

# Mixture Models for Multi-Species and Environmental Data

<sup>1,2</sup>Scott Foster, <sup>1,3</sup>Piers Dunstan, <sup>4</sup>David Warton, <sup>2,4</sup>Francis Hui, <sup>2</sup>Ross Darnell, <sup>5</sup>Geof Givens, <sup>6</sup>Grant Dornan

<sup>1</sup>CSIRO Wealth from Oceans Flagship, <sup>2</sup>CSIRO Division of Computational Informatics, <sup>3</sup>CSIRO Division of Marine and Atmospheric Science, <sup>4</sup>University of New South Wales, School of Mathematics and Statistics, <sup>5</sup>Colorado State University, Department of Statistics, <sup>6</sup>Steadman Philippon Research Institute

- Ecological inference and management decisions often depend on data from many species.
- A proper and useful statistical analysis quantifies the important patterns of variation, whilst reducing the complexity in multi-species data.
- Currently, analysis is frequently done by: 1) performing species-by-species analyses (e.g. univariate regression and extensions) and then combining results, or 2) by combining data (clustering) and then performing a group-by-group analysis.
- Neither of the standard approaches are entirely satisfactory as important aspects of the variance in the data can be lost when moving from step to step. Also, the propagation of uncertainty is difficult and is subsequently (often) ignored.
- We introduce two models, based on mixture models, that address these issues. One model type, *species archetype models (SAMs)* exploits similarities in individual species' responses to the environment. The second type, *regions of common profile (RCP) models*, exploits similarities in the assemblage patterns at each site.

## Models

We use 2 variants of mixture models to represent variation in data from many species. Let the species' data be given by  $\{y_{ij}\}$  and the environmental data be  $\{w_i, x_i\}$ , where  $i = 1 \dots N$  index sites and  $j = 1 \dots S$  index species. We have used these methods for  $S \approx 300$  species and  $N \approx 1200$  - limited by size of survey data at the moment.

### Species Archetype Models (SAMs, for when inference on species is required)

Mixture of regressions model to group individual species' responses to environmental gradients into *archetypal* responses. Only  $K \ll S$  archetypal responses are interpreted, instead of  $S$  models from a species-by-species analysis. The model for the expectation is:

$$E[y_{ij}] = \sum_{k=1}^K \pi_k E[y_{ij} | \text{archetype } k], \quad \text{where}$$

$$g(E[y_{ij} | \text{archetype } k]) = \alpha_i + \mathbf{w}_i^\top \boldsymbol{\tau}_j + \mathbf{x}_i^\top \boldsymbol{\beta}_k$$

and the distribution of each species' data may have species-specific parameters [1, 2, 3, 4, 5].

See Figure 1 for some archetypal environmental responses of fish off the south-eastern coast of mainland Australia. Analysis performed using physical environmental covariates, biomass data and a Tweedie model.

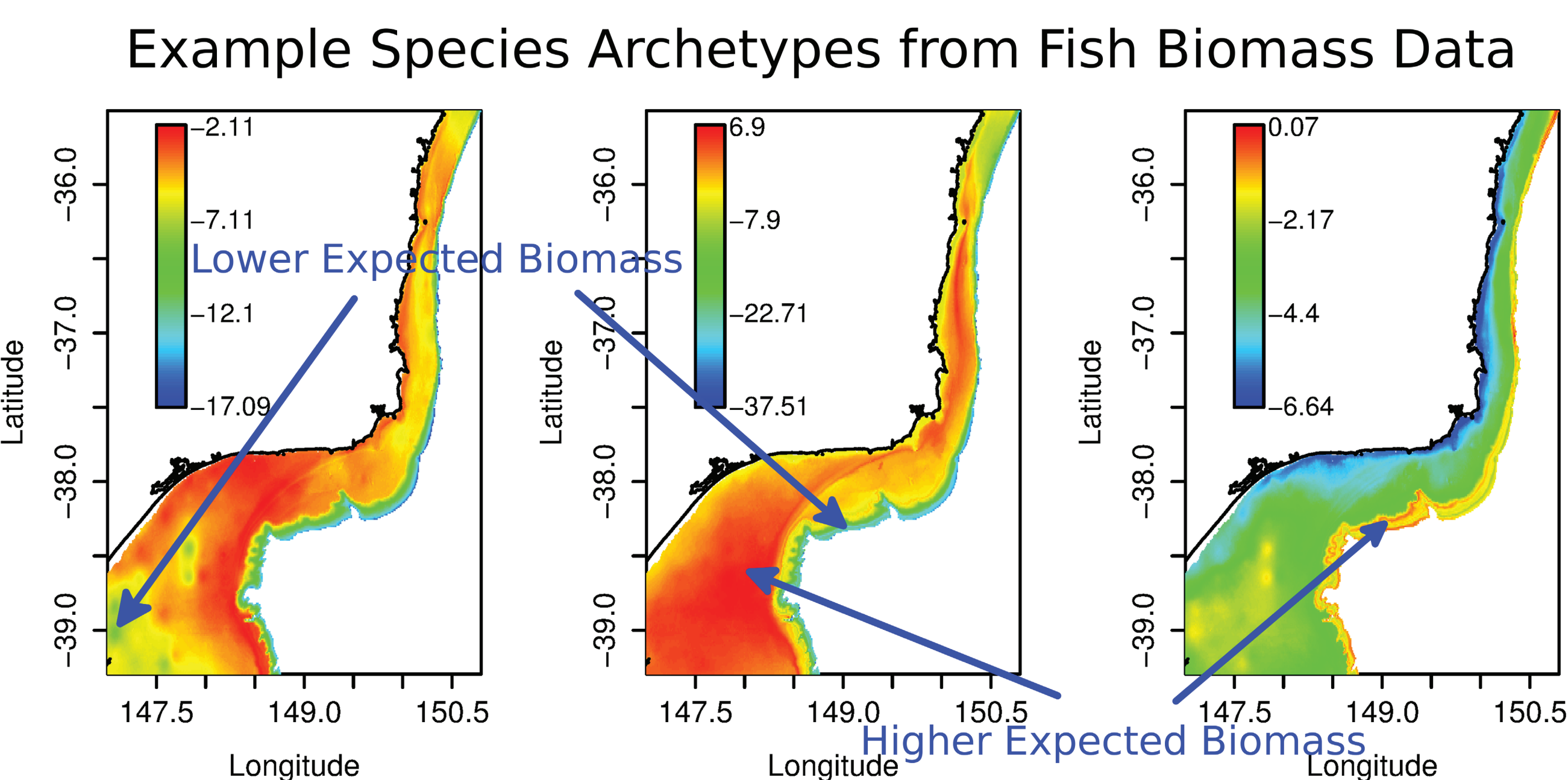


Figure 1: Example results from a SAM analysis. Maps are for expected biomass, on the log scale.

### Regions of Common Profiles (RCP, for when inference on sites is required)

Mixture of experts model to group sites based on its species profile in relation to environment. This allows interpretation of only  $H \ll N$  RCPs to study assemblage-environment relationships. No species, nor site, is hard-clustered into a particular RCP. The model for expectation is:

$$E[y_{ij}] = \sum_{h=1}^H \pi_h(\mathbf{w}_i, \mathbf{x}_i) E[y_{ij} | \text{RCP } h],$$

where  $\pi_h(\mathbf{w}_i, \mathbf{x}_i)$  is a multinomial logit link function [6].

See Figure 2 for the predicted probability of each region-type, based on fish assemblages for the North West Shelf of Australia. Analysis performed using physical environmental covariates and fish presence/absence data.

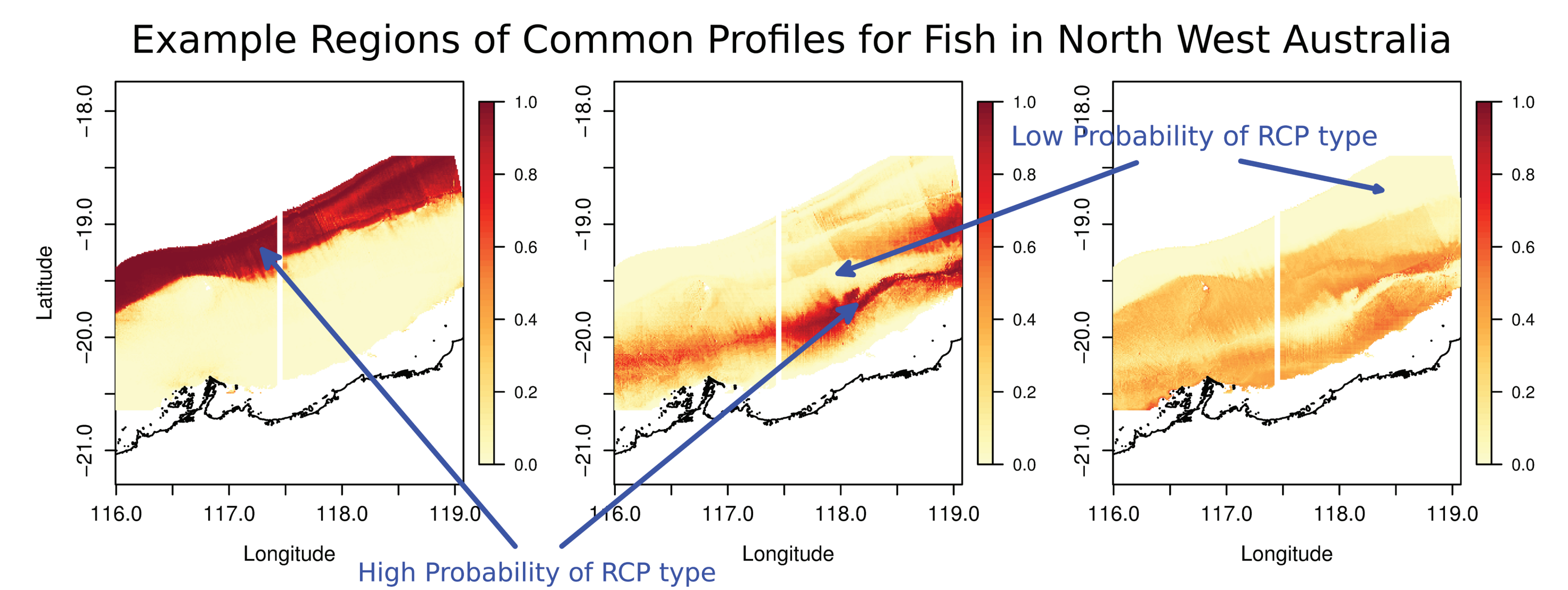


Figure 2: Example results from an RCP analysis. Maps are for the probability of each region type occurring at each location.

## Benefits of Using SAMs and RCPs

**Interpretability** - complexity of multi-species data reduced by clustering species (SAMs) or sites (RCPs) [1, 6, 4].

**Prediction Performance** - rarer species 'borrow strength' from common species [3].

**Diagnostics** - checking model adequacy, e.g. residual plots [2, 5].

**Model Selection** - using variants of common methods [7, 8].

**Flexibility** - many choices to match model to data [2, 5].

**Transparency** - model is formally specified.

**Efficiency** - proper statistical inference [1, 7].

## References

- [1] P.K. Dunstan, S.D. Foster, and R. Darnell. Model based grouping of species across environmental gradients. *Ecological Modelling*, 222:955 - 963, 2011.
- [2] P.K. Dunstan, S.D. Foster, F.K.C. Hui, and D.I. Warton. Finite mixture of regression modeling for high-dimensional count and biomass data in ecology. *Journal of Agricultural, Biological, and Environmental Statistics*, 18:357 - 375, 2013.
- [3] F.K.C. Hui, D.I. Warton, S.D. Foster, and P.K. Dunstan. To mix or not to mix: comparing the predictive performance of mixture models vs. separate species distribution models. *Ecology*, 94:1913 - 1919, 2013.
- [4] R. Leaper, P.K. Dunstan, S.D. Foster, N.S. Barrett, and G.J. Edgar. Do communities exist? complex patterns of overlapping marine species distributions. *Ecology*, in press.
- [5] S.D. Foster, P.K. Dunstan, F. Althaus, and A. Williams. The cumulative effect of trawling on a multi-species fish assemblage off south eastern Australia. In review.
- [6] S.D. Foster, G.H. Givens, G.J. Dornan, P.K. Dunstan, and R. Darnell. Modelling biological regions from multi-species and environmental data. *Environmetrics*, 24(7):489 - 499, 2013.
- [7] F.K.C. Hui, D.I. Warton, and S.D. Foster. An akaike information criterion for mixture models. In review.
- [8] F.K.C. Hui, D.I. Warton, and S.D. Foster. The group lasso for finite mixture of regression models, with applications to species distribution modelling. In review.

