



Inserm



**SORBONNE
UNIVERSITÉ**

Sorbonne University

Master 2 BMC - Systems Immunology

**Characterization of the inflammatory status in osteosarcoma by
in silico RNA-Seq analysis**

By

HUYNH Minh-Anh

September 2022

Table of Contents

List of Abbreviations

Acknowledgements

1	Introduction	1
2	Objective and experimental principle	2
3	Methods	4
3.1	Data collection	4
3.2	Normalization	4
3.3	Construction of gene signatures	4
3.4	Construction of inflammatory groups	5
3.5	MDS Clustering	5
3.6	Determination of cell population abundance	5
3.6.1	MCP Counter	5
3.6.2	CIBERSORTx.	5
3.7	Differential Expression Gene analysis	5
3.7.1	Differentially Expressed Genes analysis (DEG)	5
3.7.2	Over Representation Analysis (ORA)	5
3.8	Statistical Testing	6
3.9	Gene Set Enrichment Analysis (GSEA)	6
4	Results	7
4.1	Determination of inflammatory groups representing intratumor inflammatory status	7
4.2	Characterization of osteosarcomas associated to inflammatory status	8
4.3	Characterization of intra-tumor inflammation associated to inflammatory status . .	8
4.3.1	General relationship of inflammatory status with immune response	8
4.3.2	Similarity of inflammatory signatures	13
4.4	Biological mechanisms underlying the inflammatory groups	14
4.4.1	Differential Gene Expression Analysis	14
4.4.2	Gene Set Enrichment Analysis (GSEA)	14
5	Discussion	18
5.1	Limitations	18
6	Bibliography	20
7	Appendix: code	

List of Abbreviations

- DEG : Differentially Expressed Genes
- MDS : MultiDimensional Scaling
- MSigDB : Molecular Signature Database
- scRNA-seq : single-cell RNA sequencing
- TANs: Tumor-associated Neutrophils
- TME: Tumor MicroEnvironment
- TPM : Transcripts Per Million
- UMAP : Uniform Manifold Approximation and Projection for Dimension Reduction

Acknowledgements

Thank you for following this tutorial!

I hope you'll find it useful to write a very professional dissertation.

1 Introduction

- introduce the reader to the subject area and clarify the knowledge gap that the dissertation research will fill.
- set the context for the dissertation by reviewing the relevant literature.
- include relevant references to general (theoretical papers and reviews) and specific (specific to the particular question addressed) literature, to justify the research that has been undertaken and define the questions being addressed.
- state the primary research questions and hypotheses in the final paragraph.
- follow an ‘inverted triangle’ format, progressing from general scientific ideas and why they matter to the specific research questions addressed in the dissertation project.

The introduction should not be just a ‘Literature Review’.

Osteosarcomas are one of the most aggressive diseases, characterized by high tumor heterogeneity and numerous genomic instabilities. They usually occur in long bones and represent one-fifth of all primary malignant bone tumors and 2.4% of pediatric cancers (Fu et al., 2022). Osteosarcomas occur most frequently in children aged 14-18 years, and more than 30% of cases, within adults over 40 years of age (Mirabello, 2009). The combination of surgery and chemotherapy has increased the 5-year survival rate of patients to over 70%, however metastatic osteosarcoma, frequently occurring in the lungs, drops the survival rate to less than 20%. (Fu et al., 2022).

The difficulty in treating these tumors is largely due to the genomic instability that characterizes them and that leads to a great inter and intra tumor heterogeneity. Thus, the cancer cells can have various phenotypes. Some tumors are osteoblastic, chondroblastic or fibroblastic dominant.

It is established that tumor cells recruit numerous accessory cell types that will form a stroma that supports tumor growth, metastatic progression and resistance to treatment. Cellular (immune cells) and molecular (cytokines, metabolites) inflammatory mediators play a key role in the tumor stroma. The main immune component has been shown to be M2 macrophages. However, scRNA-seq analyses have revealed other immune cell populations such as neutrophils that are also important mediators of inflammation. Bone tumors are a source of inflammatory factors that serve to recruit and educate the stroma.

An issue that remains poorly documented concerns the heterogeneity of the inflammatory landscape of osteosarcomas and the possible relationships between tumor inflammatory status (TIS)

and phenotypic and/or clinical characteristics of these tumors. Similarly, the characteristics of neutrophils associated with bone tumors and their role in tumor inflammation have not been explored to date.

Recently, immunotherapy has proven to be particularly effective in treating previously difficult and deadly cancers, and has been promoted as a medical care therapy in many of them. However, osteosarcomas have not withstood immunotherapy and have instead turned to conventional chemotherapy, which is not very effective and has a plethora of side effects. In order to find a cure for these difficult-to-treat cancers, a better understanding of the tumor microenvironment (TME) could potentially lead to effective and novel targeted TME therapies, especially since osteosarcomas exhibit high heterogeneity and the immunological mechanisms of immune resistance are still unknown. Current advances in bulk RNA sequencing (RNA-Seq) and single-cell RNA sequencing (scRNA-seq) have shown their potential in exploring the tumor microenvironment (TME) to explore intra-tumor heterogeneity and cellular dialogue between tumor and inflammatory cells.

Neutrophils are the primary innate immune cells recruited during an inflammation and multiple papers have already elucidated that tumor-associated neutrophils (TANs) or circulating neutrophils are related to worse patient survival therapy and chemoresistance (Faget et al., 2021; Long et al., 2021). Neutrophils support cancer development through 3 pathways : they're able to promote cancer initiation, assistance of metastasis and increase of tumor growth (Faget et al., 2021).

- Tumeur ne répondent pas à l'immunothérapie
- Tumeur source d'inflammation
- Inflammation et TME non étudié dans ostéosarcomes

2 Objective and experimental principle

The study aims to analyze public dataset and to make a custom analysis using bioinformatic tools in order to understand the cellular crosstalk and interactions in the TME, between neutrophils and tumor cells and their inflammatory status.

Collaborating with Dr Jean-Marc Schwartz's team in Manchester and Sophie Khakoo, a Master student in Systems Biology, collecting data from a scRNA-Seq dataset, a customized downstream analysis will be performed on GSE87686 and Target-OS, and a scRNAseq dataset in order to identify the relationships between neutrophils and osteosarcoma through the use of genetic signatures

available from MSigDB and customized gene lists.

3 Methods

3.1 Data collection

Analyses were performed on two bulk RNA-seq osteosarcoma datasets, TARGET-OS Osteosarcoma and GSE87686. Data was collected from TARGET-OS using GDC Data Transfer Tool UI (v1.0.0), returning 19493 protein coding genes and 88 samples for TARGET-OS, containing both raw data and TPM data. GSE87686 data was obtained through the lab's previously pre-processed kallisto files, downloaded through SRA Run Selector to obtain SRA run files. Data was then imported from kallisto files via *tximport* (v1.22.0) R package. Genes were converted to ENST and ENSG and finally HUGO gene symbols through *biomaRt* (v2.50.3).

3.2 Normalization

Raw data was converted to TPM as it is the best performing normalization method than FPKM or RPKM, based on its preservation of biological signal as compared to other methods (Abrams et al., 2019). Calculation was performed using counts and lengths for each gene, returned from the *tximport* to kallisto process.

Normalization by Z-score with the *scale* function was used to normalize the data for each gene, in order to be able to visualize the data in corresponding heatmaps.

Z-score of the mean of TPM data was used to construct grouped heatmaps for a given gene or gene signature.

Inflammatory groups characterizing the intensity of inflammatory status in tumors were created by choosing the lowest and highest mean of Z-score of the Hallmark Inflammatory Response signature from MSigDB, containing 200 genes. ICAM4 was notably not present in the dataset in TARGET-OS cohort. The groups were cut off evenly using the *ntile* function in *dplyr* (v1.0.8) R package.

3.3 Construction of gene signatures

Gene signatures relevant to the research topic were obtained from MSigDB via *msigdbR* (v7.5.1) HAY_BONE_MARROW_NEUTROPHIL.v7.5.1

Canonical markers for osteosarcoma markers were adapted from Zhou et al. (2020) 's analysis

and their canonical markers generated from the literature in their Supplementary Table 2.

3.4 Construction of inflammatory groups

The groups were cut off evenly using the *ntile* function in *dplyr* R package.

3.5 MDS Clustering

3.6 Determination of cell population abundance

3.6.1 MCP Counter

Estimation of tumor immune infiltration using bulk RNA-seq can be estimated using Microenvironment Cell Populations-counter (MCP-counter) (Becht et al., 2016). Non-metric Kruskal-Wallis testing was used to determine significant differences between cell populations between Low, Medium and High groups with $P < 0.05$ were considered as statistically significant.

3.6.2 CIBERSORTx.

CIBERSORTx algorithm was also used for immune cell deconvolution, from the CIBERSORTx platform (<https://cibersortx.stanford.edu/>) to generate abundance for 22 immune cells.

3.7 Differential Expression Gene analysis

3.7.1 Differentially Expressed Genes analysis (DEG)

Using *DESeq2* (v1.34.0) (Love et al., 2014), standard DEG pipeline using raw data was performed between inflammatory groups (Low, Medium, High). DEGs were identified with *adjusted p.value* (or *FDR*) ≤ 0.05 and $|\log_2\text{FoldChange}| \geq 1$ or $\log_2\text{FoldChange} \leq -1$. An interpretation of the FDR/Benjamini-Hochberg method for controlling the FDR is implemented in *DESeq2* in which we rank the genes by p-value, then multiply each ranked p-value by m/rank (m = total number of tests).

3.7.2 Over Representation Analysis (ORA)

enrichR v(3.0) was used (Xie et al., 2021; **R-enrichR?**) for functional gene enrichment/pathway analysis. Following database were queried : “*GO_Molecular_Function_2021*”, “*Human_Gene_Atlas*”, “*BioPlanet_2019*”, “*GO_Biological_Process_2021*”, “*GO_Cellular_Component_2021*”

3.8 Statistical Testing

Statistical testing was performed using *R software 4.2.0 (2022-04-22)*. Kruskal-Wallis testing was performed for non-metric comparative analysis between groups. Post-hoc analysis was performed using Dunn's test as opposed to Wilcoxon due to the test taking into account Kruskal-Wallis's rank. $P < 0.05$ and $P_{adj} < 0.05$ was considered statistically significant.

3.9 Gene Set Enrichment Analysis (GSEA)

4 Results

Some more guidelines from the School of Geosciences.

This section should summarize the findings of the research referring to all figures, tables and statistical results (some of which may be placed in appendices).

- include the primary results, ordered logically - it is often useful to follow the same order as presented in the methods.
- alternatively, you may find that ordering the results from the most important to the least important works better for your project.
- data should only be presented in the main text once, either in tables or figures; if presented in figures, data can be tabulated in appendices and referred to at the appropriate point in the main text.

Often, it is recommended that you write the results section first, so that you can write the methods that are appropriate to describe the results presented. Then you can write the discussion next, then the introduction which includes the relevant literature for the scientific story that you are telling and finally the conclusions and abstract – this approach is called writing backwards.

4.1 Determination of inflammatory groups representing intratumor inflammatory status

Functional clusters of osteosarcoma samples were created using k-means clustering of from MDS visualization of the Hallmark Inflammatory Response signature.

Inflammatory groups characterizing the intensity of inflammatory status in tumors were created by first creating groups of inflammatory status. The number of groups were tried using k-means clustering algorithm on a MDS visualization of the scaled by Z-score of the 88 tumor samples for the hallmark inflammatory response signature from MSigDB, containing 200 genes. ICAM4 was notably not detected in the dataset in the TARGET-OS cohort. Thus, 199 genes from the signature were used.

The groups are chosen from Figure 1.C as they are relevant as they correspond fairly well to functional groups, defined by k-means clustering based on MDS visualization (**Figure 1**). Each

sample is thus attributed to its inflammatory status and this group will be subsequently used for the following results.

4.2 Characterization of osteosarcomas associated to inflammatory status

In order to characterize osteosarcomas, the *Hp Osteosarcoma* gene signature from MSigDB was used to see whether the inflammatory groups can be related to gene expression from this signature. Visually, the heatmap representing the gene signature, annotated with the inflammatory groups, does seem to indicate that the samples express different genes between Low, Medium and High inflammation group (**Fig. 2A**). However, the dendrogram clustering the samples indicates that the gene signature does not cluster well with the inflammatory annotations. However, despite high heterogeneity between samples and inflammatory status,

Comparison of the expression of specific osteosarcoma markers relating to osteoblastic, chondroblastic, fibroblastic markers has also been done, through a heatmap representation. Hierarchical clustering of the samples does not appear to be associated with corresponding inflammatory status. However it does reveal that there are groups of osteoblastic, chondroblastic and fibroblastic osteosarcomas which is expected (**Fig. 2C**).

The mean of markers of proliferation (MKI67, PCNA, TOP2A) associated to osteosarcomas have been compared to inflammatory status, along with the mean of the three markers. Kruskal-Wallis testing is significant ($p = 0.00968$) and post-hoc Dunn analysis reveals that the mean of the proliferation markers between low and high group is significantly different ($p = 0.016$). The data suggests that proliferation is hindered when inflammatory status is high in the osteosarcoma samples.

4.3 Characterization of intra-tumor inflammation associated to inflammatory status

4.3.1 General relationship of inflammatory status with immune response

4.3.1.1 Relationship with ESTIMATE and inflammatory signatures Using ESTIMATE algorithm from *tidyestimate* R package, an immune score has been calculated for each sample which reflects the immune infiltration in a given tumor sample. The violin plot represents the values obtained for each inflammatory group (**Fig. 3**).

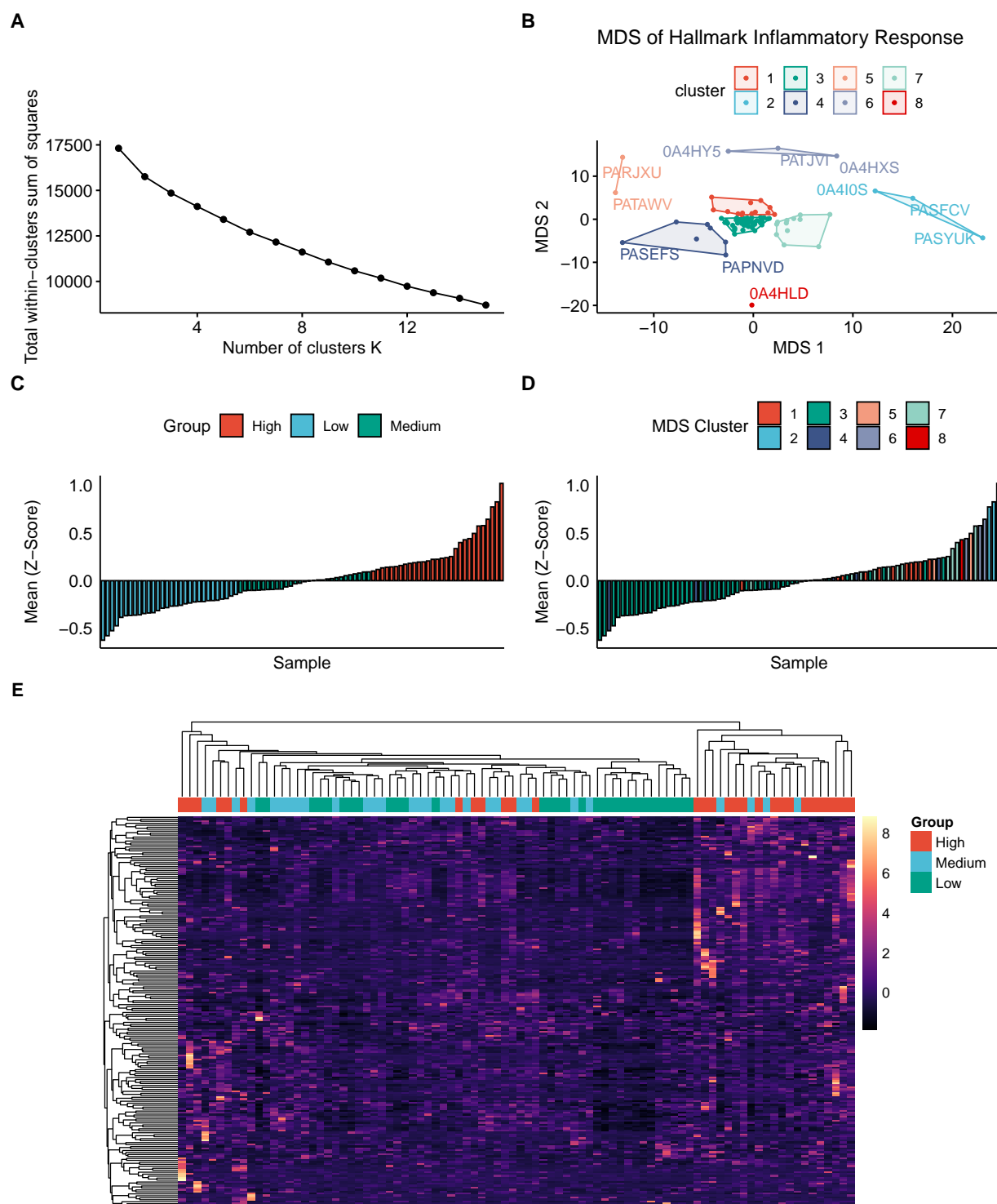


Figure 1: Construction of inflammatory groups. E. Heatmap of the Hallmark Inflammatory Response signature with samples annotated to their corresponding inflammatory groups. Rows correspond to genes and columns correspond to samples.

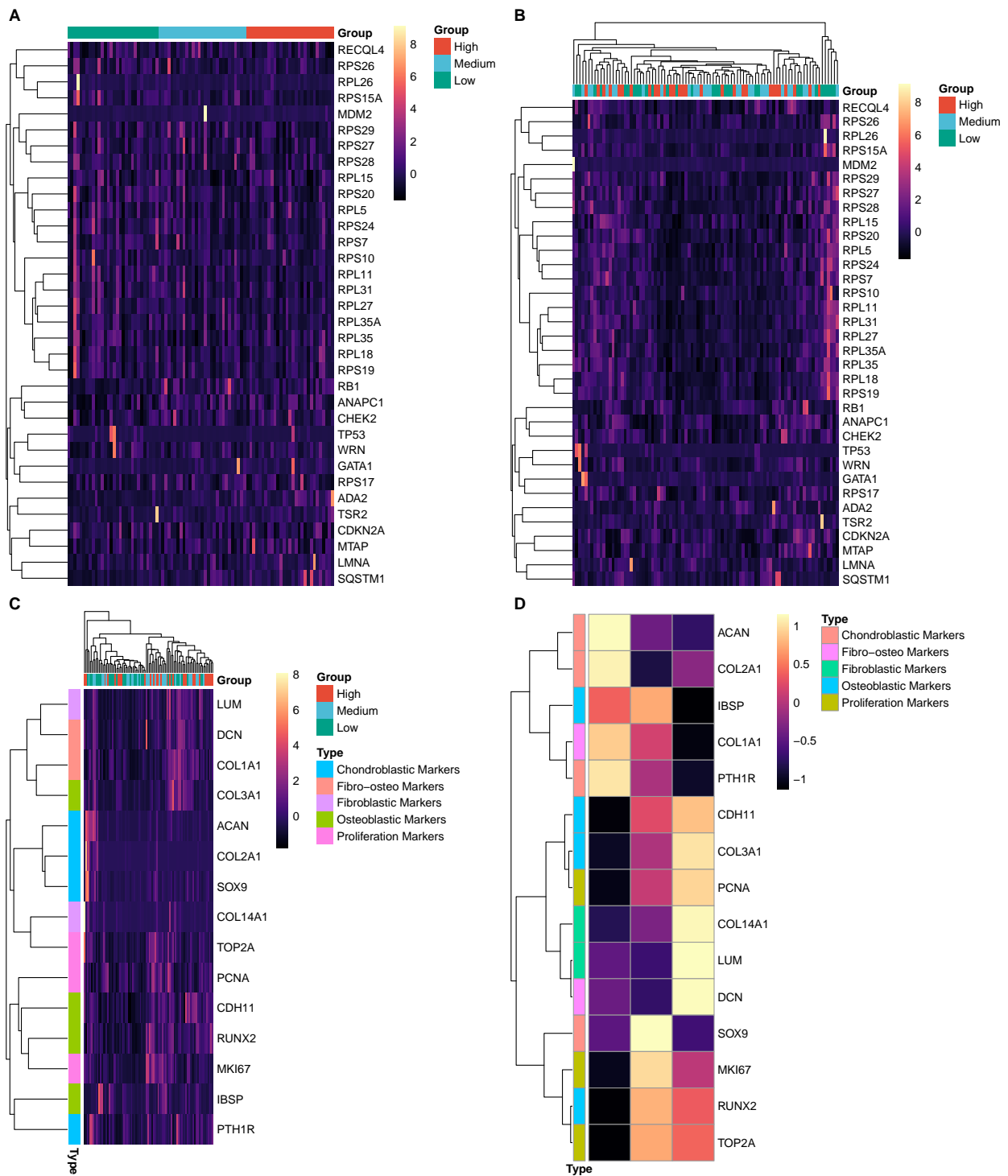


Figure 2: Heatmap of HP Osteosarcoma signature with samples annotated to their respective inflammatory groups.

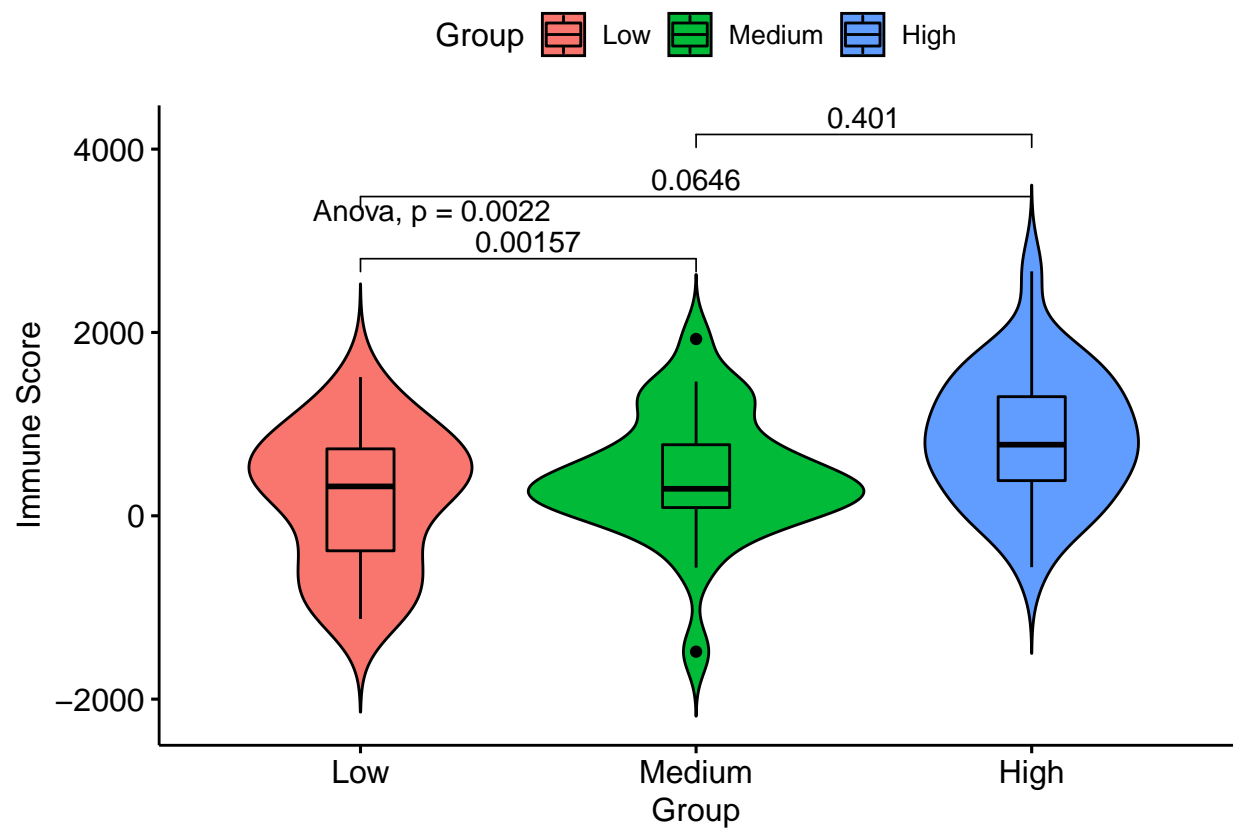


Figure 3: Violin plot of ESTIMATE score for each inflammatory group. ANOVA was performed followed by Tukey's post-hoc analysis. $P < 0.05$ is considered statistically significant.

4.3.1.2 Immune abundance by immune deconvolution algorithm Immune cell abundance can be determined thanks to various immune deconvolution algorithm on TPM normalized bulk-RNASeq. Here, MCP-counter, CIBERSORTx and xCell have been tried. In MCP-counter, the abundance of neutrophil is increased in the High group compared to the Low (P =)

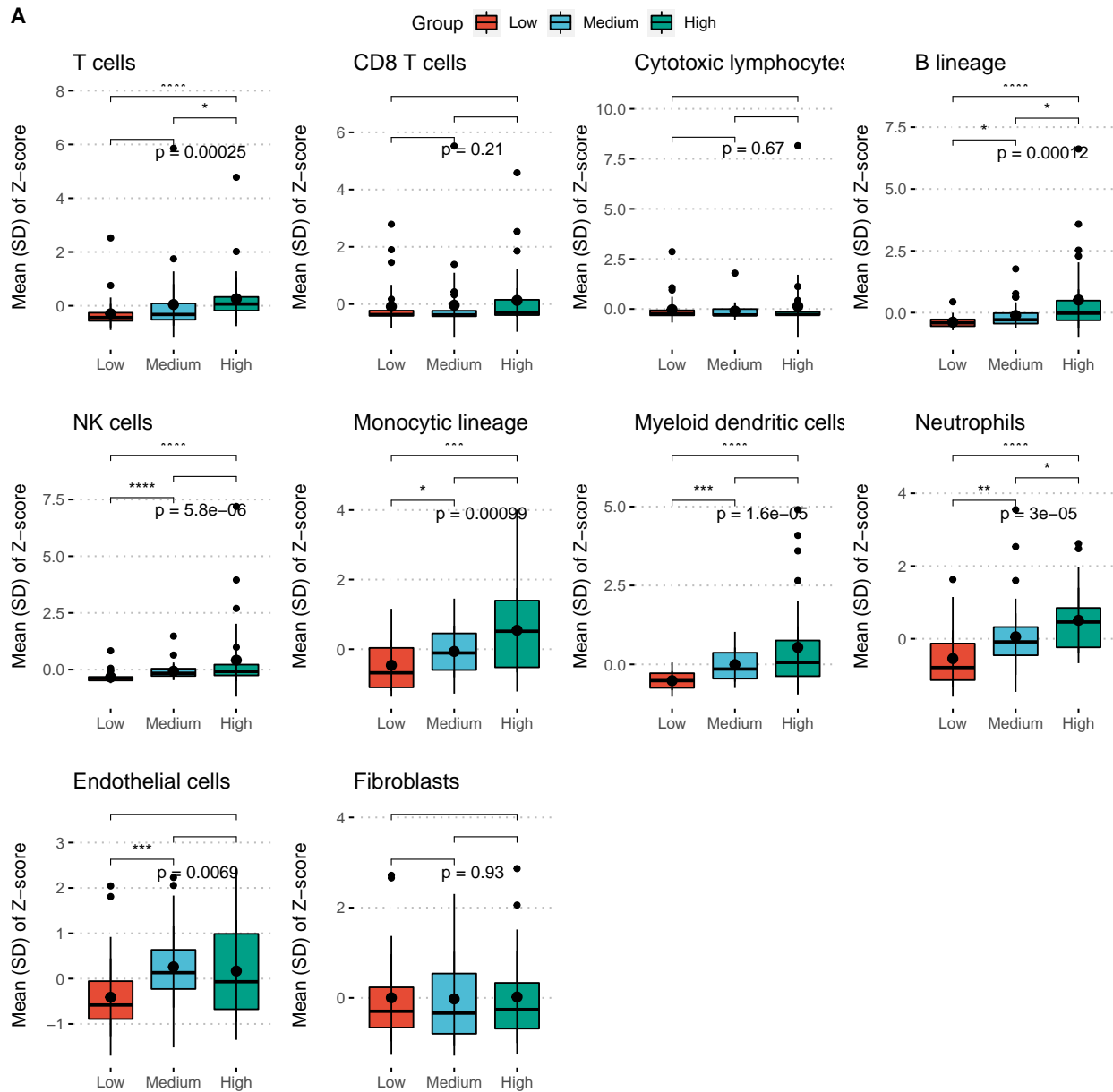


Figure 4: Immune deconvolution plot. (A) MCP-counter

4.3.1.3 Survival curve In order to determine if inflammation groups are linked with increased survival odds, a Kaplan-Meier survival analysis was performed between low and high group of inflammation (**Figure 5**). Log-rank analysis shows that the survival curve between low and high

group of inflammation are not statistically different ($P = 0.16$). However the P value is low enough that it could be considered a trend. Estimation considers 29 patients in the low group and 28 patients in the high group, and estimates that only 9 patients remain alive in the low group while 20 patients are alive in the high group.

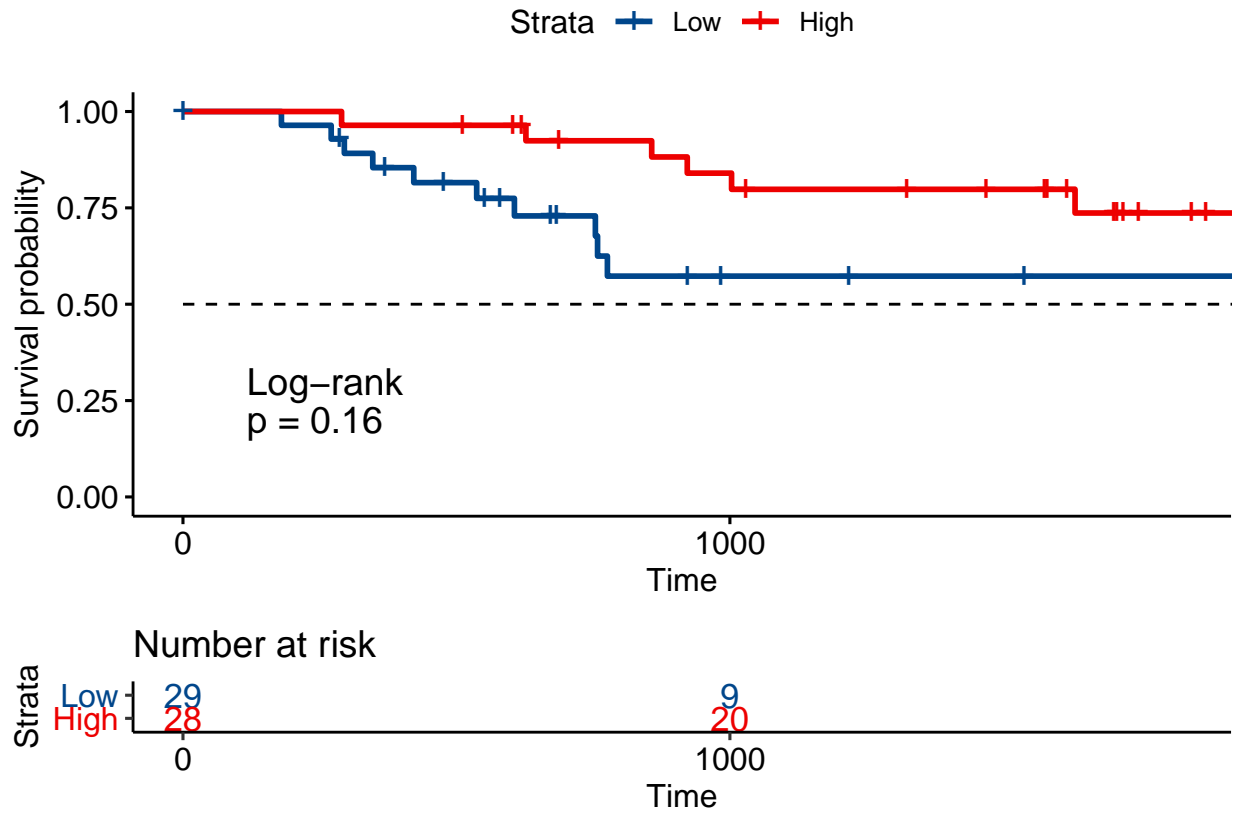


Figure 5: Kaplan-Meier survival analysis plot of low vs high group indicate a trend associating inflammatory status with survival prognostic. Horizontal and vertical axes represent survival times and rates, respectively. Blue and red colors represents low and high inflammation group respectively. Plus signs indicate censored values. Depicted P-values were obtained by the log-rank test. The diagram is cut at 1825 days (5 years). $P < 0.05$ is considered statistically significant.

4.3.2 Similarity of inflammatory signatures

In order to see if the inflammatory signatures used have overlapping genes or not, a venn diagram was constructed from the signatures (**Figure 6**). Overall, the three immune signatures depicting immune infiltration (immune score), inflammation response and TNFa pathway have low overlapping genes and can be considered different.

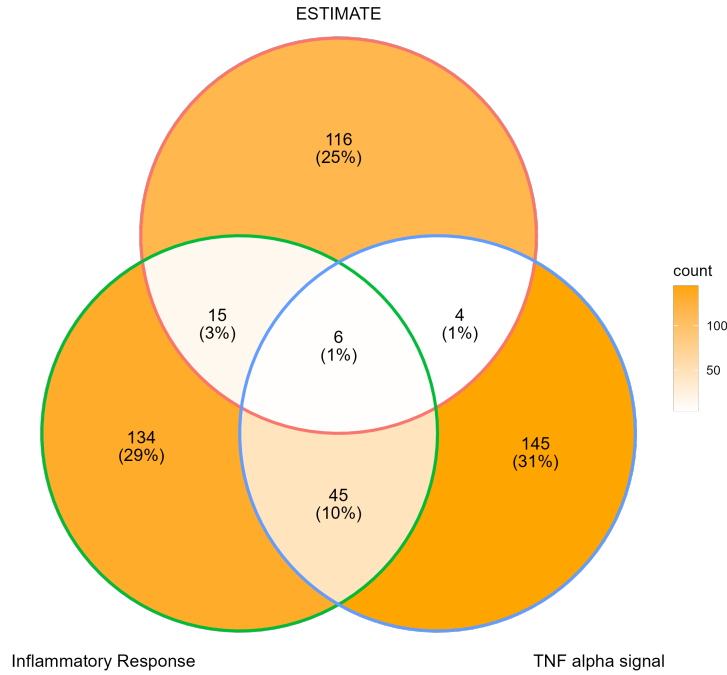


Figure 6: Venn diagram comparing three inflammatory signatures.

4.4 Biological mechanisms underlying the inflammatory groups

4.4.1 Differential Gene Expression Analysis

A differential gene expression analysis was performed in Low versus High group.

4.4.2 Gene Set Enrichment Analysis (GSEA)

GSEA was performed on 95 selected signatures from MSigDB, based on the theme of osteosarcoma, angiogenesis, neutrophils, migration, hypoxia and inflammation. Only significant pathways are reported in the table. Certain pathways are not significant but relevant, such as “Gobp Positive Regulation of Vascular Associated Smooth Muscle Cell Proliferation” ($P = 0.13$), “Gobp Positive Regulation of Inflammatory Response”, and all other pathways concerning dendritic, lymphocyte, monocyte, macrophage chemotaxis and migration.

Interestingly, there is no found significant pathway relating to positive regulation of inflammation, but there are significant pathways relating to negative regulation of inflammation, indicating that in the high group, there is negative regulation of inflammation. Similarly, the only signature with a negative enrichment score is “Gobp Positive Regulation of Neutrophil Activation”. It is also notable that the signature *Hp Osteosarcoma* in GSEA is not significantly different ($P = 0.38$).

Low vs High

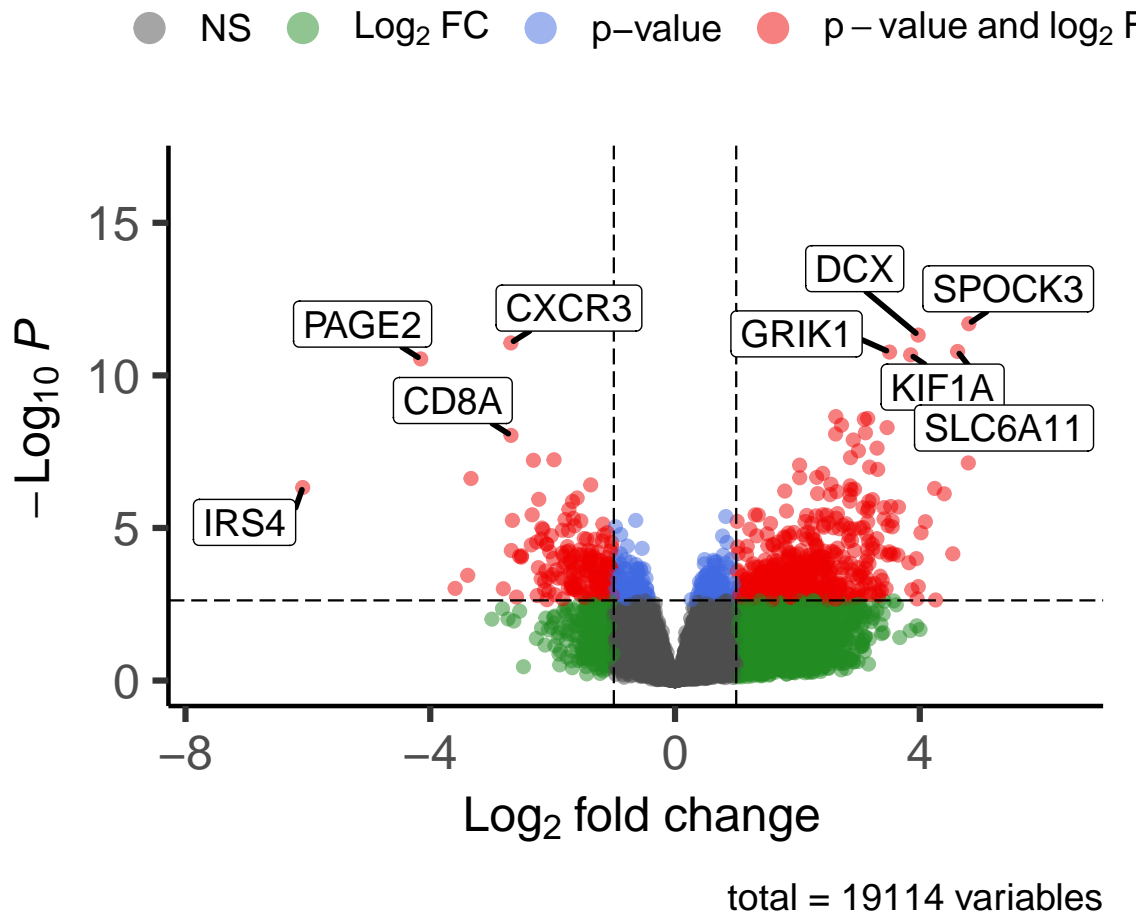


Figure 7: Volcano plot of the differentially expressed genes in Low vs High condition. Negative Log_2 of fold change represent a gene that is more differentially expressed Low condition, and *vice versa*. $P \leq 0.05$ and Log_2 fold change ≤ 1 or ≥ 1 was used for cutoff.

Table 1: Results of GSEA by clusterProfiler, pre-ranked by descending order of log2 of fold change for genes differentially expressed in low vs high condition. Genes at the top of the list are thus differentially expressed in the high group, and vice versa. Only significant pathways are shown, out of 95 selected pathways. P.adjust ≤ 0.10 is considered significant.

Pathway	Set Size	NES	p.adjust	p.sign
Gobp Cell Migration	1431	1.40	0.00	***
Reactome Neutrophil Degranulation	475	1.49	0.00	***
Hallmark Hypoxia	199	1.56	0.00	***
Theilgaard Neutrophil at Skin Wound Dn	226	1.56	0.00	***
Hay Bone Marrow Neutrophil	439	1.38	0.00	***
Hallmark Inflammatory Response	199	1.50	0.00	***
Hp Abnormality of Neutrophils	267	1.46	0.00	***
Gobp Negative Regulation of Inflammatory Response	144	1.51	0.00	**
Hp Abnormal Neutrophil Count	189	1.46	0.00	**
Gobp Endothelial Cell Migration	216	1.43	0.00	**
Kegg Leukocyte Transendothelial Migration	113	1.50	0.00	**
Gobp Negative Regulation of Neuroinflammatory Response	14	1.66	0.01	**
Wp Neovascularisation Processes	37	1.57	0.02	*
Gobp Sprouting Angiogenesis	127	1.42	0.02	*
Hallmark Tnfa Signaling via Nfkb	200	1.34	0.02	*
Theilgaard Neutrophil at Skin Wound Up	75	1.48	0.03	*
Travaglini Lung Neutrophil Cell	332	1.23	0.04	*
Gobp Negative Regulation of Neutrophil Activation	5	1.50	0.04	*
Sig Chemotaxis	45	1.48	0.05	*
Reactome Cellular Response to Hypoxia	75	1.44	0.06	*
Hay Bone Marrow Immature Neutrophil	189	1.28	0.07	*
Gobp Positive Regulation of Neutrophil Activation	6	-1.69	0.09	*
Wp Angiogenesis	24	1.46	0.09	*

```
## Warning in utils::citation(..., lib.loc = lib.loc): pas de champ date dans le  
## fichier DESCRIPTION du package 'bookdown'
```

5 Discussion

the purpose of the discussion is to summarize your major findings and place them in the context of the current state of knowledge in the literature. When you discuss your own work and that of others, back up your statements with evidence and citations.

- The first part of the discussion should contain a summary of your major findings (usually 2 – 4 points) and a brief summary of the implications of your findings. Ideally, it should make reference to whether you found support for your hypotheses or answered your questions that were placed at the end of the introduction.
- The following paragraphs will then usually describe each of these findings in greater detail, making reference to previous studies.
- Often the discussion will include one or a few paragraphs describing the limitations of your study and the potential for future research.
- Subheadings within the discussion can be useful for orienting the reader to the major themes that are addressed.

5.1 Limitations

- Immune deconvolutions methods do not give the same results, thus they are unclear. Only MCP Counter reports significantly different neutrophil abundance between low and high group.
- Immune deconv paper has 7 limitations (Sturm et al., 2019).
- Z-score and equal mean of groups may not be the best normalization method and way of making groups. Geometric mean (inspired from a package) could have been used but returned 0 value despite all values being positive.
- Literature analysis comparing normalization methods reveals library size normalization (TPM, FPKM, RPKM) methods perform worse than distribution normalization methods (DESeq2, TMM), which are recommended by Zhao et al. (2021). Notably, library size normalization methods assume that the total amount of mRNA per cell is identical in every condition.
- DESeq2's normalized count could have then been used, as it is possible to get it from raw-counts after performing DESeq2 analysis, using two functions, `estimateSizeFactors` and

counts(normalized = TRUE). However some immune deconvolution methods still expect TPM as input.

- Using DESeq2's normalized raw count as input, it is inconsistent to then use TPM in order to perform other analysis.
- Perhaps some outliers could have been removed in order to better identify clusters.
- In single-cell, one potential reason that we do not find neutrophils is because chemotherapy induces aplasia which kills off all short lived immune cells.
- Données de littérature peu nombreuses sur ostéosarcomes car cancer rare, cohorte petite et plus petite dans GSE87686
- Nouvelle cohorte de IGR
- Using mean could have resulted in a less precise formation of groups as the distribution of our samples regarding the inflammatory score was not normal. Thus groups based on the median or the quartiles of the median values of inflammatory samples could have been more appropriate.

6 Bibliography

- Abrams, Z.B., Johnson, T.S., Huang, K., Payne, P.R.O., and Coombes, K. (2019). A protocol to evaluate RNA sequencing normalization methods. *BMC Bioinformatics* 20, 679.
- Becht, E., Giraldo, N.A., Lacroix, L., Buttard, B., Elarouci, N., Petitprez, F., Selves, J., Laurent-Puig, P., Sautès-Fridman, C., Fridman, W.H., et al. (2016). Estimating the population abundance of tissue-infiltrating immune and stromal cell populations using gene expression. *Genome Biology* 17, 218.
- Faget, J., Peters, S., Quantin, X., Meylan, E., and Bonnefoy, N. (2021). Neutrophils in the era of immune checkpoint blockade. *Journal for Immunotherapy of Cancer* 9, e002242.
- Fu, Y., Jin, Z., Shen, Y., Zhang, Z., Li, M., Liu, Z., He, G., Wu, J., Wen, J., Bao, Q., et al. (2022). Construction and validation of a novel apoptosis-associated prognostic signature related to osteosarcoma metastasis and immune infiltration. *Translational Oncology* 22, 101452.
- Long, W., Chen, J., Gao, C., Lin, Z., Xie, X., and Dai, H. (2021). Brief review on the roles of neutrophils in cancer development. *Journal of Leukocyte Biology* 109, 407–413.
- Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology* 15, 550.
- Sturm, G., Finotello, F., Petitprez, F., Zhang, J.D., Baumbach, J., Fridman, W.H., List, M., and Aneichyk, T. (2019). Comprehensive evaluation of transcriptome-based cell-type quantification methods for immuno-oncology. *Bioinformatics* 35, i436–i445.
- Xie, Z., Bailey, A., Kuleshov, M.V., Clarke, D.J.B., Evangelista, J.E., Jenkins, S.L., Lachmann, A., Wojciechowicz, M.L., Kropiwnicki, E., Jagodnik, K.M., et al. (2021). Gene Set Knowledge Discovery with Enrichr. *Current Protocols* 1, e90.
- Zhao, Y., Li, M.-C., Konaté, M.M., Chen, L., Das, B., Karlovich, C., Williams, P.M., Evrard, Y.A., Doroshow, J.H., and McShane, L.M. (2021). TPM, FPKM, or Normalized Counts? A Comparative Study of Quantification Measures for the Analysis of RNA-seq Data from the NCI Patient-Derived Models Repository. *Journal of Translational Medicine* 19, 269.
- Zhou, Y., Yang, D., Yang, Q., Lv, X., Huang, W., Zhou, Z., Wang, Y., Zhang, Z., Yuan, T., Ding, X., et al. (2020). Single-cell RNA landscape of intratumoral heterogeneity and immunosuppressive microenvironment in advanced osteosarcoma. *Nature Communications* 11, 6322.

7 Appendix: code

Analyses were conducted using the R Statistical language (version 4.1.2; R Core Team, 2021) on Windows 10 x64 (build 22000). Code is available at <https://github.com/Minh-AnhHuynh/Osteosarcoma-Project>.

Code philosophy follows “The tidyverse style guide” written by Hadley Wickham. and usage of tidyverse functions, in order to ensure consistent code readability, usage and naming.

```
librarian::shelf(tidyverse, magrittr, survival, survminer, snakecase, glue, vroom)

survival_data <-
  vroom("../TargetOS-Osteosarcoma/02_data/GDC/GDC_clinical_data/clinical.tsv") %>%
  setNames(to_snake_case(colnames(.))) %>%
  mutate(
    "sample" = map_chr(
      case_submitter_id,
      ~ str_replace(., pattern = "TARGET-40-", replacement = "")
    ),
    .before = 1
  ) %>%
  select(
    sample,
    vital_status,
    days_to_death,
    age_at_diagnosis,
    days_to_last_follow_up
  )
hallmark_mean_groups <-
  vroom("../TargetOS-Osteosarcoma/03_results/hallmark_mean_groups.tsv")

# From ?Surv : The status indicator, normally 0=alive, 1=dead. Other choices are
# TRUE/FALSE (TRUE = death) or 1/2 (2=death)
survival_data <-
  left_join(hallmark_mean_groups, survival_data) %>%
  dplyr::mutate(
    vital_status = case_when(
      vital_status == "Alive" ~ "0",
      vital_status == "Dead" ~ "1"
    ),
    days_to_death = str_replace(days_to_death, "'--", NA_character_),
    days_to_death = as.numeric(days_to_death),
    group_mean = case_when(
      group_mean == "Low" ~ "1",
      group_mean == "Medium" ~ "2",
```

```

    group_mean == "High" ~ "3"
  ),
  group_mean = as.numeric(group_mean),
  vital_status = as.numeric(vital_status),
  days_to_last_follow_up = as.numeric(days_to_last_follow_up)
) %>%
drop_na(days_to_last_follow_up) %>%
filter(!group_mean == 2) %>%
write_tsv("../TargetOS-Osteosarcoma/03_results/survival_inflammatory_data.tsv")

# make_survplot <- function(quantile)

fit <-
  survfit(Surv(days_to_last_follow_up, event = vital_status) ~ group_mean,
    data = survival_data
  )
summary(fit)
summary(fit)$table

# Change color, linetype by strata, risk.table color by strata
# Testing for trend log rank produces p = 0.20, instead of regular log rank =
# 0.40. Plotting for only two groups produce log rank = 0.16.
# Chi-2 square testing produces p = 0.20.

ggsurv <- ggsurvplot(
  fit,
  pval = TRUE,
  pval.method = TRUE,
  risk.table = TRUE,
  risk.table.col = "group_mean",
  surv.median.line = "hv",
  legend.labs = c("Low", "High"),
  ggtheme = theme_pubr(),
  # fun = "cumhaz",
  palette = "lancet",
  xlim = c(0, 1825)
)
# Save Plot and Table
ggsave(
  "../TargetOS-Osteosarcoma/04_figures/ggplot/survival_plot-inflammation-group_mean.png"
  plot = survminer:::build_ggsurvplot(ggsurv)
)
# Save plot only
ggsave(

```

```

glue(
  "../TargetOS-Osteosarcoma/04_figures/ggplot/survival_plot-inflammation-group_mean-pl
),
plot = gg surv$plot,
dpi = "retina"
)
# Chi-2 square testing
surv_diff <-
  survdiff(Surv(days_to_last_follow_up, event = vital_status) ~ group_mean,
    data = survival_data
  )

```

```

##
## The 'cran_repo' argument in shelf() was not set, so it will use
## cran_repo = 'https://cran.r-project.org' by default.
##
## To avoid this message, set the 'cran_repo' argument to a CRAN
## mirror URL (see https://cran.r-project.org/mirrors.html) or set
## 'quiet = TRUE'.

```

package	version	source
assertthat	0.2.1	CRAN (R 4.2.1)
Biobase	2.56.0	Bioconductor
BiocGenerics	0.42.0	Bioconductor
BiocManager	1.30.18	CRAN (R 4.2.0)
bookdown	0.27.3	Github (rstudio/bookdown@900f92106bda673ac6ad8d76317c9017c1b17528)
car	3.1.0	CRAN (R 4.2.1)
carData	3.0.5	CRAN (R 4.2.1)
data.table	1.14.2	CRAN (R 4.2.1)
DESeq2	1.36.0	Bioconductor
dplyr	1.0.9	CRAN (R 4.2.0)
EnhancedVolcano	1.14.0	Bioconductor
forcats	0.5.1	CRAN (R 4.2.1)
GenomeInfoDb	1.32.2	Bioconductor
GenomicRanges	1.48.0	Bioconductor
ggplot2	3.3.6	CRAN (R 4.2.1)
ggplotify	0.1.0	CRAN (R 4.2.1)
ggpmisc	0.4.7	CRAN (R 4.2.1)
ggpp	0.4.4	CRAN (R 4.2.1)
ggpubr	0.4.0	CRAN (R 4.2.1)
ggrepel	0.9.1	CRAN (R 4.2.1)
ggVennDiagram	1.2.0	CRAN (R 4.2.1)
glue	1.6.2	CRAN (R 4.2.0)
gridExtra	2.3	CRAN (R 4.2.1)

package	version	source
IRanges	2.30.0	Bioconductor
kableExtra	1.3.4	CRAN (R 4.2.1)
knitr	1.39	CRAN (R 4.2.0)
magrittr	2.0.3	CRAN (R 4.2.0)
Matrix	1.4.1	CRAN (R 4.2.1)
MatrixGenerics	1.8.1	Bioconductor
matrixStats	0.62.0	CRAN (R 4.2.1)
moments	0.14.1	CRAN (R 4.2.0)
msigdbr	7.5.1	CRAN (R 4.2.1)
patchwork	1.1.1	CRAN (R 4.2.1)
pheatmap	1.0.12	CRAN (R 4.2.1)
purrr	0.3.4	CRAN (R 4.2.0)
readr	2.1.2	CRAN (R 4.2.1)
rmarkdown	2.14	CRAN (R 4.2.0)
RSQLite	2.2.14	CRAN (R 4.2.1)
rstatix	0.7.0	CRAN (R 4.2.1)
S4Vectors	0.34.0	Bioconductor
sessioninfo	1.2.2	CRAN (R 4.2.1)
snakecase	0.11.0	CRAN (R 4.2.1)
stringr	1.4.0	CRAN (R 4.2.0)
SummarizedExperiment	26.1	Bioconductor
tibble	3.1.7	CRAN (R 4.2.0)
tidyr	1.2.0	CRAN (R 4.2.1)
tidyverse	1.3.1	CRAN (R 4.2.1)
uwot	0.1.11	CRAN (R 4.2.1)
viridis	0.6.2	CRAN (R 4.2.1)
viridisLite	0.4.0	CRAN (R 4.2.1)
vroom	1.5.7	CRAN (R 4.2.1)

Summary

English Summary

no more than 250 words for the abstract

- a description of the research question/knowledge gap – what we know and what we don't know
- how your research has attempted to fill this gap
- a brief description of the methods
- brief results
- key conclusions that put the research into a larger context