

# Small nucleolar RNAs in breast cancer

Laboratory of Biological Data Mining

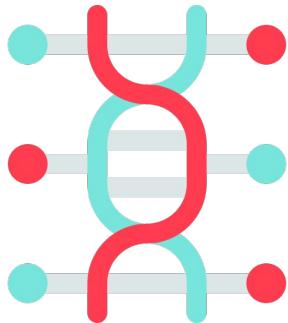
Group 1: Lorenzo Bocchi, Valentino Frasnelli, Erich Robbi, Annalisa Xamin



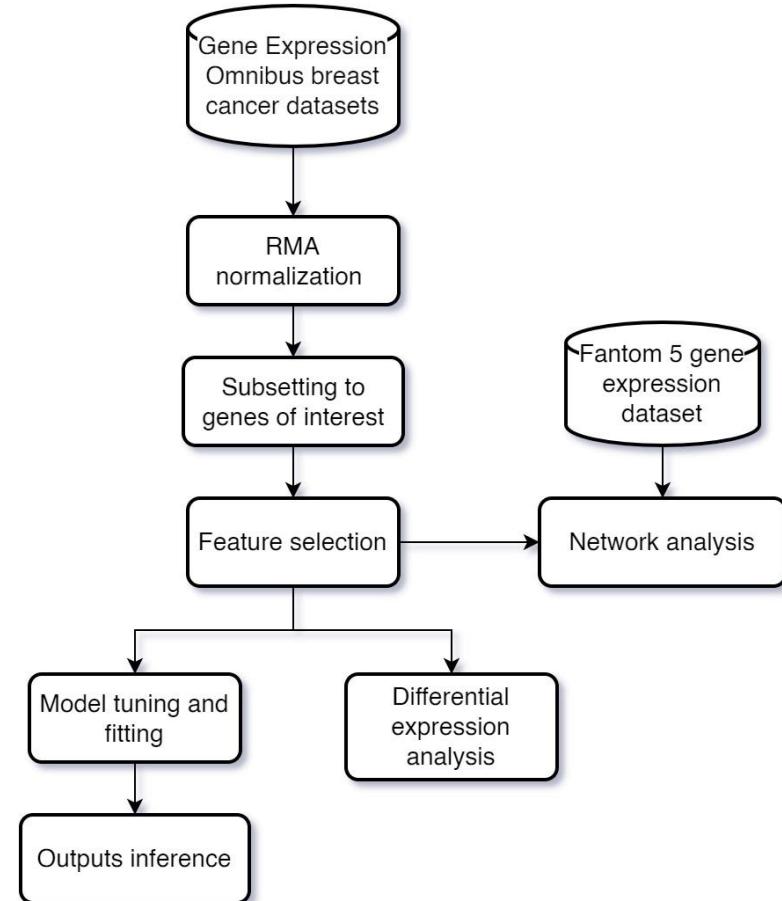
# Introduction



# Problem definition



# Project Pipeline



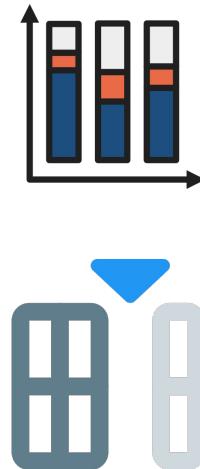


## Data collection





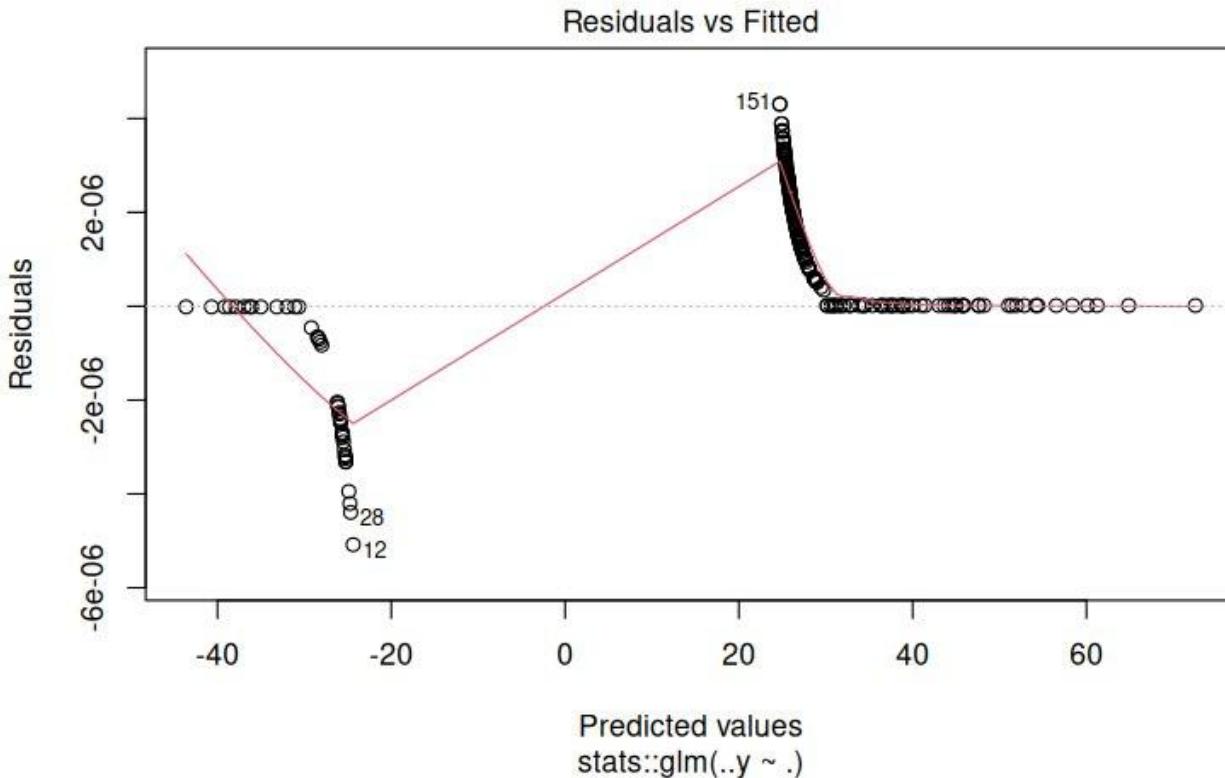
# Data preparation



# Model fitting



# First dataset - suspicious results





## Random forest performance measures on our genes of interest (old data)

	Accuracy	Precision	Recall	F1
No Cancer	0.9615	1	0.80	0.89
Cancer	0.9615	0.95	1	0.98



## Random forest fit on housekeeping genes (old data)

	Accuracy	Precision	Recall	F1
No Cancer	0.9221	1	0.60	0.75
Cancer	0.9221	0.91	1	0.95



## Random forest performance measures on our genes of interest (new data)

	Accuracy	Precision	Recall	F1
Basal	0.8571	0.75	0.82	0.78
HER	0.8286	0.63	0.63	0.63
L. A	0.8857	0.67	0.86	0.75
L. B	0.8833	0.60	0.75	0.67
Normal	1	1	1	1

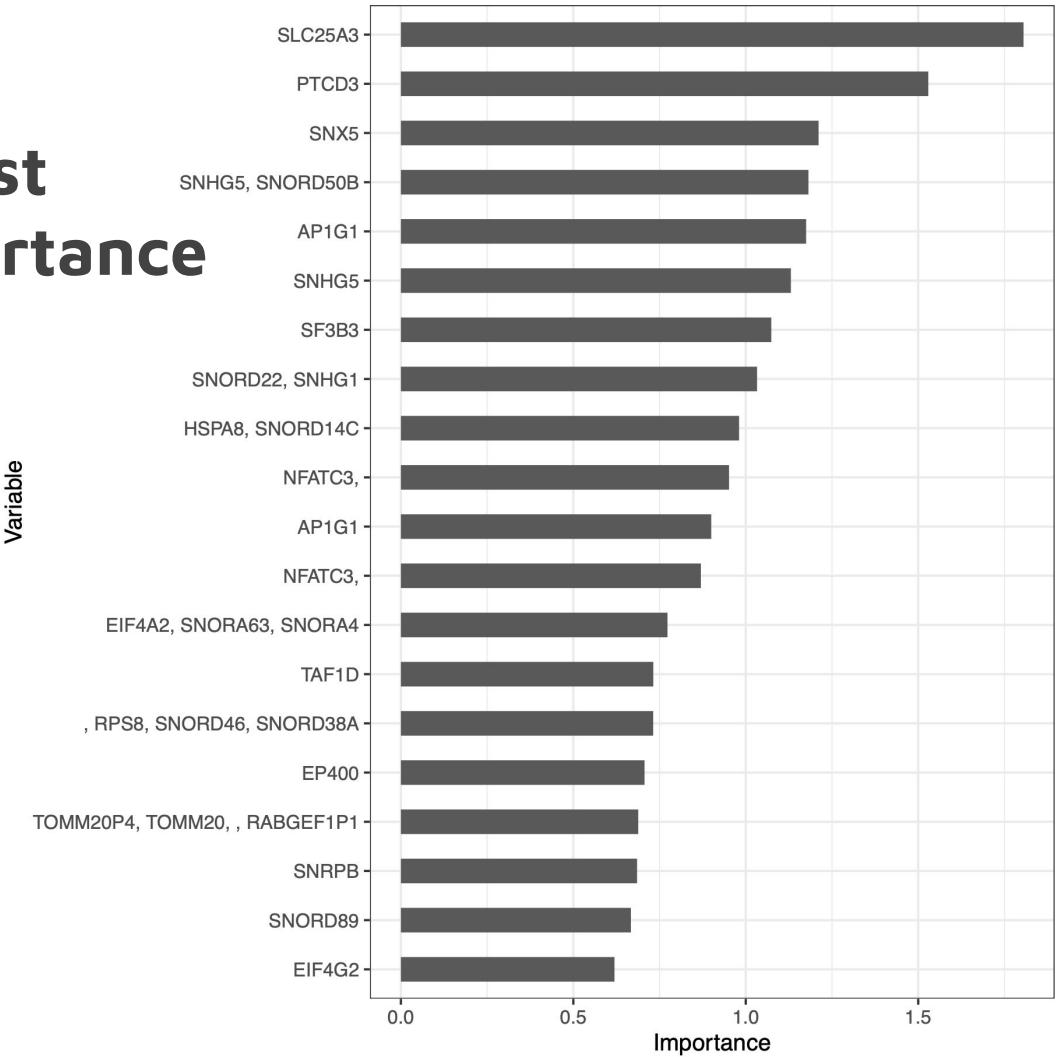


## Random forest performance metrics on lesser variant genes (new data)

	Accuracy	Precision	Recall	F1
Basal	0.6	0.41	0.64	0.5
HER	0.65	0	0	0
L. A	0.74	0.24	0.14	0.18
L. B	0.71	0.40	0.5	0.44
Normal	0	0	0	0

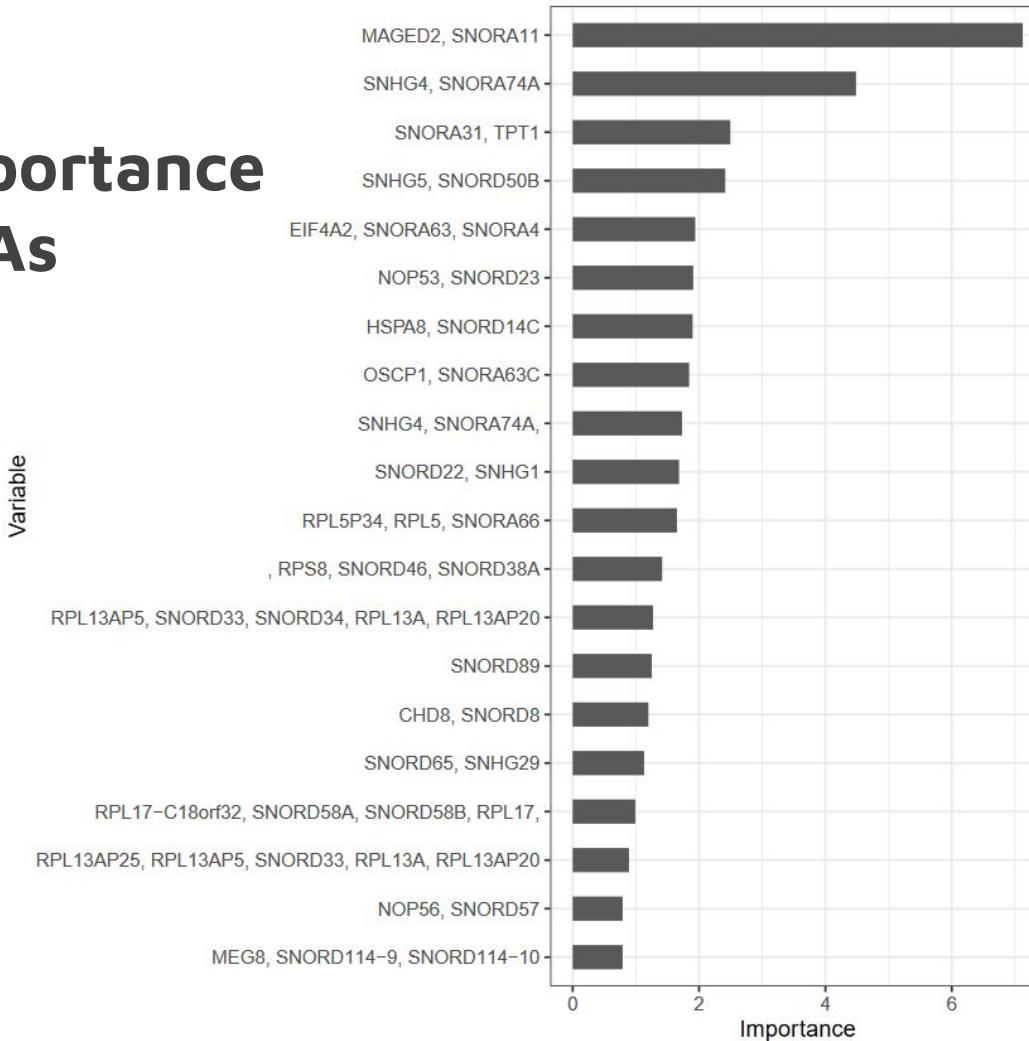


# Random Forest variable importance





# Variable importance only snoRNAs





# Final SVM model

Property	Value
Cost	0.005332322
Support Vectors	92/102
Training Error	0.019608



## Final SVM model performance measures

	Accuracy	Precision	Recall	F1
Basal	0.8857	0.77	0.91	0.83
HER	0.8857	0.75	0.75	0.75
L. A	0.9429	0.78	1	0.88
L. B	0.8857	1	0.5	0.67
Normal	1	1	1	1

	Accuracy	Precision	Recall	F1	
Basal	0.8571	0.75	0.82	0.78	Random forest
HER	0.8286	0.63	0.63	0.63	
L. A	0.8857	0.67	0.86	0.75	
L. B	0.8833	0.60	0.75	0.67	
Normal	1	1	1	1	

	Accuracy	Precision	Recall	F1	
SVM	Basal	0.8857	0.77	0.91	0.83
	HER	0.8857	0.75	0.75	0.75
	L. A	0.9429	0.78	1	0.88
	L. B	0.8857	1	0.5	0.67
	Normal	1	1	1	1



## Logistic regression with L1 penalty performance measures

	Accuracy	Precision	Recall	F1
Basal	0.8857	0.82	0.82	0.82
HER	0.8571	0.71	0.63	0.67
L. A	0.8571	0.60	0.86	0.71
L. B	0.8286	.67	0.5	0.57
Normal	1	1	1	1

# Most important variables

Coefficient (Probe id)	Estimate	Gene Symbols
214744_s_at	-5.807828e-01	RPL23
228971_at	-5.005921e-01	SLC25A3
232355_at	-4.876120e-01	MEG8, SNORD114-3
235102_x_at	4.474731e-01	SNORD3B-1, SNORD3B-2, SNORD3D, SNORD3A, SNORD3C
223666_at	4.409927e-01	SNX5
231096_at	-4.260707e-01	PCAT4
242856_at	-4.194403e-01	MEG8, SNORD114-9, SNORD114-10
221621_at	3.693238e-01	SNHG20
228879_at	-3.637258e-01	SNORD104, SNORA50C
241448_at	2.737583e-01	TOMM20



# Least important variables

Coefficient (Probe id)	Estimate	Gene Symbols
1555177_at	0	PRKAA1
215011_at	0	SNHG3
229038_at	0	CWF19L1
226428_at	0	TNPO2
240083_at	0	MEG8
200716_x_at	-8.590805e-05	RPL13AP5, RPL13A
200031_s_at	-1.930631e-04	RPS11
224741_x_at	-2.245289e-04	GAS5
200692_s_at	2.554860e-04	HSPA9
209476_at	-2.689458e-04	TMX1

# Network Analysis



## List of genes not yet expanded by gene@home

PCAT4	SNORD114-14	RPL23
SNORD114-1	RPS3	SNORD114-19
RACK1	SNORD113-4	NAN
SNORD114-13	RPS11	RPS2
SNORD114-20	SNHG20	SNORD115-23
SNORD116-13	CDKN2B-AS1	SNORD113-3
SNORD115-32	SLC25A3	ZFAS1
HIF1A-AS2	SNORD114-21	



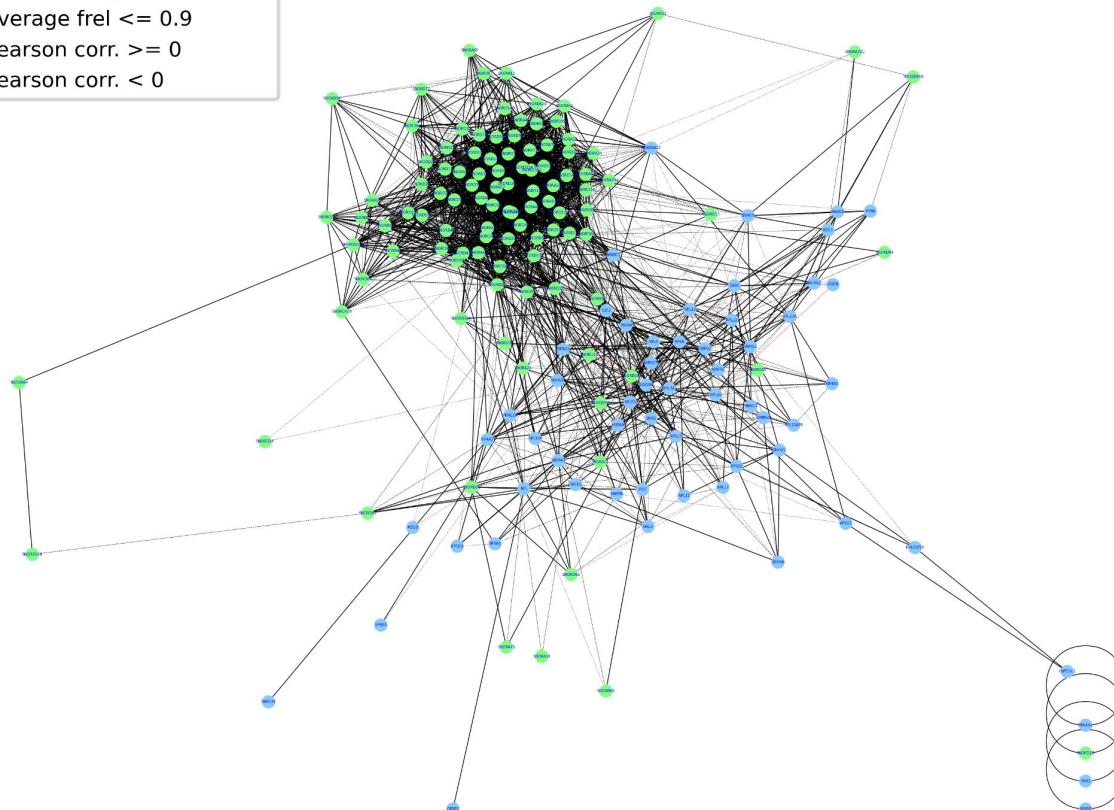
## List of genes not yet expanded by gene@home

PCAT4	<del>SNORD114-14</del>	RPL23
<del>SNORD114-1</del>	RPS3	<del>SNORD114-19</del>
RACK1	<del>SNORD113-4</del>	NAN
<del>SNORD114-13</del>	RPS11	RPS2
<del>SNORD114-20</del>	SNHG20	<del>SNORD115-23</del>
<del>SNORD116-13</del>	CDKN2B-AS1	<del>SNORD113-3</del>
<del>SNORD115-32</del>	SLC25A3	ZFAS1
<del>HIF1A-AS2</del>	<del>SNORD114-21</del>	

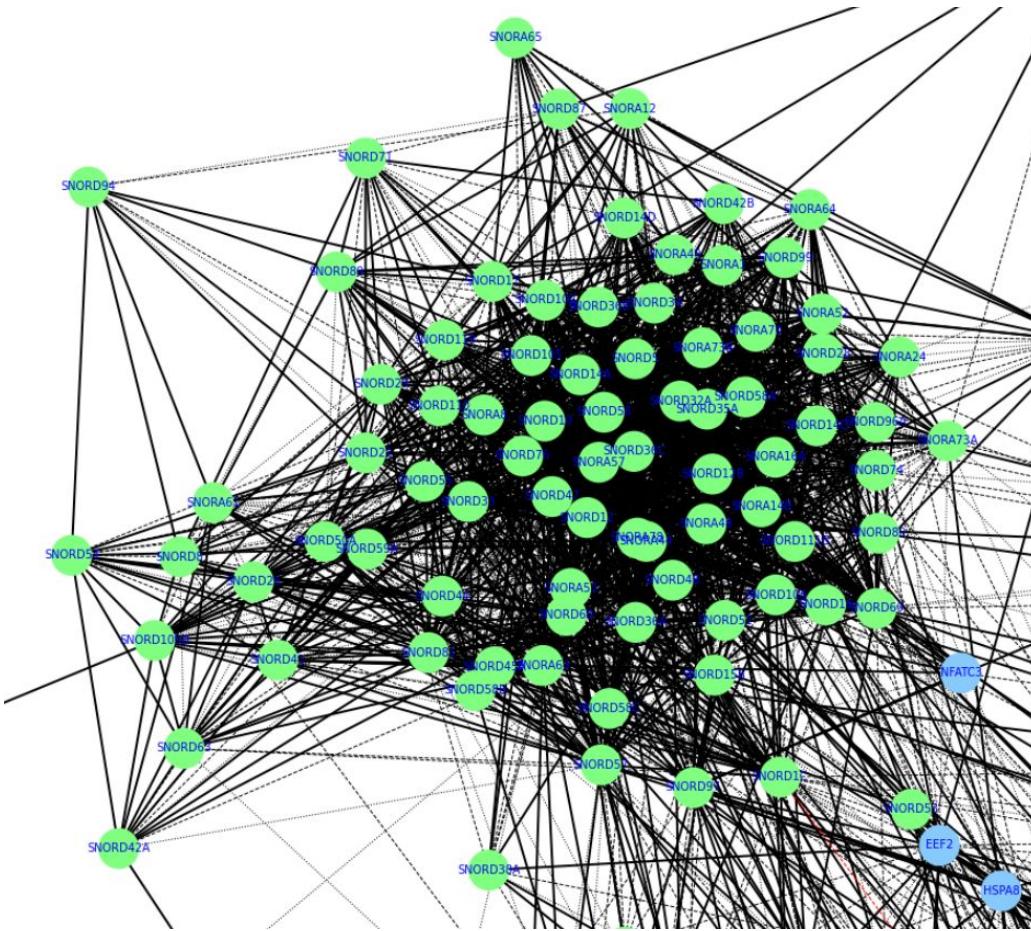


# Complete graph

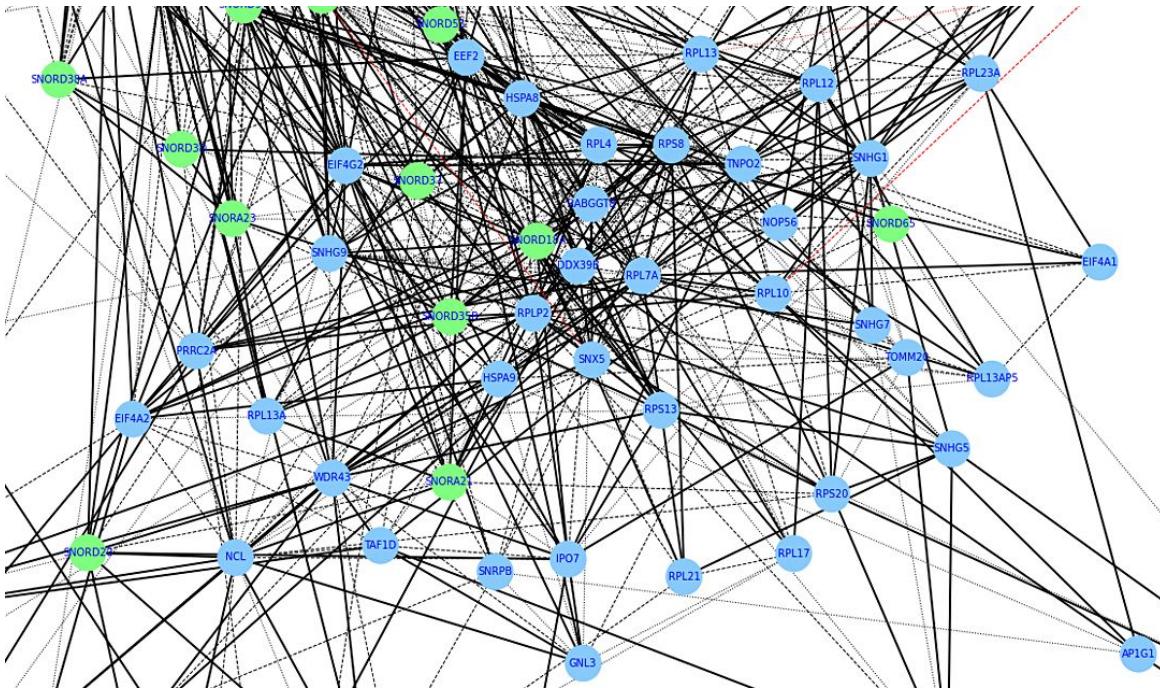
- average frel > 0.95
- - - 0.9 < average frel <= 0.95
- .... average frel <= 0.9
- Pearson corr. >= 0
- Pearson corr. < 0



# Complete graph

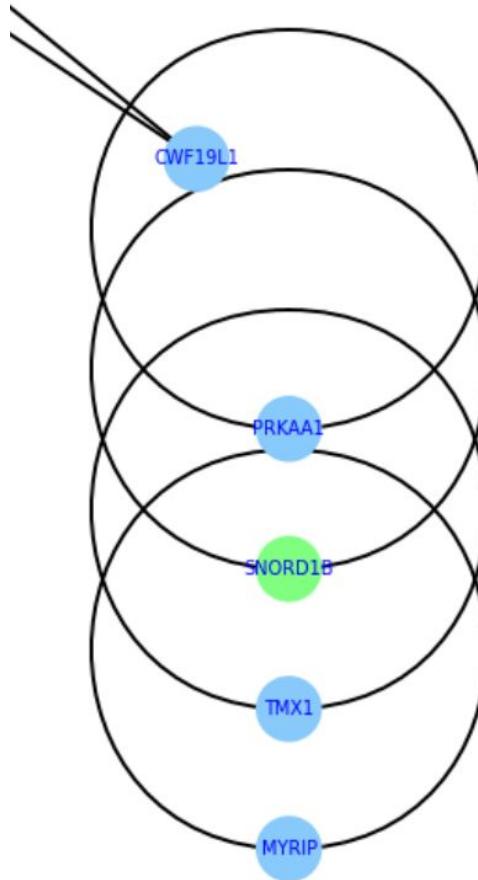


# Complete graph





# Complete graph





# Validation @0.85

SNORD41 → SNORD60

SNORA73B → SNORA73A

SNORA78 → SNORD28

SNORD50A → SNORD59A

SNORD20 → SNORA75

SNORD108 → SNORD64

SNORD12 → SNORD12B

SNORD47 → SNORD59A

SNORD47 → SNORD60

SNORA70 → SNORD52

SNORD94 → SNORD22



# Validation @0.95

SNORA73B → SNORA73A

SNORD108 → SNORD64

SNORD47 → SNORD59A

SNORD47 → SNORD60

SNORD20 → SNORA75

SNORD12 → SNORD12B

SNORD94 → SNORD22



# Connections between snoRNAs and their respective host genes

<b>snoRNA</b>	<b>Host Gene</b>
SNORD18A	RPL4
SNORD20	NCL
SNORA73A	SNHG3
SNORA14B	TOMM20
SNORA38	PRRC2A
SNORD82	NCL
SNORA78	SNHG9
SNORA48	NFATC3
SNORD68	RPL13
SNORD97	EIF4G2
SNORA75	NCL
SNORD37	EEF2

# JOINED RESULTS



# TOMM20

- Important predictor for both the random forest models
- When compared to the other predictors, its coefficient has experienced the least shortening.

snoRNA	freq. rel.
SNORA14B	0.93
SNORD13	0.86



# SNX5

- Important predictor for the lasso model
- Important predictor for the RF

snoRNA	freq. rel.
SNORD104	0.98
SNORD1C	0.92
SNORD53	0.94
SNORD58C	0.91
<b>TOMM20</b>	0.8584



# SNHG1

- Important predictor for RF

snoRNA	freq. rel.
SNORD104	0.99
SNORD57	0.99
GAS5	0.97



# HSPA8

- Important predictor for RF
- SNORA16A - 0.94
- SNORA44 - 0.97
- SNORA61 - 0.86
- SNORD119 - 0.98
- SNORD13 - 0.88
- SNORD15B - 0.95
- SNORD18A - 0.98
- SNORD1C - 0.87
- SNORD34 - 0.92
- SNORD35B - 0.95
- SNORD37 - 0.88
- SNORD57 - 1
- SNORD68 - 0.88
- SNORD84 - 0.95
- SNORD97 - 0.97



## NFATC3

- Important predictor for RF
- SNORA14B - 0.95
- SNORA16A - 0.99
- SNORA44 - 0.86
- **SNORA48** - 0.98
- SNORA53 - 0.98
- SNORA57 - 0.92
- SNORA63 - 0.98
- SNORA73B - 0.94
- SNORA78 - 0.93
- SNORD10 - 0.89
- SNORD104 - 0.90
- SNORD14A - 0.90
- SNORD22 - 0.90
- SNORD56 - 0.95
- SNORD57 - 0.89
- SNORD97 - 0.97
- **SNX5** - 0.89

# Differential Gene Expression

# ***limma* powers differential expression analyses for RNA-sequencing and microarray studies**

**Matthew E. Ritchie<sup>1,2</sup>, Belinda Phipson<sup>3</sup>, Di Wu<sup>4</sup>, Yifang Hu<sup>5</sup>, Charity W. Law<sup>6</sup>, Wei Shi<sup>5,7</sup> and Gordon K. Smyth<sup>2,5,\*</sup>**

<sup>1</sup>Molecular Medicine Division, The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville, Victoria 3052, Australia, <sup>2</sup>Department of Mathematics and Statistics, The University of Melbourne, Parkville, Victoria 3010, Australia, <sup>3</sup>Murdoch Childrens Research Institute, Royal Children's Hospital, 50 Flemington Road, Parkville, Victoria 3052, Australia, <sup>4</sup>Department of Statistics, Harvard University, 1 Oxford Street, Cambridge, MA 02138-2901, USA, <sup>5</sup>Bioinformatics Division, The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville, Victoria 3052, Australia, <sup>6</sup>Institute of Molecular Life Sciences, University of Zurich, Winterthurerstrasse 190, Zurich 8057, Switzerland and <sup>7</sup>Department of Computing and Information Systems, The University of Melbourne, Parkville, Victoria 3010, Australia

Received November 09, 2014; Revised January 04, 2015; Accepted January 06, 2015

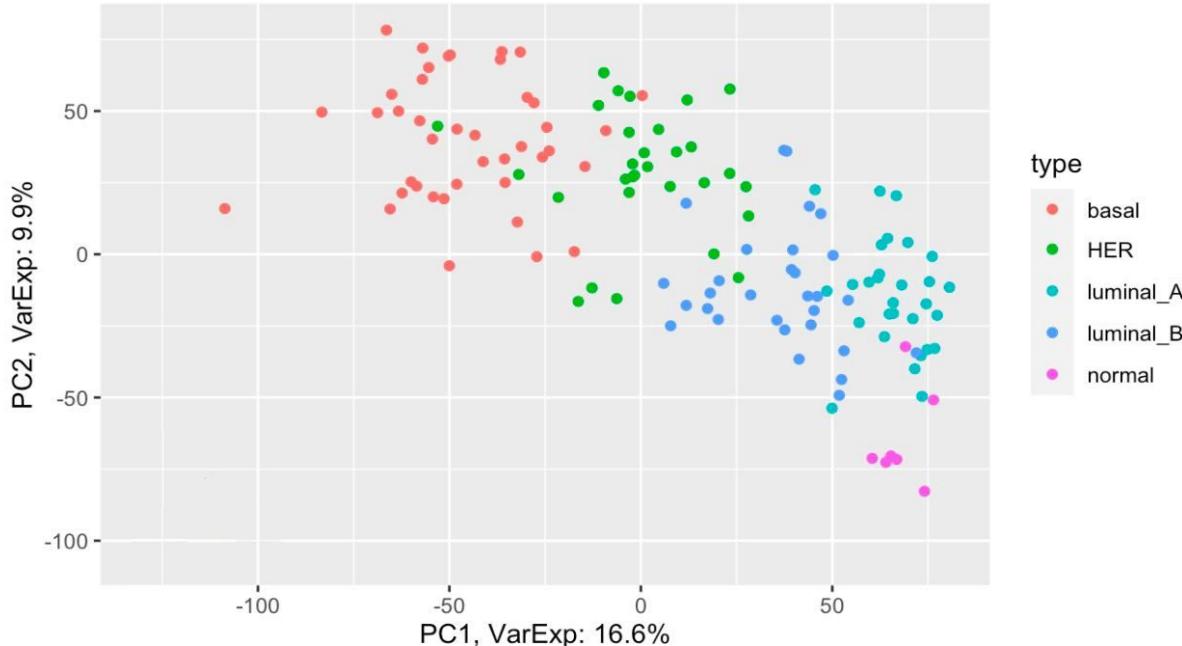
<https://doi.org/10.1093/nar/gkv007>

<https://bioconductor.org/packages/release/bioc/html/limma.html>



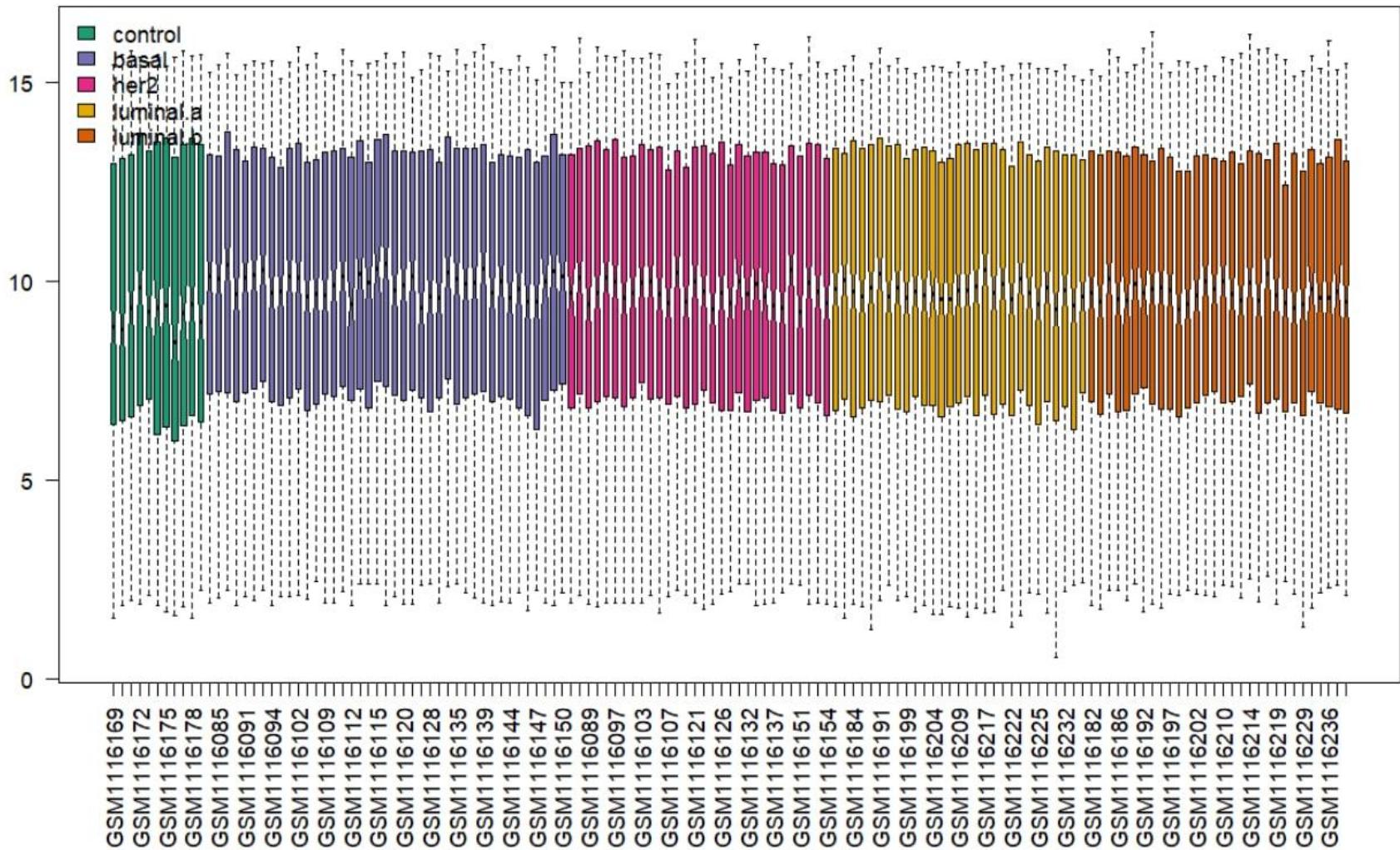
# Exploratory Data Analysis

PCA plot of the log-transformed raw expression data



PCA without cell lines

# GSE45827/GPL570



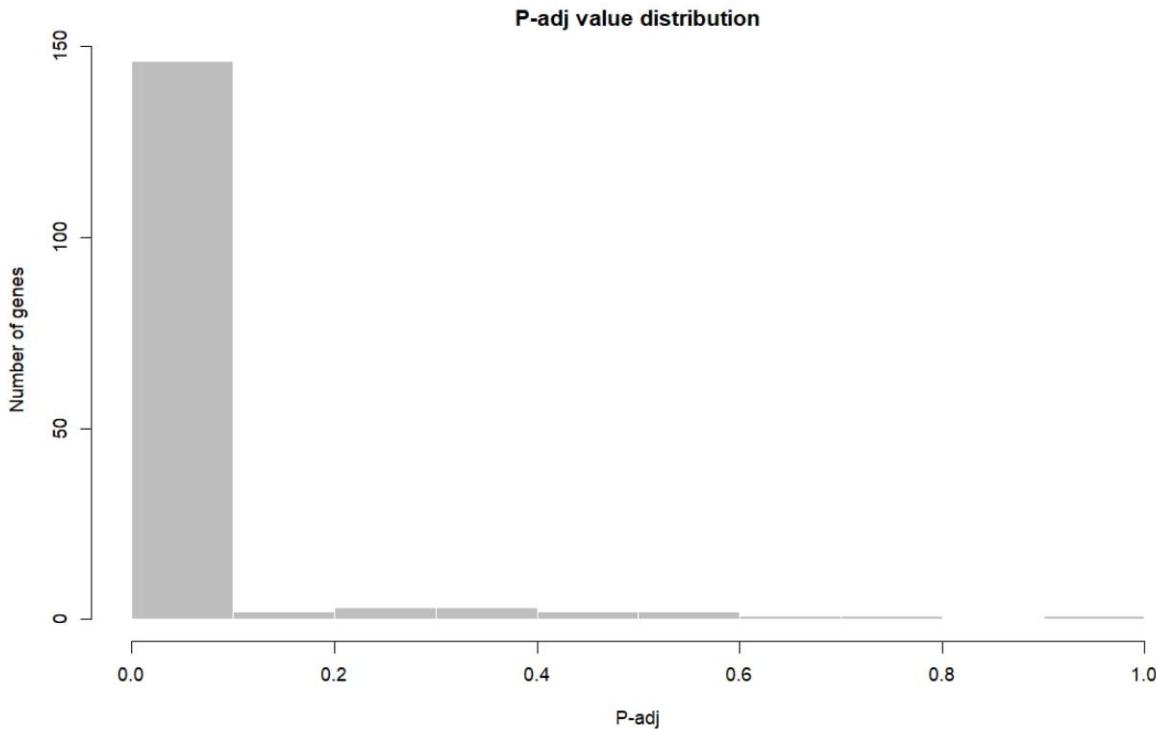


# Estimating the proportion of true null hypotheses

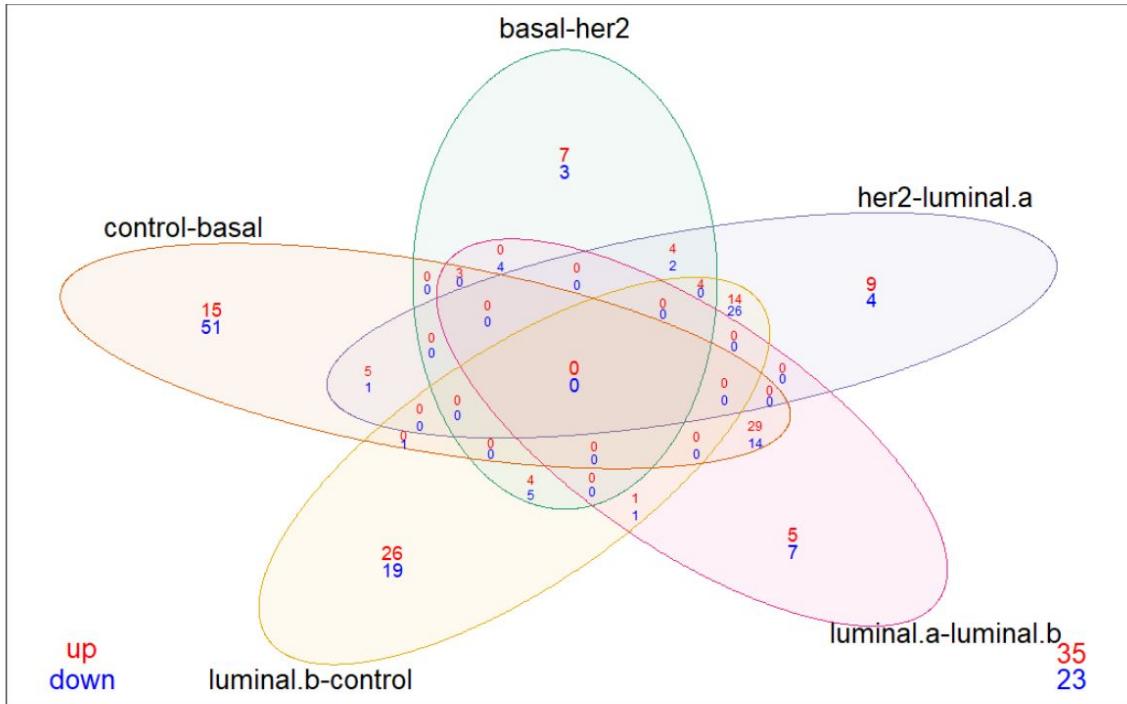
We build a histogram of P-values for all genes.

The null hypothesis is that most genes are not differentially expressed.

We reject the null hypothesis.

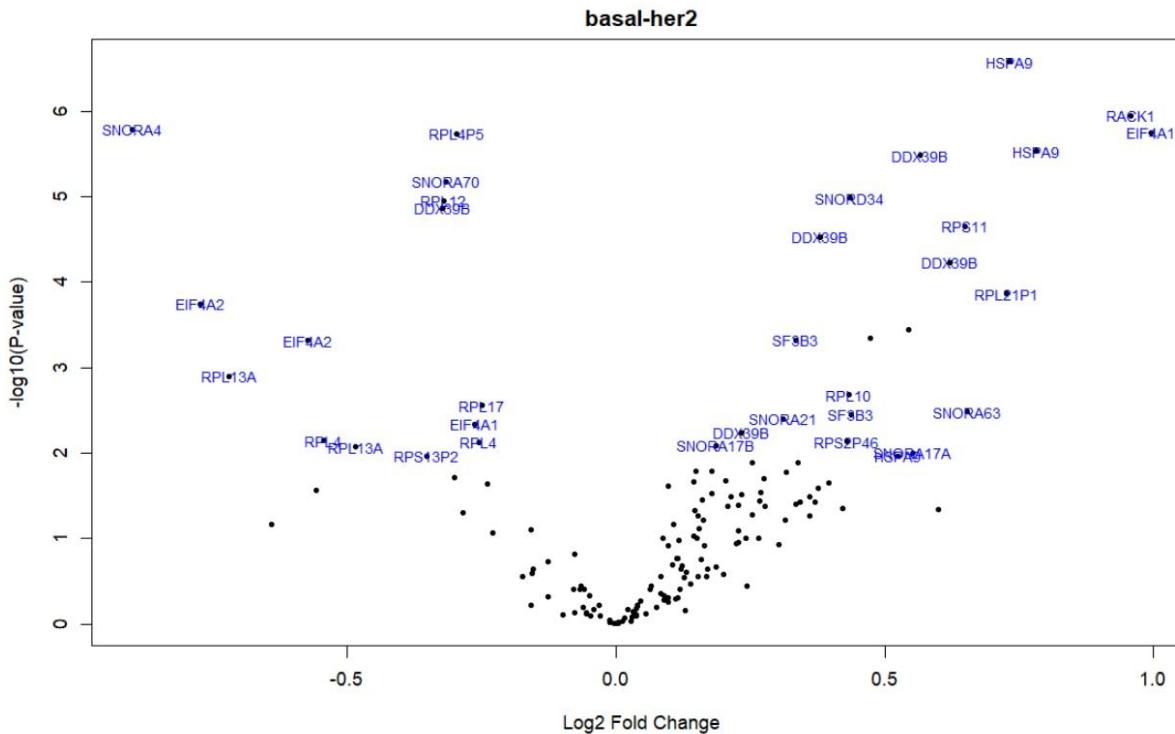


# Venn diagram with both up and down-regulated genes

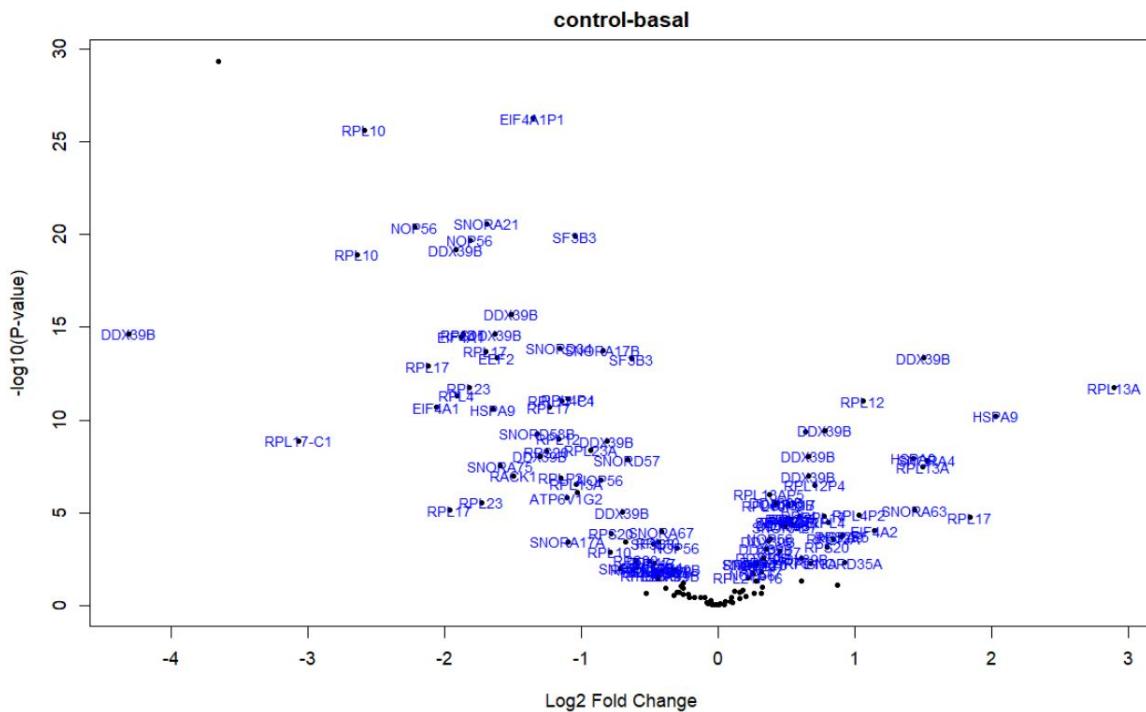


Venn diagram showing overlap in the number of DE genes for the 5 comparisons, generated by the *vennDiagram* function.

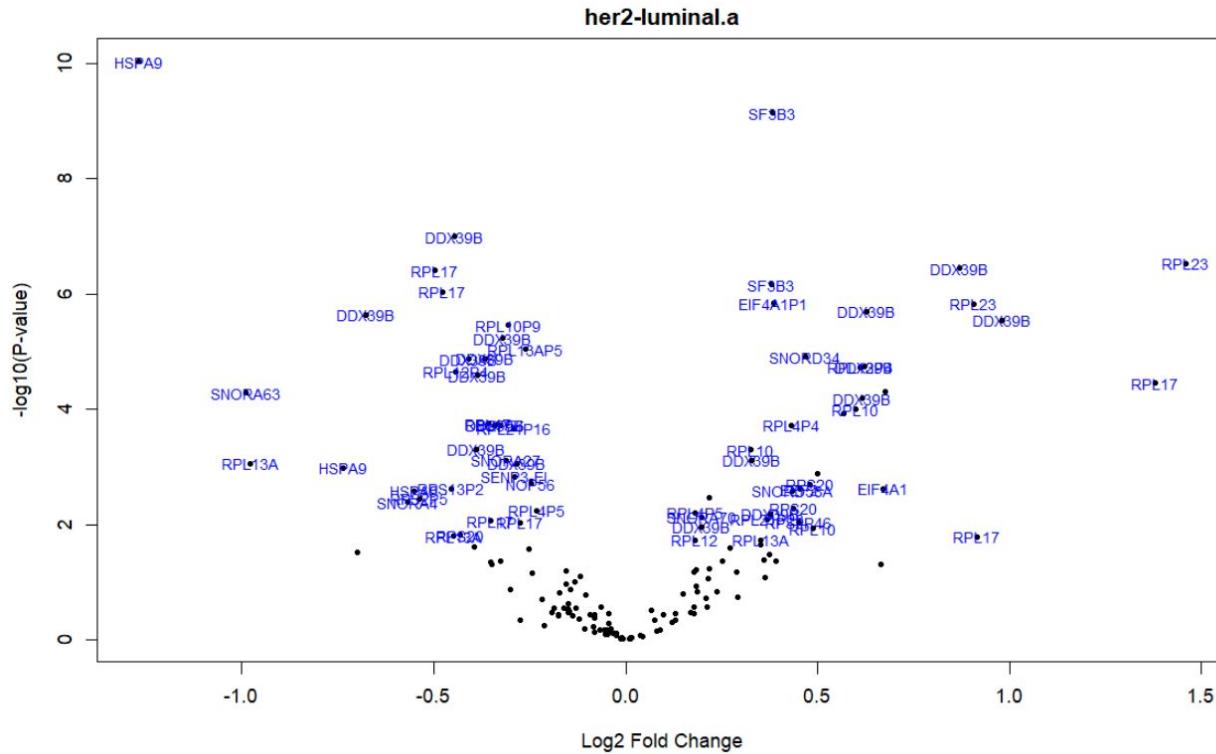
# Volcano plot (1/5) - Basal VS HER2



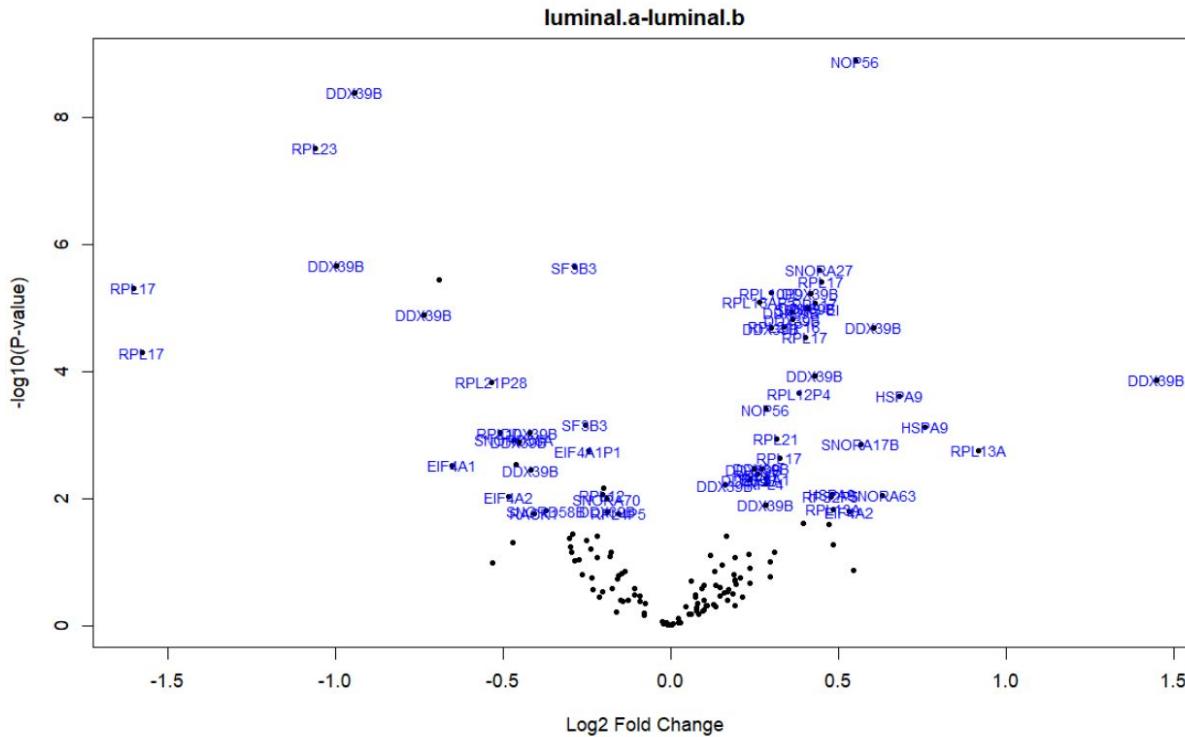
# Volcano plot (2/5) - Control VS Basal



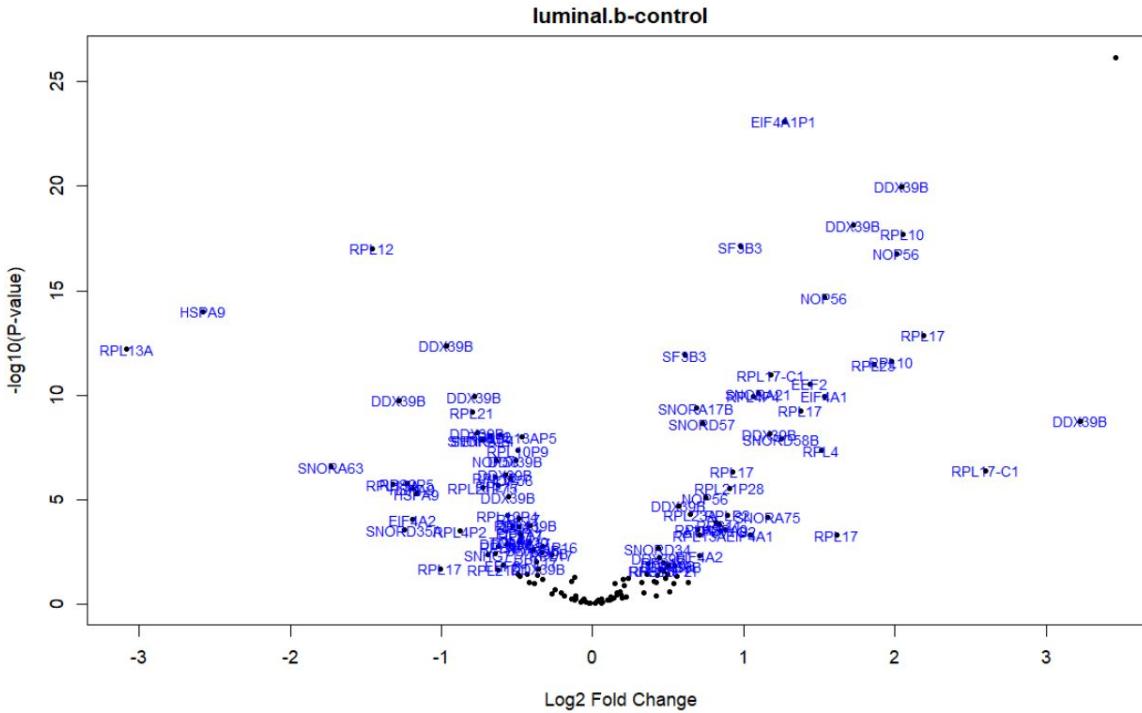
# Volcano plot (3/5) - HER2 VS Luminal A



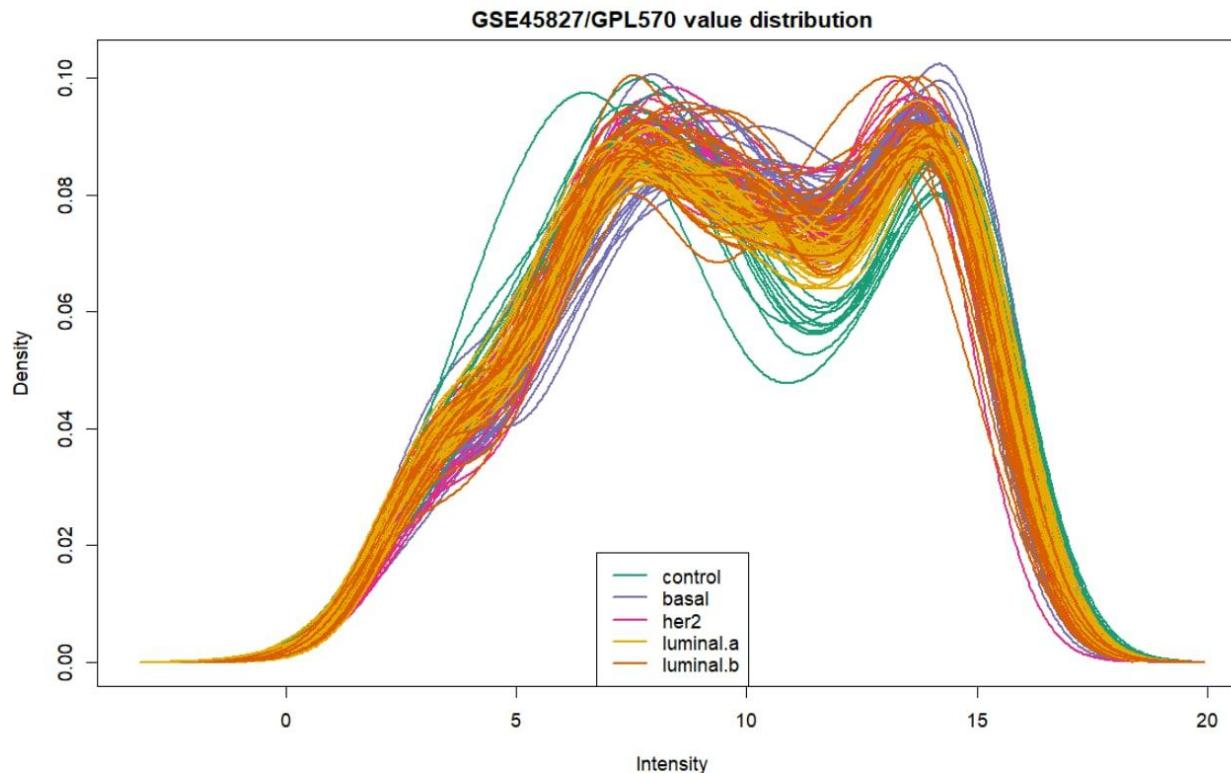
# Volcano plot (4/5) - Luminal A VS Luminal B



# Volcano plot (5/5) - Luminal B VS Control

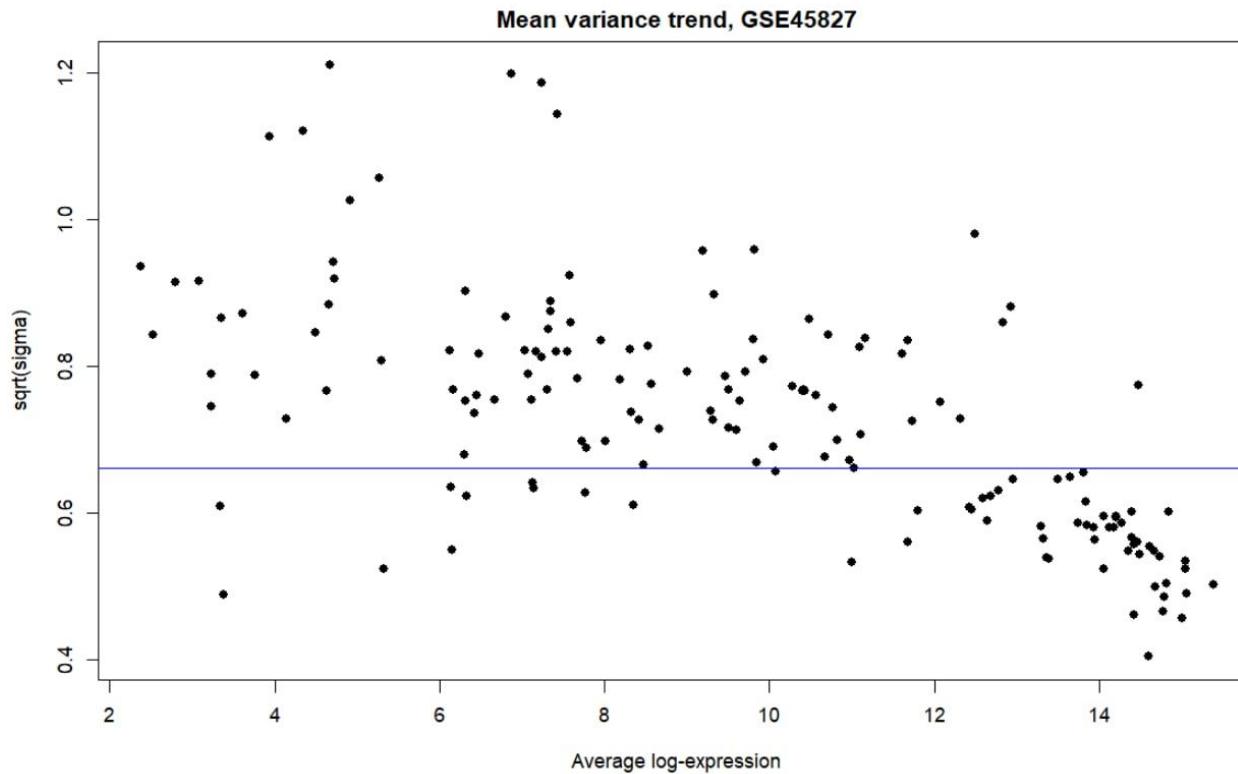


# Expression value distribution





# Mean-variance trend

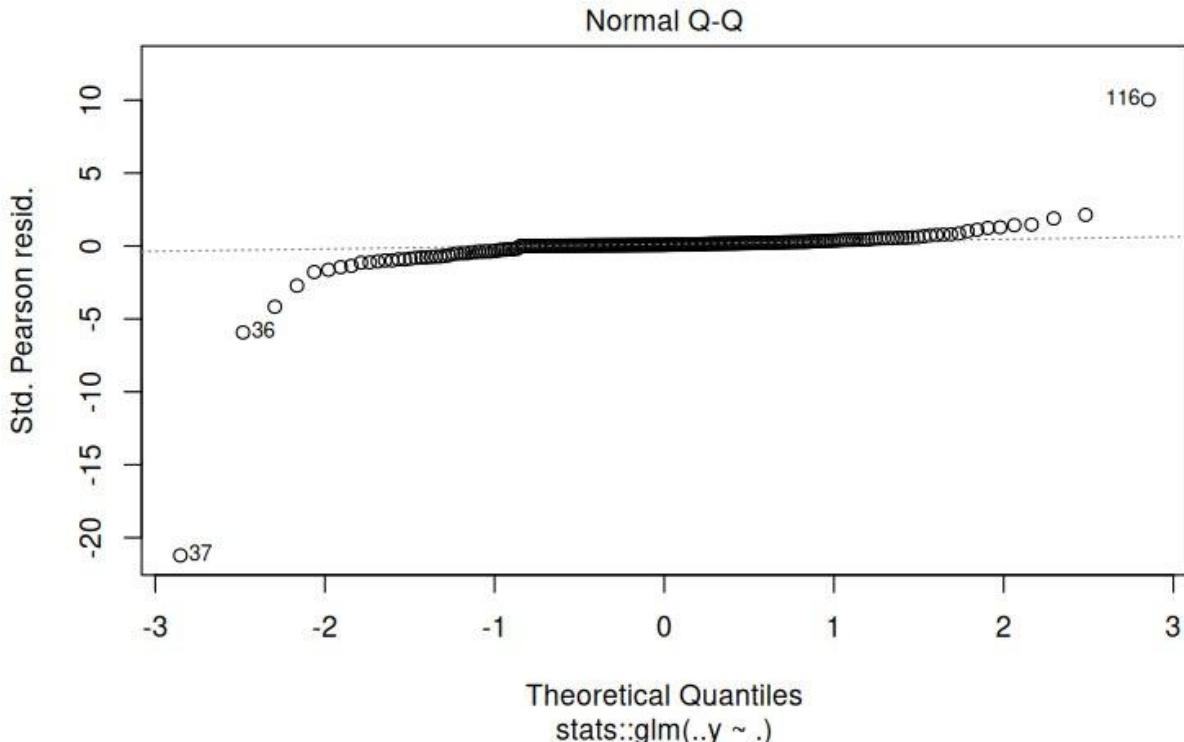


# Thanks for your attention

# Appendix



# Normal QQ plot

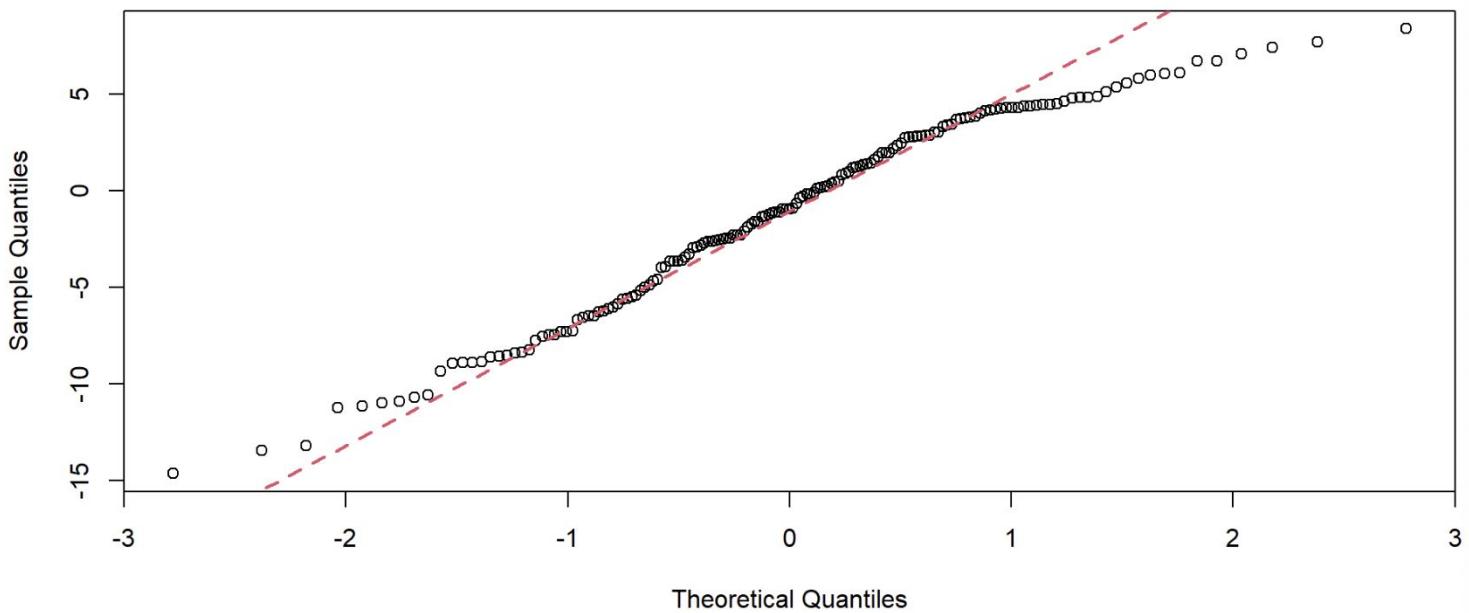


# Differential Gene Expression Analysis



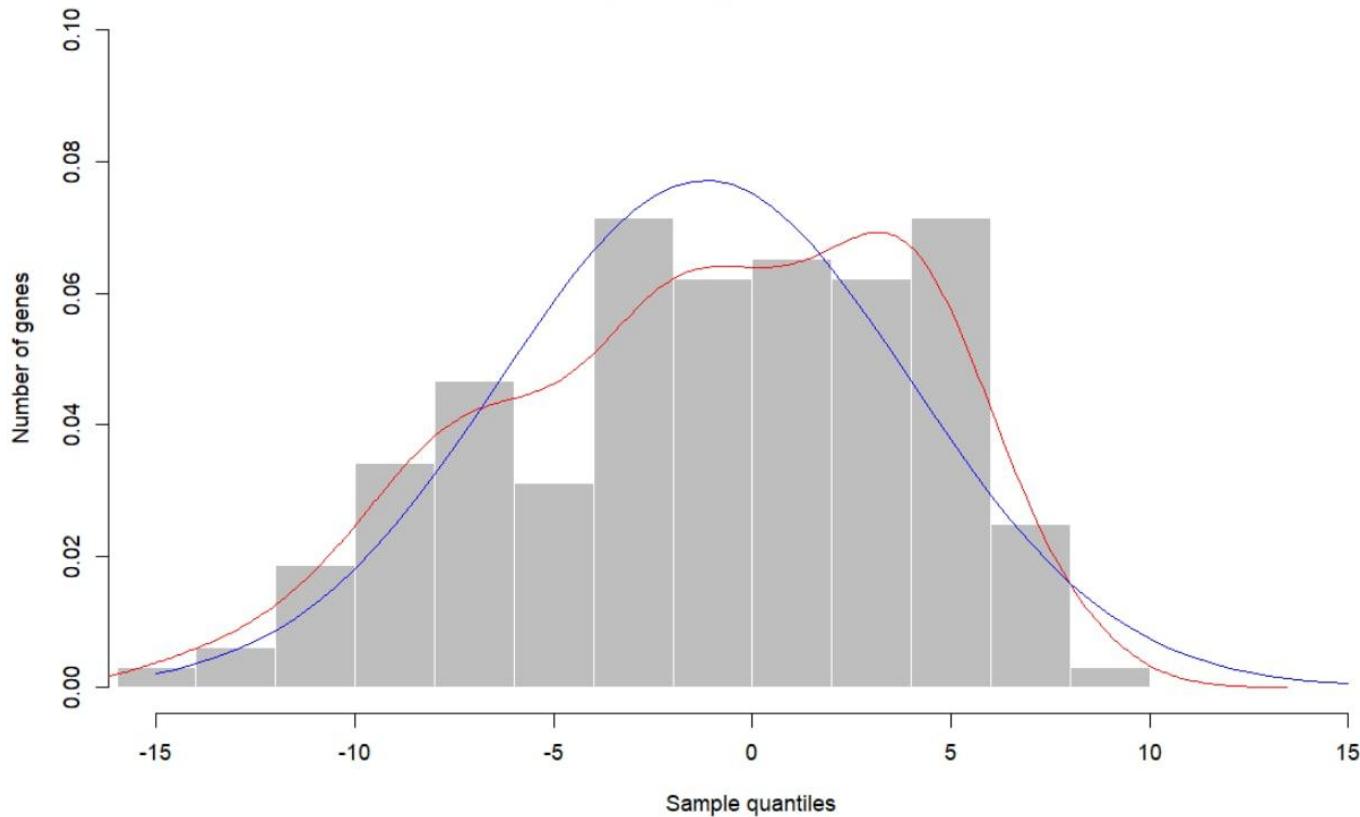
# Q-Q plot for t-statistic

Moderated t statistic



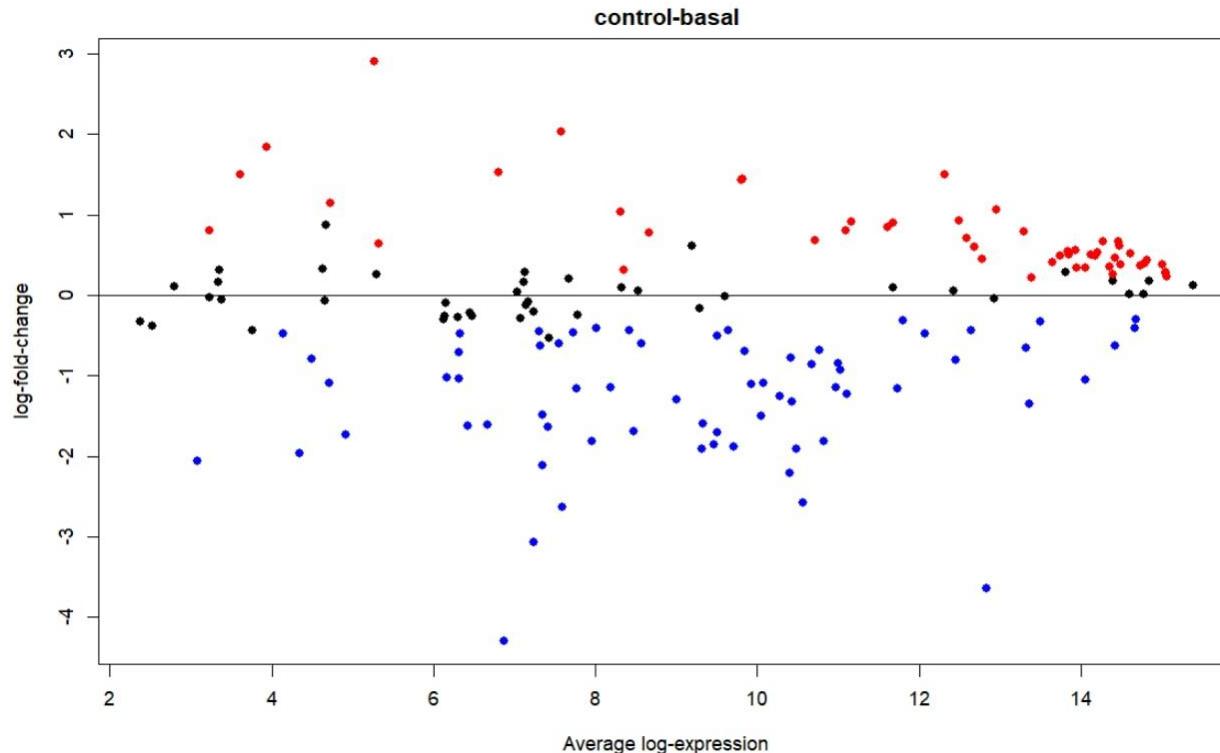


Sample quantiles probability distribution



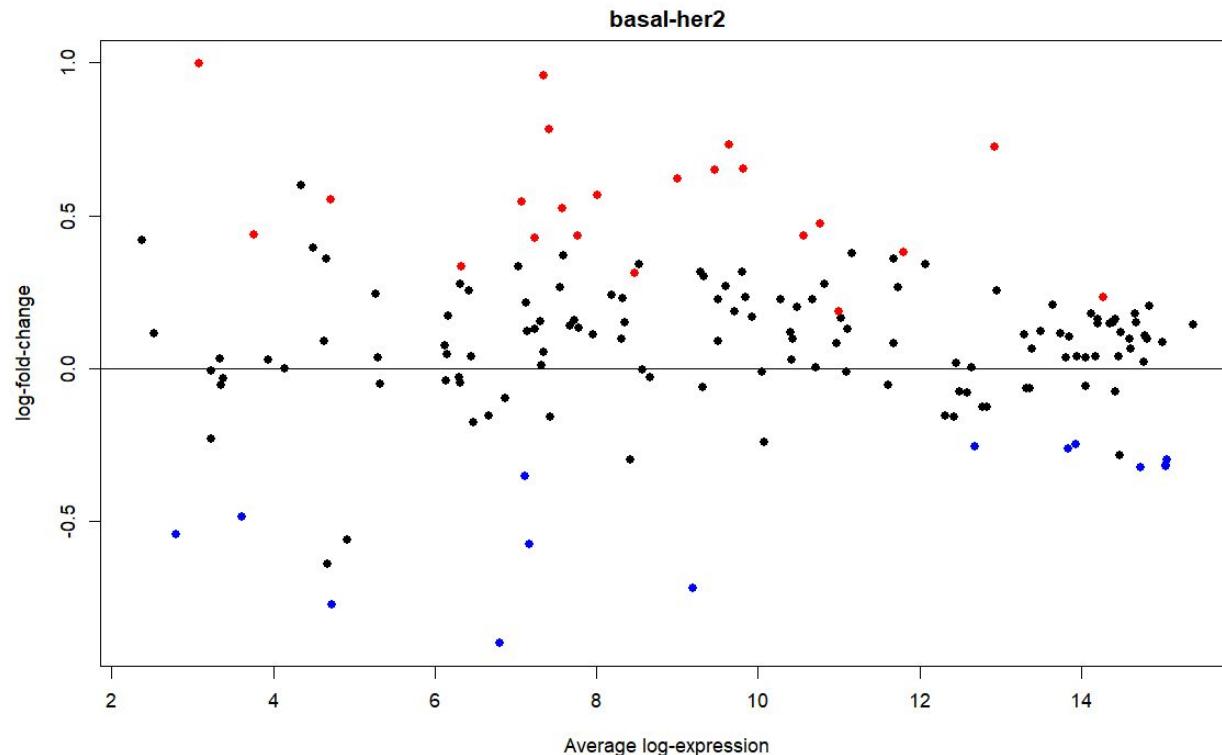


# Mean-Difference plot (1/5)



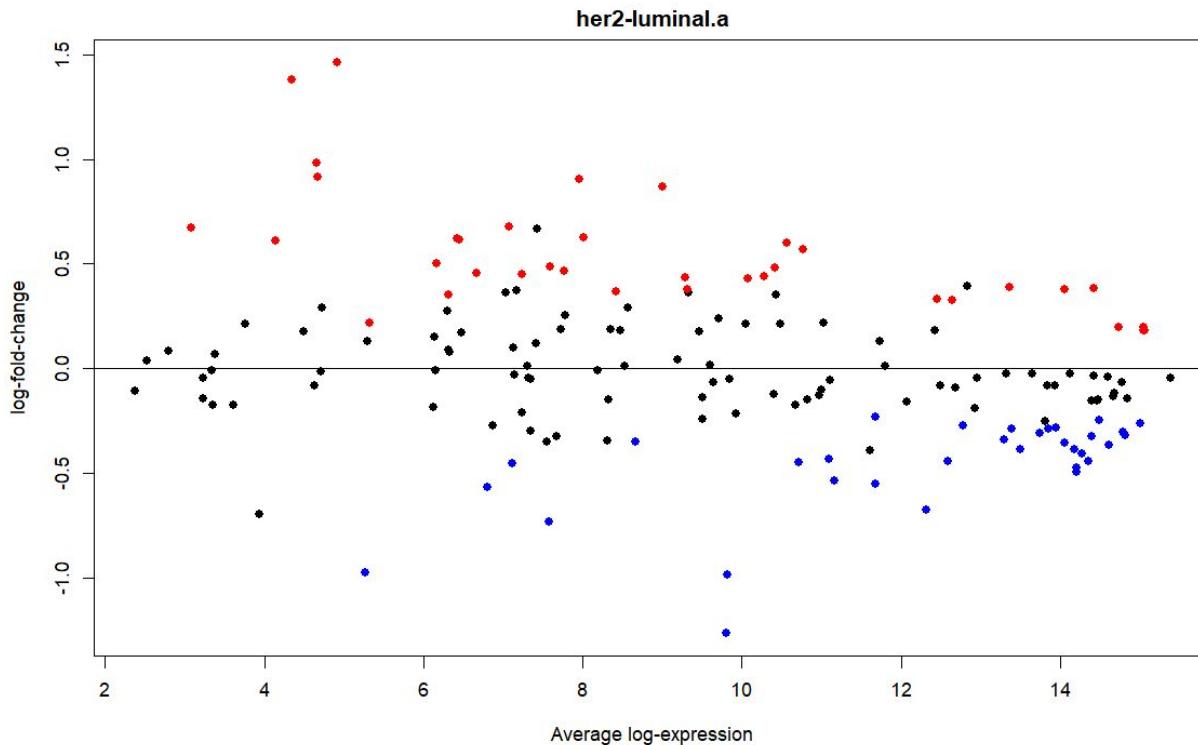


# Mean-Difference plot (2/5)



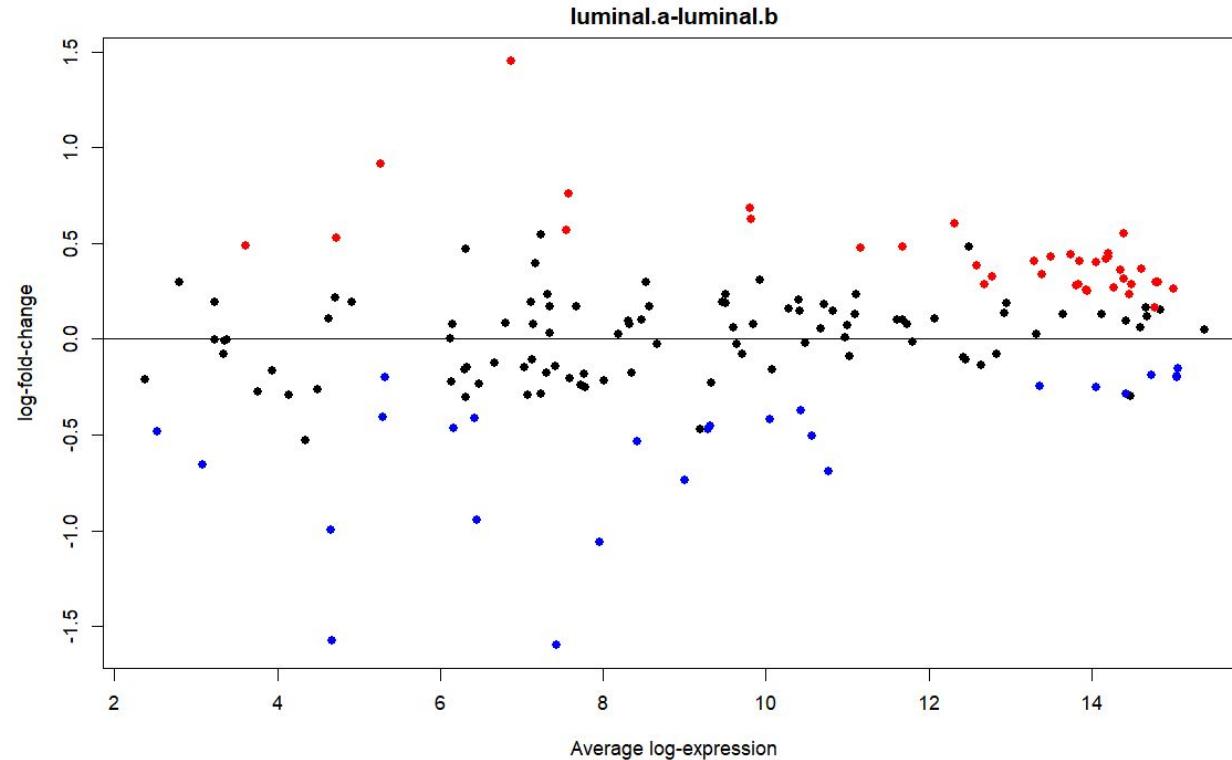


# Mean-Difference plot (3/5)



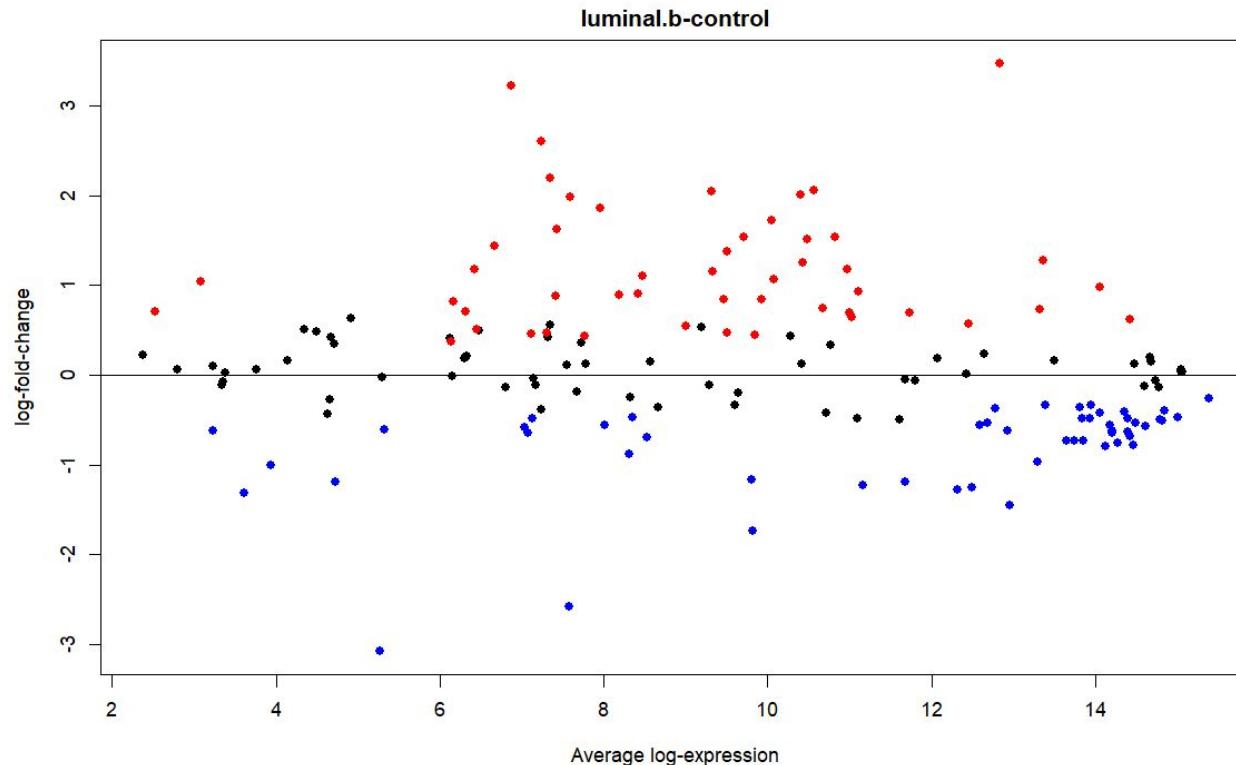


# Mean-Difference plot (4/5)



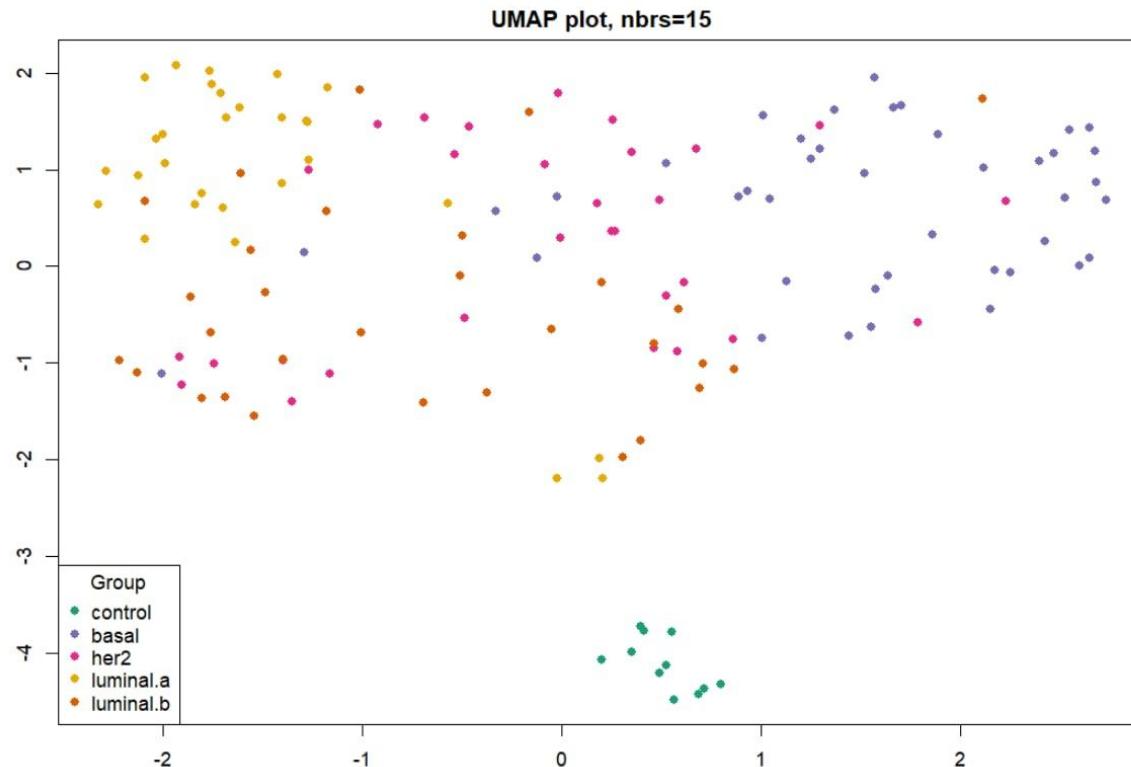


# Mean-Difference plot (5/5)

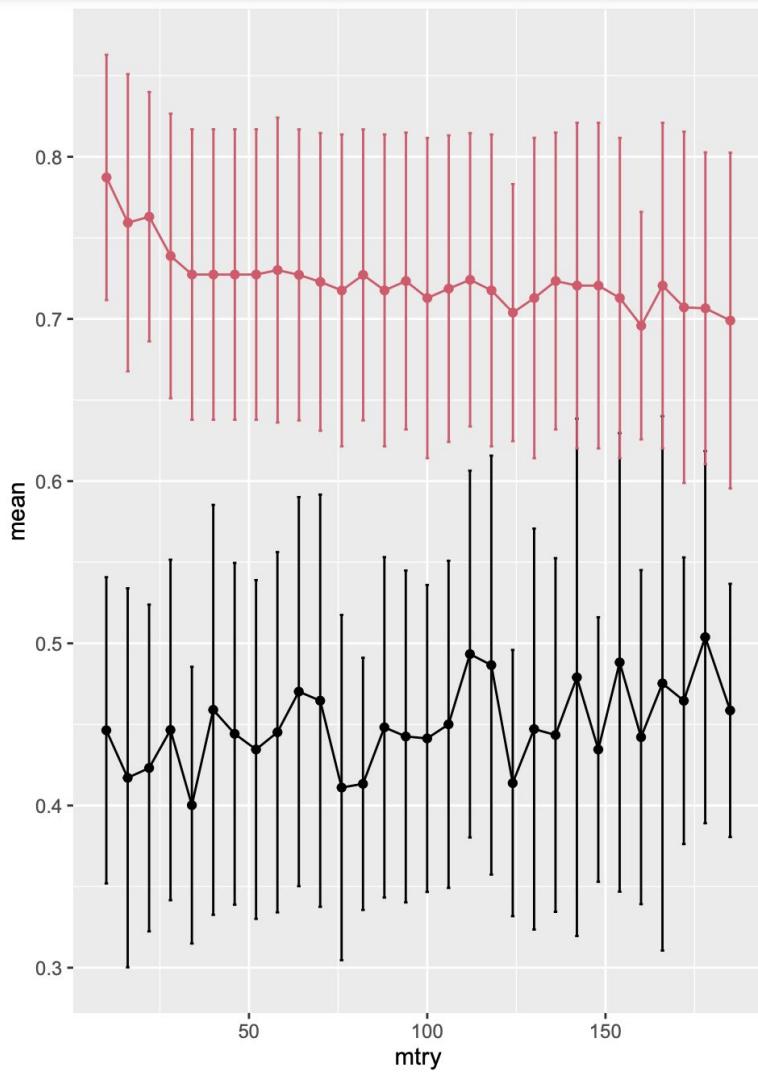




# UMAP



# Random Forest



## Comparison random forest performance

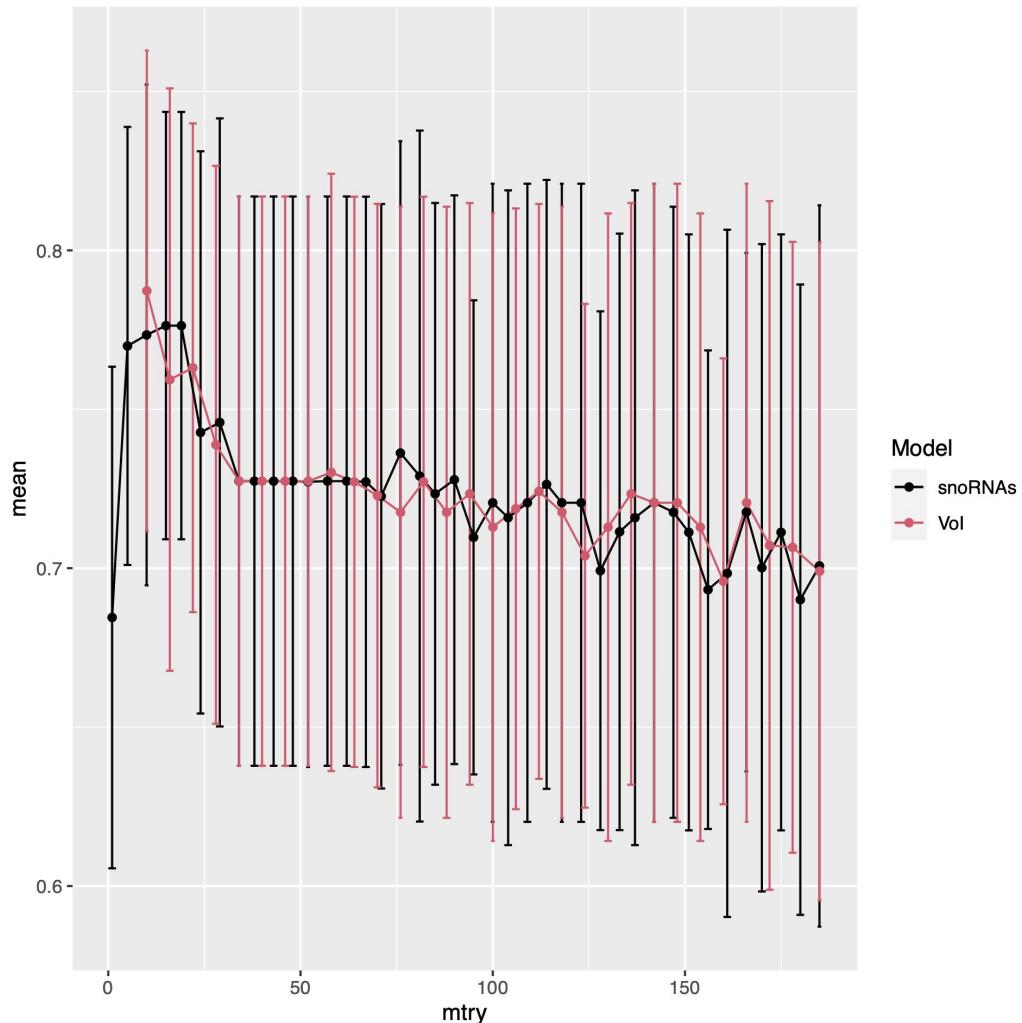
Model

- Low Variance Predictors
- Vol

Vol = Variables of Interest



# RF snoRNAs vs variables of interest





## Random forest fit on our genes of interest (old data)

		Truth	
		No Cancer	Cancer
Predicted	No Cancer	12	0
	Cancer	3	62



<b>mtry*</b>	<b>Measure</b>	<b>Mean</b>	<b>N. folds</b>	<b>Std. error</b>
10	F1	0.7872726	5	0.03857438
22	F1	0.7630541	5	0.03922816
16	F1	0.7593441	5	0.04674310
28	F1	0.7388160	5	0.04477342
58	F1	0.7301482	5	0.04794219

\*mtry is the number of predictors used for the splits



		Truth				
		Basal	HER	L. A	L. B	Normal
Predicted	Basal	9	2	0	1	0
	HER	2	5	0	1	0
	L. A	0	1	6	2	0
	L. B	0	0	1	4	0
	Normal	0	0	0	0	1

<b>mtry*</b>	<b>Measure</b>	<b>Mean</b>	<b>N. folds</b>	<b>Std. error</b>
178	F1	0.5038095	5	0.05853361
112	F1	0.4933707	5	0.05767972
154	F1	0.4882505	5	0.07213041
118	F1	0.4865629	5	0.06589169
142	F1	0.4790185	5	0.08137526

\*mtry is the number of predictors used for the splits



# Random forest confusion matrix on lesser variant genes

		Truth				
		Basal	HER	L. A	L. B	Normal
Predicted	Basal	7	6	2	2	0
	HER	2	0	1	1	0
	L. A	1	1	1	1	0
	L. B	1	1	3	4	1
	Normal	0	0	0	0	0

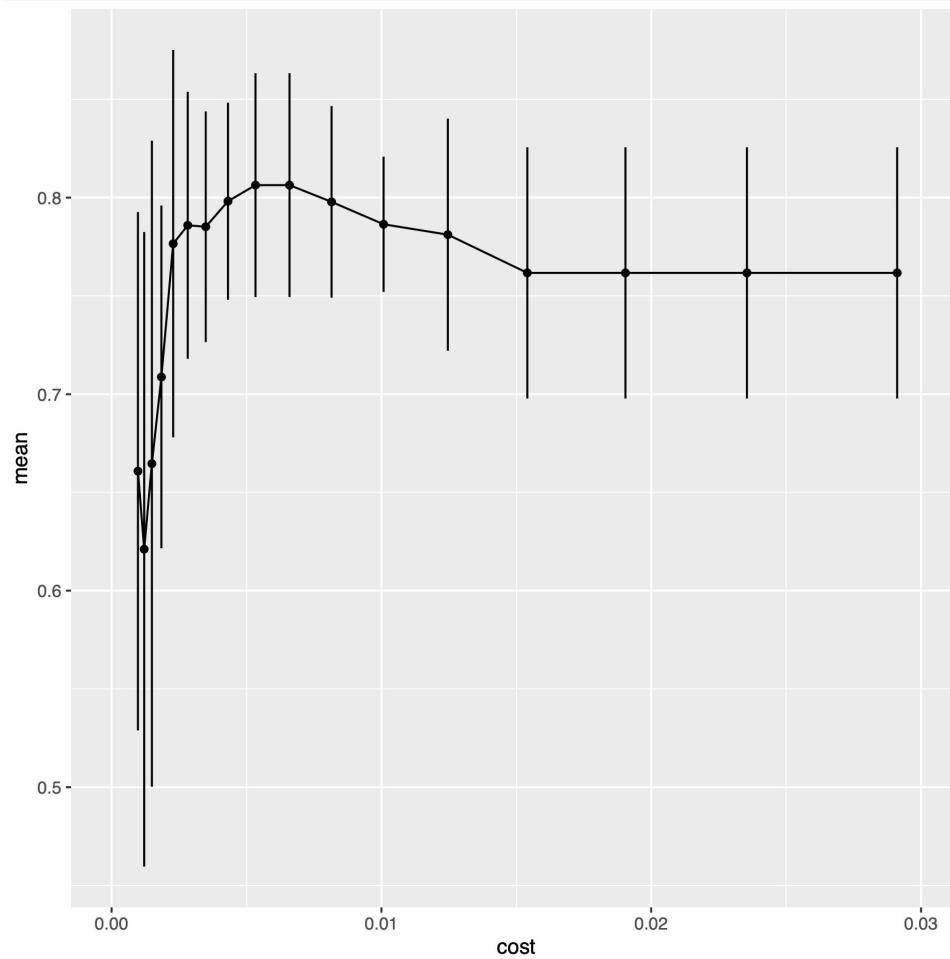


<b>cost</b>	<b>Measure</b>	<b>Mean</b>	<b>N. folds</b>	<b>Std. error</b>
0.005332322	F1	0.8063554	5	0.02904139
0.006592777	F1	0.8063554	5	0.02904139
0.004312850	F1	0.7981699	5	0.02556259
0.008151179	F1	0.7978359	5	0.02487277
0.010077956	F1	0.7864484	5	0.01755194

# SVM



# SVM cost





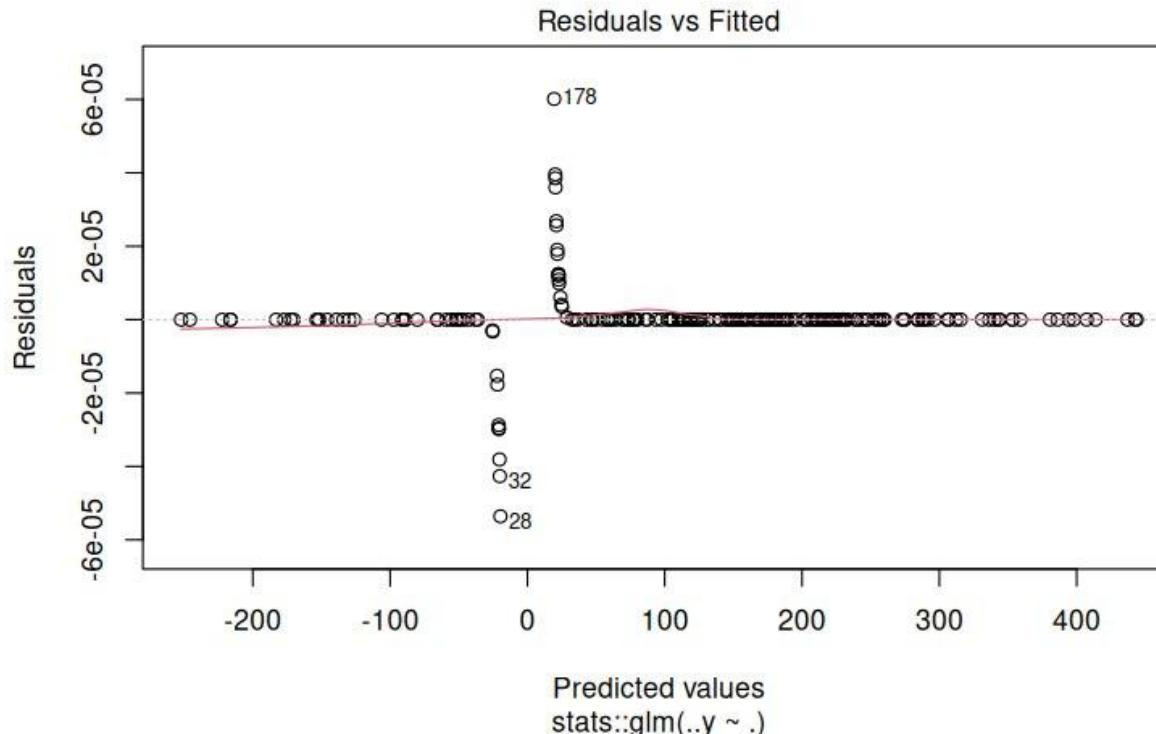
# Final SVM model confusion matrix

Predicted	Truth				
	Basal	HER	L. A	L. B	Normal
Basal	10	2	0	1	0
HER	1	6	0	1	0
L. A	0	0	7	2	0
L. B	0	0	0	4	0
Normal	0	0	0	0	1

# Logistic Regression

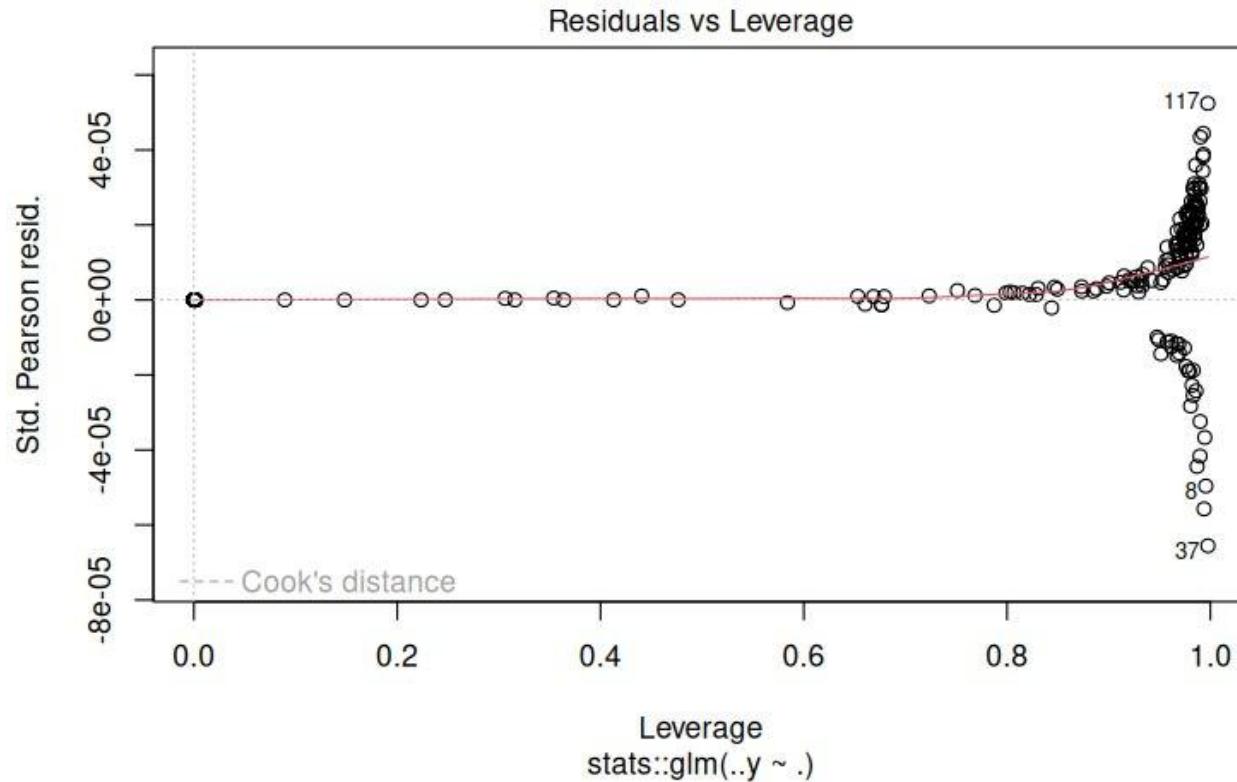


# Residuals vs fitted old



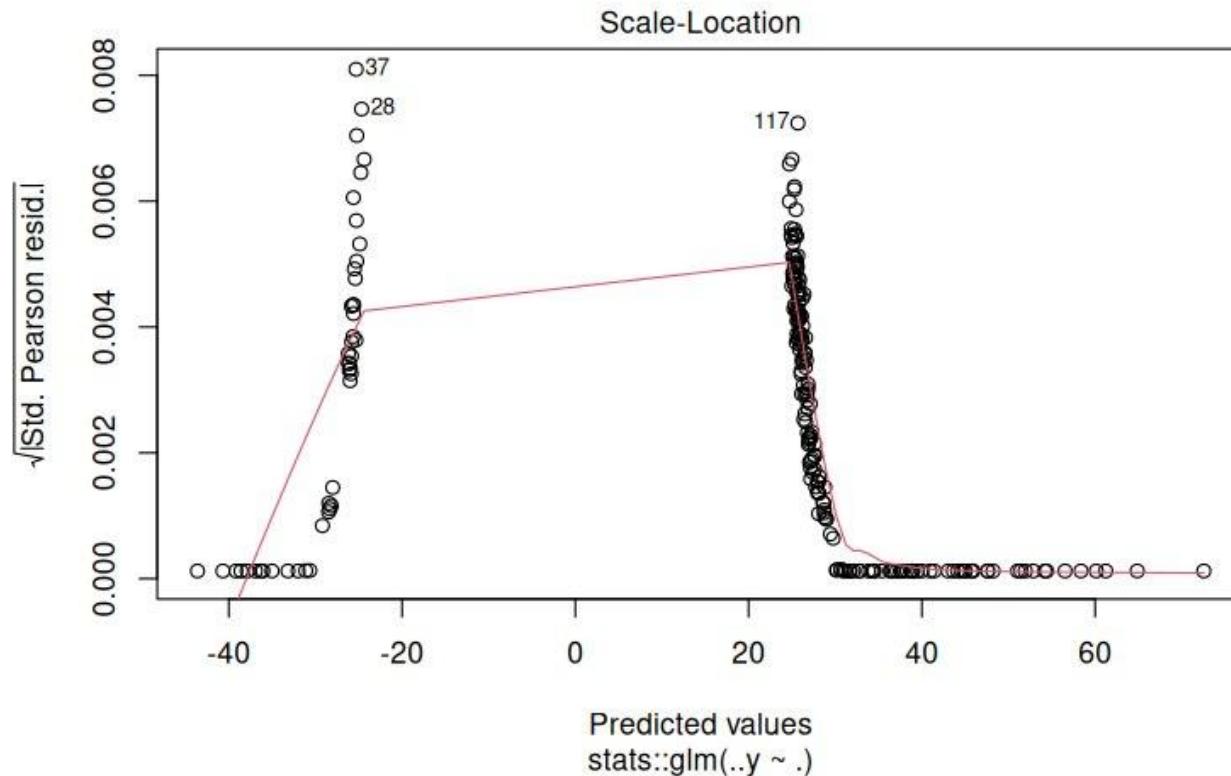


# Reslev old





# Scaleloc old





## Logistic regression with L1 penalty crossvalidation results

penalty	Measure	Mean	N. folds	Std. error
6.250552e-01	F1	0.8195279	5	0.04381708
1.264855e-05	F1	0.8137864	5	0.03345629
1.206793e-06	F1	0.8099060	5	0.04936689
1.757511e-08	F1	0.8039975	5	0.04661561
6.866488e-09	F1	0.8030488	5	0.04050916



# Logistic regression with L1 penalty confusion matrix

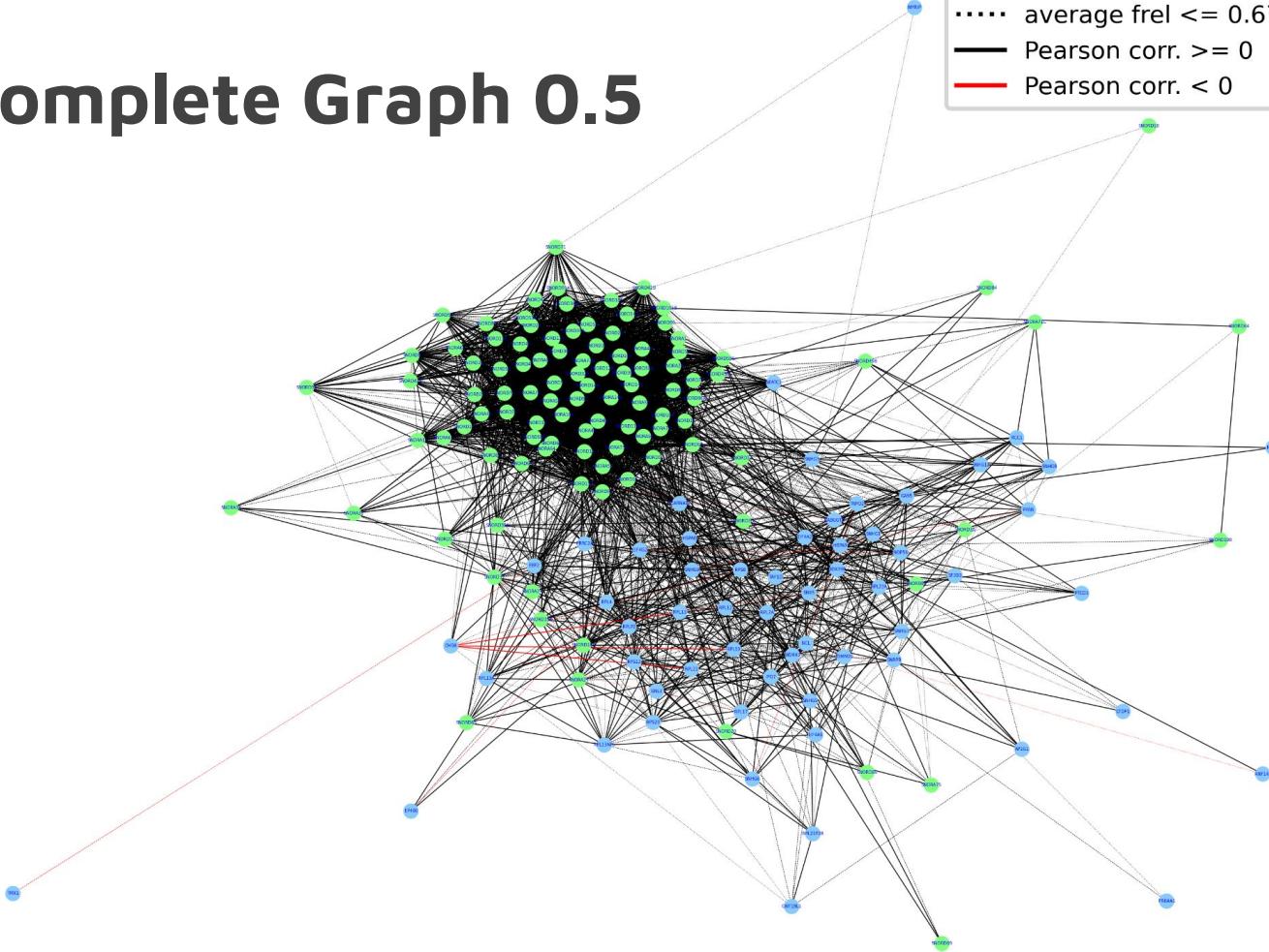
		Truth				
		Basal	HER	L. A	L. B	Normal
Predicted	Basal	9	1	0	1	0
	HER	2	5	0	0	0
	L. A	0	1	6	3	0
	L. B	0	1	1	4	0
	Normal	0	0	0	0	1

# Network analysis



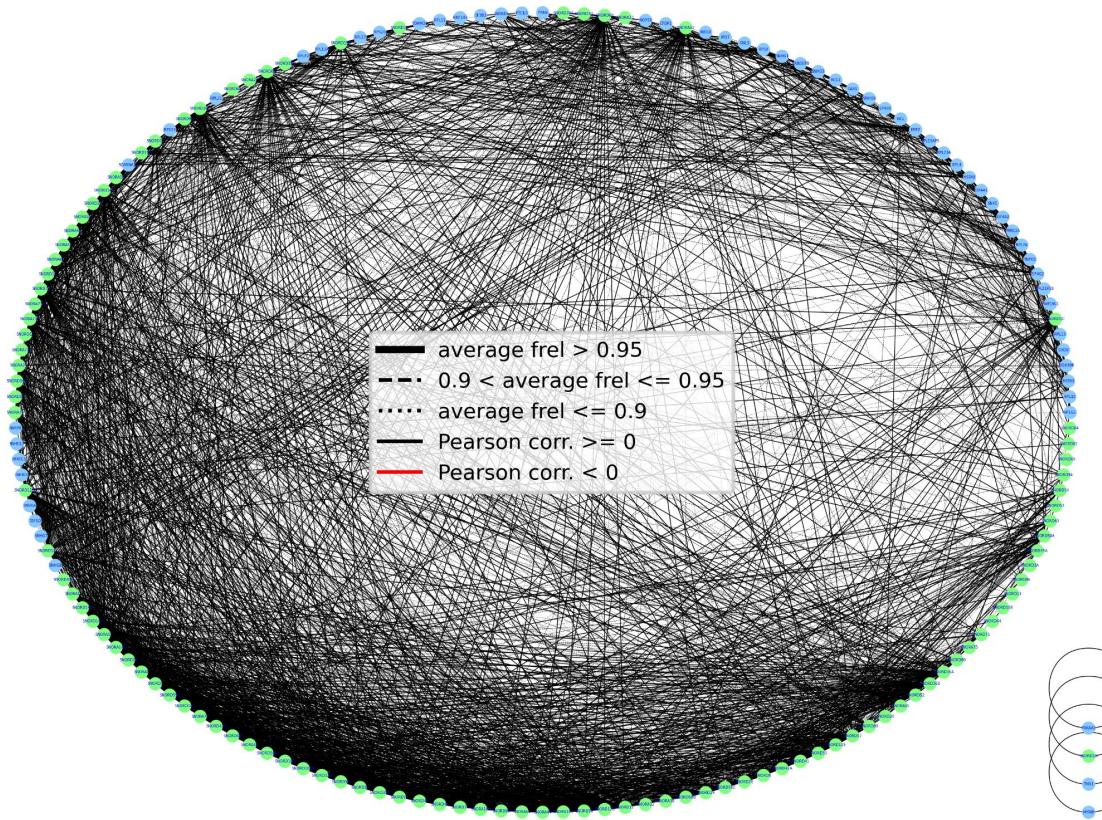
# Complete Graph 0.5

- average frel > 0.83
- - - 0.67 < average frel <= 0.83
- .... average frel <= 0.67
- Pearson corr. >= 0
- Pearson corr. < 0

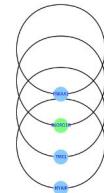
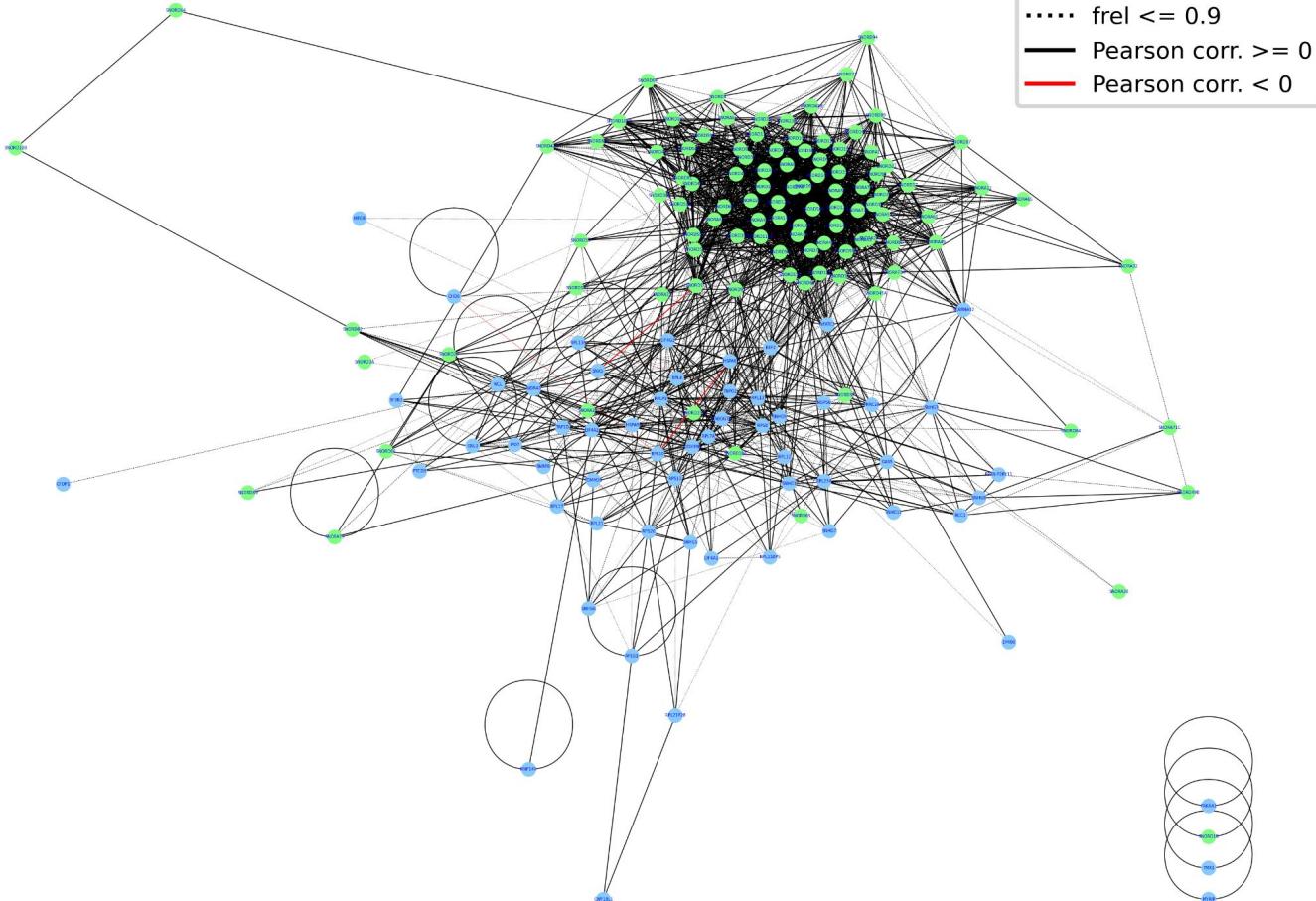




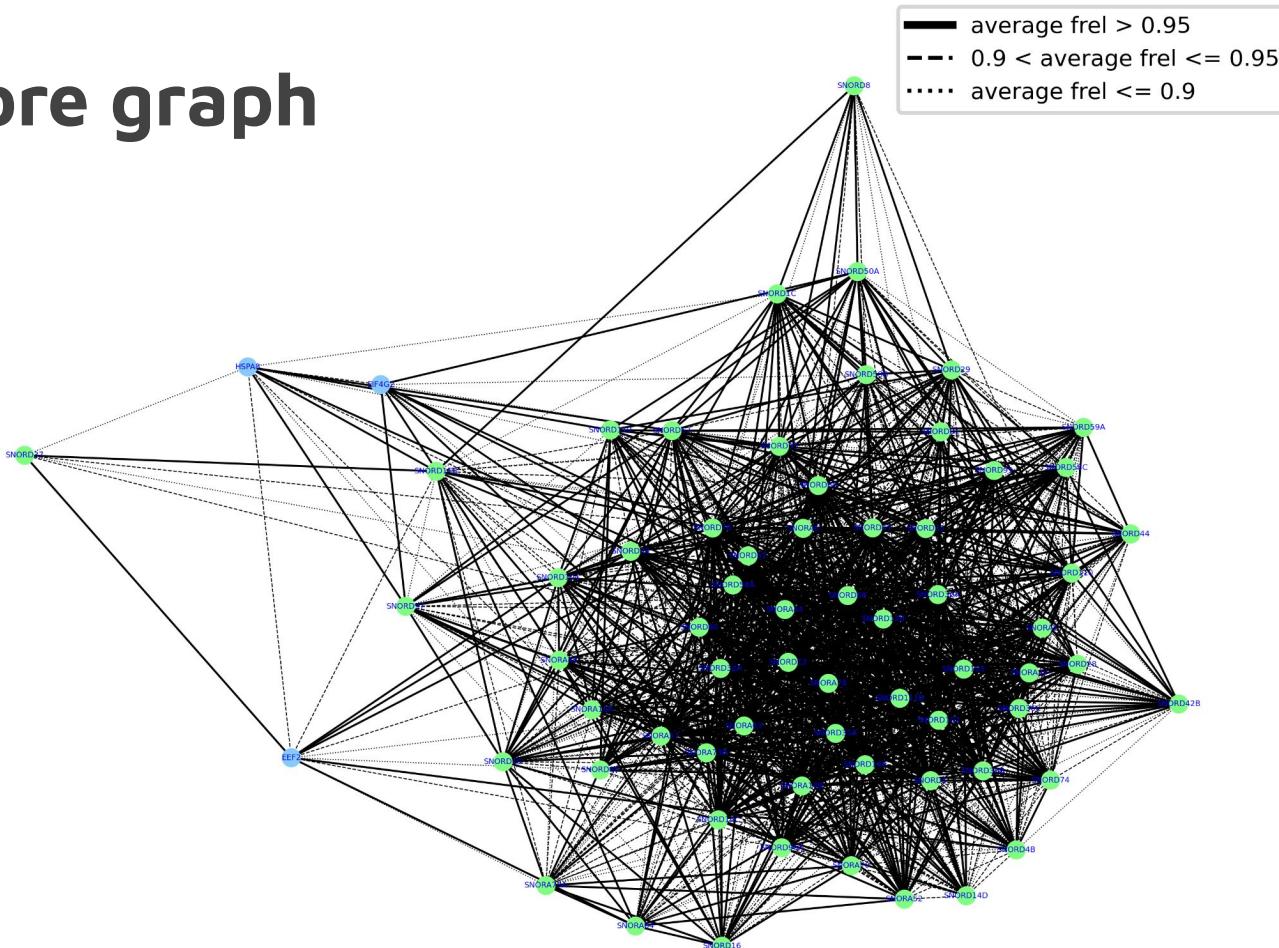
# Complete Graph Circular



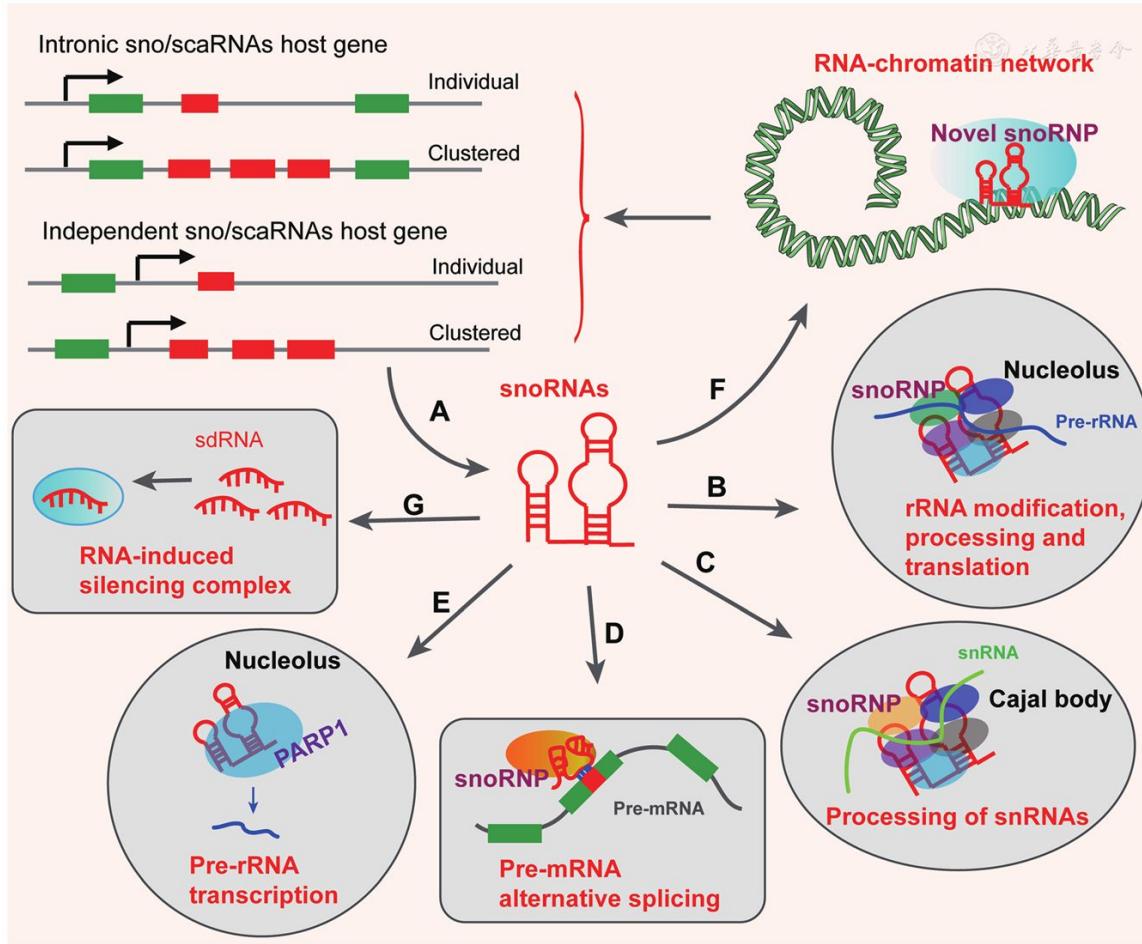
# Complete Graph not ignored



# Core graph



# Theory part





# Breast cancer subtypes

**Luminal A breast cancer** is estrogen receptor-positive and progesterone receptor-positive, HER2-negative, and has low levels of the protein Ki-67, which helps control how fast cancer cells grow. Luminal A cancers tend to grow more slowly than other cancers, be lower grade, and have a good prognosis.

**Luminal B breast cancer** is estrogen receptor-positive and HER2-negative, and also has either high levels of Ki-67 (which indicate faster growth of cancer cells) or is progesterone receptor-negative.



# Breast cancer subtypes

**HER2-enriched breast cancer** is estrogen receptor-negative and progesterone receptor-negative and HER2-positive. HER2-enriched cancers tend to grow faster than luminal cancers and can have a worse prognosis, but are usually successfully treated with targeted therapy medicines aimed at the HER2 protein.

**Triple-negative or basal-like breast cancer** is estrogen receptor-negative, progesterone receptor-negative, and HER2-negative.

Triple-negative breast cancer is considered more aggressive than either luminal A or luminal B breast cancer.

Font: <https://www.breastcancer.org/types/molecular-subtypes>