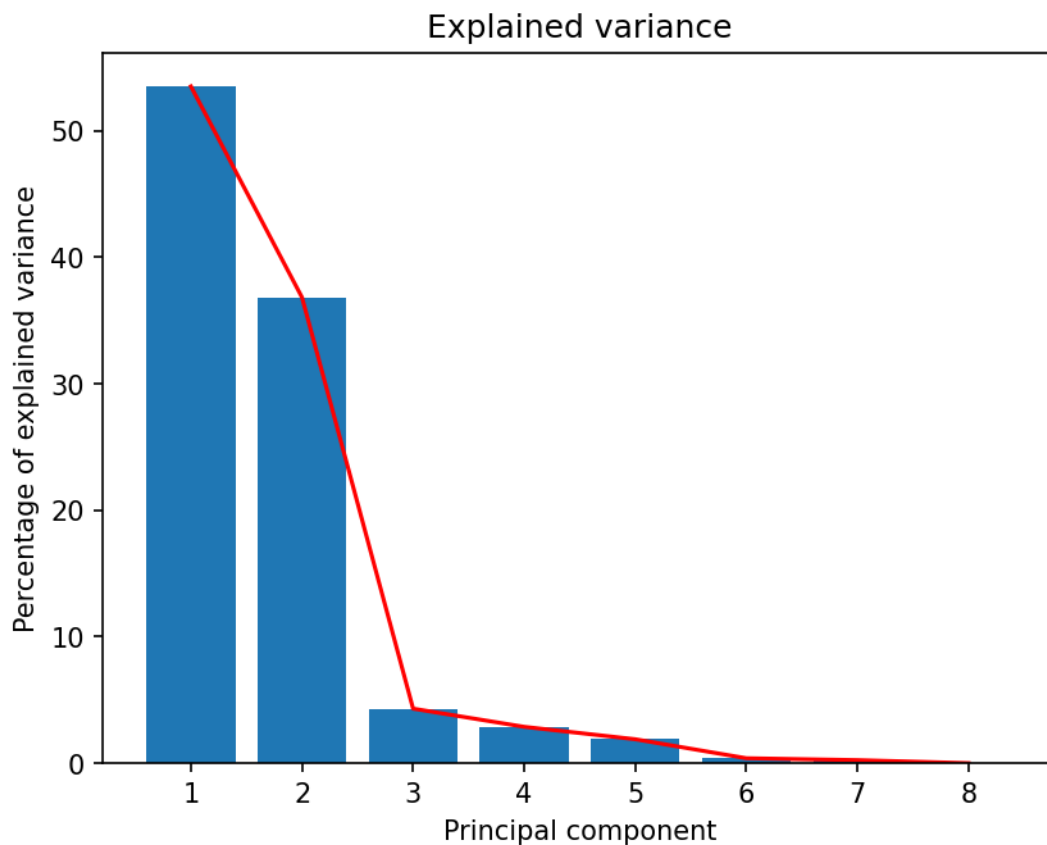


BT4110 – Computational Biology Lab  
**Transcriptome Data Analysis Assignment**

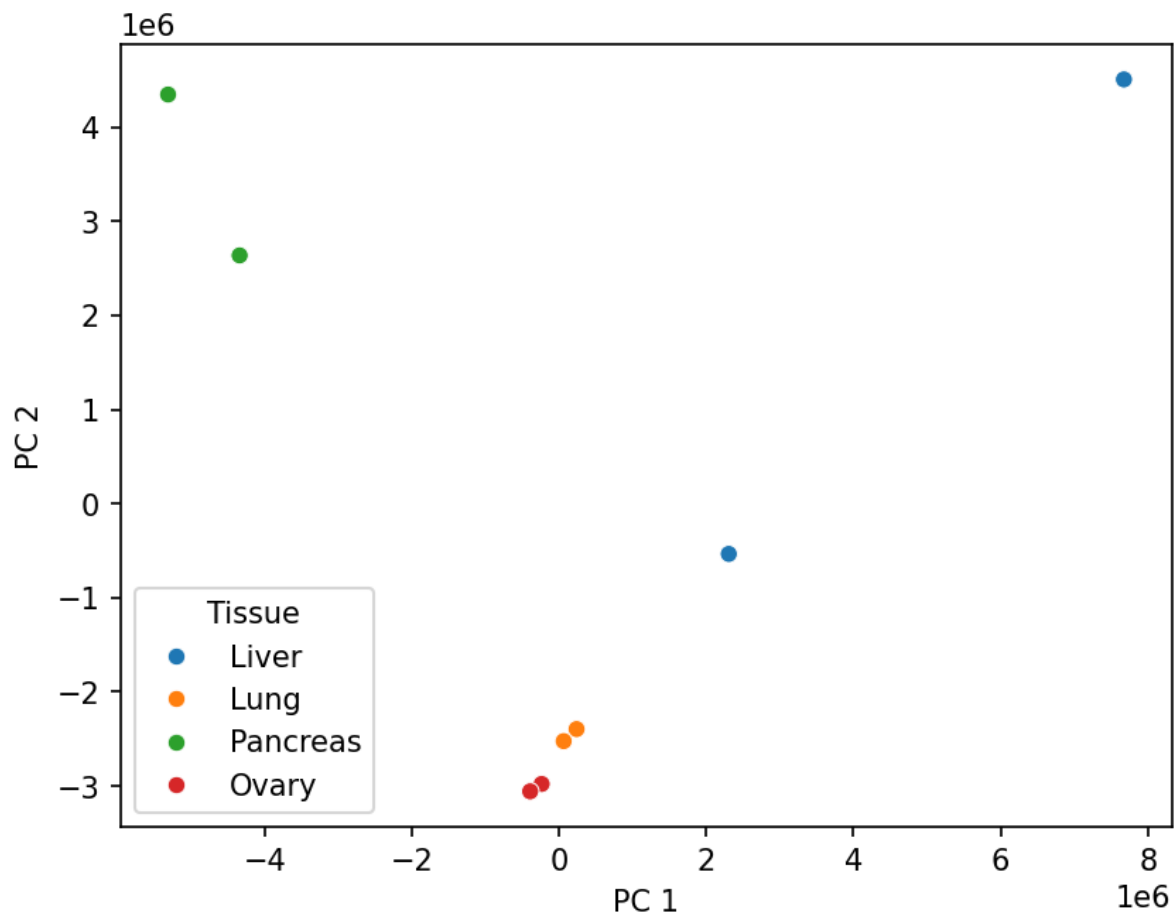
1

After scaling the gene expression data, PCA was performed using sklearn. Based on this, 8 principal components were obtained. The percentage of variance explained by each of the components is shown below:

Principal component	Explained variance (%)
PC 1	5.35E+01
PC 2	3.68E+01
PC 3	4.31E+00
PC 4	2.87E+00
PC 5	1.89E+00
PC 6	3.86E-01
PC 7	2.38E-01
PC 8	4.43E-28



When the first two PCs are plotted against each other, clustering of tissue types can be observed.



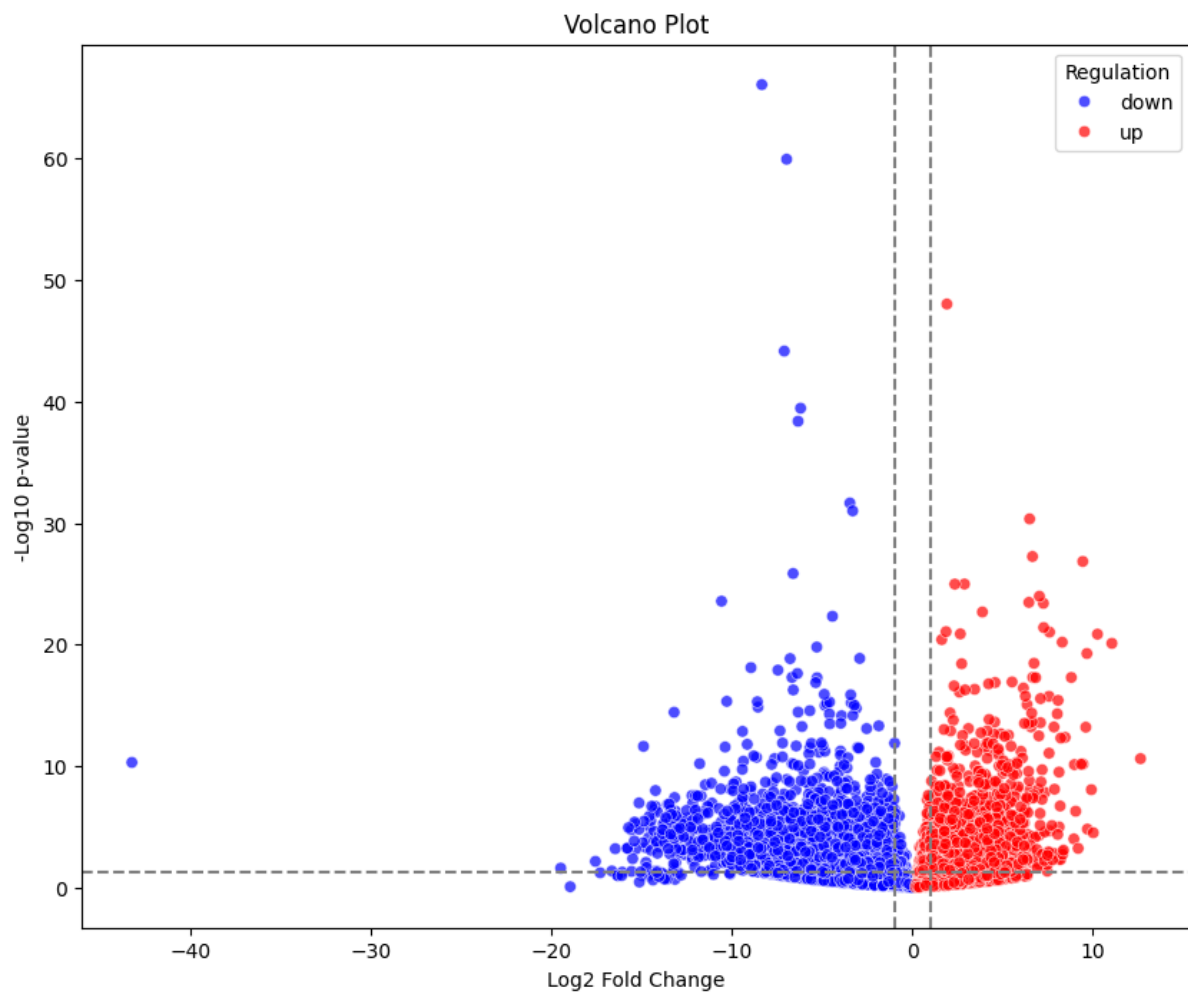
We can see that the ovary samples and lung samples are quite close to each other. The pancreas samples are also relatively close to each other. However, the liver samples are quite distant from each other. Overall, the PCA effectively captures the tissue types from the gene expression data.

## 2

PyDESeq2 was used to identify upregulated and downregulated genes in the ovary samples as compared to those from other tissues. This was done by suitably modifying the given metadata file.

Tissue	group
X5EGGE	Other
X5EQMM	Other
X5GICA	Other
X5N9D6	Other
X5P9GB	Other
X5N9ES	Other
X5P9GS	Ovary
X5N9E1	Ovary

Based on the output of DESeq2, 17355 genes are upregulated and 37237 genes are downregulated in the ovary samples as compared to the other tissues. This can be visualised using a volcano plot.

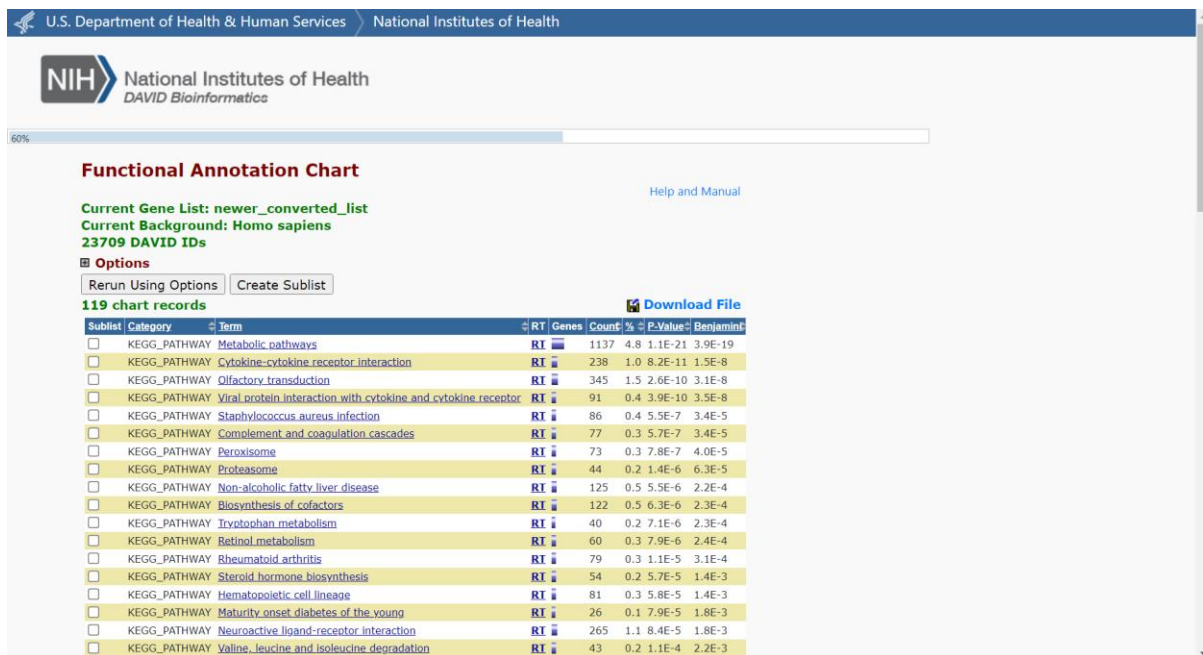


The list of upregulated genes was fed into DAVID for functional enrichment analysis. After mapping the genes to their names in DAVID DB, their pathway functional annotation was performed using the KEGG Pathway database. Of the 17355 upregulated genes, 10885 were available in the KEGG database.

Sublist	Category	Term	RT	Genes	Count	%	P-Value	Benjamini
<input type="checkbox"/>	KEGG_PATHWAY	Herpes simplex virus 1 infection	RT	304	2.8	5.0E-38	1.8E-35	
<input type="checkbox"/>	KEGG_PATHWAY	Cell cycle	RT	90	0.8	8.2E-11	1.4E-8	
<input type="checkbox"/>	KEGG_PATHWAY	Ribosome	RT	89	0.8	5.8E-8	6.8E-6	
<input type="checkbox"/>	KEGG_PATHWAY	Coronavirus disease - COVID-19	RT	115	1.1	7.8E-8	6.8E-6	
<input type="checkbox"/>	KEGG_PATHWAY	TGF-beta signaling pathway	RT	57	0.5	8.7E-6	6.1E-4	
<input type="checkbox"/>	KEGG_PATHWAY	Polycomb repressive complex	RT	46	0.4	1.5E-5	8.8E-4	
<input type="checkbox"/>	KEGG_PATHWAY	Nucleocytoplasmic transport	RT	55	0.5	4.8E-5	2.4E-3	
<input type="checkbox"/>	KEGG_PATHWAY	Wnt signaling pathway	RT	79	0.7	1.5E-4	6.6E-3	
<input type="checkbox"/>	KEGG_PATHWAY	Hedgehog signaling pathway	RT	31	0.3	5.1E-4	1.9E-2	
<input type="checkbox"/>	KEGG_PATHWAY	Dopaminergic synapse	RT	61	0.6	5.4E-4	1.9E-2	
<input type="checkbox"/>	KEGG_PATHWAY	Fanconi anemia pathway	RT	30	0.3	8.9E-4	2.8E-2	
<input type="checkbox"/>	KEGG_PATHWAY	Homologous recombination	RT	24	0.2	9.8E-4	2.9E-2	
<input type="checkbox"/>	KEGG_PATHWAY	Melanogenesis	RT	48	0.4	1.1E-3	2.9E-2	
<input type="checkbox"/>	KEGG_PATHWAY	Thyroid hormone signaling pathway	RT	56	0.5	1.1E-3	2.9E-2	
<input type="checkbox"/>	KEGG_PATHWAY	Hippo signaling pathway	RT	69	0.6	1.2E-3	2.9E-2	
<input type="checkbox"/>	KEGG_PATHWAY	Signaling pathways regulating pluripotency of stem cells	RT	64	0.6	1.3E-3	2.9E-2	
<input type="checkbox"/>	KEGG_PATHWAY	Hepatocellular carcinoma	RT	73	0.7	1.9E-3	3.8E-2	

- The herpes simplex virus has been [shown](#) to infect the ovary.
- The ovary is the site of gametogenesis via meiosis, which employs the cell cycle pathway.
- Meiosis also requires homologous recombination, which is shown to be enriched.
- The [TGF- \$\beta\$  signalling pathway](#), the [Wnt pathway](#), and the [hedgehog pathway](#), have all been shown to play a role in ovarian follicle development.
- The Polycomb repressive complex has been [shown](#) to play an important role in the epigenetic regulation of ovary function.
- The proteins associated with the Fanconi anaemia pathway are also [known](#) to be involved in follicle development in the ovaries.

The list of downregulated genes was then fed into DAVID. After mapping the genes to their names in DAVID DB, their pathway functional annotation was performed using the KEGG Pathway database. Of the 37237 upregulated genes, 23709 were available in the KEGG database.



None of the pathways identified seem to be directly and specifically associated with the ovaries. Thus, DESeq2 has accurately identified biologically relevant upregulated and downregulated genes in the ovary compared to other tissue types.