

# **Perform differential expression analysis on fibrotic and non-fibrotic patients under 4 different treatments on HPC**

## **Author**

Sam (Cheng-Hsiang) Lu

Email: [Cheng-Hsiang.Lu@cshs.org](mailto:Cheng-Hsiang.Lu@cshs.org)

## **Mentor**

David Casero

Email: [David.Casero@cshs.org](mailto:David.Casero@cshs.org)

## **Background**

### **Inflammatory bowel diseases(IBD)**

Inflammatory bowel diseases (IBD) are a group of chronic conditions that cause inflammation and damage to the digestive tract. The two main types of IBD are Crohn's Disease and Ulcerative Colitis.

Crohn's disease can affect any part of the digestive tract, like small or large intestine, and can cause symptoms such as abdominal pain, diarrhea, weight loss, and fatigue. It can also cause complications such as fistulas (abnormal

connections between different parts of the intestine) and strictures (narrowing of the intestine).

Ulcerative colitis, on the other hand, affects only the colon and rectum and causes symptoms such as bloody diarrhea, abdominal pain, and a frequent need to pass stools. It can also lead to complications such as inflammation of the skin, eyes, and joints.

Both Crohn's disease and ulcerative colitis are chronic conditions, meaning they can last for a lifetime and require ongoing treatment to manage symptoms and prevent complications.

## **Induced Pluripotent Stem Cells(iPSCs)**

Induced pluripotent stem cells (iPSCs) are a type of stem cell that are generated in the laboratory by reprogramming adult cells, such as skin or blood cells, to a pluripotent state. A pluripotent state means that the cells have the potential to develop into any type of cell in the body, just like embryonic stem cells.

iPSCs offer several advantages as they can be generated from the patient's own cells, avoiding issues with immune rejection, ethical concerns and the need for embryos.

iPSCs can be used to study the underlying causes of diseases, test new drugs and therapies, and potentially generate replacement tissues or organs for transplantation.

## **Aim**

In this project, we will be analyzing RNA-seq data from 19 samples, comprising of 10 samples with fibrotic complications and 9 non-fibrotic samples. Each sample has undergone two runs and 4 different treatments(untreated, TGF-B, TNF-A, and TGF-B+TNF-A), resulting in a total of 151 samples(1 library failed). We used induced pluripotent stem cells (iPSC) to differentiate into myofibroblasts and stimulated the system with different signals to observe its development. The objective is to investigate the effect of four different treatments: untreated, TGF-B, TNF-A, and TGF-B+TNF-A on the development of the system. In the end, we will perform differential expression analysis to identify the genes that are differentially expressed in fibrotic and non-fibrotic samples under 4 treatments.

# Pipelines

## In HPC

### Convert 151 Fastq to Fasta files

First, I put all fastq.gz files in one folder and list all fastq files' name in fastqfiles.txt

```
ls *.q.gz > fastqfiles.txt
```

Cut redundant suffix “\_R1\_trimmed” and list all fastq files' name in libraryname.txt and prefix.txt

```
ls *.q.gz | cut -f 1 -d '.' | sed 's/_R1_trimmed//g' >libraryname.txt
```

```
ls *.q.gz | cut -f 1 -d '.' | sed 's/_R1_trimmed//g' > prefix.txt
```

Form a table with 3 columns: fastqfiles.txt libraryname.txt prefix.txt

```
paste fastqfiles.txt libraryname.txt prefix.txt > tofastatable.txt
```

Create small-sized fasta-formatted files. To submit this job to the cluster on HPC, you need to read the file, library, and prefix. Once you have done that, run the script “generatefastaFromFastaqz” which will combine the script “DCfastaqTofastaLibraryId.pl”. This results in small-sized fasta-formatted files contain only one header and one sequence per read. You can find all scripts in the “scripts” folder.

```
cat tofastatable.txt | awk {print}' | while read file library prefix ;  
do qsub -cwd -o $PWD -e $PWD -l h_data=2048M,h_rt=8:00:00  
$HOME/scripts/generatefastaFromFastaqz $file $library $prefix
```

done

This is what each fasta-formatted file would look like

```

>417iP34untreatedM_S96_1
TNGGCTGGCTAAAGAAGTGAGTATGACCCCAGAGGCCAGAGAGGGCAGGGAGAGAATGCCTGGCCACTT
>417iP34untreatedM_S96_2
CNCTTCTTTTCGGCGTAGCTCATCAACCTCATATGGTGTGAGCCTTGCTACTTCCGGATGTTCCACATAA
>417iP34untreatedM_S96_3
CNTCCGCCTTTCTGTTCTGTTTTGCCGTCCCCTCCCGCCTGCTGCTGGCCCTGCCTCGCCCCGGGCGGC
>417iP34untreatedM_S96_4
ANCCATGTGAAAAACAGAACTGTAGGGGTTTTTTGTTTGTGTTTTTTGAGATGGGGTCTTGCTCTGT
>417iP34untreatedM_S96_5
ANTCCGATGTTCTGCAATTTTCTGTGCCCTAGGTTGAACTTCTTTCAGCATTGCACTAGCATCTTCATC
>417iP34untreatedM_S96_6
CNCGACATTAGAAGGTTTTCTGTGGATGGATCGGGCACCGTCTTCCTCATATTCCTTTTGGAGACCC
>417iP34untreatedM_S96_7
CNCCTTGCGCACCTCCACCATGACCAGCCCTTCCTTACCAAGCCCAGCCCCACATCGCCCTTGAATC
>417iP34untreatedM_S96_8
TNAGGTCTCACTGTCCTGGCCTGACCTTCAGCTCTCCAACACTGGGCCCCGCCGGGTCTCCGGGAGCCA
>417iP34untreatedM_S96_9
CTGGGGTTATGAGTTTATAGTTGGGAACTTCCTTACAGAGTTTATCATAGGTAGCTTTGTCAAACAAGA
>417iP34untreatedM_S96_10
CCAGTGAGAAGACAGCTTTGCAGTCACACTGGAGATCAGAGTTCCAGGCTGCAGCATGTCACCAACGCC

```

## Generate auxiliary files and directories for each sample

Put a list of names of all fasta files in the directory and save them in a text file named “fastafiles.txt”

```
ls *fasta.gz > fastafiles.txt
```

Cut the redundant suffix “.fasta.gz” from the names of all fasta files and generate a new list of file names with the suffix removed in a text file named “targetdirectories.GTF.txt”

```
cat fastafiles.txt |sed 's/.fasta.gz//g' > targetdirectories.GTF.txt
```

Create a separate directory for each sample listed in “fastafiles.txt”

```
cat fastafiles.txt |while read line ; do mkdir
${line/.fasta.gz/GTFpass1/} ; done
```

## Form the submission script called “sendmyof”

Add a shebang line at the beginning of your script file named “sendmyof” to indicate the interpreter that should be used to execute the script

```
echo '#!/bin/bash/' > sendmyof
```

The command below runs the “generatesendscriptSingleGTFParam” script with several input parameters to map the RNA-seq data with STAR. The input parameters include the list of target directories containing the input data (“targetdirectories.GTF.txt”), the subdirectory name (“GTFpass1”), a parameter file containing settings for STAR alignment (“Parameters.txt”), a prefix for output files (“myof”), the path to the STAR index directory (“/home/luc/RNASEQ\_MASTER/Hsapiens/GRC38/INDEXES/GRCh38.primary.33.basicselected.STAR2.7.3a/”), the path to the input data directory (“/home/luc/iPSC/MYOFIBROBLAST/”), the amount of free memory to use (“mem\_free=32G”), and the number of threads to use (“8”). In the end, it will generate a “processLaneSingleGTFParam” file and run the STAR package in each sample’s folder.

```
./generatesendscriptSingleGTFParam targetdirectories.GTF.txt GTFpass1
Parameters.txt myof
/home/luc/RNASEQ_MASTER/Hsapiens/GRC38/INDEXES/GRCh38.primary.33.basic
selected.STAR2.7.3a/ /home/luc/iPSC/MYOFIBROBLAST/ mem_free=32G 8 >>
sendmyof
```

Change sendmyof into executable mode and run sendmyof

```
chmod a+x sendmyof
. sendmyof
```

It will take less than one day to run through 151 samples and generate each sample a folder which contain every output from STAR

## Create a table summarizing the mapping statistics for each sample

Change directory into one sample file which ends with “GTFpass1”. Extract the first column from the mapping statistics file and store it in “temp2.txt”

```
grep "|" 008iP22TGFbM_S71GTFpass1/008iP22TGFbM_S71GTFpass1Log.final.out
| cut -f 1 -d "|" | sed 's/^ *///g' | awk 'NR>3 {print}' > temp2.txt
```

The first column from the mapping statistics file

```

                Started job on
                Started mapping on
                Finished on
Mapping speed, Million of reads per hour

                Number of input reads
                Average input read length
                UNIQUE READS:
                Uniquely mapped reads number
                Uniquely mapped reads %
                Average mapped length
                Number of splices: Total
Number of splices: Annotated (sjdb)
                Number of splices: GT/AG
                Number of splices: GC/AG
                Number of splices: AT/AC
                Number of splices: Non-canonical
                Mismatch rate per base, %
                Deletion rate per base
                Deletion average length
                Insertion rate per base
                Insertion average length
                MULTI-MAPPING READS:
                Number of reads mapped to multiple loci
                % of reads mapped to multiple loci
                Number of reads mapped to too many loci
                % of reads mapped to too many loci
                UNMAPPED READS:
Number of reads unmapped: too many mismatches
                % of reads unmapped: too many mismatches
                Number of reads unmapped: too short
                % of reads unmapped: too short
                Number of reads unmapped: other
                % of reads unmapped: other
                CHIMERIC READS:
                Number of chimeric reads
                % of chimeric reads

```

Create an empty temporary file for storing intermediate results

```
rm tempprev.txt
touch tempprev.txt
```

Extract the total mapped reads from each subsequent mapping statistics file and combine with previous results

```
ls *pass1/*final.out | while read line ; do
grep "|" $line | cut -f 2 > temp.txt
paste tempprev.txt temp.txt > tempnew.txt
mv tempnew.txt tempprev.txt
done
```

Remove the first column and write the final results to a file called “mappingstatsFirstpass.txt”

```
cut -f 2- tempprev.txt | awk 'NR>3 {print}' > tempnew.txt
mv tempnew.txt tempprev.txt
paste temp2.txt tempprev.txt > mappingstatsFirstpass.txt
```

The mappingstatsFirstpass.txt would look like this

Mapping speed, Million of reads per hour	452.92	344.45	460.17	487.96	331.34	441.99	560.78	487.93	532.33	351.69	381.99	447.59	463.27
Number of input reads	46046803	48055908	43588067	47982561	39392754	33886137	72122988	53266027	75	75	75	75	75
Average input read length	75	75	75	75	75	75	75	75	75	75	75	75	75
Uniquely mapped reads number	42647348	36668193	39982295	43570276	36214207	31289999	66061689	48190947	92.62%	89.75%	91.73%	90.80%	91.93%
Uniquely mapped reads %	92.62%	89.75%	91.73%	90.80%	91.93%	92.34%	91.60%	90.47%	93.77%	90.29%	90.70%	92.04%	91.95%
Average mapped length	75.21	75.21	75.20	75.20	75.19	75.21	75.20	75.20	75.20	75.24	75.24	75.23	75.22
Number of splices: Total	14107761	11743321	12025357	12627415	11851597	10196588	28128263	13942120	14029848	11673663	11928403	11945482	12536747
Number of splices: Annotated (sjdb)	14015583	11657575	11928403	12523568	11777228	11787471	10140891	19995363	18015583	14368	87757	62434	59661
Number of splices: GT/AG	78050	74078	84368	87757	62434	59661	145139	96183	139381	88755	60185	79428	94829
Number of splices: AT/AC	7048	6261	6785	8291	6179	5252	12249	9152	15682	8499	6008	7801	9790
Number of splices: Non-canonical	7080	5407	5801	7799	5756	4355	9871	8901	10933	5153	3764	6609	7788
Mismatch rate per base, %	0.27%	0.28%	0.29%	0.28%	0.28%	0.28%	0.31%	0.28%	0.29%	0.29%	0.28%	0.27%	0.28%
Deletion rate per base	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%	0.02%
Deletion average length	1.47	1.47	1.47	1.47	1.48	1.47	1.46	1.47	1.48	1.50	1.48	1.48	1.47
Insertion rate per base	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%	0.01%
Insertion average length	1.61	1.52	1.43	1.43	1.63	1.60	1.47	1.43	1.62	1.42	1.38	1.44	1.59
Number of reads mapped to multiple loci	2398878	3287402	2498197	3093230	2029231	1702575	4120922	3568843	3990452	2580456	1853980	3009727	2483263
% of reads mapped to multiple loci	5.21%	8.05%	5.73%	6.45%	5.15%	5.02%	5.71%	6.70%	4.80%	5.64%	6.11%	6.49%	4.84%
Number of reads mapped to too many loci	222693	207795	314715	343251	343822	264218	498200	397171	262921	547542	254116	316223	407743
% of reads mapped to too many loci	0.48%	0.51%	0.72%	0.72%	0.87%	0.78%	0.69%	0.75%	0.32%	1.20%	0.84%	0.68%	0.79%
Number of reads unmapped: too many mismatches	1269	1809	1322	1446	1105	1801	2452	1565	1799	1824	1173	1452	1228
% of reads unmapped: too many mismatches	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
Number of reads unmapped: too short	759644	674719	771365	950053	788663	614921	1410179	1081555	902701	1286763	815091	960438	1181777
% of reads unmapped: too short	1.65%	1.65%	1.77%	1.98%	2.00%	1.81%	1.96%	2.03%	1.09%	2.81%	2.69%	2.07%	2.30%
Number of reads unmapped: other	16971	16710	20173	24305	15726	13423	29546	25946	17237	24268	22820	24661	12953
% of reads unmapped: other	0.04%	0.04%	0.05%	0.05%	0.04%	0.04%	0.05%	0.02%	0.05%	0.08%	0.05%	0.03%	0.03%
Number of chimeric reads	0	0	0	0	0	0	0	0	0	0	0	0	0
% of chimeric reads	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%

## Compile Counts

Generate a directory called “COUNTS” and copy all gene count files to this folder.

```
mkdir COUNTS
cp *pass1/*PerGene* COUNTS/
```

Clean filenames

```
ls *tab|while read line ; do mv $line ${line/GTFpass1ReadsPerGene.out/}
; done
```

For each count file, extract and create the five count tables

```
ls *.tab | while read line ; do
echo $line
cat $line | awk 'NR==3{print}' | cut -f 2- > ${line/tab/nofeature\.tab}
cat $line | awk 'NR==4{print}' | cut -f 2- > ${line/tab/ambiguous\.tab}
cat $line | awk 'NR>4{print}' | cut -f 2 > ${line/tab/nostrand\.tab}
cat $line | awk 'NR>4{print}' | cut -f 3 > ${line/tab/sense\.tab}
cat $line | awk 'NR>4{print}' | cut -f 4 > ${line/tab/antisense\.tab}
done
```

Make a Geneid list from one of the count tables as  
“countsannot\_GRCh38.primary.Selected.Geneid.txt”

```
ls 008iP22TGFbM_S71.tab | head -1 | while read line; do
cut -f 1 $line | awk 'NR>4{print}' >
countsannot_GRCh38.primary.Selected.Geneid.txt
done
```

Create a file listing the names of all samples as “RBarretTNFATGFBsamples.txt”

```
ls *.sense.tab | sed 's/\.sense.tab//g' | tr -s " " "\n" | sed 's/_1//g'
> RBarretTNFATGFBsamples.txt
```

Make count tables for sense, anti-sense, nostrand, ambiguous, and nofeature reads

```
# Combine all sense counts into RBarretTNFATGFB_sense.ALL.cnt
paste *.sense.tab > RBarretTNFATGFB_sense.ALL.cnt
```

```
# Combine all antisense counts into RBarretTNFATGFB_antisense.ALL.cnt
paste *.antisense.tab > RBarretTNFATGFB_antisense.ALL.cnt
```

```
# Combine all nostrand counts into RBarretTNFATGFB_nostrand.ALL.cnt
paste *.nostrand.tab > RBarretTNFATGFB_nostrand.ALL.cnt
```

```
# Combine all ambiguous counts into RBarretTNFATGFB_ambiguous.cnt
cat *ambiguous.tab > RBarretTNFATGFB_ambiguous.cnt
```

```
# Combine all nofeature counts into RBarretTNFATGFB_nofeature.cnt
cat *nofeature.tab > RBarretTNFATGFB_nofeature.cnt
```

I am going to use “RBarretTNFATGFB\_antisense.ALL.cnt” file for the further analysis



## In MATLAB

### Transfer data to your local laptop

Read counts, annotation, and mappability

```
RBarretTNFATGFBcnt = textread('RBarretTNFATGFB_antisense.ALL.cnt','');  
RBarretsamplesTNFATGFB = textread('RBarretTNFATGFBsamples.txt','%s');  
RBarretsampleskeysTNFATGFB = textread('samplekeys_Sam.txt','%s');  
  
% Calculates the sum of the counts in RBarretTNFATGFBcnt, divides the  
result by 1000000, and rounds the result to the nearest integer.  
RBarretTNFATGFBmeta_seqdepth=round(sum(RBarretTNFATGFBcnt)/1000000);
```

### Import annotation

You can find these annotation in the “mappability and R code” folder

```
Gencode_33_Selected_MappSS=textread('mappability and R  
code/gencode.v33.Selected.ReadsPerGene.out.MappSS.txt','');  
Gencode_33_Selected_MappUS=textread('mappability and R  
code/gencode.v33.Selected.ReadsPerGene.out.MappUS.txt','');  
Gencode_33_Selected_Geneid=textread('mappability and R  
code/gencode.v33.annotation.Selected.geneid.txt','%s\n');  
Gencode_33_Selected_Biotype=textread('mappability and R  
code/gencode.v33.annotation.Selected.biotype.txt','%s\n');  
Gencode_33_Selected_Genename=textread('mappability and R  
code/gencode.v33.annotation.Selected.genename.txt','%s\n');
```

### Compile counts

First, initializes a new variable called RBarretTNFATGFBTPM with the same count data as RBarretTNFATGFBcnt. Then, iterates over each gene in the count data matrix. For each gene, the corresponding row in RBarretTNFATGFBTPM is updated by dividing the count data by the read counts from the “Gencode\_33\_Selected\_MappSS”, multiplying by 1000, and storing the result in RBarretTNFATGFBTPM.

Finally, iterates over each sample in the TPM data matrix. For each sample, the corresponding column in RBarretTNFATGFBTPM is updated by dividing the

values in the column by the sum of the values in the column, multiplying by 1,000,000, and storing the result in RBarretTNFATGFBTPM. This step **normalizes the TPM values** across samples and scales the resulting values to TPM.

```
RBarretTNFATGFBTPM = RBarretTNFATGFBcnt;
for i=1:size(RBarretTNFATGFBcnt,1)
% Divids the gene count matrix RBarretTNFATGFBcnt by
Gencode_33_Selected_MappSS matrix, which is the sum of the transcript
length of each gene
RBarretTNFATGFBTPM(i,:) =
RBarretTNFATGFBcnt(i,+)/Gencode_33_Selected_MappSS(i)*1000;
end
% Sets any NaN or Inf values resulting from the normalization process to
0
RBarretTNFATGFBTPM(isnan(RBarretTNFATGFBTPM)) = 0;
RBarretTNFATGFBTPM(isinf(RBarretTNFATGFBTPM)) = 0;
for i=1:size(RBarretTNFATGFBTPM,2)
% Scale the TPM values so that the sum of expression values across each
sample of the matrix is equal to 1,000,000. This ensures that the
expression values are comparable across different samples and allows
meaningful comparisons of gene expression levels between different
samples.
RBarretTNFATGFBTPM(:,i) =
RBarretTNFATGFBTPM(:,i)/sum(RBarretTNFATGFBTPM(:,i))*1000000;
end
```

## Make my first dendrogram

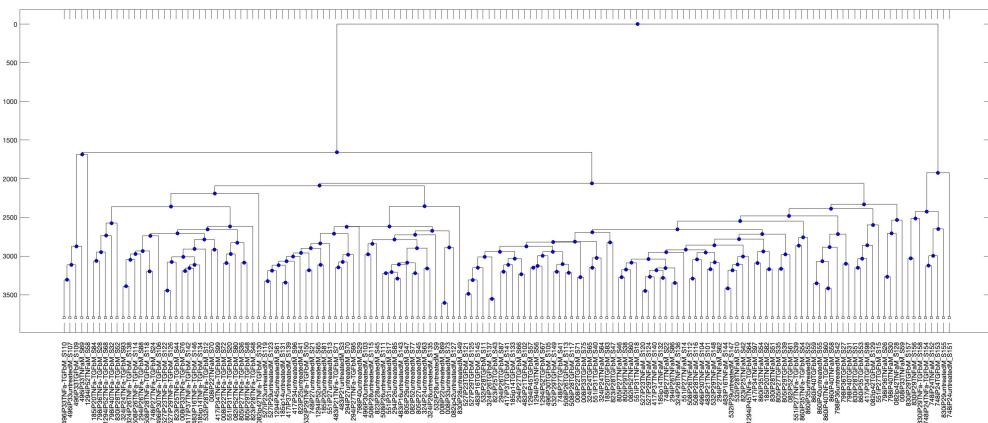
Make a dendrogram to visualize the relationships among samples in the RBarretTNFATGFB dataset based on their gene expression profiles.

1. I generates a random selection of 1,000 genes from the TPM data matrix.
2. Calculates the pairwise distances between the selected genes.
3. Creates a for loop that iterates 9,999 times. For each iteration, a new random selection of 1,000 genes is generated, and the pairwise distances between these genes are added to the previous 'thisdist' calculation.
4. Converts the one-dimensional distance vector 'thisdist' into a distance matrix 'thisdistmat' using the 'squareform' function.
5. Generates a hierarchical clustering tree based on the distance matrix

'thisdistmat'.

Overall, I perform a clustering analysis on a subset of genes in the RBarretTNFATGFB dataset to visualize the relationships among samples based on their gene expression profiles.

```
thisrand = unique(randi([1 size(RBarretTNFATGFBTPM,1)],1,1000));  
thisdist = pdist(RBarretTNFATGFBTPM(thisrand,:));  
for i=1:9999  
    thisrand = unique(randi([1 size(RBarretTNFATGFBTPM,1)],1,1000));  
    thisdist = thisdist+pdist(RBarretTNFATGFBTPM(thisrand,:));  
end  
thisdistmat = squareform(thisdist/10000);  
thistree = seqlinkage(thisdistmat,'average', RBarretsamplesTNFATGFB)  
plot(thistree, 'ORIENTATION', 'top')
```



## All biotypes counts percents

This part of codes is performing all biotypes counts percents across multiple samples.

```
% Used the unique function and stored the allbiotypes variable.  
allbiotypes = unique(Gencode_33_Selected_Biotype);  
  
% Created A cell array allbiotypeslength to store the lengths of each  
% biotype name.  
allbiotypeslength = cell(length(allbiotypes),1);
```

```

% Two new matrices, allbiotypescounts and allbiotypescountspereents, are
initialized with zeros. These matrices have dimensions (number of unique
biotypes) x (number of samples in the TPM data). They will be used to
store the number of reads (counts) and the percentage of total reads
(%TPM) for each biotype in each sample.
allbiotypescounts =
zeros(length(allbiotypes),size(RBarretTNFATGFBCnt,2));
allbiotypescountspereents =
zeros(length(allbiotypes),size(RBarretTNFATGFBCnt,2));

for i=1:length(allbiotypes)
% Found all the indices of Gencode_33_Selected_Biotype that match the
current biotype. Then, returned a vector of **indices** where the
biotype occurs in Gencode_33_Selected_Biotype.
temp = strmatch(allbiotypes{i}, Gencode_33_Selected_Biotype);
% Stored the length of the temp vector represents the number of genes
with the current biotype in the Gencode_33_Selected_Biotype.
allbiotypeslength{i} = length(temp);
if length(temp)>1
% Sum the expression values for all genes with the current biotype
across all samples. The resulting sums are stored in the corresponding
row of the allbiotypescounts matrix.
allbiotypescounts(i,:) =
sum(RBarretTNFATGFBCnt(strmatch(allbiotypes{i},
Gencode_33_Selected_Biotype),:));
allbiotypescountspereents(i,:) =
allbiotypescounts(i,:)./sum(RBarretTNFATGFBCnt)*100;
end
end

dlmwrite('allbiotypescountspereents.txt',
allbiotypescountspereents,'delimiter','\t')
writetable(cell2table(allbiotypes),'allbiotypes.txt','WriteVariableName
s',0)

```

I Formulated a spreadsheet by using “allbiotypes.txt” and “allbiotypescountspereents.txt” on Excel and add min, max, and average for each biotype. You can find my spreadsheet [here](#). As you can see, protein\_coding genes have the average of 98.65% among other biotypes, which is what we want.

Biotypes	min	max	average
IG_C_gene	0.000000000	0.000446580	0.000013340
IG_D_gene	0.000000000	0.000000000	0.000000000
IG_J_gene	0.000000000	0.000013813	0.000000235
IG_V_gene	0.000000000	0.000073058	0.000010604
Mt_rRNA	0.228350000	1.165600000	0.487934371
Mt_tRNA	0.000141280	0.002049800	0.000492376
TEC	0.011838000	0.029033000	0.018976556
TR_C_gene	0.000000000	0.002667100	0.000287637
TR_D_gene	0.000000000	0.000011542	0.000000580
TR_J_gene	0.000000000	0.000076445	0.000009219
TR_V_gene	0.000000000	0.000061956	0.000011904
lncRNA	0.537580000	1.612500000	0.825408079
miRNA	0.002352100	0.014811000	0.004202784
misc_RNA	0.000641580	0.002443400	0.001048973
<b>protein_coding</b>	<b>97.675000000</b>	<b>99.144000000</b>	<b>98.658132450</b>
rRNA	0.000000000	0.000028828	0.000005066
ribozyme	0.000000000	0.000007357	0.000000293
sRNA	0.000000000	0.000004680	0.000000135
scRNA	0.000000000	0.000000000	0.000000000
scaRNA	0.000002616	0.000324330	0.000036518
snRNA	0.000262070	0.001464200	0.000543351
snoRNA	0.001727500	0.004335800	0.002883131
vaultRNA	0.000000000	0.000000000	0.000000000

### Keep only protein coding genes

We are going to keep only protein coding genes for the next step.

```

allbiotypes=unique(Gencode_33_Selected_Biotype);
% Found the 15th unique value, which is protein_coding, of
Gencode_33_Selected_Biotype in the array proteincodingindx.
proteincodingindx = strmatch(allbiotypes{15},
Gencode_33_Selected_Biotype);
biotypeindx = proteincodingindx;

% Created an array additionalgenes contains the indices of genes that
have certain prefixes such as 'MT-', 'H1', 'H2', 'H3', 'H4', 'RPL', or
'RPS' in their names.
additionalgenes = [strmatch('MT-',Gencode_33_Selected_Genename) ;
strmatch('H1',Gencode_33_Selected_Genename);
strmatch('H2',Gencode_33_Selected_Genename);
strmatch('H3',Gencode_33_Selected_Genename);
strmatch('H4',Gencode_33_Selected_Genename) ;
strmatch('RPL',Gencode_33_Selected_Genename) ;
strmatch('RPS',Gencode_33_Selected_Genename)];

% Created an array nonadditionalgenes with the same length as the
Gencode_33_Selected_Genename array.
nonadditionalgenes = 1:length(Gencode_33_Selected_Genename);
% Removed the indices of genes in additionalgenes from the
nonadditionalgenes array.
nonadditionalgenes(additionalgenes) = [];

% mappableindx containing the indices of elements in the
Gencode_33_Selected_MappSS array that are greater than 50.
mappableindx = find(Gencode_33_Selected_MappSS>50);

% a new variable finalIndexGeneric which is the intersection of three
other variables: biotypeindx, nonadditionalgenes, and mappableindx.
finalIndexGeneric =
intersect(biotypeindx,intersect(nonadditionalgenes,mappableindx));
% finds the indices of rows in RBarretTNFATGFBCnt that have a sum
greater than 150.
countindx = find(sum(RBarretTNFATGFBCnt')>150);

% updates finalIndexGeneric to be the intersection of finalIndexGeneric
and countindx.
finalIndexGeneric=intersect(finalIndexGeneric,countindx);

```

```
% creates a new variable RBarretTNFATGFBCnt_GMask which is a subset of  
RBarretTNFATGFBCnt corresponding to the rows indexed by  
finalIndexGeneric.
```

```
RBarretTNFATGFBCnt_GMask = RBarretTNFATGFBCnt(finalIndexGeneric,:);
```

```
Gencode_33_Selected_Geneid_GMask =
```

```
Gencode_33_Selected_Geneid(finalIndexGeneric);
```

```
Gencode_33_Selected_Genename_GMask =
```

```
Gencode_33_Selected_Genename(finalIndexGeneric);
```

```
Gencode_33_Selected_MappSS_GMask =
```

```
Gencode_33_Selected_MappSS(finalIndexGeneric);
```

```
Gencode_33_Selected_MappUS_GMask =
```

```
Gencode_33_Selected_MappUS(finalIndexGeneric);
```

```

% normalizes the expression data like we do previously
RBarretTNFATGFBExpression_GMask = RBarretTNFATGFBCnt_GMask;
for i=1:size(RBarretTNFATGFBExpression_GMask,2)
RBarretTNFATGFBExpression_GMask(:,i) =
RBarretTNFATGFBCnt_GMask(:,i)/sum(RBarretTNFATGFBCnt_GMask(:,i))*100000
0;
end
for i=1:size(RBarretTNFATGFBExpression_GMask)
RBarretTNFATGFBExpression_GMask(i,:) =
RBarretTNFATGFBExpression_GMask(i,+)/Gencode_33_Selected_MappSS_GMask(i
)*1000;
end
RBarretTNFATGFBExpression_GMask(isnan(RBarretTNFATGFBExpression_GMask))
= 0;
RBarretTNFATGFBExpression_GMask(isinf(RBarretTNFATGFBExpression_GMask))
= 0;

% RBarretTNFATGFB_CPM_GMask contains the expression data normalized only
by CPM, using the same normalization method as the code above.
RBarretTNFATGFB_CPM_GMask = zeros(size(RBarretTNFATGFBCnt_GMask));
for i=1:size(RBarretTNFATGFBCnt_GMask,2)
RBarretTNFATGFB_CPM_GMask(:,i) =
RBarretTNFATGFBCnt_GMask(:,i)/sum(RBarretTNFATGFBCnt_GMask(:,i))*100000
0;
end

% RBarretTNFATGFBTPM_GMask contains the expression data normalized only
by TPM.
RBarretTNFATGFBTPM_GMask = RBarretTNFATGFBCnt_GMask;
for i=1:size(RBarretTNFATGFBCnt_GMask,1)
RBarretTNFATGFBTPM_GMask(i,:) =
RBarretTNFATGFBCnt_GMask(i,+)/Gencode_33_Selected_MappSS_GMask(i)*1000;
end
RBarretTNFATGFBTPM_GMask(isnan(RBarretTNFATGFBTPM_GMask)) = 0;
RBarretTNFATGFBTPM_GMask(isinf(RBarretTNFATGFBTPM_GMask)) = 0;
for i=1:size(RBarretTNFATGFBTPM_GMask,2)
RBarretTNFATGFBTPM_GMask(:,i) =
RBarretTNFATGFBTPM_GMask(:,i)/sum(RBarretTNFATGFBTPM_GMask(:,i))*100000
0;
end

```

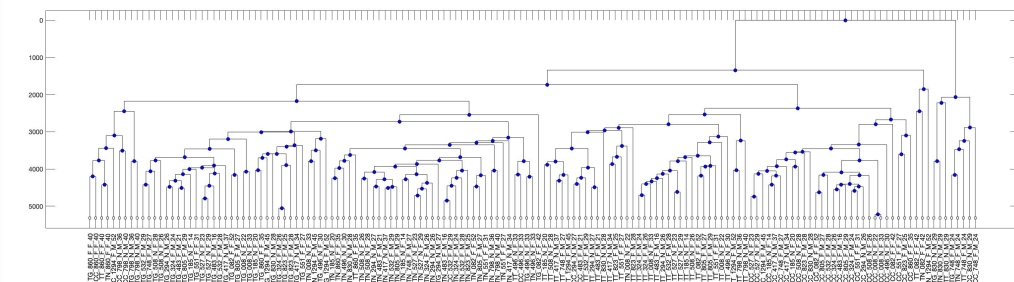


## Make dendrogram with only protein coding genes

Perform hierarchical clustering on a subset of the gene expression data stored in the variable `RBarretTNFATGFBTPM_GMask` with only protein coding genes.

```
thisrand = unique(randi([1 size(RBarretTNFATGFBTPM_GMask,1)],1,1000));
thisdist = pdist(RBarretTNFATGFBTPM_GMask(thisrand,:));
for i=1:9999
    thisrand = unique(randi([1 size(RBarretTNFATGFBTPM_GMask,1)],1,1000));
    thisdist = thisdist+pdist(RBarretTNFATGFBTPM_GMask(thisrand,:));
end
thisdistmat = squareform(thisdist/10000);
thistree = seqlinkage(thisdistmat,'average',
RBarretsampleskeysTNFATGFB)
plot(thistree,'ORIENTATION','top')
```

We can tell that basically the plot is clustered by their treatments.



## What is the percent of the top 100 genes

Calculates the top 100 expressed genes in each sample based on their transcript per million (TPM) values in the `RBarretTNFATGFBTPM_GMask` matrix.

```
% iterates over 151 samples, it first sorts the TPM values of all genes
in descending order and stores the indices of the sorted genes in y. The
top 100 expressed genes in the sample are obtained by selecting the
first 100 indices in y, and these indices are appended to a running list
of all top 100 indices yall.
yall=[];
for i=1:151
[x y]=sort(RBarretTNFATGFBTPM_GMask(:,i),'descend');
yall=unique([y(1:100); yall]);
top100percent(i)=sum(RBarretTNFATGFBTPM_GMask(y(1:100),i))/1000000;
end
```

After calculating the sum of the percent of the top 100 genes, it is 58.25%.

In the end, we store the **Gencode\_33\_Selected\_Geneid\_GMask.txt**, **Gencode\_33\_Selected\_Genename\_GMask.txt**, **Gencode\_33\_Selected\_MappSS\_GMask.txt**, **RBarretTNFATGFBTPM\_GMask.txt**, and **RBarretTNFATGFBCnt\_GMask.txt** for our further analysis in R.

```
writetable(cell2table(Gencode_33_Selected_Geneid_GMask), 'Gencode_33_Selected_Geneid_GMask.txt', 'WriteVariableNames', 0)
writetable(cell2table(Gencode_33_Selected_Genename_GMask), 'Gencode_33_Selected_Genename_GMask.txt', 'WriteVariableNames', 0)
dlmwrite('Gencode_33_Selected_MappSS_GMask.txt',
Gencode_33_Selected_MappSS_GMask, 'delimiter', '\t')
dlmwrite('RBarretTNFATGFBTPM_GMask.txt',
RBarretTNFATGFBTPM_GMask, 'delimiter', '\t')
dlmwrite('RBarretTNFATGFBCnt_GMask.txt',
RBarretTNFATGFBCnt_GMask, 'delimiter', '\t')
```

## In R

### Install packages

We have to install BiocManager, BiocLite, IHW, DESeq2, and ggplot2. Then, read in RBarretTNFATGFBCnt\_GMask.txt, RBarretTNFATGFBsamples.txt, samplekeys\_Sam.txt, and Gencode\_33\_Selected\_Genename\_GMask.txt.

```

setwd("/Users/LuC/Desktop/Cedars-Sinai/PROJECTS/IBD_RNASeq/RBARRETTNFATGFB/")#setwd("/Users/samuellu/Desktop/Cedars-Sinai/PROJECTS/IBD_RNASeq/RBARRETTNFATGFB/")if
(!requireNamespace("BiocManager", quietly = TRUE))
install.packages("BiocManager")#BiocManager::install("BiocLite")#BiocManager::install("IHW")#BiocManager::install("DESeq2")#install.packages("ggplot2")library(DESeq2)library(IHW)library(ggplot2)library(ggrepel)RBarretTNFATGFBcntGMask =
as.matrix(read.table("RBarretTNFATGFBcnt_GMask.txt"))sampleNameTNFATGFB
= as.matrix(read.table("RBarretTNFATGFBsamples.txt"))sampleKeyTNFATGFB
= as.matrix(read.table("samplekeys_Sam.txt"))genenames =
as.matrix(read.table("Gencode_33_Selected_Genename_GMask.txt"))``

#### Generate samplekeys_Sam.tab

```

I have to separate samplekeys\_Sam.txt by "\_" to get samplekeys\\_Sam.tab before next step. Here are my code in terminal.

## In terminal

## Create an empty file to store the output

```
touch samplekeys\_Sam.tab
```

## Loop over the sample names and split them by "\_"

```
for sample in $(cat samplekeys\_Sam.txt); do IFS=_ read -r col1 col2 col3 col4 col5 <<< "$sample" echo -e "$col1\t$col2\t$col3\t$col4\t$col5" >> samplekeys\_Sam.tab done
```

```
#### Generate a sampleTableTNFATGFB
```

```
##### Generate a sampleTableTNFATGFB
```

The sampleTableTNFATGFB contains Treatment, Line, Pheno, Sex, Pass, Factor, and Batch.

```
```sampleTableTNFATGFB =  
read.table("samplekeys_Sam.tab") rownames(sampleTableTNFATGFB)<-  
sampleKeyTNFATGFBcolnames(sampleTableTNFATGFB)<-  
c("Treatment","Line","Pheno","Sex","Pass") sampleTableTNFATGFB$Factor <-  
paste(sampleTableTNFATGFB$Treatment,sampleTableTNFATGFB$Pheno,sep="_")#  
concatenating the "Line" and "Pass" columns with an underscore  
separatorsampleTableTNFATGFB$Batch <-  
paste(sampleTableTNFATGFB$Line,sampleTableTNFATGFB$Pass,sep="_") colname  
s(RBarretTNFATGFBcntGMask) <-  
sampleKeyTNFATGFBwrite.table(sampleTableTNFATGFB,file="sampleTableTNFAT  
GFB.txt", sep = "\t", col.names = FALSE)````  

```

<br>

```
#### DESeq2 package
```

The RBarretTNFATGFBcntGMaskBatch is created with the **\*\*Batch\*\*** information specified in the design formula, while the RBarretTNFATGFBcntGMaskFactor is created with the **\*\*treatment** and **phenotype\*\*** information specified in the design formula.

Next, the DESeq function is used to estimate size factors and dispersion values for the DESeqDataSet objects.

Finally, the varianceStabilizingTransformation function is used to perform variance stabilizing transformation on the DESeqDataSet objects. This transformation is important for reducing the effect of noise and heteroscedasticity in the data, making it more suitable for downstream analyses such as differential gene expression analysis.

```
```RBarretTNFATGFBcntGMaskBatch <-  
DESeqDataSetFromMatrix(RBarretTNFATGFBcntGMask, colData=  
sampleTableTNFATGFB,design= ~Batch)RBarretTNFATGFBcntGMaskBatch <-  
DESeq(RBarretTNFATGFBcntGMaskBatch)RBarretTNFATGFBcntGMaskFactor <-  
DESeqDataSetFromMatrix(RBarretTNFATGFBcntGMask, colData=
```

```
sampleTableTNFATGFB,design= ~Factor)RBarretTNFATGFBcntGMaskFactor <-
DESeq(RBarretTNFATGFBcntGMaskFactor)RBarretTNFATGFBcntGMaskBatch_vsd <-
varianceStabilizingTransformation(RBarretTNFATGFBcntGMaskBatch,blind=FA
LSE)RBarretTNFATGFBcntGMaskFactor_vsd <-
varianceStabilizingTransformation(RBarretTNFATGFBcntGMaskFactor,blind=F
ALSE)``
```

#### #### PCA

```
``# performs principal component analysis (PCA) on the variance-
stabilized counts data pcabatch <-
prcomp(t(assay(RBarretTNFATGFBcntGMaskBatch_vsd)))# gives the
percentage of variance explained by each principal
componentpercentVarbatch <-
round(100*pcabatch$sdev^2/sum(pcabatch$sdev^2))# pcabatch$rotation is a
matrix containing the loadings of the principal components.aloadbatch <-
abs(pcabatch$rotation)
# normalizes the loadings in aloadbatch so that each column (i.e., PC)
sums to 1. aloadrelativebatch <- sweep(aloadbatch, 2,
colSums(aloadbatch), "/")# pcabatch$x is a matrix containing each
sample's coordinate on each principal componentpcabatchALL <-
pcabatch$xpcabatchR<- cbind(pcabatchALL,sampleTableTNFATGFB)# centers
the PC1 scores in pcabatchR to have a mean of 0. This is done so that
the PC1 variable can be used as a covariate in the subsequent
differential expression analysis.pcabatchR$PC1 <- scale(pcabatchR$PC1,
center = TRUE)RBarretTNFATGFBcntGMaskPC1 <-
DESeqDataSetFromMatrix(RBarretTNFATGFBcntGMask, colData=
pcabatchR,design= ~PC1)RBarretTNFATGFBcntGMaskPC1 <-
DESeq(RBarretTNFATGFBcntGMaskPC1)RBarretTNFATGFBcntGMaskPC1_vsd <-
varianceStabilizingTransformation(RBarretTNFATGFBcntGMaskPC1,blind=FALS
E)``
```

#### #### PCA plots

The first plot shows the relationship between **PC1** and **PC2** colored by **Pheno** variable, with the point size indicating the **Treatment** variable.

The second plot is similar to the first, but the data points are colored by the **Sex** variable, and the point size indicates the **Treatment** variable.

The third plot shows the relationship between PC2 and PC9 colored by Pheno variable, with the point size indicating the Treatment variable.

```
```ggplot(pcabatchR, aes(PC1, PC2, color= Pheno)) +
geom_point(aes(size= Treatment),alpha=0.6,stroke =
3)+geom_point(aes(size= Pheno),color="black",alpha=0.2) +
xlab(paste0("PC1: ",percentVarbatch[1],"% variance")) +
ylab(paste0("PC2: ",percentVarbatch[2],"% variance")) +
geom_text_repel(aes(label = sampleKeyTNFATGFB),size=4,box.padding =
0.35, point.padding = 0.5,segment.color = 'grey50')+
theme_bw()ggplot(pcabatchR, aes(PC1, PC2, color= Sex)) +
geom_point(aes(size= Treatment),alpha=0.6,stroke =
3)+geom_point(aes(size= Treatment),color="black",alpha=0.2) +
xlab(paste0("PC1: ",percentVarbatch[1],"% variance")) +
ylab(paste0("PC2: ",percentVarbatch[2],"% variance")) +
theme_bw()ggplot(pcabatchR, aes(PC2, PC9, color= Pheno)) +
geom_point(aes(size= Treatment),alpha=0.6,stroke =
3)+geom_point(aes(size= Treatment),color="black",alpha=0.2) +
xlab(paste0("PC2: ",percentVarbatch[2],"% variance")) +
ylab(paste0("PC9: ",percentVarbatch[9],"% variance")) +
geom_text_repel(aes(label = sampleKeyTNFATGFB),size=4,box.padding =
0.35, point.padding = 0.5,segment.color = 'grey50')+ theme_bw()```
```

In the PCA plot, there are 4 different groups separate by treatment. We have the CC group(untreated) at the right corner, the TG group(TGFB) at the bottom, TN group(TNFA) at the right corner, and the TT group(TGFB + TNFA) at the top. There are separate by PC1(19%).![]

(/Users/samuellu/Desktop/Cedars-

Sinai/PROJECTS/GitHub/Pics/PCA\_Pheno\_Treatment.jpeg)<br>#### A series of bar plots (one for each principal component)

Each bar plot represents the loadings of all samples on a given principal component.

```
```# The resulting vector coul will contain 12 colors from the "Set3"
palette.library(RColorBrewer)coul <- brewer.pal(12, "Set3")# generates
colors for a plot based on the batch
variablecolors=pcabatchR$Batchallbatches<-unique(pcabatchR$Batch)for (i
in 1:38){ colors[pcabatchR$Batch==allbatches[i]]<-
coul[i%12+1]}thinlines=c(seq(4,72,8),75,seq(83,151,8))thicklines=c(seq
(8,72,8),79,seq(87,151,8))# first half of the barplot would be the non-
```

```

fibrotic group and the second part would be the fibrotic group# The
order would be CC, TG, TN,
TTsamplesorder=c(4,1,3,2,8,5,7,6,28,25,27,26,32,29,31,30,36,33,35,34,40
,37,39,38,52,49,51,50,55,53,55,54,68,65,67,66,72,69,71,70,76,73,75,74,8
0,77,79,78,84,81,83,82,88,85,87,86,116,113,115,114,120,117,119,118,139,
136,138,137,143,140,142,141,12,9,11,10,16,13,15,14,20,17,19,18,24,21,23
,22,44,41,43,42,48,45,47,46,60,57,59,58,64,61,63,62,92,89,91,90,96,93,9
5,94,100,97,99,98,104,101,103,102,108,105,107,106,112,109,111,110,124,1
21,123,122,128,125,127,126,132,129,131,130,135,133,134,147,144,146,145,
151,148,150,149)#create 38 barplots and saving each of them as a PNG
filefor (i in 1:38) { filename = paste("PC_",i,".png", sep = "")
png(filename)
barplot(pcabatchALL[samplesorder,i],col=colors[samplesorder],las=2,xaxt
='n',space=0) for (i in 1:length(thinlines)) { abline(v =
thinlines[i], col = "black",lty = 3) } for (i in 1:length(thicklines))
{ abline(v = thicklines[i], col = "black",lty = 1) } abline(v = 72,
col = "red",lty = 1)
dev.off()}write.csv(aloadrelativebatch,file="aloadrelativeMask_batchmod
el_filtered.csv")write.csv(pcabatch$x,file="pca_batchmodel_x.csv")```

```

A red line is drawn at position 72 in order to separate the non-fibrotic group and the fibrotic group. Within each patient, the treatment order would be CC, TG, TN, TT. The color of each bar represents the batch of the sample, with a unique color assigned to each batch. The vertical lines on the plot indicate the position of specific loadings, with thin and thick lines indicating different positions.

For example, this is PC\_1.png. From this plot, you can see that the highest expression is closely related with TNF- $\alpha$ . As for the first patient, compare to the control(untreated), the TNF- $\alpha$  group is much higher and the TT group (TGF- $\beta$ +TNF- $\alpha$ ) is not that high.

```



```

#### PCA rank matrix

Takes csv files and converts it to the txt files with the second column onwards. It does this by first removing the first row using awk, replacing multiple commas with tabs using tr, and removing the first column using cut.

## In Mac terminal

```
cat aloadrelativeMask_batchmodel_filtered.csv | awk 'NR>1{print}' | tr -s “,” “ |  
cut -f 2- > aloadrelativeMask_batchmodel_filtered.clean.txt cat  
pca_batchmodel_x.csv | awk 'NR>1{print}' | tr -s”, ” | cut -f 2- >  
pca_batchmodel_x.clean.txt
```

Read in the preprocessed data files created in the previous steps and store them in variables `pcabatch_samples` and `pca_loadings`, respectively.

## In Matlab

```
pcabatch_samples = textread('pca_batchmodel_x.clean.txt',''); pca_loadings =  
textread('aloadrelativeMask_batchmodel_filtered.clean.txt','');
```

Sort the three columns of `pca_loadings` in descending order and store the sorted values in variables `x1`, `x2`, and `x3`, and the corresponding indices in `y1`, `y2`, and `y3`.

```
%%%%%%%%%% PCA SUPP %x = pca_loading number, y = its index [x1  
y1]=sort(pca_loadings(:,1),'descend'); [x2 y2]=sort(pca_loadings(:,2),'descend');  
[x3 y3]=sort(pca_loadings(:,3),'descend');
```

Determine the rank of each row in the original order for the first three principal components and store the ranks in a matrix `pcarankmatrix`.

```
[x y z]=intersect(1:length(y1),y1); pcarankmatrix(:,1)=z; [x y  
z]=intersect(1:length(y2),y2); pcarankmatrix(:,2)=z; [x y  
z]=intersect(1:length(y3),y3); pcarankmatrix(:,3)=z;
```

```
% contains the rank of each feature in the original order for the first three  
principal components dlmwrite('pcarankmatrix.txt', pcarankmatrix,'delimiter',')
```



#### #### Spreadsheet

By making a spread sheet, we can easily manage our data by a single glance. I used command line, Excel and R to orginize it.

The first sheet built on Excel contains patients order, patients id, phenotypes, and sex. You can visit the sheet by clicking [here] (/spreadsheet/Barret\_Myofibroblast\_TGFTNF\_MASTER.xlsx).

```
![(/Users/samuellu/Desktop/Cedars-Sinai/PROJECTS/GitHub/Pics/spreadsheet_patient.png)
```

<br>

The second sheet includes the names and percentages of all biotypes, along with their respective minimum, maximum, and average values, providing us with a comprehensive overview.

## In terminal

```
paste allbiotypes allbiotypescountspercent > combine_allbiotypes_percent.txt
```

## In R#

```
allbiotypes_percentsheet2_1 <-  
list("Biotypes")sheet2_2 <-  
sampleKeyTNFATGFBcombined_s  
heet2 <- c(sheet2_1,  
sheet2_2)combined_spreadsheet2  
<-
```

```
as.matrix(read.table("combine_all  
biotypes_percent.txt"))colnames(  
combined_spreadsheet2) <-  
combined_sheet2write.table(com  
bined_spreadsheet2,file="combin  
ed_spreadsheet2.txt", sep = "  
row.names = FALSE)
```

**add their respective minimum,  
maximum, and average values on  
Excel**

```
![(/Users/samuellu/Desktop/Cedars-  
Sinai/PROJECTS/GitHub/Pics/spreadsheet_allbiotypes_percent.png)
```

<br>

The third sheet includes Genename, Geneid, Mapp, PC1, PC2, PC3, and patient's TPM values.

## **In terminal**

```
paste Gencode_33_Selected_Genename_GMask.txt  
Gencode_33_Selected_Genename_GMask.txt  
Gencode_33_Selected_Geneid_GMask.txt  
Gencode_33_Selected_MappSS_GMask.txt pcarankmatrix.txt >  
combine_test.txt # In R # spreadsheetsheet3_1 <-  
list("Genename","Genename","Geneid","Mapp","PC1","PC2","PC3")sheet3_2<-
```

```
sampleKeyTNFATGFBcombined_headers <- c(sheet3_1,
sheet3_2)combined_spreadsheet <-
as.matrix(read.table("combine_test.txt"))colnames(combined_spreadsheet) <-
combined_headerscombined_spreadsheet <-
combined_spreadsheet[order(combined_spreadsheet[,1]),] #sort by the first
columnwrite.table(combined_spreadsheet,file="combined_spreadsheet.txt", sep
= ", row.names = FALSE)```
```

In Excel, we can sort the spreadsheet with PC1, PC2, and so on to see the correlation between the treatment and the expression level in each gene.

[illegible]

## Future works

## References