

36-315: Statistical Graphics and Visualization

Handout 22

Date: April 14, 2003

Finding the variables that explain outliers

Useful functions for examining outliers: `locate.bbox`, `identify`, `points`

Example of explaining outliers with maps:

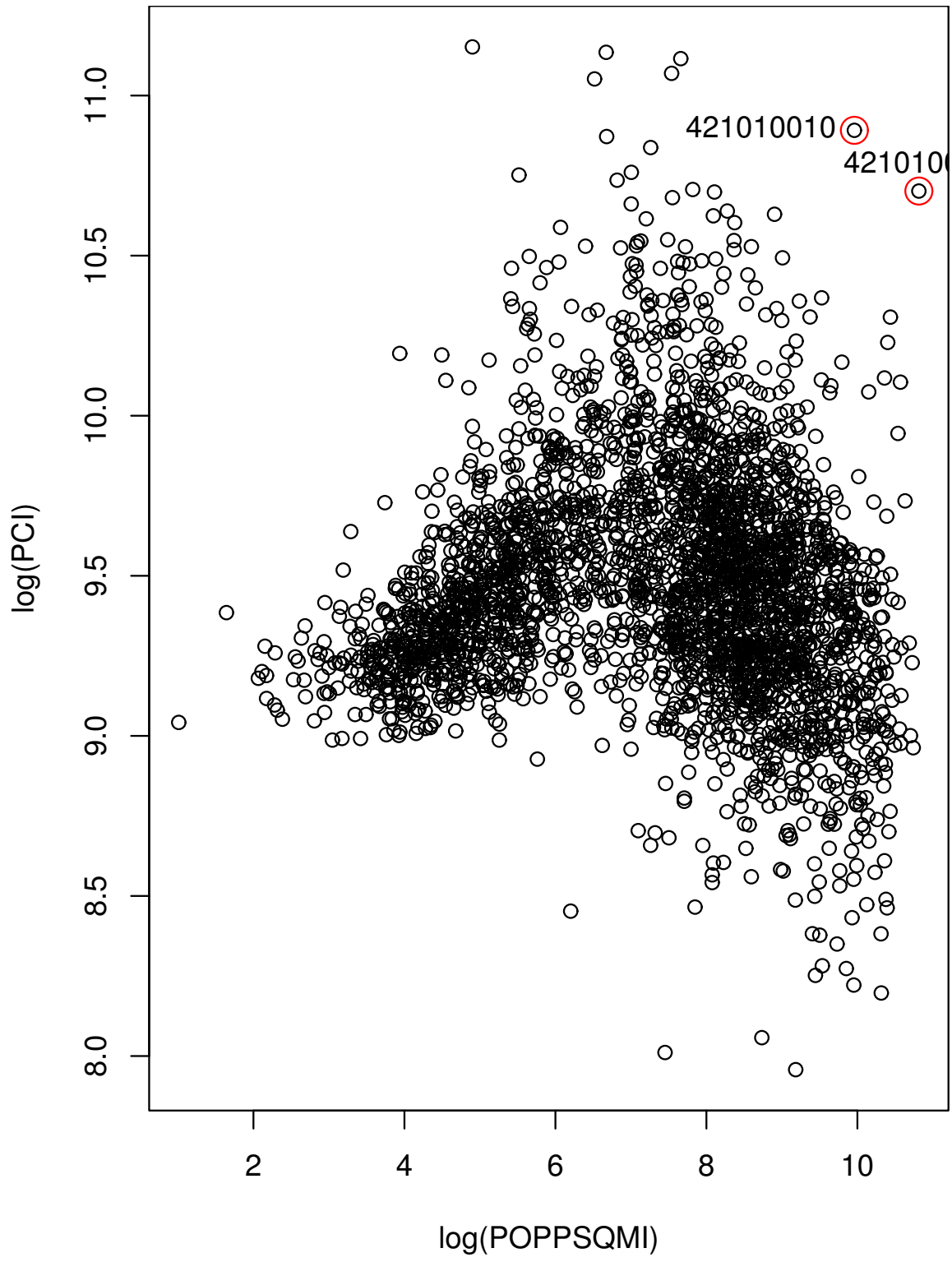
<http://www.mimas.ac.uk/argus/Tutorials/CartoViz/PopViz/Figs/C20Outliers.html>

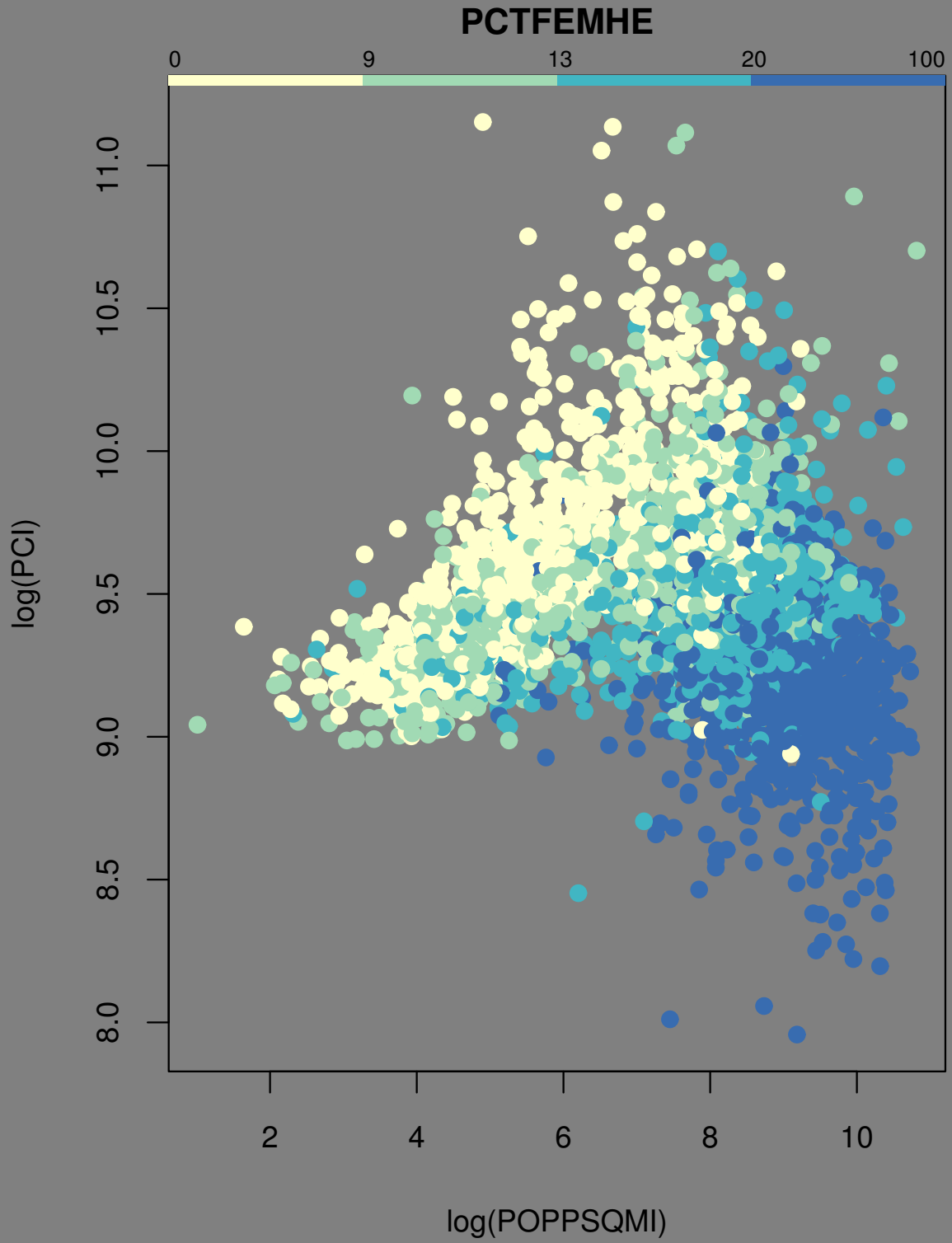
Hypervariate—Having four or more dimensions.

Projection—A method for reducing a dataset of many dimensions to fewer (usually two) dimensions, by taking a weighted sum. Geometrically, it is like taking a picture of a high-dimensional point cloud onto two-dimensional film, where the dataset is rotated and then all dimensions beyond the first two are dropped.

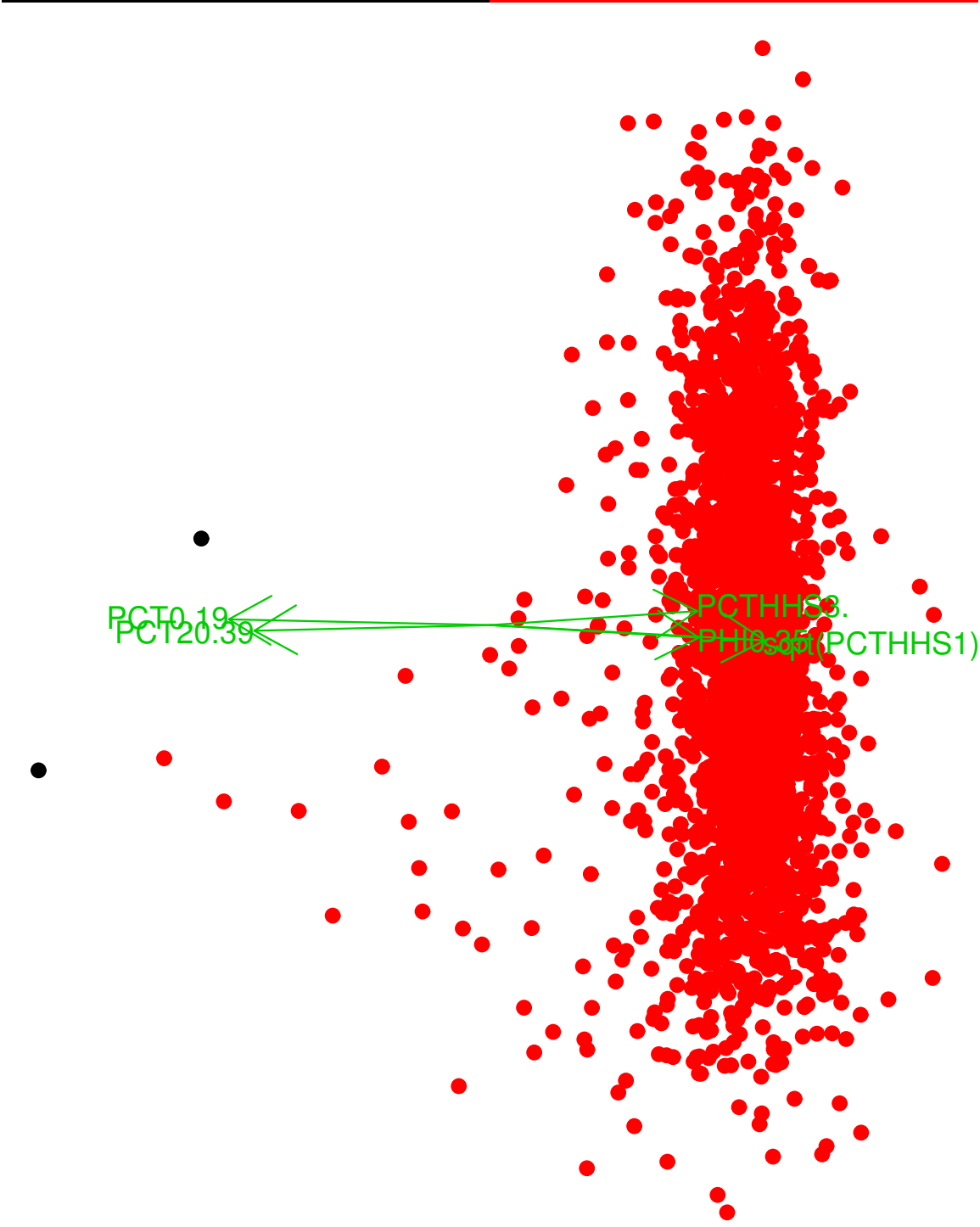
Discriminative projection—A projection which tries to separate data groups as much as possible. Typically we assume the groups have a normal distribution, so “separation” amounts to the distance between the group means, divided by the group standard deviations. This uniquely defines the projection.

To find the variables that make a data point unusual, make a discriminative projection which separates the point from the rest of the data. For explaining an outlier, first exclude the variables that were originally used to find the outlier.





Can we explain the outliers better?

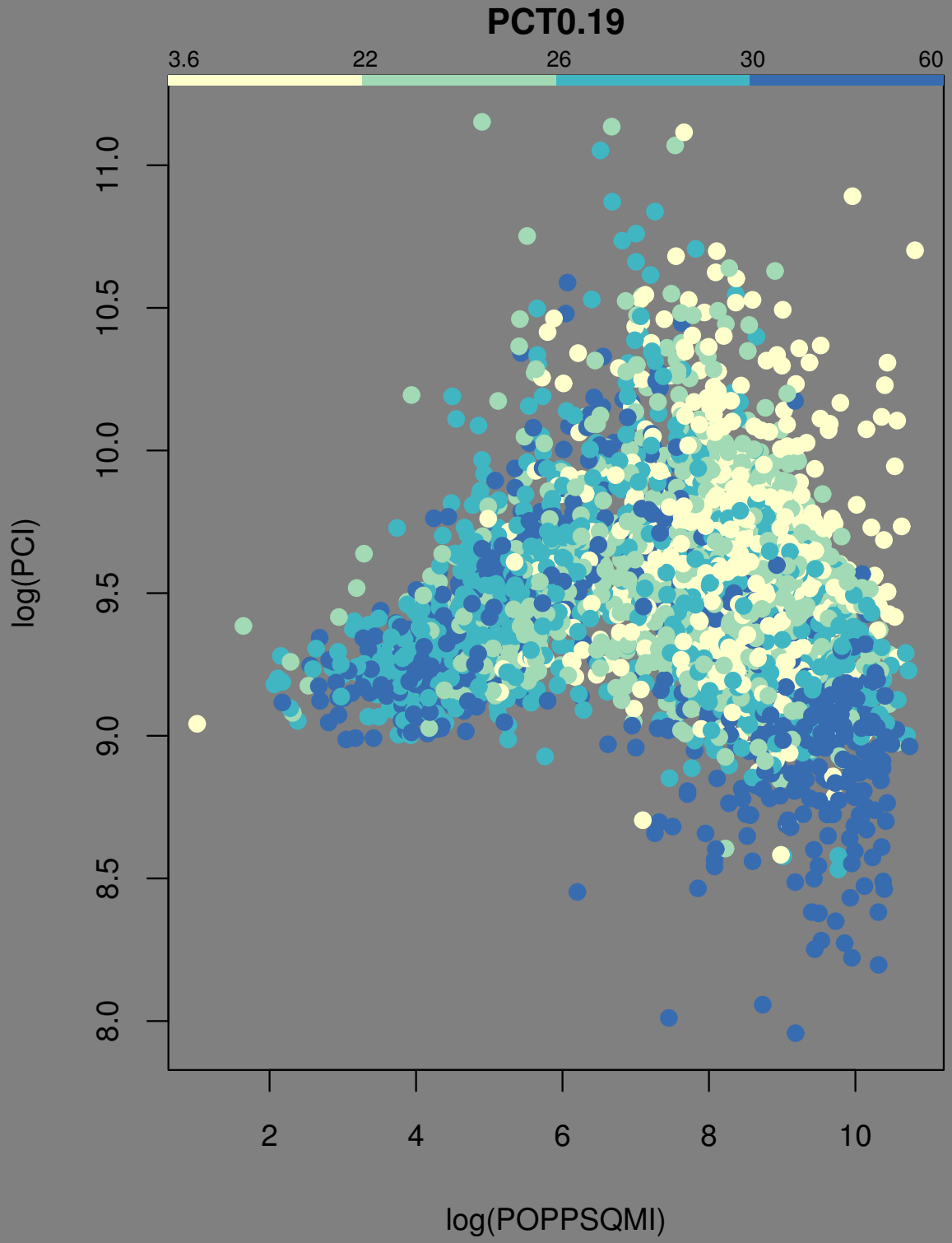


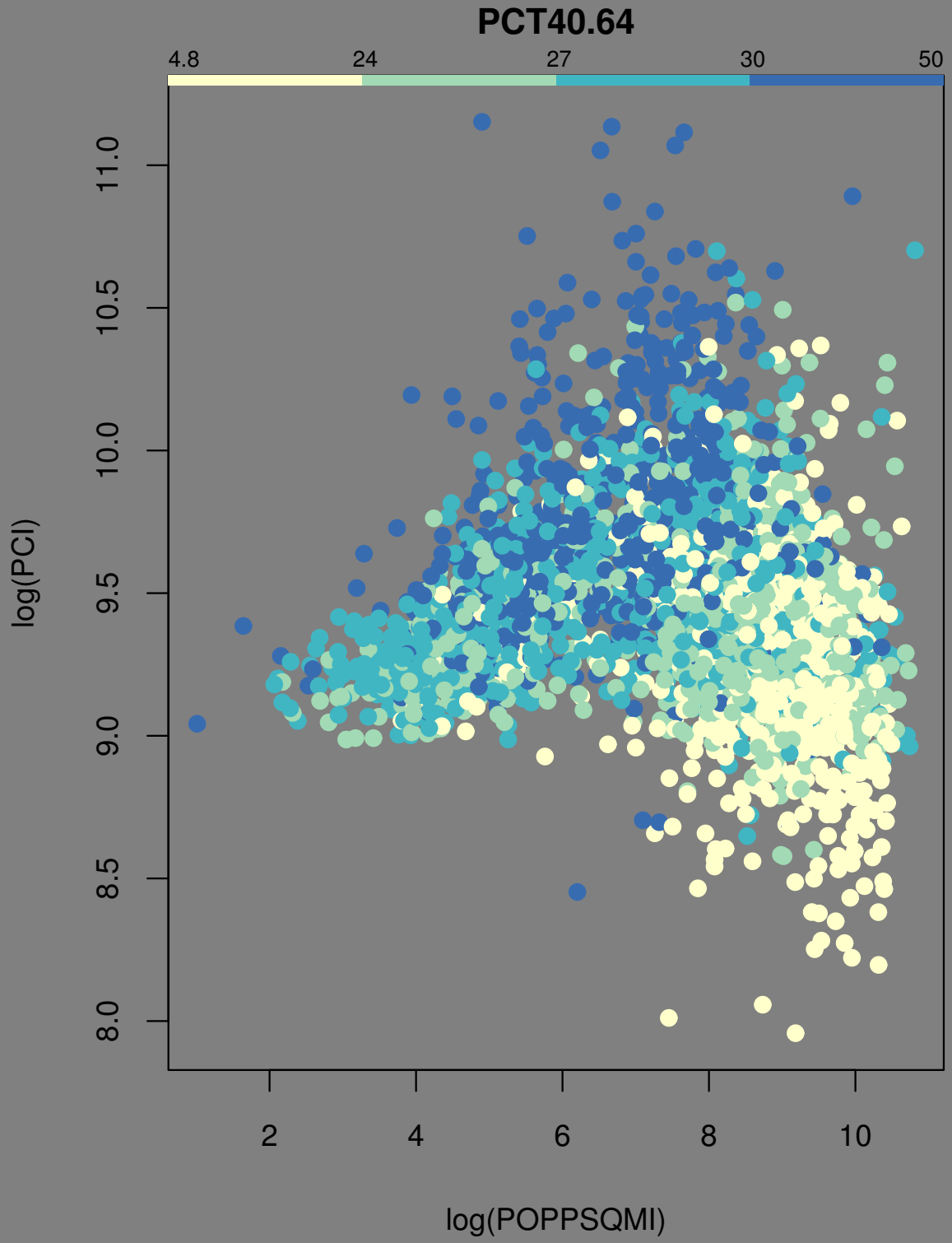
Discriminative projection

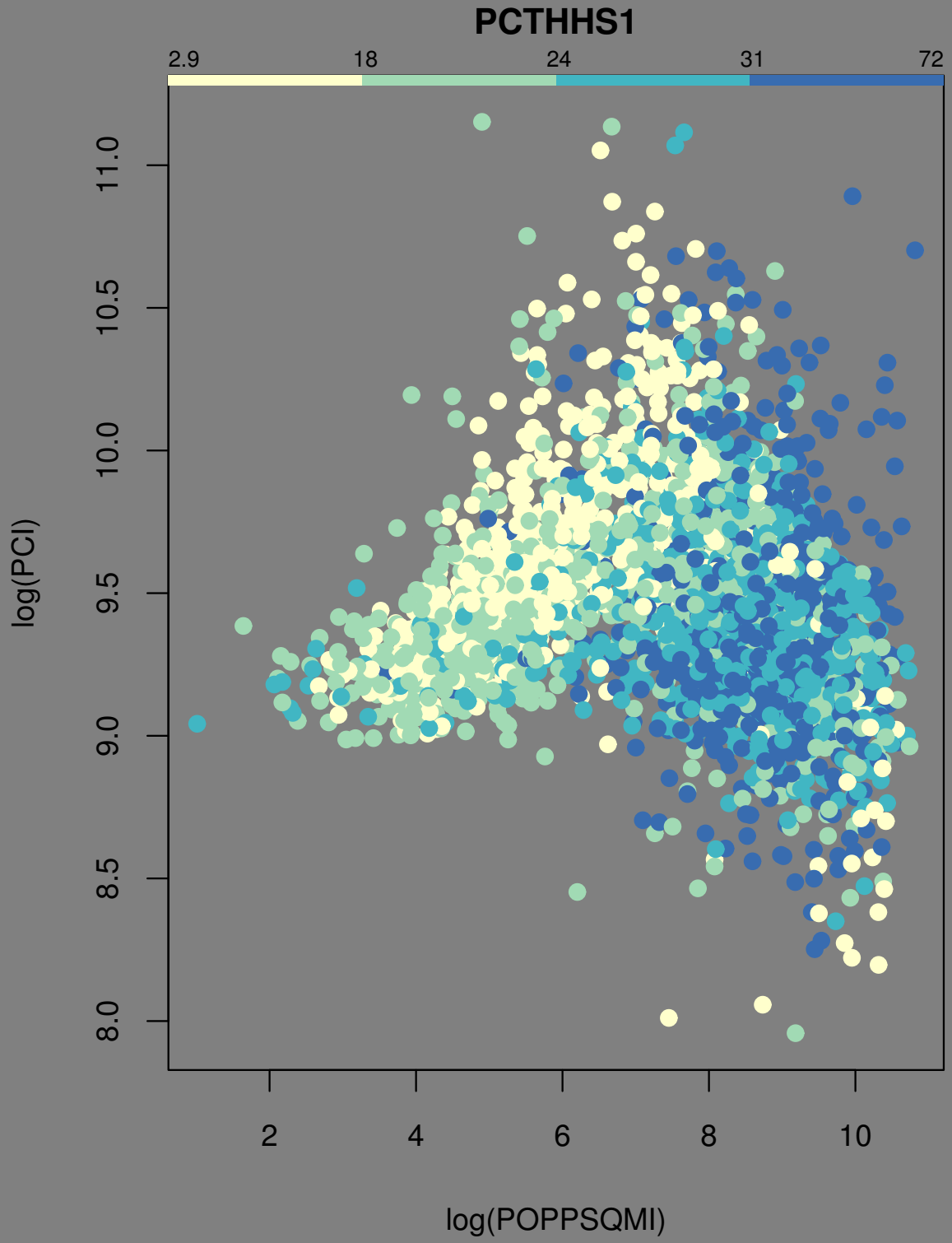
Variable weights used in the projection:

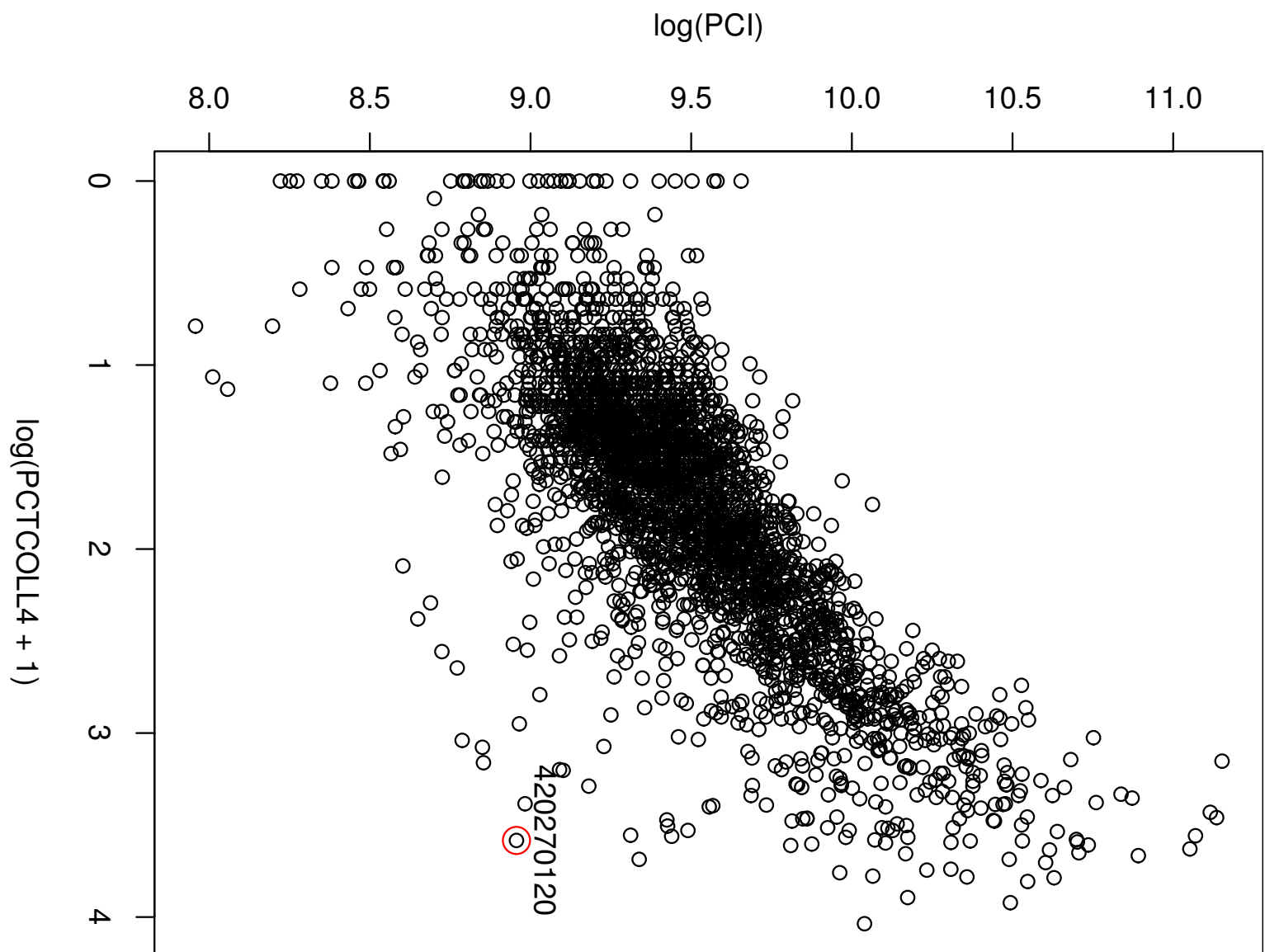
PCT0.19	-0.354221126
PCT20.39	-0.319680787
PCT40.64	-0.230885360
PHI0.75	-0.165780582
PCTFAMHH	-0.149873844
sqrt(PCTNFHHS)	-0.130067171
...	
sqrt(PCTHHS5.)	0.156165703
sqrt(PHI0.15)	0.162485771
PCT18.24	0.228470905
PCTHHS2	0.240949045
PCTHHS3.	0.280324301
PHI0.35	0.282467035
sqrt(PCTHHS1)	0.371919142

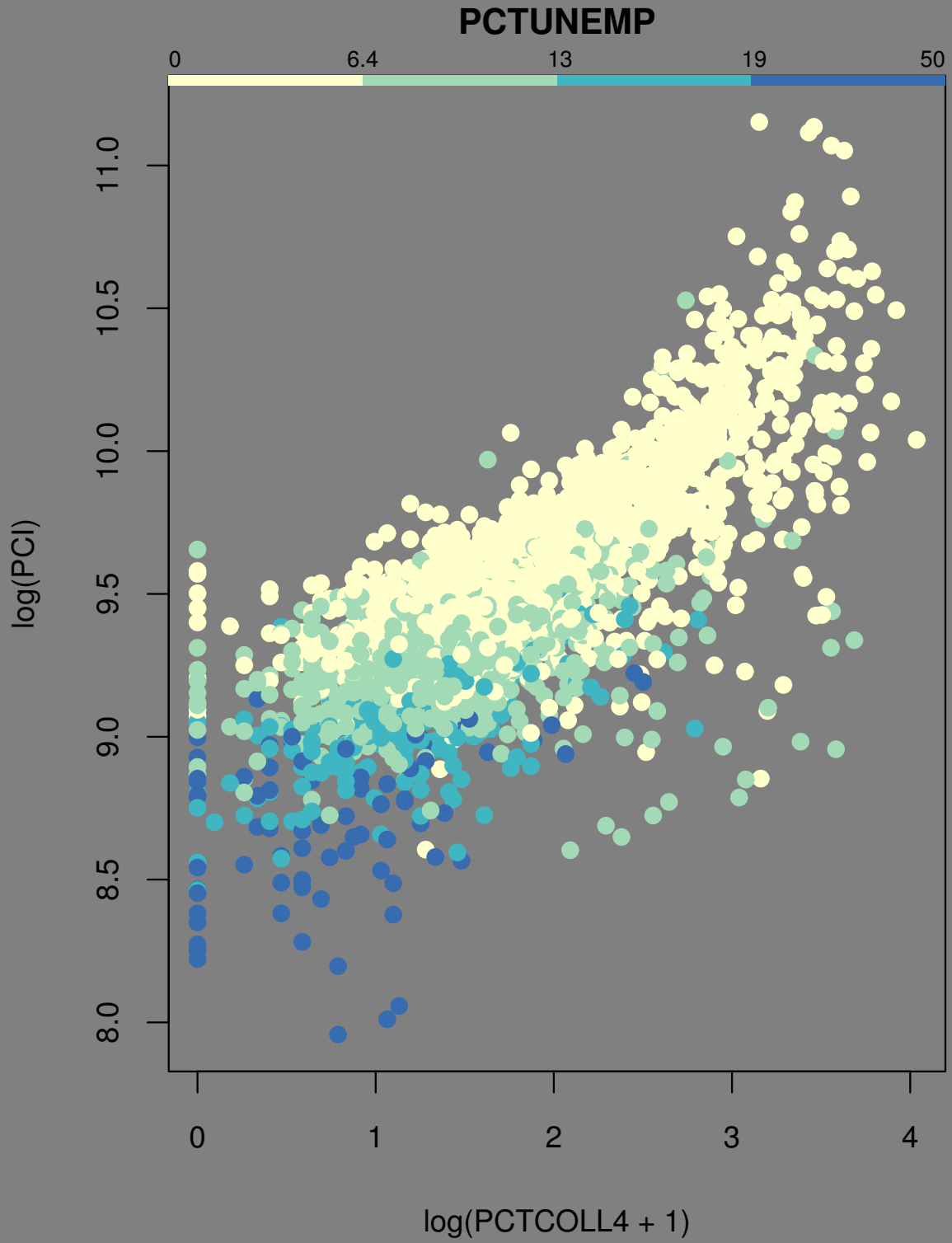
Largest weights (in magnitude) best explain the outliers







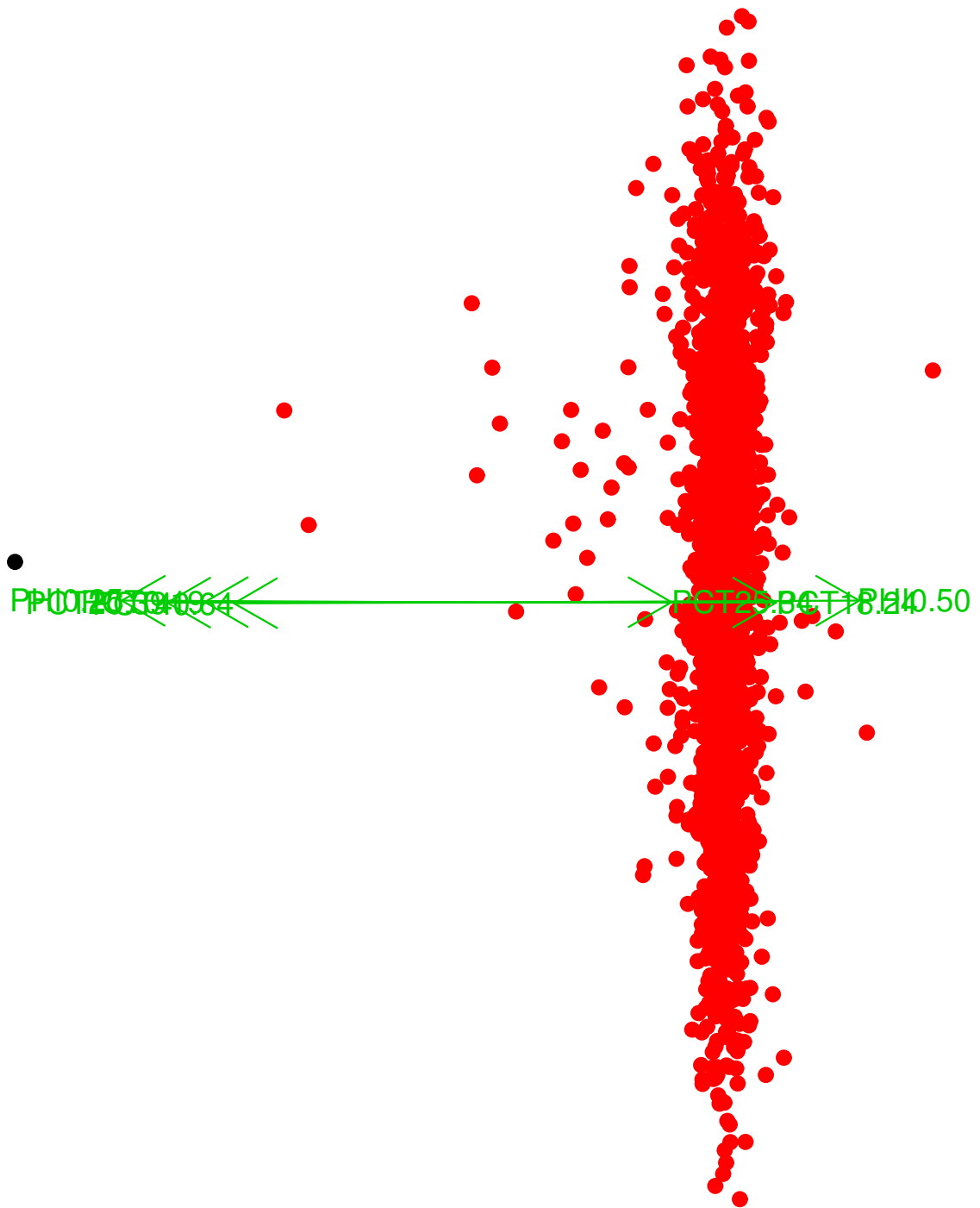




Can we explain the outliers better?

420270120

NOT 420270120



Discriminative projection

Variable weights used in the projection:

PHI0.25	-4.089384e-01
PCT20.39	-3.562167e-01
PCT0.19	-3.125527e-01
PCT40.64	-2.785998e-01
PHI35.50	-1.634679e-01
PHI0.35	-1.021199e-01
PHI25.35	-1.018871e-01
PCTHHS3.	-1.213641e-02
...	
PCT60.64	1.006973e-01
PCT5.9	1.044438e-01
PCT0.4	1.058806e-01
PCT45.54	1.480137e-01
PCT35.44	1.749425e-01
PCT25.34	2.294755e-01
PCT18.24	3.509350e-01
PHI0.50	4.474431e-01

Largest weights (in magnitude) best explain the outliers

