

# Alternative estimates of the hake species splitting model biased towards observer data

By OLRAC SPS  
Silvermine House  
Steenberg Office Park  
Tokai 7945

16 September 2013

## Background

At a DWG meeting held on 27 February 2013 OLRAC SPS was asked to carry out further work on the hake species splitting model, in particular to produce a variant which was biased towards the observer data. Further to this, at a DWG meeting of 28<sup>th</sup> August 2013, there was a request to present the standard errors of the GLM estimates presented in FISHERIES/2013/AUG/SWG-DEM/20. This document therefore presents results on the differences in the species splitting GLM between the observer and survey data, as well as the standard errors associated with the GLM parameter estimates. However, since these results raise some questions about the agreement between different approaches, further work is carried out here, revisiting the merits of year and spatial effects in the species splitting model.

## Methods

The methods adopted are as follow:

The model utilises a GLM with a binomial distribution and a logit link function. Model effects are additive in logit space, via an equation of the following form:

$$P = \frac{e^{\Psi}}{1 + e^{\Psi}} = \frac{1}{1 + e^{-\Psi}}$$

with  $\Psi = \mu + \lambda_{sizeclass} + \gamma_{depth} + \phi_{latitude} + \alpha_y$

where: P is the proportion of *Merluccius paradoxus*;

$\mu$  is the intercept;

$\alpha_y$  is the year parameter for year y;

$\phi_{latitude}$  is the latitude (West Coast) and longitude (East Coast) parameter;

$\lambda_{sizeclass}$  is the size class specific parameter;

$\gamma$  is the covariate parameter for depth;

The target variable for this GLM (dependent variable) is the number of *Merluccius paradoxus* which occur in a sample, and the number of trials is the total number of hake in the sample.

The above is the GLM model used in FISHERIES/2013/FEB/SWG-DEM/12. The GLM was run for

- Observer data only

- Survey data only,
- A combination of Observer and Survey data.

In the case of the combined observer/survey data, the datasets were combined using a so-called 'balance reduction' method, so that there are an equal number of records from the observer and the survey datasets. Balance reduction involves reducing the number of records for the more plentiful data source by a random selection process.

After fitting the GLM to the three different data combinations, three methods were used to determine the importance of the model effects as well as the standard errors for the parameter estimates. These are:

1. **Method 1a:** Use of a cross validation approach where the prediction errors for a model run on a random 50% selection of the data is used to predict the other 50% of outcomes. Method 1a uses the **proportion** of *M. paradoxus* in samples as the outcome to be predicted.
2. **Method 1b:** Also uses a cross validation approach where the prediction errors for a model run on a random 50% selection of the data is used to predict the other 50% of outcomes. Method 1a uses the **number** of *M. paradoxus* in samples as the outcome to be predicted.
3. **Method 2:** Use of a group jackknife approach where the data are partitioned randomly into 12 exclusive groups, and each group is left out. The acceptability of models is based on standard errors of the effects that are considered (depth, size, year, latitude, longitude).

## Results and Discussion

The original sample sizes in terms of the total number of fish sampled is given as follows (prior to balance reduction):

- West Coast: Research survey data: n = 10193 fish
- West Coast: OROP+SADSTIA observer data: n = 10248 fish
- South Coast: Research survey data: n = 8531 fish
- South Coast: OROP+SADSTIA observer data: n = 2552 fish

Table 1 shows the results using Methods 1a and b.

Table 2 and 4 show the GLM model parameter estimates as well as results from the Jackknife (Method 2) with the inclusion of various model effects for the east coast for the different data sources.

Table 5 and 7 show the GLM model parameter estimates as well as results from the Jackknife with the inclusion of various model effects for the west coast for the different data sources.

Figure 1 and 2 show the GLM predicted proportion of *M. paradoxus* at different depths for the different size classes on the west and east coast respectively.

This work has highlighted that different approaches to model selection are not necessarily compatible with each other. Method 2, for example, leads to different conclusions about an acceptable model structure compared to Method 1a. The jackknife results suggest that a year effect and a longitude/latitude effect is statistically significant in some cases (involving a particular coast and combination of observer and survey data). Method 1a favours the simpler models involving only depth and size effects. Method 1b is an attempt to see whether the two can be reconciled by using cross validation at the level of number of *M. paradoxus* predicted per sample. Note that Method 1a would be biased in favour of models which avoid errors in the larger proportions of *M. paradoxus*.

As noted in previous submissions, Method 1a supports simpler models with the inclusion of depth and size class as model parameters.

Considering the S.E. associated with the jackknife estimates it is clear that the observer and survey data are significantly different when considered for the species splitting algorithm.

## References

- Gaylard J.D. and M.O. Bergh. 2003. An investigation into the procedure used to split commercial catches of hake on the South African South Coast into *Merluccius paradoxus* and *Merluccius capensis*. BEN/JAN04/SAH2b
- Gaylard J.D. and M.O. Bergh .2004a. A size-dependent species splitting mechanism applied to hake catches off the South African West Coast. WG/08/04/D:H:13
- Gaylard J.D. and M.O. Bergh M. 2004b. A species splitting mechanism for application to the commercial hake catch data 1978 to 2003. Marine Coastal Management Document WG/09/04/D: H: 21. 8 pp.
- Gaylard J.D. and M.O. Bergh. 2007. Further comparison of hake species splits from observer data with the survey-generated splitting algorithm.
- Gaylard J.D. and M.O. Bergh. 2009. Update of the hake species split models in the light of more recent survey data and a revision of the large/medium/small size classification. MCM/2009/NOVEMBER/SWG-DEM/...
- Leslie R.W. B. Rose and J. Scholte. 1998. Hake grading by Irvin & Johnson and by Sea Harvest. WG/01/98/D:H:03
- OLRAC SPS. 2013a. Addendum #1 to “[FISHERIES/2013/FEB/SWG-DEM/12] A further update of the hake species splitting model. OL-SPS.” – Updated parameter estimates for the hake species split model.
- OLRAC SPS. 2013b. [FISHERIES/2013/FEB/SWG-DEM/12]. A further update of the hake species splitting model. February 2013.
- OLRAC SPS. 2013c. Addendum #2 to “[FISHERIES/2013/FEB/SWG-DEM/12] A further update of the hake species splitting model. OL-SPS.”

**Table 1. Results of the cross validation approach where the prediction method is either Method 1a (proportion of *Merluccius paradoxus* in the sample) or Method 1b (number of *Merluccius paradoxus* in the sample) to establish the importance of different GLM parameters on the species splitting algorithm.**

<i>Data Source</i>	Method 1a		Method 1b
	<i>Survey</i>	<i>Observer</i>	<i>Observer</i>
West Coast			
depth	0.117	0.061	411.4
depth lat	0.117	0.061	404.3
depth size class	0.05	0.05	298.2
depth year	0.117	0.061	402.1
East Coast			
depth	0.027	0.064	342
depth long	0.027	0.064	314.8
depth size class	0.018	0.056	309.1
depth year	0.027	0.064	277.5

**Table 2. Model parameter estimates of the species split GLM (depth and Size Class stratified) for the East Coast on the SADSTIA/OROP data, the survey data and the combined data. The Jackknife GLM estimates and SE as described in Method 2 is also given.**

Stratum GLM parameters Data Source	East Coast								
	Size Class & Depth								
	<i>Both</i>			<i>Survey</i>			<i>Observer</i>		
	<u>Complete GLM</u>	<u>Mean of Jackknife</u>	<u>SE Jackknife</u>	<u>Complete GLM</u>	<u>Mean of Jackknife</u>	<u>SE Jackknife</u>	<u>Complete GLM</u>	<u>Mean of Jackknife</u>	<u>SE Jackknife</u>
Intercept	-6.05	-6.051	0.163	-23.178	-23.184	1.794	-1.866	-1.866	0.101
SizeClass 2	2.292	2.293	0.119	10.998	10.992	1.433	2.889	2.889	0.086
Sizeclass 3	1.437	1.437	0.106	7.362	7.352	1.405	1.356	1.356	0.090
SizeClass 4	0	0	0	0	0	0	0	0	0
depth	0.022	0.022	0.001	0.073	0.073	0.004	0.01	0.010	0.001

**Table 3. Model parameter estimates of the species split GLM (depth, Size Class and year stratified) for the East Coast on the SADSTIA/OROP data, the survey data and the combined data. The Jackknife GLM estimates and SE as described in Method 2 is also given.**

Stratum GLM parameters Data Source	East Coast								
	Size Class & Depth & Year								
	<i>Both</i>			<i>Survey</i>			<i>Observer</i>		
	<u>Complete GLM</u>	<u>Mean of Jackknife</u>	<u>SE Jackknife</u>	<u>Complete GLM</u>	<u>Mean of Jackknife</u>	<u>SE Jackknife</u>	<u>Complete GLM</u>	<u>Mean of Jackknife</u>	<u>SE Jackknife</u>
Intercept	-4.989	-4.986	0.223	-38.464	-38.663	2.263	-1.900	-1.896	0.235
SizeClass 2	2.745	2.745	0.150	11.181	11.179	1.440	2.806	2.807	0.100
Sizeclass 3	1.597	1.596	0.103	7.436	7.430	1.406	1.356	1.356	0.084
SizeClass 4	0	0	0	0	0	0	0	0	0
depth	0.022	0.022	1.151E-17	0.073	0.073	0.004	0.010	0.010	5.753E-18
1986	-2.239	-2.247	0.459	14.549	14.734	1.378			
1987	-2.207	-2.213	0.372	15.668	15.868	1.224			
1988	-2.928	-2.935	0.479	14.653	14.807	1.515			
1989	-4.737	-4.768	0.921	13.378	13.534	1.428			
1990	-2.594	-2.605	0.387	15.02	15.207	1.154			
1991	<b>-3.142</b>	<b>-3.315</b>	<b>2.560</b>	14.905	14.617	5.116			
1992	-3.661	-3.673	0.343	13.719	13.881	1.592			
1993	-2.06	-2.068	0.307	15.483	15.666	1.339			
1994	-2.282	-2.288	0.270	15.201	15.362	1.432			
1995	-3.394	-3.400	0.354	14.318	14.476	1.542			
1996	-2.087	-2.094	0.281	16.219	16.396	1.263			
1997	-1.525	-1.533	0.446	16.302	16.484	1.384			
1999	-2.007	-2.013	0.249	15.672	15.851	1.626			
2000	-3.193	-3.202	0.411	13.732	13.914	1.160			
2001	-2.916	-2.924	0.297	13.807	13.987	1.348			
2002	2.957	2.946	0.468	15.944	16.095	1.504	<b>0.621</b>	<b>0.608</b>	<b>0.468</b>
2003	<b>-0.351</b>	<b>-0.355</b>	<b>0.290</b>				<b>-0.331</b>	<b>-0.335</b>	<b>0.237</b>
2004	-1.252	-1.256	0.201	14.706	14.872	1.392	<b>0.144</b>	<b>0.139</b>	<b>0.273</b>
2005	<b>-0.416</b>	<b>-0.420</b>	<b>0.313</b>	15.51	15.675	1.005	<b>0.066</b>	<b>0.061</b>	<b>0.218</b>
2006	-0.457	-0.460	0.223	15.317	15.491	1.354	<b>0.301</b>	<b>0.295</b>	<b>0.276</b>
2007	-1.991	-1.995	0.283	14.79	14.968	1.211	<b>0.510</b>	<b>0.512</b>	<b>0.493</b>
2008	-1.867	-1.871	0.259	14.87	15.050	1.258	<b>0.035</b>	<b>0.030</b>	<b>0.235</b>
2009	-1.337	-1.341	0.290	14.808	14.989	1.380	<b>-0.232</b>	<b>-0.237</b>	<b>0.237</b>
2010	-0.863	-0.866	0.259	14.101	14.295	1.287	<b>0.313</b>	<b>0.308</b>	<b>0.322</b>
2011	-2.494	-2.500	0.308	13.98	14.163	1.167	-1.632	-1.638	0.371
2012	0	0	0	0	0	0	0	0	0

**Table 4. Model parameter estimates of the species split GLM (depth, Size Class, year and longitude stratified) for the East Coast on the SADSTIA/OROP data, the survey data and the combined data. The Jackknife GLM estimates and SE as described in Method 2 is also given.**

Stratum GLM parameters Data Source	East Coast								
	Size Class & Depth & Year & Long								
	<i>Both</i>			<i>Survey</i>			<i>Observer</i>		
	Complete GLM	Mean of Jackknife	SE Jackknife	Complete GLM	Mean of Jackknife	SE Jackknife	Complete GLM	Mean of Jackknife	SE Jackknife
Intercept	-4.57	-4.579	0.404	-15.112	-15.131	0.895	1.305	1.297	0.339
SizeClass 4	-2.88	-2.879	0.144	-11.658	-11.652	1.254	-2.837	-2.839	0.103
Sizeclass 3	-1.32	-1.315	0.122	-4.569	-4.570	0.455	-1.494	-1.497	0.119
SizeClass 2	0	0	0	0	0	0	0	0	0
depth	0.02	0.020	1.151E-17	-5.106	-5.367	2.758	0.009	0.009	5.753E-18
2012	2.13	2.137	0.462	0.078	0.078	0.005	<b>-0.544</b>	<b>-0.531</b>	<b>0.419</b>
2011	<b>-0.45</b>	<b>-0.441</b>	<b>0.405</b>	<b>-0.724</b>	<b>-0.716</b>	<b>0.668</b>	-2.357	-2.348	0.566
2010	1.04	1.045	0.304	<b>0.203</b>	<b>0.209</b>	<b>0.915</b>	<b>-0.470</b>	<b>-0.460</b>	<b>0.473</b>
2009	0.70	0.701	0.349	<b>0.241</b>	<b>0.254</b>	<b>0.765</b>	-0.870	-0.861	0.418
2008	<b>0.29</b>	<b>0.296</b>	<b>0.402</b>	<b>0.693</b>	<b>0.701</b>	<b>0.410</b>	<b>-0.618</b>	<b>-0.609</b>	<b>0.460</b>
2007	<b>0.43</b>	<b>0.430</b>	<b>0.350</b>	<b>0.354</b>	<b>0.350</b>	<b>0.466</b>	<b>0.174</b>	<b>0.190</b>	<b>0.748</b>
2006	1.67	1.677	0.332	1.138	1.134	0.469	<b>-0.232</b>	<b>-0.224</b>	<b>0.416</b>
2005	1.90	1.908	0.469	1.465	1.458	0.647	<b>-0.457</b>	<b>-0.448</b>	<b>0.424</b>
2004	0.79	0.795	0.359	<b>0.246</b>	<b>0.230</b>	<b>0.423</b>	<b>-0.509</b>	<b>-0.500</b>	<b>0.441</b>
2003	1.99	1.999	0.413	1.844	1.827	0.554	-0.816	-0.806	0.402
2002	5.55	5.545	0.661	<b>0.175</b>	<b>0.168</b>	<b>0.573</b>	0	0	0
2001	<b>-0.41</b>	<b>-0.413</b>	<b>0.590</b>						
2000	<b>-0.79</b>	<b>-0.791</b>	<b>0.400</b>	<b>-0.213</b>	<b>-0.208</b>	<b>0.440</b>			
1999	<b>0.28</b>	<b>0.282</b>	<b>0.427</b>	1.524	1.514	0.744			
1997	<b>0.81</b>	<b>0.805</b>	<b>0.553</b>	1.949	1.963	0.679			
1996	<b>0.19</b>	<b>0.187</b>	<b>0.480</b>	1.835	1.826	0.777			
1995	-0.99	-0.992	0.482	<b>0.679</b>	<b>0.661</b>	<b>0.878</b>			
1994	<b>-0.09</b>	<b>-0.082</b>	<b>0.448</b>	<b>0.431</b>	<b>0.430</b>	<b>0.364</b>			
1993	<b>0.23</b>	<b>0.231</b>	<b>0.527</b>	1.765	1.765	0.534			
1992	-1.65	-1.655	0.474	<b>-1.086</b>	<b>-1.103</b>	<b>0.852</b>			
1991	<b>-1.24</b>	<b>-1.376</b>	<b>2.169</b>	<b>-0.665</b>	<b>-1.008</b>	<b>3.241</b>			
1990	<b>-0.71</b>	<b>-0.713</b>	<b>0.493</b>	<b>-0.009</b>	<b>-0.003</b>	<b>0.672</b>			
1989	-2.51	-2.527	1.004	<b>-0.192</b>	<b>-0.205</b>	<b>0.920</b>			
1988	<b>-0.77</b>	<b>-0.766</b>	<b>0.485</b>	<b>-0.142</b>	<b>-0.173</b>	<b>0.963</b>			
1987	<b>-0.28</b>	<b>-0.272</b>	<b>0.355</b>	<b>0.306</b>	<b>0.332</b>	<b>0.534</b>			
1986	0	0	0	0	0	0			
27	-21.58	-21.437	1.981	-12.798	-12.527	2.785			
26	<b>0.03</b>	<b>0.034</b>	<b>0.156</b>	1.104	1.095	0.511	<b>-0.180</b>	<b>-0.180</b>	<b>0.117</b>
25	1.22	1.221	0.111	1.699	1.702	0.442	0.580	0.580	0.101
24	1.52	1.519	0.112	3.079	3.082	0.510	0.783	0.782	0.075
23	1.15	1.149	0.121	3.656	3.661	0.328	0.320	0.319	0.121
22	<b>-0.09</b>	<b>-0.091</b>	<b>0.201</b>	2.72	2.714	0.373	<b>0.530</b>	<b>0.528</b>	<b>0.498</b>
21	0	0	0	0	0	0	0	0	0

**Table 5. Model parameter estimates of the species split GLM (depth and Size Class stratified) for the West Coast on the SADSTIA/OROP data, the survey data and the combined data. The Jackknife GLM estimates and SE as described in Method 2 is also given.**

Stratum GLM parameters Data Source	West Coast								
	Size Class & Depth								
	<i>Both</i>			<i>Survey</i>			<i>Observer</i>		
	<u>Complete GLM</u>	<u>Mean of Jackknife</u>	<u>SE Jackknife</u>	<u>Complete GLM</u>	<u>Mean of Jackknife</u>	<u>SE Jackknife</u>	<u>Complete GLM</u>	<u>Mean of Jackknife</u>	<u>SE Jackknife</u>
Intercept	-5.743	-5.743	0.120	-12.852	-12.853	0.293	-2.687	-2.687	0.071
SizeClass 2	2.688	2.689	0.041	5.783	5.784	0.155	3.058	3.058	0.053
Sizeclass 3	1.457	1.457	0.052	2.045	2.044	0.141	1.41	1.410	0.045
SizeClass 4	0	0	0	0	0	0	0	0	0
depth	0.019	0.019	1.151E- 17	0.037	0.037	0.001	0.01	0.010	5.753E- 18

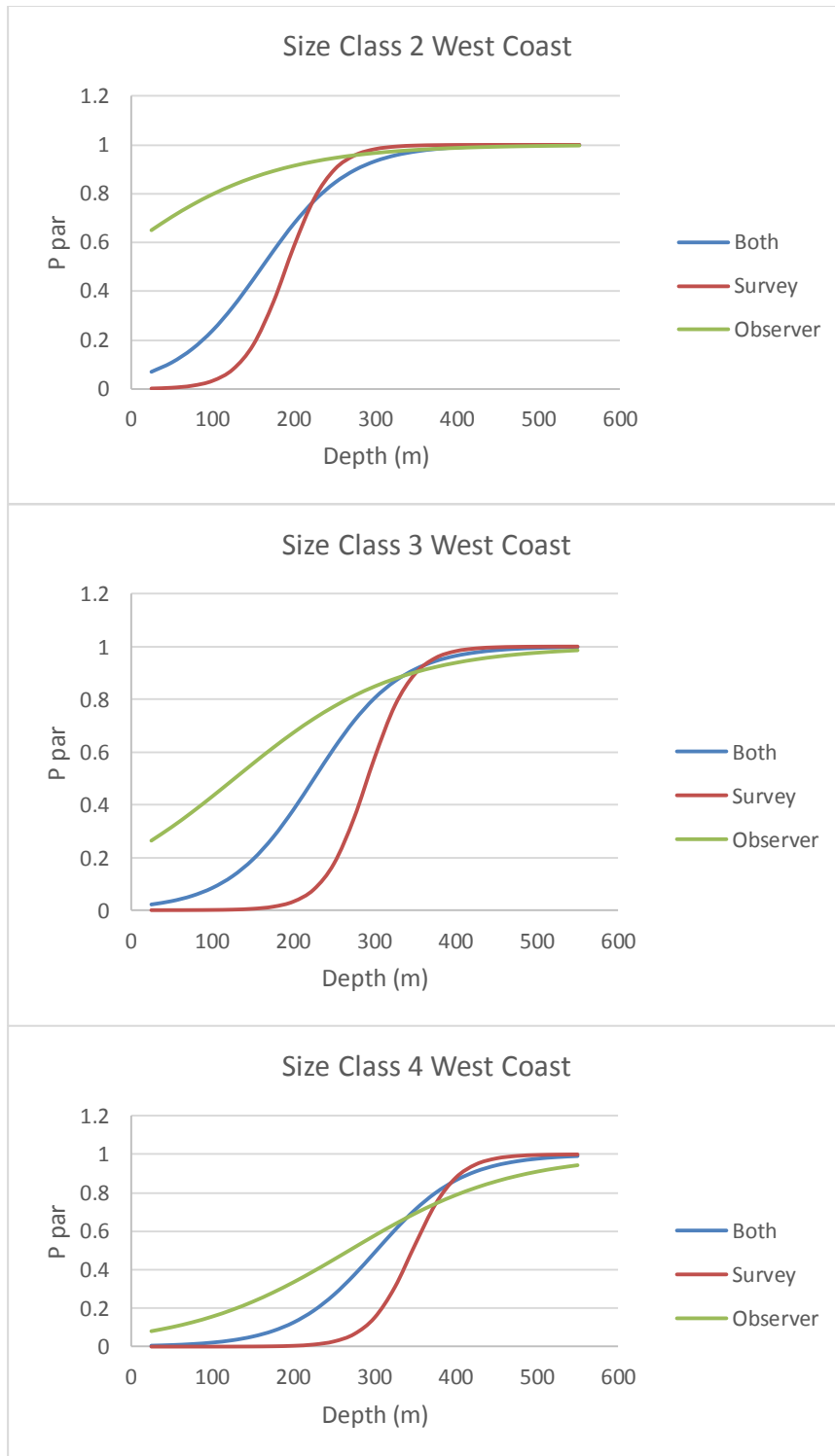
**Table 6. Model parameter estimates of the species split GLM (depth, Size Class and year stratified) for the West Coast on the SADSTIA/OROP data, the survey data and the combined data. The Jackknife GLM estimates and SE as described in Method 2 is also given.**

Stratum GLM parameters Data Source	West Coast								
	Size Class & Depth & Year								
	<i>Both</i>			<i>Survey</i>			<i>Observer</i>		
	<u>Complete</u> GLM	<u>Mean of</u> Jackknife	<u>SE</u> Jackknife	<u>Complete</u> GLM	<u>Mean of</u> Jackknife	<u>SE</u> Jackknife	<u>Complete</u> GLM	<u>Mean of</u> Jackknife	<u>SE</u> Jackknife
Intercept	-5.698	-5.698	0.179	-13.498	-13.506	0.243	-2.937	-2.937	0.111
SizeClass 2	2.945	2.945	0.038	6.407	6.412	0.138	3.015	3.015	0.052
Sizeclass 3	1.485	1.485	0.048	2.321	2.322	0.168	1.393	1.394	0.043
SizeClass 4	0	0	0	0	0	0	0	0	0
depth	0.018	0.018	0.001	0.042	0.042	0.002	0.010	0.010	5.753E-18
1984	-0.97	-0.970	0.169	-1.613	-1.614	0.278			
1985	-1.275	-1.276	0.149	-1.965	-1.965	0.297			
1986	-1.188	-1.188	0.172	-2.042	-2.043	0.208			
1987	-0.563	-0.564	0.128	-0.858	-0.859	0.194			
1988	-0.401	-0.401	0.148	-0.655	-0.655	0.315			
1989	-1.112	-1.113	0.112	-2.009	-2.011	0.456			
1990	-1.091	-1.091	0.090	-1.733	-1.734	0.266			
1991	<b>-0.298</b>	<b>-0.298</b>	<b>0.219</b>	<b>-0.378</b>	<b>-0.378</b>	<b>0.372</b>			
1992	-0.667	-0.667	0.154	-1.034	-1.035	0.261			
1993	-0.284	-0.284	0.123	-0.708	-0.711	0.345			
1994	<b>-0.045</b>	<b>-0.045</b>	<b>0.149</b>	<b>-0.092</b>	<b>-0.091</b>	<b>0.197</b>			
1995	-0.571	-0.571	0.126	-0.57	-0.571	0.164			
1996	<b>-0.051</b>	<b>-0.051</b>	<b>0.188</b>	<b>0.01</b>	<b>0.010</b>	<b>0.253</b>			
1997	<b>-0.384</b>	<b>-0.385</b>	<b>0.199</b>	<b>-0.305</b>	<b>-0.306</b>	<b>0.273</b>			
1999	<b>-0.242</b>	<b>-0.242</b>	<b>0.144</b>	<b>-0.121</b>	<b>-0.122</b>	<b>0.213</b>			
2000	-1.352	-1.350	0.208	-1.983	-1.986	0.241			
2001	-1.033	-1.034	0.149	-1.944	-1.945	0.235			
2002	0.311	0.311	0.151	<b>-0.419</b>	<b>-0.420</b>	<b>0.212</b>	<b>0.240</b>	<b>0.240</b>	<b>0.192</b>
2003	0.428	0.428	0.112	<b>0.309</b>	<b>0.309</b>	<b>0.234</b>	<b>0.029</b>	<b>0.029</b>	<b>0.109</b>
2004	0.27	0.270	0.126	-1.541	-1.542	0.234	0.529	0.529	0.121
2005	0.685	0.685	0.116	<b>-0.031</b>	<b>-0.032</b>	<b>0.180</b>	0.511	0.511	0.105
2006	0.666	0.666	0.082	<b>0.511</b>	<b>0.509</b>	<b>0.317</b>	0.343	0.342	0.079
2007	0.429	0.429	0.082	-1.465	-1.465	0.318	0.420	0.420	0.122
2008	0.281	0.281	0.072	<b>-0.325</b>	<b>-0.327</b>	<b>0.363</b>	<b>0.114</b>	<b>0.114</b>	<b>0.097</b>
2009	0.547	0.547	0.100	<b>-0.252</b>	<b>-0.252</b>	<b>0.273</b>	0.348	0.348	0.070
2010	<b>-0.072</b>	<b>-0.072</b>	<b>0.114</b>	<b>-0.065</b>	<b>-0.065</b>	<b>0.245</b>	<b>0.007</b>	<b>0.007</b>	<b>0.101</b>
2011	-0.199	-0.199	0.093	<b>0.349</b>	<b>0.349</b>	<b>0.348</b>	-0.446	-0.446	0.075
2012	0	0	0	0	0	0	0	0	0

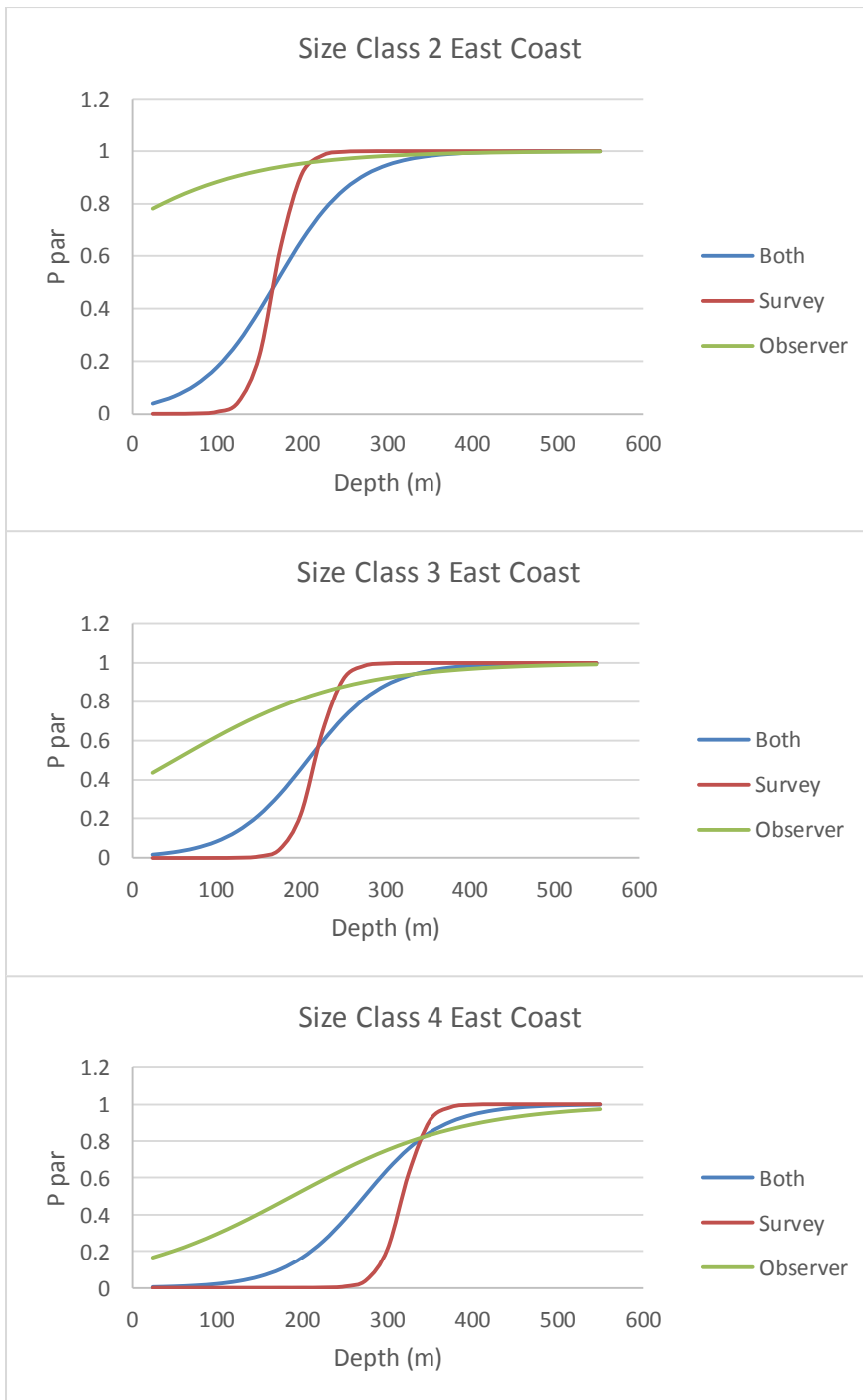


**Table 7. Model parameter estimates of the species split GLM (depth, Size Class, year and latitude stratified) for the West Coast on the SADSTIA/OROP data, the survey data and the combined data. The Jackknife GLM estimates and SE as described in Method 2 is also given.**

Stratum GLM parameters Data Source	West Coast								
	Size Class & Depth & Year & Lat								
	<i>Both</i>			<i>Survey</i>			<i>Observer</i>		
	Complete GLM	Mean of Jackknife	SE Jackknife	Complete GLM	Mean of Jackknife	SE Jackknife	Complete GLM	Mean of Jackknife	SE Jackknife
Intercept	-4.32	-4.322	0.179	-8.601	-8.607	0.266	4.138	4.066	0.562
SizeClass 4	-3.12	-3.117	0.051	-6.479	-6.485	0.158	-2.988	-2.986	0.048
Sizeclass 3	-1.59	-1.591	0.041	-4.158	-4.162	0.206	-1.617	-1.616	0.065
SizeClass 2	0	0	0	0	0	0	0	0	0
depth	0.02	0.018	1.151E-17	0.043	0.043	0.001	0.010	0.010	5.486E-18
2012	0.92	0.920	0.165	1.552	1.553	0.290	<b>-0.067</b>	<b>-0.068</b>	<b>0.194</b>
2011	0.70	0.701	0.130	1.925	1.925	0.207	-0.518	-0.521	0.154
2010	0.85	0.849	0.131	1.418	1.419	0.197	<b>-0.117</b>	<b>-0.116</b>	<b>0.196</b>
2009	1.41	1.406	0.137	1.196	1.197	0.335	<b>0.162</b>	<b>0.162</b>	<b>0.167</b>
2008	1.12	1.120	0.113	1.167	1.166	0.316	<b>-0.033</b>	<b>-0.033</b>	<b>0.164</b>
2007	1.18	1.180	0.137	<b>0.059</b>	<b>0.059</b>	<b>0.238</b>	0.341	0.335	0.163
2006	1.42	1.423	0.104	1.955	1.953	0.221	<b>0.217</b>	<b>0.217</b>	<b>0.157</b>
2005	1.56	1.564	0.101	1.468	1.468	0.242	<b>0.253</b>	<b>0.253</b>	<b>0.163</b>
2004	1.05	1.045	0.089	<b>-0.034</b>	<b>-0.034</b>	<b>0.195</b>	0.379	0.379	0.174
2003	1.26	1.261	0.134	1.922	1.923	0.163	<b>-0.141</b>	<b>-0.139</b>	<b>0.126</b>
2002	1.13	1.128	0.156	1.105	1.105	0.242	0	0	0
2001	<b>-0.03</b>	<b>-0.032</b>	<b>0.195</b>	<b>-0.278</b>	<b>-0.280</b>	<b>0.193</b>			
2000	<b>-0.38</b>	<b>-0.382</b>	<b>0.192</b>	<b>-0.418</b>	<b>-0.421</b>	<b>0.364</b>			
1999	0.72	0.715	0.162	1.436	1.437	0.221			
1997	0.56	0.559	0.232	1.165	1.165	0.146			
1996	0.99	0.992	0.186	1.65	1.651	0.183			
1995	0.41	0.409	0.137	0.982	0.981	0.212			
1994	0.94	0.939	0.191	1.462	1.462	0.155			
1993	0.73	0.725	0.189	0.857	0.855	0.287			
1992	<b>0.28</b>	<b>0.283</b>	<b>0.200</b>	0.496	0.494	0.234			
1991	0.69	0.685	0.196	1.193	1.195	0.222			
1990	<b>-0.14</b>	<b>-0.135</b>	<b>0.153</b>	<b>-0.22</b>	<b>-0.221</b>	<b>0.269</b>			
1989	<b>-0.14</b>	<b>-0.139</b>	<b>0.169</b>	<b>-0.494</b>	<b>-0.495</b>	<b>0.363</b>			
1988	0.59	0.587	0.197	0.893	0.894	0.209			
1987	0.44	0.441	0.124	0.666	0.666	0.198			
1986	<b>-0.22</b>	<b>-0.223</b>	<b>0.160</b>	<b>-0.494</b>	<b>-0.494</b>	<b>0.275</b>			
1985	<b>-0.33</b>	<b>-0.329</b>	<b>0.189</b>	<b>-0.389</b>	<b>-0.389</b>	<b>0.224</b>			
1984	0	0	0	0	0	0			
36	1.34	1.339	0.067	0.701	0.704	0.235	-3.836	-3.764	0.525
35	0.42	0.419	0.067	-0.586	-0.585	0.115	-4.155	-4.083	0.540
34	0.35	0.353	0.080	-0.859	-0.860	0.189	-3.969	-3.898	0.550
33	0.47	0.473	0.057	-0.324	-0.324	0.154	-3.818	-3.748	0.558
32	0.72	0.717	0.051	<b>-0.291</b>	<b>-0.291</b>	<b>0.157</b>	-3.191	-3.120	0.569
31	0.53	0.533	0.049	<b>-0.063</b>	<b>-0.062</b>	<b>0.114</b>	-2.859	-2.786	0.545
30	0.85	0.851	0.069	0.28	0.281	0.134	-3.624	-3.552	0.524
29	0	0	0	0	0	0	0	0	0



**Figure 1.** The average predicted proportion of *M. paradoxus* at different depths on the West Coast for the survey data, SADSTIA/OROP data and a combination of both data sources with depth and size class as model effects.



**Figure 2.** The average predicted proportion of *M. paradoxus* for different depths on the West Coast for the survey data, SADSTIA/OROP data and a combination of both data sources with depth and size class as model effects.