# Quantification of Individual Rugby Player Performance through Multivariate Analysis and Data Mining

A thesis presented for the fulfilment
of the requirements for the degree of
Doctor of Philosophy
at Massey University, Albany,
New Zealand.

Paul J. Bracewell B.Sc M.Appl.Stat(Hons)

2003

# MASSEY UNIVERSITY
## APPLICATION FOR APPROVAL OF REQUEST TO EMBARGO A THESIS
### (Pursuant to AC 98/168 (Revised 2), Approved by Academic Board 16.02.99)

Name of Candidate:    Paul J. Bracewell          ID Number:    95052126

Degree:         PhD          Dept/Institute/School:     Statistics/IIMS

Thesis Title:    Quantification of Individual Rugby Player Performance Through Multivariate Analysis

       and Data Mining

Name of Chief Supervisor:    Denny H. Meyer          Telephone Extn:    9495

As author of the above named thesis, I request that my thesis be embargoed from public access

until (date)    28/02/05          for the following reasons:

☒    Thesis contains commercially sensitive information.

☐    Thesis contains information which is personal or private and/or which was given on the basis that it not be disclosed.

☐    Immediate disclosure of thesis contents would not allow the author a reasonable opportunity to publish all or part of the thesis.

☐    Other (specify): _____

_____

Please explain here why you think this request is justified:

This thesis details the construction of a commercial rating system and discusses the business

processes that contribute to the calculation of such a rating system.  Public dissemination of this

information would remove the competitive edge the sponsoring company would otherwise possess in

the commercial sector.

_____

_____

Signed (Candidate):    _____    Date:    22/2/02

Endorsed (Chief Supervisor):    _____    Date:    22/2/02

Approved/~~Not Approved~~ (Representative of VC):    _____    Date:    24/2/02

---

*Note: Copies of this form, once approved by the representative of the Vice-Chancellor, must be bound into every copy of the thesis.*

# Abstract

This doctoral thesis examines the multivariate nature of performance to develop a contextual rating system for individual rugby players on a match-by-match basis.

The data, provided by Eagle Sports, is a summary of the physical tasks completed by the individual in a match, such as the number of tackles, metres run and number of kicks made. More than 130 variables were available for analysis. Assuming that the successful completion of observed tasks are an expression of ability enables the extraction of the latent dimensionality of the data, or key performance indicators (KPI), which are the core components of an individual's skill-set.

Multivariate techniques (factor analysis) and data mining techniques (self-organising maps and self-supervising feed-forward neural networks) are employed to reduce the dimensionality of match performance data and create KPI's. For this rating system to be meaningful, the underlying model must use suitable data, and the end model itself must be transparent, contextual and robust.

The half-moon statistic was developed to promote transparency, understanding and interpretation of dimension reduction neural networks. This novel non-parametric multivariate method is a tool for determining the strength of a relationship between input variables and a single output variable, whilst not requiring prior knowledge of the relationship between the input and output variables. This resolves the issue of transparency, which is necessary to ensure the rating system is contextual.

A hybrid methodology is developed to combine the most appropriate KPI's into a contextual, robust and transparent univariate measure for individual performance. The KPI's are collapsed to a single performance measure using an adaptation of quality control ideology where observations are compared with perfection rather than the average to suit the circumstances presented in sport.

The use of this performance rating and the underlying key performance indicators is demonstrated in a coaching setting. Individual performance is monitored with the use of control charts enabling changes in form to be identified. This enables the detection of strengths/weakness in the individual's underlying skill-set (KPI's) and skills.

This process is not restricted to rugby or sports data and is applicable in any field where a summary of multivariate data is required to understand performance.

# Acknowledgements

During the course of this thesis I have received support, encouragement and advice from my supervisors, academics, sports-people, fellow postgraduate students and business associates. I thank all these people whole-heartedly for their assistance, understanding and contributions allowing this thesis to be completed.

I am especially grateful to the input provided by my supervisor, Associate Professor Denny Meyer, who inspired my pursuit of statistics, nurtured my passion for sport statistics, encouraged me to chase my ideas and provided me with the necessary direction and support required to complete this research. I cannot speak highly enough of Denny's wonderful influence upon my research and development as a student. Additionally, her assistance with defining the theoretical variance associated with the univariate parametric half-moon statistic was most helpful.

The encouragement and guidance provided by my co-supervisor, Dr Siva Ganesh, is also greatly appreciated. Additionally, I wish to thank Professor Jeff Hunter and Dr Paul Cowpertwait for their constructive comments for improving the heuristic test of independence introduced in Chapter Five

The input of business mentor, Mr Chris Lines, of the Eagle Technology Group was also extremely valuable. His assistance was instrumental in obtaining funding via the Graduate in Industry Fellowship from Technology New Zealand. Further, Chris provided excellent support relating to the feasibility of the models created. Additional to the academic and business contribution was the beneficial assistance provided by the numerous top-level coaches, selectors and players who made themselves available for discussion and freely voiced their opinions.

I am also grateful to the support provided by my family. Some of the philosophies pursued in this thesis were sown at a young age by listening to my father and his brothers debating their theories on sports coaching and performance – all with first class playing experience, in cricket and/or rugby. Given that they are still involved in high level coaching, their ideas proved to be relevant, challenging and thought provoking.

# Table of Contents

## Part One: Introducing The Eagle Rating

**Part Two: Improving The Eagle Rating**

# List of Illustrations

# List of Tables