exist, either form exists, or performance is random. If autocorrelation is present, then form exists. Intuitively performance would be considered random, due to the apparent lack of predictability of such a sport. "Uncertainty plays a large role in sports, and one can argue that the uncertainty associated with sports outcomes is one reason that sports are so popular (Stern, 1997, p19)." It has been shown that baseball is a game of chance (Cook, 1977). An analysis of team tactics as related to the game of baseball and analysis of the annual World Series competition revealed that results were subject more to the laws of chance than the relative calibre of the competing teams. Taking a simplistic view of competitive sports suggests this may also be the case in cricket (Assuming everyone is equally able to compete, and that natural ability will differ, dependent on the pool of talent available). Logistically it would be ideal if performance is random. If this is the case then it is a relatively simple task to select the best individuals, provided that the sampling distributions to which the data belong are known.

In order to fully understand the summary statistics presented, and make effective use of the available information, the statistical distribution for each of the performance outputs needs to be known. Fulfilment of this requirement and that of randomness satisfies the most important assumptions regarding inference and quality control.

An overview of previous research on performance output measures in cricket is presented in the next section.

2.1 Literature Review

Cricket is a fascinating game. The outcome of every ball delivered is recorded. Given the wealth of statistics it comes as a surprise then to note the relative lack of literature based upon the statistics of this summer sport. Present research in this field favours the one-day game and the batting perspective. Very little literature exists on the other major aspects of individual involvement, Bowling, Fielding and Wicketkeeping. Perhaps this is due to the fact that the data required for an effective analysis in these fields are rather specialised and messy and not readily available. As mentioned in Chapter One, Performance Indicators are used to quantify performance and require additional input from each ball bowled, in a form that has no direct impact on the match result. So it is no surprise to note the lack of literature dealing with performance indicators. This thesis avoids performance indicators based upon the assumption that individual performance is expressed via performance outputs.

In this section the statistical literature is discussed in relation to random performance and distribution fitting in cricket. Firstly the performance measures are defined.

2.1.0 Performance Output Measures

Four performance output measures will be defined for batsmen, six for bowlers.

a) Batting Performance for Individual Batsmen

A single value, either individual score, or individual contribution to the team total most easily determines the performance of a Batsman in an innings. Contribution is the percentage of the team total an individual has contributed.

$$Contribution = \frac{Individual\ Score}{Team\ Total} \times 100\%$$

Individual Score is the number of runs credited to an individual during an innings and Team Total is simply the total number of runs amassed by that team whilst batting in that innings.

Over a period of time, a batsman's worth is investigated via their batting average. The traditional batting average is expressed below.

$$\text{Traditional Batting Average} = \frac{\sum \text{Individual Scores}}{\text{dismissals}}$$

However, our interest is with what an individual is expected to score in a given innings, as measured by a batting average.

$$\text{Batting Average} = \frac{\sum \text{Individual Scores}}{\text{innings}}$$

Thus the aggregate score is divided by the total number of times batted to provide the batting average. This performance measure circumvents the debate surrounding the handling of 'not outs' by only considering the average score per innings. This is different from the traditional batting average, shown above, which estimates the runs scored between dismissals, however this method seems redundant due to the time constraints placed upon the game, especially as we are to consider the expected number of runs in an innings. A further discussion is included in Appendix B.

Finally, average contribution can be used as a measure of batting performance as shown below.

$$\text{Average Contribution} = \frac{\sum \text{Contribution}}{\text{innings}}$$

It is defined in a nature similar to that of batting average. The sum of individual contributions is divided by the total number of innings.

b) Bowling Performance for Individual Bowlers

Of the individual disciplines, bowling is perhaps the hardest to evaluate quantitatively. A typical bowling analysis consists of four variables, Runs conceded, Maidens bowled, Overs bowled and Wickets taken. There is no easy way of interpreting these values independently. History plays a large part of how these statistics are perceived as does the game situation. This section briefly reviews the statistical methods for evaluating an individual's bowling performance.

Kimber (1993) proposed a two-dimensional graphical display for comparing bowlers in cricket based on strike rate (SR) and economy rate (ER), taking advantage of the relationship that these two values have with the Bowling Average (AV).

$$SR \times ER = 100AV$$

These values are traditionally calculated as follows: the Economy Rate (ER) is defined as the runs conceded per ball.

$$\text{Economy Rate} = \frac{\text{Total Runs Conceded}}{\text{Total Balls Bowled}}$$

The Strike Rate (SR) is defined as the number of balls bowled per wicket taken.

$$\text{Strike Rate} = \frac{\text{Total Balls Bowled}}{\text{Wickets Taken}}$$

(Kimber, 1993)

However, this relationship does not take into consideration the team situation, and other confounding variables that confront a bowler, such as the state of the game, combined with environmental factors, as these can have an impact on how the specific individual's involved, batsman and bowler, approach each delivery. As a brief example, a batsman is more likely to attack the bowler towards the end of the innings, with wickets remaining in a run chase, than a batsman trying to save the match by remaining not out in a last wicket partnership when a run chase is no longer viable. Strike Rate has an additional problem, in that if a player fails to take a wicket, a value for SR is not returned as the divisor is zero.

Thus SR is not suitable for evaluation on an innings by innings basis. This measure could be calculated using all the match results for a season, but in terms of selection and monitoring a player's performance, it is too late to address an individual's worth at the end of the season. Thus only players who have taken wickets can have strike rate as a performance measure.

Bracewell (3)(1998) detailed a novel way to evaluate individual bowling performance, incorporating SR and ER into two separate indices that considered relative performance to the team. This involved an attempt to form ratio's that took into account an individual's performance in relation to the team performance. The Attack Ratio involved inverting SR for both team and individual so that wickets taken was no longer the denominator.

The ratios are defined as follows:

$$\text{Economy Ratio} = [\text{Opposition Total} / \text{Total Overs} - \text{Runs Conceded} / \text{Overs}]$$

$$\text{Attack Ratio} = [\text{Wickets/Overs} - \text{Total wickets/Total Overs}]$$

(Bracewell, (3) 1998)

However it was found that as the number of overs bowled by an individual increased, the score for both indices tended to zero. This was because as a player bowls more and more overs (approaches 50%) this player is having a huge influence on the team performance. His performance therefore reflects the team performance very closely.

The final evolution of performance measures for bowlers involved multiplying the ratio's by a weighting factor related to time (overs). The problem described earlier was removed in this way.

In addition it was found the Attack index needed to be multiplied by a wicket weighting factor, defined in terms of w , the number of wickets taken in any innings. This index is therefore innings specific whereas the other measures are more general.

The wicket weighting factor in the Attack Index is given by $p(w)$, the probability of taking a certain number of wickets in an innings. Standardisation allows the indices to be compared on similar scales.

The indices are therefore defined as follows:

$$\text{ECONOMY INDEX} = [\text{Economy Ratio} \times \sqrt{\text{Overs}}]$$

$$\text{ATTACK INDEX} = [(\text{Attack Ratio} \times \sqrt{\text{Overs}}) / (1 - p(w))]$$

(Bracewell, (3) 1998)

2.1.1 Investigation of Bowling Measures

Of the statistical analyses performed using cricket data, bowling is an area deficient in research. Only Kimber (1993) and Bracewell (3)(1998) have examined how to measure an individual's bowling performance. Kimber sought to do this via a graphical display based on Strike Rate and Economy Rate, whereas Bracewell tried extending these values relative to the team.

2.1.2 Randomness

Very little research has been done on the aspect of randomness in an individual's performance in cricket. A distantly related team sport, baseball, was found to be essentially random (Cook, 1977). There is anecdotal evidence supporting the claim that the role of an individual within a game is random, generally commenting on the apparent lack of predictability of cricket. Berkman, (1990), Brittenden (1994) and Turner (1998) are just a small selection of cricket observers that subscribe to the unpredictability of cricket view. Hunting through player biographies also reveals that those who play the game express this view.

Danaher (1989) applied a Run's test to 6 English County Cricketers and found that none showed a significant runs pattern at the 5% significance level. The batsmen chosen were of varying batting ability but chosen because they were either top, or close to the top, of their team's batting averages list.

Kumar (1996) suggested that cricket is not by chance. However, this assertion was based solely upon over run rates in one-day cricket. The implications of this are manifested in the troublesome interrupted match rules. If over rates were random, then the simple Average Run Rate (ARR) rule would suffice, as this is based on the assumption that run rate of the batting side does not change during the innings. Instead, the resources available to a team play an important role in determining the outcome of a one-day match. One only needs to look to the Duckworth-Lewis model (1996) to see the effect that time (overs in hand) and wickets in hand have in determining a batting side's capacity for team total. Team strategies also illustrate this point. As a simple illustration of batting capacity, this model accepts the fact that a side is more capable (or daring) of scoring runs when only 2 wickets have been lost, as opposed to being 8 down, with 10 overs remaining. The reasoning behind this is; with the loss of only 2 wickets, presumably the better batsmen are still available, and there are plenty of individual's remaining. Thus batsmen are more able to go after their shots, as the consequences to the team of their dismissal are not as great. Whereas, a batting side with 8 wickets down needs to adopt a more cautious approach, as once a team is dismissed, there is no further chance of adding to the team total.

2.1.3 Distribution of Performance Measures

a) Batting Performance for Individuals

In the late 1970's no satisfactory distribution for individual batting scores had been found, although there had been several unsuccessful attempts (Pollard, 1977).

The first attempts to model the frequencies of individual batting scores began when Elderton (1927) used a Type X Pearson curve, which is exponential, to model the scores of the Yorkshire batsman Tunnnicliffe. Later Elderton (1945) concluded that individual scores were based upon a geometric distribution from the observation of 4 English county cricketers. Wood (1945, p13) also explored the use of a geometric distribution by considering "the massed results of the scores of large groups of batsmen over very many innings." From the evidence gleaned from this study Wood stated there was "good and solid ground for saying that a Geometric Progression applies wholly or mainly to batting scores at cricket (p14)." However, he encountered problems with this series at each end, more so when dealing with the commencement of an individual's innings.

More recently Danaher (1989) tried fitting an exponential distribution to the batting frequencies of 114 English county cricket players batting frequencies over one season and found only 32.5% fitted an exponential distribution. The level of significance used was not detailed.

Reep, Benjamin and Pollard (1971) showed that for some players a negative binomial distribution fitted well. This distribution was also used to model the scores from batting partnerships. Ganesalingam, Kumar and Ganeshanandam (1994) cite the use of the geometric, exponential and negative binomial distributions to model individual scores but once again mention the unsatisfactory attempts at distribution fitting.

Two World-Class players, Geoff Boycott and Ian Botham, were the centre of Burrows and Talbot's (1985) study regarding the exponential distribution. They found an adequate fit to Boycott's 77 innings and the 50 played by Botham. This study also considered the handling of 'Not Out's'. It was found that by adding the mean of the exponential distribution to a not out score an estimate is found for what the individual was likely to score in that particular innings. As exponential random variables have no memory, this is a valid estimate. Furthermore, using this information to establish a player's compensated batting average through a set of iterative equations and solving the first order difference equation, simply resulted in the traditional batting average normally quoted; the total number of runs scored divided by the number of times dismissed. However, the nature of the competitive game is glossed over. This study "excluded limited overs games since innings in these games are necessarily restricted (p46)." To some extent all cricket played has some time restriction, whether 50 overs, 3 days or 5 days, hence the need for declarations, in the pursuit of victory. Further discussion on the effect of time limitations is provided in Appendix B.

Pollard (1977) conceded "that a more elaborate model needs to be developed to describe the distribution of batsman's scores (p129)." This was due to the fact that previous results did not cater for the higher than expected frequencies of failures to score, compared to the theoretical models.

Bracewell (2)(1998) suggested a discrete version of a mixed exponential distribution for score and a relatively new concept in cricket statistics, contribution, based upon 5609 observations of individuals in the top 6 of the batting order from New Zealand domestic first class cricket. This involved separating the occurrence of zero and recalculating the mean to find the parameters of the distribution involving the non-zero values.

2.2.0 Data

The data used in this study refers to the full score cards of New Zealand first-class cricket obtained from The Shell Cricket Almanack's (Payne and Smith). A first-class match, defined by the Imperial Cricket Conference, 19 May 1947, is a match of three or more day's duration between two sides of eleven players officially adjudged first-class (Payne, F., & Smith I., 1996). All matches considered were played between the commencement of the 1993-94 season and the completion of the 1997-98 season. Hence all matches considered had a maximum duration of four days. Only domestic games were considered. Thus the results only apply to matches played in New Zealand conditions, and individual's competing in the Shell Trophy competition.

The Data was entered into Microsoft Excel, relevant sections were then transferred to MINITAB for the statistical analysis, under the categories listed in Appendix A. Whilst some of the columns are redundant for this study, they were entered for possible use in future studies. Appendix A lists a full description for each column. Appendix B discusses the handling of not outs for batting scores.

2.2.1 Calculation of Batting Measures

Batting performance can be evaluated on a univariate level. Generally speaking an average is the basis for any inference regarding an individual's ability.

Bracewell (2)(1998), proposed taking into account the team situation when examining batting performance. Thus batting contribution was defined and briefly analysed. To fit the distributions of both contribution and runs scored, a discrete form of the exponential distribution was applied.

If the probability of not scoring is p_0 , and the likelihood of an individual's score, or level of contribution, is exponentially distributed with mean β , then the mixture probability density function for x is given on the following page.

$$f_X(x) = \begin{cases} p_0 & \text{if } x = 0; \\ (1 - p_0) \times 1/\beta \times e^{-x/\beta} & \text{if } x > 0. \end{cases} \quad (\text{Smith, 1993})$$

The probability of a certain score is given by the area of the corresponding interval of the probability density function. Considering individual scores and contribution, “not scoring” is the failure to score a run.

Analysing scores from 5609 individuals from only the top six of the batting order yielded the following models.

1) The fitted model for individual scores:

$$\bar{f}_X(x) = \begin{cases} 0.075 & \text{if } x = 0; \\ 0.032 e^{-x/28.955} & \text{if } x > 0. \end{cases}$$

2) The fitted model for individual contribution:

$$f_X(x) = \begin{cases} 0.092 & \text{if } x = 0; \\ 0.066 e^{-x/13.686} & \text{if } x > 0. \end{cases}$$

A chi-square goodness-of fit test indicated both models were of significantly good fit at the 5% significance level.

Obviously the ability of batsmen in the top six differs. Thus, the suggested models are contaminated. However, the nature of the distributions give an insight into how the individual performance outputs for batting are distributed.

2.2.2 Calculation of Bowling Indices

In this study we shall use the measure of performance for the bowling indices developed by Bracewell (3)(1998). These indices provide a relative measure appropriate for use on an innings by innings basis. As Bracewell's indices were created by examining first class matches with a maximum duration of three days, they need to be re-evaluated for the data used in this study, because the data extended to four-day matches in many cases. Indices are calculated for all players for each innings played.

a) Wicket Weighting Factor

An important element of the Attack index is the probability of taking w wickets in an innings. The initial study by Bracewell (1998) found that the number of wickets taken by a bowler in an innings could be presented by a model in the following form.

$$w = a - b.p(w)^c$$

Where w represents wickets per bowler per innings, $p(w)$ indicates relative frequency for w , and a, b, c are constants. These constants were estimated, with 95% confidence intervals, as $a = 11.9$ (11.42, 12.38) and $b = 14.3$ (13.43, 15.17) and $c = 0.2$ by Bracewell (3)(1998).

A χ^2 test involving the new data with the above equation and parameters estimates yielded an observed χ^2 value of 32.952. The critical χ^2 value at a 5% level of significance with 11 degrees of freedom was 19.675. As the observed value exceeded the critical value the null hypothesis that the data is from the given distribution must be rejected in favour of the alternative, that there is no fit at the 5% level of significance.

Thus the parameters needed re-estimation. It was likely that the value for c would be close to 0.2, providing an ideal starting point.

Using an iterative approach it was found that c was 0.25. Using this factor to linearize the equation a regression analysis was performed allowing the final estimation for a and b . The adjusted r^2 value of the regression was 99.5, indicating the proposed regression line explained almost all of the variation in the data.

Estimates for both a and b were acquired, and are given below, with 95% confidence intervals, $a = 11.1$ (10.61, 11.59) and $b = 14.2$ (13.19, 15.21).

The value for a in the model must always be greater than ten. If it is equal to ten, then this assumes that the probability of taking all ten in an innings is impossible. Evidence proves this not to be the case: A.E. Moss, in the 1889/90 season took all 10 wickets in an innings for Canterbury against Wellington at Christchurch on debut (Payne & Smith, 1996).

The interval for ' a ' does not contain Bracewell's (3)(1998) estimate. This suggests that the probabilities change slightly given the extra allowable day. However the changes provide only minimal difference. When 5 or more wickets are taken in an innings the probabilities are approximately equal. The most noticeable difference is the probability that no wickets are taken in an innings. The revised estimate is less than the initial probability from the three-day game. This suggests that in an extended match a player is less likely to go without a wicket. This is possibly due to the chance of having a prolonged bowling spell.

However, the interval for ' b ' contains both the original estimate and its confidence interval.

The resultant χ^2 value of 12.855 indicated a suitable fit, as it is less than the critical χ^2 value of 18.307 at the 5% level of significance with 10 degrees of freedom.

This information suggests a new model for wickets (w) as follows.

$$w = 11.1 - 14.2p(w)^{1/4}$$

with the probability of a bowler taking a certain number of wickets in an innings as given by:

$$p(w) = \left(\frac{11.1 - w}{14.2} \right)^4 \quad w = 0, 1, \dots, 10$$

This probability function is used in the attack index to describe a bowler's attacking abilities relative to that of the team.

b) Attack Index

As indicated in 2.1.0 Bracewell's (3)(1998) attack index is

$$\text{ATTACK INDEX} = [(\text{Attack Ratio} \times \sqrt{\text{Overs}}) / (1 - p(w))]$$

where Attack Ratio and $p(w)$ have been defined previously and overs is the number of overs bowled by the individual in the innings.

Standardising the above equation by first subtracting the mean (-0.07883) and then dividing by the standard deviation (0.37461) gives a value comparable to the standardised economy index.

A 95% Confidence interval for the population attack mean suggests that it probably lies between -0.0541 and -0.0912. As zero is not contained within this confidence interval indicating that the value for the sample mean is significantly different from zero. This interval also includes the original parameter estimates for standardisation.

Utilising the fact that

$$\frac{(n-1)s^2}{\sigma^2}$$

is χ^2 with $n-1$ degrees of freedom allows a confidence to be formed interval for the standard deviation. As a result a 95% confidence interval for the standard deviation shows it probably falls between 0.3855 and 0.3858. Due to 1 not falling in this interval, the value given for standard deviation needs to be used to standardise the attack index.

Both these intervals contain the corresponding Bracewell (3)(1998) estimates.

This suggests that the attack index provides similar results for 3 and 4 day matches. Most importantly this index is not sensitive to match duration.

The final formula for the Attack Index is given below.

$$\text{STANDARDISED ATTACK INDEX} = \frac{[(\text{Attack Ratio} \times \sqrt{\text{Overs}}) / (1 - p(w))] + 0.07883}{0.37461}$$

c) Economy Index

As indicated in 2.1.0 Bracewell's (2)(1998) economy index is:

$$\text{ECONOMY INDEX} = [\text{ECONOMY RATIO} \times \sqrt{\text{OVERS}}]$$

Where Economy Ratio has been defined previously and overs is the number of overs bowled by the individual in the innings.

Standardising the above equation by first subtracting the mean (0.3657) and then dividing by the standard deviation (3.2016) gives a value to comparable to the standardised attack index. A 95% Confidence interval for the mean reveals that the population mean probably lies between 0.3851 and 0.3463. As zero is not contained within this confidence interval the value for the mean can not be ignored and must be included in the standardisation. Similarly a 95% confidence interval for the population standard deviation shows that it probably falls between 3.2001 and 3.2035. Due to 1 not falling in this interval, the value given for standard deviation needs to be used to standardise the economy index.

The final formula is given below.

$$\text{STANDARDISED ECONOMY INDEX} = \frac{[\text{ECONOMY RATIO} \times \sqrt{\text{OVERS}}] - 0.3657}{3.2016}$$

For the Strike Rate defined in 2.1.0, a low value is most desirable. This suggests that the bowler has good attacking abilities by taking wickets quickly. However, a problem arises when no wickets are taken, and there is a high chance of this happening in an innings (0.4) (Thus this value is only useful at the end of a season, provided a wicket has been taken). For this reason an alternative strike rate definition is used in this paper, namely the inverse of the conventional strike rate.

Both these indices were found to be approximately normal in distribution. Normality on an individual player level is investigated later in the results section.

The performance of these indices, standardised economy and standardised attack, is best expressed graphically. Both of the following graphs indicate a one-day match type scenario. For this simulation the opposition scored 220 runs from 50 overs and the bowler in question delivered 10 overs. These parameter values are also printed immediately below the graphics that follow.

Firstly the Attack Index is examined given the specified conditions and is shown on the next page. As the 'raw' measure of an individual's ability is the number of wickets taken, a direct comparison between wickets taken and index score is useful.

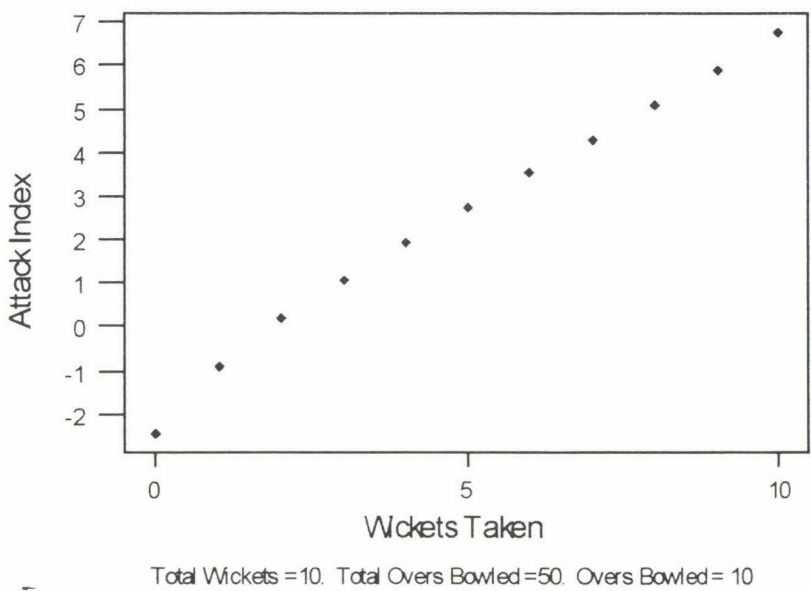
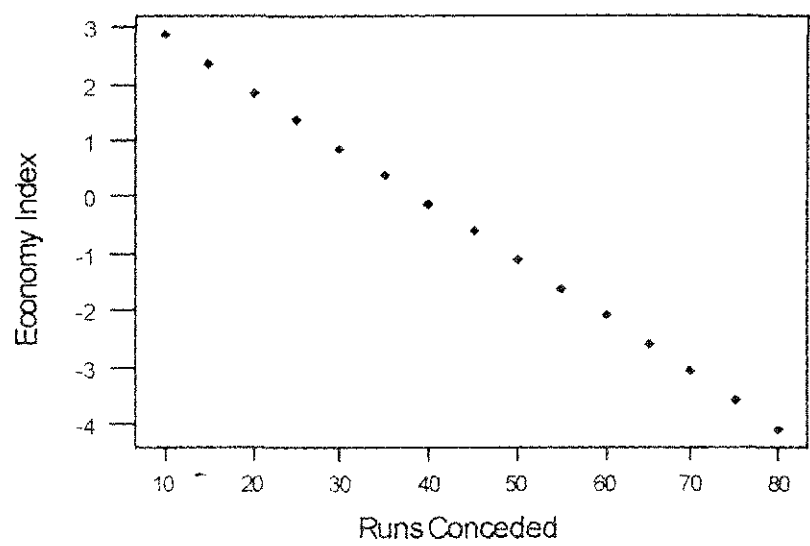


Figure 1. *Relative Effectiveness of the Attack Index*

The first graph, showing the resultant value of the Attack index for various quantities of wickets taken in an innings, shows the immediate value of the attack index. The number of wickets taken correlates positively with the index score, as expected. It can be seen that in the given circumstances, taking 5 wickets in an innings corresponds to an index score of approximately 3. When examined in context the value of the indices is strengthened. Five wickets were taken in an innings only twice in the 97/98 Shell Cup season, which included 34 matches (Blake (Central Districts vs Northern Districts) and Maxwell (Canterbury vs Auckland)). Five was also the maximum number of wickets taken by an individual in the 97/98 Shell Cup competition. This puts the y-axis parameters in perspective. Attack indices of above 3 are relatively rare.

Below the Economy index is evaluated similarly.



Opposition Total = 220. Total Overs Bowled = 50. Overs Bowled = 10

Figure 2. *Relative Effectiveness of the Economy Index*

The y-axis depicts the relative score for the individual's performance given differing values of runs conceded. Obviously the fewer runs conceded the better, and thus this corresponds to higher scores for the economy index. The above graph is much simpler to interpret. In this case there exists a negative correlation between runs conceded and Index score, as expected.

Essentially, given an opposition total of 220 (run rate of 4.4), an individual will concede between 10 (1 run per over) and 70 (7 runs per over) almost all of the time. That is an individual is most unlikely to concede more than 70 runs in a 10 over spell.

The preceding graphs depict the value of the indices obtained from the varying number of wickets taken, and runs conceded. Having shown that the indices perform to expectation, they can be used effectively as a performance measure for bowling outputs.

2.3 Methods

2.3.0 Introduction

“The objective of statistical inference is to draw conclusions or make decisions about a population based on a sample selected from the population (Montgomery, 1997, p78).” Considering cricket, each time an individual participates in a match, a sample of their true ability is revealed. Can a series of scores be used to make inferences about an individual’s ability? If the probability distribution of a population from which the sample is gathered is known, then the probability distribution of the various statistics computed from the sample data can be determined (Montgomery, 1997). More importantly, it can be established what a player is expected to score and thus their progress can be monitored, which is especially relevant for team selection.

A population is a set of measurements that can be described by a set of numerical measures called parameters (Ott, Mendenhall, 1985). In most applications of statistics the parameters are not known but inferences about them are made using information contained in a sample.

For time series analysis it is assumed that for each time point t , Z_t is a random variable. Thus the behaviour of Z_t will be determined by a probability distribution (Cryer). In this instance time t , refers to each innings and Z_t refers to a performance output.

Previous studies have assumed that the data are independent, in that for each individual bowler the previous match result does not have a direct impact on the following match result. At first class level this is a safe assumption as it is presumed that players who reach this level have developed the necessary mental skills.

As the level gets higher, this assumption would hold more strongly. If the data were dependent it would make the job of prediction and thus selection a very simple job. "Cricket, fortunately, is less predictable than that (Berkmann, 1990)." This section attempts to find if this assumption of independence, or randomness holds.

2.3.1 Tests for Randomness

a) Runs Test

A popular technique for testing the randomness of observed data is the runs test. This is a non-parametric method based on the theory of runs, where a run is a succession of identical values which is immediately neighboured by different values (Freund, 1992). The total number of runs appearing in an arrangement is often a good indication of a possible lack of randomness.

The data involved must be binary. In most instances, the data is converted to such a state by reclassifying the data using ones for values greater than the median and zeros for values less than the median.

The runs test is based on three values n_1 , n_2 and u . They are as follows:

n_1 = the number of data values below median

n_2 = the number of data values above median

u = number of runs.

The null hypothesis of randomness is rejected at the α level of significance if

$$u \leq 'u_{\alpha/2} \quad \text{or} \quad u \geq u_{\alpha/2}$$

Where $'u_{\alpha/2}$ is the maximum value for which the probability of getting a value less than or equal to it does not exceed $\alpha/2$ and $u_{\alpha/2}$ is the minimum value for which the probability of getting a value greater than or equal to it does not exceed $\alpha/2$ when the null hypothesis is true (Freund, 1992).

Values for $u_{\alpha/2}$, and $u_{\alpha/2}$ are to be found in table XI of Freund (1992). If n_1, n_2 are both greater than 15 then u is approximately normally distributed with:

$$\text{Mean} = \frac{2n_1n_2}{n_1 + n_2} + 1$$

$$\text{Variance} = \frac{2n_1n_2(2n_1n_2 - n_1 - n_2)}{(n_1 + n_2)^2(n_1 + n_2 + 1)}$$

Consider the scores of Adam C. Parore in the data set: 6, 8, 87, 4, 6, 40, 26, 133, 0, 84, 14, 91, 0, 63, 111, 87.

This series of scores has a median of 33. Thus the re-coded binary data is as follows: 0, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 1.

In the given example $n_1 = n_2 = 8$, and $u = 12$. Table XI of Freund indicates $u_{0.025} = 4$ and $u_{0.025} = 14$. Thus the null hypothesis of randomness is rejected if $u \geq 14$ or $u \leq 4$ at the 5% level of significance. As the obtained value of u does not violate the limits, there is insufficient evidence to reject the hypothesis of randomness

MINITAB performs this calculation, requesting only the median and the column in which the data is stored.

In this study we need to test for randomness using two measures of performance simultaneously, that is, Attack and Economy Indices for bowling in order to determine whether performance is consistent (not form dependent).

This means we need to extend the runs test to a 2-dimensional test. Consider a 2-dimensional graph showing the performance of a player for a series of innings using two appropriate indices.

Considering the example given by evaluating a bowler's performance where 2 independent standardised normal populations are presented in a two-dimensional array. Intuitively if a player is performing consistently relative to the team then the points will tend to be clustered close together. Hence the distance between neighbouring points is minimal. As a consequence of this, by considering the distance from the origin of immediate neighbouring performances it is possible to test for randomness in performance. When standard normal populations are used, the squared distance from the origin is examined. Essentially a multidimensional problem is collapsed to one dimension. As the populations being dealt with are independent standard normal it is the T^2 statistic that is being dealt with. The runs test can then be applied to the distances by complementing either the mean or median of the independent dimensional distances. The median is suggested, as this is not prone to influence by outliers.

b) Autocorrelation

Autocorrelation usually presents itself as correlation between consecutive points in time series. However an autocorrelation function measures the correlation between points lagged by 1,2,...,k periods. In our analysis the time series points correspond to the performance of a given player in a series of innings (one period = 1 innings).

At this stage the interest is not the nature of any patterns in the series, but merely the existence of any patterns in individual performance. If the majority of individuals do not display any performance autocorrelation, this will be indicative of randomness in individual performance.

The Autocorrelation function of MINITAB generates a detailed plot of the autocorrelation function, including the 95% confidence limits for the correlations at lag k. If these confidence limits are crossed then the assumption that there is nil autocorrelation is abandoned. Essentially this means testing whether all autocorrelations for all lags up to lag k equal zero (MINITAB, 1996).

The default lag ($n/4$) was used for k , where n is the number of observations in the series. The Ljung-Box Q statistics acts as safeguard against the explosion of the probability of Type I errors by testing the null hypothesis that the autocorrelations for all lags up to k equal zero (MINITAB, 1996).

If the Majority of time series do exhibit autocorrelation, then the job of the selector is much harder. No longer is the average estimate of an individual's ability sufficient. From the historical data, performance predictions need to be made using prior performance. In terms of recording results from this analysis, when the 95% confidence limits have been crossed, the player in question is noted as displaying significant autocorrelation, thus failing the assumption of randomness.

2.3.2 Distribution Fitting

In the next section we determine which of the following distributions best fit the performance output data for an individual player.

Three distributions are investigated for the batting data, Exponential, Negative Binomial and Geometric. As discussed in 2.1.0, these distributions have been used to model individual player performance by other authors. Below the properties of each distribution is listed. In these formulae p denotes the probability of a success and $(1-p)$ denotes the probability of a failure for independent Bernoulli trials.

Exponential

<i>Pdf</i>	$f(x \beta) = (1/\beta).e^{-x/\beta}, \quad 0 \leq x \leq \infty, \quad \beta > 0$
<i>Mean and variance</i>	$E(X) = \beta, \quad \text{VAR}(X) = \beta^2$

Negative Binomial

<i>Pmf</i>	$P(X=x r,p) = p^r(1-p)^x; \quad x = 0,1,2,\dots; 0 \leq p \leq 1$
<i>Mean and variance</i>	$E(X) = \binom{r+x-1}{x} r(1-p)/p, \quad \text{VAR}(X) = r(1-p)/p^2$
For $r = 1$	$E(X) = (1-p)/p, \quad \text{VAR}(X) = (1-p)/p^2$

Geometric

<i>Pmf</i>	$P(X=x p) = p(1-p)^{x-1}; \quad x = 1,2,\dots; \quad 0 \leq p \leq 1$
<i>Mean and variance</i>	$E(X) = 1/p, \quad \text{VAR}(X) = (1-p)/p^2$

(Casella, Berger, 1990)

These distribution are waiting time distributions. These are suitable as our interest is with the number of runs scored till completion of the innings. The geometric distribution is the simplest of the waiting time distributions, and is also a special case of the negative binomial distribution when r is set at 1 (Casella, Berger, 1990). Note that for the formulae above X denotes the number of failures before the r th success for the negative binomial, while X denotes the trial corresponding to the first success.

The negative binomial, and thus the geometric, are discrete versions of the exponential function. The exponential distribution is a continuous distribution whereas the data here is discrete.

Initially it may seem redundant including both the geometric and negative binomial with r set at 1. However, there is a key difference in the estimator used for the sample mean, as seen in the above table. Previous research suggests that these are the most likely distributions to model individual batting scores. Hence the inclusion of the negative binomial with $r = 1$.

In this study parameter estimation is done using the method of moments. The method of moments is one of the oldest methods for parameter estimation. This method consists of equating the first few moments of a population to the comparable moments of a sample, obtaining the required number of equations needed to solve for the unknown parameters of the population (Freund, 1992).

Given a population has r parameters, the method of moments consists of solving the system of equations

$$m_k' = \mu_k' \quad k = 1, 2, \dots, r \text{ for the } r \text{ parameters}$$

$$m_1' = \frac{1}{n} \sum_{i=1}^n x_i^k$$

All three distributions being dealt with here require only one parameter to be estimated (Negative Binomial set $r=1$).

Thus $m_1' = \mu_1'$ is used.

Exponential

The mean is given by $E(X) = \beta = \mu_1'$, and the expected value from the sample is the mean $= m_1'$.

Therefore setting the method of moments estimator for the exponential parameter is simply:

$$\hat{\beta} = \bar{x}$$

Negative Binomial

$E(X) = (1-p)/p = \mu_1'$ and $\bar{X} = m_1'$. Therefore setting:

$$\bar{x} = \frac{1-p}{p}$$

Rearranging give the method of moments estimator for p for the Negative Binomial Distribution:

$$\hat{p} = \frac{1}{\bar{x} + 1}$$

Geometric.

$E(X) = 1/p = \mu_1'$,

Setting $m_1' = \mu_1'$ provides: $\bar{x} = \frac{1}{p}$

Subsequent rearrangement yields a method of moment estimator as follows:

$$\hat{p} = \frac{1}{\bar{x}}$$

Using the estimates for the parameters of the given probability distributions, the data can be modelled and the fit evaluated using the chi-square Goodness-of-Fit test.

Fitting a Mixed Distribution

As previously discussed in 2.1.2 a mixed model may be a better fit, due to the higher than expected number of zeros (Smith, 1993) (Bracewell, (2) 1998). As a result a separate component needs to be built into the probability model to cater for the number of zeros. The second component of the 'ducks and runs' distribution deals with the non-zero portion. Let p_0 be the probability of a zero score. In order to fit a mixed distribution it is necessary to multiply the probability model by $(1-p_0)$ so that the area under the probability model is equal to one, that is $(1-p_0)p_x(x)$ for $x > 0$. For a geometric distribution, the sum of the probabilities for all possible scores, shown below, clearly converges to one as n approaches infinity:

$$p_0 + (1-p_0) \sum_{i=1}^n p(1-p)^{i-1}$$

For calculation of the parameters of a mixed distribution (p_0 and β or p), the fraction of data set at 0 is separated out and the mean recalculated. The new sample mean is then used as the parameter estimate for the probability distribution (β or p).

The probability mass function of an individual batsman's contribution can be presented in the following form:

$$\text{Pmf} \quad P(X=x|p) = p_c(1-p_c)^x; \quad 0 \leq x \leq 100; \quad 0 \leq p_c \leq 1$$

Where p_c represents the reciprocal of the mean contribution and x corresponds to the random variable for percentage of the team total. It then follows that the probability mass function of individual batsmen scores can be represented as follows: –

$$\text{Pmf} \quad P(X=x|p_c, p_s) = \begin{cases} p_c & \text{if } x = 0; \quad 0 \leq p_c \leq 1 \\ (1-p_c)p_s(1-p_s)^{x-1}; & x = 1, 2, \dots; \quad 0 \leq p_s \leq 1 \end{cases}$$

Once again, p_c represents the reciprocal of the mean contribution and x corresponds to the random variable for individual total.

Similarly, p_s is the mean score inverted.

Normality Test

Bracewell (3) (1998) hypothesized the bowling indices for individuals are normally distributed. To test this hypothesis a normality test needs to be performed. The normality test for the bowling indices involved the generation of a normal probability plot. The probability for the x -values (index) is calculated then plotted against a standard normal probability score. A least-squares line is fitted to the points. This forms an estimate for the cumulative distribution function from which the data for the population is drawn. The Anderson-Darling test for normality is used, which is an ECDF (empirical cumulative distribution function) based test.

2.4 Results

2.4.0 Introduction

To enable a chi-square goodness-of-fit test to be performed on the batting outputs, the data needed to be broken into manageable segments, displayed as follows:

Score	0	1-10	11-20	21-30	31-40	41-50	51-100	100+
Contribution	0	1- 5	6-10	11-15	16-20	21-25	26-30	31+

All tests were performed at a 5% significance level.

2.4.1 Batting Results

Considering only individuals who batted in 20 or more innings yielded 66 individuals for analysis. A brief summary of the results follows. Appendix D contains the full results.

		Pass	Fail
Autocorrelation	Score	62/66	4/66
	Contribution	62/66	4/66
Runs Test	Score	64/66	2/66
	Contribution	64/66	2/66

Table 1: Tests for randomness in Batting Performance Measures

Score	Pass	Fail	Obtained χ^2	Critical χ^2
Exponential	49/66	17/66	739.03	581.51
Geometric	52/66	14/66	705.04	581.51
Negative Binomial	53/66	13/66	711.85	581.51
Mixed Exponential	65/66	1/66	336.47	512.06
Mixed Geometric	65/66	1/66	335.04	512.06
Mixed Negative Binomial	65/66	1/66	397.14	512.06

Table 2: *Distribution Fitting for Individual Batting Scores*

The obtained χ^2 and critical values shown in tables 2 and 3 refer to the fit of the model over the entire population. That is the χ^2 values for all individuals are summed and compared to the critical χ^2 value. For the standard distributions this was 527 degrees of freedom ($66 \times 8 - 1$) and 461 degrees of freedom for the mixed distributions ($66 \times 7 - 1$). This helped confirm the best model.

Contribution	Pass	Fail	Obtained χ^2	Critical χ^2
Exponential	62/66	4/66	482.38	581.51
Geometric	56/66	10/66	645.29	581.51
Negative Binomial	62/66	4/66	449.44	581.51
Mixed Exponential	57/66	9/66	470.55	512.06

Table 3: *Distribution Fitting for Individual Batting Contribution*

2.4.2 Bowling

Similarly only individuals who bowled in 20 or more innings were considered, providing 35 individuals for examination. A brief overview of the results follows. Appendix C gives the results in full.

		Pass	Fail
Autocorrelation	Economy	34/35	1/35
	Attack	32/35	3/35
Runs Test	Economy	32/35	3/35
	Attack	34/35	1/35
	Bivariate	32/35	3/35

Table 4: *Tests for randomness in Batting Performance Measures*

		Pass	Fail
Normality	Economy	31/35	4/35
	Attack	33/35	2/35

Table 5: *Normality test for Individual Bowling Indices*

2.5 Discussion

The evidence provided from the analyses performed in this chapter clearly suggests that individual performance in the primary disciplines of first class cricket in New Zealand is random.

Considering first the case of batting. Only 2 individuals from the sample of 66 failed the runs test. These two individuals, Mark Haslam and Shayne O'Connor, failed the test for both contribution and score. As both are primarily selected as bowlers (Haslam SLA and O'Connor LFM) and have low medians it could be argued that the basis behind the non-random behaviour is that their skill level with the bat is not sufficient. Four individuals failed the test for Autocorrelation. Due to the low numbers violating the assumption of randomness, expressed through the runs test and the test for autocorrelation, it is considered sufficient evidence to claim that batting is random in New Zealand first class cricket.

A similar situation applies to the bowling results. From the sample of 35, at most 3 failed the runs-test, or the test for autocorrelation. Once more the majority exhibit random behaviour and this is taken as sufficient evidence for stating that bowling performance outputs are random in New Zealand first class cricket.

According to the analyses performed individual batting scores are best modelled by a mixed geometric distribution, mixed in two parts, the zero portion and non-zero portion. Batting contribution is best modelled by the negative binomial distribution (with r set equal to 1). The zero component of this distribution represents the zero portion of the score distribution.

It is important to note that the negative binomial distribution, and the parameter from the contribution distribution, model the occurrence of zero amongst individual scores. This is an interesting phenomenon.

Obviously either score or contribution has to be mixed as both share the same number of zeros; unless the team continually scores exactly 100 in each innings, in which case the two distributions will be equivalent. That is the score is effectively the percentage contribution, as score is continually divided by 100 runs (the team total). Considering the case of scoring 0, the probability of this occurring is the same for both contribution and score. This is because for contribution $0/y = 0$, where y is the team score. Due to the sample mean representing the shape parameter, there is a difference between the shapes for the score and contribution distributions.

This is shown in the graph on the following page detailing the probability mass function for differing values of score and contribution. The means for contribution and score are 14% and 29 respectively. These were the population values of top 6 batsmen obtained from Bracewell (2) (1998).

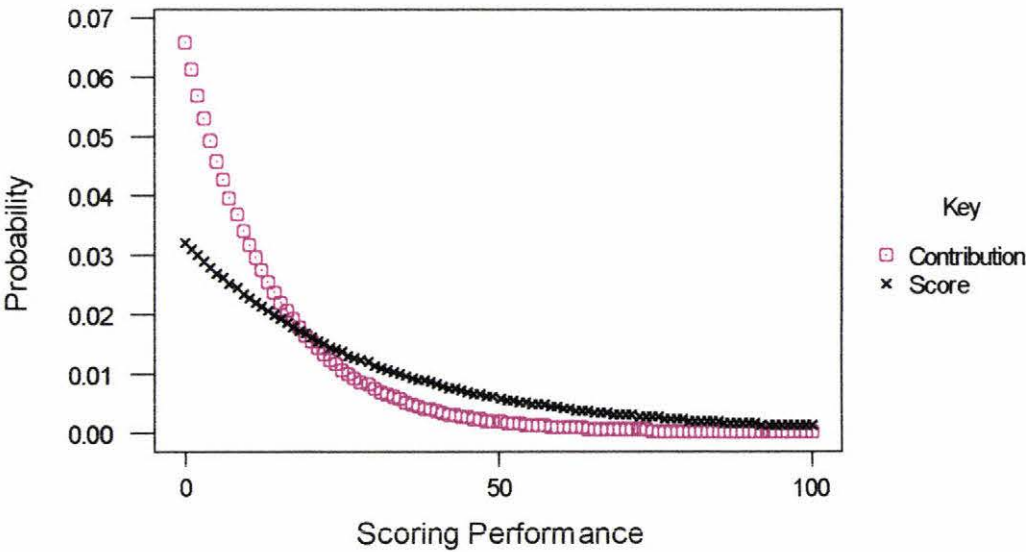


Figure 3. *Distributional Comparison of Contribution and Score*

The distribution for scores confirms the high likelihood of being dismissed early. The fact that the distributions involved are memory-less, as discussed in Appendix B, is also of interest.

This harks back to the old adage; *it only takes one ball*, referring to the fact that only one ball is needed to dismiss a batsman, no matter what score the individual is on.

The results from the normality test showed that Bracewell's (3) (1998) initial hypothesis of normality for the Bowling Indices is correct as an overwhelming majority exhibited this property (33/35 for Attack and 31/35 for Economy).

The above results confirm the beliefs discussed in the literature review. Having proved that performance outputs are random, and gained knowledge of the distributions adhered to by the data, this enables sound statistical methodology to be applied to the performance outputs. A natural extension of this knowledge is to monitor individual ability through statistical process control. This is approached in the next chapter.

Chapter 3

Monitoring Player Performance

3.1 Introduction

The implementation of quality control procedures is useful in that it allows the identification of the most statistically successful players, which is of obvious use in ranking. Furthermore, it enables the monitoring of an individual's performance, identifying any change in the estimate of a player's natural ability. This has direct relevance in revealing the impact of certain coaching and training regimes. Thus, this chapter provides the most extensive section of the thesis, looking at varying control chart methods. Unlike earlier work in this area, parametric tests are used whenever possible. This is vital because non-parametric methods avoid the influence of outliers. An attempt is made to include these influential points due to their direct bearing on our estimate of the individual's ability and their importance. Outliers in this case give an insight into the capability of an individual. In addition data per innings is analysed individually rather than creating artificial subgroups.

A unique situation arises with the study of sports data in relation to quality control. For any given standard it is preferable to be better than the average. It is therefore desirable to have an 'out-of-control' situation on the higher side of the given standard. It is important to note what the charts are being used for in this chapter. Control charts require some target value (or centre line). The quality control tests described in this chapter determine how good an estimate of an individual's ability this target represents.

Before standardising the data it is important to note the estimated parameter values (mean and standard deviation), as these are the estimates of the player's latent ability.

Performance measures are initially standardised with respect to the population, that is subtracting the population mean and then dividing by the population standard deviation, giving estimates of the individuals ability relative to those competing in the same competition. In order for charts to be based upon the standard normal distribution, these indices are standardised using personal means and standard deviations. When the data is standardised and tested with the quality control tools, the test is for how reliable our estimate of the individual's ability is.

Chapter Two sought to prove the fundamental assumptions involved with statistical process control, namely independence and normality, in the context of performance evaluation in cricket. Following on from the findings of chapter two, it is relatively easy to apply control charts to the bowling indices and contribution as the assumptions of normality, and independence are upheld (normality for contribution is achieved by transformation). Thus, the application of conventional parametric quality control methods is assessed in this chapter.

Initially, possible techniques for the situation presented by individual scores are reviewed and a selection applied to real data. Standard procedures for use with normal data are then discussed. For the univariate case, three methods will be examined, namely the shewhart control chart for individual observations with run rules, CUSUM and EWMA. Then a new type of non-parametric control chart based on quartiles is proposed to deal with the mixed distribution of 'ducks and runs' presented by individual batting scores for an innings. The theoretical quartiles of this distribution are used, maintaining the integrity of the distribution. In the entire scheme of things we are interested in a control method that picks up a change in performance within a season. Thus the control chart must pick up changes rapidly. In designing charts and rules consideration must be given to the number of 'samples' per season. The basis behind the need for a short nominal ARL is determined by the structure of the Shell Trophy competition and also the relative lack of sampling opportunities supplied by cricket in general.

Effectively, if five rounds are played an individual has the chance of participating in a maximum of 10 innings. Hence short nominal ARL's are implemented.

Using short ARL's the performance of various quality control procedures is compared with the new quartile chart.

Multivariate procedures are then examined at the end of the chapter for combined bowling and all-round efforts. To adequately explain all-round performance in a match the maximum and minimum values are incorporated into the formation of Hotelling's T^2 control statistic.

3.2 Previous Studies

Very little work has been conducted in the area of quality control in relation to cricket. However, techniques that may be of use to the sporting context are discussed. Firstly techniques for monitoring individual scores are considered. The initial discussion concerns non-parametric methods.

Most Non-parametric Quality Control Techniques are based on signs and ranks. Bakir, S.T. Reynolds M.R. (1979) proposed a method that sought to quickly detect any significant change in mean process in non-normal populations. This method was based around the Wilcoxon signed rank statistic, where the ranking is within groups. A CUSUM type scheme was then applied to the Wilcoxon statistics. It is mentioned in this study that within group ranking is useful when groups of observations are taken at each time point. If single time points are taken then the observations must be divided artificially into groups. For the example presented by first class cricket a case can be put forward for a subgroup of size 2. This is because in a first class match conceivably there can be two observations at each time point (first and second innings). However there is no guarantee that an individual will bat (or bowl) in a game, as a direct consequence of how well the individual's team plays. Whilst scores are independent, random variables (Chapter 2), it makes no sense to artificially group different innings outside of a single match.

As cricket is played under varying conditions against varying opposition it is not suitable to artificially create subgroups. Information relating to an individual's tendency to struggle under differing conditions is potentially lost. Also, the special nature of an outlying score can be forfeited. It is well known that an outlier can influence our estimate of an individual's ability. In this case, the effect will be to inflate the mean. However, we want to retain that influence as it indicates the individual is more capable than suggested by the bulk of the data. By the reduction to ranks the nature of very large scores is removed. Our estimate of an individual's performance indicates what a person is capable of scoring. The presence of an outlier can reveal that our estimate is possibly wrong and the player in question is capable of much more. Outliers are usually regarded as aberrations or errors. But in cricket outliers must be regarded in a totally different light. High scores are valuable observations for performance measures. Another problem arising with the use of ranks dwells with the presence of ties. This makes the use of ranks inexact (Rossini, 1997). The probability of this occurring is quite high due to the likelihood of an individual not scoring.

Hackl and Ledolter (1992) proposed a non-parametric technique utilising sequential ranking. This method involved using the sequential ranks of observations in association with an EWMA control chart. The method is outlier resistant, as all rank based charts are. Hence, the importance of the extreme score described previously is ignored.

After their early attempt with the Wilcoxon signed rank statistic Reynolds and Bakir joined Amin (1995) in investigating a method based on the sign statistic gathered from within artificially created group. As discussed previously it is inappropriate to deal with subgroups in first class cricket as they are not always present.

McGilchrist and Woodyer (1975) applied a Distribution-free CUSUM procedure. However, this method also reduced scores to being above or below median thus ignored the impact of high scores.

Park and Reynolds (1987) implemented a location parameter obtained from linear placement statistics. Once again these values were obtained using within group rankings. Problems arise, as with scores we are interested in the influence of extreme outliers (very high scores) and their effect on the estimate of an individual's ability – thus we want to maintain the original values, if possible. The obtained statistics are applied to standard versions of Shewhart and CUSUM charts. The linear placement statistics are reduced to ranked data.

Non-parametric methods are not assumption free. The assumption of independence is critical (Rossini, 1997). When parametric assumptions are met, non-parametric methods are less powerful. However, when these assumptions fail, non-parametric methods can be more powerful. For this reason a new non-parametric method is proposed for handling the mixed distribution presented by individual batting scores for an innings. This method differs from the above non-parametric methods in that it can easily be adapted to a parametric method, considers individual observations, and doesn't disregard the influence of outliers. As a result this method should be more powerful and appropriate.

3.2.1 Applying the Methodology of Previous Studies

Only Bracewell (3)(1998) has applied general quality control methodology to individual performance data. A difficult situation arises if we wish to monitor performance via individual score.

As outlined in the literature review of this chapter most non-parametric Quality Control techniques are unsuitable due to the fact that they are resistant to outliers. As mentioned earlier, sport presents us with a novel application of quality control procedures. In this situation outliers are important as they shape the estimate of the individual's natural ability and give an insight into what a specific player is capable of.

In addition previous work has considered subgroups of data, where natural subgroups do not occur the authors recommend the creation of artificial groupings of innings results for an individual player. In this thesis it is argued that this approach is also not appropriate because:

- Considering the match context, subgroups of size 2 would be reasonable except that a player may bat in only one innings.
- Cricket is played in variable conditions, playing surfaces or environmental conditions may vary and as a result artificial subgroups are meaningless.
- Outliers are lost

In this section three of the methods outlined in the literature review will be applied to the batting scores of Hartland and Horne and compared to the Quartile chart later.

a) **Non-Parametric EWMA** (Hackl, Ledolter, 1992)

This control chart is based on an EWMA of sequential ranks. Where the sequential rank, $R_t^{(g)}$, is an observation's rank amongst the most recent g observations. This chart performs well with slowly trending process levels. Once again a short ARL is required and this is obtained approximately from Figure 1 and Table 1 (Hackl et al, 1992). An immediate problem arises with the selection of g . As we require a short nominal ARL, we also require a relatively small g , which can lead to correlations among successive ranks. The control statistic is defined as follows:

$$T_t = (1 - \lambda)T_{t-1} + \lambda R_t^{(g)}, \quad t = 1, 2, \dots,$$

The initial value, T_0 , is set at zero. Three parameters are required, obtained from tabulated values where the group size is taken at the smallest available value.

Thus the parameters are set as follows to give a nominal ARL ≈ 19.8 , $\lambda=0.25$, $g=4$, $h=0.2980$.

The resultant statistics, T_t , can be examined via normal chart form, where an alarm for an out-of-control situation is signalled if $|T_t| > h$. As this is the most recent non-parametric technique it is also applied to the simulation data in section 3.4.

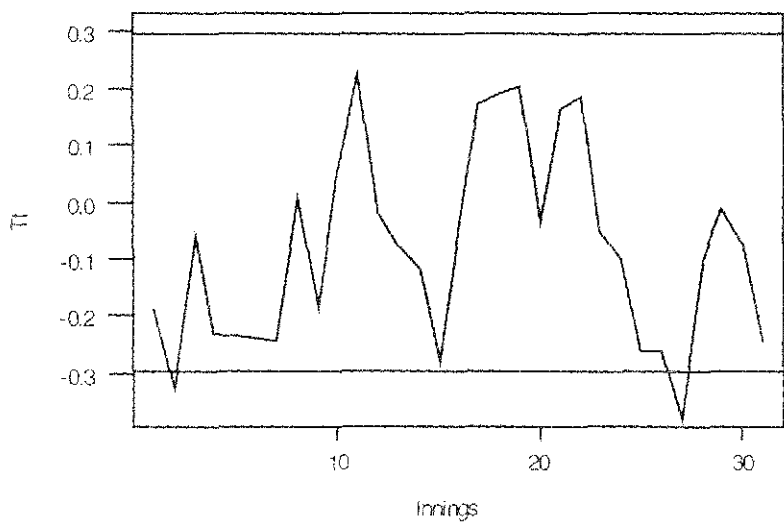


Figure 4. *Non-parametric EWMA for B.R. Hartland*

Two Alarms are signalled for Hartland, one very early on and the other towards the end of the series. Both are associated with signals for inferior form.

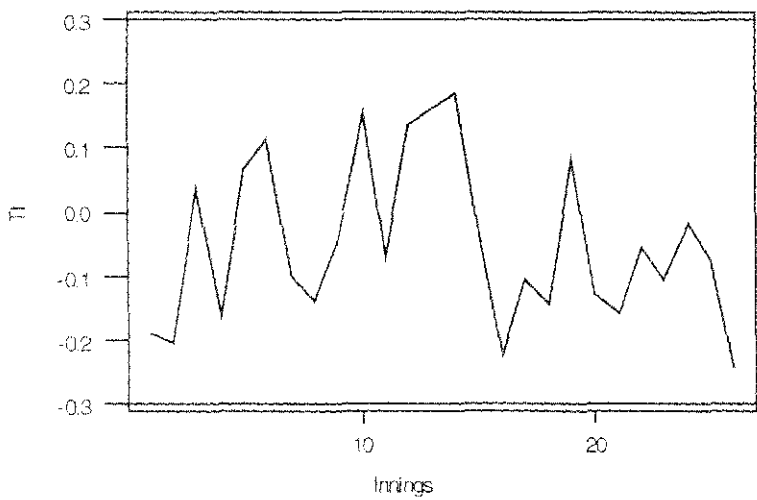


Figure 5. *Non-parametric EWMA for M.J. Horne*

No alarms are signalled for Horne, suggesting that he is consistently playing to the level his natural ability implies.

b) Non-Parametric CUSUM (McGilchrist, Woodyer, 1975)

To allow detection of changes in the extreme distribution posed in hydrology

McGilchrist et al developed a distribution free CUSUM. Using an even number of observations the control statistic V_i is defined as follows:

$$V_i = \sum_{j=1}^i q(X_j - k)$$

Where $q(x) = 1, x \geq 0,$
 $-1, x < 0,$

In order to make $V_n = 0$, k is the sample median.

The series V_0, V_1, \dots, V_n is then a one-dimensional random walk. When the observations in question are independent then all paths from $(0,0)$ to $(n,0)$ are equally likely. Two methods are available to evaluate if the process is in control. The first sets control limits, the second considers the number of returns to zero and is effectively a test for randomness and will not be pursued further.

The control limits for this control chart are evaluated from table 16 in Conover (1971). When these control limits are crossed, this is evidence of an out-of-control process. To have a comparable nominal run length for the non-parametric EWMA the level of significance is set at 6% (ARL=16.7). For a 6% significance level the critical level for V_i is approximately $1.882 \times \sqrt{m}$ where m is half the sample size of the individual in question. Thus for Horne, $m=12$, the critical value becomes 6.52 at the 6% level of significance. For Hartland, $m=16$, the corresponding value is 7.528.

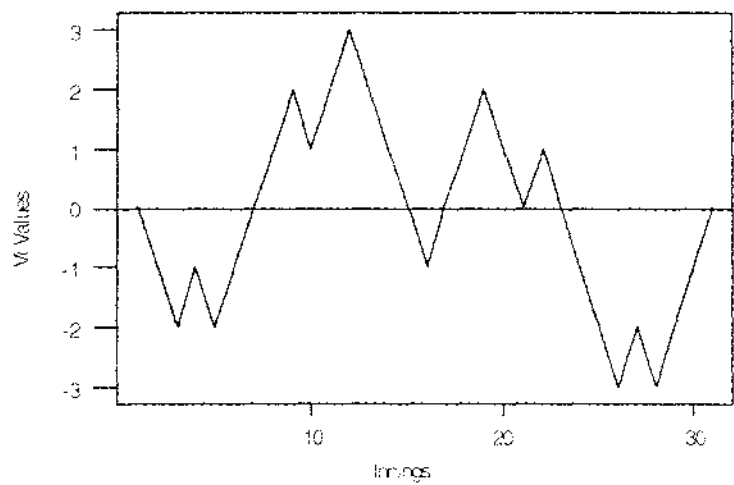


Figure 6. *Distribution-free CUSUM for B.R. Hartland*

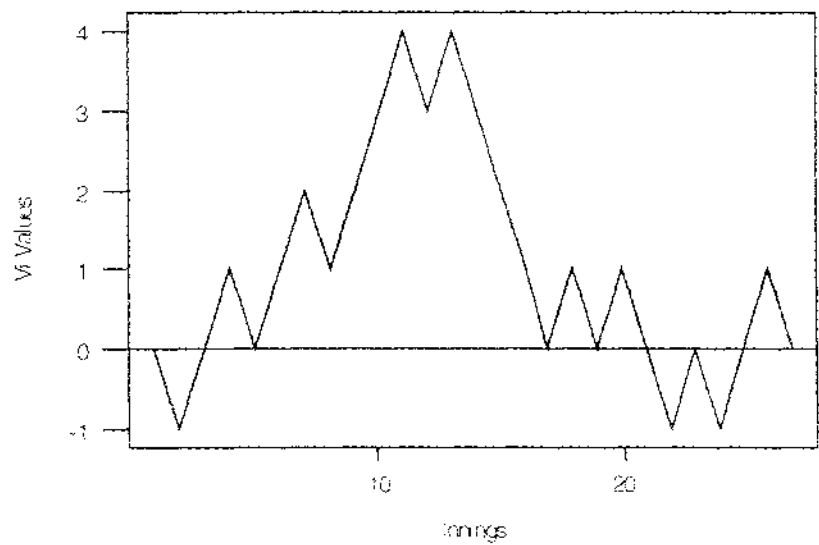


Figure 7. *Distribution-free CUSUM for M.J. Horne*

As the control limits are not breached this implies that the estimates for both player s have been consistently achieved.

c) A non-parametric CSCC procedure based on Within Group Rankings

(Bakir, Reynolds, 1979)

This non-parametric procedure was developed to quickly detect any shift in the mean process level. Using Wilcoxon signed rank statistics and within group ranking, a CUSUM type procedure is implemented. CSCC, as defined by Bakir et al, given in the sub-title is better known as the CUSUM (Cumulative Sum Control Chart). For the within group ranking subgroups are required. Where these do not occur naturally, they must be created artificially. In this instance artificial subgroups of size 4 are created. To detect shifts in the process level on the positive side the following control value is applied.

$$\sum_{i=1}^m (SR_i - k) - \min_{0 \leq m \leq n} \sum_{i=1}^m (SR_i - k) \geq h$$

Similarly, to detect deviations in the process level on the negative side the formula below is implemented.

$$\max_{0 \leq m \leq n} \sum_{i=1}^m (SR_i - k) - \sum_{i=1}^m (SR_i - k) \geq h$$

The Wilcoxon signed rank statistic is defined as follows:

$$SR = \sum_{X_i - \mu_0 > 0} \text{rank}(|X_i - \mu_0|) \quad (\text{Smith, P. 1993})$$

Signed rank has the advantage of considering the relative rankings of the magnitudes of the data points. A problem associated with the signed rank test in this instance is the assumption of a symmetrical distribution. Clearly the mixed distribution presented by individual batting scores is not symmetrical. As a result the median is used to circumvent this problem. In this case the mean μ_0 has been replaced by the median. Effectively the median is subtracted, then the values ranked. Finally, those ranks associated with points greater than or equal to the median are added to provide SR.

From the tabulated values of Bakir et al the parameters are chosen to provide a short nominal ARL in the vicinity of 17. Choosing $k=0$, $h=6$, and $g=4$ provides a nominal $ARL=17.07$. Applying the above procedure to the score series of Horne and Hartland an alarm is signalled in both cases after 8 innings (2 groups) indicating that the estimate of batting ability for both individuals has changed. For each individual this alarm was associated with a decrease in performance.

Bracewell (3)(1998) showed how shewhart control charts could be used to monitor bowling performance. It was also shown that the interpretation of zone run rules could be modified to accommodate the example presented by cricket. A brief demonstration of how multivariate control charts using Hotelling's T^2 statistic was also provided in the same study. These approaches are described and extended in the rest of this chapter.

3.3 Univariate Quality Control Methodology

It is appropriate to monitor an individual's performance with control charts. Provided the measurements of the 'product' are reflective of quality, function, or performance then the nature of the 'product' has no bearing on the general applicability of control charts (Montgomery, 1997).

The control chart is a useful tool in statistical process control. First developed by W.A. Shewhart, the Shewhart charts are widely accepted as standard tools for monitoring process of univariate independent and nearly normal measurements (Liu & Tang, 1996).

Control charts have three fundamental uses:

- 1 Reduction of process variability
- 2 Monitoring and surveillance of a process
- 3 Estimation of product or process parameters

(Montgomery, 1997).

It is the second use that is of the essence in the application to cricket, and possibly other sports. Process in industry is the parallel term to player performance.

Control charts have found frequent applications in both manufacturing and non-manufacturing settings (Montgomery, 1997).

The third use is also of relevance when dealing with team selection. This is a result of the interest in the estimate of an individual's ability in relation to other player's available for selection.

Before standardising the data it is important to note the parameter values as these are the estimates of the player's latent ability. In particular those applied to the bowling indices. These are initially standardised with respect to the population, by the nature of the indices, giving estimates of the individuals ability relative to those competing in the same competition. For charts that assume a standard normal distribution, these indices are standardised again, for within person evaluation, using individual means and standard deviations. When the data is standardised and tested with the quality control tools, the test is for how reliable the mean is as our estimate of the individual's ability.

3.3.1 Shewhart Control Charts

Shewhart charts are strongly dependent on the assumption of normality and independence. Also assumed is the absence of between subgroup variation when the process is in control (61.325 Study Guide). However, this statement is irrelevant in this study, as only individual innings observations are taken, that is only subgroups of size one exist.

The operation of a shewhart control chart with only action limits is slow detecting small shifts in process level (61.325 Study Guide). However, the Shewhart chart can be sensitised by utilising zone rules.

Shewhart Charts detect large shifts faster than EWMA & CUSUM (Montgomery, 1997). A major positive for using this type of scheme in cricket is the simplicity of implementation and interpretation.

The bowling indices are both normally distributed as determined in Chapter 2. Contribution, whilst being from a negative binomial distribution can be transformed to normality. As the negative binomial distribution is the discrete equivalent of the exponential function by raising to the power of 1/3.6 the transformation yields a Weibull Distribution, which is well approximated by the normal distribution (Montgomery, 1997).

Having acknowledged the fact that the data is normal, limits based on the normal distribution are appropriate. The Shewhart model for control charts is given as:

$$\mu_M \pm 3\sigma_M$$

where M is a given quality measure, referred to as the quality statistic. For appropriate use it is assumed that M is both independent and normally distributed with mean, μ_M , and standard deviation, σ_M . The central line is set at μ_M with the upper limits at $\mu_M + 3\sigma_M$ and lower limits at $\mu_M - 3\sigma_M$.

The control chart procedure is actually a sequence of hypothesis tests, with the control limits corresponding to a confidence interval (99.73% using the 3 sigma limits).

Having found the mean, \bar{x} , the standard deviation s and the distribution for both bowling indices (Attack index and Economy index), for the player under review, it is relatively straightforward to apply statistical quality control methodology. In particular control charts are ideal to monitor the relevant processes, which in this case relate to bowling performance. The nature of the indices sets it up ideally for use with Zone Rules for Control Chart Interpretation.

However minor alterations need to be made to the labelling of the zones and to the rules to make it compatible with the evaluation of an individual's performance.

Whilst being similar, there are fundamental differences in testing for quality in a product and sport. Typically quality control monitors the maintenance of certain control limits and deviation from a common mean.

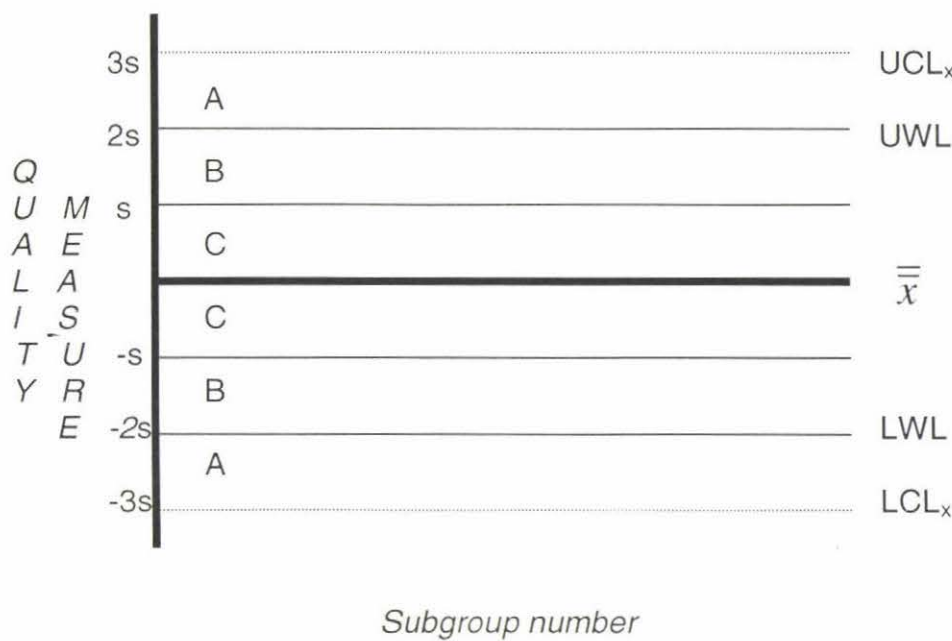


Figure 8. Control Chart with Warning Lines

Where UCL= Upper Control Limit
 LCL= Lower Control Limit
 LWL= Lower Warning Limit
 UWL= Upper Warning Limit

In a sporting context, which side of the mean a point falls is important and needs to be built into any control chart.

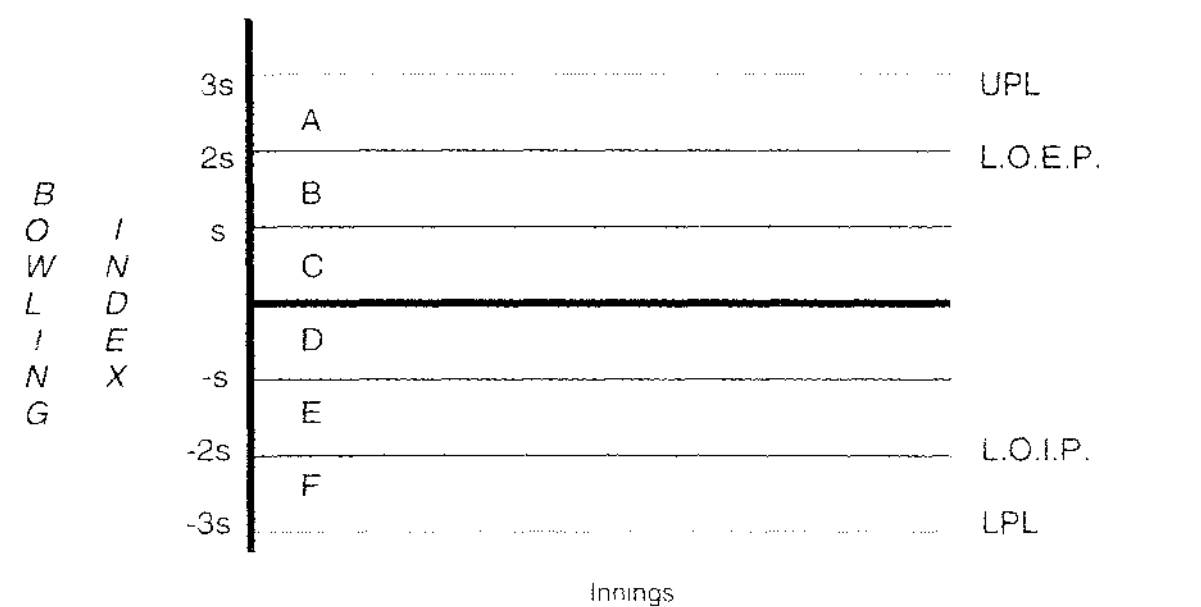


Figure 9. *Control Chart with Warning Lines for Cricket*

- L.O.E.P.
- Level of excellent performance
- L.O.I.P.
- Level of inferior performance
- UPL
- Upper performance level
- LPL
- Lower performance level

As discussed earlier, the probability that a point will fall into each of the zones can be established. The distribution for the respective indices is assumed to be standard normal. As estimates of μ and σ are unbiased the control limits are given by:

$$\hat{\mu} \pm 3\hat{\sigma}$$

Where σ is the standard deviation and μ is the mean. If a standard shewhart chart is applied, without zone rules, then to achieve a desired ARL of approximately 17 the control limits need to be altered. A significance level of 6% presents the appropriate ARL (making use of the fact that $ARL = 1/p$, $1/17 \approx 0.06$). Thus the control limits for a shewhart control chart are set approximately 2 standard deviations from the mean.. However, if all the zone rules are implemented, retaining the 3σ limits provides an ARL of approximately 17.

The performance of the Shewhart charts can be improved by the implementation of zone rules. These are listed in a cricketing context next.

Zone Rules for Cricket

Bracewell (3)(1998) proposed the following interpretations on the Zone Rules in a cricketing context. Understanding the intrinsic differences between the product attributes and sports performance allows the zone rules for Shewhart Control Charts to be manipulated to indicate 'Out-of-Control' conditions. This effectively means that the individual's performance relative to the team is no longer random in most instances. The tests given by Montgomery (1997) are modified to suit the situations presented in cricket.

Test 1. Extreme Points

Points that fall outside the control limits. Falling outside of zone A indicates an exceptionally brilliant performance relative to the team. Conversely a point falling outside zone F indicates an exceptionally awful performance relative to the team.

Test 2. Two Out of Three Points in Zones A or F and Beyond

Two of three performances in and beyond A shows continued excellent performance. Whereas, two of three performances in and beyond f shows continued inferior performance.

Test 3. Four Out of Five Points in Zones B or E and Beyond

Four out of five successive points in zone B or beyond indicates continued good performance. The same situation for zone E and beyond reveals continued bad performance.

Test 4. Runs above or Below the Centreline

This test considers long runs (eight or more successive innings) either strictly above or below the centreline. This indicates either consistently above expected performance (above the centreline) or consistently below expected performances (below the centreline).

Test 5. Linear Trend Identification

Following the continued increase or decrease of six successive points an alarm is given indicating the presence of a systematic trend. An increasing trend indicates the player is performing better relative to the team in each outing. The opposite applies to a decreasing trend, a relative worsening in performance.

Test 6. Oscillatory Trend Identification

When 14 successive points oscillate up and down an alarm is signalled. This may indicate a player who performs better in the first innings than the second or vice versa.

Test 7. Avoidance of Zones C and D Test

An alarm is signalled when eight successive points fail to fall in zones C and D. This could suggest an individual is either at the top of their game or fails.

Test 8. Run in Zones C and D Test

When 15 successive points fall in only zones C and D an alarm is sounded. This shows continued relatively average performance.

If no signal is given, this indicates the individual has performed as expected given the estimate of their ability.

These tests are useful in that they allow the identification of trends in performance. Each test caters for a relevant set of circumstances that can be useful to the selector, critic and player.

The numbers identifying the tests differ from those used in MINITAB. As mentioned earlier, the shewhart chart can be set up in two ways. The first uses only the control limits that are set at 2 standard deviations away from the process mean. Secondly all the run rules are applied along with the traditional 3-sigma limits. Both cases can be shown simultaneously.

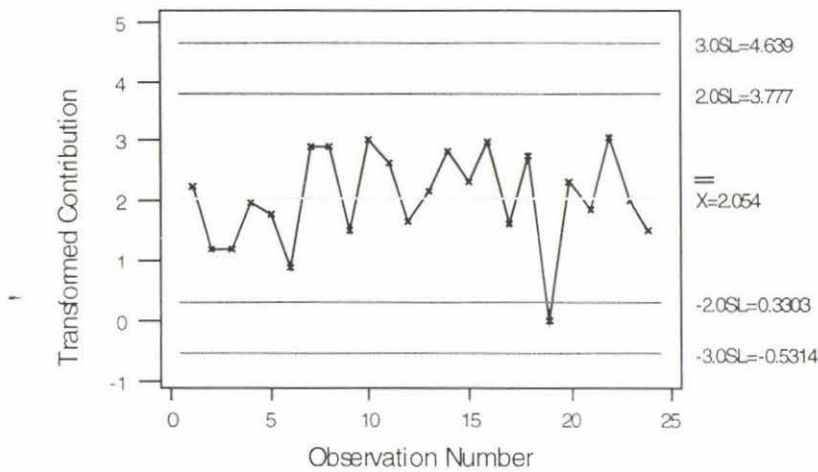


Figure 10. *Shewhart Control Chart of M.J. Horne's Transformed Batting Contribution With Zone Run Rules.*

The initial graph (zone rules, 3-sigma limits) reveals no signals, indicating that Horne is performing to expectations in terms of contributing to the team total. This also indicates that the impressive estimates for his natural ability are significantly adhered to.

If the 2-sigma limits are applied, an alarm is signalled at point 19. However, this corresponds to a score of zero, which is not necessarily an indication of form change.

As it is impossible to achieve a negative score with cricket data as applied to this type of control chart, negative control limits should be set at zero. Nevertheless, it is useful to see the impact the effect control limits have in this type of situation. Effectively every time an individual fails to score an alarm is signalled.

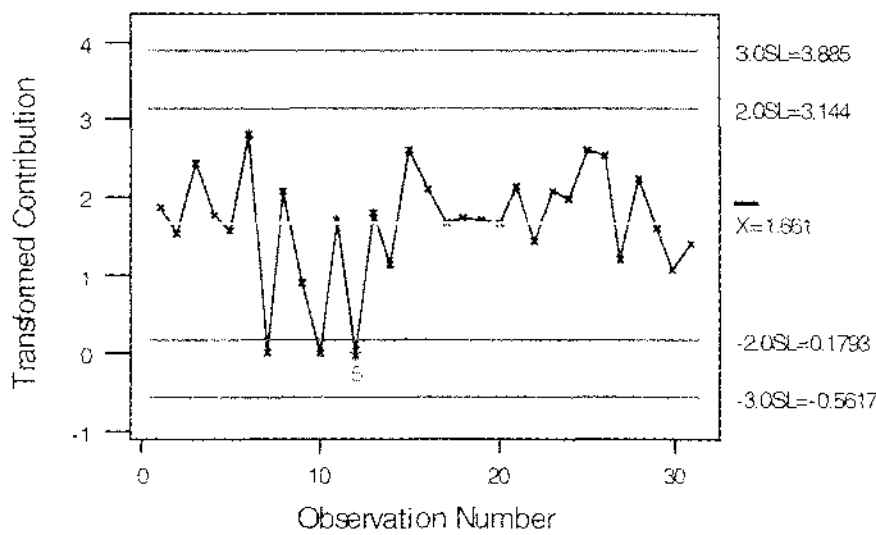


Figure 11. *Shewhart Control Chart of B.R. Hartland's Transformed Batting Contribution With Zone Run Rules.*

A violation of zone rule 5 (Two Out of Three Points in Zone F) at point 12 indicates a loss of form. This means Hartland has not been performing as well as expected given our estimate of his ability. This effectively means that Hartland's contribution to the team total is less than expected. Considering the 2-sigma limits, an alarm is given at point 7. Once more this corresponds to a zero score. This suggests that operating the 2-sigma limit boundaries is not as useful as the 3-sigma limits and zone rules.

Unless the player has modified their technique or developed new concentration skills, one would not expect a large shift in the estimate of natural ability. As the Shewhart does not detect small shifts in process as quickly as the CUSUM or EWMA schemes, they are investigated next.

3.3.2 CUSUM

The British Statistician Page first proposed cumulative sum (CUSUM) charts in 1954. This procedure involves cumulating sums, such that past values have an impact on the control statistic.

All measures of performance involved in this study are standardised to make computations somewhat easier for setting up general standards and relating these in terms of player performances. Also the implementation of the charts is designed to monitor the given estimate of ability. It is preferable to work with one-sided standardised CUSUMs for the case presented by cricket. As mentioned earlier, a special situation arises in the application of quality control methodology to sport. If a player's performance process changes such that an alarm is signalled, it is necessary to note if the alarm was due to superior or inferior performance. Obviously if an individual's performance is improving the desired situation has occurred. Standardised values are used to compare within player performance on relative scales.

As we are dealing with individual observations the statistic required for the CUSUM scheme is given as:

$$z_i = \frac{X_i - \bar{X}}{\sigma_X}$$

For detecting shifts on the upper side of the mean the procedure is defined as

$$S_{Hi} = \max \{0, (z_i - k) + S_{Hi-1}\}$$

With S_{H0} set at zero. The slack constant, k , must be less than 3 or a situation akin to a shewhart chart is implemented, which detects only large shifts. Generally this value is 0.5, designed to detect smaller shifts.

The next part of the chart is the adoption of some threshold value, for which when crossed, indicates an out of control situation. This value is referred to as h .

A similar procedure is used for detecting shifts on the lower side of the mean:

$$S_{Li} = \max \{0, (-z_i - k) + S_{Li+1}\}, \text{ with } S_{H0} \text{ set at zero.}$$

The optimal values for the parameters of the CUSUM procedure can be found from Gan's (1991) nomograph. To achieve a small nominal ARL of approximately 17 h is set a 0.25 and k to 1.7. Whenever a signal is given, this implies an out of control situation, that is the estimate of the player's natural ability has changed. If a cause for the alarm is found the CUSUM is reset to zero.

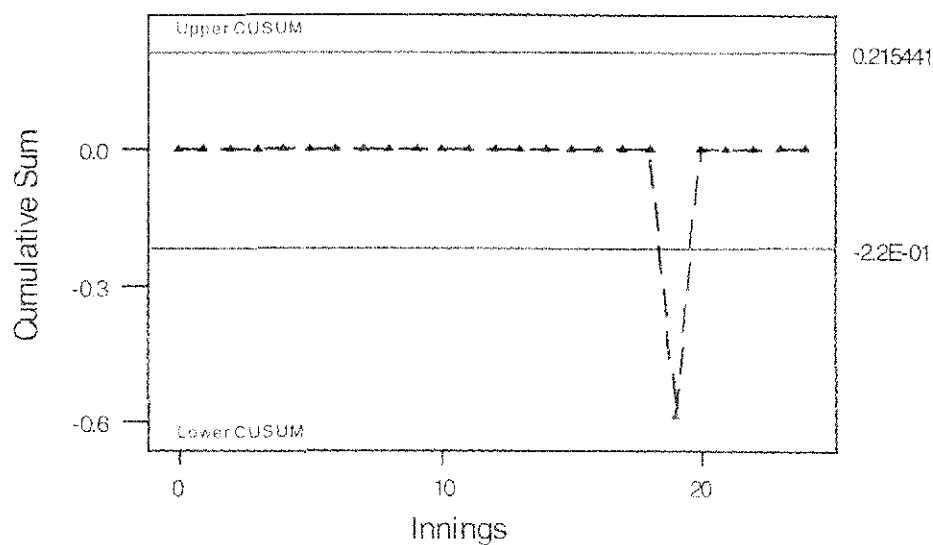


Figure 12. CUSUM Control Chart of M.J. Horne's Transformed Batting Contribution.

A signal is given at point 19, once again corresponding to a score of zero. Apart from the one 'duck', no other signals are given, indicating that Horne is performing to expectations in terms of contributing to the team total.

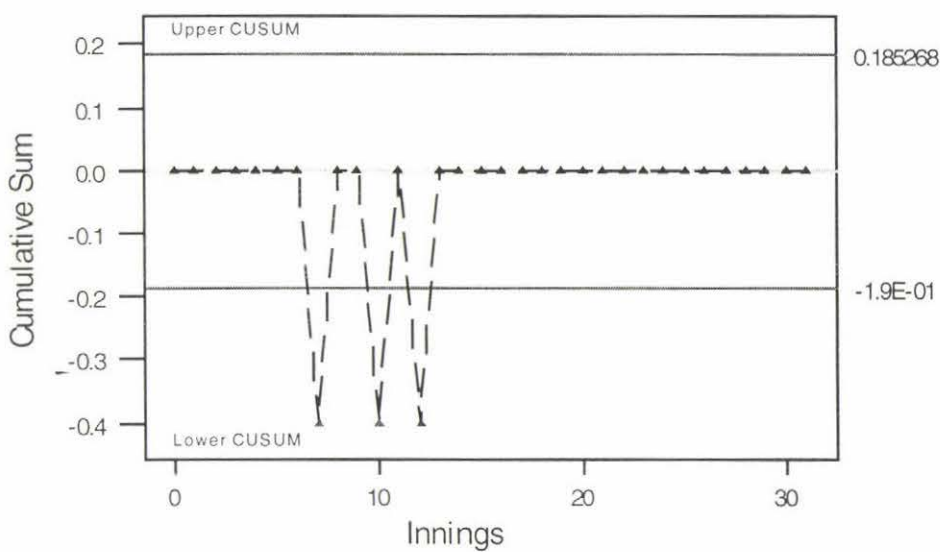


Figure 13. *CUSUM Control Chart of B.R. Hartland's Transformed Batting Contribution.*

A signal is given at each point Hartland failed to score. This is the same as the Shewhart chart with only the 2-sigma limits in operation. It is not necessarily confirmation of deterioration in form. However, with three alarms in a short space of time, this suggests that the true estimate of mean batting ability for Blair Hartland is actually significantly less than noted.

3.3.3 EWMA

An alternative to the Shewhart control chart is the Exponentially Weighted Moving Average control chart developed by Roberts (1959). The performance of the EWMA chart is similar to the CUSUM scheme, but easier to set up and operate (Montgomery, 1997). The EWMA chart proves to be useful when it is not practical to take more than a single observation per sample, as is the situation presented by cricket. An advantage is the effect of averaging to detect process level shifts and damping out some effect of random errors on individual observations, due to the reliance on past observations.

The exponentially weighted moving average is defined as follows:

$$z_i = \lambda x_i + (1-\lambda) z_{i-1}$$

(Montgomery, 1997)

Where λ is a constant greater than zero, but no larger than one. The starting value (when time, i , is one) is the process target and hence equal to the population mean. The control limits are defined as follows

$$UCL = \mu_0 + L\sigma \sqrt{\frac{\lambda}{(2-\lambda)}[1-(1-\lambda)^{2i}]}$$

$$LCL = \mu_0 - L\sigma \sqrt{\frac{\lambda}{(2-\lambda)}[1-(1-\lambda)^{2i}]}$$

$$CL = \mu_0$$

The factor L represents the width of the control limits. λ indicates the weighting placed on previous values. To obtain an appropriate ARL, of approximately 17, values are taken from Crowder's (1989) nomograph to detect a shift of one standard deviation. L and λ are set at 2 and 0.25 respectively.

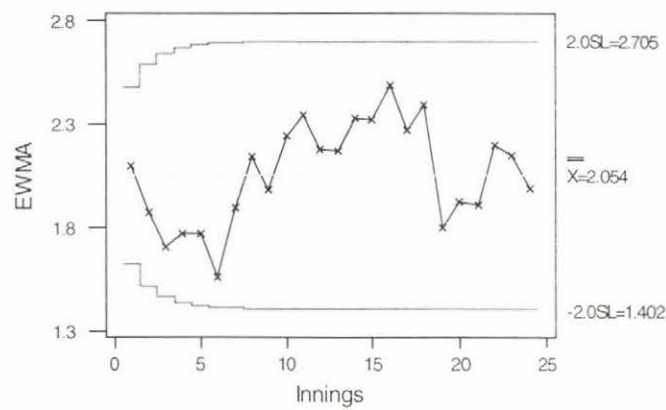


Figure 14. *EWMA Control Chart of M.J. Horne's Transformed Batting Contribution With Zone Run Rules.*

As no signals are given, further confirmation that Horne is performing to expectations in terms of contributing to the team total is found in the above EWMA chart.

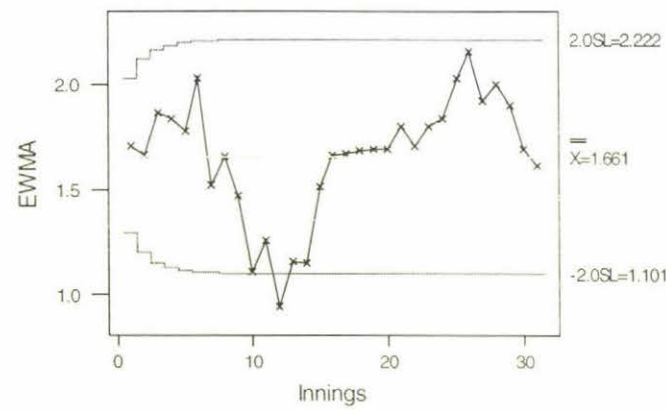


Figure 15. *EWMA Control Chart of B.R. Hartland's Transformed Batting Contribution With Zone Run Rules.*

In contrast to the Shewhart and CUSUM scheme, the EWMA chart only signals after the third zero score, suggesting that there is no significant change in form. As the EWMA scheme is not readily influenced by the zero scores, this is the preferred option, for assessing batting contribution.

3.4 Proposed Control Chart Based on Quartiles.

The performance of a control chart in the traditional sense is very sensitive to the assumptions when dealing with individual scores due to the relatively high likelihood of failure. When the assumptions of either independence or normality fail, traditional methods introduce high probabilities for false alarms (Type I errors) (Vasilopoulos, Stamboulis, 1978). Problems arise in that the distribution for individual batting scores is mixed geometric, as indicated in Chapter 2, and thus transformation will not yield an approximately normal distribution. However, we are wanting to preserve the influence of extreme scores as these are an indication of an individual's capability, which non-parametric methods negate. A potential problem with standard control charts is that a single extreme outlier can trigger an out-of-control situation. We want a technique that is devoid of a distribution, or enhances, or is adaptable to a given distribution, and makes use of extreme values for individual observations. An approach is introduced here based around the simple concept of quartiles.

The proposed method is based on Quartiles. As sample sizes are generally small, theoretical quartiles are used rather than the observed values, which also allows parametric influences. To maintain the attachment to a given distribution, the theoretical quartiles are gained using the estimate of the mean. Hence, outliers potentially influence these values and as a result information pertaining to outliers is not lost. How these values were obtained is discussed later. Essentially, due to the nature of the 'ducks and runs' distribution, specifically the number of zero's occurring, the use of a shewhart type scheme is inappropriate. Moreover, the LCL is incompatible with the number of zeros that are generated, effectively causing an alarm every time a zero is recorded.

From Chapter 2 there is sufficient evidence to imply that individual batting scores are from a mixed geometric distribution. As this distribution is discrete, to find the theoretical quartiles involved investigating the cumulative probabilities. Values for the theoretical quartiles were computed based on the mean score and mean contribution.

Using EXCEL a table was produced listing the probability of a given score, given varying batting parameters, mean score and mean contribution. A scatterplot of the mean contribution and mean score revealed a linear relationship between score and contribution. Consequently a 99% Prediction interval, shown below, from a simple linear regression gives the most likely region of interest, giving general bounds on which to form the table.

The spreadsheet was designed to sum all previous values. Where the cumulative probability was closest to the values of the quartiles (0.25, 0.5, 0.75), the corresponding score was taken. Invariably all the values fell between the positive integer values. As individual scores are represented by only positive integers and zero, the precise score at which they were obtained is not needed, so values given are of the form $(m+n)/2$ where m and n are the two immediately neighbouring points that surround the quartile value of interest.

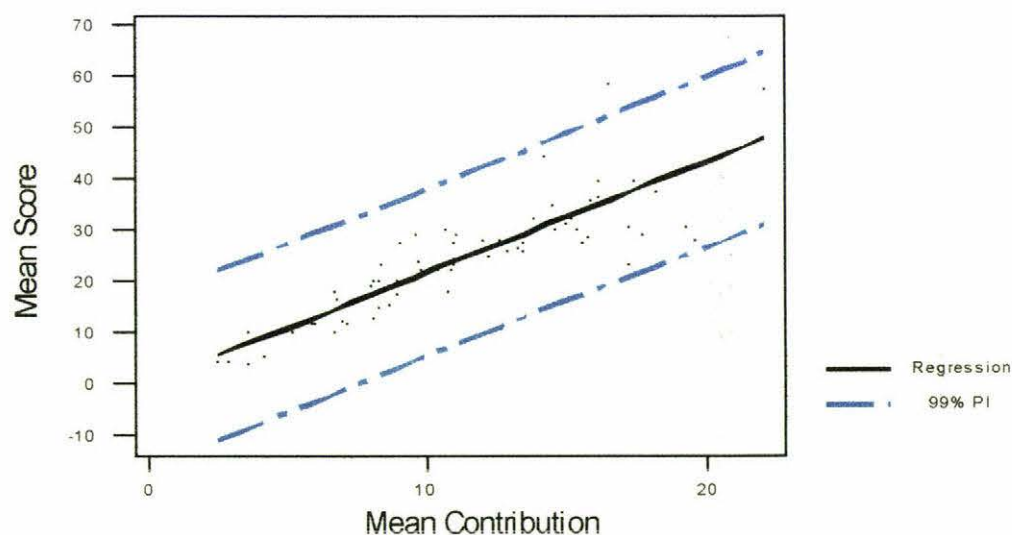


Figure 16. *Fitted Line Plot of Mean Contribution Vs Mean Score*

Three control lines are drawn corresponding to the lower quartile, median and upper quartile. These theoretical values are observed from tables listed in Appendix E, utilising the average score of the non-zero component and mean contribution. The exact reasoning behind this is given in the development of the mixed model described earlier. The mean score and mean contribution are then used to evaluate the theoretical median and Quartiles. Thus, this method has adapted from being distribution free to being based upon mixed geometric.

To signal an alarm rules are developed based on the probabilities of a point appearing in a certain zone or pattern of zones which can be used to identify performance. These rules try to emulate the zone rules that supplement the Shewhart charts.

Due to the nature of cricket, it takes a number of observations to be able to effectively estimate a player's ability. It is logical to also infer this holds for a change in ability. The quartile chart with zone rules described below will give a false signal on average every 16.8 innings. If five rounds of Shell Trophy matches (five matches, ten innings) are played, this equates to approximately one signal in two seasons. For an out of control situation, either through loss of form or an improvement a signal will be picked up approximately in one season. The run rules designed for use with this chart are set at a 95% confidence limit and are relatively loose. Altering this to a stricter plan 99.5% would correspond to larger run lengths which is obviously unacceptable because of the relative lack of sampling opportunities.

The figure below details the positioning of the zones about the quartiles.

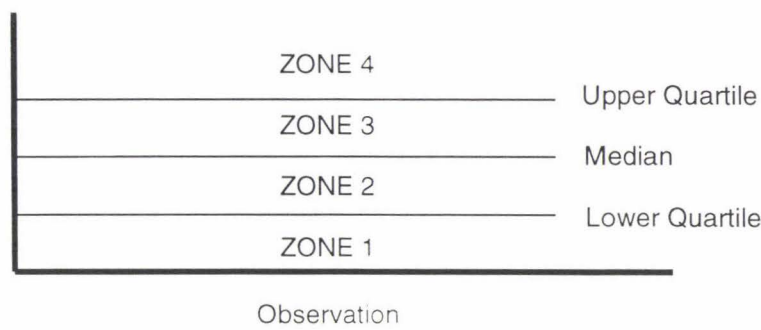


Figure 17. *Quartile Control Chart*

In creating a set of zone rules for the Quartile chart it was attempted to have run rules comparable to those implemented by the Shewhart Control Chart. Obviously shewhart charts assume normality, which is not found in this situation. A situation akin to a uniform distribution is created. Four zones with equal probability (0.25) are created based upon the Median, Upper Quartile and lower Quartile. Knowing the probability that a point falls into a certain zone enables improbable run lengths to be established. That is certain patterns of data – similar to that used in the zone rules for Shewhart Control Charts – that are unlikely to occur in an in-control situation can be established. The run lengths are given where the probability for a given string of values drops below a given error setting. In this instance run lengths are given for an error setting of 0.05. That is there is less than a 95% chance of a given pattern occurring. This setting is chosen as it is sufficiently tight for the context presented. 6 rules are proposed as follows.

Zone Rules for Quartile Chart

An 'alarm' in this instance refers to a change in form signal.

H_0 : Playing to natural ability

1 Runs in Extremities

Involves a run of points exclusively in zones 1 and 4.

As there is a 0.5 probability of a score occurring in these zones, under H_0 an alarm will be sounded after 5 points

$$P[\text{False Alarm}] = 0.5 \times 0.5 \times 0.5 \times 0.5 \times 0.5 = 0.03125$$

In terms of a player's performance this signal suggests that a player is either failing or going on to make a big total, relative to the estimate of their natural ability.

2 Runs in Central Zones

Similarly, this involves a run of points exclusively in zones 2 and 3.

As there is a 0.5 probability of a score occurring in these zones, under H_0 an alarm will be sounded after 5 points

$$P[\text{False Alarm}] = 0.5 \times 0.5 \times 0.5 \times 0.5 \times 0.5 = 0.03125$$

In terms of a player's performance this signal suggests that a player is getting starts but not going on to make a big total, relative to the estimate of their natural ability.

3 Runs in one Zone

As there is a 0.25 chance of falling in any give zone under H_0

$$P[\text{False Alarm}] = 0.25 \times 0.25 \times 0.25 = 0.015625$$

3 runs in one zone results in an alarm. Depending on which zone the alarm originates, influences the interpretation of the signal. An alarm from zone 4 indicates an extremely good series of scores, whereas 3 points in zone 1 shows poor form.

4 Runs above or below the median

This involves a series of scores in either Zones 1 and 2 or in Zones 3 and 4. As there is a 0.5 probability of a score occurring in these zones, under H_0 an alarm will be sounded after 5 points

$$P[\text{False Alarm}] = 0.5 \times 0.5 \times 0.5 \times 0.5 \times 0.5 = 0.03125$$

Again depending on which region has caused the alarm sways the interpretation. Runs above the median are highly favourable and thus indicate good form. Conversely a signal below the median indicates inadequate performance.

5 Points increasing /decreasing

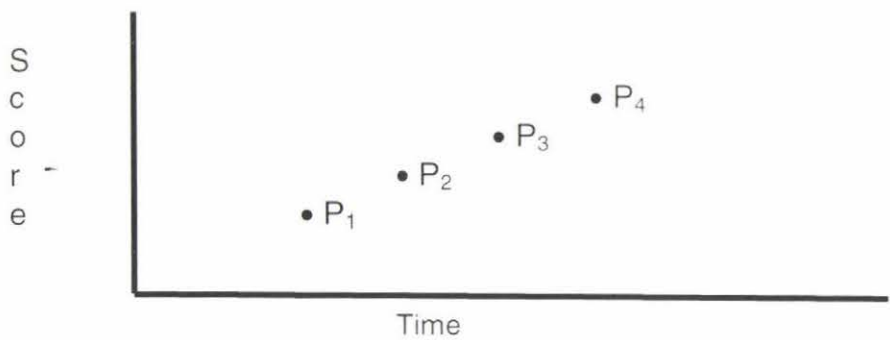


Figure 18. Establishing Number of Consecutive Increasing Points for Alarm

The number of strictly increasing points can be found by solving an integral of the following form such that the resultant probability is below 0.05.

$$P[P_4 > P_3 > P_2 > P_1] = \int_0^1 \int_0^{P_4} \int_0^{P_3} \int_0^{P_2} 1 .dP_1 .dP_2 .dP_3 .dP_4$$

This is because for the sequence to be increasing P_1 must lie in the range $(0, P_2)$; P_2 must fall in the range $(0, P_3)$ and so forth. Solving the above integral yields a probability of 0.0417 (4d.p.), the first sequence (4 consecutive points) to fall below the defined 5% Type I error.

A similar situation exists for points strictly decreasing except the integration takes place over the ranges $(p_n, 1)$ instead of $(0, p_n)$.

Therefore for increasing or decreasing points $P[\text{False Alarm}] < 0.05$ occurs after 4 points (0.0417). If the points are increasing this suggests a player is improving, conversely a decrease suggests poor form.

6 Runs outside of zone 4 or any specified zone

$P[\text{False Alarm}] = 0.75^{11} = 0.0422$

Runs outside of zone 4 indicate a player has not gone on to record a successful score as frequently as expected according to the estimate of their ability.

Zone rules for Quartile control chart

Rule	Number of points for signal	
	P[False Alarm] 0.05	0.005
1 Extreme Points	5	8
2 Central Points	5	8
3 Points in One Zone	3	4
4 Points Above / Below Median	5	8
5 Points Increasing / Decreasing	4	5
6 Points Avoiding One Zone	11	19

Table 6. Signal Length in Consecutive Points.

Situations 1 and 2 signal performance patterns in performance that can be related to tactics.

- 1 Player either scores large or is dismissed cheaply. May indicate they are a shaky starter, but once underway, are set to push on for a relatively productive innings
- 2 Indicates a player is dealing with mediocrity. Suggests that they are easily settled but do not kick on to score the large scores that are expected.

The zones can be thought of in a more familiar cricket thinking – Zone 1 is a failure, zones 2 and 3 are ‘starts’ and zone 4 is a success.

The rules can be divided up into favourable responses, unfavourable responses and indifferent responses. These can be used similar to the CUSUM scheme. How these are handled will affect run length.

Rule	Improvement	Worsening
3	points in zone 4	points in zone 3
4	above median	below median
5	increasing	decreasing
6	points avoiding zone 1	points avoiding zone 4

Whilst the run length appears to be small, when thinking of the amount of cricket played, it equates to a reasonable period of time.

3.5 Comparison of Univariate Control Chart Performance

To enable comparisons of the previously discussed univariate control charts an exponential type population was assumed for the scores. These were than transformed by a power of 1/3.6, to introduce normality, as discussed earlier. A number of simulations were then performed and the result displayed below. In each case alarms were reset after a signal.

In all cases an ARL of approximately 17 was sought for the in-control situation and the charting parameters were chosen accordingly. The first column, parameters, refers to the values that define the ‘ducks and runs’ distribution. First value quoted is the mean score, followed by the mean contribution. Corresponding values for the control chart are found in Appendix E. A mean score of 30 and a mean contribution of 16 represent the in-control situation. The out-of control situations change by increments of five runs and one percent respectively. The charting parameters for each of the methods are given as follows for the mean score and mean contribution:

- CUSUM: $h=0.25$, $h=1.7$
- EWMA: $\lambda=0.25$, $L=2$
- Shewhart: $k=2$, All zone rules used.
- Non-parametric EWMA (Hackl, Ledolter, 1992): $\lambda=0.25$, $g=4$, $h=\pm 0.2980$.

These values were confirmed through trial and error from initial values obtained approximately from tables mentioned previously. The ARL for each situation was established from 10 samples of 500 individual innings.

Parameters	Quartile	CUSUM	EWMA	Shewhart	Non-parametric EWMA
30,16	16.8	15.3	16.7	13.4	19.8
35,17	9.9	25.0	11.5	14.1	16.5
25,15	10.8	15.6	14.2	11.7	18.9
40,18	8.7	13.4	6.8	11.5	16.5
20,14	8.6	17.2	8.8	11.1	19.4

Table 7. Comparative ARL's.

From the run lengths the functionality of the Quartile Chart is seen. Whilst it is not 'quick' in picking up changes in performance the cautious approach is especially relevant in sport statistics where perseverance with a selected individual can reap rewards. The EWMA chart is the next best performer, however a decrease in ability by 5 runs and 1% is picked up effectively. The EWMA is chosen as the next best performer as it is the only other control chart that has an ARL for the out-of-control situation less than the in-control ARL, as well as having shorter ARL's for the more extreme changes (40,18; 20,14) than the moderate changes (35,17; 25,15).

In addition, Table 8, clarifies the first table by providing a more equivalent comparison. This is achieved by considering the ratio of the out-of-control ARL to that of the in-control ARL. That is:

$$Ratio = \frac{\text{Out - of - control ARL}}{\text{In - control ARL}}$$

Parameters	Quartile	CUSUM	EWMA	Shewhart	Non-parametric EWMA
30,16	1	1	1	1	1
35,17	0.59	1.03	0.69	1.05	0.83
25,15	0.64	1.03	0.85	0.87	0.95
40,18	0.52	0.92	0.41	0.86	0.83
20,14	0.51	0.93	0.52	0.83	0.98

Table 8. Ratio's Comparing Run Length From In-Control to Out-of-Control.

Table 8 largely confirms the results of Table 7, that the Quartile chart out performs the other methods for detecting shifts in the natural ability of an individual.

As found in Chapter 2, the distribution of individual scores can not be assumed to be of an exponential nature. The bizarre nature of the EWMA, CUSUM and shewhart_ARL's is easily explained by the assumption of normality. Clearly scores are non-normal and not transformable to normality. The mixed distribution for scores is required to handle the larger than expected number of zero's, which are modelled by contribution. Mean contribution is almost always going to be less than the mean score (unless the team total averages less than 100). Due to the nature of the distribution involved this has an immediate impact on the shape of the distribution, as shown previously in chapter 2, figure 3.

Thus the increased number of zeros has a huge impact on the estimated standard deviation in the transformed data. Also as the non-zero portion is exponentially based, the higher the mean the wider the range of values. Because the standard deviation is enlarged due to the number of zeros and the distance to the transformed non-zero portion, a moderate shift of a positive increment of five runs coupled with a 1% increase in average contribution has the longest ARL. The 1% increase in average contribution results in a lower spread as the number of zero's increases meaning the results are more compact as the shift in mean is only moderate. Obviously in a situation where average contribution equals average mean then the transformation to normality is stable and normality can be inferred. The effect of the transformation to introduce normality is shown graphically below. Figure 19 details the transformation on randomly generated exponential distributed (discrete) score data_ (8500 rows) with a mean score of 30 and mean contribution of 16. The disproportionate number of zeros is easily seen in the transformed histogram. Whilst a similar problem exists with the contribution data shown in figure 20 (8500 rows, mean contribution = 16), the approximation of normality is clearly better.

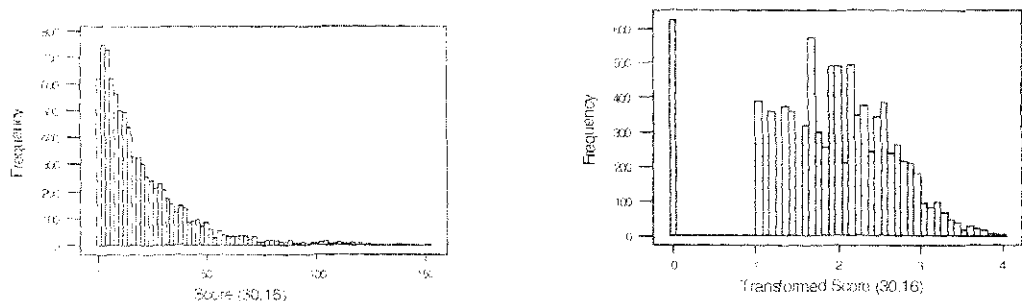


Figure 19. *Histograms of simulated Score data with and without Transformation*

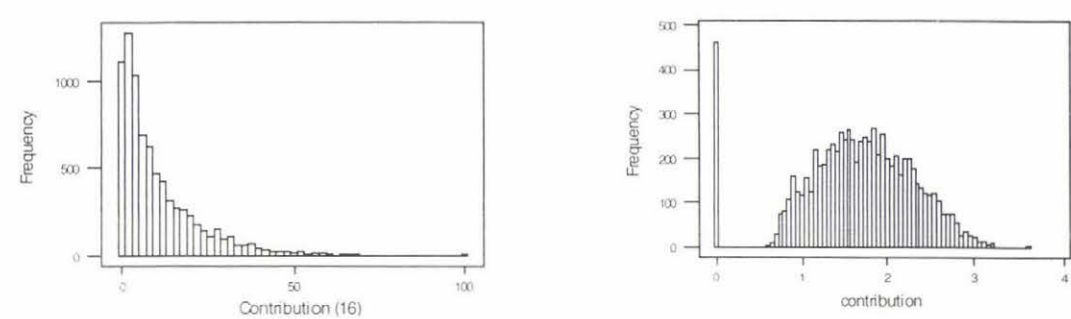


Figure 20. *Histograms of simulated Contribution data with and without Transformation*

The theoretical quartiles are obtained by finding the mean contribution, then finding the mean of the non-zero scores. These values are then used to obtain the corresponding quartiles given in Appendix E.

Application of Quartile Chart

From the data, Horne has a mean contribution of 21.20, which is set at 21, and a mean non-zero score of 59.96, set to 60. Using the tables in Appendix E the quartile values as shown below. The same procedure is performed with Hartland mean score of 30.54 (31) and mean contribution of 10.77 (11).

	Horne	Hartland
Upper Quartile	82.5	42.5
Median	39.5	14.5
Lower Quartile	14.5	6.5

Table 9. *Theoretical Quartile Limits for Horne and Hartland*

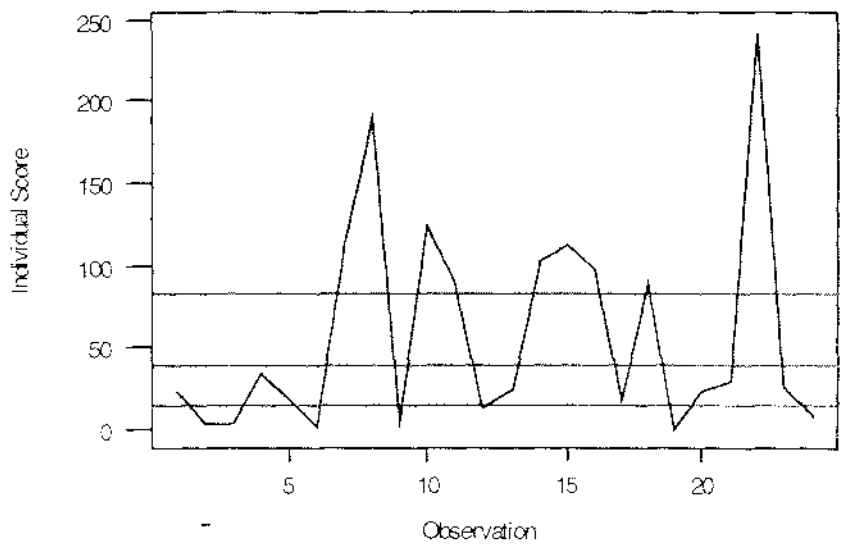


Figure 21. *Quartile Control Chart for M. J. Horne*

From the above plot 2 alarms are signalled for the individual scores of M.J. Horne. These are at points 5 (5 points below median) and 16 (3 points in same zone). The previous charts gave no signals: however, they are not directly comparable as they deal with contribution as opposed to score. It is of interest to note that Horne only played 8 innings for Auckland before moving south to Otago. If this is taken into account with our interpretation, this suggests that Horne under performed for Auckland, yet upon moving to Otago blossomed, with the only alarm coming from three successive points in zone 4. It becomes obvious that the southern form is behind his selection in the New Zealand Cricket side.

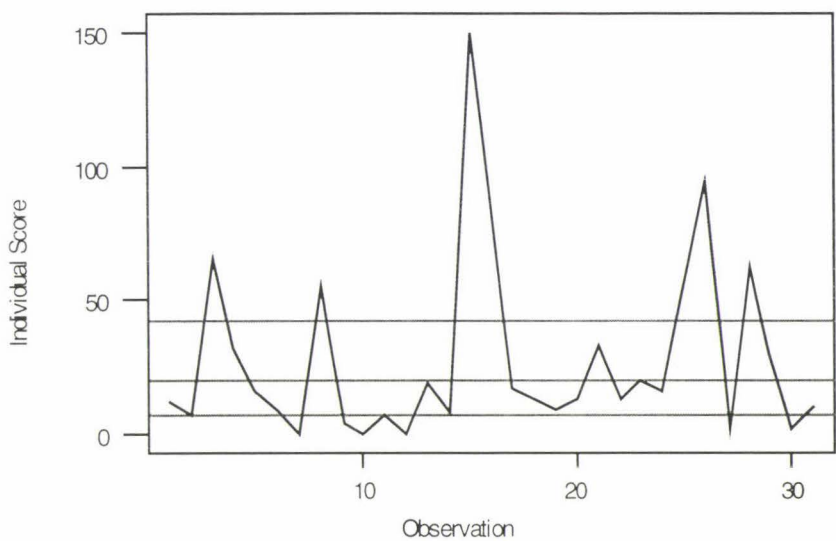


Figure 22. *Quartile Control Chart for B. R. Hartland*

Former New Zealand representative, B.R. Hartland’s quartile chart portrays a different example. Previously the CUSUM and Shewhart charts picked up a change for the worse in ability at point 12 (3-sigma limits, zone rules) as did the EWMA. This loss of form is pre-empted by the Quartile Chart by a signal at point 6 (4 decreasing points).

3.6 Multivariate Charts

Multivariate control charts are a useful extension of the univariate tests discussed previously. They enable the analysis of several different performance measures simultaneously. This is especially relevant when dealing with the performance data generated by bowlers. Several methods exist.

Lowry et al (1992) notes the MEWMA scheme performs better than MCUSUM (Crosier, 1988) when the process is initially out of control. They also note that performances are nearly matched if a shift in the mean vector is delayed. As a consequence only the multivariate EWMA scheme is considered in this study.

Hotelling T^2 is effectively an extension of the Shewhart Chart (Lowry, Montgomery, 1995). Consequently normality is assumed. Based on the most recent observation and is insensitive to small and moderate shifts in the mean vector.

Montgomery mentions that a dimension reduction method, such as principal component analysis can be of use in multivariate control charting procedures. This is particularly useful when interpreting an out-of-control situation. However the principal components do not always provide meaningful indices. Indeed, Bracewell (1)(1998) found that principle component analysis provided ineffectual results when dealing with the bowling data.

A problem arises in dealing with the Multivariate Allrounder Data. In a single match a player has the opportunity to both bat and bowl at most twice. It is also feasible for an individual in the playing eleven to do neither. However, it is proposed that this can be dealt with if managed on an entire match basis, as opposed to innings by innings.

Each match then becomes the sample and a varying number of subgroup sizes (1 or 2) and dimensions (1,2 or 3) are presented.

Thus the control limits must be dynamic to accommodate an individual's input in a game. This is achieved by taking the player's maximum value for each of the p dimensions involved in any innings. p is given by one of three states 1, 2 or 3 and control limits are plotted accordingly. These dimensions are as follows:

1 = Batting (Transformed Contribution)

2 = Bowling (Attack index)

3 = Bowling (Economy Index)

Maximum represents the player's most favourable performances with respect to the team. However, if we are interested in an individual under performing then the minimum must be dealt with. It is possible to plot both simultaneously, much the same as a CUSUM scheme.

3.6.1 Hotelling T^2 Chart

A further extension of the methods outlined above is to consider the two bowling indices at once using multivariate control methods. An ideal method for this is a multivariate T^2 chart as explained by Bracewell (3)(1998).

In this section multivariate quality control procedures are utilised in two situations. The first deals with bowling where 2 independent, standard normal populations are presented. Secondly, All-round performance is analysed where 3 independent factors are considered the 2 bowling indices and the normalised and standardised batting contribution data.

The general convention is to consider p correlated characteristics, which are assumed to be from a multivariate normal distribution. The i th individual observation is then given as:

$$X_i' = (X_{i1}, \dots, X_{ip})$$

Which leads on to the estimated mean vector as follows:

$$\bar{X}_i = \frac{1}{m} \sum_{i=1}^m X_{ij} \quad \text{where} \quad \bar{X}'_m = (\bar{X}_1, \dots, \bar{X}_p)$$

The estimated covariance matrix is then:

$$S_{m_i} := \frac{1}{m-1} \sum_{j=1}^m (X_{ij} - \bar{X}_{m_i})(X_{ij} - \bar{X}_{m_i})'$$

This is then combined to give the control statistic:

$$Q_i = (X_i - \bar{X}_{m_i})' S_{m_i}^{-1} (X_i - \bar{X}_{m_i}) \tag{61.725 Study Guide}$$

Due to the large sample sizes involved in preliminary research, and they manner in which they are dealt with, it is assumed that the mean and standard deviations are the true population values. Hence the control limits are found using the chi-squared distribution, where p is the number of characteristics involved. Thus the control limits can be expressed as follows:

$$\begin{aligned} \text{LCL} &= \chi^2_{1-\alpha/2,p} \\ \text{UCL} &= \chi^2_{\alpha/2,p} \end{aligned}$$

(Tracey, N.D., Young, J.C., & Mason, R.L., 1992)

Our first interest is with how well a bowler complies with the estimate of ability. In this instance there are two variables Attack Index (AI) and Economy Index (EI). Each individual observation is considered as a single subgroup. In the following analysis it is assumed that the individual's results have been standardised. The estimated mean vector contains the grand mean for attack (0) and the grand mean for economy presented as follows:

$$\bar{X}_i = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

Then the estimated covariance matrix is found using the standard deviation for AI (1) and EI (1) and the correlation between EI, AI (0) and AI, EI (0).

$$S_{m_i} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

This information is combined to give the control statistic for the observation X_i :

$$T^2_K = (X_i - \bar{X})' S_m^{-1} (X_i - \bar{X})$$

Due to the standardised nature of the indices and the absence of any correlation, the covariance matrix is simply the identity matrix. As a result it is displayed in a simplified form below:

$$T^2_K = (X_{I\text{ (ATTACK)}})^2 + (X_{I\text{ (ECONOMY)}})^2$$

In order to have an ARL comparable to that used in the univariate examples the significance level has to be calculated. For an ARL of 16.8 the significance level is simply $1/16.8 = 0.06$, indicating a significance of 94%. The control limits are then found at a 6% level of significance:

$$LCL = \chi^2_{1-\alpha/2;p} = \chi^2_{0.97;2} = 0.0609$$

$$UCL = \chi^2_{\alpha/2;p} = \chi^2_{0.03;2} = 7.0131$$

(Tracey, N.D., Young, J.C., & Mason, R.L. ,1992)

A chart is then drawn to illustrate this. These examples follow. Chris Brown and Paul Wiseman are used as examples. Both are specialist bowlers, Brown, right arm medium fast and Wiseman, right arm off-spin.

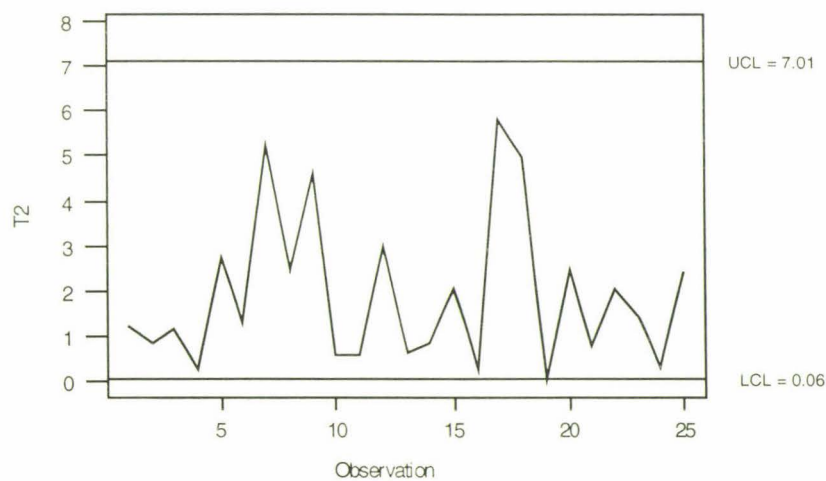


Figure 23. T^2 Control Chart of C.M. Brown for Bowling Indices

The previous plot indicates no violation of the control limits for performances by Chris Brown. This indicates that the estimate for his natural ability is correct, and that the observed performances fit within the expected scope.

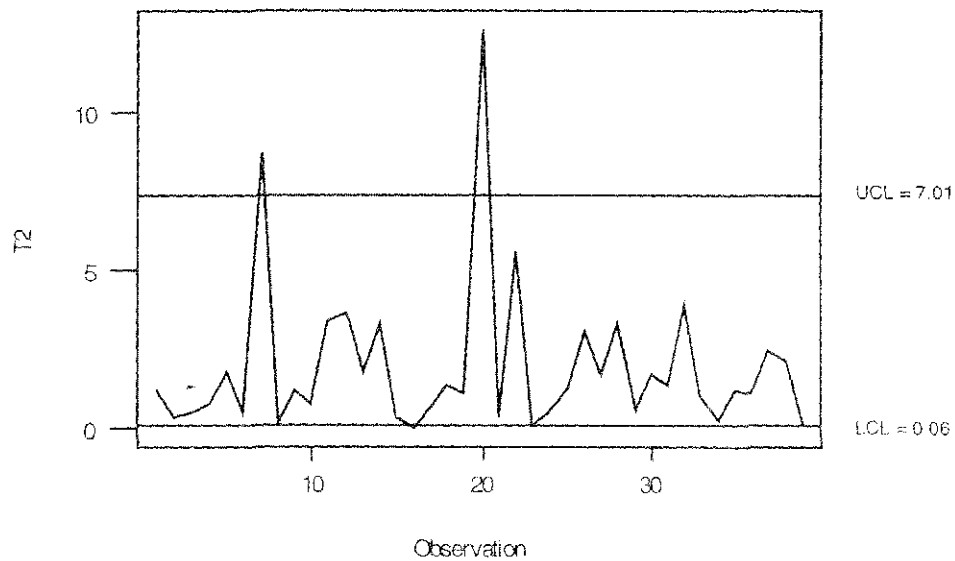


Figure 24. T^2 Control Chart of P.J. Wiseman for Bowling Indices

Two signals are given for Wiseman. However, due to the use of the squared distance from the origin used in the construction of the T^2 Chart, we are unable to detect which index is responsible and if the alarm is due to improved or decreased individual performance.

As $p=2$, the relationship $T^2_k = UCL$ leads to a bivariate control chart with a circular control region. Corresponding index scores are then plotted as shown on the next page.

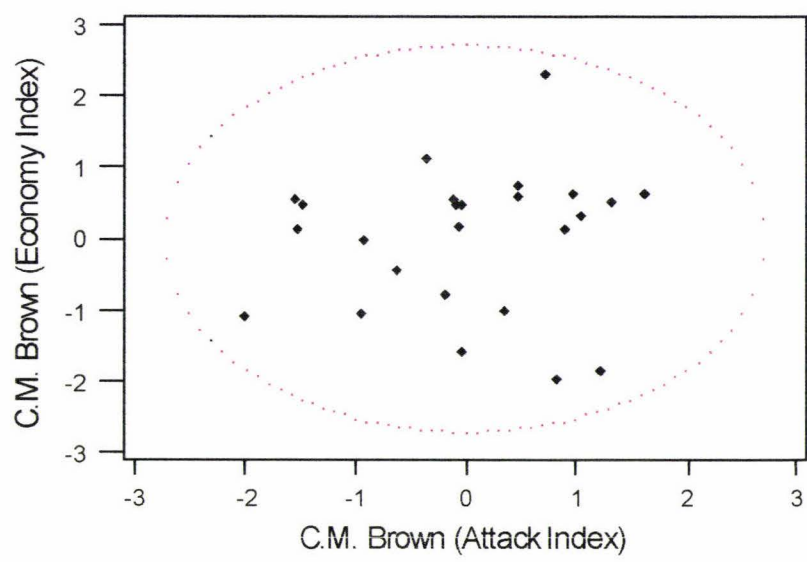


Figure 25. *Bivariate Control Chart of C.M. Brown for Bowling Indices*

For Brown, as with the T^2 Control Chart, all points fall inside the control region.

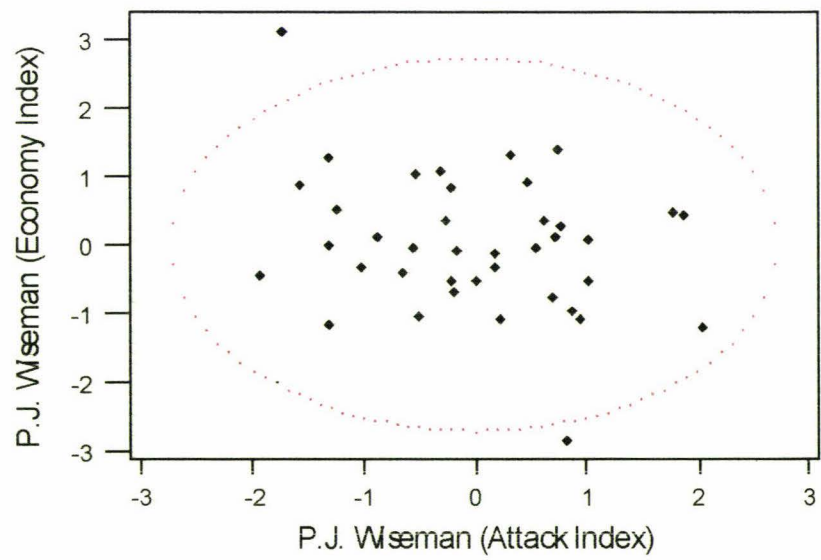


Figure 26. *Bivariate Control Chart of P.J. Wiseman for Bowling Indices*

The bivariate chart reveals the exact reasoning behind the alarms for Wiseman. Both correspond to large scores on the Economy Index. One corresponds to a highly economical performance and the other to an expensive spell of bowling.

The bivariate example above can be expanded further to include as many dimensions as deemed necessary. The following is an application to allround performance, which considers the 3 normally distributed, independent aspects of an individual's ability. A problem arises due to the non-normality of the batting variables. Thus it is inappropriate to include individual batting scores. However, as relative contributions to team performance are considered, score is redundant to this chart.

Secondly we deal with the all-round facets of an individual's play. Once again as independent, standard normal populations are being dealt with, the quality statistic Q_i is simplified to give the result below:

$$Q = (X_A \ X_E \ X_B) \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} X_A \\ X_E \\ X_B \end{pmatrix}$$

$$= (X_A)^2 + (X_E)^2 + (X_B)^2$$

Where X_A = Attack Index
 X_E = Economy Index
 X_B = Transformed Batting contribution

Next the control limits are expanded to cater for the additional characteristic to give limits at a 5% significance level of:

$$LCL = X^2_{0.97, 3} = 0.2451$$

$$UCL = X^2_{0.03, 3} = 8.9473$$

As standardised, independent values are used, the three-way chart has a spherical control region, which due to the limitations of conventional graphical displays is not shown here.

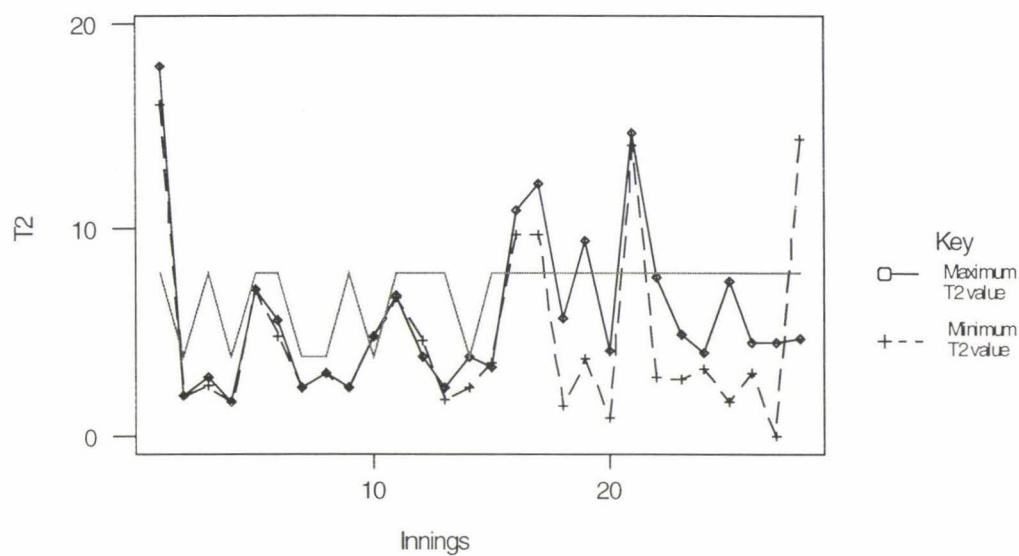


Figure 27. T^2 Control Chart of A.C. Barnes

Both maximum and minimum values are plotted simultaneously in the above graph. The maximum value is of the most use (due to the squaring of negative values – creating a positive numbers). This enables the variation in a match performance to be seen. Several alarms are signalled, but to understand the reason for the alarm and the performance aspect influencing the alarm either the raw data or a three dimensional plot need to be examined.

A problem with the T^2 Control chart is the use of data only from the current sample which means that the impact of previous values is ignored (Montgomery, 1997). This poses an interesting question in terms of how much a prior performance affects the current performance of an individual. From Chapter 2 it was found that performance was random, suggesting that there is no significant impact. This is followed up with an investigation of the MEWMA scheme.

3.6.2 MEWMA

Developed by Lowry et al (1992), the MEWMA chart is the multivariate extension of the EWMA chart.

$$Z_t = \lambda x_t + (1 - \lambda)Z_{t-1}$$

The MEWMA is defined as follows

Where, λ is once again limited to values greater than 0 but no larger than 1.

An out of control signal is given if

$$Q_t = Z_t \Sigma_z^{-1} Z_t > H \quad \text{where} \quad \Sigma_z = \left(\frac{\lambda}{2 - \lambda} \right) \Sigma$$

H is the control limit and set according to the desired ARL. In this example it is set at 5. Similar to the univariate case $\lambda=0.25$. As the values are standardised and normally distributed the control statistic, Q_t becomes:

$$Q_t = \left(\frac{2 - \lambda}{\lambda} \right) (Z_{Ai}^2 + Z_{Tb}^2 + Z_{Rb}^2)$$

Where Z_{Ai} is the i th value of the attack bowling index and so forth.

In this form it is easily manipulated using spreadsheet programmes such as EXCEL. This information is then presented in a graphical form, which is then used to interpret the inputted values.

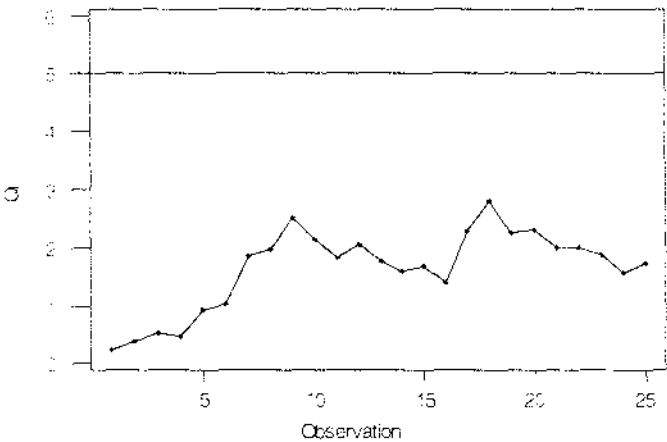


Figure 28. MEWMA Control Chart of Bowling Indices for C.M. Brown

Similar to the previous control charting methods, no alarm is signalled for Brown's performance.

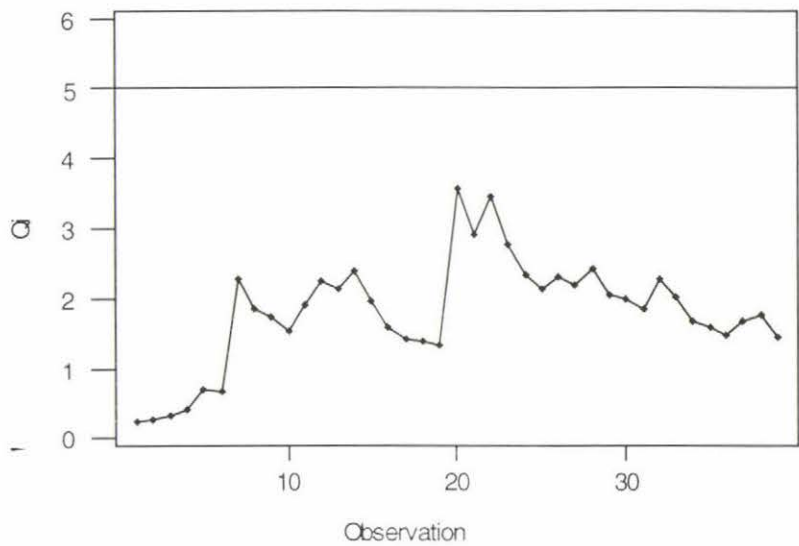


Figure 29. *MEWMA Control Chart of Bowling Indices for P.J. Wiseman*

The above chart for Wiseman indicates no change in performance parameters. As previous values are taken into account, the effect of the outliers that caused alarms previously are dampened. From a cricket point of view the nature of the alarms comes as no surprise. The role of a spin bowler can vary depending on the state of the game, nature of the wicket and so forth. Specifically, a spinner generally fills two roles, either that of a container, or an attacker. From the bivariate case, it was seen that the highly expensive innings also corresponded to an attack index score of approximately 1. From this it can be inferred that Wiseman was in an attacking capacity in that innings. A further explanation is contained in Chapter 4 with reference to homogeneity of variance. A similar argument can be applied to the innings that resulted in a highly economical index score. This implied that Wiseman played a containing role.

Further cricket debate, regarding the role that certain types of bowlers play in a match, can be applied to the alarms signalled. This further strengthens the use of control charting procedures in cricket with special regard to selection .

If there are no signals then this suggests that the estimate for the player’s ability is significantly accurate. Unfortunately due to the variable number of parameters in the all-rounder situation, the MEWMA is only suited to the bowling data.

3.7 Summary

		Number of Points to First Alarm	
Score	Method	Horne	Hartland
	Quartile	5	6
	Non-parametric EWMA	-	2
	Non-parametric CUSUM	-	-
	Within Group Ranks CUSUM	8	8
Contribution	Shewhart (Transformed)	19	7
	Shewhart (Transformed)	19	7
	CUSUM (Transformed)	19	8
	EWMA (Transformed)	-	12

Table 10. Comparison of the Number of Points to First Alarm for Horne and Hartland for Different Methods.

The above table shows the comparison of different methods for assessing batting prowess. The number of innings to the first Alarm is given. Both players appear to have runs of below average performance – Horne’s first 6 innings fall below the expected median and Hartland fails to score in three innings over a short period of six innings. As a result, alarms are expected early on. Thus the best overall performance comes from the new Quartile chart, because it gives an early signal for Horne. When alarms were signalled in the other methods, this tended to correspond with failures to score, which does not necessarily equate to a change in performance level. As charts dealing with individual values (as opposed to subgroup or weighted data) sudden shocks in the system can cause alarms, whilst not necessarily denoting performance changes.

Thus it is recommended that the new quartile chart (individual scores) is operated in conjunction with the EWMA scheme (individual contribution) for batting data. This is because one shows instantaneous results, and the other indicates evolution over time.

A similar situation applies to the bowling data. Whilst univariate procedures can be applied to the bowling indices, it is recommended that the multivariate techniques are applied. Thus the Hotelling T^2 chart along with the MEWMA are suggested for application. A bivariate chart is also necessary as a diagnostic tool.

For the all-rounder data, only the Hotelling T^2 chart is appropriate due to the potential variation of the number of parameters.

Hotelling's T^2 Statistic also makes an ideal tool for selecting the man of the match. This can be achieved by taking the highest value associated with positive parameter inputs.

Chapter 4

Ranking and Selection as Applied to Cricket

Whilst quality control procedures are implemented in this study, the interest is purely with the monitoring of an individual's performance, and consequently our estimate of the player's natural ability. This estimate is used to rank individuals so that a statistically optimal team can be selected. Obviously many factors need consideration when selecting an individual; technique, mental concentration, team dynamics and the like. As mentioned earlier, it is assumed that a player's ability is reflected in the statistics that are obtained and thus the immediate usefulness of these statistics becomes apparent.

Selecting a side solely from statistics obviously has limitations. There are human characteristics that are difficult to evaluate quantitatively, such as competitiveness, that are desirable in a team situation. We can infer the existence of hypothetical constructs involving human traits through the statistics that are expressed at the completion of a game. As with most sampling situations, our estimate of natural ability is subject to sampling error. We must also understand that as more games are played, this inherent error is reduced as a result of increased sample size. Ultimately, following the conclusion of a player's career, all that is testament to that individual's ability is the summary of performance outputs.

Having quantified the estimate of an individual's ability it is of use to rank these estimates and select the relevant number of individuals for a team. Having selected a side it is then of interest to note the probability of correct selection. We use the ranking and selection component to find the best players statistically, then use the quality control procedures described in Chapter 3 to monitor the estimate of the individual's ability.

There are a number of ways in which a statistical selection procedure could be set up, but this depends directly on the individuals available for selection and the needs of the team to be selected. Several methods exist for selection, depending on the type of distribution, what is known about the parameters of that distribution, and the number of parameters that describe a distribution. An attempt is made to keep this section as simple as possible to enable use by non-statisticians. As a result it is necessary to establish if the variance is constant amongst bowlers so that a one-stage procedure may be implemented.

For batsmen, scores are reduced to binary data, enabling the proportion of scores to be established, and hence a one-stage procedure is also applied.

Setting up this sort of a procedure has immediate appeal for choosing age group sides, where a large number of candidates need sorting through, provided the measures used relate accurately to player performance. The task of selecting paper sides for tournament teams or other result based selections, also becomes simpler.

A template for the make-up of the team to be selected is needed. This is determined by a number of factors including environmental conditions, specific personalities available and skill level of opposition.

A NZ Shell Trophy Statistical XI is chosen as an example, listed later. The selection criteria required direct participation in at least 20 innings from the duration of this study. Available players are listed in Appendix F. The selection template for this side revolved around a squad of 12; 5 batsmen, 1 batsman/keeper, 3 Fast Bowlers, 2 Spin Bowlers and 1 medium pace bowler.

Depending on the requirement of the team chosen and the variation necessitated in the bowling department, the bowling categories can be broken down further to differentiate between left arm/right arm, offspin/legspin, inswing/outswing and so forth.

If two players are equally matched, other comparative skills need to be evaluated, such as batting ability or fielding ability. This may also be done with specialist batting positions such as the top order to provide a stable batting order.

Having selected the side, each individual is monitored through the quality control procedures outlined in chapter three. Also contenders for positions in the team need monitoring in case of injury or loss of "ability" by an incumbent.

The first area investigated is that of ranking the estimates of an individual's natural ability.

4.1 Ranking the Estimate of an Individual's Ability

The first area investigated is that of ranking the estimates of an individual's natural ability. Ranking is simply ordering the data numerically; smallest to largest and assigning an ordered whole number sequentially to that series, similar to assigned placings in a running race. For population parameter θ , this is defined as follows: $\theta_1 < \theta_2 < \dots < \theta_k$ for k individuals (Mukhopadhyay, Solanky 1994). Thus, θ_k is the population parameter associated with the best individual.

The simplest method is to select the population (individual) associated with the highest mean individuals (Mukhopadhyay, Solanky 1994). However to do this we must first establish if the variance is equal between individuals (populations).

4.2 Homogeneity of Variance

The simplest method is to select the population associated with the highest mean (Mukhopadhyay, Solanky 1994). However to do this we must first establish if the variance is equal between populations, which in this situation refers to equal variance between individual players. Two different situations need to be assessed for the cases presented by batting and bowling.

a) Bowlers

As found previously in Chapter 2 the data for the bowling indices is normally distributed and random. As a result Bartlett's test statistic is used to test for homogeneity of variance. This is defined as follows:

$$B = \frac{\left(\sum v_i\right) \log \left(\sum v_i S_i^2 / \sum v_i\right) - \sum v_i \log S_i^2}{1 + \left\{ \sum (1/v_i) - 1 / \sum v_i \right\} / \{3(k-1)\}}$$

Where $S_i^2 = \sum_{j=1}^{n_i} (X_{ij} - \bar{X})^2 / (n_i - 1)$ the sum of the individual sample variance, k = number of samples and $v_i = n_i - 1$ (MINITAB, 1996).

MINITAB outputs B and also a p -value that lists the probability of rejecting the null hypothesis; that the variance is equal. If the p -value is greater than 0.05 the null hypothesis is accepted. This test is then applied to the bowling data and the results recorded in the table below.

TYPE	Single Indices		Combination Indices		
	Economy	Attack	0.5A+0.5E	0.75A+0.25E	0.25A+0.75E
ALL	0.000	0.001	0.000	0.000	0.000
FAST	0.001	0.419	0.093	0.448	0.003
SLOW	0.000	0.190	0.193	0.790	0.000
MEDIUM	0.374	0.000	0.000	0.000	0.597

Table 11. *P-Values from Test for Homogeneity of Variance for Bowlers*

Reviewing table 10 yields some fascinating results. Obviously the variances are not the same for all individuals for economy or attack ($p < 0.05$). However, breaking down the population into the specific types of bowlers defined by Smith and Payne (1994-98), an insight into the roles of these differing styles is witnessed.

For both the fast bowlers and spinners, who are generally perceived to be attacking bowlers, the variance is approximately equal for the attack index, but not for the economy index. As the job of these two general groups is to take wickets, this makes the task of selecting the best individuals easier. Reflecting on the cricketing context, when attacking, a wide range of run conceding situations occur. Often these situations are shaped by how well the bowling side has batted. Spinners are at their best when giving the ball 'flight', this creates 'loop' and 'turn', key aspects in beating a batsmen. However because it is slower through the air, batsmen can score runs relatively quickly – at a risk. Consequently the bowler is trying to deceive the batsman into playing a false stroke, leading to a dismissal. Also if the fielding side is dominating, this can see the introduction of numerous fielders clustered around the bat, as the batsmen fight to survive. However, this leaves holes elsewhere in the field that can be exploited. Many other situations are presented, due to the dynamic nature of the game, and thus presents variations in the runs conceded. It also stands to reason that the strength of the batting in a team from which a bowler comes from, will have an influence on the variation of the indices.

A similar situation occurs with the quick bowlers, with batsmen being able to use the pace of the ball to work runs. Also with attacking fields, false shots, such as edges, can be profitable for the batting team.

While these two groups look for wickets, the runs conceded varies greatly between each individual, dependent on a vast range of situations, tactics and conditions, some mentioned earlier.

Conversely, medium pacers have homogeneity of variance for the economy index, but not for attack. By their very nature, medium pace bowlers are stock bowlers, used to tie up an end, conceding as few runs as possible. Their job is to restrict run rate.

Because they are not quick, batsmen cannot use the pace of the ball to score runs, and yet because the ball is coming through quicker than a spinner, batsmen are less inclined to use their feet to get to the pitch of the ball and work it around. Due to their restrictive nature, wickets can fall through the batting side falling to frustration or the need to try and increase the run rate. Also as the pace is less than that of the quicks, the ball is in the air longer and thus has a better chance of moving in the air, especially in New Zealand conditions, creating the chance of beating the bat. However, as medium pace bowlers lack the sheer pace of the quicks or the guile of a spinner, they can be relatively easy to defend. Once again, the variation is dependent on the situation a bowler is confronted with. Some teams may find themselves in attacking situations.

Obviously when assessing a bowler we want to consider both attack and economy simultaneously. However as discussed previously, different types of bowlers have different roles. Thus we need different weightings on the indices to get an overall performance measure.

We can define a univariate performance measure (β) as a combination of the two independent indices, Attack (A) and Economy (E), using a weighting factor (ω).

$$\beta = \omega A + (1-\omega)E, \quad \text{Where } 0 \leq \omega \leq 1.$$

For Spin Bowlers and Fast bowlers ω is set at 0.75. Similarly ω is set at 0.25 for medium pacers.

Obviously the weightings are flexible depending on the game situation with regard to game plan or conditions. For a one-day type set up, ω may be set at zero, placing the entire emphasis on Economy.

b) Batsmen

With the batting measures a difficult situation arises, we have a mixed distribution. For batting two parameters define the probability of a certain score, mean score and mean contribution. In the distribution model, mean contribution inversed, denotes the probability of failure to score. As our initial interest is with what an individual is expected to score in any given innings, the two parameters must be combined to give a single expected value. Thus the expected score for an individual in an innings is given as the modified mean:

$$\text{Modified Mean} = \frac{100 - \left(\frac{100}{\bar{C}} \right)}{100} \times \bar{X}$$

Where \bar{X} is the mean of the non zero scores and \bar{C} is the mean contribution. This is necessary, as the probability of failure to score must be built into our expected value. This leaves a single value to assess. Obviously the variances are not going to be equal, but if we assume that the modified mean is approximately Geometric, then the variance is dependent on the mean. Thus we can take a single stage selection process and chose the individual with the highest mean.

Following the justification for using only the mean as an estimate for an individual's natural ability in a one-stage selection procedure, this information is combined with the ranked data to calculate the probability of correct selection.

4.3 Establishing the Probability of Correct Selection

Having established constant variance, the population (individual) with the largest mean is selected. Using a natural selection rule this is defined as

$$\theta_{jn} = \text{Max}_{1 \leq i \leq k} \theta_{in}$$

(Mukhopadhyay, Solanky, 1994)

Where θ is the population parameter for the mean and θ_{jn} is associated with the best individual. To find the probability of correct selection $P(\text{CS})$ the batting data was reduced to a binary situation. The mean expected score for the population was found and then each individual score coded depending on which side of the mean it fell (1 for above, 0 for below). This then enabled the use of proportions when evaluating the $P(\text{CS})$, based on the assumption that the sample size is sufficiently large for the normal approximation to the binomial distribution to hold.

A similar situation is applied to the bowling data, except that we can that the data comes from an independent normally distributed population with constant variance.

We consider confidence intervals for $\mu_1 - \mu$ where μ_1 is the individual mean and μ is the comparable population mean, and we look for the presence of zero in this interval to determine if there is a significant difference or not. In this example one population is the top individual, the second population is the remainder or the comparable individuals. Obviously if zero is contained there is no significant difference. We want to find the critical cut-off point at which zero is no longer included. This is effectively the largest one-sided confidence interval that does not contain zero. From this we gather the significance of the difference which enables the $P(\text{CS})$ to be determined. In a roundabout-way we are finding the most confidence we have in zero not being contained by our confidence interval.

Normally this is set up using the following equations with the significance level α , having previously been set. In this instance we are trying to find α and thus find $P(\text{CS})$. For the normal population, involved with bowling performance, we use the formula:

$$z = \frac{\bar{x}_1 - \bar{x}_2 - \delta}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \quad (\text{Smith, 1993})$$

To establishing the probability of correct selection for batsmen, the following formula is applied:

$$z = \frac{\hat{p}_1 - \hat{p}_2 - \delta}{\sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}} \quad (\text{Smith, 1993})$$

Instead of setting the critical value we use the obtained z value to solve for the significance level α , which relates directly to the $P(\text{CS})$, and relating this value to the cumulative probability density.

This is effectively creating a one-sided confidence interval, giving the probability at which the difference between the two populations becomes significant.

We are in fact testing the null hypothesis that there is no difference between populations against the alternative that there is a difference. As a result of using ordered data, a one sided test is used to test if $\theta_1 > \theta_2$. However instead of commencing with a defined level of significance, α , we solve the above equations to find α , which leads directly to $P(\text{CS})$.

As the null hypothesis is no difference, δ is zero. Thus the rearranged equation for the $P(\text{CS})$ for bowlers becomes:

$$P \left[z = \frac{\bar{x}_1 - \bar{x}_2 - \delta}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \right] = 1 - \alpha = P(\text{CS})$$

Having found z , this value is used to find the cumulative function, which enables us to solve for α . However this is not necessary as our interest is with the $P(\text{CS})$ which is $1-\alpha$. The identical situation applies to the proportion formula for evaluating the $P(\text{CS})$ for batsmen.

Having selected an individual three questions are asked.

1. Is the individual better than the population?
2. Is the individual better than the next person?
3. Is the individual better than the best individual not selected?

The probability of correct selection can easily be identified for each of these cases.

4.4 New Zealand Statistical XI 1994-98

This section selects a side chosen from eligible players from the Shell Trophy competition, given a suitable team template. Then the probability of correct selection for each of the individuals involved is evaluated, with respect to the rest of the population, the next best individual and the best non-selected individual. Finally a discussion detailing the relevant cricket thinking behind the selection finishes the chapter. The New Zealand Statistical XI for the seasons 1994-98 is given as follows:

Qualification 20 innings

M.J. Horne
C.Z. Harris
M.J. Greatbatch
C.D. McMillan
R.G. Twose
M.D. Bell
A.J. Penn
C.M. Brown
S.B. O'Connor
S.W. Duff
M.W. Priest
N.J. Astle

For this example, the selection template utilised is based around, 6 Batsmen, 3 Fast or Fast Mediums, 2 Spinners and 1 medium pacer. A list of all the players considered for the side is given in Appendix F. One will note the absence of a number of incumbent New Zealand Representatives. Due to the nature of the international programme and injuries, many of the elite players did not fulfil the requirement of direct participation in 20 innings of the Shell Trophy. Thus the players selected are the form players of the domestic competition and not necessarily the best players in the country during the given time frame. This is not a problem, instead it allows the next tier to push through. Estimates for the top players exist and the tools to monitor the maintenance of this estimate over time are also available.

A divulgence of the player's selected and the merits for their selection could be entered into here, but that is more appropriate in a sports magazine. However, it is necessary to make a few observations and justifications.

Mark Greatbatch started off his cricket as a wicketkeeper and in the Almanacks (94-98) is listed as such. Methods for selecting wicketkeepers are not covered as this is a much more involved process that requires additional information to that which is supplied in a post match analysis.

Interestingly Nathan Astle bowled in 20 innings yet did not bat in 20, an indication of Canterbury's dominance during this period, A rather handy player to be selecting purely as a bowler.

Instead we have to look at what the information tells us about the state of the game and how can this be combined with the previous methods.

To establish P(CS), a coding system based on proportions was implemented.

Batting scores are reduced to binary data by recoding the scores; 1, if greater than the population mean or 0 if below. For this data set the population mean was 26.41.

This then enabled the proportion of scores by an individual above the population mean to be assessed by the following fraction:

$$\frac{\text{Number of 1's/}}{\text{Number of innings.}}$$

This makes use of the ‘ducks and runs’ distribution that batting scores conform to. As a player’s mean increases, it is expected that more scores be above the population mean. The proportion of scores above the population mean allows for a comparison of proportions, which enables the estimation of P(CS).

As mentioned earlier, the selection template utilised for this team is based around, 6 Batsmen, 3 Fast or Fast Mediums, 2 Spinners and 1 medium pacer.

The table below defines P(CS) for each category.

- Population* compares the individual concerned with the rest of the population in the defined role.
- Next Best* compares the individual concerned with the next best ranked individual in the defined role.
- Best Non-Selected* compares the individual concerned with the best individual not selected in the defined role.

Name	Role	Population	Next Best	Best Non-Selected
M.J. Horne	BAT	92.42	95.84	92.71
C.Z. Harris	BAT	16.69	17.16	34.71
M.J. Greatbatch	BAT/KEEP	69.70	16.12	74.69
C.D. McMillian	BAT	99.42	36.66	98.39
R.G. Twose	BAT	98.84	93.17	98.07
M.D. Bell	BAT	68.57	69.98	69.98
A.J. Penn	FAST	61.42	62.51	65.16
C.M. Brown	FAST	58.15	54.54	54.54
S.B. O'Connor	FAST	57.40	50.28	50.28
S.W. Duff	SPIN	74.30	69.82	69.82
M.W. Priest	SPIN	69.00	52.99	52.99
N.J. Astle	MEDIUM	76.80	67.20	67.20

Table 12. *Probability of Correct Selection for NZ Statistical XI*

The three numerical columns detail the probability of correct selection, as a percentage, from three different comparisons.

Firstly, *Population*, gives the probability of the estimate of an individual's ability being greater than population mean. The *Next Best* lists the probability that the selected individual's ability is greater than the estimate of natural ability for the very next ranked person. Finally, *Best Non-Selected* gives the probability the selected individual's ability is greater than the estimate of natural ability for the best individual not selected in the defined role in each of those cases.

Through the use of the proportions we make the assumption that batting scores follow the distribution defined in Chapter Two. Any departure from this mixed distribution will have an impact on the results. This is evident with the results displayed for C.Z. Harris. In the 96/97 Season Harris just failed to become the first player to score three double centuries in a New Zealand Season when dismissed for 198 in the Shell Trophy Final (Payne, Smith 1997). It becomes obvious that Harris' average is heavily influenced by the presence of 3 extreme scores. However, by the very nature of the game, it is necessary to include these values. Extreme scores show an individual's ability to bat for long periods of time. This not only depends on a player's technique and other physical attributes, but also on the mental skills a player possesses. Thus while the $P(CS)$ for Harris is only 17%, we can draw inferences about his mode of play.

The table also reveals that the selection for most of the batsmen is clear-cut, with the percentages detailing the probability of correct selection for three batsmen in the 90's and two in the 70's. The bowling results are not as decisive, with percentages ranging between 50 and mid 70.

Having discussed methods for monitoring the estimate of an individual's ability and for selecting the best individuals, these procedures need to be combined in an effective and meaningful manner. Chapter 5 attempts to do this.

Chapter 5

Implementation of Statistics in Selection Methodology

After showing how cricket statistics relating to performance outputs can be used with a number of different statistical techniques in the previous chapters, this chapter demonstrates how this information can be combined to select a cricket team.

As mentioned repeatedly throughout this thesis, the main assumption rests with the belief that an individual's natural ability is expressed over time through statistics representing performance outputs. These measures are directly obtainable from the official match records and as a result extremely easy to obtain. Thus they can be applied at different levels.

Proving the previous assumption of randomness in performance outputs enabled the application of general statistical methodology. The applicability to statistical theory was further confirmed after the distributional properties of the different performance measures were established. This information then enabled the implementation of quality control procedures. Having found a method for monitoring the natural ability of an individual, ranking and selection methodologies were applied to find the statistically "best" individuals.

Building upon the knowledge generated in this study, it is very simple to put in place general guidelines for the use of performance output statistics in cricket team selection. Due to homogeneity of variance, the individual's associated with the highest means are chosen. This allows the identification and quantification of potential talent. Quality control procedures are used to establish the accuracy of this measure.

A general outline of procedure, as developed within the thesis, is detailed below:

- Estimate natural ability from performance outputs
- Check estimate with control procedures
- Rank estimates of individual ability
- Assess qualitative information
- Create team selection template
- Make selection decisions
- Establish probability of correct selection
- Combine with qualitative information
- Confirm selection
- Monitor estimate with control procedures
- Repeat process following next sampling opportunity

Due to the dynamic nature of both cricket and people the process must be continually updated, incorporating the most recent information.

Statistics are not the sole tool in forming a collective unit to play cricket. However, used correctly, sport statistics can aid the selector in the decision making process. The techniques discussed here are not confined to selection procedures, but can also act as a diagnostic tool, identifying the impact of specialist training and coaching on the estimate of ability.

The potential for the methods applied in this thesis are huge, provided the assumption between natural ability and sport statistics can be proved.

Conformation of the relationship between performance indicators and performance outputs would directly strengthen the value and merit of sport statistics in team selection.

Appendix A

Data Description

Year: Refers to the season; 1994 (1993/94), 1995 (1994/95), 1996 (1995/96), 1997 (1996/97), 1998 (1997/98).

Game: A sequential positive integer value assigned to each match. Combined with the year, enabled data checking and correction.

Province: Otago (O), Canterbury (C), Central Districts (CD), Wellington (W), Northern Districts (N), Auckland (A).

Toss: Indicates which side won the toss, 1 = win, 0 = loss.

Innings: 1 = first batting innings (Team 1 first innings)

2 = second batting innings (Team 2 first innings)

3 = third batting innings

4 = fourth batting innings

In most instances the third batting innings refers to team 1's second innings.

However, when the follow on was enforced, it becomes team 2's second innings.

In this situation, the fourth batting innings becomes team 1's second innings.

Position: Details the position occupied in the batting order. An integer value between 1 and 11 inclusive.

Name: An identity label for each player formed using the initials of their name. For example, Aaron Craig Barnes = ACB, and Campbell James Marie Furlong = CJMF. Where players shared the same initials an additional vowel from the last name of one player was added as shown here, Mark James Haslam = MJH, and Matthew Jeffrey Horne = MJHo.

Result: Indicates the final result of the match.

1 = First Innings Win

2 = First Innings Loss

3 = Outright Win following First Innings Win

4 = Outright Loss following First Innings Loss

5 = Outright Win following First Innings Loss

6 = Outright Loss following First Innings Win

0 = No result.

The following labels refer directly to batting variables. When a player was not required to bat, this was denoted in the appropriate columns as a missing value (*)

Type: Batting Style, Right Hand Bat (RH) or Left Hand Bat (LH).

Score: Runs scored by the individual. An integer value with a minimum of zero

Mins: Time spent at the crease during the individual's batting innings.

Balls: The number of deliveries faced by an individual whilst batting.

6's: Number of '6's' in the batsman's score

4's: Number of '4's' in the batsman's score

Team Score: The team total in an innings

Dismissal: The manner by which the batsman was dismissed. There are 10 potential ways to be dismissed in cricket, but not all were present in this study. Caught (C), Bowled (B), LBW Leg Before Wicket (L), Run Out (R), Stumped (S), Hit Wicket (H). The following dismissal types were not present in the time frame studied; Handled Ball, Obstructing the Field, Hit the Ball Twice and Timed Out. Also listed in this column were; Not Out (N), DNB Did Not Bat (D), Absent Hurt (A), Retired (R). For general statistical purposes, Retired is classified as Not Out and Absent Hurt is reclassified as Did Not Bat.

The following labels refer directly to bowling variables.

Innings: Refers to the innings in which a bowler operated

- 1=first bowling innings
- 2= second bowling innings
- 3 =third bowling innings
- 4 = fourth bowling innings

If the Follow On was enforced the same situation detailed in the batting variables was applied.

Type: The type of stock delivery a Bowler is rated as from one of the following.

- Right Medium (RM)
- Left Medium (LM)
- Slow Left-Arm (SLA)
- Right Leg-break (RLB)
- Right Off-break (ROB)
- Right Fast Medium (RFM)
- Right Slow Medium (RSM)
- Left Fast Medium (LFM)
- Left Slow Medium (LSM)

These rankings were obtained directly from the section entitled '*Register of Players*' from the Cricket Almanack's (Smith, Payne). Other categories apply such as Right Fast and Right Medium Fast, however Smith and Payne deemed that no one analysed during this time period were of sufficient pace.

Overs: The number of overs bowled by an individual and given as a positive rational number. Where part overs were given, these were converted to a decimal fraction. For example, a score card lists a bowler has delivered 6.2 overs, this refers to 6 completed overs and two balls out of a possible six, which is entered as 6.3333 (6 1/3). This allows for the calculation of the number of balls bowled by an individual by multiplying the new value by six.

Maidens: The number of completed overs bowled by an individual in which no runs were penalised against the bowler.

Runs: Runs conceded by the bowler.

Wickets: Number of dismissals credited to the bowler.

OppRuns: The opposition Team Total.

OppWick: Number of dismissals taken by the team

Totov: The total number of overs bowled by the bowling team.

Appendix B

Not Out: The Eleventh Form of Dismissal?

The handling of not out's is important as it impacts directly on the estimate of an individual's ability. The batting average is a useful estimate of a player's natural ability, as discussed in Chapter 1. This thesis has sought to find the best methods for monitoring the estimate of an individual's natural ability based on the expected score per innings. However the traditional average defines the expected score between dismissals. As a result the handling of Not Out's is a contentious issue in evaluating batting statistics.

The Traditional Average is deemed redundant due to the handling of Not Out scores. Discussion follows detailing the merits of considering only the Expected score per innings.

Each time a batter goes to the crease, the opposition has ten different ways of dismissing that individual. If at the completion of the innings the individual has not fallen victim to one of these methods, then that individual is deemed not out. Also if a player retires during their innings and does not return to continue the innings then they are also deemed not out.

"The usual estimate of batting average can grossly overstate a cricketer's batting ability (Danaher, p5, 1989)." This is due solely to the use of the number of dismissals as the denominator. "The underlying theory is, perhaps, that from the batsman's point of view he was 'not out' and as it were, continued his innings at the next opportunity (Wood, 1945). This thinking ignores the situations imposed by the game. In his quest for the distribution of individual batting scores Wood (1945) conceded that "I had to desert the old idea that a 'not out' innings had not been completed, which must, I think, be regarded as a pleasant fiction (p3)." As a result Wood treated a Not Out innings as a completed innings. This brings us to look at the time limitations imposed on the game of cricket.

Time Limitations

Cricket is a team game played over a specified period of time. The most popular time frame is the one-day game where 50 overs are allocated to each side for batting. First class cricket is played over a minimum of three days, with each team having the opportunity to bat twice. Limits are imposed on how the game is played by the time assigned to the match. At the commencement of each match, each side enters the fray with the specific intent of finishing the game with the most favourable result, which is generally winning. When this option is exhausted the next most favourable option is pursued and so forth. Depending on the competition being played and the advantages available for winning can influence the course of a match.

Cricket has in place certain 'boundaries' that restrict the span of a match. Specifically in a first class match the following situations impose limits:

- 1) *Specified time limit.* Shell trophy matches have a maximum of 4 days duration.
- 2) *Team limits.* Ten wickets are allocated to a batting side per innings. Once ten individuals are dismissed in a given innings, that innings is completed.
- 3) *Pursuit of winning limits.* In order to win a match sufficient time must be allocated to bowling a side out. Declarations may be required to see that the team has the best chance of a favourable result.

If we consider the case of a timeless first class match, that is no time limits are imposed, thus weather and other factors do not have the same influence. Each team will bat as long as possible, until the team is dismissed or the match has been won. Potentially one player will be left not out at the end of each innings. We come to the first limit on a batsman's score. This is imposed by the rules of the game. The opposition only need dismiss 10 of the 11 batsmen to terminate that innings. An individual cannot score runs without a batting partner. Every player who participates understands this fact. Thus it is an accepted limit imposed by the rules of the game of cricket.

Therefore the debate over what a player may have scored becomes redundant, as the scope and rules of the match do not allow the continuation of an individual innings given that the "time" boundaries have lapsed.

If we now look at the specific instances an individual may be left not out, we see that a completed not out innings is equivalent to an 'out' innings.

1) TEAM DISMISSED.

The player's opportunity to score runs has expired once the team is dismissed. There is no further chance to add to the individual total. Considering what a person may have scored is therefore irrelevant. Although the opposition did not dismiss the player they have been given the marching orders by their team-mates.

2) DECLARATION

It must not be forgotten that cricket is a team game with the ultimate goal being victory. Generally declaring the innings closed is made in the best interest of winning. In pursuit of the 'ultimate goal' it is sometimes necessary to sacrifice a player's innings. Thus the captain is responsible for this dismissal.

3) RESULT

Once a result is achieved, the job is done. There is no point considering what a batter would have scored if it does not fit with the context of the game.

4) TIME

Two factors are out of the hands of the players, time and weather. Where no result is achieved and the game is forced to end, the not out individuals have run out of time. The opportunity to score further runs has finished. It does not make sense to consider what might have been when it falls outside the scope of the match. The hands of time have effected this dismissal.

In considering all the above factors it makes sense to treat not out, as out. As all instances mean that no further addition to the individual score is possible, the thought of what might have been is irrelevant. However, as it indicates the opposition has been unable to dismiss the individual it is still important to record not outs.

This debate can be extended with a brief statistical analysis. Reducing the dismissal data to a binary situation, either not out or dismissed, enabled a One-way analysis of variance, to investigate any differences in the number of not out's by position. Assuming the approximation of Normality from the binomial distribution as n is large and np is small.

Analysis of Variance for Not Outs by Batting Position

Source	DF	SS	MS	F	P
Position	10	38.905	3.890	38.13	0.000
Error	3370	343.834	0.102		
Total	3380	382.739			

Level	N	Mean	StDev	Individual 95% CIs For Mean Based on Pooled StDev
1	350	0.0629	0.2431	(-*-)
2	350	0.0686	0.2531	(--*-)
3	336	0.0923	0.2998	(-*-)
4	326	0.0675	0.2512	(-***)
5	318	0.0629	0.2432	(-***)
6	307	0.0879	0.2837	(--*-)
7	297	0.0808	0.2730	(-***)
8	291	0.1203	0.3258	(-**-)
9	282	0.1738	0.3796	(---*-)
10	271	0.2657	0.4425	(---*-)
11	253	0.4506	0.4985	(---*-)

Pooled StDev =	0.3194	0.15	0.30	0.45
----------------	--------	------	------	------

These results reveal there is a great deal of overlap in the 95% confidence intervals for the number of not out's for the batting positions 1 to 8, inclusive.

This implies that the treatment of not out's must be the same for the first eight batsmen. As a result over time, as the relative proportions of not out's are the same, there is no significant impact on the estimation of the natural ability. Thus if the average becomes the aggregate divided by the number of innings batted in, although the final value will be less than or equal to the previous value, a result is gained that is relative to the first eight batsmen.

Furthermore, the graph below indicates the 95% confidence intervals for batting score by position and clearly shows the ranked nature of a team's batting order. Outside the non-specialist batsmen the order is clearly ranked, with player's higher in the order better than those who follow. Thus it can be implied that if number 9 is left not out when the team is dismissed, then the opposition is well capable of removing number 9, as the preceding eight are regarded as better players in the given situation and they had succumbed. The same argument can be applied to numbers 10 and 11.

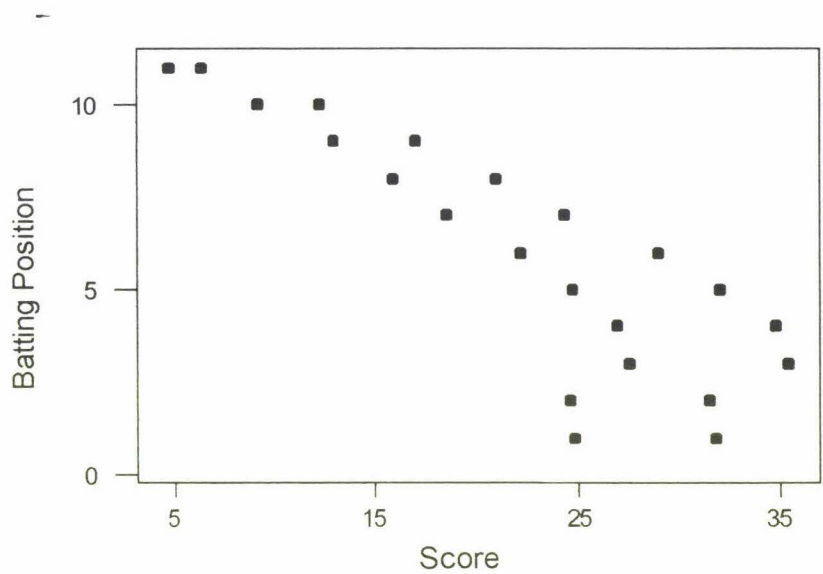


Figure 30. *Plot Showing Ranked Nature of Population Batting Order*

Conversely it could be argued that if the openers are left not out then it could be argued that the opposition are not capable of dismissing them, however, the context and time boundaries of the game must be considered.

The first section found that the distribution of scores and contribution are based on discrete forms of the exponential distribution. A special property of this distribution is the fact that it is a memory-less distribution; that is what has happened in the

past does not affect the future. In this case the present total of an individual does not affect how many runs a player goes on to score.

The proof for this is shown in Lindgren (1993, p168) and is shown by finding the tail probability of the conditional distribution and showing that it is the same as the unconditional distribution. As a direct result Burrows and Talbot add the mean to compensate a batsman's batting average for a Not Out. This may be further adapted to score, which is a mixed distribution and detailed below.

$$\begin{aligned}
 P(T > r + t | T > r) &= \frac{P(T > r + t)}{P(T > r)} \\
 &= \frac{(1 - P_0)\lambda e^{-\lambda(r+t)}}{(1 - P_0)\lambda e^{-\lambda r}} \\
 &= e^{-\lambda t}, t > 0
 \end{aligned}$$

This proof implies that the first run is the hardest to obtain. That is once a player is off the mark the probability of scoring no further runs is expressed by the standard distribution rather than the zero portion of the mixed distribution. This is obviously true as when a batsmen first gets to the crease they need to adapt themselves to the conditions, pace of the wicket, movement in the air and so forth. As a result of the distribution being memoryless the prior score does not have an effect on what may be scored, given the boundaries placed on the game. However, referring back to the limits placed upon a batting innings by the nature of cricket, it is not appropriate to add the mean score to the not out score, to correct an individual's average.

The basis of the argument presented here is that an individual's innings has been viewed in its entirety, as a direct result of the boundaries imposed by both the co-operative and competitive nature of first class cricket, hence the title of this segment; *Not Out, the eleventh form of dismissal?* Thus statistical methods relating to censored data have not been included (Smith, 1993).

Appendix C

Bowling Results

		Distribution Fitting				Tests for Randomness							
		economy	attack	Normality (p-values)		Autocorrelation		Runs Test (p-values)					
Name	Innings	mean	mean	Economy	Attack	E	A	Median	Economy	Median	Attack	Median	Bivariate
ACB	32	-0.212	-0.963	0.086	0.002			-0.58	0.7194	-0.281	0.4724	1.7000	0.1509
AJG	41	0.328	-0.299	0.383	0.754			-0.218	0.1558	0.3430	0.0403	1.3600	0.0210
AJP	27	-0.153	0.329	0.146	0.730			0.203	0.3296	-0.3970	0.8385	1.3870	0.8946
ART	36	0.460	0.072	0.447	0.472			0.353	0.0911	0.2240	0.1764	1.0710	0.7515
CJD	20	0.141	-0.147	0.458	0.627			-0.328	0.3583	0.3000	0.6460	1.8610	0.3583
CJMF	31	-0.181	-0.322	0.863	0.470	Y		-0.346	0.8503	-0.1660	0.5874	1.7590	0.8596
CMB	25	-0.339	0.285	0.008	0.393			0.246	0.5448	-0.0310	0.5448	1.2820	0.0241
CP	20	0.174	-0.081	0.780	0.032			-0.049	0.6460	0.1410	0.6460	2.3240	0.6460
CZH	26	0.409	-0.360	0.247	0.142			-0.585	0.6890	0.4890	0.1096	1.2850	0.1096
DGS	24	-0.586	0.100	0.307	0.115			0.206	0.4040	-0.7840	0.4040	2.3590	0.2082
DKM	20	0.093	0.096	0.685	0.246			0.355	0.1684	-0.0260	0.1684	1.3640	0.6460
EJM	21	-0.404	0.479	0.393	0.205			0.224	0.8142	-0.2880	0.2667	1.0710	0.8308
GEB	45	0.243	-0.227	0.170	0.146		Y	-0.228	0.1783	0.3800	0.1759	1.2530	0.0980
GPS	21	0.228	-0.029	0.027	0.902			0.073	0.3286	0.3360	0.1429	1.1350	0.4404
GRJ	47	-0.007	0.203	0.079	0.758			0.2559	0.1042	0.105	0.3005	1.342	0.4786
HTD	47	-0.989	0.197	0.246	0.864			0.149	0.8852	-0.9040	0.6604	2.7000	0.6780
JDW	31	-0.548	-0.391	0.631	0.133			-0.444	0.1987	-0.3200	0.8596	1.3600	0.4750
JTCV	37	0.257	0.251	0.002	0.164			0.076	0.8711	0.3480	0.6220	0.9060	0.1347
MHA	26	0.248	-0.034	0.523	0.716			-0.038	0.4235	0.3110	0.4235	1.8400	0.6890
MHR	25	0.112	-0.387	0.110	0.779			-0.485	0.8443	0.1760	0.8443	0.6740	0.1547
MJH	33	0.294	-0.168	0.539	0.335			-0.168	0.8637	0.3940	0.8637	1.0250	0.3734
MNH	31	0.271	-0.140	0.925	0.241			-0.238	0.2029	0.2400	0.3639	1.1610	0.5793
MWP	51	0.483	-0.031	0.660	0.781			-0.101	0.0343	0.5720	0.8853	1.0930	0.6693
NJA	20	0.638	0.019	0.059	0.566			-0.416	0.1680	0.4290	1.0000	1.7300	0.6460
PJW	39	0.427	-0.033	0.356	0.566			-0.02	0.4194	0.3960	0.4194	1.3930	0.6292
RGP	35	0.163	-0.092	0.007	0.304			-0.12	0.8676	0.3190	0.8676	1.4810	0.8676
RGT	30	0.367	-0.085	0.559	0.186			-0.332	0.0260	0.3080	0.7103	1.4030	0.2652
RJK	28	0.102	-0.159	0.177	0.378			-0.288	0.1237	-0.0400	0.7148	0.9800	0.4413
RLH	22	-0.146	0.050	0.405	0.200			0.343	0.0291	-0.0380	0.3824	1.1580	0.0291
SBO	28	-0.342	0.244	0.727	0.410		Y	0.256	0.7002	-0.3070	0.7002	1.2700	0.7002
SBS	25	-0.615	-0.061	0.240	0.894		Y	0.035	0.0646	-0.3360	0.3103	1.6350	0.8315
SJR	28	-0.862	0.338	0.702	0.332			0.375	1.0000	-0.9240	0.2482	1.6670	0.4413
SWD	23	1.109	-0.102	0.029	0.791			-0.021	0.8380	0.9430	0.8235	1.3840	0.2894
WAW	49	-0.205	0.201	0.106	0.689			0.188	0.8875	-0.3640	0.8829	1.2350	0.3108
WW	22	0.483	-0.092	0.772	0.523			-0.024	1.0000	0.4530	1.0000	0.9820	0.6623
Failures				4	2	1	3		3		1		3

Appendix D Batting Results

					Distribution Fitting							
					Score							
					excluding zero portion	MIXED EXPONENTIAL	MIXED GEOMETRIC	MIXED NEGATIVE BINOMIAL	including zero portion	GEOMETRIC	EXPONENTIAL	NEGATIVE BINOMIAL
Identity	Innings	N	N*	Mean	chi-square value a=0.05, 6df(12.592)				Mean	chi-square value a=0.05, 7df(14.067)		
ACB	45	41	4	24.1	3.0664	2.9867	3.2192	21.96	5.2366	5.1892	5.2456	
AJG	39	31	8	14.74	2.3559	2.6327	2.2357	11.72	15.3289	13.5621	13.5596	
AJP	21	17	4	12.35	4.5824	4.6442	4.7001	10	7.0192	9.1022	8.422	
ART	28	25	3	27.24	4.9278	4.8144	5.0684	24.32	8.1182	7.7147	7.9357	
ATR	40	38	2	31.76	7.596	7.5793	7.6447	30.18	7.6702	8.0092	8.0043	
BAP	23	21	2	32.52	8.7225	8.6587	8.8123	29.7	10.2837	10.7259	10.7828	
BRH	31	28	3	30.54	7.6256	7.4596	7.8203	27.58	10.6463	10.3063	10.5391	
CBG	35	28	7	34.39	2.286	2.2476	2.3476	27.51	36.9578	30.2727	31.2247	
CDC	25	24	1	30.88	9.1994	9.1294	9.2908	29.64	8.8334	9.1053	9.1482	
CDI	24	21	3	32.24	1.8625	1.8692	1.8753	28.21	8.4835	8.0077	8.1749	
COM	34	33	1	45.42	2.767	2.73	2.8335	44.09	2.728	2.7726	2.7944	
CJMF	27	25	2	17.64	4.6667	4.6907	5.2104	16.33	4.5791	4.9904	4.8757	
CJN	23	22	1	22.77	6.2882	6.2085	6.5365	21.78	5.9848	6.179	6.2097	
CMB	20	14	6	5.71	0.5889	0.3355	1.33	4	2.7591	1.7311	1.691	
CMS	44	39	5	34.08	3.9165	3.8201	4.3177	30.2	14.0553	12.8988	13.217	
CZH	23	19	4	70.5	1.4154	1.4362	1.7721	58.3	43.3387	36.0814	36.6444	
DJM	35	33	2	30.48	6.9484	6.7585	6.9826	28.74	7.038	6.9945	7.139	
DJN	20	18	2	29.78	5.0016	4.9027	5.2268	26.8	7.0593	6.9072	7.0445	
GEB	47	40	7	32.35	8.7387	8.5833	9.7994	27.53	28.5064	25.6014	26.3097	
GPB	22	21	1	36.48	5.1803	5.1177	5.3167	34.82	5.1558	5.248	5.3017	
GPS	28	26	2	24.54	10.2718	10.5894	12.172	22.79	10.5273	12.3932	11.8582	
GRJ	33	23	10	5.739	1.1013	0.9445	3.4522	4	5.0127	4.2682	3.794	
GRS	40	37	3	30.95	9.2984	9.3619	10.569	28.62	10.951	11.9108	11.8155	
HDB	23	19	4	18.26	7.0987	7.1792	10.4437	15.09	11.9657	13.6766	13.4916	
HTD	30	24	6	12.5	3.3303	3.3719	5.3713	10	8.3063	8.297	8.4425	
ISB	39	36	3	29.28	0.6861	0.6582	0.8177	27.03	2.7381	2.5823	2.6356	
JDW	44	40	4	24.2	0.8959	0.8329	1.0205	22	3.4521	3.1326	3.235	
JIP	24	19	5	18.47	3.7521	3.5483	3.899	14.63	13.4223	10.6805	11.372	
JMM	28	24	4	20.29	1.9232	1.8924	2.5048	17.39	6.6675	6.1309	6.3455	
JTCV	41	39	2	29.41	1.222	1.1586	1.2407	27.98	1.4107	1.3553	1.3818	
JVWV	20	18	2	28.33	3.6651	3.6755	4.2857	25.5	5.7029	5.9197	5.9668	
LGH	47	44	3	34.09	6.1859	6.0643	6.3301	31.91	7.6104	7.5218	7.6593	
LKG	27	23	4	31.87	5.2785	5.2081	6.0978	27.15	16.1845	14.7734	15.1593	
MAS	24	20	4	12	3.664	4.0166	6.7175	10	5.529	8.0423	7.1131	
MDBa	49	44	5	38.5	1.9642	1.9482	2.5695	34.57	12.6498	11.7973	12.0454	
MDBe	39	35	4	43.63	7.9015	7.8414	8.6337	39.15	17.778	17.0678	17.3648	
MEL	20	18	2	19.61	7.4476	7.4351	8.4803	17.65	7.6882	8.5177	8.4309	
MEP	44	38	6	35.97	4.8183	4.8007	6.4724	31.07	22.7945	21.2429	21.6512	
MGC	31	28	3	21.04	5.5253	5.7829	7.4496	19	6.8782	8.5726	8.0897	
MHA	45	41	4	33	13.0825	12.9444	14.0997	30.07	16.9171	17.2863	17.4407	
MHR	46	40	6	41.2	4.2972	4.2874	5.7207	35.83	25.0336	23.0191	23.4806	
MJG	30	28	2	60.5	6.7457	6.743	7.4551	56.5	10.7264	11.0365	11.0989	
MNH	33	30	3	25.17	3.9609	3.9221	4.4386	22.88	5.7937	5.9064	5.9615	
MSS	29	25	4	43.28	8.5011	8.3669	8.7244	37.31	23.7079	21.2414	21.7353	
MWD	23	21	2	39.76	3.5339	3.5352	4.062	36.3	6.7125	6.8057	6.8675	
MWP	43	40	3	29.4	8.5929	8.4349	8.9012	27.35	9.5083	9.5932	9.7315	
PJBC	32	29	3	25.28	9.7849	9.5664	10.2064	22.91	11.0958	11.1589	11.3536	
PJW	42	35	7	13.89	4.749	4.6106	6.5113	11.57	9.3055	9.6387	9.5781	
PWD	20	19	1	30.26	3.253	3.2098	3.3732	28.75	3.2435	3.3197	3.3468	
RAJ	48	44	4	28.34	6.7265	6.7257	7.6202	25.98	9.3618	9.7701	9.8029	
RAL	47	45	2	26.8	1.034	0.8569	0.9212	25.66	0.9211	0.7362	0.8503	
RGH	46	39	7	23.62	6.6139	6.5018	7.8435	20.02	18.7042	17.0941	17.6783	
RGP	40	33	7	30.12	5.2319	5.146	6.4209	24.85	28.5012	24.4009	25.2371	
RGT	43	40	3	42.2	4.551	4.4472	4.565	39.26	7.9654	7.6057	7.7696	
RPW	21	18	3	13.06	6.963	7.0646	8.9265	11.19	7.3938	8.6648	8.5698	
RTK	24	22	2	21.64	5.3227	5.5664	6.7899	19.83	5.8717	7.2896	6.8304	
SBO	28	21	7	14.95	6.1002	6.1325	10.8599	11.21	18.7691	18.1692	18.5524	
SBS	23	21	2	19.71	1.6683	1.6745	1.9714	18	2.1859	2.3154	2.2874	
SJR	25	22	3	14.45	9.3268	9.4241	11.2057	12.72	9.354	10.8061	10.5713	
SML	25	23	2	34.65	7.1709	7.1606	7.8924	31.88	8.8637	9.2717	9.3176	
SRM	35	32	3	25.97	7.491	7.559	8.8118	23.74	8.9886	10.0315	9.8565	
SWD	26	22	4	29.73	3.6824	3.6984	5.0742	25.15	14.3219	13.3942	13.7102	
WAW	37	31	6	23.84	3.7155	3.6499	4.9279	19.97	15.8716	14.2845	14.743	
MJH	35	23	12	8.09	0.8454	0.5296	2.3281	5.31	14.5691	8.6589	10.0006	
MJHo	24	23	1	60	8.3752	8.4852	8.4086	57.5	9.837	9.8124	9.7072	
RTK	21	11	10	6.27	1.3832	1.8793	8.2082	3.29	18.4545	14.1967	13.6155	
FAILURES					1	1	1		17	14	13	

Distribution Fitting						Tests for Randomness					
Contribution						Autocorrelation		Runs Test			
excluding zero portion	MIXED EXPONENTIAL	including zero portion	GEOMETRIC	EXPONENTIAL	NEGATIVE BINOMIAL			Score		Contribution	
Mean	6df(12.592)	Mean	chi-square value $\alpha=0.05$, 7df(14.067)			S	C	Median	p-value	Median	p-value
10.43	3.7887	9.5	4.9299	4.0569	3.833			11.00	0.1637	5.030	0.4310
6.89	10.1005	5.48	16.1099	10.8751	8.9823			6.00	0.8743	3.070	0.8743
3.584	5.0856	2.901	34.6081	12.3347	6.7206			4.00	0.2667	1.639	0.8960
11	7.966	9.82	9.882	9.2298	8.9263			21.00	0.4413	8.330	1.0000
19.33	1.7551	18.37	2.4738	2.1579	1.9589			29.50	0.5313	15.300	0.2003
10.63	7.4193	9.7	9.0858	8.4834	8.2603			11.00	0.5274	3.540	0.5274
19.61	2.4057	10.77	3.348	2.798	2.5828			13.00	0.5874	7.030	0.1987
15.6	15.5915	12.48	12.6838	13.0689	13.7025			16.00	0.0586	7.780	0.8676
15.44	3.9136	14.82	3.7639	3.7738	3.8955			28.00	0.8315	9.860	0.8443
15.86	12.8668	13.88	13.1401	13.2423	13.4755			14.50	0.6765	12.570	0.6765
14.3	7.5595	13.88	7.7329	7.6651	7.7725			34.00	0.4719	11.480	0.7277
6.74	4.7996	6.24	7.3421	5.7367	5.2255			11.00	0.5961	3.160	0.3224
9.76	1.916	9.34	2.6428	2.1769	2.0214			16.00	0.8380	6.960	0.8380
2.566	1.3187	1.796	18.13	2.5209	0.6021			3.00	0.6761	1.710	0.6460
17.31	8.8036	15.35	8.1036	8.121	8.3361		Y	16.00	0.5419	10.850	0.2226
16.64	9.4508	13.75	8.3304	8.2823	8.3809			29.00	0.8235	10.110	0.8235
9.58	18.589	9.03	22.1167	20.8488	20.2897	Y		12.00	0.2315	4.230	0.8598
12.35	2.2972	11.12	3.1102	2.7206	2.5439			18.00	0.6096	8.180	0.6460
12.57	7.3603	10.7	6.5128	6.2067	6.3754			17.00	0.4422	7.510	0.4589
16.3	3.749	15.66	4.2132	4.0659	4.0087			32.50	0.6623	13.860	0.1902
8.35	3.6755	7.76	5.3887	4.4085	4.0336			10.50	0.2482	5.190	0.7002
2.873	3.9	2.002	29.0132	5.4699	2.1231			3.00	0.3993	1.200	0.1126
11.06	5.639	10.23	5.3743	5.0483	5.1599			20.50	0.5218	6.960	0.0549
8.62	7.0216	7.12	12.2583	10.1531	9.2064			6.00	0.4514	3.670	0.0977
6.59	28.1067	5.28	6.8154	5.5665	5.9003			8.50	1.0000	3.280	0.4575
13.51	9.1123	12.48	6.4181	5.7105	5.3057			15.00	0.6292	8.810	0.2577
10.88	3.276	9.89	1.4715	0.7854	0.5949			12.50	0.1275	6.830	0.7604
8.26	5.4502	6.54	8.0061	6.4725	5.8963			9.50	0.4040	4.700	0.2108
8.93	3.2192	7.66	4.6256	3.747	3.4621			14.00	0.7002	5.370	0.7002
12.06	6.7167	11.47	7.435	7.0483	6.9714			19.00	0.6584	10.870	0.6274
13.47	6.1359	12.13	6.519	6.4078	6.4578			20.00	0.6460	7.740	0.6460
13.84	5.6406	12.96	6.5655	6.1202	5.957			22.00	0.8600	11.350	0.6559
10.93	5.7224	9.31	6.0576	5.7637	5.8196			16.00	0.0784	8.530	0.3296
5.11	4.9775	4.25	16.6595	9.6615	6.9873			6.00	1.0000	2.330	0.4040
14.56	4.4747	13.07	5.0323	4.5936	4.4821			23.00	0.3137	9.550	0.6670
17.44	5.7257	15.65	6.3916	6.1238	6.0398		Y	36.00	0.2577	10.620	0.0750
6.6	2.6305	5.94	4.8722	3.5735	3.0602			7.50	0.3583	4.490	0.3583
15.02	15.0595	12.97	22.4489	20.7833	19.5641			16.50	0.0675	7.640	0.5419
7.959	7.5901	7.189	9.4935	8.593	8.3601			14.00	0.9710	7.273	0.8596
14.65	19.8698	13.35	19.6648	19.8144	20.2906			12.00	0.4588	5.100	0.2897
15.88	11.0093	13.81	10.1375	10.1757	10.4641			25.00	0.5416	11.030	0.0713
18.42	4.7407	17.19	5.3588	5.17	5.0836			29.00	0.2652	16.210	0.0635
10.95	14.1397	9.96	16.5742	15.9585	15.7915			14.00	0.5992	6.450	0.5992
18.24	8.8185	15.73	9.3289	9.2265	9.2631		Y	14.00	0.5657	12.680	0.8550
16.2	1.9503	14.79	2.2845	2.1078	2.0485			19.00	0.1323	9.150	0.1323
9.04	6.6811	8.41	8.3257	7.3001	7.0168	Y		16.00	0.8801	4.850	0.8705
17.32	8.0115	15.7	7.3766	7.528	7.8205			8.50	0.4724	6.630	0.1509
5.89	4.5094	4.91	13.8157	7.5247	5.3763			6.50	1.0000	2.896	1.0000
17.73	4.234	16.84	4.4895	4.4125	4.4099			17.50	0.6460	12.370	0.6460
13.28	8.1137	12.17	9.1558	8.6363	8.5173			14.50	1.0000	9.630	0.5596
12.87	3.9479	12.32	3.9014	3.8	4.0096			15.00	0.9047	7.870	0.4628
8.974	9.5697	7.608	10.9102	9.9479	9.8877			16.50	0.7656	6.338	0.7656
12.27	13.1433	10.12	13.5791	13.3195	13.5029			21.50	0.7488	8.250	1.0000
16.24	6.9687	15.11	7.448	7.2275	7.2149	Y		24.00	0.1660	11.320	0.1660
7.05	4.6341	6.05	8.7205	6.5219	5.5694			7.00	0.1140	3.680	0.3847
8.26	3.9622	7.57	5.6054	4.65	4.2686			7.50	1.0000	3.400	0.4040
5.82	5.2931	4.367	10.8095	6.6885	5.4654	Y		6.50	0.0210	2.582	0.0210
10.73	1.5279	9.8	2.3574	1.867	1.6587			12.00	0.2984	6.810	0.2984
8.03	2.7152	7.07	4.6515	3.3853	2.9412			11.00	0.3103	4.230	0.8315
15.27	2.2366	14.05	3.5337	3.039	2.7095			15.00	0.8443	10.660	0.5340
9.7	3.7187	8.86	4.9747	4.175	3.9328			16.00	0.8676	6.870	0.6031
11.72	10.3502	9.92	10.3571	10.2077	10.3943		Y	15.00	1.0000	5.420	0.6890
8.09	5.0735	6.78	7.4928	5.689	5.1561			10.00	0.8711	3.830	0.6200
4.206	7.5143	2.764	36.0077	11.2436	6.7138			3.00	0.0022	0.909	0.0103
22.12	7.4013	21.2	8.117	7.6088	7.1718			25.00	1.0000	14.000	0.6765
3.58	19.6073	1.878	31.2068	10.7566	9.4896			1.00	0.7438	0.235	0.5078
	9		10	4	4	4	4		2		2

Ducks and Runs Distribution Theoretical Quartiles

[illegible]

0.5		MEAN CONTRIBUTION												
16	17	18	19	20	21	22	23	24	25	26	27	28		
													5	
													6	
													7	
													8	
													9	
													10	
													11	
													12	
													13	
													14	
													15	
													16	
													17	
													18	
													19	
													20	
													21	
14.5													22	
15.5													23	
15.5													24	
16.5	16.5												25	
16.5	17.5												26	
17.5	17.5												27	
18.5	18.5	18.5											28	
18.5	18.5	19.5											29	
19.5	19.5	19.5											30	
20.5	20.5	20.5											31	
20.5	20.5	20.5	21.5										32	
21.5	21.5	21.5	21.5										33	
21.5	22.5	22.5	22.5										34	
22.5	22.5	22.5	22.5	23.5									35	
23.5	23.5	23.5	23.5	23.5									36	
23.5	23.5	24.5	24.5	24.5									37	
24.5	24.5	24.5	24.5	24.5	25.5								38	
25.5	25.5	25.5	25.5	25.5	25.5								39	
25.5	25.5	26.5	26.5	26.5	26.5								40	
26.5	26.5	26.5	26.5	26.5	27.5	27.5							41	
26.5	27.5	27.5	27.5	27.5	27.5	27.5							42	
27.5	27.5	27.5	28.5	28.5	28.5	28.5							43	
28.5	28.5	28.5	28.5	28.5	28.5	29.5	29.5						44	
28.5	29.5	29.5	29.5	29.5	29.5	29.5	29.5						45	
29.5	29.5	29.5	29.5	30.5	30.5	30.5	30.5						46	
30.5	30.5	30.5	30.5	30.5	30.5	30.5	31.5	31.5					47	
30.5	30.5	31.5	31.5	31.5	31.5	31.5	31.5	31.5					48	
31.5	31.5	31.5	31.5	32.5	32.5	32.5	32.5	32.5					49	
32.5	32.5	32.5	32.5	32.5	32.5	32.5	33.5	33.5	33.5				50	
32.5	32.5	33.5	33.5	33.5	33.5	33.5	33.5	33.5	33.5				51	
33.5	33.5	33.5	33.5	33.5	34.5	34.5	34.5	34.5	34.5				52	
33.5	34.5	34.5	34.5	34.5	34.5	34.5	34.5	35.5	35.5	35.5			53	
34.5	34.5	34.5	35.5	35.5	35.5	35.5	35.5	35.5	35.5	35.5			54	
35.5	35.5	35.5	35.5	35.5	36.5	36.5	36.5	36.5	36.5	36.5			55	
35.5	36.5	36.5	36.5	36.5	36.5	36.5	36.5	37.5	37.5	37.5	37.5		56	
36.5	36.5	36.5	37.5	37.5	37.5	37.5	37.5	37.5	37.5	37.5	37.5		57	
37.5	37.5	37.5	37.5	37.5	37.5	38.5	38.5	38.5	38.5	38.5	38.5		58	
37.5	37.5	38.5	38.5	38.5	38.5	38.5	38.5	38.5	39.5	39.5	39.5	39.5	59	
38.5	38.5	38.5	38.5	39.5	39.5	39.5	39.5	39.5	39.5	39.5	39.5	40.5	60	

0.75	MEAN CONTRIBUTION													
16	17	18	19	20	21	22	23	24	25	26	27	28		
													5	
													6	
													7	
													8	
													9	
													10	
													11	
													12	
													13	
													14	
													15	
													16	
													17	
													18	
													19	
													20	
													21	
31.5													22	
32.5													23	
34.5													24	
35.5	35.5												25	
36.5	36.5												26	
37.5	38.5												27	
39.5	39.5	39.5											28	
40.5	40.5	40.5											29	
41.5	42.5	42.5											30	
43.5	43.5	43.5											31	
44.5	44.5	44.5	44.5										32	
45.5	46.5	46.5	46.5										33	
47.5	47.5	47.5	47.5										34	
48.5	48.5	48.5	48.5	49.5									35	
49.5	49.5	50.5	50.5	50.5									36	
51.5	51.5	51.5	51.5	51.5									37	
53.5	53.5	53.5	53.5	53.5	53.5								38	
53.5	53.5	54.5	54.5	54.5	54.5								39	
55.5	55.5	55.5	55.5	55.5	55.5								40	
56.5	56.5	56.5	56.5	57.5	57.5	57.5							41	
57.5	57.5	58.5	58.5	58.5	58.5	58.5							42	
59.5	59.5	59.5	59.5	59.5	59.5	59.5							43	
60.5	60.5	60.5	60.5	60.5	61.5	61.5	61.5						44	
61.5	61.5	62.5	62.5	62.5	62.5	62.5	62.5						45	
63.5	63.5	63.5	63.5	63.5	63.5	63.5	64.5						46	
64.5	64.5	64.5	64.5	64.5	65.5	65.5	65.5	65.5					47	
66.5	66.5	66.5	66.5	66.5	66.5	66.5	66.5	66.5					48	
66.5	67.5	67.5	67.5	67.5	67.5	67.5	68.5	68.5					49	
68.5	68.5	68.5	68.5	68.5	69.5	69.5	69.5	69.5	69.5				50	
69.5	69.5	70.5	70.5	70.5	70.5	70.5	70.5	70.5	70.5				51	
70.5	71.5	71.5	71.5	71.5	71.5	71.5	72.5	72.5	72.5				52	
72.5	72.5	72.5	72.5	72.5	73.5	73.5	73.5	73.5	73.5	73.5			53	
73.5	73.5	73.5	74.5	74.5	74.5	74.5	74.5	74.5	74.5	74.5	75.5		54	
74.5	75.5	75.5	75.5	75.5	75.5	75.5	76.5	76.5	76.5	76.5			55	
76.5	76.5	76.5	76.5	76.5	77.5	77.5	77.5	77.5	77.5	77.5	77.5	77.5	56	
77.5	77.5	77.5	78.5	78.5	78.5	78.5	78.5	78.5	78.5	78.5	79.5	79.5	57	
78.5	79.5	79.5	79.5	79.5	79.5	79.5	80.5	80.5	80.5	80.5	80.5		58	
80.5	80.5	80.5	80.5	80.5	81.5	81.5	81.5	81.5	81.5	81.5	81.5	81.5	59	
81.5	81.5	81.5	82.5	82.5	82.5	82.5	82.5	82.5	82.5	82.5	83.5	83.5	82.5	60

Appendix F

Eligible Players for NZ Statistical XI

Batsmen

M.H. Austen
 M.D. Bailey
 A.C. Barnes
 H.D. Barton
 M.D. Bell
 I.S. Bildiff
 G.E. Bradburn
 C.M. Brown
 G.P. Burnett
 P.J.B. Chandler
 M.G. Croy
 C.D. Cumming
 H.T. Davis
 P.W. Dobbs
 M.W. Douglas
 S.W. Duff
 C.J.M. Furlong
 C.B. Gaffaney
 A.J. Gale
 L.K. Gemon
 M.J. Greatbatch
 C.Z. Harris
 M.N. Hart
 R.G. Hart
 B.R. Hartland
 M.J. Haslam
 M.J. Horne
 L.G. Howell
 C.D. Inham
 G.R. Jonas
 R.A. Jones
 R.J. Kennedy
 R.T. King
 M.E. Lane
 R.A. Lawson
 S.M. Lynch
 S.R. Mather
 C.D. McMillan
 J.M. Mills
 D.J. Murray
 D.J. Nash
 C.J. Nevin
 S.B. O'Connor
 J.I. Pamment
 M.E. Parlane
 A.J. Penn
 R.G. Petrie
 B.A. Pocock
 M.W. Priest
 A.T. Reinholds

M.H. Richardson
 S.J. Roberts
 M.A. Sigley
 M.S. Sinclair
 C.M. Spearman
 G.R. Stead
 S.B. Styris
 G.P. Sulzbeger
 A.R. Tait
 R.G. Twose
 J.T.C. Vaughan
 J.D. Wells
 J.W. Wilson
 P.J. Wiseman
 W.A. Wisneski
 R.P. Wixon

Bowlers

N.J. Aistle
 M.H. Austen
 A.C. Barnes
 G.E. Bradburn
 C.M. Brown
 H.T. Davis
 C.J. Drum
 S.W. Duff
 C.J.M. Furlong
 A.J. Gale
 C.Z. Harris
 M.N. Hart
 M.J. Haslam
 R.L. Hayes
 G.R. Jonas
 R.J. Kennedy
 E.J. Marshall
 D.K. Morrison
 S.B. O'Connor
 A.J. Penn
 R.G. Petrie
 M.W. Priest
 C. Pringle
 M.H. Richardson
 S.J. Roberts
 D.G. Sewell
 S.B. Styris
 G.P. Sulzbeger
 A.R. Tait
 R.G. Twose
 J.T.C. Vaughan
 W. Watson
 J.D. Wells
 P.J. Wiseman
 W.A. Wisneski

References

- Amin, R.W, & Reynolds, M.R.Jr; Bakir, S; (1995). Non-parametric Quality Control Charts Based on the Sign Statistic. *Communications in Statistics – Theory and Methods*. 24(6); 1597.
- Anton, H. (1995). *Calculus with analytic geometry*. (5th ed.). New York: John Wiley & Sons.
- Bakir, S.T, & Reynolds, M.R.Jr; (1979). A Non-Parametric Procedure for Process Control Based on Within-Group Ranking. *Technometrics*. 21(2); 175-183.
- Benaud, R. (1995). *The appeal of cricket*. London: Hodder and Stoughton.
- Berkmann, M. (1990). *The complete guide to test cricket in the eighties*. Avon: The Bath Press.
- Bracewell, P.J. (1) (1998). 61.772: Assignment Four. Unpublished Assignment, Massey University, New Zealand.
- Bracewell, P.J. (2) (1998). Evaluating a batman's performance. Unpublished project paper, Massey University, New Zealand.
- Bracewell, P.J. (3) (1998). Evaluating a bowler's performance, the application of statistical quality control to cricket. Unpublished project paper, Massey University, New Zealand.
- Brittenden, R. (1991). *Smithy, Just a Drummer in the Band*. Auckland: Moa Beckett.
- Burrows, B.L, & Talbot, R.F; (1985). Boycott Botham and The Exponential Distribution. *Teaching Statistics*. 7; 42-48.
- Casella, G., & Berger, R.L., (1990). *Statistical inference*. California: Duxbury Press.
- Champ, C.W; Woodall, W.H; (1987). Exact results for Shewart Control Charts with Supplementary Runs Rules. *Technometrics*. 29(4); 393-399.
- Chatfield, C. (1989). *The analysis of time series, an introduction*. (4th ed). London: Chapman & Hall.
- Conover, W.J. (1971). *Practical non-parametric statistics*. New York: John Wiley & Sons.
- Cook, E. (1977). *An analysis of baseball as a game of chance by the monte carlo method*. Pp 50-54 in: Optimal Strategies in Sport, Ladany S.P.: Machol R.E. ed. Amsterdam: North Holland.
- Crosier, R.B; (1988). Multivariate Generalizations of Cumulative Sum Quality-Control Schemes. *Technometrics*. 30 (3); 291-303.

-
- Crowder, S.V. (1989). Design of Exponentially Weighted Moving Average Schemes. *Journal of Quality Technology*. 21 (3); 155-162
- Cryer, J.D; 61.342 *Introductory time series analysis, Course Textbook*. Boston: PWS-Kent Publishing Company.
- Danaher, P.J. (1989). Estimating a cricketers batting average using the product limit estimator. *New Zealand Statistician* 24: 2-5
- De Silva, B.M., & Swartz, T.B. (1997). Winning the coin toss and the home advantage in one-day international cricket matches. *The New Zealand Statistician* 32: 16-22.
- Doganaksoy, N, & Faltin, F.W; Tucker, W.T. (1991). Identification of Out of Control Quality Characteristics in a Multivariate Manufacturing Environment. *Communication in Statistics; Theory and Methods*. 20(9); 2775-2790.
- Draper, N.R., & Smith, H. (1981). *Applied regression analysis*. (2nd ed.). New York: John Wiley & Sons
- Duckworth, F., Lewis. T, (1996). *A fair method for resetting the target in interrupted one-day cricket matches*. Presented at Mathematics and Computers in Sport, Gold Coast, Australia.
- Elderton, W.P. (1945). Cricket scores and some skew correlation distribution. *Journal of the Royal Statistical Society Series A* 108: 1-11.
- Franks, B.D., Deutsch, H. (1973). *Evaluating performance in physical education*. New York: Academic Press
- Freund, J.E; (1992). *Mathematical statistics*. (5th ed). New Jersey: Prentice-Hall International, Inc.
- Gan, F.F. (1991) An Optimal Design of CUSUM Control Charts. *Journal of Quality Technology*. 23 (4); 279-286.
- Ganesalingam, S., Ganesanandam, S. & Kumar, K. (1994) *A Statistical look at cricket data*. Presented at Mathematics and Computers in Sport, Queensland, Bond University, Australia.
- Greenberg, J., & Baron, R.A. (1997). *Behaviour in organizations*. (6th ed.). New Jersey: Prentice-Hall.
- Griffith, G.K. (1996). *Statistical process control methods for long and short runs*. (2 ed.). ASQC: Wisconsin.
- Hackl, P. & Ledolter, J; (1992). A New Non-parametric Quality Control Technique. *Communications in Statistics – Simulation and Computation*. 21(2); 423-444.
- Hayter, A.J; Tsui, K-L. (1994). Identification and Quantification in Multivariate Quality Control Problems. *Journal of Quality Technology*. 26(3); 197-208.
- Hawkins, J.M. (1990). *Oxford senior dictionary*. (7th ed.). Oxford University Press: London.

-
- Kimber, A. (1993). A graphical display for comparing bowlers in cricket. *Teaching Statistics* 15: 84-86.
- Kumar, K., (1996). *Is Cricket Really by Chance?* Presented at Mathematics and Computers in Sport, Gold Coast, Australia.
- Lindgren, B.W. (1993). *Statistical theory*. (4th ed.). New York: Chapman & Hall.
- Liu, R.Y., & Tang, J. (1996). Control charts for dependent and independent measurements based on bootstrap methods. *Journal of the American Statistical Association* 91: 1694-1700.
- Lowry, C.A; Montgomery, D.C. (1995). A Review of Multivariate Control Charts. *IIE Transactions*. 27; 800-810.
- Lowry, C.A. & Woodall, W.H; Champ, C.W; Rigdon, S.E; (1992). A Multivariate Exponentially Weighted Moving Average Control Chart. *Technometrics*. 34(1); 46-53.
- McGilchrist, C.A, & Woodyer, K.D; (1975). Note on a Distribution-Free Cusum Technique. *Technometrics*. 17(3); 321-325.
- Meyer, D; (1995). Disadvantages of Moving Averages. *The New Zealand Statistician*. 30(1); 9-22.
- MINITAB (1996). *MINITAB reference manual*. Release 11. USA
- Montgomery, D.C. (1997). *Introduction to statistical quality control*. (3rd ed.). New York: John Wiley & Sons.
- Mukhopadhyay, N; Solanky, T.K.S. (1994). *Multistage selection and ranking procedures*. New York: Marcel Dekker Inc.
- Ott, L., Mendenhall, W., (1985). *Understanding Statistics*. (4th ed.). Boston: Duxbury Press.
- Park, C. & Reynolds, M.R.Jr; (1987). Non-parametric Procedures for Monitoring a Location Parameter Based on Linear Placement Statistics. *Sequential Analysis*. 6; 303-323.
- Payne, F. & Smith, I. (1998). *1998 New Zealand Cricket Almanack*. (51st ed). Hodder Moa Beckett, Auckland.
- Payne, F. & Smith, I. (1997). *The 1997 Shell Cricket Almanack of New Zealand*. (50th ed). Hodder Moa Beckett, Auckland.
- Payne, F. & Smith, I. (1994). *The 1994 Shell Cricket Almanack of New Zealand*. (47th ed). Moa Beckett, Auckland.

-
- Payne, F. & Smith, I. (1995). *Cricket Almanack of New Zealand*. (48th ed). Hodder Moa Beckett, Auckland.
- Payne, F. & Smith, I. (1996). *The 1996 Shell Cricket Almanack of New Zealand*. (49th ed). Hodder Moa Beckett, Auckland.
- Pollard, R. (1977). *Cricket and statistics*. Pp 129-130 in: Optimal Strategies in Sport, Ladany S.P.: Machol R.E. ed. Amsterdam: North Holland.
- Reep, C., Benjamin, B., & Pollard, R. (1971). Sport and the Negative Binomial Distribution. Pp 188-195 in: Optimal Strategies in Sport, Ladany S.P.: Machol R.E. ed. Amsterdam: North Holland.
- Roberts, S.W. (1959). Control Chart Tests based on Geometric Moving Averages. *Technometrics*. 1: 239-250.
- Rossini, A., (1997). *General critique of rank methods*. http://franz.stat.wisc.edu/~rossini/courses/intro-biomed/text/General_Critique_of_Rank_Methods.html
- Schweighert, W.A., (1994). *Research methods & statistics for psychology*. California: Brooks/Cole.
- Smith, N. (1994). *Kiwis declare*. Auckland: Random House.
- Smith, P.J. (1993). *Into Statistics*. Melbourne: Nelson
- Sprent, P. (1993). *Applied nonparametric statistical methods*. (2nd ed). London: Chapman & Hall.
- Stern, H.S. (1997). A Statistician Reads the Sports Pages. *Chance* 10: 19-23.
- Tracey, N.D., Young, J.C., & Mason, R.L., (1992). Multivariate control charts for individual observations. *Journal of Quality Technology* 24: 88-95.
- Turner, G., & Turner, B. (1998). *Lifting the covers*. Wellington: GP Print.
- Vasilopolous, A.V., & Stamboulis, A.P. (1978). Modification of Control Chart Limits in the Presence of Data Correlation. *Journal of Quality Technology*. 30(1); 20-30.
- Wolfowitz, J; (1943). On The Theory of Runs With Application to Quality Control. *Annals of Mathematical Statistics*. 14; 280-288.
- Wood, G.H. (1945). Cricket scores and geometric progression. *Journal of the Royal Statistical Society Series A* 108: 12-22.
- 61.325 Study Guide (1997). *Statistical methods for quality improvement*. Massey University.
- 61.725 Study Guide (1998). *Statistical quality control*. Massey University.