

# Cluster Results

## Concepts that need to be investigated

Much has been done thus far in regards to the research, and a clear picture of what is further required now exists. The work with SOMs is likely sound, but does require that a few more variables be tested. Furthermore, it is not yet certain that SOMs will be the best clustering technique. To that end, K-means clustering and hierarchical clustering have also been identified as alternative techniques. This now gives rise to the possibility that two papers could emerge from this work. One on the resultant clustering of synoptic air-sea states during coastal MHWs, and another paper that discusses the strengths of these various clustering techniques.

## The metrics for each MHW in each cluster

This requires that once the different events have been clustered, regardless of the technique used, or the variables controlled for within (see below), a summary of the event metrics must also be provided. These then will allow for the second more meaningful round of the interpretation of the results.

Table 1: The possible metrics that may be of interest for summarising the events clustered into each node. Node numbers given here correspond to Figure 1 and 2. (continued below)

node	count	summer	autumn	winter	spring	west	south	east
1	19	3	7	4	5	6	11	2
2	5	2	2	1	0	3	2	0
3	12	0	0	0	12	0	12	0
4	15	1	0	6	8	2	12	1
5	2	0	0	2	0	2	0	0
6	8	0	8	0	0	1	7	0
7	5	0	0	5	0	2	3	0
8	17	4	5	4	4	4	12	1
9	12	0	5	6	1	2	10	0
NA	95	10	27	28	30	22	69	4

Table 2: Table continues below

duration_min	duration_mean	duration_max	int_cum_min	int_cum_mean
15	24.4	65	20.57	68.242
19	25.2	43	45.97	70.895
15	32.0	47	45.26	89.687
15	24.8	35	23.77	59.027
27	27.0	27	41.49	47.452
15	19.2	25	23.73	51.747
18	26.0	31	33.67	47.867
16	27.9	98	24.49	75.261
15	18.8	27	23.87	43.202
15	25.1	98	20.57	64.829

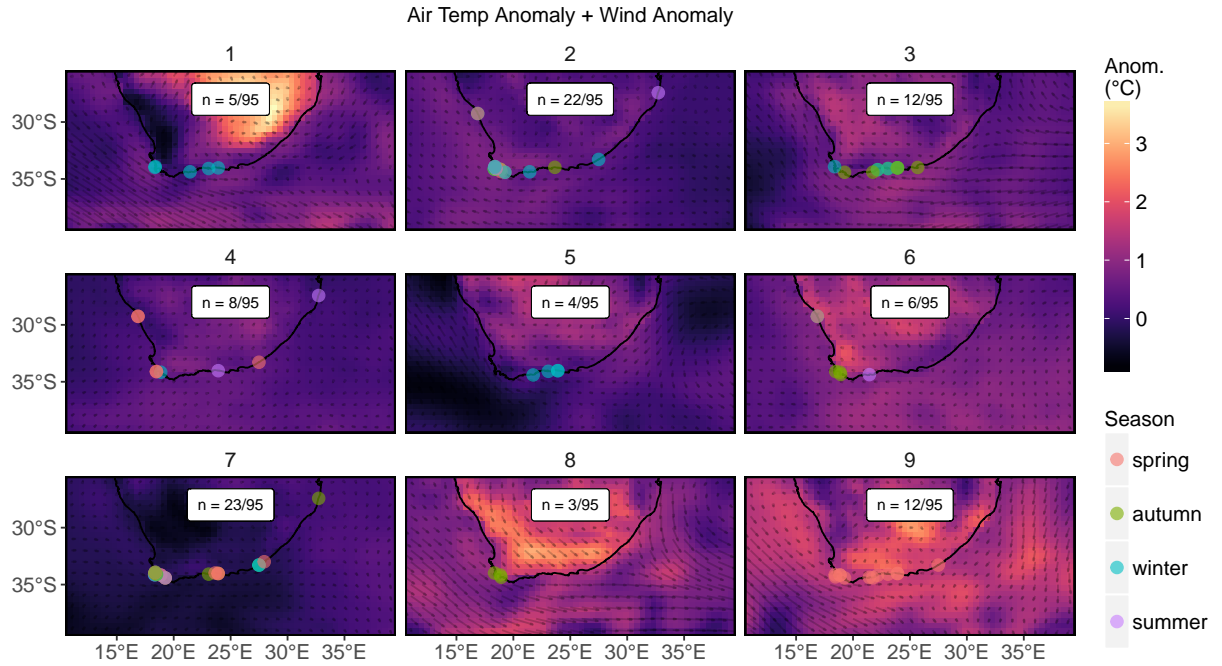
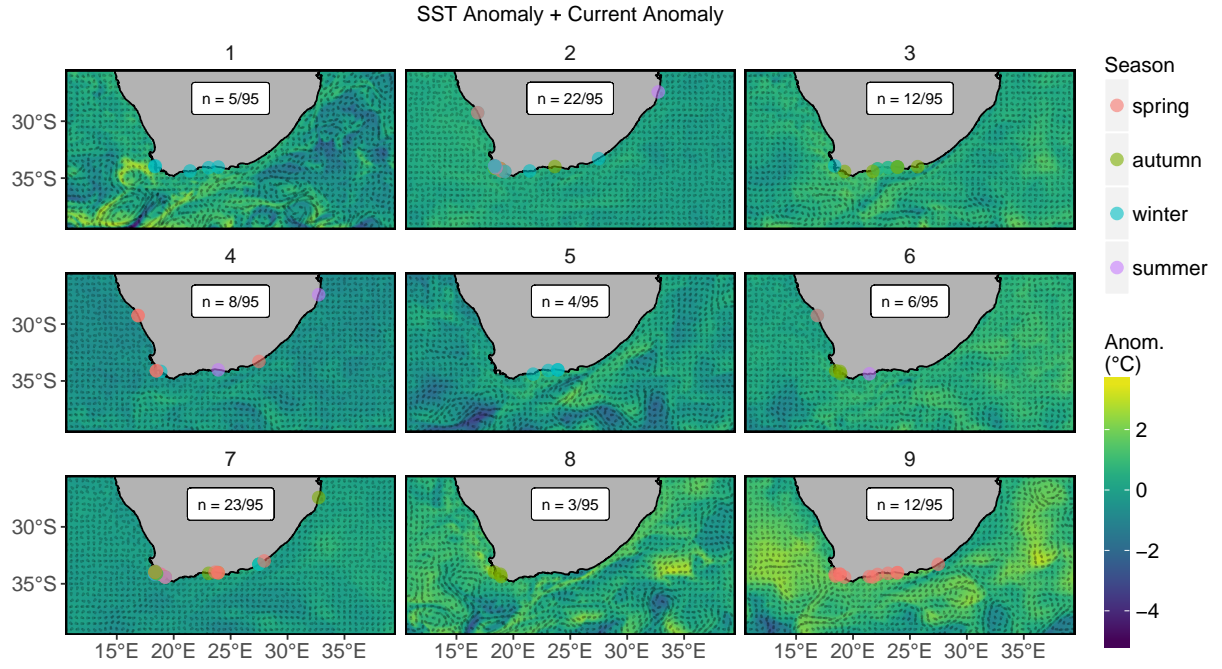


Figure 1: The results of a SOM clustering of the syoptic air-sea anaomaly data during coastal MHWs. The clusters shown here correspond to the following table and figure.

int_cum_max	int_max_min	int_max_mean	int_max_max
160.80	1.63	3.935	7.34
91.43	2.77	4.479	6.94
137.08	2.36	4.021	5.18
158.12	2.16	3.578	7.66
53.42	2.28	2.381	2.48
92.14	2.01	3.499	7.37
62.49	2.10	2.669	4.03
308.20	1.68	3.891	6.90
69.02	1.86	3.083	4.80
308.20	1.63	3.666	7.66

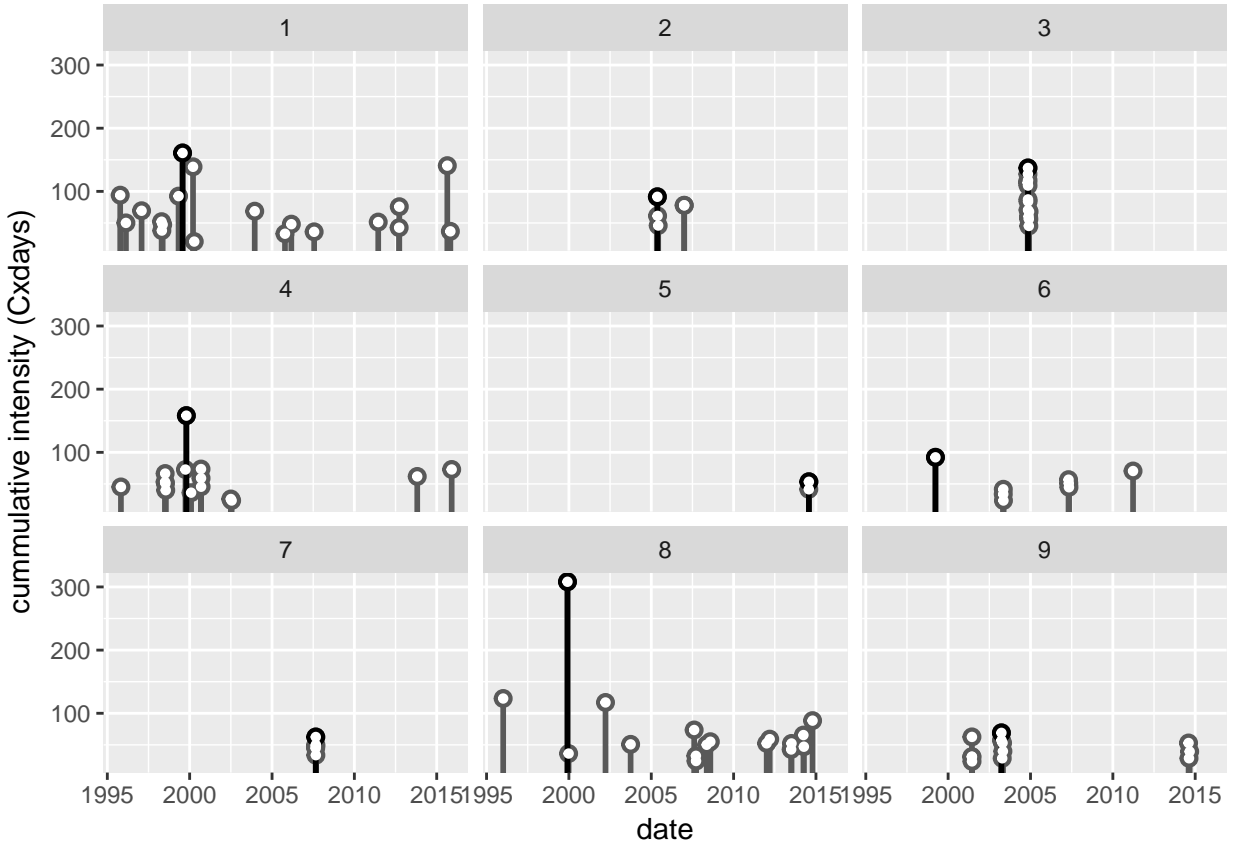


Figure 2: Lolliplots showing start date and cummulative intensity of each event by each node, as seen in Figure 1.

Figure 2 most clearly demonstrates that the SOM nodes in Figure 1 generally consist of either several events that occurred at the same time, or a blend of many disparate events. The nodes that contain more events are therefore much more ‘neutral’, meaning there is very little visible by way of air-sea anomalies. This is as I had feared and is currently my largest criticism of this technique. If one were to add more SOM nodes (a bad idea given that there are only 95 data vectors to be clustered), the new nodes will only allow some of the events within the larger clusters to break away and form their own node. The complexity of the air-sea states is such that any consistent patterns that may occur are obfuscated by everything else happening in the study area. Therefore it is necessary to reduce the dimensions of the input data by drawing a more narrow box around the study area to investigate the effect this may have. But first, the BRAN data need to be appropriately rounded down to a resolution of 0.5 degrees in order to match the ERA-Interim data.

## Effect of pixel resolution on clustering

The more dimensions/ variables one introduces to a cluster analysis, the more stress will exist in the results. As large stress values are generally considered to be a negative result in clustering, it is best to attempt to reduce it where possible. One way of doing that for this research is by reducing the pixel resolution of the reanalysis products. There are two reasons that this cannot simply be done out of hand. The first is that the reduction in resolution may affect the clustering of the events. So this must be documented. The other problem this presents is that the reduction of pixel resolution would require that any results produced be shown at this same reduced resolution. And because the goal is to show meso-scale forcing on the coast, higher pixel resolutions would be preferable. Regardless, the ERA-Interim data are at a resolution of 0.5 degrees, which requires that the BRAN data be reduced to this same resolution for appropriate cluster comparison. Beyond this initial required reduction in resolution, the question then is what effect does the further rounding of the data produce? Here we look at three resolutions: 0.5, 1.0 and 2.0 degree lon/ lat.

Table 4: Table showing the number of events clustered into which of the 9 SOM nodes. The different columns show the effect that reducing the resolution of the data has on the clustering.

res_all	res_0.5	res_1.0	res_2.0
19	5	12	7
5	13	9	12
12	12	18	16
15	7	4	4
2	17	4	4
8	8	25	10
5	10	10	12
17	7	9	8
12	16	4	22

The results table generated from the clustering of events into different nodes shown above is not very informative because it is known that the SOM algorithm always reshuffles these data into different nodes. Due to the very high dimensions of the data, the algorithm is not able to find a single best answer. Therefore, it is more informative to further order each column from highest to lowest so as to see if the general clustering of events is similar.

Table 5: Table showing the number of events within a node scored in descending order. The different columns show the effect that reducing the resolution of the data has on the clustering.

res_all	res_0.5	res_1.0	res_2.0
2	5	4	4
5	7	4	4
5	7	4	7
8	8	9	8
12	10	9	10
12	12	10	12
15	13	12	12
17	16	18	16
19	17	25	22

When the data are reshuffled into descending order based on the number of events clustered into each node

we see that the results are much more similar than they first appeared. The next step in this portion of the analysis is not the creation of figures, but rather clustering events that are generally clustered together. This then would allow for the comparison to be scaled up so it could be replicated 1,000 times to be more thorough.

## Comparing large replication sets

Because the SOM nodes are shuffled during each run, comparing the results directly becomes difficult. We can however iterate the SOM analysis many times and save a vector of results for each event showing into which node it was cast for each run. This then creates a unique information vector for each event that can then be used to further cluster the events into ‘mean’ clusters. This helps to address the issue of SOM node ‘drift’. The following two tables show the results of a SOM run on the data 10 times.

Table 6: Table showing the number of events within the same node over 10 runs on the same data. Note how variable the node assignments may be.

1r	2r	3r	4r	5r	6r	7r	8r	9r	10r
16	7	12	18	12	17	18	13	16	10
18	9	8	16	7	5	9	18	9	7
6	16	6	9	6	6	12	6	12	12
11	4	8	12	8	9	6	10	11	6
6	4	7	4	6	17	21	7	10	5
12	17	9	9	6	8	7	7	8	28
12	12	13	6	17	13	5	12	18	12
7	9	26	6	17	8	10	10	5	7
7	17	6	15	16	12	7	12	6	8

Table 7: Table showing the number of events within a node scored in descending order. The different columns show each run. Note that the results are generally consistent, but not completely.

1r	2r	3r	4r	5r	6r	7r	8r	9r	10r
6	4	6	4	6	5	5	6	5	5
6	4	6	6	6	6	6	7	6	6
7	7	7	6	6	8	7	7	8	7
7	9	8	9	7	8	7	10	9	7
11	9	8	9	8	9	9	10	10	8
12	12	9	12	12	12	10	12	11	10
12	16	12	15	16	13	12	12	12	12
16	17	13	16	17	17	18	13	16	12
18	17	26	18	17	17	21	18	18	28

Now let’s up the anti a bit and run the SOM 100 times. Because the resultant data frame will become unwieldy, we will create means across the number of clustered events as coerced into descending order. These then will be compared to the mean for the 1 run and 10 run data frames.

Table 8: Table showing the number of events within a node scored in descending order. The different columns show the mean of 1, 10, and 100 SOM runs. Note that the results are remarkably similar.

1r	10r	100r
6	5	5
6	6	6
7	7	7
7	8	8
11	9	10
12	11	11
12	13	13
16	16	15
18	20	20

As we may see from the table above, the ‘wobble’ present in the SOM analysis is not great. I considered running the analysis 1,000 times but this would take several hours on my computer and I think that the consistency shown between 1, 10, and 100 runs is sufficient to put this question to rest.

## Effect of lat/ lon extent on clustering

With more traditional cluster analyses, the values being compared would have far fewer dimensions. In this regard one would endeavour to only include variables that seem relevant to the question being asked. For example, if clustering different rock pools by the species found within them, one would likely create better results by not including any anomalous species found in the results (like a cow fish). In regards to this work, it is best to include only the pixels that are likely relevant to the meso-scale features that may be impacting the coast. More specifically, cutting out the Agulhas retroflection above the Southern Ocean will prevent any behaviour there from affecting the clustering of events that are occurring along the coastline of South Africa.

That all being said, it may indeed be relevant, at least in the sense of potential teleconnections, to include the Agulhas retroflections over the Southern Ocean. Therefore, I have decided to trim the total study area by 1 degree on the East, South and West extents, rather than just the South extent. The North border of the study area is left unchanged as this would begin to remove some of the coastal stations. The effect of trimming the study area 1 degree at a time is presented below. A SOM is run on each trimmed data frame 10 times to produce smoother results. The inquiry into the SOM ‘drift’ above shows that 10 iterations return comparable results to 100 iterations, so 10 are used here in the interest of speed. A total of 5 degrees are iteratively trimmed.

Table 9: Table showing the number of events within a node scored in descending order. The different columns show the mean of 10 SOM runs on differing amounts of spatial reduction in the extent of the study area. Note that the results remain similar, but there appears to be break in the results once 3 degrees of lon & lat have been removed from the East, South, and West edges of the study area.

trim_0	trim_1	trim_2	trim_3	trim_4	trim_5
5	4	5	6	6	6
6	6	7	7	7	7
7	7	7	7	8	7
8	8	8	8	9	9
9	9	10	10	10	10

trim_0	trim_1	trim_2	trim_3	trim_4	trim_5
11	11	12	12	11	12
13	12	13	13	12	12
16	15	14	14	15	15
20	22	20	18	18	18

The table above shows that removing portions of the study area does have an effect on the SOM clustering of the data. Leaving the study area ‘as is’ is shown in the column labeled ‘trim\_0’. These are the results shown in the previous tables and figures. Removing one degree of lat/ lon appears to cluster the events more into the one large cluster, with fewer events in the smaller clusters. Removing 2 degrees of lat/ lon reverses this trend. Trimming 3 to 5 degrees of lat/ lon then all appears to produce very similar results. The difference between these different extents is not large, but I think it does show that the edges of the study area are having an effect on the clustering. And that reducing the area does allow for a more even clustering of the events into nodes. It may be best to use a study area with 3 degrees trimmed from the East, South, and West edges. But for now I shall continue to use the full study area.

## Effect of running air and sea variables separately

It may be that air and sea values work in tandem with one another to force MHWs, but it is more likely that they do not. Except for perhaps VERY extreme situations (e.g. once per decade). Therefore it is necessary to run all clustering techniques on air-sea values combined, as well as separately, in order to quantify the potential effect they have on clustering. Up until this point all variables have been run together, here we pull them apart to see how the results may differ. We do this by running all air or sea variables together and then by running air or sea temperature and wind/ current vectors separately.

Table 10: Table showing the number of events within a node scored in descending order. The different columns show the mean of 10 SOM runs on different subsets of the data. The first column ‘all\_all’, shows the previous results of 10 SOM runs on the full, unsubsetted, untrimmed, etc. data. Note that the full air data (i.e. temp, U and V) appears the most similar to the normal results, and that the sea current data provides the most similar results.

all_all	air_all	air_temp	air_uv	sea_all	sea_temp	sea_uv
5	4	5	4	4	4	4
6	6	7	4	5	6	5
7	7	8	6	6	8	5
8	8	10	7	8	9	6
9	11	11	8	9	11	8
11	12	12	9	11	12	9
13	14	12	14	13	13	10
16	15	13	20	17	14	14
20	18	18	23	21	17	34

As we see in the table above, separating the variables from one another has a large effect on the clustering of the data. Unexpectedly, the air data appear to produce results the most similar to the results generated by running the SOM analysis on the full data set. This suggests that the air data (possibly the UV more than the temperature) are driving the SOM clusters, and not the sea data. My general perception of all of the tests and results thus far has been that it is the sea, and not the air, that is usually driving events. Which makes this result a bit trickier to deal with. The importance of air data on clustering the data must

be acknowledged, but I also think that this provides a clear argument against clustering the data using all of the variables at once because we generally want to see what the effect the ocean temperature is having on the coastal MHWs. One reason why air temperature and winds may be having a stronger effect on the clustering is that they are much more even, allowing for easier clustering of the results. This issue may require further analysis but I will leave it as is for now.

## Do normal days cluster apart from events?

The idea here is to include the 366 daily synoptic climatology values in with the synoptic MHW values to see if they cluster differently. From previous SOM and K-means cluster analyses on these daily climatology values it was determined that 2 - 4 nodes/ clusters was optimal. Anything more than that was over fitting. I attribute this to the much more even nature of the climatological data. Therefore I hypothesised that if one were to cluster the daily clims with the event data they would be placed into different clusters. There are 95 events and 366 daily clims, so this should allow for a pattern to be seen if one does exist. Up until this point SOMs have been the primary tool of investigation into clustering, but for this last step I opt to use HCA to produce a dendrogram of the results. I do this because it provides a better visual index of the difference between the points. It also allows for a nice visual representation of the two different classes of data being compared (i.e. event data and daily clim data). Additionally, an MDS is run to allow for another dimension of spatial difference to be inferred.

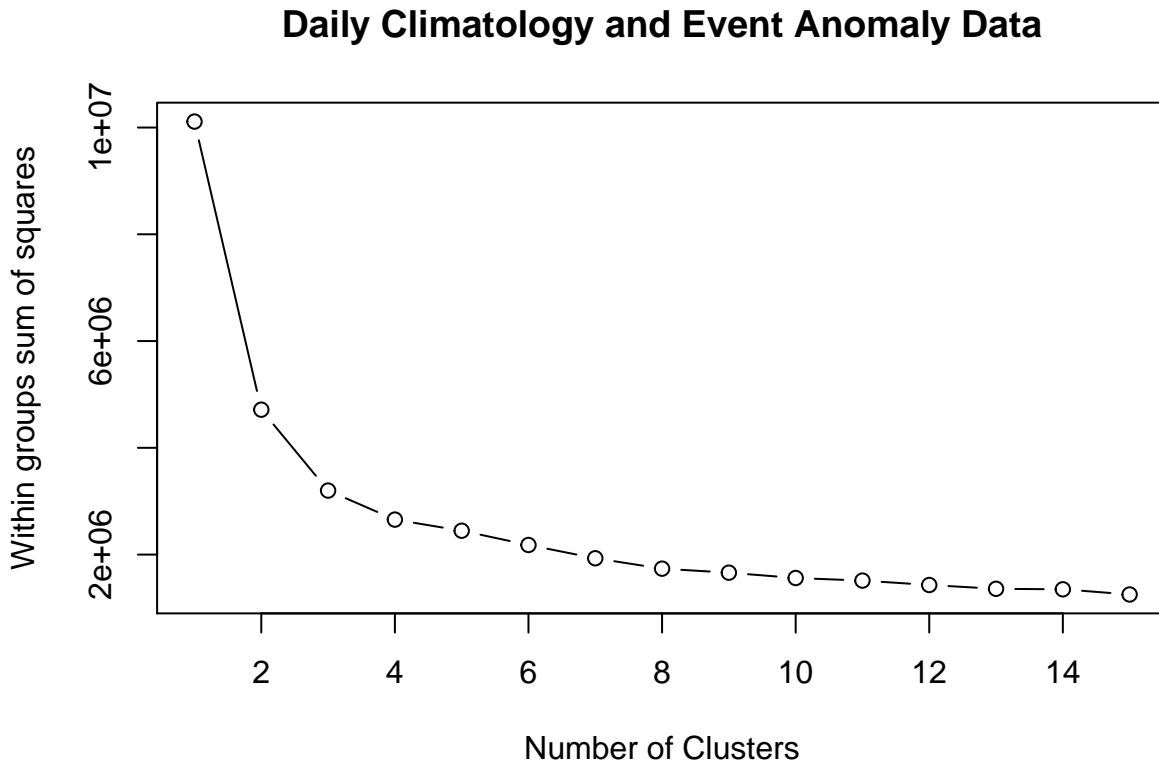


Figure 3: Plot showing the decrease in WGSS as more clusters are used for the results of an HCA on the anomaly values for synoptic air-sea states during events and daily climatologies.

Before running HCA and any other cluster or ordination techniques on these data we want to see how many clusters would be reasonable to use. The figure above shows that 4 to 6 is a good choice. With that known, we now run HCA, create a dendrogram and overlay some clusters.

The dendrogram from the HCA very clearly shows that the event data separate out from the daily clim data almost perfectly. Besides the overlayed clusters, one may also see that the final branch on which each event



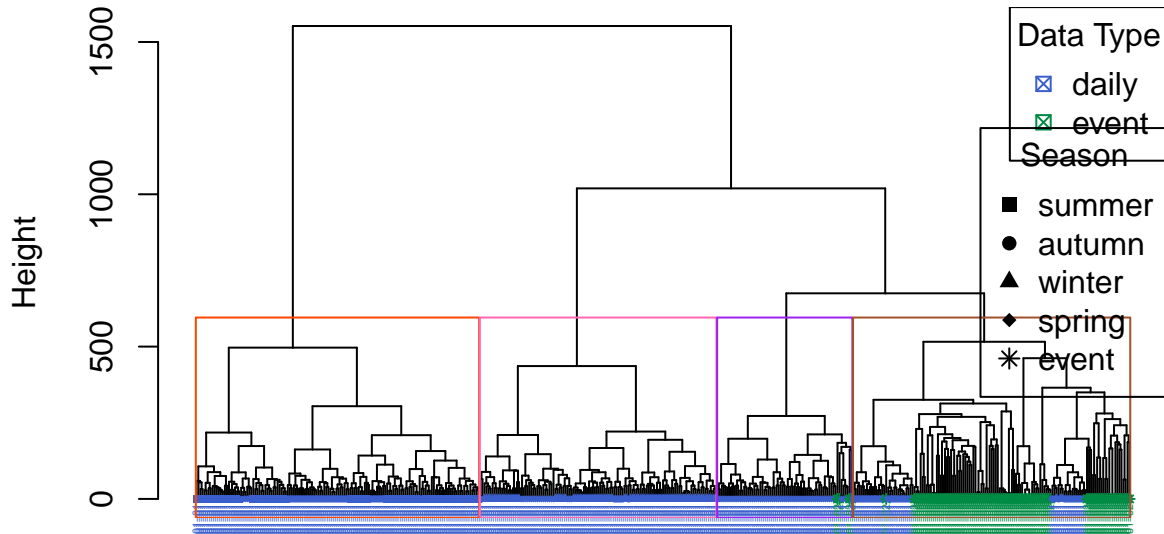


Figure 4: Dendrogram showing the results of an HCA on the anomaly values for synoptic air-sea states during events and daily climatologies. The daily clims and event data are shown with different colours. The dates or event names are included but are not legible.

sits is much longer than the daily clims. This means that not only are the events clustered apart from the daily clim data, but also that the individual events are also much more dissimilar from any of the other data points than the daily clims. But let's not stop there. The dendrogram makes a very convincing case for the dissimilarity between event and daily clim data, but let's add another dimension by creating an ordiplot via MDS.

I find these results very exciting. I think this ordiplot shows very clearly that the synoptic air sea states during the 366 daily climatologies are different from almost all of the synoptic air-sea states during coastal MHWs. As one may see from the flat ellipse of blue squares (the daily clim points), the variance represented in the x axis is seasonality. Indeed, if the dates are included in the figure above they are in a contiguous state. With January 1st in the top left edge of the ellipse of blue squares and the dates then move clockwise. So May is roughly in the middle of the top of the ellipse and October in the middle on the bottom. The synoptic states during events appear to be controlled by the variance represented by the y axis. This must be some sort of variance that is aseasonal. Likely the anomalous characteristics of air and or sea that occur during the events. This will require further investigation but I think it will prove to be a very strong result. Even if it isn't central to the question of what are the air-sea states during extreme events, it certainly helps to show that whatever those states may be, they are different from the common air-sea states.

This dot plots shows the seasonality of the clustering in a chronological order. The colours of the dots show if they are daily climatologies or event data. The x axis shows the date or event name, but there are too many to read. The take away message from this is to see how the clustering very clearly progresses throughout the year in a very even fashion. With cluster 1 representing summer, 2 shows Autumn, 3 winter, and 4 is the spring days. Beyond that, we see that almost all of the event data is clustered in with the Autumn data. This means that conditions during autumn most closely resemble the air-sea state during an extreme event. This could be taken to mean that ecosystems are naturally at more risk during this time of year. Or, perhaps, due to the consistency of the seasonality, that species would be more prepared for these conditions during this time of year and therefore less susceptible. One would need to do more research to say. But that isn't the focus of this work anyway. When this figure is taken in conjunction with the MDS plot above we are able to say that the event data is not most like the autumn daily clims, but rather they are the least dissimilar to these data. Because they are very different from the daily clims.

I understand that this may not look as clear to the reader as it does to me, so please let me know in what ways I am failing to communicate the patterns I see in these data so I can better delve deeper into them in

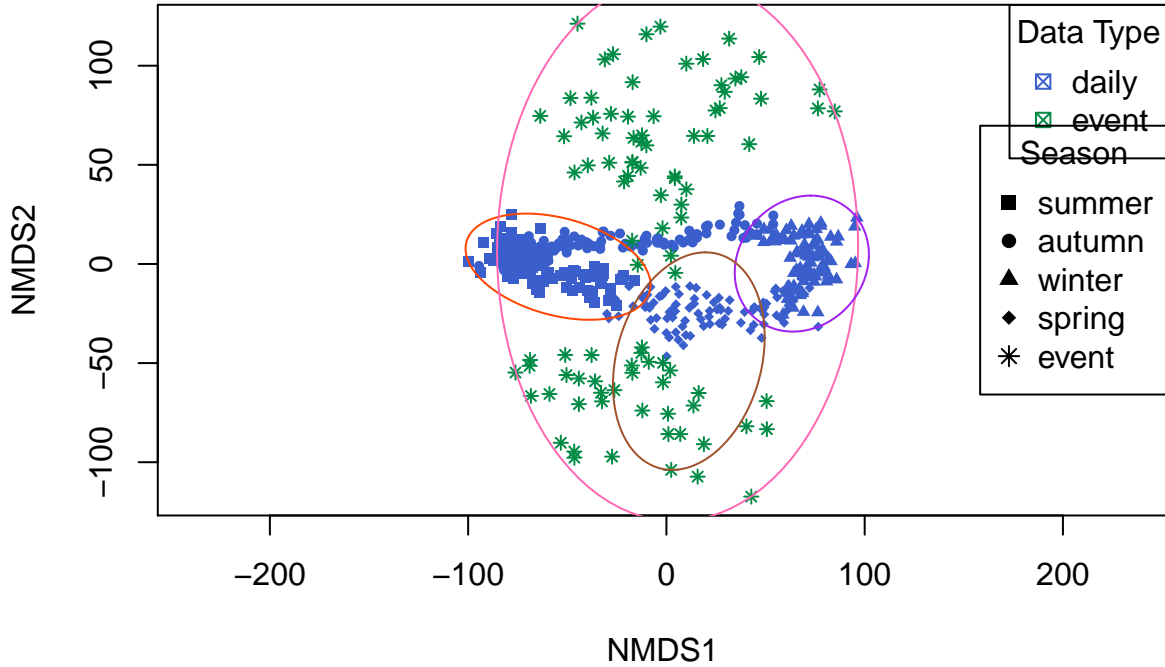


Figure 5: Ordiplo showing the results of an MDS on the anomaly values for synoptic air-sea states during events and daily climatologies. The daily clims and event data are shown with different colours. The dates or event names have not been included but are available. The clusters from the dendrogram are shown here but the colours do not correspond.

order to paint a clear picture for the publication.

## Clustering techniques that need to be investigated

### Hierarchical clustering

HCA differs from the other two techniques outlined below in that it does not cluster the data simultaneously, based on the least stress that can be found between data vectors. Rather it iteratively divides (or combines) data vectors as the algorithm moves down (or up) a classification tree. Always looking for the point at which clusters of data may be split (or combined). This method may benefit this research

### Dendrograms

### K-means clustering

The simplest method of clustering, and for that reason still one of the best. This is a basic algorithm that takes all data vectors and positions them in a 2D space. It then picks K points and sees, given the best possible fit of all dimensions being used, which data vectors are closest to which of the K points. This process is then repeated x number of times until a best fit is found. The data vectors are classified into the cluster centroid to which they are closest.

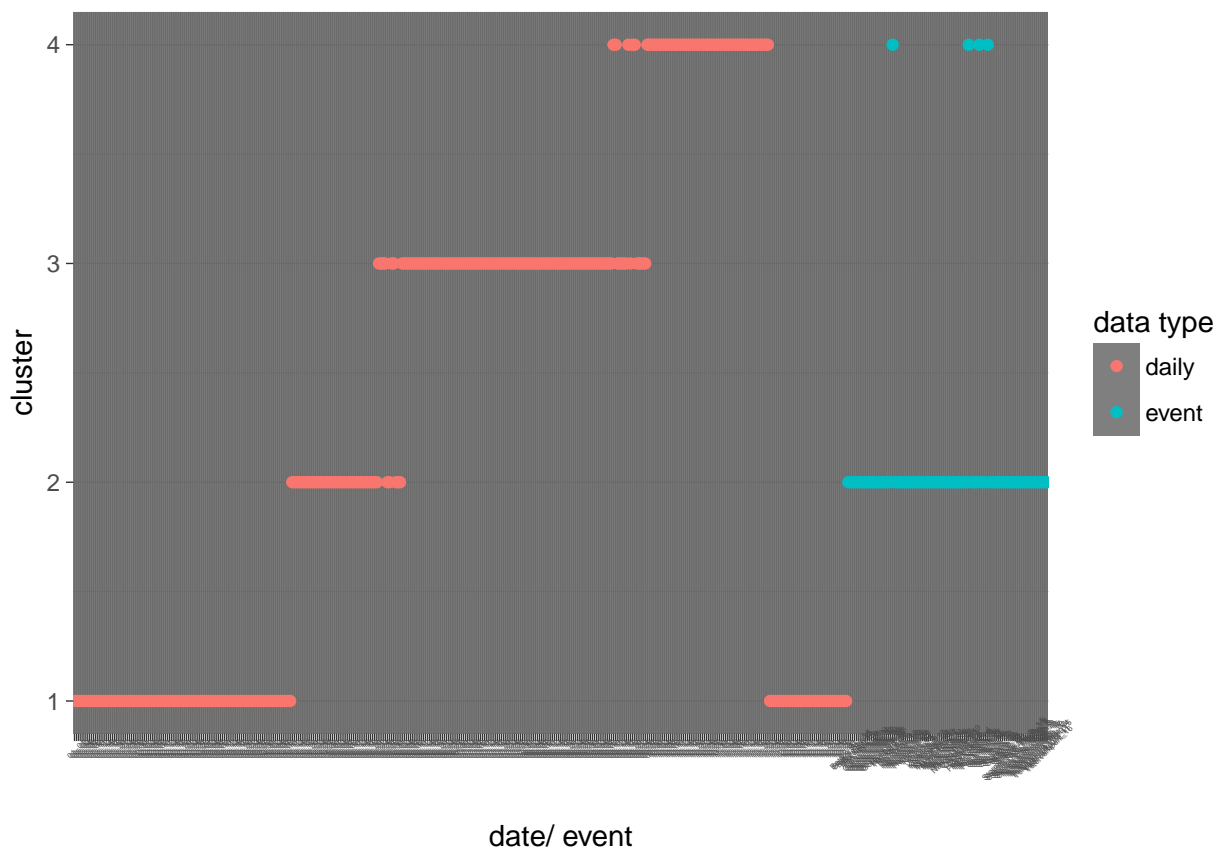


Figure 6: Dotplot showing the clustering of the anomaly values for synoptic air-sea states during events and daily climatologies as shown in the dendrogram and ordiplot. The daily climis and event data are shown with different colours. The dates or event names are shown on the x axis but are illegible.

## **Ordiplots**

## **SOMs**

The originally proposed technique and perhaps, once this dust settles, the reigning champion. The SOM technique is apart from the previous two methods in that it accounts for the gradient that exists between the nodes it clusters the given data into. Meaning that the positions of the nodes in 2D space is relevant, unlike HCA and K-means.

## **SOM nodes**

## **MDS**

Multi-dimensional scaling provides another possible layer of interpretation of these data. By highlighting which pixels on the map belong to which meso-scale properties (e.g. Agulhas, Benguela, Agulhas retroflexion) it is then possible to overlay the effect of these pixels, and therefore meso-scale features, on top of the ordiplots generated by MDS.

## **Ordiplots**

## **References**