

Figure 1

For each sample,  
 $s = 1, \dots, S$ , a list  
 containing CDR3  
 amino acid sequences  
 from the t cell recep-  
 tor repertoire,  $\mathbf{x}_s$ .

The parameter value  
 of  $k$ ; on current data  
 sets this has been op-  
 timised at  $k = 5$ .

Value of parameter  $\mathbf{p}$ ;  
 on current data sets  
 this has been opti-  
 mised as  $\mathbf{p} = 2, 3, 4$ ,  
 note the first principl-  
 al component is ex-  
 cluded since this cap-  
 tures between batch  
 variability.

Identify the set of all k-mers contained in  $x_1, \dots, x_S$ ;  
 this set has  $K$  elements and is denoted  $\mathcal{K} = \{k_1, \dots, k_K\}$

Calculate the  $S \times K$  k-mer frequency matrix,  $M$ , where

$$M_{ij} = \frac{\sum_t I_{x_{st}=k_j}}{\sum_t \sum_m I_{x_{st}=k_m}}$$

and  $I$  is the indicator function.

Calculate the principal components of  $M$ .  
 If  $\mathbf{w}_i = w_{i,1}, \dots, w_{i,K}$  is the  $i^{th}$  eigenvector,  
 then for sample  $s$ , the  $i^{th}$  principal component  
 is given by  $M_s \cdot \mathbf{w}_i$ , where  $M_s$  is  
 the vector obtained from the  $s^{th}$  row of  $M$ .

Hierarchically cluster the samples using  
 the principal components indexed by  $\mathbf{p}$  using  
 Ward's mumimum variance method.

Classify the samples as either Coeliac or non-coeliac  
 based on clustering. Coeliac clusters will be larger  
 and consisting of highly similar samples whilst  
 non-coeliac clusters will be smaller more diverse

The lists  $\mathbf{x}_1, \dots, \mathbf{x}_S$   
 will contain repeat se-  
 quences and will vary  
 in length, whereas  
 $\mathcal{K}$  contains only  
 unique substrings of  
 length  $k$  contained in  
 $\mathbf{x}_1, \dots, \mathbf{x}_S$

Scale the number of  
 times the k-mer is ob-  
 served in the sample  
 by the total number  
 of k-mers observed in  
 that sample, this ac-  
 counts for different  
 sequencing depths per  
 sample.

Transforms the data  
 onto a new coordinate  
 system, whereby the  
 first principal compo-  
 nent contains the  
 greatest variance

Only using the sub-  
 set of the PCs given  
 by  $\mathbf{p}$  to reduce di-  
 mensionality whilst  
 retaining the main  
 variation. This clus-  
 tering method builds  
 up a tree of clusters  
 by iteratively merg-  
 ing those which have  
 the smallest increased  
 variance.

Figure 2a

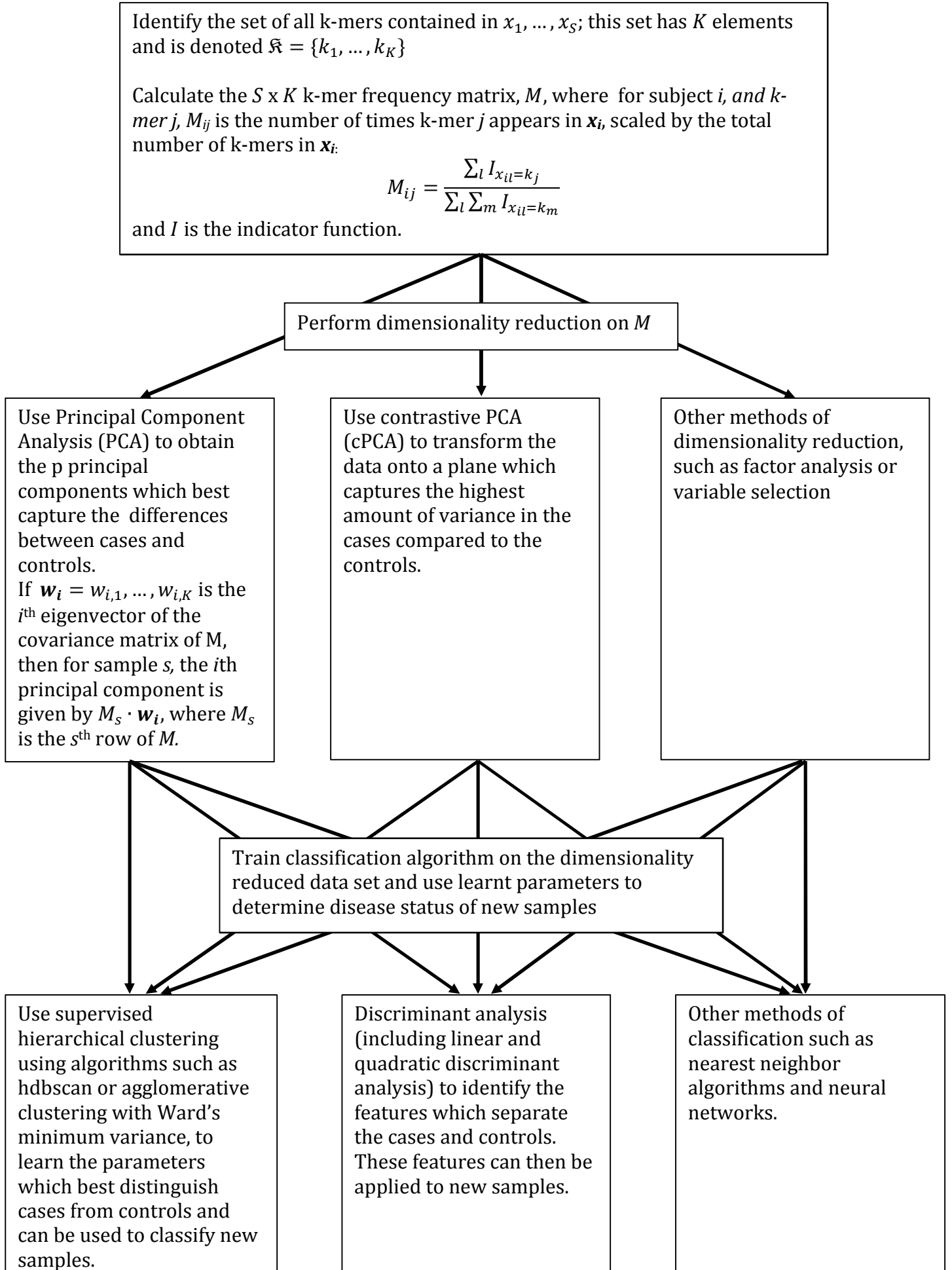


Figure 2 b

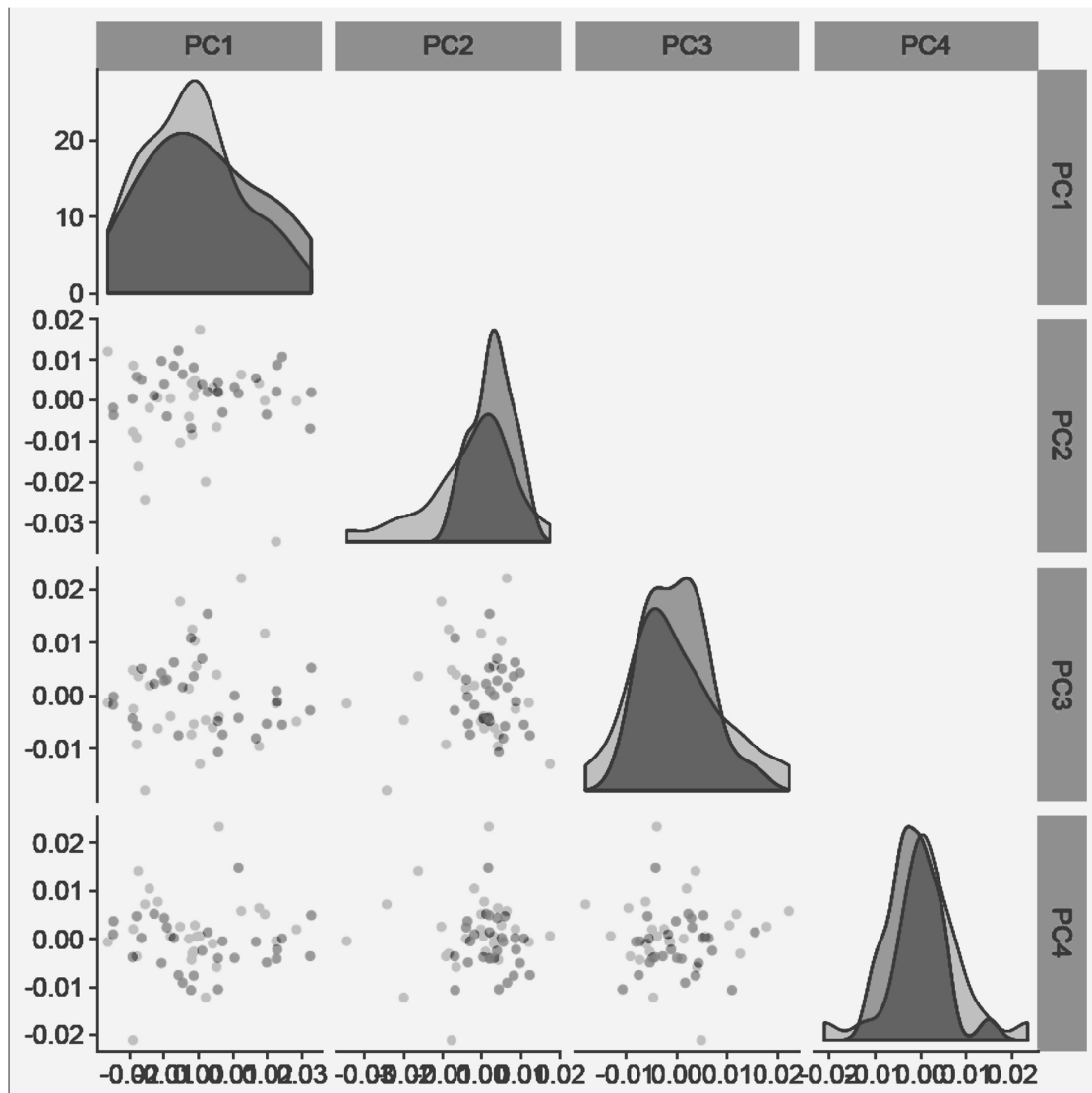


Figure 2c

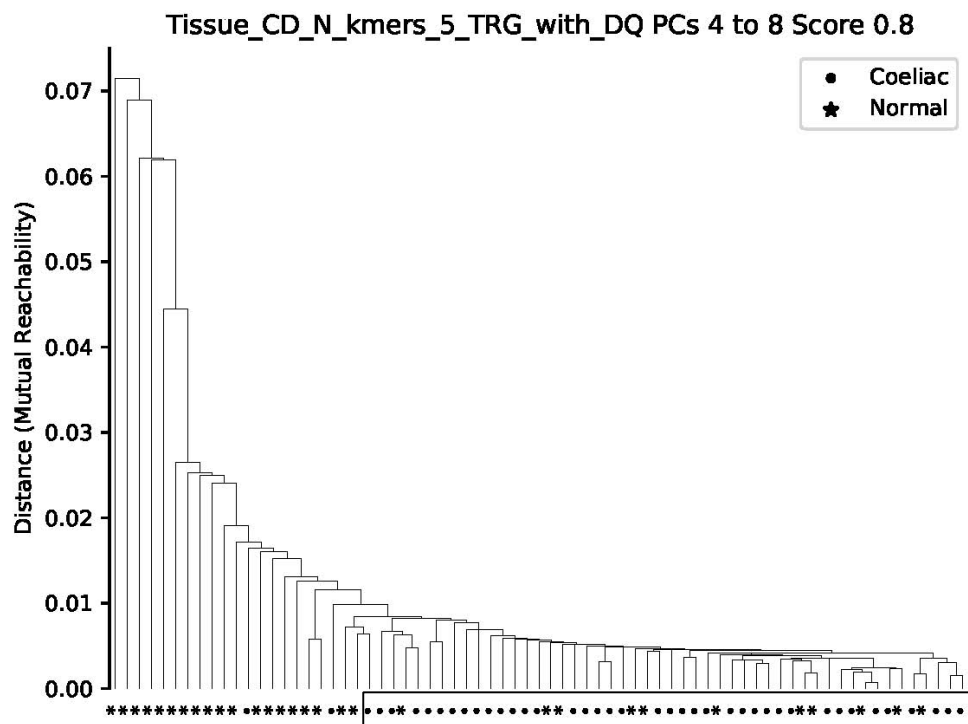


Figure 3a

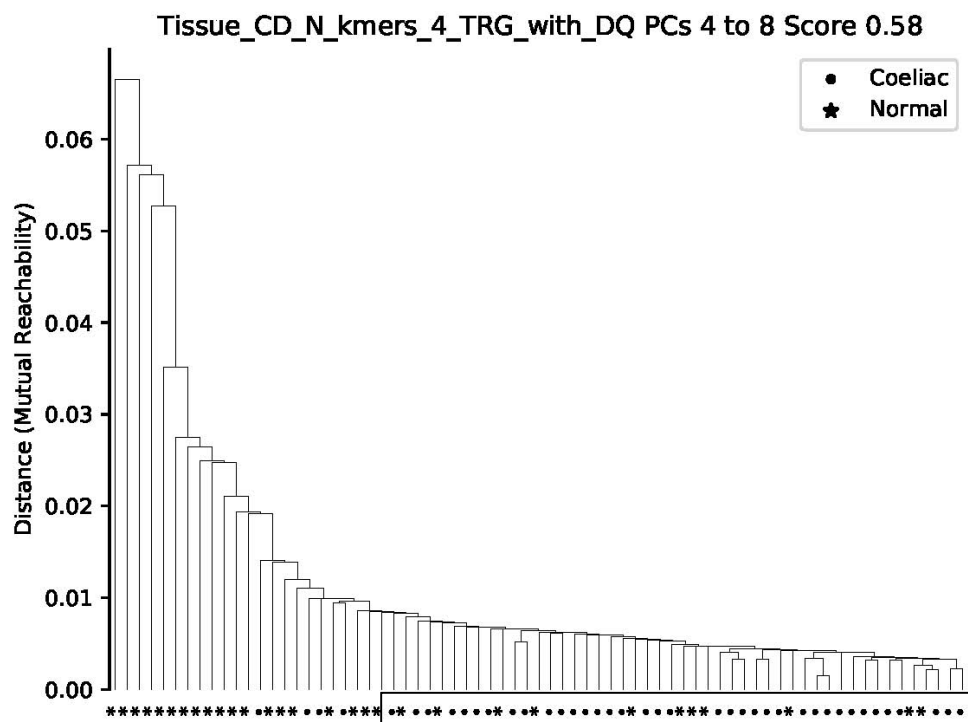


Figure 3b

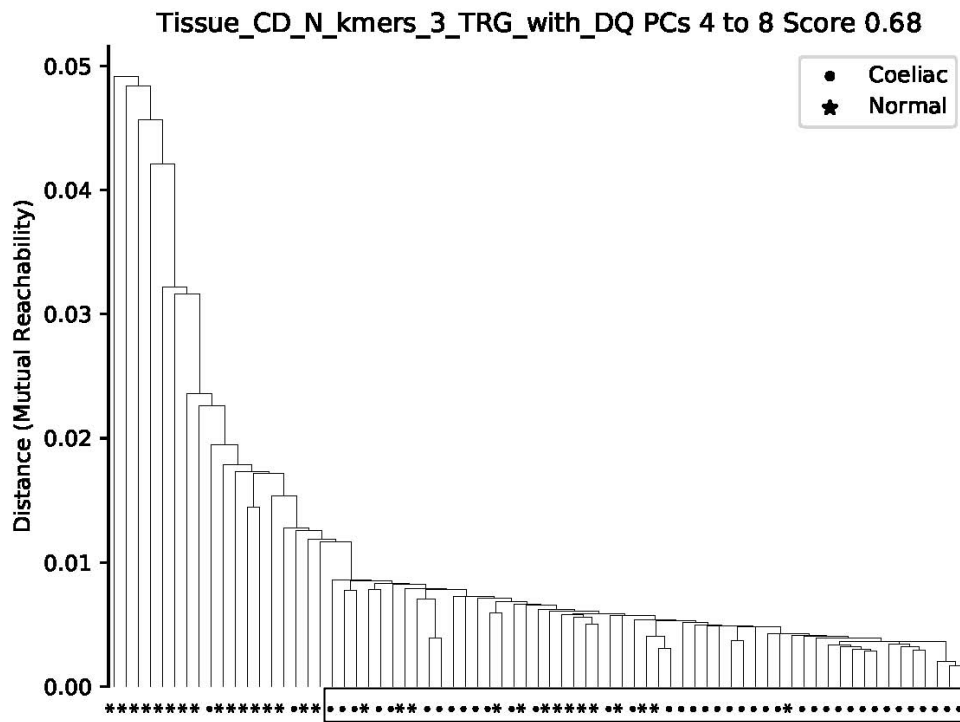


Figure 3c

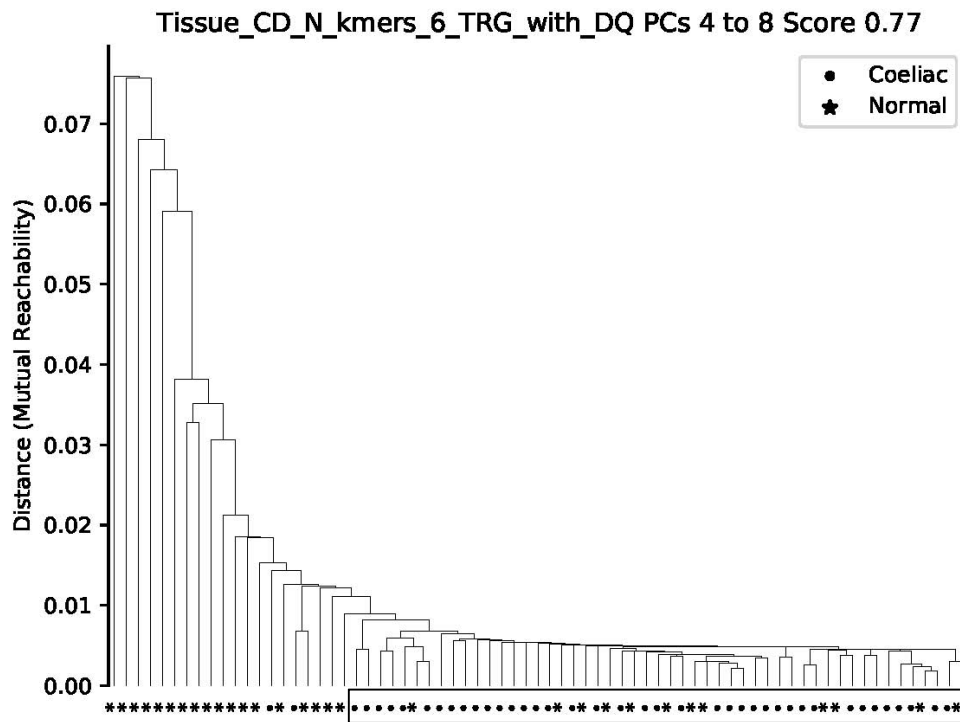


Figure 3d

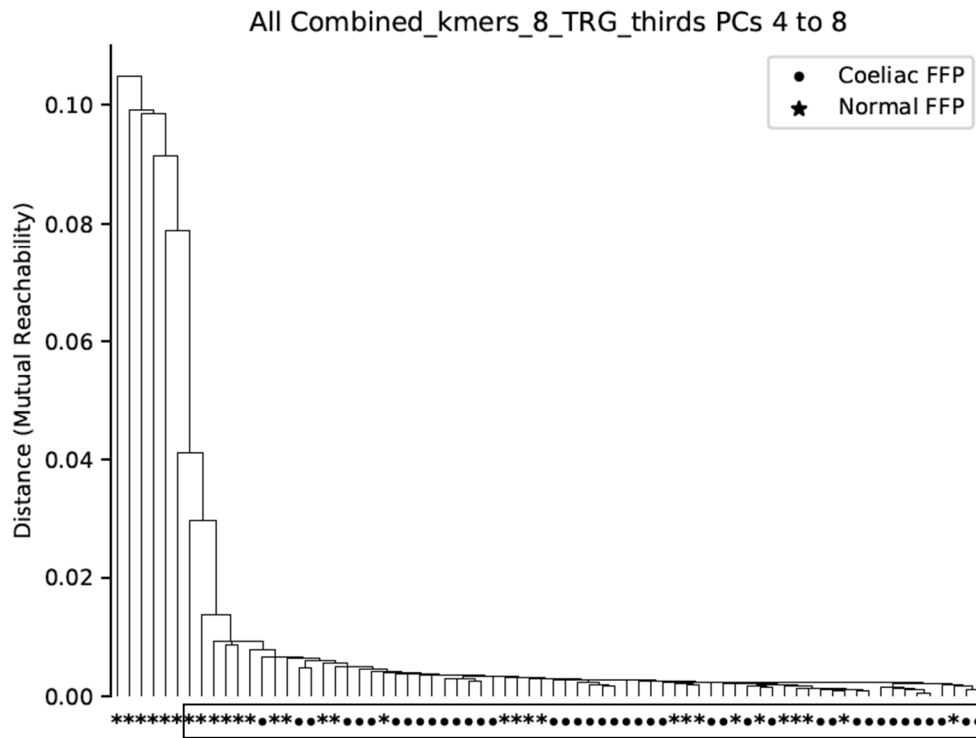


Figure 3e

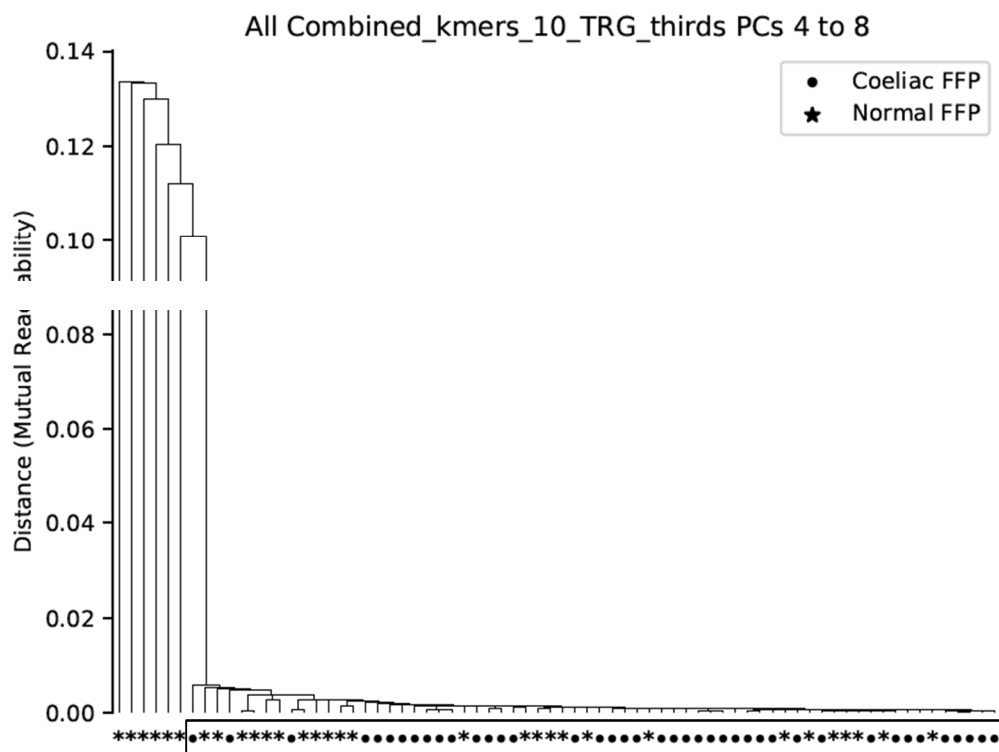


Figure 3f

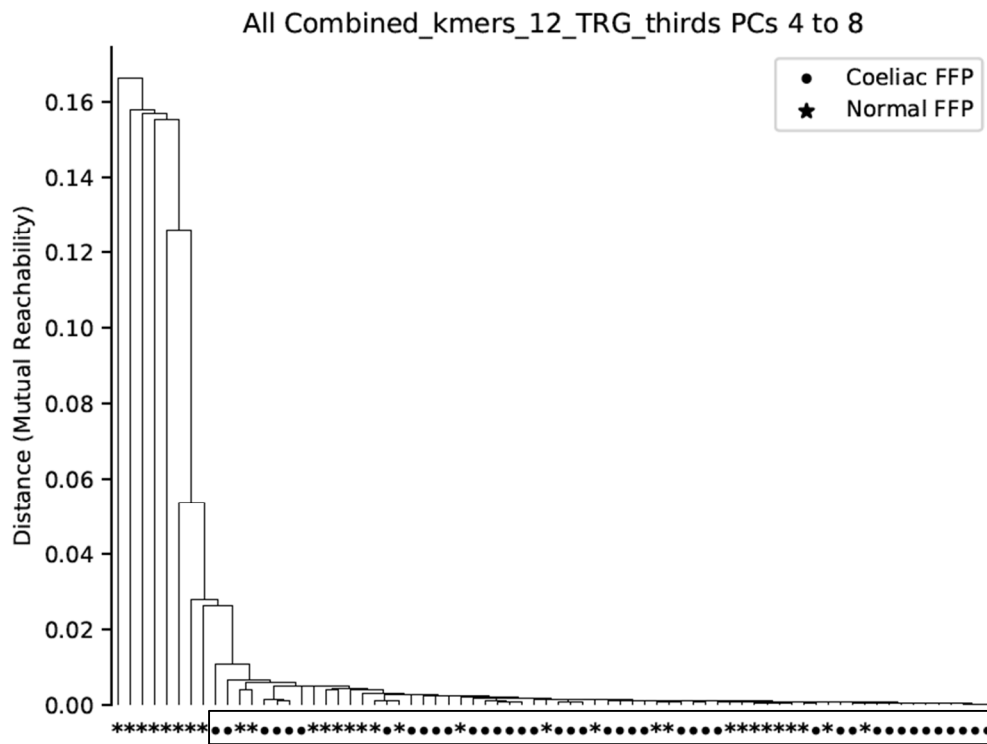


Figure 3g



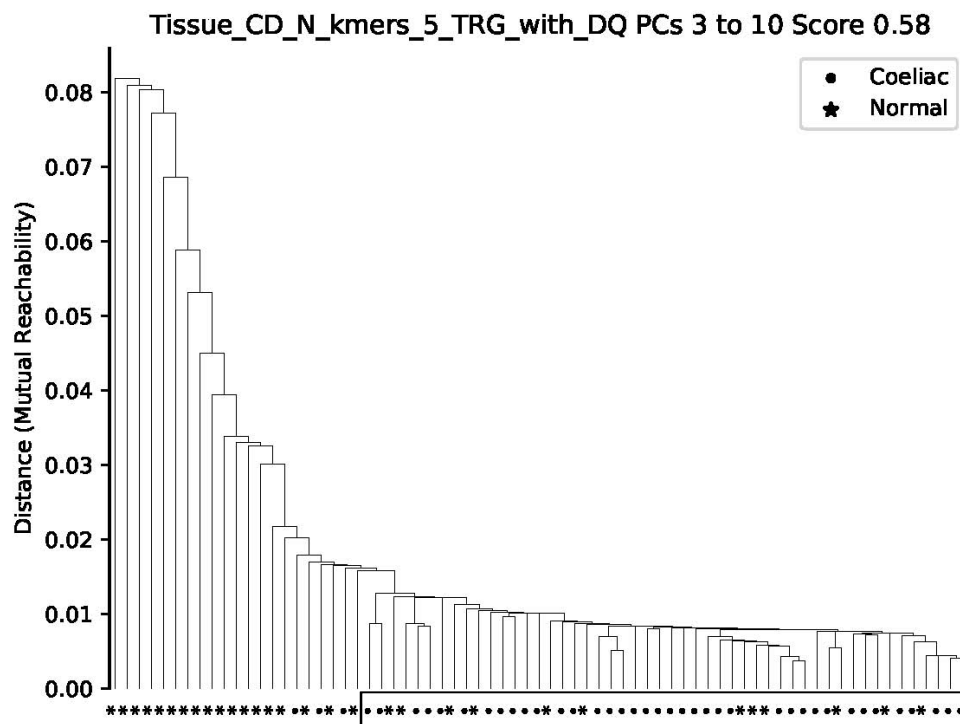


Figure 4a

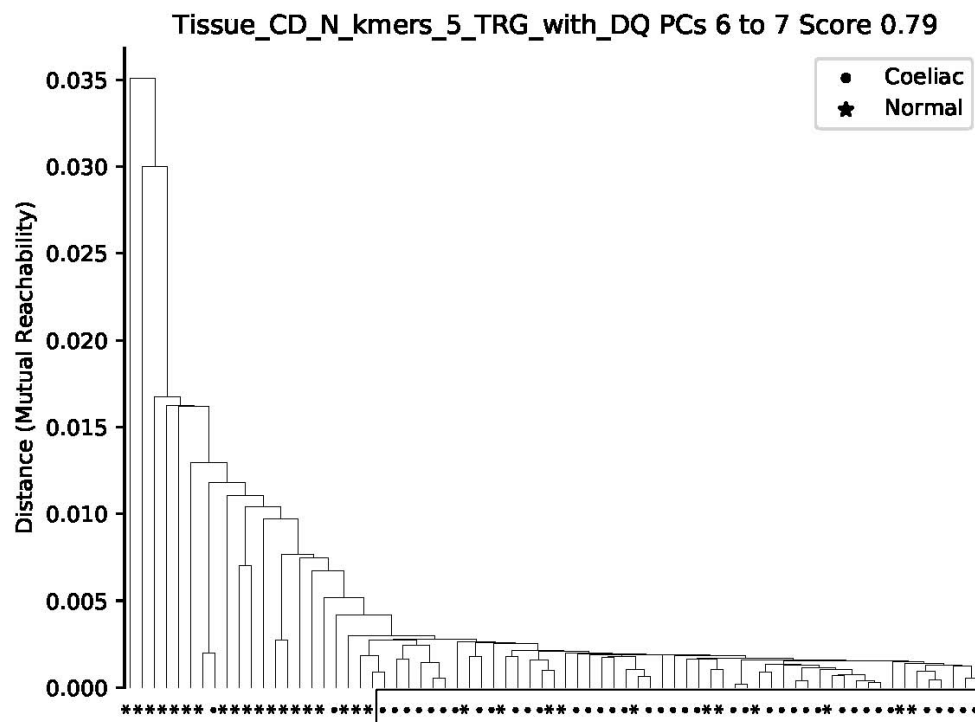


Figure 4b

kmer\_matrix\_5mers PCs 4 to 8

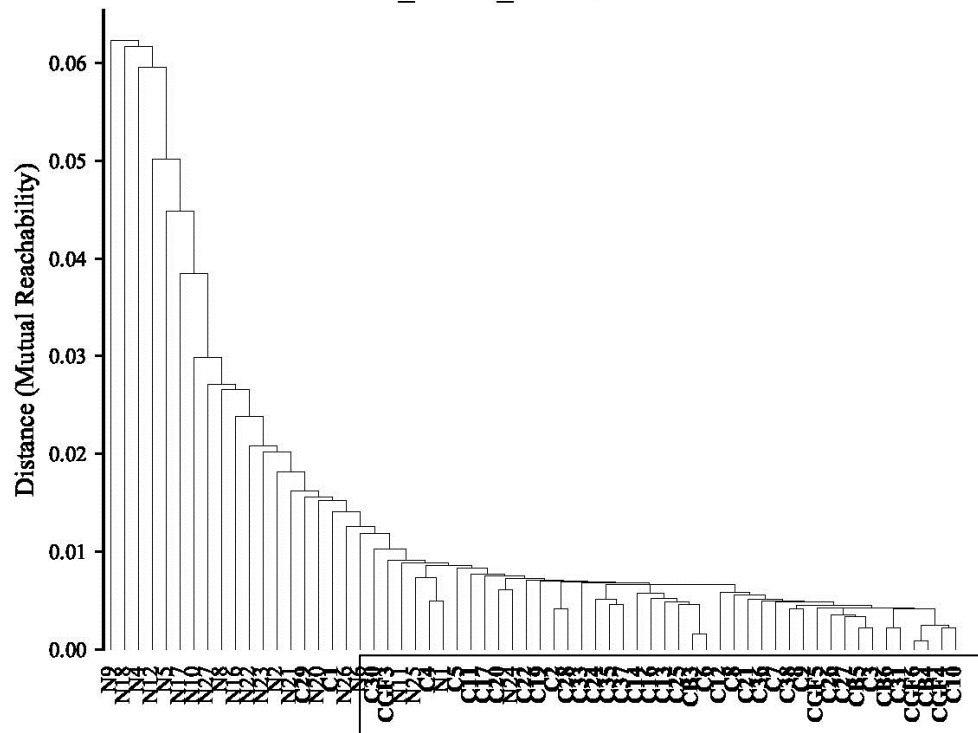


Figure 5

All Combined\_kmers\_5\_TRG PCs 4 to 8

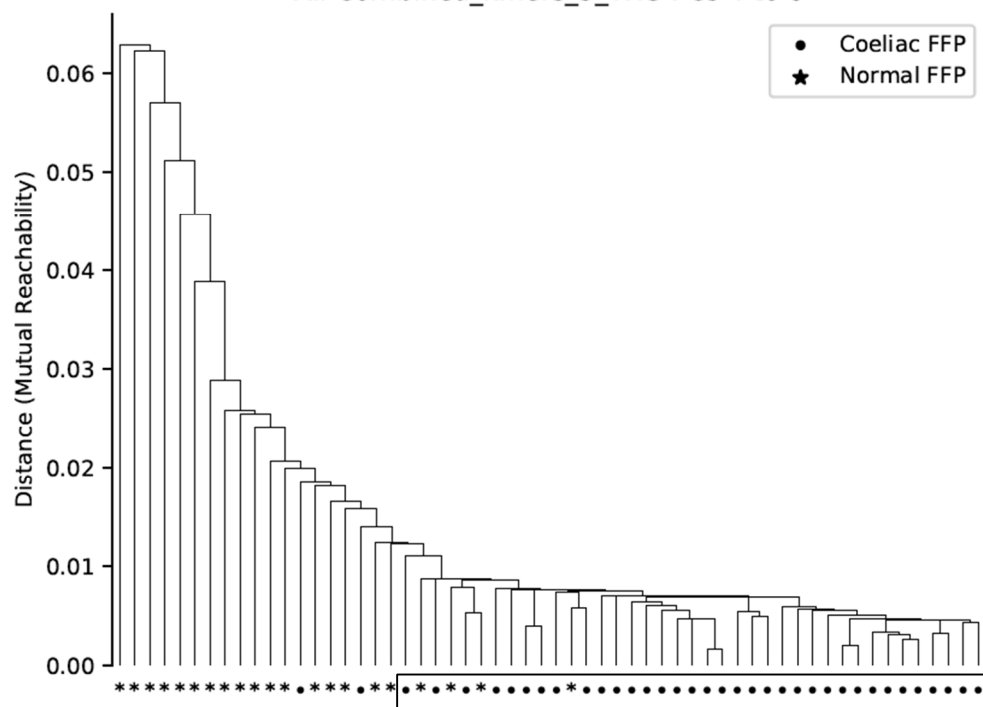


Figure 6a

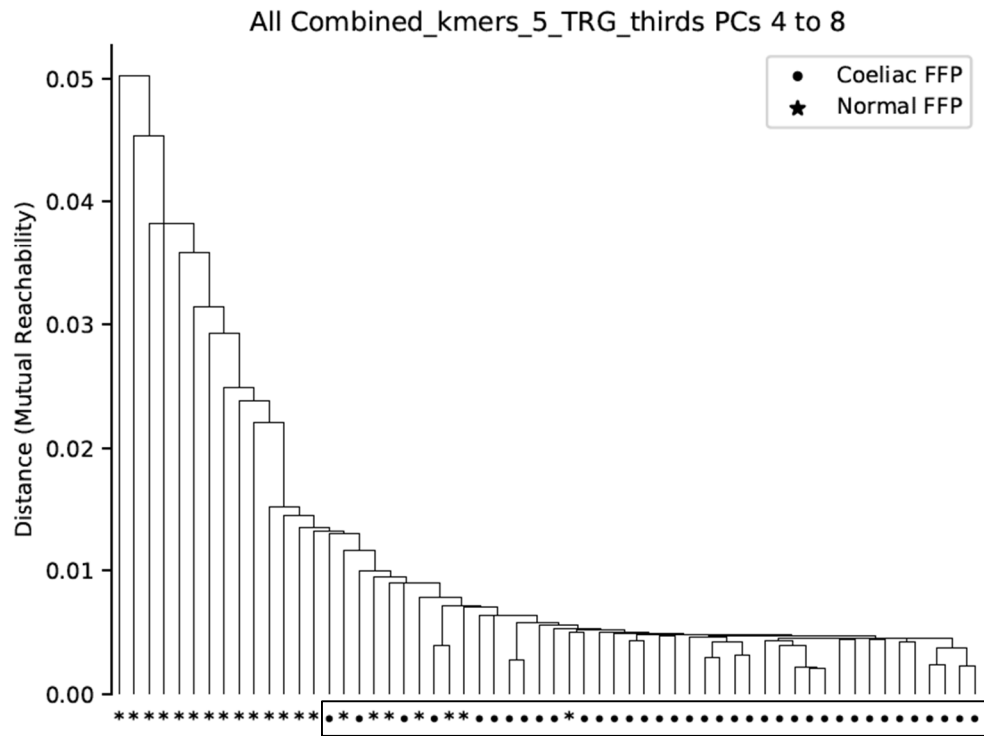


Figure 6b

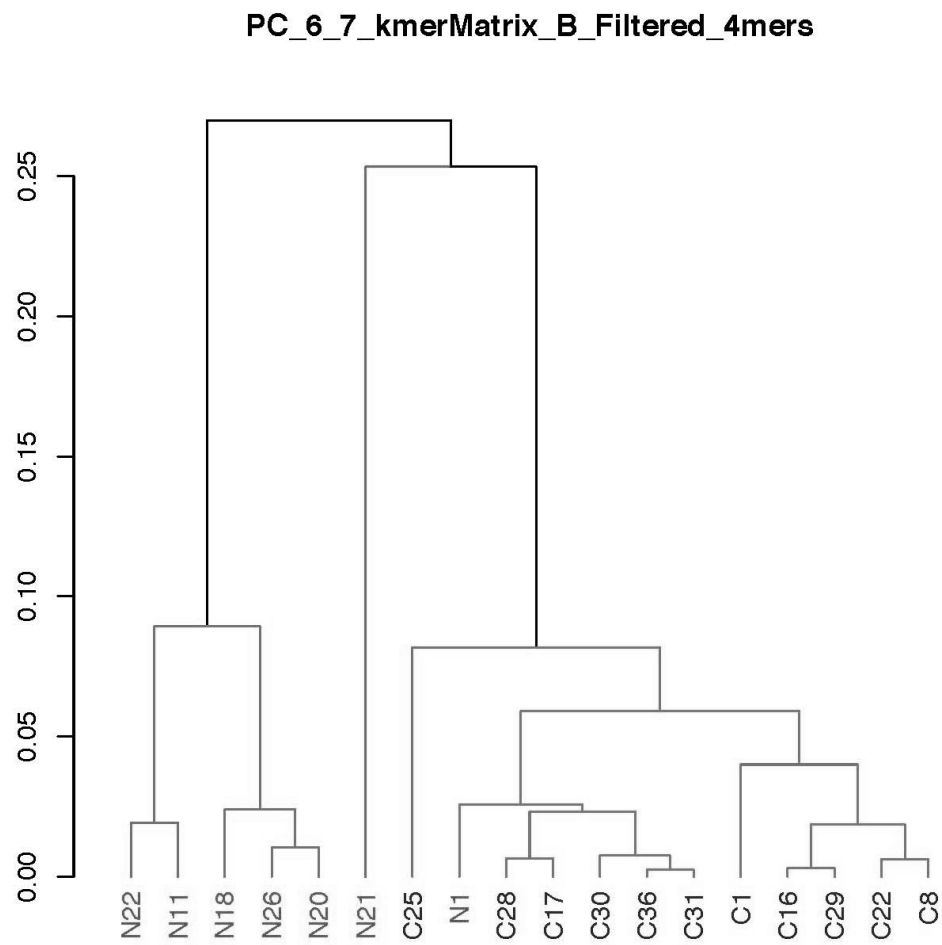


Figure 7a

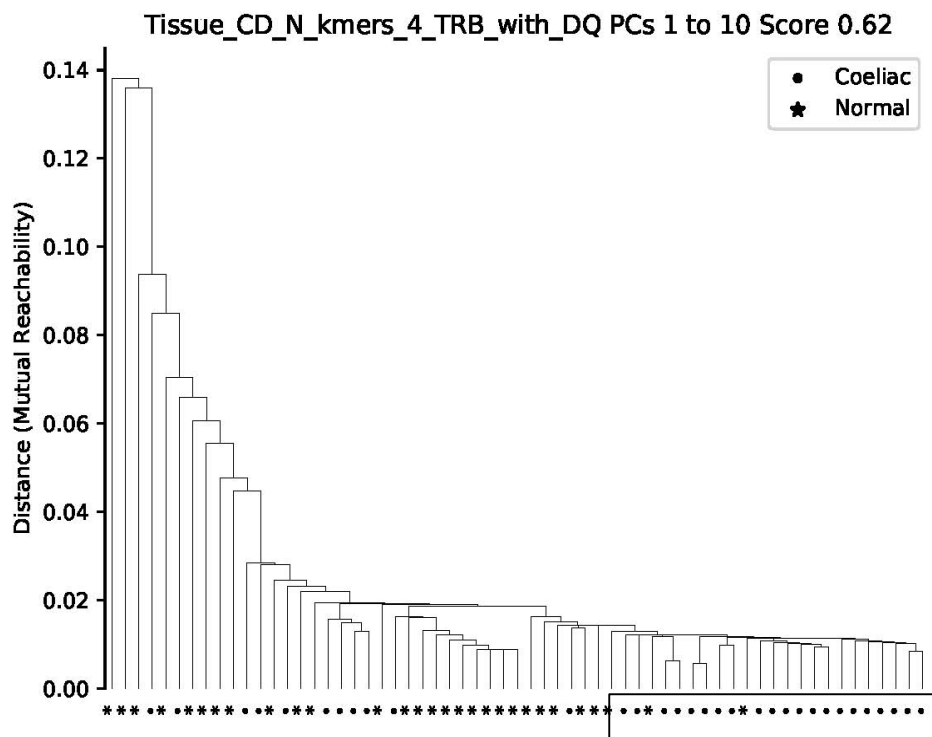


Figure 7b

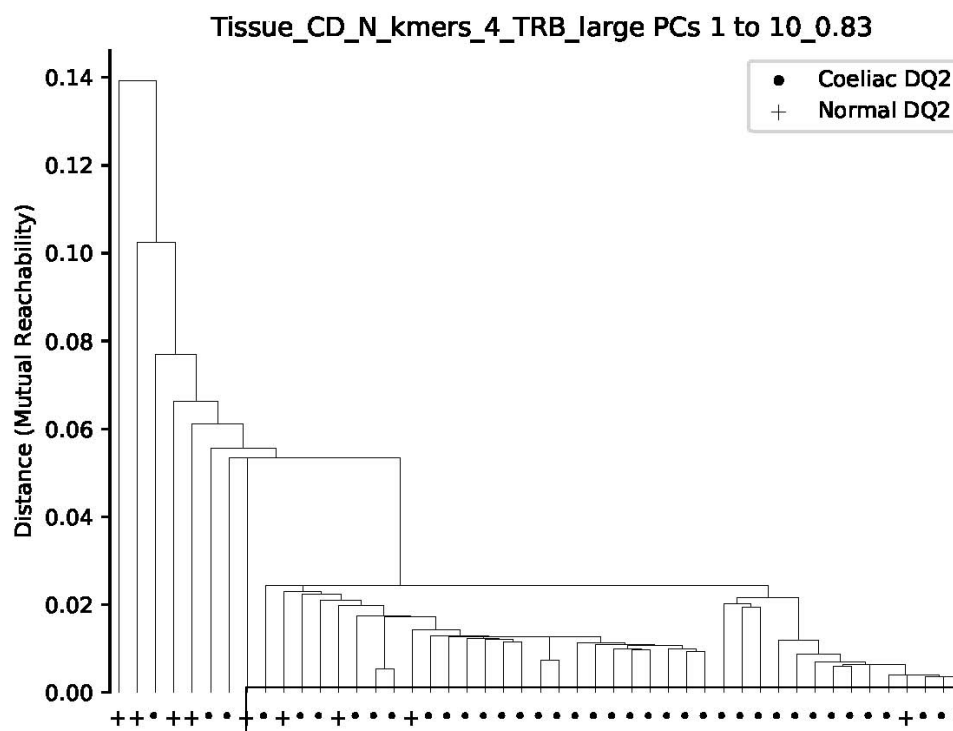


Figure 7c

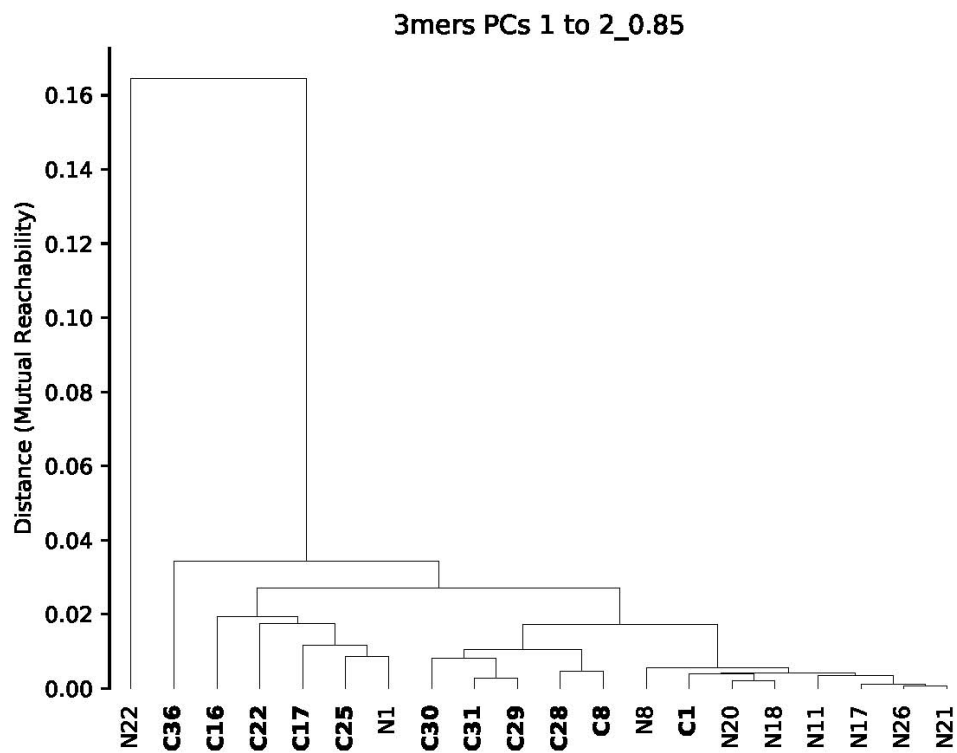


Figure 7d

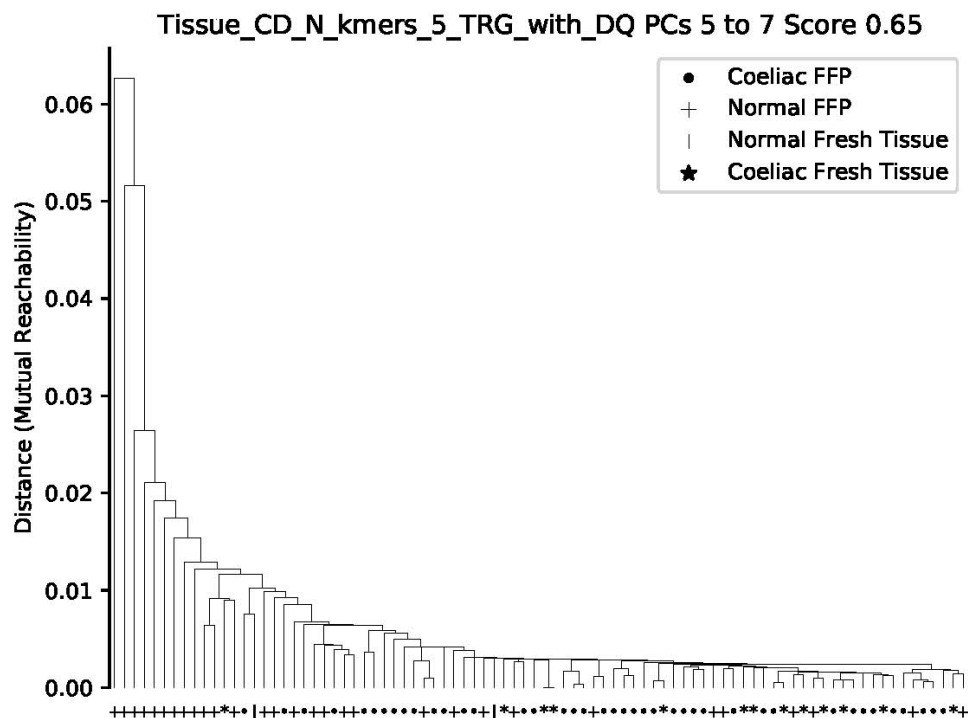


Figure 8

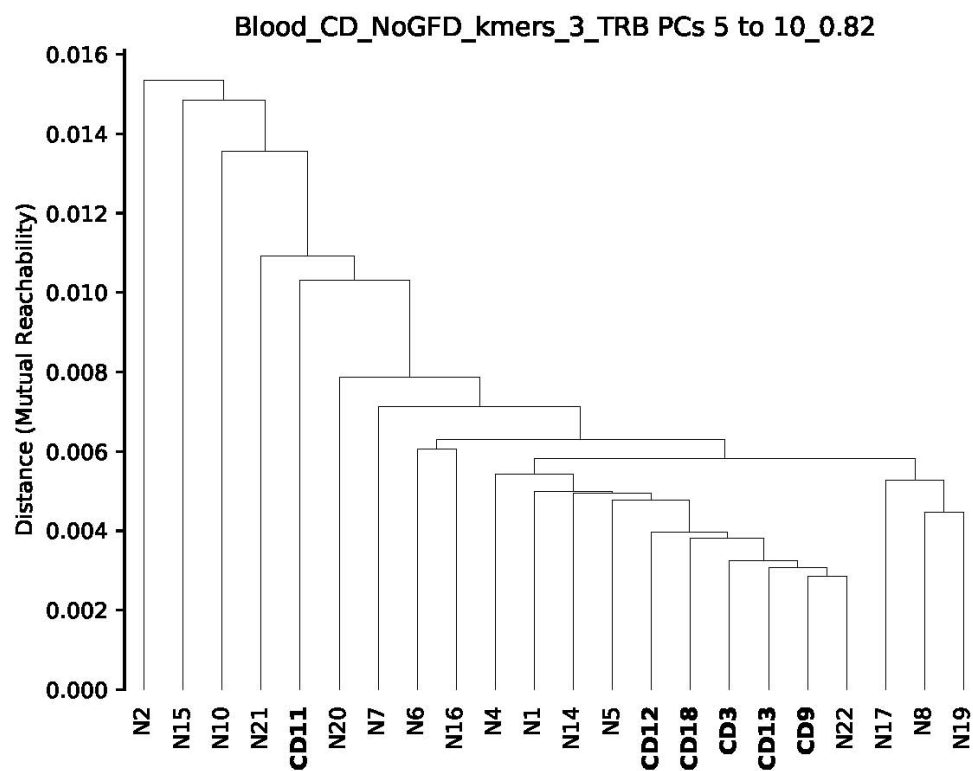


Figure 9

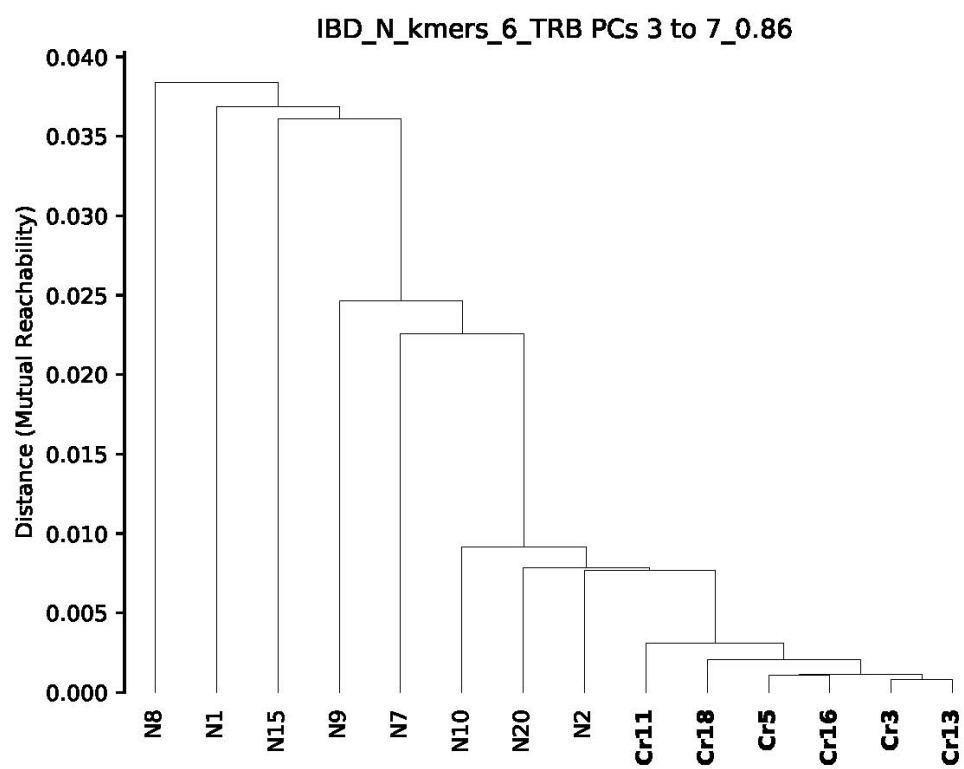


Figure 10

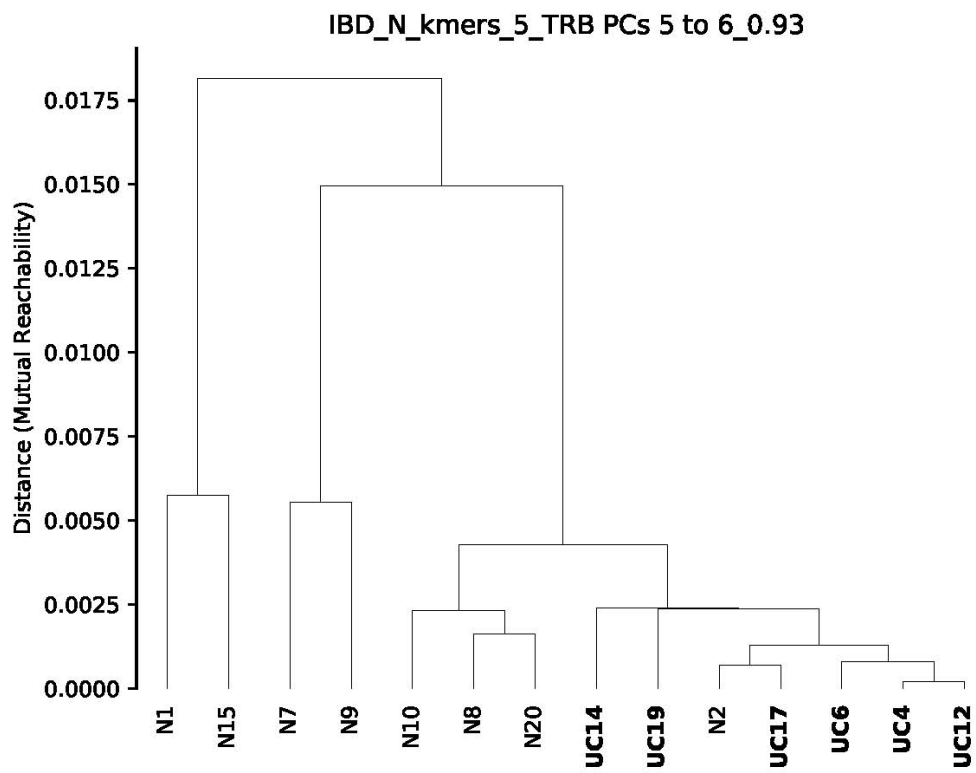


Figure 11a.

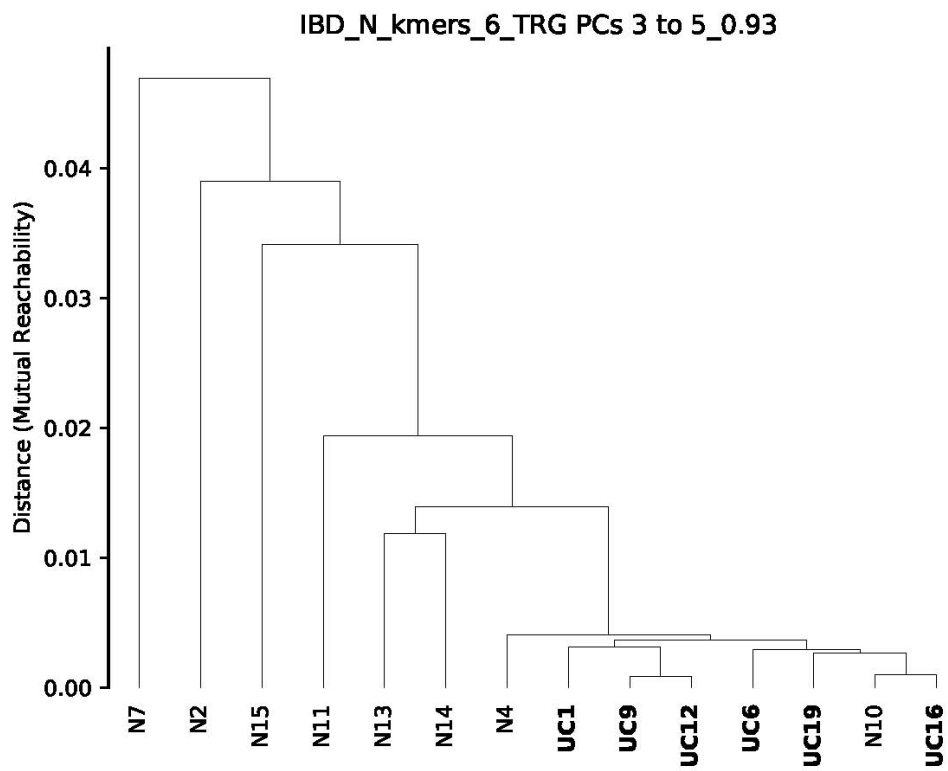


Figure 11b



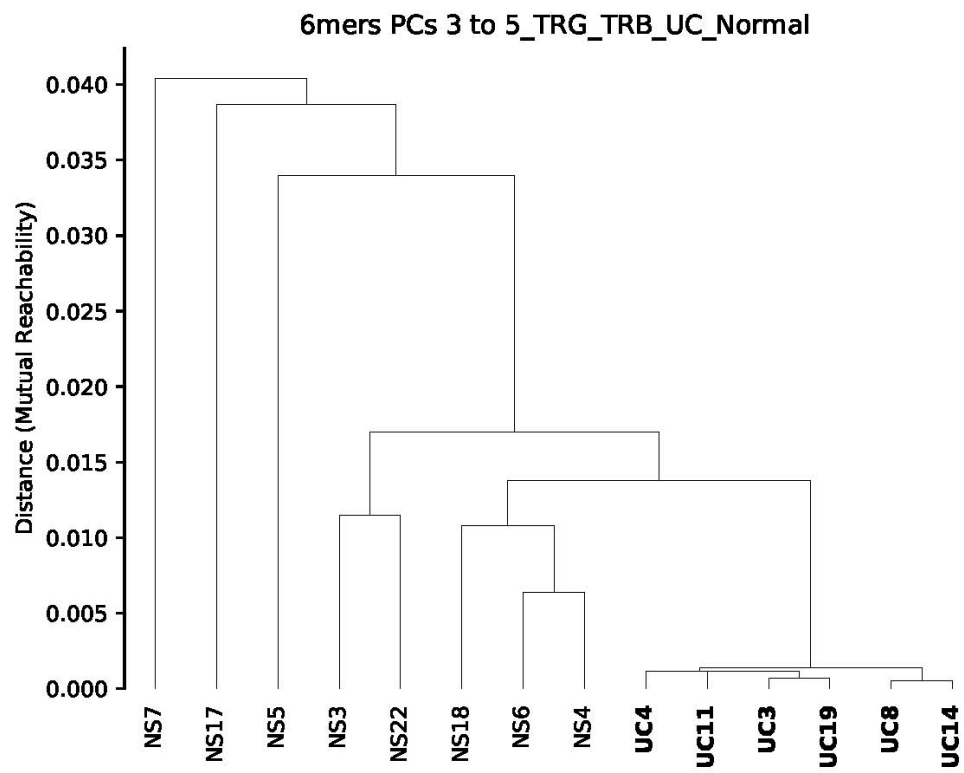


Figure 11c

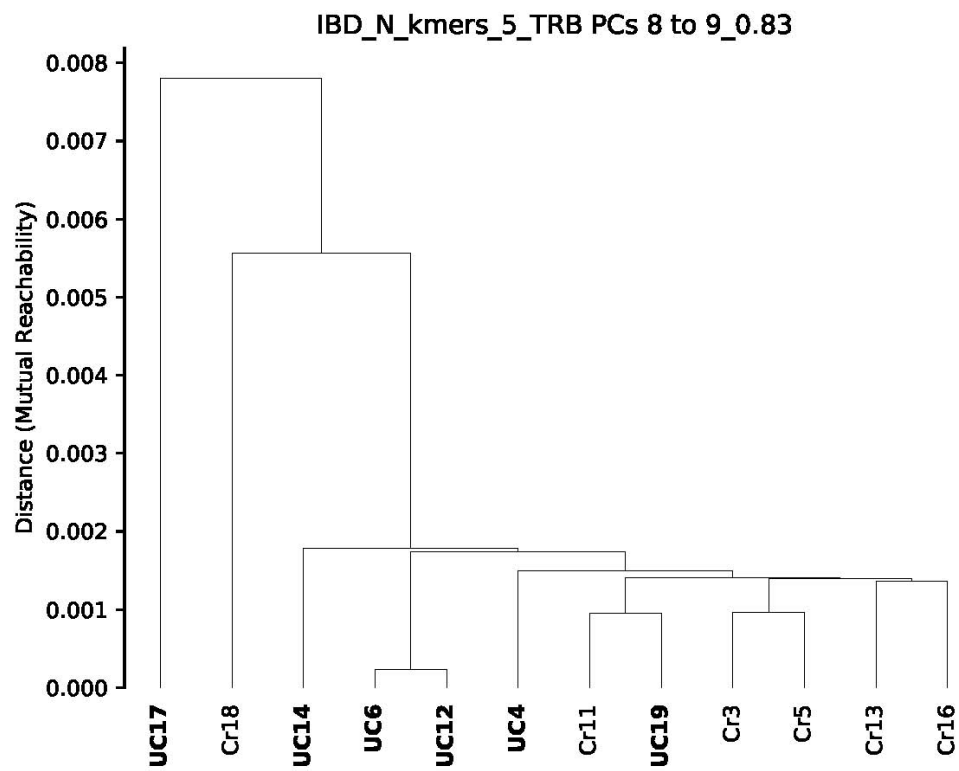


Figure 12

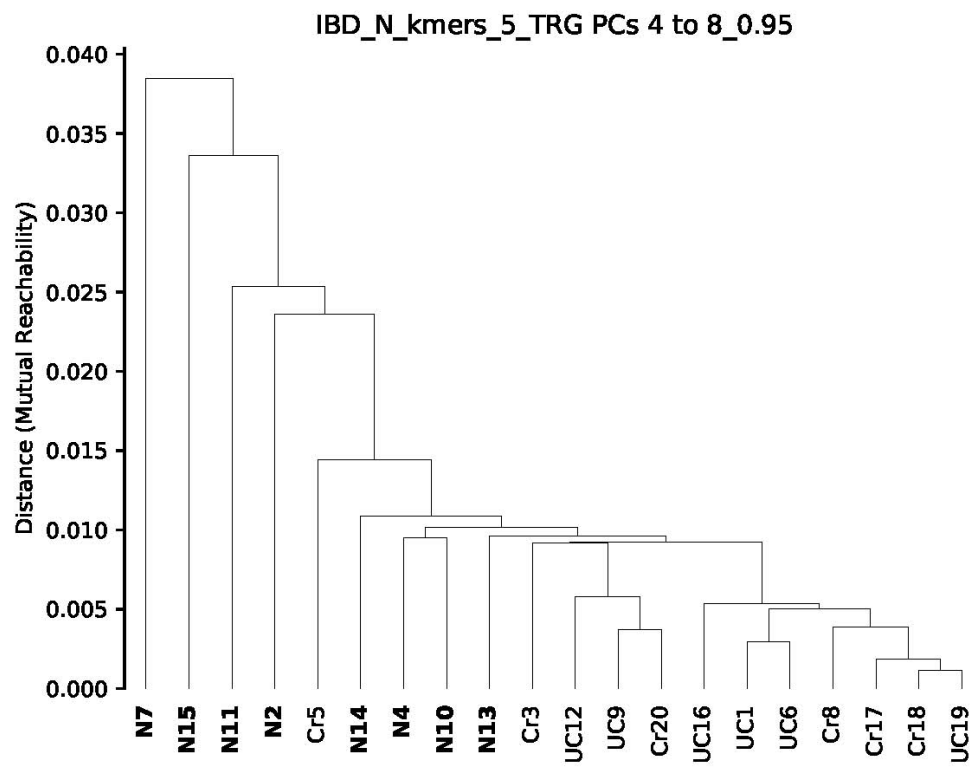


Figure 13

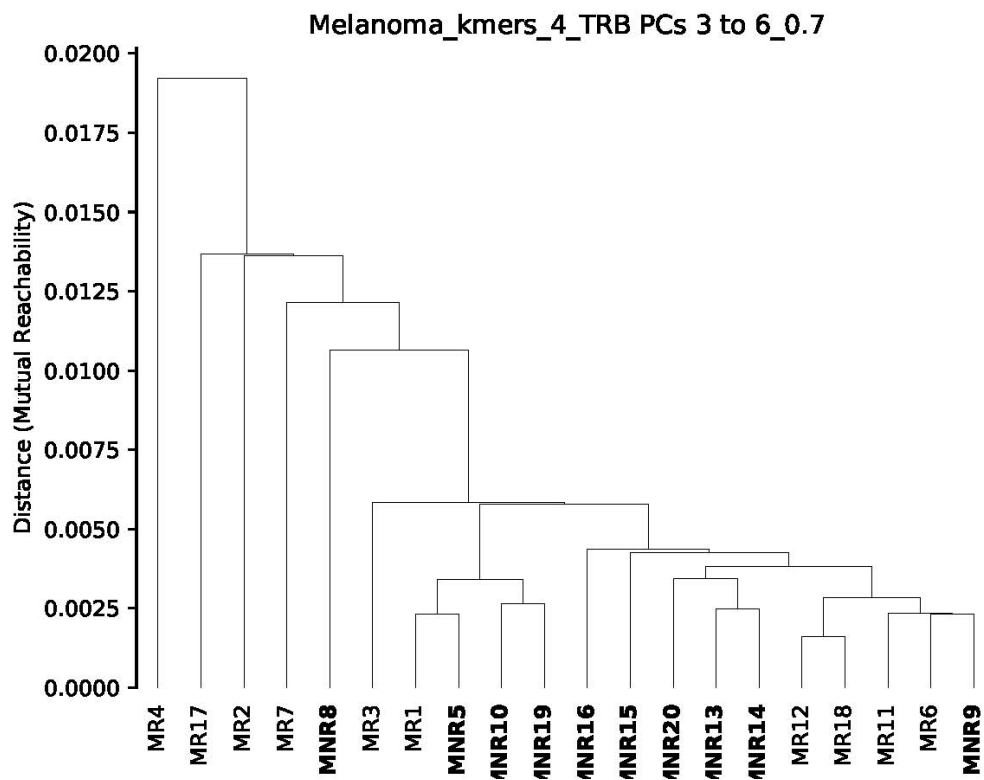


Figure 14

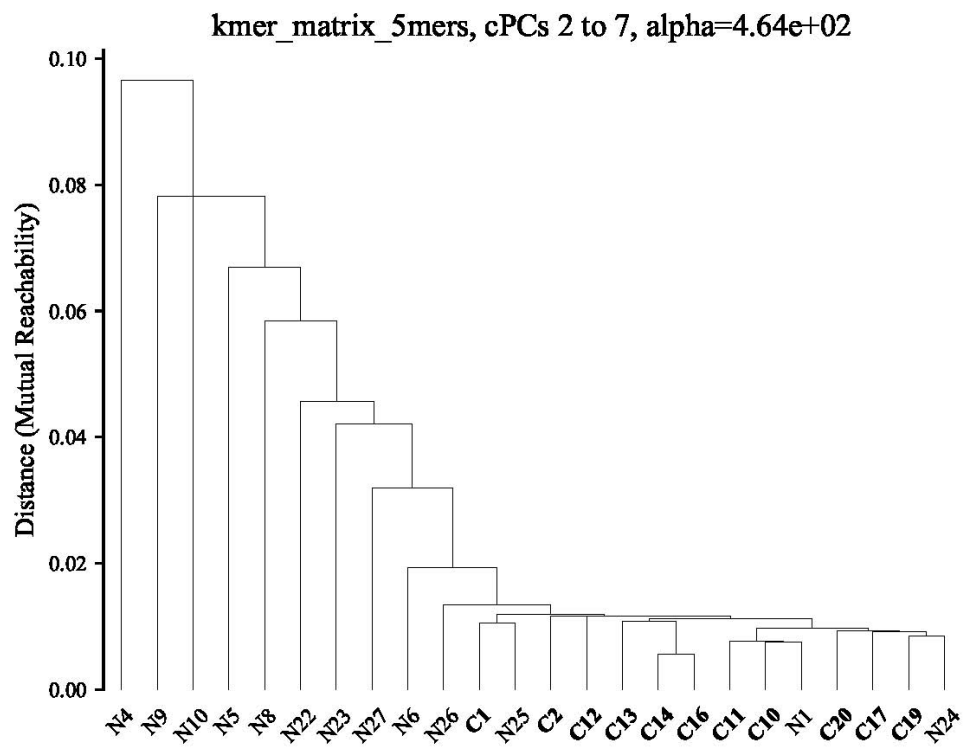


Figure 15.

| Principal components | Score (k-mers without positional annotation) | New sample prediction accuracy (without positional annotation) | Overall Accuracy (k-mers without positional annotation) | Score (k-mers with positional annotation) | New sample prediction accuracy (with positional annotation) | Overall Accuracy (k-mers with positional annotation) |
|----------------------|--|--|---|---|---|--|
| 1 to 2               | 0.778  | 1.00   | 0.806   | 0.815                                     | 1.00  | 0.839  |
| 1 to 3               | 0.833  | 1.00   | 0.855   | 0.833                                     | 1.00  | 0.855  |
| 1 to 4               | 0.870  | 1.00   | 0.887   | 0.833                                     | 1.00  | 0.855  |
| 1 to 5               | 0.815  | 1.00   | 0.839   | 0.852                                     | 1.00  | 0.871  |
| 1 to 6               | 0.889  | 1.00   | 0.903   | 0.889                                     | 1.00  | 0.903  |
| 1 to 7               | 0.926  | 1.00   | 0.935   | 0.889                                     | 1.00  | 0.903  |
| 1 to 8               | 0.926  | 1.00   | 0.935   | 0.907                                     | 1.00  | 0.919  |
| 1 to 9               | 0.889  | 1.00   | 0.903   | 0.907                                     | 1.00  | 0.919  |
| 1 to 10              | 0.907  | 1.00   | 0.919   | 0.907                                     | 1.00  | 0.919  |
| 2 to 3               | 0.759  | 1.00   | 0.790   | 0.796                                     | 1.00  | 0.823  |
| 2 to 4               | 0.796  | 1.00   | 0.823   | 0.833                                     | 1.00  | 0.855  |
| 2 to 5               | 0.796  | 1.00   | 0.823   | 0.852                                     | 1.00  | 0.871  |
| 2 to 6               | 0.852  | 1.00   | 0.871   | 0.870                                     | 1.00  | 0.887  |
| 2 to 7               | 0.889  | 1.00   | 0.903   | 0.870                                     | 1.00  | 0.887  |
| 2 to 8               | 0.926  | 1.00   | 0.935   | 0.907                                     | 1.00  | 0.919  |
| 2 to 9               | 0.926  | 1.00   | 0.935   | 0.907                                     | 1.00  | 0.919  |
| 2 to 10              | 0.926  | 1.00   | 0.935   | 0.907                                     | 1.00  | 0.919  |
| 3 to 4               | 0.759  | 1.00   | 0.790   | 0.759                                     | 1.00  | 0.790  |
| 3 to 5               | 0.815  | 1.00   | 0.839   | 0.852                                     | 1.00  | 0.871  |
| 3 to 6               | 0.852  | 1.00   | 0.871   | 0.833                                     | 1.00  | 0.855  |
| 3 to 7               | 0.889  | 1.00   | 0.903   | 0.852                                     | 1.00  | 0.871  |
| 3 to 8               | 0.926  | 1.00   | 0.935   | 0.926                                     | 1.00  | 0.935  |
| 3 to 9               | 0.926  | 1.00   | 0.935   | 0.907                                     | 1.00  | 0.919  |
| 3 to 10              | 0.944  | 1.00   | 0.952   | 0.907                                     | 1.00  | 0.919  |
| 4 to 5               | 0.796  | 1.00   | 0.823   | 0.852                                     | 1.00  | 0.871  |
| 4 to 6               | 0.870  | 1.00   | 0.887   | 0.852                                     | 1.00  | 0.871  |
| 4 to 7               | 0.870  | 1.00   | 0.887   | 0.852                                     | 1.00  | 0.871  |
| 4 to 8               | 0.944  | 1.00   | 0.952   | 0.926                                     | 1.00  | 0.935  |
| 4 to 9               | 0.926  | 1.00   | 0.935   | 0.926                                     | 1.00  | 0.935  |
| 4 to 10              | 0.907  | 1.00   | 0.919   | 0.907                                     | 1.00  | 0.919  |
| 5 to 6               | 0.796  | 1.00   | 0.823   | 0.815                                     | 1.00  | 0.839  |
| 5 to 7               | 0.833  | 1.00   | 0.855   | 0.815                                     | 1.00  | 0.839  |
| 5 to 8               | 0.907  | 1.00   | 0.919   | 0.907                                     | 1.00  | 0.919  |
| 5 to 9               | 0.907  | 1.00   | 0.919   | 0.907                                     | 1.00  | 0.919  |
| 5 to 10              | 0.907  | 1.00   | 0.919   | 0.926                                     | 1.00  | 0.935  |
| 6 to 7               | 0.833  | 1.00   | 0.855   | 0.852                                     | 1.00  | 0.871  |
| 6 to 8               | 0.907  | 1.00   | 0.919   | 0.907                                     | 1.00  | 0.919  |
| 6 to 9               | 0.907  | 1.00   | 0.919   | 0.907                                     | 1.00  | 0.919  |
| 6 to 10              | 0.889  | 1.00   | 0.903   | 0.889                                     | 1.00  | 0.903  |
| 7 to 8               | 0.833  | 1.00   | 0.855   | 0.852                                     | 1.00  | 0.871  |
| 7 to 9               | 0.852  | 1.00   | 0.871   | 0.889                                     | 1.00  | 0.903  |
| 7 to 10              | 0.852  | 1.00   | 0.871   | 0.852                                     | 1.00  | 0.871  |
| 8 to 9               | 0.815  | 1.00   | 0.839   | 0.852                                     | 1.00  | 0.871  |
| 8 to 10              | 0.852  | 1.00   | 0.871   | 0.833                                     | 1.00  | 0.855  |
| 9 to 10              | 0.796  | 1.00   | 0.823   | 0.796                                     | 1.00  | 0.823  |

Figure 16.

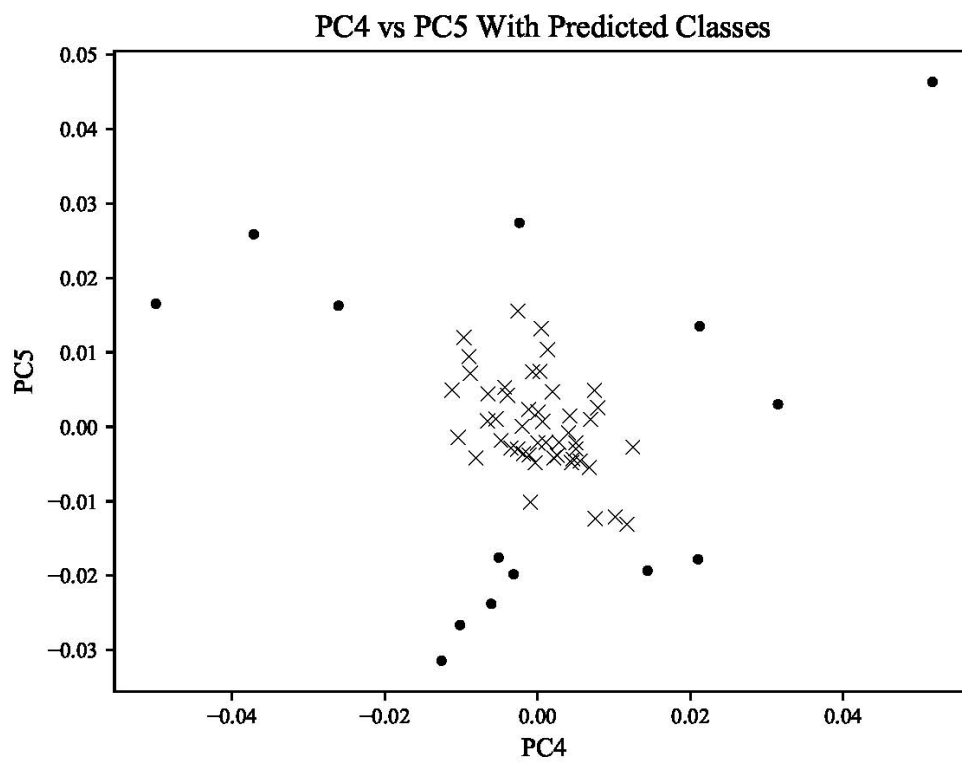


Figure 17a.

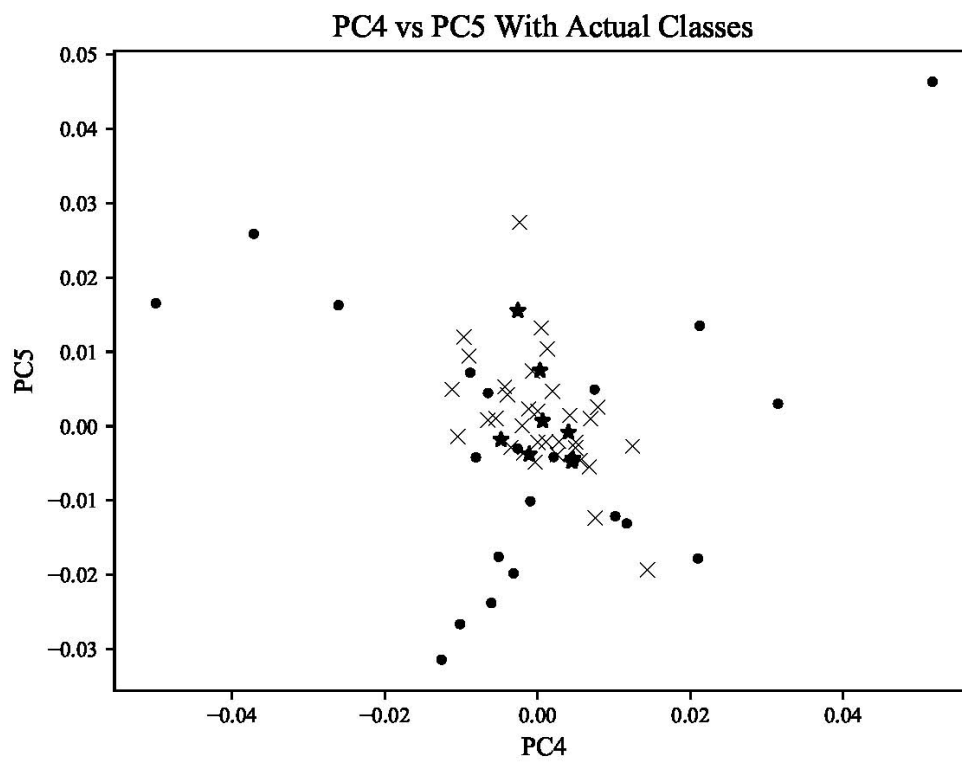


Figure 17b.

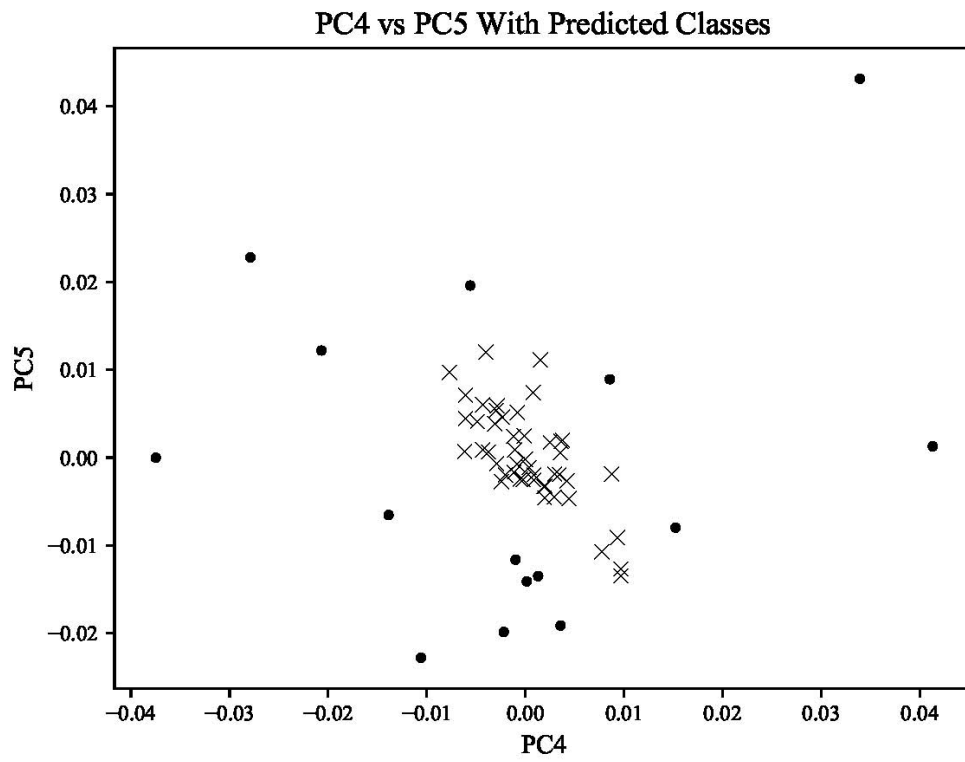


Figure 17c.

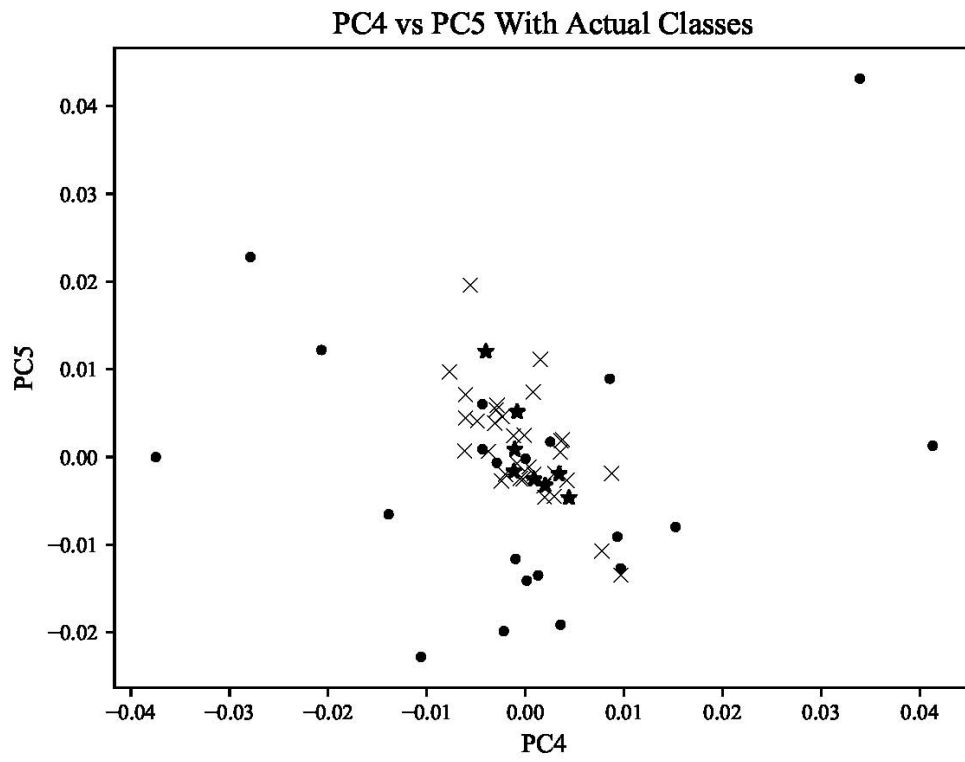


Figure 17d.

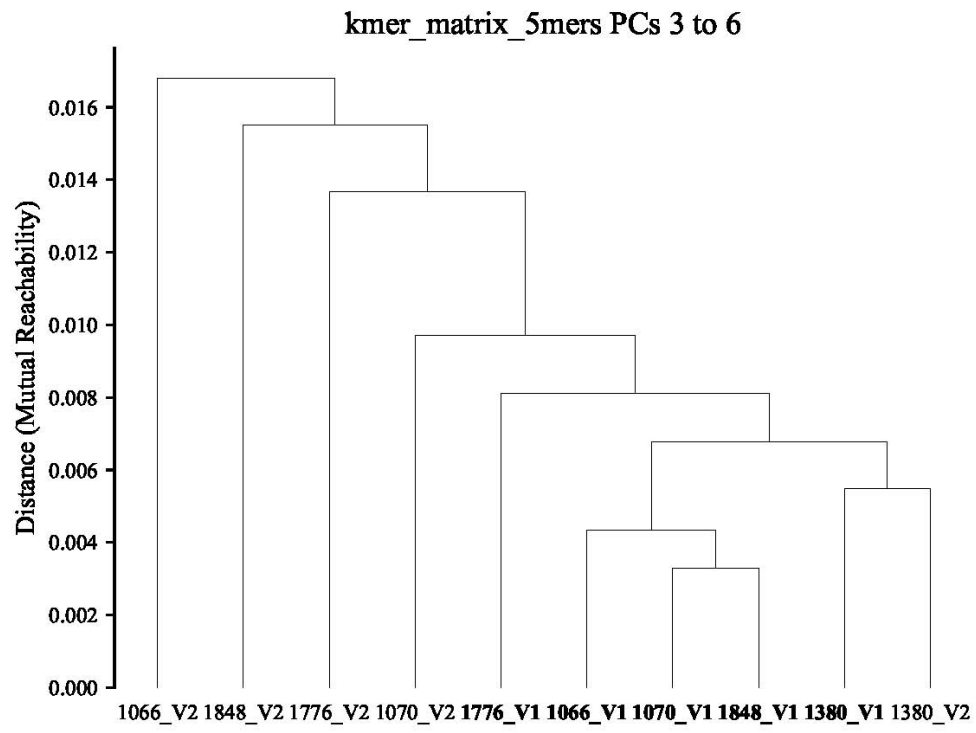


Figure 18