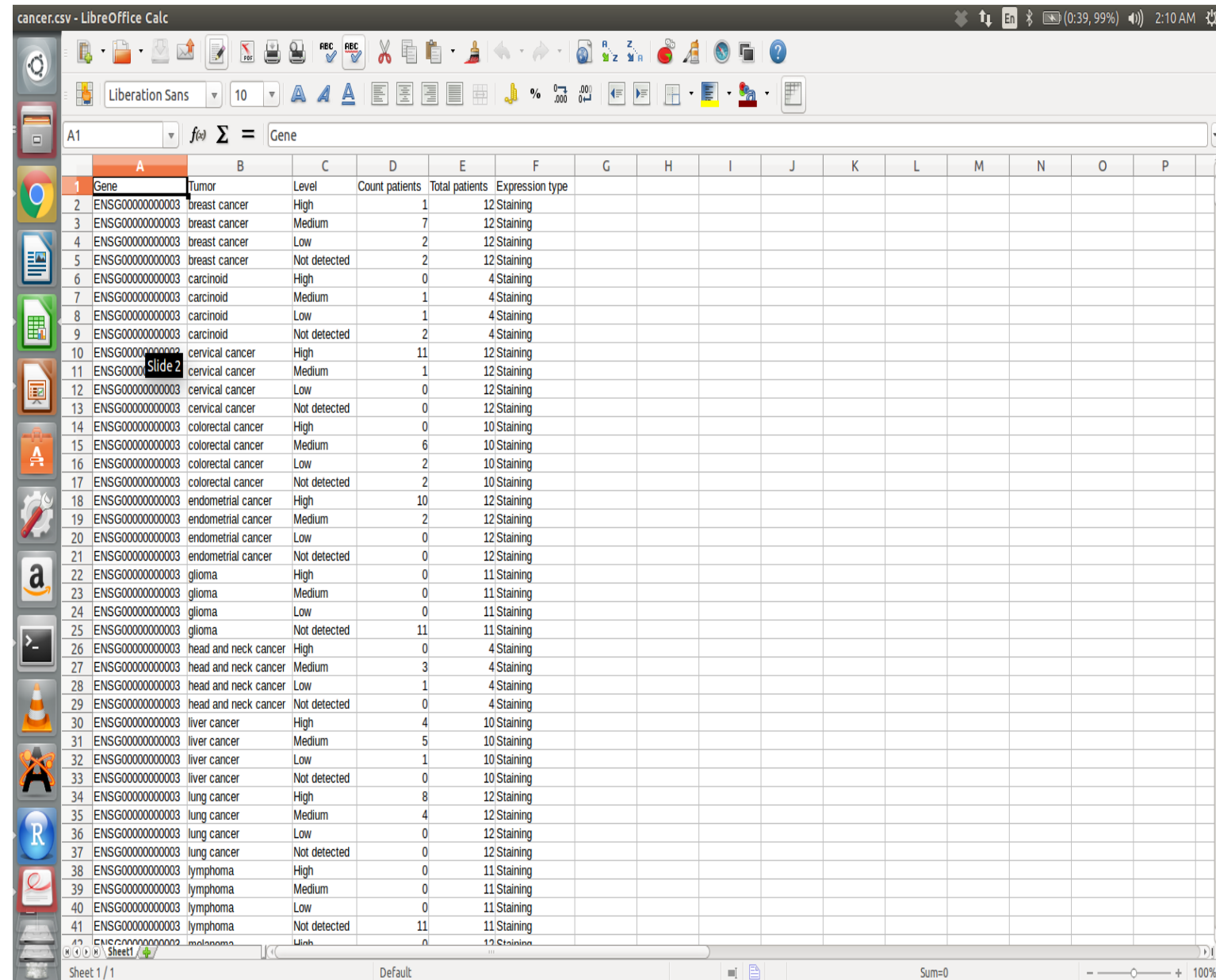# COMPARATIVE ANALYSIS ACROSS CANCERS USING PROTEIN DATA
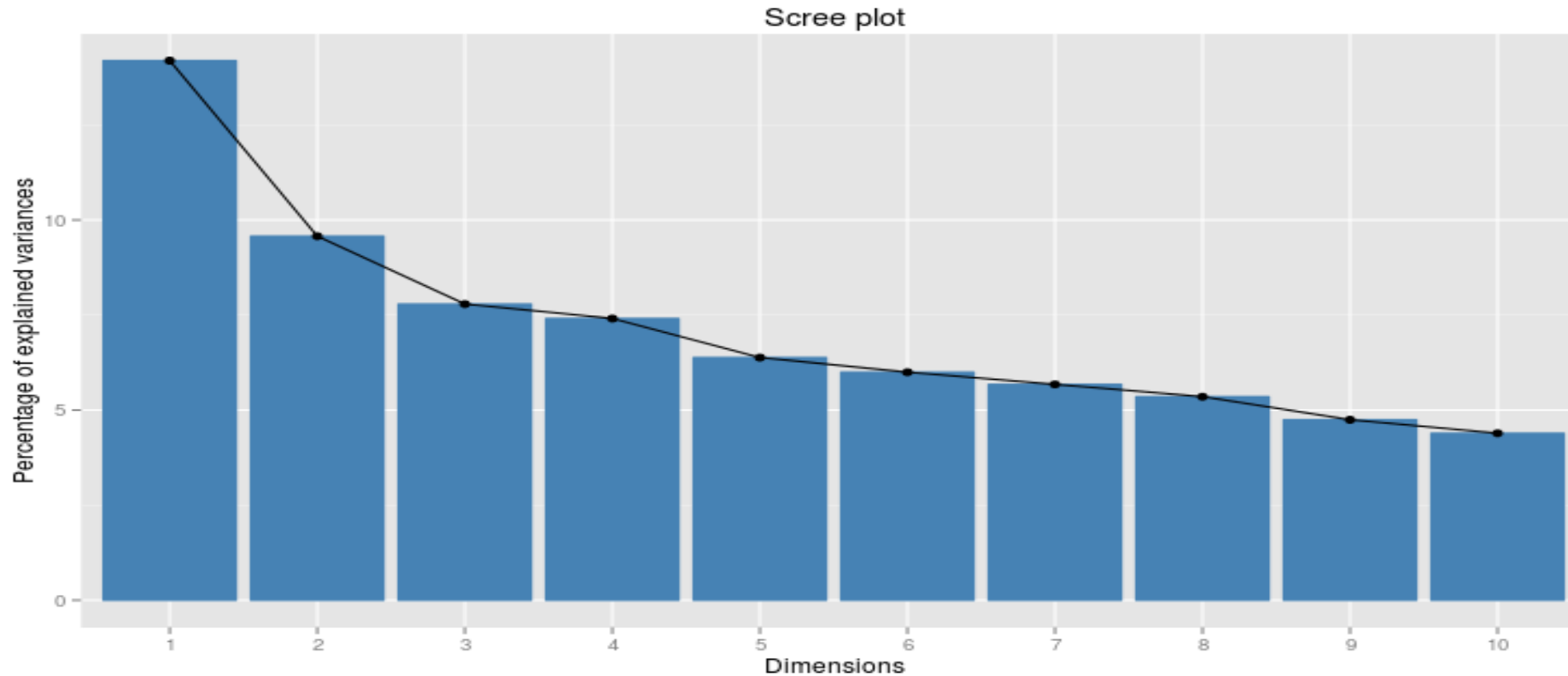
Noor Pratap Singh

# DATA

- We worked with cancer data and tissue data.

- Taken from Human Protein Atlas(HPA).

- Cancer data contains a total of 16613 unique Ensembl Ids.

- We also took specific patients' data from http://msb.embopress.org/content/10/3/721

# MCA Analysis

- We did MCA Analysis on the entire set of genes.



We get the above contribution along the respective dimensions.

# Overall MCA Plot

The picture is clumsy and not much can be interpreted from it.

But it plots the entire variables as well as the individuals and the variables being huge messes it up.



MCA factor map - Biplot

# Some Cancers are clearly separated !

- We plot an individual wise plot.

- We can clearly see not all cancers have uniform distribution with respect to gene levels.



Individuals factor map - MCA

Other Plots

Variable categories- MCA

All variables

MCA factor map - Biplot

Asymmetric Plot

MCA factor map - Biplot

Biplot with cos2 > 0.6

# Cluster Analyis

- We still don't know the exact dimensions to chose.
- So we do cluster analysis to actually find out which cancers can be actually clubed together.
- We get a total of 7 clusters.

# MCA for metabolic genes

- We do the same analysis on the same data but only for metabolic genes.

# Cluster Analysis

- Overall individual wise both look very similar.

- The cluster analysis brings a different story.

- There are 7 clusters in former while 6 clusters in latter.

- Also position of cancers in cluster change.

- So far we have considered only metabolic genes.

- We plan to perform similar analysis based on other proccesess such as cell cycle etc.

# Hepatocellular Carcinoma

Comparing Normal and Cancer Tissue data

**Problems**
- Cancer data associated with each gene for each cancer 4 levels and their counts.
- Normal tissue data on the other hand had associated with it each a cell type with only *one* level and no count.
- So it made comparing the two datasets difficult.

**Solution**
- So in order to fix a level with respect to gene for a cancer we make 4 basis of comparison(50,75,90, 100) based on count clubbing the 4 levels into 2.

# Identification of Differentially Expressed Genes(DEG)

- Once the level of liver cancer w.r.t to each gene is calculated based on type of mode(50,75,90,100) we compare it with normal tissue data of hepatocytes.

- The genes which do not match contribute to the differentially expressed genes(DEG).

- As per our calculation we get the following number of DEG:
    1. **50% -** 3547
    2. **75% -** 1532
    3. **90% -** 697
    4. **100%-** 310

# Functional Annotation

- After identification of DEGs for functional annotation of DEGS we tried different tools like David, Genecodis and GOSim package in R.

- For DAVID analysis we used BP 5 and we get following results:
  - **50%** 84 GoTerms with 34 having p value < 0.05
  - **75%** 53 GoTerms with 25 having p value < 0.05
  - **90%** 18 GoTerms with 9 having p value < 0.05
  - **100%** 4 GoTerms with 1 having p value < 0.05

- The results were not very satisfactory.

# GOSim

- We used 75 % differentially expressed genes.
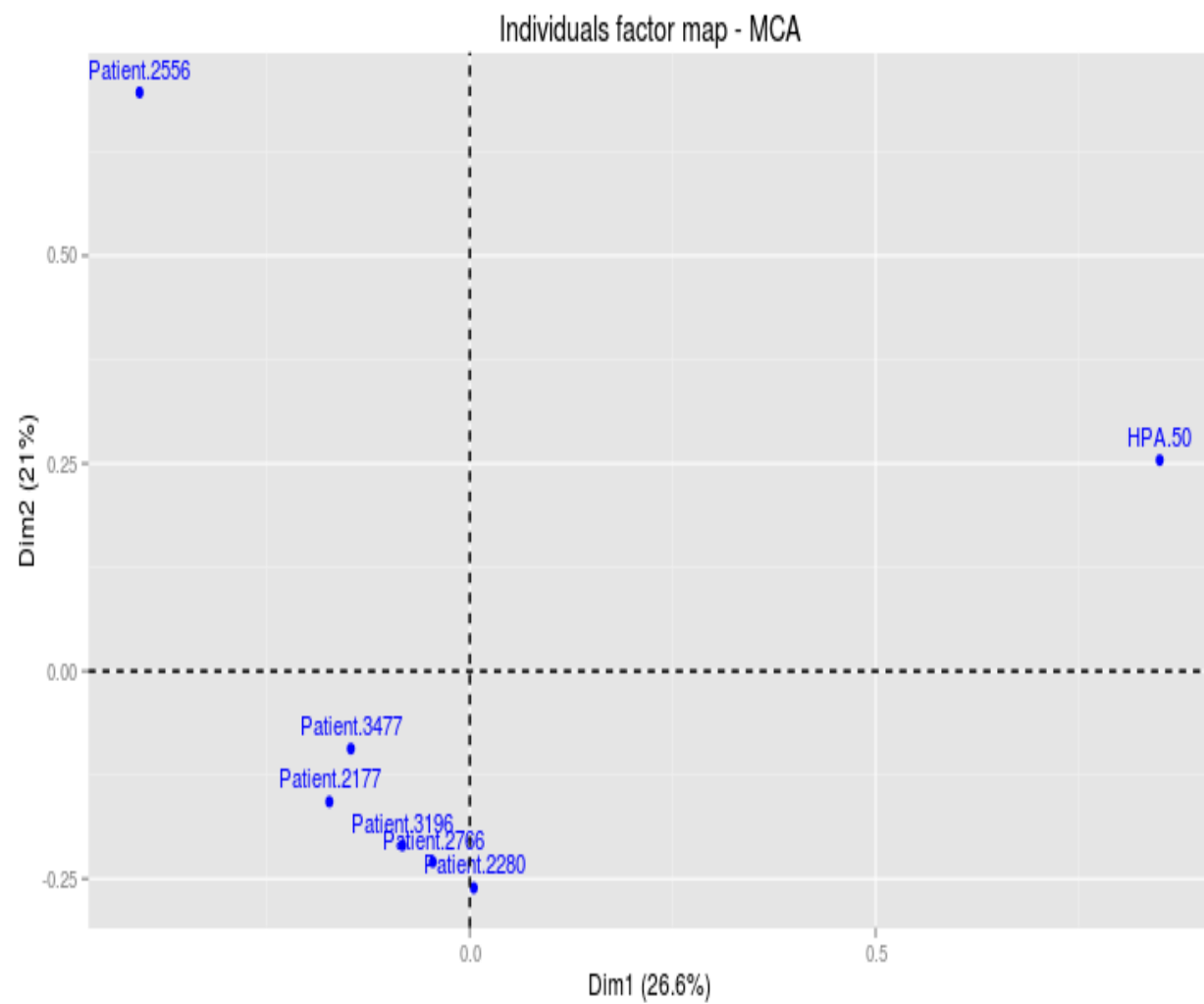
- We did GOenrichment analysis on the above set.

- We get a total of 32 different GOTerms which do not match DAVID.

- We also found the similarity matrix using Jiang and Conrath method.

- Based on above similarity matrix we use k mediod clustering with k = 8 (maximum silhouette width tried for different values of k)

- We again did GOenrichment analysis on individual clusters.

# Patient Wise Comparison

- In order to verify that whether we have correctly picked DEGs we also take patient data.(The number of genes are less than HPA and 6 patients)

- We also do a patient wise protein expression comparison with HPA at 50% since it is the most generic.

- We get the following number of proteins with same expression level across patients: 11946, 12414, 11442, 12206, 12258, 12151.

- A question might arise that IDs that do not match may be the DEGS but if you match DEGS at 50% with above genes you get the following:
  - 813, 1198, 666, 1130, 979, 856.

- A point worth noting that 50 % is the most generic case and the patients might be in different stages of cancer.

# MCA Analysis



Individuals factor map - MCA

Clustering

Cluster Dendrogram