

NCBI DATASET ON T-ALL CANCER

Exploring the dataset:

GSE63602 is a dataset from the national center for biotechnology information (NCBI)

*. The dataset is about a disease called T-cell acute lymphoblastic leukemia (T-ALL)

*. It is a type of cancer of the blood & bone marrow. The spongy tissue inside bones where blood cells are made.

*. This develops from white blood cells which are called lymphoblasts

*. The primary cause of the disease is unknown but reports suggest that it is caused by genetic mutation.

How was the dataset created?

The researchers brought a total of 9 T-ALL patients & 4 normal persons. They collected the blood & bone marrow samples from each patient.

From the collected cells, scientists extract total RNA (The RNA includes mRNA, miRNA). These were then sent to a high-throughput sequencing & finally got the dataset.

*. There are a total of 14,000 rows & 4 columns in the dataset. Each row is a kind of RNA that is present).

• The first column is sample ID & the last four columns are the different types of RNA.

What can be found inside the dataset?

- i) ID: These are unique names for the miRNA's.
Each identifies which miRNA's the sample is represented.
- ii) Base mean: This is the average normalized expression of the miRNA across all samples.
- iii) Base mean A: Average normalized expression of a miRNA in group A.
- iv) Base mean B: Average normalized expression of a miRNA in group B.
- v) Fold change: The ratio of expression between two groups (often TACO vs normal).
- vi) log₂ fold change: Logarithmic base 2 of fold change.
$$\log_2 FC = \log_2(\text{fold change})$$
- vii) Pvalue: Probability that the observed difference happened by chance, assuming there is no real difference between groups.
- viii) Padj: Corrected P-value for multiple testing, using FDR.

If Padj is closer to 1 (the cells are more active), closer to 0 (the cells are less active).

Padj tells if the change is caused by the gene activity or is it just a random noise.

Some insightful research paper level questions that can be solved with this dataset.

- i) Can we build a model to distinguish T-ALL samples from healthy control samples based on gene or mRNA expression profiles?
- ii) Are there molecular subtypes of T-ALL that can be identified using unsupervised ML techniques?

How ML models have affected impact on predicting diseases.

There are few challenges that are faced by ML models:

- i) Heterogeneous datasets
- ii) Class imbalance
- iii) Scalability

Project continuation:

Task in ML project competition held at IIT-B.

1. Types of Questions:

"Predicting Leukemia subtypes & discovering Big marker miRNAs using ML on RNA-seq data".

Dataset used is Leukemia.

Final Project Report submitted by team A22.

"Developed an end-to-end ML pipeline to classify T-cell leukemia using RNA-seq expression data, identified candidate microRNA biomarker regulating pathway."

- i) Identifying key RNA Biomarkers linked to leukemia.
- ii) Distinguish cancer vs normal Expression signatures.
- iii) Link dysregulated RNA's to cancer pathways
- iv) construct RNA co-expression networks.
- v) Develop a minimal biomarker signature panel.

STEPS:

1. Data Acquisition & quality control.
2. Exploratory Data Analysis
3. Feature Selection & Statistical Testing
4. Build interpretable models.
5. Biological interpretation.
6. Validation & visualization
7. Documentation & sharing.

1. DATA ACQUISITION & QUALITY CONTROL:

There are totally 14000 rows in the file. In these 14000 rows there is a lot of unwanted genes which has to be removed based on certain conditions. The conditions that I used to remove the genes are:

- i) If $P_{adj} > 0.1$ - we should remove these because these are not statistically significant.

8) i) Filtering genes with $\text{mean} < 10$,
why do we have to do that.

ii) ~~out~~ genes with very low counts
produce unstable dispersion estimates
which can distort the results.

This 2 are the main methods that I used to remove the data from the dataset.

2. Exploratory DATA Analysis:

EDA gives us a sense of how the RNA expression data behave across samples & can reveal biological structure, batch effects, or outliers before you start the modeling.

i) check expression distributions across samples

Best EDA practices for biological data:

A) summary statistics

B) Data transformation & normalization

C) visualizing distributions

D) sample clustering

E) gene correlation

F) outlier detection

G) functional annotation

What is Box Plot?

Box Plots are used in statistics to graphically display various parameters at a glance. In a boxplot the median, the interquartile range and the outliers can be read.

Concept of upregulation & downregulation

If the log₂ fold change > 0 - the gene is more active in T-ALL (up regulation)

If the log₂ fold change < 0 - the gene is less active in T-ALL (down regulation)

Using Interquartile Range:

Formula:

$$\text{Upper bound} = Q_3 + 1.5 * \text{IQR}$$

Max upper outliers = sig_genes [sig_genes

: log₁₀ [log₂ fold change] > upper bound]

This above formula gives the top 10 genes that are responsible for the cancer.

Why volcano plot?

We can use volcano plot to identify the genes which are very highly significant & also probabilistically

It takes two columns for comparison

Two parameters - log₁₀ Padj & log₂ fold change

(9)

If the log₂ Fold change is above 0, that means a gene is upregulated in leukemia samples compared to control.

From the Volcano Plot we can find that the cells which are more active in control & cells which are more active in Patient, & also what is the probability of these cells are showing actual difference rather than being a noise.

Next we can work with PCA:

what is PCA? PCA is a dimensionality reduction technique.

PCA is a powerful dimensionality reduction techniques widely used in bioinformatics.

But PCA cannot be done in this dataset.

because unlike traditional dataset where we have a lot of samples. In this dataset we do not have those samples.

So, I have decided to skip PCA.

In the future works that I do I would love to use PCA, because PCA gives a really good visualization for a large amount of samples.

10

Heatmaps:

A heatmap is a visual way to display how strongly genes are expressed across samples or conditions.

What we can tell from a heatmap:

• which genes are upregulated & which genes are downregulated.

• visual grouping or clustering.

• magnitude of change.

MA Plot:

MA plot is a scatter plot which is commonly used in gene expression analysis to visualize the relationship between the average expression level and the magnitude of expression change between two conditions.

MA plots visualize the relationship

between log₂ foldchange (y values) & mean of normalized counts (x values)

for differential expression analysis

results.

Blue = up, red = down, grey = no change.

Open = not significant, full = significant.

Blue = up, red = down, grey = no change.

Open = not significant, full = significant.

NEXT IS FUNCTIONAL ENRICHMENT & PATHWAY

ANALYSIS: ~~and I left yesterday 29.4.21~~

~~ESOP. tomorrow will do it in~~

Gene ontology (GO) & other enrichment analysis.

~~inherent weight given Biological framework used to~~

~~determine or describe the functions of~~

~~genes and their products - linear, consistent,~~

~~structured, hierarchical lineage across~~

~~different species. (Int'l. R. Consortium)~~

~~based on it is DFBP out for now~~

~~GO Gene ontology provides a standardized~~

~~vocabulary to describe gene function. There~~

~~are 3 main ontology categories:~~

i) Biological Process

ii) molecular function

iii) cellular component

Gene ontology is a pre-built framework.

What in general Functional Enrichment analysis & Pathway analysis mean:

Functional enrichment & Pathway analysis are very useful for interpreting which biological process, molecular functions, or signaling pathways are affected.

12

VALUING & FILTERING INFORMATION IN THE DATA

By performing this I have found out that in GO-cellular component-2023

billions of genes are present in the genome and of those collagen-containing extracellular matrix is the most abundant. This collagen-containing extracellular matrix has really a lot of genes with representation. A statistically significant portion of the genes in the filtered list (the ones that were selected based on p-value & foldchange) are annotated to the "collagen-containing extracellular matrix" pathway.

KEGG Pathway enrichment: 1.0.1.1

KEGG is a database of manually curated pathways, covering metabolism, cell signaling, disease mechanisms & more.

RESULTS THAT I GOT FROM THE FOLLOWING

1. GO RESULTS: i) GO Biological Process-2023

i.) GO Biological Process-2023

from filtering → collagen-containing Extracellular matrix

ii.) GO-Cellular-Component-2023

→ collagen-containing Extracellular matrix

iii.) GO Molecular Function-2023

→ collagen-containing Extracellular matrix

(13)

2. KEGG-2021_Human nt and I test with

i) Extracellular matrix organization

3. Reactome-2022_Protein-DNA interaction

ii) Extracellular matrix organization.

method not yet used previously etc.

changes to and even revert to previous to

What can we do next?

What I did is that I have downloaded a csv file which has the gene list which are common in my csv file & in GO database, KEGG Database

standard PPI for interact. & seen event (i)

i) Statistical strength genes 2nd 3rd 4th

ii) Separately upregulated vs downregulated genes

iii) Functional network / Pathway visualization.

iv) validate biological relevance.

v) ML / Predictive modeling.

TOOLS TO USE instead of R script FT (ii)

swi 2nd result & merge with 2nd

* Cytoscape. rich tool, slow n.s

* string-db for protein-protein interaction

bio2rdf * ClusterProfiler/EnrichmentNetwork plots.

bioperlweb to interact with swi n.s

1A

Now that I have the genes which are specific to ECM.

Based on the GO Interpretation most of the genes belong to ECM category & all the remaining genes are in the category of granule & there are lot of granule types.

Plan:

Phase 1:

- i) There are a total of 19 granule types & some genes are repeating in granule types. So the first process is to find all the unique Granule genes & keep them as one master granule list.
- ii) If there is a particular granule gene which appears 5 times then we can note that down.
- iii) We have to determine if the granule genes are upregulated or downregulated.

Phase 2: (using ENTREZ)

In Phase 2 we will dive into the visualization of the files, here are the files which are going to be used for visualization

- i) AT_biotypes.csv & step AT
- ii) matched_ecm-genes - They are the file which have ecm-genes out of other 4k genes
- iii) Granule unique-genes - This file contains the unique granule genes. Out of 9 granule file.

(entrez) S.51 ← N.A.T

I have performed the visualization & here are the results that I found.

Step 1: AT_biotypes.csv, matched AT

Image 1: AT-ALL vs Normal.Png

• AT-ALL - All samples measured two

This is a heatmap representation of the

genes which are downregulated &

Upregulated.

Legend: Yellow (4.0): very low expression.

" Light Green (8.0): moderate expression

" Dark Blue (14.0): extremely high expression

Ex: COL3A1 (The top): here most of the

AT-ALL "most intense" cell was step

"970" AT-ALL → 14.7 (Dark Blue) is much

Normal → 10.4 (medium blue) in

This is a massive upregulation. The leukaemia cells are producing huge amounts of Type 3 collagen.

16

ELANE (The Enzyme):

With others with T-ALL, gene is silent.
But with T-ALL \rightarrow 17.6 (Green)
With others normal \rightarrow 12.7 (Dark Blue)

This gene is downregulated. In a healthy person, ELANE is active to protect the body.

COL4A2 (The structural switch):

T-ALL \rightarrow 12.2 (Blue)

With others normal \rightarrow 4.4 (Yellow)

In healthy people, the gene is nearly silent (Yellow). In T-ALL it "wakes up" and becomes highly active. This is a perfect oncogene.

Image 2: Clusters of granules, comparing

Normal (Yellow): gene at normal avg level

Positive (Red): The gene is "high"

Negative (Blue): The gene is "low"

The bottom red block: These are the genes we call the "leukemia team". They are always on in T-ALL. It is always "ON" in (normal) A.O.C. Leukemia.

Normal cell. Not different from B cell

Leukemia cell: it turns ON the leukemic genes

The red Blue Block, These are usually "Healthy team". The leukemia has actively silenced these genes. Since these genes (like MMP8) usually destroy collagen, silencing them allows the "Red Block" to grow out of control.