

# SEAK Pink Salmon 2023 Forecast Process

Sara Miller

September 12, 2022

## Objective

To forecast the Southeast Alaska (SEAK) pink salmon commercial harvest in 2023.

## Executive Summary

Forecasts were developed using an approach originally described in Wertheimer et al. (2006), and modified in Orsi et al. (2016) and Murphy et al. (2019). We used a similar approach to Murphy et al. (2019), but assumed a log-normal error. This approach is based on a multiple regression model with juvenile pink salmon catch-per-unit-effort (CPUE) and temperature data from the Southeast Alaska Coastal Monitoring Survey (SECM; Piston et al. 2021) or from satellite sea surface temperature data (SST and SST Anomaly, NOAA Global Coral Bleaching Monitoring, ([https://coastwatch.pfeg.noaa.gov/erddap/griddap/NOAA\\_DHW\\_monthly.html](https://coastwatch.pfeg.noaa.gov/erddap/griddap/NOAA_DHW_monthly.html); [https://coastwatch.pfeg.noaa.gov/erddap/griddap/NOAA\\_DHW.html](https://coastwatch.pfeg.noaa.gov/erddap/griddap/NOAA_DHW.html))). See the document `satellite_SST_process-September_2022` for details about the temperature variables. Based on prior discussions, the index of juvenile abundance (i.e., CPUE) was based on the pooled-species vessel calibration coefficient.

The model performance metric one-step ahead MAPE for the last five years (years 2018 through 2022) and for the last ten years (years 2013 through 2022) was used to evaluate the forecast accuracy of the 18 individual models. The 2023 forecast was based on a . . . .

## Analysis

### Model data

The data used in the model are shown in table 1.

Table 1: Model data. This does not include the temperature data.

Year	Harvest	CPUE
1998	42.50	2.48
1999	77.80	5.62
2000	20.30	1.60
2001	67.00	3.73
2002	45.30	2.87
2003	52.50	2.78
2004	45.30	3.08

Year	Harvest	CPUE
2005	59.10	3.90
2006	11.60	2.04
2007	44.80	2.58
2008	15.90	1.17
2009	38.00	2.32
2010	24.00	2.33
2011	58.90	4.11
2012	21.30	1.51
2013	94.70	3.52
2014	37.20	2.14
2015	35.10	3.80
2016	18.40	2.45
2017	34.70	4.35
2018	8.10	0.35
2019	21.10	1.17
2020	8.07	1.14
2021	48.40	2.15
2022	16.00	0.88
2023	NA	1.45

## Individual, multiple linear regression models

Biophysical variables based on data from Southeast Alaska were used to forecast the harvest of adult pink salmon in Southeast Alaska, one year in advance, using individual, multiple linear regression models (models m1–m18; Table 2). The simplest regression model (model m1) consisted of only the predictor variable juvenile pink salmon CPUE ( $X_1$ ), while the other 17 regression models consisted of the predictor variable juvenile pink salmon CPUE and a temperature index ( $X_2$ ),

$$E(Y_i) = \hat{\alpha}_i + \hat{\beta}_{1i}X_1 + \hat{\beta}_{2i}X_2.$$

The temperature index was either the SECM survey ISTI temperature data (Murphy et al. 2019) or one of the 16 satellite-derived SST data (Huang et al. 2017). Although the simplest model only contained CPUE, including temperature data with CPUE is likely a more accurate measure of juvenile abundance if temperature affects the proportion of juveniles that migrate through Icy Strait in a given year (Murphy et al. 2019). The response variable ( $Y$ ; Southeast Alaska adult pink salmon harvest in millions) and CPUE data were natural log transformed in the model, but temperature data were not. The forecast ( $\hat{Y}_i$ ), and 80% prediction intervals (based on output from program R; R Core Team 2021) from the 18 regression models were exponentiated and bias-corrected (Miller 1984),

$$\hat{F}_i = \exp(\hat{Y}_i + \frac{\sigma_i^2}{2})$$

where  $\hat{F}_i$  is the preseason forecast (for each model  $i$ ) in millions of fish, and  $\sigma_i$  is the variance (for each model  $i$ ).

## Performance metrics

The model summary results using the performance metric one-step ahead MAPE are shown in table 2; the smallest value is the preferred model. The performance metric one-step ahead MAPE was calculated as:

1. Estimate the regression parameters at time  $t-1$  from data up to time  $t-1$ .
2. Make a prediction of  $\hat{Y}_t$  at time  $t$  based on the predictor variables at time  $t$  and the estimate of the regression parameters at time  $t-1$  (i.e., the fitted regression equation).
3. Calculate the MAPE based on the prediction of  $\hat{Y}_t$  at time  $t$  and the observed value of  $Y_t$  at time  $t$ ,

$$\text{MAPE} = \left| \frac{\exp(Y_t) - \exp(\hat{Y}_t + \frac{\sigma_t^2}{2})}{\exp(Y_t)} \right|.$$

4. For each individual model, average the MAPEs calculated from the forecasts,

$$\text{MAPE} = \frac{1}{n} \sum_{t=1}^n \left| \frac{\exp(Y_t) - \exp(\hat{Y}_t + \frac{\sigma_t^2}{2})}{\exp(Y_t)} \right|.$$

Table 2: Summary of the adjusted R squared value and the 5-year and 10-year one step ahead MAPEs.

model	terms	AdjR2	5-year MAPE	10-year MAPE
m1	CPUE	0.612	57%	63%
m2	CPUE + ISTI20_MJJ	0.817	37%	35%
m3	CPUE + Chatham_SST_May	0.799	30%	25%
m4	CPUE + Chatham_SST_MJJ	0.748	42%	38%
m5	CPUE + Chatham_SST_AMJ	0.801	31%	28%
m6	CPUE + Chatham_SST_AMJJ	0.778	36%	32%
m7	CPUE + Icy_Strait_SST_May	0.783	31%	24%
m8	CPUE + Icy_Strait_SST_MJJ	0.738	43%	38%
m9	CPUE + Icy_Strait_SST_AMJ	0.772	33%	28%
m10	CPUE + Icy_Strait_SST_AMJJ	0.756	37%	33%
m11	CPUE + NSEAK_SST_May	0.788	28%	24%
m12	CPUE + NSEAK_SST_MJJ	0.754	38%	33%
m13	CPUE + NSEAK_SST_AMJ	0.789	29%	27%
m14	CPUE + NSEAK_SST_AMJJ	0.771	32%	28%
m15	CPUE + SEAK_SST_May	0.770	32%	28%
m16	CPUE + SEAK_SST_MJJ	0.740	40%	36%
m17	CPUE + SEAK_SST_AMJ	0.776	32%	30%
m18	CPUE + SEAK_SST_AMJJ	0.756	35%	31%

Figure 1: The 2023 SEAK pink salmon harvest (millions) forecast by model with 80% prediction intervals (corrected for log transformation bias in a linear-model) around each forecast.

## Results:

### Model Diagnostics

Model diagnostics for model m11 included residual plots, the curvature test, and influential observation diagnostics using Cook’s distance (Cook 1977), the Bonferroni outlier test, and leverage plots. Model diagnostics were used to identify observations that were potential outliers, had high leverage, or were influential (Zhang 2016). These observations may have significant impact on model fitting and may need to be excluded.

Table 3: Detailed output for model m11. Juvenile years 1998, 1999, 2005, 2012, and 2019, and 2020 (years 1999, 2000, 2006, 2013, 2020, and 2021) show the largest standardized residual. Year refers to the forecast year. Fitted values are log-transformed.

year	SEAKCatch	CPUE	temp	resid	hat_values	Cooks_distance	std_resid	fit_bias_corrected
1998	42.5	2.48	7.35	0.29	0.04	0.01	0.94	33.51
1999	77.8	5.62	7.65	-0.45	0.29	0.40	-1.71	127.46
2000	20.3	1.60	6.70	-0.30	0.11	0.04	-1.03	28.86
2001	67.0	3.73	7.23	0.12	0.09	0.00	0.39	62.64
2002	45.3	2.87	6.66	-0.11	0.11	0.01	-0.36	52.88
2003	52.5	2.78	6.39	-0.02	0.15	0.00	-0.08	56.42
2004	45.3	3.08	7.57	0.16	0.05	0.00	0.53	40.48
2005	59.1	3.90	7.89	0.17	0.09	0.01	0.59	52.10
2006	11.6	2.04	8.42	-0.38	0.12	0.08	-1.31	17.80
2007	44.8	2.58	6.98	0.15	0.06	0.01	0.48	40.66
2008	15.9	1.17	6.90	-0.27	0.11	0.03	-0.92	21.82
2009	38.0	2.32	6.64	-0.04	0.10	0.00	-0.12	41.33
2010	24.0	2.33	7.32	-0.23	0.04	0.01	-0.75	31.59
2011	58.9	4.11	7.76	0.02	0.11	0.00	0.06	60.67
2012	21.3	1.51	7.25	0.01	0.07	0.00	0.02	22.19
2013	94.7	3.52	6.95	0.45	0.10	0.08	1.51	63.62
2014	37.2	2.14	6.59	0.01	0.11	0.00	0.02	38.76
2015	35.1	3.80	8.15	-0.20	0.11	0.02	-0.67	44.80
2016	18.4	2.45	8.92	0.09	0.21	0.01	0.33	17.63
2017	34.7	4.35	8.92	-0.16	0.26	0.04	-0.59	42.62
2018	8.1	0.35	7.75	-0.22	0.19	0.05	-0.79	10.61
2019	21.1	1.17	7.53	0.27	0.09	0.03	0.91	16.93
2020	8.1	1.14	8.42	-0.32	0.19	0.10	-1.15	11.70
2021	48.4	2.15	8.26	0.94	0.10	0.37	3.18	19.90
2022	16.0	0.88	7.29	0.03	0.12	0.00	0.11	16.25

### Cook’s distance

Cook’s distance is a measure of influence, or the product of both leverage and outlier. Cook’s distance,

$$D_i = \frac{e_{PSi}^2}{k+1} * \frac{h_i}{1-h_i},$$

where  $e_{PSi}^2$  is the standardized Pearson residuals,  $h_i$  are the hat values (measure of leverage), and  $k$  is the number of predictor variables in the model, is a measure of overall influence of the  $i_{th}$  data point on all  $n$  fitted values (Fox and Weisburg 2019). A large value of Cook’s distance indicates that the data point is an influential observation. Cook’s distance values greater than  $4/(n-k-1)$ , where  $n$  is the number of

observations (i.e., 25), was used as a benchmark for identifying the subset of influential observations (Ren et al. 2016). Therefore, a Cook's distance cut-off of 0.18 was used; observations with a Cook's distance greater than 0.18 were investigated further (Figure 1a).

### Leverage

An observation that is distant from the average covariate pattern is considered to have high leverage or hat-value. If an individual observation has a leverage value  $h_i$  greater than 2 or 3 times  $p/n$  (Ren et al. 2016), it may be a concern (where  $p$  is the number of parameters in the model including the intercept (i.e., 3), and  $n$  is the number of observations in the model (i.e., 25);  $p/n = 3/25 = 0.12$  for this study). Therefore, a leverage cut-off of 0.24 was used; observations with a leverage value greater than 0.24 were investigated further (Figure 1b).

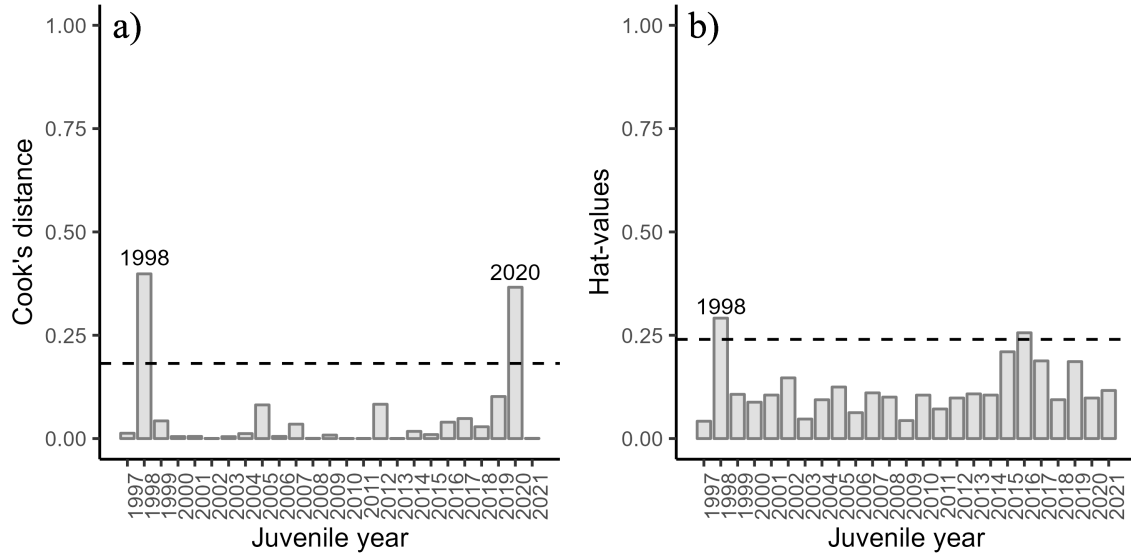


Figure 1: Diagnostics plots of influential observations including a) Cook's Distance (with a cut-off value of 0.18), and b) leverage values (with a cut-off value of 0.24) from model m11.

### Influential datapoints

To determine if a variable has a relationship with residuals, a lack-of fit curvature test was performed. In this test, terms that are non-significant suggest a properly specified model. The CPUE term was significant in the lack-of-fit curvature test ( $P < 0.05$ ), suggesting some lack of fit for this term (Figure 2a). Diagnostics indicated that two of the data points were above the cut-off value for the Cook's distance (Figure 1a). One observation had a high leverage value (Figure 1b). Based on the Bonferroni outlier test, one of the data points had a studentized residual with a significant Bonferroni  $P$ -value suggesting that one of the data points impacted the model fitting; although observations 2, 3, 9, 16, 23, and 24 and were the most extreme (juvenile years 1998, 1999, 2005, 2012, and 2019, and 2020) based on standardized residuals (Figure 3a; Table 5). Based on the lightly curved fitted lines in the residual versus fitted plot (Figure 3b), the fitted plot shows some lack of fit of the model.

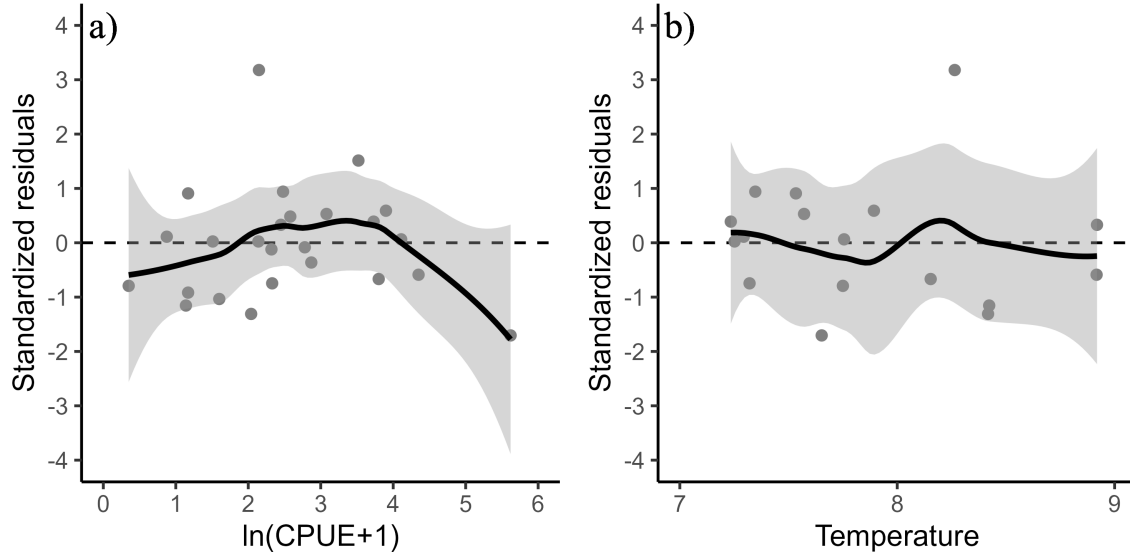


Figure 2: Standardized residuals versus predicted plots for a) CPUE and b) temperature.

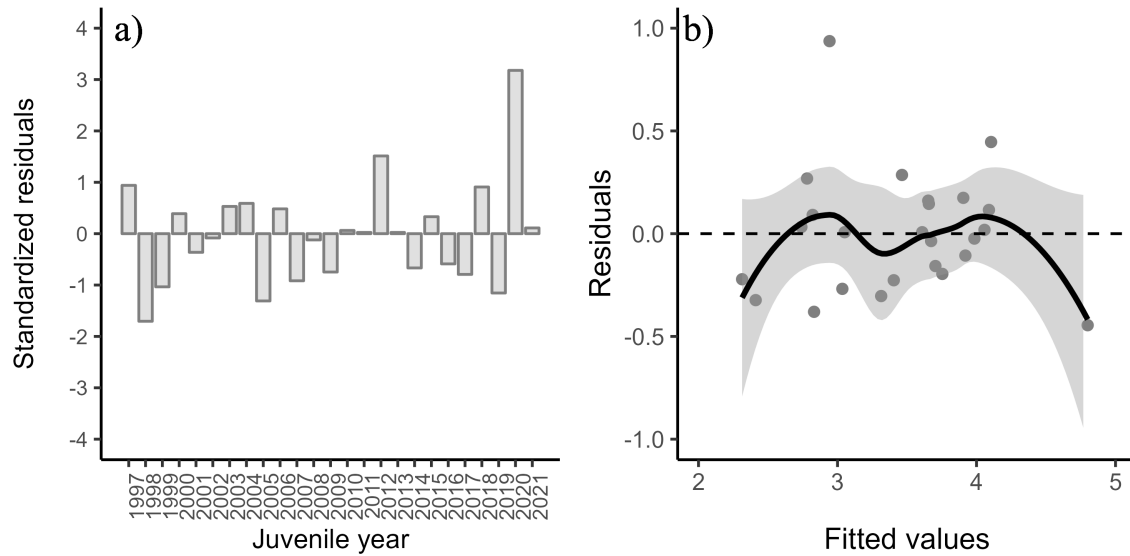


Figure 3: a) Standardized residuals versus juvenile year and b) residuals versus fitted values for model m11. Positive residuals indicate that the observed harvest was larger than predicted by the model.

The best regression model based on the performance metric one-step ahead MAPE was model m11 (i.e., the model containing CPUE, and a May NSEAK SST data). The adjusted  $R^2$  value for model m11 was 0.79 (Table 3) indicating overall a good model fit. Based upon a model that includes juvenile pink salmon CPUE and May NSEAK SST (model m11), the 2023 SEAK pink salmon harvest is predicted to be in the weak range with a point estimate of xxxx million fish (80% prediction interval: xxxx to xxxx million fish).

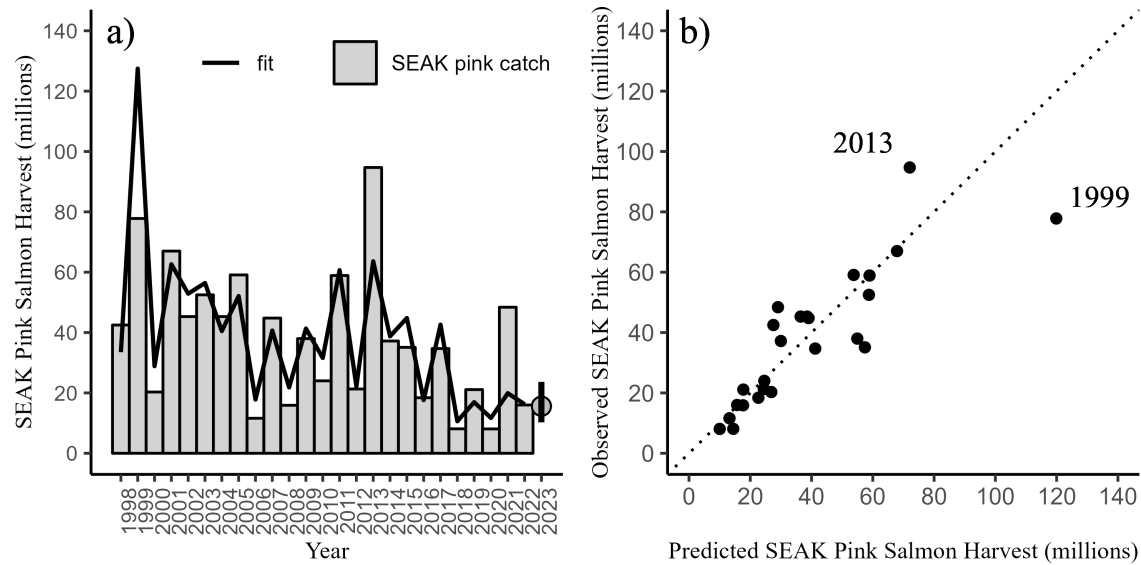


Figure 4: A) SEAK pink salmon harvest (millions) by year with the model fit (line). The predicted 2023 forecast is symbolized as a grey circle with an 80% prediction interval (xxx to xxx million fish). B) SEAK pink salmon harvest (millions) against the fitted values from model m11 by year. The dotted line is a one to one line.

## References

- Cook, R. D. 1977. Detection of influential observations in linear regression. *Technometrics* 19: 15-18.
- Fox, J. and S. Weisburg. 2019. *An R Companion to Applied Regression*, Third Edition. Thousand Oaks CA: Sage Publications, Inc.
- Huang, B., P. W. Thorne, V. F. Banzon,, T. Boyer, G. Chepurin, J. H. Lawrimore, M. J. Menne, T. M. Smith, R. S. Vose, and H. M. Zhang. 2017. Extended reconstructed sea surface temperature, version 5 (ERSSTv5): upgrades, validations, and intercomparisons. *Journal of Climate* 30:8179–8205.
- Miller, D. M. 1984. Reducing transformation bias in curve fitting. *The American Statistician* 38: 124-126.
- Murphy, J. M., E. A. Fergusson, A. Piston, A. Gray, and E. Farley. 2019. Southeast Alaska pink salmon growth and harvest forecast models. North Pacific Anadromous Fish Commission Technical Report No. 15: 75-81.
- NOAA Coral Reef Watch (NOAA\_DHW\_monthly dataset). 2021, updated daily. NOAA Coral Reef Watch Version 3.1 Monthly 5km SST and SST Anomaly, NOAA Global Coral Bleaching Monitoring Time Series Data, May 1997-June 2021. College Park, Maryland, USA: NOAA/NESDIS/STAR Coral Reef Watch program. Data set accessed 2021-10-01 at [https://coastwatch.pfeg.noaa.gov/erddap/griddap/NOAA\\_DHW\\_monthly.html](https://coastwatch.pfeg.noaa.gov/erddap/griddap/NOAA_DHW_monthly.html).
- NOAA Coral Reef Watch (NOAA\_DHW dataset). 2021, updated daily. NOAA Coral Reef Watch Daily Near-real-Time Global 5km SST and SST Anomaly, NOAA Global Coral Bleaching Monitoring Time Series Data, July 2021. College Park, Maryland, USA: NOAA/NESDIS/STAR Coral Reef Watch program. Data set accessed 2021-10-01 at [https://coastwatch.pfeg.noaa.gov/erddap/griddap/NOAA\\_DHW.html](https://coastwatch.pfeg.noaa.gov/erddap/griddap/NOAA_DHW.html).
- Orsi, J. A., E. A. Fergusson, A. C. Wertheimer, E. V. Farley, and P. R. Mundy. 2016. Forecasting pink salmon production in Southeast Alaska using ecosystem indicators in times of climate change. *N. Pac. Anadr. Fish Comm. Bull.* 6: 483–499. (Available at <https://npafc.org>)

Piston, A. W., J. Murphy, J. Moss, W. Strasburger, S. C. Heintz, E. Fergusson, S. Miller, A. Gray, and C. Waters. 2021. Operational Plan: Southeast coastal monitoring, 2021. ADF&G, Regional Operational Plan No. ROP.CF.1J.2021.02, Douglas.

R Core Team. 2021. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL: <http://www.r-project.org/index.html>

Ren et al. reference

Wertheimer A. C., J. A. Orsi, M. V. Sturdevant, and E. A. Fergusson. 2006. Forecasting pink salmon harvest in Southeast Alaska from juvenile salmon abundance and associated environmental parameters. In Proceedings of the 22nd Northeast Pacific Pink and Chum Workshop. Edited by H. Geiger (Rapporteur). Pac. Salmon Comm. Vancouver, British Columbia. pp. 65–72.

Zhang, Z. 2016. Residuals and regression diagnostics: focusing on logistic regression. *Annals of Translational Medicine* 4: 195.



# Appendix

Table 4: Parameter estimates for the 18 individual models.

model	term	estimate	std.error	statistic	p.value
m1	(Intercept)	2.3446817	0.195	12.039	0.000
m1	CPUE	0.4277626	0.069	6.236	0.000
m2	(Intercept)	7.2759139	0.961	7.570	0.000
m2	CPUE	0.4900832	0.049	10.088	0.000
m2	ISTI20_MJJ	-0.5620524	0.108	-5.181	0.000
m3	(Intercept)	5.5746972	0.696	8.007	0.000
m3	CPUE	0.4872575	0.051	9.569	0.000
m3	Chatham_SST_May	-0.4480194	0.095	-4.736	0.000
m4	(Intercept)	6.6338302	1.181	5.617	0.000
m4	CPUE	0.4533798	0.056	8.138	0.000
m4	Chatham_SST_MJJ	-0.4441657	0.121	-3.664	0.001
m5	(Intercept)	6.3136306	0.844	7.480	0.000
m5	CPUE	0.4737561	0.050	9.452	0.000
m5	Chatham_SST_AMJ	-0.5325521	0.112	-4.768	0.000
m6	(Intercept)	6.7390100	1.043	6.463	0.000
m6	CPUE	0.4625582	0.053	8.799	0.000
m6	Chatham_SST_AMJJ	-0.5140892	0.121	-4.257	0.000
m7	(Intercept)	5.1869991	0.667	7.772	0.000
m7	CPUE	0.4996308	0.054	9.266	0.000
m7	Icy_Strait_SST_May	-0.4202629	0.096	-4.364	0.000
m8	(Intercept)	6.1748137	1.118	5.524	0.000
m8	CPUE	0.4584470	0.057	8.026	0.000
m8	Icy_Strait_SST_MJJ	-0.3889864	0.112	-3.462	0.002
m9	(Intercept)	5.8747371	0.868	6.769	0.000
m9	CPUE	0.4823046	0.054	8.886	0.000
m9	Icy_Strait_SST_AMJ	-0.4910139	0.119	-4.129	0.000
m10	(Intercept)	6.2146681	1.024	6.069	0.000
m10	CPUE	0.4686008	0.055	8.458	0.000
m10	Icy_Strait_SST_AMJJ	-0.4537813	0.119	-3.823	0.001
m11	(Intercept)	5.2506033	0.666	7.885	0.000
m11	CPUE	0.4643401	0.051	9.030	0.000
m11	NSEAK_SST_May	-0.3999366	0.089	-4.470	0.000
m12	(Intercept)	6.4231506	1.094	5.872	0.000
m12	CPUE	0.4372503	0.055	7.988	0.000
m12	NSEAK_SST_MJJ	-0.4158686	0.110	-3.767	0.001
m13	(Intercept)	6.0130213	0.828	7.263	0.000
m13	CPUE	0.4545559	0.051	8.920	0.000
m13	NSEAK_SST_AMJ	-0.4940128	0.110	-4.499	0.000
m14	(Intercept)	6.4439839	1.006	6.406	0.000
m14	CPUE	0.4465449	0.053	8.444	0.000
m14	NSEAK_SST_AMJJ	-0.4775477	0.116	-4.121	0.000
m15	(Intercept)	5.2527324	0.725	7.248	0.000
m15	CPUE	0.4619235	0.053	8.642	0.000
m15	SEAK_SST_May	-0.3699946	0.090	-4.101	0.000
m16	(Intercept)	6.1864825	1.106	5.592	0.000
m16	CPUE	0.4296916	0.056	7.651	0.000
m16	SEAK_SST_MJJ	-0.3713206	0.106	-3.509	0.002
m17	(Intercept)	5.9806776	0.872	6.857	0.000

model	term	estimate	std.error	statistic	p.value
m17	CPUE	0.4521473	0.052	8.629	0.000
m17	SEAK_SST_AMJ	-0.4543651	0.107	-4.230	0.000
m18	(Intercept)	6.2475961	1.035	6.035	0.000
m18	CPUE	0.4394329	0.055	8.062	0.000
m18	SEAK_SST_AMJJ	-0.4266927	0.112	-3.813	0.001