# A New Cross-Platform Multi-Signature Classifier Approach To Predict Neuroblastoma Patients Outcome

Acquaviva M., Cornero A., Fardin P., Blengio F., Belli M.L., Varesio L.

*Laboratory of Molecular Biology, Gaslini Institute, Genoa, Italy;*

NETTAB Meeting 2011

# Background

- Neuroblastoma is the most common pediatric solid tumor of the sympathetic nervous system

- High variability in clinical behavior

- Reliable patients outcome predictions are often difficult to assess

Development of new predictive tools to assist established Neuroblastoma risk factors is mandatory

# Background

## Generation of a Multi-Signature Ensemble classifier
## for NeuroBlastoma patients outcome prediction (NB-MuSE-classifier)

➢  Ability to take into account the biological and prognostic information derived from a-priori knowledge (gene expression signatures).

➢ Possibility to combine different machine learning algorithms prediction power.

- 182 Neuroblastoma patients: U133Plus2 Gene Expression Profiles

- 35 Neuroblastoma related gene signatures from literature

- 22 Machine Learning algorithms tested

NB-MuSE-classifier

# Background

Previous  Work Results

| Single Signature Classifier | External validation Accuracy (%)^ | Paradigm |
|---|---|---|
| Chen 1 | 85 | BayesNet |
| Di Pietro 1 | 83 | BayesNet |
| Fredlund 1 | 80 | ClassificationViaRegression |
| Asgharzadeh 1 | 83 | ComplementNaiveBayes |
| Fransson 1 | 85 | ComplementNaiveBayes |
| De Preter 2 | 87 | IBk |
| Wei 1 | 83 | IBk |
| De Preter 1 | 83 | KStar |
| Oberthuer 1 | 87 | Logistic |
| Hahn 1 | 82 | MultiLayerPerceptron |
| McArdle 1 | 80 | MultiLayerPerceptron |
| Oe 1 | 80 | MultiLayerPerceptron |
| Nevo 2 | 87 | NaiveBayes |
| Shimada 1 | 80 | NBTree |
| Vermeulen 1 | 85 | NBTree |
| Ohira 1 | 85 | RandomForest |
| Fischer 1 | 81 | SimpleLogistic |
| Fardin 1 | 83 | Voted Perceptron |
| Nevo 1 | 80 | Voted Perceptron |
| **NB-MuSE** | **94** | **DecisonTable** |

Promising preliminary results

but

Evaluation Limited by relatively small test dataset

# Background

The availability of samples is one important limiting factor in developing reliable prognostic classifiers, especially for rare tumors such as neuroblastoma.

Neuroblastoma  repositories are often characterized by heterogeneous  high-througput datasets

A  multi-signature  classification  framework  which  can  use  different  array datasets (different gene expression platforms, arrayCGH, etc.)   to:

- improve biological and prognostic   a-priori  information
- extend the sample size  used for validation

## Generation of a cross-platform  NB-MuSE-classifier
### Exon 1.0 ST ⟺ U133Plus2 Data

➢ Ability to take into account the biological and prognostic information derived from a-priori knowledge (gene expression signatures).

➢ Possibility to combine different machine learning algorithms prediction power.

➢ Ability to be trained and tested on different type of high-throughput datasets (cross-platform feature), such as different gene expression arrays.
This feature permits the integration of heterogeneous datasets and the extension of sample size used for validation.

# Workflow for cross-platform NB-MuSE-classifier construction
## Exon 1.0 ST ⟺ U133Plus2 Data

**Phase 1**

| Signature 1 | Signature 2 | Signature 3 | Signature ... | Signature n |

Classifiers training & testing on DS1 (53 patients) on Alive/Dead labels
(KNIME-WEKA, 19 algorithms tested, Leave One Out C.V.)

| **Classifier** 1.1 <br> **Classifier** 1.2 <br> .... <br> **Classifier** 1.19 | **Classifier** 2.1 <br> **Classifier** 2.2 <br> .... <br> **Classifier** 2.19 | **Classifier** 3.1 <br> **Classifier** 3.2 <br> .... <br> **Classifier** 3.19 | **Classifier**...1 <br> **Classifier**...2 <br> .... <br> **Classifier**...19 | **Classifier** n.1 <br> **Classifier** n.2 <br> .... <br> **Classifier** n.19 |

Classifiers validation on external dataset DS2 (53 patients)
(cut-off on 80% prediction accuracy and best classifiers selection)

| Classifier 1 | **Classifier 2** | Classifier 3 | **Classifier ...** | **Classifier n** |

| | Predictions 2 | | Predictions... | Predictions n |

**Phase 2**

NB-MuSE-classifiers training & testing on predictions made on DS2 on Alive/Dead labels
(KNIME-WEKA, 19 algorithms tested, Leave One Out C.V.)

NB-MuSE-classifiers validation on predictions made on external dataset DS3
(56 patients)
(KNIME-WEKA)

**NB-MuSE-classifier Exon**
(best performing model)

*Cross-platform evaluation*
*Selection of the most reliable classifier*

**NB-MuSE-classifier U133Plus2**
(best performing model)

- 162 new neuroblastoma patients: Affymetrix Exon 1.0 ST Array (gene level).

| Clinical Characteristics | DS1 | DS2 | DS3 | Global Dataset |
|---|---|---|---|---|
| NB stage | % | % | % | % |
| st4s | 7.55 | 9.43 | 12.50 | 9.88 |
| st4 | 35.85 | 47.17 | 50.00 | 44.44 |
| st3 | 20.75 | 20.75 | 12.50 | 17.90 |
| st2 | 13.21 | 11.32 | 10.71 | 11.73 |
| st1 | 22.64 | 11.32 | 14.29 | 16.05 |
| age at diagnosis | | | | |
| <=1 y.o.a. | 43.40 | 28.30 | 42.86 | 38.27 |
| >1 y.o.a | 56.60 | 71.70 | 57.14 | 61.73 |
| mycn amplification | | | | |
| yes | 21.15 | 25.00 | 20.00 | 22.01 |
| no | 78.85 | 75.00 | 80.00 | 77.99 |
| overall survival | | | | |
| alive | 71.70 | 83.02 | 80.36 | 78.40 |
| dead | 28.30 | 16.98 | 19.64 | 21.60 |
| number of patients | 53 | 53 | 56 | 112 |

- 20 neuroblastoma related gene signatures

- 19 Machine Learning algorithms

| Machine learning algorithm | Cathegory |
|---|---|
| Bayes Logistic regression | Bayesan |
| BayesNet | Bayesan |
| Complement Naive Bayes | Bayesan |
| Naive Bayes | Bayesan |
| Logistic | Functions |
| Multi Layer Perceptron | Functions |
| Simple Logistic | Functions |
| Voted perceptron | Functions |
| IB1 | Lazy |
| IBK | Lazy |
| Kstar | Lazy |
| Bagging | Meta-learner |
| Classification via regression | Meta-learner |
| Decision Table | rules |
| Zero R | rules |
| J48 | tree |
| NBTree | tree |
| Random Forest | tree |
| Random Tree | tree |

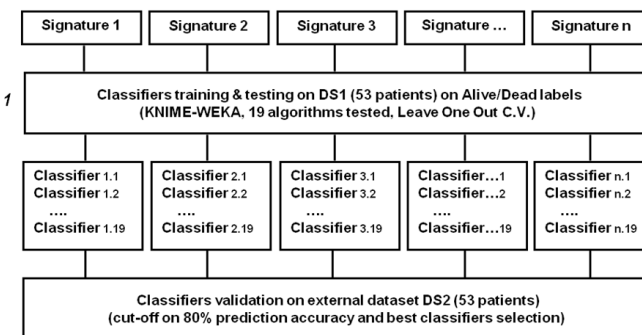# DS1  trainingdataset (53 patients): one for each signature
Expression dataset  used to train the single signature classifiers to predict patients overall  survival (Alive/Dead labels)

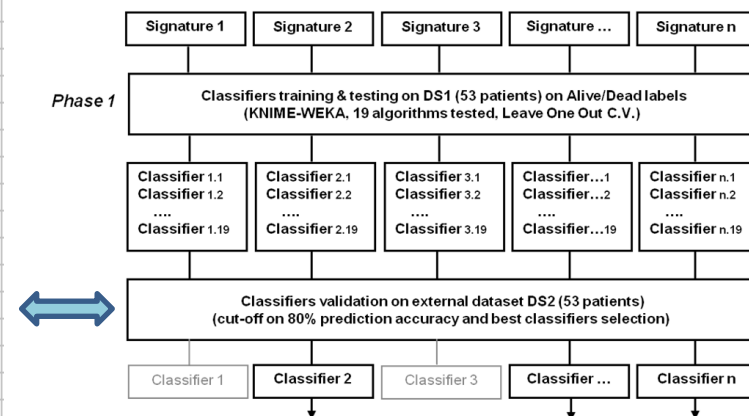| Patient ID | Label | 2319423 | 2319802 | 2319881 | 2395146 | 2395245 | 2395564 | 2395890 | 2395965 |
|---|---|---|---|---|---|---|---|---|---|
| nrc0002 | A | 136.2 | 848.9 | 105 | 69.1 | 402.3 | 31.2 | 218.5 | 336.9 |
| nrc0003 | A | 105.5 | 1062.1 | 136.2 | 75 | 359.8 | 35.1 | 283 | 470.9 |
| nrc0004 | A | 66.3 | 788.2 | 87.9 | 32.6 | 252.5 | 36.2 | 159.4 | 212 |
| nrc0005 | A | 85.1 | 813.2 | 93.1 | 32.4 | 173.7 | 37.8 | 87.8 | 286.5 |
| nrc0006 | A | 100.1 | 749.3 | 111.5 | 53 | 274.5 | 34.2 | 223.4 | 343.8 |
| nrc0007 | D | 79 | 1019 | 75.4 | 116.8 | 292.6 | 31 | 228.4 | 412.4 |
| nrc0010 | A | 74.2 | 1023.7 | 115.9 | 37.8 | 232.1 | 38.1 | 134.8 | 448.4 |
| nrc2536 | A | 81.7 | 1011.3 | 127.1 | 289.6 | 339.6 | 41.4 | 238.5 | 479.7 |
| nrc2537 | D | 64.9 | 783.4 | 112.2 | 36.9 | 229.1 | 39 | 156.9 | 128.4 |
| nrc2538 | A | 110.2 | 685.8 | 108.5 | 41.4 | 425 | 36.9 | 365.1 | 233 |
| nrc2541 | A | 97.3 | 1049.1 | 108.4 | 64.1 | 337.9 | 35.9 | 164.3 | 323.5 |
| nrc2542 | D | 112.1 | 372.4 | 87.5 | 61.5 | 127.5 | 38.9 | 75.3 | 208.1 |
| nrc2544 | A | 85.5 | 544.6 | 102.4 | 41 | 213.9 | 48.6 | 109.2 | 216.6 |
| nrc2545 | A | 97.9 | 1039.8 | 117.1 | 170.3 | 257.7 | 50.6 | 141.2 | 445.4 |
| nrc2546 | A | 96.8 | 763.1 | 84.2 | 104.3 | 310.5 | 63.4 | 211.9 | 402.6 |
| nrc2549 | A | 145.7 | 506.4 | 74 | 48.3 | 230.2 | 27.4 | 118 | 248.9 |
| nrc2550 | D | 57.4 | 788.5 | 99.1 | 42 | 182.3 | 32.6 | 118.9 | 184.2 |
| nrc2552 | A | 70.5 | 358.5 | 97.1 | 146.3 | 266.5 | 48.1 | 274.1 | 303.8 |
| nrc2555 | A | 110.4 | 646.7 | 91.4 | 68.8 | 267.9 | 43.9 | 160.8 | 292.8 |
| nrc2556 | A | 120.6 | 644.9 | 118.6 | 59.7 | 280.9 | 158.3 | 257.4 | 279.6 |
| nrc2557 | D | 58 | 434.4 | 78.5 | 35.2 | 210.7 | 40 | 138.4 | 199.3 |
| nrc2558 | D | 75.1 | 667.4 | 99.7 | 68.8 | 237.7 | 39.1 | 146.1 | 266.9 |
| nrc4001 | A | 72 | 919.6 | 136.6 | 30.2 | 352.1 | 48.4 | 229 | 401.2 |
| nrc4002 | A | 66.8 | 633.6 | 75.9 | 38.7 | 249.8 | 27.1 | 151.5 | 236.7 |
| nrc4005 | A | 91.7 | 1032 | 105.7 | 93.2 | 343.6 | 31.3 | 220.7 | 516 |
| nrc4006 | A | 106.7 | 738 | 109.4 | 31.8 | 383.5 | 26.7 | 214.1 | 386 |
| nrc4007 | A | 61.3 | 1169.8 | 131.3 | 46.1 | 305.3 | 25.1 | 262.3 | 555.2 |
| nrc4008 | A | 119 | 1035.3 | 146.6 | 129.8 | 406.3 | 30.4 | 324.7 | 432 |
| nrc4009 | A | 167 | 606.3 | 129.2 | 150.3 | 308.6 | 53.6 | 185.5 | 258.8 |
| nrc4010 | D | 114.8 | 1067.4 | 105.7 | 80.4 | 383.2 | 30.8 | 298.6 | 408 |

.... **Gene IDs**

**DS2 external validation dataset (53 patients): one for each signature**
Expression dataset used to validate the single signature classifiers trained on DS1

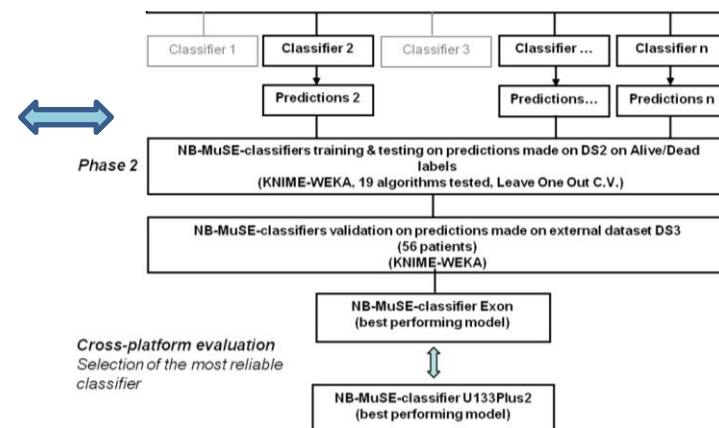| Patient ID | Label | 2320581 | 2339786 | 2340315 | 2356115 | 2460551 | 2492064 | 2515707 | 2528347 | .... **Gene IDs** |
|---|---|---|---|---|---|---|---|---|---|---|
| nrc6082 | A | 201.3 | 106.3 | 122.5 | 1271 | 239.3 | 118.7 | 131.8 | 91.1 | |
| nrc6083 | A | 161.9 | 90.3 | 131 | 2099.8 | 244.6 | 105.2 | 170.2 | 126.2 | |
| nrc6085 | A | 146.6 | 103.6 | 115.8 | 816.4 | 237.5 | 162.4 | 280.1 | 121.1 | |
| nrc6086 | A | 121.5 | 88.2 | 95.4 | 1109.3 | 207.5 | 156.9 | 185.6 | 111.9 | |
| nrc6087 | A | 122.7 | 159.2 | 133.9 | 745.2 | 246.9 | 64.2 | 194 | 89.1 | |
| nrc6088 | A | 112.5 | 88.8 | 134.3 | 681.6 | 222.9 | 103.9 | 81.9 | 133 | |
| nrc6090 | A | 161.7 | 68.5 | 183.7 | 652 | 264 | 101.8 | 123 | 113.9 | |
| nrc6091 | A | 209.9 | 125.8 | 174.9 | 1058.8 | 294.8 | 165.8 | 755.6 | 91.8 | |
| nrc6098 | A | 196 | 102.7 | 106.2 | 1002.6 | 223.5 | 155 | 161.9 | 151 | |
| nrc6103 | A | 278.3 | 96.9 | 137.2 | 1849.8 | 260.1 | 107.2 | 225.9 | 100.3 | |
| nrc6106 | D | 516.6 | 205 | 113.6 | 1313.8 | 264 | 143.8 | 257.2 | 79.4 | |
| nrc8113 | A | 164.8 | 107.5 | 114.5 | 1833.7 | 236.8 | 189.2 | 238.3 | 138.8 | |
| nrc8114 | A | 189 | 118.6 | 135.5 | 1091.6 | 214.9 | 186.7 | 732.1 | 132.1 | |
| nrc8119 | A | 290.7 | 119.2 | 121.5 | 1078.4 | 261.5 | 127.1 | 159.6 | 102.8 | |
| nrc8127 | A | 188 | 121.5 | 118.3 | 1669.1 | 282.3 | 132.4 | 262.1 | 116.6 | |
| nrc8134 | A | 126.4 | 95.6 | 137.8 | 1009.1 | 235.9 | 156.8 | 186.9 | 136.6 | |
| nrc8135 | D | 112.3 | 50.4 | 201.6 | 543.6 | 235.4 | 74.3 | 129 | 82.6 | |
| nrc8137 | A | 157.6 | 92.3 | 167.9 | 747.1 | 196.8 | 125.9 | 406.6 | 88.1 | |
| nrc8142 | D | 412.6 | 188.7 | 120.4 | 1098.6 | 259.1 | 140.8 | 447.5 | 75 | |
| nrc8145 | A | 224.9 | 126.9 | 110.1 | 1584.6 | 236.5 | 153.3 | 169.6 | 119.8 | |
| nrc8151 | A | 123.5 | 115.9 | 99.2 | 1368.2 | 249.7 | 132.2 | 133.6 | 132.7 | |
| nrc8152 | A | 64.2 | 54.8 | 178.4 | 240.2 | 161.8 | 70.2 | 64.9 | 70.9 | |
| nrc8153 | A | 115.2 | 75.1 | 145.2 | 1403.8 | 238.9 | 139.4 | 141.7 | 108.1 | |
| nrc8155 | A | 141.6 | 92 | 85.3 | 2073.3 | 204.7 | 124.4 | 298.1 | 111.8 | |
| nrc8157 | A | 251.8 | 135.8 | 194.3 | 1440.8 | 229.2 | 82.8 | 156.5 | 76.3 | |
| nrc8159 | A | 286.3 | 175.5 | 148 | 1465.2 | 376 | 153.3 | 102.1 | 96.5 | |
| nrc8160 | A | 327.2 | 234.6 | 139.6 | 1604.2 | 256 | 139.1 | 262.4 | 84 | |
| nrc8161 | A | 230.6 | 100.7 | 100.6 | 1375.5 | 253.4 | 130.8 | 302.1 | 100.5 | |
| nrc8167 | A | 261.6 | 108.3 | 119.9 | 1903.1 | 223.9 | 116.9 | 93.2 | 82.7 | |
| nrc8172 | D | 166.2 | 100.6 | 126.4 | 1185.1 | 270.9 | 263.2 | 457.4 | 125.9 | |

⋮



- Selection of the best single signature classifiers: evaluation prediction accuracy (>80%), sensitivity, specificity and recall.

**DS2 transformed in Prediction Matrix (NB-MusE-classifier training set, 53 patients)**
Prediction matrix assembled from the predictions performed on DS2 by the best single-signature classifiers selected during the first phase

| Patient ID | Label | Asgharza | Asgharza | Chen_1 N | DePreter_ | DePreter_ | DiPietro_ | Fardin_1_ | Fisher_1 |
|---|---|---|---|---|---|---|---|---|---|
| nrc6082 | A | A | A | A | A | A | A | A | A |
| nrc6083 | A | A | A | A | A | A | A | A | A |
| nrc6085 | A | A | A | A | A | A | A | A | A |
| nrc6086 | A | A | A | A | A | A | A | A | A |
| nrc6087 | A | A | A | A | A | A | A | A | A |
| nrc6088 | A | A | A | A | A | A | A | A | A |
| nrc6090 | A | A | A | A | A | A | A | A | A |
| nrc6091 | A | A | A | D | D | D | D | A | A |
| nrc6098 | A | A | A | A | A | A | A | A | A |
| nrc6103 | A | A | A | A | A | A | A | A | A |
| nrc6106 | D | A | A | A | A | A | A | D | A |
| nrc8113 | A | A | A | A | A | A | A | A | A |
| nrc8114 | A | A | A | A | D | D | D | A | A |
| nrc8119 | A | A | A | A | A | A | A | A | A |
| nrc8127 | A | A | A | A | A | A | A | A | A |
| nrc8134 | A | A | A | A | A | A | A | A | A |
| nrc8135 | D | A | A | D | D | D | D | A | A |
| nrc8137 | A | A | A | A | D | D | D | A | A |
| nrc8142 | D | A | A | D | D | A | D | A | D |
| nrc8145 | A | A | A | A | A | A | A | A | A |
| nrc8151 | A | A | A | A | A | A | A | A | A |
| nrc8152 | A | D | A | A | A | A | A | A | A |
| nrc8153 | A | A | A | A | A | A | A | A | A |
| nrc8155 | A | A | A | A | A | A | A | A | A |
| nrc8157 | A | D | A | A | A | A | A | A | A |
| nrc8159 | A | A | A | A | A | A | A | A | A |
| nrc8160 | A | A | A | A | A | A | A | D | A |
| nrc8161 | A | A | D | A | D | A | D | A | A |
| nrc8167 | A | A | A | A | A | A | A | A | A |
| nrc8172 | D | A | A | A | D | A | D | A | A |

.... **Signature IDs**



Classifier 1  Classifier 2  Classifier 3  Classifier ...  Classifier n
Predictions 2  Predictions...  Predictions n

*Phase 2* — NB-MuSE-classifiers training & testing on predictions made on DS2 on Alive/Dead labels (KNIME-WEKA, 19 algorithms tested, Leave One Out C.V.)

NB-MuSE-classifiers validation on predictions made on external dataset DS3 (56 patients) (KNIME-WEKA)

NB-MuSE-classifier Exon (best performing model)

*Cross-platform evaluation*
*Selection of the most reliable classifier*
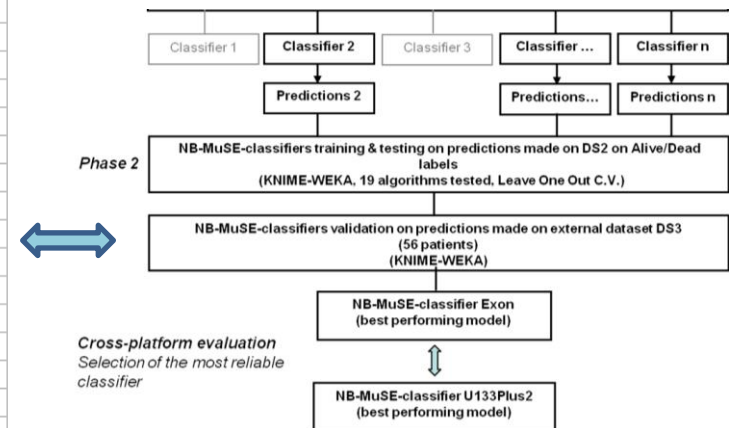
NB-MuSE-classifier U133Plus2 (best performing model)

**DS3 transformed in Prediction Matrix (NB-MusE-classifier validation set, 56 patients)**
Prediction matrix assembled from the predictions performed on DS3 by the best single-signature classifiers selected during the first phase

| Patient ID | Label | Asgharza | Asgharza | Chen_1_N | DePreter_ | DePreter_ | DiPietro_' | Fardin_1_ | Fisher_1_ |
|---|---|---|---|---|---|---|---|---|---|
| nrc4025 | A | A | A | A | A | A | A | A | A |
| nrc4026 | A | A | A | A | A | A | A | A | A |
| nrc4027 | A | A | A | A | A | A | A | A | A |
| nrc4034 | A | A | A | A | A | A | A | D | A |
| nrc4036 | A | A | A | A | A | A | A | A | A |
| nrc4037 | D | A | A | A | A | A | A | A | A |
| nrc4038 | D | A | A | A | A | A | A | A | A |
| nrc4039 | A | D | A | D | D | A | D | A | A |
| nrc4040 | A | D | A | D | D | D | D | A | A |
| nrc4041 | A | A | A | A | A | A | A | A | A |
| nrc4042 | D | D | A | D | D | D | D | D | A |
| nrc4043 | A | A | A | A | A | A | A | A | A |
| nrc4044 | A | D | D | D | D | D | D | D | A |
| nrc4045 | A | A | A | A | D | A | A | A | A |
| nrc4046 | D | A | A | A | D | D | D | A | D |
| nrc4047 | D | D | A | D | D | D | D | D | A |
| nrc4052 | A | A | A | A | A | A | A | A | A |
| nrc4053 | A | A | A | A | A | A | A | A | A |
| nrc4055 | A | A | A | A | A | A | A | A | A |
| nrc4056 | A | A | A | A | A | A | A | A | A |
| nrc4058 | A | A | A | D | A | A | A | A | A |
| nrc4059 | A | A | A | A | D | A | A | A | A |
| nrc4060 | A | A | A | D | A | A | A | A | A |
| nrc4062 | D | A | A | A | D | D | D | A | A |
| nrc4064 | D | D | A | A | D | A | A | A | A |
| nrc4066 | D | A | A | D | D | D | D | D | A |
| nrc4076 | A | A | A | A | A | A | A | A | A |
| nrc4077 | A | A | A | A | A | A | A | A | A |
| nrc4079 | A | A | A | A | A | A | A | A | A |
| nrc4081 | A | A | A | A | A | A | A | A | A |
| nrc4084 | A | A | A | A | A | A | A | A | A |
| nrc4087 | A | A | A | A | A | A | A | A | A |
| nrc4090 | A | A | A | A | D | A | D | D | A |

.... **Signature IDs**



Phase 2

Cross-platform evaluation
*Selection of the most reliable classifier*

DS2 and DS3 transformation steps are the core of the Cross Platform feature

# Preliminary Results

**Cross-platform evaluation of Multi-Signature Classifiers performance.**
The resulting multi-signature classifiers have been cross-tested on the external datasets and the relative performance has been evaluated. The Exon based NB-MuSE-classifier showed higher stability and reliability across the test datasets.

| | Label | TruePositives | FalsePositives | TrueNegatives | FalseNegatives | Recall | Precision | Sensitivity | Specifity | Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| **NB-MuSE-classifier Exon** | Alive | 42 | 6 | 5 | 3 | 0.933 | 0.875 | 0.933 | 0.455 | 0.839 |
| **K-Star** | Dead | 5 | 3 | 42 | 6 | 0.455 | 0.625 | 0.455 | 0.933 | |
| Test on U133Plus2 DS2 | Alive | 43 | 5 | 9 | 3 | 0.935 | 0.896 | 0.935 | 0.643 | 0.867 |
| 60 patients | Dead | 9 | 3 | 43 | 5 | 0.643 | 0.750 | 0.643 | 0.935 | |
| Test on U133Plus2 DS3 | Alive | 41 | 4 | 14 | 3 | 0.932 | 0.911 | 0.932 | 0.778 | 0.887 |
| 62 patients | Dead | 14 | 3 | 41 | 4 | 0.778 | 0.824 | 0.778 | 0.932 | |
| **NB-MuSE-classifier U133Plus2** | Alive | 43 | 3 | 15 | 1 | 0.977 | 0.935 | 0.977 | 0.833 | 0.940 |
| **Decision-Table** | Dead | 15 | 1 | 43 | 3 | 0.833 | 0.938 | 0.833 | 0.977 | |
| Test on Exon DS2 | Alive | 44 | 9 | 0 | 0 | 1 | 0.830 | 1 | 0 | 0.830 |
| 53 patients | Dead | 0 | 0 | 44 | 9 | 0 | | 0 | 1 | |
| Test on Exon DS3 | Alive | 44 | 10 | 1 | 1 | 0.978 | 0.815 | 0.978 | 0.091 | 0.804 |
| 56 patients | Dead | 1 | 1 | 44 | 10 | 0.091 | 0.500 | 0.091 | 0.978 | |

## Conclusions

- We developed a new classification model based on Exon expression data testable on the prediction matrices previously assembled from U133Plus2 data

- We successfully tested the cross-platform feature of NB-MuSE-classifier

- We have been able to evaluate and compare the two classifiers performance on respectively 109 and 122 (DS2+DS3)  new neuroblastoma patients.

**Future Directions**

- Optimization of classifiers learning parameters and cross-validation set-ups

- Optimization of a-priori information selection (NB-related signatures)

- Test on randomized datasets

- Integration of arrayCGH and miRNA datasets