

Cluster Analysis of
pavement deterioration
among climatic zones

Climatic zones of India

Objective: To divide the country into regions based on Climatic factors that affect pavement deterioration.

Features used for clustering:

- Average precipitation
 - no. of wet days
 - no. of high precipitation days
 - Average daily maximum temperature
 - no. of hot days
 - Yearly maximum
 - Average daily minimum temperature
 - no. of cold days
 - yearly minimum temperature
 - average relative humidity



Data Collection

Required Climatic data of the last 10 years was downloaded from the Indian Meteorological department, Pune website and NASA Data access portal.

The data for rainfall, minimum and maximum temperature was downloaded from [IMD Pune](#) website in the form of .grd files.

These files were converted to usable .csv files using imdlib python library through the [linked code](#) found on [Youtube](#)

The humidity data was downloaded from the [NASA Data access portal](#) in the form of 6 .csv files with each having data of area of $10^\circ \times 10^\circ$ lat*lon. The data points were then adjusted to be similar with all other climatic data That was collected before from IMD.

The raw data was then transformed and calculated to get final features for clustering.

A	B	C	D	E
X	Y	1	2	
37	73.5	8.5	32.87	32.52
40	76.5	8.5	31.04	30.94
41	77.5	8.5	32.04	31.97
42	78.5	8.5	31.11	30.93
71	76.5	9.5	31.27	31.18
72	77.5	9.5	30.67	30.67
73	78.5	9.5	29.44	29.3
74	79.5	9.5	30.61	30.16



X	Y	average	total	wetd
68.5	23.5	8.018193	2981.727	
69.5	21.5	7.921131	2944.821	
69.5	22.5	8.272256	3075.683	
69.5	23.5	8.361742	3109.377	
69.5	24.5	9.085993	3378.911	
69.5	27.5	10.88134	4047.351	
70.5	21.5	10.28052	3821.359	
70.5	22.5	8.719271	3241.752	

	X	Y	avghum
87	68.5	23.5	61.67446
69	69.5	21.5	67.25214
78	69.5	22.5	63.77786
88	69.5	23.5	51.87132
97	69.5	24.5	44.21612
122	69.5	27.5	33.35617
70	70.5	21.5	56.50129



Spatially Connected Clusters: Uniting Data Points with Agglomerative Clustering

The final collected data was shorted down to 309 spatial points

The spatial weights of these points along with the climatic factors were inputed to a hierachial clustering algorithm to form clusters

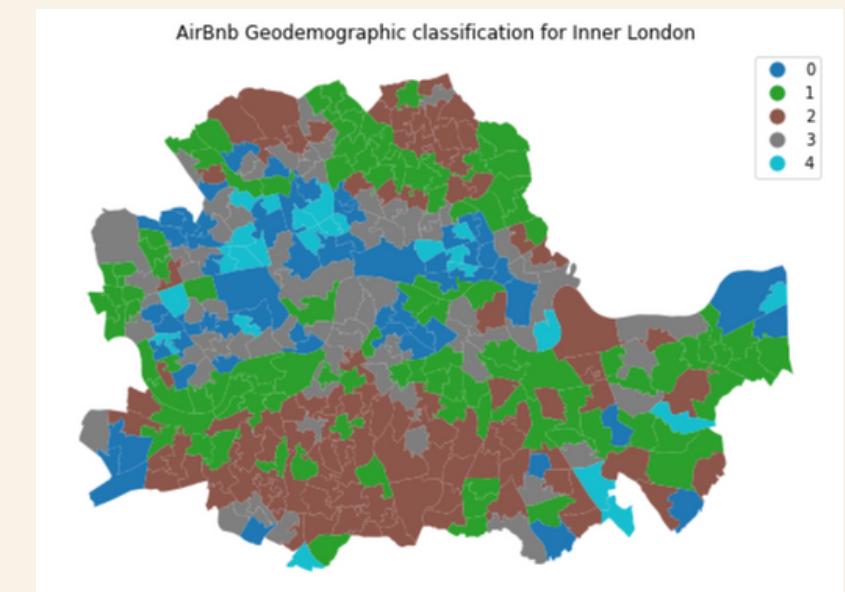
Regionalisation algorithm:

Regionalization is the subset of clustering techniques that impose a spatial constraint on the classification. In other words, the result of a regionalization algorithm contains areas that are spatially contiguous.

Agglomerative clustering with spatial connectivity constraint is one such regionalisation algorithm. It uses either spatial weights matrix or shape files to aggregate areas that are neighbours to each other.

Here, **in our data**, we do not have spatial weights matrices or shape files. Rather, we have climatic data at known points. Hence, we use the geopandas and pysal libraries to convert latitudes and longitudes to geodataframe and then spatial weights with queen criteria.

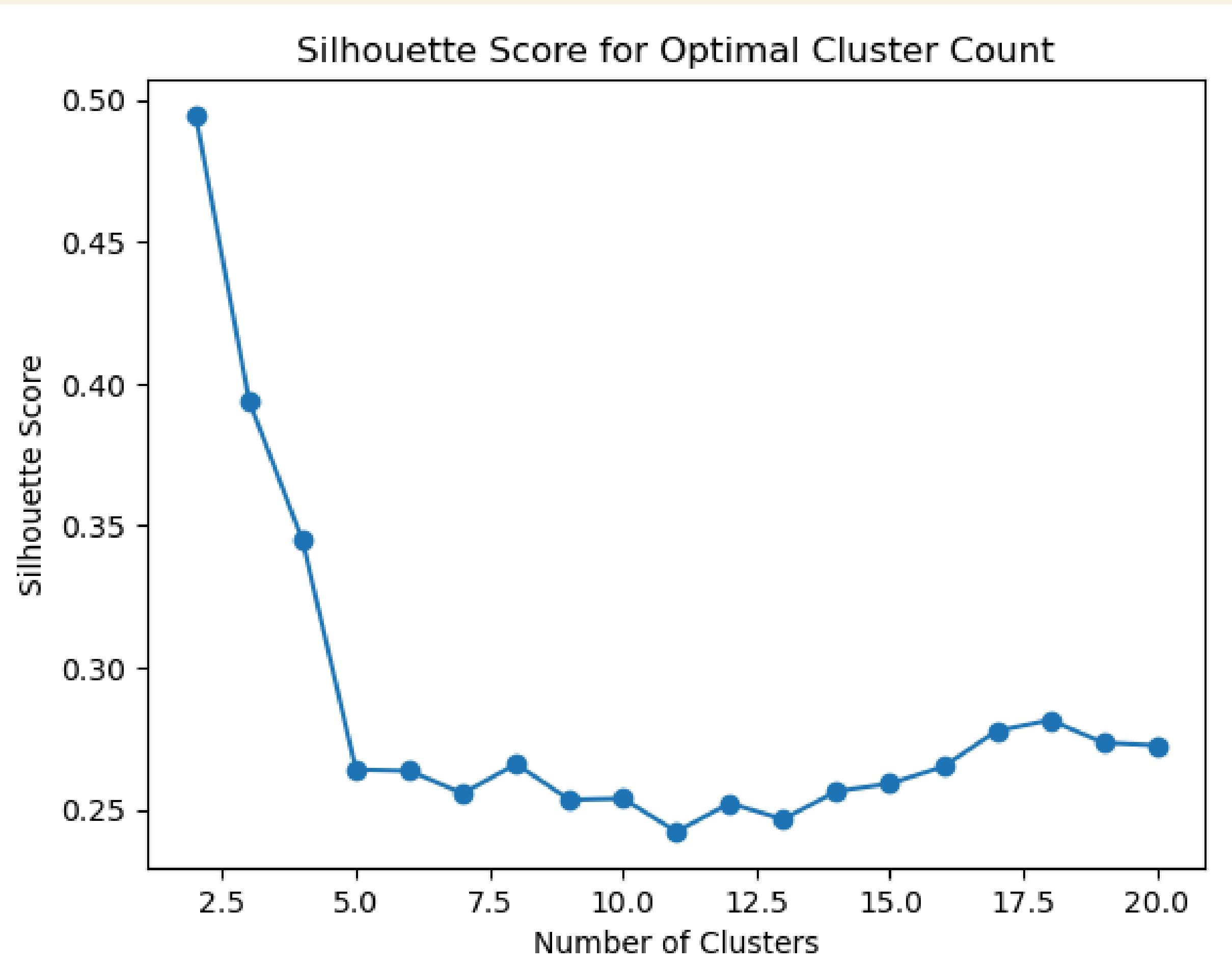
[Reference here.](#)



To find the optimal number of clusters, the elbow method using **silhouette scores** as evaluation metrics is used.

The silhouette scores of all the climatic variables was calculated and the elbow was found at **5 number of clusters**.

The silhouette value is a measure of how similar an object is to its own cluster (cohesion) compared to other clusters (separation).

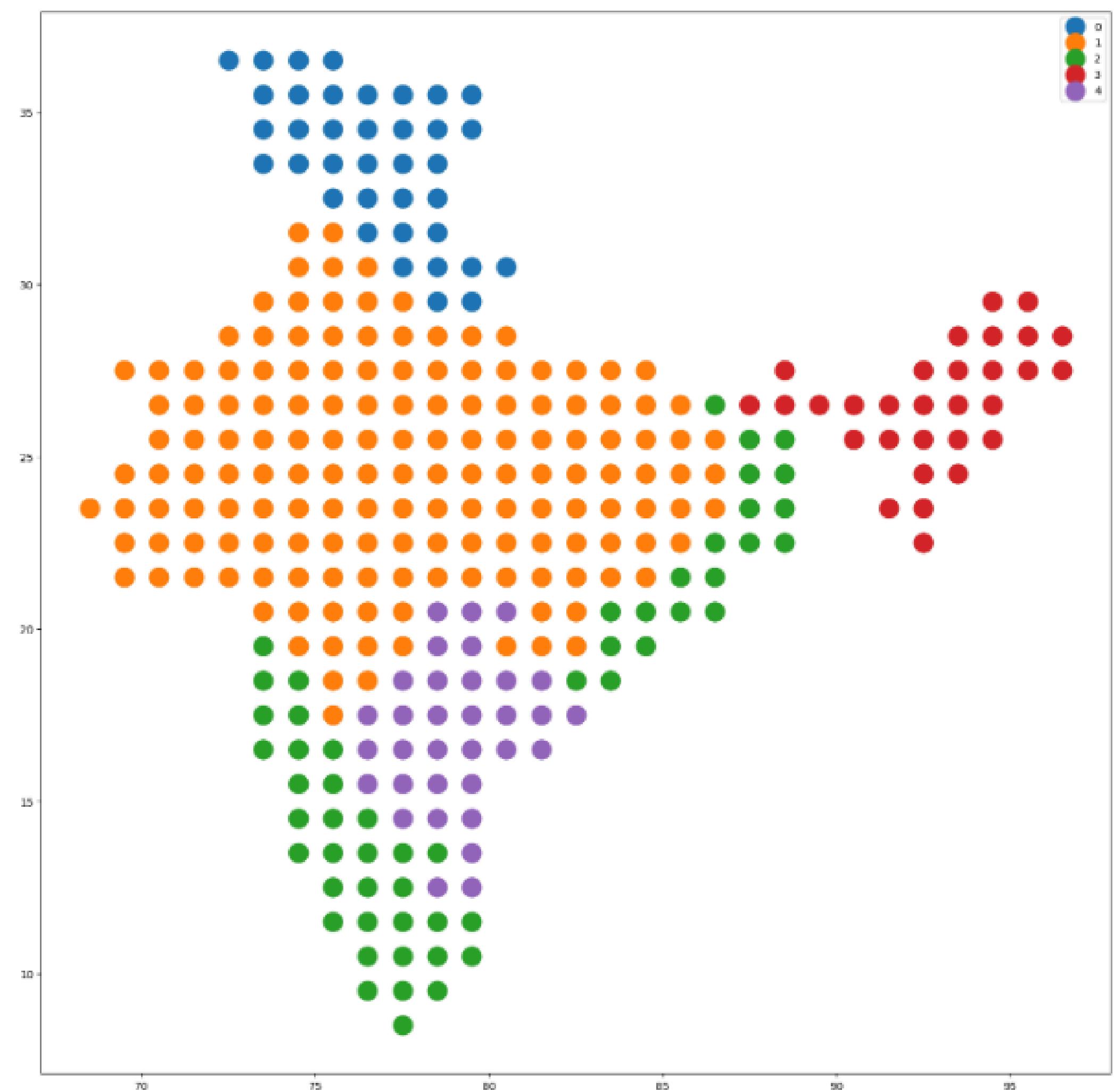


The following picture shows the cluster formation when rainfall, temperature and humidity are taken as features.

The Davies-Bouldin Index and the Silhouette score came out to be **1.1648** and **0.2661** for the above results.

This Davies Bouldin index represents the average 'similarity' of clusters, where similarity is a measure that relates cluster distance to cluster size. A model with a lower Davies-Bouldin index has a better separation between the clusters.

In our case, since the climatic data is continuous and irregular, the resulting Davies Bouldin Index is reasonable to be a little more than 1 and the silhouette score to be less than 0.5.

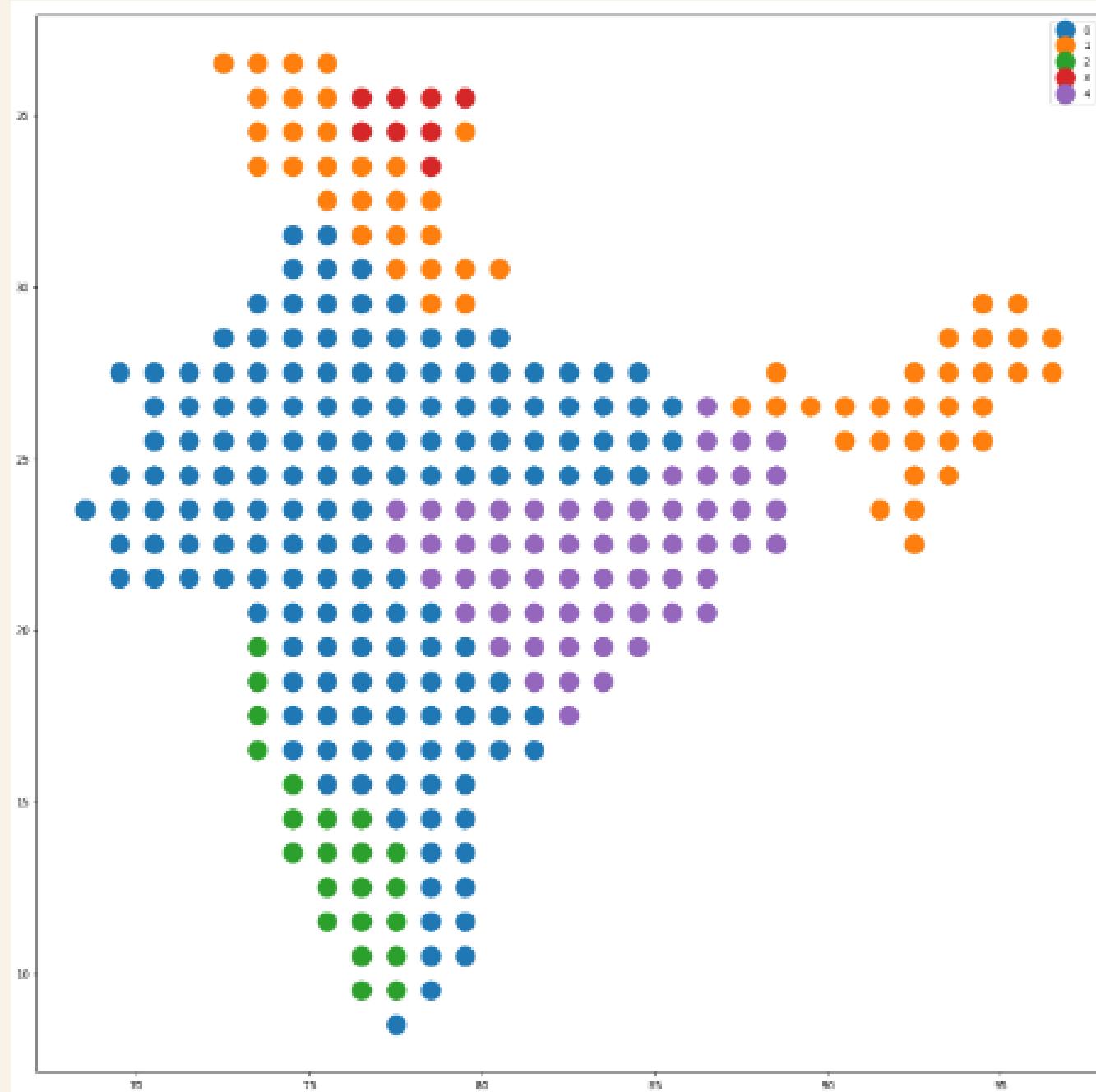


Preprocessing and Results

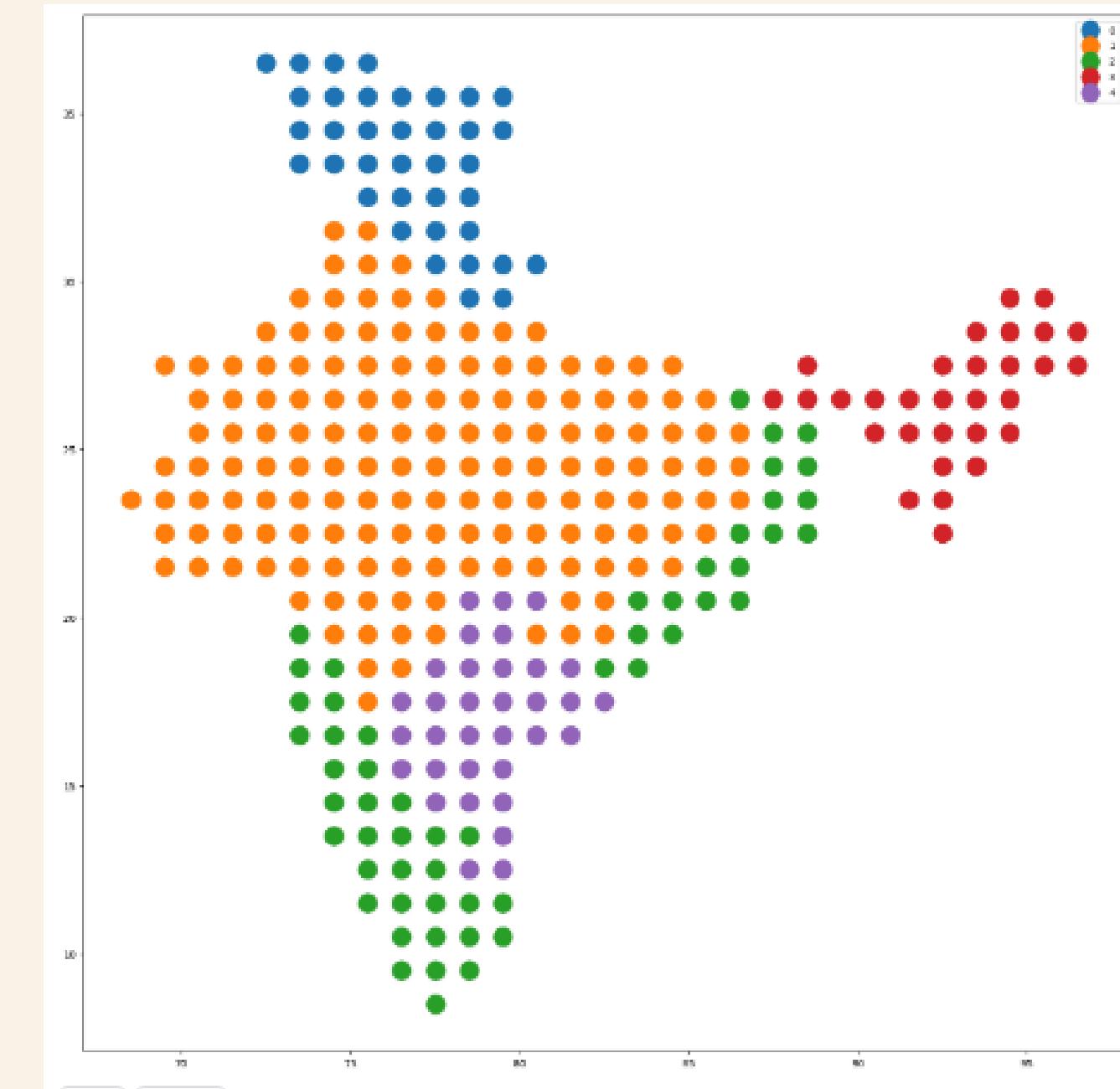
The features were standardized before being passed through clustering algorithm in order to ensure that there was no favour towards the features that naturally have larger magnitude like humidity. The difference in results with and without standardisation are shown in the next slide.

	X	Y	average	wetdays	high_precip	maxaverage	yearlymax	hotdays	minaverage	yearlymin	avghum		average	wetdays	high_precip	maxaverage	yearlymax	hotdays	minaverage	yearlymin	avghum	sagg13cls	
0	count	37.000000	37.000000	37.000000	37.000000	37.000000	37.000000	37.000000	37.000000	37.000000	37.000000		155.000000	155.000000	155.000000	155.000000	155.000000	155.000000	155.000000	155.000000	155.000000	155.0	
	mean	76.689189	33.581081	15.056938	112.162162	28.567568	22.672967	34.060847	19.297297	12.962146	-1.239456	62.546940		11.806554	68.709677	20.83871	32.510431	43.776995	212.303226	21.073187	6.274766	50.996887	1.0
	std	2.092816	2.073354	3.371626	49.347134	18.403714	2.364418	2.350089	24.921950	2.516910	2.581675	6.803490		4.091051	23.339384	8.18653	1.013583	1.723511	27.229275	1.037090	2.227239	7.789554	0.0
	min	72.500000	29.500000	9.974820	17.000000	3.000000	19.233430	30.714497	0.000000	8.990645	-4.422795	51.871158		5.139942	24.000000	4.00000	28.694117	38.417915	146.000000	18.074120	2.841949	33.356171	1.0
	25%	75.500000	32.500000	11.153087	93.000000	4.000000	20.813515	32.218369	2.000000	11.226514	-3.058757	57.331925		8.871994	52.500000	15.00000	31.917870	42.856134	198.000000	20.502214	4.655874	44.890640	1.0
	50%	76.500000	33.500000	16.090432	124.000000	35.000000	21.630200	33.220177	4.000000	11.871618	-2.317508	61.434274		10.311999	68.000000	20.00000	32.603501	44.330505	212.000000	21.117941	5.663892	51.349397	1.0
	75%	78.500000	35.500000	17.924352	151.000000	44.000000	24.863992	35.806114	33.000000	15.220402	1.148496	67.624952		14.990567	83.500000	28.00000	33.160722	44.906582	226.000000	21.511383	7.792997	56.406764	1.0
	max	80.500000	36.500000	19.287686	167.000000	54.000000	27.002304	38.610123	82.000000	17.715110	3.603627	76.215904		21.085709	125.000000	45.00000	34.666935	46.454216	300.000000	23.591546	12.159250	68.526247	1.0
1	average	wetdays	high_precip	maxaverage	yearlymax	hotdays	minaverage	yearlymin	avghum														
2	average	wetdays	high_precip	maxaverage	yearlymax	hotdays	minaverage	yearlymin	avghum														
3	average	wetdays	high_precip	maxaverage	yearlymax	hotdays	minaverage	yearlymin	avghum														
4	average	wetdays	high_precip	maxaverage	yearlymax	hotdays	minaverage	yearlymin	avghum														

0-blue
1-orange
2-green
3=red
4-violet



without standardising



with standardising

The after results seem to be favourable when compared with the IRI clusters as their distribution was also similar in the mid-region.

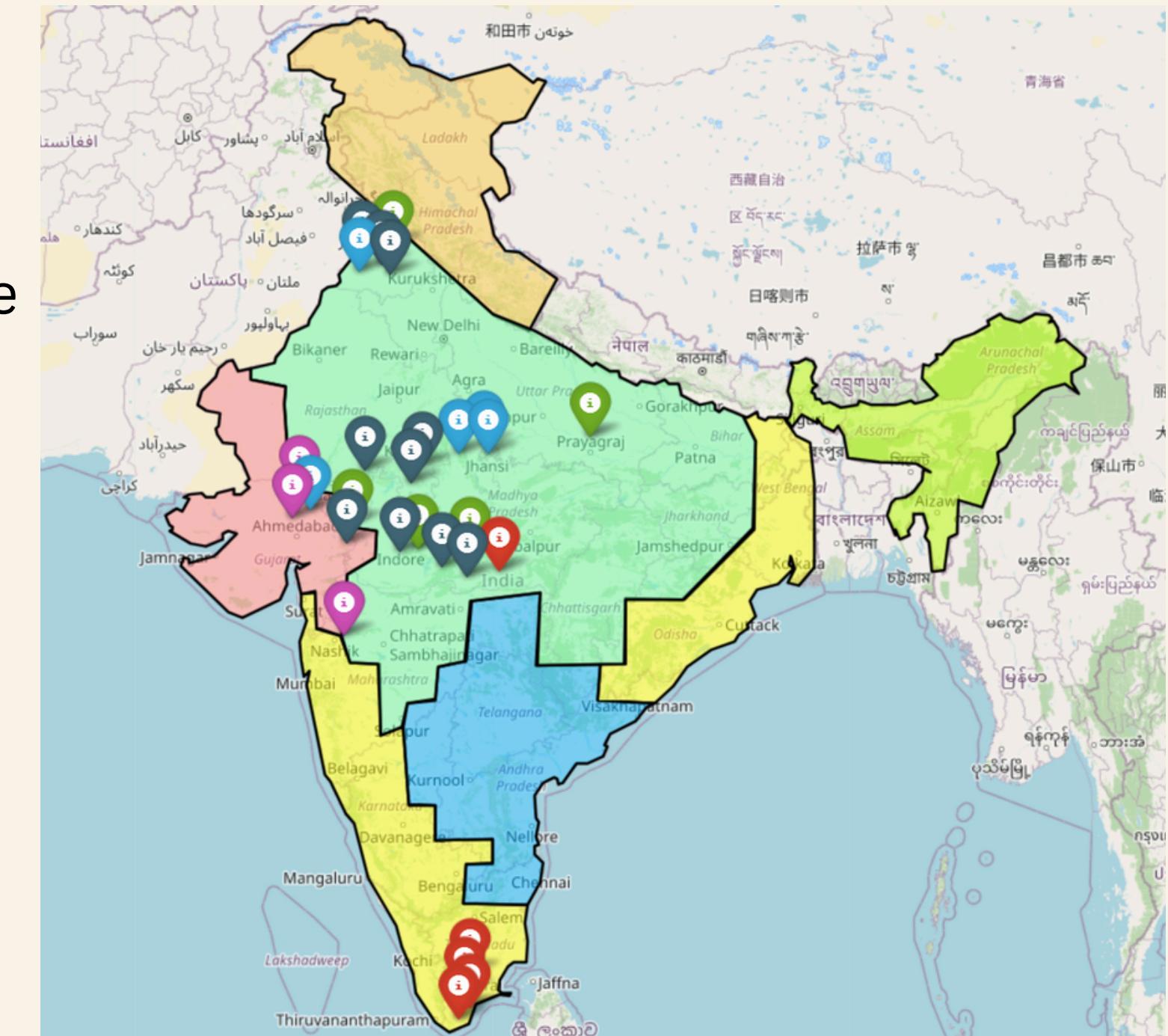
Validation of IRI difference with pavements climatic zone

The next part of the project deals with analysing how IRI is behaving with respect to the climatic zone that pavement falls in.

To find out the patterns of pavement deterioration with respect to the climate, we first cluster all the IRI data we have into different ranges. The process is elaborated in the upcoming slides.

Steps followed:

- Cleaning and transforming raw IRI data
- Assigning the IRI data to categorical variables
- Performing statistical analysis to find out patterns
- Clustering the categories together according similar IRI-difference range
- Comparing the locations of similar categories with obtained climatic zones



Data Preprocessing

The data was provided by MORTH, India and NHAI. The data consisted of section-wise IRI measurements for a set of national highways of India over multiple dates.

The IRI measures cannot be directly compared to find out patterns in pavement deterioration. Hence the difference of the IRI measurements of two dates for the same section of road is considered to be a factor. the change observed over an year of time in different regions is compared for our analysis.

C	NH Number(New)	District	Start Chainage(km)	End Chainage(km)	Lane Number	IRI-15	IRI-16	IRI-17	16-15diff	16-17diff
0	NH0752	Jhalawar	89.517	89.6	L1	2.41	5.19	5.98	2.78	0.79
1	NH0752	Jhalawar	89.600	89.7	L1	2.41	4.29	5.87	1.88	1.58
2	NH0752	Jhalawar	89.700	89.8	L1	2.48	4.46	4.56	1.98	0.10
3	NH0752	Jhalawar	89.900	90.0	L1	3.24	3.50	3.78	0.26	0.28
4	NH0752	Jhalawar	90.200	90.3	L1	3.25	3.31	3.41	0.06	0.10
...
282	NH0138	Tuticorin	18.900	18.8	R2	1.52	2.62	2.79	1.10	0.17
283	NH0138	Tuticorin	22.100	22.0	R2	2.00	2.52	3.01	0.52	0.49
284	NH0138	Tuticorin	38.100	38.0	R2	2.59	2.68	2.82	0.09	0.14
285	NH0138	Tuticorin	39.700	39.6	R2	2.26	2.38	2.99	0.12	0.61
286	NH0138	Tuticorin	32.800	32.9	L2	2.30	2.51	3.40	0.21	0.89

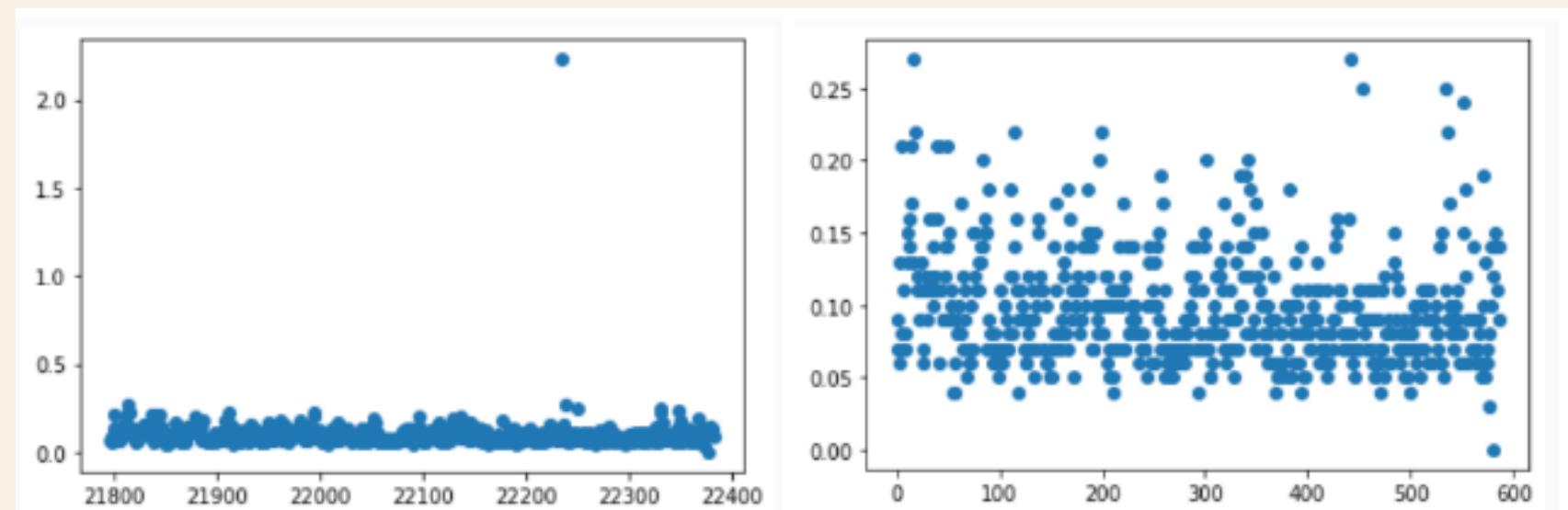
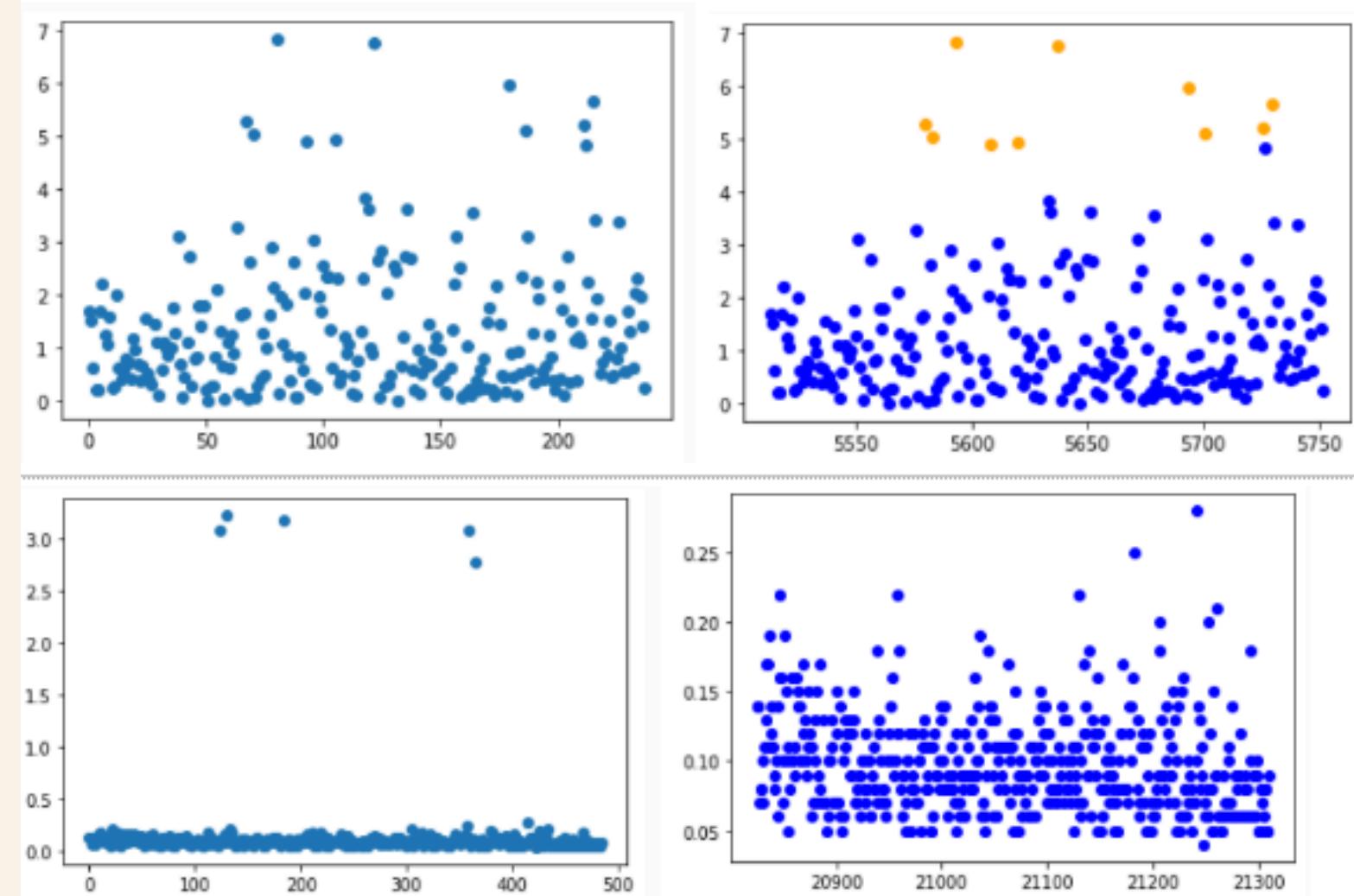
Data cleaning

All the measurements of IRI difference greater than 3 are ignored as they are practically rare and might be wrong measurements.

Outliers are detected and removed using DBSCAN clustering algorithm

The figure shows the type of points that are being considered outliers. (orange ones are outliers and fig2 shows the data after removal).

The range of IRI-difference where the section doesn't have much points are being considered outliers. Such range is found out for every section and removed.



Statistical understanding of the data

In order to cluster the data in a useful way, we have to consider the geographical locations too.

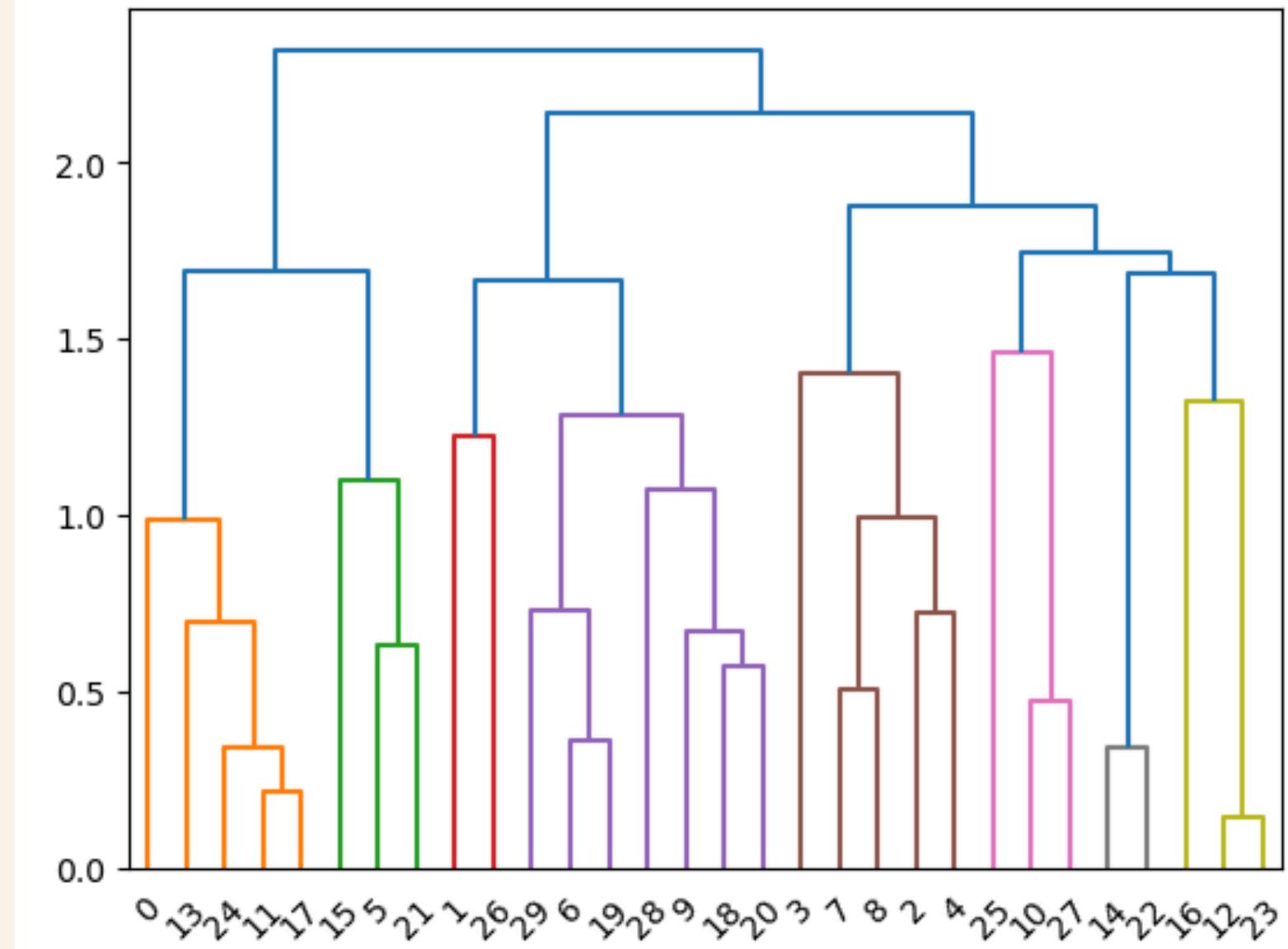
To choose from the above two, similarity is being checked in the two cases. A **null hypothesis** is proposed that there are no significant differences between the means of the groups. The null hypothesis is rejected when the p-value is less than 0.05

The Kruskal test(alternative to ANOVA test) was conducted for two cases

- single Highway, multiple districts. data of each $NH \cap Di$ is taken as a group and there are j such groups ($0 < i < j$) {p-value average =}
- One district, multiple highways, data of each $D \cap NH_i$ is taken as a group and there are j such groups ($0 < i < j$)

The results showed that the data within a district are more similar than in the first case. So, district wise clustering is chosen

Clustering districts with similar ΔIRI ranges



Cluster 1: [0, 5, 11, 13, 15, 17, 21, 24]

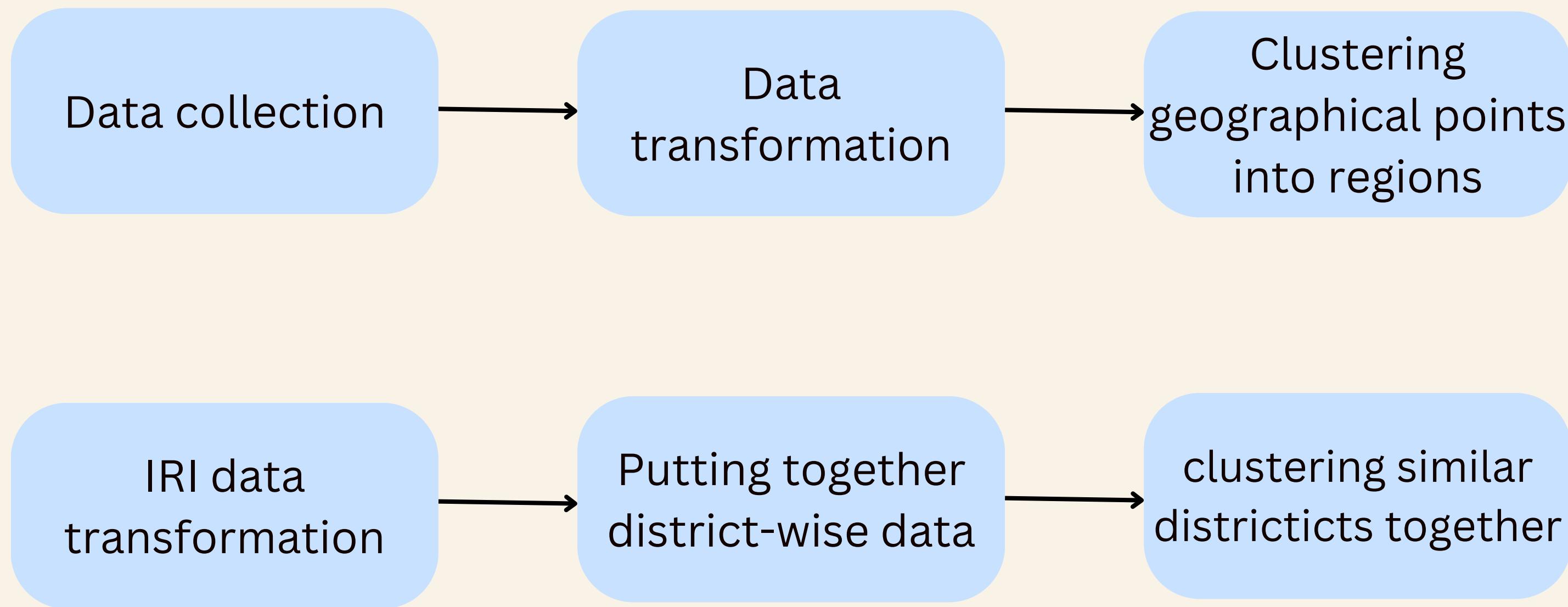
Cluster 2: [1, 6, 9, 18, 19, 20, 26, 28, 29]

Cluster 3: [2, 3, 4, 7, 8]

Cluster 4: [10, 25, 27]

Cluster 5: [12, 14, 16, 22, 23]

- Kruskal test is performed on all the districts' data and found out that they are all significantly different
- A pairwise t-test which is Conover test is performed for all the datasets and the p-value matrix is obtained
- The p-value matrix is being passed into the hierachial clustering algorithm with multiple types of linkage methods.
- The average of p-values of each cluster obtained from different methods is used as evaluation metrics.
- It is found out that the highest p-value is obtained in case of average linkage method.
- The datasets in each cluster are visualised and finally choosen for comparing with climatic zones.

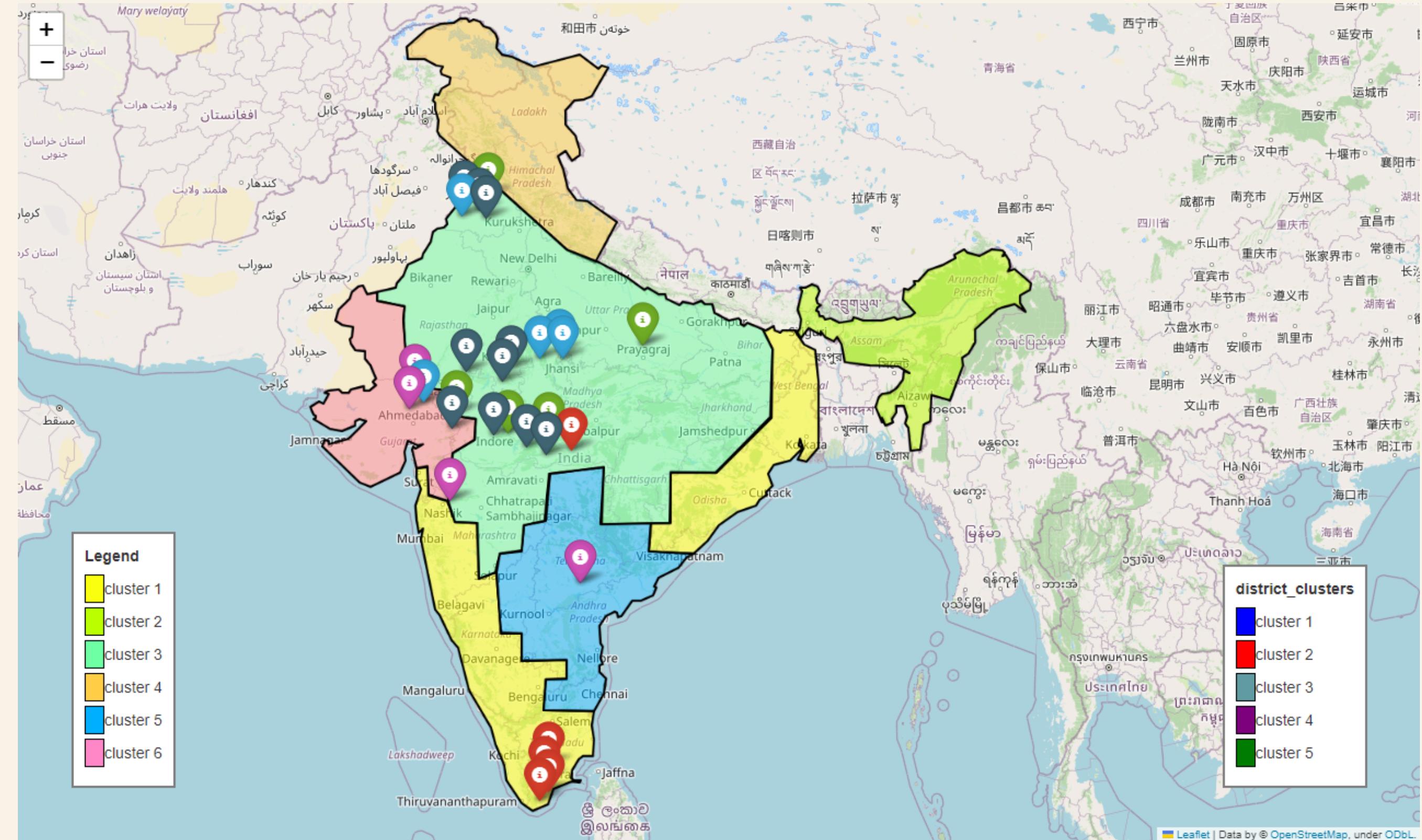


Plotting the clusters and climatic zones together

The districts are plotted with the colour corresponding to the cluster they belong on the map that contains climatic zones.

A pattern can be seen that the districts that belong in cluster 1,3 and 5 are lying in the north and eastern parts of the country and others in southern parts.

It is to be noted that both the clusters 2 and 4 show very less deterioration compared to other clusters.





Project Proposal

- Borcelle Studio -

