

Bivariate Models for the Analysis of Internal Nitrogen Use
Efficiency: Mixture Models as an Exploratory Tool

Isabel Munoz Santa (Masters of Applied Science in Biometrics)

A thesis submitted for
the degree of Masters of Applied Science in Biometrics
in the University of Adelaide

School of Agriculture Food and Wine

July 2014

Acknowledgments

I would like to express my deepest gratitude to my advisor Dr. Olena Kravchuck for her expertise, guidance, support and encouragement. This thesis would not have been possible without her. Thank you for your efforts to make me a better professional and give me the opportunity to study here!

I would also like to thank my second supervisor Dr. Petra Marschner for her guidance in the biological aspects of my thesis and Dr. Stephan Haefele who kindly provided the data for the case study of this thesis and provided support in the interpretation of the results.

I would like to acknowledge the Faculty of Science for providing the Turner Family Scholarship which supported this research and provided travel assistance to attend the Australian Statistical Conference, July 2014 Adelaide, Young Statistician Conference, February 2013 Melbourne, International Biometrics Society Conference, December 2014 Mandurah and Australian Statistical Conference in conjunction with the Institute of Mathematical Statistics Annual Meeting, July 2014 Sydney.

I would like to thank all the people at the Biometry Hub: Bev, David, Jules, Paul, Stephen and Wayne for their big smiles and for creating a positive working environment as well as the wonderful group of statisticians at Waite, meeting regularly for the professional development discussions.

Thank you to all the friends I have met in Adelaide, where I have had one of the best professional, cultural and personal experiences of my life. Thanks to Negar, Mohsen, Casey, Amanda, Rodrigo, Mariana, Fien, Daniela, Diego, Kanch, Antonio, Maria, Ruben, Alfonso, Lidia, Pablo, Ana, Antonija, Roey, Chris, Konrad, Diana, Luis and Lorinda.

Finally, with deep love and admiration, I thank my family for their love and support

from thousands of kilometres away and Martin for his support, love and contagious positive vision of life.

Contents

1	Motivations and thesis outline	1
1.1	Feeding the world requires an efficient use of nitrogen fertilisers	1
1.2	Nitrogen efficiency measures	5
1.3	Strategies for improving the uptake and utilisation of nitrogen by cereals	8
1.4	Review of grain yield and nitrogen uptake analyses in agricultural research	9
1.4.1	Studies selected for the review	9
1.4.2	Amount and format of grain yield and nitrogen uptake data	12
1.4.3	Relationship between grain yield and nitrogen uptake	12
1.4.4	Common methods of analysis of grain yield and nitrogen uptake field data and their limitations	13
1.5	Objectives of the thesis	16
1.6	Thesis outline	18
2	Ratios of jointly normal variables	21
2.1	Introduction to the ratio of jointly normal variables	21
2.2	Distribution of the ratio: history and properties	23
2.2.1	Geary (1930) and Fieller (1932) expressions of the pdf	23
2.2.2	Marsaglia (1965, 2006) expression of the pdf	24
2.2.3	Pham-Gia et al. (2006) expression of the pdf	27
2.3	Normal approximation of the pdf of the ratio	30
2.4	Estimators of the ratio	33
2.4.1	Point estimators: average of ratios, ratio of averages	33
2.4.2	Confidence sets of the ratio of expected values	35
2.5	On the distributional properties of internal nitrogen use efficiency in rice.	42

2.6	Summary	45
3	Fundamentals of finite mixture models	47
3.1	Non-technical introduction	47
3.2	Common use of mixture models	51
3.3	Mathematical definition	52
3.4	Classifying data into groups: label random vectors and posterior probabilities	53
3.5	Maximum likelihood estimation of the mixture parameters	54
3.6	The EM algorithm for the estimation of mixture parameters	55
3.7	The EM algorithm for the estimation of parameters of mixtures of multivariate Gaussian distributions	57
3.7.1	Illustration of the EM algorithm on simulated data	57
3.8	Difficulties in selecting the MLE of mixture models of Gaussian distributions with heteroscedastic components	59
3.8.1	Unboundedness of the likelihood function	60
3.8.2	Multiple local maxima	61
3.8.3	Spuriousities	62
3.9	Strategy to select the MLE of mixtures with heteroscedastic Gaussian components	64
3.9.1	Starting strategies for the EM algorithm	66
3.10	Bayesian approach to estimating parameters of mixture models of multivariate Gaussian distributions	67
3.10.1	The Gibbs sampler	69
3.10.2	The Gibbs sampler for a mixture of multivariate Gaussian distributions	69
3.10.3	Label switching problem	72
3.11	Selecting the number of mixture components	74
3.11.1	Information Criteria	74
3.11.2	Likelihood ratio test for selecting the number of clusters	76
3.12	Summary	78

4	Bivariate models for internal nitrogen use efficiency: mixture models as an exploratory tool	81
5	Conclusions and future lines of research	111
	Appendix A List of studies in the review	117
	Appendix B Journal information of the studies in the review	131
	Appendix C Equivalence between Pham-Gia et al. (2006) and (Marsaglia, 1965, 2006) expressions of the pdf of the ratio	133
	Appendix D Application of the EM algorithm for estimating the parameters of a mixture of multivariate Gaussian distributions	135
	Appendix E R code for fitting mixtures models of univariate Gaussian distributions	141
	Appendix F R code for fitting mixture models of bivariate Gaussian distributions	153

Abstract

Ratios are commonly used among plant and soil scientists, in particular to express the plant nutrient utilisation efficiency of macro- and micro-nutrients. The internal nutrient efficiency can be understood in terms of maximising yield per a unit of nutrient in the plant. At present, IE_N data are usually collected from designed field trials where different treatments are applied (e.g. fertiliser treatments) and analysed by univariate linear mixed models. However, univariate linear models on the ratio do not maintain information on the original traits, including their correlation, which presents a challenge when interpreting the effect of agronomic practices or environmental conditions on the process of nutrient conversion into grain. Moreover, the distributional properties of ratios do not comply with the assumptions of these linear models favoured in the area of soil and plant science research. A more suitable approach is to collect the traits of interest and to use bivariate analyses. These analyses preserve the information on the original traits and avoid issues associated with the ratio distributional properties.

If the data comes from field studies, different experimental and environmental conditions may lead to the presence of patterns (groups) in the data in addition or concurrently with designed treatments. Researchers in plant and soil sciences may be interested in identifying those conditions, for example to understand the nature of genotype-by-environment interactions. The inspection of the groups may reveal the factors defining them, thus gaining insight into the experimental or environmental drivers of the biological traits. Among bivariate analyses, bivariate mixture models of Gaussian distributions are an appropriate methodology for identifying clusters in the nutrient efficiency data, assuming that the traits are jointly normal. Studying this methodology for the analysis of the internal nitrogen use efficiency traits is the focus

of the present thesis.

The application of bivariate mixture models is suggested here as a complementary analysis to bivariate mixed models in designed field trials and for exploratory purposes only. The exploratory and supplementary character of the mixture analysis is due to the potential violation of the independence assumption when the data are collected from designed field trials.

In this project, bivariate mixed and mixture models are applied to a real-life designed field trial on non-irrigated rice in Thailand for the analysis of grain yield (*GY*) and plant nitrogen uptake (*NU*) data. The univariate counterparts of these analyses are also applied on the ratio of these two traits (the internal nitrogen use efficiency). The advantages of the bivariate analyses are discussed in comparison to the univariate analyses on the ratio. In this case study, the bivariate mixture approach revealed that soil water availability post-flowering and N supply in soil are the potential factors defining the mixture groups.

The present work can be readily extended to the analysis of other similar traits in agriculture when the objective is to explore potential environmental conditions affecting the traits under study. In order to fully exploit the proposed methodology, field survey is suggested as a more appropriate sampling procedure for the application of mixture models than collecting data from designed field trials.

Declaration of originality

I certify that this work contains no material which has been accepted for the award of any other degree or diploma in my name in any university or other tertiary institution and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where reference has been made in the text. In addition, I certify that no part of this work will, in the future, be used in a submission in my name for any other degree or diploma in any university or other tertiary institution without the prior approval of the University of Adelaide and where applicable, any partner institution responsible for the joint award of this degree.

I give consent to this copy of my thesis, when deposited in the University Library, being made available for loan and photocopying, subject to the provisions of the Copyright Act 1968.

I also give permission for the digital version of my thesis to be made available on the web, via the University digital research repository, the Library Search and also through web search engines, unless permission has been granted by the University to restrict access for a period of time.

Common abbreviations in this thesis

BVN	bivariate normal distribution
χ^2	chi-square distribution
CV	coefficient of variation
ρ	correlation coefficient
σ_{xy}	covariance between random variables X and Y
$Cov()$	covariance operator
cdf	cumulative distribution function
CV_x	CV of a random variable X
p	dimension of a random vector
D	Dirichlet distribution
\sim	distributed as
EM	Expectation and Maximisation
μ_x	expected value of a random variable X
$\boldsymbol{\mu}$	expected value of a random vector
$\boldsymbol{\mu}_i$	expected value of the i -th component of a mixture of multivariate normal distributions
$E()$	expected value operator
$\exp\{\}$	exponential function
F	F-distribution
Γ	gamma function
GY	grain yield
$\leftrightarrow \Leftrightarrow$	if only if
G_i	i -th group of a mixture of distributions

i.i.d.	independent and identically distributed
IE_N	internal nitrogen use efficiency
$L()$	log likelihood function
MLE	maximum likelihood estimate
π	mixing proportions
ψ	mixture parameters
$Mult$	multinomial distribution
MVN	multivariate normal distribution
NU	nitrogen uptake
g	number of groups in a mixture
τ_{ij}	posterior probabilities of the j -th observation in G_i
pdf	probability density function
\propto	proportional
tr	trace of a matrix
n	sample size
Φ	standard univariate normal cdf
φ	standard univariate normal pdf
S	straw yield
T	student's t-distribution
N	univariate normal
σ	variance
σ_x	variance of a random variable X
Σ	variance-covariance matrix
Σ_i	variance-covariance matrix of the i -th component of a multivariate normal distributions
$Var()$	variance operator
θ_i	vector of parameters of the i -th component
W	Wishart distribution

List of Tables

1.1	Causes of N losses and associated environmental impacts.	3
1.2	Main Nitrogen Use Efficiency (NUE) indices.	7
2.1	Conditions for the four scenarios of Fieller's confidence sets	38
B.1	Journals and their impact factor for studies in the review	131

List of Figures

1.1	Nitrogen pathway from soil to grain	4
1.2	Distribution of the studies in the review by continents	11
1.3	Distribution of the studies in the review by the year of publication (left) and the paper citation index (right).	12
1.4	Typical scatter plots of grain yield and nitrogen uptake in wheat	13
2.1	Different shapes of the probability density function of the ratio of two jointly normal variables	27
2.2	Probability density function of the ratio of two jointly normal variables for different values of the coefficient of variation of the denominator (CV_x)	32
2.3	Probability density function of the ratio of two jointly normal variables for different values of the coefficient of variation of the numerator (CV_y) but same CV_x	32
2.4	Bootstrap distribution of the estimators defined in Eq. 2.11 (left) and Eq. 2.12 (right)	33
2.5	Feasible cases of Fieller's confidence set of the ratio of expected values .	37
2.6	Construction of a wedge given the confidence interval of the ratio of means (left) and vice versa (right)	40
2.7	Confidence sets of the ratio of expected values in Von Luxburg & Franz (2009).	41
2.8	Grain yield versus nitrogen uptake from a sample of non-irrigated rice in northeast Thailand (Naklang et al., 2006)	43
2.9	Probability density function of the ratio of two jointly normal variables with parameters given in Eq. 2.21	44

2.10	Bootstrap distribution of the estimators defined in Eq. 2.11 (left) and Eq. 2.12 (right) for the data in Fig. 2.8	44
3.1	Solutions of three clusters obtained by the EM algorithm for the case study (Chapter 4)	50
3.2	Spurious solution obtained by the EM algorithm when fitting 7 components to the case study data (Chapter 4).	50
3.3	Bimodal distribution generated from a mixture of three univariate normal components	52
3.4	Scatter plot of a sample from a mixture two bivariate normal	58
E.1	Histogram of the internal nitrogen use efficiency data and the mixture components found by the EM algorithm initiated from random starts .	144
E.2	Histogram of the internal nitrogen use efficiency data and the mixture components found by the EM algorithm initiated from a partition provided by the K-means algorithm	145
E.3	Histogram of the internal nitrogen use efficiency data and the mixture components found by the EM algorithm initiated on a random subsample	147
E.4	Histogram of the internal nitrogen use efficiency data and the mixture components found after running several short runs of the EM algorithm	149
F.1	Cluster partition found by the EM algorithm initiated from random starts	159
F.2	Cluster partition found by the EM algorithm initiated from the partition obtained by the K-means algorithm	162
F.3	Cluster partition found by the EM algorithm initiated from simulated means	165
F.4	Cluster partition found by the EM algorithm initiated from the mixture estimates obtained by running the EM algorithm on a random subsample of 200 observations	170
F.5	Cluster partition found by the EM algorithm initiated from the mixture estimates after running several short runs of the EM algorithm	179