

Riffle: an R Package for Nonmetric Clustering

Geoffrey B. Matthews and Robin A. Matthews

Western Washington University

Bellingham, WA, USA



- Dissimilar data types
 - Chemical
 - * ph, alkalinity
 - Physical
 - * temperature, percent canopy cover, sediment size, land use classes
 - Biological
 - * chlorophyll, sex (male, female, juvenile)
 - * rare species (counts 1-2)
 - * common species (counts 10,000-100,000)

Problems for Multivariate Data Analysis

- Censored data.
 - Tied ranks and reduced variance when “<5” \Rightarrow “5”.
 - Systematic bias when omitted.
- Missing data.
 - Omit entire row when one variable column is missing?
- Noisy, “useless” parameters.
 - Measured anyway.
 - Can be unrelated to major patterns.

Riffle: an R Package for Nonmetric Clustering



Riffle

Matthews & Hearne, *IEEE PAMI*, 1991

A clustering algorithm:

- group similar points into clusters.

A nonmetric algorithm:

- uses only order statistics for continuous data
- can handle both continuous and categorical data together

Uses variables independently:

- ignores scattered missing values
- uses incommensurable variables without normalizing



Proportional Reduction in Error

- Measuring Predictability for Categorical Variables

	red	green	blue	Errors
A	5	8	2	7
B	2	3	9	5
C	8	1	0	1
Totals	15	12	11	13

Errors predicting (red, green, blue) *a priori*: $12 + 11 = 23$

Errors predicting (red, green, blue) *given* (A, B, C): $7 + 5 + 1 = 13$

Proportional reduction in error: $\frac{23-13}{23} = \frac{10}{23}$

- More meaningful and robust than, e.g., χ^2

Riffler: an R Package for Nonmetric Clustering



Proportional Reduction in Error

Independent variables:

	red	green	blue	Errors
A	6	3	9	9
B	4	2	6	6
C	2	1	3	3
Totals	12	6	18	18

Minimum:
0% reduction

$\frac{0}{18}$

Perfectly predictable variables:

	red	green	blue	Errors
A	12	0	0	0
B	0	0	18	0
C	0	6	0	0
Totals	12	6	18	0

Maximum:
100% reduction

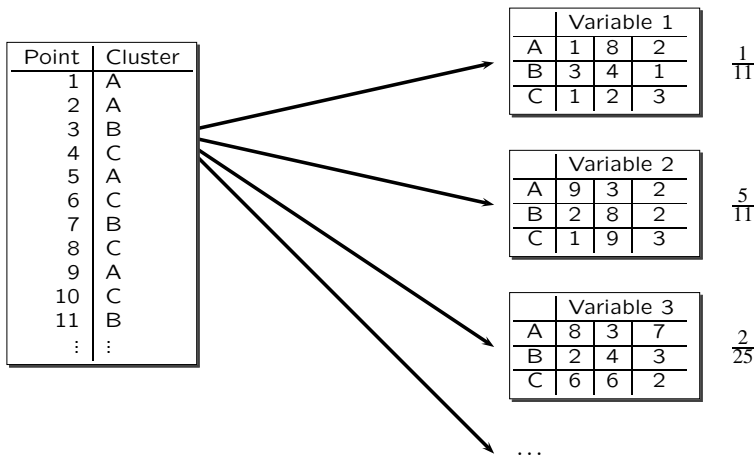
$\frac{18}{18}$

Riffler: an R Package for Nonmetric Clustering



Clustering with categorical variables

- Assign clusters to maximize predictability over other variables.

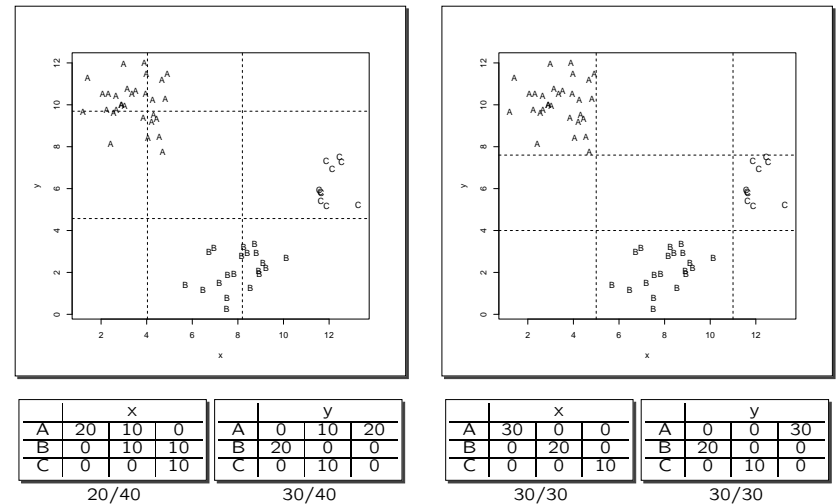


Riffler: an R Package for Nonmetric Clustering



Handling ordered variables

- Cuts adjusted to maximize predictability of clusters

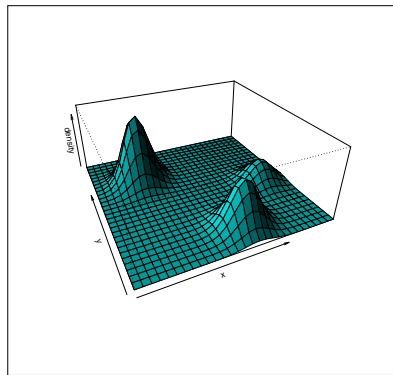
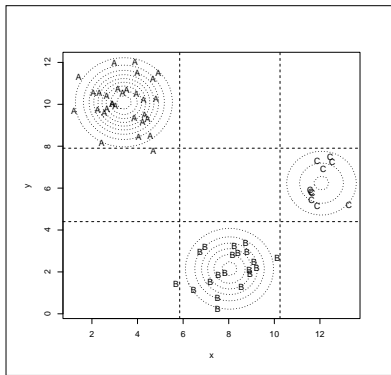


Riffler: an R Package for Nonmetric Clustering



Cutting Gaussian variables

- Generate independent Gaussians from cluster statistics μ_i, σ_i
- Cut where max likelihood changes from one to another.



	x			y		
A	30	0	0	A	0	1
B	1	19	0	B	20	0
C	0	0	10	C	0	10

28/29

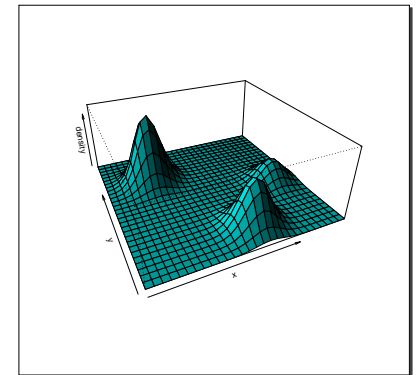
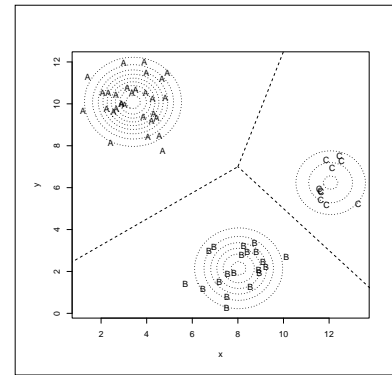
30/31

Riffle: an R Package for Nonmetric Clustering



Alternative handling of Gaussian variables (EM)

- Assign to most likely group, instead of max predictability.
- Not used in Riffle.



Riffle: an R Package for Nonmetric Clustering



Essential Algorithm

```

variables <- quantile.cuts(data)
clusters <- seed.clusters(variables)
score <- reduction.in.error(variables, clusters)

while (improving(score)) {
  variables <- best.cuts(variables, clusters)
  clusters <- best.clusters(variables, clusters)
  score <- reduction.in.error(variables, clusters)
}

return (clusters, variables)

```

Riffle: an R Package for Nonmetric Clustering



Getting things started

To find initial cuts for variables:

- Use quantiles for cut points.
- Use quantiles for μ_i , overall σ for σ_i .

To find initial clusters, given cut variables:

- Select one point randomly as seed.
- Find other seeds by selecting points as different as possible.
- Assign each seed to a different cluster.
- Assign all other points to cluster of most similar seed.

Riffle: an R Package for Nonmetric Clustering



Embellishments

- Each variable is dealt with independently.
- Each variable has a score (predictability vs. cluster).
- Use score to eliminate variables, or rank them in importance.
- We use this to handle the curse of dimensionality and find a small set of critical variables.
- Data reduction

Riffle: an R Package for Nonmetric Clustering



Data Exploration vs. Confirmation

- Clustering in general is exploratory.
- Clustering data with known groups:
 - correlation between clusters and groups measures significance.
 - identifies important variables as the ones with high predictability.
 - determine not only significance of effect, but also which variables are affected the most.
 - we have used this to chart seasonal effects.

Riffle: an R Package for Nonmetric Clustering



Conclusion

- We have used **Riffle** successfully for over 10 years for ecological and toxicological data analysis.
- **Riffle** can cluster using incommensurate variables.
- **Riffle** handles censored data and missing data with few assumptions.
- **Riffle** can reduce complexity in highly multivariate datasets.
- **R** package available 2006.

Riffle: an R Package for Nonmetric Clustering

