

Cluster Analysis: Identifying Parkinson's Disease Subtypes

Jesse Mu

Wednesday, June 10

1 Preprocessing

1.1 Dataset Description

951 subjects, 145 metrics, collected 15-4-2012. From Pablo Martinez Martín. 170 subjects with missing values (brought down to 781); these were removed automatically, even if the missing values were not included in the selected features below. This will need to be changed later on, by keeping those removed that still have all selected features and perhaps with some compensation for missing values.

1.2 Selected Features

Combination of non-motor scale (NMS) symptoms and standard motor symptoms.

Name	Type	Format	Description
nms_d1	byte	%8.0g	cardiovascular
nms_d2	byte	%8.0g	sleep/fatigue
nms_d3	byte	%8.0g	mood/cognition
nms_d4	byte	%8.0g	percep/hallucinations
nms_d5	byte	%8.0g	attention/memory
nms_d6	byte	%8.0g	gastrointestinal
nms_d7	byte	%8.0g	urinary
nms_d8	byte	%8.0g	sexual function
nms_d9	byte	%8.0g	miscellaneous
tremor	float	%9.0g	tremor
bradykin	float	%9.0g	bradykinesia ¹
rigidity	float	%9.0g	rigidity
axial	float	%9.0g	axial ²
pigd	float	%9.0g	postural instability and gait difficulty

Table 1: Selected Features and Details

Name	μ	σ	min-max
nms_d1	1.76	3.32	0-24
nms_d2	8.71	8.76	0-48
nms_d3	8.70	11.83	0-60
nms_d4	1.65	3.94	0-33
nms_d5	5.22	7.44	0-36
nms_d6	5.67	6.92	0-36
nms_d7	8.02	9.09	0-36
nms_d8	3.57	5.97	0-24
nms_d9	6.99	7.74	0-48
tremor	2.59	2.63	0-12
bradykin	2.49	1.39	0-6
rigidity	2.34	1.36	0-6
axial	3.28	2.75	0-12
pigd	3.36	2.77	0-12

Table 2: Descriptive Statistics

1.3 Dimensionality Reduction: PCA

May not be useful? If we're trying to identify *clinically* relevant features, merging them may not be a good idea.

Figure 1 shows scree test elbow occurs around 2 or 3. Also, eigenvalues 1 and 2 $>$ 1, while 3 is around .9

2 k -means

2.1 Identifying optimal number of clusters

2.1.1 WSS Error Scree Test

Figure 2 shows no optimal elbow in scree test! Maybe 2-3?

2.1.2 Gap Statistic

Optimal cluster is the local maximum of the gap statistic, but it appears to be consistently increasing in Figure 3.

2.1.3 Average Silhouette Width

Figure 4 shows average silhouette width as being consistently under 0.25 for all clusters, implying the data is not well structured.

¹Impaired ability to adjust the body's position.

²Issues affecting the middle of the body.

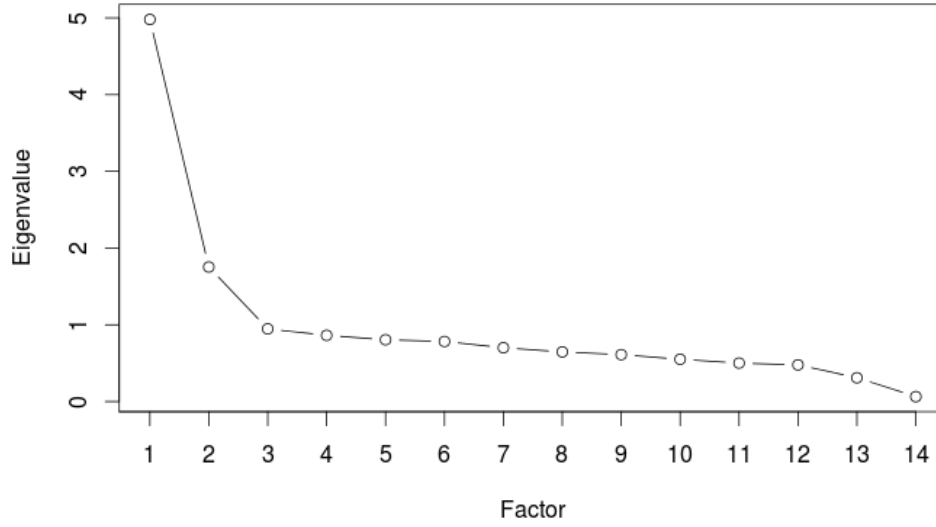


Figure 1: Scree test: eigenvalues by factor

2.2 Cluster statistics

k	n	Within SS	sum(Within SS)
2	201/580	4248.585/4132.434	8381.019
3	420/231/130	2618.368/1973.82/3076.542	7668.73
4	61/372/145/203	1481.25/1845.389/2147.988/1609.555	7084.183

Table 3: Cluster statistics

2.3 Centers

```

CLUSTERS: 2
=====
> 1   nms_d1   nms_d2   nms_d3   nms_d4   nms_d5   nms_d6
      nms_d7   nms_d8   nms_d9   tremor   bradykin
0.7328282 0.9345720 0.9810287 0.8195052 0.7908599 0.8551764
0.8493069 0.6376458 0.6289463 0.2008254 0.8019272
rigidity   axial     pigd
0.6840038 1.0834910 1.0634082
> 2   nms_d1   nms_d2   nms_d3   nms_d4   nms_d5
      nms_d6   nms_d7   nms_d8   nms_d9   tremor
-0.2539629 -0.3238776 -0.3399772 -0.2840009 -0.2740739
-0.2963629 -0.2943288 -0.2209772 -0.2179624 -0.0695964

```

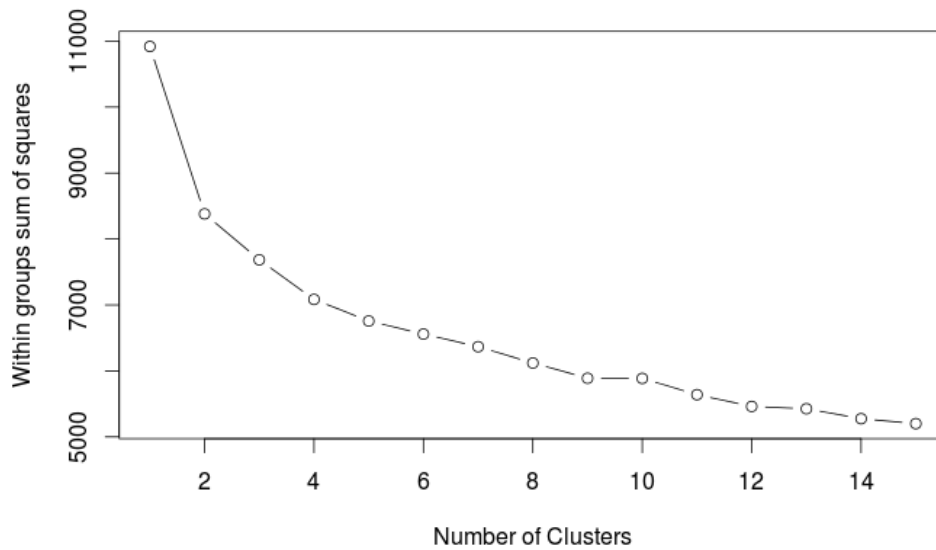


Figure 2: Scree test: WSS error by cluster size

```

    bradykin    rigidity      axial      pigd
-0.2779093 -0.2370427 -0.3754857 -0.3685259
CLUSTERS: 3
=====
> 1    nms_d1      nms_d2      nms_d3      nms_d4      nms_d5
      nms_d6      nms_d7      nms_d8      nms_d9      tremor
-0.2699345 -0.3571672 -0.3574942 -0.2776501 -0.2579928
      -0.3084614 -0.3030016 -0.2270260 -0.1867338 -0.2531402
    bradykin    rigidity      axial      pigd
-0.6091393 -0.5542033 -0.5792769 -0.5775312
> 2    nms_d1      nms_d2      nms_d3      nms_d4      nms_d5
      nms_d6      nms_d7      nms_d8      nms_d9
-0.13369057 -0.07938936 -0.05832302 -0.22742096 -0.18981647
      -0.02540265 -0.15964421 -0.08973885 -0.16993231
      tremor    bradykin    rigidity      axial      pigd
0.39256552 0.69210577 0.63956741 0.43974428 0.44762184
> 3    nms_d1      nms_d2      nms_d3      nms_d4      nms_d5      nms_d6
      nms_d7      nms_d8      nms_d9      tremor    bradykin
1.1096539 1.2949935 1.2586166 1.3011330 1.1708044 1.0417061
      1.2626038 0.8929277 0.9052504 0.1202787 0.7381697
    rigidity      axial      pigd
0.6540408 1.0901181 1.0704805
CLUSTERS: 4
=====

```

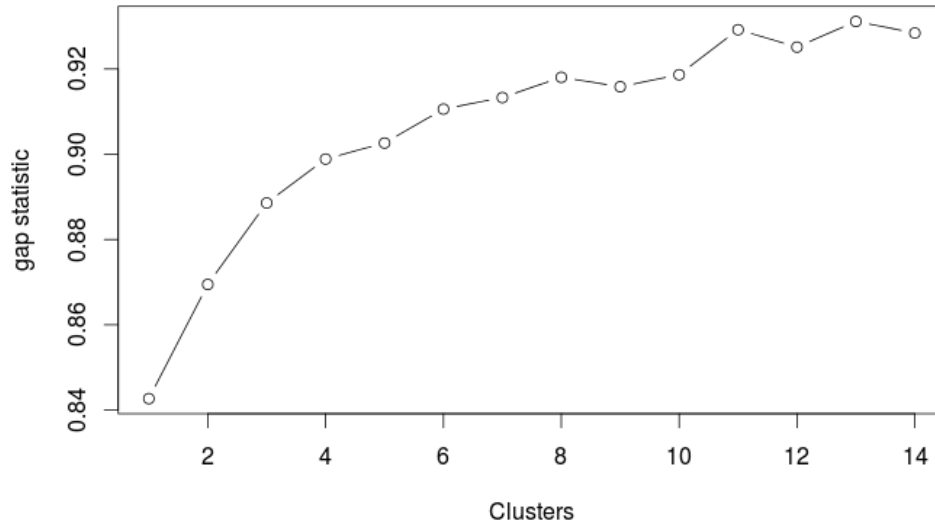


Figure 3: Gap statistic by cluster size

```

> 1      nms_d1      nms_d2      nms_d3      nms_d4      nms_d5      nms_d6
      nms_d7      nms_d8      nms_d9      tremor      bradykin
1.5558981 1.4133002 1.1184738 1.9877193 1.2550406 1.5834288
      1.4194860 0.7793719 0.8480830 0.4984324 1.5155436
      rigidity      axial      pigd
1.5055810 1.9495563 1.9242566
> 2      nms_d1      nms_d2      nms_d3      nms_d4      nms_d5
      nms_d6      nms_d7      nms_d8      nms_d9      tremor
-0.3206187 -0.4993585 -0.4670797 -0.3084269 -0.3380326
      -0.3885648 -0.3868053 -0.2817080 -0.3620072 -0.2526726
      bradykin      rigidity      axial      pigd
-0.5841212 -0.5460419 -0.5968519 -0.5927227
> 3      nms_d1      nms_d2      nms_d3      nms_d4      nms_d5
      nms_d6      nms_d7      nms_d8      nms_d9
0.35964469 0.82182933 0.92114183 0.30981044 0.75178256
      0.42320932 0.66431905 0.63701383 0.85893238
      tremor      bradykin      rigidity      axial      pigd
-0.34050774 -0.15173201 -0.20269121 0.04852427 0.04032785
> 4      nms_d1      nms_d2      nms_d3      nms_d4      nms_d5
      nms_d6      nms_d7      nms_d8      nms_d9
-0.13688715 -0.09662665 -0.13812235 -0.25339210 -0.29466914
      -0.06605132 -0.19223309 -0.17297198 -0.20498319
      tremor      bradykin      rigidity      axial      pigd
0.55647025 0.72337965 0.69299204 0.47325105 0.47914113

```

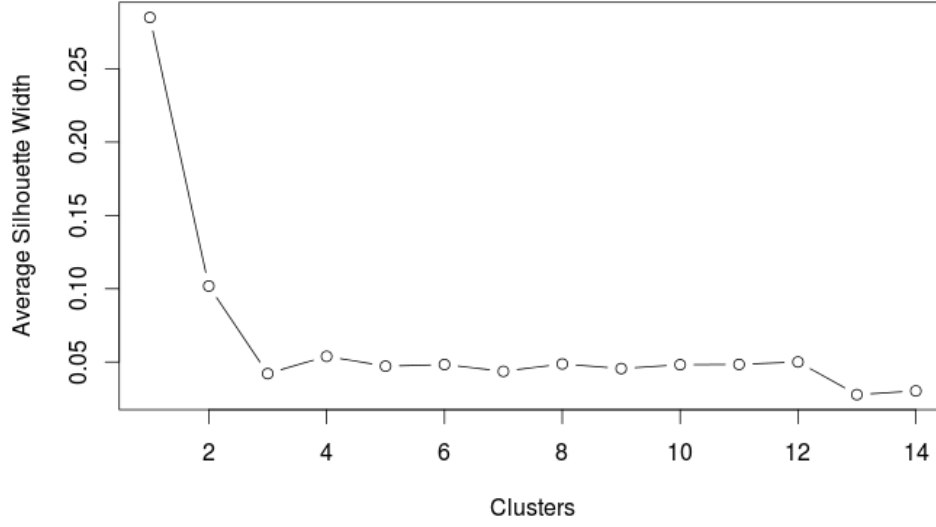


Figure 4: Average silhouette width by cluster size

2.4 Decision tree classifier based on clusters

k	CP ³	CV Xerror ⁴	Root Feature	Root Error	Figure
2	0.0348	0.134	$\text{axial} \geq 0.44$	0.257	Figure 5
3	0.0100	0.194	$\text{bradykin} < 0.0041$	0.462	Figure 6
4	0.0100	0.248	$\text{bradykin} < 0.0041$	0.523	Figure 7

Table 4: k -kmeans decision trees statistics

3 Biclustering

4 Subspace clustering

5 Bayesian Networks

³Complexity Parameter

⁴10-fold cross validation

Pruned Tree, 2 clusters

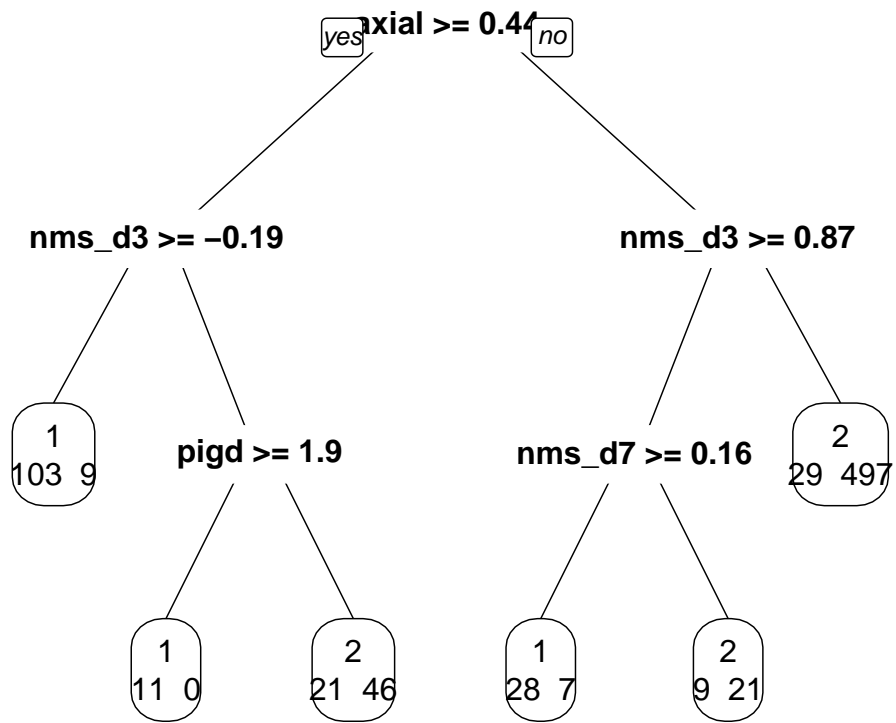


Figure 5: Decision Tree from k -means clustering, 2 clusters

Pruned Tree, 3 clusters

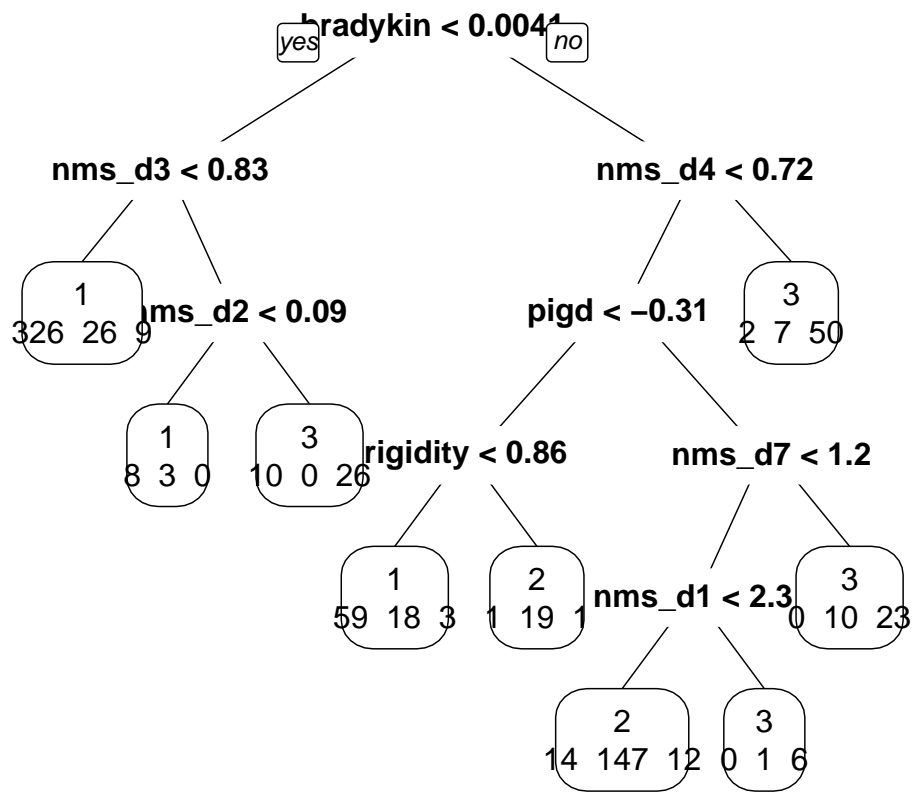


Figure 6: Decision Tree from k -means clustering, 3 clusters

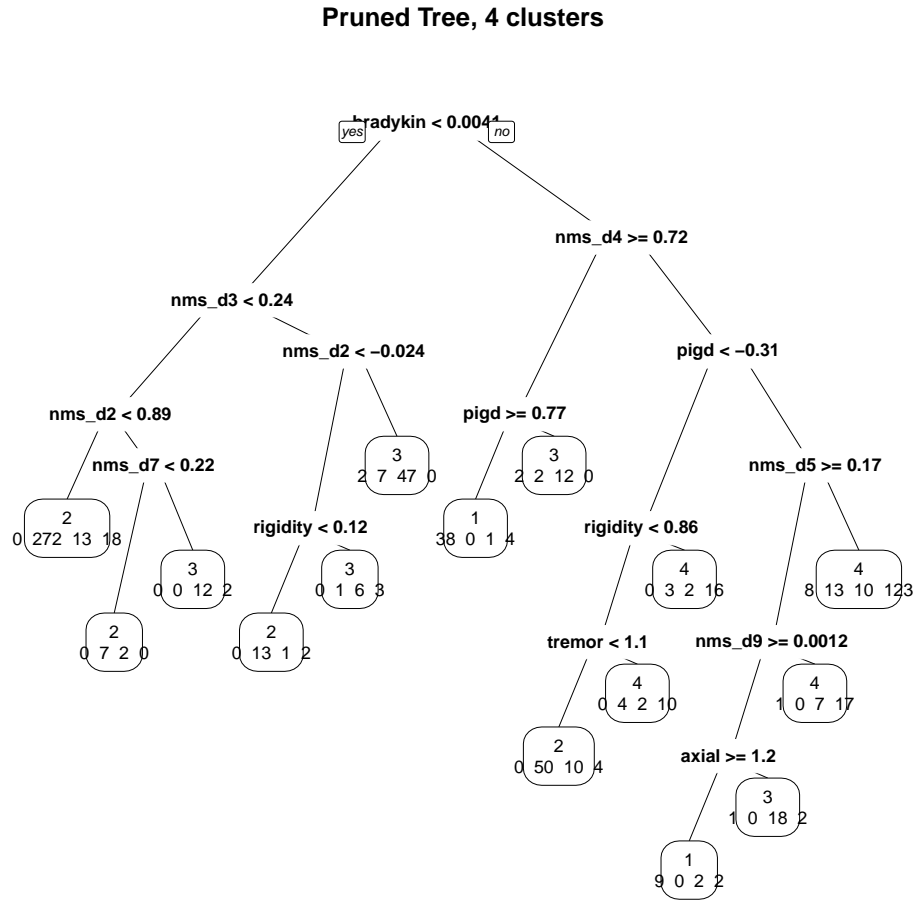


Figure 7: Decision Tree from k -means clustering, 4 clusters