

Standardization of bigeye and yellowfin tuna CPUE by Japanese longline in the Indian Ocean which includes cluster analysis

Takayuki Matsumoto¹, Keisuke Satoh¹ and Simon Hoyle²

¹*National Research Institute of Far Seas Fisheries (NRIFSF), Fisheries Research Agency (FRA), 5-7-1, Orido, Shimizu, Shizuoka, 424-8633, Japan*

²*IOTC consultant*

Abstract

Standardizations of Japanese longline CPUE for bigeye and yellowfin tuna in multiple Indian Ocean regions were conducted using generalized linear models (GLM) with log normal errors. The models incorporated fishing power based on vessel ID where available, and used cluster analysis to account for targeting. The variables year-quarter, vessel ID, latlong5 (five degree latitude-longitude block), cluster and number of hooks were used in the standardization. The numbers of clusters selected varied among regions and species, but in all cases were either 4 or 5. Dominant species differed depending on clusters. The effects of each covariate differed depending on species and region. The CPUE trends were similar to those estimated last year, though with some differences due to the inclusion of vessel effects and cluster variables.

1. Introduction

To date, national scientists have mainly standardized bigeye and yellowfin tuna Japanese longline CPUE using generalized linear models (GLM), with log normal errors and either operational or aggregated catch and effort data (e.g. Matsumoto et al., 2016a, b). The standardizations have incorporated the effects of fishing season, area, fishing gear (number of hooks between floats and gear material) and an environmental factor (sea surface temperature). These may be termed ‘simple’ and ‘traditional’ methods.

In 2016, IOTC joint CPUE analysis (CPUE workshop) was conducted and ‘joint CPUEs’ were created for bigeye and yellowfin tuna, based on Japanese, Taiwanese and Korean longline operational data (Hoyle et al., 2016). These models account for fishing power based on vessel ID where available, and use cluster analysis to incorporate targeting. Joint CPUEs were considered to be more representative of status of the stocks and so were used for base models of stock assessment. At that time fleet-specific CPUE indices were prepared for Japanese longline using the same methods, but were not presented, so it was not possible to compare the joint and Japanese-only longline CPUE indices. This year the joint CPUE analysis workshop was again held and CPUE indices for each fleet as well as joint CPUE were created. The workshop also provided training to the participants on standardization methods. This document reports the standardization of bigeye and yellowfin tuna Japanese longline CPUE conducted at this year’s joint CPUE analysis, using the same methods as that for joint CPUE.

2. Materials and methods

Data

Operational level (set by set) Japanese longline logbook data were used. The data were available for 1952-2016 (data for 2016 were preliminary), with the fields year, month and day of operation, location to 1° of

latitude and longitude, vessel call sign, no. of hooks between floats (HBF), number of hooks per set, date of the start of the fishing cruise, logbook identifier, and catch in number of each species. Vessel call signs were available from 1979 onward and were used for the vessel identifier. The operations with hooks per set above 5000 and less than 200 were removed. Sets after 1975 with HBF missing or > 25 were removed. Sets before 1975 with missing HBF were allocated HBF of 5, according to standard practice with Japanese longline data (e.g. Langley et al. 2005; Hoyle et al. 2013; Ochi et al. 2014).

Each set was allocated to bigeye and yellowfin regions (Fig. 1). These regions are the same as those in Hoyle et al. (2016).

Cluster analysis

We clustered the data using the approach applied by Hoyle et al. (2015). We removed all sets with no catch of any of the species, and then aggregated by vessel-month. Set level data contains variability in species composition due to the randomness of chance encounters between fishing gear and schools of fish. This variability leads to some misallocation of sets using different fishing strategies. Aggregating the data tends to reduce the variability, and therefore reduce misallocation of sets. For these analyses we aggregated the data by vessel-month, assuming that individual vessels tend to follow a consistent fishing strategy through time. One trade-off with aggregation in this way is that vessels may change their fishing strategy within a month, which will result in misallocation of sets. For the purposes of this paper we refer to aggregation by vessel-month as trip-level aggregation, although the time scale is (for distant water vessels) in most cases shorter than a fishing trip. In the data prior to 1979 vessel id was not available, but we were able to cluster them by vessel-month because the logbook id, available for the first time in the current data set, could be used to identify sets on the same vessel-trip.

We calculated proportional species composition by dividing the catch in numbers of each species by catch in numbers of all species in the vessel-month. Thus the species composition values of each vessel-month summed to 1, ensuring that large catches and small catches were given equivalent weight. The data were transformed by centering and scaling, so as to reduce the dominance of species with higher average catches. Centering was performed by subtracting the column (species) mean from each column, and scaling was performed by dividing the centered columns by their standard deviations.

We clustered the data using the hierarchical Ward hclust method, implemented with function hclust in R, option ‘Ward.D’, after generating a Euclidean dissimilarity structure with function ‘dist’. This approach differs from the standard Ward D method which can be implemented by either taking the square of the dissimilarity matrix or using method ‘ward.D2’ (Murtagh & Legendre 2014). However in practice the method gives similar patterns of clusters to other methods, more reliably than ward.D2 (Hoyle et al 2015).

Data were also clustered using the kmeans method, which minimises the sum of squares from points to the cluster centres, using the algorithm of Hartigan and Wong (1979). It was implemented using function kmeans in the R stats package (R Core Team 2014).

Selecting the number of groups

We used several subjective approaches to select the appropriate number of clusters. In most cases the approaches suggested the same or similar numbers of groups. First, we applied hclust to transformed trip-

level data and examined the hierarchical trees, subjectively estimating the number of distinct branches. Second, we ran kmeans analyses on untransformed trip-level data with number of groups k ranging from 2 to 25, and plotted the deviance against k . The optimal group number was the lowest value of k after which the rate of decline of deviance became slower and smoother. Third, following Winker et al (2014) we applied the `nScree()` function from the R `nFactors` package (Raiche & Magis 2010), which uses various approaches (Scree test, Kaiser rule, parallel analysis, optimal coordinates, acceleration factor) to estimate the number of components to retain in an exploratory PCA. Where there was uncertainty about the number of clusters, we selected the option with more clusters.

We plotted the `hclust` clusters to explore the relationships between them and the species composition and other variables, such as HBF, number of hooks, year, and set location. Plots included boxplots of a) proportion of each species in the catch, by cluster; b) the distributions of variables by cluster; and c) maps of the spatial distribution of clusters, one map for each cluster.

In some analyses clusters that caught very few of the species of interest were omitted, because they provide little relevant information and may cause analysis problems due to large numbers of zeroes, and memory problems due to large sample sizes. Cluster selection was based on review and discussion of the plots of covariates and species compositions by cluster. Analyses were run both with and without these clusters – see the ‘Models and datasets’ section.

For standardization of each region, data were selected for vessels that had fished for at least N1 quarters in that region. The standard level of N1 was 8 quarters in the equatorial regions and 2 quarters in the southern regions. Subsequently, vessels, 5° cells, and year-quarters were included if they had at least 100 sets. For analyses of the 1952-1979 period this criterion was reduced to 50 sets, to increase the size of the dataset. For datasets with more than 60,000 sets the number of sets in each stratum (5° square * year-quarter) was limited by randomly selecting 60 sets without replacement from strata with more than this number of sets. Testing suggested that this approach did not cause bias, and the effects on trends of random variation were reduced to very low levels at 30 sets per stratum (Hoyle & Okamoto 2011), suggesting that 60 sets was more than adequate.

CPUE standardization, and fleet efficiency analyses

CPUE standardization methods generally followed the approaches used by Hoyle et al. (2015). The operational data were standardized using generalized linear models in R.

GLM (generalized linear models) that assumed a lognormal and delta lognormal distribution was conducted, and in this report only the methods and results for lognormal distributions are shown. In this approach the response variable $\log(\text{CPUE}+k)$ was used, and a Normal distribution assumed. The constant k , added to allow for modelling sets with zero catches of the species of interest, was 10% of the mean CPUE for all sets. The following models were used:

$$\ln(\text{CPUEs}+k) \sim \text{yrqtr} + \text{vessid} + \text{latlong5} + \text{cluster} + f(\text{hooks}) + \epsilon$$

$$\ln(\text{CPUEs}+k) \sim \text{yrqtr} + \text{vessid} + \text{latlong5} + g(\text{HBF}) + f(\text{hooks}) + \epsilon$$

where *yrqtr*: year and quarter; *vessid*: effect of vessel ID; *latlong5*: effect of five degree latitude and

longitude; *cluster*: effect of cluster; *f(hooks)*: function of number of hooks modelled with a cubic spline; *g(HBF)*: function of the number of hooks between floats modelled with a cubic spline; ϵ : error term.

Data periods

Vessel identity information was only available from 1979, so could not be applied uniformly across all years. The discontinuity in 1979 could be addressed in several different ways. We therefore analyzed the data in several ways so as to provide the assessment scientists with appropriate data. For each of the approaches above, four analyses were carried out as shown below.

Analysis	Years	Vessel effects
1	1952-1979	No
2	1979-2016	Yes
3	1952-2016	No
4	1952-2016	Yes

It is possible to standardize the time series with vessel effects by assigning an identical dummy value to all vessels without vessel identity information. This was done for analysis 3). However using a dummy value introduces several problems. First, not all vessels begin to report their call sign at once in 1979, and those that do are self-selected and not randomly selected from the vessel population. Therefore it cannot be assumed that fishing power remains constant after 1979 for the dummy vessel id, so the transition in 1979 may introduce a discontinuity into the time series. The discontinuity can be limited in scope by restricting the overlap between dummy and real vessel IDs to one year – 1979 – and removing sets with missing vessel IDs after this time. Secondly, residuals may be more variable before 1979, without a true vessel ID in the model, which can introduce bias into the standardization.

One approach for addressing the discontinuity in analysis 3) is to adjust the time period 1952-1978 so that the relative averages in 1978 and 1979 are the same as they are in analysis 4), without vessel effects. However we considered that a better approach may be to estimate two time series 1952-1979 without vessel effects, and a second time series 1979-2015 with vessel effects (omitting all sets without vessel IDs). These are analyses 1) and 2) above. Subsequently the analyst can use them as desired, for example concatenating them after adjusting the averages so that the estimates for 1979 are the same.

Indices of abundance

Indices of abundance were obtained by applying the R function `predict.glm` to model objects. Binomial time effects were obtained by generating time effects from the glm and adjusting them so that their mean was the proportion of positive sets across the whole dataset. The main aim with this approach is to obtain a CPUE that varies appropriately, since variability for a binomial is greater when the mean is at 0.5 than at 0.02 or 0.98, and the multiplicative effect of the variability is greater when the mean is lower. The outcomes were normalised and reported as relative CPUE with mean of 1.

Uncertainty estimates were provided by applying the R function `predict.glm` with `type = "terms"` and `se.fit=TRUE`, and taking the standard error of the year-quarter effect. For the delta lognormal models we used

only the uncertainty in the positive component. Uncertainty estimates from standardizing commercial logbook data are in general biased low and often ignored by assessment scientists, since they assume independence and ignore autocorrelation associated with (for example) consecutive sets by the same vessels in the same areas. There may be a very large mismatch between the observation error in CPUE indices and the process error in the indices that is estimated in the assessment. This is particularly true for distant water longline CPUE, where very large sample sizes generate small observation errors.

Residual distributions and Q-Q plots were produced for all but the binomial analyses. For the lognormal positive analyses that included cluster in the model, median residuals were plotted by cluster. For all lognormal positive analyses, residuals by year-quarter were plotted by flag; median residuals by year-quarter were plotted by flag; and median residuals by 5° cell were mapped onto a contour plot for each flag.

We compared the indices with the area-specific Japanese bigeye (Matsumoto et al. 2016a) and yellowfin (Matsumoto et al. 2016b) indices from 2016. For each comparison, each dataset was first normalized by dividing through by its mean for 1980-2000, and the datasets plotted on the same axes. Secondly, the joint indices were divided by the matching year-quarter values from the Japanese indices, and these ratios were plotted to show the relative trends of the two time series.

3. Results and discussion

Cluster analysis

The aim of the cluster analysis was first to identify separate fishing strategies in the data for each species, regional structure, fleet, and region, and so to better understand the fishing practices; and second to assign each unit of fishing effort to a particular fishing strategy, so that the clusters could be used in standardization.

Species compositions were plotted by cluster for each region and fleet, as were the relative distributions of covariates (**Fig. 2- Fig. 5**). Dominant species differed depending on clusters. Clusters with low levels of the target species were excluded from standardization datasets. Numbers of clusters were 4 or 5.

Fig. 6 and **Fig. 7** show the effect of each covariate for bigeye and yellowfin region, respectively (from 1979 onward with vessel ID). For bigeye tuna, the vessel effect usually increased with time, and the HBF effect increased with the number of HFB. As for yellowfin tuna, vessel effect increased with time in several regions, and HBF effect peaked in the lowest or middle number of HFB.

Fig. 8 and **Fig. 9** show the trend of standardized CPUE for bigeye and yellowfin, respectively, without and with vessel effects. The trend differs between species and among regions, but CPUE usually shows decreasing trend with part of them increasing in the recent years.

Fig. 10 shows comparison of bigeye and yellowfin CPUE with those created last year based on ‘traditional’ method (without cluster analysis and vessel effect). The trend of both CPUEs was mostly similar, but there are some differences especially in the early period. Decline in CPUE is slower for new CPUE, probably because of the results of incorporating vessel effect and/or targeting. As for yellowfin region 4, there are many missing values in the new CPUE. This is partly because some data were eliminated based on cluster analysis, but also because of the settings that eliminated strata with too few sets per vessel, quarter, or 5

degree square. .

Fig. 11 and **Fig. 12** show distribution of standardized residuals and QQ plots for bigeye and yellowfin, respectively.

4. References

- Hoyle SD, Okamoto H 2011. Analyses of Japanese longline operational catch and effort for bigeye and yellowfin tuna in the WCPO, WCPFC-SC7-SA-IP-01. Western and Central Pacific Fisheries Commission, 9th Scientific Committee. Pohnpei, Federated States of Micronesia.
- Hoyle SD, Okamoto H, Yeh Y-m, Kim ZG, Lee SI, Sharma R 2015. IOTC-CPUEWS02 2015: Report of the 2nd CPUE Workshop on Longline Fisheries, 30 April – 2 May 2015. 126 p.
- Hoyle, S., Chang, Y., Kim, D. N., Lee, S., Matsumoto, T., Satoh, K. and Yeh, Y. (2016) Collaborative study of tropical tuna CPUE from multiple Indian Ocean longline fleets in 2016. IOTC-2016-WPTT18-14.
- Matsumoto, T., Nishida, T., Satoh, K. and Kitakado, T. (2016a) Japanese longline CPUE for bigeye tuna in the Indian Ocean standardized by GLM. IOTC-2016-WPTT18-13. pp 17.
- Matsumoto, T., Nishida, T., Satoh, K. and Kitakado, T. (2016b) Japanese longline CPUE for yellowfin tuna in the Indian Ocean standardized by generalized linear model. IOTC-2016-WPTT18-25. pp 22.
- Murtagh F, Legendre P 2014. Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward's Criterion? *Journal of Classification* 31(3): 274-295.
- Winker H, Kerwath SE, Attwood CG 2014. Proof of concept for a novel procedure to standardize multispecies catch and effort data. *Fisheries Research* 155: 149-159.

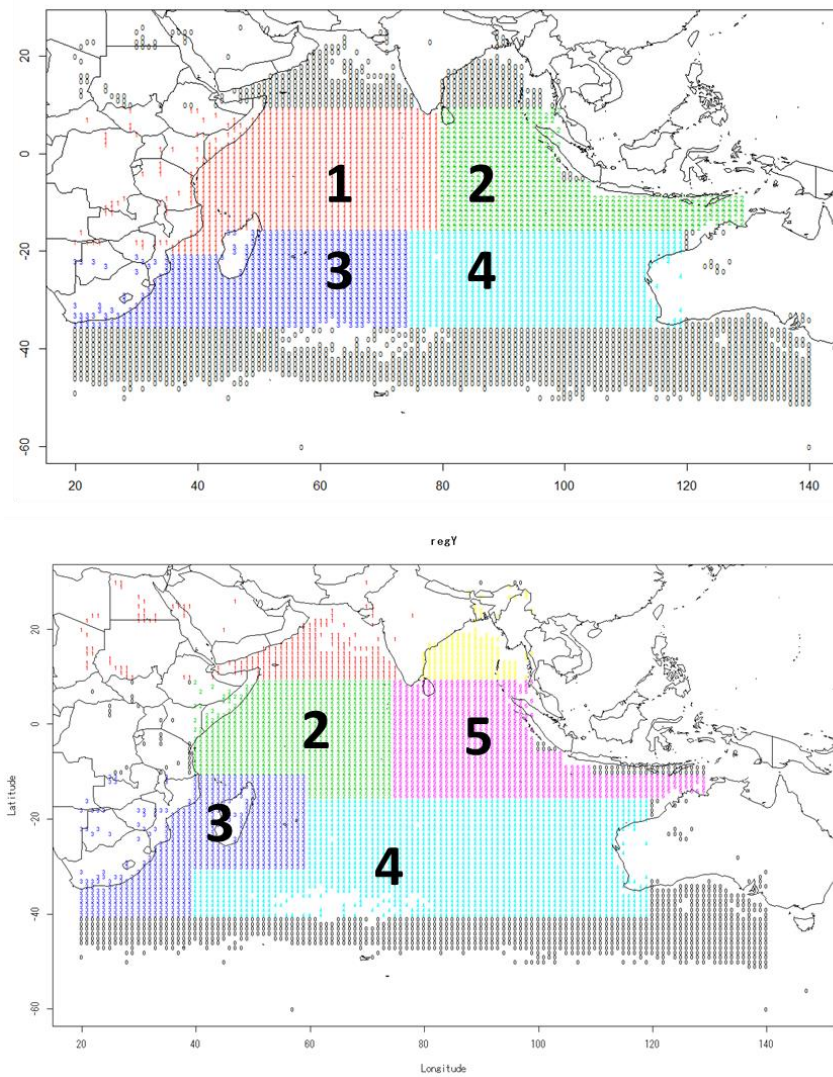


Fig. 1. Maps of the regional structures used to estimate bigeye (top) and yellowfin (bottom) CPUE indices.

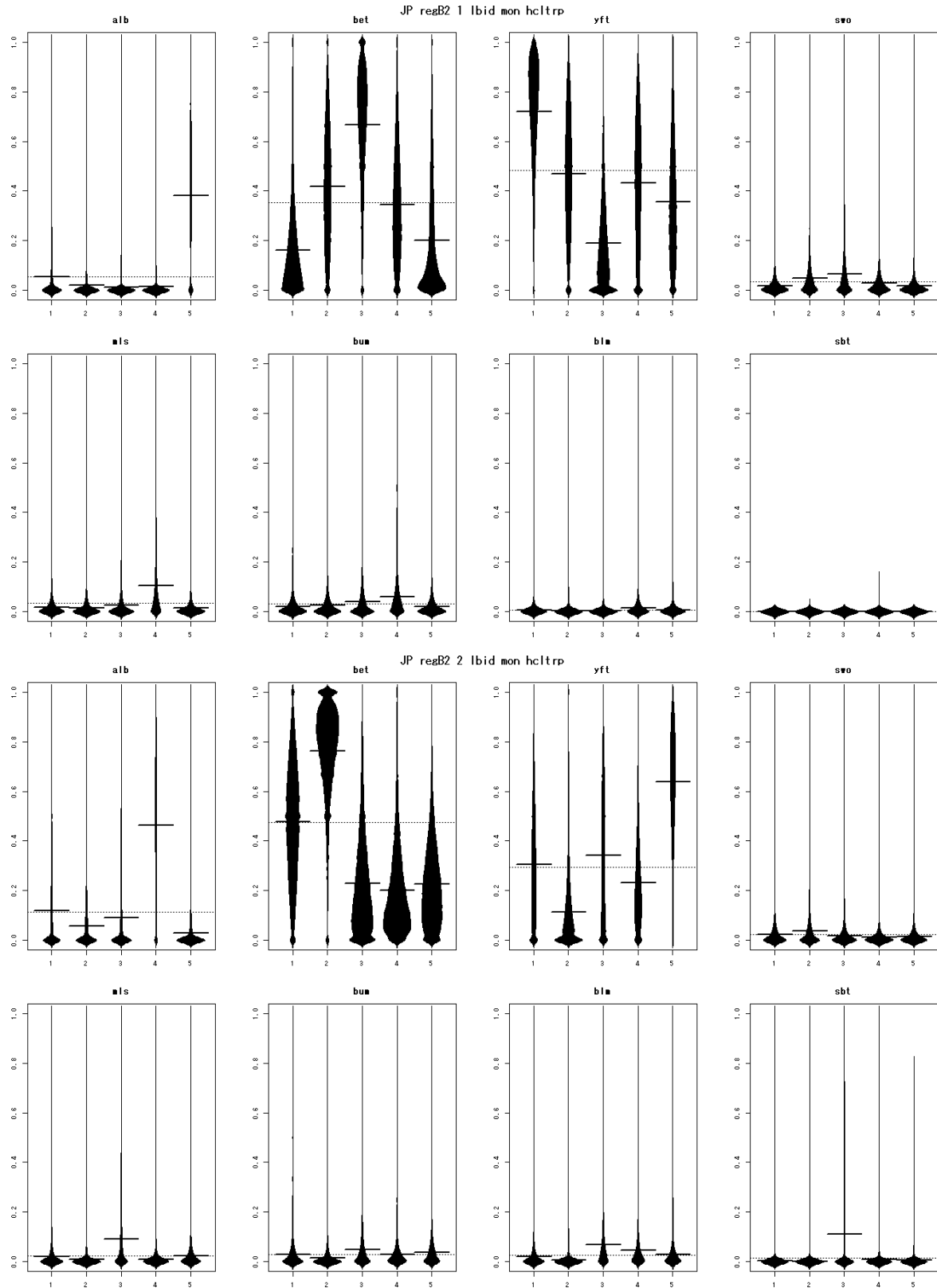


Fig. 2. Beanplots for bigeye region showing species composition by cluster. The horizontal bars indicate the medians.

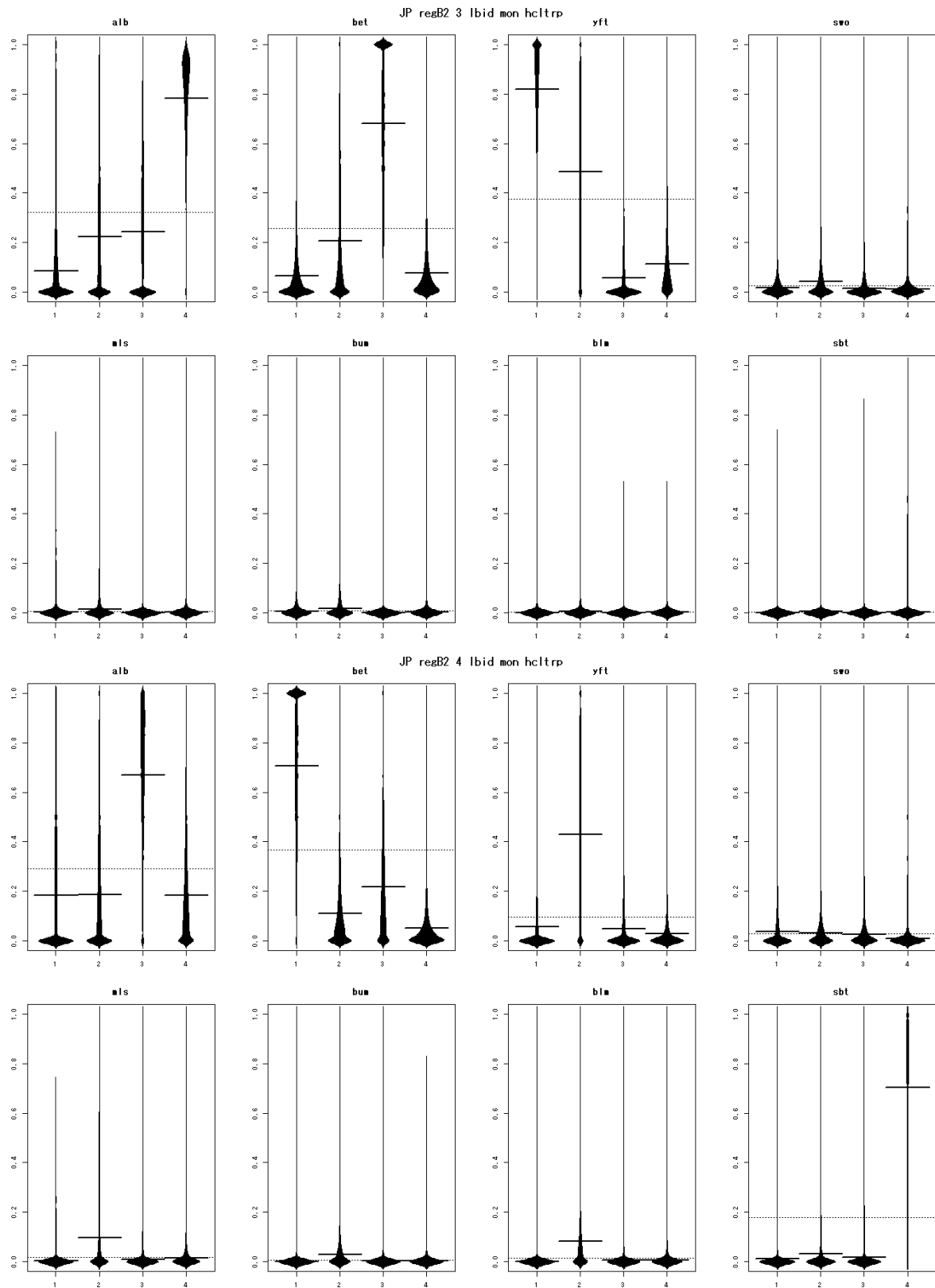


Fig. 2. Beanplots for bigeye region showing species composition by cluster. The horizontal bars indicate the medians. (continued)

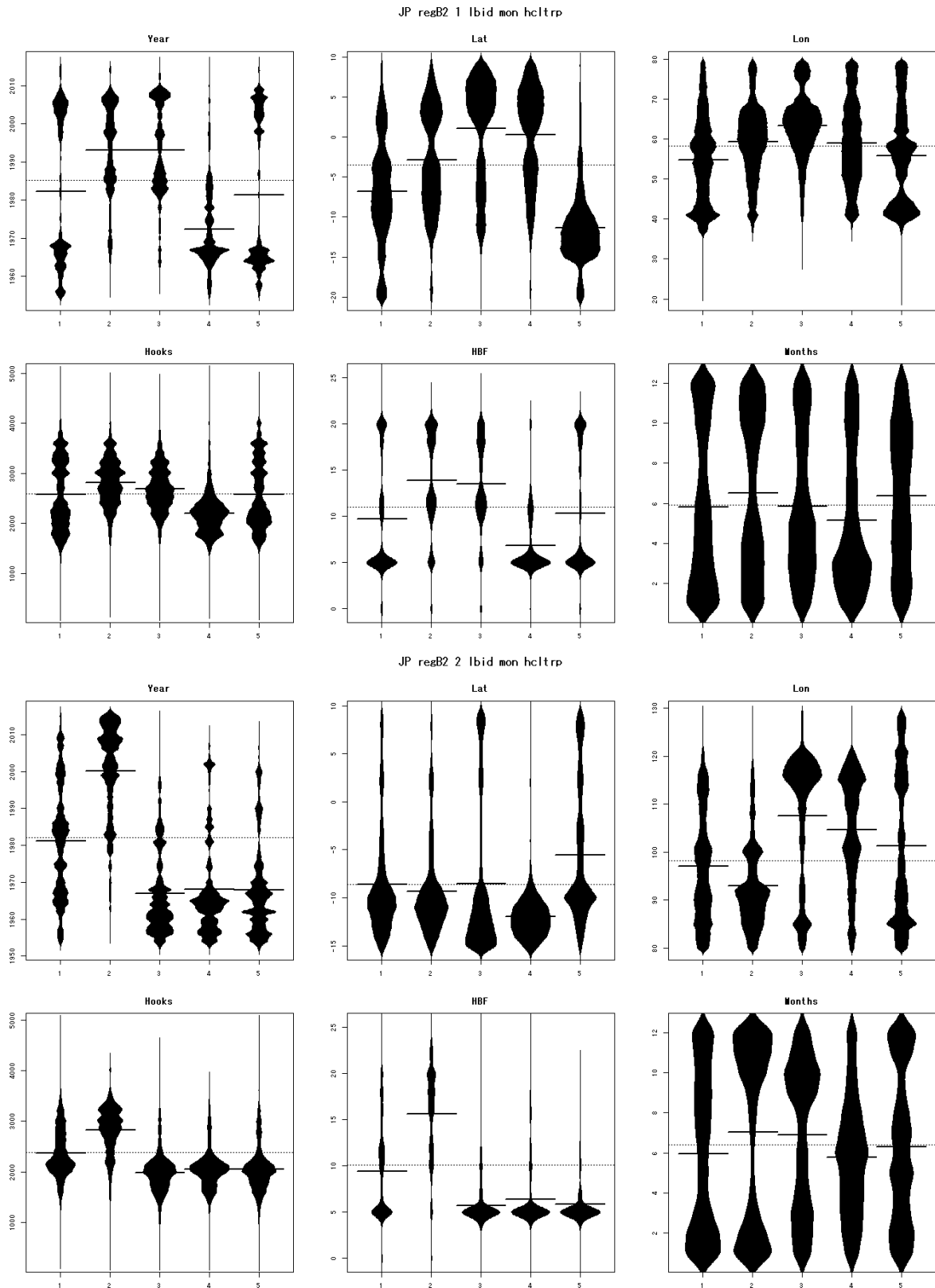


Fig. 3. Beanplots for bigeye region showing number of sets versus covariate by cluster. The horizontal bars indicate the medians.

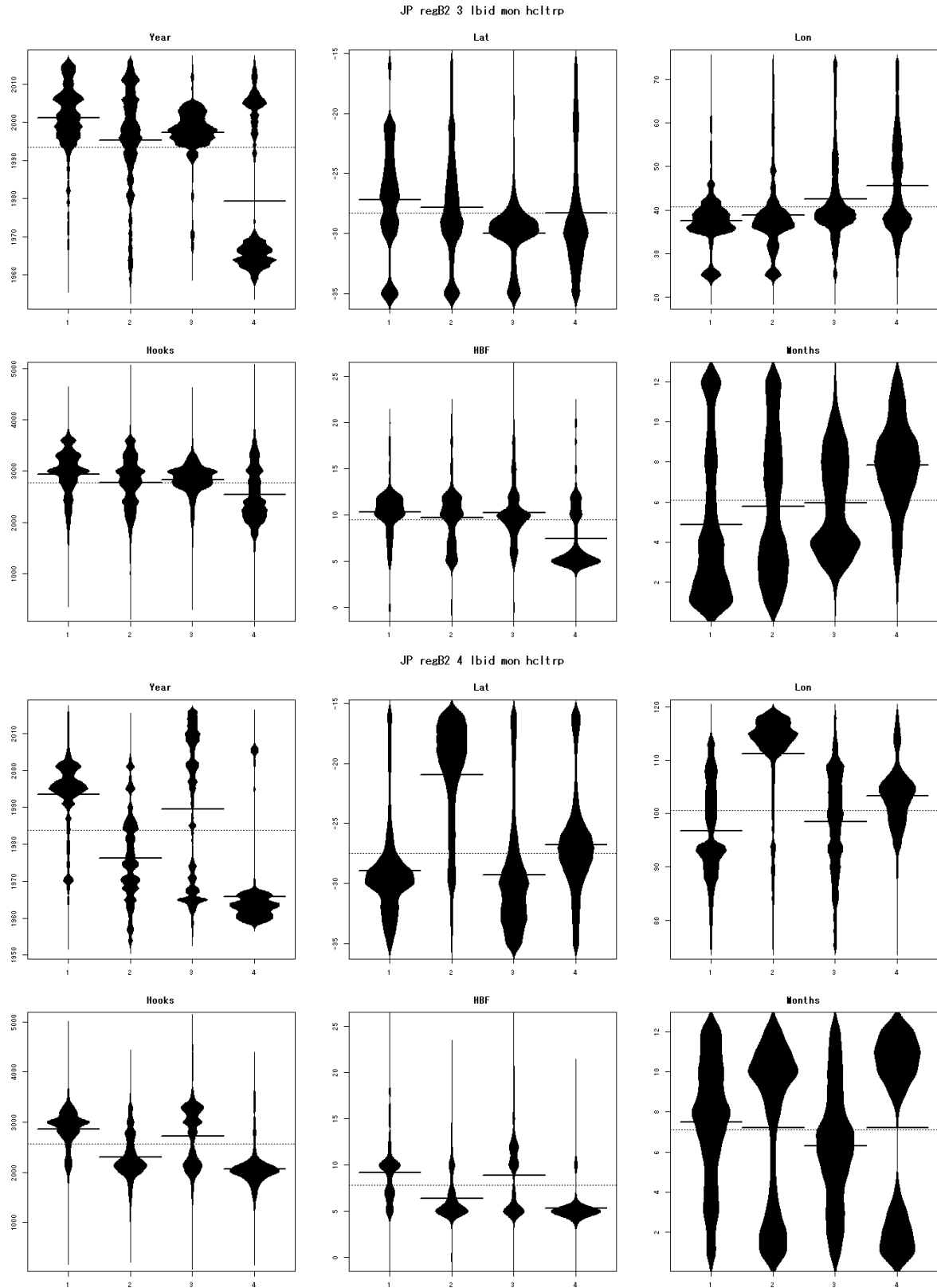


Fig. 3. Beanplots for bigeye region showing number of sets versus covariate by cluster. The horizontal bars indicate the medians. (continued)

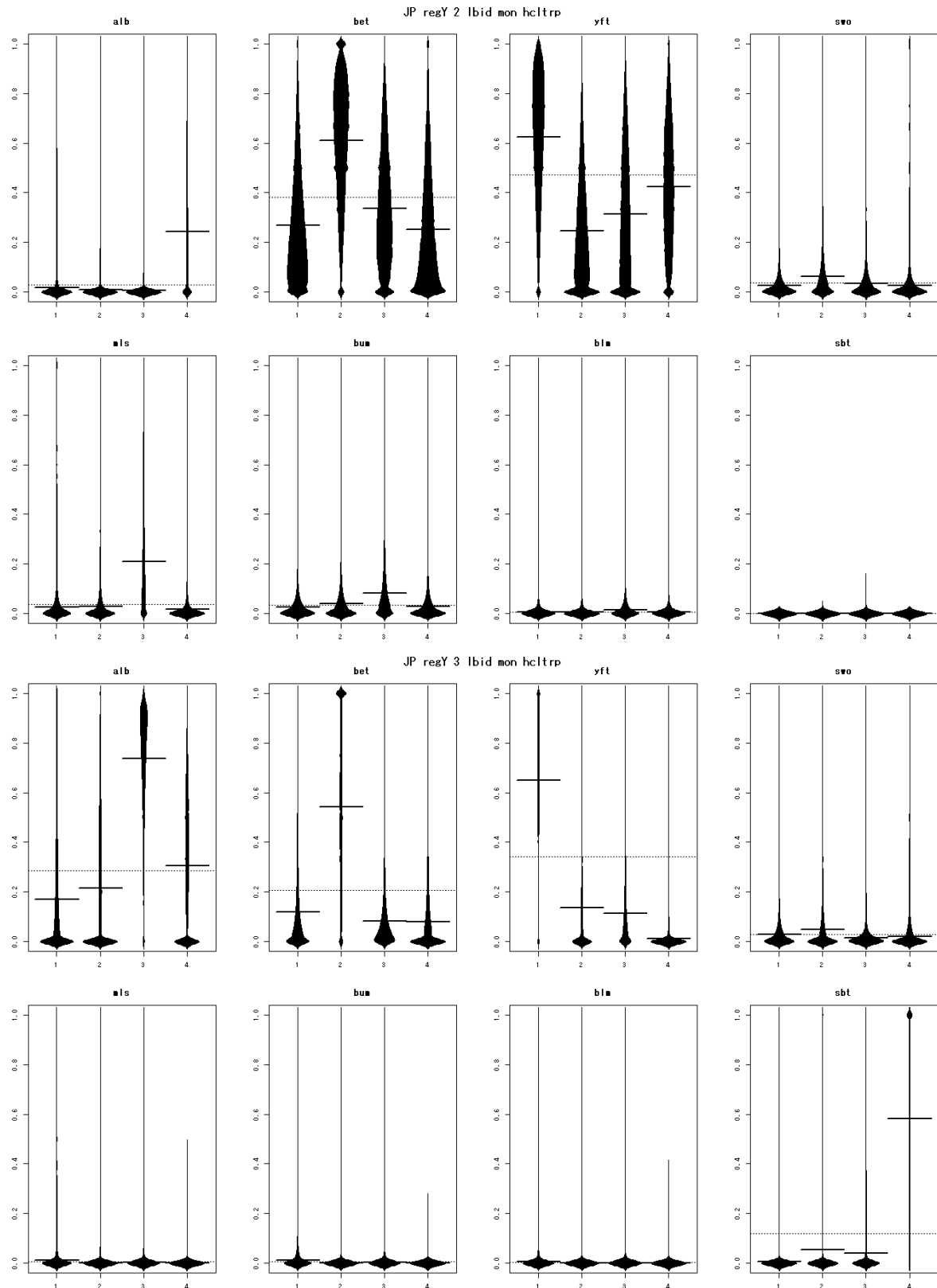


Fig. 4. Beanplots for yellowfin region showing species composition by cluster. The horizontal bars indicate the medians.

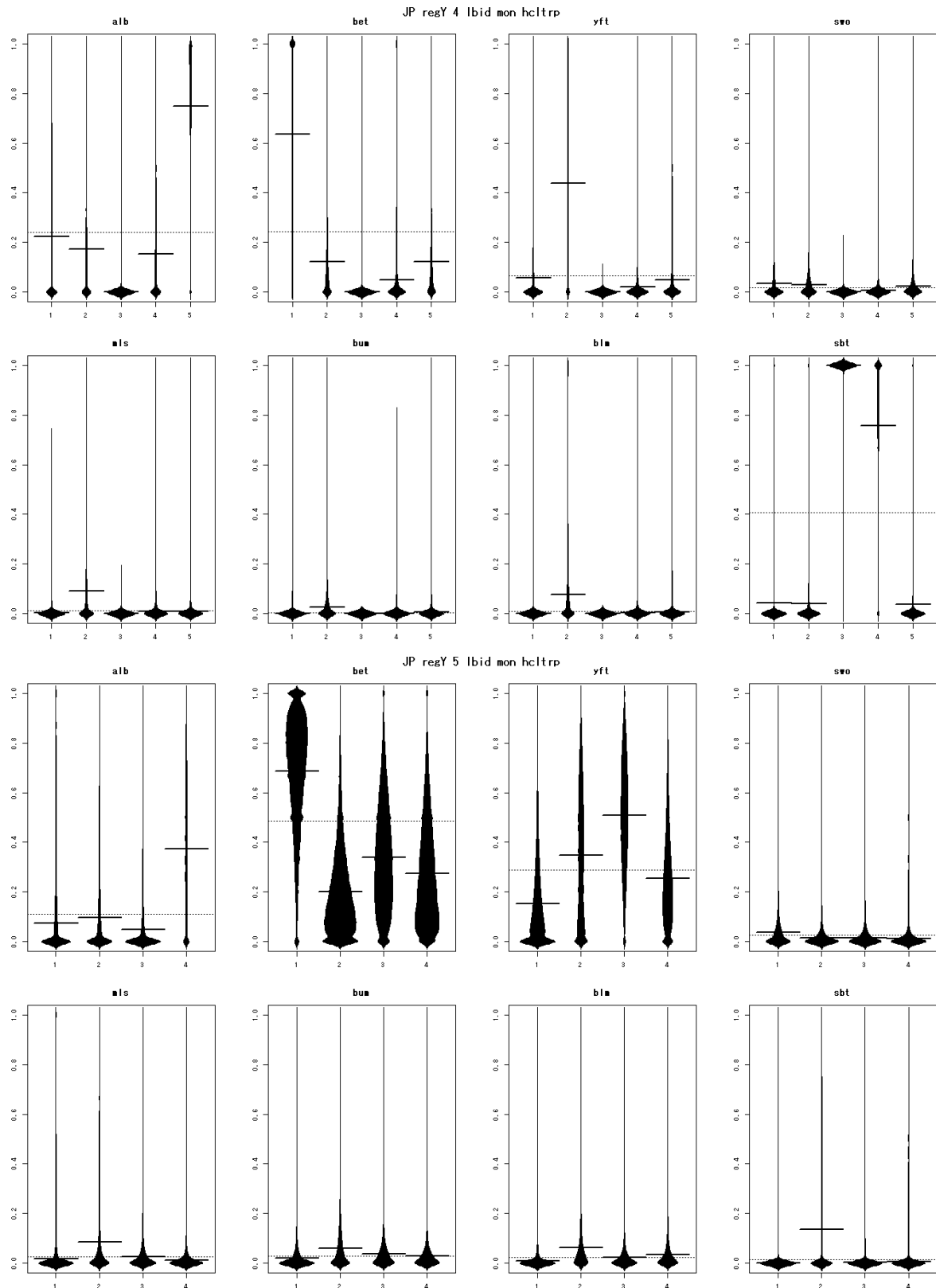


Fig. 4. Beanplots for yellowfin region showing species composition by cluster. The horizontal bars indicate the medians. (continued)

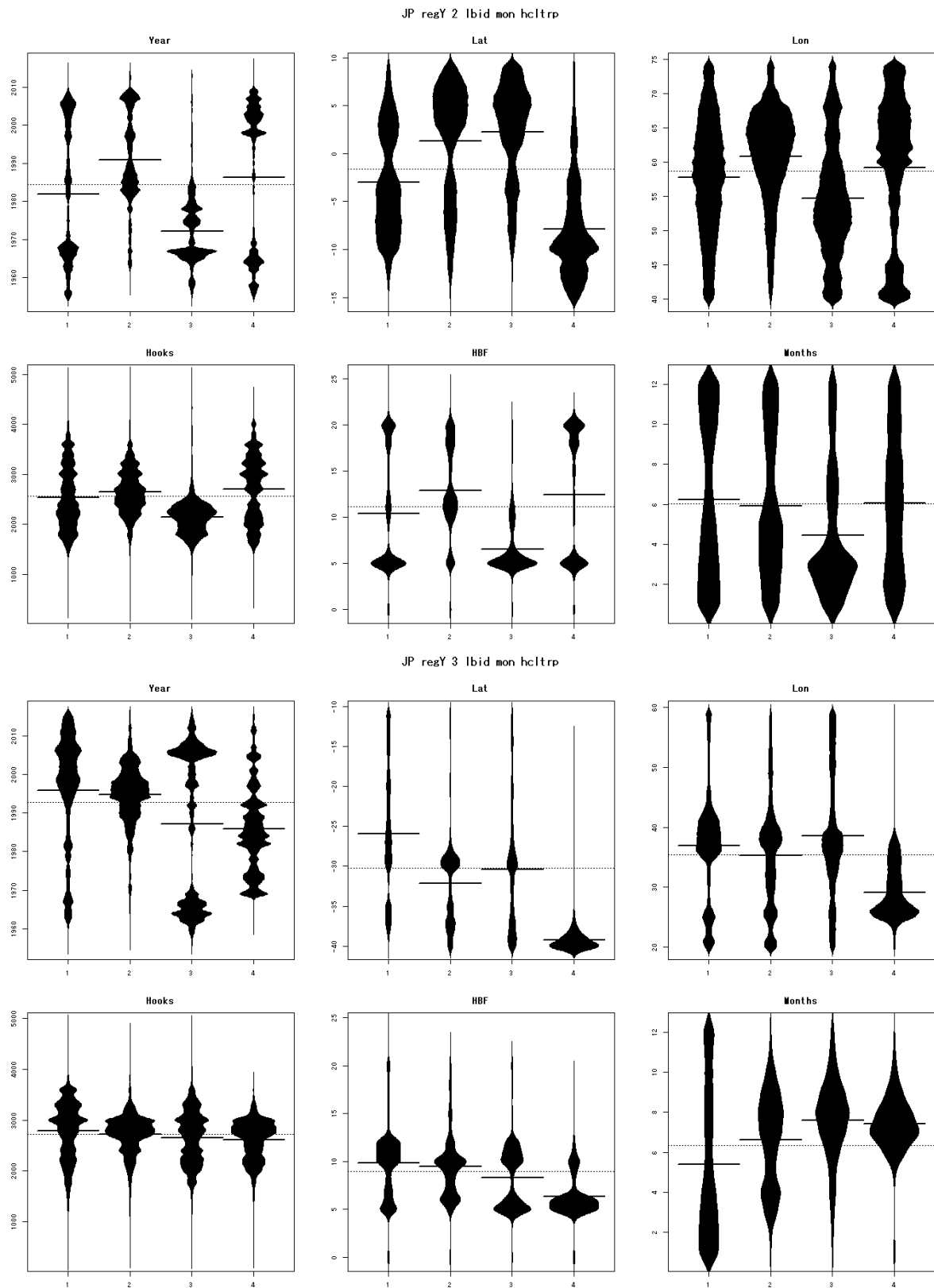


Fig. 5. Beanplots for yellowfin region showing number of sets versus covariate by cluster. The horizontal bars indicate the medians.

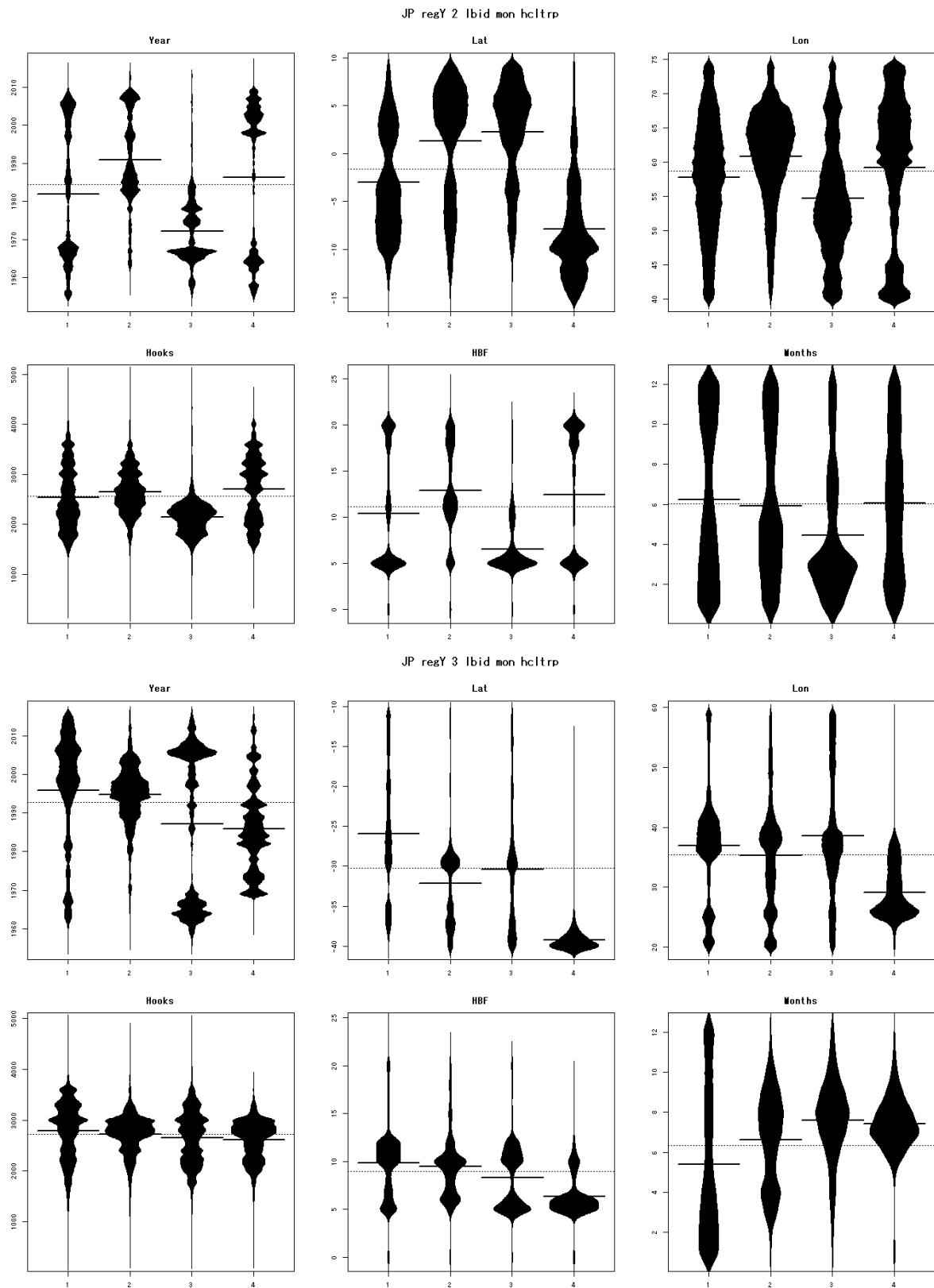
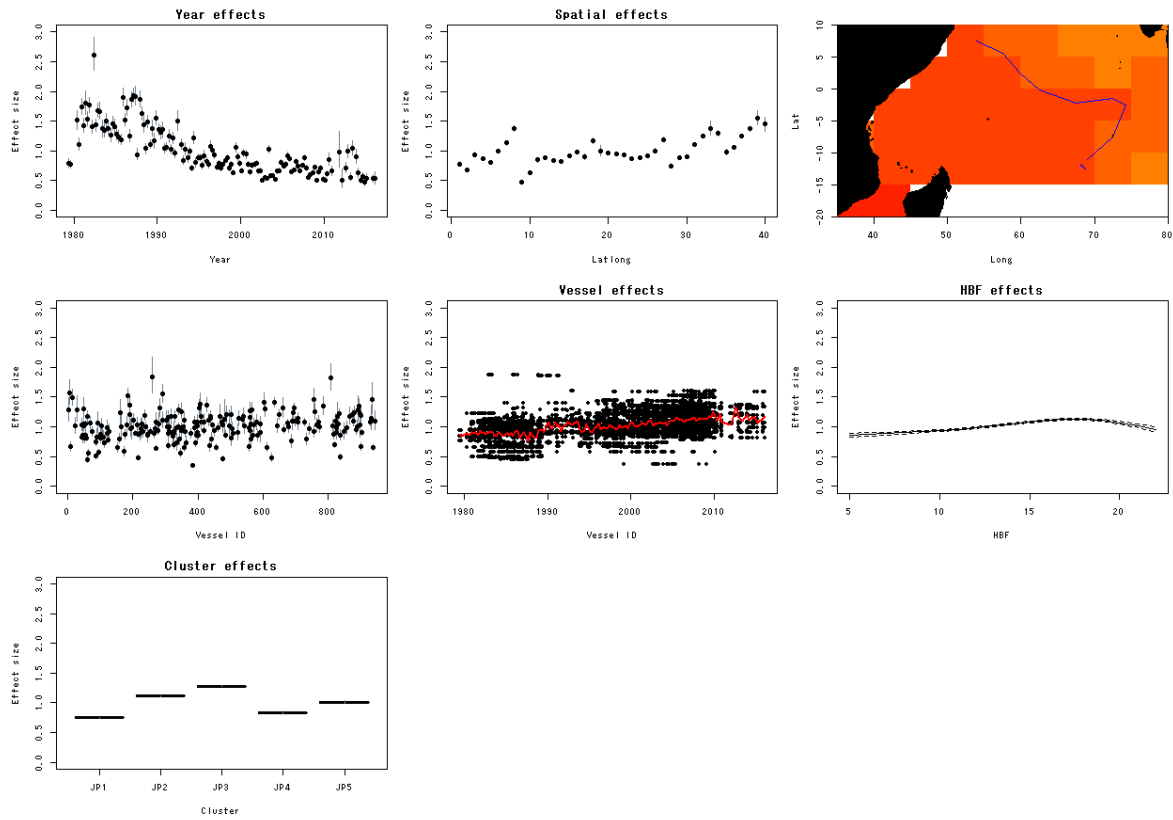


Fig. 5. Beanplots for yellowfin region showing number of sets versus covariate by cluster. The horizontal bars indicate the medians. (continued)

Region 1



Region 2

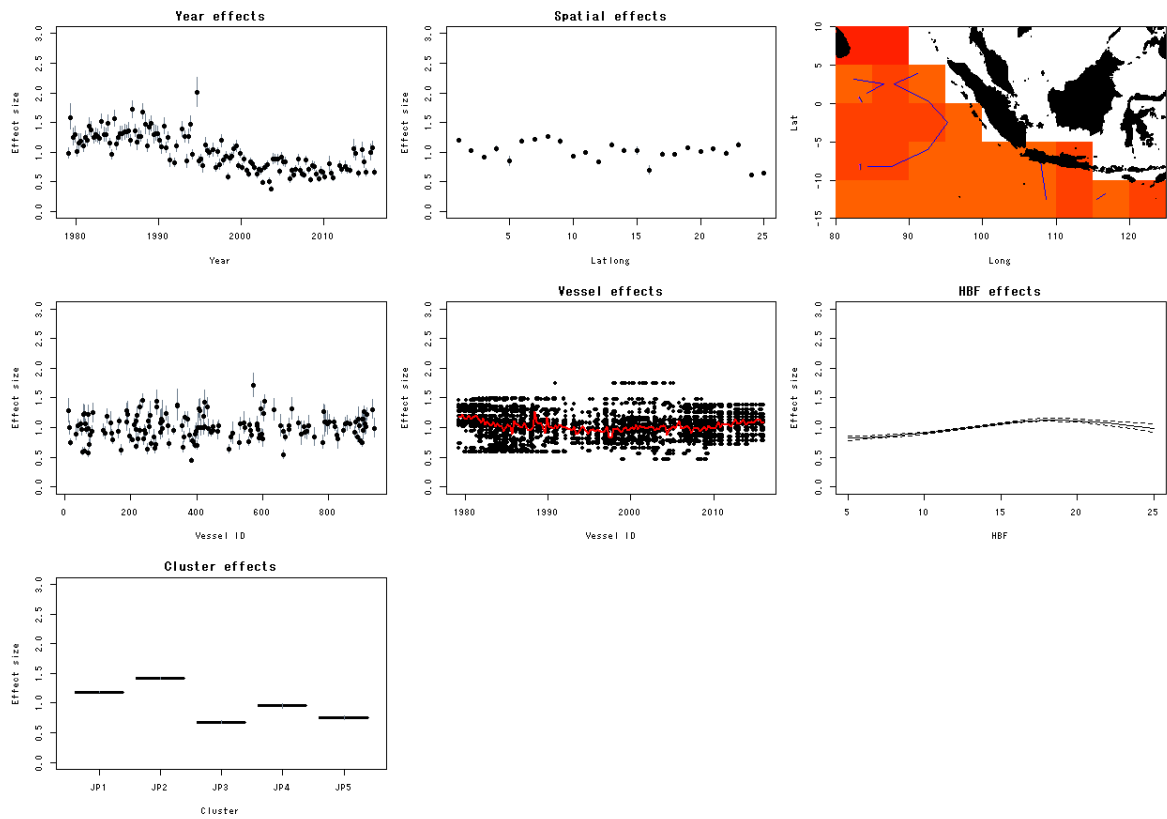
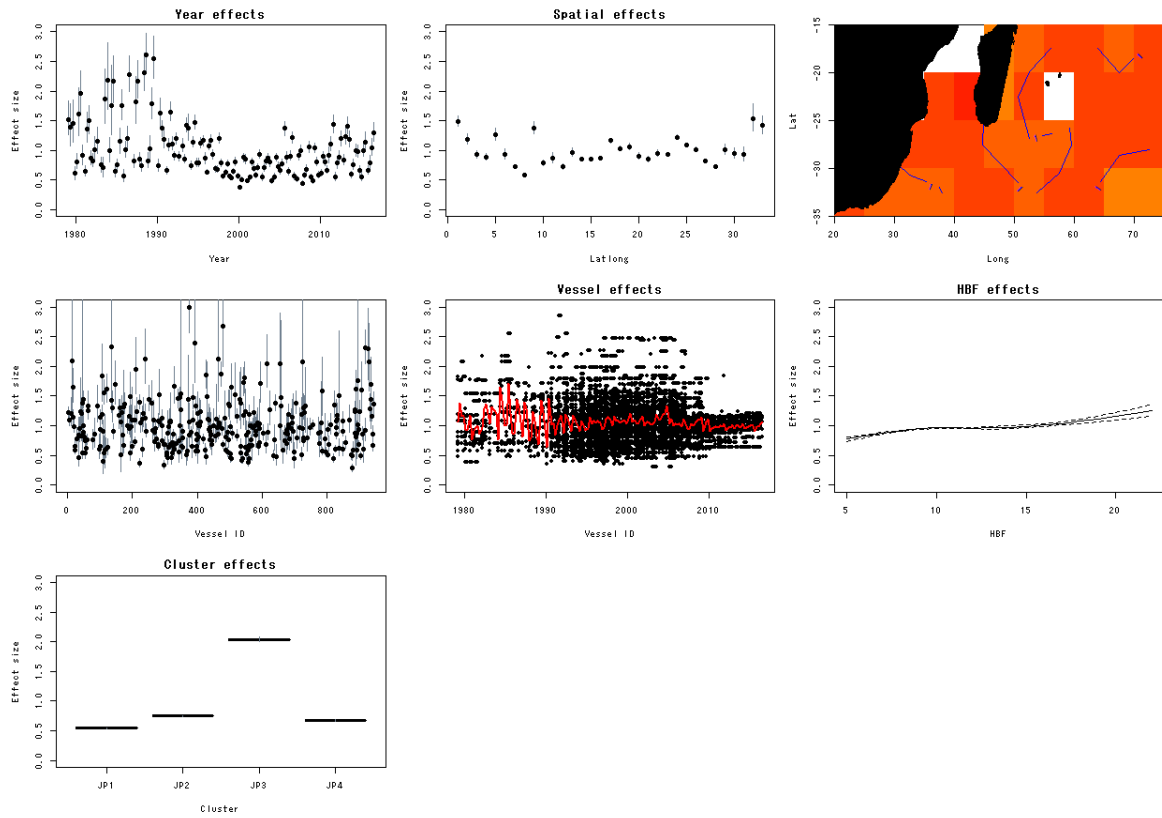
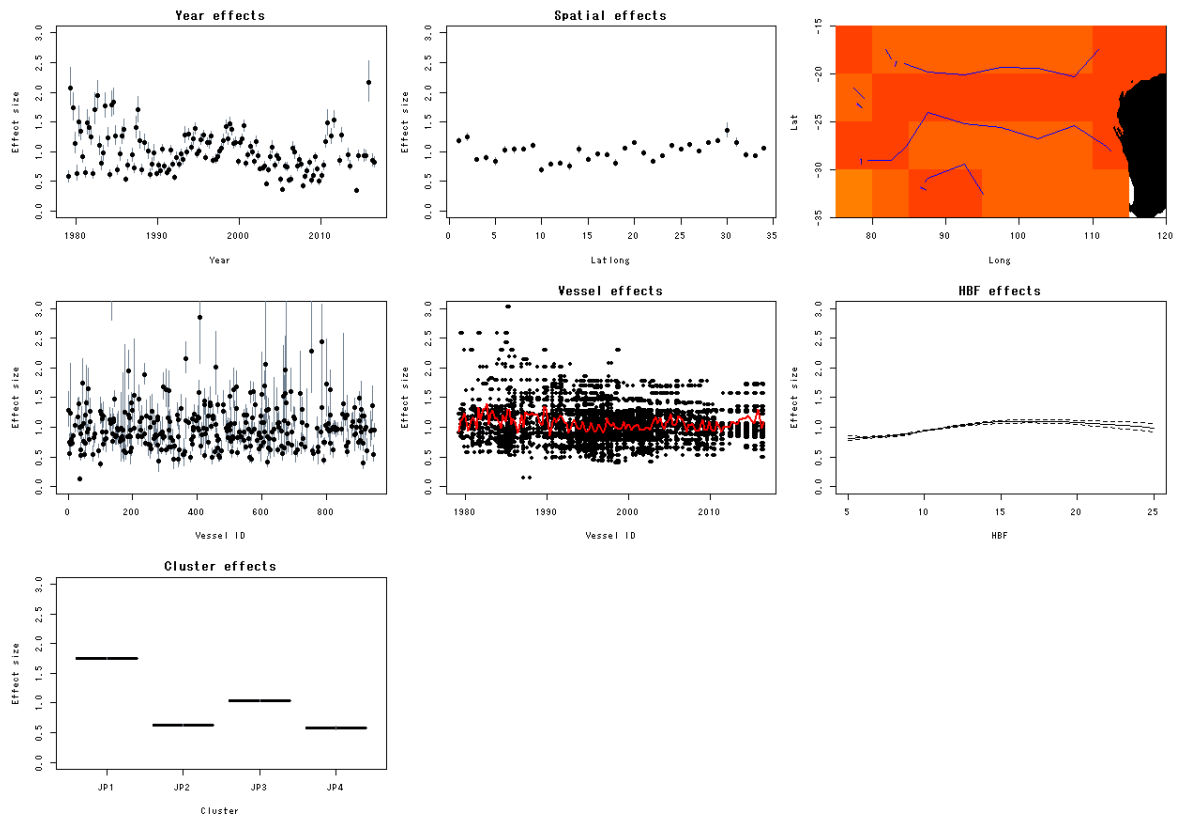


Fig. 6. Effect of each covariate for bigeye region (from 1979 onward with vessel ID).

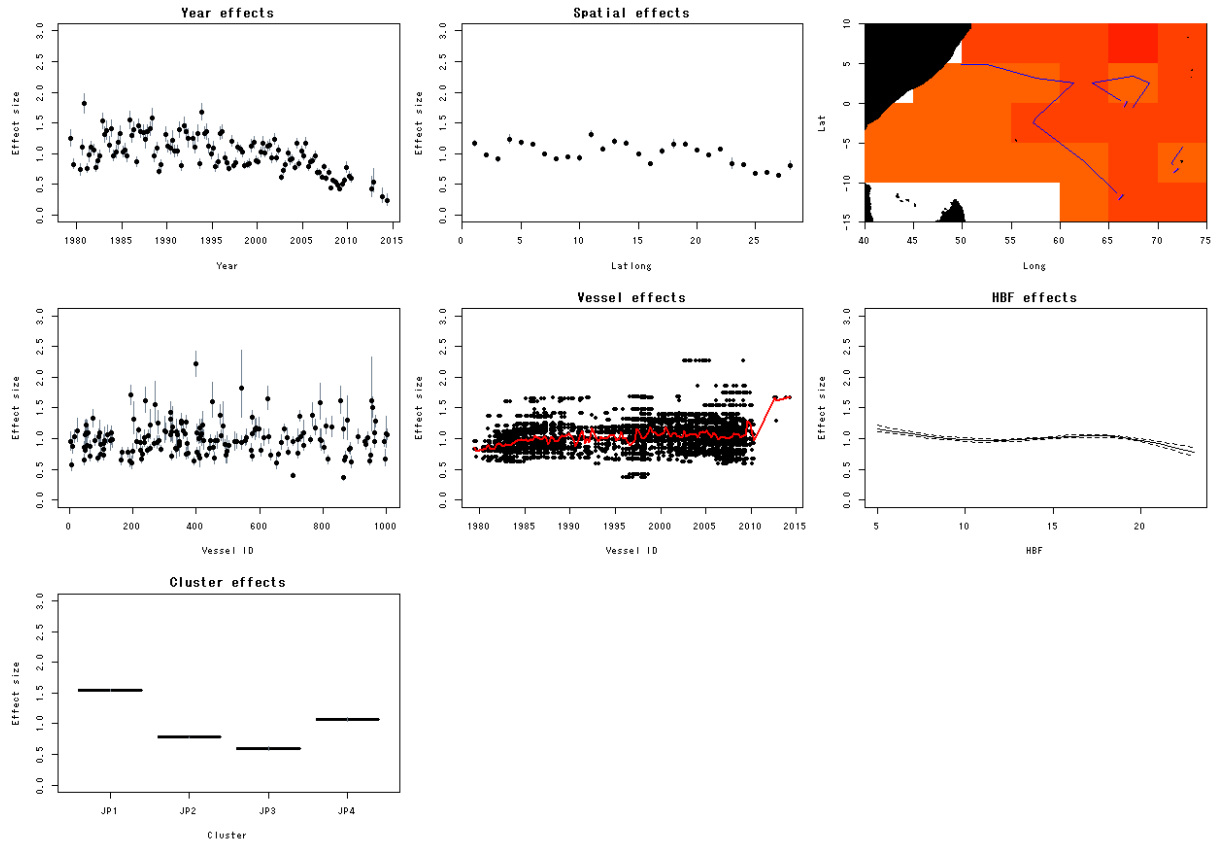
Region 3



Region 4

**Fig. 6.** Effect of each covariate for bigeye region (from 1979 onward with vessel ID). (continued)

Region 2



Region 3

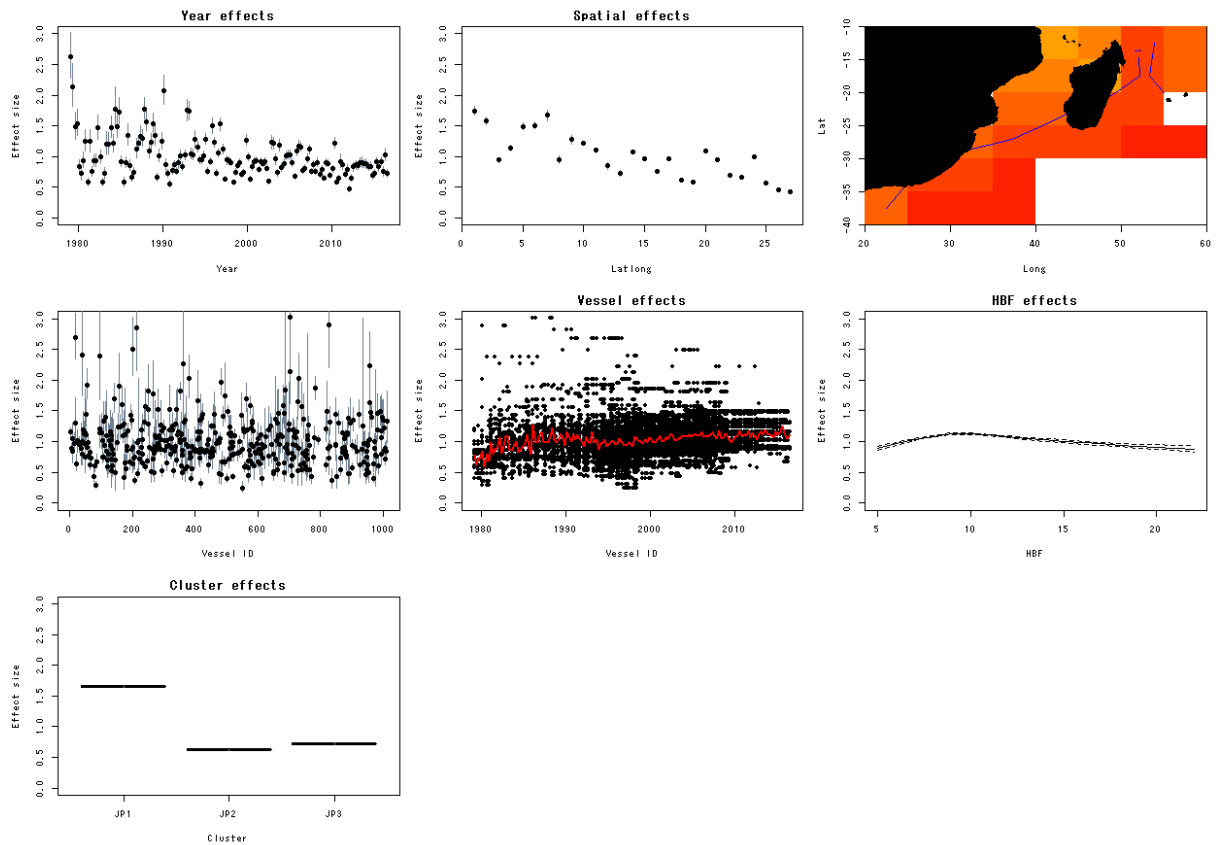
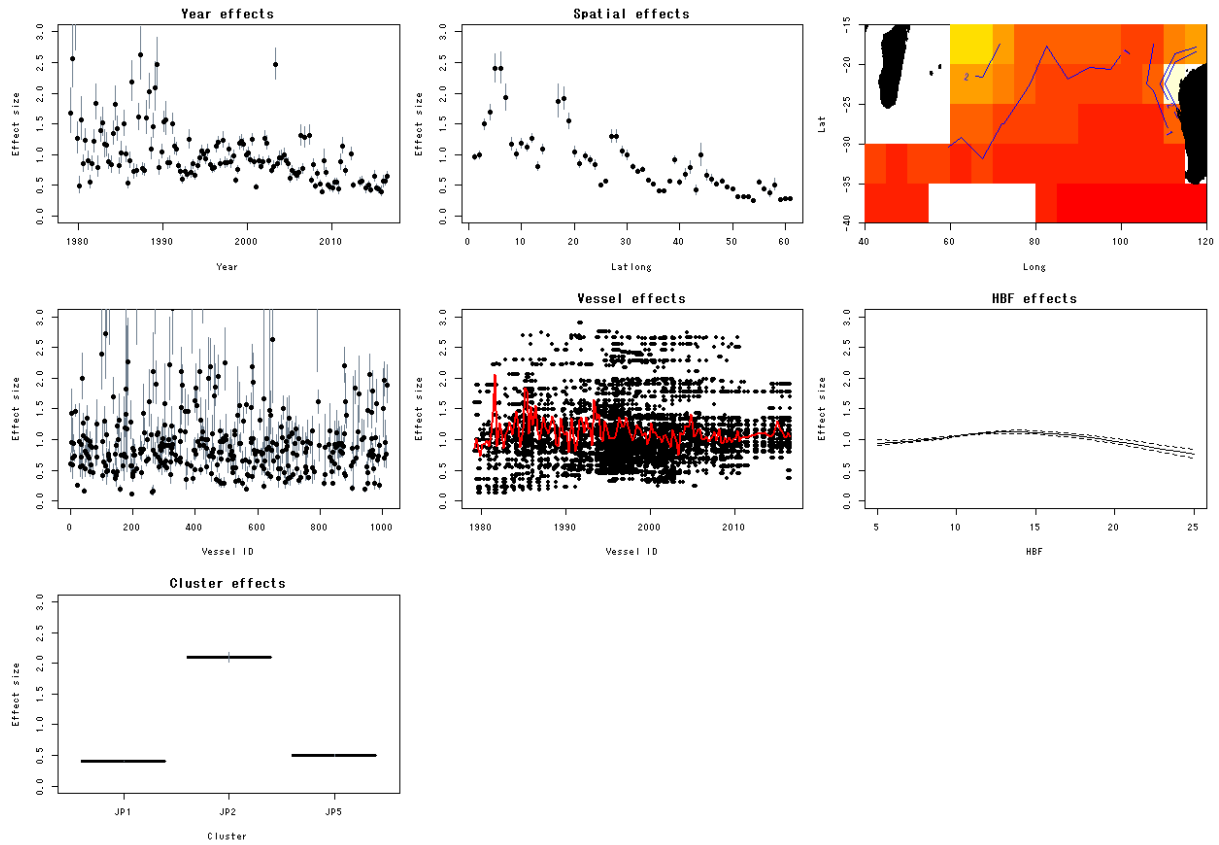
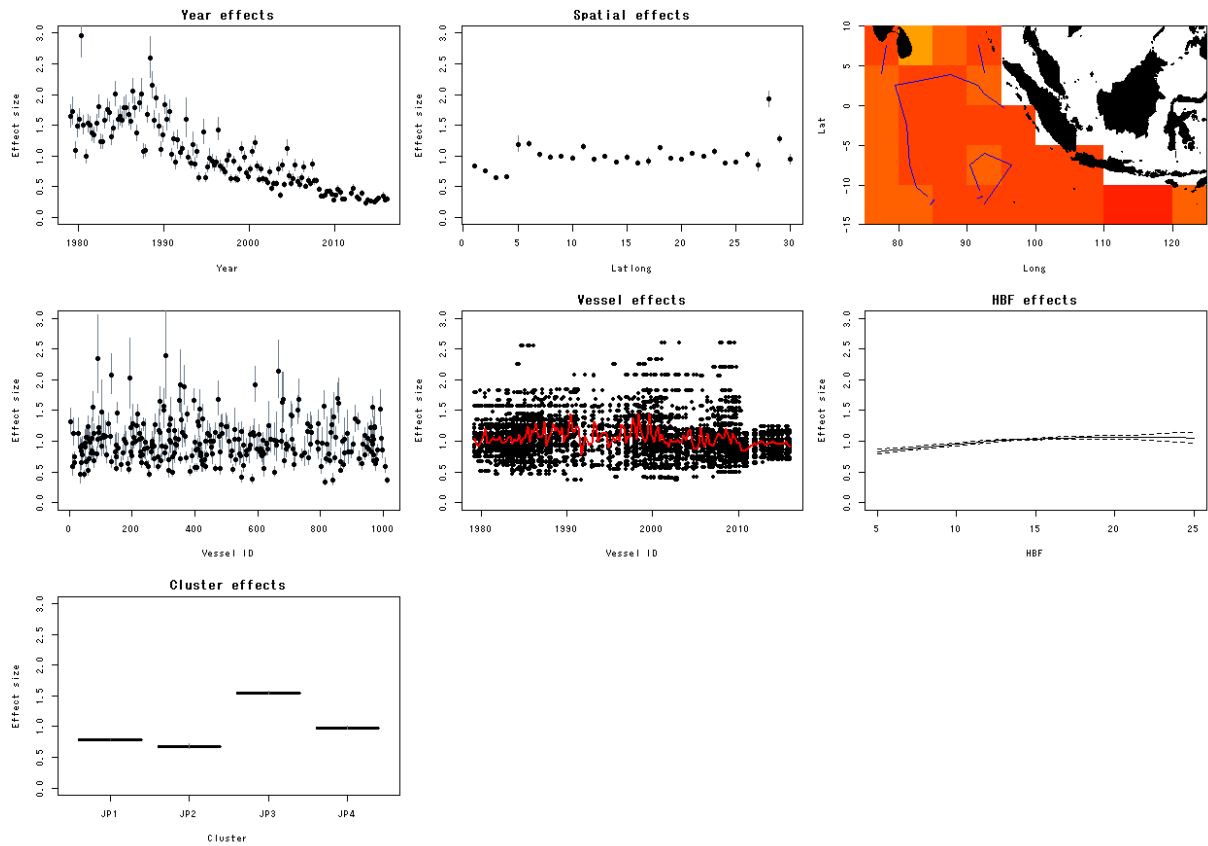


Fig. 7. Effect of each covariate for yellowfin region (from 1979 onward with vessel ID).

Region 4



Region 5

**Fig. 7.** Effect of each covariate for yellowfin region (from 1979 onward with vessel ID). (continued)

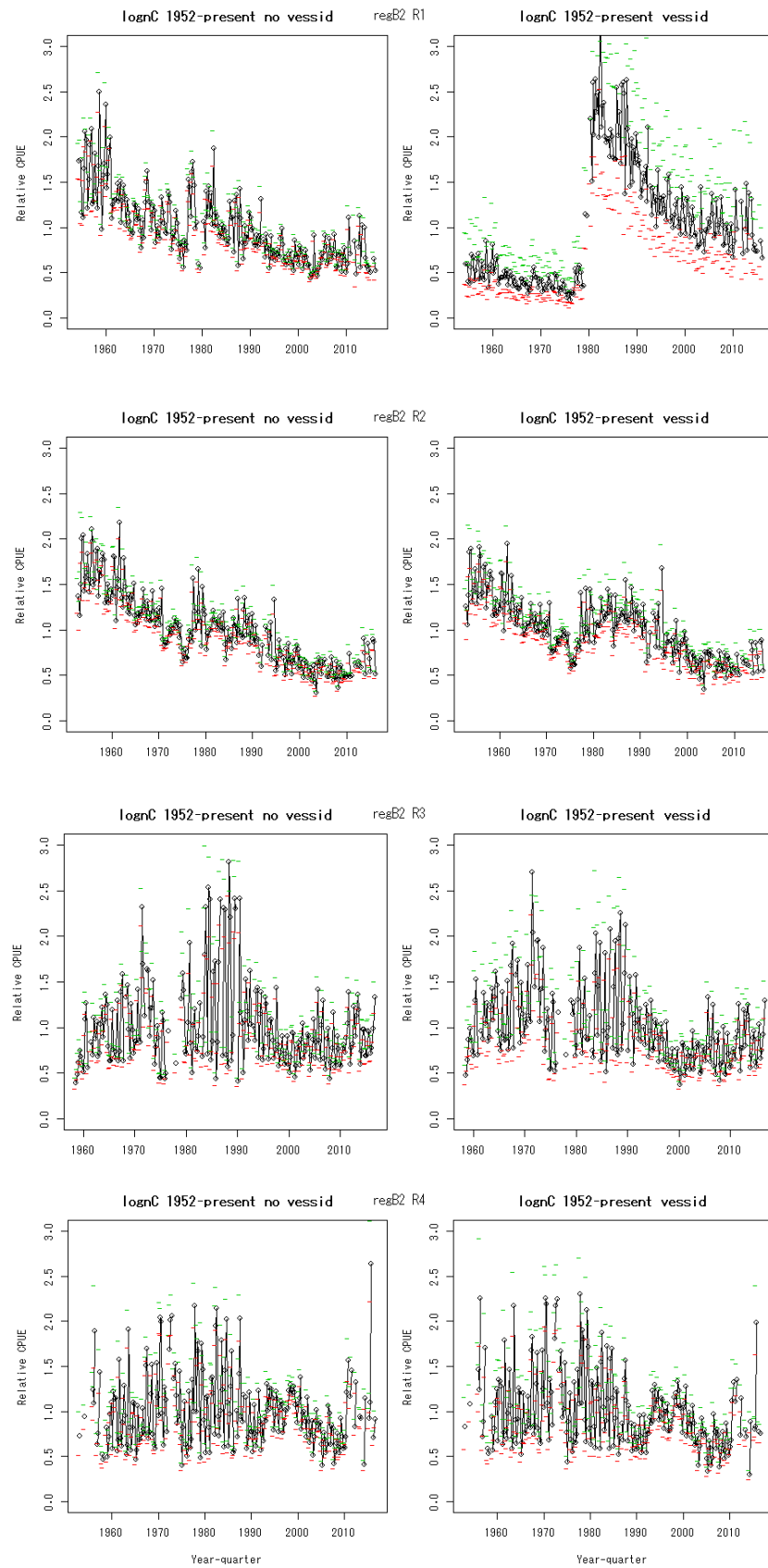


Fig. 8. Trend of quarterly CPUE (left: without vessel effects, right: with vessel effects) of bigeye.

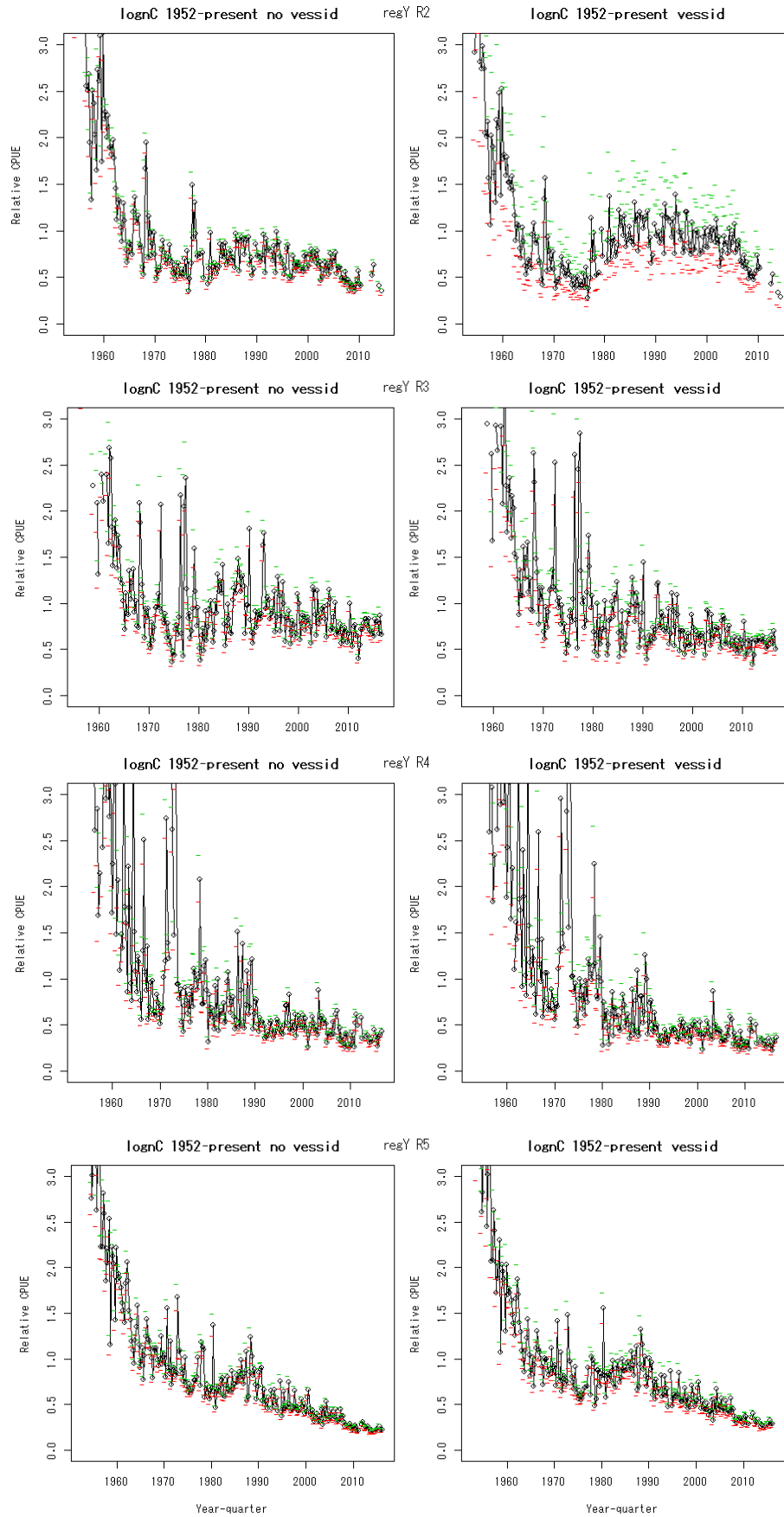


Fig. 9. Trend of quarterly CPUE (left: without vessel effects, right: with vessel effects) of yellowfin.

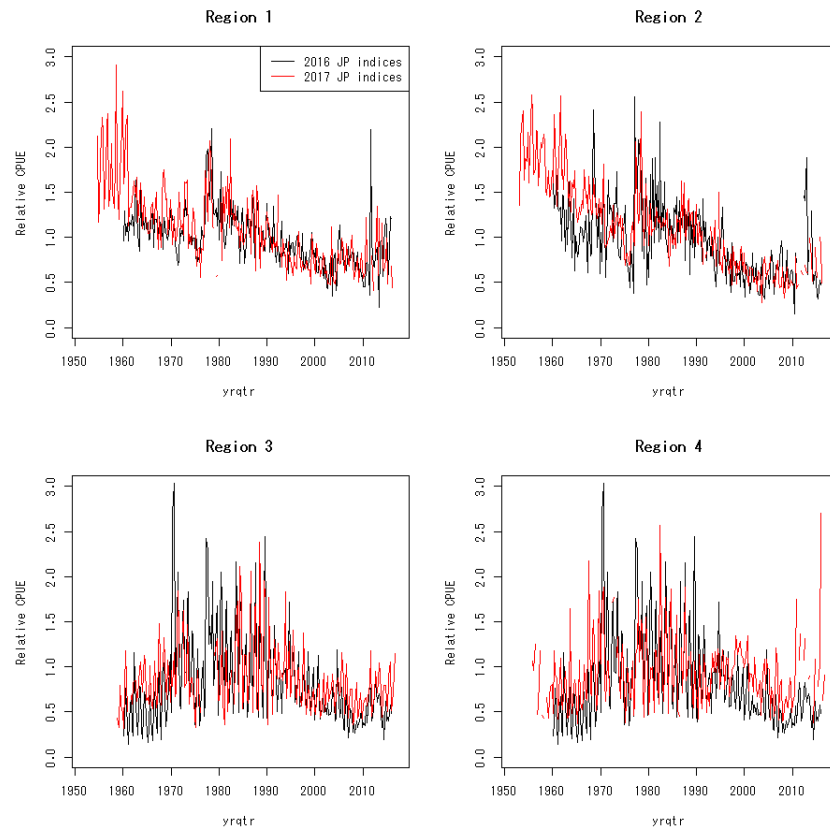
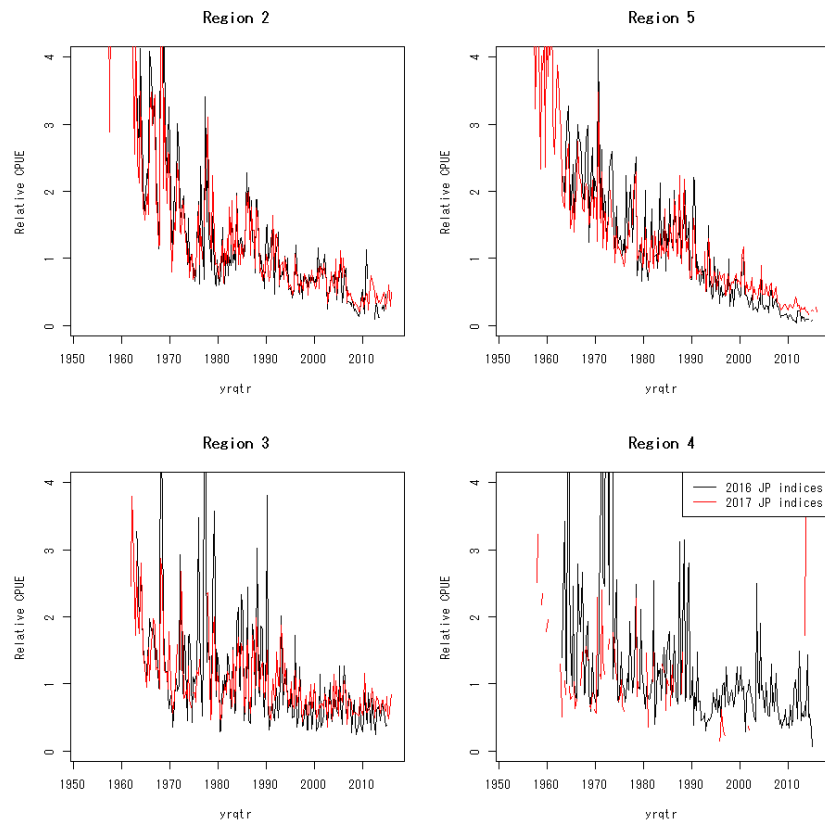
BET**YFT**

Fig. 10. Comparison of CPUE series of bigeye and yellowfin tuna with those reported in 2016.

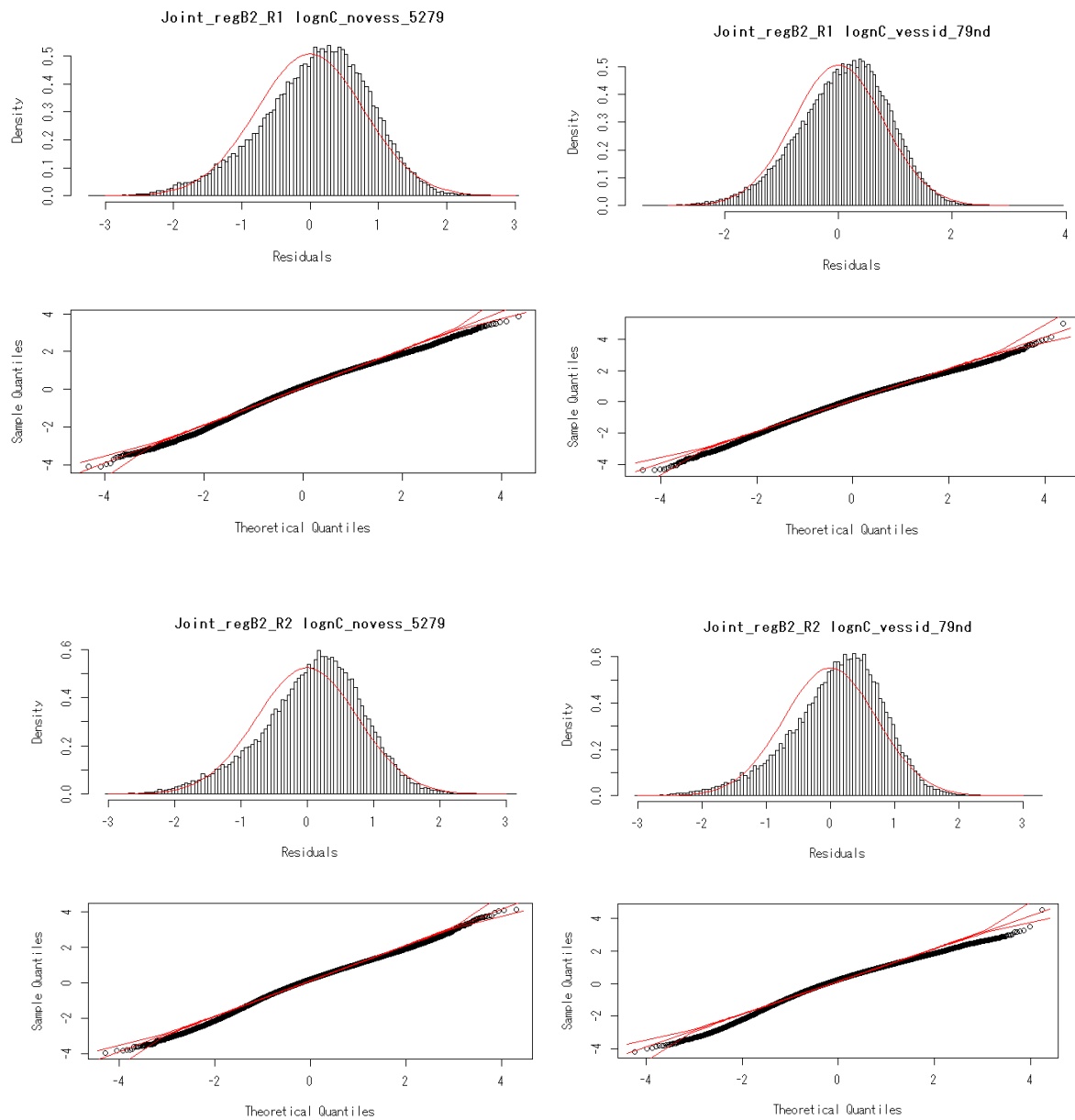


Fig. 11. Standardized residuals of CPUE standardization for bigeye.

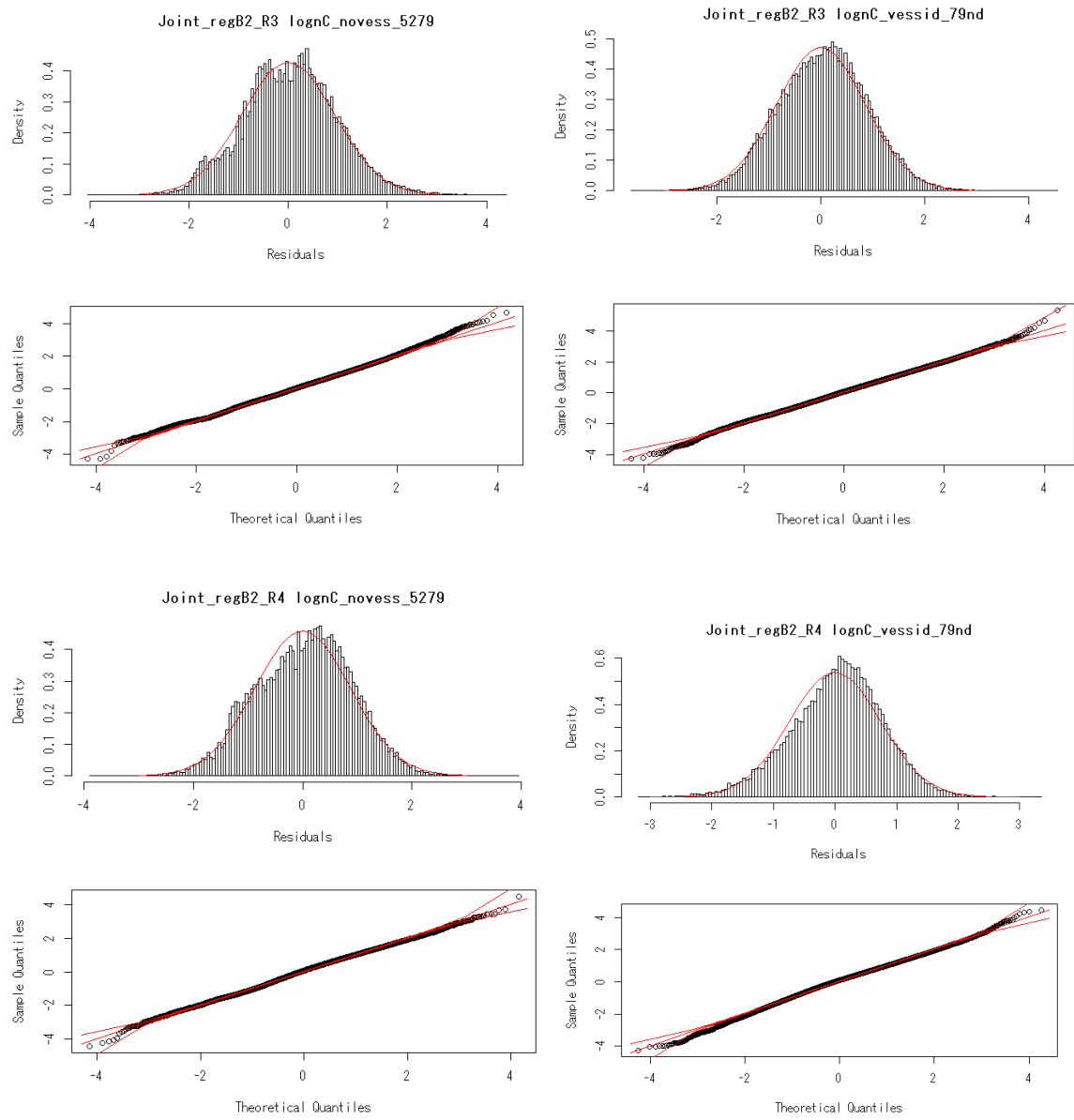


Fig. 11. Standardized residuals of CPUE standardization for bigeye. (continued)

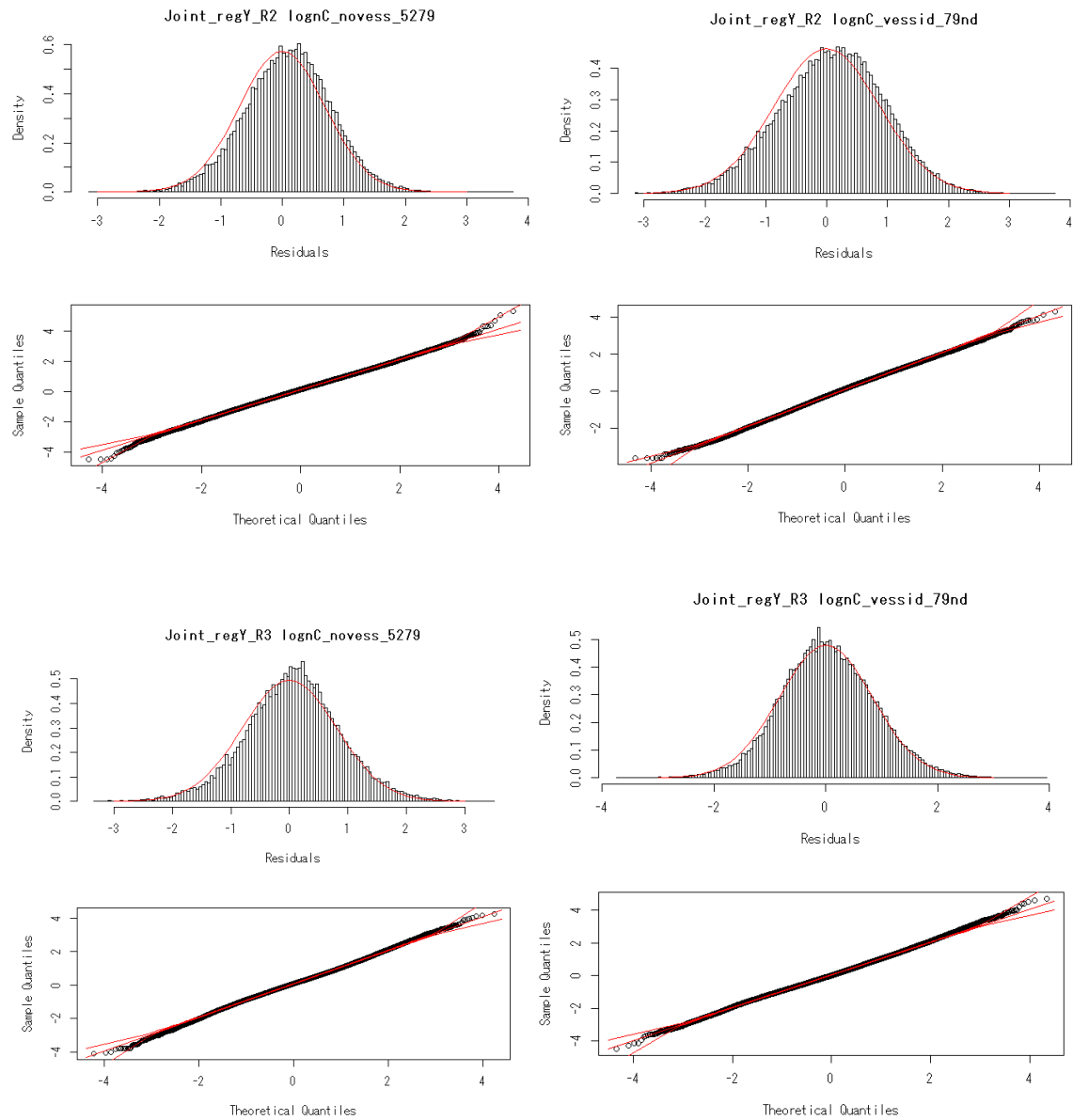


Fig. 12. Standardized residuals of CPUE standardization for yellowfin.

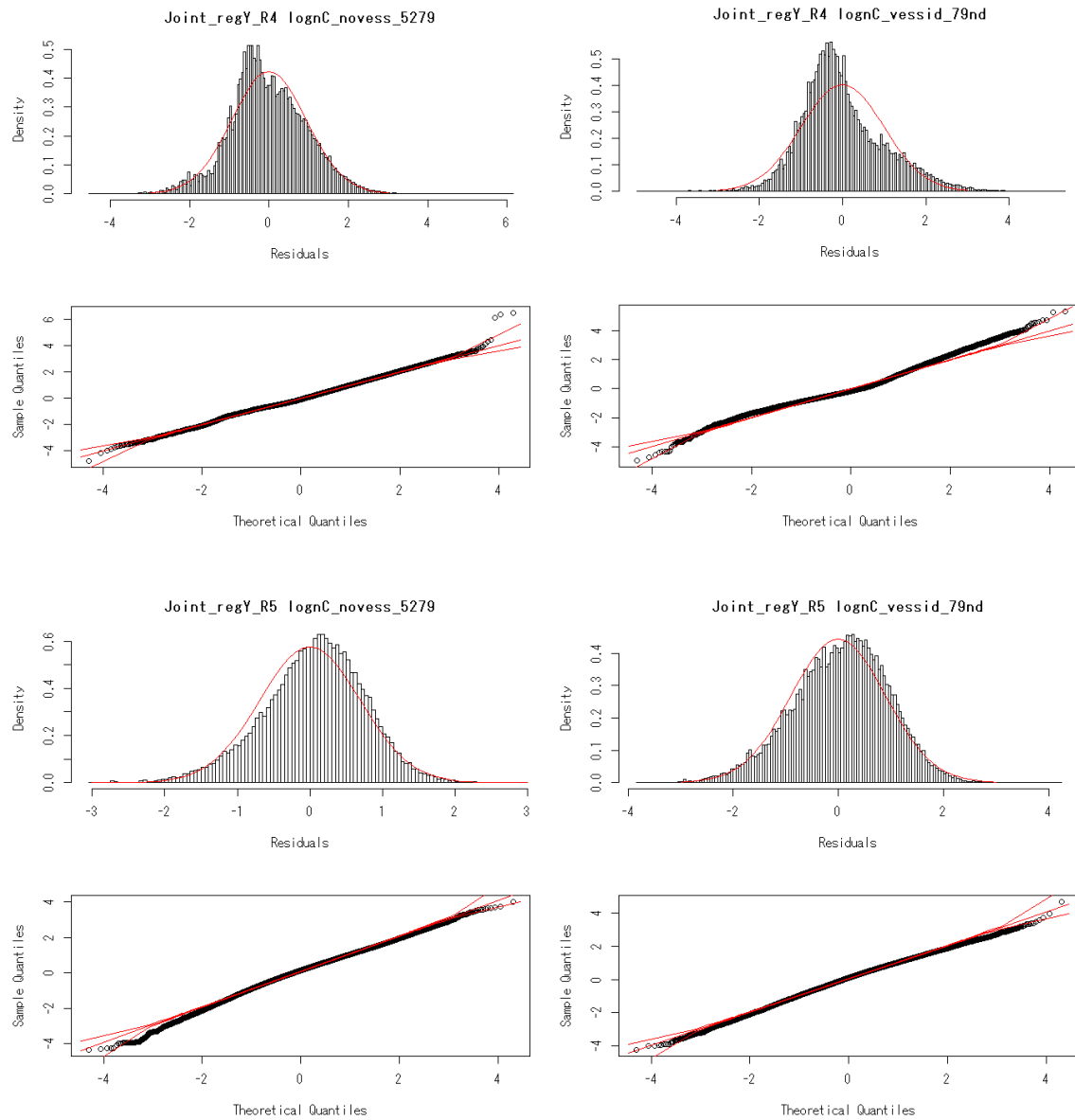


Fig. 12. Standardized residuals of CPUE standardization for yellowfin. (continued)