# Microarray vs RNA-seq data
# for chronic pancreatitis in humans:
# a comparison of differential gene expression
# results

Video tutorial:

https://youtu.be/S4-i-sPLI2M

GitHub repository:

https://github.com/lalagkaspn/chronic_pancreatitis_omics.git
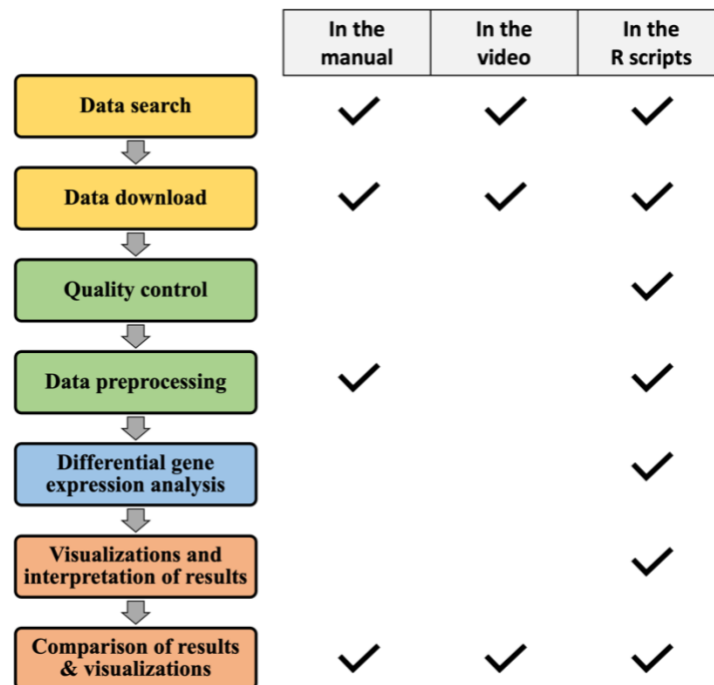
Panagiotis N. Lalagkas

December 2023

# TABLE OF CONTENTS

# 1. Purpose statement

This manual is divided into three parts. First, it demonstrates how to search the [Gene Expression Omnibus](#) (GEO) database for publicly available microarray and RNA-sequencing (RNA-seq) data and how to download it in R. Second, it emphasizes the importance of data preprocessing through specific examples in R. Third, it provides a recommended workflow and visualization techniques to compare two lists of significantly differentially expressed genes obtained after the analysis of microarray or RNA-seq data for the same disease.

Other intermediate steps, such as differential gene expression analysis, are also of high importance. However, they are not the primary focus of this manual and are not mentioned in detail. But, to build this manual, all these steps had to be performed. Therefore, if you want to explore them, you can find all the R scripts used for downloading and preprocessing the gene expression data, conducting the differential gene expression analysis and comparing two lists of differentially expressed genes in this [GitHub repository](#).

The following graphical abstract provides an overview of all the steps typically involved in a differential gene expression analysis. Its purpose is to give you an overall perspective before delving deeper into the present manual. In this illustration, you can also see which of these steps are covered in the manual, the tutorial video, or the R scripts.

| | In the manual | In the video | In the R scripts |
|---|:---:|:---:|:---:|
| **Data search** | ✓ | ✓ | ✓ |
| **Data download** | ✓ | ✓ | ✓ |
| **Quality control** | | | ✓ |
| **Data preprocessing** | ✓ | | ✓ |
| **Differential gene expression analysis** | | | ✓ |
| **Visualizations and interpretation of results** | | | ✓ |
| **Comparison of results & visualizations** | ✓ | ✓ | ✓ |

## 2. Approach

We will search the GEO database for publicly available microarray and RNA-seq data for pancreatitis and normal samples. Then, we will download this data in R, an open-source, powerful programming language designed for statistical computing and data analysis. R offers a wide variety of packages, available in both CRAN and Bioconductor repositories, which we can use for data manipulation, gene expression analysis and visualization.

After downloading the microarray and RNA-seq data, we will preprocess it and perform a differential gene expression analysis by comparing pancreatitis samples to normal samples. Eventually, we will compare two lists of differentially expressed genes derived from the analysis of either microarray or RNA-seq data, create insightful figures and run simple statistical tests.

## 3. Background

The analysis of gene expression data has significantly helped scientists to gain a better understanding of the underlying biology of various diseases[1]. Two of the most widely used techniques for quantifying the expression of genes are microarray and RNA-seq, with the latter being the more modern approach. Consequently, knowing how to search, download, preprocess and analyze these two types of data is a valuable skill for every bioinformatician.

Both microarray and RNA-seq techniques can simultaneously measure the expression of thousands of genes in a single experiment. However, there is one important difference between them. Microarrays rely on pre-designed probes, which can introduce bias to the results as only sequences that hybridize with the probes can be detected. On the other hand, RNA-seq is a sequence-agnostic technique. This means that it can capture the expression of a broader range of genes, and it can detect low-abundance transcripts and novel isoforms[2].

Nevertheless, both microarrays and RNA-seq techniques have been extensively used to improve our understanding of biological processes and the biology of complex diseases. Thus, a plethora of microarray and RNA-seq data is publicly available for many complex diseases. Given that both techniques measure the expression of genes, logical questions are raised: do they give the same results or not? If not, are the differences important and can they be traced back to their potential cause? Previous studies have shown a high correlation between the gene expression profiles generated by the two techniques. Additionally, RNA-seq was found to be better at

detecting a greater number of differentially expressed genes with higher fold-change[3]. But this might differ based on the question under study and the analyzed datasets.

Therefore, the purpose of this manual is to analyze microarray and RNA-seq data of pancreatitis and healthy samples to identify statistically significantly dysregulated genes and compare the results obtained from analyzing microarray and RNA-seq data.

## 4. Installation of R packages

Throughout this manual, we will be using various R packages (from either CRAN or Bioconductor repositories). To install all the necessary packages, please refer to the "packages_installation.R" script on this GitHub repository. Below, you find a list of these packages, their intended use, and hyperlinks to their respective websites.

Packages annotated with "*" are CRAN packages.

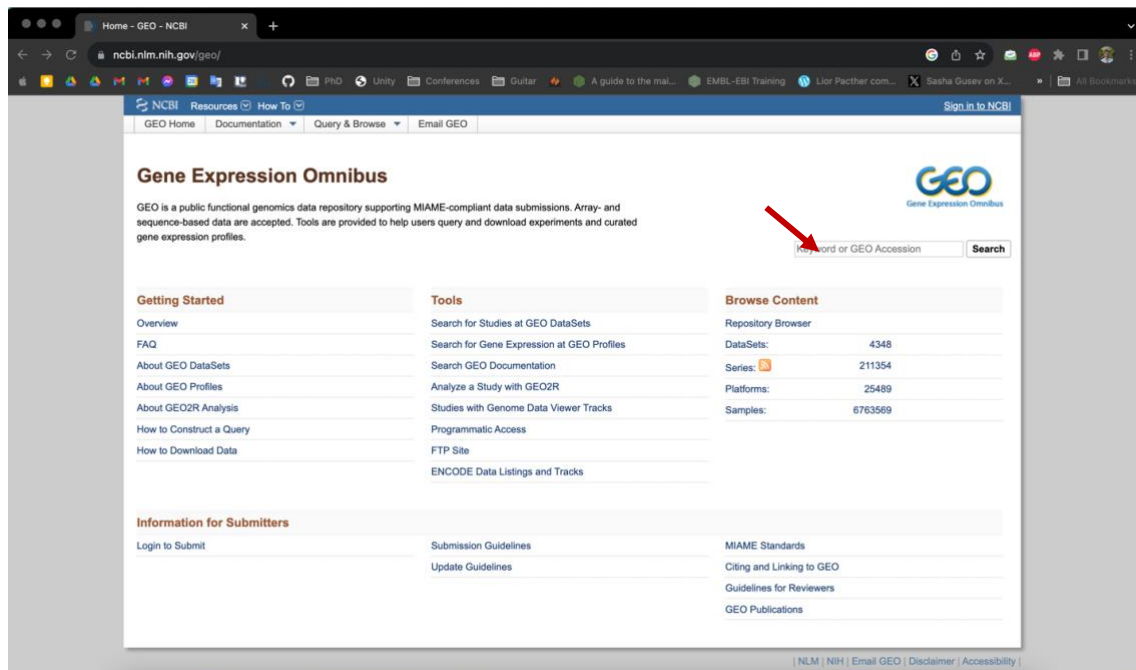Packages annotated with "†" are Bioconductor packages.

- Data manipulation
    - * dplyr: https://cran.r-project.org/web/packages/dplyr/index.html
    - * matrixStats: https://cran.rstudio.com/web/packages/matrixStats/index.html
    - * reshape2: https://cran.r-project.org/web/packages/reshape2/index.html
- Save data into excel file
    - * openxlsx: https://cran.r-project.org/web/packages/openxlsx/index.html
- Download microarray and RNA-seq data from the GEO database
    - † GEOquery:
      https://bioconductor.org/packages/release/bioc/html/GEOquery.html
- Imputation of missing values in gene expression datasets
    - † impute: https://bioconductor.org/packages/release/bioc/html/impute.html
- Conversion of human gene IDs to Entrez IDs:
    - † org.Hs.eg.db:
      https://bioconductor.org/packages/release/data/annotation/html/org.Hs.eg.db.html
- Differential gene expression analysis:
    - † limma: https://bioconductor.org/packages/release/bioc/html/limma.html
- Visualizations:
    - * ggplot2 (various plots): https://ggplot2.tidyverse.org/

- o  * pheatmap (heatmap): https://cran.r-project.org/web/packages/pheatmap/index.html
- o  † EnhancedVolcano (volcano plot): https://bioconductor.org/packages/release/bioc/html/EnhancedVolcano.html
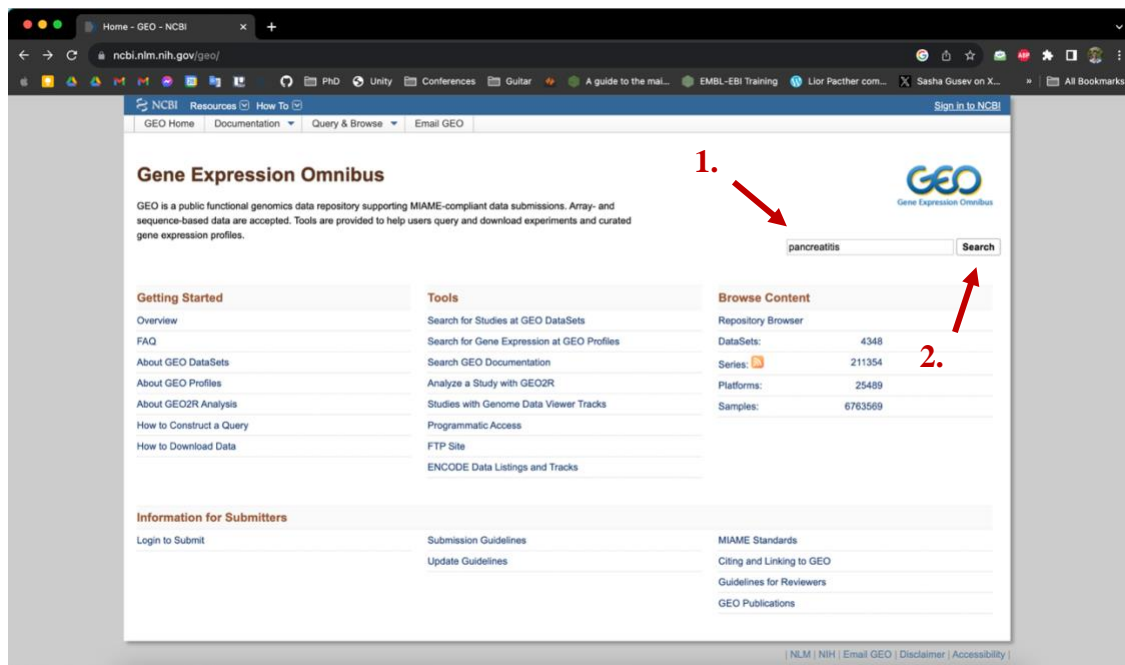
# 5. Gene expression data search in GEO

GEO is a repository for gene expression data from array- and sequencing-based experiments. All the data in this database is publicly available. Below, you can see a screenshot of the GEO landing page (https://www.ncbi.nlm.nih.gov/geo/).

To search for gene expression data, you must provide keywords in the white rectangle next to the "Search" button (as indicated by the red arrow in the screenshot). Examples of keywords that can be used include specific disease terms, the species of interest or the desired method used to generate the gene expression data. In our case, we are interested in searching for gene expression data related to "pancreatitis" in humans and generated using either microarray or RNA-seq method.

Let's begin by searching for "pancreatitis" data. We enter "pancreatitis in the white rectangle box and click "Search".
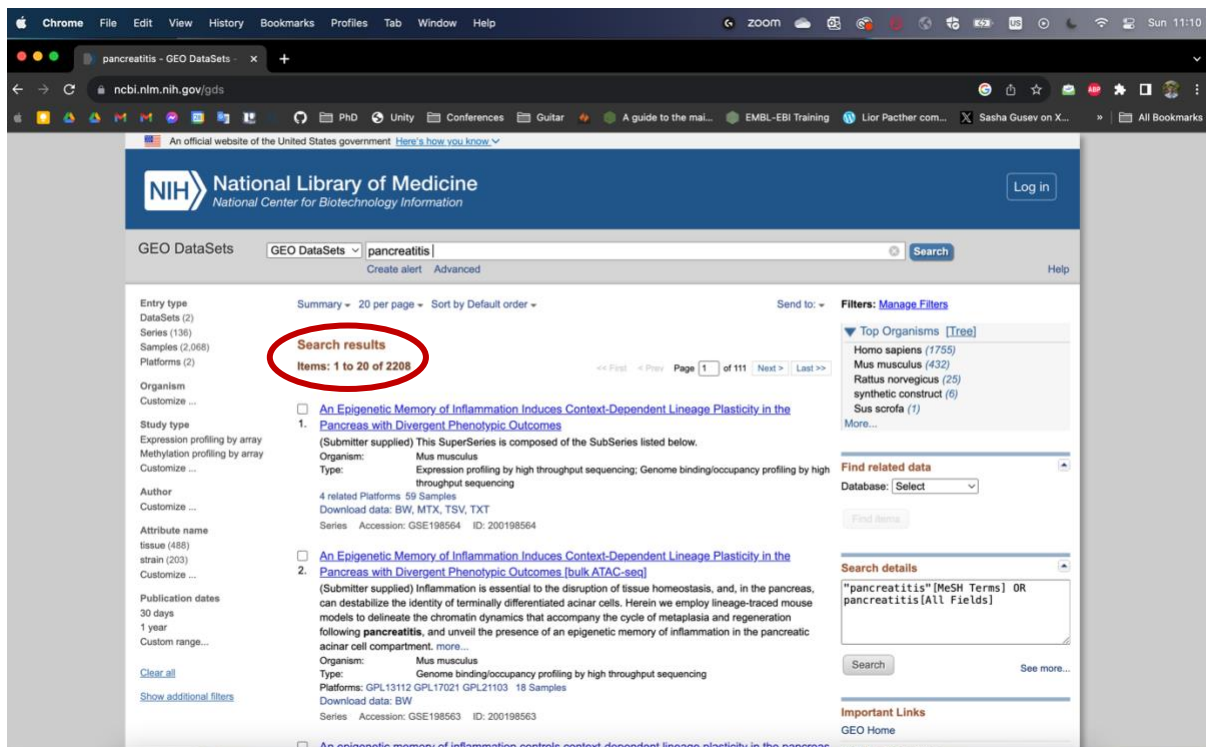


A white box appears (see screenshot below). As you can see, we have the option to search the GEO DataSets or GEO Profiles database. A GEO DataSet is a collection of biologically and statistically comparable samples that were analyzed using the same process and contains the expression of many genes. In contrast, a GEO Profile includes expression measurements for an individual gene across many samples in a GEO DataSet. In our case, we want to look for GEO DataSets. So, we click on the hyperlink for the GEO DataSets (the hyperlink is usually the number of results; so, we click on the "2208"). Note that the number of returned results might change in the future as more datasets will be added from future studies.
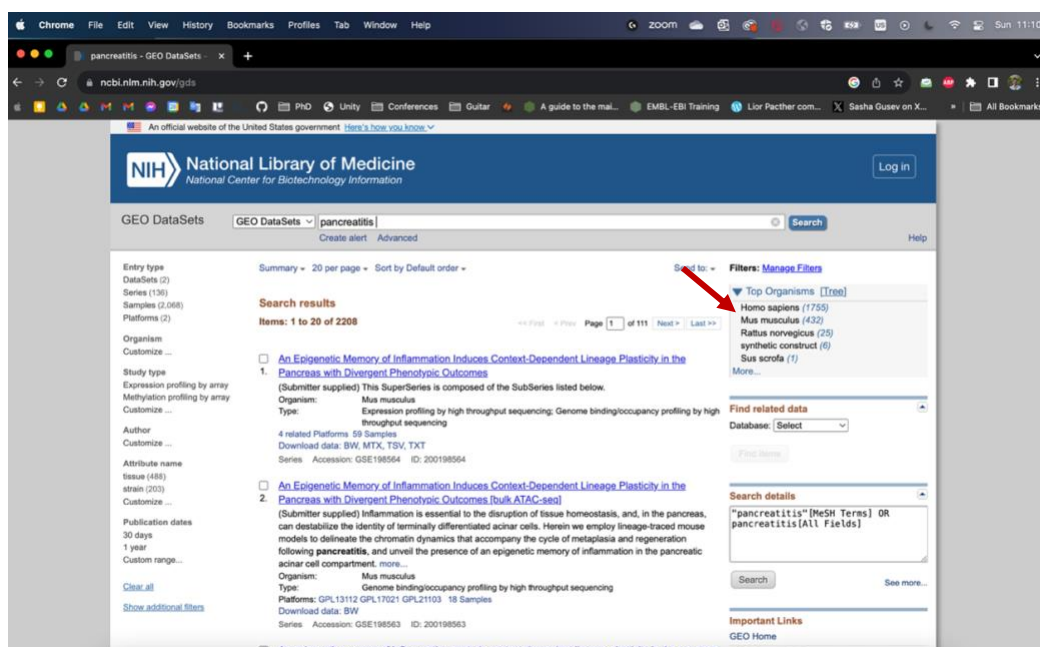


There are 2208 results for "pancreatitis" in the GEO DataSets Database.

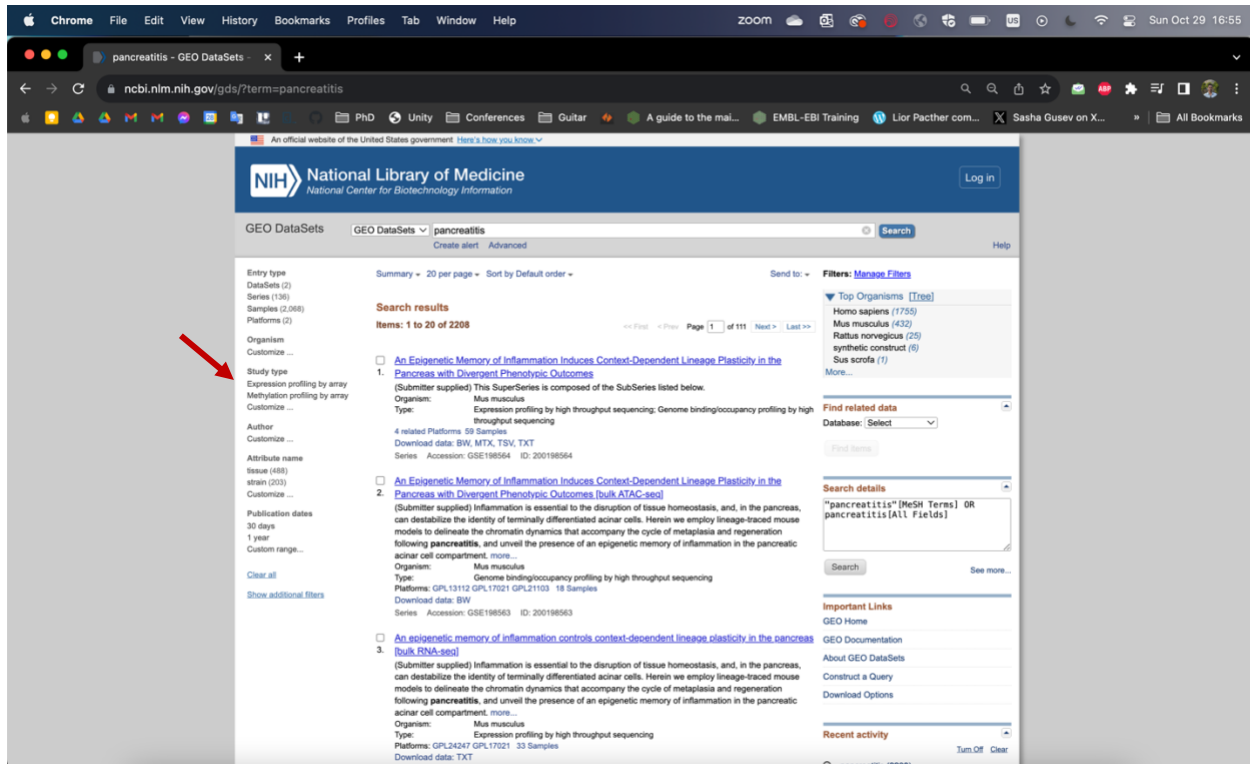There are 22690 results for "pancreatitis" in the GEO Profiles Database.

Now, we can see all the GEO DataSets available in the GEO database and are related to pancreatitis. We have 2,208 results (see red circle below).



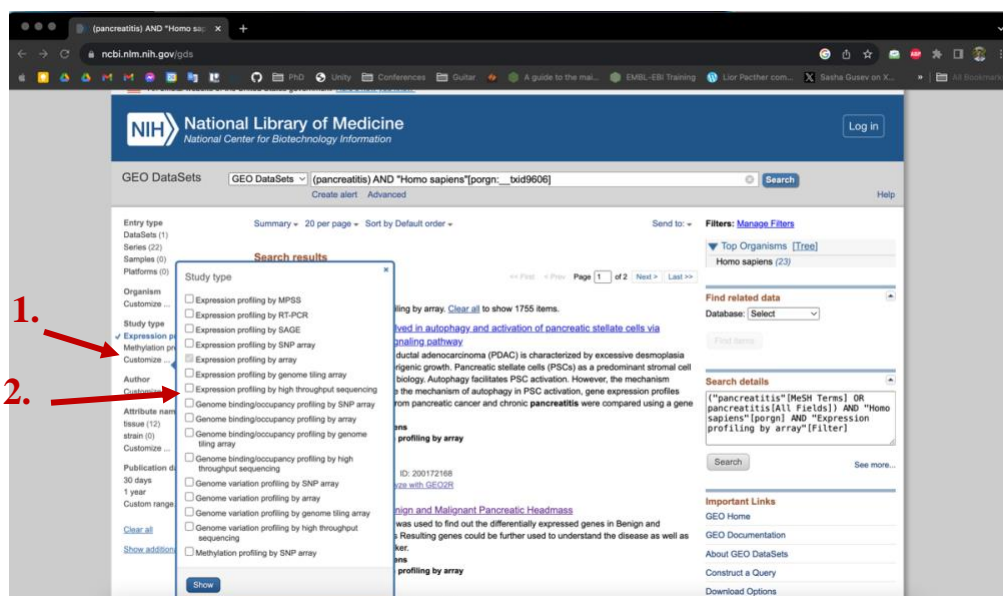Next, we filter the results for those related to humans (Top Organisms = Homo sapiens; see right red arrow below).

To search for **microarray data**, we select the study type to be "Expression profiling by array" (see red arrow below).



To search for **RNA-seq data**, we select the "Study type" to be "Expression profiling by high throughput sequencing" instead. We find this option if we click on "Customize" under the Study type.

Eventually, we find 23 GEO DataSets for pancreatitis in humans generated by microarrays. However, it is possible that some of the returned GEO DataSets do not contain pancreatitis or normal samples. After manual inspection, we keep the following 5 microarray datasets:

- GSE143754: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE143754

- GSE61166: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE61166

- GSE71989: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE71989

- GSE101462: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE101462

- GSE77858: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE77858

Similar to microarrays, we find 59 GEO DataSets for pancreatitis in humans generated by RNA-seq. Again, after manual inspection, we keep the following 2 RNA-seq datasets that contain gene expression data from both pancreatitis patients and normal samples:

- GSE194331: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE194331

- GSE133684: https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE133684

At this point, it is important to acknowledge that our upcoming comparison involves samples from 5 microarray studies and samples from 2 RNA-seq studies. This difference in the number of studies will potentially introduce bias to our analysis. Additionally, these samples stem from different individuals, which might introduce biological variation. These factors are essential to consider when interpreting our analysis results at the end of this manual. Ideally, we would like to have same samples from the same individuals processed using both microarray and RNA-seq methods to minimize bias.

## 6. Gene expression data download in R

### 6.1. Microarray data

To download GEO microarray expression data in R, we use the R package "GEOquery". Given that you have already installed it by following the "package_installation.R" script from this GitHub repository, we load it into R.

```
12  ## Load GEOquery
13  library(GEOquery)
```

We create a vector containing the accession numbers of the GEO microarray DataSets we want to download (the 5 datasets mentioned earlier).

```
16  # Vector with datasets to download from GEO
17  datasets = c("GSE143754", "GSE61166", "GSE71989", "GSE101462", "GSE77858")
```

As we want to download more than one GEO DataSets, we create an empty list named "GEOsets" (row 20), which we populate later with the downloaded data. Using the *getGEO()* function in a for loop, we download all the available data for each GEO DataSet (rows 21-24). Eventually, we name the elements of the list using the corresponding GEO accession numbers (row 26).

```
19  # Download data
20  GEOsets = list()
21  for (i in 1:length(datasets)){
22    GEOsets[[i]] = getGEO(datasets[i])
23  }; rm(i)
24  GEOsets = unlist(GEOsets)
25  # Name the elements of the list with the GEO accession numbers
26  names(GEOsets) = datasets; rm(datasets)
```

After downloading the data, we can click on the "GEOSets" variable in our R environment to inspect it. We can see that our list has 5 elements, each named after a GEO accession number.



Let's investigate one element to see what information is included. As you can see below, it contains information about the experiment (experimentData), the expression counts per probe (assayData), and metadata about the samples (phenoData) and each probe (featureData). Feel free to further explore the data.

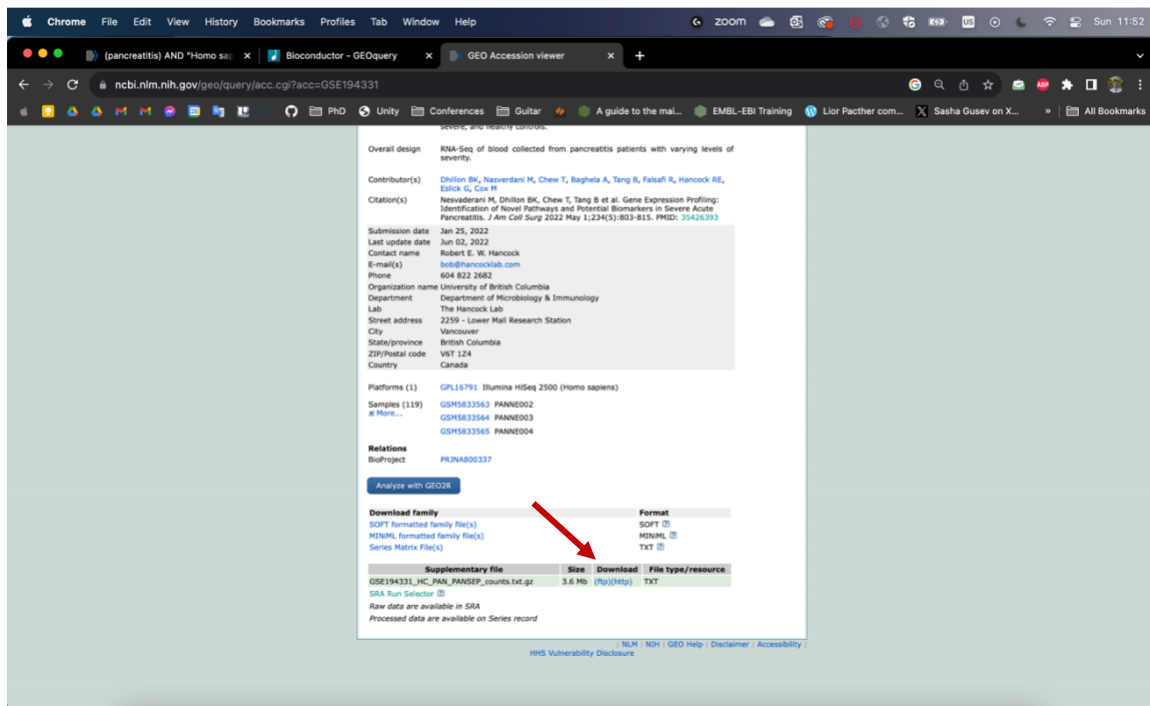| Name | Type | Value | |
|------|------|-------|---|
| GEOsets | list [5] | List of length 5 | |
| GSE143754 | S4 [70523 x 26] (Biobase::Exp | S4 object of class ExpressionSet | |
| experimentData | S4 (Biobase::MIAME) | S4 object of class MIAME | Experiment data |
| name | character [1] | 'Bishnupriya,,Chhatriya' | |
| lab | character [1] | '' | |
| contact | character [1] | 'chhatriya.bishnupriya@gmail.com' | |
| title | character [1] | 'Transcriptome profiling of Benign and Malignant Pancreatic Headmass' | |
| abstract | character [1] | 'Microarray was used to find out the differentially expressed genes in Benign a ... | |
| url | character [1] | 'https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE143754' | |
| pubMedIds | character [1] | '33160365' | |
| samples | list [0] | List of length 0 | |
| hybridizations | list [0] | List of length 0 | |
| normControls | list [0] | List of length 0 | |
| preprocessing | list [0] | List of length 0 | |
| other | list [24] | List of length 24 | |
| .__classVersion__ | list [2] (Biobase::Versions) | List of length 2 | |
| assayData | environment [1] | <environment: 0x28948a318> | Expression counts |
| exprs | double [70523 x 26] | 14.17 13.63 13.45 13.47 12.21 10.27 14.08 13.70 13.54 13.59 12.56 10.88 13.50 11 ... | |
| phenoData | S4 [26 x 41] (Biobase::Annotat | S4 object of class AnnotatedDataFrame | |
| varMetadata | list [41 x 1] (S3: data.frame) | A data.frame with 41 rows and 1 column | |
| data | list [26 x 41] (S3: data.frame) | A data.frame with 26 rows and 41 columns | Phenotypic data |
| dimLabels | character [2] | 'sampleNames' 'sampleColumns' | |
| .__classVersion__ | list [1] (Biobase::Versions) | List of length 1 | |
| featureData | S4 [70523 x 15] (Biobase::Ann | S4 object of class AnnotatedDataFrame | |
| annotation | character [1] | 'GPL17586' | Feature data |
| protocolData | S4 [26 x 0] (Biobase::Annotate | S4 object of class AnnotatedDataFrame | |
| .__classVersion__ | list [4] (Biobase::Versions) | List of length 4 | |

TIP: to access an element of a list in R, you can use the list name followed by double square brackets and the name (or index) of the desired element. To access a specific component of an element, you can use the "@". For example, if you want to access the expression data for the GEO DataSet "GSE143754", you can use the following command:

GEOsets[["GSE143754"]]@**assayData**

## 6.2. RNA-seq data

If we follow the same process as above to download gene expression data for RNA-seq experiments in R, we will notice that the assayData element, which contains the expression values, is empty. This happens because the getGEO function is designed to download expression data from microarrays but not from RNA-seq experiments. Instead, we can download RNA-seq expression data directly from the webpage of each GEO DataSet. For example, to download the RNA-seq expression data for the GEO DataSet "GSE194331", we search the GEO database for "GSE194331" (at the same white box that we searched for pancreatitis earlier). After clicking search, we arrive at the page shown below. This is the webpage of the study we searched for. If

we scroll to the bottom of this webpage, we find the expression data in a "txt" file. We download and save it locally in our laptop by clicking on the hyperlink at the third column (see red arrow).



NOTE: every GEO DataSet has a dedicated webpage similar to the one shown above. This webpage includes all the information related to the study/dataset. All this information is also included in the list we previously viewed about the microarray data (the one we downloaded in R using the getGEO() function). For example, experimentData contains the data presented above the "gray" box in the screenshot. phenoData contains all the information found under the "Samples" section of the webpage. featureData contains all the data provided under the "Platforms" section of the webpage.

## 7. Gene expression data preprocessing in R

Data preprocessing is a crucial step in any data analysis workflow. However, the specific preprocessing steps to apply depend on the data at hand and may vary from case to case. In the context of the present manual, all the preprocessing steps applied to the data can be found in detail in the corresponding github scripts (GSE_microarrays.R rows 44-597; GSE_rnaseq.R rows 49-346). Below, the most important ones are briefly mentioned.

### 7.1. Microarray data

#### 7.1.1. Dealing with missing expression data

Some of the probes in microarray DataSets have missing expression data for some samples (due to technical issues). In such cases, we can impute the missing values by leveraging information from the expression of a gene in other samples. However, when a gene has missing values for many samples, imputation may not be accurate. Therefore, we first remove the probes with missing values in more than 25% of the samples. The 25% threshold was chosen arbitrarily, as no gold-standard guidelines exist. After filtering them out, we use a method called k-NN imputation to impute the missing values in the remaining probes (GSE_microarrays.R rows 301-330).

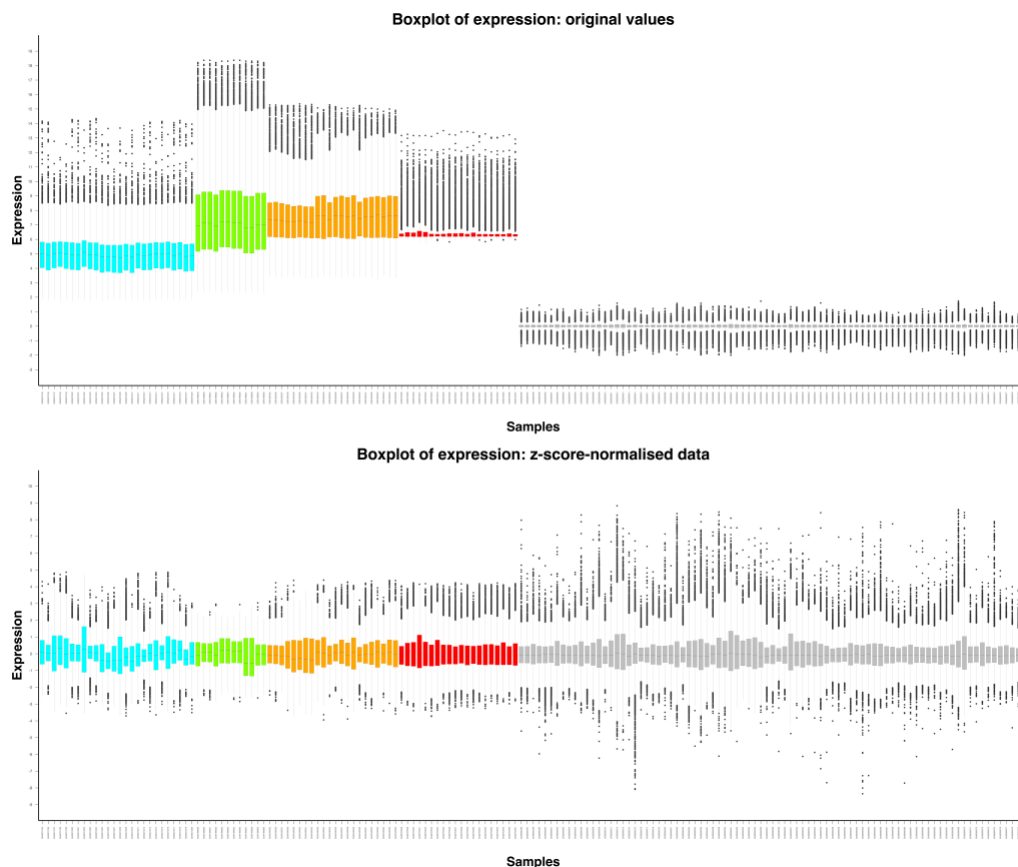#### 7.1.2. Dealing with multiple probes matching to the same gene

Some microarrays include multiple probes that match the same gene. This is intentional and done to capture different parts of a gene as it improves coverage and the accuracy of the calls. However, in downstream analysis we only need one expression value per gene per sample. To address this, we need to combine the information from all the probes that match the same gene. Therefore, for each gene, we find all the matched probes and calculate the variance of their expression values. Eventually, for each gene, we keep the matched probe with the highest variance (GSE_microarrays.R: rows 345-357, 450-522). Please note that various methods have been proposed in the past for this preprocessing step, but this is out of the scope of this manual.

#### 7.1.3. Matching microarray probe sequences to genes

In some cases, such as for the "GSE61166" dataset, we might only have access to the sequence of the probes instead of the gene IDs matching to each probe. To obtain the gene IDs for each probe, so that we know the identity of genes and their expression, we used the BLASTN online tool (https://blast.ncbi.nlm.nih.gov/Blast.cgi). This tool accepts nucleotide sequences as input and outputs RefSeq IDs of genes with sequences that match the input sequences. We only keep RefSeq IDs of genes with sequences that had a 100% match with the provided nucleotide sequences. You can find the code for preparing the BLASTN input data and processing of the returned outputs from BLASTN in the rows 359-449 of the GSE_microarrays.R.

### 7.1.4. Normalization of gene expression counts

The final preprocessing step, before conducting the differential gene expression analysis, typically involves the normalization of the expression data. This step is crucial to ensure that all expression data, as obtained from various experiments, are on the same scale. While there are several normalization methods available, here we employ one of the most widely used: the z-normalization method. This method scales the expression data of all genes per sample to have a mean value of zero and a standard deviation of one. The code for applying the z-normalization method can be found in the rows 539-553 of the GSE_microarrays.R script. To get a better understanding of why normalization of gene expression counts is important, please see below two boxplot figures that show the expression data for each gene per sample before and after the z-normalization method:



Boxplot of expression: original values

Boxplot of expression: z-score-normalised data

As you can see, after the z-normalization (bottom boxplot), the expression values of all genes per sample are centered around zero (same scale), which was not the case before the normalization (top boxplot). This plot, along with other quality control plots can be found in the

[QC/GSE_microarrays subdirectory on GitHub](). All these plots were created using the GSE_microarrays.R script, rows 555-903.
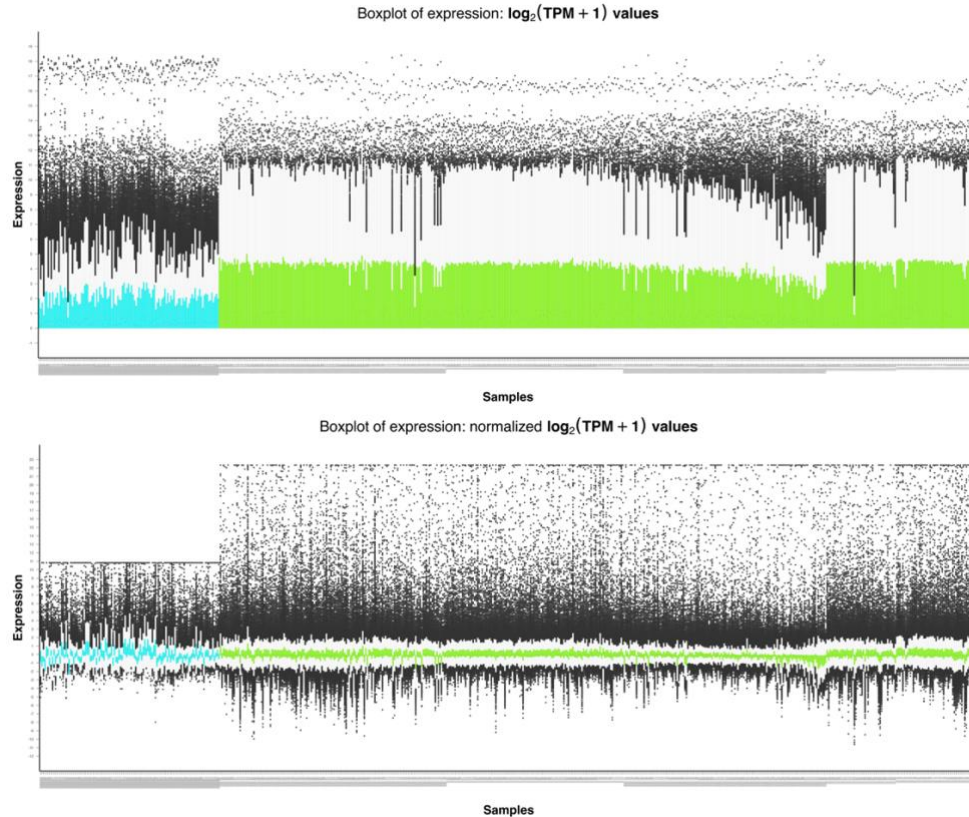
## 7.2.   RNA-seq data

### 7.2.1.   Convert raw counts to TPM

"GSE133684" provided gene expression data in transcripts per million (TPM), while "GSE194331" made available the raw gene expression data (raw counts). Consequently, we need to convert raw counts to TPM, so both datasets are on the same scale. To accomplish this, we first compute the Reads Per Kilobase (RPK) for each gene. This involves obtaining information about gene length (GSE_rnaseq.R rows: 159-221). Then, we calculate the RPK values for each gene by dividing the expression counts by the transcript length (in kilobases). Using the calculated RPK values, we compute the TPM. For those interested, the whole process is implemented and described in detail in rows 223-242 of the GSE_rnaseq.R script.

### 7.2.2.   Normalization of gene expression counts

As with microarray data, the final preprocessing step before conducting the differential gene expression analysis in RNA-seq data is the normalization of the expression data. However, in the case of RNA-seq data, we are working with TPM values rather than counts. Therefore, we first log2 transform the TPM values and then apply z-normalization (GSE_rnaseq.R rows: 280-326). If you want to learn about the rationale behind the log2 transformation, you can read this [thread](). Below, you can see two boxplots showing the expression data of genes per sample before (top boxplot) and after (bottom boxplot) the log2 transformation and z-normalization:

Boxplot of expression: $\log_2(\text{TPM} + 1)$ values



Boxplot of expression: normalized $\log_2(\text{TPM} + 1)$ values



Again, we can see that after the normalization (bottom boxplot), the expression values of all genes per sample are centered around zero, which is not the case before normalization (top boxplot). This plot, along with other quality control plots, can be found in the QC/GSE_RNAseq subdirectory on GitHub. All the plots can be generated using the GSE_rnaseq.R script rows 348-644.

## 8. Differential gene expression analysis

After finishing the data preprocessing for both microarray and RNA-seq data, we are ready to conduct the differential gene expression analysis. To do so, we use the R package "limma". Limma employs a linear model for each gene and assesses whether its expression is statistically significantly different between two groups, such as pancreatitis and normal control groups. All the steps to perform the differential gene expression analysis using microarray data can be found in the GSE_microarrays.R rows 905-1012, and using RNA-seq data in the GSE_rnaseq.R rows 647-752. The output of these scripts is two excel files containing different statistics per gene. You can

access them [here](#). Below, we see the first rows of the one for the microarray data, but the same columns are found in the file for the RNA-seq data.



Next, we use these two excel files to compare the results obtained by analyzing microarray and RNA-seq data.

## 9. Comparison of results

### 9.1. Expected results

Before conducting any analysis, it is useful to think about what you expect the results to be. In our experiment, we anticipate observing differences in the two lists of differentially expressed genes for several reasons. First, the microarray and RNA-seq data we analyzed do not originate from the same individuals and different number of samples were included in each data. Ideally, we would want to have gene expression microarray and RNA-seq data from the same individuals to reduce the inter-individual variance. Second, microarray uses pre-defined probes, while RNA-seq does not. Consequently, we anticipate finding more differentially expressed genes in RNA-seq data

compared to microarray data. Third, according to existing literature, although we expect most of the genes to be found dysregulated (downregulated or upregulated) in the same direction, we expect to observe higher fold-changes in the differentially expressed genes from RNA-seq data compared to microarray data. Having all these in mind, let's now compare the two lists of differentially expressed genes we have.

## 9.2. Description of the data

First, we load the two excel files containing the results from the differential gene expression analysis in R. As they are excel files, we use the function read.xlsx() from the "openxlsx" R package to load them. With the function head(), we can inspect the files.

```
3  ## load libraries
4  library(dplyr)
5  library(openxlsx)
6
7  ## load data
8  cp_vs_normal_microarrays = read.xlsx("DGEA/GSE_microarrays/DGEA_results.xlsx", sheet = 2)
9  cp_vs_normal_rnaseq = read.xlsx("DGEA/GSE_RNAseq/DGEA_results.xlsx", sheet = 2)
```

In both files, each row represents a gene. Therefore, the number of rows in each file tells us how many genes were tested for differential gene expression in each case.

```
> nr_genes_microarrays = nrow(cp_vs_normal_microarrays)
> nr_genes_rnaseq = nrow(cp_vs_normal_rnaseq)
> nr_genes_microarrays
[1] 13744
> nr_genes_rnaseq
[1] 20550
```

We have 13,744 genes for microarray and 20,550 genes for RNA-seq. This larger number of genes tested in the RNA-seq data compared to microarray data aligns with our expectations. As mentioned earlier, this happens because RNA-seq does not rely on pre-defined probes and can capture a larger number of genes compared to microarrays.

Next, we aim to determine how many genes were tested in both cases. In other words, we want to find the overlap of the genes in the two files. To do this, we use the intersect() function.

```
> ## find the overlap of genes in both files
> nr_genes_overlap = length(intersect(cp_vs_normal_microarrays$EntrezGene.ID,
+                                       cp_vs_normal_rnaseq$EntrezGene.ID))
> nr_genes_overlap
[1] 13029
```
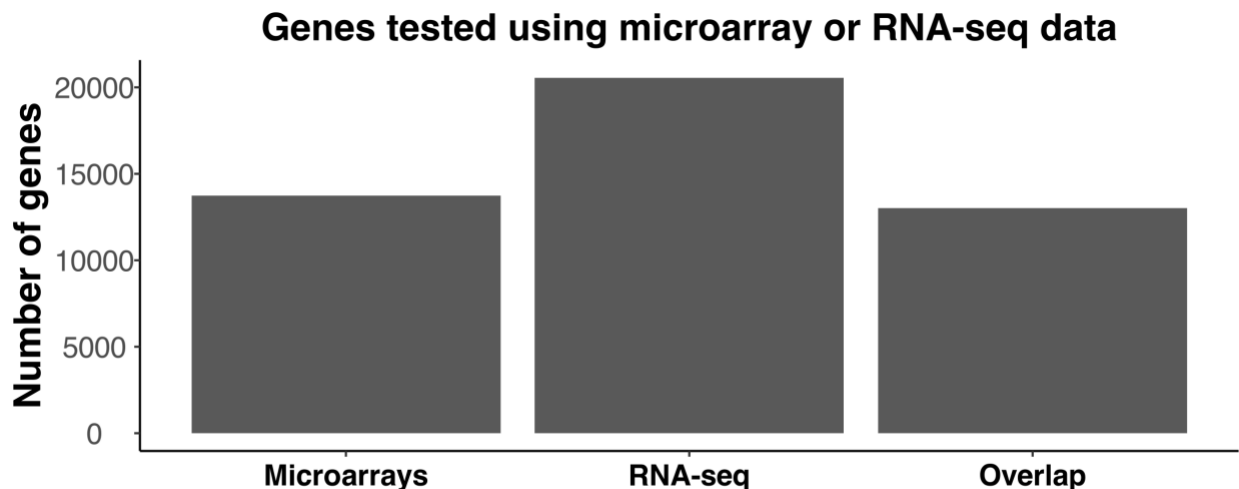
We find that 13,029 genes appear in both files. Let's create a barplot to visualize this information using ggplot2 and the following code:

```
38  ## barplot of genes overlap
39  # create data frame with numbers to be visualized
40  barplot_geneoverlap = data.frame(method = c("Microarrays", "RNA-seq", "Overlap"),
41                                   values = c(nr_genes_microarrays, nr_genes_rnaseq, nr_genes_overlap))
42  # convert "method" to be a factor --> this will help us visualize the data in the row they appear in the barplot_geneoverlap
43  barplot_geneoverlap$method = factor(barplot_geneoverlap$method, levels = barplot_geneoverlap$method, labels = barplot_geneoverlap$method)
44  # create barplot
45  ggplot(barplot_geneoverlap, aes(x = method, y = values)) +
46      geom_col() +
47      labs(title = "Genes tested using microarray or RNA-seq data") +
48      xlab("") + # empty title for x axis
49      ylab("Number of genes") + # name for y axis
50      theme_classic() +
51      theme(plot.title = element_text(face = "bold", size = 18, hjust = 0.5), # bold and center title
52            axis.text.y = element_text(angle = 0, hjust = 0.5, vjust = 0.5, # change font size and position of y axis text
53                                       margin = margin(t = 1, unit = "cm"),
54                                       size = 14),
55            axis.text.x = element_text(face = "bold", color = "black", # change font size and position of x axis text
56                                       angle = 0, hjust = 0.5, vjust = 0.5,
57                                       size = 14),
58            axis.title = element_text(angle = 0, hjust = 0.5, face = "bold", # change font size and location of axes titles
59                                      margin = margin(t = 3, unit = "cm"),
60                                      size = 18))
61  # save plot
62  ggsave(filename = "Number_of_genes_overlap.tiff",
63         path = "DGEA/results_comparison/",
64         width = 2500, height = 1080, device = 'tiff', units = "px",
65         dpi = 300, compression = "lzw", type = type_compression)
```

The output is this:



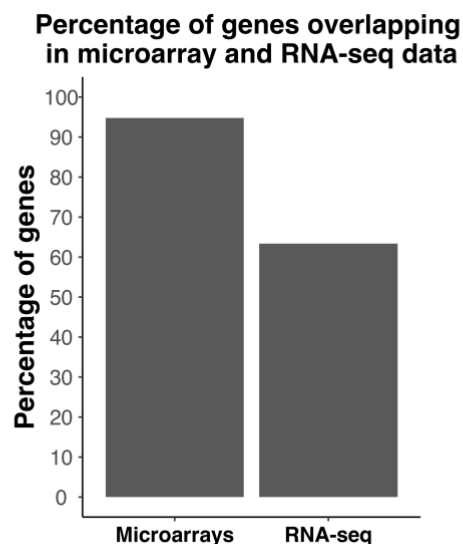**Genes tested using microarray or RNA-seq data**

We can see that almost all the genes tested in microarrays were also tested in RNA-seq. The few genes that were not tested in RNA-seq might be genes with low expression values that happened not to be captured by RNA-seq, just by chance. This is attributed to the bias of both methods, which tend to capture transcripts that are more abundant in samples. Let's now create a barplot that shows the percentages of overlapping genes in each method, instead of the number of genes.

The code is:

```
67  ## create a barplot that shows the percentages of overlapping genes in each method rather than the actual number of genes
68  # data
69  barplot_geneoverlap_percentages = data.frame(method = c("Microarrays", "RNA-seq"),
70                                               values = c((nr_genes_overlap / nr_genes_microarrays) * 100, # calculate percentages
71                                                          (nr_genes_overlap / nr_genes_rnaseq) * 100))
72  # create barplot
73  ggplot(barplot_geneoverlap_percentages, aes(x = method, y = values)) +
74    geom_col() +
75    labs(title = "Percentage of genes overlapping \nin microarray and RNA-seq data") +
76    xlab("") + # empty title for x axis
77    ylab("Percentage of genes") + # name for y axis
78    scale_y_continuous(breaks = seq(0, 100, 10), limits = c(0, 100)) + # set value breaks and limits for y axis
79    theme_classic() +
80    theme(plot.title = element_text(face = "bold", size = 18, hjust = 0.5), # bold and center title
81          axis.text.y = element_text(angle = 0, hjust = 0.5, vjust = 0.5, # change font size and position of y axis text
82                                     margin = margin(t = 1, unit = "cm"),
83                                     size = 14),
84          axis.text.x = element_text(face = "bold", color = "black", # change font size and position of x axis text
85                                     angle = 0, hjust = 0.5, vjust = 0.5,
86                                     size = 14),
87          axis.title = element_text(angle = 0, hjust = 0.5, face = "bold", # change font size and location of axes titles
88                                    margin = margin(t = 3, unit = "cm"),
89                                    size = 18), aspect.ratio = 1.3)
90  # save plot
91  ggsave(filename = "Number_of_genes_overlap_percentages.tiff",
92         path = "DGEA/results_comparison/",
93         width = 2500, height = 1600, device = 'tiff', units = "px",
94         dpi = 300, compression = "lzw", type = type_compression)
```

The output is:



21

Again, we can see that almost all the microarray genes are present in the RNA-seq data. However, approximately 35% of the RNA-seq genes were not found in the microarray data.

## 9.3.    Overlap of significantly dysregulated genes

Let's now test if both methods identify genes to be dysregulated in the same direction. To do so, we first need to filter for the overlapping genes. We can use the intersect() function to find the overlapping genes and then use dplyr syntax to filter both variables for the overlapping genes:

```
96  ## filter for overlapping genes
97  overlapping_genes = intersect(cp_vs_normal_microarrays$EntrezGene.ID, cp_vs_normal_rnaseq$EntrezGene.ID)
98  cp_vs_normal_microarrays_common_genes = cp_vs_normal_microarrays %>% filter(EntrezGene.ID %in% overlapping_genes)
99  cp_vs_normal_rnaseq_common_genes = cp_vs_normal_rnaseq %>% filter(EntrezGene.ID %in% overlapping_genes)
```

We then filter for the statistically significantly dysregulated genes. We do this by filtering both data frames for adjusted p.value < 0.05 (this is a standard threshold of significance used in research):

```
101  ## filter for significantly dysregulated genes
102  cp_vs_normal_microarrays_common_genes_sig = cp_vs_normal_microarrays_common_genes %>% filter(adj.P.Val < 0.05)
103  cp_vs_normal_rnaseq_common_genes = cp_vs_normal_rnaseq_common_genes %>% filter(adj.P.Val < 0.05)
```

Let's see how many genes were found to be dysregulated in each case:

```
> ## number of significantly dysregulated genes
> nr_dysregulated_genes_microarrays = nrow(cp_vs_normal_microarrays_common_genes_sig)
> nr_dysregulated_genes_microarrays
[1] 1761
> nr_dysregulated_genes_rnaseq = nrow(cp_vs_normal_rnaseq_common_genes_sig)
> nr_dysregulated_genes_rnaseq
[1] 1655
```

A total of 1,761 genes were found to be dysregulated using microarray data, while 1,655 genes were found to be dysregulated using RNA-seq data. Out of these, only 258 were significantly dysregulated in both experiments. This is a surprisingly low number, as we had expected most significantly dysregulated genes to be identified by both methods. If we apply a second filter requiring the absolute log fold-change to be greater than one ($|logFC| > 1$), we find 829 differentially expressed genes for microarray data but 0 for the RNA-seq data. Several reasons could be responsible for this discrepancy. Inter-individual variability might be one of them since

we did not have microarray and RNA-seq data from the same individuals, or the fewer replicates for the RNA-seq data. Another potential reason could be related to the normalization method we employed for RNA-seq. It is important to note that different normalization methods can yield different results. Additionally, the imputation step with the microarray data could also play a role. To gain further insights, we could repeat the analysis by changing the aforementioned preprocessing methods. However, this is outside of the scope of this manual.
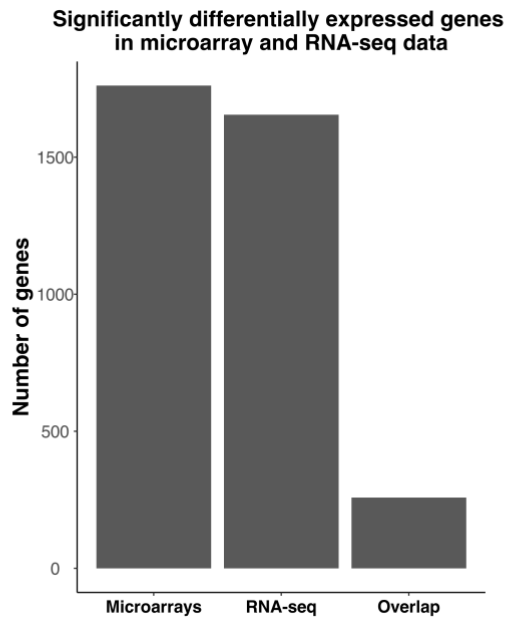
As previously done, let's create a barplot using ggplot2 that summarizes these results. The code is the following:

```r
## barplot to summarize the overlap of dysregulated genes in both data
# data
barplot_dys_geneoverlap = data.frame(method = c("Microarrays", "RNA-seq", "Overlap"),
                                     values = c(nr_dysregulated_genes_microarrays, nr_dysregulated_genes_rnaseq, length(overlap_sig_dys_genes)))
barplot_dys_geneoverlap$method = factor(barplot_dys_geneoverlap$method, levels = barplot_dys_geneoverlap$method,
                                        labels = barplot_dys_geneoverlap$method)
# create barplot
ggplot(barplot_dys_geneoverlap, aes(x = method, y = values)) +
  geom_col() +
  labs(title = "Significantly differentially expressed genes \nin microarray and RNA-seq data") +
  xlab("") + # empty title for x axis
  ylab("Number of genes") + # name for y axis
  theme_classic() +
  theme(plot.title = element_text(face = "bold", size = 18, hjust = 0.5), # bold and center title
        axis.text.y = element_text(angle = 0, hjust = 0.5, vjust = 0.5, # change font size and position of y axis text
                                   margin = margin(t = 1, unit = "cm"),
                                   size = 14),
        axis.text.x = element_text(face = "bold", color = "black", # change font size and position of x axis text
                                   angle = 0, hjust = 0.5, vjust = 0.5,
                                   size = 14),
        axis.title = element_text(angle = 0, hjust = 0.5, face = "bold", # change font size and location of axes titles
                                  margin = margin(t = 3, unit = "cm"),
                                  size = 18), aspect.ratio = 1.3)
# save plot
ggsave(filename = "Number_of_sig_dysregulated_genes_overlap.tiff",
       path = "DGEA/results_comparison/",
       width = 2800, height = 1600, device = 'tiff', units = "px",
       dpi = 300, compression = "lzw", type = type_compression)
```

The output is:

**Significantly differentially expressed genes in microarray and RNA-seq data**

We see that the number of commonly significantly dysregulated genes in the two datasets is very low (most right bar). However, from this visualization we cannot tell if the overlap is significant or not. To test this, we will run a Fisher's exact test (hypergeometric test). In this test, we will ask if the overlap we observe is greater than expected by chance. To do so, we need to create a 2x2 table which will contain the information about how many genes were found to be significantly dysregulated in both studies, only in microarray, only in RNA-seq or in none of them. The code we will use is this:

```
## check the significance of the overlap
# the background genes are the genes tested in both microarry and RNA-seq data (n=13,029)
# create 2x2 table
contingency_table = matrix(ncol = 2, nrow = 2)
colnames(contingency_table) = c("in RNA-seq", "not in RNA-seq")
rownames(contingency_table) = c("in microarray", "not in microarray")
# populate the 2x2 table
# genes found to be significantly dysregulated in both data
contingency_table[1, 1] = length(overlap_sig_dys_genes)
# genes found to be significantly dysregulated in RNA-seq but not in microarray
contingency_table[2, 1] = length(setdiff(cp_vs_normal_rnaseq_common_genes_sig$EntrezGene.ID, cp_vs_normal_microarrays_common_genes_sig$EntrezGene.ID))
# genes found to be significantly dysregulated in microarray but not in RNA-seq
contingency_table[1, 2] = length(setdiff(cp_vs_normal_microarrays_common_genes_sig$EntrezGene.ID, cp_vs_normal_rnaseq_common_genes_sig$EntrezGene.ID))
# genes not found to be significantly dysregulated in any data
contingency_table[2, 2] = length(overlapping_genes) - contingency_table[1, 1] - contingency_table[2, 1] - contingency_table[1, 2]
```

The created 2x2 table should look like this:

```
> contingency_table
                   in RNA-seq not in RNA-seq
in microarray             258           1503
not in microarray        1397           9871
```

Then, we run the Fisher's exact test using the fisher.test() function:

```
# assess the significance of the overlap using the Fisher's exact test (hypergeometric test)
fisher.test(contingency_table, alternative = "greater")
```

```
> fisher.test(contingency_table, alternative = "greater")

        Fisher's Exact Test for Count Data

data:  contingency_table
p-value = 0.00517
alternative hypothesis: true odds ratio is greater than 1
95 percent confidence interval:
 1.071554      Inf
sample estimates:
odds ratio
  1.212882
```

As we can see, the overlap seems to be significant (p = 0.005 < 0.05).


## 9.4.    Direction of dysregulation in overlapping genes

Next, let's investigate whether the overlapping significantly differentially expressed genes exhibit dysregulation in the same direction in both microarray and RNA-seq data. To accomplish this, we first need to filter our two gene lists for these genes:

```
144  ## filter for overlapping significantly dysregualted genes
145  cp_vs_normal_microarrays_common_genes_sig_ovelap = cp_vs_normal_microarrays_common_genes_sig %>% filter(EntrezGene.ID %in% overlap_sig_dys_genes)
146  cp_vs_normal_rnaseq_common_genes_sig_overlap = cp_vs_normal_rnaseq_common_genes_sig %>% filter(EntrezGene.ID %in% overlap_sig_dys_genes)
```

Now, we use the "logFC" column (log-fold change) to create a scatterplot. In this scatterplot, each point represents a gene. The x-axis depicts the log fold-change of each gene as found by the microarray data, while the y-axis represents the log fold-change of each gene as found by the RNA-seq data. To do this, we first should create a data frame with three columns. The first column should be the gene identifiers (EntrezID), the second the log-fold change in microarrays and the third the log-fold change in the RNA-seq data.

```
152  ## create a data frame that contains the logFC from each data type
153  logfc = data.frame(EntrezGene.ID = cp_vs_normal_microarrays_common_genes_sig_ovelap$EntrezGene.ID, # first column contains the gene IDs
154                     microarrays_logFC = cp_vs_normal_microarrays_common_genes_sig_ovelap$logFC, # second column contains the logFC from microarrays
155                     rnaseq_logFC = cp_vs_normal_rnaseq_common_genes_sig_overlap$logFC) # third column contains the logFC from RNA-seq
156
```

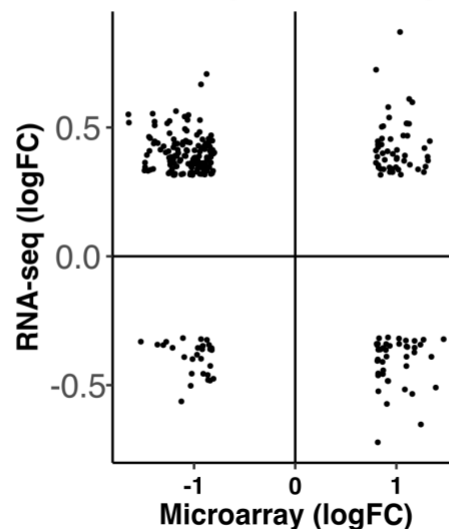Next, we create the scatterplot using ggplot2. The code is:

```
157  ## create scatterplot
158  ggplot(logfc, aes(x = microarrays_logFC, y = rnaseq_logFC)) +
159    geom_point(size = 0.5) +
160    labs(title = "logFC of overlapping significantly dysregulated genes \n in microarry and RNA-seq data") +
161    xlab("Microarray (logFC)") + # empty title for x axis
162    ylab("RNA-seq (logFC)") + # name for y axis
163    geom_vline(xintercept = 0) + # add vertical line
164    geom_hline(yintercept = 0) + # add horizontal line
165    theme_classic() +
166    theme(plot.title = element_text(face = "bold", size = 12, hjust = 0.5), # bold and center title
167          axis.text.y = element_text(angle = 0, hjust = 0.5, vjust = 0.5, # change font size and position of y axis text
168                                     margin = margin(t = 1, unit = "cm"),
169                                     size = 14),
170          axis.text.x = element_text(face = "bold", color = "black", # change font size and position of x axis text
171                                     angle = 0, hjust = 0.5, vjust = 0.5,
172                                     size = 10),
173          axis.title = element_text(angle = 0, hjust = 0.5, face = "bold", # change font size and location of axes titles
174                                    margin = margin(t = 3, unit = "cm"),
175                                    size = 12), aspect.ratio = 1.3)
176  # save plot
177  ggsave(filename = "logFC_of_sig_dysregulated_genes_overlap.tiff",
178         path = "DGEA/results_comparison/",
179         width = 1920, height = 1080, device = 'tiff', units = "px",
180         dpi = 300, compression = "lzw", type = type_compression)
```

The output is:



logFC of overlapping significantly dysregulated genes
in microarry and RNA-seq data

If there was complete concordance between the results from microarray and RNA-seq data, we would expect all the points/genes to fall in the top-right or bottom-left quartiles. However, we observe many points in the top-left and bottom-right quartiles, indicating that there are numerous genes found to have opposite log-fold change values in microarrays compared to RNA-seq data.

With this in mind, let's test if there is a significant correlation between the log-fold change of dysregulated genes in microarray and RNA-seq data. To do this, we can run a Pearson correlation test using the cor.test() function and defining the parameter method to be "pearson":

```
182  ## Pearson correlation analysis for logFC in microarray and RNA-seq data
183  pearson_correlation = cor.test(logfc$microarrays_logFC, logfc$rnaseq_logFC, method = "pearson")
```

The output:

```
> pearson_correlation

        Pearson's product-moment correlation

data:  logfc$microarrays_logFC and logfc$rnaseq_logFC
t = -4.9175, df = 256, p-value = 1.569e-06
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 -0.4014997 -0.1780423
sample estimates:
      cor
-0.2937796
```

We observe a negative correlation (cor = -0.29). This suggests that the majority of significantly differentially expressed genes between chronic pancreatitis and normal samples were found to be dysregulated in opposite directions in microarray and RNA-seq data. Additionally, we find that this correlation is statistically significant (p-value = 1.569e-06 < 0.05). This finding is surprising, as we initially expected the log-fold change to be in the same direction for the same genes in both datasets.

# 10.    Conclusion

While there is a significant overlap in the significantly dysregulated genes identified from the analysis of both microarray and RNA-seq data, there are significant discrepancies in the direction of dysregulation. However, given that the analyzed microarray and RNA-seq data were not obtained from the same individuals, we cannot draw definitive conclusions. Furthermore, some of the decisions made during the preprocessing steps, such as selection of imputation and normalization methods, may impact the results. Nevertheless, it appears that the method of analysis and expression inference can have an impact on the findings and researchers should keep this in mind when they compare differential gene expression results from different methods and studies.

# 11.    References

1. Kaklamani, V.G. and Gradishar, W.J. (2006) 'Gene expression in breast cancer', *Current Treatment Options in Oncology*, 7(2), pp. 123–128. doi:10.1007/s11864-006-0047-0.

2. Wang, C. *et al.* (2014) 'The concordance between RNA-seq and microarray data depends on chemical treatment and transcript abundance', *Nature Biotechnology*, 32(9), pp. 926 932. doi:10.1038/nbt.3001.

3. Zhao, S. *et al.* (2014) 'Comparison of RNA-seq and microarray in transcriptome profiling of activated T cells', *PLoS ONE*, 9(1). doi:10.1371/journal.pone.0078644