

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

Predicting albumen gland length of grain crop pest snails

Kathy Ruggiero

Department of Statistics
The University of Auckland

28 November 2023

Email: k.ruggiero@auckland.ac.nz

Outline

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- 1 The problem
- 2 Reproductive activity of snails
- 3 Study data
- 4 Some findings
- 5 Objectives
- 6 Methods and results

Controlling snail populations

Requires a combination of methods:

- Cultural
 - 🍷 Field hygiene, e.g. weed control and removal of refuges
- Biological
 - 🍷 Predators, e.g. beetles, lizards, and birds (ducks, chickens or guinea fowl)
- Chemical
 - 🍷 Baits

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

Baiting efficacy

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- Limited field life
 - 🍷 Rainfall affects physical integrity and dilutes concentration of active ingredients
 - 🍷 Temperature-related degradation
- Requires re-application every 2–4 weeks
- Timing is critical
 - 🍷 Rule of thumb: Bait in the autumn when snails are actively feeding and prior to egg laying
 - 🍷 Idea: Preventing egg laying to mitigate risks of harvest contamination by juveniles (due to their small size)

Baiting efficacy

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

Can efficacy be improved by optimising timing of bait application?

What triggers reproduction activity?

Snail reproduction activity

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- Hermaphrodites but “prefer” to mate with another
- Mature in 1-2 years
- Aestivate through summer, generally dormant
- Rain triggers snail activity, temporary feeding
- Autumn rains prompt feeding, mating, and egg-laying
- Lay up to 6 batches of 80 eggs each, hatching after 2 weeks

Snail reproduction activity

- Albumen gland swells when reproductively active

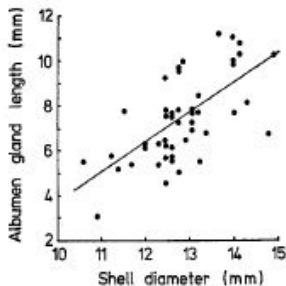


Fig. 7. Length of albumen gland as a function of shell diameter for *Cernuella virgata* collected in the pasture at Mt Benson during March 1985. Equation of the regression line is: $y = 1.32x - 9.37$, $r_{30} = 0.681$, $P < 0.05$.

Image source: Baker, G. H. (1988). The life-history, population-dynamics and polymorphism of *Cernuella virgata* (Mollusca, Helicidae). *Australian Journal of Zoology*, 36, 497-512.

Study: Biology and ecology of pest snails

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- GRDC¹-funded project conducted by SARDI²
- Aimed at understanding pest snails in agricultural regions of southern and western Australia
- Four snail species and two slug species
- Several sites across southern and western Australia
- The rest of this talk will focus on a single snail species from a single site
- The methods are applicable to all four species across all sites

¹Grain Research and Development Corporation

²South Australian Research and Development Institute

The data

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

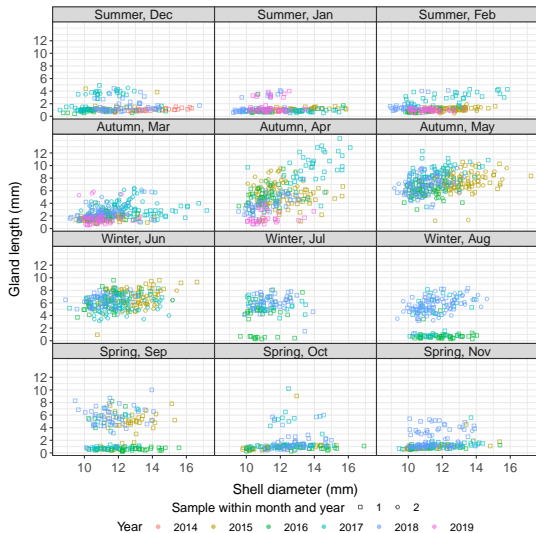
Objectives

Methods and
results

- Quota sampling
- Approximately monthly samples gathered for size and reproductive trait analysis
- Longest and most comprehensive dataset:
 - 📎 Dec 2014 to Apr 2019, covering 64 sampling occasions
 - 📎 Sample sizes ranged from 17 to 45 snails per sample
 - 📎 Total 2498 animals during the study period
- Daily Australian Bureau of Meteorology data
- Micro-climate data logger (30-minute intervals)

SARDI snail data: Single species, single site

December 2014 to April 2019



Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

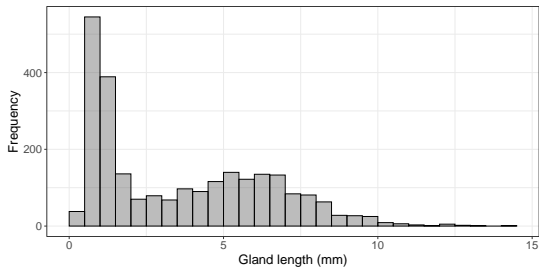
Some
findings

Objectives

Methods and
results

SARDI snail data: Single species, single site

Distribution of albumen gland length



- Gland lengths primarily cluster around two modes: around 1.0 mm and 5.8 mm.
- Bimodality plausibly explained by two reproductive states:
 - 📌 Higher mode (swollen glands) indicates reproductive activity (State A).
 - 📌 Lower mode indicates reproductive inactivity (State I).

Snail Data Analysis: Gland Length Modeling

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- Baker (1988) modelled gland length as a function of shell diameter
- SARDI data, exhibiting bimodal gland length distribution, indicates we need a model accommodating two reproductive states (A and I)
- Unsupervised clustering required due to the lack of state labels

Objectives

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- Separate snails into reproductively active and inactive states (today)
- Identify environmental (climate) variables which trigger reproductive activity

Unsupervised clustering of reproductive state

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

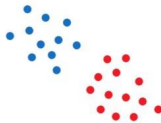
Some
findings

Objectives

Methods and
results

k-Means Clustering

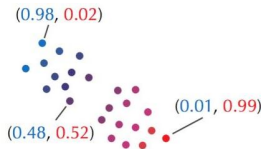
- Each snail assigned to one cluster (hard clustering)



- Tends to produce spherical and equally sized clusters
- Sensitive to outliers and initial centroids

Gaussian mixture models

- Assigns probabilities to cluster membership (soft clustering)



- Accommodates varying cluster shapes and sizes

Log-transformed data characteristics

Snail gland
length
prediction

Kathy
Ruggiero

The problem

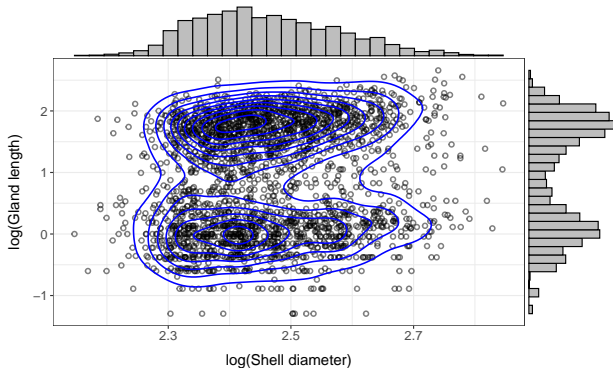
Reproductive
activity

Study data

Some
findings

Objectives

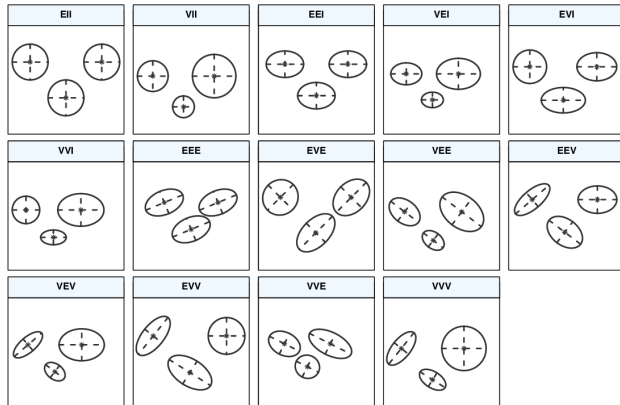
Methods and
results



- Mixture of bivariate normals
- Ellipsoidal covariance matrix; variable orientation along coordinate axes
- Moderate (State A) and weak (State I) correlation between variables

Multi-dimensional Gaussian Mixture Model

Parametrizations of covariance matrices³ (volume, shape, and orientation)



³Implemented in the Mclust (version 5) R package

Model selection

Multi-dimensional Gaussian Mixture Model

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- How many components should be included in the mixture?
- Which covariance matrix should be adopted?

Information criteria for model selection

Multi-dimensional Gaussian Mixture Model

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

Bayesian Information Criterion (BIC)

- Penalises model complexity using the number of parameters and sample size
- Tends to favour simpler models and may overlook certain complex structures

Integrated Complete-data Likelihood (ICL) criterion

- Penalises BIC by incorporating an entropy term which quantifies the overlap of observations between clusters
- Tends to favour solutions with clearly separated clusters

Information criteria for model selection

Multi-dimensional Gaussian Mixture Model

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

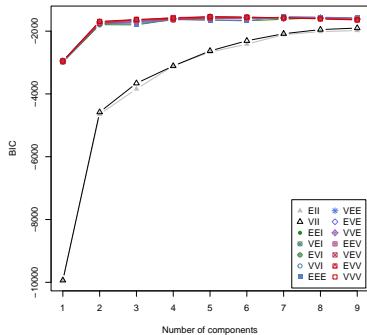
Study data

Some
findings

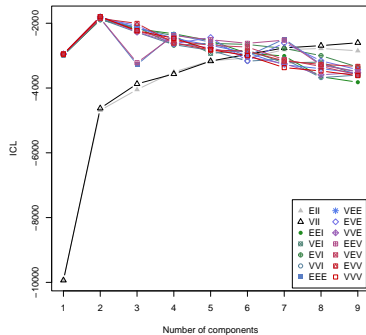
Objectives

Methods and
results

BIC



ICL



2-Cluster uncertainty plot

Multi-dimensional Gaussian Mixture Model

Snail gland
length
prediction

Kathy
Ruggiero

The problem

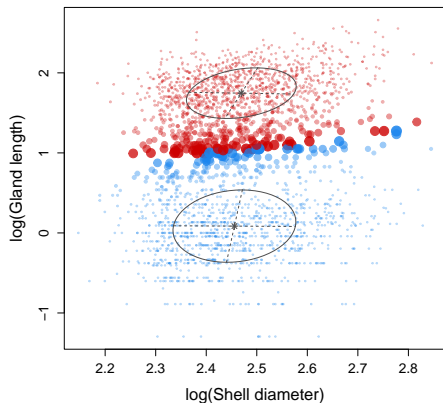
Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results



- Uncertainty is given by size of point
- 94.7% of cases have probability >0.9 of belonging to the assigned cluster

```
null device
```

```
1
```

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

Cluster reliability

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- Cannot assess cluster *accuracy* (correctness of cluster assignment) because true reproductive state is unknown
- Can use *bagging*⁴ to assess cluster consistency, i.e. sensitivity to small changes in the input data

⁴BagClust2, Dudoit and Fridlyand (2003) *Bioinformatics*, 19

Bag clustering 2 algorithm

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- 1 Initialise $A = [a_{ij}]$ and $M = [m_{ij}]$ matrices to zeros.
 - 📌 A records concurrences of observations in the same cluster across bootstrap samples
 - 📌 M records total occurrences of observations in the same bootstrap sample
- 2 Form the b th bootstrap sample $L_b = (x_{b1}, \dots, x_{bn})$.
- 3 Apply clustering procedure P to L_b and obtain cluster labels $P(x_{bi}; L_b)$.
- 4 Update matrices A and M for each pair of observations based on cluster concurrence.
 - 📌 $a_{ij} \leftarrow a_{ij} + I[x_i \in L_b, x_j \in L_b, P(x_i; L_b) = P(x_j; L_b)]$
 - 📌 $m_{ij} \leftarrow m_{ij} + I[x_i \in L_b, x_j \in L_b]$
- 5 Repeat Steps 2–4 B times and compute dissimilarity matrix $D = [d_{ij}]$, where $d_{ij} = 1 - \frac{a_{ij}}{m_{ij}}$
- 6 Cluster the n original observations based on the dissimilarity matrix D

Cluster reliability

Snail gland
length
prediction

Kathy
Ruggiero

The problem

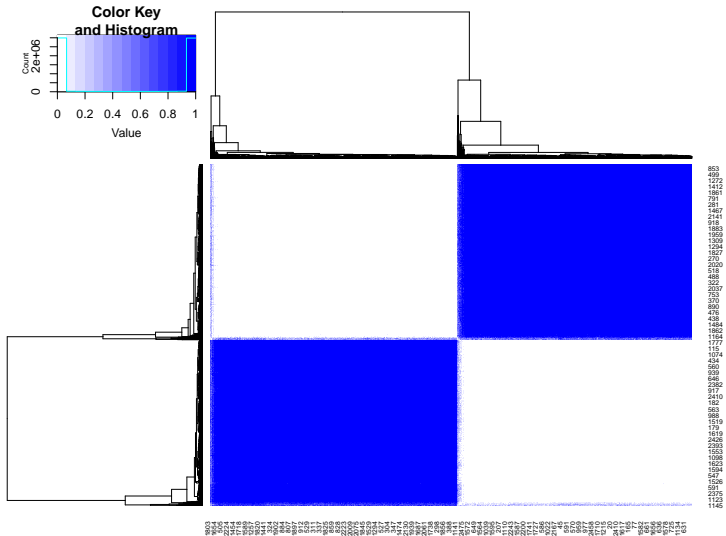
Reproductive
activity

Study data

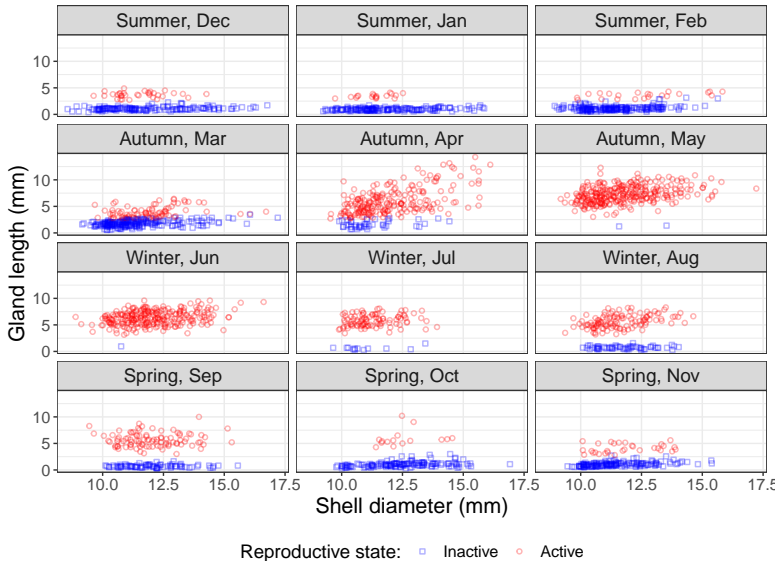
Some
findings

Objectives

Methods and
results



GMM clustering results



Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

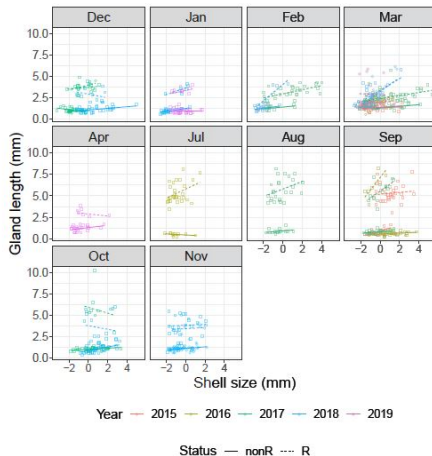
Some
findings

Objectives

Methods and
results

Next steps

- Predict gland length for “standard” sized snail by sample (month/year) and state



Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

Next steps

Snail gland
length
prediction

Kathy
Ruggiero

The problem

Reproductive
activity

Study data

Some
findings

Objectives

Methods and
results

- Correlate gland length with BOM and micro-climate data
 - 👉 Tree-based approach: binary (reproductive state) vs continuous ($\log(\text{GL}/\text{SD})$) response
- Single-step method, using a Bayesian approach?