



Universidade de Brasília
Instituto de Ciências Exatas - IE
Departamento de Estatística - EST

**Análise de Expressões Gênicas
em Síndrome de Sjögren**

Bárbara Jacqueline Cruvinel Matos

16/0002796

Trabalho de Conclusão de Curso

Orientadora: Prof. Dra. Joanlise Marco de Leon Andrade

Brasília
2020

Bárbara Jacqueline Cruvinel Matos

Análise de Expressões Gênicas em Síndrome de Sjögren

Trabalho de Conclusão de Curso apresentado ao Departamento de Estatística da Universidade de Brasília como parte dos requisitos necessários à obtenção do título de Bacharel em Estatística.

Orientador: Prof. Dra. Joanlise Marco de Leon Andrade

Universidade de Brasília – UnB

Instituto de Matemática

Departamento de Estatística

Brasília

14 de Dezembro de 2020

Bárbara Jacqueline Cruvinel Matos

Análise de Expressões Gênicas em Síndrome de Sjögren/ Bárbara Jacqueline
Cruvinel Matos. – Brasília, 14 de Dezembro de 2020-

81 p. : il. (algumas cores.) ; 30 cm.

Orientador: Prof. Dra. Joanlise Marco de Leon Andrade

Trabalho de Conclusão de Curso – Universidade de Brasília – UnB

Instituto de Matemática

Departamento de Estatística, 14 de Dezembro de 2020.

1. Análise de Expressões Gênicas. 2. Síndrome de Sjögren. I. Orientador. II.
Universidade de Brasília. III. Faculdade de Ciências Exatas. IV. Departamento
de Estatística.

CDU 02:141:005.7

Bárbara Jacqueline Cruvinel Matos

Análise de Expressões Gênicas em Síndrome de Sjögren

Trabalho de Conclusão de Curso apresentado ao Departamento de Estatística da Universidade de Brasília como parte dos requisitos necessários à obtenção do título de Bacharel em Estatística.

Trabalho aprovado. Brasília, 14 de Dezembro de 2020:

Prof. Dra. Joanlise Marco de Leon Andrade

Prof. Dr. George Freitas von Borries

Prof. Dra. Thais Carvalho Valadares Rodrigues

Brasília
14 de Dezembro de 2020

Agradecimentos

Agradeço a Deus, primeiramente. Sem Ele nada em minha vida seria possível. Agradeço também à minha família, em especial a minha mãe que sempre me apoiou em tudo que eu precisava durante a minha vida e sempre esteve disposta a me ajudar em tudo o que foi preciso.

Eternos agradecimentos à Universidade de Brasília que abriu tantas portas e me ensinou tanto, em especial a minha professora e orientadora, Dra. Joanlise Marco de Leon Andrade, pelas diretrizes, disponibilidade, ensinamentos, compreensão, incentivo, paciência, além de seu grande desprendimento em ajudar-me sempre, viabilizar este projeto e acreditar em minha capacidade. Dr. George Freitas Von Borries, agradeço pelo apoio, ajudas, incentivos e inúmeras oportunidades de conhecimento.

Resumo

Neste estudo foram analisados dados de expressão gênica de um grupo de casos de Síndrome de Sjögren Primária e um grupo controle, de indivíduos saudáveis. O principal objetivo foi a realização de análises de perfil de expressão gênica, com base na identificação de genes diferencialmente expressos entre os casos e controles. Três abordagens foram utilizadas para esse fim: testes de comparação de médias, testes não paramétricos e o algoritmo PPCLUST (Particionamento com uma robusta medida de distância para dados com altas dimensões e um tamanho amostral baixo). Análises de agrupamento por técnicas não supervisionadas dos transcritos selecionados foram realizadas para a visualização dos perfis de expressão. Destaca-se a superexpressão de genes ligados ou induzidos por interferons. Também foram identificados sistemas biológicos por Anotação Genômica (*Annotation*), a maior parte dos quais também associados a interferons. Os resultados obtidos são consistentes com os descritos em outros estudos de Síndrome de Sjögren.

Palavras-chaves: Análise de Perfil de Expressão Gênica. Síndrome de Sjögren. Interferons. Análise de Agrupamento. Análise de Cluster. Correção por testes múltiplos. FDR. PPCLUST.

Abstract

In this study, gene expression data from Primary Sjögren's Syndrome cases and a group of healthy controls were analyzed. The main objective was to perform gene expression profiling on the identified differentially expressed genes between cases and controls. Three different approaches were employed for this purpose: parametric tests, non-parametric tests and PPCLUST algorithm (Partition clustering of high dimensional low sample size data based on p-values). Unsupervised Cluster Analyses were performed for profiling visualization. Several interferon genes and interferon inducible genes were identified. Important biological pathways were identified by Annotation, most of which were related to interferon pathways. Results were consistent with other Sjögren's Syndrome genetic studies.

Key-words: Gene expression profiling. Sjögren's Syndrome. Interferon genes. Cluster analysis. Multiple Testing Correction. FDR. PPCLUST.

Lista de ilustrações

Figura 1 – Resumo esquemático dos mecanismos genéticos e epigenéticos associados com a suscetibilidade à SSp	21
Figura 2 – Resumo esquemático do processo de obtenção das intensidades de expressão gênica	27
Figura 3 – Processo da síntese proteica (ou expressão gênica)	28
Figura 4 – Ilustração de agrupamento hierárquico aglomerativo e divisivo.	36
Figura 5 – Exemplo de <i>Heatmap</i> de perfil de expressão gênica.	37
Figura 6 – Distribuição de expressões gênicas de casos e controles em escala original.	41
Figura 7 – Distribuição de expressões gênicas de casos e controles em escala log2.	42
Figura 8 – Distribuição das expressões médias de casos e controles em escala original.	42
Figura 9 – Distribuição das expressões médias de casos e controles em escala log2.	43
Figura 10 – Distribuição de FC por médias e de FC por medianas.	43
Figura 11 – Distribuição dos valores utilizados para o PPCLUST.	44
Figura 12 – Distribuição de P-valores do Testes de normalidade Shapiro-Wilk.	44
Figura 13 – Número de grupos definidos pelo algoritmo PPCLUST por P-valor.	45
Figura 14 – Estrutura dos grupos obtidos pelo algoritmo PPCLUST.	45
Figura 15 – Diagrama de Venn das 3 listas de transcritos.	47
Figura 16 – <i>Heatmap</i> das expressões gênicas em escala log2 para a Lista L1	48
Figura 17 – <i>Heatmap</i> das expressões gênicas em escala log2 para a Lista L2	48
Figura 18 – <i>Heatmap</i> das expressões gênicas em escala log2 para a Lista L3	49
Figura 19 – <i>Heatmap</i> das expressões gênicas em escala log2 para a Lista L4	49
Figura 20 – <i>Heatmap</i> das expressões gênicas em escala log2 para a Lista L5	50
Figura 21 – Distribuição do número de genes da L1 detectados nos sistemas biológicos mais significantes	52
Figura 22 – Distribuição do número de genes da L2 detectados nos sistemas biológicos mais significantes	53
Figura 23 – Distribuição do número de genes da L3 detectados nos sistemas biológicos mais significantes	53
Figura 24 – Distribuição do número de genes da L4 detectados nos sistemas biológicos mais significantes	54
Figura 25 – Distribuição do número de genes da L5 detectados nos sistemas biológicos mais significantes	54
Figura 26 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista L1	56
Figura 27 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista L2	56

Figura 28 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista L3	57
Figura 29 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista L4	57
Figura 30 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista L5	58
Figura 31 – Diagrama das ligações entre os sistemas biológicos identificados na Lista L1	59
Figura 32 – Diagrama das ligações entre os sistemas biológicos identificados para a Lista L2.	59
Figura 33 – Diagrama das ligações entre os sistemas biológicos identificados para a Lista L3.	60
Figura 34 – Diagrama das ligações entre os sistemas biológicos identificados para a Lista L4.	60
Figura 35 – Diagrama das ligações entre os sistemas biológicos identificados para a Lista L5.	61
Figura 36 – Gráfico de redes para os genes e seus sistemas biológicos relacionados identificados na Lista L1	62
Figura 37 – Cnetplot - WILCOX.	62
Figura 38 – Cnetplot - PPCLUST.	63
Figura 39 – Cnetplot - interseção dos transcritos entre PPCLUST, Teste T e Wilcoxon.	63
Figura 40 – Cnetplot - União dos transcritos entre PPCLUST, Teste T e Wilcoxon.	64
Figura 41 – Ilustração da ação de interferons.	66

Lista de tabelas

Tabela 1 – Tabela do método FDR	34
Tabela 2 – Número de transcritos selecionados para cada etapa das abordagens de identificação de TDEs.	46
Tabela 3 – Genes mais sub/super-expressos na L5	51
Tabela 4 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE 1	76
Tabela 5 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE 2	77
Tabela 6 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE 3	78

Tabela 7 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE	
4	79
Tabela 8 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE	
5	80

Lista de abreviaturas e siglas

AECG	Grupo de Consenso Americano-Europeu, do inglês <i>American-European Consensus Group</i> ;
ANA	Anticorpos Antinucleares;
ANTI-RO	Anticorpos contra o antígeno Ro (proteína citoplasmática);
CBC	Hemograma Completo, do inglês <i>Complete blood count</i> ;
Células NK	Células Exterminadoras Naturais , do inglês <i>Natural Killer Cell</i> ;
DHEA-S	Sulfato de Deidroepiandrosterona, do inglês <i>Dehydroepiandrosterona Sulfate</i> ;
DHEA	Deidroepiandrosterona, do inglês <i>Dehydroepiandrosterone</i> ;
DNA	Ácido Desoxirribonucleico, do inglês <i>Deoxyribonucleic acid</i> ;
EAGC ou GWAs	Estudos de Associação de Genoma Completo, do inglês <i>Genome-wide association studies</i> ;
FDR	Descobertas Taxa de Falsas, do inglês <i>False Discovery Rate</i> ;
HLA	Sistema Antígeno Leucocitário Humano, do inglês <i>Human leukocyte antigen</i> ;
IFN	Interferon;
IRF5	Fator Regulador do Interferão 5, do inglês <i>Interferon regulatory factor 5</i> ;
meQTL	Loci Controladores de Traços Quantitativos da Metilação, do inglês <i>Methylation Quantitative Trait Loci</i> ;
miRNA	Micro Ácido Ribonucleico, do inglês <i>Micro Ribonucleic acid</i> ;
NF- κ B	Fator Nuclear κ B, do inglês <i>Factor Nuclear κB</i> ;
NOD2	Domínio de Oligomerização Nucleotídica 2, do inglês <i>Nucleotide-binding Oligomerization Domain Containing 2</i> ;
OMRF	Fundação de Pesquisa Médica de Oklahoma, do inglês <i>Oklahoma Medical Research Foundation</i> ;

OR	Razão de Chances, do inglês <i>Odds Ratio</i> ;
PPCLUST	Cluster de partição de dados de tamanho amostral alto e dimensional baixo com base em p-valores, do inglês <i>Partition clustering of high dimensional low sample size data based on p-values</i> .
RNA	Ácido Ribonucleico, do inglês <i>Ribonucleic acid</i> ;
RNA _m ou mRNA	Ácido Ribonucleico Mensageiro, do inglês <i>Messenger Ribonucleic Acid</i> ;
SNP	Polimorfismo de Nucleotídeo Simples, do inglês <i>Single Nucleotide Polymorphism</i> ;
STAT4	Transdutor de Sinal e Ativador da Transcrição 4, do inglês <i>Signal transducer and activator of transcription 4</i> ;
TDEs	Transcritos diferencialmente expressos;
TLR	Receptores do tipo Toll, do inglês <i>Toll-like receptors</i> ;
TNF	Fatores de Necrose Tumoral, do inglês <i>Tumor Necrosis Factor</i> ;
SS	Síndrome de Sjögren.

Sumário

	Lista de ilustrações	11
	Lista de tabelas	12
1	INTRODUÇÃO	19
2	METODOLOGIA	25
	2.1 Conjunto de dados	25
	2.2 Estimaco de expresses gnicas	27
	2.2.1 Procedimentos de controle de qualidade	29
	2.2.2 Correo de <i>Background</i> , Normalizao e Sumarizao	29
	2.3 Anlises de expresses gnicas	30
	2.3.1 Anlise exploratria	30
	2.3.2 Mtodos de filtraem	30
	2.3.3 Filtraem por magnitude de efeito	35
	2.3.4 Anlises de agrupamento	35
	2.3.5 Anotao genmica	38
3	RESULTADOS	41
	3.1 Anlise exploratria	41
	3.2 Anlise de perfil de expresso gnica	45
	3.3 Identificao de sistemas biolgicos e funes genticas	51
4	DISCUSSO E CONCLUSOES	65
	REFERNCIAS	67
	REFERNCIAS	69
	APNDICES	71
	APNDICE A – DESCRIO DOS SISTEMAS BIOLGICOS IDENTIFICADOS EM ANLISES DE ANOTAO	73
	APNDICE B – LISTA DE TRANSCRITOS ESTATISTICAMENTE SIGNIFICANTES	76

**APÊNDICE C – INFORMAÇÕES RELACIONADAS AO SOFTWARE
R. E AO BANCO DE DADOS. 81**

1 Introdução

A Síndrome de Sjögren (SS) é uma doença autoimune sistêmica caracterizada pela infiltração de linfócitos em glândulas exócrinas, causando processos inflamatórios e eventual diminuição em suas funções normais (Ramos-Casals e outros, 2012). As glândulas lacrimais e salivares são as mais comumente afetadas, o que resulta em sintomas e sinais de *secura* intensa na boca (xerostomia) e nos olhos (ceratoconjuntivite seca) devido à diminuição na produção de saliva e lágrima. Indivíduos portadores de SS podem apresentar outras complicações graves incluindo fadiga profunda, dor crônica e a infiltração em outras glândulas exócrinas (como o pâncreas, glândulas sudoríparas, glândulas mucosas dos tratos respiratório, gastrointestinal e uro-genital) e apresentar ainda neuropatias e risco aumentado de linfomas (Tapinos e outros, 1999).

A *secura* pode afetar profundamente a qualidade de vida dos portadores de SS, interferindo em funções diárias básicas, como comer, falar e dormir, dado que a redução do volume de saliva e a subsequente perda das propriedades antibacterianas da saliva podem acelerar infecções, como a cárie dentária e a doença periodontal. Glândulas mucosas do trato respiratório superior e inferior também podem ser afetadas, levando a um ressecamento do nariz, garganta e traqueia, resultando muitas vezes em tosse seca crônica (Kassan e Moutsopoulos, 2004).

A SS pode ser classificada em dois tipos: Síndrome de Sjögren Primária (SSp), que ocorre na ausência de outro distúrbio autoimune subjacente, e em Síndrome de Sjögren Secundária (SSs), que está associada a outros distúrbios autoimunes, principalmente lúpus eritematoso sistêmico (LES), artrite reumatoide e esclerose sistêmica (Pasoto, Martins e Bonfa, 2019).

Estima-se que a prevalência mundial de Síndrome de Sjögren esteja em torno de 0,5% a 1,0% da população adulta, ocorrendo em mulheres em 90% dos casos (Carsons e Patel, 2019), em sua maioria caucasianas (Patel e Shahane, 2014) em idades entre os 40 e 60 anos (Qin e outros, 2015). Estudos sugerem que o risco aumentado nessa faixa etária esteja relacionado à ocorrência da menopausa, que causa uma alteração na razão andrógeno-estrogênio e a diminuição de ambos os hormônios (Brandt e outros, 2015). Há evidências de que baixos níveis de andrógenos estão associados à SS. Estudos reportaram concentrações do andrógeno deidroepiandrosterona (DHEA) 40 a 50% mais baixas em pacientes com SS do que em controles pareados por idade e sexo (Valtysdóttir, Wide e Hällgren, 2001). Além disso, foi descoberto que, em glândulas salivares de pacientes com SS, o DHEA foi efetivamente convertido em testosterona e não em di-hidrotestosterona (o que seria o processo normal), provavelmente devido à expressão diminuída das principais

enzimas esteroideogênicas e sua localização subcelular anormal (Nikolov e Illei, 2009).

Fatores ambientais, como infecções virais prévias (a vírus Epstein-Barr, citomegalovírus, vírus herpes humano, vírus da hepatite C, dentre outros) ou bacterianas (a *Helicobacter pylori*) já foram apontados como potenciais desencadeadores da resposta imune ao tecido glandular, devido à sua frequente concomitância em pacientes com SS (Felberg e Dantas, 2006). Além disso, segundo pesquisa realizada (Theander e outros, 2015), autoanticorpos, que são anticorpos produzidos pelo sistema imune que atuam contra uma ou mais proteínas do próprio indivíduo, podem estar presentes 18 a 20 anos antes do diagnóstico de SS primária, principalmente em famílias com vários casos de doenças autoimunes, em que o perfil de autoanticorpos, juntamente com a avaliação do risco genético, permite a identificação de indivíduos suscetíveis em um estado pré-doença (Theander e outros, 2015).

Vários estudos de genes candidatos e, mais recentemente, estudos de associação de genoma completo (EAGC, do inglês *Genome-wide association studies* ou GWAs) reportaram associações de SSp com muitos fatores genéticos, destacando-se os maiores efeitos com genes *HLA* e efeitos moderados com genes *IRF5* e *STAT4*, consistentes em várias etnias (Imgenberg-Kreuz e outros, 2019), como ilustrado na Figura 1. Os genes incluídos nesta Figura, foram identificados com significância de estudos de associação de genoma completo (GWAS), após correção para testes múltiplos que envolvem centenas de milhares a milhões de testes realizados para todos os marcadores separadamente.

Além de associações com os genes *HLA*, *IRF5* e *STAT4*, estudos de genes candidatos para a SSp (Imgenberg-Kreuz e outros, 2019) identificaram outros *loci* associados a maiores riscos para o desenvolvimento da síndrome (Quadro 1 que contém alguns deles assim como a função de cada um).

Também foram observadas em estudos de perfil de expressão gênica associações de SSp com muitos desses genes candidatos, como *IFN*, *EPSTI1*, *STAT1* e *IFI44L*, que são relacionados sobretudo à atividade da doença e à presença de autoanticorpos (Yao e outros, 2019). Tais estudos utilizaram plataformas chamadas de chips de *microarray* (ou microarranjos em português), que permitem a estimação (por quantificação de RNA mensageiro) de expressão gênica de dezenas de milhares de transcritos (trecho do DNA já transcrito em uma molécula de RNA) simultaneamente. Representam uma das abordagens mais amplamente utilizadas para se determinar diferenças globais de transcriptoma (medição em larga escala de transcritos) entre casos e controles saudáveis. Estudos de perfil de expressão gênica permitiram a identificação das denominadas “assinaturas” de expressão da doença. Uma assinatura representa um grupo de genes co-expressos (e diferentemente expressos entre casos e controles) que pertencem a sistemas biológicos específicos. Os genes ou os sistemas identificados podem servir como biomarcadores para diagnóstico, classificação e previsão de resposta a medicamentos (Li e outros, 2013).

Nesse contexto, o presente trabalho teve como objetivo a realização de uma análise de perfil de expressões gênicas com base na identificação de transcritos diferencialmente expressos (TDEs) entre casos com SSp e controles saudáveis.

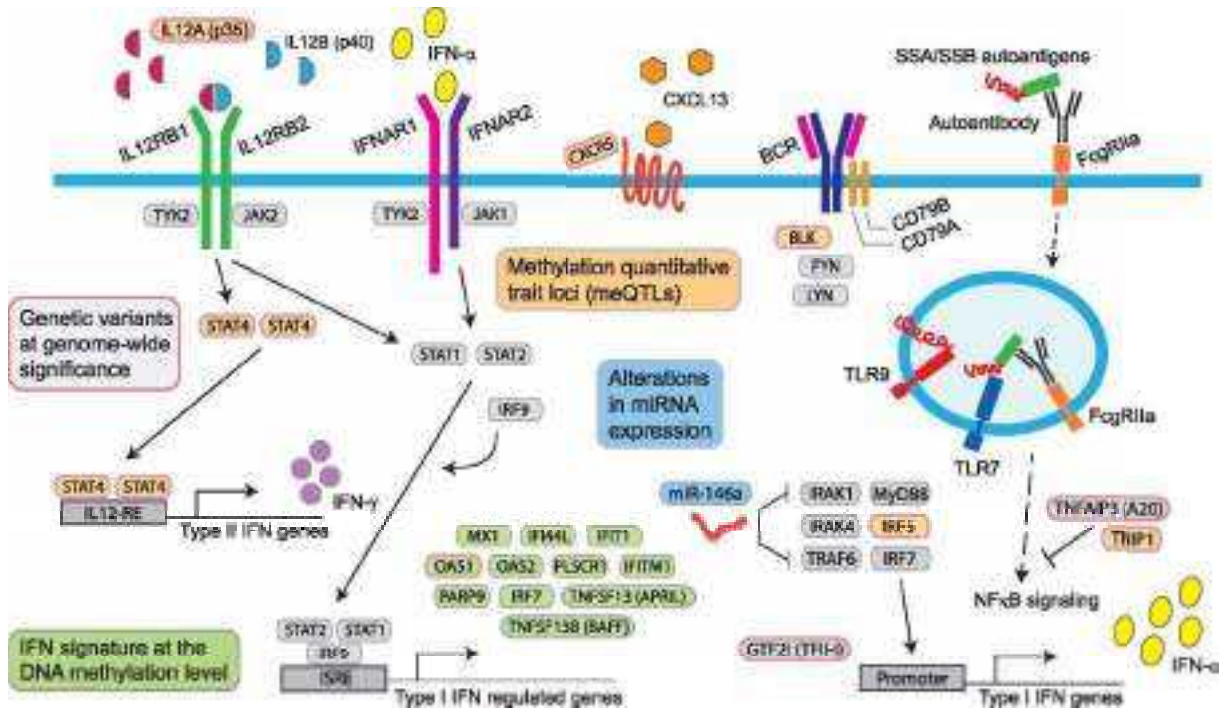


Figura 1 – Resumo esquemático dos mecanismos genéticos e epigenéticos associados com a suscetibilidade à SSp. Caixas com bordas em rosa representam genes localizados fora da região HLA com SNPs associados à suscetibilidade de SSp com significância de estudos de associação de genoma completo (GWAS) após correção para testes múltiplos e incluem *BLK*, *CXCR5*, *GTF2I*, *IL12A*, *IRF5*, *OAS1*, *STAT4*, *TNFAIP3* e *TNIP1*. Caixas verdes representam genes induzidos por interferon (IFN), como *MX1*, *IFI44L*, *OAS1*, *OAS2*, *TNFSF13B* e *IRF7*, com locais CpG hipometilados em SSp. Caixas laranjas incluem meQTLs (loci controladores de traços quantitativos da metilação), referentes à associação entre um variante genético (SNP) e o nível de metilação em um local CpG próximo como *BLK*, *CXCR5*, *IL12A*, *IRF5-TNPO3*, *STAT4* e *TNIP1*. Caixas azuis representam alterações na expressão do miRNA, incluindo o *miRNA-146*, foram identificadas em SSp (Imgenberg-Kreuz e outros, 2019).

Quadro 1: Listagem de genes que conferem maior risco para a SS identificados em estudos de genes candidatos (p-valores corrigidos por testes múltiplos inferiores a 5%).

Gene	SNP	OR (95% CI)	Função
IRF5-TNPO3	rs2004640	1.93 (1.15,3.42)	Tipo I de IFN e sinalização TLR , produção de citocina produção (por exemplo, IL-6, IL-12),apoptose.
	rs10488631	1.57 (1.23,1.99)	
	CGGGGindel	2.00 (1.50,2.70)	
		1.49 (1.24, 1.79)	
STAT4	rs7574865	1.46 (1.09, 1.97)	Tipo I e II de IFN, sinalização NOD2 .
		1.40 (1.21, 1.62)	
	rs7582694	1.41 (1.14,1.73)	
		1.40 (1.15,1.70)	
BLK- FAM167A	rs12549796	1.37 (1.15,1.63)	Sinalização e desenvolvimento de células B. timopoiese , secreção e síntese de insulina e sinalização de NF-kB .
EBF1	rs3843489	1.68 (1.29,2.18)	Potenciador da atividade transcricional durante o desenvolvimento de células B.
TNSF4/OX40L	rs1234315	1.34 (1.14,1.64)	Co-estimulação, proliferação de células T e produção de citocinas.
BAFF/INFSF 13B	5 haplotype	2.60 (1.70,4.10)	Maturação, proliferação e sobrevivência de células B. além da sobrevivência de células epiteliais.
	5 haplotype (TATT,GTTC)	-	
NCR3/NKp30	rs11575837	0.48	Regulação das comunicações dendríticas e celulares NK e regulação das citocinas Th1 (IL-12 e IFN- γ).
	rs2736191	0.56	
PTPN22	rs2476601	2.42 (1.24,4.75)	Ativação de células T.
TNIP1	rs3792783	1.33 (1.16,1.52)	Sinalização de NF-kB e EGF/ERK e apoptose induzida por TNF .
	rs7708392	1.21 (1.08,1.36)	
TNFAIP3	rs2230926	3.26(1.31,8.12)	Sinalização de NF-kB (repressor), apoptose induzida por TNF , sinalização de TLR4 , produção de citocinas (por exemplo, IL-1B).
BAFF-R	His159Tyr	4.13 (1.19,14.3)	Sinalização de NF-kB2.
LTA/LTB/INF	rs1800629	2.00 (1.61,2.49)	Organogênese linfoide, manutenção do tecido linfoide terciário.
	rs909253	1.59 (1.34,1.89)	
MECP2	rs17435	1.33 (1.12,1.59)	Silenciamento transcricional mediado por metilação.

Adaptado e traduzido de: (Imgenberg-Kreuz, et al.,2019)

OR: *odds ratio*.

TLR: toll-like receptors, receptores do tipo Toll, que são uma família de proteínas transmembrânicas de tipo I que formam uma parte do sistema imunológico inato.

NOD2: (em inglês: nucleotide-binding oligomerization domain containing 2) é uma proteína que desempenha um importante papel no sistema imune.

Timopoiese: processo em que dentro do timo os timócitos entram em um processo de seleção e maturação.

NF- κ B: (factor nuclear κ B) é um complexo proteico que desempenha funções como fator de transcrição e um papel fundamental na regulação da resposta imunitária à infecção.

Células dendríticas: nomeadas por suas formas de sondagem, “tipo árvore” ou dendríticas, são responsáveis pelo início das respostas imunes adaptativas e, portanto, funcionam como as “sentinelas” do sistema imunológico.

Células NK: células exterminadoras naturais (do inglês Natural Killer Cell) são um tipo de linfócitos citotóxicos necessários para o funcionamento do sistema imunitário inato.

TNF: (do inglês: Tumor necrosis factor), são os fatores de necrose tumoral e este gene codifica uma citocina pró-inflamatória importante, produzida principalmente por macrófagos, que tem sido relacionada com diversas doenças, como as autoimunes e as degenerativas. Estudos indicam que polimorfismos do TNF podem influenciar seu nível de expressão e atividade e, portanto, poderiam influenciar a susceptibilidade a tumores (Lopes e outros, 2014).

2 Metodologia

2.1 Conjunto de dados

Neste trabalho foram analisados dados de expressões gênicas de 222 indivíduos de descendência europeia, sendo 190 diagnosticados com SSp (grupo de casos) e 32 indivíduos saudáveis (grupo de controle) de um estudo de caso-controle (Lessard e outros, 2013) com 15.063 transcritos e aprovado pelos devidos comitês de ética¹ e consentimento informado (documento assinado) pelos participantes.

Os dados são públicos e fazem parte do repositório funcional de dados genômicos *Gene Expression Omnibus (GEO)* (Edgar, Domrachev e Lash, 2002).

Os pacientes com SS foram avaliados e recrutados por médicos especialistas na Universidade de Minnesota ou na *Oklahoma Medical Research Foundation (OMRF)*. O critério adotado para a classificação dos casos foi o da American-European Consensus Group (AECG) de 2002 (Vitali e outros, 2002), apresentado no Quadro 2.

Segundo o critério AECG, para ser diagnosticado com Síndrome de Sjögren Primária um indivíduo deve apresentar:

- Quaisquer 4 dos 6 critérios, incluindo obrigatoriamente o item IV (Histopatologia) ou o VI (Autoanticorpos) ou
- Quaisquer 3 dos 4 itens de critérios objetivos (III, IV, V, VI).

Tendo como critérios de exclusão o tratamento anterior de radiação de cabeça e pescoço, infecção por hepatite C, AIDS, linfoma preexistente, sarcoidose, enxerto por doença do hospedeiro ou o uso atual de drogas anticolinérgicas.

¹ Oklahoma Medical Research Foundation Review Boards e University of Minnesota Institutional Review Board.

Quadro 2: Critério adotado para a classificação dos casos da (AECG) de 2002.

I. Sintomas oculares - resposta positiva a pelo menos uma das seguintes perguntas:	<ol style="list-style-type: none"> 1. Tem problemas oculares diários e persistentes, relacionados a quadro de olho seco há mais de 3 meses? 2. Tem sensação de areia ou queimação ocular? 3. Usa colírios lubrificantes mais de 3 vezes ao dia?
II. Sintomas orais - resposta positiva a pelo menos uma das seguintes perguntas:	<ol style="list-style-type: none"> 1. Tem uma sensação diária de boca seca há mais de 3 meses? 2. Tem inchaço recorrente ou persistente das glândulas salivares, na idade adulta? 3. Sente necessidade frequente de ingerir líquidos para ajudar na deglutição de alimentos sólidos?
III. Sinais oculares - evidência objetiva de envolvimento ocular definido como resultado positivo para, pelo menos, um dos dois testes a seguir:	<ol style="list-style-type: none"> 1. Teste de Schirmer I, realizado sem anestesia (<5 mm em 5 minutos). 2. Pontuação da Rosa Bengala ou outra pontuação de corante ocular (>4 dna escala de van Bijsterveld's).
IV. Histopatologia: achados histopatológicos nas glândulas salivares menores (obtidas através da mucosa de aparência normal) sia-loadenite linfocítica focal:	Avaliada por um histopatologista especialista, com um escore de foco >1, definido como um número de focos linfocitários adjacentes aos ácinos mucosos de aparência normal com pelo menos 50 linfócitos por 4 mm ² de tecido da glândula salivar.
V. Comprometimento da glândula salivar: evidência objetiva, definida por um resultado positivo para, pelo menos, um dos seguintes testes de diagnóstico:	<ol style="list-style-type: none"> 1. Fluxo salivar total não estimulado (<1,5 ml em 15 minutos). 2. Sialografia parotídea mostrando a presença de sialectasias difusas (padrão pontuado, cavitário ou destrutivo), sem evidência de obstrução nos ductos principais. 3. Cintilografia salivar mostrando captação tardia, concentração reduzida e/ou excreção retardada do traçador.
VI. Autoanticorpos: presença de pelo menos um dos seguintes autoanticorpos séricos:	Anticorpos contra os antígenos Ro (SSA) ou La (SSB), ou ambos.

Os dados de expressões gênicas utilizados são resultantes de dois experimentos de *microarray* separados. O RNA total foi obtido pela coleta de sangue periférico em tubos PAXGene (BD Company) e extraídos seguindo protocolos do fabricante (Qiagen). Concentrações de RNA foram determinadas utilizando-se um espectômetro NanoDrop (Thermo Scientific). A qualidade de RNA foi avaliada pelo Agilent 2100 Bioanalyzer baseada na razão ribossômica 28S/18S e no número de integridade de RNA. cDNA de dupla hélice foi sintetizado usando um promotor T7 e cRNA biotin-labeled foi transcrito utilizando o sistema de amplificação TotalPrep de RNA da Illumina (Ambion). Amostras foram hibridizadas ao chip Human WG-6 v3.0 BeadChip microarrays (da empresa Illu-

mina), que contém 48.803 seqüências de mRNA distintas. Os *microarrays* foram lavados para remoção dos “alvos” excedentes com alto rigor e marcados com estreptavidina-Cy3. Os dados de expressão gênica com base na intensidade fluorescente foram obtidos pela leitura dos chips no scanner BeadStation 500 da Illumina ou o iScan (Lessard e outros, 2013), conforme ilustrado na Figura 2.

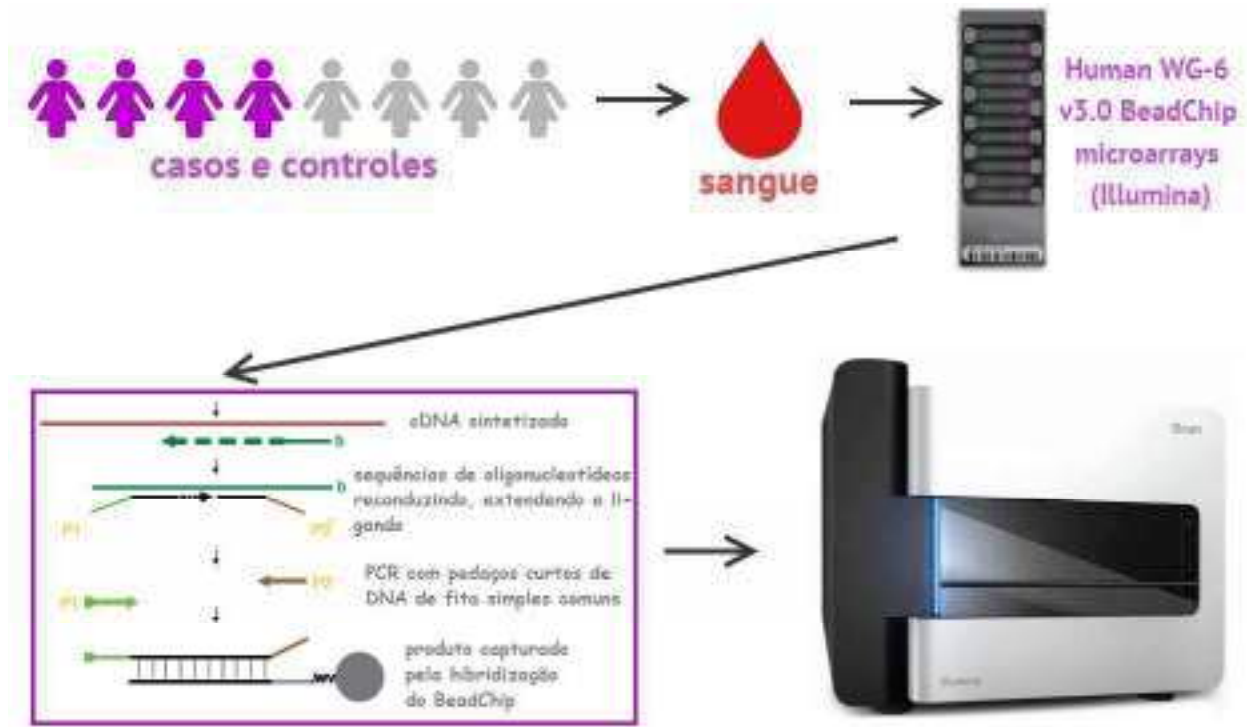


Figura 2 – Resumo esquemático do processo de obtenção das intensidades de expressão gênica. <<https://www.illumina.com/>>

2.2 Estimação de expressões gênicas

A expressão gênica é o processo pelo qual a informação codificada em uma seqüência de DNA é utilizada para a formação de um produto genético funcional. Quando o produto é uma molécula de proteína, o processo é denominado síntese de proteína ou síntese protéica e ocorre quando RNA mensageiro (RNAm) é produzido (processo chamado de transcrição) e traduzido em proteína (processo chamado de tradução). Sequências de DNA não-codificante produzem RNA funcional, que inclui elementos regulatórios de atividade gênica, como por exemplo RNA transportador (tRNAs), RNA ribossômico (rRNAs), microRNAs (miRNAs), RNAs não-codificante longo (lncRNAs). A expressão gênica é um processo que regula rigidamente se proteínas devem ser produzidas ou não e em que quantidades, permitindo que a célula responda a estímulos diversos (Griffiths e outros, 2005).

Desafios associados à medição das expressões de proteínas, como a imprecisão na quantificação de algumas moléculas de proteína, diferentes alternativas de transcrição e variações no nível da proteína, maiores que no nível do RNA, estimularam a utilização da quantidade de mRNA produzido (transcrito) para se estimar a expressão gênica (Arivaradajan e Misra, 2019). Com isso, para as análises de transcriptoma² utiliza-se a quantificação dos mRNAs em um dado organismo ou tecido. O perfil de expressão gênica em diferentes condições ambientais, estados patológicos, fisiológicos ou de desenvolvimento e finalmente a caracterização de polimorfismos associados aos genes transcritos (SNPs³ e formas alternativas de *splicing*) (Bravim, 2013).

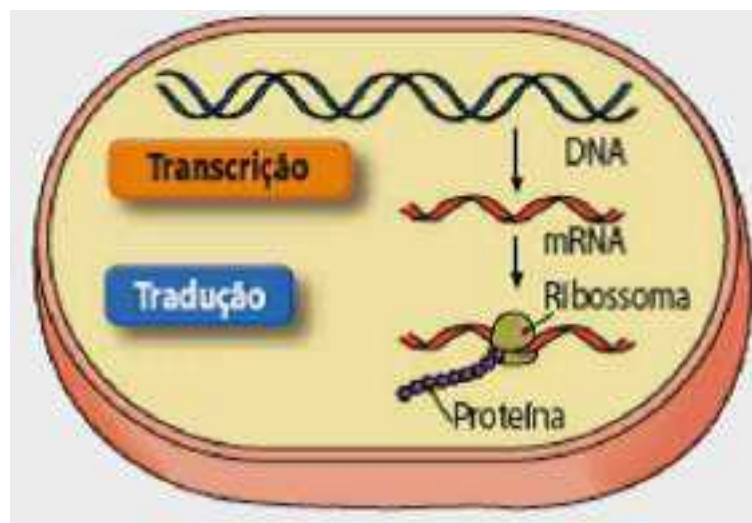


Figura 3 – Processo da síntese proteica (ou expressão gênica)

<<https://i.pinimg.com/originals/96/6d/07/966d07e9195385b3c50d294e23ef44ff.png>>

Com o objetivo de estimar os transcritos, utilizaram o chip de alto rendimento Human WG-6 v3.0 BeadChip (Illumina, 2008) da empresa Illumina para a análise das amostras hibridizadas que permite a medição de expressão do genoma inteiro e estudos integrados de genotipagem. Sua cobertura inclui mais de 48.000 transcritos escolhidos com base em dois diferentes bancos de dados que agrupam diversos genes que estão ligados a diferentes funções biológicas, são eles: *Reference Sequence (RefSeq)* e *UniGene*.

A análise de perfil de expressão gênica envolve uma sequência de etapas que podem ser divididas em **análises no nível de probe**, realizadas previamente por outros pesquisadores (Lessard e outros, 2013), e em **análises no nível de expressão gênica**, foco do presente trabalho.

² Transcriptoma é o conjunto completo (ou em larga escala) de transcritos (mRNAs, rRNAs, entre outros) de um dado organismo, órgão, tecido ou linhagem celular.

³ Variações na sequência de DNA dos humanos.

As análises no nível de probe ⁴ que foram realizadas (Lessard e outros, 2013) incluíram em sequência: procedimentos de controle de qualidade, correção de *background*, normalização e sumarização, descritas a seguir.

2.2.1 Procedimentos de controle de qualidade

O conjunto de dados utilizado envolve dados de dois experimentos de *microarray* separados e por isso procedimentos de controle de qualidade foram realizados separadamente para cada subconjunto.

Com o intuito de identificar e eliminar dados de baixa qualidade, foram excluídos os transcritos expressos em menos de 10% dos indivíduos (limiar de detecção: P-valor < 0,05) por subconjunto de dados. Foram ainda excluídos os transcritos com percentuais de valores faltantes estatisticamente diferentes entre os dois subconjuntos de dados (P-valor < 0,001 pelo teste exato de Fisher (Agresti, 2019)). Os probes relativos à genes não anotados ou com resultados inconsistentes entre os dois bancos de dados também foram excluídos.

2.2.2 Correção de *Background*, Normalização e Sumarização

O procedimento de correção de *background* tem por objetivo estimar o verdadeiro nível de hibridização do cDNA ⁵, removendo ruídos (sinais de intensidade não específicas ou *non-specific binding*) do sinal de intensidade total, detectado por leitura em scanner (Figura 2).

A etapa seguinte, de normalização, é essencial para garantir que as expressões gênicas das diferentes observações (dados de *microarray* para cada indivíduo) sejam comparáveis, removendo potenciais erros sistemáticos de medição associados aos processos de hibridização e leitura (por scanner) de imagens, problemas físicos nos chips, efeito de lote de reagentes, condições laboratoriais entre outras fontes de variação. A normalização quantílica tem por objetivo fazer duas distribuições idênticas em propriedades estatísticas, pra isso, transforma os dados impondo a mesma distribuição empírica da média obtida pela média de cada quantil nas amostras. Além disso, os dados são transformados para a escala logarítma de base 2.

Métodos de sumarização resumem os dados à nível de *probe* em uma única medida (por transcrito) que representa o nível estimado de expressão para cada transcrito. O método da mediana polida é comumente utilizado para esse fim dado que ele protege contra sondas outlier, extraindo efeitos de linhas e colunas através da subtração recursiva

⁴ Fragmento de RNA usado em genética para pesquisar um determinado gene ou outra sequência de DNA

⁵ cDNA ou DNA complementar, é o DNA sintetizado a partir de uma molécula de RNA mensageiro

das medianas de linhas e de colunas.

2.3 Análises de expressões gênicas

A avaliação dos dados de expressão gênica para os transcritos incluídos nesta etapa (após as exclusões de controle de qualidade) foi iniciada por uma análise descritiva dos dados. Em seguida, testes estatísticos e algoritmos multivariados foram utilizados para a identificação de TDEs entre grupos do estudo. Análises de agrupamento foram utilizadas para a visualização de perfis de expressão diferenciados entre os grupos de interesse. Além disso, procedimentos de anotação genômica foram realizados para identificar informações de genes (dos transcritos diferencialmente expressos que se ligam a algum gene) relacionados a sistemas biológicos associadas a Síndrome de Sjögren.

2.3.1 Análise exploratória

Gráficos e tabelas das distribuições (na escala original e em escala \log_2) foram elaborados para avaliar distribuições das expressões gênicas e identificar possíveis irregularidades e padrões nas observações. Testes de Shapiro-Wilk (Shapiro e Wilk, 1965) foram utilizados para avaliar a normalidade (hipótese nula) das distribuições de expressão gênica de todos os transcritos para casos e controles. Este teste foi escolhido por possuir maior poder para detectar desvios de normalidade por assimetria, padrão esperado para alguns transcritos.

2.3.2 Métodos de filtragem

Alguns métodos podem ser utilizados para a identificação dos chamados transcritos diferencialmente expressos (TDEs), cujas diferenças entre níveis de expressão são estatisticamente significantes. Comumente, são identificados por testes de comparação de médias ou de medianas, com correção por testes múltiplos. Modelos de regressão são utilizados quando covariáveis clínicas ou moleculares estão disponíveis.

Neste trabalho foram utilizadas três abordagens para a identificação de TDEs entre casos e controles para cada transcrito, listadas abaixo e descritas em sequência:

1. Testes T com correção de Welsh (para variâncias desiguais) para comparação de médias, seguidos por correção por testes múltiplos;
2. Testes não paramétricos de Wilcoxon-Mann-Whitney para comparação de medianas, seguidos por correção por testes múltiplos;
3. O Algoritmo PPCLUST (Partition clustering of high dimensional low sample size data based on p-values).

Sendo a primeira delas, o Teste T, a mais utilizada na literatura mas que depende da suposição de normalidade dos transcritos, que muitas vezes não é atendida, além da dependência de uma correção por testes múltiplos. A segunda, Wilcoxon-Mann-Whitney, um teste não paramétrico com os mesmos objetivos do Teste T, mas sem a necessidade de uma suposição de normalidade, mantendo-se ainda a necessidade de correção por testes múltiplos e, o PPCLUST, que mesmo sendo um algoritmo de agrupamento tem sua finalidade adaptada e é uma técnica invariante a transformações monótonas que não depende de correção por testes múltiplos.

1. Teste T - Welch de comparação de médias

Com o objetivos de identificar transcritos diferencialmente expressos, o Teste T foi utilizado. Seguindo para isso as seguintes hipóteses para cada um dos k transcritos (Morettin e Bussab, 2017):

$$H_0 : \mu_{(k\text{-ésimo caso})} = \mu_{(k\text{-ésimo controle})}$$

$$H_A : \mu_{(k\text{-ésimo caso})} \neq \mu_{(k\text{-ésimo controle})}$$

A estatística do teste é dada por:

$$T = \frac{\bar{X}_{casos} - \bar{X}_{controles}}{\sqrt{S_{casos}^2/n + S_{controles}^2/m}} \sim t(v), \quad (2.1)$$

em que \bar{X}_{casos} e $\bar{X}_{controles}$ representam médias amostrais; S_{casos}^2 e $S_{controles}^2$, variâncias amostrais e n e m , número de observações para os grupos casos e controles, respectivamente.

O número de graus de liberdade é aproximado pela expressão de Welch-Satterthwaite:

$$v = \frac{(A + B)^2}{(A^2/(n - 1) + B^2/(m - 1))}, \quad (2.2)$$

em que $A = S_{casos}^2/n$ e $B = S_{controles}^2/m$.

2. Teste de Wilcoxon-Mann-Whitney

Tendo em vista que alguns dos transcritos podem apresentar distribuições acentuadamente assimétricas, utilizou-se o teste não paramétrico de Wilcoxon-Mann-Whitney para a identificação de TDEs, cujas hipóteses são:

$$H_0 : \theta_{casos} = \theta_{controles};$$

$$H_A : \theta_{casos} \neq \theta_{controles},$$

em que θ_i representa a mediana de i . A estatística do teste é dada por:

$$W = S_n - \frac{1}{2}n(n + 1), \quad (2.3)$$

em que S_i é somas dos postos dos elementos de i , $S_n = \frac{1}{2}(m+n)(m+n+1) - S_m$ e n e m o tamanho amostral dos casos e dos controles, respectivamente.

O p-valor do teste é calculado da seguinte maneira:

$$p - \text{valor} = \begin{cases} 2P(W > W_{obs} - 1) & , \text{ se } W_{obs} > \frac{mn}{2}; \\ 2P(W > W_{obs}) & , \text{ se } W_{obs} \leq \frac{mn}{2}, \end{cases} \quad (2.4)$$

onde segue-se uma distribuição normal com média $\frac{n(n+m+1)}{2}$ e variância $\frac{nm(n+m+1)}{12}$.

3. PPCLUST

O algoritmo *Partition clustering of high dimensional low sample size data based on p-values* (PPCLUST) (von Borries e Wang, 2009) também foi utilizado para a identificação de TDEs. Trata-se de um algoritmo de agrupamento de particionamento com uma robusta medida de distância que pode agrupar variáveis de um banco de dados com altas dimensões e um tamanho de amostra baixo (HDLSS - *High dimensional low sample size data based on p-values*). É um algoritmo não hierárquico baseado em modelagem de dados por mistura de distribuições não especificadas (abordagem não paramétrica), usando P-valores como medida de similaridade e regra automática de parada.

O PPCLUST tem como base a metodologia de ANOVA não-paramétrica e tem como funcionalidade agrupar variáveis com a mesma distribuição de dados, resultando ao final do procedimento em grupos com uma variância alta entre si e pequena dentro do grupo, sem que para isso seja preciso a pré-especificação do número de grupos a serem formados.

A primeira etapa envolve o cálculo das diferenças entre expressões dos casos e a mediana dos controles.

O algoritmo é invariante a transformações monótonas mas posto que esses dados foram disponibilizados na escala \log_2 , essa escala foi mantida para o cálculo dessas diferenças, como na expressão:

$$Y_{ij} = \log_2(X_{ij}^{\text{casos}}) - \log_2(\theta^{\text{controles}}), \quad (2.5)$$

em que X_{ij}^{casos} denota o valor da i -ésima variável (transcrito) para o j -ésimo caso e $\theta^{\text{controles}}$ representa a mediana de expressão para o grupo controle. $\{Y_{ij}, 1 \leq j \leq n_i\}$ são consideradas observações independentes de alguma distribuição desconhecida $F_i(x)$, $i : 1, 2, \dots, n$. Quer-se testar a seguinte hipótese nula:

$$H_0 : F_1(x) = \dots = F_2(x) = \dots = F_n(x). \quad (2.6)$$

A Estatística do Teste se dá pela razão entre o quadrado médio do tratamento (MST_R) e o, do erro (MSE_R):

$$F_R = \frac{MST_R}{MSE_R} = \frac{\frac{1}{n-1} \sum_{i=1}^n (\bar{R}_i - \tilde{R}_{..})^2}{\frac{1}{n} \sum_{i=1}^n \frac{1}{n_i} S_{R,i}^2}, \quad (2.7)$$

em que R_{ij} representa o rank da observação Y_{ij} no conjunto de todas as observações $n_1 + n_2 + \dots + n_a$ e R_{ij} é uniformemente distribuído entre 1 e $\sum_{i=1}^a n_i$; \bar{R}_i representa a classificação média das observações para o i -ésimo transcrito; $\tilde{R}_{..}$ é definido pela média geral não ponderada de classificações de todos os transcritos e $S_{R,i}^2$ é a variância da amostra calculada usando o rank das observações do i -ésimo transcrito.

Além disso, em relação a distribuição assintótica segue que:

Teorema 1 ((Wang e Akritas, 2004) *Teste de ausência de efeito de grupo*)

Seja $H_0 : F_1(x) = \dots = F_a(x)$ atendido, com $F_i(x)$ arbitrário. Se $n_i \geq 2$ fixo, assumindo que as observações são independentes, os seguintes limites existem:

$$v_2^2 = \lim_{a \rightarrow \infty} \frac{1}{a} \sum_{i=1}^a \frac{1}{n_i} \sigma_i^2 > 0; \quad (2.8)$$

e

$$\tau_2 = \lim_{a \rightarrow \infty} \frac{1}{a} \sum_{i=1}^a \frac{2\sigma_i^4}{n_i(n_i - 1)}. \quad (2.9)$$

Então, se $a \rightarrow \infty$ tem-se que:

$$\sqrt{a}(F_R - 1) \xrightarrow{d} N(0, \tau_2/v_2^4). \quad (2.10)$$

Nesse algoritmo, para obter uma medida de similaridade com objetivo de agrupar os dados, utiliza-se o P-valor obtido com a estatística $\sqrt{a}(F_R - 1)$ e, então, compara-se o P-valor com o nível de significância α pré-especificado. O valor desse nível α pré-especificado geralmente é testado com diversos valores com o objetivo de verificar se os agrupamentos são consistentes e/ou quando vai ocorrer uma estabilização no número de grupos.

Ademais, um nível de significância usado que leva a clusters muito pequenos indica que o nível α não é pequeno o suficiente e os resultados de cluster obtidos não são confiáveis.

Além disso, quando esse nível α estável é elevado as variáveis testadas são dadas como distribuições similares e, quando pequeno, há evidências contra H_0 de que pelo menos uma variável tem distribuição diferente das demais.

Para informações sobre a implementação do PPCLUST no R ver: Lins (2018)

Correção por Testes Múltiplos

Estudos de *microarray* comumente envolvem a realização de dezenas de milhares de testes estatísticos (um para cada transcrito) para a identificação de TDEs. Se cada teste é realizado a um nível de significância α , a probabilidade de se rejeitar incorretamente pelo menos umas das hipóteses nulas será muito maior que o nível α . Métodos de correção por testes múltiplos são utilizados para controlar a probabilidade da ocorrência de erros do tipo I. O método de Benjamini e Hochberg (1995), abordagem menos conservadora para a correção de múltiplas, foi utilizado para ajustar os p-valores obtidos nas abordagens 1. e 2. Foram calculadas taxas de descobertas falsas, do inglês *False Discovery Rates* (FDRs) para controlar a proporção de falsos positivos entre o conjunto de hipóteses rejeitadas nesses conjuntos de testes.

A Tabela 1 apresenta a distribuição de testes múltiplos por resultado (significante ou não) e situação real (H_0 verdadeira ou não) (Garcia, 2019).

Tabela 1 – Tabela do método FDR - False Discovery Rate.

	Significante	Não significativo	Total
H_0 verdadeiro	F	$m_0 - F$	m_0
H_A verdadeiro	T	$m_1 - T$	m_1
Total	S	m - S	m

A proporção de resultados incorretamente declarados como significantes entre todos os resultados significantes é dada por:

$$\frac{n^\circ \text{ falsos positivos}}{n^\circ \text{ testes significativos}} = \frac{F}{S}.$$

O FDR(T) é calculado através da seguinte esperança:

$$FDR(T) = E \left[\frac{F(T)}{S(T)} \right].$$

Aqui, T representa o limiar usado para considerar os P-valores como significativos ($0 < T \leq 1$).

Sendo esse cálculo de falsos positivos feito com os P-valores obtidos abaixo do α especificado, utilizando para isso uma distribuição binomial(r, α), onde r representa a quantidade de P-valores rejeitados pelo α especificado antes da correção. Ajustando assim a distribuição real dos P-valores dos dados.

2.3.3 Filtragem por magnitude de efeito

Geralmente, seleciona-se para as análises subsequentes, apenas os TDEs com maiores diferenças entre grupos, que corresponderiam aos maiores efeitos (positivos ou negativos). Para tanto, pode-se utilizar uma métrica denominada *fold-change*, definida como a razão de médias (ou de medianas) de dois grupos. Neste trabalho, foram considerados TDEs os transcritos com *fold-change* (FC), definido pela razão entre as médias (ou medianas) de casos e controles, superior a 1,5 ou inferior a 0,67.

Para os transcritos selecionados pela abordagem 1. (por testes t ajustados por testes múltiplos) FC foi calculado por:

$$FC_i = \frac{\bar{X}_i^{casos}}{\bar{X}_i^{controles}}, \quad (2.11)$$

em que \bar{X}_i^l representa a média para o i-ésimo transcrito de t (casos ou controles).

Já para as abordagens 2. (por testes de Wilcoxon-Mann-witney) e 3. (pelo algoritmo PPLUST) bem como para a interseção e união das abordagens 1. a 3., FC foi calculado por:

$$FC_i = \frac{\theta_i^{casos}}{\theta_i^{controles}}, \quad (2.12)$$

em que θ_i^l representa a mediana para o i-ésimo transcrito de t (casos ou controles).

2.3.4 Análises de agrupamento

Análises de Agrupamento foram realizadas objetivando a definição de padrões de expressão gênica entre casos e controles com base nos transcritos selecionados, separando eles em grupos, baseando-se nas características que possuem. A idéia básica consiste em colocar em um mesmo grupo o que é similar de acordo com algum critério pré-determinado (Linden, 2009).

Além disso, as Análises de agrupamento são exemplos de técnicas de aprendizado não supervisionado em que as categorias de classificação (os agrupamentos) são desconhecidas, mas observa-se recursos relacionados às categorias não observadas (Agresti, 2019). Portanto, essas análises visam encontrar algoritmos ótimos para agrupar grupos que são semelhantes de acordo com algum critério, não existindo grupos pré-definidos nem quantos grupos se formarão.

Além disso, podem ser classificadas como de particionamento e métodos hierárquicos. Um método de particionamento constrói k clusters, ou seja, classifica os dados em k grupos, que juntos atendem aos requisitos de uma partição (cada grupo deve conter pelo menos um objeto e cada objeto deve pertencer a exatamente um grupo) (Kaufman e

Rousseeuw, 2019), um exemplo de um algoritmo desse tipo é o próprio PPCLUST. Já os algoritmos hierárquicos não constroem uma única partição com *clusters* k , mas lidam com todos os valores de k na mesma execução, ou seja, a partição com $k = 1$ (todos os objetos estão juntos no mesmo *cluster*) faz parte dos resultados e também a situação com $k = n$ (cada objeto forma um *cluster* separado com apenas um elemento). Intermediariamente, todos os valores de $k = 2, 3, \dots, n - 1$ são abordados em um tipo de transição gradual, com a diferença que entre $k = r$ e $k = r + 1$ os agrupamentos de r se dividem para obter os agrupamentos de $r + 1$ (ou, de maneira diferente), dois dos *clusters* $r + 1$ se combinam para produzir o *cluster* r (Kaufman e Rousseeuw, 2019). Neste trabalho foi utilizado um algoritmo hierárquico, que pode ser subdividido de duas maneiras, descritas a seguir e ilustradas na Figura 4:

- Aglomerativos: Cada item de dados é considerado como um grupo individual, e grupos são recursivamente fundidos de acordo com algum critério de proximidade até produzir um bom agrupamento no final do processo (Miranda, 2012).
- Divisivos: Inicialmente, o conjunto de todos os dados é considerado como sendo um único grupo e, em seguida ele é recursivamente dividido de acordo com algum critério de proximidade para produzir um bom agrupamento final (Miranda, 2012).

Técnicas Aglomerativas e Divisivas constroem sua hierarquia em direção oposta e podem produzir resultados distintos (Kaufman e Rousseeuw, 2019).

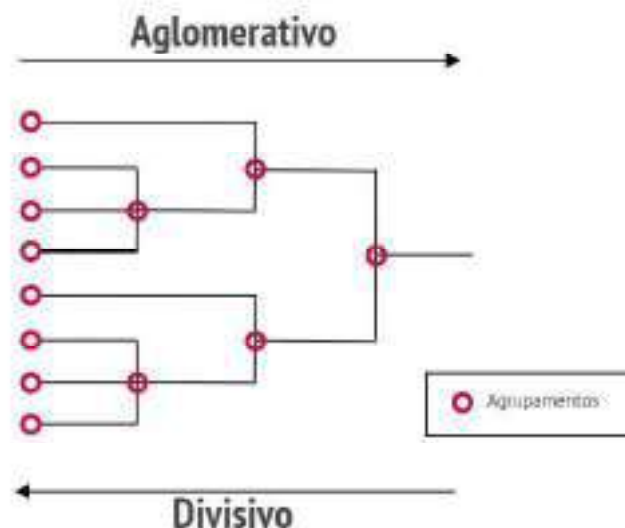


Figura 4 – Ilustração de agrupamento hierárquico aglomerativo e divisivo.

Neste trabalho foi utilizado um algoritmo Hierárquico Divisivo. A medida de similaridade utilizada foi a distância euclidiana, que além de ser a mais utilizada em dados

em escala LOG, também reduz a variabilidade desses dados por se tratar de uma medida quadrática.

Distância Euclidiana:

$$d(x, y) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2}, \quad (2.13)$$

p : número de variáveis;

x : i-esimo transcrito para o agrupamento dos transcritos ou o j-ésimo caso ou controle para o agrupamento de casos ou controles;

y : i-esimo transcrito para o agrupamento dos transcritos ou o j-ésimo caso ou controle para o agrupamento de casos ou controles, diferente de x .

O agrupamento hierárquico é normalmente apresentado em um dendrograma e para uma visualização ainda melhor desses resultados, em genética principalmente, é comumente usado a apresentação desses resultados em *heatmaps* com dendrogramas (na parte superior e ao lado) geralmente, como na Figura 5, retratando o processo de mesclagem dos *clusters* em função da métrica de proximidade. A associação dessas duas ferramentas gráficas possibilita melhor visualização das magnitudes dos efeitos. Além disso, uma separação por grupos, também é possível, como é o caso da separação entre casos e controles.

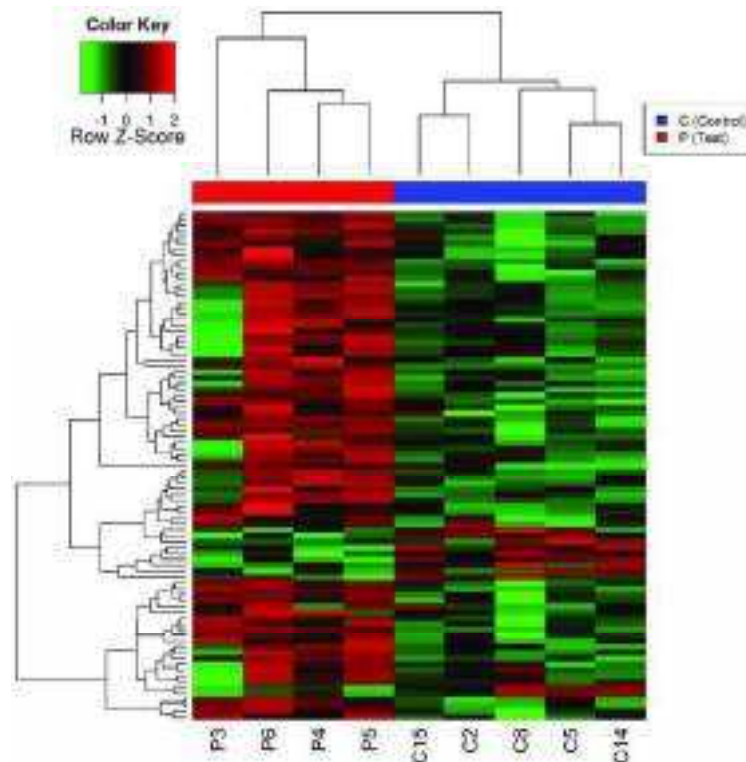


Figura 5 – Exemplo de *Heatmap* de perfil de expressão gênica.

2.3.5 Anotação genômica

A anotação genômica, do inglês *genomic annotation*, pode ser definida como o processo de descoberta de componentes importantes do genoma, principalmente genes e seus produtos, adicionando a eles análises e interpretações necessárias para se extrair sua importância biológica e colocando-os no contexto dos processos biológicos, pressupondo para isso, uma forte correlação entre coexpressão e a relação funcional (Stein, 2001).

Tipicamente, essa análise envolvem o denominado *pathway analysis* (análise de caminhos), que tem como objetivo identificar proteínas vinculadas a transcritos que são relacionadas dentro de uma função biológica ou que ligam uma função biológica a outra. Esse procedimento permite encontrar processos celulares distintos, doenças ou vias de sinalização estatisticamente associadas à seleção de genes diferencialmente expressos entre as duas amostras (Aguiar e Severino, 2010).

A análise de caminhos precisa de uma base de conhecimento com redes de coleta e interação de caminhos. O conteúdo, a estrutura e a funcionalidade das coleções de caminhos geralmente variam em diferentes bases de dados. Sendo aqui utilizada a coleção de caminhos da base **Reactome** (reactome.org, 2020).

A base **Reactome** é um recurso com curadoria on-line para dados de caminhos humanos, com base no qual pode-se inferir reações equivalentes em várias espécies não humanas. É utilizado também como recurso de aprendizado além de ser uma ferramenta computacional para auxiliar na interpretação de microarranjos e conjuntos de dados semelhantes em larga escala (Vastrik e outros, 2007) (reactome.org, 2020).

A base de dados **Reactome** tem como principal característica que as identificações de cada gene (colunas IDs, UniProt, símbolos genéticos ou IDs ChEBI) são mapeados para a análise de caminhos, análises de super-representação e da análise de conectividade entre moléculas.

A super-representação diz respeito ao cálculo de uma lista de proteínas, obtidas através de algumas das identificações supracitadas que contém mais anotações em cada caminho da base **Reactome** do que seria esperado. A análise de super-representação envolve um teste estatístico que utiliza a distribuição hipergeométrica e determina se certas vias da base **Reactome** estão super-representadas (enriquecidas) nos dados enviados.

Por meio desse cálculo da distribuição hipergométrica, um P-valor é obtido e é calculado um ajuste para comparações múltiplas por Dunn-Sidak (Abdi, 2007) ⁶ que se assemelha bastante com o método por FDR, visando controlar o Erro do Tipo I.

⁶ Ajuste derivado do método de Bonferroni (Bland e Altman, 1995)

Esse ajuste se dá seguindo a seguinte equação:

$$p' = 1 - (1 - p)^r, \quad (2.14)$$

em que p' representa o P-valor ajustado por Dunn-Sidak, p representa o P-valor obtido pela distribuição hipergeométrica e r é igual ao número de hipóteses (Biocomputing, 2007).⁷. Além disso, a medida *gene ratio* (razão de genes ou razão gênica) definida pela porcentagem do total de genes diferencialmente expressos é também utilizada na análise de anotação.

⁷ Dunn-Sidak é derivado da correção de Bonferroni, este que possui $p' = \frac{p}{k}$

3 Resultados

Procedimentos no nível de *probe*, conduzidos por (Lessard e outros, 2013), envolveram a eliminação de dados de baixa qualidade com a filtragem dos transcritos que foram expressos em menos de 10% dos indivíduos (limiar de detecção: $P < 0,05$) e probes com taxas diferenciais ausentes entre casos e controles ($P < 0,001$ pelo teste exato de Fisher). Em seguida, os dados foram transformados para uma escala logarítmica de base 2 e foi realizada uma Normalização Quantílica. Tais procedimentos levaram à seleção de 15.063 (dos 48.803 transcritos incluídos no chip) para os 190 casos de pSS e 32, do grupo controle utilizados nas análises no nível de expressão gênica, cujos resultados são apresentados a seguir.

3.1 Análise exploratória

Os boxplots das Figuras 6 e 7 apresentam os dados de expressão gênica em escala original e em escala \log_2 por grupo. Observa-se acentuada assimetria na escala original em ambos os grupos, o que evidencia a necessidade da transformação em escala logarítmica. Como tais gráficos incluem expressões dos 15.063 transcritos selecionados não é possível se observar diferenças nas distribuições dos grupos de casos e controles.

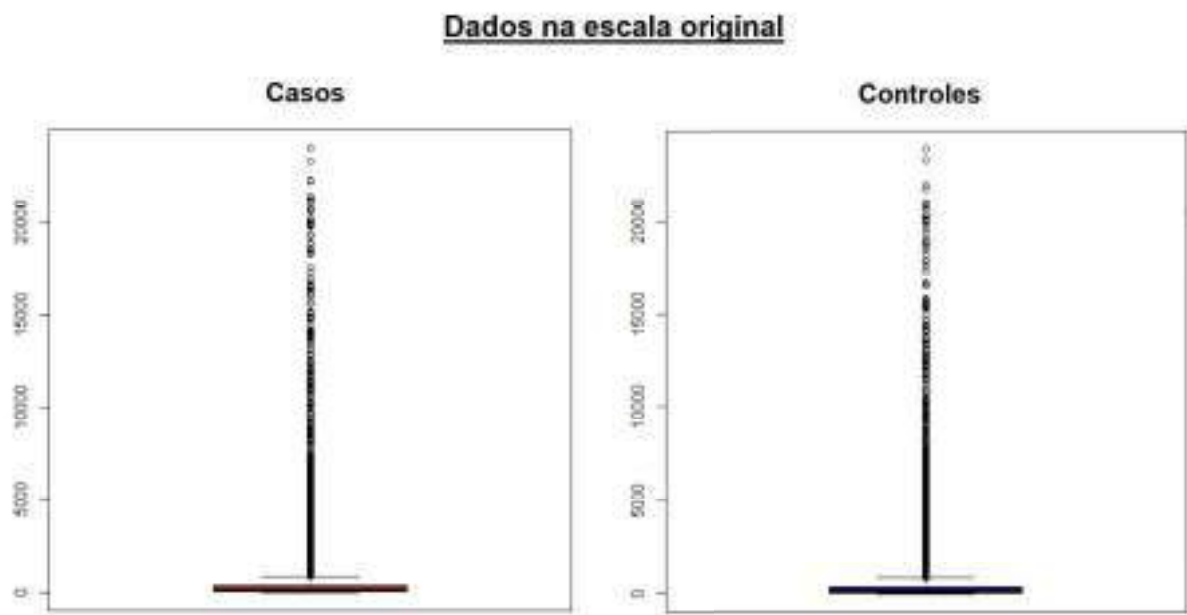


Figura 6 – Distribuição de expressões gênicas de casos e controles em escala original.

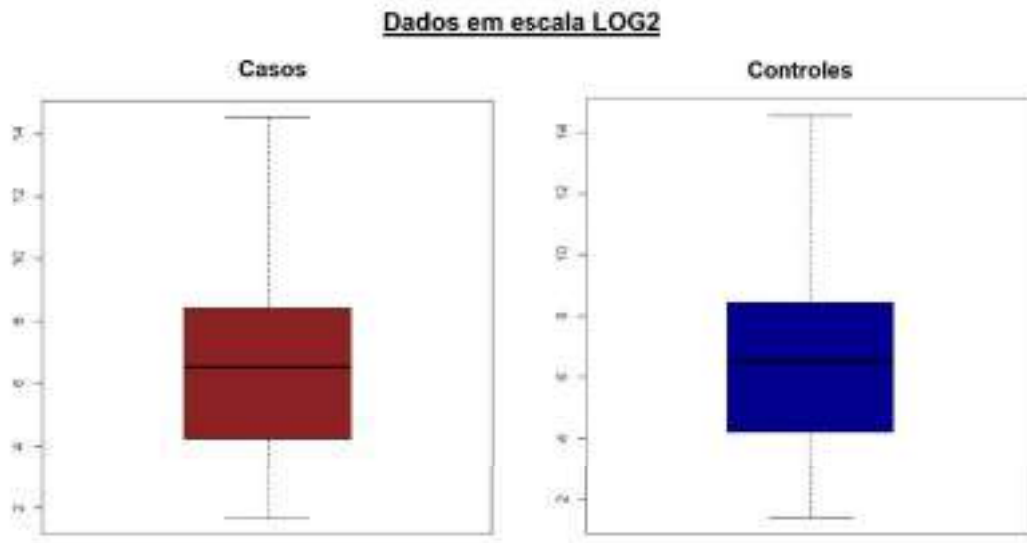


Figura 7 – Distribuição de expressões gênicas de casos e controles em escala log2.

A distribuição das médias dos transcritos (Figuras 8 e 9) destaca que apenas uma pequena parcela dos transcritos possui expressões médias elevadas em ambos os grupos.

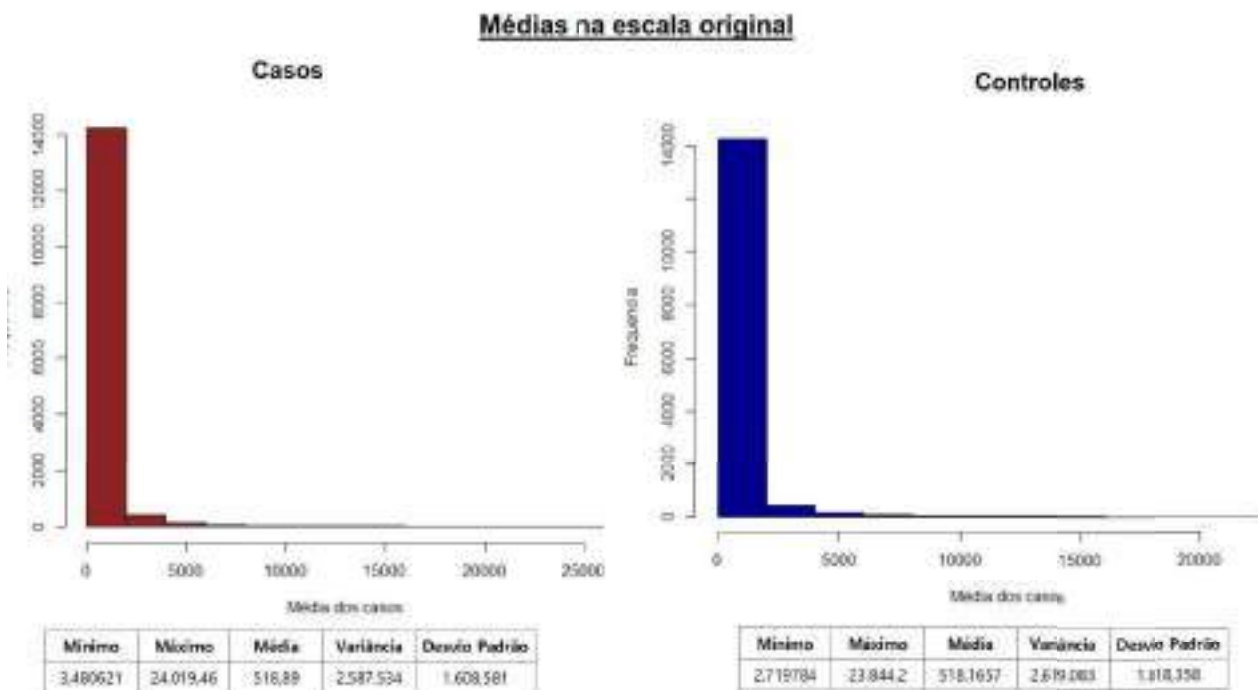


Figura 8 – Distribuição das expressões médias (de todos os transcritos) de casos e controles em escala original.

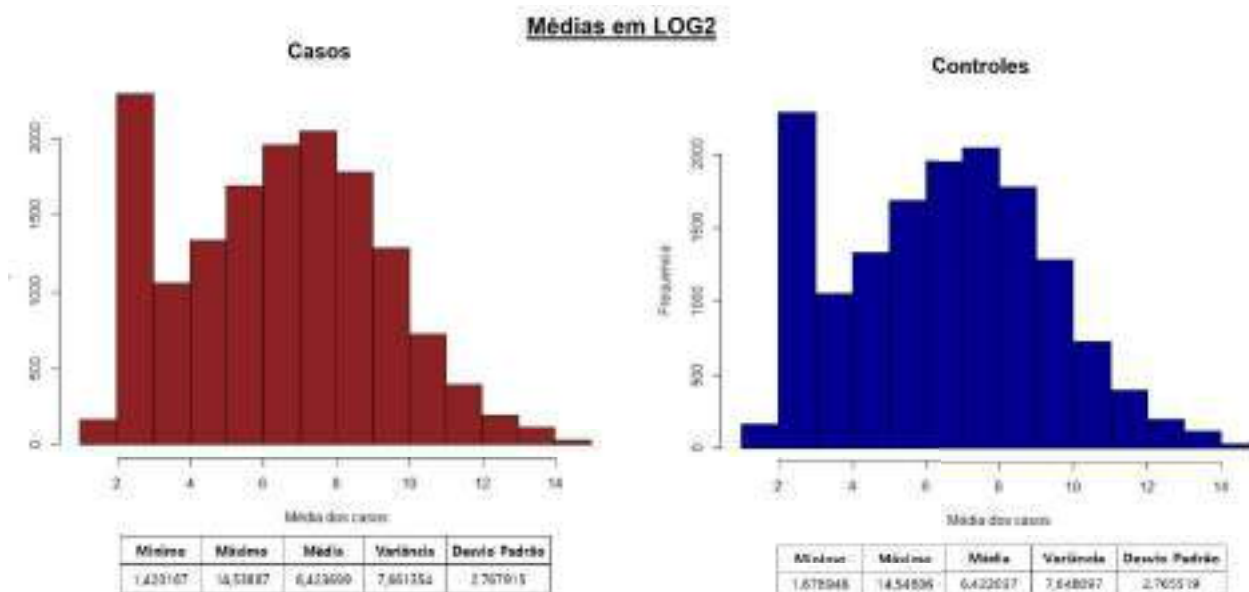


Figura 9 – Distribuição das expressões médias (por transcrito) de casos e controles em escala log2.

A Figura 10 apresenta as distribuições de FC dos 15.063 transcritos, definidas pela razão de médias e de medianas (de casos por controles), respectivamente. Para ambas as medidas, observa-se que as razões FC se concentram em valores muito próximos de 1 para a vasta maioria dos transcritos.

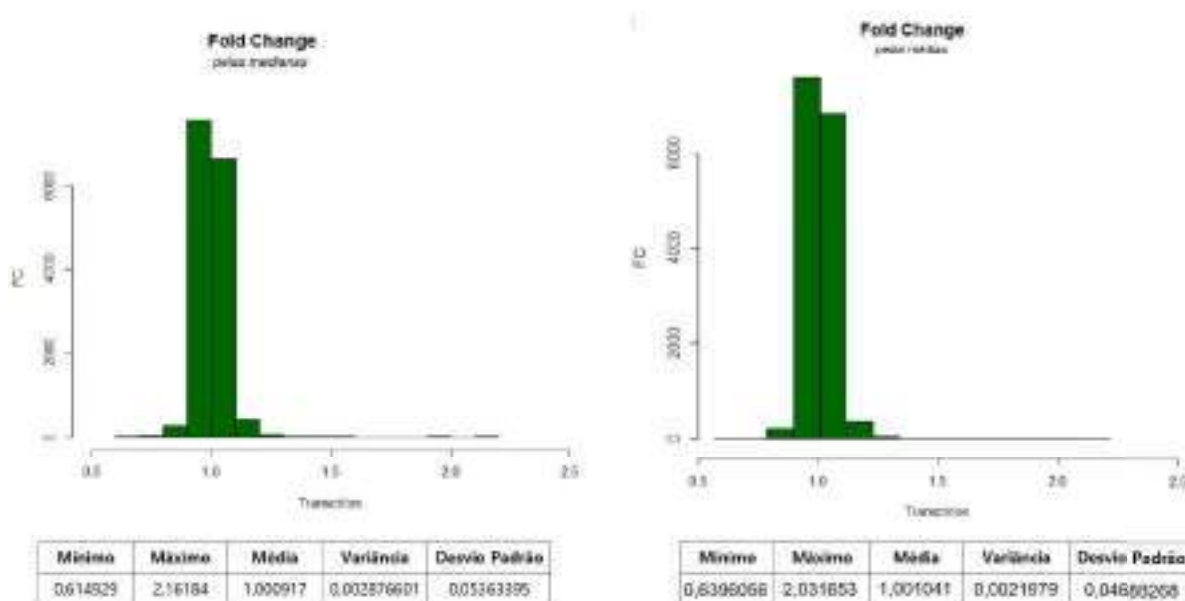


Figura 10 – Distribuição de FC por médias e de FC por medianas.

Para a implementação do algoritmo PPCLUST, para cada transcrito, foram utilizadas as expressões gênicas dos casos subtraídas pela mediana das expressões do grupo controle em escala log2 (escala que os dados foram disponibilizados), chegando então a: $(\log_2(X^{casos_i}) - \log_2(\theta_{controles}))$. A Figura 11 apresenta a distribuição das médias desses valores.

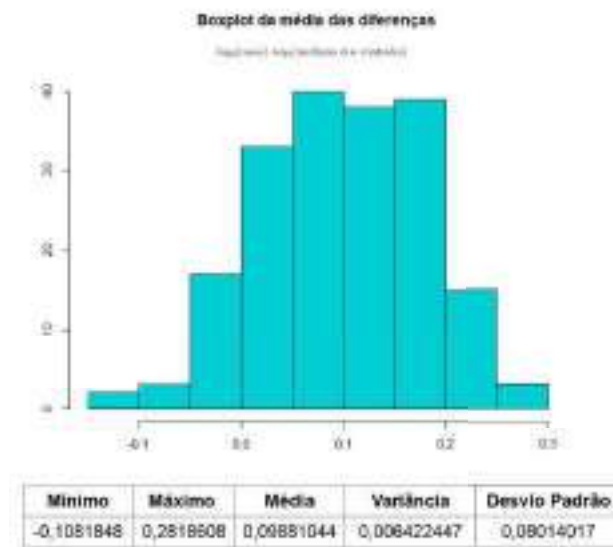


Figura 11 – Distribuição dos valores utilizados para o PPCLUST.

A Figura 12 apresenta a distribuição de p-valores de testes de Normalidade (Shapiro-Wilks) para os 15.063 transcritos em log2. Considerando-se um nível $\alpha=0,05$, observou-se que 8.469 (aproximadamente 44%) dos transcritos apresentaram evidência contra a hipótese nula de normalidade. Tais resultados apontam para a necessidade da utilização das abordagens por Wilcoxon-Mann-Whitney e PPCLUST na identificação de TDEs.

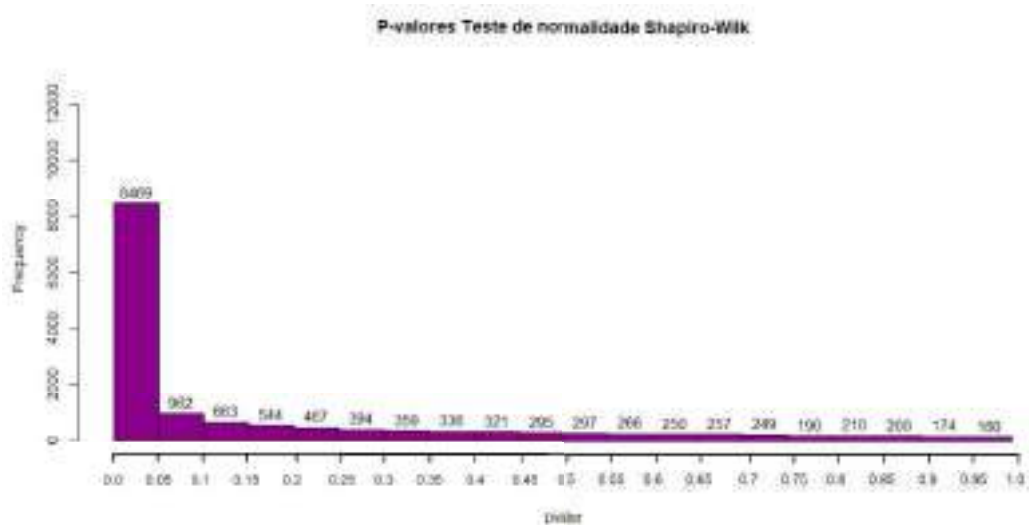


Figura 12 – Distribuição de P-valores do Testes de normalidade Shapiro-Wilk.

3.2 Análise de perfil de expressão gênica

Após as análises exploratórias preliminares iniciaram-se os procedimentos para a identificação de TDEs sob as três diferentes abordagens (Tabela 2).

Para a escolha do nível α adequado no algoritmo PPCLUST, de modo que o número de grupos definidos fosse estável, foram testados níveis α de $10e^{-1}$, $10e^{-4}$, $10e^{-8}$ e $10e^{-12}$, que resultaram, respectivamente, em 13,13,14 e 13 grupos (Figura 13).

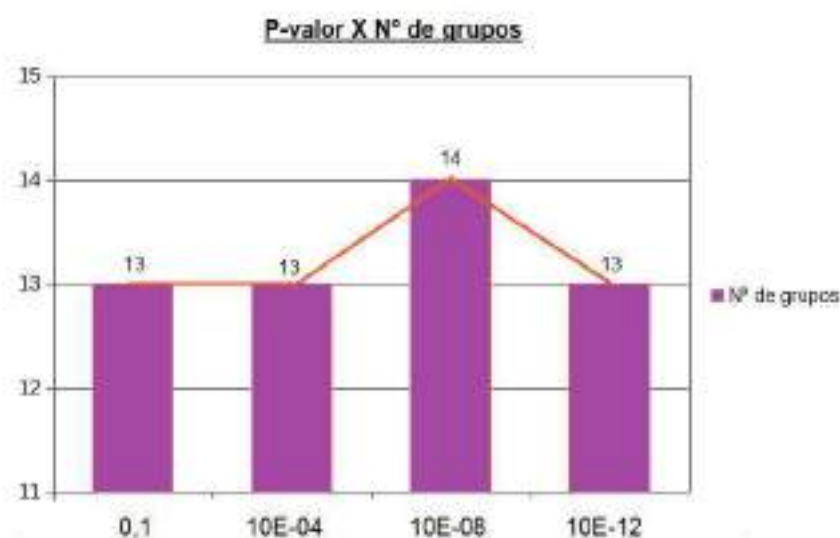


Figura 13 – Número de grupos definidos pelo algoritmo PPCLUST por P-valor.

Atingiu-se uma certa estabilidade em 13 grupos com o menor dos níveis α 's testados de $10e^{-1}$. O agrupamento resultante é apresentado na Figura 14. Destaca-se que o grupo 0 foi criado para indicar transcritos que não puderam ser incluídos em nenhum dos grupos criados (são sobras do processo de agrupamento).

Grupo	0	1	2	3	4	5	6	7	8	9	10	11	12	13
Quantidade transcritos	35	4	92	647	17	1865	2616	3586	2964	1976	891	277	89	4

↑ Transcritos que não se encaixam em nenhum grupo

Grupos extremos

Figura 14 – Estrutura dos grupos obtidos pelo algoritmo PPCLUST.

Após o processo de agrupamento, transcritos dos grupos extremos 1 a 4 e 10 a 11 foram considerados TDEs, totalizando 2.021 transcritos.

A separação em grupos extremos e não extremos foi arbitrária e baseada na visualização da distribuição dos transcritos que apresentassem menor quantidade nos grupos considerados extremos que nos demais grupos, sugerindo que esses seriam TDEs.

Para efeitos de consistência com as outras abordagens, uma filtragem adicional por magnitude de efeito foi efetuada, excluindo-se transcritos com FC de medianas entre $(1.5)^{-1}$ e 1.5.

Os números de TDEs identificados por cada uma das abordagens, incluindo seus respectivos cortes por magnitude e correções por testes múltiplos são apresentados na Tabela 2. Destaca-se que, cerca de 50 % dos transcritos apresentou significância genômica (p-valores corrigidos pela taxa de FDR $< 0,05$) e apenas 10 % apresentou magnitudes de efeitos maiores, definidos por $FC < (1.5)^{-1}$ ou $FC > 1.5$.

Tabela 2 – Número de transcritos selecionados para cada etapa das abordagens de identificação de TDEs.

Etapa/Abordagem	Teste t	Wilcoxon-Mann-Whitney	PPCLUST
Após a realização do teste	4.214	4.340	2.021
Com correção pela taxa FDR	2.133	2.165	-
Com corte por magnitude (<i>fold change</i>)	232	257	295

O diagrama de Venn (Figura 15) abaixo apresenta as relações entre os conjuntos de TDEs, identificados pelas 3 abordagens utilizadas. Observa-se que 318 transcritos foram identificados por qualquer das 3 técnicas (União dos conjuntos). Destaca-se que 50% destes (158 transcritos) foram detectados pelas 3 diferentes abordagens, enquanto que 35% foram detectados por apenas uma. Apesar da abordagem baseada no teste t detectar um menor número de TDEs que a abordagem baseada no teste de Wilcoxon-Mann-Witney, apenas 3 transcritos foram identificados unicamente pelo último. Cinquenta e seis transcritos foram detectados unicamente pela abordagem PPCLUST e 53, unicamente pela abordagem baseada no teste t.

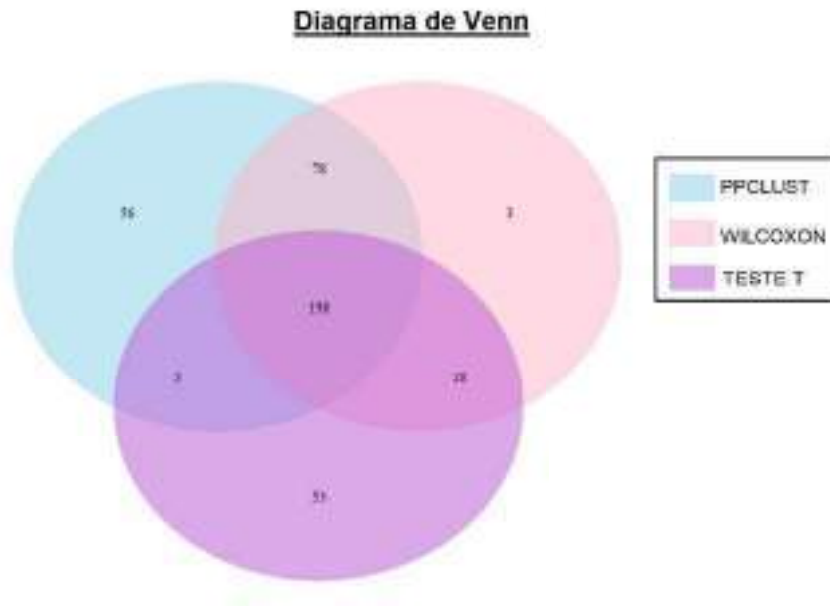


Figura 15 – Diagrama de Venn das 3 listas de transcritos (resultantes de testes paramétricos, não paramétricos e PPCLUST).

Em sequência, são apresentados conjuntos de gráficos para cada uma das cinco listas de TDEs, descritas a seguir (completar):

- Lista **L1**: TDEs identificados por teste t, corrigidos pela taxa de FDR e com $FC_{medias} > 1,5$ ou $FC_{medias} < 1,5^{(-1)}$;
- Lista **L2**: TDEs identificados por teste Wilcoxon-Mann-Witney, corrigidos pela taxa de FDR e com $FC_{medianas} > 1,5$ ou $FC_{medianas} < 1,5^{(-1)}$;
- Lista **L3**: TDEs identificados por algoritmo PPCLUST e com corte pelo $FC_{medianas} > 1,5$ ou $FC_{medianas} < 1,5^{(-1)}$;
- Lista **L4**: Interseção das listas L1, L2 e L3;
- Lista **L5**: União das listas L1, L2 e L3.

As Figuras 16 a 20 apresentam os perfis de expressão gênica para os TDEs selecionados para as listas **L1** a **L5**. Tais figuras evidenciam nitidamente a existência de grupos de transcritos super-expressos (em vermelho) e sub-expressos (em verde) no grupo de casos, comparado ao grupo de controles. Pode-se observar uma concentração maior de transcritos super-expressos entre os detectados como TDEs pelas abordagens 1 e 2 e pela intersecção (**L4**). A abordagem que envolve o algoritmo PPCLUST detectou mais transcritos sub-expressos que as outras abordagens.

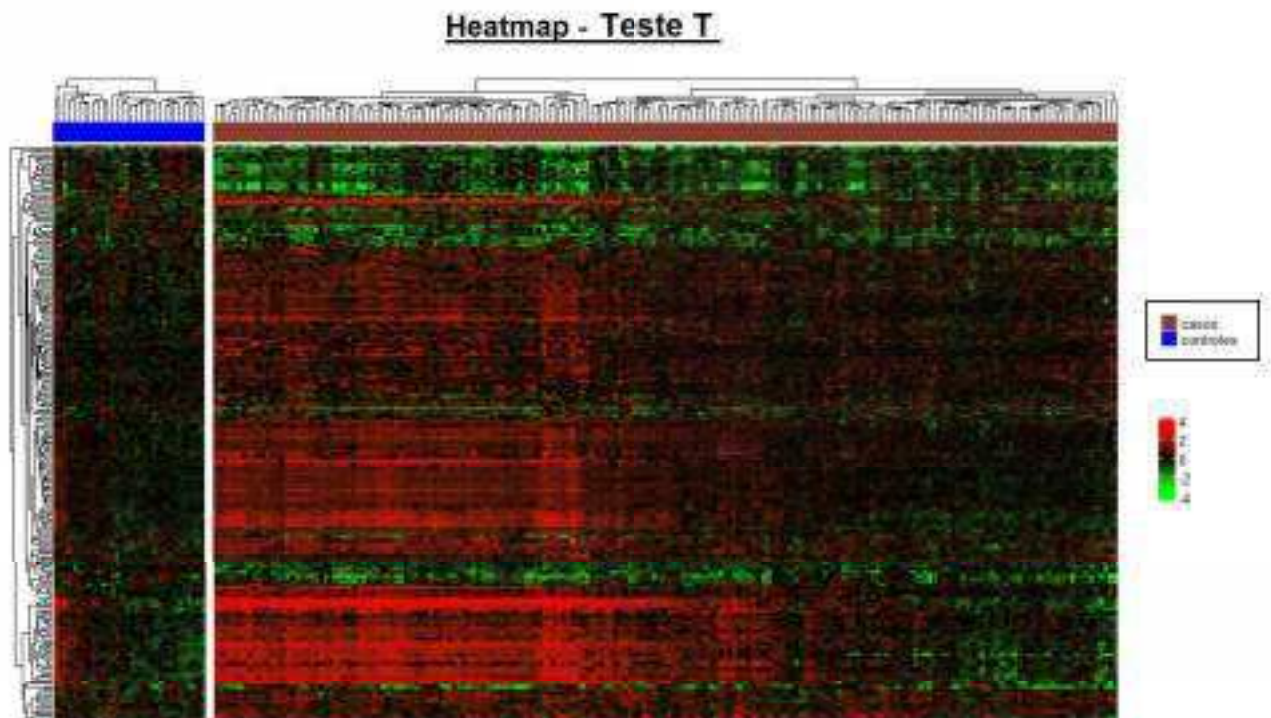


Figura 16 – *Heatmap* das expressões gênicas em escala log2 para a Lista L1 (testes t corrigidos por FDR e FCs por médias)

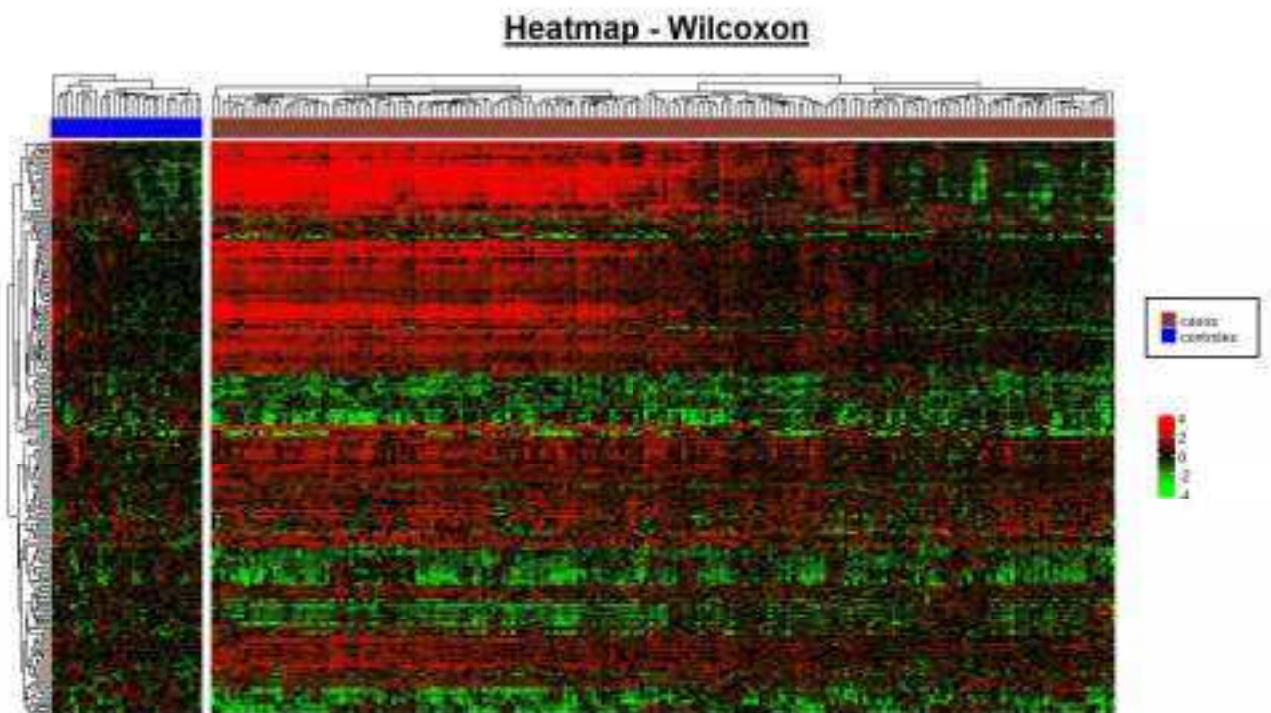


Figura 17 – *Heatmap* das expressões gênicas em escala log2 para a Lista L2 (testes Wilcoxon-Mann-Whitney corrigidos por FDR e FCs por medianas)

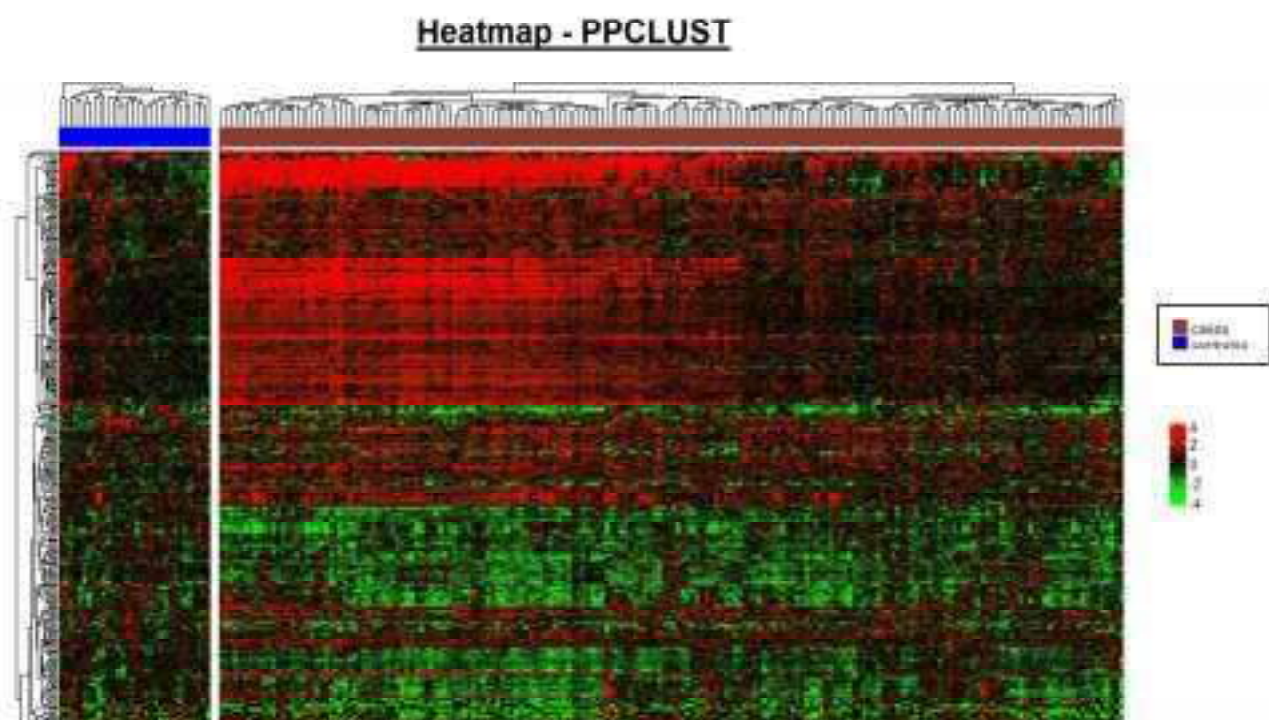


Figura 18 – *Heatmap* das expressões gênicas em escala log2 para a Lista L3 (algoritmo PPCLUST e FCs por medianas)

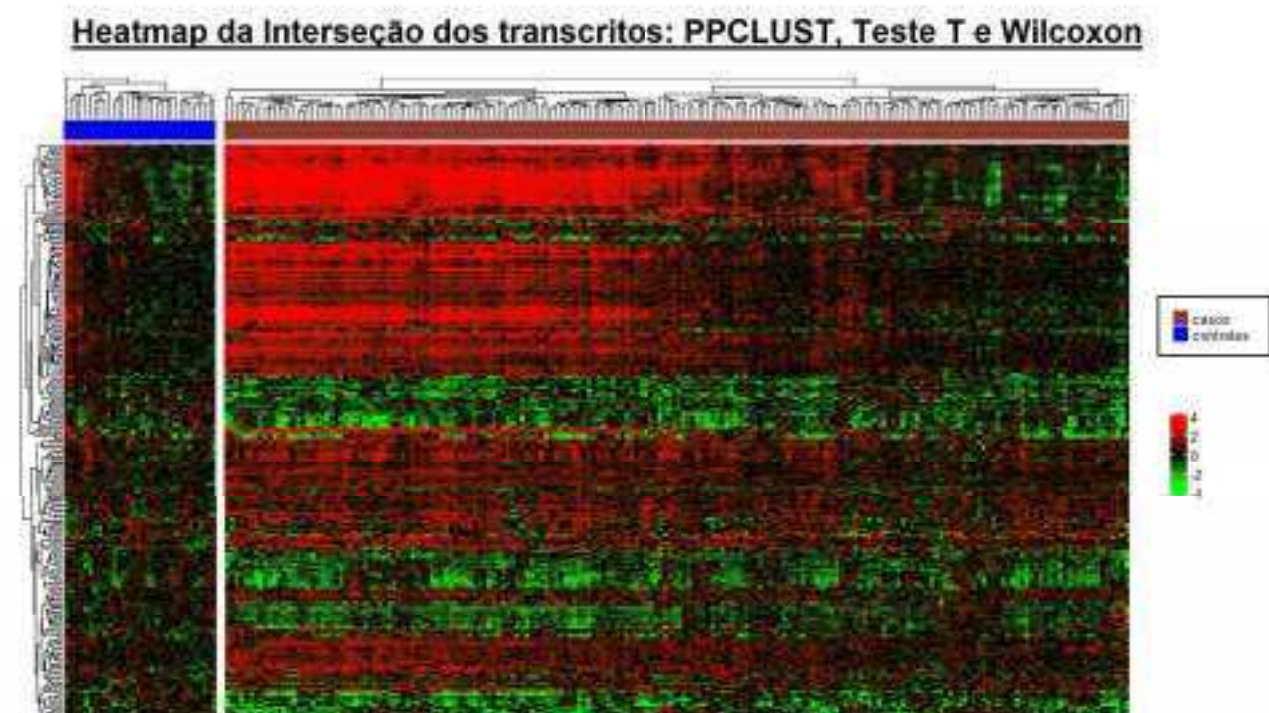


Figura 19 – *Heatmap* das expressões gênicas em escala log2 para a Lista L4 (Intersecção das abordagens)

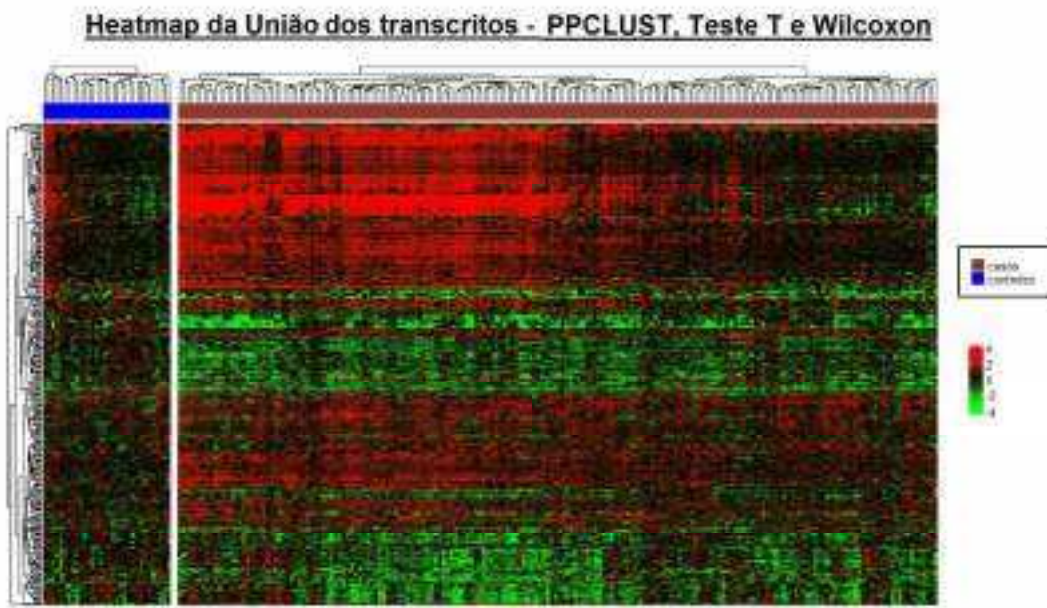


Figura 20 – *Heatmap* das expressões gênicas em escala \log_2 para a Lista L5 (União das abordagens)

A Tabela 3 apresenta os 10 genes para maiores FCs entre a união das 3 abordagens (**L5**). Destaca-se entre estes, genes de interteron (*IFI44L*, *IFI27*, *IFI44*, *IFIT1*) ou induzidos por diferentes tipos de interferons bem como outros genes, relacionados a função de imunidade e defesa contra vírus e bactérias.

Tabela 3 – Genes mais sub/super-expressos na lista de união dos TDEs detectados pelas 3 abordagens (L5).

GENE	FUNÇÃO	SUB/ SUPER- EXPRESSO NOS CASOS
<i>CMPK2</i>	Ligada a resposta celular ao lipopolissacarídeo. O lipopolissacarídeo é um dos componentes principais da membrana exterior da bactéria gram-negativa .	Super-expresso
<i>RSAD2</i>	É induzível por intelferon tipo I e tipo II e desempenha um papel importante no estado antiviral celular. Pode inibir uma ampla gama de vírus de DNA e RNA, incluindo citomegalovírus humano, vírus da hepatite C, vírus do Nilo Ocidental, vírus da dengue, vírus sindbis, vírus influenza A, vírus sendai, vírus estomatite vesicular e vírus da imunodeficiência humana (HIV-1).	Super-expresso
<i>IFI44L</i>	Um gene comum de ser encontrado em pacientes com doenças autoimunes que também é estimulado por interferon e é ligado a resposta de defesa imune e a resposta de defesa à vírus.	Super-expresso
<i>IFI27</i>	Pode desempenhar um papel na resposta vascular à lesão e é um novo modulador de respostas imunes inatas que regulam os receptores nucleares anti-inflamatórios.	Super-expresso
<i>LAMP3</i>	A expressão elevada de LAMP3 está associada à transição epitelial-mesenquimal (as células epiteliais perdem sua polaridade celular e adesão celular e ganham propriedades migratórias) e potencial metastático no câncer de esôfago.	Super-expresso
<i>OTOF</i>	Mutações nesse gene são uma causa de surdez recessiva não síndrômica neurossensorial.	Super-expresso
<i>IFI44</i>	É um paralogue do IFI44L que é estimulado por intelferon e ligado a respostas imunes, resposta às bactérias e a vírus.	Super-expresso
<i>TUBB2A</i>	Principal constituinte dos microtúbulos, que são estruturas proteicas que fazem parte do citoesqueleto nas células eucarióticas. Além de estar ligado ao processo do ócio celular mitótico e a migração de neurônios.	Super-expresso
<i>EPSTI1</i>	Gene de resposta ao interferon que foi originalmente identificado como um gene induzido por fibroblastos estromais no câncer de mama, mas novos estudos têm mostrado que ele está relacionado também a suscetibilidade à doença inflamatória crônica, lúpus eritematoso sistêmico.	Super-expresso
<i>IFIT1</i>	Estimulado por interferon e ligado a resposta de defesa imune e a resposta de defesa à vírus.	Super-expresso

(Omim,2020) (Uniprot, 2020) (Genecards, 2020)

3.3 Identificação de sistemas biológicos e funções genéticas

Os resultados das análises de anotação são apresentados em conjuntos de 5 gráficos (para as listas L1 a L5).

Gráficos de barras da razão gênica:

Os gráficos de barras do número de genes para cada uma das listas (Figuras 21 a 25) apresentam os sistemas mais significantes identificados com base no número de genes de cada lista (**L1** a **L5**). Utilizam gradação de cor de acordo com o P-valor ajustado que advém da hipergeométrica do teste que é realizado.

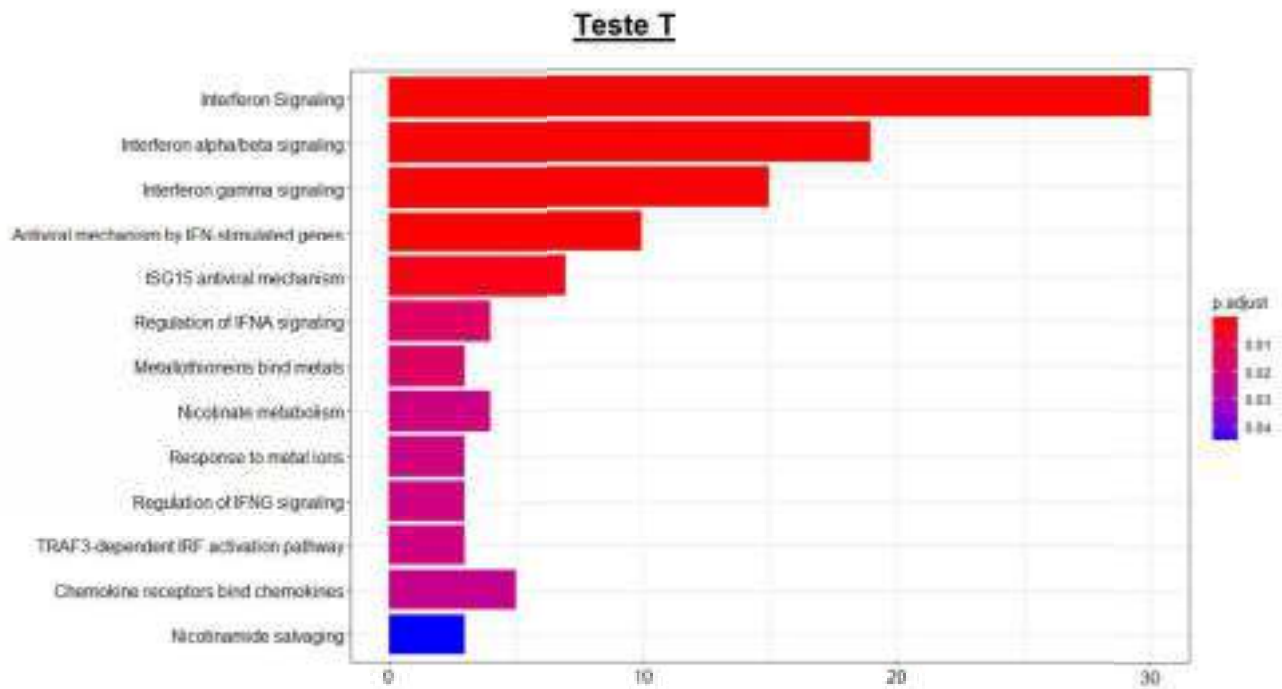


Figura 21 – Distribuição do número de genes da Lista **L1** detectados nos sistemas biológicos mais significantes.

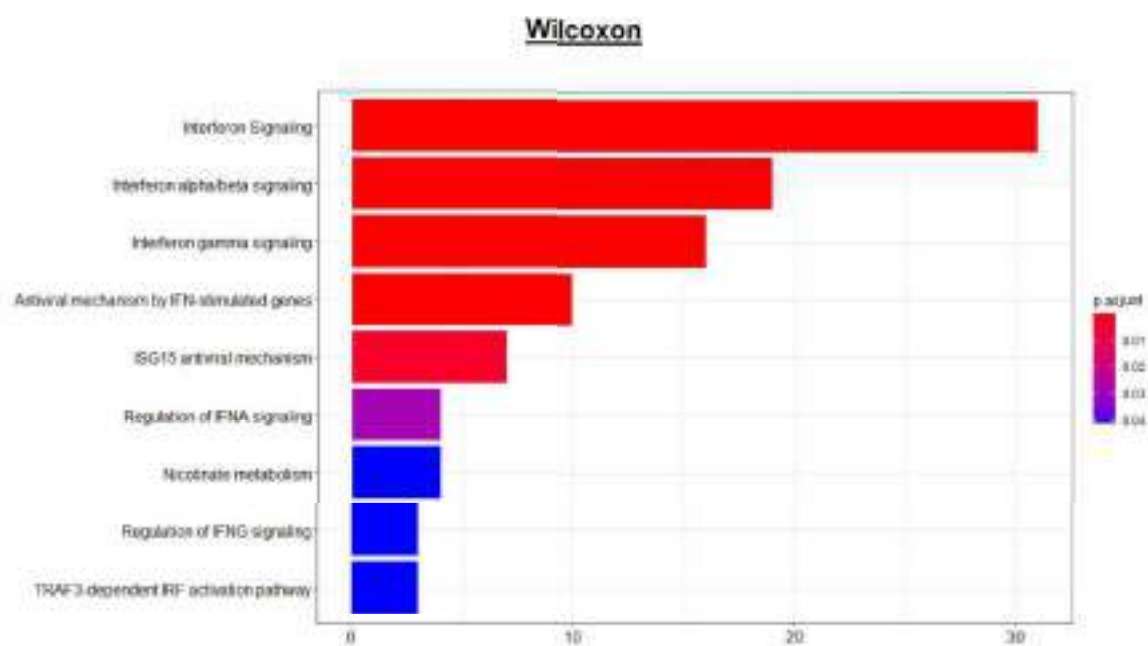


Figura 22 – Distribuição do número de genes da Lista **L2** detectados nos sistemas biológicos mais significantes.

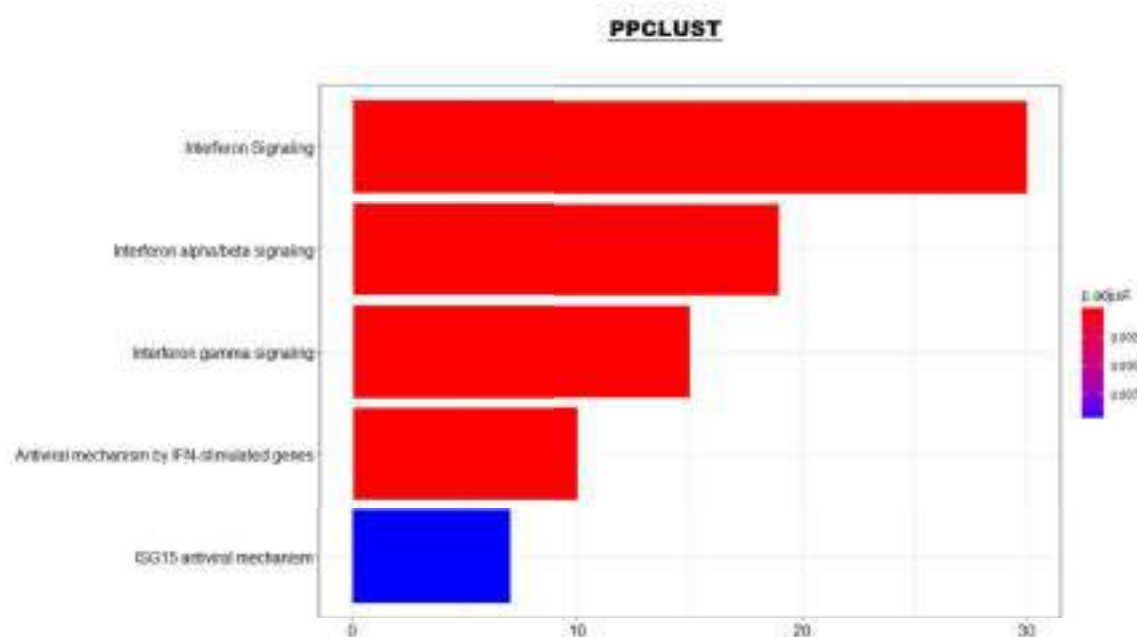


Figura 23 – Distribuição do número de genes da Lista **L3** detectados nos sistemas biológicos mais significantes.

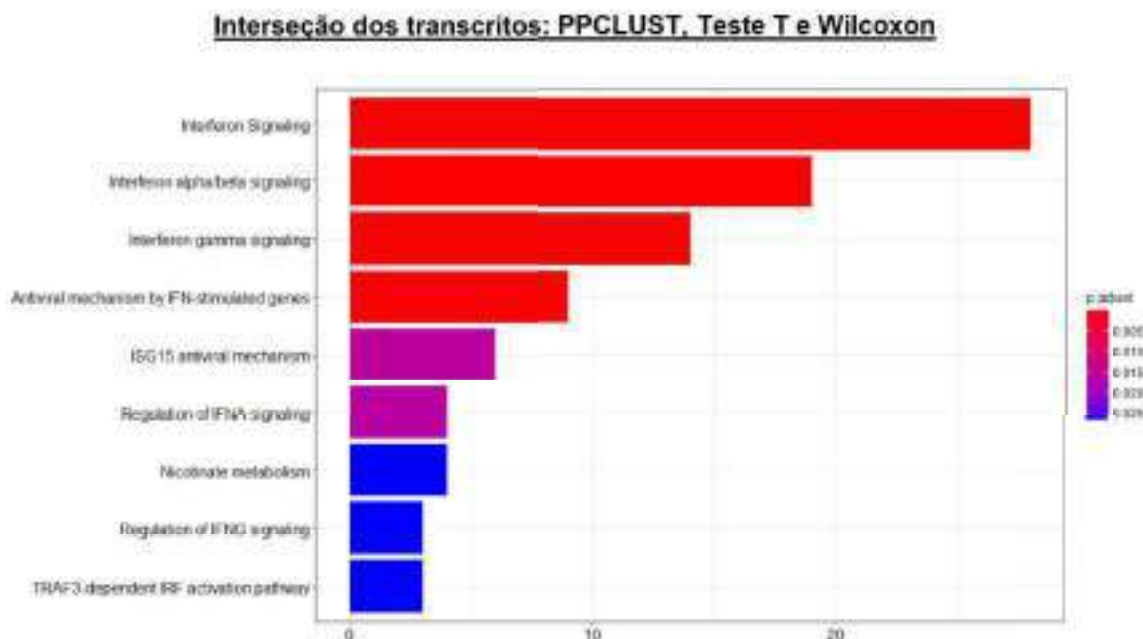


Figura 24 – Distribuição do número de genes da Lista **L4** detectados nos sistemas biológicos mais significantes.

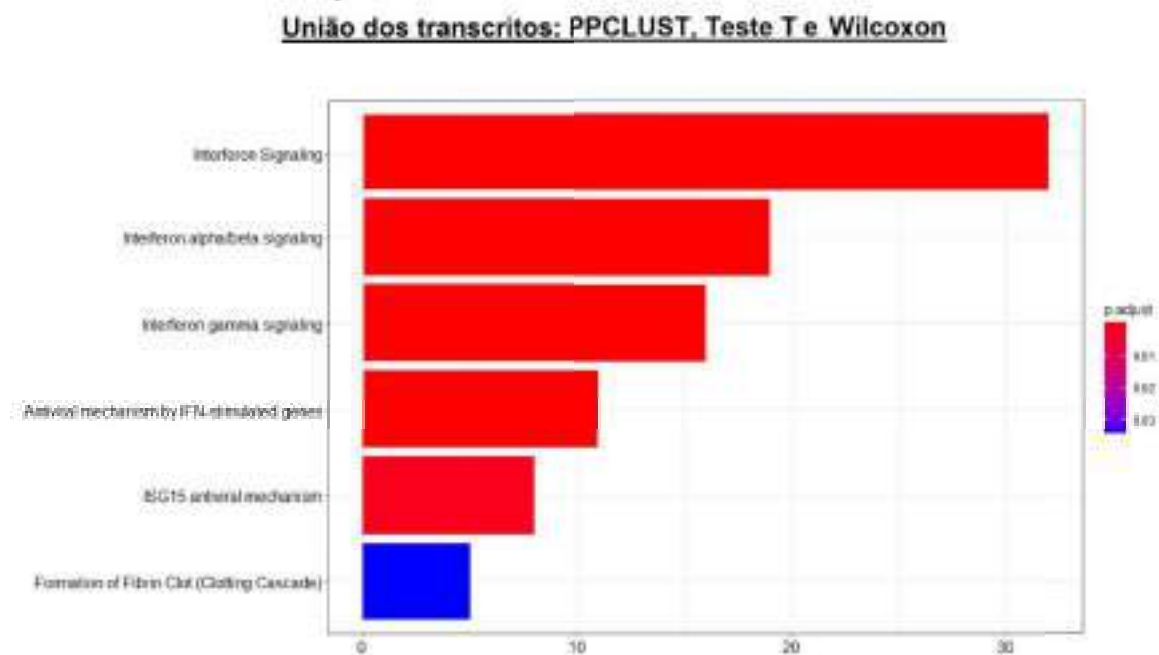


Figura 25 – Distribuição do número de genes da Lista **L5** detectados nos sistemas biológicos mais significantes.

Observa-se que a abordagem baseada no Teste T (**L1**) foi a que identificou o maior número de sistemas biológicos e a lista baseada no PPCLUST (**L3**) a que menos identificou, porém esta última com P-valores mais baixos que nas demais listas, indicando então uma maior consistência para esses sistemas nela identificados.

Todos os gráficos apresentam os sistemas de sinalização de interferon, de interferon alfa/beta e de interferon gama bem como o sistema de mecanismo antiviral por genes induzidos por interferon.

Além disso, nota-se que os sistemas identificados pela lista por Wilcoxon-Mann-Whitney (**L2**) e pela Interseção (**L4**) identificaram exatamente os mesmos sistemas, se diferenciando apenas em seu P-valor ajustado.

Ademais, partes dos sistemas identificados pela Lista **L1** foram visualizadas em todas as demais listas, sendo a lista da União, a Lista **L5**, a única em que foi identificado um sistema que não estava na **L1**, a formação de coágulos.

A descrição dos sistemas biológicos identificados pelos procedimentos das listas L1 a L5 estão na Tabela do Apêndice A.

Gráficos de pontos da razão gênica:

Os resultados fornecidos nas Figuras 21 a 25 podem ser alternativamente apresentados como nas Figuras 26 a 30, que apresentam os sistemas biológicos detectados pela medida razão gênica (*gene ratio*¹). A graduação de cores indica a magnitude do P-valor e o tamanho do círculo, a quantidade de genes identificados como parte de cada um desses sistemas biológicos.

¹ Porcentagem do total de genes diferencialmente expressos.

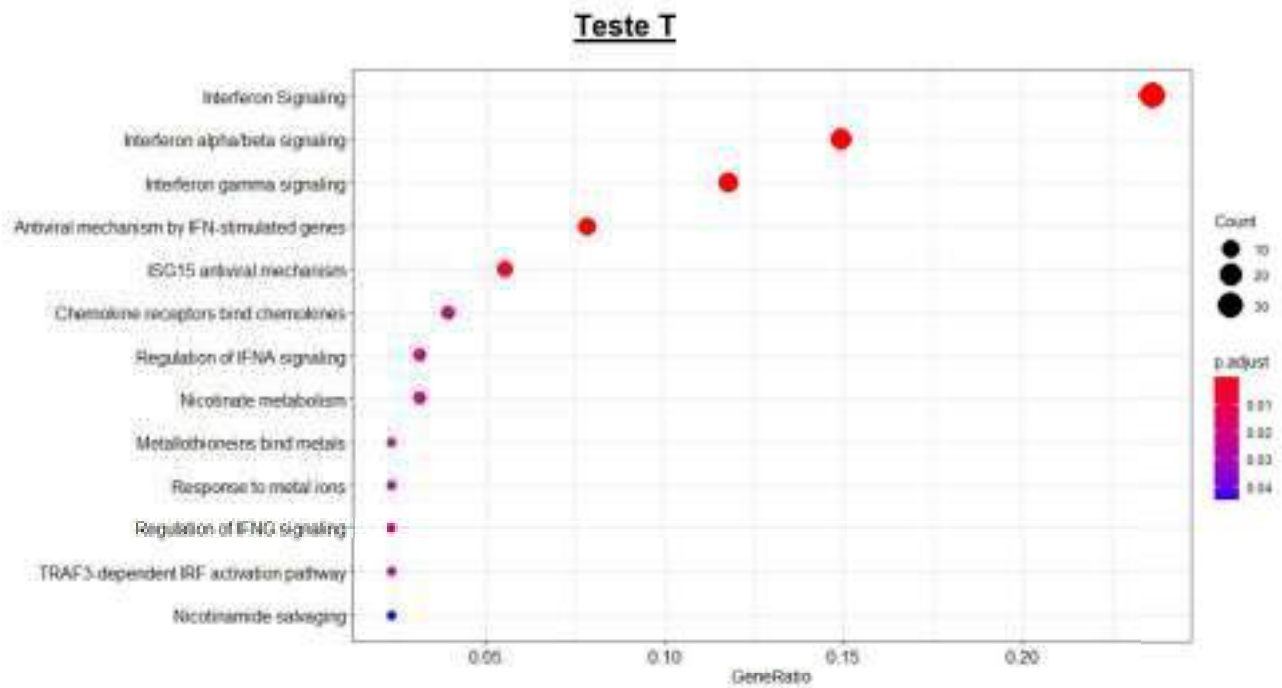


Figura 26 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista **L1**.

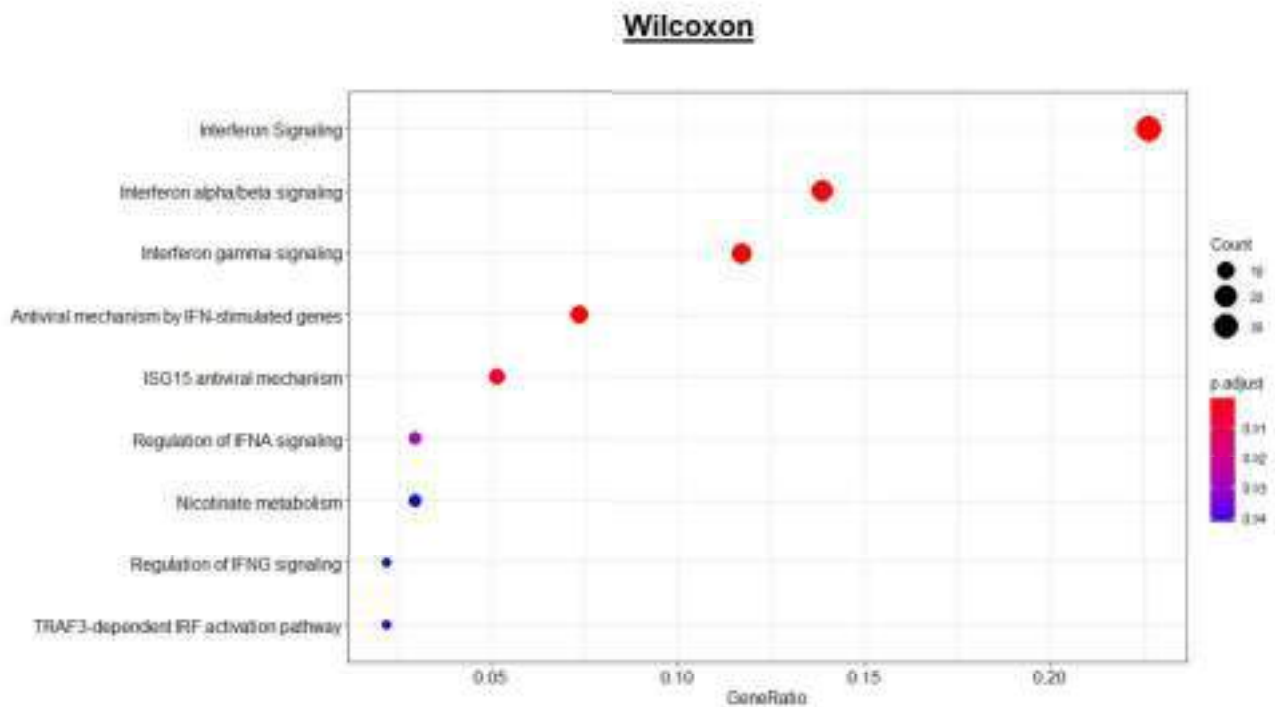


Figura 27 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista **L2**.

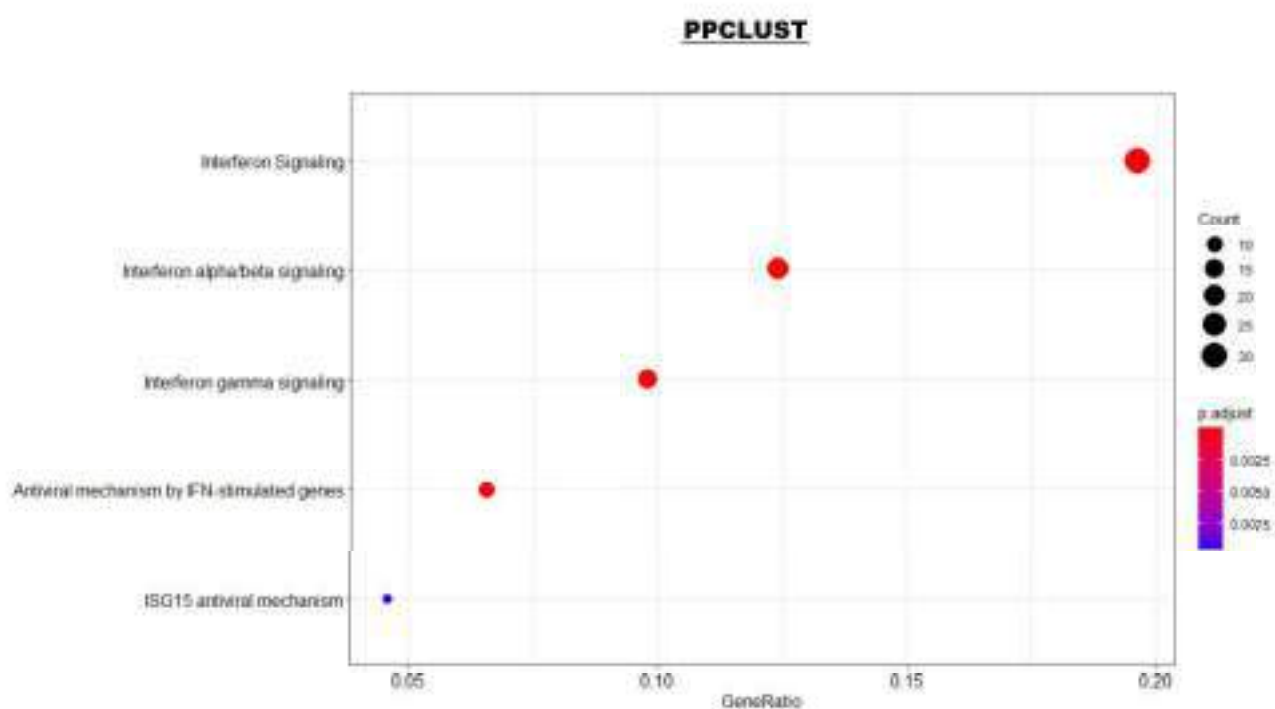


Figura 28 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista **L3**.

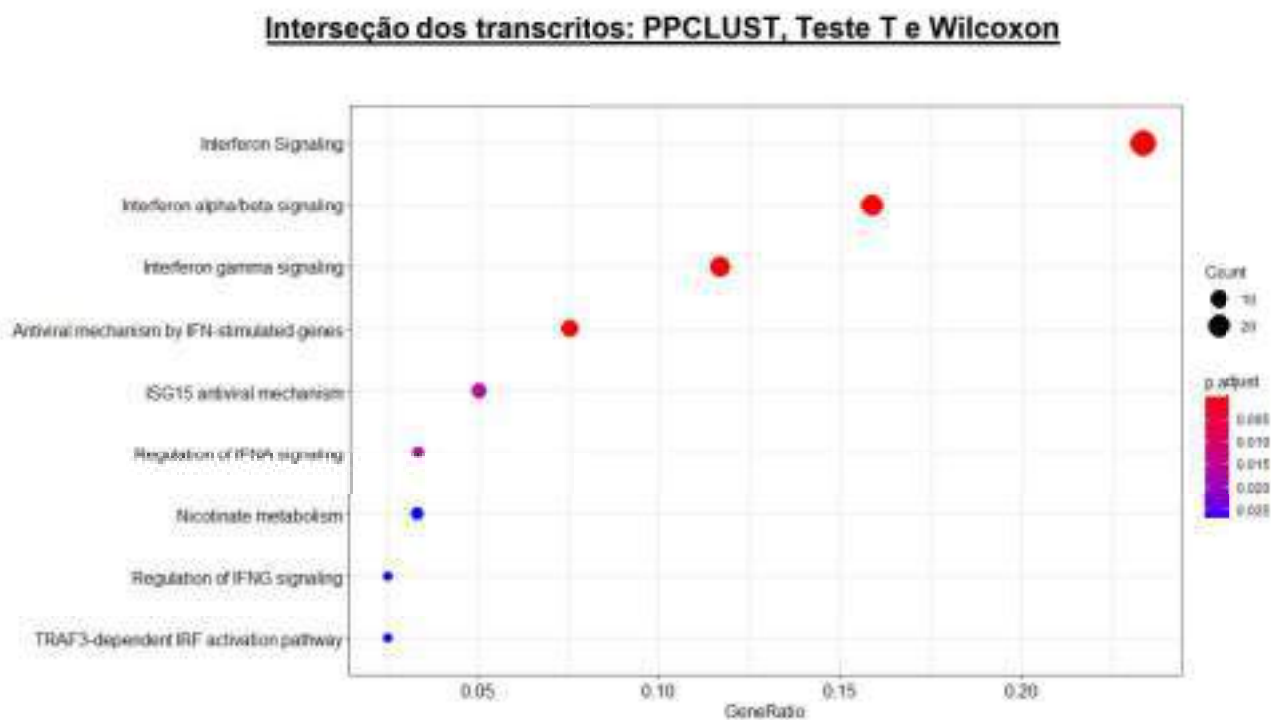


Figura 29 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista **L4**.

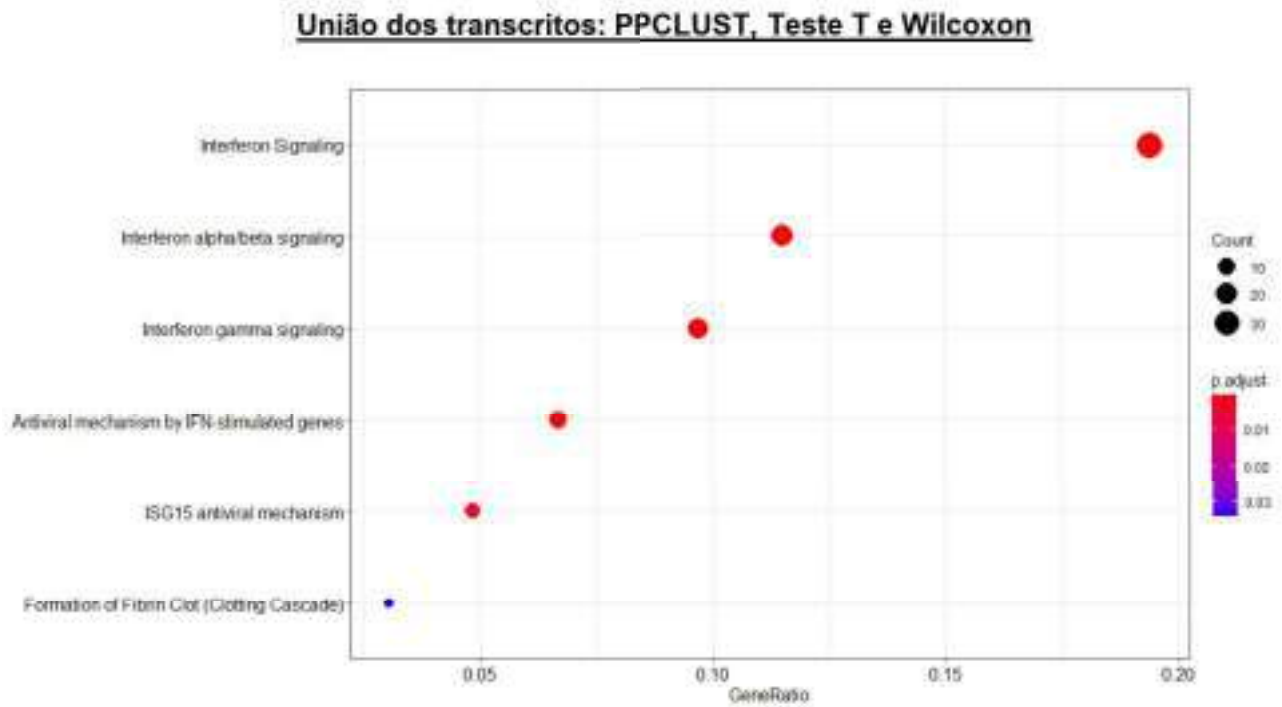


Figura 30 – Distribuição da razão gênica dos sistemas biológicos mais significantes identificados na Lista **L5**.

Diagramas de conexão de sistemas (*emaplots*):

São mapas de anotação que organizam os sistemas biológicos em uma rede com bordas conectando o conjunto desses sistemas quando relacionados. As Figuras 31 a 35 apresentam ligações entre sistemas biológicos mais significantes identificados nas 5 listas de TDEs. A gradação de cores indica a magnitude do P-valor e o tamanho do círculo, a quantidade de genes identificados como parte de cada um desses sistemas biológicos.

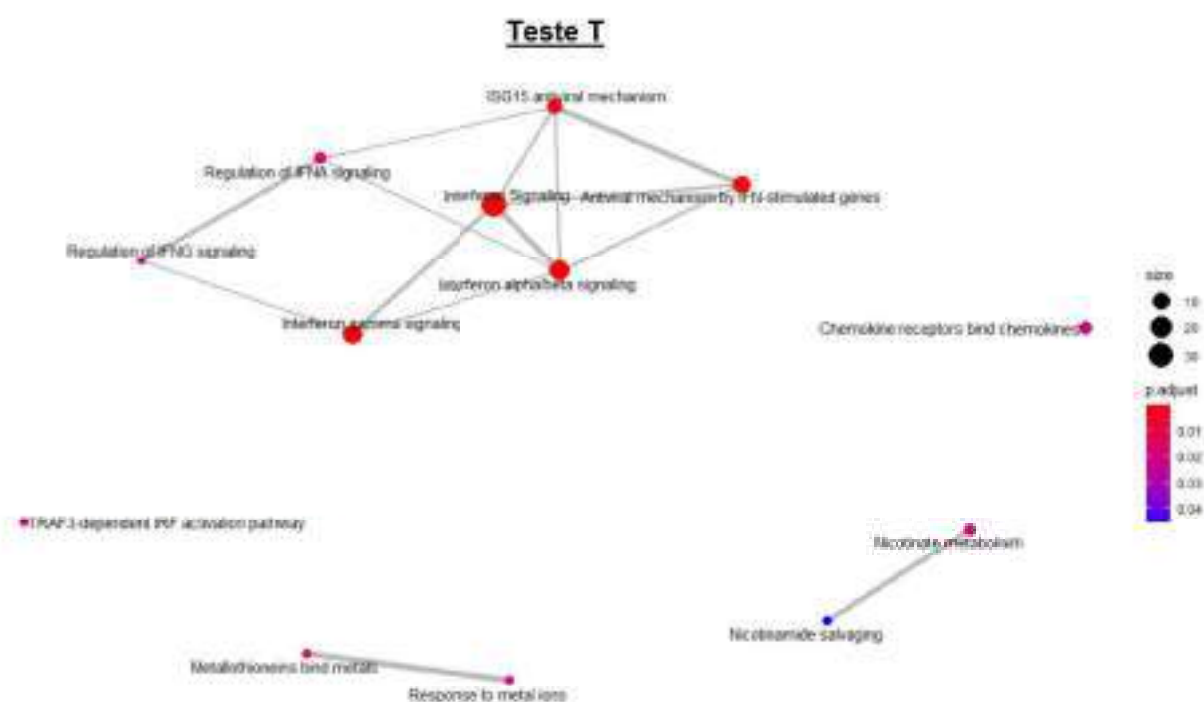


Figura 31 – Diagrama das ligações entre os sistemas biológicos identificados na Lista L1.

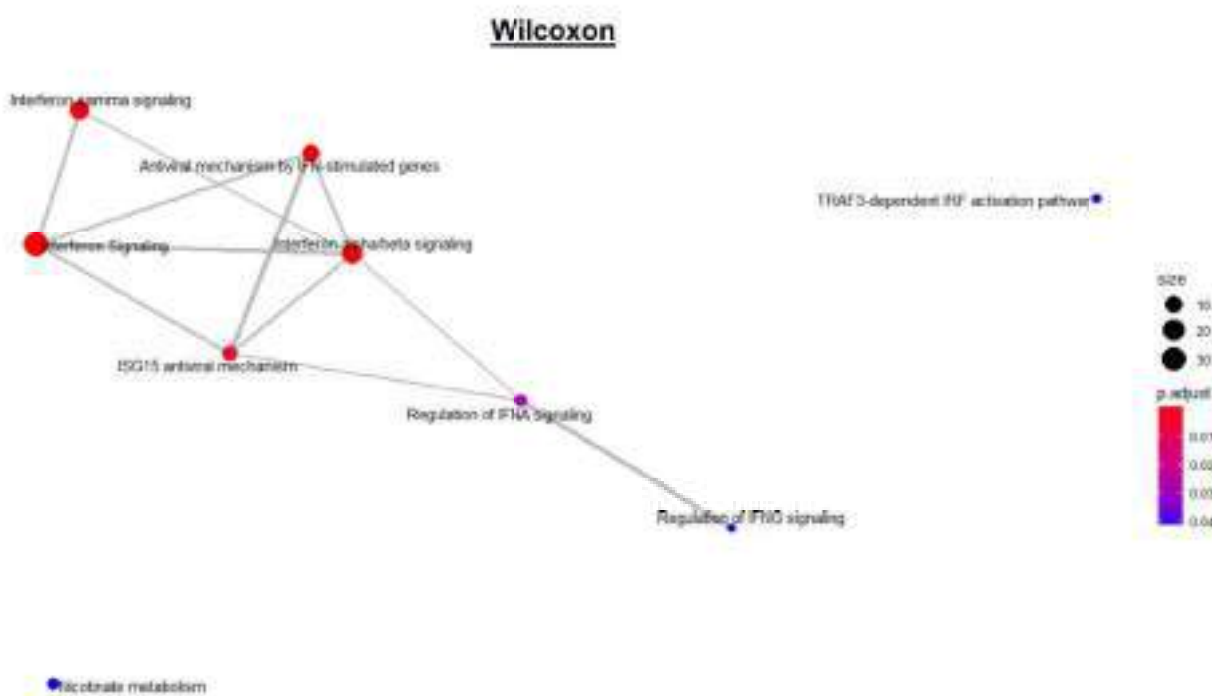


Figura 32 – Diagrama das ligações entre os sistemas biológicos identificados para a Lista L2.

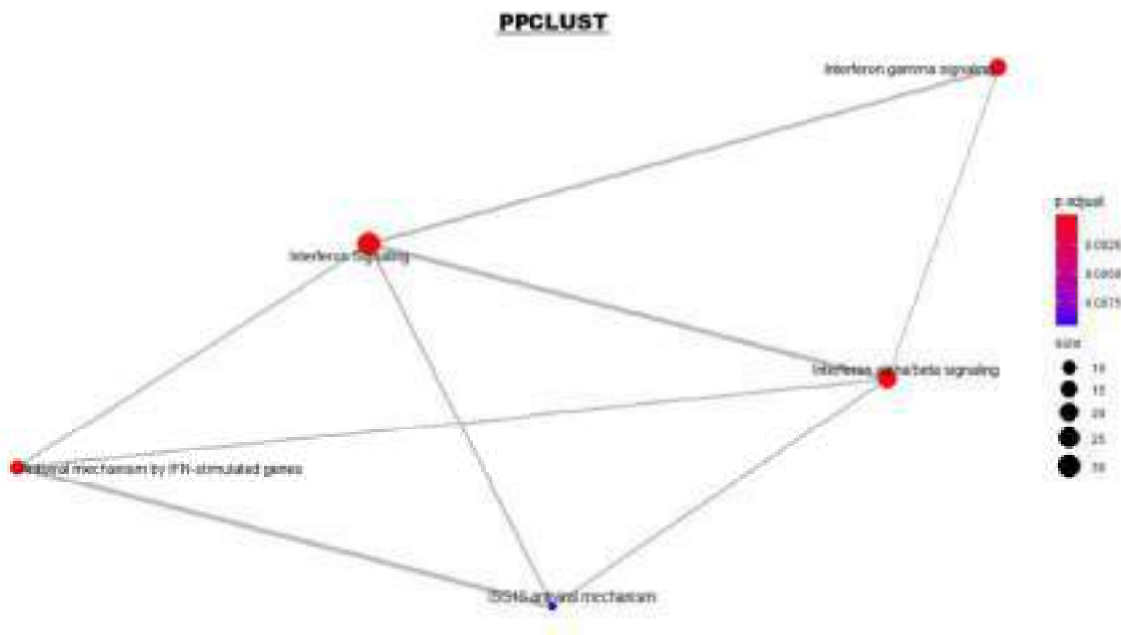


Figura 33 – Diagrama das ligações entre os sistemas biológicos identificados para a Lista L3.

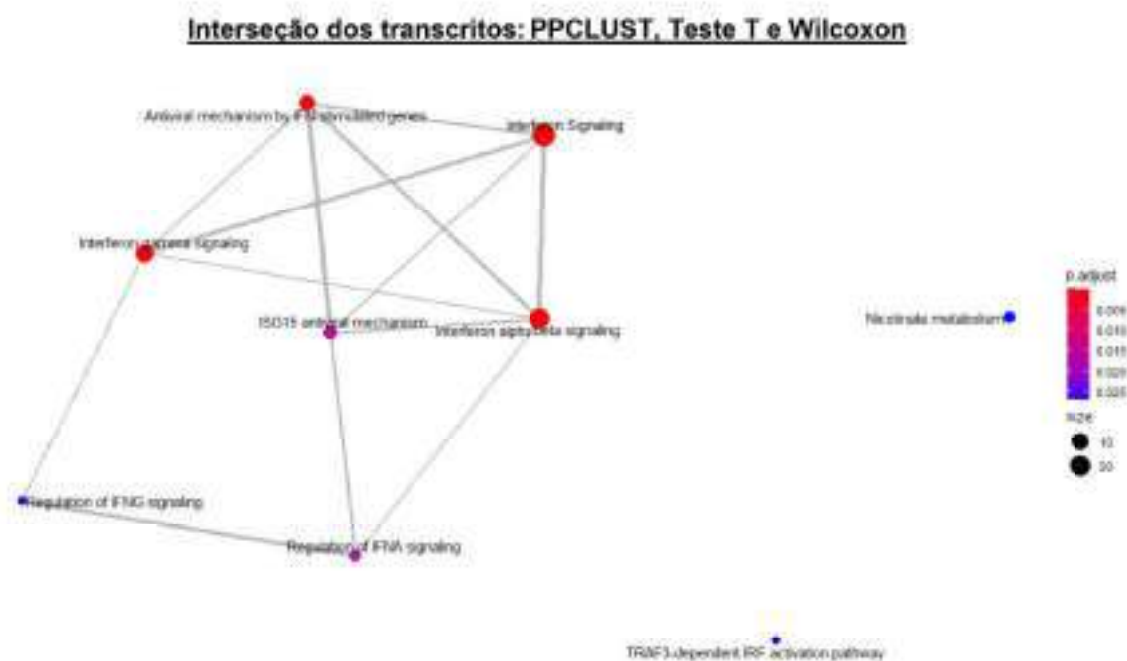


Figura 34 – Diagrama das ligações entre os sistemas biológicos identificados para a Lista L4.

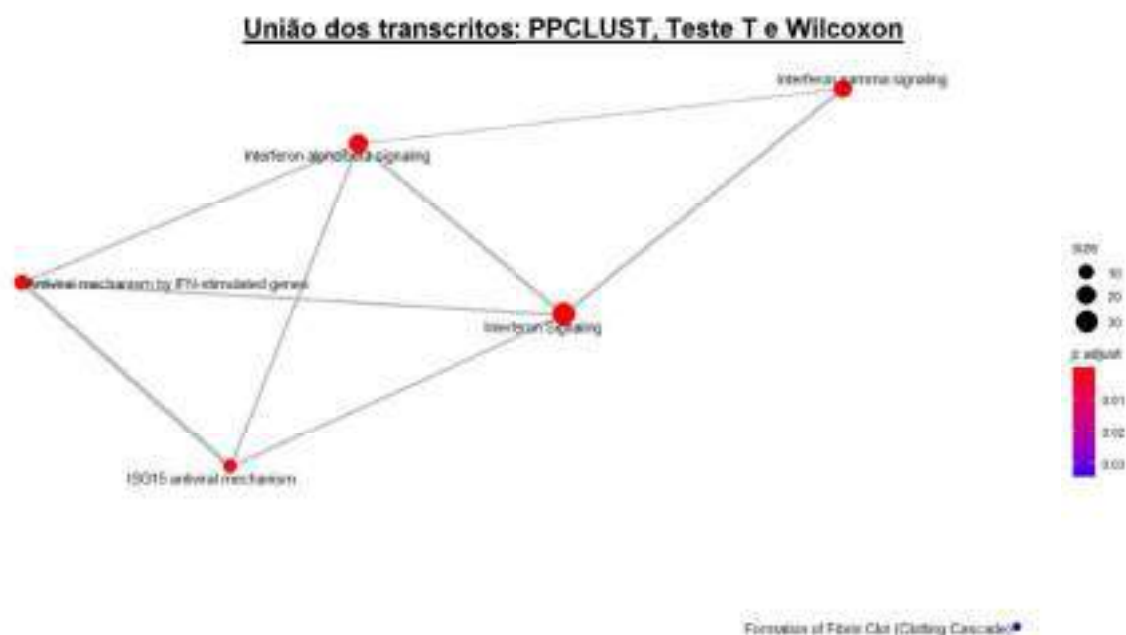


Figura 35 – Diagrama das ligações entre os sistemas biológicos identificados para a Lista L5.

Observa-se que os sistemas detectados para cada uma das listas estão bastante interligados, sendo o diagrama referente à Lista **L1** o que apresentou menos conexões, provavelmente porque apresentou mais sistemas. Sendo o referente a Lista **L3** o único que apresentou todos esses sistemas interligados.

Gráficos de rede, ou *cnetplots*:

Nas Figuras 36 a 40 são apresentadas ligações entre sistemas e seus respectivos genes dentre os que pertencem a uma das listas de transcritos.

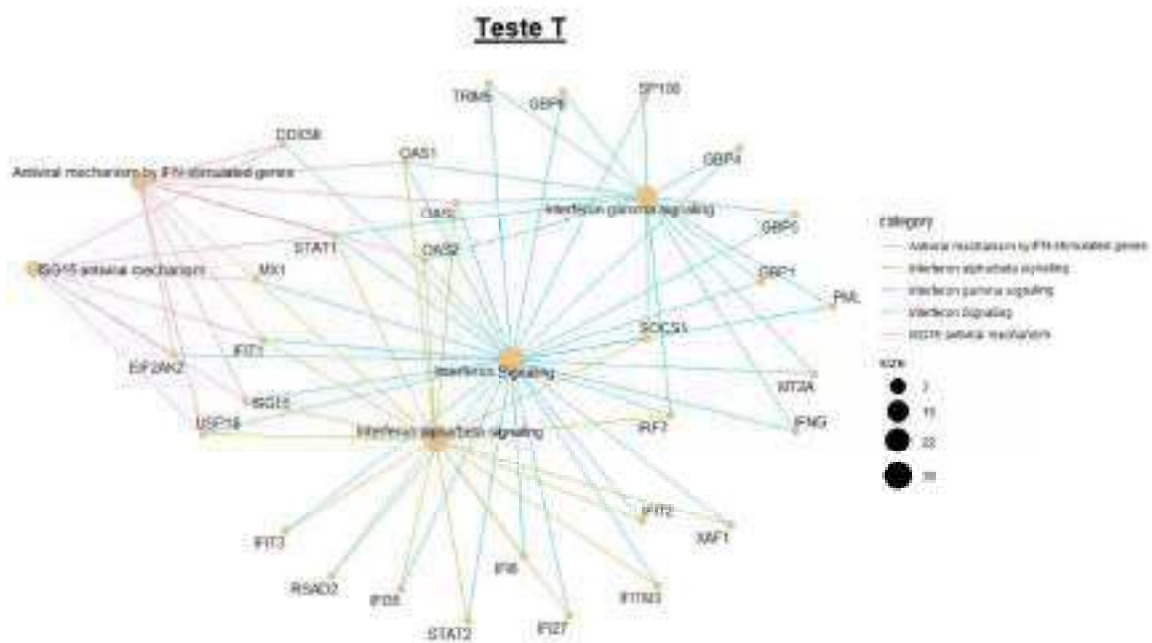


Figura 36 – Gráfico de redes para os genes e seus sistemas biológicos relacionados identificados na Lista L1.

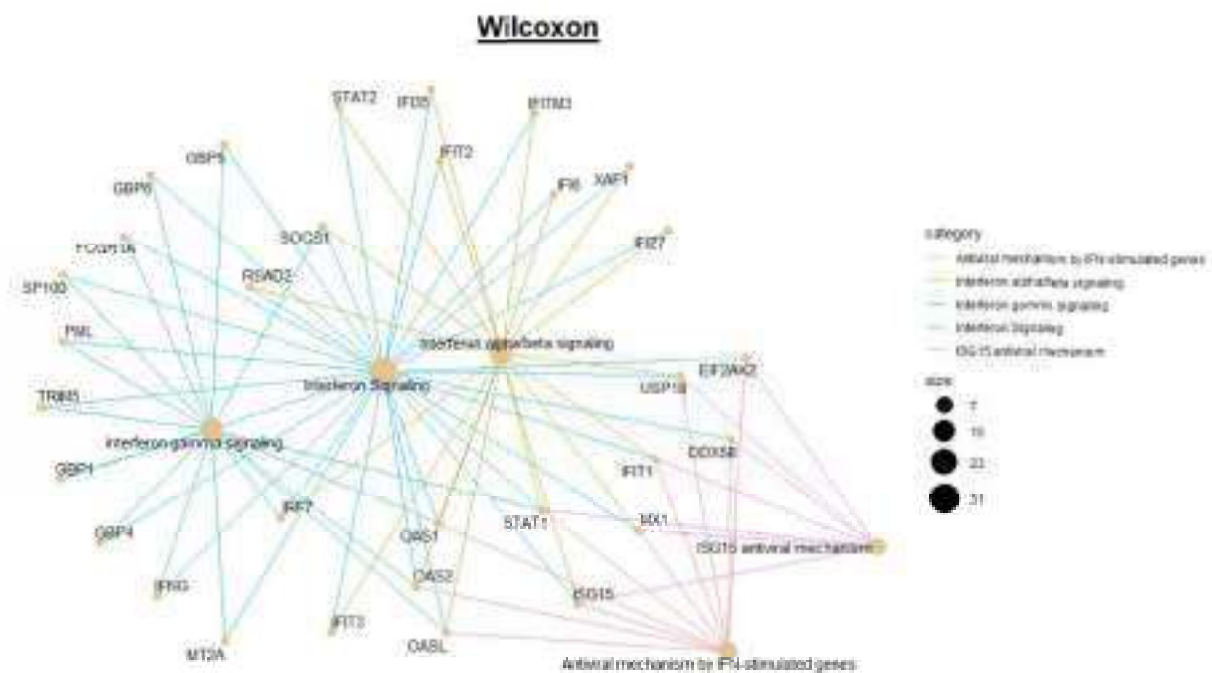


Figura 37 – Gráfico de redes para os genes e seus sistemas biológicos relacionados identificados para a Lista L2.

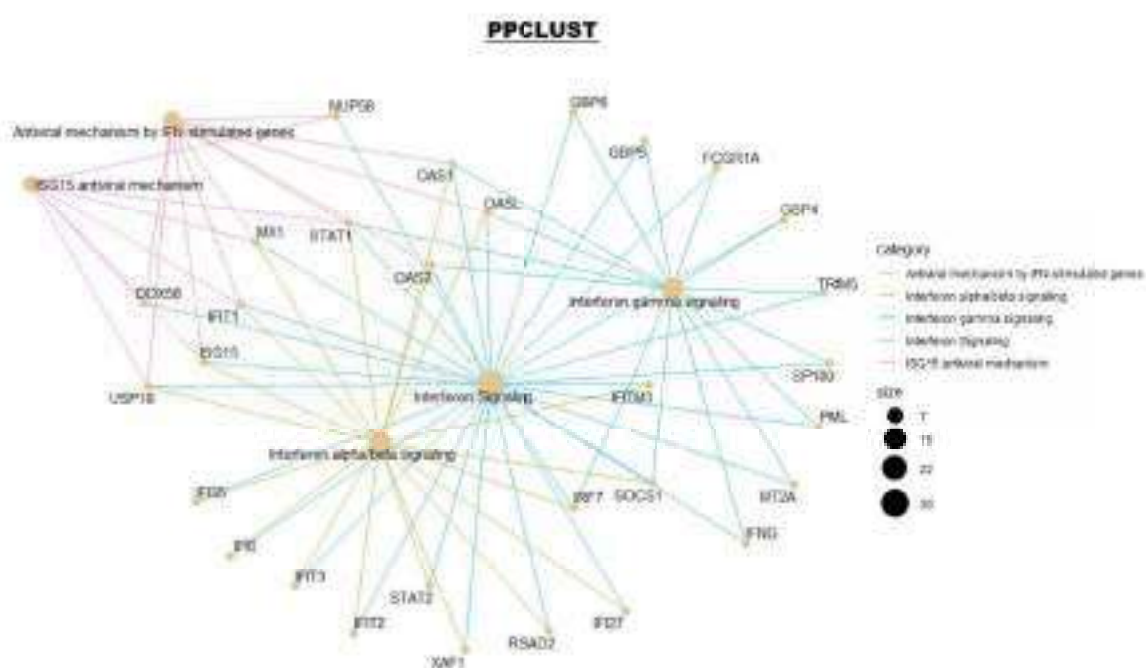


Figura 38 – Gráfico de redes para os genes e seus sistemas biológicos relacionados identificados para a Lista L3.

Interseção dos transcritos: PPCLUST, Teste T e Wilcoxon

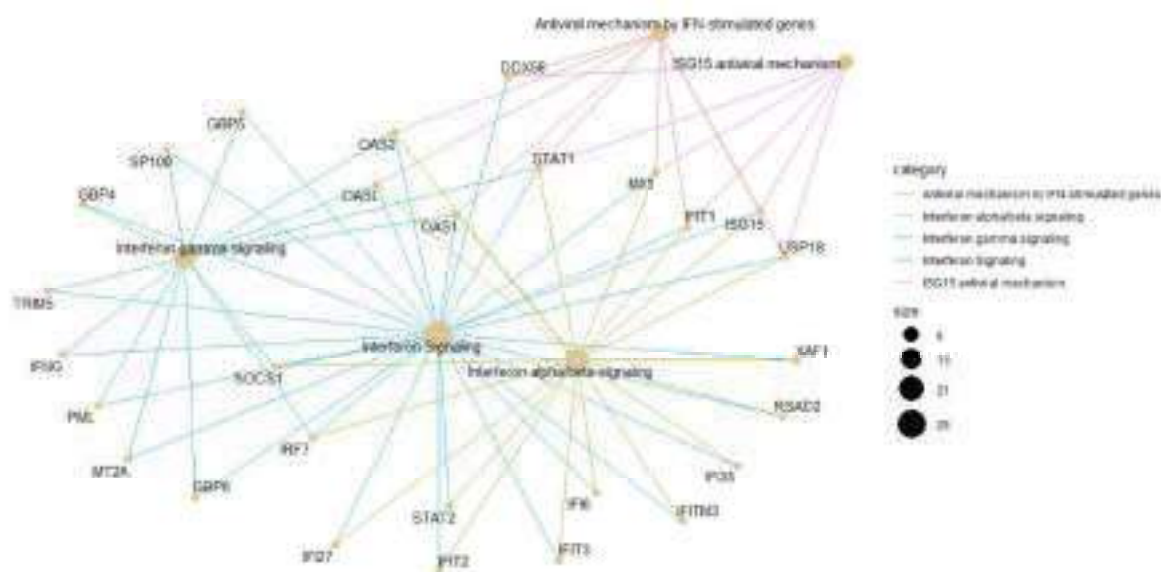


Figura 39 – Gráfico de redes para os genes e seus sistemas biológicos relacionados identificados para a Lista L4.

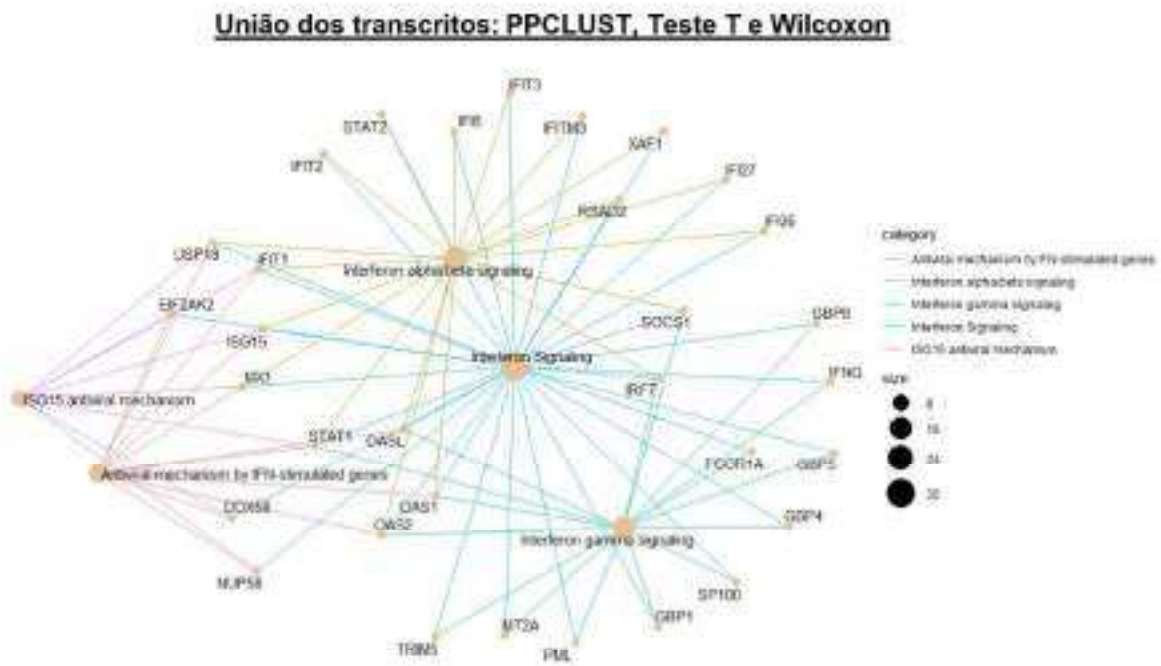


Figura 40 – Gráfico de redes para os genes e seus sistemas biológicos relacionados identificados para a Lista L5.

Observa-se que muitos genes pertencem a mais de um sistema biológico, sendo possível observar a complexidade nas relações entre genes e sistemas biológicos em todas as 5 listas.

4 Discussão e Conclusões

O presente trabalho envolveu uma análise de perfil de expressão gênica, comparando-se dados de casos de pSS e de controles saudáveis. Três abordagens distintas foram utilizadas para a identificação de transcritos diferencialmente expressos (TDEs), quais sejam, 1. Por p-valores de testes t, corrigidos por FDR $< 0,05$ e *Fold Change* (FC) de médias; 2. Por p-valores de testes Wilcoxon-Mann-Whitney, corrigidos por FDR e FC de medianas; 3. Por resultados do algoritmo PPCLUST e FC de medianas. Além das listas de TDEs obtidas por tais abordagens (**L1**, **L2** e **L3**, respectivamente), utilizou-se a lista composta pela intersecção (**L4**) e outra pela união (**L5**) de **L1**, **L2** e **L3**.

Observou-se que aproximadamente 50% dos TDEs estavam presentes em todas as técnicas. Em todas as listas de TDEs, foram identificados genes associados a SS também observados em outros estudos (em (Imgenberg-Kreuz e outros, 2019) e (Yao e outros, 2019)). Cinco dos 10 genes identificados como diferencialmente expressos estavam também presentes na lista dos 20 mais associados no estudo de YAO e outros (2019) e os 10 genes dados como mais diferencialmente expressos (Tabela 3) são também observados na Figura 1 (Imgenberg-Kreuz e outros, 2019), resultante de análises de estudos de genes candidatos.

Destaca-se entre os resultados obtidos, um número considerável de genes de interferons ou induzidos por interferons super-expressos. Interferons são citocinas que desempenham um papel central nas respostas imunes iniciais, com funções antivirais e antitumorais (Reactome, 2020) (Figura 41).

Entre os genes sub-expressos, foi identificado um gene relacionado a coagulação. Indivíduos com SS podem apresentar problemas de coagulação (Oimim,2020) (Uniprot, 2020) (Genecards, 2020).

Com respeito às análise de anotação, houve bastante concordância entre os sistemas identificados nas Listas **L1** a **L5**, sendo **L1** a que identificou mais sistemas, seguida da **L2** e então **L3**, com apenas 5 sistemas, confirmados pelos resultados obtidos pelas outras duas técnicas. Os sistemas biológicos identificados já foram reportados em outros estudos genéticos de SS e algumas outras doenças autoimunes (Imgenberg-Kreuz e outros, 2019).

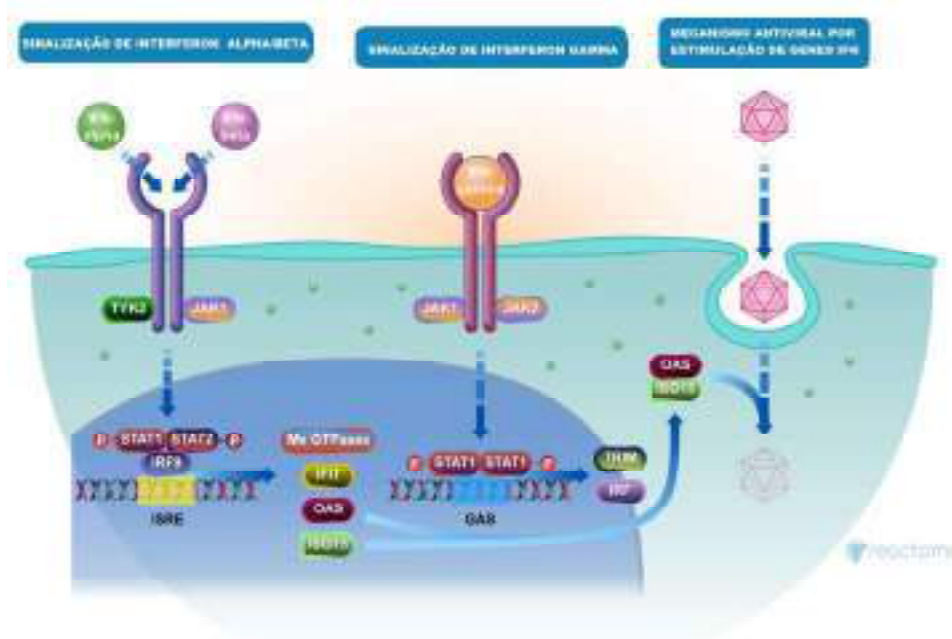


Figura 41 – Ilustração da ação de interferons. (Reactome, 2020)

Apesar de haver bastante concordância entre as listas (**L1**, **L2** e **L3**) e os sistemas biológicos identificados pelas 3 abordagens, tal comparação evidencia que a escolha do método de filtragem pode afetar substancialmente os resultados. Portanto, uma cuidadosa análise exploratória dos dados se faz necessária. Além disso, apesar de não ter sido tratado neste trabalho, procedimentos de controle de qualidade são essenciais e diferenças entre técnicas de normalização devem ser observadas.

Outras técnicas poderiam ter sido utilizadas buscando encontrar TDEs, como é o caso de uma Análise de Discriminante em que um gene teórico, super-expresso em todas as amostras do grupo 1 e sub-expresso em todas as amostras do grupo 2, seleciona genes que possuam a correlação de Pearson mais alta com um discriminador ideal.

Outra opção seria o método paramétrico de Bayes Empíricos para TDEs (**EBarrays**) que tem como base o cálculo de probabilidades a posteriori de padrões de expressão diferencial em várias condições.

O presente estudo permitiu uma análise global das expressões gênicas em SS. Mais estudos são necessários para um melhor entendimento dos processos moleculares que levam ao desenvolvimento e progressão da síndrome e a identificação de biomarcadores moleculares. Em particular, há o interesse em se identificar subtipos de pacientes e avaliar o perfil de expressão gênica por complicações clínicas, severidade da síndrome e resposta à medicamentos.

Referências

- ABDI, H. Bonferroni and šidák corrections for multiple comparisons. *Encyclopedia of measurement and statistics*, Sage Thousand Oaks, CA, v. 3, p. 103–107, 2007. Citado na página 38.
- AGRESTI, A. *An introduction to categorical data analysis third edition*. [S.l.]: Wiley New York, 2019. Citado 2 vezes nas páginas 29 e 35.
- AGUIAR, P. M. d. C.; SEVERINO, P. Biomarkers in parkinson disease: global gene expression analysis in peripheral blood from patients with and without mutations in park2 and park8. *Einstein (Sao Paulo)*, SciELO Brasil, v. 8, n. 3, p. 291–297, 2010. Citado na página 38.
- ARIVARADARAJAN, P.; MISRA, G. *Omics approaches, technologies and applications: integrative approaches for understanding OMICS data*. [S.l.]: Springer, 2019. Citado na página 28.
- BENJAMINI, Y.; HOCHBERG, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, Wiley Online Library, v. 57, n. 1, p. 289–300, 1995. Citado na página 34.
- BIOCOMPUTING, P. S. on. *Pacific Symposium on Biocomputing 2007: Maui, Hawaii, 3-7 January 2007*. [S.l.]: World Scientific Publishing Company, 2007. Citado na página 39.
- BIOCONDUTOR. *bioconductor.org*. 2020. <<https://bioconductor.org/>>. [Online; accessed 22-June-2020]. Citado na página 81.
- BLAND, J. M.; ALTMAN, D. G. Multiple significance tests: the bonferroni method. *Bmj*, British Medical Journal Publishing Group, v. 310, n. 6973, p. 170, 1995. Citado na página 38.
- BRANDT, J. E. et al. Sex differences in sjögren’s syndrome: a comprehensive review of immune mechanisms. *Biology of sex differences*, Springer, v. 6, n. 1, p. 19, 2015. Citado na página 19.
- BRAVIM, F. Análise de expressão gênica. In: . [S.l.]: Universidade Federal do Espírito Santo, Laboratório de Biotecnologia Aplicado ao Agronegócio, Brasil, 2013. Citado na página 28.
- CARSONS, S. E.; PATEL, B. C. Sjogren syndrome. In: *StatPearls [Internet]*. [S.l.]: StatPearls Publishing, 2019. Citado na página 19.
- EDGAR, R.; DOMRACHEV, M.; LASH, A. E. Gene expression omnibus: Ncbi gene expression and hybridization array data repository. *Nucleic acids research*, Oxford University Press, v. 30, n. 1, p. 207–210, 2002. Citado na página 25.

- FELBERG, S.; DANTAS, P. E. C. Diagnóstico e tratamento da síndrome de sjögren. *Arquivos brasileiros de oftalmologia*, Conselho Brasileiro de Oftalmologia, 2006. Citado na página 20.
- GARCIA, A. A. F. Biometria de marcadores genéticos - tópico 3: Mapas genéticos i segregação mendeliana. In: . [S.l.]: Departamento de Genética - ESALQ/USP, 2019. Citado na página 34.
- GENECARDS. *genecards.org*. 2020. <<https://genecards.org/>>. [Online; accessed 28-June-2020]. Citado 2 vezes nas páginas 51 e 65.
- GRIFFITHS, A. J. et al. *An introduction to genetic analysis*. [S.l.]: Macmillan, 2005. Citado na página 27.
- ILLUMINA, I. *HumanHT-12 v3 Expression BeadChip*. 2008. <https://www.illumina.com/Documents/products/datasheets/datasheet_humanht_12.pdf>. [Online; accessed 18-July-2020]. Citado na página 28.
- IMGENBERG-KREUZ, J. et al. Genetics and epigenetics in primary sjögren's syndrome. *Rheumatology*, 2019. Citado 4 vezes nas páginas 20, 21, 22 e 65.
- KASSAN, S. S.; MOUTSOPOULOS, H. M. Clinical manifestations and early diagnosis of sjögren syndrome. *Archives of internal medicine*, American Medical Association, v. 164, n. 12, p. 1275–1284, 2004. Citado na página 19.
- KAUFMAN, L.; ROUSSEEUW, P. J. *Finding groups in data: an introduction to cluster analysis*. [S.l.]: John Wiley & Sons, 2009. v. 344. Citado 2 vezes nas páginas 35 e 36.
- LESSARD, C. J. et al. Variants at multiple loci implicated in both innate and adaptive immune responses are associated with sjögren's syndrome. *Nature genetics*, Nature Publishing Group, v. 45, n. 11, p. 1284–1292, 2013. Citado 4 vezes nas páginas 25, 27, 28 e 41.
- LI, H. et al. Interferons in sjögren's syndrome: genes, mechanisms, and effects. *Frontiers in immunology*, Frontiers, v. 4, p. 290, 2013. Citado na página 21.
- LINDEN, R. Técnicas de agrupamento. *Revista de Sistemas de Informação da FSMA*, n. 4, v. 4, n. 4, p. 18–36, 2009. Citado na página 35.
- LINS, R. d. S. Implementação computacional de algoritmos para agrupamento de dados hdlls e hdlss. 2018. Citado na página 33.
- LOPES, M. P. et al. Estrutura e variabilidade do promotor do gene do fator de necrose tumoral humano (tnf). Universidade Federal de Goiás, 2014. Citado na página 23.
- MIOZZI, L. et al. Functional annotation and identification of candidate disease genes by computational analysis of normal tissue gene expression data. *PLoS One*, Public Library of Science, v. 3, n. 6, 2008. Citado na página 38.
- MIRANDA, P. A. V. de. Métodos de agrupamento (clustering) - aula 18. In: . [S.l.]: Universidade de São Paulo (USP), 2012. Citado na página 36.
- MORETTIN, P. A.; BUSSAB, W. O. *Estatística básica*. [S.l.]: Saraiva Educação SA, 2017. Citado na página 31.

- NIKOLOV, N. P.; ILLEI, G. G. Pathogenesis of sjögren's syndrome. *Current opinion in rheumatology*, NIH Public Access, v. 21, n. 5, p. 465, 2009. Citado na página 20.
- OMIM. *omim.org*. 2020. <<https://omim.org/>>. [Online; accessed 28-June-2020]. Citado 2 vezes nas páginas 51 e 65.
- PASOTO, S. G.; MARTINS, V. A. de O.; BONFA, E. Sjögren's syndrome and systemic lupus erythematosus: links and risks. *Open access rheumatology: research and reviews*, Dove Press, v. 11, p. 33, 2019. Citado na página 19.
- PATEL, R.; SHAHANE, A. The epidemiology of sjögren's syndrome. *Clinical epidemiology*, Dove Press, v. 6, p. 247, 2014. Citado na página 19.
- QIN, B. et al. Epidemiology of primary sjögren's syndrome: a systematic review and meta-analysis. *Annals of the rheumatic diseases*, BMJ Publishing Group Ltd, v. 74, n. 11, p. 1983–1989, 2015. Citado na página 19.
- RAMOS-CASALS, M. et al. Primary sjögren syndrome. *British Medical Journal Publishing Group*, v. 344, p. e3821, 2012. Citado na página 19.
- REACTOME. *reactome.org*. 2020. <<https://reactome.org/>>. [Online; accessed 28-June-2020]. Citado 8 vezes nas páginas 65, 66, 73, 76, 77, 78, 79 e 80.
- REACTOME.ORG. *reactome.org*. 2020. <<https://reactome.org/>>. [Online; accessed 22-June-2020]. Citado na página 38.
- SHAPIRO, S. S.; WILK, M. B. An analysis of variance test for normality (complete samples). *Biometrika*, JSTOR, v. 52, n. 3/4, p. 591–611, 1965. Citado na página 30.
- STEIN, L. Genome annotation: from sequence to biology. *Nature reviews genetics*, Nature Publishing Group, v. 2, n. 7, p. 493–503, 2001. Citado na página 37.
- TAPINOS, N. I. et al. *Sjögren's Syndrome*. Boston, MA: Springer US, 1999. 127–134 p. Citado na página 19.
- THEANDER, E. et al. Prediction of sjögren's syndrome years before diagnosis and identification of patients with early onset and severe disease course by autoantibody profiling. *Arthritis & rheumatology*, Wiley Online Library, v. 67, n. 9, p. 2427–2436, 2015. Citado na página 20.
- UNIPROT. *uniprot.org*. 2020. <<https://uniprot.org/>>. [Online; accessed 28-June-2020]. Citado 2 vezes nas páginas 51 e 65.
- VALTYSDÓTTIR, S. T.; WIDE, L.; HÄLLGREN, R. Low serum dehydroepiandrosterone sulfate in women with primary sjögren's syndrome as an isolated sign of impaired hpa axis function. *The Journal of rheumatology*, The Journal of Rheumatology, v. 28, n. 6, p. 1259–1265, 2001. Citado na página 19.
- VASTRIK, I. et al. Reactome: a knowledge base of biologic pathways and processes. *Genome biology*, BioMed Central, v. 8, n. 3, p. 1–13, 2007. Citado na página 38.
- VITALI, C. et al. Classification criteria for sjögren's syndrome: a revised version of the european criteria proposed by the american-european consensus group. *Annals of the rheumatic diseases*, BMJ Publishing Group Ltd, v. 61, n. 6, p. 554–558, 2002. Citado na página 25.

VON BORRIES, G.; WANG, H. Partition clustering of high dimensional low sample size data based on p-values. *Computational statistics & data analysis*, Elsevier, v. 53, n. 12, p. 3987–3998, 2009. Citado na página 32.

WANG, H.; AKRITAS, M. G. Rank tests for anova with large number of factor levels. *Journal of Nonparametric Statistics*, Taylor & Francis, v. 16, n. 3-4, p. 563–589, 2004. Citado na página 33.

YAO, Q. et al. Identifying key genes and functionally enriched pathways in sjögren's syndrome by weighted gene co-expression network analysis. *Frontiers in genetics*, Frontiers, v. 10, p. 1142, 2019. Citado 2 vezes nas páginas 20 e 65.

Apêndices

APÊNDICE A – Descrição dos Sistemas Biológicos Identificados em Análises de Anotação

Descrição dos Sistemas Biológicos Identificados em Análises de Anotação - Parte 1. (Reactome, 2020)

- ***Interferon Signaling* | Sinal de Interferon:** Citocinas que desempenham um papel central no início de respostas imunes, especialmente efeitos antivirais e antitumorais;
- ***Interferon alpha/beta signaling* | Sinal de Interferon alfa/beta:** Tipo de interferon, citocinas, que desempenham um papel central no início de respostas imunes, especialmente efeitos antivirais e antitumorais;
- ***Interferon gamma signaling* | Sinal de Interferon gama:** Tipo de interferon, citocinas, que desempenham um papel central no início de respostas imunes, especialmente efeitos antivirais e antitumorais;
- ***Antiviral mechanism by IFN-stimulated genes* | Mecanismo antiviral por genes estimulados por Interferons:** Os interferons ativam a sinalização que leva à indução transcricional de centenas de genes estimulados por interferon. As proteínas codificadas por inteferon incluem efetores diretos que inibem a infecção viral através de diversos mecanismos, bem como fatores que promovem respostas imunes adaptativas. As proteínas geradas pelas vias interferon desempenham papéis-chave na indução de respostas imunes inatas e adaptativas;
- ***ISG15 antiviral mechanism* | Mecanismo antiviral ISG15:** É fortemente induzido após a exposição a interferons tipo I, vírus, lipopolissacarídeo bacteriano e outros estresses. Uma vez liberado, o ISG15 maduro se conjuga com uma série de proteínas alvo, um processo denominado ISGylation;
- ***Regulation of IFNA signaling* | Regulação da sinalização IFNA:** Várias proteínas e mecanismos envolvidos no controle da extensão da estimulação do ligante da sinalização de interferons alpha e beta;
- ***Regulation of IFNG signaling* | Regulação da sinalização IFNG:** Várias proteínas e mecanismos envolvidos no controle da extensão da estimulação do ligante da sinalização de interferons alpha e beta;

- ***Nicotinamide salvaging* | Recuperação de nicotinamida:** Modulador da ação de três enzimas envolvidas na recuperação da nicotinamida;
- ***Nicotinate metabolism* | Metabolismo de nicotinato:** Modulador da ação de três enzimas envolvidas na recuperação da nicotinamida;
- ***Nicotinate metabolism* | Metabolismo de nicotinato:** São cofatores importantes em várias centenas de reações redox (oxidação-redução);
- ***Metallothioneins bind metals* | Metalotioneínas ligam metais:** Embora as funções das metalotioneínas não tenham sido totalmente elucidadas, elas parecem participar na desintoxicação de metais pesados, no armazenamento e transporte de zinco e na bioquímica redox;
- ***Response to metal ions* | Resposta a íons metálicos:** Embora metais como zinco, cobre e ferro sejam necessários como cofatores para enzimas celulares, eles também podem catalisar substituições de metais prejudiciais ou reações redox inespecíficas, se não forem retirados;
- ***Chemokine receptors bind chemokines* | Receptores de quimiocinas ligam quimiocinas:** As quimiocinas podem ser divididas em dois grupos funcionais, inflamatório e homeostático, mas a discriminação não é estrita e algumas sobreposições são encontradas. As quimiocinas inflamatórias são produzidas em condições inflamatórias por infiltração e células residentes em resposta a mediadores pró-inflamatórios (IL-1 e TNF- α), produtos bacterianos (LPS) e agentes infecciosos (vírus);
- ***Formation of Fibrin Clot* | Formação de coágulos de fibrina:** Ligados a cicatrização, formação de coágulos, para a cicatrização.

APÊNDICE B – Lista de transcritos estatisticamente significantes

Tabela 4 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE 1. (Reactome, 2020)

ID	LÔC FC	GENE SYMBOL	GENE TITLE
TLNW_1713393	0.97755056*ZNF604	ZNF604	Zinc finger protein 604
TLNW_1714402	1.42041920*ZC3H3C2	ZC3H3C2	Zinc finger CCHC-type-containing 2
TLNW_2492017	-0.68589465*ZBTB10	ZBTB10	Zinc finger and BTB domain containing 10
TLNW_1675596	-0.65022736*ZKRB1	ZKRB1	ZK-related 1
TLNW_1742618	1.48077070*ZNF1	ZNF1	ZNF-associated factor 1
TLNW_1671200	-0.72693660*ZNF5	ZNF5	Zinc finger WW domain-containing protein 5
TLNW_2793095	-0.62613443*ZNF157	ZNF157	Zinc finger WW domain-containing protein 157
TLNW_1807375	0.02079869*ZNF501	ZNF501	Zinc finger WW domain-containing protein 501
TLNW_1757755	-0.65748030*ZNF1	ZNF1	Zinc finger WW domain-containing protein 1
TLNW_1740000	0.00027050*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2095012	0.03271312*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1003040	0.72610013*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2838775	2.14912125*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2044013	1.04230027*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1692738	-0.71873160*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1713003	-1.04405085*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1715002	1.04941583*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1803000	0.00013400*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1737595	0.72591502*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1797250	0.68291002*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1693000	0.70940299*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1693000	1.07357004*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1795008	0.74892000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1693000	-0.68510000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1759075	-0.58680125*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1795750	0.64400000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1770007	-0.69030015*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2195577	-0.64754000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1700100	0.02007000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2240000	0.63000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1700000	0.72000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2350000	0.77231000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1693000	-0.64000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1750707	-0.60070000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1795000	0.75700000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1700000	0.64700000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1750000	0.67070000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2403000	0.60430000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1695000	0.76300000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1690100	0.77240000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1691300	0.66300000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1770000	0.61000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2204000	0.00000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1700000	-0.70000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1691300	-0.71000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2074400	-0.68000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1693000	1.51000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1703000	0.72540000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2240000	0.70000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2000000	0.66000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1774700	0.73000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1775300	0.62500000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1714400	-1.00000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1703000	0.84650000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2084200	0.81400000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1703000	-0.60000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2211700	0.84710000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1703000	0.85000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1770007	0.58000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1670000	1.66070000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1693000	-0.68000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1670000	0.68000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1805000	0.73957000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1700000	-1.00000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2173075	1.10100000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1693000	0.00000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1657075	-0.07000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1775125	-0.73000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1707055	-0.65010000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2084200	-0.77400000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1697000	0.82500000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2017050	0.82000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1800000	0.82100000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1740000	0.82000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1733045	0.62140000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2090024	-0.76000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_2000000	0.80000000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108
TLNW_1810000	0.64157000*ZNF108	ZNF108	Zinc finger WW domain-containing protein 108

Tabela 5 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE 2. (Reactome, 2020)

ID	LOG FC	GENE SYMBOL	GENE TITLE
TLMN_1695840*	0,827894185	PTEN22*	protein tyrosine phosphatase, non-receptor type 22*
TLMN_1798233*	0,588210022	PSMD9*	proteasome subunit beta 9*
TLMN_1810620*	1,288680843	PRUNE2*	prune homolog 2*
TLMN_1600017*	-0,609510043	PRRS-ARHGAP8*	PRRS-ARHGAP8 readthrough*
TLMN_1810600*	0,813404084	PNP11*	polynucleotide nucleotidyltransferase 1*
TLMN_1728019*	0,834910235	PML*	promyelocytic leukemia*
TLMN_1745242*	0,96368423	PLSCR1*	phospholipid scramblase 1*
TLMN_1656857*	-0,581758804	PLAU*	plasminogen activator, urokinase*
TLMN_1704537*	-0,808271130	PHGDH*	phosphoglycerate dehydrogenase*
TLMN_2178873*	0,862194971	PHG3*	polycomb target homolog 3*
TLMN_1704870*	-0,808300540	PGCYP1*	peptidoglycan recognition protein 1*
TLMN_1745522*	-0,862132807	PI-AVI*	platelet factor 4 variant 1*
TLMN_1706580*	-0,847064455	POZK1IP1*	POZK1 interacting protein 1*
TLMN_1684862*	0,503515805	PKC4*	pyruvate dehydrogenase kinase 4*
TLMN_1815057*	1,082168247	PDGFRB*	platelet derived growth factor receptor beta*
TLMN_1731224*	0,807282865	PARP9*	poly(ADP-ribose) polymerase family member 9*
TLMN_2053527*	0,876619457	PARP9*	poly(ADP-ribose) polymerase family member 9*
TLMN_1698731*	1,014181128	PARP14*	poly(ADP-ribose) polymerase family member 14*
TLMN_1718558*	0,867677074	PARP12*	poly(ADP-ribose) polymerase family member 12*
TLMN_1721411*	0,73489482	PARP10*	poly(ADP-ribose) polymerase family member 10*
TLMN_1657373*	-0,803697125	PGH2*	prolyl 3-hydroxylase 2*
TLMN_1603847*	2,293122274	PI3K*	piasterin*
TLMN_1732768*	0,814657798	PIR2AG1*	pi3kcyt receptor family 2 subfamily AG member 1 (gene/protein)?
TLMN_1674811*	1,422484828	OAS1*	2'-5'-oligoadenylate synthetase like*
TLMN_1681721*	1,181873523	OAS1*	2'-5'-oligoadenylate synthetase like*
TLMN_1736729*	1,180695905	OAS2*	2'-5'-oligoadenylate synthetase 2*
TLMN_1656291*	0,955271848	OAS1*	2'-5'-oligoadenylate synthetase 1*
TLMN_1675040*	1,203441803	OAS1*	2'-5'-oligoadenylate synthetase 1*
TLMN_2410826*	1,303529656	OAS1*	2'-5'-oligoadenylate synthetase 1*
TLMN_1806448*	0,635703141	NTNG2*	netrin G2*
TLMN_1772627*	-0,812712746	NSG1*	neuron specific gene family member 1*
TLMN_2411238*	-0,803002708	NRCAM*	neuronal cell adhesion molecule*
TLMN_2048768*	0,858472388	NFE2L3*	nuclear factor, erythroid 2 like 3*
TLMN_1693497*	-0,818143450	NECB2*	N-terminal EF-hand calcium binding protein 2*
TLMN_1727618*	0,813125036	NDUFAB1*	NADH dehydrogenase (ubiquinone) complex assembly factor 6*
TLMN_1788689*	-0,583244172	NBL1*	neuroblastoma 1, DAN family BMP antagonist*
TLMN_1716733*	-1,528175883	MYOM2*	myomesin 2*
TLMN_1716062*	-0,888940427	MYL4*	myosin light chain 4*
TLMN_1705297*	-0,597345444	MYBPH*	myosin binding protein H*
TLMN_1733915*	-0,894882286	MX11*	MAX interactor 1, dimerization protein*
TLMN_1662358*	1,640079431	MX1*	MX dynamin like GTPase 1*
TLMN_1688684*	0,781161796	MT2A*	metallothionein 2A*
TLMN_1691156*	0,631734978	MT1A*	metallothionein 1A*
TLMN_2370836*	0,950822802	MS4AA4*	membrane spanning 4-domains 4A*
TLMN_1776039*	0,897224883	MS4A1*	membrane spanning 4-domains 4A*
TLMN_1775016*	0,825836208	MPZL2*	myelin protein zero like 2*
TLMN_1717046*	0,508040002	MOC3B*	MOC3 kinase activator 3B*
TLMN_2398916*	-1,589206834	MMP28*	matrix metalloproteinase 28*
TLMN_1803673*	-0,862274548	MISP3*	MISP family member 3*
TLMN_1738588*	-0,840271481	MGL1*	monoglyceride lipase*
TLMN_2138589*	0,645268976	MERTK*	MER proto-oncogene, tyrosine kinase*
TLMN_1669321*	-0,897264607	MATK*	megakaryocyte-associated tyrosine kinase*
TLMN_2319000*	-0,891689262	MATK*	megakaryocyte-associated tyrosine kinase*
TLMN_1662336*	0,769064357	MASTL*	microtubule associated serine/threonine kinase like*
TLMN_1684404*	1,500611262	LY8E*	lymphocyte antigen 8 complex, locus E*
TLMN_2228242*	0,813632148	LSM3*	LSM3 homolog, U6 small nuclear RNA and mRNA degradation associated*
TLMN_1773858*	-1,105251381	LRRN3*	leucine rich repeat neuronal 3*
TLMN_2048581*	-0,921944621	LRRN3*	leucine rich repeat neuronal 3*
TLMN_1696004*	0,841754228	LRRK1*	leucine rich repeat kinase 1*
TLMN_2249018*	-1,261858321	LRRK2*	leucine rich repeat containing 2*
TLMN_1877186*	0,659291140	LRRK3*	leucine rich repeats and calpain homology domain containing 3*
TLMN_1865736*	-0,968868894	LOC729451*	uncharacterized LOC729451*
TLMN_1715780*	0,743247756	LCAL5*	galectin 5*
TLMN_1654880*	0,769300508	LENG8*	leukocyte receptor cluster member 8*
TLMN_1683792*	1,054086743	LAP3*	leucine aminopeptidase 3*
TLMN_1713581*	0,824866134	LAMP5*	lysosomal associated membrane protein family member 5*
TLMN_2170814*	2,811254393	LAMP3*	lysosomal associated membrane protein 3*
TLMN_2386798*	-0,910980148	KLRC3*	koller cell lectin like receptor C3*
TLMN_1748123*	0,899101032	KLHL14*	kelch like family member 14*
TLMN_1673521*	-0,860008988	KISS1R*	KISS1 receptor*
TLMN_1771482*	-0,682017346	KIAA1324*	KIAA1324*
TLMN_1738488*	-0,587034393	KEL*	Kell blood group, metallo-endopeptidase*
TLMN_1798224*	0,807880028	KDM5C*	lysine demethylase 5C*
TLMN_1733811*	0,847581470	JUP*	junction plakoglobin*
TLMN_1741870*	0,881844543	JUP*	junction plakoglobin*
TLMN_2165993*	-0,827628834	ITLN1*	intelectin 1*
TLMN_1796755*	-0,592288813	ITGB5*	integrin subunit beta 5*
TLMN_2054019*	1,801347541	ISG15*	ISG15 ubiquitin-like modifier*

Tabela 6 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE 3. (Reactome, 2020)

ID	LOG FC	GENE SYMBOL	GENE TITLE
TILMN_1811488	4,813193324	*IRX3*	*trojanus homeobox 3*
TILMN_2349961	0,710597537	*IRF7*	*interferon regulatory factor 7*
TILMN_2340643	-0,69407111	*INSC*	*inscutable homolog (Drosophila)*
TILMN_1007710	0,614631629	*IL37*	*interleukin 37*
TILMN_2207291	1,054802508	*IFNG*	*interferon gamma*
TILMN_1005750	1,63215301	*IFITM3*	*interferon induced transmembrane protein 3*
TILMN_1664543	1,737878725	*IFIT3*	*interferon induced protein with tetratricopeptide repeats 3*
TILMN_1701789	1,552718774	*IFIT3*	*interferon induced protein with tetratricopeptide repeats 3*
TILMN_2239754	1,108060128	*IFIT3*	*interferon induced protein with tetratricopeptide repeats 3*
TILMN_1738428	1,024439098	*IFIT2*	*interferon induced protein with tetratricopeptide repeats 2*
TILMN_1707695	1,972678302	*IFIT1*	*interferon induced protein with tetratricopeptide repeats 1*
TILMN_1781373	0,940745037	*IFIH1*	*interferon induced with helicase C domain 1*
TILMN_2347700	1,729058907	*IFI6*	*interferon alpha inducible protein 6*
TILMN_1723812	2,836717413	*IFI44L*	*interferon induced protein 44 like*
TILMN_1760062	2,188853084	*IFI44*	*interferon induced protein 44*
TILMN_1745374	0,664145947	*IFI35*	*interferon induced protein 35*
TILMN_2058782	2,834600357	*IFI27*	*interferon alpha inducible protein 27*
TILMN_1656310	0,842928639	*IDO1*	*indoleamine 2,3-dioxygenase 1*
TILMN_1671089	-0,955395475	*IOMER2*	*iomer scaffolding protein 2*
TILMN_1705984	0,684548429	*HNMT*	*histamine N-methyltransferase*
TILMN_2066050	0,930573654	*HLA-DRA*	*major histocompatibility complex, class II, DR beta 6 (pseudogene)*
TILMN_2198239	-0,94809176	*HGD*	*homogentisase 1,2-dioxygenase*
TILMN_1653480	1,476989795	*HES4*	*hes family BHLH transcription factor 4*
TILMN_1787509	0,377920187	*HEL22*	*helicase with zinc finger 2*
TILMN_2318568	-1,194683200	*HCC-CIR1*	*host cell factor C1 regulator 1*
TILMN_2121400	0,664608641	*HDCGF*	*heparin binding EGF like growth factor*
TILMN_1723139	1,037879173	*HGD2*	*glycerol-3-phosphate dehydrogenase 2*
TILMN_1005221	-0,633841260	*GOLGA7*	*golgin A7*
TILMN_1723049	-0,989166211	*GJC2*	*gap junction protein gamma 2*
TILMN_1708002	0,613528923	*GCH1*	*GTP cyclohydrolase 1*
TILMN_1758053	1,508310949	*GBP6*	*guanylate binding protein family member 6*
TILMN_2114568	1,004074476	*GBP5*	*guanylate binding protein 5*
TILMN_1771385	0,758321776	*GBP4*	*guanylate binding protein 4*
TILMN_1701114	1,396427094	*GBP3*	*guanylate binding protein 3*
TILMN_2148789	1,182911945	*GBP1*	*guanylate binding protein 1*
TILMN_1698725	0,631434156	*FERMD3*	*FERM domain containing 3*
TILMN_1655325	-0,840661787	*FLOC*	*foliculin*
TILMN_1814952	-0,786884214	*FLOC*	*foliculin*
TILMN_1762531	-0,717990803	*FGF9*	*fibroblast growth factor 9*
TILMN_1727982	-0,68496644	*FEZ1*	*fertilization and elongation protein zeta 1*
TILMN_1779031	0,65829578	*FEZ1*	*fertilization and elongation protein zeta 1*
TILMN_1654389	0,650096054	*FCG3RA*	*Fc fragment of IgG receptor 2a*
TILMN_2302757	-0,959266278	*FCG3BP*	*Fc fragment of IgG binding protein*
TILMN_1717083	-0,640348463	*FBXO9*	*F-box protein 9*
TILMN_1701955	0,712298529	*FBXO9*	*F-box protein 9*
TILMN_2380919	-0,69407602	*FBXN2*	*fublin 2*
TILMN_1701012	0,782069512	*FAM720*	*family with sequence similarity 72 member C*
TILMN_1797301	-1,083752628	*FAM712B*	*family with sequence similarity 712 member B*
TILMN_2401253	0,602542867	*FAM33A*	*family with sequence similarity 33 member A*
TILMN_1708105	0,594883345	*EZH2*	*enhancer of zeste 2 polycomb repressive complex 2 subunit*
TILMN_1708747	1,030497058	*EXO3*	*endonuclease (5'-3'), endonuclease G-like*
TILMN_1709032	-0,598656142	*EVA1B*	*eva-1 homolog B*
TILMN_1700671	1,255586923	*ETV7*	*ETS variant 7*
TILMN_2388547	2,146902559	*EPST11*	*epithelial stromal interaction 1 (broad)*
TILMN_2323427	-0,635481464	*EPI641*	*erythrocyte membrane protein band 4.1*
TILMN_1708502	1,239439307	*EIF2AK2*	*eukaryotic translation initiation factor 2 alpha kinase 2*
TILMN_1767322	-0,589674342	*EDAR*	*ectodysplasin A receptor*
TILMN_1730284	-0,945551939	*DSC1*	*desmoglein 1*
TILMN_2402640	-0,757310987	*DSC1*	*desmoglein 1*
TILMN_2047112	-0,604385499	*DPCD*	*deleted in primary ciliary dyskinesia homolog (mouse)*
TILMN_1678422	0,994053554	*DDX56*	*DEAD-box helicase 56*
TILMN_1795181	1,354480159	*DDX60*	*DEAD-box helicase 60*
TILMN_1797001	0,85279248	*DDX58*	*DEAD-box helicase 58*
TILMN_1818304	-0,750823795	*DDX57*	*DEAD-box helicase 57*
TILMN_1801008	0,606586642	*DAC1*	*deshvechled binding antagonist of beta casein 1*
TILMN_1708303	1,13091854	*CYP4F2*	*cytochrome P450 family 4 subfamily F member 22*
TILMN_1705403	0,683477791	*CYP2S1*	*cytochrome P450 family 2 subfamily S member 1*
TILMN_1651752	0,732825373	*CXOR21*	*chromosome X open reading frame 21*
TILMN_1674640	-0,751589724	*CXCR6*	*CX-C motif chemokine receptor 6*
TILMN_1752562	-0,716741967	*CXCL5*	*CX-C motif chemokine ligand 5*
TILMN_1791759	1,772280073	*CXCL10*	*CX-C motif chemokine ligand 10*
TILMN_1812995	0,689917682	*CTSL*	*cathepsin L*
TILMN_2374036	0,604399096	*CTSL*	*cathepsin L*
TILMN_1737932	-0,592405952	*CRAT*	*camitine O-acetyltransferase*
TILMN_1684724	0,795998783	*CR2*	*complement component 3d receptor 2*
TILMN_1788551	-0,691475495	*CPAS*	*carboxypeptidase A3*
TILMN_2184049	0,657009139	*COX7B*	*cytochrome c oxidase subunit 7b*
TILMN_1783909	-1,052385664	*COL6A2*	*collagen type VI alpha 2 chain*

Tabela 7 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE 4. (Reactome, 2020)

ID	LOG FC	GENE SYMBOL	GENE TITLE
TLMN_1711514*	-0,900645828	COCH*	"cochlin"
TLMN_1703821*	3,120357909	CMPK2*	"cytidine/uridine monophosphate kinase 2"
TLMN_2066274*	-0,04540804	CLEC9A*	"C-type lectin domain family 9 member A"
TLMN_1662259*	-0,0007875	CLEC4E*	"C-type lectin domain family 4 member C"
TLMN_1756328*	0,756285526	CKS2*	"CDC28 protein kinase regulatory subunit 2"
TLMN_2072290*	0,587892781	CKS2*	"CDC28 protein kinase regulatory subunit 2"
TLMN_1732534*	0,859739044	CHMP5*	"charged multivesicular body protein 5"
TLMN_2270100*	0,770493447	CEP85L*	"centrosomal protein 85 kDa"
TLMN_1716815*	-0,75466514	CEACAM1*	"carcinoembryonic antigen related cell adhesion molecule 1"
TLMN_2371172*	0,879664781	CEACAM1*	"carcinoembryonic antigen related cell adhesion molecule 1"
TLMN_1663390*	0,790686172	CDC20*	"cell division cycle 20"
TLMN_1651316*	0,865871514	CD69*	"CD69 molecule"
TLMN_2233783*	0,791955854	CD38*	"CD38 molecule"
TLMN_1701814*	0,591561057	CD274*	"CD274 molecule"
TLMN_1726589*	-1,868645392	CD248*	"CD248 molecule"
TLMN_1678833*	0,582756103	CCR1*	"C-C motif chemokine receptor 1"
TLMN_1786125*	0,753668939	CCNA2*	"cyclin A2"
TLMN_2107088*	0,612738046	CCNA1*	"cyclin A1"
TLMN_1772864*	0,865724277	CCL8*	"C-C motif chemokine ligand 8"
TLMN_1764038*	0,886797514	CCL23*	"C-C motif chemokine ligand 23"
TLMN_1720448*	1,090663678	CCL7*	"C-C motif chemokine ligand 7"
TLMN_1707979*	1,228384848	CARD17*	"caspase recruitment domain family member 17"
TLMN_1737088*	-0,880789382	CAPN5*	"calpain 5"
TLMN_1678873*	-0,75098788	CABP5*	"calcium binding protein 5"
TLMN_1668317*	-0,587152051	C5AR2*	"complement component 5a receptor 2"
TLMN_1810752*	-0,68210541	C4BPA*	"complement component 4 binding protein alpha"
TLMN_2224486*	0,840252227	C3orf14*	"chromosome 3 open reading frame 14"
TLMN_1710746*	0,618423176	C2*	"complement component 2"
TLMN_1662578*	0,89788903	C1GALT1*	"core 1 synthase, glycoprotein-N-acetylgalactosamine 2-beta-galactosyltransferase 1"
TLMN_1690241*	-1,60551053	BATF2*	"basic leucine zipper ATF-like transcription factor 2"
TLMN_1734655*	-0,505477781	ATG9B*	"mitochondry related 9B"
TLMN_2374005*	1,394750440	ATF3*	"activating transcription factor 3"
TLMN_1815184*	0,864788708	ASPM*	"abnormal spindle microtubule assembly"
TLMN_1695414*	-0,704874268	ASF1B*	"anti-silencing function 1B histone chaperon"
TLMN_2186878*	0,808467828	ARMT1*	"arabic rosidae methyltransferase 1"
TLMN_1668853*	0,657397512	ARHGAP25*	"Rho GTPase activating protein 25"
TLMN_1800182*	0,773601001	APOBEC3A*	"apolipoprotein B mRNA editing enzyme catalytic subunit 3A"
TLMN_1708368*	1,066788955	ANKRD22*	"ankyrin repeat domain 22"
TLMN_2132539*	1,03008408	ANKRD22*	"ankyrin repeat domain 22"
TLMN_1783443*	0,587794208	ALOX15*	"arachidonate 15-lipoxygenase"
TLMN_2387712*	-0,812661559	AK5*	"adenylyate kinase 5"
TLMN_1681301*	0,625281321	AIM2*	"absent in melanoma 2"
TLMN_1770454*	0,912973843	AGRN*	"agrin"
TLMN_1707925*	-0,853824337	ABHD12B*	"aldehyde dehydrogenase domain containing 12B"
TLMN_1692417*	0,884394504		
TLMN_1694638*	0,76439301		
TLMN_1654892*	0,610390546		
TLMN_1657932*	-1,077861336		
TLMN_1658578*	0,60943002		
TLMN_1667444*	-1,287874744		
TLMN_1676708*	-0,581046645		
TLMN_1680757*	-1,115620056		
TLMN_1680772*	-0,749120346		
TLMN_1681728*	-0,868155826		
TLMN_1688172*	-0,751942346		
TLMN_1690365*	0,778591736		
TLMN_1693838*	-0,79789858		
TLMN_1694127*	-1,423217589		
TLMN_1695485*	0,843233478		
TLMN_1708004*	0,919373435		
TLMN_1710844*	0,628541854		
TLMN_1723433*	-0,615810815		
TLMN_1730491*	0,645850005		
TLMN_1735364*	0,620213353		
TLMN_1740317*	0,669679185		
TLMN_1751289*	-0,653068883		
TLMN_1756992*	0,624357768		
TLMN_1788551*	0,76327392		
TLMN_1770785*	-0,77428903		
TLMN_1773887*	-0,817404837		
TLMN_1782407*	1,7087001		
TLMN_1798101*	0,845883858		
TLMN_1804350*	-1,061582287		
TLMN_1808768*	0,679412781		
TLMN_1822871*	-2,31239811		
TLMN_1834888*	-0,625708477		
TLMN_1860003*	0,638125151		
TLMN_1882108*	-0,631667878		

Tabela 8 – APÊNDICE: Lista de transcritos estatisticamente significantes - PARTE 5.
(Reactome, 2020)

ID	LOG FC	GENE SYMBOL	GENE TITLE
"ILMN_1870005"	+1,000109527		
"ILMN_1887011"	-0,714887909		
"ILMN_1889950"	0,707957407		
"ILMN_1888691"	0,614521730		
"ILMN_2286080"	0,586810841		
"ILMN_2377829"	0,712354649		

APÊNDICE C – Informações relacionadas ao software R e ao banco de dados.

- Fonte dos dados: <<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE51092>>
 - Dados brutos: `GSE51092RAW.tar`
 - Dados não normalizados: `GSE51092nonnormalized.txt.gz`
 - Utilizou-se a opção: "Analyze with GEO2R"
- *Heatmap*: pacote *ComplexHeatmap* do projeto *Bioconductor*¹.
- Anotação genômica: "*ReactomePA*" do projeto *Bioconductor*.

¹ Bioconductor é um projeto de software com código aberto para bioinformática, que fornece ferramentas para a análise e compreensão de dados genômicos de alto rendimento utilizando linguagem de programação estatística R (BIOCONDUTOR, 2020) com o intuito de promover a compreensão dos ensaios biológicos atuais e emergentes de alto rendimento.