# Assignment 3

1. First, I headed to NCBI website to get info about the data and get SRA ids so that I can import then to galaxy portal.
   https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE29968
   https://www.ncbi.nlm.nih.gov/Traces/study/?acc=PRJNA140847&o=acc_s%3Aa



NCBI Data base page



SRA page with annotation of dataset

2. After I headed to galaxy portal at https://usegalaxy.org/ to import data from option **Faster Download and Extract Reads in FASTQ**.



As you can see in the history data is downloaded after typing SRR ids from SRA page in NCBI.
Approx. 4 GB data was downloaded from NCBI which took around 3 hours

3. After downloading datasets in fastq files, quality report is generated using **FastQC.**



**Inference**



## ✅ Basic Statistics

| Measure | Value |
|---|---|
| Filename | SRR278174.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 15589765 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 38-58 |
| %GC | 54 |

6 reports were made as our dataset contain 3 dataset of normal tissue and 3 for infected tissue.
Basic Statistics part of fastqc report contains basic info like what is the size of reads, GC percentage in all reads, how reads were recorded etc. Here in our example there were about 15million read of size from 38 to 58 which has 54% of GC content as it is recorded with illumina it also contains quality score and 0 is of poor quality.

## Per base sequence quality



Per base sequence quality part of fastqc report contains info of box plot of quality score at every base position of read and our read length is max 58 so x axis of 58 length, there is a trend in all report that per base quality score decreases as position of read increases

## Per sequence quality scores



Per sequence quality score part of fastqc report contains frequency of mean sequence phred score i.e. average of phred score is taken and then these scores are grouped by intervals to get frequency of each grp. In our example we can see that about 8M sequence or mostly has phred score around 33.

## Per base sequence content



Per base sequence content part of fastqc report contains info of base at position of read. It gives us percentage of base at every position of sequence or reads. We can see that in the initial content are mixed up which also shows that there is bias that why there is warning sign before heading in graph.

**Per sequence GC content**



Per sequecnce GC content part of fastqc report shows that what should be distribution of GC across all sequence and what is in the data. Blue line represent theoretical data line or we can ideal situation and red line represent our data from both line it shown that most of the sequence should contain about 50% percent of GC content. In our example about 1.2 M sequence have 50% GC content.

**Sequence Length Distribution**



Sequence length distribution part of fastqc contain info of length of reads or sequence. It shows frequency of length of read in data. Here data either has length of 38 or 58. This may be because of primer selection.

**Sequence Duplication Levels**



Sequence duplication levels contains Info of duplication in the data, red line represent number of distinct sequence that are duplicated similarly blue line represent count of all sequence.

4. After downloading datasets in fastq files, reads were aligned using **Bowtie2** using inbuild reference genome hg19 version as it is suitable for the study.





Bowtie2 took around 6 hours to complete

Default parameters were selected in the bowtie2 except reference genome and history visibility. Bowtie2 is fast and take low memory as compared to other alignment tools like hisat2 etc. In these reads were align to reference genome to get count matrix in the subsequent steps.

```
25541142 reads; of these:
  25541142 (100.00%) were unpaired; of these:
    2254424 (8.83%) aligned 0 times
    5827495 (22.82%) aligned exactly 1 time
    17459223 (68.36%) aligned >1 times
91.17% overall alignment rate
```

This is stat of one of the data after alignment it shows that around 68% of read aligned more than 1 time and 22% of sequence exactly one time and 8% of sequence are not aligned, these sequences can be adapter sequence which is used in illumina.

5. Download annotation file from ensembl.
   Data from https://www.gencodegenes.org/human/release_19.html
   First data was download into computer then on galaxy portal



6. Using htseq-count tool count matrix was created and annotation or gtf file was upload in the previous step.
   Union mode was selected and min 10 quality was selected. BAM file from bowtie2 result was selected as input.

7. **Metadata checklist**
   a) Reads originate from **Homo sapiens.**
   b) Yes, reference genome is available from Homo sapiens that is download from ensembl which was hg19 version.
   c) GPL10999 Illumina Genome Analyzer IIx (Homo sapiens) platform was used.
   d) These are short reads
   e) **ID list**

| SRA ID | GEO ID | Sample ID | Type |
|---|---|---|---|
| SRR278173 | GSM741690 | 16N | Non-Tumour |
| SRR278174 | GSM741691 | 17N | Non-Tumour |
| SRR278175 | GSM741692 | 19N | Non-Tumour |
| SRR278176 | GSM741693 | 16T | Tumour |
| SRR278177 | GSM741694 | 17T | Tumour |
| SRR278178 | GSM741695 | 19T | Tumour |

8. From this step and onward all steps were taken in R code. For Differential Gene Expression analysis DESeq2 was used in R. Before coding in R I downloaded count file from galaxy in tabular form in my PC.

```r
#####################
getwd()
# change the working directory
setwd("E:/Project/RNA-seq")
getwd()
#####################




#####################
# TO do filtering and DE analysis on HTseq data
library(DESeq2)
#####################



#####################
#To take files from HTseq as count matix and covert to DESeq data type
#Take all files name
CM_Files <- grep("Tumour",
                 list.files("Count_matrix/"),
                 value=TRUE)

CM_Samples <- c("Non_Tumour_1","Non_Tumour_2","Non_Tumour_3","Tumour_1","
Tumour_2","Tumour_3")

#Take all samples here Tumour and Non-Tumour
CM_Condition <- sub("(.*Tumour).*",
                    "\\1",
                    CM_Files)

#Make table which is in the form of DE Seq data
Table <- data.frame(sampleName = CM_Files,
                    fileName = CM_Files,
                    condition = CM_Condition)
Table$condition <- factor(Table$condition)

#Make DESeqDataSet data
DESeq_Data <- DESeqDataSetFromHTSeqCount(sampleTable = Table,
                                         directory = "Count_matrix/",
                                         design= ~ condition)
#####################
```

**Directory E:/Project/RNA-seq** working site
**Directory E:/Project/RNA-seq/Count_matrix** count tables were stored

Files were download and renamed as Tumour_1 that represent tumour count matrix in first patient similarly Non-Tumout_1 represent first patient non tumour rna-seq count matrix from galaxy which is generated after htseq-count.

Then in the previous code condition and files were import with htseq function in deseq2

After getting data model from DESeq2 was fitted to get DESeq type file, which also make count file in single file which make it easy to calculate dispersion, size factors etc.
Base mean, LogFC(fold change), Pvalue, FDR value was calculated b passing DESeq type data in results function.
As data contains ID from ensembl we have to convert it into normal gene symbol by using org.Hs.eg.db library.
LogCPM or log of count per million is not calculated using results so it calculated using edgeR library.

```
    #Differenctial gene expressoion
    DESeq <- DESeq(DESeq_Data)
    DESeq$type<-c('single-read','single-read','single-read','single-
read','single-read','single-read')
    DESeq$type<-factor(c('single-read'))


    DESeq_result <- results(DESeq,
                    pAdjustMethod = "BH",
                    alpha = 0.1)


    ####################
    #convert gene names
    library("org.Hs.eg.db")
    DESeq_result$hgnc_symbol <- mapIds(org.Hs.eg.db,
                              keys=gsub("\\..*","",row.names(DESeq_r
esult)),

                              column="SYMBOL",
                              keytype="ENSEMBL",
                              multiVals="first")
    DESeq_result$entrezid <- mapIds(org.Hs.eg.db,
                              keys=gsub("\\..*","",row.names(DESeq_result
)),

                              column="ENTREZID",
                              keytype="ENSEMBL",
                              multiVals="first")
```

Head of DESeq_result matrix

```
> head(DESeq_result)
log2 fold change (MLE): condition Tumour vs Non Tumour
Wald test p-value: condition Tumour vs Non Tumour
DataFrame with 6 rows and 10 columns
                  baseMean log2FoldChange    lfcSE      stat    pvalue        padj hgnc_symbol  entrezid    logCPM       SD
                 <numeric>      <numeric> <numeric> <numeric> <numeric>   <numeric> <character> <character> <numeric> <numeric>
ENSG00000000003.10  647.0206      -2.710219  0.504565 -5.371394 7.81304e-08 1.87102e-06       TSPAN6        7105   4.86742  0.768782
ENSG00000000005.5     0.0000            NA        NA        NA        NA          NA         TNMD       64102      -Inf      NaN
ENSG00000000419.8    95.9206      -0.258414  0.385290 -0.670701 5.02411e-01 6.87636e-01         DPM1        8813   3.41146  0.400908
ENSG00000000457.9    18.2084      -0.354243  0.655616 -0.540321 5.88976e-01 7.52701e-01        SCYL3       57147   1.69829  0.488435
ENSG00000000460.12   20.3621       1.159525  0.674801  1.718322 8.57378e-02 2.09620e-01     C1orf112       55732   1.71792  1.017300
ENSG00000000938.8    20.8955       1.779838  0.925642  1.922815 5.45033e-02 1.51882e-01          FGR        2268   1.43100  1.228274
```

Following graph are created in R and saved in PNG file which are then attached here

**Standard deviation vs average logCPM distribution of data**



As we can see that most of the gene are highly expressed and logCPM as mostly gene lie around 2 and 4.

**MA-plot**



This graph shows us dispersion of fold change ie if logFC is 2 then gene expression is 4 times when condition is changed i.e. from non-tumour to tumour and x axis contain normalised count of count of genes.

**PCA plot**



PCA plot is calculated using regularized log transformation is DESeq2 function PC1 is good in clustering group in condition while PC2 is not. As left cluster is of non-tumour read and right is of tumour read, data is separated greatly in x axis but no in y axis.

**Sample to sample distance**



This heatmap shows the relative distance between 2 sample, it show how much 2 sample are similar if square is dark blue then distance is 0 it means they are 100% similar, as similarity decreases square colours become light blue.

**Dispersion Estimation**



This dispersion graph shows how gene should be align, red line is the expected dispersion value for genes of a given expression strength and black dots is a gene with maximum likelihood estimation (MLE) of dispersion and blue dots are in the expected dispersion value to filter data we can make data so that every dot is red or blue dots.

**Histogram of pvalue vs frequency**



X axis is divided into 50 division and each bar represent frequency of that particular p value and following graph represent data set which are select (p-value < 0.05)

This graph represent mean of normalised of selected dataset i.e. which have p value less than 0.05 against ratio of small p value.

Following code snippet was used to filter out genes

```
#gene selection
filtered_genes <- as.data.frame(DESeq_result[order(DESeq_result$padj),])

filtered_genes <- filtered_genes[!(filtered_genes$baseMean==0),]
filtered_genes <- filtered_genes[!is.na(filtered_genes$pvalue),]
filtered_genes <- filtered_genes[!is.na(filtered_genes$padj),]
filtered_genes <- filtered_genes[!is.na(filtered_genes$hgnc_symbol),]
filtered_genes <- filtered_genes[!(filtered_genes$logCPM==-Inf),]

filtered_genes <- filtered_genes[(filtered_genes$pvalue<0.05),]
filtered_genes <- filtered_genes[(filtered_genes$padj<0.1),]

filtered_genes <- filtered_genes[(filtered_genes$logCPM>2),]

filtered_genes <- filtered_genes[(filtered_genes$log2FoldChange>4 | fil-
tered_genes$log2FoldChange< -4),]
```

Conditions were
1. P-value < 0.05
2. FDR value using Benjamini Hochberg method < 0.1
3. Log(CPM) >2
4. Log(FC) >4 or <-4

Dataset was sorted by FDR value and row with NA value was removed.

This gives us 184 genes

9. Gene ontology was performed using GoFuncR function.

```
#gene ontology
library(GOfuncR)
library(Homo.sapiens)


gene_ids = c(filtered_genes$hgnc_symbol)
input_hyper = data.frame(gene_ids, is_candidate=1)
res_hyper = go_enrich(input_hyper, n_randset=100)
```

Above code is executed to generate data in GO id terms

```
> ontology
stats$ontology: biological_process
          ontology    node_id                          node_name raw_p_underrep raw_p_overrep FWER_underrep FWER_overrep
6   biological_process GO:0030198      extracellular matrix organization           1   1.459300e-20             1            0
7   biological_process GO:0043062 extracellular structure organization           1   1.577141e-20             1            0
8   biological_process GO:0009888                    tissue development           1   3.924169e-19             1            0
12  biological_process GO:0018149                 peptide cross-linking           1   1.236495e-16             1            0
14  biological_process GO:0043588                       skin development           1   4.981828e-15             1            0
16  biological_process GO:0030154                  cell differentiation           1   2.856382e-13             1            0
17  biological_process GO:0008544                  epidermis development           1   4.229357e-13             1            0
18  biological_process GO:0048869          cellular developmental process           1   8.026372e-13             1            0
19  biological_process GO:0048856          anatomical structure development           1   1.249367e-12             1            0
20  biological_process GO:0007275  multicellular organism development           1   1.297294e-12             1            0
-----------------------------------------------------------------------------------------------------------
stats$ontology: cellular_component
          ontology    node_id                              node_name raw_p_underrep raw_p_overrep FWER_underrep FWER_overrep
1   cellular_component GO:0005615                      extracellular space           1   2.681641e-26             1            0
2   cellular_component GO:0005576                     extracellular region           1   4.431553e-26             1            0
3   cellular_component GO:0062023 collagen-containing extracellular matrix           1   6.423385e-26             1            0
5   cellular_component GO:0031012                     extracellular matrix           1   4.544422e-22             1            0
9   cellular_component GO:0070062                    extracellular exosome           1   4.090721e-18             1            0
10  cellular_component GO:1903561                    extracellular vesicle           1   6.897469e-18             1            0
11  cellular_component GO:0043230                  extracellular organelle           1   7.719615e-18             1            0
15  cellular_component GO:0031982                                  vesicle           1   2.688858e-14             1            0
21  cellular_component GO:0098644            complex of collagen trimers           1   2.078136e-12             1            0
23  cellular_component GO:0001533                         cornified envelope           1   1.023321e-11             1            0
-----------------------------------------------------------------------------------------------------------
stats$ontology: molecular_function
          ontology    node_id                                                  node_name raw_p_underrep raw_p_overrep FWER_underrep FWER_overrep
4   molecular_function GO:0005201                    extracellular matrix structural constituent           1   4.350024e-22             1            0
13  molecular_function GO:0005198                           structural molecule activity           1   1.787739e-15             1            0
36  molecular_function GO:0030020 extracellular matrix structural constituent conferring tensile strength   1   1.224640e-10             1            0
40  molecular_function GO:0061134                      peptidase regulator activity           1   2.560436e-10             1            0
49  molecular_function GO:0004866                     endopeptidase inhibitor activity           1   1.083178e-09             1            0
51  molecular_function GO:0004867          serine-type endopeptidase inhibitor activity           1   1.250417e-09             1            0
54  molecular_function GO:0030414                       peptidase inhibitor activity           1   1.810398e-09             1            0
55  molecular_function GO:0061135                     endopeptidase regulator activity           1   1.810398e-09             1            0
59  molecular_function GO:0005518                              collagen binding           1   1.213656e-08             1            0
64  molecular_function GO:0048407            platelet-derived growth factor binding           1   3.017867e-08             1            0
```

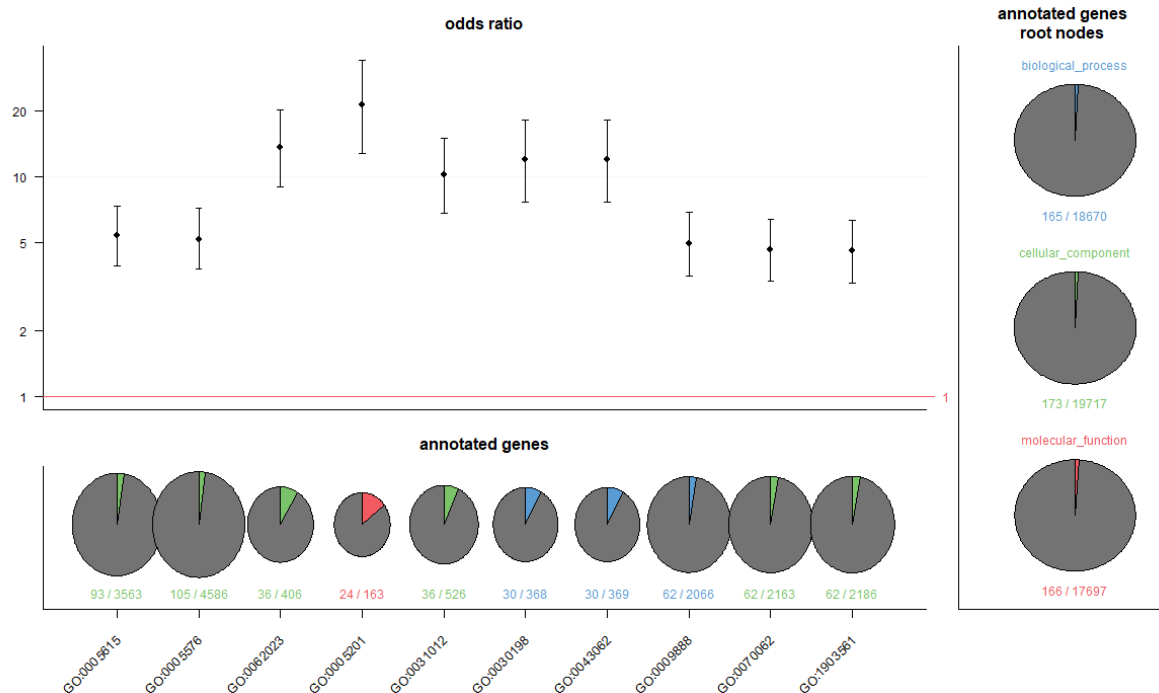Above is the top 10 ontology result which is grouped by gene ontology types

Most of the top gene ontology result shows that genes were in the path of extracellular space and region of skin of vesicle and that is true as we have taken sample from esophageal squamous cell.

Biological processes in this ontology mainly consist of cellular matrix and structure, tissue and skin development processes lead to the reason that these genes can be cancer promoter.

Cellular component of ontology also mainly consists of cellular region and space, as sample is taken from infected cell and normal cell so ontology will be focused of cell area and parts.

Molecular function in these ontologies mainly consist of endopeptidase or peptidase activity which show high activity of breakage of peptide bonds in the cellular area or space.

If we confine all the info then top ontology show that there is high activity of breakage of peptide bond cell area which leads to abnormal development of tissue and cell.
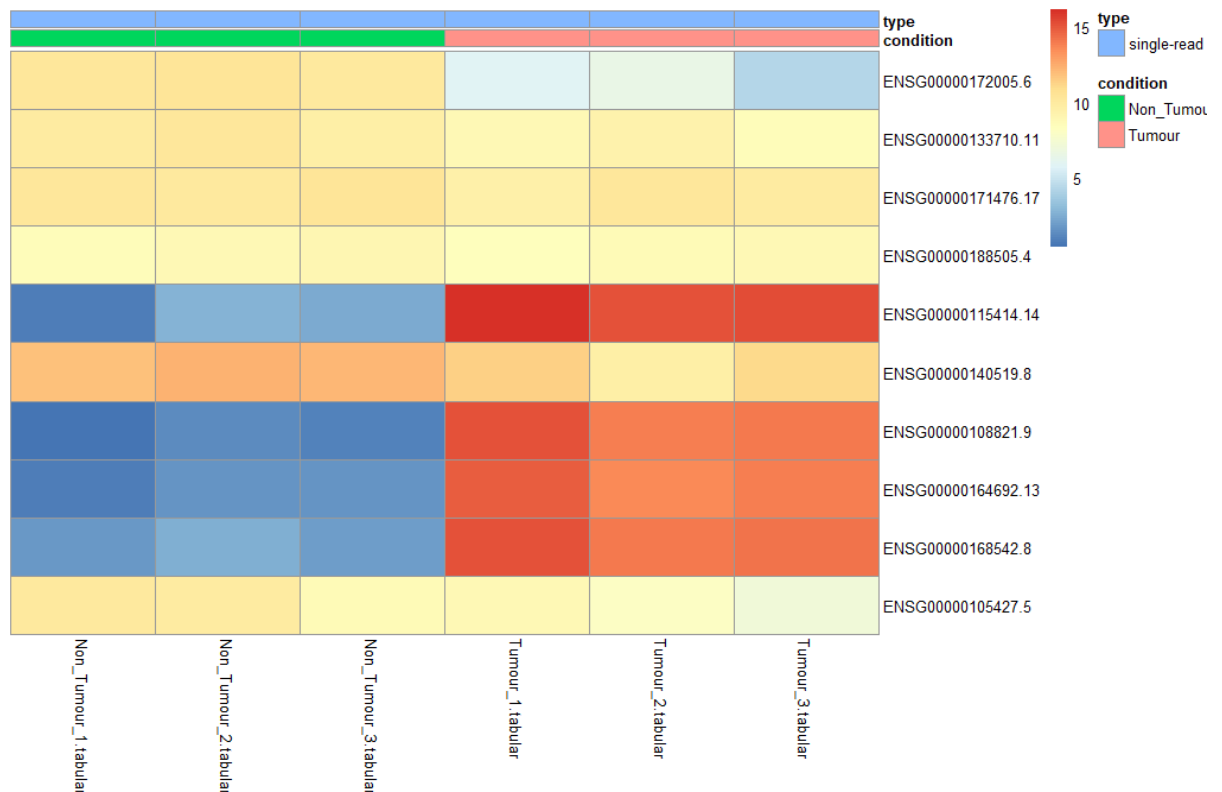
This is graph generated from GoFuncR function, x axis of annotated genes contains name of GO ID which are top of ontology data when sorted by p value.

For eg in first pic chart annotation gene section the amounts of candidate and background genes that are annotated to the GO-categories and the root nodes so 93 candidate gene out of 3563 background genes in cellular component root.

The top panel shows the odds-ratio comparing the GO-categories with their root nodes.

10. Heatmap



This graph was made using top 10 expressing which were calculated when we did ontology process.

From the above heatmap we can clearly see that most of the gene here are highly expressed in tumour cell as compared to non-tumour cell

11. Microarray vs RNA-seq data

| Microarray | RNA-seq |
|---|---|
| Old technology, and it contained non discrete data | New technology, and it contains discrete data that is reads have finite number |
| Most of the old is in the microarray form | Almost all data in the recent is done using RNA-seq |
| Microarray can be done for sample which can be fluorescently tagged | RNA-seq can be done for many sample |
| It requires transcript specific probes | It doesn't require any transcript specific probes |
| It has intensity of light that is captured by camera | It has base with quality score that is captured by camera |
| For high quality of data high number of RNA are required | High quality data can be generated using even low number of RNA |
| In this lots of manual work requires so degree of error by human is very large | While most the work is done by computer so degree of error by human is very low |