On automating assessment model diagnostics and the need for simulation testing

Henning Winker^{1*}, Massimiliano Cardinale², Francesco Masnadi^{3,4}, Laurence Kell⁵, Iago Mosqueira⁶, Felipe Carvalho⁷

*Henning.Winker@ec.europa.eu

CAPAM, January 2022

¹Joint Research Centre (JRC), European Commission, TP 051, Via Enrico Fermi 2749, 21027 Ispra (VA), Italy ²Swedish University of Agricultural Sciences, Department of Aquatic Resources, Institute of Marine Research, Lysekil, Sweden ³National Research Council, Institute for Marine Biological Resources and Biotechnology (CNR IRBIM), Ancona, Italy ⁴Department of Biological, Geological, and Environmental Sciences (BiGeA), Alma Mater Studiorum—University of Bologna, Bologna, Italy ⁵Centre for Environmental Policy, Imperial College London, Weeks Building, 16-18 Princes Gardens, London SW7 1NE, UK ⁶Wageningen Marine Research, Haringkade 1, 1976CP IJmuiden, the Netherlands ⁷NOAA Fisheries, Pacific Islands Fisheries Science Center, Honolulu, HI, United States

Outline

- A very brief evolution of model diagnostic approaches
- Some real world examples
- Towards automation of model diagnostic in benchmarks
- The need for simulation testing
- A simplified trial how a simulation could look like
- Remarks on model plausibility







"Modern": Model diagnostics

+ Data conflicts

+ Prediction Skill

M. Maunder, K.R. Piner/Fisheries Research 192 (2017) 16–27

Fig. 5. Flow chart depiction of a simplification of the recommended modelling approach.







'ss3diags': Stock Synthesis

R package hosted on Github

https://github.com/JABBAmodel/ss3diags

Winker H, Carvalho F, Cardinale M, Kell LT

The R package `ss3diags` automates a set of model diagnostics tools plots and ensemble weighting functions for Stock Synthesis

User guidelines for Advanced Model Diagnostics with ss3diags

Felipe Carvalho (NOAA) Hennig Winker (JRC-EC) Massimiliano Cardinale (SLU)

Laurence Kell (Sea++)

20 August, 2021

Contents

1	Getting started 1.1 Installation 1.2 Loading built-in example data	$\frac{1}{2}$
2	Model Diagnostics with ss3diags 2.1 Residual diagnostics 2.2 Retrospective and Forecast bias 2.3 Hindeast Cross-Validation and prediction skill	3 3 8 9
3	Model uncertainty 1	14
4	Cookbook Recipies 2 4.1 Retrospectives with hindcasts 2 4.2 R0 profiling 2 4.3 ASPM diagnostic 2 4.4 Jittering 2	21 21 24 24 24

a ...



'JABBA': First was JABBA

R package hosted on Github

https://github.com/JABBAmodel/JABBA

Winker H, Carvalho F, Kapur M

iabbamodel normindex utils		58ec143 on Nov 24 2021 102 commits	Just Another Bayesian Biomass Assessment		
Example Example R Version1.1_files data man .gitignore Convergence Vignette.Rmd DESCRIPTION	test push normindex utils delete oldfiles Making a JABBA package normindex utils Provide Tutorial add code and figures for manuscript normindex utils	4 months ago 2 months ago 2 years ago 2 years ago 2 years ago 4 years ago 4 years ago 4 years ago 2 months ago 2 months ago 2 months ago	Assessment Releases No releases published Packages		
 JABBAmodel.Rproj Model Diagnostics Vignette Rmd 	Making a JABBA package	2 years ago	No packages published		
NAMESPACE NAMESPACE README.md Tutorial Vignette.Rmd Tutorial Vignette.md	normindex utils Update README.md Update Tutorial Vignette.Rmd Add MD Tutorial File	2 months ago 2 years ago 4 years ago 3 years ago	Contributors 2 Jabbamodel Just Another Bayesian Bio.		
JABBA: Just And	ther Bayesian Biomass	Assessment	Languages R 100.0%		



'a4adiags': Statistical-Catch-Age

R package hosted on Github

https://github.com/flr/a4adiags

Mosqueira I & Winker H

✓ сэчастт on Dec 3, 2021 ③ 41 commi 13 months ag 2 months ag 13 months ag 13 months ag 13 months ag 12 months ag 2 months ag 13 months ag
13 months ag 2 months ag 13 months ag 13 months ag 13 months ag 13 months ag 12 months ag 2 months ag 13 months ag 13 months ag
2 months ag 13 months ag 13 months ag 13 months ag 12 months ag 2 months ag 13 months ag 13 months ag 13 months ag
13 months ag 13 months ag 13 months ag 12 months ag 2 months ag 13 months ag 13 months ag 13 months ag
13 months ag 13 months ag 12 months ag 2 months ag 13 months ag hchxval 12 months ag
13 months ag 12 months ag 2 months ag 13 months ag hchxval 12 months ag
12 months ag 2 months ag 13 months ag hchxval 12 months ag
2 months ag 13 months ag hchxval 12 months ag
13 months ag hchxval 12 months ag
hchxval 12 months ag
13 months ag
2 months ag
13 months ag
12 months ag
13 months ag
13 months ag
13 months ag







'AAP': Statistical-Catch-Age (Aaarts & Poos)

R package hosted on Github

https://github.com/iagomosqueira/AAP

Mosqueira I & Poos JJ

C	iagomosqueira/AAP: Aart>	× +	_ 0 😣
$\leftarrow \rightarrow$	0 6 0 6 0	ೆ https://github.com/iagomosqueira/AAP/ 🖣 🏠 💷 🔕 🌒 🛷	<i>⊡ m</i> § w ≡
()	Search or jump to	Pull requests Issues Marketplace Explore Public Q Pin Q Unwatch 1 • ¥ Fork 0	↓ + • 5 • • • • • • • • • • • • • • • • •
4	Code 🕢 Issues 👫	th Pull requests ⊙ Actions ⊞ Projects □ Wiki ⑦ Security ∠ Insights	
	iagomosqueira v.0.2.8.90	Aarts and Poos Stock that Estimates Bycato	Assessment Model
	R	v.0.2.8.9001 2 months ago	
	data-raw	Missing pin set as NULL 9 months ago	,
	data	Small updates, new Makefile 8 months ago 1 watching	
	inst	Recompiled exe 9 months ago 9 0 forks	
	man	New icon, plus aap.sa 2 months ago	
	tests	No tests yet 2 years ago Releases	
	vignettes	Added Sage.knots 2 years ago No releases published	
D	.Rbuildignore	v0.2.5 2 years ago	
D	CHECKLIST.md	Initial commit 2 years ago	
۵	DESCRIPTION	v.0.2.8.9001 2 months ago	
۵	Makefile	Small updates, new Makefile 8 months ago Publish your first package	
۵	NAMESPACE	xval added 9 months ago	
D	NEWS	Initial commit 2 years ago	
R	NEWS md	Initial commit 2 years and	P



1980

2000

2020

0.2

1960



2015 - 2017 - 2019 2016 - 2018 - 2020 (ref

BTS	SNS
ک۵ ^{۳۵} کورو _ا کو	$ \frac{1}{2}$ $ \frac{1}{2}$ $\frac{1}{2}$ $\frac{1}{2$
	$\sim - \frac{1}{2} \sum_{i=1}^{2} \sum_{j=1}^{2} \sum_{i=1}^{2} \sum_{i$
᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆᠆	$-\frac{1}{2} \omega^2 \omega^2 \omega^2 \omega^2 \omega^2 \omega^2 \omega^2 \omega^2 \omega^2 \omega^2$
	$\mathbf{v} = -\frac{\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}}\mathbf{u}_{\mathbf{v}$
- 	ဗိန္ကအုတ္ပိတိုက္ရာ ဆိုင္နဲ႔အတိုက္တာ ကို က က က က က က က က က က က က က က က က က က
	- <mark>,648,000,000,000,000,000,000,000,000,000,0</mark>
	7
$-\cdots - \mathbf{l}_{\mathbf{c}} \mathbf{l}} \mathbf{l}_{\mathbf{c}} \mathbf{l}_{\mathbf{c}} \mathbf{l}_{\mathbf{c}} \mathbf{l}_$	
⁻	ω
1	
1970 1980 1990 2000 2010 2020	1970 1980 1990 2000 2010 2020

Cookbook in action

- 1. ICCAT Shortfin make 2019 (SS3: Benchmark evaluation)
- 2. GFCM Sicilian Hake 2019 (SS3: Benchmark)
- 3. ICES Eastern Baltic Cod 2019 (SS3: Benchmark)
- 4. ICES Bothnian Herring 2019 (SS3: Benchmark)
- 5. IOTC Skipjack tuna 2020 (SS3: Benchmark Ensemble)
- 6. GFCM Adriatic Sole 2020 (SS3: Benchmark Ensemble)
- 7. ICES Iberian sardine 2020 (SS3: Base-Case development)
- 8. ICCAT Med Swordfish 2020 (JABBA, Benchmark)
- 9. ICCAT S.Atl. Albacore 2020 (JABBA, Benchmark)
- 10. IOTC Yellowfin tuna **2021** (SS3: Benchmark Ensemble)
- 11. ICCAT Bigeye tuna 2021 (SS3 JABBA, mpb: Benchmark Ensemble)
- 12. WCPFC Blue marlin 2021 (SS3, Benchmark Ensemble)
- 13. IOTC Blue shark 2021 (SS3, JABBA: Benchmark)
- 14. IOTC Albacore 2021 (SS3, MSE grid)
- 15. IOTC Swordfish 2021 (SS3, MSE grid)
- 16. SEDAR 2021 Gulf of Mexico Scamp Grouper (SS3, Benchmark)
- 17. ICCAT Mediterranean Albacore 2021 (JABBA, Benchmark)
- 18. Sweden Vendace 2022 (SS3: Benchmark Ensemble)
- 19. ICES Pandalus shrimp 2022 (SS3, Benchmark Ensemble)
- 20. GFCM Adriatic Hake 2022 (SS3: Updated Assessment)
- 21. GFCM Mantis Shrimp 2022 (SS3: Advice Assessment)



A cookbook for using model diagnostics in integrated stock assessments

Felipe Carvalho^{a, g,1}, Henning Winker^{b,1}, Dean Courtney^c, Maia Kapur^d, Laurence Kell^e, Massimiliano Cardinale^f, Michael Schirripa⁸, Toshihide Kitakado^h, Dawit Yemaneⁱ, Kevin R. Piner^j, Mark N. Maunder^{k,1}, Ian Taylor^m, Chantel R. Wetzel^m, Kathryn Doeringⁿ, Kelli F. Johnson^m, Richard D. Methot^m















Cookbook in action (2019)

General Fisheries Commission of the Mediterranean (GFCM) Stock: Sicilian Hake Model: Stock Synthesis

Diagnostics:

- Runs Tests
- Joint-Residual Plots (RMSE)
- ASPM
- R0 Profiling
- Retrospective Analysis
- Hindcast cross-validations



European Commission

Cookbook in action (2019) GFCM Stock: Sicilian Hake Model: Stock Synthesis

ALK

h = 0.99

E1

9.74E-05

1289

2653

Model

AIC

Convergence

Likelihood

ALK

h = 0.88

E2

2.02E-04

1304

2737



Cookbook in action (2019) GFCM Stock: Sicilian Hake Model: Stock Synthesis

Model	E1	E2	E3
Convergence	9.74E-05	2.02E-04	9.41E-05
Likelihood	1289	1304	1761
AIC	2653	2737	3001
RunsTest Index	Р	F	F
Runs Test Length	Р	Р	Р
Restrospective bias	F	Р	F
Forecast bias	F	Р	F
MASE	F	Р	F
Weighting Score	0.4	0.8	0.2



Cookbook in action (2020) **ICCAT** Stock: South Atlantic Albacore Model: JABBA



Hindcast Cross-Validation CTPLL: MASE = 0.6



2010

202

0.2

0.1

0.0

0.1

-0.2

1960

1970

1980

1990

Year

2000

Process Error Deviates





Towards automation of model diagnostics





FAO-GFCM WGSAD 24-28/01/2022



Stock assessment of *Solea solea* in GSA 17 Benchmark update

Masnadi F., Cardinale M., Carbonara P., Donato F., Sabatini L., Pellini G., Scanu M., Dragičević B., Armelloni E. N., Martinelli M., Domenichetti F., Santojanni A., Colella S., Giovanardi O., Raicevich S., Fabi G., Marceta B., Vrgoč N., Isajlovic I., Milone N., Arneri E., Scarcella G.









ALMA MATER STUDIORUM Università di Bologna



Swedish University of Agricultural Sciences



Stock Assessment workflow for Benchmark of Solea solea in GSA 17





Stock Synthesis 3 + Ensemble approach (delta-MVLN) •

Vear

Diagnostic - Convergence & stability (ref run 1)



Diagnostic - Goodness of the fit (ref run 1)



Diagnostic - Consistency (ref run 1)

A cookbook for using model diagnostics in integrated stock assessments

Felipe Carvalho^{41,4,1}, Henning Winker^{b,1}, Dean Courtney^e, Maia Kapur^d, Laurence Kell⁹, Massimiliano Cardinale^f, Michael Schirripa⁸, Toshihide Kitakado^h, Dawit Yemaneⁱ, Kevin R. Piner^j, Mark N. Maunder^{k,1}, Ian Taylor^m, Chantel R. Wetzel^m, Kathryn Doeringⁿ, Kelli F. Johnson^m, Richard D. Methot^m

Jabbamodel / ss3diags Public

Mohn's indices (both ρ M and ρ F) limit value:

- Long lived species: 0.20 ; -0.15
- Short lived species: 0.30 ; -0.22



 Hindcasting (CrossValidation) Diagnostic - Prediction skills (ref run 1)



Model weighting (diagnostic scores)

A cookbook for using model diagnostics in integrated stock assessments

Felipe Carvalho^{n,e,1}, Henning Winker^{h,1}, Dean Courtney^c, Maia Kapur^d, Laurence Kell^e, Massimiliano Cardinale^f, Michael Schirripa⁸, Toshihide Kitakado^h, Dawit Yemane^f, Kevin R. Piner^f, Mark N. Maunder^{k,1}, Ian Taylor^m, Chantel R. Wetzel^m, Kathryn Doeringⁿ, Kelli F. Johnson^m, Richard D. Methot^m $\frac{W(Diagnostics):}{W(Diags 1) + W(Diags 2) + W(Diags 3) \dots + W(Diags N)}$ Num of W(Diags)

Jabbamodel / ss3diags

1. Convergence & stability		Convergenc	e and stability	Goodness of the fit						Consistency			Prediction skills						
		Positive	Uthering			R	un test			Joint-	residuals		Retrospective	analysis		H	indcasting (M	ASE)	
- Positive Hessian	Run name	Hessian	Jittering	Index	lenGNS_ITA	lenTBB_ITA	lenGTR_HRV	lenOTB_ITA	lenSoleMon	Index	Length	Retro_SSB	Forecast_SSB	Retro_F	Forecast_I	Index	SurveyLen	COMfleet	W(Diagnostics)
littoring	Run1	Passed		Passed	Passed	Passed	Passed	Passed	Passed	15.2	3.1	-0.083	-0.070	0.021	0.035	0.726	0.399	0.320	1.00
- Jittering	Run2	Passed		Passed	Passed	Passed	Passed	Passed	Passed	14.7	3.1	-0.058	-0.054	0.026	0.052	0.863	0.363	0.312	1.00
	Run3	Passed		Passed	Passed	Passed	Passed	Passed	Passed	14.9	3.1	-0.061	-0.053	0.016	0.036	0.766	0.382	0.316	1.00 1.00 1.00 1.00 1.00 1.00
2. Goodness of the fit	Run4	Passed		Passed	Passed	Passed	Passed	Passed	Passed	15.4	3.1	-0.074	-0.059	0.018	0.029	0.714	0.407	0.319	1.00
	Run5	Passed		Passed	Passed	Passed	Passed	Passed	Passed	14.7	3.1	-0.040	-0.036	0.014	0.040	0.842	0.370	0.312	1.00
	Run6	Passed		Passed	Passed	Passed	Passed	Passed	Passed	14.9	3.1	-0.036	-0.030	0.008	0.026	0.743	0.334	0.316	1.00
- Joint-residuals	Run7	Passed		Passed	Passed	Passed	Passed	Passed	Passed	15	3.1	-0.078	-0.064	0.034	0.047	0.744	0.410	0.317	1.00
- Runs tests	Run8	Passed		Passed	Passed	Passed	Passed	Passed	Passed	14.4	3.1	-0.037	-0.033	0.017	0.042	0.825	0.377	0.312	1.00
- Kuns tests	Run9	Passed	Barrad	Passed	Passed	Passed	Passed	Passed	Passed	14.7	3.1	-0.054	-0.044	0.021	0.040	0.750	0.396	0.315	1.00
3. Consistency	Run10	Passed	rasseu	Passed	Passed	Passed	Passed	Passed	Passed	21.2	3.3	0.126	0.157	analysis Forecast_ Index Index SurveyLen COMfleet W(Diagnostics) 0.021 0.035 0.726 0.399 0.320 1.00 0.026 0.052 0.863 0.363 0.312 1.00 0.016 0.036 0.766 0.382 0.316 1.00 0.016 0.036 0.766 0.382 0.316 1.00 0.018 0.029 0.714 0.407 0.319 1.00 0.014 0.040 0.842 0.370 0.312 1.00 0.008 0.026 0.743 0.334 0.316 1.00 0.008 0.026 0.743 0.334 0.316 1.00 0.017 0.042 0.825 0.377 0.312 1.00 0.017 0.042 0.825 0.377 0.312 1.00 0.017 0.042 0.826 0.375 1.00 -0.016 -0.072 0.967 0.455 0.351 1.00 -0.					
	Run11	Passed		Passed	Passed	Passed	Passed	Passed	Passed	20.1	3.6	0.013	0.003	-0.009	0.041	1.362	0.450	0.351	0.93
5. Consistency	Run12	Passed		Passed	Passed	Passed	Passed	ALREVIenOTB_ITAIenSoleMonIndexsedPassedPassed15.2sedPassedPassed14.7sedPassedPassed14.9sedPassedPassed15.4sedPassedPassed14.7sedPassedPassed14.7sedPassedPassed14.7sedPassedPassed14.9sedPassedPassed15.5sedPassedPassed15.5sedPassedPassed14.4sedPassedPassed20.1sedPassedPassed20.1sedPassedPassed20.2sedPassedPassed19.4sedPassedPassed20.1sedPassedPassed20.1sedPassedPassed20.1sedPassedPassed20.1sedPassedPassed16.7sedPassedPassedPassedsedPassedPassed16.6sedPassedPassedPassedsedPassedPassed16.6sedPassedPassed16.7	3.4	0.083	0.092	-0.067	-0.037	1.166	0.388	0.367	0.93		
	Run13	Passed		Passed	Passed	Passed	Passed	Passed	Passed	20.2	3.2	0.123	0.162	-0.113	-0.087	0.796	0.472	on skills W(Diagnostics) ig (MASE) W(Diagnostics) 99 0.320 1.00 63 0.312 1.00 63 0.312 1.00 62 0.316 1.00 70 0.319 1.00 70 0.312 1.00 34 0.316 1.00 10 0.317 1.00 96 0.315 1.00 95 0.375 1.00 50 0.351 0.93 88 0.367 0.93 88 0.367 1.00 164 0.344 0.93 163 0.354 1.00 123 0.344 1.00	
- Retrospective analysis	Run14	Passed		Passed	Passed	Passed	Passed	Passed	Passed	19.4	3.4	0.042	0.043	-0.040	-0.001	1.098	0.464	0.344	0.93
- Retrospective analysis	Run15	Passed		Passed	Passed	Passed	Passed	Passed	Passed	20.1	3.2	0.086	0.102	-0.078	-0.044	0.957	0.463	0.354	1.00
	Run16	Passed		Passed	Passed	Passed	Passed	Passed	Passed	16.7	3.1	0.070	0.081	-0.067	-0.024	0.777	0.423	0.346	1.00
4 Due disting abills	Run17	Passed		Passed	Passed	Passed	Passed	Passed	Passed	16.6	3.1	0.049	0.051	-0.045	0.001	0.887	0.421	0.340	1.00
4. Prediction Skills	Run18	Passed		Passed	Passed	Passed	Passed	Passed	Passed	16.7	3.1	0.062	0.070	-0.058	-0.014	0.810	0.423	0.344	1.00

- Hindcasting (CrossValidation)

13.33 x 7.50 in

Ensemble model Results



Commission

Ensemble model Results



Survey Index increasing & -25% landings compare to 2019 !!!

Cookbook in action (2022)

ICES

Stock: Northern shrimp (*Pandalus borealis*)

Model: Stock Synthesis – Length-based 2 Area Model

Max Cardinale, Francesco Masnadi, Alessandro Orio, Mikaela Bergenius, Katja Nören, Christopher Griffiths





ICES Benchmark of Northern shrimp in 3a and 4a east









Joint MASE < 1



The need for simulation testing



The need for more simulation testing

- 1. Identify Causes of Misspecifications
 - Observation process or system dynamics?

2. Specificity/Sensitivity of Diagnostic tests?

False Negative vs False Positive

3. Implications on status/advice?

- Not all misspecifications may be equally consequential
- Is there asymmetric risk among diagnostic?



Original Article

Looking in the rear-view mirror: bias and retrospective patterns in integrated, age-structured stock assessment models

Felipe Hurtado-Ferro^{1*}, Cody S. Szuwalski^{1,2}, Juan L. Valero³, Sean C. Anderson⁴, Curry J. Cunningham¹, Kelli F. Johnson¹, Roberto Licandeo⁵, Carey R. McGilliard^{6‡}, Cole C. Monnahan⁷, Melissa L. Muradian⁷, Kotaro Ono¹, Katyana A. Vert-Pre⁸, Athol R. Whitten¹, and André E. Punt¹



Figure 2. General design of the simulation study.



Figure 4. Experimental design showing the time-varying processes (a - c), fishing patterns (d - f), and time-varying patterns for the processes (g). For the patterns of time variance, only some cases are shown to reduce clutter (solid and dotted lines showold and recent timing, respectively; black and grey lines show gradual and sudden patterns, respectively). See text for further explanation.

European Commission

Causes for retrospective patterns with focus on unaccounted time-varying growth, M and selectivity Rule-of-thumb Mohn's ρ ranges: Low-medium productivity: **-0.15 – 0.2** High productivity: **-0.22 - 0.30**



Full length article

Can diagnostic tests help identify model misspecification in integrated stock assessments?



Felipe Carvalho^{a,b,*}, André E. Punt^c, Yi-Jay Chang^d, Mark N. Maunder^{e,f}, Kevin R. Piner^g

32



Fig. 3. General design of the simulation study.





ASPMs

- Promising candidate to identify misspecifications in the system dynamics (Carvalho et al. 2016)
- Apart from overprecision, ASPMs appear to perform well in simulations to capture the dynamics



Emerging ASPM questions:

- When do correctly specified ASPMs fail in simulations? e.g. combinations of σ_R , ρ and generation time
- How to treat time-varying empirical input data of weight-at-age, maturity-at-age and M-at-age?
- Are there any data-rich age-based examples where the ASPM was consistent with the full model?



Focus Question: Hind casting

Is there a management strategy that relates closely to what we observer (and can then test predictions for)?

Perhaps we should go one step back first and ask:

Are latent management quantities more reliable from models that can be cross-validated using observations that are unknown to the model (i.e. have prediction skill)?

We can test this with simulations!



ICES Journal of Marine Science (2021), 78(6), 2244–2255. https://doi.org/10.1093/icesjms/fsab104

Validation of stock assessment methods: is it me or my model talking?

Laurence T. Kell ^{1,*}, Rishi Sharma², Toshihide Kitakado³, Henning Winker⁴, Iago Mosqueira ⁵, Massimiliano Cardinale ⁶, and Dan Fu⁷



Prediction Skill with Extension: The Vantage Point Approach

ARTICLE

Retrospective forecasting — evaluating performance of stock projections for New England groundfish stocks Elizabeth N. Brooks and Christopher M. Legault

Fig. 1. Illustration of "retrospective forecasting". (a) Retrospective models are created by removing 1, 2, ..., 7 years of data from the full time series (current model). (b) Forecasts are then made from each retrospective model to the end of year 2007.



- Increased prediction residuals (1+2+3,...)
- More contrast?
- How far can we forecast?

Another look at measures of forecast accuracy Rob J. Hyndman^{a,*}, Anne B. Koehler^{b,1} ^a Department of Econometrics and Business Statistics, Monash University, VIC 3800, Australia ^b Department of Decision Sciences and Management Information Systems, Miami University, Oxford, Ohio 45056, USA SSmase(re)

Index	Season	MASE	MAE.PR	MAE.base	MASE.adj	n.eval
Survey	1	1.065	0.302	0.284	1.065	4

Toshihide Kitakado ICCAT Bigeye 2018





Prediction Skill with Extension: The Vantage Point Approach



The options for Simulation testing are infinite

OM generation?

- 1. Condition + Simplify
- 2. Off the shelf OMs
- 3. Coverage + variation in F-trajectory
- 4. N Stock examples
- 5. Time-varying processes (e.g. Selectivity)

Misspecifications EM?

- 1. Natural Mortality M
- 2. Steepness h
- 3. Somatic Growth
- 4. Maturity
- 5. Selectivity
- 6. Catchability
- 7. Weighting: CV, ESS
- 8. Catch history (underreporting)
- 9. Recruitment function
- 10. Spatial

Influences?

- 1. Life history (productivity, generation time)
- 2. Process error (sigmaR, rho)
- 3. Time series length (observation horizon)
- 4. Sampling (e.g. ESS, CV, AR1)
- 5. Stock depletion (left or right of Bmsy?)
- 6. Contrast in exploitation history
- 7. Data (survey vs CPUE; age vs length)

Model complexity?

- 1. N Indices
- 2. N Size/Age comps setup
- 3. Historical vs Recent time series
- 4. N Fleets
- 5. Sex-structured?
- 6. Seasonal?
- 7. Spatial?



Proposal: Start with simplified "off the shelf" OMs for the system dynamics (1 Fleet)

- Cover more life histories, exploitation patterns and process error dynamics, data-rich vs data moderate
- Minimize run time to 1-3 min (1 OM = ~ 8 EMs x 5 hindcasts x 100 iterations = 1 hour on 75 cores)
- Easier interpretation of main effects



Strawman: Simulation design



Sensitivity: the ability of a test to correctly identify a misspecification **Specificity**: the ability of a test to identify a correctly specified model (Model)





Evaluating impacts on status/advice?

A simple simulation trial





https://github.com/henning-winker/spmpriors

stk = flmvn_traits(Genus="Sparus",Species="aurata",M=c(0.43,0.1),h=c(0.6,0.9),Plot=T)

European Commission

 Received: 17 October 2018
 Revised: 21 October 2019
 Accepted: 31 October 2019

 DOI: 10.1111/faf.12427

ORIGINAL ARTICLE

FISH and FISHERIE

Predicting recruitment density dependence and intrinsic growth rate for all fishes worldwide using a data-integrated life-history model

James T. Thorson^{1,2}

- SSR BevHold with h = 0.87
- M = 0.43



Simulation Run with LIME: Sparus aurata

Simulation Run with LIME: Sparus aurata



Time

Time





















SSplotHCxval()







SSplotHCxval()









Plausibility

			Model
Diagnostics	Data	BevHolt	Ricker
Joint RMSE	Survey Index	18.5	26.2
	Length Comps	4	5.3
Runs Test	Survey Index	Passed	Passed
	Survey Mean Length	Passed	Passed
	Fishery Mean Length	Passed	Passed
Mohn's Rho	SSB	-0.04	-0.02
Forecast Bias	SSB	-0.07	-0.11
HCxval MASE	Survey Index	0.84	1.32
	Survey Mean Length	0.88	1.38
	Fishery Mean Length	1.04	1.93

Diagnostics to add: ASPM class

Statistics to note: LL, AIC, convergence













Model:BevHolt







Remarks on model plausibility

"The likelihood of a scenario considered in simulation trials representing reality, relative to other scenarios also under consideration. Plausibility may be estimated formally based on **some statistical approach**, or specified based on expert judgement, and can be used to weight performance statistics when integrating over results for different scenarios. [...] The aim of conditioning is to select those OMs consistent with the data and reject OMs that do not fit these data satisfactorily and, as such, are considered implausible."

RFMO MSE Glossary 2018

"The **SC AGREED** that it is useful to develop a set of generic criteria for model plausibility, utilising best practice in evaluating model **convergence** and **data fits**, **retrospective pattern** and **forecast bias**, and **prediction skill**, as well as other potential aspects of model diagnostics. The **SC NOTED** that establishing such guidance and criteria can help ensure that the stock assessments are transparent and comprehensive and allow stakeholders to have a good grasp of the scientific process. Stock specific plausibility criteria can also be considered to evaluate if assessment results are consistent with prior knowledge about the exploitation history and population biology.

Report of IOTC Scientific Committee 2020

"A highly **plausible** scenario is one that fits prior knowledge, with many sources of corroboration, without the complexity of explanation, and with minimal conjecture (Connell, 2006). Plausibility may be determined formally, based on a statistical approach to determine whether a **system equivalent to the model generated the data** or based on expert judgement."

Laurence T. Kell, Coilin Minto, and Hans D. Gerritsen. IJMS. Evaluation of the skill of length-based indicators to identify stock status and trends, accepted 2022.



