

Análise de Expressão Diferencial

Pablo Rodrigo Sanches

Departamento de Genética – FMRP/USP

psanches@usp.br

Roteiro de análise

1. DESeq2 (normalização e cálculo da expressão diferencial)
 - a) tratamento vs. controle
2. Obter os resultados:
 - a) PCA
 - b) Distância entre amostras
 - c) Dispersão
 - d) Histograma de p valor
 - e) Genes significativamente DE
 - f) Tabelas de genes significativamente DE

Normalização

- Descobrir como os genes se expressam diferencialmente em diferentes condições ou tecidos.
- Número de reads é correlacionado com o nível de expressão do gene
- RNA-Seq oferece uma aproximação quantitativa da abundância dos transcritos na forma de contagens.
- Contagens precisam ser normalizadas
 - Comprimento das diferentes moléculas de RNAs
 - Profundidade do sequenciamento de diferentes bibliotecas

Normalização

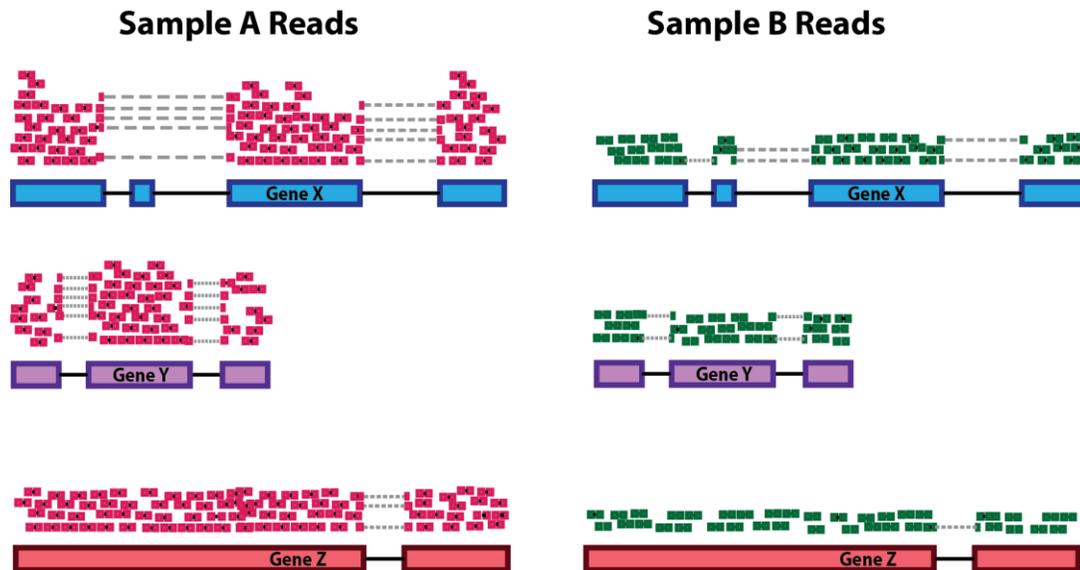
- Exemplo:
 - Biblioteca 1 → 12 milhões de reads mapeados
 - Biblioteca 2 → 16 milhões de reads mapeados

Loco	Tamanho loco (pb)	Nro. Reads Biblioteca 1	Nro. Reads Biblioteca 2
GeneA	800	24	38

É possível afirmar que temos maior expressão do GeneA na Biblioteca 2 quando comparada à Biblioteca 1 ???

Normalização

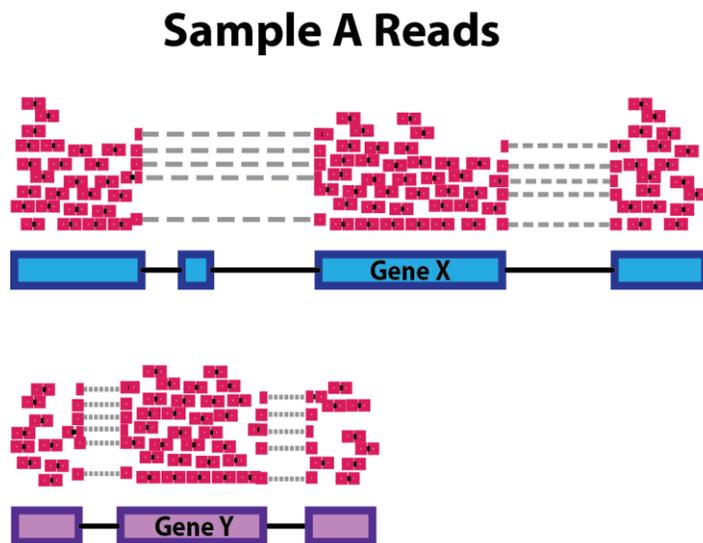
Profundidade do sequenciamento:



- No exemplo, cada gene parece ter dobrado a expressão na Amostra A em relação à Amostra B, no entanto, isso é uma consequência da Amostra A ter o dobro da profundidade de sequenciamento.

Normalização

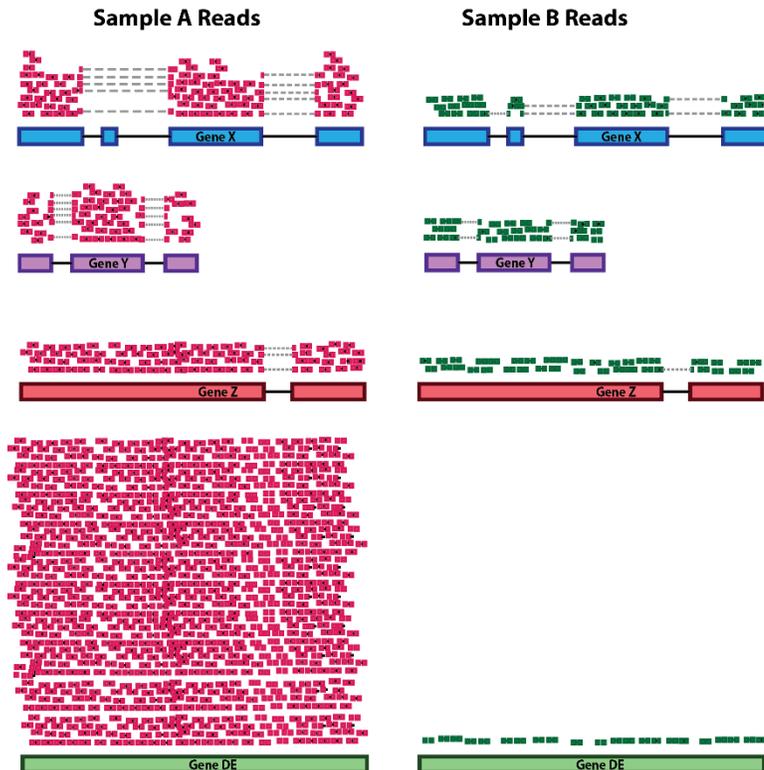
Comprimento do gene:



- No exemplo, *Gene X* e *Gene Y* têm níveis semelhantes de expressão, mas o número de leituras mapeadas para *Gene X* seria muito maior do que o número mapeado para *Gene Y* porque *Gene X* é mais longo.

Normalização

Composição de RNA:



- No exemplo, se fosse para dividir cada uma das amostras pelo número total de contagens para normalização, as contagens seriam muito enviesadas pelo Gene DE (ocupa a maior parte das contagens para uma das amostras).

DESeq2 – Mediana das razões

- O DESeq2 executa uma normalização onde a média geométrica é calculada para cada gene em todas as amostras. A contagem de um gene em cada amostra é então dividida por essa média. A mediana dessas proporções em uma amostra é o fator de tamanho dessa amostra. Este procedimento corrige o tamanho da biblioteca e o viés da composição do RNA, que pode surgir, por exemplo, quando apenas um pequeno número de genes são altamente expressos em uma condição de experimento, mas não na outra.

Exemplo de aplicação

	amostraA	amostraB	referência	amostraA/referência	amostraB/referência	normalizadoA	normalizadoB	FC	log2FC
geneA	2800	1000	1673,32	1,67	0,60	3429,29	816,50	0,24	-2,07
geneB	30	15	21,21	1,41	0,71	36,74	12,25	0,33	-1,58
geneC	500	750	612,37	0,82	1,22	612,37	612,37	1,00	0,00
geneD	40	80	56,57	0,71	1,41	48,99	65,32	1,33	0,42
geneE	500	1200	774,60	0,65	1,55	612,37	979,80	1,60	0,68
			Mediana	0,82	1,22				

Renomear arquivos Gene counts no Galaxy (saída do software StringTie)

The screenshot displays the Galaxy web interface for editing dataset attributes. The main content area shows the 'Edit dataset attributes' form with the following fields:

- Name:** Oh II StringTie on data 106 and data 109: Gene counts (highlighted in yellow)
- Info:** (Empty text area)
- Annotation:** (Empty text area)
- Database/Build:** ----- Additional Species Are Below -----
- Number of comment lines:** (Empty text area)

A black text box is overlaid on the 'Name' field, containing the text: `Oh II RNA STAR...; Oh III RNA STAR...; 3h II RNA STAR...; ...`. A yellow arrow points to the 'Save' button in the top right corner of the form.

The 'History' panel on the right side of the interface shows a list of datasets. The dataset '126: StringTie on data 106 and data 109: Gene counts' is highlighted in yellow, and a black arrow points to its 'x' icon.

DESeq2 – Via Galaxy (cont.)

The screenshot shows the Galaxy web interface for the DESeq2 tool. The central configuration panel includes the following fields and options:

- Specify a factor level, typical values could be 'tumor', 'normal', 'treated' or 'control'**: 3h
- Counts file(s)**: A list of files is shown, with the following files selected (highlighted in blue):
 - 129: 0h III StringTie on data 106 and data 112: Gene counts
 - 128: StringTie on data 106 and data 112: Assembled transcripts
 - 127: StringTie on data 106 and data 109: Transcript counts
 - 126: 0h II StringTie on data 106 and data 109: Gene counts
- 3: Factor level**: 0h
- Specify a factor level, typical values could be 'tumor', 'normal', 'treated' or 'control'**: 0h
- Counts file(s)**: A list of files is shown, with the following files selected (highlighted in blue):
 - 130: StringTie on data 106 and data 112: Transcript counts
 - 129: 0h III StringTie on data 106 and data 112: Gene counts
 - 128: StringTie on data 106 and data 112: Assembled transcripts
 - 127: StringTie on data 106 and data 109: Transcript counts
 - 126: 0h II StringTie on data 106 and data 109: Gene counts
 - 125: StringTie on data 106 and data 109: Assembled transcripts
- (Optional) provide a tabular file with additional batch factors to include in the model.**: Nothing selected
- Files have header?**: Yes
- Choice of Input data**: Count data (e.g. from HTSeq-count, featureCounts or StringTie)
- Visualising the analysis results**: Yes
- Output normalized counts table**: No

Annotations in Portuguese provide additional context:

- Nome do tratamento**: Points to the factor level field.
- Selecione os arquivos de contagem das réplicas do tratamento**: Points to the selected counts files.
- Arquivos de contagem possuem cabeçalho**: Points to the "Files have header?" option.
- Arquivos de entrada foram gerados pelo StringTie**: Points to the "Choice of Input data" dropdown.

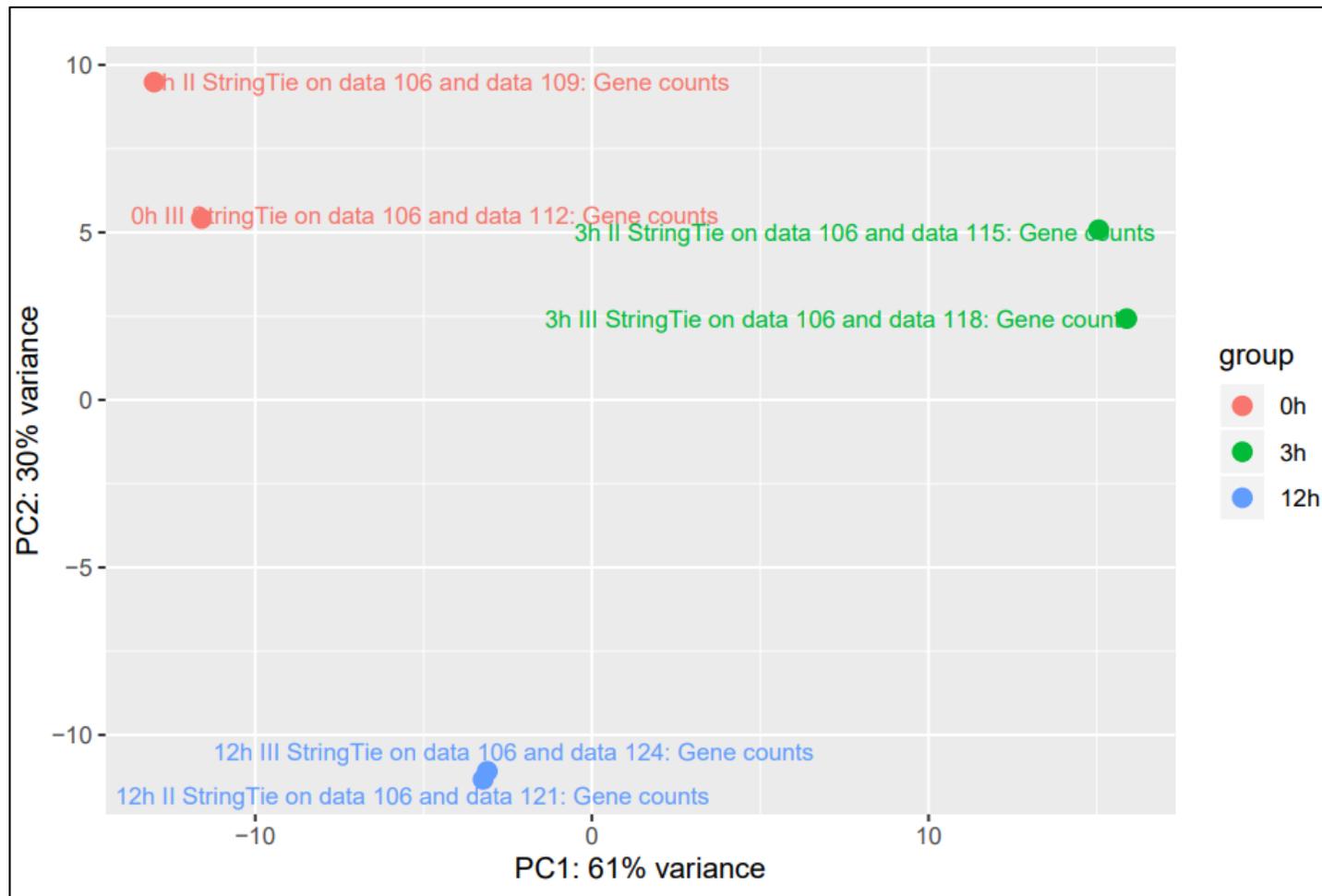
The right sidebar shows a history list of datasets, including:

- 142: StringTie on data 106 and data 124: Transcript counts
- 141: 12h III StringTie on data 106 and data 124: Gene counts
- 140: StringTie on data 106 and data 124: Assembled transcripts
- 139: StringTie on data 106 and data 121: Transcript counts
- 138: 12h II StringTie on data 106 and data 121: Gene counts
- 137: StringTie on data 106 and data 121: Assembled transcripts
- 136: StringTie on data 106 and data 118: Transcript counts
- 135: 3h III StringTie on data 106 and data 118: Gene counts
- 134: StringTie on data 106 and data 118: Assembled transcripts
- 133: StringTie on data 106 and data 115: Transcript counts
- 132: 3h II StringTie on data 106 and data 115: Gene counts

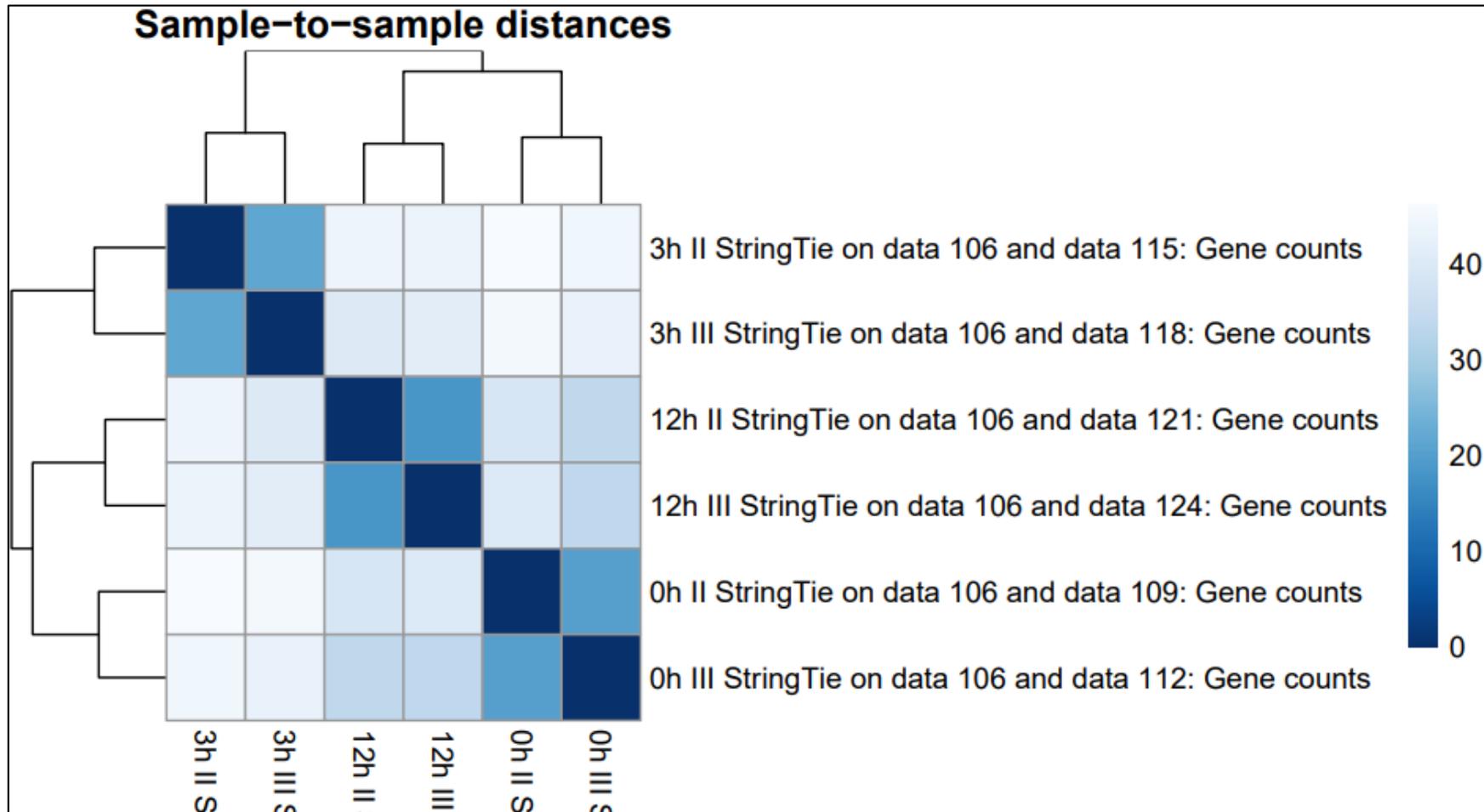
DESeq2 – Via Galaxy (cont.)

The screenshot displays the Galaxy web interface for the DESeq2 tool. The left sidebar lists various tools, with 'deseq2' selected. The main panel shows the tool's configuration options, including 'Output rLog normalized table', 'Output VST normalized table', and 'Output all levels vs all levels of primary factor'. The 'Output all levels vs all levels of primary factor' option is checked, and a black box with white text 'Gerar resultados para -> 12h vs. 0h; 3h vs. 0h; 12h vs. 3h' is overlaid on the 'Yes' radio button. The right sidebar shows a history of jobs, including '142: StringTie on data 106 and data 124: Transcript counts', '141: 12h III StringTie on data 106 and data 124: Gene counts', '140: StringTie on data 106 and data 124: Assembled transcripts', '139: StringTie on data 106 and data 121: Transcript counts', '138: 12h II StringTie on data 106 and data 121: Gene counts', '137: StringTie on data 106 and data 121: Assembled transcripts', '136: StringTie on data 106 and data 118: Transcript counts', '135: 3h III StringTie on data 106 and data 118: Gene counts', '134: StringTie on data 106 and data 118: Assembled transcripts', '133: StringTie on data 106 and data 115: Transcript counts', and '132: 3h II StringTie on d...'. The top navigation bar includes 'Analyze Data', 'Workflow', 'Visualize', 'Shared Data', 'Help', 'User', and 'Using 30%'.

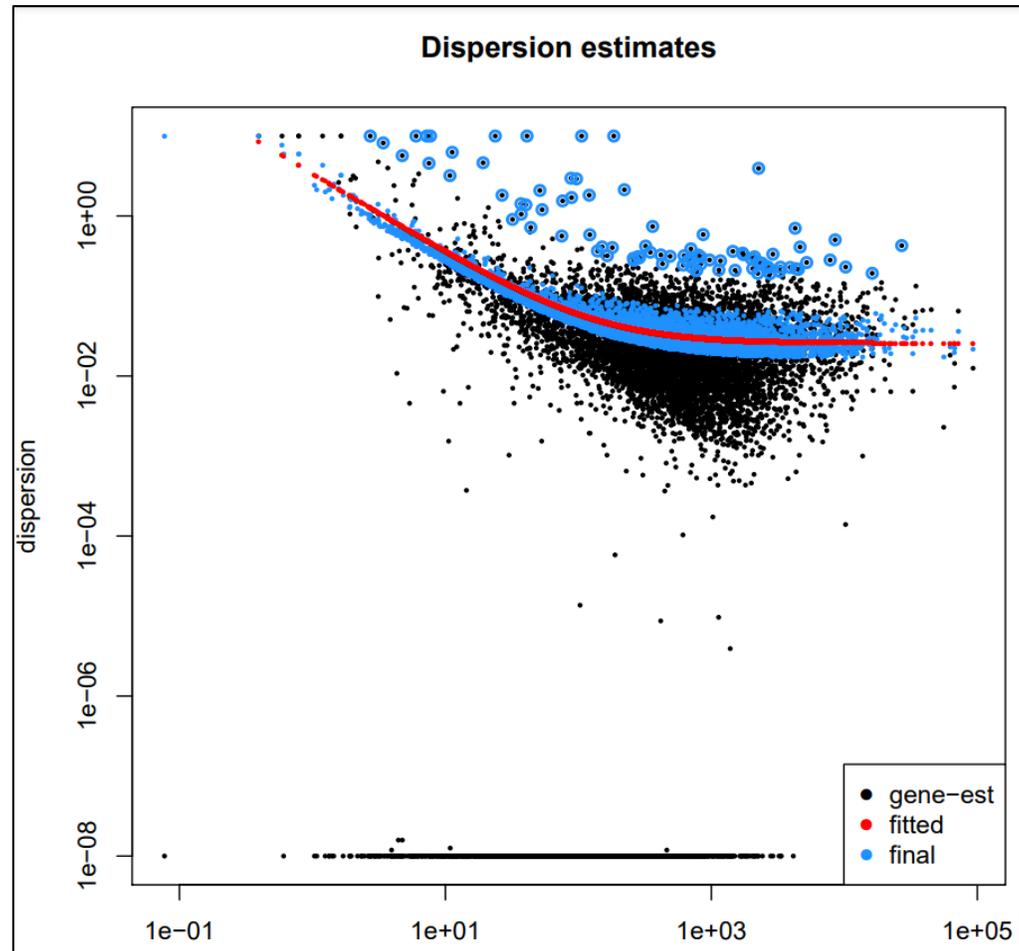
Exemplo de Resultado DESeq2 (PCA)



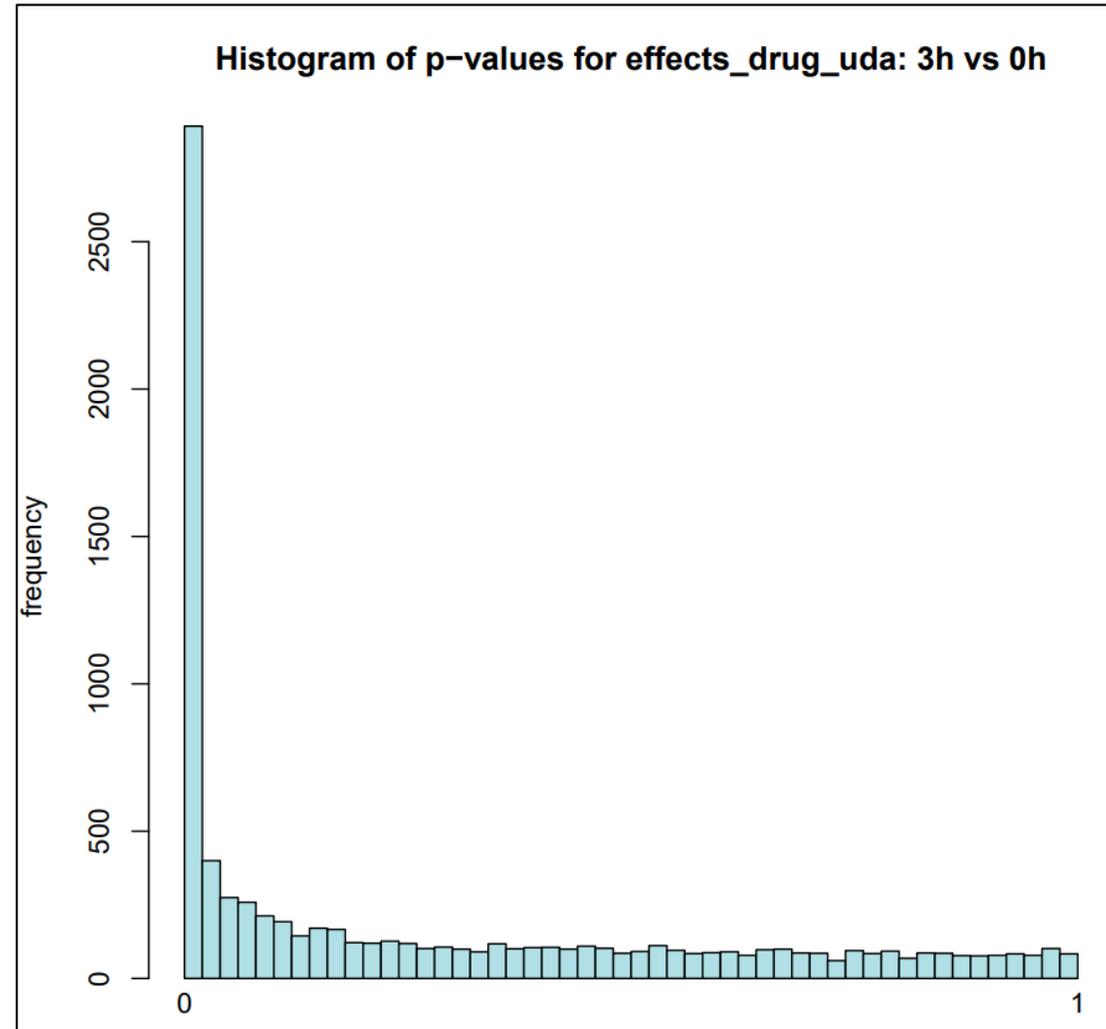
Exemplo de Resultado DESeq2 (Distância entre amostras)



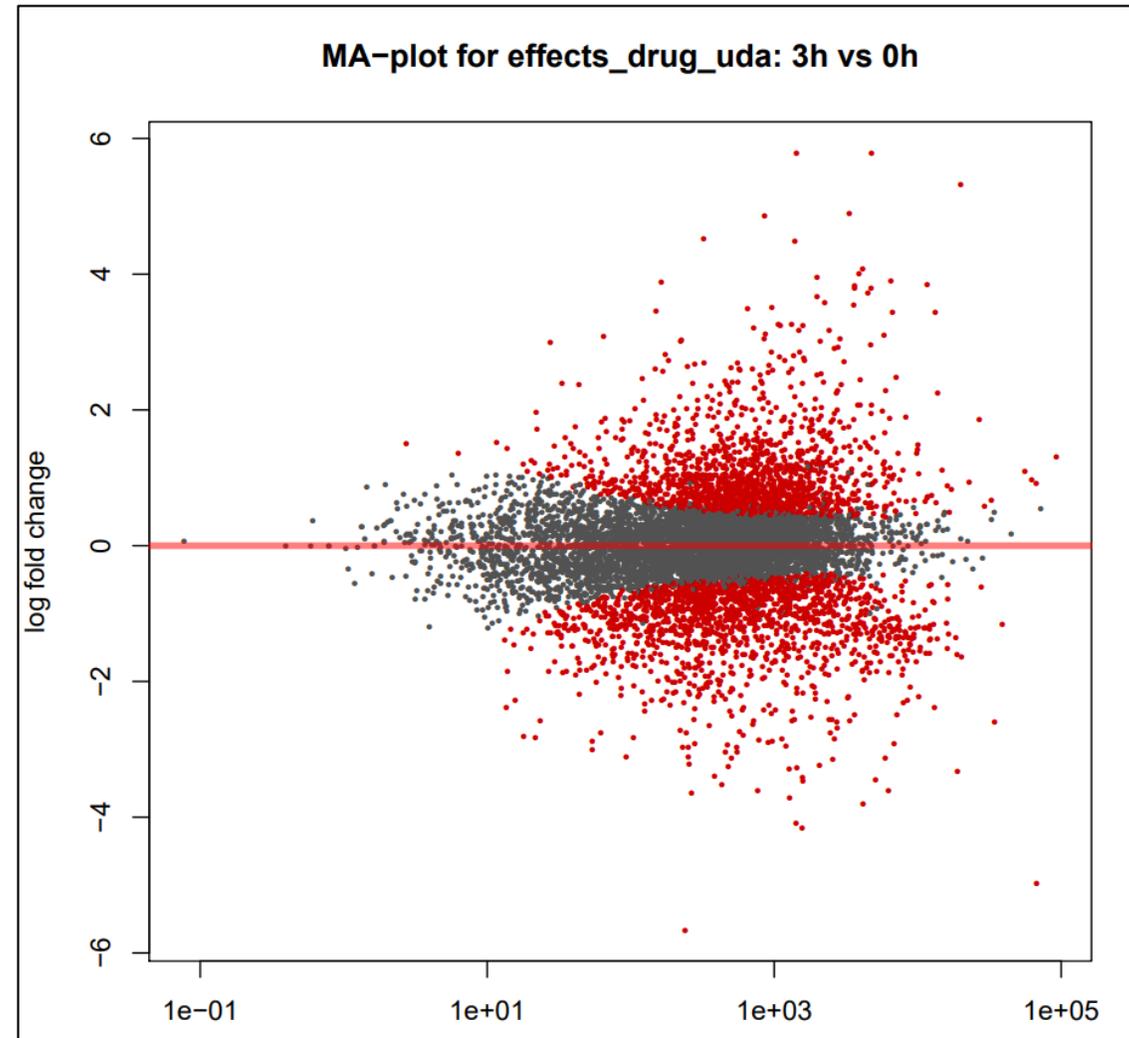
Exemplo de Resultado DESeq2 (Dispersão)



Exemplo de Resultado DESeq2 (Histograma de p valor)

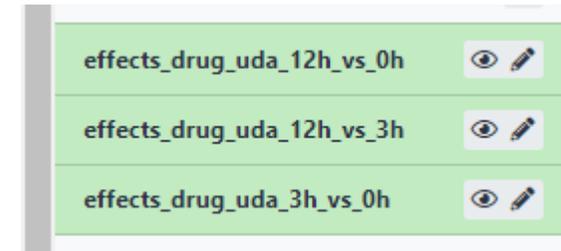
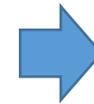
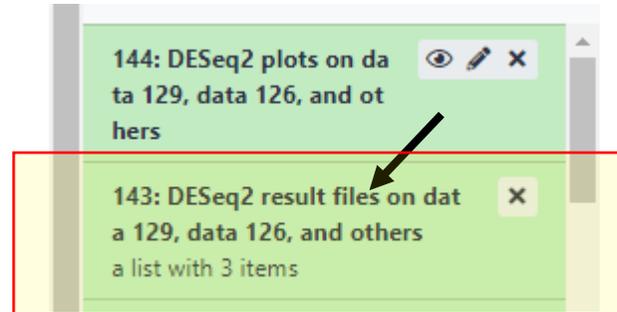


Exemplo de Resultado DESeq2 (DEG)



Exemplo de Resultado DESeq2 (Tabelas DEG)

- 12h vs 0h
- 12h vs 3h
- 3h vs 0h



Exemplo de Resultado DESeq2 (Tabelas DEG)

The screenshot displays the Galaxy web interface for a DESeq2 analysis. The main content area shows a table of differentially expressed genes (DEGs) with columns for gene ID, log2 fold change, and adjusted p-value. The right sidebar shows a 'History' panel with a list of analysis steps, including 'DESeq2 result files on data 129, data 126, and others'. A red box highlights the entry 'effects drug_uda 3h vs 0h' in the history panel, with an arrow pointing to it.

Gene ID	log2 Fold Change	Adjusted P-value
TERG_04234	67473.0935325852	-4.96773316638587
TERG_00823	4132.88765837273	4.06899139801621
TERG_04960	4770.59303749011	5.78488740348485
TERG_01937	857.116355441076	4.85761395468055
TERG_00254	1398.00852381666	4.48077471620284
TERG_08077	3340.49253836308	4.89681919985616
TERG_05484	6513.77903450769	3.90227713189188
TERG_04232	3638.99663086208	3.79821989809618
TERG_03078	19983.8897966338	5.32312644340679
TERG_01512	2259.41269421887	3.57475612970528
TERG_12530	3911.67702738133	4.00392296410416
TERG_07143	4763.34375494071	3.78624381937886
TERG_01405	1437.7200729774	5.78750556309314
TERG_05621	11715.207019069	3.85582247514341
TERG_03483	6693.75976184962	3.43140366618469
TERG_03343	1442.27418969119	-3.26396704750338
TERG_02747	1322.89955674594	3.26315654491761
TERG_04038	3596.14089525809	3.55334456602225
TERG_02023	4161.01592916159	-7.79981320939872
TERG_03154	13288.0374927027	3.43367452666396
TERG_01742	1595.02730840687	3.23975585567722
TERG_01784	2790.09704216066	2.916537623235819
TERG_07516	240.81849145649	-5.66325605793869
TERG_07830	324.252793994781	4.51423142835636
TERG_04500	1492.1736296593	3.16538719947916
TERG_07659	4705.29046710537	2.96071296180051
TERG_06509	3643.64855241637	3.82830305162481
TERG_02909	2000.63666467812	3.6658077338571
TERG_00162	4498.93228100454	3.73028325741556
TERG_08004	1420.88488847186	-4.08815590174838
TERG_05363	1991.57594480271	3.95886120136674
TERG_02423	722.981589991999	3.21161474145848
TERG_01233	1072.35006917928	3.26609711136417
TERG_01599	1213.92544601065	-2.9545252360843
TERG_12351	871.792090012748	3.1266671836548
TERG_06883	2564.7662443609	-3.15312741426125
TERG_00911	1272.05775029198	-3.29654407042257
TERG_03832	957.665313231076	2.84658814047705
TERG_11538	1636.25129894608	2.72093724828664
TERG_11928	966.524531905092	3.50670334044019
TERG_07919	1511.93960325045	2.85020305994746
TERG_07721	1107.15482889117	3.23681520258725
TERG_08363	973.856612973878	-2.88644254490119
TERG_02918	2101.26616512965	3.01383486413062
TERG_07691	3073.43754501511	2.71307769118404
TERG_04714	2417.35366459197	3.17347079867012

Sobre as colunas da tabela de resultados DESeq2

Coluna	Descrição
1	geneid = Identificador do gene
2	baseMean = média das contagens normalizadas tomadas em todas as amostras
3	log2FoldChange = mudança log2 vezes entre os grupos. Por exemplo, valor 2 significa que a expressão aumentou 4 vezes
4	lfcSE = erro padrão da estimativa log2FoldChange
5	stat = estatística de Wald
6	pvalue = valor p do teste de Wald
7	padj = valor p ajustado de Benjamini-Hochberg (Taxa de Falsas Descobertas - FDR)

Download dos resultados

The screenshot displays the Galaxy web interface. The main area shows a list of data files with columns for file name, size, and other metadata. The right-hand panel shows the 'History' section, which lists recent jobs. The job 'effects_drug_uda_3h_vs_0h' is highlighted with a red box, and a black arrow points to the 'Download' button below it. Another black arrow points to the 'Download' button in the history panel.

File Name	Size	Other Info
TERG_04234	67473.0935325852	-4.96773316638587
TERG_00823	4132.88765837273	4.06899139801621
TERG_04960	4770.59303749011	5.78488740348485
TERG_01937	857.116355441076	4.85761395468055
TERG_00254	1398.00852381666	4.48077471620284
TERG_08077	3340.49253836308	4.89681919985616
TERG_05484	6513.77903450769	3.90227713189188
TERG_04232	3638.99663086208	3.79821989809618
TERG_03078	19983.8897966338	5.32312644340679
TERG_01512	2259.41269421887	3.57475612970528
TERG_12530	3911.67702738133	4.00392296410416
TERG_07143	4763.34375494071	3.78624381937886
TERG_01405	1437.7200729774	5.78750556309314
TERG_05621	11715.207019069	3.85582247514341
TERG_03483	6693.75976184962	3.43140366618469
TERG_03343	1442.27418969119	-3.26396704750338
TERG_02747	1322.89955674594	3.26315654491761
TERG_04038	3596.14089525809	3.55334456602225
TERG_02023	4161.01592916159	-3.79981320939872
TERG_03154	13288.0374927027	3.43367452666396
TERG_01742	1595.02730840687	3.23975585567722
TERG_01784	2790.09704216066	2.91653762325819
TERG_07516	240.81849145649	-5.66325605793863
TERG_07830	324.252793994781	-4.51423142835636
TERG_04500	1492.1736296593	3.16538719947916
TERG_07659	4705.29046710537	2.96071296180051
TERG_06509	3643.64855241637	3.82830305162481
TERG_02909	2000.63666467812	3.6658077338571
TERG_00162	4498.93228100454	3.73028325741556
TERG_08004	1420.88488847186	-4.08815590174838
TERG_05363	1991.57594480271	3.95886120136674
TERG_02423	722.981589991999	3.21161474145845
TERG_01233	1072.35006917928	3.26609711136417
TERG_01599	1213.92544601065	-2.9545252360843
TERG_12351	871.792090012748	3.1266671836548
TERG_06883	2564.7662443609	-3.15312741426125
TERG_00911	1272.05775029198	-3.29654407042257
TERG_03832	957.665313231076	2.84658814047705
TERG_11538	1636.25129894608	2.72093724828664
TERG_11928	966.524531905092	3.50670334044019
TERG_07919	1511.93960325045	2.85020305994746
TERG_07721	1107.15482889117	3.23681520258725
TERG_08363	973.856612973878	-2.88644254490119
TERG_02918	2101.26616512965	3.01383486413062
TERG_07691	3073.43754501511	2.71307769118404