Statistical Advances in Inferring and Forecasting Community Responses to Climate Drivers, using time-series analysis

Northwest Fisheries Science Center Elizabeth Holmes, Eric Ward, Mark Scheurell, Dan Pendleton*, Tessa Francis*

National Center for Ecological Analysis and Synthesis Stephanie Hampton, Lindsay Scheef*, Elizabeth Wolkovich*

LIMNOLOGY and OCEANOG

Assessing m monitoring a compariso sampling da

Lindsay P. Scheef¹*, D Scheuerell², and Davi ¹National Center for Ec ²Northwest Fisheries Sci ³Channel Islands Nation ⁴Sir Alister Hardy Found



CONTRIBUTED RESEARCH ARTICLES

MARSS: Multivariate Autoregressive State-space Models for Analyzing Time-series Data

Could my statistical analyses be cast as a MARSS model (a particular class of time-series model)?

What the next steps for fitting these kind of models?

eis. There are a number of existing is packages for fitting this class of models, including ssnir (Doth, 2010) to allow fixed and shared values within param-

High-resolution data collected over the past 60 years by a single family of Siberian scientists on Lake Baikal reveal significant warming of surface waters and long-term changes in the basal food web of the world's largest, most ancient lake. Attaining depths over 1.6 km, Lake Baikal is the deepest and most voluminous of the world's great lakes. Increases in average water temperature (1.21 °C since 1946), chlorophil *a* (300% since 1978) are the world's largest part of the

narine plankton

725 Montlake Blvd. E, Seattle, WA

ence Center, Hatfield Marine Science of Marine Resources Studies, Oregon

ODEUR†, PHILLIP S.

PETERSON[†]

11

ccurate predictions about ecosystem ponses in a food-web context to supwe conduct time-series analyses with in the Northern California Current on community interactions. Autorel ocean climate conditions. Negative arm phase during the time series, ankton communities. Local environouthern Oscillation) were associated wironmental correlates of zooplank-

arm phase for upwelling as a covarimeous quantitation of community community structure varies with

Using multivariate autoregressive models to infer species interaction strengths: a statistical approach



Lots of applications MAR models to analyze freshwater plankton datasets

Citation	System	
Francis et al 2012	Marine plankton	
Scheef et al 2012	The data have no	
Hall et al 2009 Hampton et al 2008	observation error or have	
Duffy 2007	known error	
Hampton et al 2006		
Huber and Gaedke 2006	Covariates are known	
Hampton and Schindler 2	ponfoctly (no onnon)	
Hampton et al 2006	perfectly (no error)	
Carpenter et al 2005		
Ives et al 2003	No missing values	nkton
Beisner et al 2003	5	
Klug and Cottingham 2001	Decemble confidence on	
Fischer et al 2001	Reasonable confidence on	
Klug et al 2000	how to group species	
Ives et al 1999	Theory + Meshwarer pla	hkton
Ives 1995	Theory	

Some problems we have with ecological data and models

Multivariate time-series data with

- Lots of gaps (missing data)
- (Unknown) observation error
- Complex (unknown) relationships between observation and underlying process trajectory
- Non-ideal covariate data --- instrumentation changed, multiple time series

Solution \rightarrow State-space models \rightarrow MARSS

MARSS models: Multivariate AutoRegressive State-Space

SOME UNDERLYING "HIDDEN" AUTOREGRESSIVE PROCESS



Random walk x(t)=x(t-1)+e(t) +u



Mean-reverting random walk x(t)=bx(t-1)+e(t) +u

+ OBSERVATION PROCESS



NAMES

Univariate: Autoregressive state-space

Mulitvariate:

Vector autoregressive SS Multivariate autoregressive SS Dynamic linear model Structural SS time-series model

MARSS model



Model with lags (lag-p models)

$$\mathbf{x}'_t = \mathbf{B}_1 \mathbf{x}'_{t-1} + \mathbf{B}_2 \mathbf{x}'_{t-2} + \mathbf{u}' + \mathbf{w}'_t, \text{ where } \mathbf{w}'_t \sim \text{ MVN}(0, \mathbf{Q}')$$

$$\times \text{ at t-2 affects x at t}$$

In MARSS form, it becomes... $\mathbf{x}_t = \mathbf{B}\mathbf{x}_{t-1} + \mathbf{u} + \mathbf{w}_t$, where $\mathbf{w}_t \sim \text{MVN}(0, \mathbf{Q})$

$$\begin{bmatrix} \mathbf{x}'_t \\ \mathbf{x}'_{t-1} \end{bmatrix} = \begin{bmatrix} \mathbf{B}_1 \ \mathbf{B}_2 \\ \mathbf{I}_m \ 0 \end{bmatrix} \begin{bmatrix} \mathbf{x}'_{t-1} \\ \mathbf{x}'_{t-2} \end{bmatrix}_{t-1} + \begin{bmatrix} \mathbf{u}' \\ 0 \end{bmatrix} + \mathbf{w}_t, \ \mathbf{w}_t \sim \text{MVN} \left(0, \begin{bmatrix} \mathbf{Q}' \ 0 \\ 0 \ 0 \end{bmatrix} \right)$$

Multivariate moving average models (autocorrelated process or observ noise)

$$\mathbf{x}'_t = \mathbf{w}'_t + \Theta_1 \mathbf{w}'_{t-1} + \Theta_2 \mathbf{w}'_{t-2}$$
, where $\mathbf{w}'_t \sim \text{MVN}(0, \mathbf{Q}')$

In MARSS form, it becomes... $\mathbf{x}_t = \mathbf{B}\mathbf{x}_{t-1} + \mathbf{u} + \mathbf{w}_t$, where $\mathbf{w}_t \sim \text{MVN}(0, \mathbf{Q})$

$$\begin{bmatrix} \mathbf{x}'_{t-2} \\ \mathbf{x}'_{t-1} \\ \mathbf{x}'_{t} \end{bmatrix} = \begin{bmatrix} 0 \ \mathbf{I}_{m} & 0 \\ 0 & 0 \ \mathbf{I}_{m} \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x}'_{t-3} \\ \mathbf{x}'_{t-2} \\ \mathbf{x}'_{t-1} \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \mathbf{w}'_{t} \end{bmatrix}, \ \mathbf{w}_{t} \sim \text{MVN} \left(0, \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & \mathbf{Q}' \end{bmatrix} \right)$$

 $\mathbf{y}_t = \begin{bmatrix} \mathbf{\Theta}_2 \ \mathbf{\Theta}_1 \ 1 \end{bmatrix} \mathbf{x}_t$

Stochastic level model (used to detect structural breaks)

$$\begin{aligned} x_t &= x_{t-1} + w_t \\ y_t &= x_t + v_t \end{aligned}$$
 The mean level is an autoregressive process



Model with covariates to look for effects of the covariate or account for effects



W. T. Edmondson dataset courtesy of D. Schindler

Model the covariates with an observation model and process model

$$\begin{bmatrix} \mathbf{x}^{(v)} \\ \mathbf{x}^{(c)} \end{bmatrix}_{t} = \begin{bmatrix} \mathbf{B}^{(v)} & \mathbf{C} \\ 0 & \mathbf{B}^{(c)} \end{bmatrix} \begin{bmatrix} \mathbf{x}^{(v)} \\ \mathbf{x}^{(c)} \end{bmatrix}_{t-1} + \begin{bmatrix} \mathbf{u}^{(v)} \\ \mathbf{u}^{(c)} \end{bmatrix} + \mathbf{w}_{t}, \ \mathbf{w}_{t} \sim \text{MVN} \left(0, \begin{bmatrix} \mathbf{Q}^{(v)} & 0 \\ 0 & \mathbf{Q}^{(c)} \end{bmatrix} \right)$$
$$\begin{bmatrix} \mathbf{y}^{(v)} \\ \mathbf{y}^{(c)} \end{bmatrix}_{t} = \begin{bmatrix} \mathbf{Z}^{(v)} & \mathbf{D} \\ 0 & \mathbf{Z}^{(c)} \end{bmatrix} \begin{bmatrix} \mathbf{x}^{(v)} \\ \mathbf{x}^{(c)} \end{bmatrix}_{t} + \begin{bmatrix} \mathbf{a}^{(v)} \\ \mathbf{a}^{(c)} \end{bmatrix} + \mathbf{v}_{t}, \ \mathbf{v}_{t} \sim \text{MVN} \left(0, \begin{bmatrix} \mathbf{R}^{(v)} & 0 \\ 0 & \mathbf{R}^{(c)} \end{bmatrix} \right)$$

The covariates can be modeled as a autoregressive process

The covariates might have an observation process (to deal with missing values, multiple time series, changing time series)

Dynamic Factor Analysis: "PCA for time series"

Lot of time series (case reproductive output of Alaska salmon stocks)



Reduce to a few hidden drivers with weighting



Heavily used in economics, finance and engineering



We set out to solve statistical problems that ecologists have but engineers and economists don't have (so much)

Multivariate time-series data with

- Lots of gaps (missing data)
- (Unknown) observation error
- Holmes, E. E. 2010, 2012. Derivation of the EM
 algorithm for constrained and unconstrained multivariate autoregressive state-space (MARSS) models.

changed, multiple time series

Constrained or shared parameters

Solution \rightarrow A general Expectation-Maximization algorithm for MARSS models

Statistical advances aren't very useful without tools to apply them...

🖉 CRAN - Package MARSS - Windows Internet Explorer

Google "MARSS cran" or "time series task view cran"

MARSS: Multivariate Autoregressive State-Space Modeling

The MARSS package provides maximum-likelihood parameter estimation for constrained and unconstrained linear multivariat data. Fitting is primarily via an Expectation-Maximization (EM) algorithm, although fitting via the BFGS algorithm (using the op model (DLM) and vector autoregressive model (VAR) model. Functions are provided for parametric and innovations bootstra (AICb), confidences intervals via the hessian approximation and via bootstrapping and calculation of auxilliary residuals for det for parameter estimation for a variety of applications, model selection, dynamic factor analysis, outlier and shock detection, an at the R command line to open the MARSS user guide.

Version:	2.7
Depends:	MASS, mvtnorm, nlme, time, KFAS
Published:	2011-10-23
Author:	Eli Holmes, Eric Ward, and Kellie Wills, NOAA, Seattle, USA
Maintainer:	Eli Holmes <eli.holmes at="" noaa.gov=""></eli.holmes>
License:	<u>GPL-2</u>
In views:	TimeSeries
CRAN checks	MARSS results

Downloads:

 Package source:
 MARSS 2.7.tar.gz

 MacOS X binary:
 MARSS 2.7.tgz

 Windows binary:
 MARSS 2.7.zip

 Reference manual:
 MARSS.pdf

 Vignettes:
 EM Derivation Quick Start Guide Changes between versions

 Old sources:
 MARSS archive

 E. E. Holmes and E. J. Ward

Analysis of multivariate timeseries using the MARSS package

version 2.7

October 21, 2011

Mathematical Biology Program Northwest Fisheries Science Center, Seattle, WA

Developed with support by the Comparative Analysis of Marine Ecosystem Organizations (CAMEO) Program

Lots of case studies and examples from workshops we (Eric Wark, Brice Semmens, Mark Scheuerell, and myself) have taught

12	Case Study 3: Using MARSS models to identify spatial
	population structure and covariance
	12.1 The problem
	12.2 How many distinct subpopulations?
	12.3 Is Hood Canal separate?
13	Case Study 4: Dynamic factor analysis (DFA) using
	MARSS
	13.1 Dynamic factor analysis 111
	13.2 The data
	13.3 Setting up the model
	13.4 Fitting the model
	13.5 Using model selection to determine the number of trends 117
	13.6 Using a varimax rotation to determine the trend loadings \dots 120
14	Case Study 5: Using state-space models to analyze noisy
	animal tracking data123
	14.1 A simple random walk model of animal movement
	14.2 The problem
	14.3 Estimate locations from bad tag data
	14.4 Comparing turtle tracks to proposed fishing areas
	14.5 Using specialized packages to analyze tag data $\dots \dots \dots \dots 129$
15	Case Study 6: Detection of outliers and structural breaks
	using MARSS
	15.1 Detection of outliers and structural breaks

Effects of observation error on estimates of species interaction strength



Lake Washington plankton community, S.E. Hampton, NCEAS, UCSB



Week

Observation error and spurious density-dependence Original interaction matrix





Week

What are the effects of observation error on estimates of large interaction matrices? Comparison using long-term plankton studies in the English Channel

"L4" (**good**)

- 1 location
- Weekly samples at standard time of day
- Individual counts
- Very few missing values

Continuous Plankton

Recorder (CPR) (**NOISY**)

- •"platforms of opportunity" = many locations
- Log10 counts
- Times of day variable
- Lots of missing values
- Some spp poorly sampled



Scheef, L. P., D. E. Pendleton, S. E. Hampton, S. L. Katz, E. E. Holmes, M. D. Scheuerell, and D.G. Johns. 2012. Assessing marine plankton community structure from long-term monitoring data with multivariate autoregressive (MAR) models: a comparison of fixed station vs. spatially distributed sampling data. Limnology & Oceanography: Methods 10: 54-64. 14 groups 12







■¤	Group¤	Propo com	ortion of munity¤		Taxa. Included¤	Proportio: group	n∙of¶ ∝	c	
■¤		L4¤	····CPR¤	×	IIICIGGCG~	L4·&·CPR·	mean¤	c	
• ¤	<u>Chaetognaths</u> ¤	0.02¤	0.07¤	×	Sagitta spp.¤	~1.00¤	¤	c	
• ¤	Pteropods¤	0.01¤	0.02¤	×	<u>Thecosomata¤</u>	>0.99¤	¤	c 💘	
• ¤	Tunicates¤	0.03¤	0.07¤	×	Appendicularian Doliolids¤	ຢູ່ໃ 099 001⊧	a		0 300 M
• ¤	Cladocerans	0.05¤	0.04¤	×	Eyadue opp.¶ Podon opp.¤	0.66¶ 0.34⊧	¤	c 🧖	
• ¤	Amphipods¤	<0.01¤	<0.01¤	×	Gammarid-ampl Hyperiid amphij Isopods¶ Mysid-shrimp¤	iģods¶ 094¶ ∞ds¶ 0.03¶ 0.02¶ 0.01⊧	¤	r /	
• ¤	Krill¤	<0.01¤	<0.01¤	×	<u>Euphausiids</u> ×	~1.00≍	¤	C	
	Large- <u>calanoids¤</u>	0.03¤	0.08¤	×	Calarus spp.¶ Metridia spp.¶ Candacia spp.¶ Eucalarus spp.¤	095¶ 003¶ 001 001 \$	×	C	
opepodsa	Small- <u>calanoids</u> ¤	0.38¤	0.45¤	×	Pseudocalanus - Acastia - spp. ¶ Temesa - spp. ¶ Paracalanus - spp Centropaces - spr Clausocalanus - sp Ctemocalanus - sp	፼·∬ 0.33¶ 0.28¶ 0.1.5¶ ∬ 0.12¶ .¶ 0.06¶ ፼·∬ 0.02¶ p.≍ 0.01≈	×	C	
0	Cyclopoids¤	0.12¤	0.02¤	×	Qithona spp.¤	~1.00¤	¤	Ţc	
	Poecilostomatoids [®]	0.19¤	0.01¤	×	Correaeus spp.¶ Oncæa spp.¤	0.51¶ 0.49≍	a	c	
	Harpacticoids ²³	0.01¤	<0.01¤	×	Euterpina spp.¶ Clytennestra sp Microsetella spp Alteutha spp.×	0.70¶ 0.23¶ .¶ 0.05¶ 0.01≍	¤	c	
	Cirripedia¤	0.08¤	0.01¤	×	Cirripede larvae	≅ 1.00≒	¤	c	
Meroplankton	Mero, grazers¶ (miscellaneous)¤	0.06¤	0.23¤	×	Echinoderm law Bivalve lavae¶ Cyphonaste lavv Polychaete lavva Gastropod lavva	ae¶ 0.66¶ 0.19¶ ae¶ 0.05¶ e¶ 0.05¶ ₂≋ 0.04≥	×	- C	
	Decapod larvae¤	0.01¤	0.01¤	×	Crab & shrimp l	arvae¤ 1.00¤	¤	c	

B matrix

 12 x 12 interaction matrix

L4 (clean data)

- Green R>0
- Yellow (ignores obs error) R=0

CPR (noisy data)

- Red
- Pink (ignores obs error)





B matrix

- 12 x 12
- 2 groups removed because they had only 3 levels in the CPR data

L4 (clean data)

- Green
- Yellow (ignores obs error)

CPR (noisy data)

- Red
- Pink (ignores obs error)





Gets back to the "unknown observation error = poor B estimation" issue



Univariate case (one spp)

- Largely solved by independent samples of same population
- Partially solved with duplicate samples of different populations with same parameters

Multivariate case (community)

- This is where our current research is focused
- Research depends on simulations (10,000s), so fast algorithms key

"We have not succeeded in answering all our problems. The answers we have found only serve to raise a whole set of new questions. In some ways we feel we are as confused as ever, but we believe we are confused on a higher level and about more important things." Earl C. Kelley, 1951, 'The Workshop Way of Learning'

Also there is a less recognized issue: the effect of unknown environmental drivers



Univariate case (one spp)

- This is bad unless you can demean your data without removing the true fluctuations.
- If you remove those in your demeaning step, $B \rightarrow 0$ (spurious DD)
- Datasets much longer than any cycles in the unknown covariate are key.

Multivariate case (community)

• It is generally accepted that inclusion of the important environmental drivers is key for good B estimation

Finally, there are way too many estimated B elements

Constraining B will vastly improve estimation

But current model selection algorithms (for MAR) require searching a huge model space and the fitting step for MARSS is too slow, i.e. model selection steps would = months of computation



Motivation: understanding plankton dynamics from long-term marine data sets



Lots of applications to freshwater plankton datasets

Interaction matrix

$$\mathbf{x}_{t} = \mathbf{B}\mathbf{x}_{t-1} + \mathbf{u} + \mathbf{C}\mathbf{f}_{t} + \mathbf{w}_{t},$$
$$\mathbf{y}_{t} = \mathbf{Z}\mathbf{x}_{t} + \mathbf{a} + \mathbf{D}\mathbf{f}_{t} + \mathbf{v}_{t},$$

 \mathbf{y}_{t}

Assume the data have no observation error

Citation	System
Hall et al 2009	Freshwater plankton
Hampton et al 2008	Freshwater plankton
Duffy 2007	Freshwater plankton
Hampton et al 2006	Freshwater plankton
Huber and Gaedke 2006	Freshwater plankton
Hampton and Schindler 2006	Freshwater plankton
Carpenter et al 2005	Freshwater plankton
Beisner et al 2003	Freshwater plankton
Klug et al 2000	Freshwater plankton
Hampton et al 2006	Freshwater plankton
Fischer et al 2001	Freshwater plankton
Ives et al 1999	Freshwater plankton
Ives et al 2003	Freshwater plankton
Klug and Cottingham 2001	Freshwater plankton

14 groups Lots of spp







■¤	Group¤	Proportion of community¤			Proportion of community ²² Included ²³		n∙of¶ ¤	c
۰¤		L4¤	····CPR¤	¤	menadea	L4·&·CPR·	mean¤	c
• ¤	Chaetognaths¤	0.02¤	0.07¤	×	Sagitta∵spp.¤	~1.00¤	¤	C
• ¤	Pteropods [©]	0.01¤	0.02¤	×	Thecosomata¤	>0.99¤	ä	C XXX
• ¤	Tunicates¤	0.03¤	0.07¤	×	Appendicularian Doliolids¤	ឲ្រ[] 0.99 0.01≍	¤	C C C C C C C C C C C C C C C C C C C
•¤	Cladocerans [®]	0.05¤	0.04¤	×	Evadue spp.¶ Podon spp.¤	0.66¶ 0.34≍	¤	c 💓 🕅
• ¤	Amphipods¤	<0.01¤	<0.01¤	×	Gammæid ampl Hyperiid amphir Isopods¶ Mysid shrimp¤	iģods¶ 0.94¶ pods¶ 0.03¶ 0.02¶ 0.01⊧	¤	c /
• ¤	Krill¤	<0.01¤	<0.01¤	×	Euphausiids¤	~1.00¤	a	C at a farmer
	Large- <u>calanoids</u> ¤	0.03¤	0.08¤	x	Calamıs spp.¶ Metridia spp.¶ Candacia spp.¶ Eucalamıs spp.¤	095¶ 003¶ 001¶ 001\$	¤	c
opepodsa	Small- <u>calanoids</u> ¤	0.38¤	0.45¤	×	Pseudocalanus -5: Acatta -5pp. ¶ Temora -5pp. ¶ Paracalanus -5pp. Centropages -5pp Clausocalanus -5p Ctenocalanus -5p	pp.¶ 0.33¶ 0.28¶ 0.15¶ .¶ 0.12¶ .¶ 0.06¶ pp.¶ 0.02¶ p.≍ 0.01≈	×	C
0	Cyclopoids¤	0.12¤	0.02¤	×	Qithona spp.¤	~1.00¤	¤	c
	<u>Poecilostomatoids¤</u>	0.19¤	0.01¤	×	Corycaeus spp.¶ Oncæa spp.¤	0.51¶ 0.49≍	¤	c
	Harpacticoids ²³	0.01¤	<0.01¤	×	Euterpina-spp.¶ Clytennestra-spj Microsetella-spp Alteutha-spp.¤	0.70¶ 0.23¶ ∿¶ 0.05¶ 0.01≍	×	c
8	Cirripedia¤	0.08¤	0.01¤	×	<u>Cirripede</u> larvae	≅ 1.00¤	¤	c
Meroplankton	Mero, grazers¶ (miscellaneous)¤	0.06¤	0.23¤	×	Echinoderm law Bivalve lavae¶ Cyphonaste lavv Polychaste lavva Gastropod lavva	rae¶ 0.66¶ 0.19¶ cae¶ 0.05¶ ce¶ 0.05¶ s≍ 0.04≥	¤	
	Decapod larvae¤	0.01¤	0.01¤	×	Crab & shrimp l	arvae¤ 1.00¤	¤	Ľ

But the talk is about the statistical methods

- What is a multivariate autoregressive state-space model (MARSS or VARSS)?
- o A tour of different classes of time series models written as MARSS (more math)
- o Estimating parameters using an EM algorithm for MARSS models with linear constraints (more math)
- o MARSS R package
- o Estimating the species interaction matrix and covariate matrices for PLANKTON (actually more math)
- o Some results from the plankton work

Finding MLE parameters for MARSS models

Joint likelihood of y(data), x(hidden states)

$$\log \mathbf{L}(\boldsymbol{y}, \boldsymbol{x}; \Theta) = -\sum_{1}^{T} \frac{1}{2} (\boldsymbol{y}_{t} - \mathbf{Z}\boldsymbol{x}_{t} - \mathbf{a})^{\mathsf{T}} \mathbf{R}^{-1} (\boldsymbol{y}_{t} - \mathbf{Z}\boldsymbol{x}_{t} - \mathbf{a}) - \sum_{1}^{T} \frac{1}{2} \log |\mathbf{R}|$$
$$-\sum_{1}^{T} \frac{1}{2} (\boldsymbol{x}_{t} - \mathbf{B}\boldsymbol{x}_{t-1} - \mathbf{u})^{\mathsf{T}} \mathbf{Q}^{-1} (\boldsymbol{x}_{t} - \mathbf{B}\boldsymbol{x}_{t-1} - \mathbf{u}) - \sum_{1}^{T} \frac{1}{2} \log |\mathbf{Q}|$$
$$-\frac{1}{2} (\boldsymbol{x}_{0} - \boldsymbol{\xi})^{\mathsf{T}} \boldsymbol{\Lambda}^{-1} (\boldsymbol{x}_{0} - \boldsymbol{\xi}) - \frac{1}{2} \log |\boldsymbol{\Lambda}| - \frac{n}{2} \log 2\pi$$

T .

• If you can compute the marginal likelihood $L(y; \Theta)$, you can maximize that (using some Newton-based method, like BFGS). The Kalman filter will give you the marginal likelihood. Works great for lots of problems. But for many big multivariate problems it doesn't work so great.

• A different approach to finding MLE parameters for problems with hidden states is the Expectation-Maximization (EM) algorithm.

EM algorithm

Joint likelihood of y and x is log $L(y,x;\Theta) = f(y,x,\Theta)$

The EM algorithm maximizes the expected value of the joint likelihood

$$\mathbf{E}_{\mathbf{X}\mathbf{Y}}[\log \mathbf{L}(\boldsymbol{Y}, \boldsymbol{X}; \Theta); \boldsymbol{Y}(1) = \boldsymbol{y}(1), \Theta_j]$$

Expected value of the "random variable LL" conditioned on the observed data and a set of parameters

$$\begin{split} \mathbf{E}_{\mathbf{X}\mathbf{Y}}[\log \mathbf{L}(\boldsymbol{Y}, \boldsymbol{X}; \Theta); \boldsymbol{Y}(1) &= \boldsymbol{y}(1), \Theta_j] = \\ \mathsf{g}(\mathsf{E}(\mathsf{Y}\mathsf{X}), \mathsf{E}(\mathsf{X}\mathsf{X}), \mathsf{E}(\mathsf{Y}\mathsf{Y}), \mathsf{E}(\mathsf{X}), \mathsf{E}(\mathsf{Y}), \Theta) \end{split}$$

The expectations in this expected joint likelihood can be computed (for MARSS models with the Kalman smoother)

We can maximize $g(..., \Theta)$ with respect to Θ to find the Θ that maximizes the expected log likelihood.

EM algorithm for MARSS models

- 1. Start with Θ_1
- 2. Compute the expectations involving X and Y conditioned on Θ_1 and the data
- 3. Put those $E_{XY}[\log L(Y, X; \Theta); Y(1) = y(1), \Theta_j]$ and maximize with respect to Θ to get Θ_2
- 4. Compute the expectations involving X and Y conditioned on Θ_2 and the data
- 5. Put those $E_{XY}[\log L(Y, X; \Theta); Y(1) = y(1), \Theta_j]$ and maximize with respect to Θ to get Θ_3
- 6. Repeat until convergence

But the talk is about the statistical methods

- o What is a multivariate autoregressive state-space model (MARSS or VARSS)?
- o A tour of different classes of time series models written as MARSS (more math)
- o Estimating parameters using an EM algorithm for MARSS models with linear constraints (more math)
- o MARSS R package
- o Estimating the species interaction matrix and covariate matrices for PLANKTON (actually more math)
- o Stabilitiy metrics (cartoons!)
- o Some results from the plankton work

Multispecies Autoregressive Models (MARs) as used in community modeling



Not a new result but perhaps not widely recognized... unknown obs error = spurious density-dependence

Ecology Letters, (2011)

doi: 10.1111/j.1461-0248.2011.01702 x

Are patterns of density dependence in the Global Population Dynamics Database driven by uncertainty about population abundance?

Abstract

Jonas Knape* and Perry de Valpine Department of Environmental Science, Policy and Management, 137 Mulford Hall 4B114, University of California, Berkeley, Berkeley CA 94720 1154 *Correspondence: E-mail: iknape@berkelex.edu

LETTER

Density dependence in population growth rates is of immense importance to ecological theory and application, but is difficult to estimate. The Global Population Dynamics Database (GPDD), one of the largest collections of population time series available, has been extensively used to study cross-taxa patterns in density dependence. A major difficulty with assessing density dependence from time series is that uncertainty in population abundance estimates can cause strong bias in both tests and estimates of strength. We analyse 627 data sets in the GPDD using Gompertz population models and account for uncertainty via the Kalman filter. Results suggest that at least 45% of the time series display density dependence, but that it is weak and difficult to detect for a large fraction. When uncertainty is ignored, magnitude of and evidence for density dependence is strong, illustrating that uncertainty in abundance estimates qualitatively changes conclusions about density dependence drawn from the GPDD.

Keywords

Density dependence, GPDD, observation error, time series,

Embg Letters (2011)

INTRODUCTION

Density dependence in population growth rates is a fundamental concept for ecological theory as well as for population management. Estimating density dependence in wild populations has, however, proved challenging. Ideally, density dependence in growth rates should be estimated directly from the effects of density acting on the traits contributing to population growth. Given current progress in statistical methods for jointly analysing data on both population size and demographic traits (Besbeas et al. 2005), and with long-term population studies involving demographic data becoming increasingly mmon, this approach holds a bright future. However, the number of such studies is currently limited and they only cover a rather narrow range of taxa. Long-term time series on population abundance are more common and can be used to estimate density dependence in population growth rates. Under this approach, density dependence is defined as a general tendency of per capita growth rates to decrease when population size is large and increase when it is small, and is identified as a statistical pattern not tied to any specific biological mechanism (Wolda & Dennis 1993).

It was noted early that estimates and tests of density dependence based on regressing log transformed current observed population size, y_s on previous log transformed observed population size, y_{s-b} are sensitive to uncertainty in the observations (St-Amant 1970; Kuno 1971; Itô 1972; Slade 1977). Similar concerns were aired about estimates from fisheries models of stock-recruitment data (Ludwig & Walters 1981; Walters & Ludwig 1981). Uncertainty inflates the Type I error rate of tests for density dependence (Shenk et al. 1998) and tends to bias estimates towards stronger density dependence if dynamics are under-compensatory and towards weaker density dependence if dynamics are over-compensatory (Benson 1973). Bulmer (1975) devised two tests for density dependence taking the time series nature of the data into account. One of those was designed to be robust

against uncertainty about population size and has been shown to perform better than density dependence tests ignoring uncertainty in estimates of population abundance (Shenk et al. 1998). Simple procedures to correct for effects of uncertainty such as the SIMEX method have been suggested (Solow 1998; Freckleton et al. 2006) but typically require that the variance of the uncertainty about population size is known. A more direct approach to account for uncertainty is provided by state space models, first used for modelling population dynamics in the fisheries literature (e.g. Mendelssohn 1988; Sullivan 1992). State space models in these cases consist of a model of a population dynamical process combined with a model of the uncertainty in the abundance estimates, sometimes termed observation, measurement or sampling error, and may be used to estimate the variance of this uncertainty as well as to filter out its effects (de Valpine & Hastings 2002; Calder et al. 2003; Buckland et al. 2004; Dennis et al. 2006). Estimates derived from state space models tend to have smaller bias than estimates ignoring uncertainty about population abundance, but can also have large variances (Knape 2008), and the statistical properties of even simple state space model estimators are not fully understood (Dennis et al. 2006; Lebreton 2009).

The Global Population Dynamics Database (GPDD), containing over 5000 time series on population abundances obtained from various forms of population surveys, has provided an opportunity for ecologists to explore population dynamical patterns over a wide range of taxa (Inchausti & Halley 2001). Analyses using data in the GPDD have focused on, e.g., extinction risks (Fagan et al. 2001; Inchausti & Halley 2003; Brook et al 2006), population cycles (Kendall et al 1998; Murdoch et al. 2002) and effects of weather (Knape & de Valpine 2011) but, arguably, the studies stirring the most attention as well as debate have addressed population regulation and density dependence. These have explored patterns in the shape of density dependence (Sibly et al. 2005; Polansky et al. 2009) and in the strength of regulation and density dependence (Brook & Bradshaw 2006; Sibly et al. 2007;

© 2011 Hackwell Publishing Lad/CNRS

Ecology, 89(11), 2008, pp. 2994-3000 © 2008 by the Ecological Society of America

ESTIMABILITY OF DENSITY DEPENDENCE IN MODELS OF TIME SERIES DATA

JONAS KNAPE

Department of Theoretical Ecology, Ecology Building, Lund University, SE-22362, Lund, Sweden

Abstract. Estimation of density dependence from time series data on population abundance is hampered in the presence of observation or measurement errors. Fitting state-space models has been proposed as a solution that reduces the bias in estimates of density dependence caused by ignoring observation errors. While this is often true, I show that, for specific parameter values, there are identifiability issues in the linear state-space model when the strength of density dependence and the observation and process error variances are all unknown. Using simulation to explore properties of the estimators, I illustrate that, unless assumptions are imposed on the process or observation error variances, the variance of the estimator of density dependence varies critically with the strength of the density dependence. Under compensatory dynamics, the stronger the density dependence the more difficult it is to estimate in the presence of observation errors. The identifiability issues disappear when density dependence is estimated from the state-space model with the observation error variance known to the correct value. Direct estimates of observation variance in abundance censuses could there fore prove helpful in estimating density dependence but care needs to be taken to assess the uncertainty in variance estimates

Key words: density dependence; state-space models; time series analysis.

INTRODUCTION

eports

Density dependence can be loosely defined as a quantitative influence of population size on some life history or population trait of interest. The concept is of central importance to population ecology since it determines both the limiting and the short time behavior of the dynamics of populations. Empirical estimates of density dependence are therefore important from a scientific as well as from a management perspective. Assessment of density dependence in the dynamics of natural populations has however proved to be challenging (Dennis et al. 2006).

When relevant data are available, effects of density dependence can be directly linked to life history traits. For instance, density dependence in recruitment (e.g., Crespin et al. 2006) and survival (e.g., Festa-Bianchet et al. 2003) have been estimated by mark-recapture analyses and density dependence in fecundity has been inferred from data on reproduction (e.g., Solbreck and Ives 2007). Density dependence in life history traits influences density dependence in population growth rate (Lande et al. 2002). It can be argued that density dependence in population growth is the most important form of density dependence for determining long-term behavior of populations. However, since the link from demographic traits to population change is almost never known with good precision, density dependence in

Manuscrint received 12 January 2008: revised 2 June 2008: accepted 12 June 2008. Corresponding Editor: M. Lavine. ¹E-mail: jonas.knape@teorekol.lu.se

population growth rate is not easily inferred from life history data even if the effects of density dependence on several life history traits are well known. Time series analysis of population abundance data provides an alternative or complementary method that ideally could serve as a more direct way of estimating density dependence in population growth rate.

Estimates of density dependence must rely on measures of population density that are usually difficult to obtain with precision (Freckleton et al. 2006). This problem is particularly relevant to estimates of density dependence in growth rate derived from time series data on population size in that both the dependent and the independent variable are measured with uncertainty. Introducing observation error to dynamical data changes its dynamical structure (Dennis et al. 2006) and estimators relating to the dynamics of the data that do not account for observation errors are therefore often biased. Specifically, tests and estimators of density dependence based on time series data are known to be biased if observation errors are present but ignored for both direct (Kuno 1971, Walters and Ludwig 1981, Shenk et al. 1998, Freckleton et al. 2006) and delayed (Solow 2001) density dependence. An appealing method for overcoming this difficulty is provided by the statespace framework (Harvey 1990), a general term for statistical models of observations of hidden state variables that are dynamically linked through time. For time series data on population abundance, statespace models can be used for explicit modeling of both the observation and the population dynamical processes (Stenseth et al. 2003, Jamieson and Brooks 2004).

2004



Lake Washington long-term plankton monitoring

- weekly plankton sampling 1960s to present
- environmental covariate data
- standardized sampling

Seattle

Lake Washington

 basis for lots of MAR-based research into plankton community dynamics



	1	0.65	-	-0.08				
MARSS	2		0.36					
R ast = 02	3			0.34				
Or so	4				0.64	-0.17		-0.16
Q est=.6	5	0.11		0.01		0.74		
01 50	6	0.11				-0.26	0.71	-0.15
	7	0.06				-0.13		0.57

Comparison of the B matrix estimates analysis of Lake WA data mid-1970s on

	1	0.52	-	-0.06					Diatom
	2		0.44						Green
ORIGINAL	3			0.92					Cyano
	4				0.66	-0.16		-0.12	Cyclops
	5	0.23		-0.11		0.57			Daphnia
	6	0.09				-0.22	0.64	-0.14	Diatopmis
	7	0.23				-0.27		0.52	Bosmina

Those results assumed we knew where the zeros were. What if we don't know?

This is the same 7x7 interaction matrix. The distributions are posterior distributions. All B elements were estimated but I blocked out the original "zeros"



A different approach to maximizing the logL of hidden states problems: EM algorithm

- 1. Start with Θ_1
- 2. Compute the expectations involving X and Y conditioned on Θ_1 and the data
- 3. Put those in $\operatorname{E}_{\mathbf{XY}}[\log \mathbf{L}(\mathbf{Y}, \mathbf{X}; \Theta); \mathbf{Y}(1) = \mathbf{y}(1), \Theta_j]$ and maximize with respect to Θ to get Θ_2
- 4. Compute the expectations involving X and Y conditioned on Θ_2 and the data
- 5. Put those in $E_{XY}[\log L(Y, X; \Theta); Y(1) = y(1), \Theta_j]$ and maximize with respect to Θ to get Θ_3
- 6. Repeat until convergence

What's the point?

- 1) It can make certain types of model fitting problems tractable by being considerably faster and more stable
- 2) For many of the problems we work on, other approaches grind to a halt

"EM algorithms sound like fun!"

Google "MARSS cran"

🖉 CRAN - Package MARSS - Windows Internet Explorer									
🚱 💿 💌 🕼 http://cran.r-project.org/web/packages/MARS5/index.html									
File Edit View Favorites Tools Help	х 🇞 -								
X 🔗 McAfee 🖌 🗸									
🚖 Favorites 🛛 🚔									
CRAN - Package MARSS									

MARSS: Multivariate Autoregressive State-Space Modeling

The MARSS package provides maximum-likelihood parameter estimation for constrained and unconstrained linear multivariate aud data. Fitting is primarily via an Expectation-Maximization (EM) algorithm, although fitting via the BFGS algorithm (using the optim f model (DLM) and vector autoregressive model (VAR) model. Functions are provided for parametric and innovations bootstrappir (AICb), confidences intervals via the hessian approximation and via bootstrapping and calculation of auxiliary residuals for detectin for parameter estimation for a variety of applications, model selection, dynamic factor analysis, outlier and shock detection, and ad at the R command line to open the MARSS user guide.

Version:	2.7
Depends:	MASS, mvtnorm, nlme, time, KFAS
Published:	2011-10-23
Author:	Eli Holmes, Eric Ward, and Kellie Wills, NOAA, Seattle, USA
Maintainer:	Eli Holmes <eli.holmes at="" noaa.gov=""></eli.holmes>
License:	<u>GPL-2</u>
In views:	TimeSeries
CRAN checks:	MARSS results

Downloads:

 Package source:
 MARSS 2.7.tar.gz

 MacOS X binary:
 MARSS 2.7.tgz

 Windows binary:
 MARSS 2.7.zip

 Reference manual:
 MARSS.pdf

 Vignettes:
 EM Derivation Quick Start Guide User Guide Changes between versions

 Old sources:
 MARSS archive

Derivation of the EM algorithm for constrained
and unconstrained multivariate autoregressive
state-space (MARSS) models
DRAFT

Elizabeth Eli Holmes Northwest Fisheries Science Center, NOAA Fisheries 2725 Montlake Blvd E., Seattle, WA 98112 eli.holmes@noaa.gov http://faculty.washington.edu/eeholmes

October 21, 2011

C	ontents		
1	Overview	2	
2	The EM algorithm	6	
3	The unconstrained update equations	10	
4	The constrained update equations	24	
5	Computing the expectations in the update equations	39	
6	Degenerate variance modifications	46	
7	Implementation comments	55	
8	MARSS R package	57	

citation. Rolmay, E. E. 2013. Derivation of the SM algorithm for constrained and unconstrained multivariate autors gravity estate-space (MARSS) models. Unpublished report. Northeres Fisheries Science Conter, NOAA Fisheries Searcie, WA, USA.

Advances in Multivariate AutoRegressive State-Space (MARSS) Models for Analysis of Ecological Data

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_t = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{t-1} + \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}_t, \quad \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}_t \sim MVN\left(\begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, \begin{bmatrix} q_{11} & q_{12} \\ q_{21} & q_{22} \end{bmatrix} \right)$$
$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}_t = \begin{bmatrix} z_{11} & z_{12} \\ z_{21} & z_{22} \\ z_{31} & z_{32} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_t + \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}_t, \quad \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}_t \sim MVN\left(\begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}, \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \right)$$

Finding MLE parameters for MARSS models

Joint likelihood of y(data), x(hidden states)

$$\log \mathbf{L}(\boldsymbol{y}, \boldsymbol{x}; \Theta) = -\sum_{1}^{T} \frac{1}{2} (\boldsymbol{y}_{t} - \mathbf{Z}\boldsymbol{x}_{t} - \mathbf{a})^{\mathsf{T}} \mathbf{R}^{-1} (\boldsymbol{y}_{t} - \mathbf{Z}\boldsymbol{x}_{t} - \mathbf{a}) - \sum_{1}^{T} \frac{1}{2} \log |\mathbf{R}|$$
$$-\sum_{1}^{T} \frac{1}{2} (\boldsymbol{x}_{t} \mathbf{H} \mathbf{O} \mathbf{G}^{-1} \mathbf{L} (\mathbf{y}, \mathbf{X}^{1}, \mathbf{O}) \mathbf{E}^{\mathsf{T}} - \mathbf{f} (\mathbf{y}, \mathbf{X}, \mathbf{X}^{\mathsf{T}}) \log |\mathbf{Q}|$$
$$-\frac{1}{2} (\boldsymbol{x}_{0} - \boldsymbol{\xi})^{\mathsf{T}} \boldsymbol{\Lambda}^{-1} (\boldsymbol{x}_{0} - \boldsymbol{\xi}) - \frac{1}{2} \log |\boldsymbol{\Lambda}| - \frac{n}{2} \log 2\pi$$

11

• If you can compute the marginal likelihood L(y; Θ), you can maximize that (using some Newton-based method, like BFGS). optim() in R.

- The Kalman filter will give you the marginal likelihood.
- Works great for lots of problems. But for many big multivariate problems it doesn't work so great.

A different approach to parameter estimation for hidden state problems: Expectation-Maximization algorithms

Holmes, E. E. 2010. Derivation of the EM algorithm for constrained and unconstrained multivariate autoregressive state-space (MARSS) models.

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_t = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{t-1} + \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}_t, \quad \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}_t \sim MVN\left(\begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, \begin{bmatrix} q_{11} & q_{12} \\ q_{21} & q_{22} \end{bmatrix}\right)$$

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}_t = \begin{bmatrix} z_{11} & z_{12} \\ z_{21} & z_{22} \\ z_{31} & z_{32} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_t + \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}_t, \quad \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}_t \sim MVN\left(\begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}, \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix}\right)$$

$$\alpha + \beta_1 p_1 + \beta_2 p_2 + \dots$$

What's the point?

- 1) It can make **multivariate** AR **state-space** model fitting problems tractable by being considerably faster and more stable
- 2) For many of the problems we work on, other approaches grind to a halt

What are the strong species interactions? How are environmental factors affecting species?



Hampton, Scheuerell, & Schindler 2006

Written in matrix form:



Autoregressive process noise

$$\mathbf{x}'_t = \mathbf{B}\mathbf{x}'_{t-1} + \mathbf{u}' + \boldsymbol{\eta}_t$$

where $\boldsymbol{\eta}_t$ is a AR-1 (or p) process.

We re-write this as a MARSS(1) model by moving the error term into the state process

$$\mathbf{x}_t = \mathbf{B}\mathbf{x}_{t-1} + \mathbf{u} + \mathbf{w}_t$$
, where $\mathbf{w}_t \sim \text{MVN}(0, \mathbf{Q})$

$$\begin{bmatrix} \mathbf{x}' \\ \boldsymbol{\eta} \end{bmatrix}_{t} = \begin{bmatrix} \mathbf{B}_{x} \ \mathbf{I}_{m} \\ 0 \ \mathbf{B}_{\eta} \end{bmatrix} \begin{bmatrix} \mathbf{x}' \\ \boldsymbol{\eta} \end{bmatrix}_{t-1} + \begin{bmatrix} \mathbf{u}' \\ 0 \end{bmatrix} + \mathbf{w}_{t}, \ \mathbf{w}_{t} \sim \text{MVN} \left(0, \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{Q}_{\eta} \end{bmatrix} \right)$$

So...if MARSS models have such a history, can't you just use finance algorithms?

Existing methods dealt with this

 Gappy data, observation error. non-ideal covariate data

But

 Parameter estimation based on Newton methods which struggle with general MARSS models.

Goal

 Develop a general robust algorithm for constrained MARSS models based on Expectation-Maximization algorithms

Holmes, E. E. 2010, 2012. Derivation of the EM algorithm for constrained and unconstrained multivariate autoregressive state-space (MARSS) models.

What are the effects of observation error on estimates of large B matrices?



Hampton, Scheuerell, & Schindler 2006