

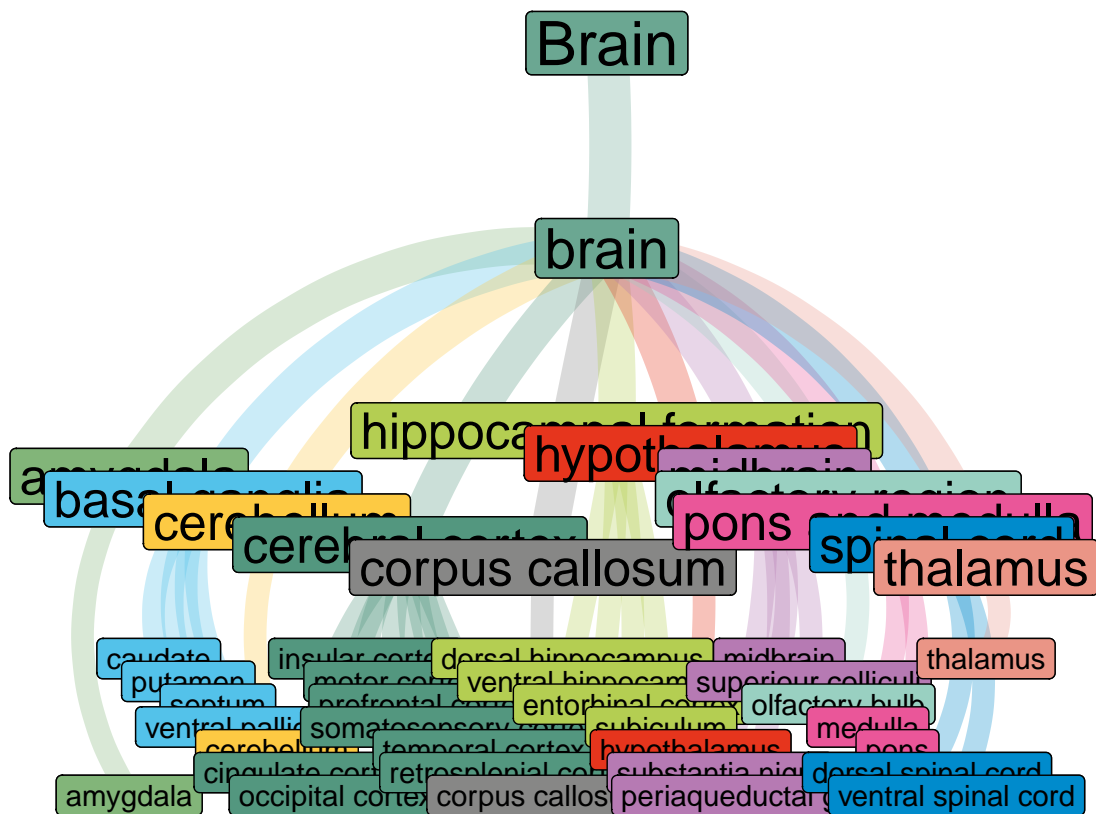
QC of pig data

Max Jonatan Karlsson

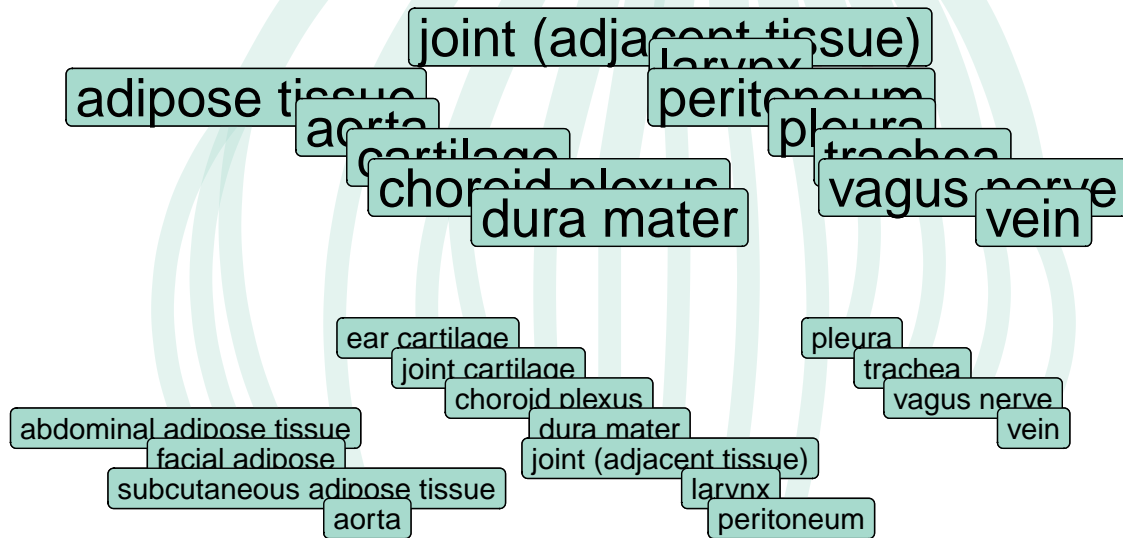
21 October 2019

Tissue grouping hierarchy

The following two plots describe the hierarchy of tissues for the Brain, and Adipose & soft tissue organs, and how they are grouped together. Similar plots for remaining organs are available in the file “Pig tissue groupings.pdf”. The levels in the plot are from top to bottom as follows: Organ, Consensus tissue, Region (Only for brain), Tissue.

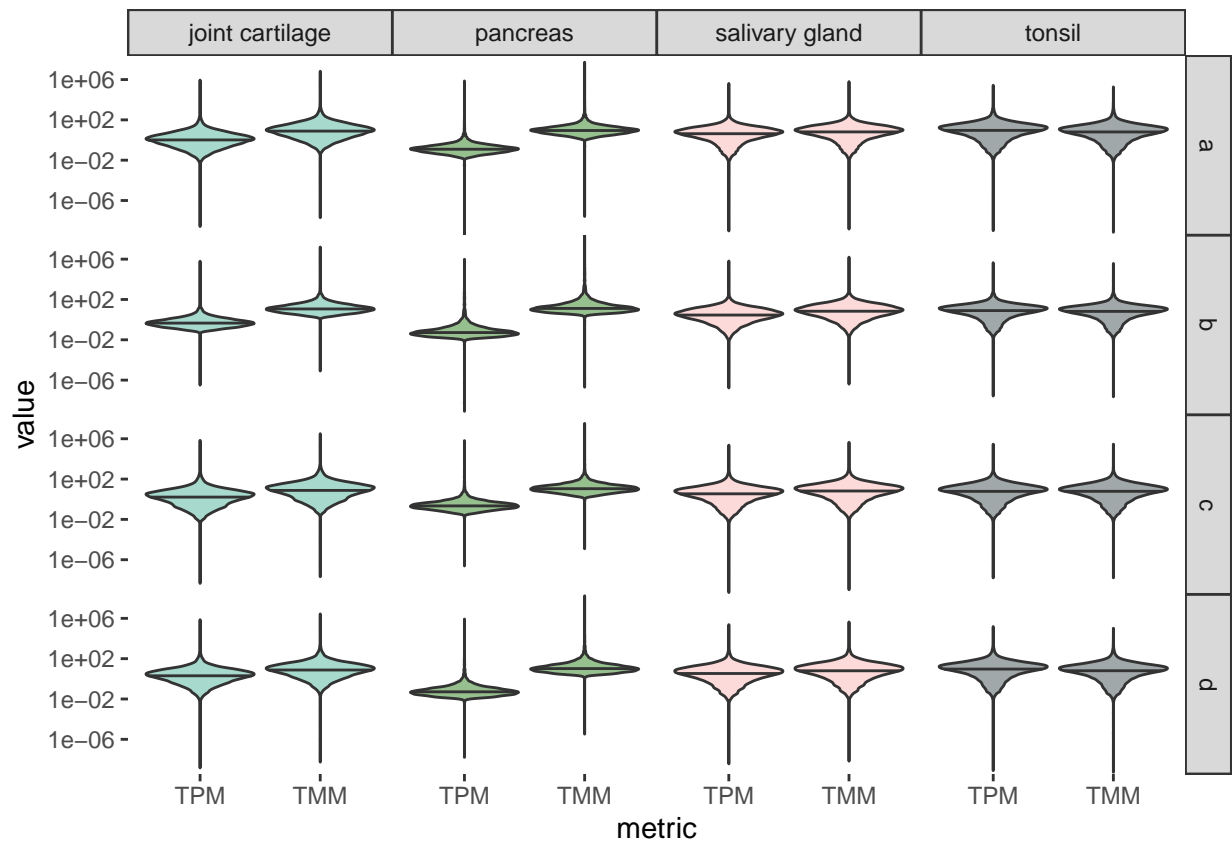


Adipose & soft tissue



Normalization

All samples were TMM normalized together with a median sample (median expression for all genes) as a reference distribution. Only protein coding genes are included. The plot below shows the distribution of expression values before and after TMM normalization for some selected tissues. Only transcripts with non-zero tpm have been included.

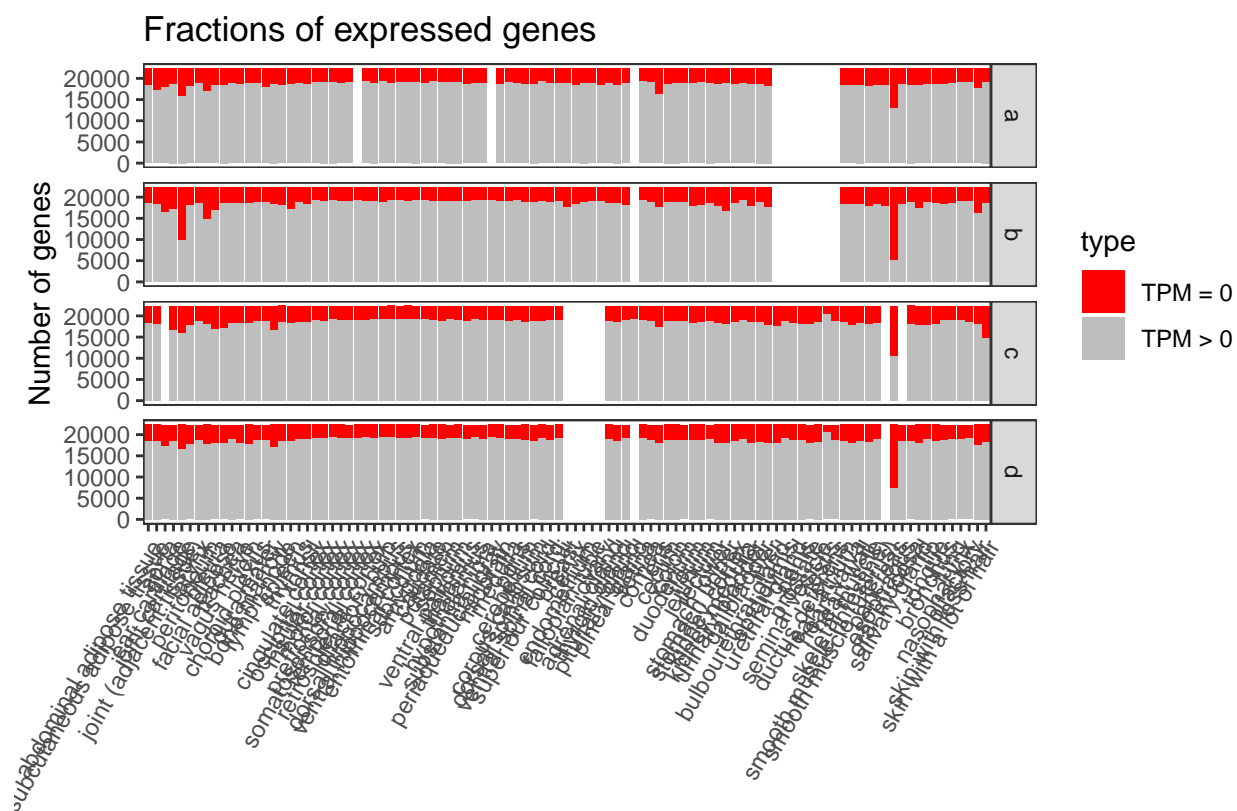


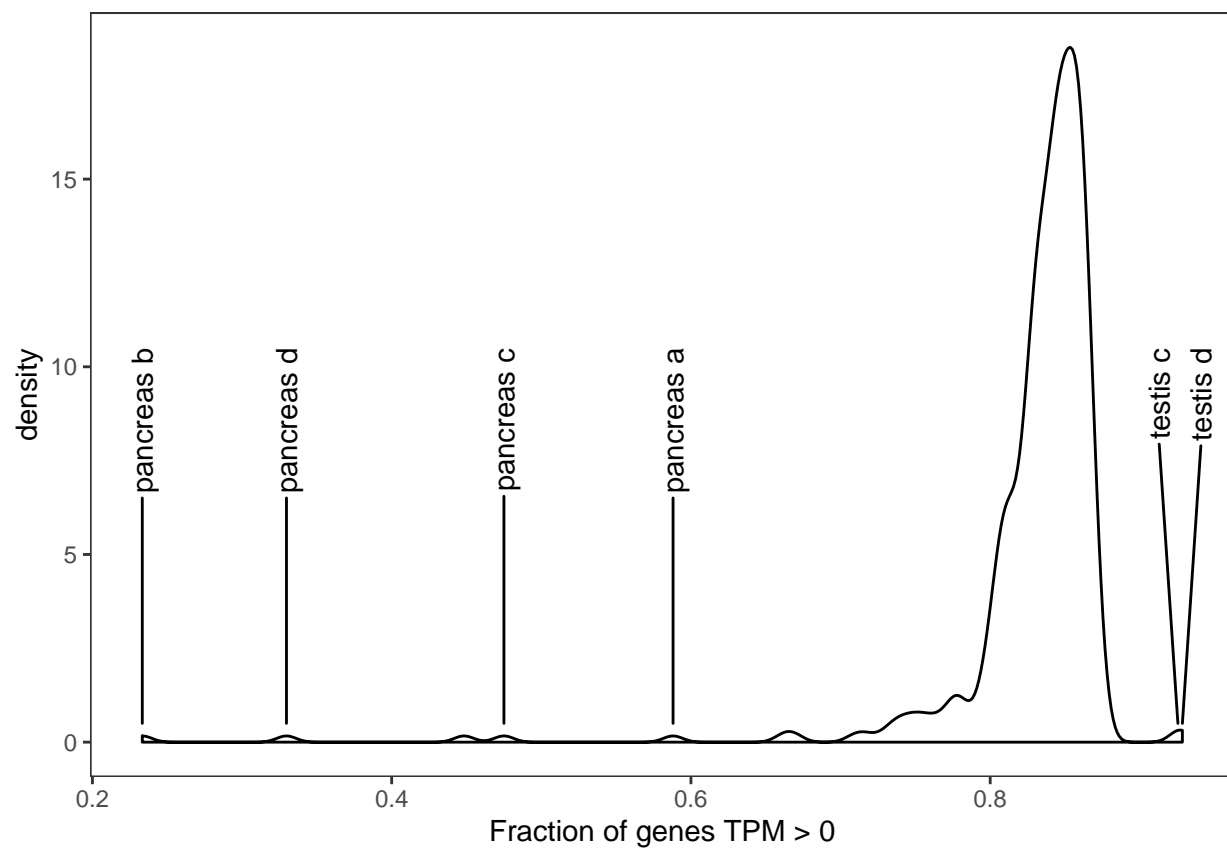
RNA and sequencing quality control

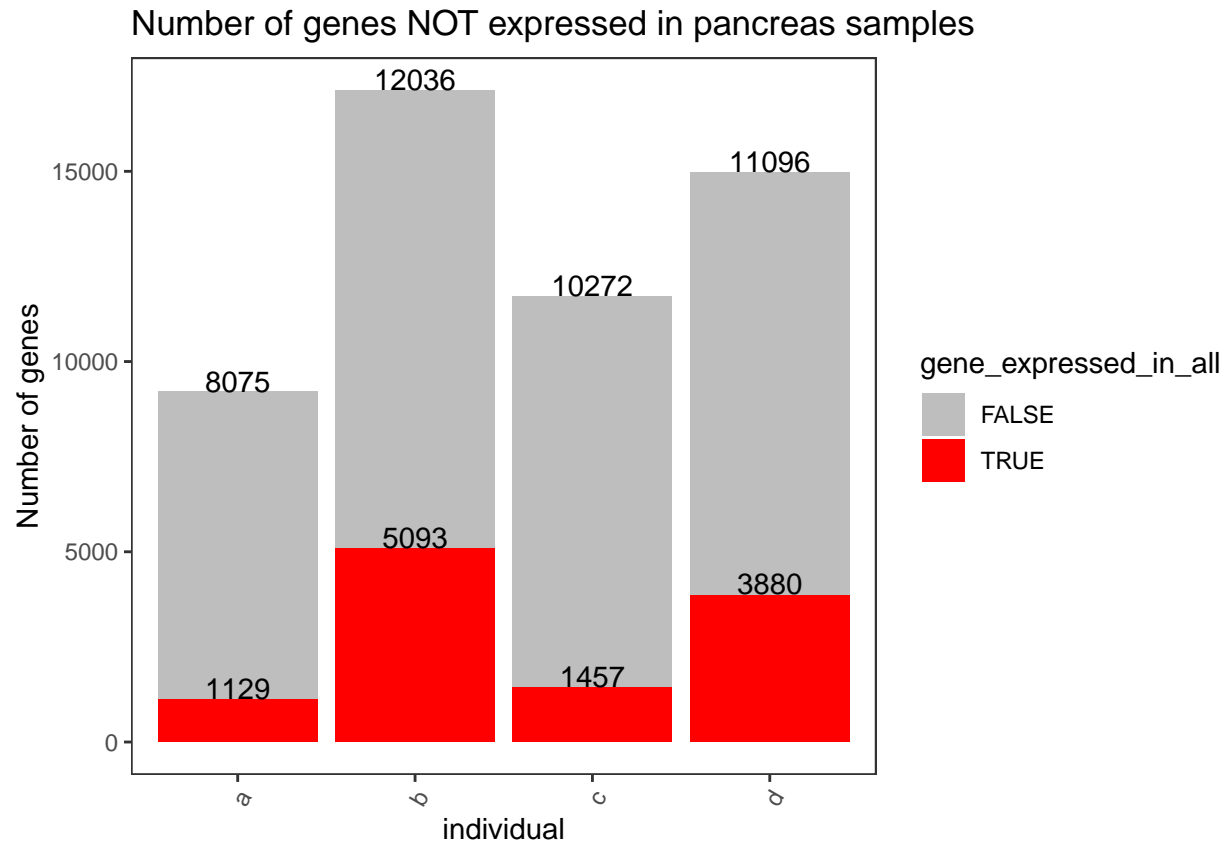


Fraction of genes with zero expression

The following plots describe the number of genes with $\text{TPM} = 0$ in each sample. We see that many genes are missing for pancreas samples - much more than for any other sample (except one joint adjacent tissue sample that will be removed further ahead). This could be due to autodigestion of the sample by the nucleases present in the pancreas. Pancreas samples b and d lack most genes, and will thus be removed.

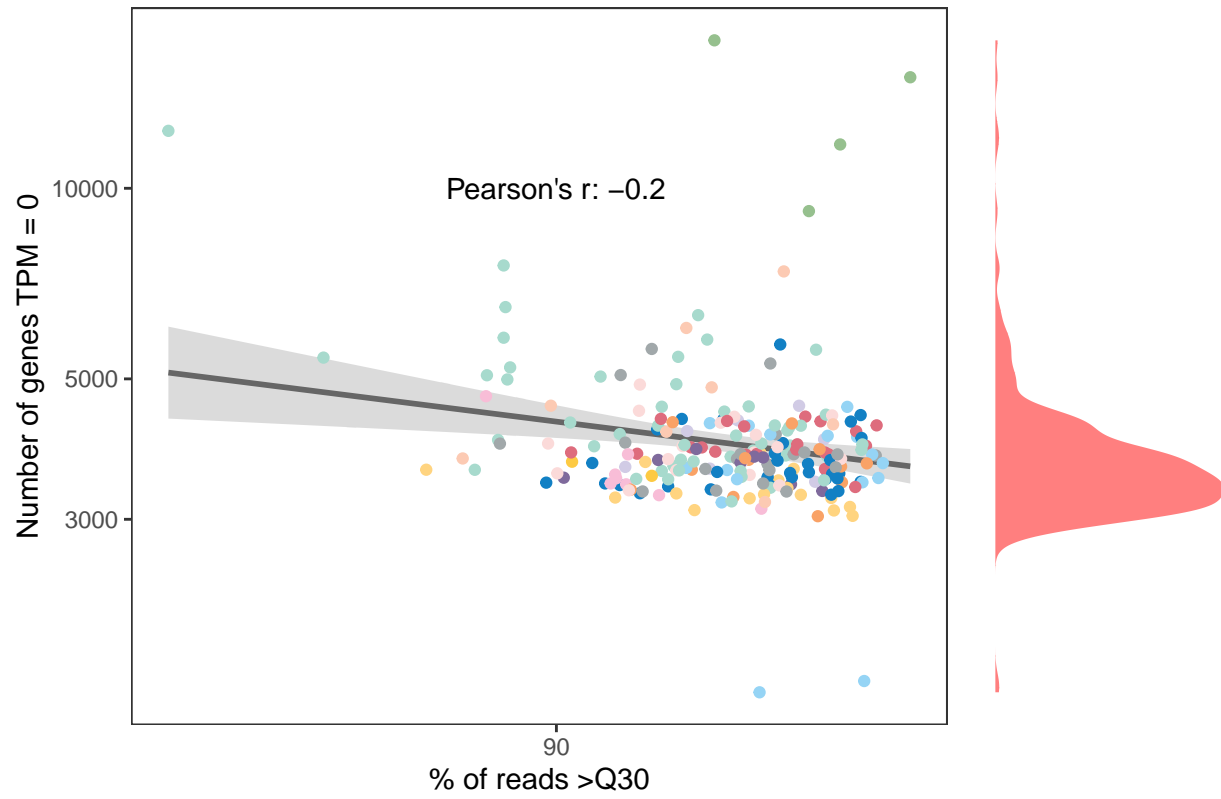


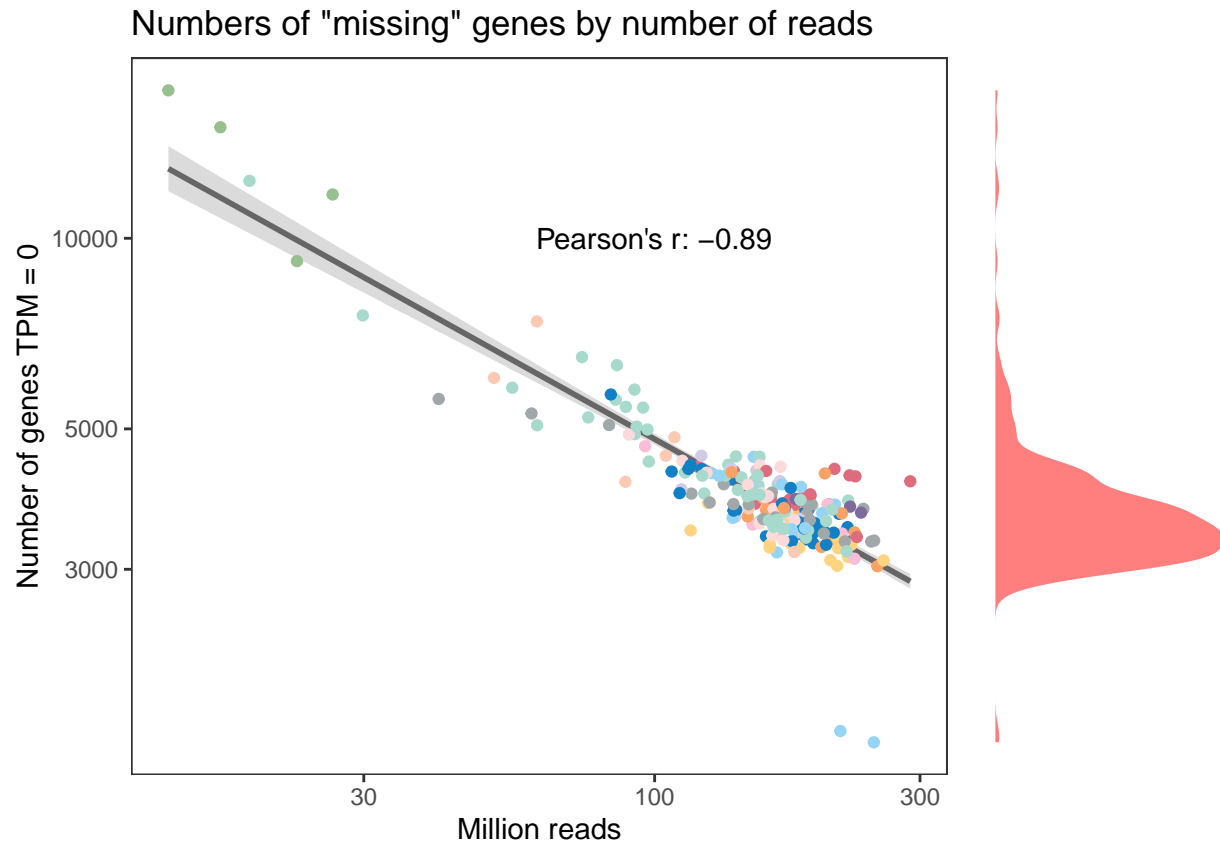




The following plots show how the number of “missing” genes in each sample is correlated with the sequencing quality metrics. We see that there is a strong correlation between number of reads and the number of genes with $TPM = 0$. This suggests that the issue with pancreas is simply sequencing depth.

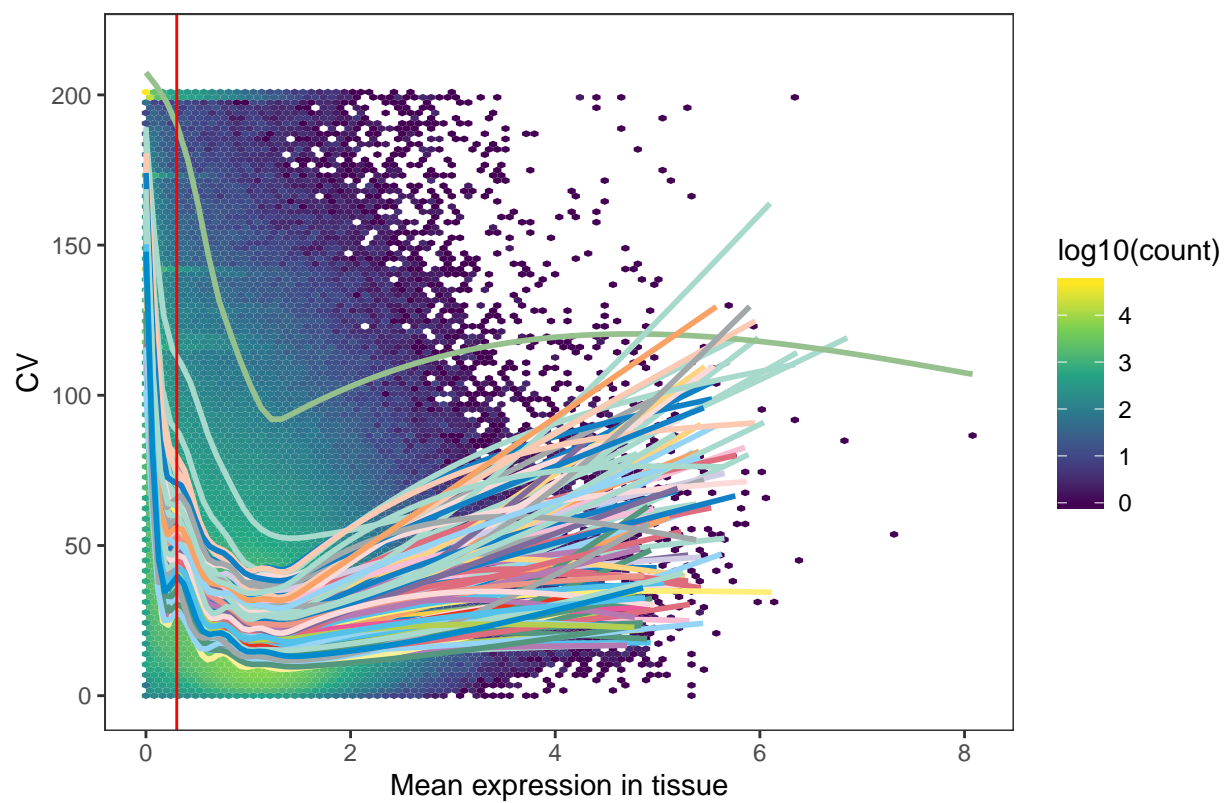
Numbers of "missing" genes by % of reads >Q30



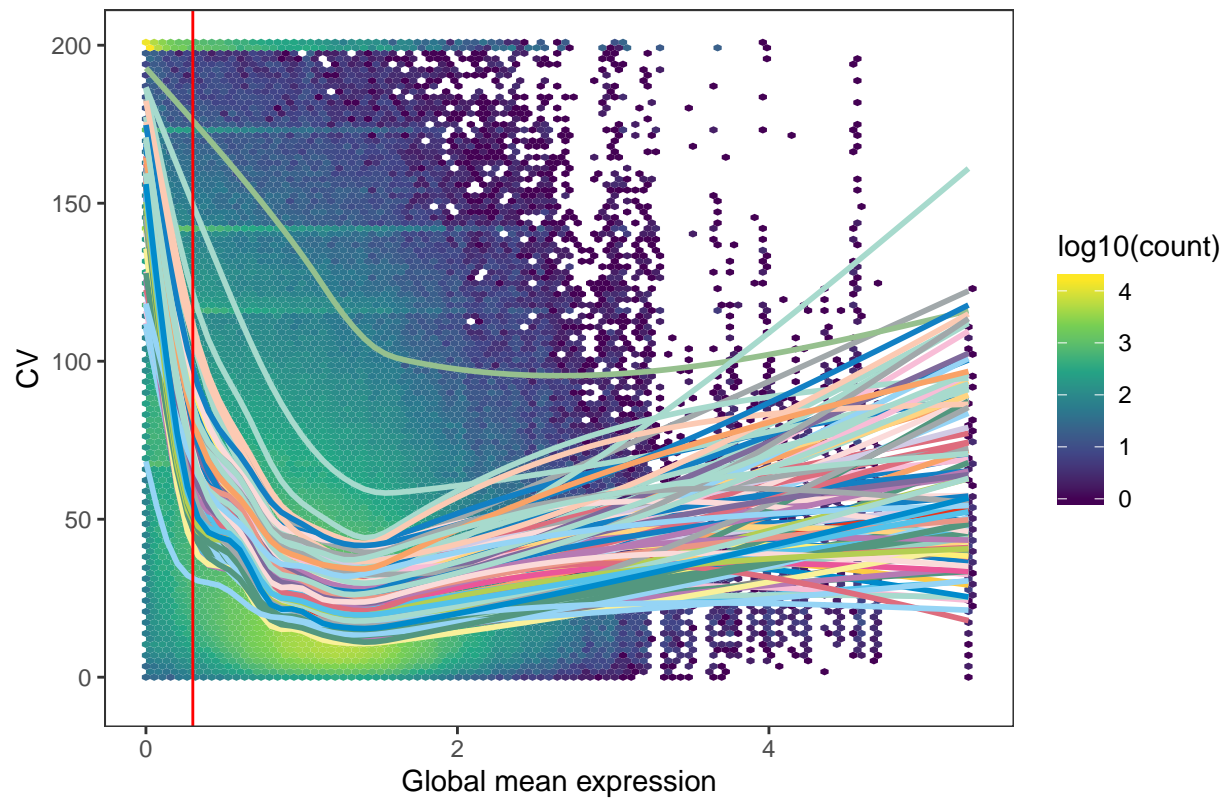


The following plots show the gene-wise variance between replicates by the genes' mean expression. We see pancreas and joint cartilage as outliers - again demonstrating that pancreas and joint cartilage b will need to be removed.

Intratissue gene variance by tissue mean expression



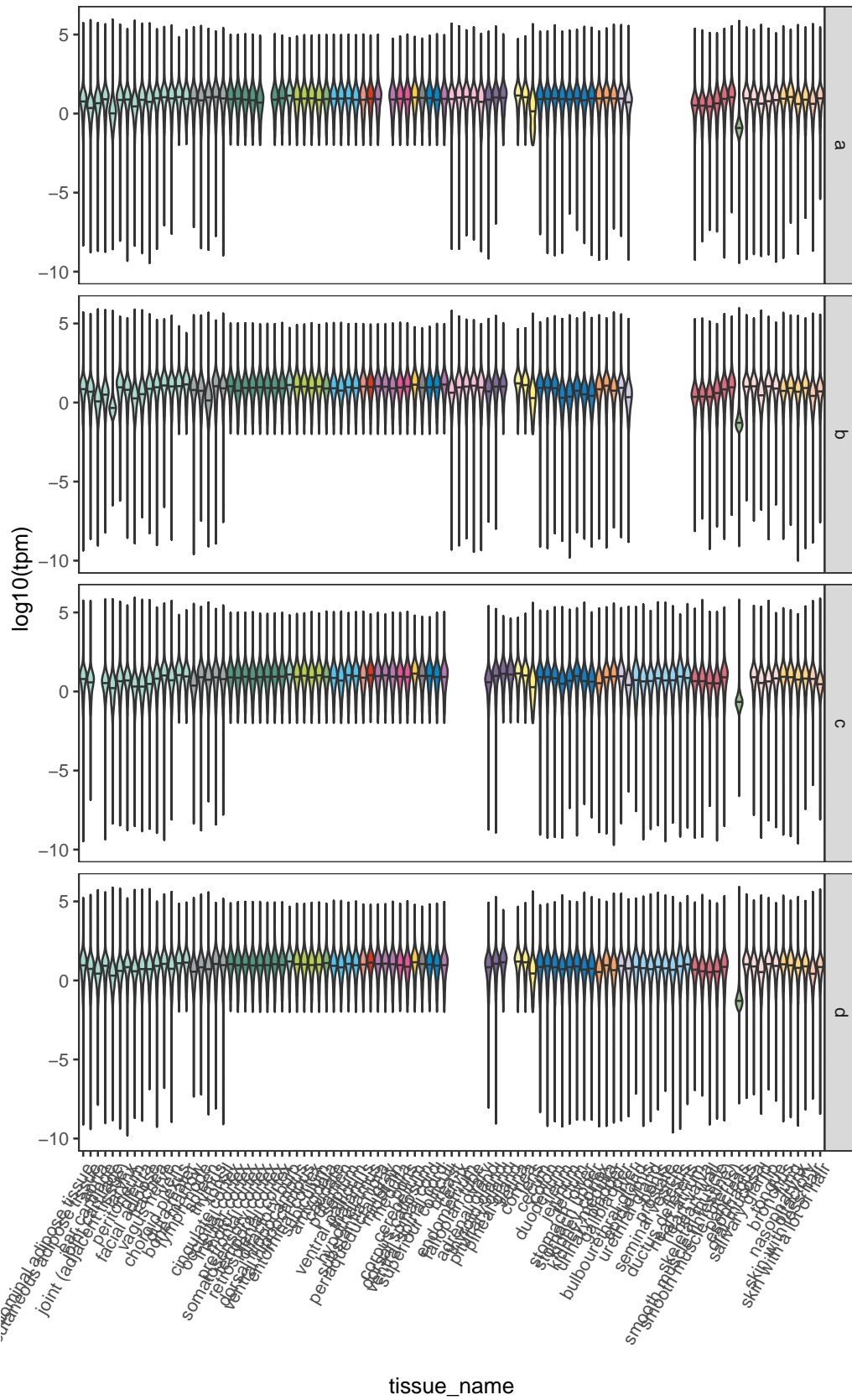
Intratissue gene variance by global mean expression



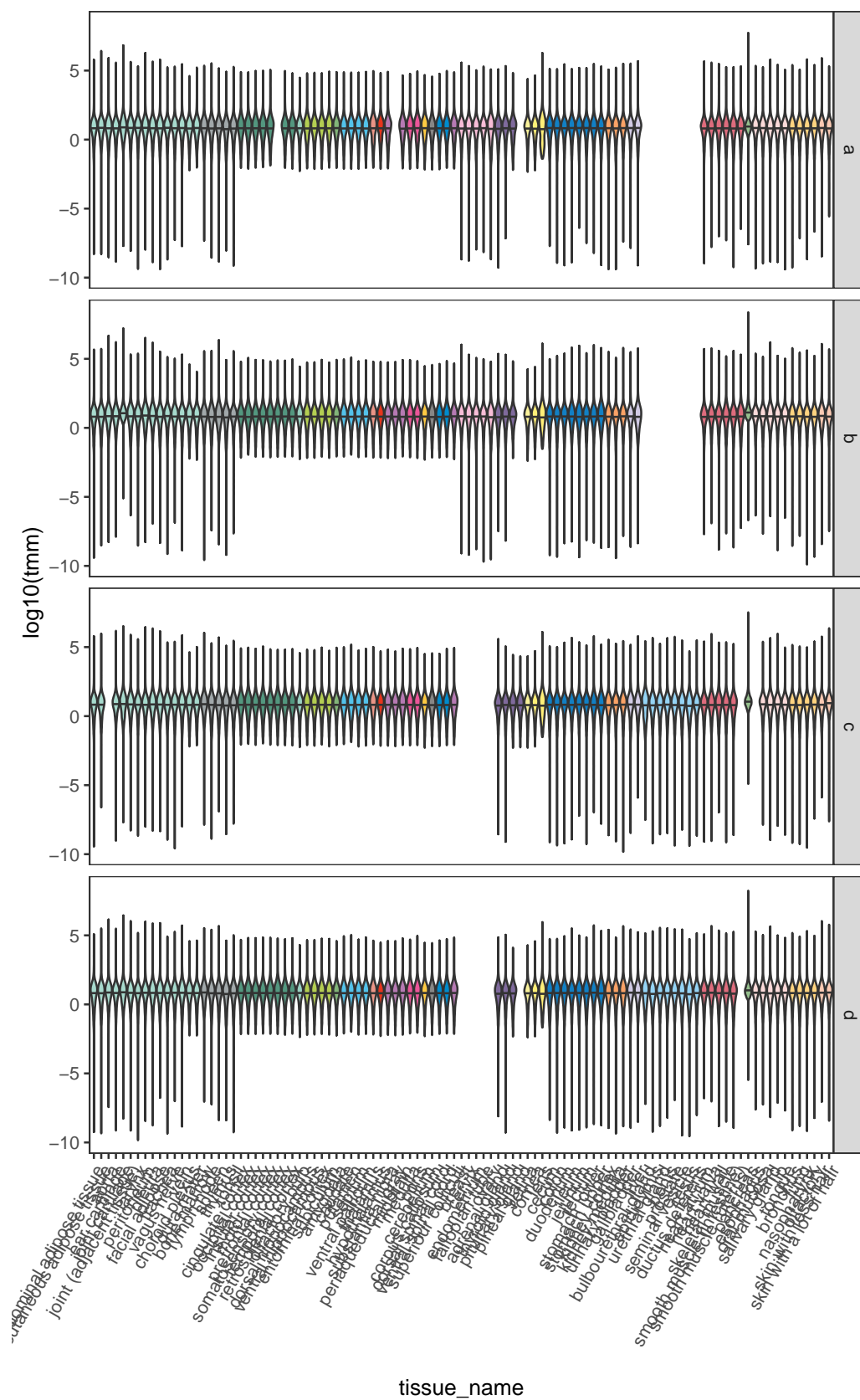
Transcript level distribution

Here we see the transcript level distribution in the four individuals (a, b, c, d) first for TPM values and then TMM values. Only transcripts with non-zero tpm have been included.

Transcript level distribution – TPM

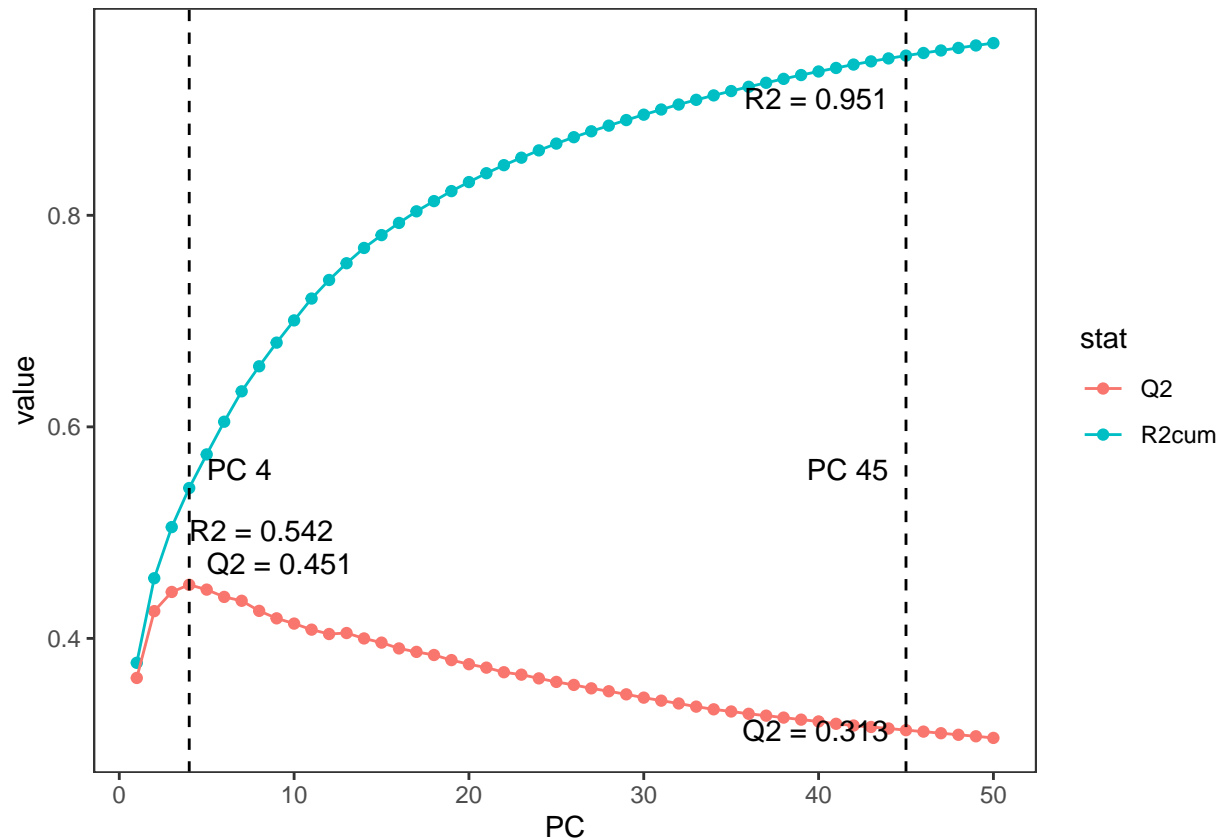


Gene level distribution – TMM normalized

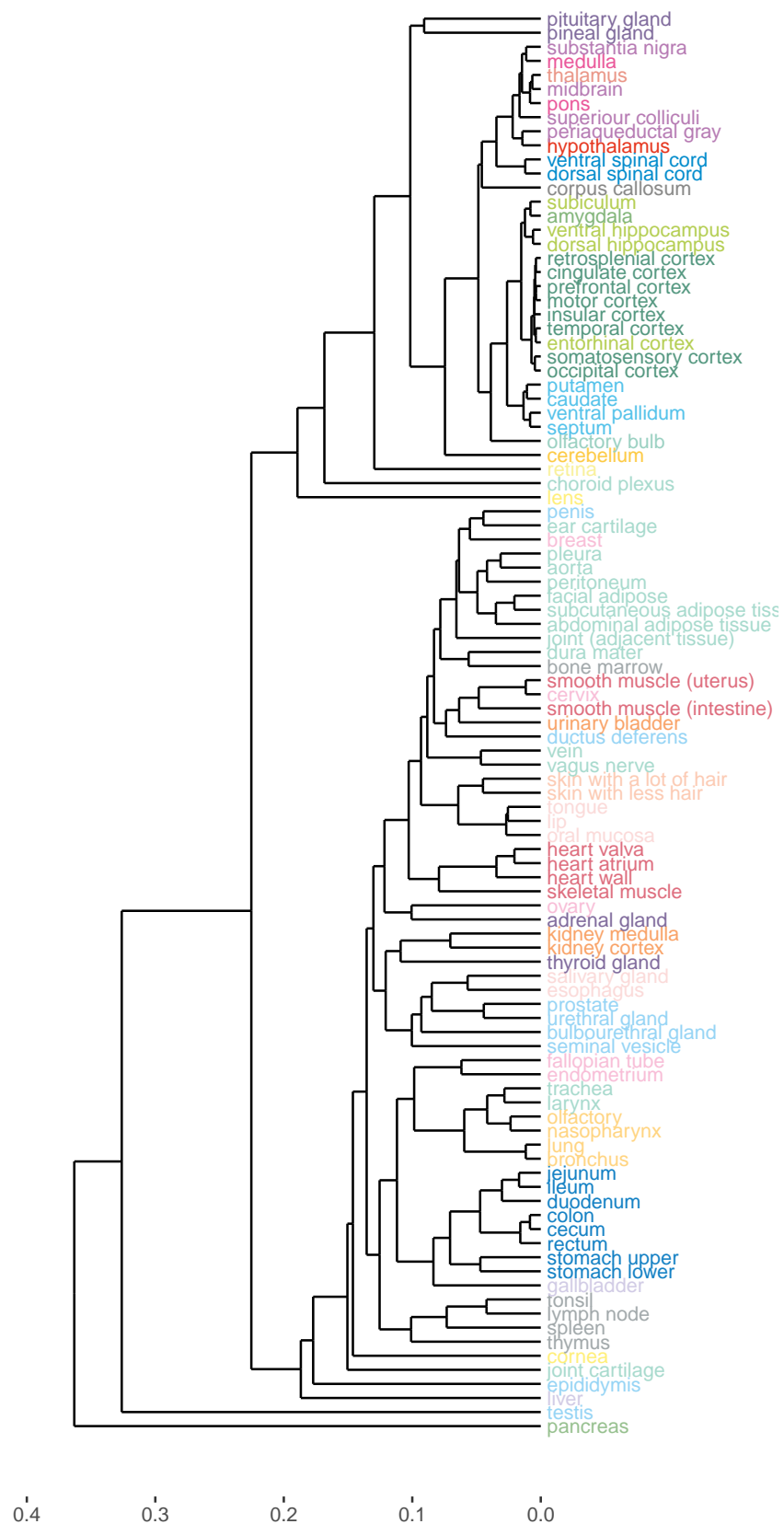


Tissue wise clustering

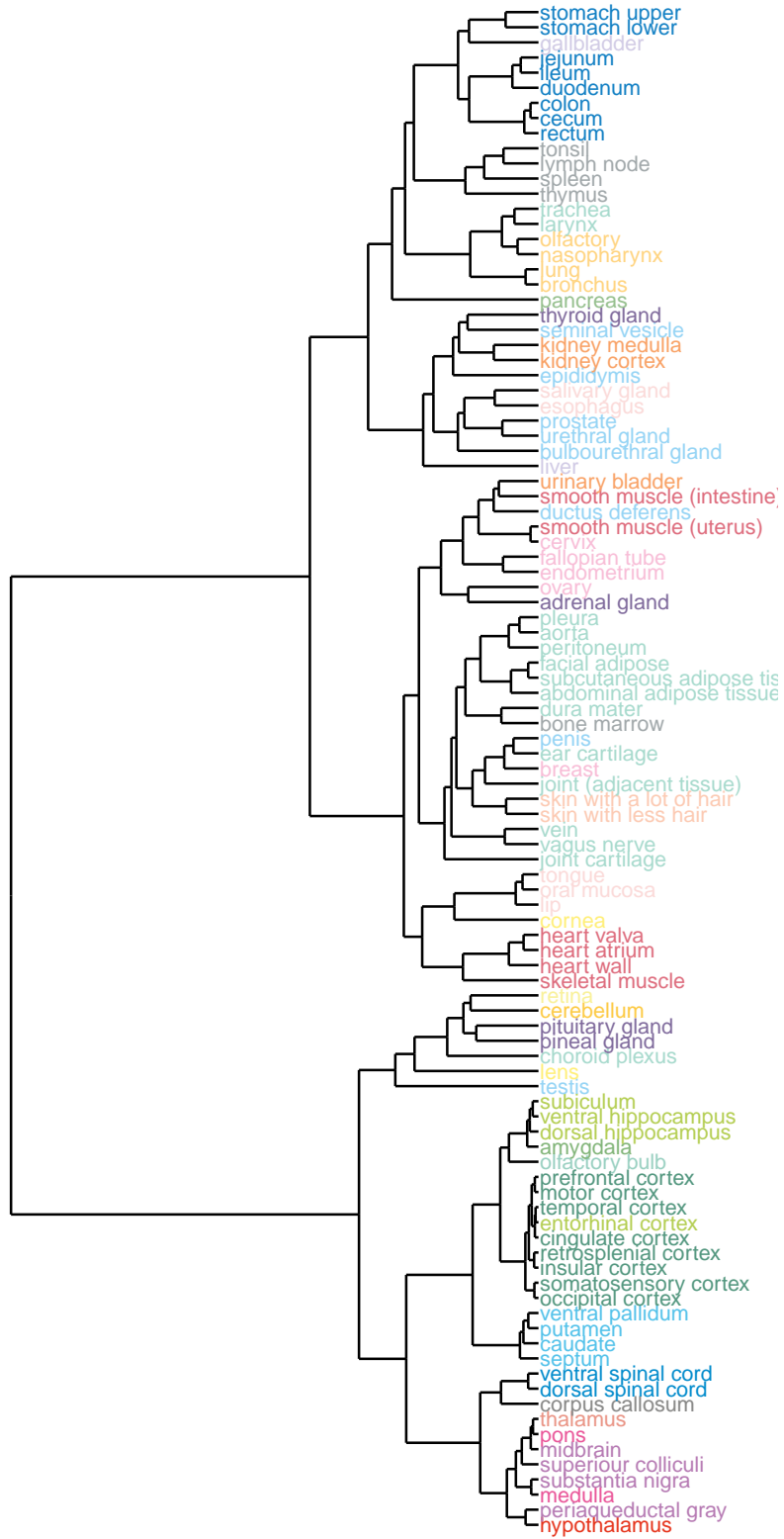
For each tissue, averages were taken for each gene and a PCA was performed on the normalized data to reduce the number of dimensions in the data. This helps to remove noise and keep only the structures in the data that separates the features of our tissues from each other. Below, we see a plot showing R2 (The fraction of explained variance) and Q2 (crossvalidated R2, can approximately be explained as the amount of information the PCA model contains). We see that 95% R2 is reached at component 39, and we will thus choose 39 components as our dimensionally reduced data. We choose this arbitrary cutoff as we want to remove noise but still keep a majority of the structure in the data.



Below we see a dendrogram built from 1 - Spearman's rho between the tissues using all genes.

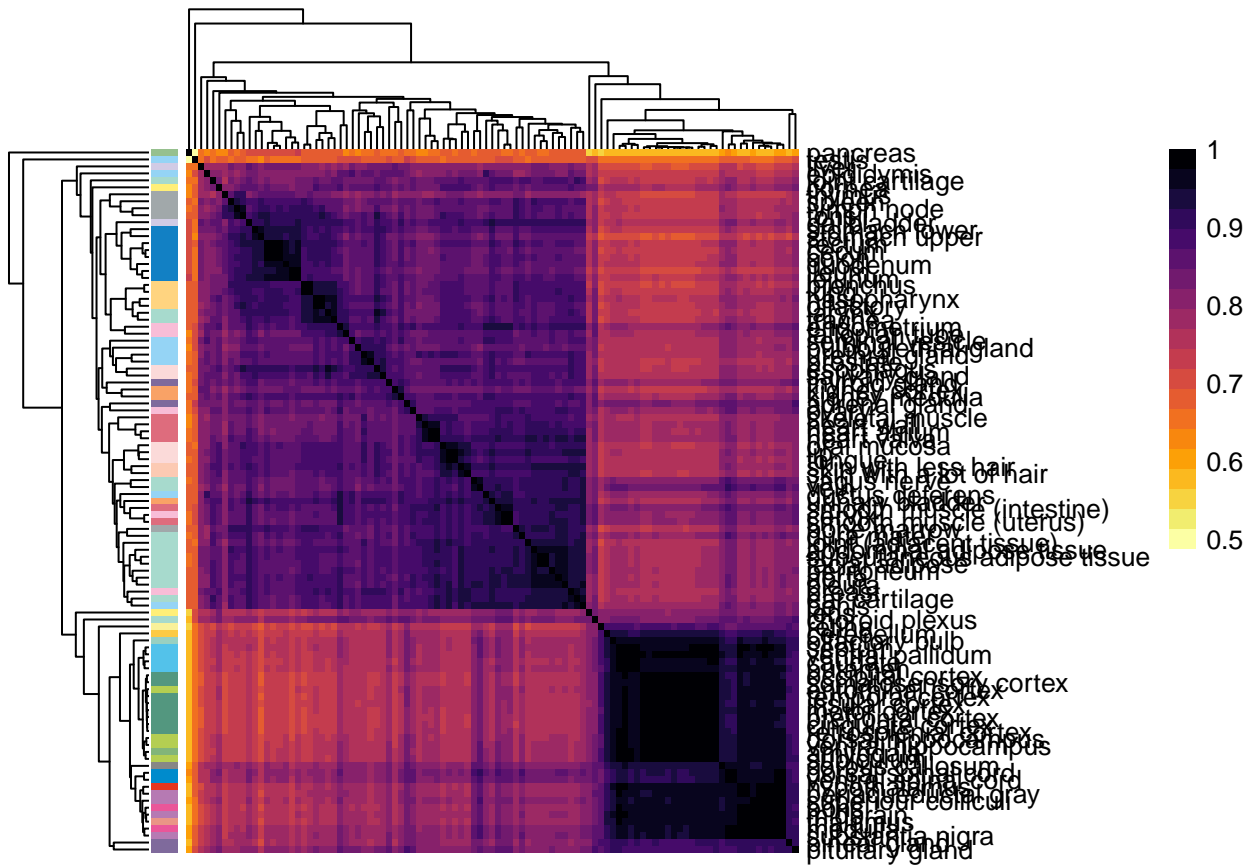


Below is an alternative (and complementary) clustering of tissues using Ward clustering of PCA scores in the 39 selected components. This type of clustering is based on minimizing the variance within each cluster and thus takes distances into account, which correlation does not.



150 100 50 0

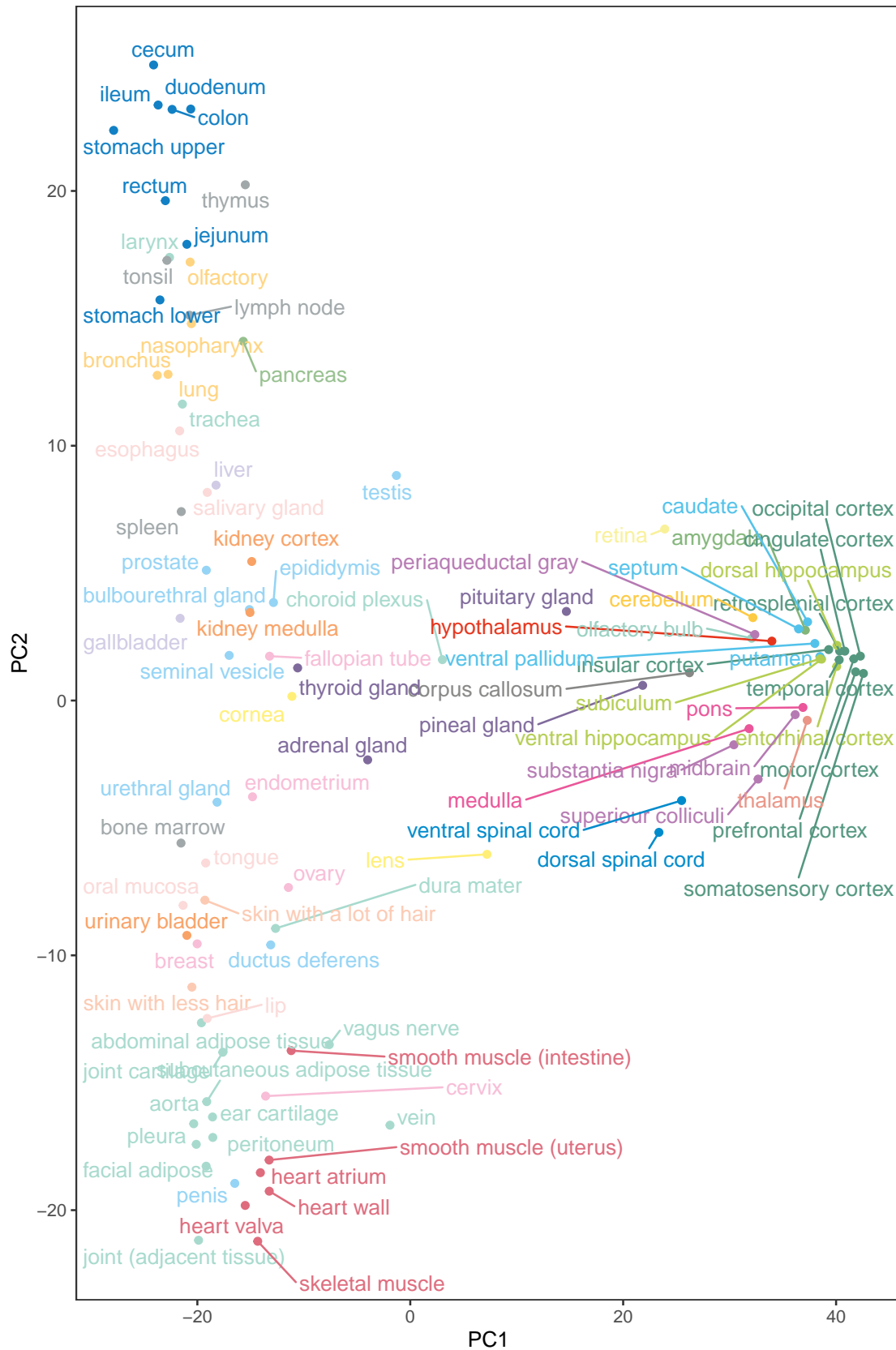
The following heatmap shows spearman correlation between tissues.



Below we see three plots showing the twodimensional projections of the data using three different strategies:

1. PCA

PCA



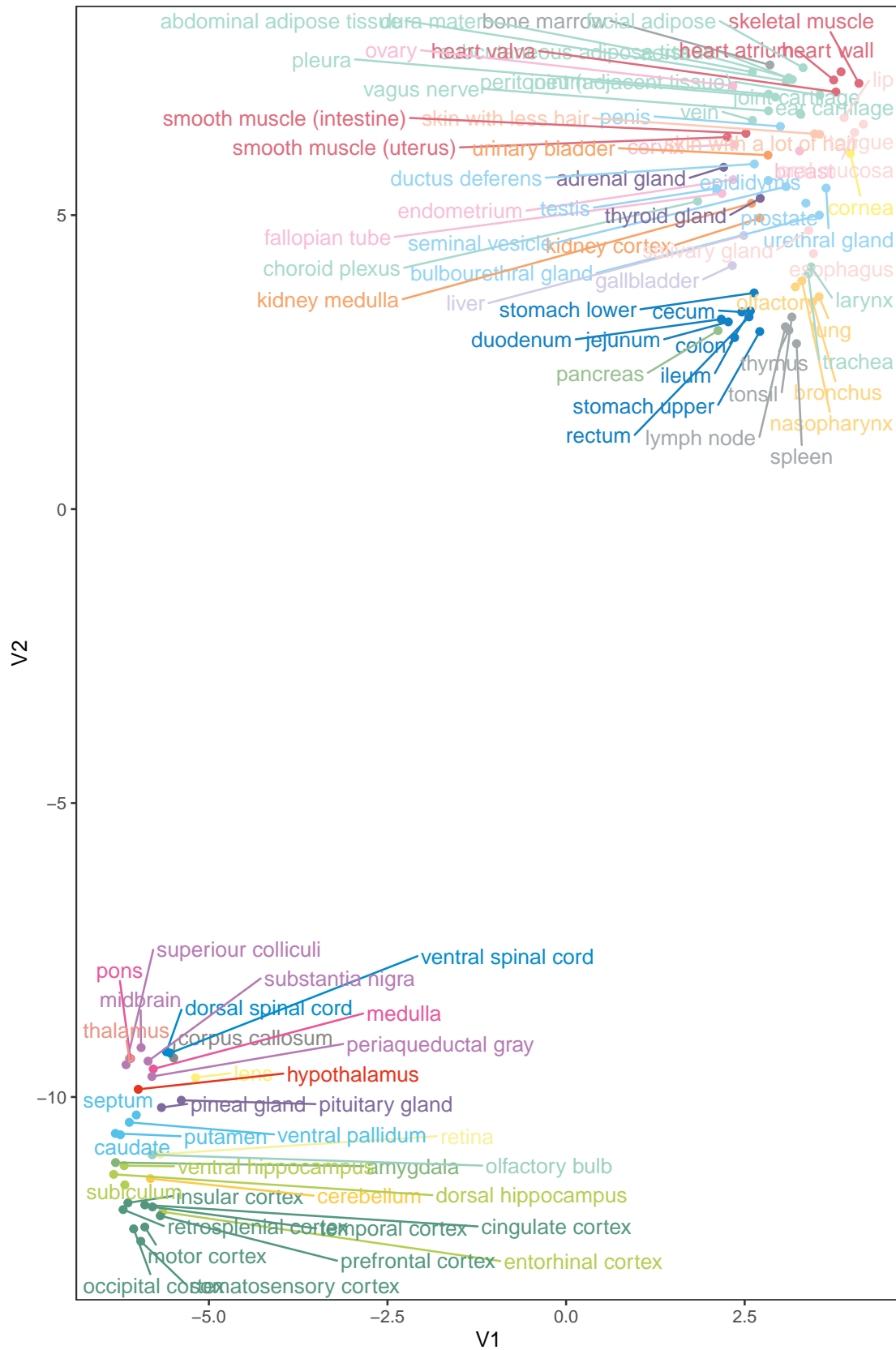
2. UMAP

subcutaneous adipose tissue oral-mucosa lining (adjacent skeletal muscle) heart valve heart wall



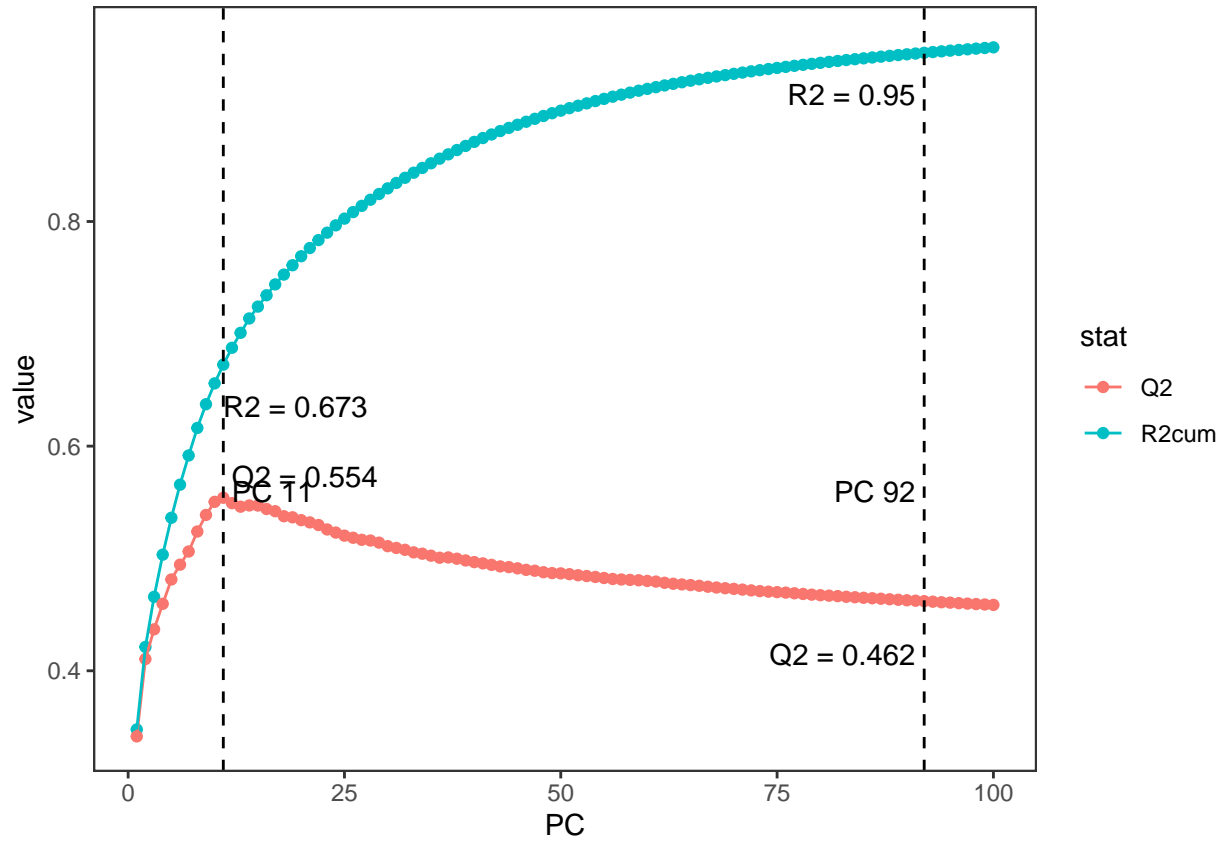
3. UMAP on 39 PCs from PCA

UMAP on PCA



Sample wise clustering

To investigate if we have any outliers or sample mixups we again perform a PCA, but this time on the full data. That is, we do not perform any averaging between samples of the same tissue. We see that 92 PCs represent 95% of the variance.



Hierarchical clustering of individual samples. Individuals are displayed to the left and the tissue name to the right for each sample. Below we see a dendrogram built from $1 - \text{Spearman's } \rho$ between the tissues using all genes.

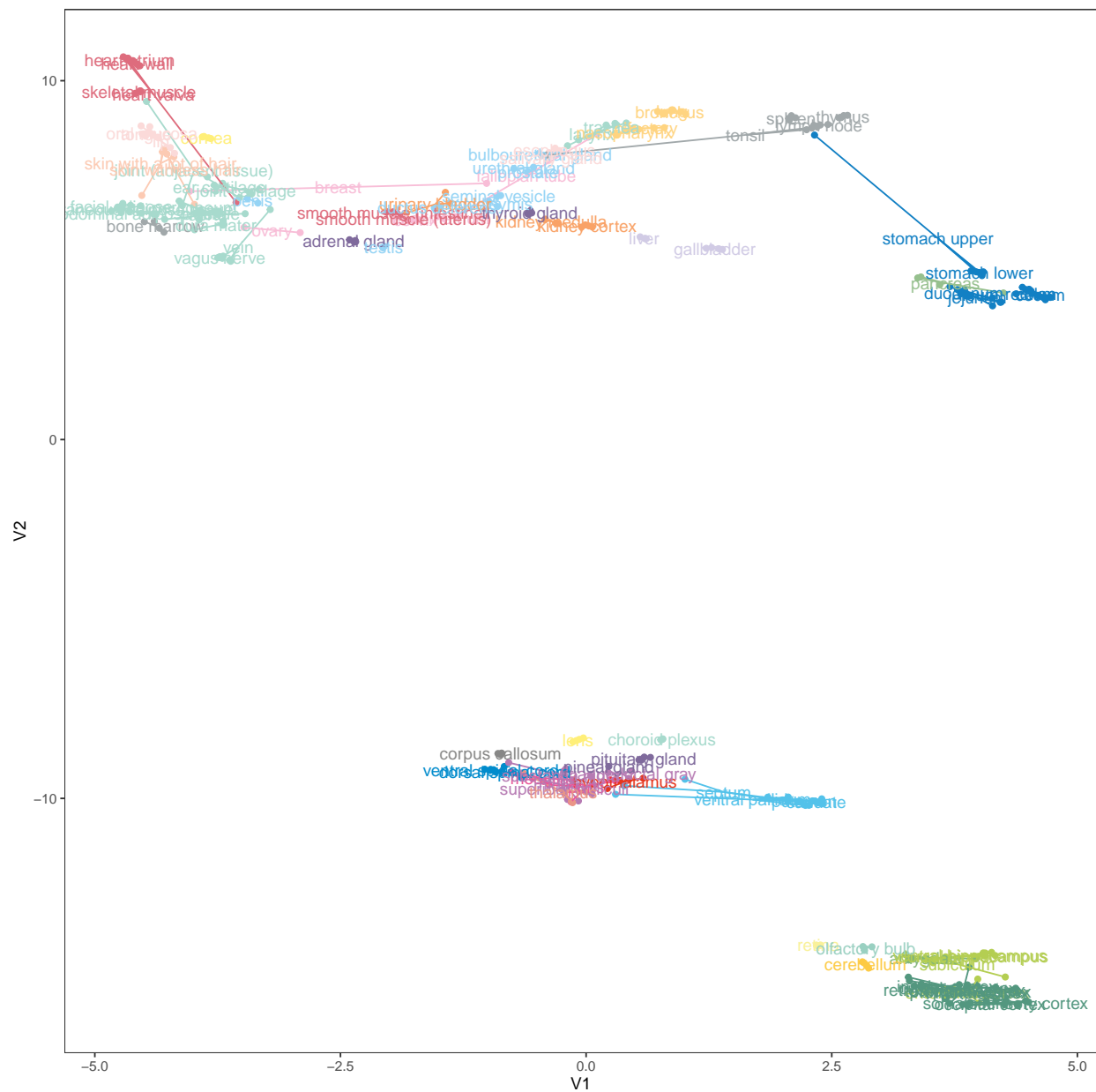
Below is an alternative (and complementary) clustering of tissues using Ward clustering of PCA scores in the 92 selected components. This type of clustering is based on minimizing the variance within each cluster and thus takes distances into account, which correlation does not.



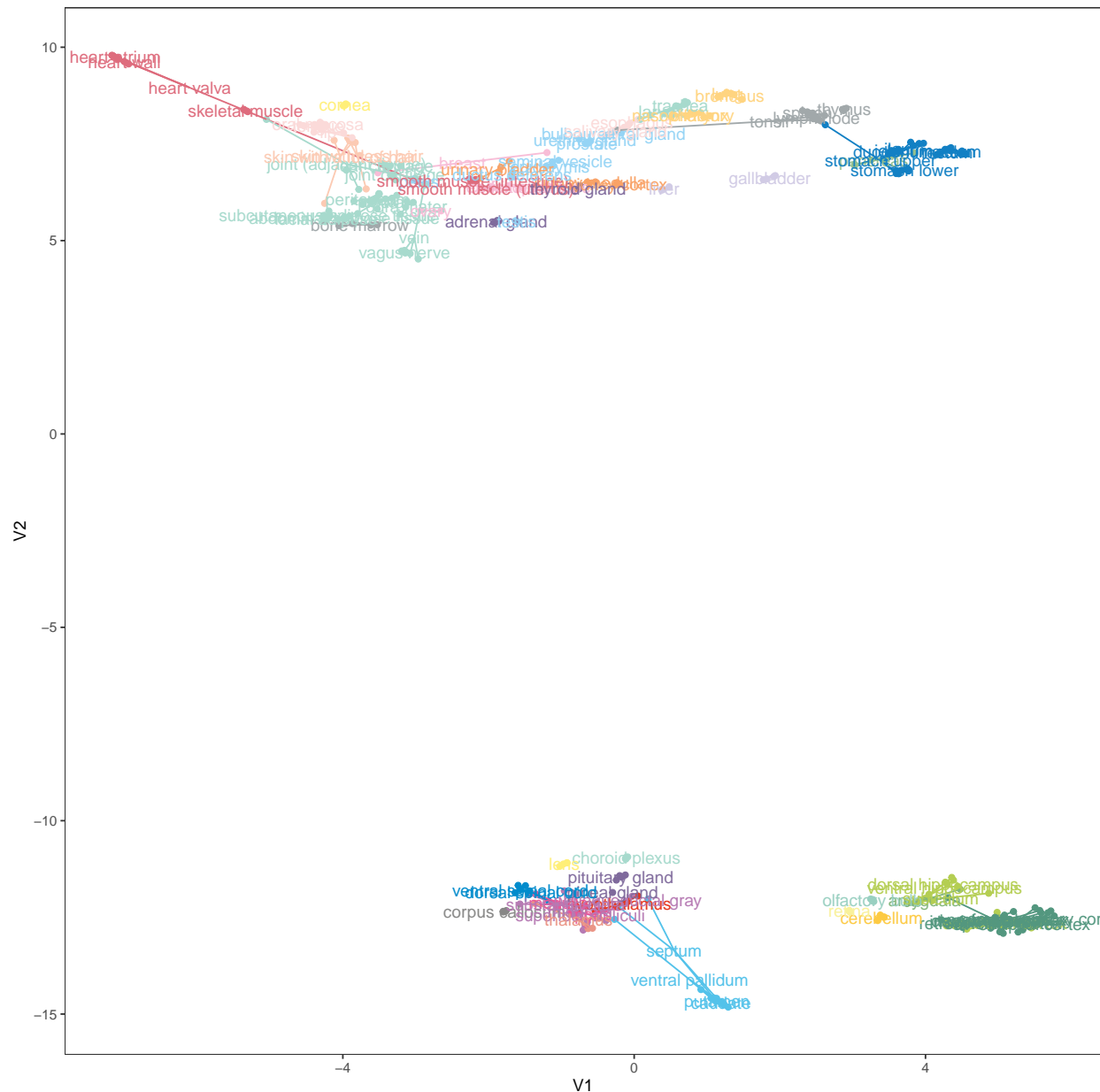
1. PCA



2. UMAP



3. UMAP on 92 PCs



Investigate sample mixup

Filtered samples

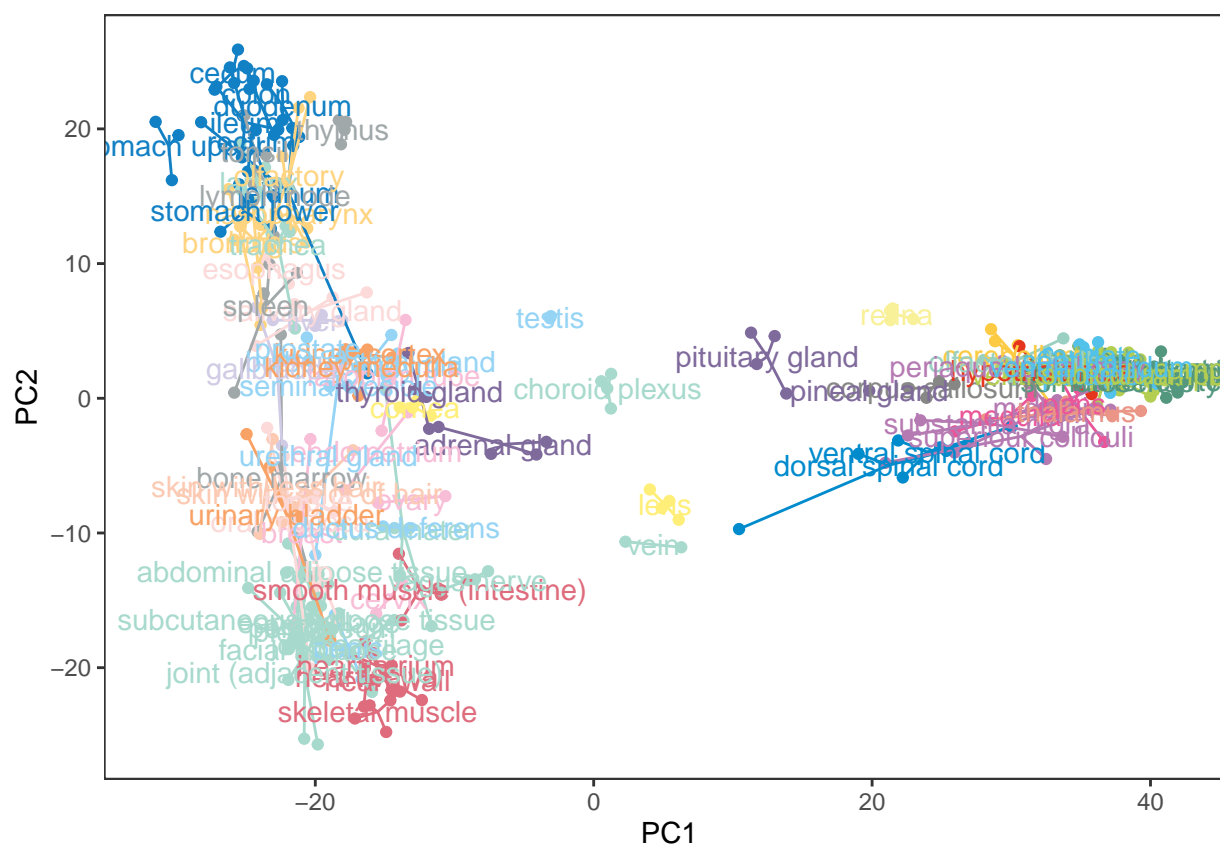
The following section describes which samples were removed from the data set. Summary plots for each individual tissue can be found in the file “Sample QC summary.pdf”. Please see this file for plots showing how samples cluster in relation to their tissue and to the dataset as a whole.

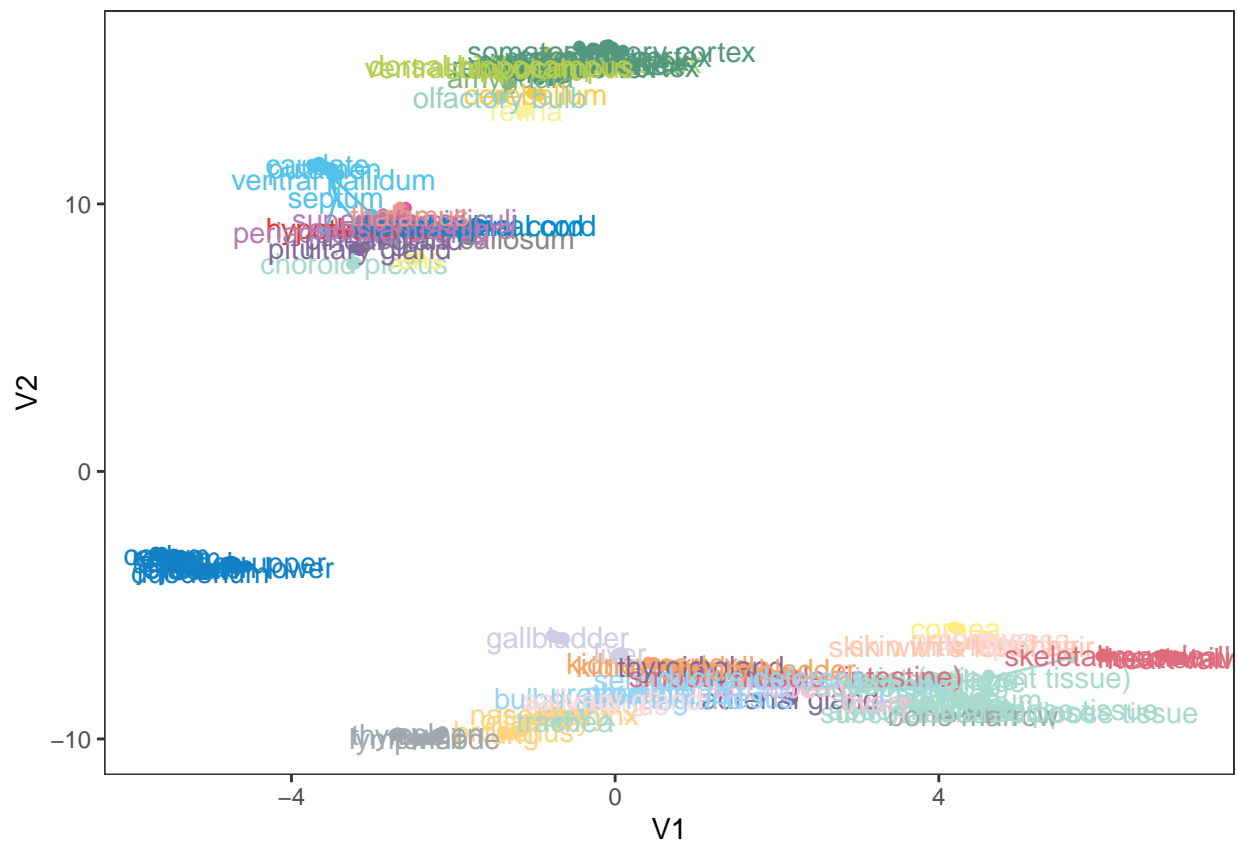
```
## Warning: Expected 2 pieces. Missing pieces filled with `NA` in 1 rows [17].
```

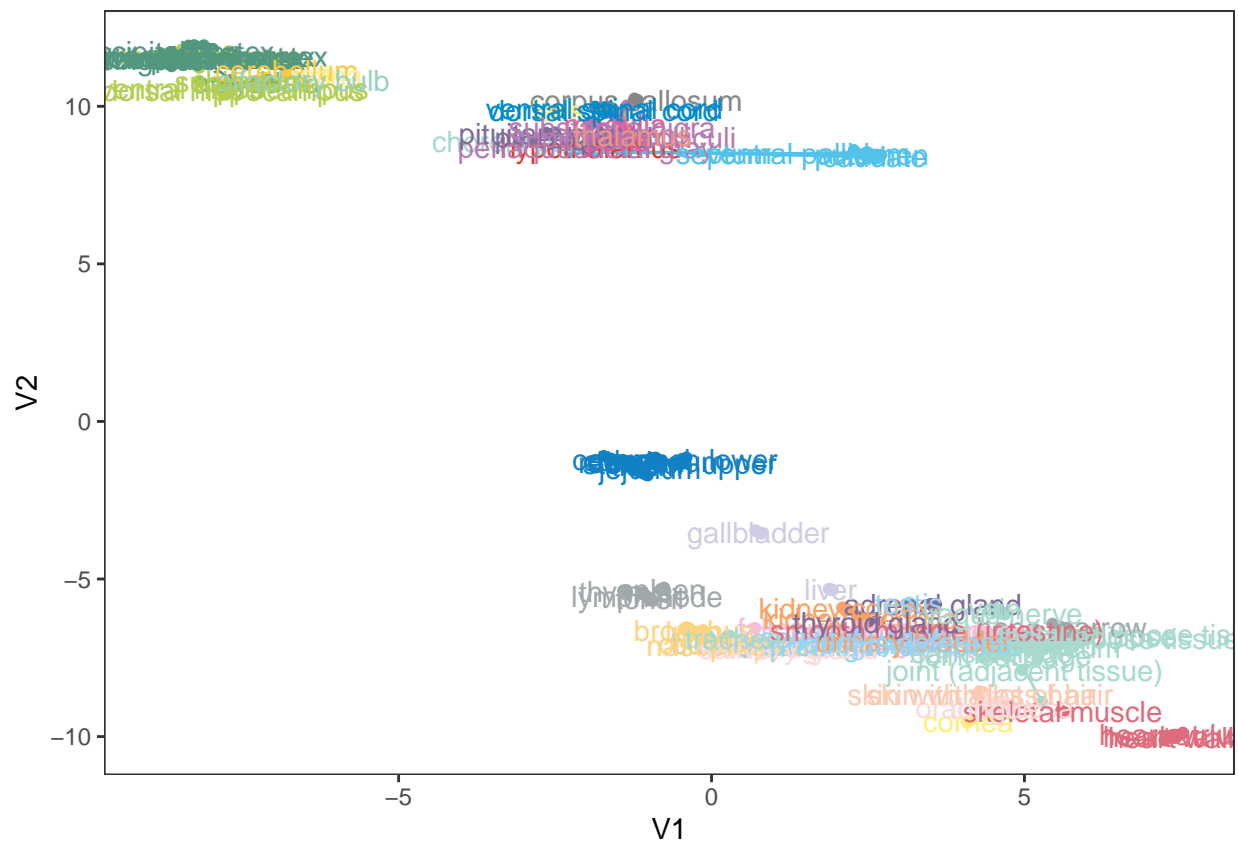
Table 1: Samples that were removed from the data set

sample_ID	tissue	individual	comment
vein_c	vein	c	Wrong tissue collected, nerve was collected instead of vein as confirmed by histology
vein_d	vein	d	Clusters far from other vein samples
skiB_a	skiB	a	Clusters far from other skin samples
skiB_b	skiB	b	Clusters far from other skin samples
skiH_b	skiH	b	Clusters far from other skin samples
pan_a	pan	a	A large proportion of genes has TPM = 0
pan_b	pan	b	A large proportion of genes has TPM = 0
pan_c	pan	c	A large proportion of genes has TPM = 0
pan_d	pan	d	A large proportion of genes has TPM = 0
carJ_b	carJ	b	A large proportion of genes has TPM = 0 (and it therefore clusters with some pancreas samples)
skiL_c	skiL	c	Histology showed the sample included a lot of hair follicles
smoU_a	smoU	a	Histology showed that the sample was mostly cervix, and clustering confirms this
smoU_b	smoU	b	Histology showed that the sample was mostly cervix, and clustering confirms this
heaV_c	heaV	c	Clusters far from other heart muscle samples
stoU_a	stoU	a	Clusters far from other GI samples and the sample included some esophagus in contrast to other samples
ton_b	ton	b	Clustering and histology shows that sample is not tonsil, but esophagus
OB-3X	OB-3X	NA	This sample was an extra “rouge sample” of olfactory bulb from one individual

The following plots show the clustering after problematic samples have been excluded.







PCA

