

InSyBio Pipelines

March 2022

Insybio Suite v3.0



InSyBio

Intelligent Systems Biology

User Manual

www.insybio.com

Insybio Pipelines

Introduction

Pipelines is a tool that provides an integrated optimized pipeline from raw data until the discovery of biosignatures, networks and predictive models. It takes as input a dataset file and after performing statistical and network analytics, it uses the dataset to train a machine learning model identifying at the same time the optimal biomarkers for the trained model.

To start the calculation:

Click in the menu “InSyBio Pipelines”, select the “Add new job” button and then:

- Upload a new Biomarkers Dataset file with features as rows and samples as columns and a corresponding Biomarkers Labels file with all labels in one row, separated by tabs. You are redirected to the Data Store where step by step instructions guide you.
- Or Select a file from the Data Store. There you can find your previously uploaded files or InSyBio pre-uploaded sample datasets.

The screenshot displays the 'InSyBio Suite - Pipeline Job' interface. On the left is a navigation menu with options: InSyBio Interact, InSyBio ncRNASeq, InSyBio Bionets, InSyBio Biomarkers, InSyBio DNA-Seq, InSyBio Pipelines, and InSyBio DataStore. The main content area is titled 'Main Inputs - Upload Files' and contains two sections for file uploads. The first section, 'Biomarkers Dataset', has a 'Title 1' field with 'healthy' and a 'Filename 1' field with 'dsfile1622716077_184.txt'. Below these are two buttons: 'Select file from Data Store' and 'Go to Data Store to Upload File'. The second section, 'Phenotypic Annotation', has a 'Title 2' field with 'softparse healthy labels' and a 'Filename 2' field with 'dsfile1647516535_7558.txt', also with 'Select file from Data Store' and 'Go to Data Store to Upload File' buttons. At the bottom, there are three questions with dropdown menus or checkboxes: 'Does your dataset have samples headers?' (No), 'Does your dataset have features headers?' (No), and 'Does the Normalization use a set of householding molecules?' (checkbox).

After that the user will have to insert the information regarding the headers. That is to inform the application if the original dataset has sample headers or feature headers. Optionally, the user will have the option to insert the names of the features that will be used for normalization.

The user should also select which Pipeline Steps he wants to be performed. There are six steps. All steps have pre-optimized parameters and steps for general purpose and pre-optimized parameters and steps for specific application domains that we use as defaults, but the user can configure these values manually. It is advised that the default values are used.

Pipeline Steps

Select preoptimized parameters:

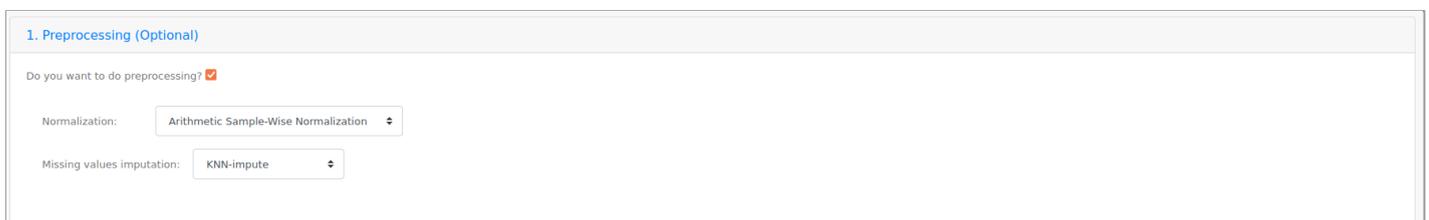
Default preoptimized parameters

1. Preprocessing (Optional)
2. Statistical Analysis
3. General Predictive Model Creation
4. Network Analysis (Optional)
5. ncRNAseq Predict (Optional)
6. Testing Multi-biomarker Predictive Analytics Model (Optional)

Submit job

Preprocessing

The first step is the Preprocessing. The user can select if preprocessing should be performed. During preprocessing we filter the dataset, perform normalization, missing values imputation, duplicate measurements averaging and outlier detection with the PCA LOF method.



The screenshot shows a user interface for the '1. Preprocessing (Optional)' step. It includes a checkbox for 'Do you want to do preprocessing?' which is checked. Below this, there are two dropdown menus: 'Normalization' is set to 'Arithmetic Sample-Wise Normalization' and 'Missing values imputation' is set to 'KNN-impute'.

More specifically, there are two kinds of normalization: arithmetic sample-wise and logarithmic. It should be noted that when the data contain negative numbers the arithmetic normalization should be chosen, since logarithmic normalization method functions only with non-negative data. If “None” is chosen, no normalization takes place.

There are two kinds of missing value imputation methods as well: average imputation and KNN imputation. Average imputation is a method in which the missing value on a certain variable is replaced by the mean of the available cases¹. On the other hand, the key idea of KNN imputation is that a point value can be approximated by the values of the points that are closest to it, based on other variables². The above missing values imputation methods are relevant only for cases where a missing value does not imply a quantification value of zero. In such cases, missing values should be replaced with zeros before uploading the dataset. If “None” is chosen then no missing value imputation takes place.

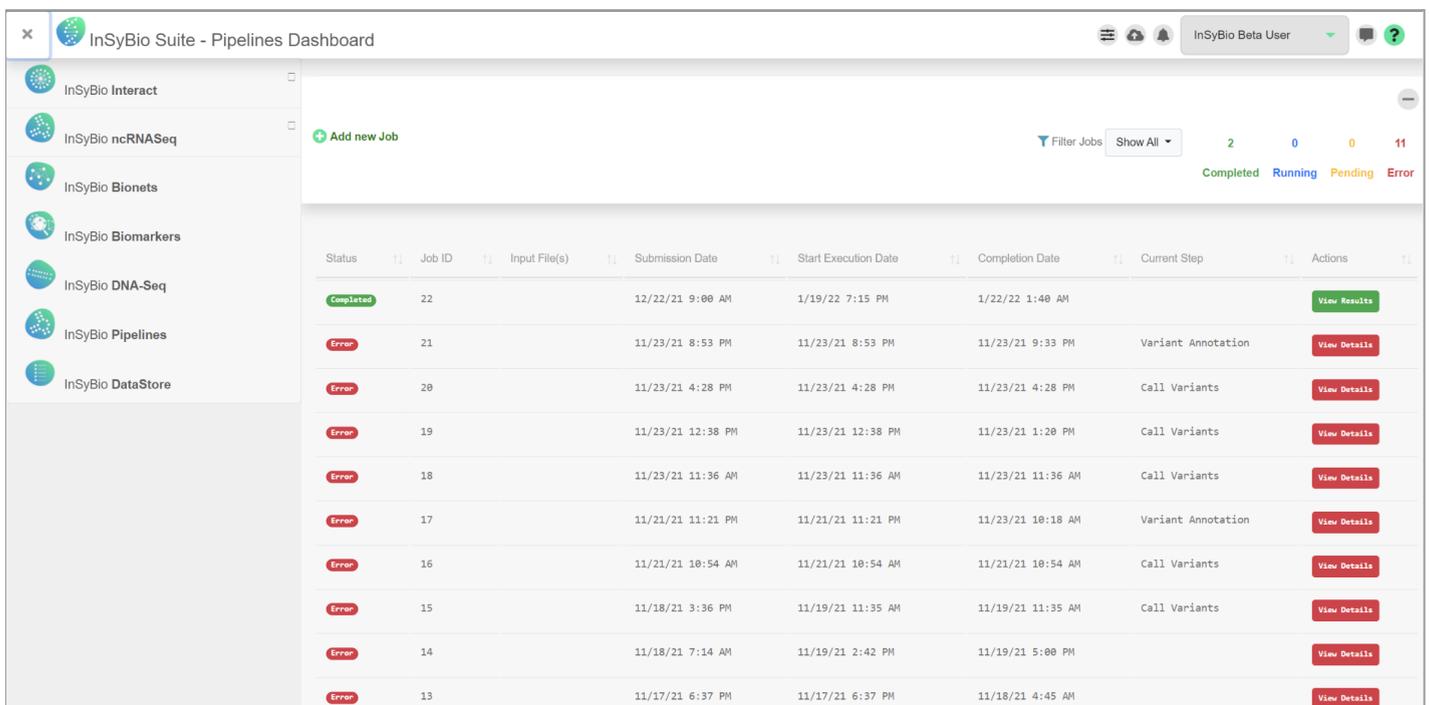
¹ [Single Imputation Methods](http://iriseekhout.com) (iriseekhout.com)

² [The use of KNN for missing values](http://towardsdatascience.com) (towardsdatascience.com)

Furthermore, if a missing values imputation method is being chosen instead of “None”, the duplicate measurements will be averaged.

To view the results:

By starting a calculation the Pipelines dashboard is updated with the submitted job. There you can view the status of your current and previous Pipelines calculations. At completion of the calculation you can select the View Details at the Actions column and view the results.



The screenshot displays the InSyBio Suite - Pipelines Dashboard. On the left, a sidebar lists various pipeline categories: InSyBio Interact, InSyBio ncRNASeq, InSyBio Bionets, InSyBio Biomarkers, InSyBio DNA-Seq, InSyBio Pipelines, and InSyBio DataStore. The main area features a table with columns for Status, Job ID, Input File(s), Submission Date, Start Execution Date, Completion Date, Current Step, and Actions. The table shows one completed job (Job ID 22) and eleven failed jobs (Job IDs 13-21). Each failed job has a 'View Details' button in the Actions column. A summary bar at the top right indicates 2 completed, 0 running, 0 pending, and 11 error jobs.

Status	Job ID	Input File(s)	Submission Date	Start Execution Date	Completion Date	Current Step	Actions
Completed	22		12/22/21 9:00 AM	1/19/22 7:15 PM	1/22/22 1:40 AM		View Results
Error	21		11/23/21 8:53 PM	11/23/21 8:53 PM	11/23/21 9:33 PM	Variant Annotation	View Details
Error	20		11/23/21 4:28 PM	11/23/21 4:28 PM	11/23/21 4:28 PM	Call Variants	View Details
Error	19		11/23/21 12:38 PM	11/23/21 12:38 PM	11/23/21 1:20 PM	Call Variants	View Details
Error	18		11/23/21 11:36 AM	11/23/21 11:36 AM	11/23/21 11:36 AM	Call Variants	View Details
Error	17		11/21/21 11:21 PM	11/21/21 11:21 PM	11/23/21 10:18 AM	Variant Annotation	View Details
Error	16		11/21/21 10:54 AM	11/21/21 10:54 AM	11/21/21 10:54 AM	Call Variants	View Details
Error	15		11/18/21 3:36 PM	11/19/21 11:35 AM	11/19/21 11:35 AM	Call Variants	View Details
Error	14		11/18/21 7:14 AM	11/19/21 2:42 PM	11/19/21 5:00 PM		View Details
Error	13		11/17/21 6:37 PM	11/17/21 6:37 PM	11/18/21 4:45 AM		View Details

In the Preprocessing tab the user will be able to download the resulting preprocessed file.

Statistical Analysis

The second step is Statistical Analysis. This step is mandatory. At this step, Differential Expression Analysis, Heatmap construction and Spearman correlation table construction are performed.

2. Statistical Analysis

Analysis Selection: Unpaired Analysis Selection

Define pvalue threshold value: 0.05

Test Selection: Automatic

InSyBio Suite - Pipelines Results

InSyBio Beta User

Job Status	Job ID	Submission Date	Execution Time	Input Data and Parameters
COMPLETED	22	Dec 22, 2021 9:00:29 AM	06 hours, 24 minutes, 52 seconds	

Preprocessing > Full Predictive Model > Full Model Testing > Statistical Analysis > Differential Expression Predictive Model > Differential Expression Model Testing > Network Analysis > Network-based Predictive Model

Download

preprocessed_data_23.txt File

Only variables annotated as genes/transcripts/proteins will be used for differential expression analysis. If a user has uploaded a phenotypic annotation file with more than two columns then multiple tasks will be created with one column of the phenotypic annotation per file. Every phenotypic column can take two or more values.

The user has to select the type of analysis to be made on the inserted dataset. There are two types of statistical analysis, paired and unpaired analysis.

Afterwards the user will have to insert the p-value threshold value, which is recommended to be 0.05.

Then, the user will choose the kind of test to be performed on the selected dataset: automatic, parametric Ebayes Test Selection, parametric 2-sided Students T-test (or

One-way ANOVA Test Selection) or non-parametric Kruskal Wallis (or Mann Whitney Test Selection) test. If the automatic version is chosen, then our algorithm will decide which test to run: the parametric or non-parametric.

To view the results:

Preprocessing	Full Predictive Model	Full Model Testing	Statistical Analysis	Differential Expression Predictive Model	Differential Expression Model Testing	Network Analysis	Network-based Predictive Model
Network-based Model Testing	miRNA Target Prediction	Enrichment Analysis					
Statistical Analysis Results	Heatmap Visualization	Volcano Plots Visualization	Significant Molecules Dataset	MQ Files	Beanplots Download	All Results Download	Run Info
Statistical Analysis Results (Top 20*)							
*You can download the full results from "All Results Download" tab.							
p-Values 0 VS 1 top20		p-Values 0 VS 2 top20		p-Values 1 VS 2 top20			
IDs	Pvalue	Adjusted Pvalue	Fold Change				
RNF212B	0.18741294693239047	0.5452433804609337	0.08475392968950879				
LOC100129198	0.9293637713589189	0.9785016692464833	0.006983114010646718				
PTPMT1	0.25870166874098385	0.6119673829551042	-0.09314611910088039				
FAM132A	0.37562148584900146	0.7011692112382196	-0.07741626791886391				
LOC157931	0.02927316214513881	0.27998063212706686	0.15770781008409007				
SLC7A11-AS1	0.11107623773055629	0.45178433046415594	-0.12683162544288085				
SIGLEC15	0.869474054318045	0.9561845021784103	0.009908001816207268				
BE327079	0.211841319811446	0.569855370076092	-0.11498243147699527				
BTBD1	0.25262234061810934	0.6065762558972843	-0.07141004166966713				
AA868500	0.7086717631489028	0.8844687572647206	0.026512553390938998				
ITGA2	0.06694479283695004	0.3752367700391344	-0.14432300651737673				

In the Statistical Analysis tab the user will be able to view the Statistical Analysis Results (the top 20 features), the Heatmaps, the Volcano plots, the significant molecules, the molecular quantification (MQ) files and he'll be able to download the Beanplots and all the resulting files. Finally, in the last tab the run information will be displayed.

General Predictive Model Creation

The third step is the General Predictive Model Creation. This step is also mandatory and allows users to train their own predictors using the biomarkers dataset, the phenotypic annotation and the parameters that they selected.

3. General Predictive Model Creation

What's your Prediction Problem?

Two-Class
 Regression
 Multi-Class

Predictor Goals

1. Selected Features Minimization	2. Classifier's Accuracy	3. F1 score	4. F2 score	5. Classifier's Precision
1	10	10	1	1
6. Classifier's Recall	7. Classifier's ROC AUC	8. Model Complexity Minimization		9. Manhattan Distance
1	1	1		1

Advanced Options

Multiobjective Optimization Framework Parameters

Population Size: 50	Arithmetic Crossover Probability: 0	Mutation Probability: 0.1
Generations: 100	Two Point Crossover Probability: 0.9	k in k-fold Cross Validation: 5

These parameters have been tested in various diagnostic and prognostic applications by InSyBio's R&D team and they have proven to provide efficient exploration and exploitation of the search space minimizing also the risk of getting trapped to local optimal solutions. If you are not a bioinformatician with expertise in machine learning, we strongly advice not to change these parameters and to contact our support team if the default values do not provide good predictive models for your dataset.

The user chooses the kind of prediction problem he has at hand. Later, he can alter the weights of the predictor goals. It is advisable to use the default values. The higher the weight, the more significant the goal.

Finally, the user can alter the multi-objective optimization framework parameters at the Advanced Options. Those are the population size, the number of generations, the arithmetic crossover probability, the two point crossover probability, the mutation probability and the number of folds k for the cross validation.

To view the results:

Preprocessing | **Full Predictive Model** | Full Model Testing | Statistical Analysis | Differential Expression Predictive Model | Differential Expression Model Testing | Network Analysis | Network-based Predictive Model

Network-based Model Testing | miRNA Target Prediction | **Enrichment Analysis**

Classification Performance

Cross validation accuracy: 74.96 %
Cross validation F1 score: 70.68 %
Cross validation Precision: 73.27 %
Cross validation Recall: 74.96 %
Cross validation F2 score: 74.62 %
Cross validation Manhattan Distance: 0.75

Training accuracy: 100.00 %
Training F1 score: 100.00 %
Training Precision: 100.00 %
Training Recall: 100.00 %
Training F2 score: 100.00 %
Training Manhattan Distance: 1.00

Model Complexity

Models

Model 1 - Number of Random Forest Trees: 52
Model 2 - Number of Random Forest Trees: 35
Model 3 - Number of Random Forest Trees: 35
Model 4 - Number of Random Forest Trees: 52
Model 5 - Number of Random Forest Trees: 52
Model 6 - Number of Support Vectors: 64
Model 7 - Number of Random Forest Trees: 35
Model 8 - Number of Random Forest Trees: 35
Model 9 - Number of Random Forest Trees: 84

Selected Inputs

ARMC7, MBOAT7, TFIP11, HTR3A, NDUVF2, KMT2E, TMEM126A, MIR1247, RNASE1, BC031588, CWC25, LINC00926, BI603728, 236523_x_at, CTBP1-A, AS2, GHDC, MED19, NXMN1, C1orf116, AI264671, EBDC1, SNHG3, AA700650, TMEM39B, UBC, MRPL14, RPS5, AL079289, LINC00445, LOC101928422, CCNB1, TMEM212, AL542578, RACGA, P1, 207371_at, LOC102723831, DNAI2, ARG1, AL353942, SFI, AA893820, STARD7, HNRNPK, INSR, HNF1A, C7orf25, PROKR2, 242436_at, LINC00942, SH3BP1, OVG1, LOC101927710, SLC36A1, LINC00607, PHK1, SNTG2, AA012953, NUSAP1, 242174_at, LOC102723927, PRRG4, PEM1, RGL4, MEZ, FAM219A, AUL45280, BQ719879, TTK, FKBP1, AI703397, RPR1, P1, FPP1R27, AW293239, JARID3, TMEM35B, ERG1, AKR7A2, DISC1, LINC00886, HSPA4, AF318321, 230663_at, LINC00939, LINC01192, EDE2A, ADAR, KITLG, AW628168, SH2D1A, ZNF445, FTL, AK093193, CXCL8, TEX30, CCNE2, BE468039, AI417657, WDR4, LOC100507506, CDK2, AF086093, COL7A1, BIRC7, LINC01304, ARNA, OSBP2, IL34, GNB2, CD37, FBXL13, FLPFR2, KIAA1324L, MIR4746, EIF3G, CHUK, HFE, N39314, AF086066, TERT5, TOP2A, ABHD14B, AI174119, DDX51, BF930294, MACROD2, GINS1, CYP2F1, GSK1, FCER2, AGO2, PITPNM2, NBSN2, FARVG, ITGAM, OPTN, IQCH-

AS1, IL22RA2, NEU3, AI598222, PTFN2, TNKI, ATP5H, PRKC2, LOC100288860, DDI1, AI698731, DOCK11, 243642_x_at, DHX15, NMB, AA909691, BC041381, CTR9, C1orf131, PROSER1, BE048525, HIST3H3, BC033361, AI024328, UBQLN4, POLR2D, YEATS2, AA648438, DAXX, PITPNB, CCNG, LINC00441, 221071_at, PRKD1, PCDH12, TMEM98, TCF7L2, RNMT, GALNT3, HT, R5A, CSNK1G3, CCL24, SDCAG8, FKHD1, PLEC, FLJ41170, ABCA17P, ULBP3, BC009884, USP54, LOC100506684, LOC101060275, 225714_s_at, FOXRED1, CAPN1, AA724992, LINC00026, 5, PLPP3, TLE2, RFTN1, 229329_s_at, CHAC2, SIVA1, CES4A, WDR64, CSRN2, GPER1, 221117_at, AW582267, KCNG2, TMEM208, MGAT4B, AA057437, ADAMTSL4, AXIL, TEX2, RGF1, CDC144A, MSS51, SF3A3, COF3, 230346_x_at, ZNF740, EMC2, SACML, NLR4, CDC42BFB, ADGRB1, HIST1H2BB, AI740763, XIAP, AF3M2, E2TB46, LOC728392, WTH3D1, CEM1, MIEF2, EV, FL, AL833053, R47946, LOC101928682, AKO26967, AVIL, LOC101927490, AI808306, Dbpht2, ASB6, HOXD-

AS2, MYO1A, GMEG, ABCB9, TBC1D23, NDUFA6, 242815_x_at, AW615179, ADAM3B, EPN3, AW298101, AL134451, BTBD7, DGCR14, NCF1C, SLC22A15, CST11, ABCCEP1, AI126321, NUDCD3, MRFS12, BF507371, BF445149, SRGN, PPP6R3, LINC00174, NUCKS1, AI820991, SRI, SRF68, BE219311, AA805239, LOC388780, MARVELD2, SPATA13, ATMIN, TEX10, CRYBG3, FAM151, A, TUBA4B, PKKX, MBDS, WT1, BC042969, AW293456, DCC2A, DEFB106B, X72882, UCHL3, AL832732, AA001390, ZMYND11, ADD3, KIAA1586, AA72874, TREML3P, 220458_at, PRRC2B, C, NOT9, C6orf47, APT, GLRA2, MRF821, MGP, ATP6V0B, RHOF, LOC101929718, 216497_at, BF510844, URI1, SEMA3F, WWC3, ELK1, NNS, SOX15, AI701857, FUS1, TRIM22, AA682265, AV6, 48843, TLR9, DUS1L, TESK1, LINC00355, KCNJB, NR3C2, PIEZO2, NFE2L1, CDC16, ADAM33, AW467070, FAM206A, AI861840, EP300, SNORD73A, TMEM41A, AF131767, ELMOD2, COX6A1P, 1, NDNF, GSDMD, LXB2-

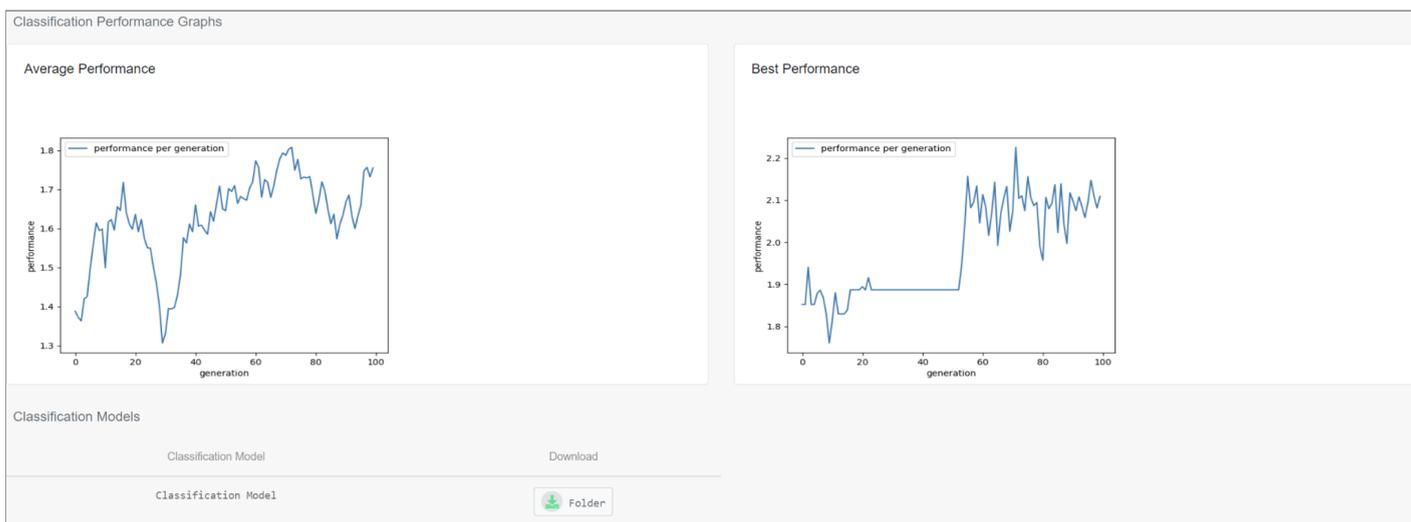
AS1, TXNIP, LOC100507006, AW367380, RPL7, 242256_x_at, LOC100652777, MRGFRX2, AI869532, ALKBH7, AW268162, 1553354_a_at, IFT20, NACAD, ABCB10, N36400, HOKX5, OSBP, KIAA0319L, C17orf64, AI911163, 234413_at, ARHGAP19, AK021486, LOC100131170, W58344, HMGAL, KLR2, RTAD5, DGKA, SNORD35B, AA707317, 243202_at, BC036639, DNAH9, RT, N2, AI963959, TMEM44-AS1, GABRR2, DUSP26, ATP2B1, MAP3K3, ELK2AP, AI921894, ZNF106, PTPRA, CCDC167, MIR6859-

1, BC001743, HUS1, IFI272, AF116671, PRKC2, ZDHHC13, 217019_at, STK31, FKBP, TLR3, MMS19, PSMB4, KLHL11, LINC00472, LRRC37A4P, RBMY1J, RORA-

AS1, AW517851, PFDNS, MIR452, AB062480, 239373_at, RA928078, BMCL1, FURIN, APTPH, KCTD13, RAB11FIP5, GTPBP6, FBLL1, GSN, 221379_at, IFI6, ZNF283, AI457965, CFP52, H, SEB6, AW450033, BC014996, LARF1B, BF511409, MILR1, GFR42, AF088044, BC012090, AFCS, AL049443, SCGB3A2, BE222450, CFP161, PFP4R4, CEP85, PPP1OC, AUL46384, SLC19A1, RPRD2, EHD2, ZMIZ2, OCIAD2, YIPF1, C15240, MRPL36, SGM3, THAP3, SH3GLIP2, LINC-

PINT, CRADD, AA477687, AA702963, MBD2, BCL11B, FOXD4, COA6, PNMA6A, DIO308, ACTA1, POLR3E, CA5B, AW303454, KLHL11, CYP2D6, KIM1B, UCAL1, RSPH10B2, GABRG2, BUB1B, AI80, 0518, TOMM5, APOC3, REG1A, CIZ1, SFRN, BC039091, EFOR, UBE2T, CARM1, PHF5A, BEX1, AK022046, CFI, PSAP, ATP11A, APT1, UCK1, UPK3A, FAB5C, LOC728095, CADM3, CMTM3, H0956, 4, TXNDC9, AV649275, ERO1A, ALDH1A2, PBK, ZFAND5, ZNF304, TMEM70, SSU72, SLC44A3, FLJ30064, RPL7AP10, ELM, SRRM5, CENPB, BF224366, ZNF146, DLX3, LOC728445, RTN4I1, AI620881, AI762446, AK025072, MICAL2, ZNF503, MXRA8, OXLD1, AUL47360, SH3BGR3, AK098328, STK35, BC034639, ILIRAPL2, RAB27A, AK022801, PRDM6, HMMB1, LOC101929762, BC027465, ORA12, HOOK1, CLS1N1, 231486_x_at, LOC100508408, LOC102724094, FAMI3A, 217329_x_at, MRPL13, COL11A2, SPATA3, MBNL3, 216494_at, CDF1, KHSPR, AK8, AF147, 397, BRPF3, BOLAA2B, HYI, HDAC8, KLK11, TMEM176B, BPNT1, SIGLEC11, PML2, 1565875_at, AW087759, ANKRD17, LOC101928284, LINC00305, HEM, NKRF, AK021967, TMEM230, RNASE, 4, AW269743, GABPB1, DPFYSL4, RALGAPB, BTNL2, LYZL1, ZYX, AI538880, IMPG1, AL080112, RAD23A, TMA16, KIFC2, IL20, RBM3, EMC8, LINC01165, FDI, SLC41A3, NUP205, ARF5, NO, D2, AI926424, AF172327, GBX1, GPANK1, FLJ33534, ERALL1, VTA1, 244582_at, USP14, FAM228A, CXorf51B, AW296118, PRSS27, LOC101930595, MRPE3, UTP11, THYN1, MRPS6, AF339, 772, HCG9, EEB41L2, LOC100996579, ADAMTSL3, ADAM3A, SNK11, AF222858, TMEM107, PDS5B, MIR4657, VANGL2, DFFA5, 243078_at, MIR3656, ACTR3C, TYR, RHHBDD3, SOX30, MAGEA4

```
.ISOC1,YFEL3,RPS24,RECK,RBM15,HID1,ZNF816-
ZNF321P,AF187554,BE71L,PNMA3,2P2,AK026701,TUBB4B,ARHGAP6,AI733177,RASA3,PRDM12,AF090925,NMRAL1P1,TIMM13,UNC5A,AQF12B,C7orf43,AV736725,AIP,AI7987
24,UBAF2L,PAFLN,BE044484,AW172407,LINC00588,LOC100133920,SNF8,FQLC2L,AK022067,CDC26,1556333_at,BC016176,ASMTL,237424_at,AI469935,LINC01535,AF086
565,AI733457,SFMBT1,TEM26-
AS1,TFAP4,240896_at,HCWC2,AI0117540,LOC102724312,RBBP5,C5orf66,CENPH,RAD54B,IFITM3,GUCY2D,BC041050,SEFN1,RFL38,HRG,BC035326,244488_at,PLEKHA6,RAB
3GAP1,AL831920,RHOT2,IQCF5,2FP62,IFNA5,PET100,SEHL1,OPN5,ATG16L1,AI457449,AL043143,LOC101928510,AW779654,MREL2L,FOLH1,MZT2A,AL075053,VSTM
4,AI022850,TAL1,228586_at,SLC12A6,AI084064,NUDT19,PABEN1,PMF1-
BGLAP,AW452392,AK022254,BE675275,MYO18A,KLF5,PDCD6IP,RHBG,EFNB3,228896_at,BE858194,CASD1,LL19,N40199,CLSTN2-
AS1,KDF1,PHLDA2,SNX5,LOC100287896,R0660,PRKV,MEIOC,MEX3C,UBL4A,KIF13B,ATG4D,PCA3,MTA1,AI074594,JUNB,NKX2-
5,MRPS17,GNL3,C4orf45,MEL,AP2L1,LOC339685,ACTL9,TSSK3,JRH4,AF339810,BC039410,AW974816,H95280,KLHDC7B,208144_s_at,1554281_at,CCDC151,EP400,ADORA2
B,AI076370,FRKD2,IFNG,LOC100506885,ARL4D,ERGIC2,X84340,FDCC2L,GFRA2,KENNA6,PTNS,STRA6,PIK3R6,UBXN11,GGH,KIARA0907,RRM2,AW082221,AK021495,ANKA2F3,
AK091415,207916_at,EMS2CL,ATAT1,EPHA5-
AS1,NSUN4,C22orf29,240654_at,TMNC2,CLC4F,N23258,AL832887,C6orf118,AJ012498,212883_at,S100A12,MICALL1,ZNF600,C9orf131,GJA1,AV700385,TRAF2,TEM10
8,TEM51,LOC100129380,PPF1R36,LTC4S,APC,SYSL,IL36A,FAM210A,AK023800,RNF152,AL050136,BC031975,LOC728805,FAE1,HUWE1,AI821649,1566268_at,CUL9,CDC42
SE1,C10orf11,NPM1,IFNA14,RASSF4,LOC101060835,AU156181,ENF831,AW237316,MARS2,FOXO3B,CEP76,TNRC6C-
AS1,EXOSC4,TMSB10,SGCD,TRIM74,LOC101927040,C3orf35,LOC155060,228549_at,LOC400620,AF147426,LOC101927929,RSLLD1,243160_at,ZNF362,CC2D2A,AF130093,D
BT,AI419968,AI923201,DHTRD1,ZDHHC4,SEPT9,BC041976,PLCL2,BC040306,RANBP1,MIEF1,AI674915,AW139719,244047_at,NDUFB8,SZT2,FAXIP1,NSMCE1,TSC22D1-
AS1,SUPT6H,IP04,LOC101927379,SLC45A2,HAT1,RNPEFL1,LOC81787,CDK5RAP3,ANGPT4,OXSM,MARK1,TEM184C,SNAPC5,228156_at,BF508839,GATAD1,BC031957,SUFU,FL
J16779,MFSDB4B,1556828_at,DKO7,EGR3,GUSB,240451_at,LOC285847,AI700768,THOC1,BC013633,POF7,CROCC,CH17-
360D5.1,NFF,PLGRN,COX8C,SHQ1,208421_at,TFAP2A,MIR4784,LRFN4,BC042181,C5orf51,234454_at,AK024973,AI733345,BE463783,1570653_at,LZIC,AI360167,WDF
Y2,TMCC3,KRT76,AW015426,LOC285902,WDRD1,MIR3917,AA010315,SMG9,LOC101928886,LOC101928787,CDBA,SPRY3,241217_x_at,LILRA2,CDT1,FLAGL2,ZNF446,C17orf8
9,BG403405,AK098724,AF085897,FAM96B,AU150817,TBC1D5,MYH14,C17orf98,BAG6,RRAS,CDCA2,ZNF24,UNG,PHF7,ZCCHC4,GJA5,SNORD18C,BLOC1S2,216943_at,BF51121
2,AL159594,PPF1R9B,AW271060,TEM171,TEM8B,SRF,AW572853,AI247365,216568_x_at,BRD3,CCDC93,CROT,FMFBF1,AI922972,LOC101929741,CNTNAP1,LOC728084,CDC
A5,ARAF1,BE71L,SRF54,U90905,HAUS1,BC037412,MIR4640,WTA6,XRCC2,BE504838,FRORSIP,SECL16A,VT1A,KDMA,FUT1,ULK3,AW205775,COLGALT2,ZNF467,GLI3,CIR,HS
D17B12,PITPNCL,MYZAP,LOC100507194,MAGEC1,AI421677,LBB1,IL10RB,SLC24A1,LOC100289045,FLAD1,1560905_at,PNMT,INPF5B,AI821694,WDR25,AI923713,SNORD3D,
MEAL,MPPED1,FLXNA3,LINC01468,242192_at,CDRI,HSD17B1,ACTL7A,TUBB2B,AF075064,LOC100310756,LOC101927770,FRK,GCCT,PLRG1,FAM110B,LINC01004,TRBC1,PTTG
2,SECISBP2,CFAP73,VEB4B,GAS1,GRIA4,217676_at,234372_at,AI377746,LOC102724782,AI760332,RBM45,SUN5,AK026468,AP1S2,AU148090,FAM86B3P,DMRTB1,PCDHGC4
,DCAD,Clorf27,MS42,HPCAL4,ESRP1,ATG4A,C10orf10,PAK6,TUBB1,ZSCAN25,MTM1,FAM155A,242223_at,RNF185-
AS1,B3GAT3,RIDA,NOB1,CAAP1,CYP19A1,SMG5,TF2,TULP3,ESENN,TMX2-
CTNND1,AA648962,RALA,MURC,TAI2,LINC01622,238586_at,LINC00427,C2orf76,UBE2C,AI138418,DDX49,AF007143,STEA3P,BE066500,LOC105370629,SCO2,AK000106,M
SH2,881578,216421_at,GRB10,AI822140,BC042815,W96141,ADD3-
AS1,GT2P7,ST18,SFTBN2,MEO,BF847120,H27618,LINC00282,AL162010,AA814006,ARHGDB,AKT2,220862_s_at,TEM258,AI280131,LOC101928521,NDUFA7,ZNF646,LNX
2,H0710,BCL7C,MSUN2,AV659223,MIR205
```



In the Full Predictive Model tab the user will be able to view the classification (or regression) performance of the cross validation and of the training set. For the two-class prediction problem the user will be able to see the accuracy, sensitivity, specificity, f1 score, f2 score, ROC AUC score. For the multi-class prediction problem the user will be able to see the accuracy, f1 score, f2 score, precision, recall and Manhattan distance and for the regression prediction problem the user will be able to view the root mean square error, the relative absolute error, the root relative squared error, the R2 (coefficient of determination) regression score, the explained variance score and the Spearman Correlation.

Additionally, the complexity of every model, which is the total number of support vector machines and the number of trees for RandomForest and also, the average and the best performance of the trained model are being displayed.

Differential Expression Predictive Model Results:

Preprocessing Full Predictive Model Full Model Testing Statistical Analysis **Differential Expression Predictive Model** Differential Expression Model Testing Network Analysis Network-based Predictive Model

Network-based Model Testing miRNA Target Prediction Enrichment Analysis

Classification Performance

Cross validation accuracy: 48.48 %
Cross validation F1 score: 47.52 %
Cross validation Precision: 50.47 %
Cross validation Recall: 48.48 %
Cross validation F2 score: 48.87 %
Cross validation Manhattan Distance: 0.48
Training accuracy: 67.19 %
Training F1 score: 65.95 %
Training Precision: 69.45 %
Training Recall: 67.19 %
Training F2 score: 67.63 %
Training Manhattan Distance: 0.67

Model Complexity

Models
Model 1 - Number of Support Vectors: 55
Model 2 - Number of Support Vectors: 55

Selected Inputs

ARM7

At the Differential Expression Predictive Model tab you can view the same results as at the Full Predictive Model tab. The difference between them is the input to this same step. For the Differential Expression prediction step the input is original dataset with only the significant biomarkers selected.

Network-based Predictive Model Results:

Preprocessing Full Predictive Model Full Model Testing Statistical Analysis Differential Expression Predictive Model Differential Expression Model Testing Network Analysis **Network-based Predictive Model**

Network-based Model Testing miRNA Target Prediction Enrichment Analysis

Classification Performance

- Cross validation accuracy:** 68.98 %
- Cross validation F1 score:** 65.44 %
- Cross validation Precision:** 67.51 %
- Cross validation Recall:** 68.98 %
- Cross validation F2 score:** 68.68 %
- Cross validation Manhattan Distance:** 0.69
- Training accuracy:** 98.44 %
- Training F1 score:** 98.42 %
- Training Precision:** 98.49 %
- Training Recall:** 98.44 %
- Training F2 score:** 98.45 %
- Training Manhattan Distance:** 0.98

Model Complexity

Models

- Model 1** - Number of Random Forest Trees: 10
- Model 2** - Number of Random Forest Trees: 13
- Model 3** - Number of Random Forest Trees: 10
- Model 4** - Number of Random Forest Trees: 10
- Model 5** - Number of Random Forest Trees: 10
- Model 6** - Number of Random Forest Trees: 10
- Model 7** - Number of Random Forest Trees: 19
- Model 8** - Number of Random Forest Trees: 19
- Model 9** - Number of Random Forest Trees: 19

At the Network-Based Predictive Model tab you can view the same results as at the Full Predictive Model tab. The difference between them is the input to this same step. For the Network-Based prediction step the inputs are the original dataset with only the significant biomarkers and the output of the Network Comparison

Network Analysis

The fourth step is Network Analysis. This step is optional and it consists of five steps:

- Bionet Create Gene Co-expression Networks.
- Biological Network Analysis
- Network Comparison Biomarkers
- Interact Enrichment Analysis
- Bionets Clustering

At each step, the default value for every parameter is selected. These values can be configured manually by the user.

4. Network Analysis (Optional)

⚠ If you want to do Biological Network Analysis or/and Network Comparison Biomarkers or/and Interact Enrichment Analysis or/and Bionets Clustering, then you should also do gene Co-expression Network Creation.

Do you want to do network analysis?

Bionets Create Gene Co-expression networks

Biological Network Analysis

Network Comparison Biomarkers

Interact Enrichment Analysis

Bionets Clustering

Gene Co-expression Network Creation

Method: Pearson

Interval of trust: 99%

Filtering parameter minimum variance:

Filtering parameter minimum average logarithmized expression:

To view the results:

In the Network Analysis tab you can view the five different tabs for each step of the

Preprocessing			Full Predictive Model			Full Model Testing			Statistical Analysis			Differential Expression Predictive Model			Differential Expression Model Testing			Network Analysis			Network-based Predictive Model								
Network-based Model Testing						miRNA Target Prediction						Enrichment Analysis																	
Co-expression Networks						Network Analysis						Network Biomarkers						Enrichment Analysis						Clustering					
Gene Co-expression Network Filename						Threshold						Download																	
mq_coexpnet_0_5_23_0.tsv						0.5						File																	
mq_coexpnet_0_5_23_2.tsv						0.5						File																	
mq_coexpnet_0_5_23_1.tsv						0.5						File																	

Network Analysis steps.

Gene co-expression network creation

Gene expression files, either uploaded directly by the users or generated through soft files parsing, can be used to generate weighted gene co-expression networks. Experienced users can tune the parameters of the algorithms used for this step and select the most suitable algorithm for them. Three options are offered:

- **Pearson Correlation:** This method adds an edge to a network if the Pearson correlation of the nodes adjacent to the edge exceeds a threshold.
- **Mutual information:** This method adds an edge to a network if the mutual information among the expression profiles of the two nodes of the edge exceeds a threshold.
- **Spearman Correlation:** This method adds an edge to a network if the Spearman correlation of the nodes adjacent to the edge exceeds a threshold.

The thresholds for adding edges are dynamically generated to alleviate problems occurring by using the same threshold for all nodes. In particular, for a single node Pearson correlations or Mutual Information or Spearman correlations between this node and all other nodes are calculated. Assuming that the Pearson correlation/Mutual Information/Spearman correlation values between a single node and all other nodes follow a normal distribution, then the threshold for adding edges is selected to be in a predefined confidence interval (90%, 95% or 99%). The confidence interval is predefined at 99% but the users can change this value in order to get denser or sparser networks. In order to force nodes to have a minimum number of edges users can also specify a minimum value for the threshold of adding an edge in the network. Experienced users can further filter nodes from the network by altering the minimum expression variance threshold and the minimum average of the logarithmized expression values threshold.

To view the results:

In the Co-expression Networks tab you can view and download the created networks.

Gene Co-expression Network Filename	Threshold	Download
mq_coexpnet_0.5_23_0.tsv	0.5	File
mq_coexpnet_0.5_23_2.tsv	0.5	File
mq_coexpnet_0.5_23_1.tsv	0.5	File

Network Comparison Biomarkers
 Interact Enrichment Analysis
 Bionets Clustering

Interval of Trust (Most Significant Nodes): 95% ▾
 Method (Most Significant Edges): Edge weight ▾
 Interval of Trust (Most Significant Edges): 95% ▾

Biological Network Analysis

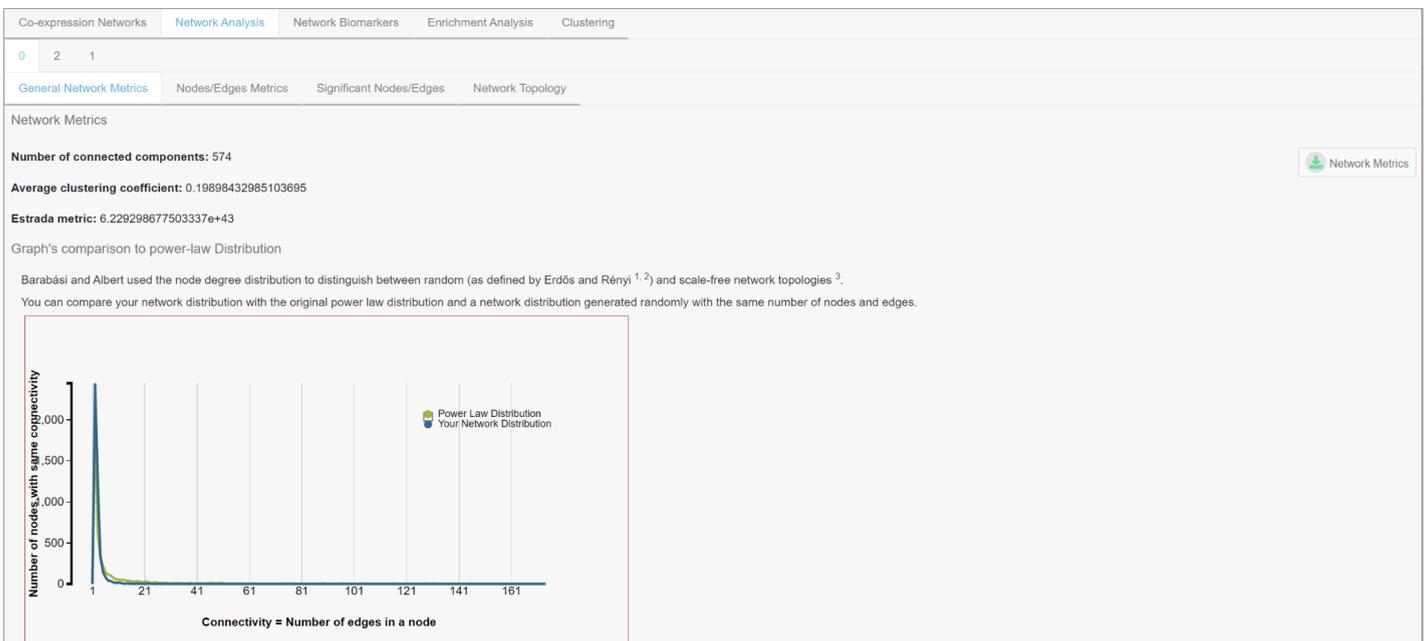
In order to analyze the biological network, the users can tune the following parameters if they are experienced:

- Method for selecting significant nodes (Pagerank (default), Clustering Coefficient, degree centrality])
- Confidence interval for locating significant nodes
- Method for selecting significant edges (Edge weight (default), Inbetweeness centrality)
- Confidence interval for locating significant edges

To view the results:

Four new tabs are generated, General Network Analysis, Node/Edges Metrics, Significant Nodes/Edges and Network Topology.

At the General Network Metrics tab users can view the most significant network metrics (clustering coefficient, Estrada index and so on) and compare the degree distribution of their network with a random network's power law distribution. Information in this tab is not available for networks with more than 225000 edges.



At the Node/Edges Metrics tab users can find the metrics for all nodes (degree centrality, clustering coefficient and pagerank centrality) and edges (edge weight and in betweenness centrality) of your network.

Co-expression Networks **Network Analysis** Network Biomarkers Enrichment Analysis Clustering

0 2 1

General Network Metrics **Nodes/Edges Metrics** Significant Nodes/Edges Network Topology

Node Metrics

Browse among metrics Node Metrics

Node	Degree Centrality	Clustering Coefficient	Pagerank Centrality
ANKRD17	0.00022341376228775692	0	0.00010364841101858703
APC	0.00044682752457551384	0	0.00013748721716292053
RANBP1	0.0013404825737265416	0.008977556109725686	0.0002606384630344983
HNF1A	0.0011170688114387846	0.005319148936170213	0.00024345433958942114
MSH2	0.00044682752457551384	0	0.00031837835861361624
PRK CZ	0.00022341376228775692	0	0.00019801492562043118
CCNB1	0.00044682752457551384	0.0	0.00023821496586395212
CDK1	0.00044682752457551384	0	0.0002643453226221462
SEPN1	0.006478999106344951	0.030863545777296785	0.0004979254986518572
ELK1	0.00022341376228775692	0	5.1240886233946666e-05
LOC100653049	0.010947274352100089	0.4298614323351367	0.0002949315974561881

Co-expression Networks **Network Analysis** Network Biomarkers Enrichment Analysis Clustering

0 2 1

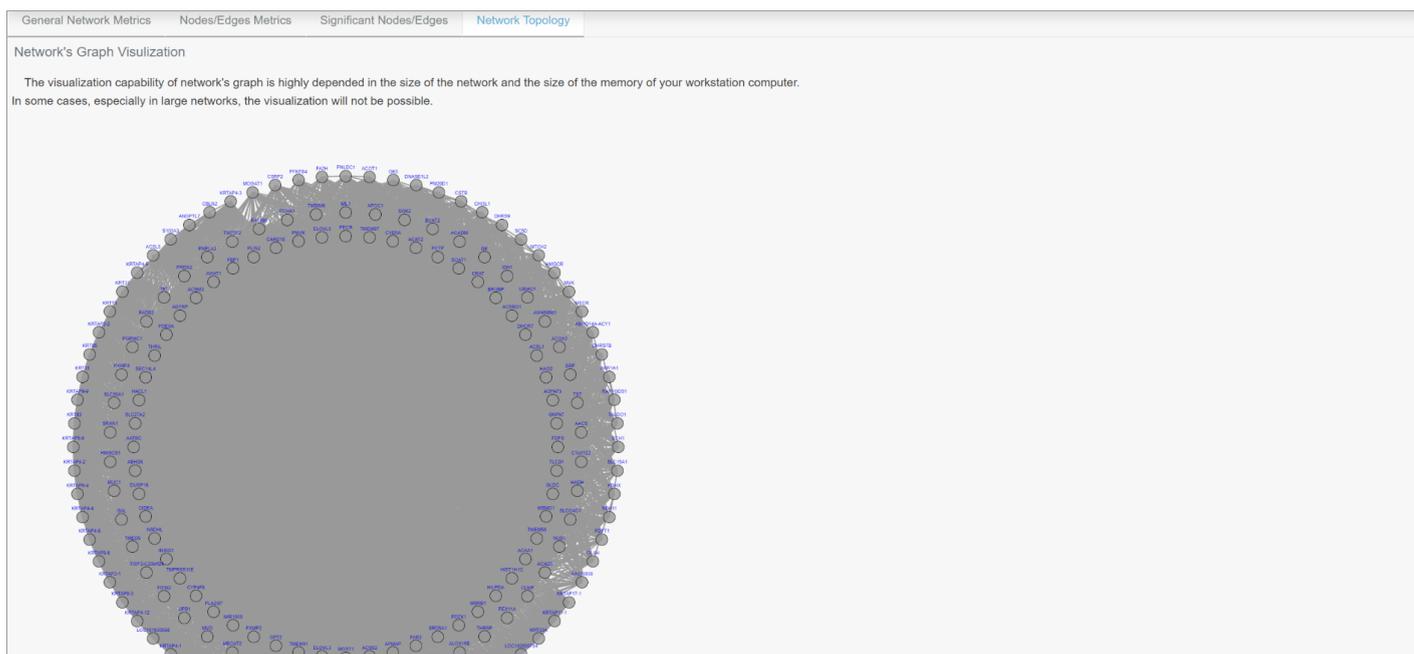
General Network Metrics Nodes/Edges Metrics **Significant Nodes/Edges** Network Topology

Most Significant Nodes

Browse among metrics Most Significant Nodes Metrics

Node	Degree Centrality	Clustering Coefficient	Pagerank Centrality	P-value
LOC100130370	0.013404825737265416	0.006368880239310772	0.0019119370202301674	5.139867197997187e-25
SARDH	0.012064343163538873	0.006093109565855748	0.0017764936692936813	1.8578161929389055e-21
HOXC10	0.02033065236818588	0.023576114267675145	0.0013400867999560791	5.695042202916822e-12
WBP2	0.022117962466487937	0.02479571364777837	0.0012866514545128629	5.134090812430496e-11
HNRNPUL1	0.021447721179624665	0.023306981599252832	0.0012729840938720332	8.861616262259541e-11
MAZ	0.021447721179624665	0.024041884981079963	0.001268088902179435	1.0757286747437055e-10
AI364876	0.007819481680071492	0.009668267675385114	0.0012634502589654794	1.291615427448428e-10
SF3A2	0.01876675603217158	0.02588632169279417	0.0012229745725881292	6.164563550787137e-10
FBXO44	0.007596067917783735	0.008500942114718898	0.0011788655361056723	3.1650058690534606e-09
CENPB	0.019660411081322608	0.026770813755190213	0.0011581736569499508	6.655405930659457e-09
PRKCSH	0.020107238605898123	0.026396007419605797	0.0011421515804376137	1.1709248358259013e-08

At the Significant Nodes/Edges tab users can access two tables including the significant nodes and edges. For each node and edge, the respective metrics and the p-values of their significance are provided. Significant edges are not available for networks with more than 225000 edges.



The Network Topology tab offers an interactive visual representation of the biological network. When networks have more than 10,000 edges, a haircut filter is applied before the visualization of the network. If the haircut filter cannot reduce the number of edges below 10,000 edges then no network visualization is provided. Networks' visualization is based on the Cytoscape plugin and it provides an interactive graphical interface. Users can retrieve information about clicked nodes and edges, export the image in different formats (a PNG, SVG, JPG), decrease opacity on mouseover and view the network using different visualization layouts (force-directed, circle or radial).

Export Network Layout

[PNG image](#) [Save Visualized Network](#)

Network Comparison Biomarkers

4. Network Analysis (Optional)

⚠ If you want to do Biological Network Analysis or/and Network Comparison Biomarkers or/ and Interact Enrichment Analysis or/and Bionets Clustering, then you should also do gene Co-expression Network Creation.

Do you want to do network analysis?

Bionets Create Gene Co-expression networks

Biological Network Analysis

Network Comparison Biomarkers

Interact Enrichment Analysis

Bionets Clustering

Network Comparison Biomarkers

Confidence interval:

Method:

It is widely accepted lately that differential expression biomarkers are large in numbers, contain a large number of false positives and mainly depict the outcome of disease mechanism and not its cause. For this reason, the current trend in biomarker discovery is to detect biomarkers by comparing biological networks. Biological networks are slightly altered in different biological conditions and changes on them are associated with the causes of disease mechanisms with high probability.

When having two biological networks of different conditions, users can use them to predict network biomarkers with an InSyBio's novel methodology. In particular, a certain network metric is selected and InSyBio BioNets attempts to detect network's nodes with significantly altered values for this network metric. Thus, our approach finds nodes whose role in the network has significantly changed among the different conditions. Experienced users can select a specific network metric among the following ones:

- Degree Centrality
- Clustering Coefficient
- Pagerank method

Pagerank method is the default one. This method triggers random walkers starting from each node. Significant nodes are collecting more information from the diffused quantities of the random walkers over time. Experienced users can also select the

confidence interval for tuning the threshold of assigning a node as biomarker. Higher confidence interval values lead to the extraction of more compact sets of biomarkers.

To view the results:

Co-expression Networks		Network Analysis	Network Biomarkers	Enrichment Analysis	Clustering			
0 vs 2	0 vs 1	2 vs 1						
Gene Expression	Confidence Score	Centrality metric in control network	Centrality metric in examined phenotype/condition network	Difference in centrality metric between examined phenotype and control networks	Database	Related Uniprot ID	Link to External Databases	
MAPK8IP3	1.0	0.000267799491643915	0.0006987786201103868	0.00018542274368304752	Gene Symbols	Q9UPT6 , E9PFH7 , B7ZMF3 , H3BN91 , G1UI24	GeneCards OMIM	
SLC12A4	0.9537835058549945	0.00030367511578434115	0.00048087304174944607	0.0002426171312549065	Gene Symbols	Q9UP95	GeneCards OMIM	
LOC100130370	0.9374170716884797	0.0007945623004465563	0.000372258856712521	9.908005172473156e-05			-	
BLZF1	0.9076678045361469	0.00022336385972749609	0.0016772649565850546	6.118801026087902e-05	Gene Symbols	Q9H269	GeneCards OMIM	
SPTBN4	0.8854177551403185	0.0003936722558950883	0.000261953139665185	0.0001216436929401671	Gene Symbols	Q9H254	GeneCards OMIM	
LRRC43	0.870937071331673	0.00010364841101858703	0.00011822655006919568	0.00022336385972749609	Gene Symbols	Q8N309	GeneCards OMIM	
GJC2	0.8674068116674456	0.00025883302930618403	0.0001191436168375295	5.288213768444955e-05	Gene Symbols	Q5T442	GeneCards OMIM	
TNK2	0.8435961158238391	0.0003960090054490455	0.0005797886505628168	0.00020149103364900262	Gene Symbols	Q07912 , C931X3 , H7C412 , C9J0G3 , H0Y5H7 , F8NER3 , H7C343 , H7B2M8 , H7B2Z3 , C931R5	GeneCards OMIM	

At the Network Biomarkers tab the results are presented in a table with Gene, Confidence Score, Centrality metric in control network, Centrality metric in examined phenotype/condition network, Difference in centrality metric between examined phenotype and control networks, Database, Related Uniprot ID, Link to External Databases columns. Clicking a Gene Expression field the user can view diseases associated with that gene. Clicking a Related Uniprot ID field the user can view the related protein in our InSyBio Interact tool. Clicking a Link to External Databases the user can view the gene in external databases.

Interact Enrichment Analysis

4. Network Analysis (Optional)

⚠ If you want to do Biological Network Analysis or/and Network Comparison Biomarkers or/ and Interact Enrichment Analysis or/and Bionets Clustering, then you should also do gene Co-expression Network Creation.

Do you want to do network analysis?

Bionets Create Gene Co-expression networks	<h3>Interact Enrichment Analysis</h3> <p>Use any known identifier for denoting your biomarkers: Uniprot IDs, gene symbols, RefSeq_id and so on. Mixed identifiers are not supported!</p> <p>Pvalue threshold ? : <input type="text" value="0.05"/></p>
Biological Network Analysis	
Network Comparison Biomarkers	
Interact Enrichment Analysis	
Bionets Clustering	

You can perform enrichment analysis with hypergeometric distribution on a given a list of proteins, genes or transcripts and produce a list of GO terms associated with the list, with their term specificity and score in the distribution. You can also provide your custom annotation, term, term type and functional annotation of molecules files, that will be appended to the default files to perform the enrichment. You can define a pvalue threshold for the biomarker to GO terms association output.

To view the results:

Network-based Model Testing		miRNA Target Prediction		Enrichment Analysis	
Network Biomarkers		Differential Expression Biomarkers			
GO Term	GO Term's Type	GO Term's Name	GO Term's Specificity	Enrichment Score	Associated Uniprot ids
GO:0034260	biological:process	negative regulation of GTPase activity	8	0.0257066828955476498	Q07912
GO:0035268	biological:process	protein mannosylation	8	0.01858060806315525	O60762,Q96E22
GO:0005881	cellular:component	cytoplasmic microtubule	6	0.01400675073029111	P38622
GO:0005643	cellular:component	nuclear pore	5	0.0018448455849128349	Q9UWD3
GO:0006417	biological:process	regulation of translation	7	0.009660680011730018	Q96EY7
GO:0006461	biological:process	protein complex assembly	5	0.012340289698257073	P11047,P29590,P63027,Q15334,Q9BMM4,Q9UQR1,Q9Y566
GO:0016311	biological:process	dephosphorylation	6	0.020318205512264802	P09467,P35813,Q9BY84
GO:0045087	biological:process	innate immune response	4	1.575854940325369e-08	O43914,O94817,P02745,P02747,P06241,P09871,P12931,P29590,Q07912,Q13263,Q8IWM3
GO:0019899	molecular:function	enzyme binding	4	0.0003491810118052967	O00267,O14975,P00167,P06241,P12931,P16157,P17535,P23378,P26368,P42224,P43246,P46108,Q03135,Q5XKE5,Q7L592,Q92802,Q99638,Q9BU72,Q9C0C2,Q9H2G9
GO:0010628	biological:process	positive regulation of gene expression	7	6.961521611935466e-06	O60895,P01127,P05549,P28906,Q03135,Q9P1Z2,Q9UK33,Q9Y6Q9
GO:0000049	molecular:function	tRNA binding	6	0.04540573670923506	Q9BV44,Q9HD40

At the Enrichment Analysis tab you can view the results that are a list of GO terms, terms type and name, specificity, enrichment score, associated Uniprot ids and input ids.

Bionets Clustering

At this step, you can analyze your Biological Network to extract complexes of similar nodes (i.e protein complexes), Weighted and unweighted Biological Networks.

For the prediction of Biological Network complexes one option is supplied:

4. Network Analysis (Optional)

⚠ If you want to do Biological Network Analysis or/and Network Comparison Biomarkers or/ and Interact Enrichment Analysis or/and Bionets Clustering, then you should also do gene Co-expression Network Creation.

Do you want to do network analysis?

Bionets Create Gene Co-expression networks	<p>Bionets Cluster</p> <p>Select Algorithm: ClusterONE - Clustering with Overlapping Neighborhood Expansion</p> <p>Algorithm parameters</p> <p>Complexes Size Threshold: <input type="text" value="3"/></p> <p>Complexes Density Threshold: <input type="text" value="0.3"/></p>
Biological Network Analysis	
Network Comparison Biomarkers	
Interact Enrichment Analysis	
Bionets Clustering	

Clustering with Overlapping Neighborhood Expansion (ClusterONE).

- Complexes size threshold: (default value 3)
- Complexes density threshold: (default value 0.3)

To view the results:

Co-expression Networks Network Analysis Network Biomarkers Enrichment Analysis **Clustering**

0 2 1

<p>cluster_1</p> <p>STAU1, TMEM70, RBM15, AI703397, RBM45, TBCB, FAM86B3P, MRO</p> <p>View Complex</p>
<p>cluster_2</p> <p>STAU1, AW205775, TMEM70, RBM15, AI703397, FKBP2, LOC105370629</p> <p>View Complex</p>
<p>cluster_3</p> <p>STAU1, TMEM70, BF510982, RBM15, AI703397, RBM45, TBCB, FAM86B3P</p> <p>View Complex</p>

At the Clustering tab you can view the different network clusters that are computed. You can also visualize them by clicking "View Complex".

The screenshot shows the InSyBio Suite interface. A modal window titled "Visualize Complex" is open, displaying a network graph with nodes labeled MRO, STAU1, FAM86B3P, AI703397, TMEM70, RBM15, TBCB, and RBM45. Below the graph, there are two buttons: "PNG image" and "Save Visualized Network". A "Close" button is located at the bottom right of the modal. The background interface shows a sidebar with cluster lists (cluster_1, cluster_2, cluster_3) and a "View Complex" button.

ncRNAseq Predict

The fifth step is the ncRNAseq Predict step. This step is also optional and with this step you can computationally predict potential miRNA targets at given Genes or Transcripts and given miRNAs. BLAST is performed in order to find potential target sites, and then the computational intelligent technique, which was applied for the prediction of miRNAs (hybrid combination of Genetic Algorithms and epsilon-SVRs), and 124 informative features are used in order to calculate a prediction score.

For this step:

- Select miRNAs and the Genes you want to search for potential targets by searching in our Database and adding them to the miRNA List and Genes List or add them manually to their Lists and separating them with commas.

The screenshot shows the "5. ncRNAseq Predict (Optional)" step. It includes a checkbox for "Do you want to do ncRNAseq target site prediction?" which is checked. Below this is a "Search miRNA" section with a dropdown menu labeled "Select miRNA" and an "Add to list" button. A "miRNAs List:" label is followed by a large empty text area for input.

To view the results:

At the tab miRNA Target Prediction the results are presented on your screen in a browse-able table, with each miRNA and gene pair in a row with their confidence score. By pressing Details at the Actions Column the specific scores between the miRNA and the gene's transcripts can be viewed. If no target sites are found "No targets found!" is presented at the score column. If one or more target sites are found you can view its UTR sequence, with the target sites of the miRNA highlighted. Multiple target sites are marked with green color and unique target sites are marked with light blue.

Testing Multi-biomarker Predictive Analytics Model

The sixth and final step is the Testing Multi-biomarker Predictive Analytics Model step. This step is also optional and allows the users to test the predictors that they have trained in the previous "training" step.

The first input file is the test dataset, which can be preprocessed or not preprocessed. The second input are the testset labels, which is optional. If the user inserts the testset labels then he'll receive as an output along with the predicted labels the performance metrics of the prediction.

It should be noted that the input dataset must have the format of the previous functionalities, that is it should have as rows the features and as columns the samples.

6. Testing Multi-biomarker Predictive Analytics Model (Optional)

Do you want to test Multi-biomarker Predictive Analytics Model?

Test set File 

Title 3:

Filename 3:

Test set Labels (Optional) 

Title 4:

Filename 4:

At the end, the user should click the Submit Job button to start the job.

To view the results:

At the Full Model Testing tab the user will be getting the following results. He'll view the predicted labels and the performance metrics. For the two class classification problem (two-class, multi-class) he'll view the test set accuracy, specificity, sensitivity, f1 score, f2 score and ROC AUC score. For the multiclass classification problem he'll view the test set accuracy, f1 score, precision, recall, f2 score and Manhattan Distance. For the regression he'll view the testset mean squared error, the test set relative absolute error and the root relative squared error.

How to get InSyBio Pipelines

To request a free one month full (evaluation) version of InSyBio Suite please email us at info@insybio.com.

To purchase InSyBio Interact commercial version 3.0 please contact us at sales@insybio.com.

About Us

InSyBio Ltd is a bioinformatics pioneer company (www.insybio.com) in precision medicine and nutrition, that focuses on developing computational frameworks and tools for the analysis of complex life-science and biological data in order to develop predictive integrated biomarkers (biomarkers of various categories) with increased prognostic and diagnostic aspects for the personalized Healthcare Industry.

InSyBio Suite consists of tools for providing integrated biological information from various sources, while at the same time it is empowered with robust, user-friendly and installation-free bioinformatics tools based on intelligent algorithms and methods.

COPYRIGHT NOTICE

External Publication of InSyBio Ltd - Any InSyBio information that is to be used in advertising, press releases, or promotional materials requires prior written approval from the InSyBio Ltd. A draft of the proposed document should accompany any such request. InSyBio Ltd reserves the right to deny approval of external usage for any reason.

Copyright 2022 InSyBio Ltd. Reproduction without written permission is completely forbidden.