

Perform Network Analysis

- Create a covariance matrix
- Find an ideal noise-removing threshold filter with Random Matrix Theory
- Analyze properties of the network

Questions to answer with alpha and beta diversity analysis:

- Does the 3G and 3M look like 1C or 2C
- How do 1G and 1M compare to 1C
- How does 2G compare to 3G, 2M compare to 3M
- How does 1C compare to 2C

```
In [1]: import pandas as pd
df_counts = pd.read_csv("data/FCF_relative_counts.csv", index_col=0)
df_annotations = pd.read_csv("data/FCF_annotations_corrected.csv", index_col=0)
```

Import Data

```
In [2]: df_counts_rel_1C = df_counts_rel.loc[df_annotations["group"] == "1C", :]
df_counts_rel_1M = df_counts_rel.loc[df_annotations["group"] == "1M", :]
df_counts_rel_2C = df_counts_rel.loc[df_annotations["group"] == "2C", :]
df_counts_rel_2M = df_counts_rel.loc[df_annotations["group"] == "2M", :]
df_counts_rel_3C = df_counts_rel.loc[df_annotations["group"] == "3C", :]
df_counts_rel_3M = df_counts_rel.loc[df_annotations["group"] == "3M", :]
```

```
In [3]: df_counts_rel
```

	opitutus spp.	paludibacter propionigenes	magnetospirillum sp.	rhodopseudomonas palustris	acetobacter spp.	bacteroides spp.	pleomorphomonas oryzae	alpha sp. acidophilus	rhodoblast
1C	0.207688	0.186672	0.111388	0.105284	0.074637	0.072618	0.061373	0.051297	0.0443
1M	0.260982	0.066341	0.188617	0.127545	0.012517	0.026242	0.050610	0.038828	0.0161
2C	0.290486	0.026069	0.162435	0.117691	0.011864	0.026510	0.041062	0.032812	0.0115
3C	0.194119	0.023516	0.234578	0.097890	0.008008	0.018172	0.031574	0.027593	0.0244
3M	0.187871	0.048505	0.132181	0.130452	0.011161	0.030600	0.049378	0.037943	0.0250
3G7A	0.104459	0.001801	0.018378	0.065455	0.067390	0.293988	0.159024	0.050529	0.0240
3G7B	0.247657	0.000743	0.022841	0.051781	0.062985	0.275727	0.206768	0.004794	0.0467
3G10A	0.221610	0.000888	0.028917	0.059782	0.080850	0.262612	0.185294	0.004707	0.0299
3G10B	0.093492	0.001327	0.090380	0.008828	0.046889	0.383666	0.146778	0.008510	0.0512
3G10C	0.102037	0.000338	0.240036	0.053780	0.061816	0.286654	0.121156	0.005679	0.0169

86 rows × 10 columns

Split Samples Into Groups

```
In [4]: groups = set(df_annotations["group"])
groups
```

```
Out[4]: {'1C', '1G', '1M', '2C', '2G', '2M', '3C', '3M', '3G'}
```

```
In [5]: df_counts_rel_1C = df_counts_rel.loc[df_annotations["group"] == "1C", :]
df_counts_rel_1M = df_counts_rel.loc[df_annotations["group"] == "1M", :]
df_counts_rel_2C = df_counts_rel.loc[df_annotations["group"] == "2C", :]
df_counts_rel_2M = df_counts_rel.loc[df_annotations["group"] == "2M", :]
df_counts_rel_3C = df_counts_rel.loc[df_annotations["group"] == "3C", :]
df_counts_rel_3M = df_counts_rel.loc[df_annotations["group"] == "3M", :]
```

```
In [6]: df_counts_rel_3M
```

	opitutus spp.	paludibacter propionigenes	magnetospirillum sp.	rhodopseudomonas palustris	acetobacter spp.	bacteroides spp.	pleomorphomonas oryzae	alpha sp. acidophilus	rhodoblast
3MA	0.174045	0.000404	0.030346	0.021195	0.000757	0.013949	0.061011	0.005201	0.0204
3MB	0.038556	0.000985	0.007862	0.021515	0.001198	0.005499	0.099063	0.007420	0.1167
3MC	0.104566	0.000922	0.500182	0.014735	0.001243	0.003717	0.044183	0.004214	0.1209
3MB	0.074561	0.003128	0.024658	0.024319	0.012068	0.090322	0.199926	0.007409	0.1384
3MC	0.149788	0.001396	0.051745	0.027828	0.001339	0.087199	0.102974	0.006518	0.4014
3MSA	0.157334	0.000190	0.061791	0.020911	0.001099	0.034584	0.130031	0.005369	0.3593
3MTA	0.203516	0.001522	0.034878	0.041955	0.007556	0.07834	0.123492	0.010883	0.2616
3MTB	0.094730	0.000746	0.024598	0.038245	0.029740	0.163735	0.192327	0.006677	0.0984
3MTC	0.195999	0.002399	0.069706	0.023421	0.023294	0.124625	0.112692	0.005527	0.1062
3M10A	0.121134	0.016020	0.113322	0.049653	0.009298	0.093914	0.127554	0.011773	0.1497
3M10B	0.125448	0.007026	0.265430	0.030463	0.030561	0.184910	0.138465	0.009670	0.0715
3M10C	0.149036	0.002733	0.201424	0.037578	0.030911	0.144593	0.084645	0.006991	0.1933

12 rows × 10 columns

Standardize Data For Correlation Analysis

Standardize relative abundance counts by OTU by subtracting the mean across all samples and dividing by variance across all samples. Use this standardized data matrix for all subsequent correlation analysis.

Source: Deng Paper

```
In [7]: df_counts_rel_1C_stan = df_counts_rel_1C.stn(df_counts_rel_1C.mean(), axis=1)
df_counts_rel_1M_stan = df_counts_rel_1M.stn(df_counts_rel_1M.mean(), axis=1)
df_counts_rel_2C_stan = df_counts_rel_2C.stn(df_counts_rel_2C.mean(), axis=1)
df_counts_rel_2M_stan = df_counts_rel_2M.stn(df_counts_rel_2M.mean(), axis=1)
df_counts_rel_3C_stan = df_counts_rel_3C.stn(df_counts_rel_3C.mean(), axis=1)
df_counts_rel_3M_stan = df_counts_rel_3M.stn(df_counts_rel_3M.mean(), axis=1)
```

```
Out[7]: df_counts_rel_1C_stan
```

	opitutus spp.	paludibacter propionigenes	magnetospirillum sp.	rhodopseudomonas palustris	acetobacter spp.	bacteroides spp.	pleomorphomonas oryzae	alpha sp. acidophilus	rhodoblast
3MA	0.737828	-0.756187	-0.716109	-0.765586	-0.727594	-1.158250	-1.118185	-0.974558	2.448
3MB	0.185156	-0.468953	-1.175262	-0.738538	-0.710066	-1.295065	-0.305736	-0.089996	-0.586
3MC	1.689498	-0.640769	2.194347	-1.334413	-0.708690	-1.32907	-1.468437	-1.368070	-0.389
3MB	0.657591	-0.148796	-0.773035	-0.492098	-0.286830	0.078213	1.831747	-1.034047	-0.500
3MC	0.079791	-0.402432	-0.971835	-0.183721	-0.704912	0.027657	-0.858348	-0.456020	1.094
3MSA	0.284504	-0.015221	-0.545445	-0.791596	-0.714287	-0.824016	0.315616	-0.019690	0.834
3MTA	1.537370	-0.506780	-0.710407	1.606705	-0.376939	-0.107779	0.212116	0.472312	0.230
3MTB	1.437370	0.030074	-0.773999	0.731798	0.401933	1.266768	1.670731	1.590231	-0.777
3MTC	1.333359	-0.313347	-0.496935	-0.571108	0.150613	0.633580	-0.016741	-0.967293	-0.729
3M10A	-0.697500	2.725996	-0.112773	1.734422	2.761679	0.136387	0.298118	1.623975	-0.400
3M10B	-0.580466	0.720491	0.702663	0.047891	0.433822	0.160482	0.008699	0.791451	-0.94
3M10C	0.059406	-0.268676	0.310365	0.673204	0.482123	0.956849	-0.613050	0.404234	-0.1

12 rows × 10 columns

```
In [9]: df_counts_rel_1C_stan.mean()
```

```
Out[9]: opitutus spp. 2.869580e-16
paludibacter propionigenes 1.663335e-16
magnetospirillum sp. -1.110223e-17
rhodopseudomonas palustris 0.609006e-09
acetobacter spp. -1.653135e-16
```

```
preotella nanensis
lactobacillus helveticus
bacillus thermoautotrophicus
corynebacterium durum
length: 141, dtype: float64
```

```
Out[10]: opitutus spp. 1.0
paludibacter propionigenes 1.0
magnetospirillum sp. 1.0
rhodopseudomonas palustris 1.0
acetobacter spp. 1.0
```

```
preotella nanensis
lactobacillus helveticus
bacillus thermoautotrophicus
corynebacterium durum
length: 141, dtype: float64
```

Network Construction

- create a pair-wise similarity (Pearson) of abundance across different samples
- determine adjacency matrix by RMT-based approach

```
In [11]: import networkx
networkx.set_random_state(1)
```

Use Random Matrix Theory To Threshold

Good overview of Random Matrix Theory by Torsten Scholok

TODO: figure out if my Poisson distribution is very often... my chi_squared value seems unreal...

```
In [12]: # nat.find_RMT_threshold(df_counts_rel_1C_stan, s_t=0.3, alpha=0.5)
```

Data Is Too Messy

Since the data is too messy, this code lets you eyeball the threshold where the data fits the Poisson distribution more than it fits the GOE distribution.

```
In [13]: # threshold = 0.49
# nat.visualize_RMT_threshold(df_counts_rel_1C_stan, threshold)
# nat.visualize_RMT_threshold(df_counts_rel_1G_stan, threshold)
# nat.visualize_RMT_threshold(df_counts_rel_1M_stan, threshold)
# nat.visualize_RMT_threshold(df_counts_rel_2C_stan, threshold)
# nat.visualize_RMT_threshold(df_counts_rel_2M_stan, threshold)
# nat.visualize_RMT_threshold(df_counts_rel_3C_stan, threshold)
# nat.visualize_RMT_threshold(df_counts_rel_3M_stan, threshold)
```

```
In [14]: threshold = 0.75
nat.visualize_RMT_threshold(df_counts_rel_1C_stan, threshold)
nat.visualize_RMT_threshold(df_counts_rel_1G_stan, threshold)
nat.visualize_RMT_threshold(df_counts_rel_1M_stan, threshold)
nat.visualize_RMT_threshold(df_counts_rel_2C_stan, threshold)
nat.visualize_RMT_threshold(df_counts_rel_2M_stan, threshold)
nat.visualize_RMT_threshold(df_counts_rel_3C_stan, threshold)
nat.visualize_RMT_threshold(df_counts_rel_3M_stan, threshold)
```

```
shape (111, 111)
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 1872621675.958992
Xsq_crit 16.25816442187412 at alpha=0.05
shape (76, 69)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 6825519284.134026
Xsq_crit 16.25816442187412 at alpha=0.05
shape (122, 121)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 264891709228.84634
Xsq_crit 16.25816442187412 at alpha=0.05
shape (122, 121)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 387623755141.487
Xsq_crit 147.03752976381804 at alpha=0.05
shape (100, 100)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 39783623 9874957
Xsq_crit 123.225224533618 at alpha=0.05
shape (184, 184)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 307530753.383557
Xsq_crit 127.68838953385 at alpha=0.05
shape (110, 110)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 387623755141.487
Xsq_crit 147.03752976381804 at alpha=0.05
shape (100, 100)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 39783623 9874957
Xsq_crit 123.225224533618 at alpha=0.05
shape (184, 184)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 387623755141.487
Xsq_crit 147.03752976381804 at alpha=0.05
shape (100, 100)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 39783623 9874957
Xsq_crit 123.225224533618 at alpha=0.05
shape (184, 184)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

Xsq_poisson 387623755141.487
Xsq_crit 147.03752976381804 at alpha=0.05
shape (100, 100)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel Density Estimate (blue line), and Empirical Spacing Distribution (grey bars).

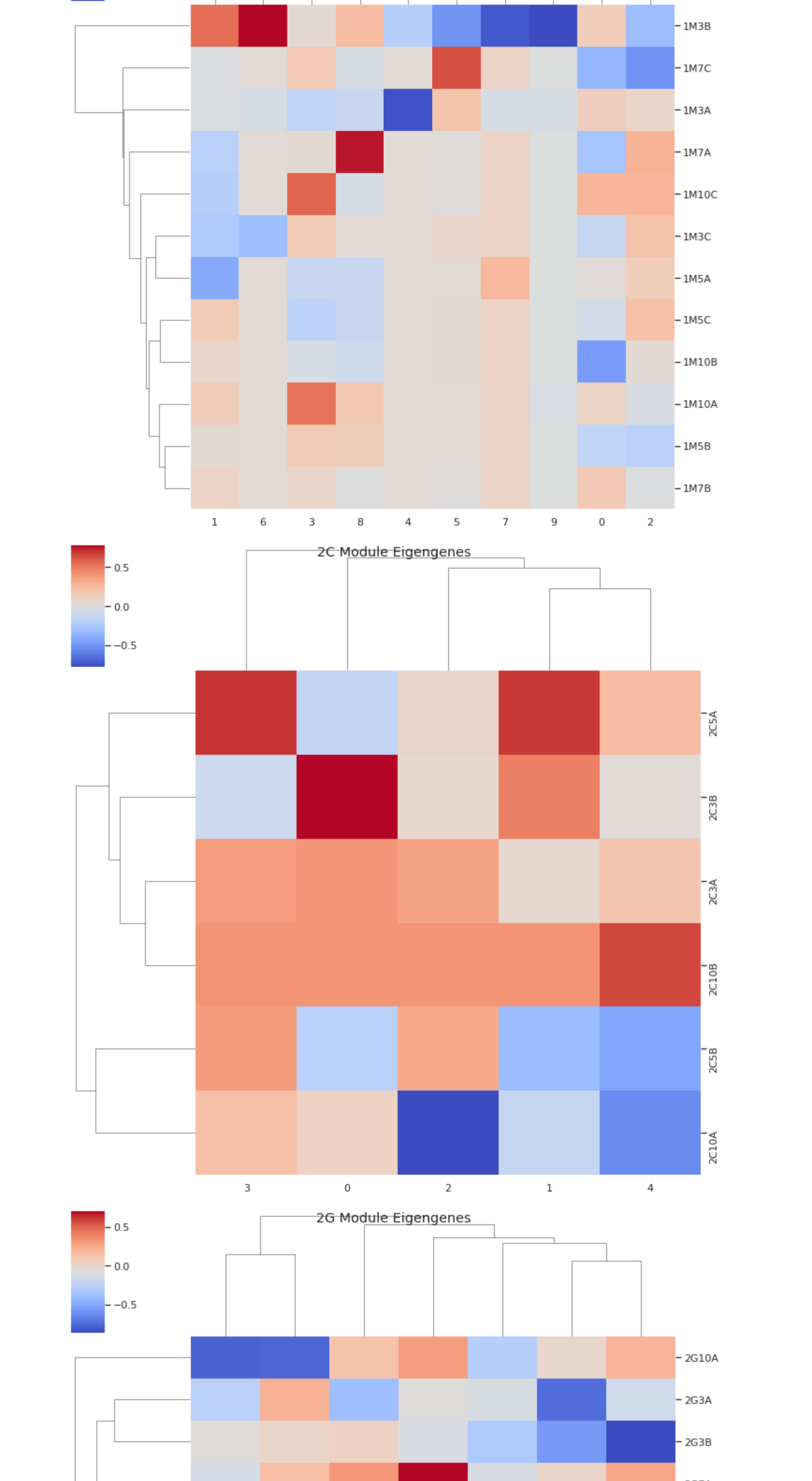
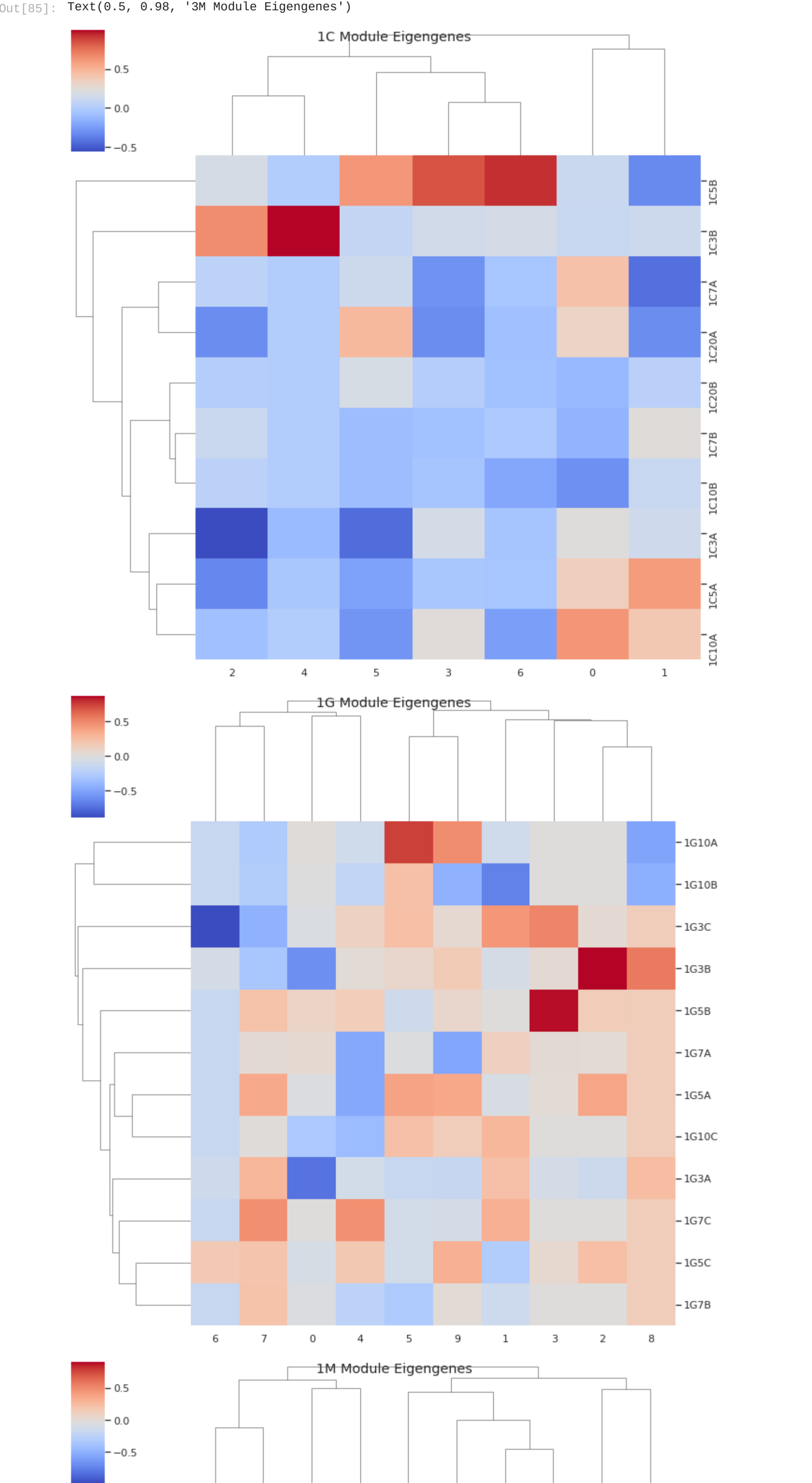
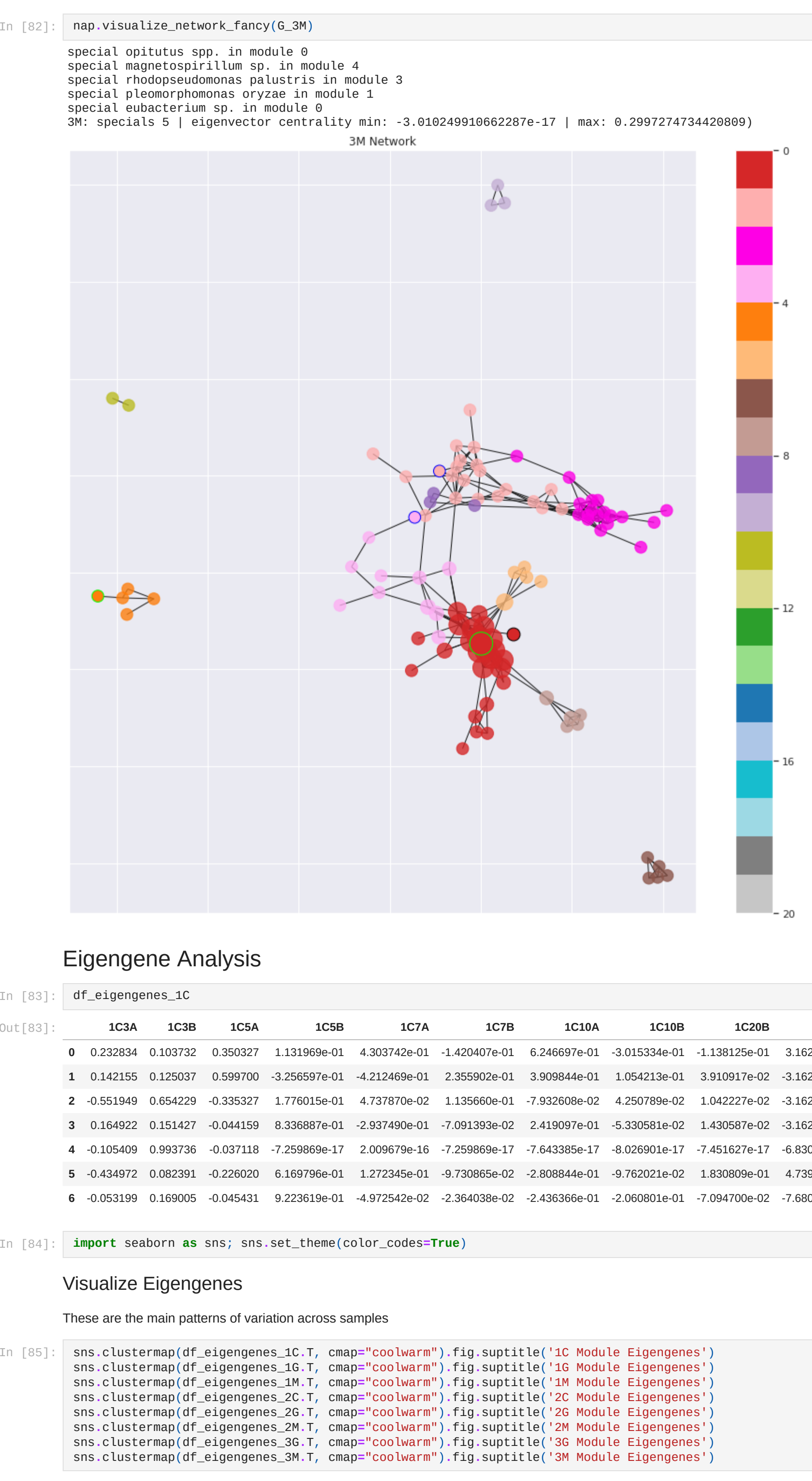
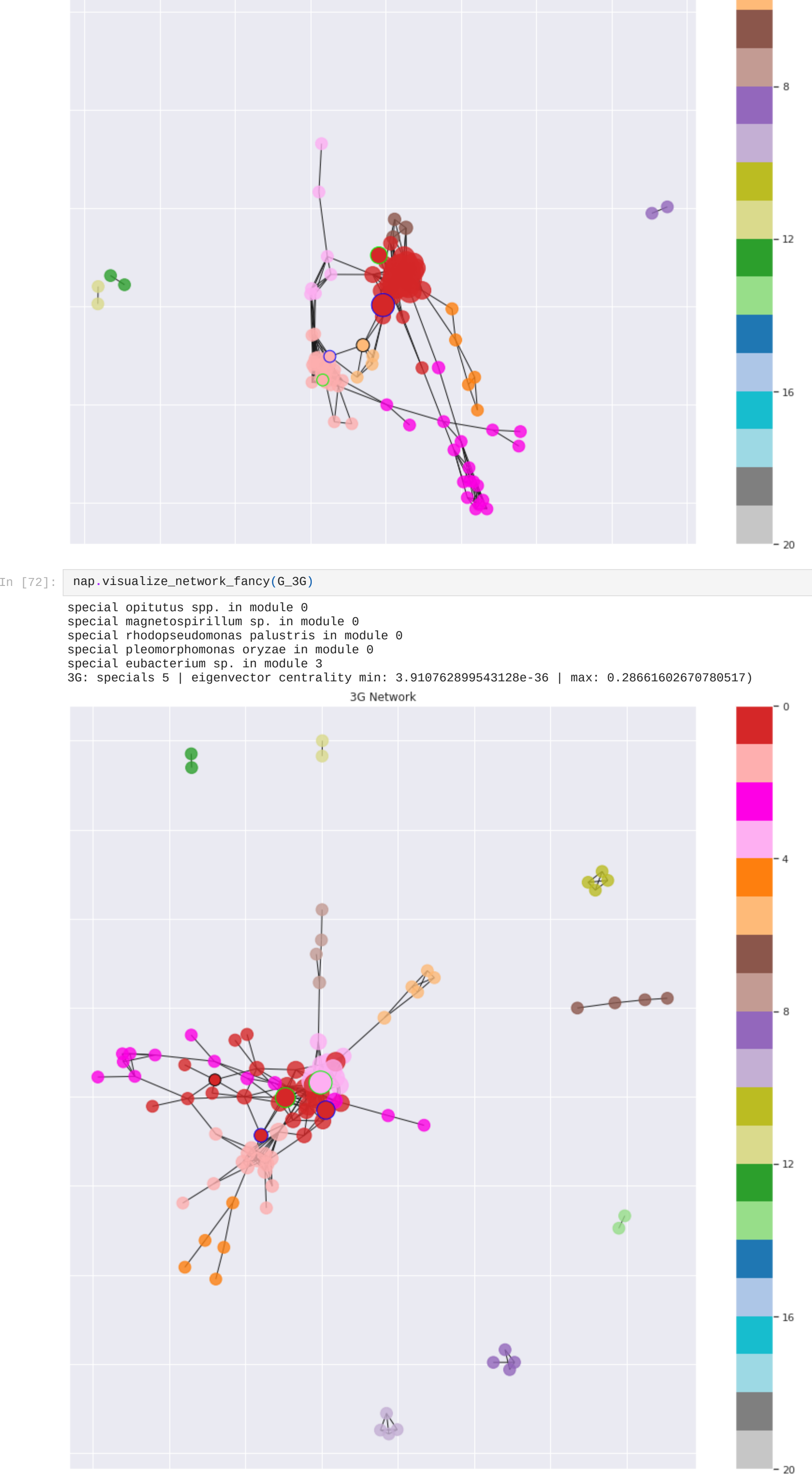
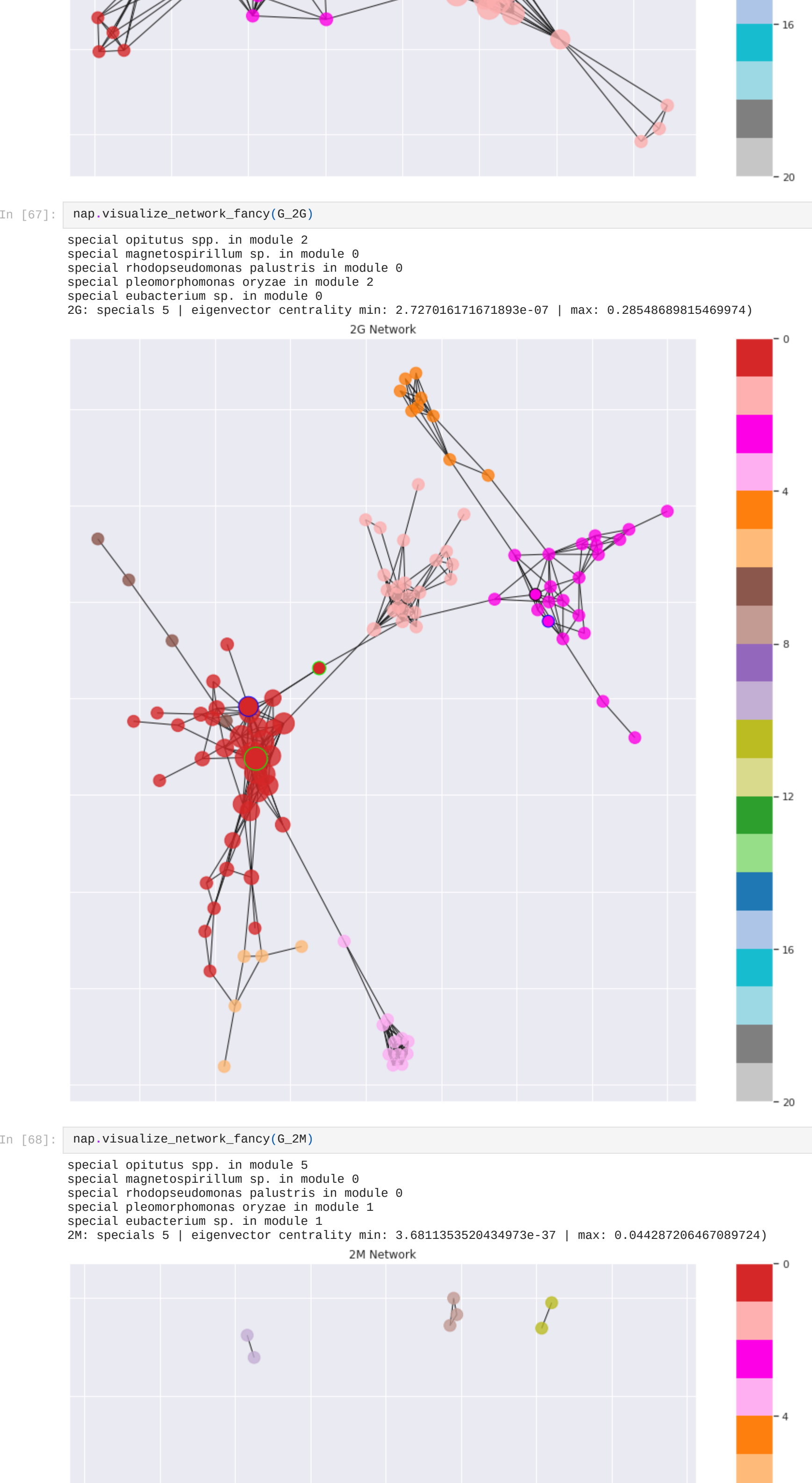
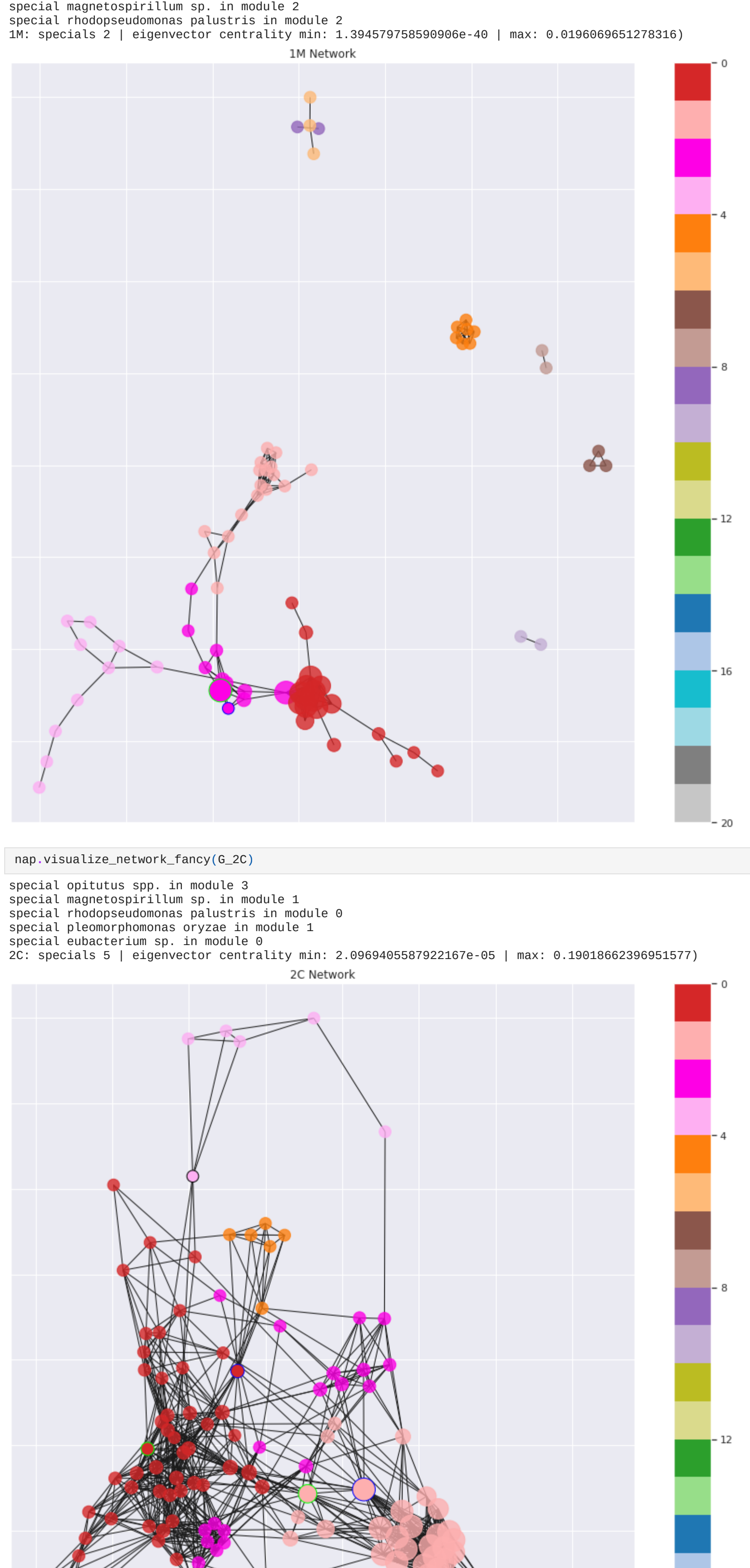
Xsq_poisson 39783623 9874957
Xsq_crit 123.225224533618 at alpha=0.05
shape (184, 184)

```
/home/gaher/miniconda3/envs/MBENV/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'di'
plot() is a deprecated function and will be removed in a future version. Please adapt your code to use either
histplot (a figure-level function with similar flexibility) or histplot (an axes-level function for histograms)
```

```
warnings.warn(msg, FutureWarning)
```

Unfolded Spacing Distribution

Density plot showing the distribution of spacing (x-axis, 0.0 to 2.5) versus density (y-axis, 0.0 to 1.4). The plot compares the Poisson distribution (green line), Gaussian Orthogonal distribution (orange line), Kernel



Eigengene Analysis

In [83]:

	df_eigengenes_1C	1C3A	1C3B	1C5A	1C5B	1C7A	1C7B	1C10A	1C10B	1C20B	1C20A
Out[83]:	0	0.232034	0.103732	0.350327	1.131909e-01	4.303742e-01	-1.420407e-01	6.246697e-01	-0.015334e-01	-1.120125e-01	3.162278e-01
	1	0.142155	0.125037	0.599700	-3.256597e-01	-4.212499e-01	2.355902e-01	3.909844e-01	1.054213e-01	3.9310917e-02	-3.162278e-01
	2	-0.551949	0.654229	-0.335327	1.776015e-01	4.737870e-02	1.135660e-01	-7.932600e-02	4.250789e-02	1.042227e-02	-3.162278e-01
	3	0.164922	0.151427	-0.044159	8.336887e-01	2.937490e-01	-7.091393e-02	2.413097e-01	-5.330581e-02	1.430587e-02	-3.162278e-01
	4	-0.105409	0.993736	-0.037118	-7.259869e-17	2.009679e-16	-7.259869e-17	-7.643385e-17	-8.020001e-17	-7.451627e-17	-6.830800e-17
	5	-0.434972	0.082391	-0.226200	6.169796e-01	1.272345e-01	-9.730865e-02	-2.808844e-01	-9.762021e-02	1.830809e-01	4.739839e-01
	6	-0.055199	0.169005	-0.045431	9.223619e-01	-4.972542e-02	-2.364038e-02	-2.436366e-01	-2.060801e-01	-7.094700e-02	-7.680426e-02

In [84]:
import seaborn as sns; sns.set_theme(color_codes=True)

Visualize Eigengenes

These are the main patterns of variation across samples

