

# **BLOBREC – a blood-based test for breast cancer**

**Analyses of test properties in the years before and after diagnosis**

**Note no.**

**SAMBA/27/19**

**Authors**

**Eiliv Lund, Marit Holden, Lars Holden**

**Date**

**September 2019**

### Authors

Eiliv Lund (1)

Marit Holden (2)

Lars Holden (2)

- 1) UiT The Arctic University of Norway, Tromsø, Norway
- 2) Norsk Regnesentral, Oslo, Norway

### Norsk Regnesentral

Norsk Regnesentral (Norwegian Computing Center, NR) is a private, independent, non-profit foundation established in 1952. NR carries out contract research and development projects in information and communication technology and applied statistical-mathematical modelling. The clients include a broad range of industrial, commercial and public service organisations in the national as well as the international market. Our scientific and technical capabilities are further developed in co-operation with The Research Council of Norway and key customers. The results of our projects may take the form of reports, software, prototypes, and short courses. A proof of the confidence and appreciation our clients have in us is given by the fact that most of our new contracts are signed with previous customers.

<b>Title</b>	<b>BLOBREC – a blood-based test for breast cancer</b>
<b>Authors</b>	<b>Eiliv Lund, Marit Holden, Lars Holden</b>
Date	September 2019
Year	2019
Publication number	SAMBA/27/19

### **Abstract**

BLOBREC is a test for distinguishing breast-cancer patients from population-based controls described by Dumeaux et al. Previously, we performed a quality control of the methods and procedures used for developing this test, and our analyses confirmed the results obtained by Dumeaux et al.

The aim of these analyses were to study the properties of the BLOBREC test in the years before or after diagnosis, and compare with time of diagnosis. In addition, we looked at effects of parity in controls, and also used a clinical stress study.

We had a case-control design with 539 pairs before diagnosis, 59 at and 429 after diagnosis. In the controls taken from the NOWAC postgenome biobank, we found no difference in percentage false positives (%P) between the pre- and postdiagnostic controls (37% and 34% respectively). The %P were similar to the case-control study at time of diagnosis; 37%. The %P for cases were except for one year before diagnosis similar to the controls pre- and postdiagnostic. Additionally, we found a weak, non-significant increase in %P for controls with many children. The stress data originated from the “second look” at one single centre in the national screening program for breast cancer. The %P (the per cent of positive tests) were lower both for cases and controls than for the original case-control series. These data were collected under more stringent conditions.

Conclusion: the BLOBREC showed higher %P for cases at diagnosis than either before or after, while the %P for controls remained identical. As previous, the %P was higher for metastatic breast cancer.

BLOBREC might be improved through a more stringent sampling of both cases and controls.

Keywords	Gene expression test; Breast cancer; Blood; Screening program; Diagnostic test; Naïve Bayes;
Target group	Clinical medicine; Systems epidemiology
Availability	Open
Project number	NR project number 220 785
Research field	MH, Bioinformatics
Number of pages	17
© Copyright	Norsk Regnesentral



# Table of Content

<b>1</b>	<b>Introduction.....</b>	<b>7</b>
<b>2</b>	<b>The BLOBREC test .....</b>	<b>7</b>
<b>3</b>	<b>Data .....</b>	<b>8</b>
3.1	Preprocessing the gene expression datasets .....	8
3.2	The prospective dataset.....	8
3.3	The postdiagnostic dataset .....	8
3.4	The stress dataset.....	9
<b>4</b>	<b>Results .....</b>	<b>9</b>
<b>5</b>	<b>Conclusion .....</b>	<b>12</b>
<b>6</b>	<b>References .....</b>	<b>12</b>
<b>7</b>	<b>Appendix .....</b>	<b>13</b>
7.1	The set of 345 differentially expressed gene .....	13
7.2	The 50 genes included in the BLOBREC test .....	13
7.3	Prediction results stratified on screening status .....	14
7.4	Selecting 10 of 50 genes for the predictor.....	15



# 1 Introduction

BLOBREC is a test for distinguishing breast-cancer patients from population-based controls described by Dumeaux et al. in [1]. In the note [2], we performed a quality control of the methods and procedures used for developing the test, and our analyses confirmed the results obtained by Dumeaux et al. The BLOBREC test was developed based on gene expression in blood at the time of diagnosis. In this note, we will use the BLOBREC test on three datasets, one with gene expression in blood before the time of diagnosis, one with gene expression in blood after the time of diagnosis, and one with gene expression in blood at the time of diagnosis.

In Section 2 we describe the BLOBREC test, while the datasets are presented in Section 3. Results are summarized in Section 4.

## 2 The BLOBREC test

In [1], three datasets with gene expression in blood at the time of diagnosis, CC1, CC2 and CC3, were used for defining the BLOBREC test. First, the CC1 and CC2 datasets were used for finding a set of 345 genes that were differentially expressed in both datasets (FDR q-value < 0.005). The defined set of differentially expressed genes consisted of 345 genes (Appendix 7.1). The CC3 dataset was then used to select 50 of the 345 genes for the predictor that separates cases from controls (Appendix 7.2).

The BLOBREC test predicts cancer or not using a Naïve Bayes method. The predictions made are based on data in the CC3 dataset that consists of  $N=118$  individuals and  $M = 50$  genes, where each individual is either a case or a control. In the CC3 dataset, there are 59 cases and 59 controls.

Let  $x_{ij}$  be the gene expression data on log-scale,  $i = 1, \dots, M$  and  $j = 1, \dots, N$ . Let group 0 consist of individuals without cancer (controls), and group 1 of individuals with cancer (cases). A new individual with data  $y_i$ ,  $i = 1, \dots, M$ , is predicted to be a case, i.e. to have cancer, if

$$\frac{p}{1-p} \prod_{i=1}^M \frac{\varphi(y_i; \mu_i^1, \sigma_i^1)}{\varphi(y_i; \mu_i^0, \sigma_i^0)} > 1,$$

and to be a control, i.e. to be without cancer, otherwise, where

- $p = \sum_{j=1, \dots, N} \frac{g_j}{N}$ , where  $g_j = 0$  if individual  $j$  of the CC3 dataset belongs to group 0, and 1 if individual  $j$  belongs to group 1.
- $\varphi$  is the probability density of the normal density.
- $\mu_i^1$  and  $\sigma_i^1$  are the mean and standard deviation computed from  $x_{ij_1}$ , for all  $j_1 \in \{1, \dots, N\}$  such that  $g_{j_1} = 1$ . Similarly,  $\mu_i^0$  and  $\sigma_i^0$  are the mean and standard deviation computed from  $x_{ij_0}$ , for all  $j_0 \in \{1, \dots, N\}$  such that  $g_{j_0} = 0$ .

Note that if some of the 50 genes are not included in the dataset with the new individual these genes are excluded from the test. Also note that before using the test on data from a new dataset, the mean and standard deviations of the test are adjusted so that the new mean and

standard deviations for the controls in the BLOBREC test becomes equal to the mean and standard deviations for the controls in the new dataset. More precisely, we compute the adjusted mean and standard deviations  $\mu_i^{0*}$ ,  $\mu_i^{1*}$ ,  $\sigma_i^{0*}$  and  $\sigma_i^{1*}$  as:

$$\mu_i^{0*} = \mu_i^0 - \mu_i^0 + \mu_i^{0N}, \quad \mu_i^{1*} = \mu_i^1 - \mu_i^0 + \mu_i^{0N}, \quad \sigma_i^{0*} = \sigma_i^0 \cdot \frac{\sigma_i^{0N}}{\sigma_i^0}, \quad \text{and} \quad \sigma_i^{1*} = \sigma_i^1 \cdot \frac{\sigma_i^{0N}}{\sigma_i^0},$$

where  $\mu_i^{0N}$  and  $\sigma_i^{0N}$  are the sample mean and standard deviation, respectively, for gene  $i$  for the controls in the new dataset.

We will also define a new test, the extended BLOBREC test, which is equal to the BLOBREC test except that it is based on the set of 345 genes instead of only 50 genes.

## 3 Data

Three breast cancer datasets are available: Two of them were based on the NOWAC postgenome biobank; the prospective dataset and the postdiagnostic dataset. The stress dataset was collected at the “second look” in one clinic. Each of these datasets will be described in more detail below.

### 3.1 Preprocessing the gene expression datasets

Each dataset was background corrected using negative control probes,  $\log_2$  transformed using a variance stabilizing technique [2], and quantile normalized. We retained probes present in at least 70% of the samples. If a gene was represented with more than one probe, the average expression of the probes was used as expression value for the gene. The probes were translated to genes using the lumiHumanIDMapping [3].

### 3.2 The prospective dataset

After removing technical outliers and controls that were diagnosed with cancer, the dataset consists of data from 539 case-control pairs that are from three different runs / batches. After preprocessing the dataset consists of 8155 genes. The cases are diagnosed with metastases, without metastases or with in-situ tumors. Note that the adjusted mean and standard deviations,  $\mu_i^{0*}$ ,  $\mu_i^{1*}$ ,  $\sigma_i^{0*}$  and  $\sigma_i^{1*}$ , are computed separately for each of the three different runs.

### 3.3 The postdiagnostic dataset

After removing technical outliers and controls that were diagnosed with cancer, the dataset consists of data from 429 case-control pairs. After preprocessing the dataset consists of 8400 genes. We updated the metastasis and follow-up time information for cases that were diagnosed with metastases (6) or new breast cancer (10) after the initial/first breast cancer diagnosis and before the blood sampling<sup>1</sup>. After removing 7 case-control pairs where the case was diagnosed with other cancers after the initial/first breast cancer diagnosis and before the

---

<sup>1</sup> For 17 of the 429 cases such information was not available.



blood sampling<sup>2</sup>, and 7 cases where the metastases status of the case was unknown, 415 cases remained in the dataset.

### 3.4 The stress dataset

The stress dataset was produced for examining the effect of stress when the women returned for a “second look” after positive findings in the ordinary mammographic screening. In this note we will use the dataset as a validation set for the BLOBREC test. The blood sampling procedure is such that all cases in the stress dataset were stressed at time of blood sampling since that was done at the time of the diagnostic biopsy, while the controls had nothing to be anxious about.

The stress dataset consisted of 40 case-control pairs and 47 323 probes. In this dataset some cases have cancer (12), while the remaining cases (28) and controls (40) are healthy. In addition to the 40 case-control pairs, the stress dataset also contains 16 samples that are selected from a pooled sample based on blood samples from 16 individuals. The data are preprocessed using the procedure described in Section 3.1, but adapted so that the probes present in the CC3 dataset were used when mapping from probes to genes. After preprocessing, the dataset consisted of the same 9 936 genes as the CC3 dataset.

## 4 Results

Table 1 below shows prediction results for the CC3 dataset. These results are taken from Table 7 in [2]. When 50 of the 341 genes are included in the test, the disease status of an individual *i* is predicted using a 50-gene best predictor that is selected using a dataset consisting of all individuals in the CC3 dataset except individual *i*.

Table 1 Prediction results for the 59 cases and 59 controls in the CC3 dataset using leave-one-out prediction both when selecting genes for the test and when computing mean and standard deviation. The results are taken from Table 7 in [2]. “N” is the number of negative tests, i.e. the number of individuals that are classified as not cancer, while “P” is the number of positive tests, i.e. the number of individuals that are classified as having cancer. “%P” is the per cent of positive tests, i.e.  $P/(P+N)$ .

	Number of genes included in test	N	P	%P
Controls	The 341 of the 345 genes present in the CC3 dataset	37	22	37
	50 of 341 genes <sup>3</sup> , simulation 1	37	22	37
	50 of 341 genes, simulation 2	36	23	39
Cases	The 341 of the 345 genes present in the CC3 dataset	15	44	75
	50 of 341 genes, simulation 1	12	47	80
	50 of 341 genes, simulation 2	15	44	75

Tables 2-6 show prediction results for the prospective, postdiagnostic and stress datasets. For each of the three datasets, we give results for tests with the 50 and 345 genes that were identified in [1]. Only genes that are included in the preprocessed dataset are used in the tests.

<sup>2</sup> For 17 of the 429 cases such information was not available.

<sup>3</sup> The 50 genes have been selected as described in [1] and [2]

Table 2 Prediction results for the 539 cases and 539 controls in the prospective dataset. "N", "P" and "%P" are defined as in Table 1. Table 6 show the prediction results stratified on screening status (screen-detected or clinically detected cancer).

Controls	N	P	%P	The BLOBREC test (41 genes included in the test)							
	338	201	37	In-situ			Without Metastases			With Metastases	
Cases	N	P	%P	N	P	%P	N	P	%P		
Year 1	17	4	19	40	24	38	10	12	55		
Year 2	7	4	36	47	20	30	19	8	30		
Year 3	12	6	33	39	27	41	19	7	27		
Year 4	10	10	50	50	17	25	24	11	31		
Year 5	4	4	50	27	14	34	15	4	21		
Year 6	1	0	0	5	8	62	3	4	57		
Year 7	0	0	0	1	3	75	1	0	0		
Year 8	0	0	0	1	0	0	0	0	0		

Controls	N	P	%P	The extended BLOBREC test (300 genes included in the test)							
	342	197	37	In-situ			Without Metastases			With Metastases	
Cases	N	P	%P	N	P	%P	N	P	%P		
Year 1	17	4	19	41	23	36	10	12	55		
Year 2	6	5	45	45	22	33	19	8	30		
Year 3	12	6	33	39	27	41	17	9	35		
Year 4	10	10	50	52	15	22	27	8	23		
Year 5	4	4	50	27	14	34	14	5	26		
Year 6	1	0	0	5	8	62	4	3	43		
Year 7	0	0	0	1	3	75	1	0	0		
Year 8	0	0	0	0	1	100	0	0	0		

The %P of controls were identical to the %P of controls in CC3. In situ had similar %P as the controls. For metastatic cases the %P were highest for metastatic cases. There was no effect of the extended BLOBREC test with regard to test properties.

Table 3 Prediction results for the 429 cases and 429 controls in the postdiagnostic dataset. "N", "P" and "%P" are defined as in Table 1. We know the metastasis status (in-situ, without metastases, with metastases) for 415 of the 429 cases.

Controls	N	P	%P	The BLOBREC test (42 genes included in the test <sup>4</sup> )							
	282	147	34	In-situ			Without Metastases			With Metastases	
Cases	N	P	%P	N	P	%P	N	P	%P		
Year 1	5	2	29	17	16	48	7	6	46		
Year 2	4	3	43	29	25	46	14	10	42		
Year 3	6	3	33	21	18	46	10	12	55		
Year 4	6	3	33	19	13	41	13	9	41		
Year 5	4	0	0	11	8	42	7	6	46		
Year 6	4	2	33	11	7	39	11	9	45		
Year 7	5	1	17	15	8	35	11	11	50		
Year 8	1	0	0	6	4	40	0	2	100		

Controls	N	P	%P	The extended BLOBREC test (301 genes included in the test <sup>5</sup> )							
	278	151	35	In-situ			Without Metastases			With Metastases	
Cases	N	P	%P	N	P	%P	N	P	%P		
Year 1	5	2	29	16	17	52	8	5	38		
Year 2	5	2	29	27	27	50	12	12	50		
Year 3	6	3	33	25	14	36	10	12	55		
Year 4	5	4	44	23	9	28	13	9	41		
Year 5	4	0	0	11	8	42	9	4	31		
Year 6	5	1	17	12	6	33	11	9	45		
Year 7	5	1	17	17	6	26	11	11	50		
Year 8	1	0	0	6	4	40	0	2	100		

The %P for controls in the postdiagnostic dataset were similar to the two previous series. In situ was similar to the controls. There was no relationship between time after diagnosis and percent positives.

<sup>4</sup> These 301 include the 42 genes for the prospective dataset and in addition the gene ZNF266.

<sup>5</sup> These 301 include the 300 genes for the prospective dataset and in addition the gene ZNF266.

Table 4 Prediction results for controls in the prospective and postdiagnostic dataset stratified on parity. "N", "P" and "%P" are defined as in Table 1.

41/42 genes included in the test	539			429			968		
	Prospective			Postdiagnostic			Both		
Parity	N	P	%P	N	P	%P	N	P	%P
0	59	24	29	34	19	36	93	43	32
1-3	259	154	37	223	116	34	482	270	36
4-6	24	19	44	25	12	32	49	31	39

The additional analyses of the %P among controls according to parity showed a weak, non-significant trend with higher parity (Fisher's test,  $p > 0.3$ ).

Table 5 Prediction results for the 12 cases with cancer, 28 individuals with biopsy, but without cancer, and the 16 pooled samples, the 40 controls, in the stress dataset. "N", "P" and "%P" are defined as in Table 1.

	The BLOBREC test (48 genes included in the test)			The extended BLOBREC test (327 genes included in the test)		
	N	P	%P	N	P	%P
Pools	15	1	6	16	0	0
Controls	34	6	15	35	5	12
With biopsy, without cancer	19	9	32	20	8	29
Cases with cancer	5	7	58	5	7	58

The stress test differs from the other series as it was collected in one surgical department. The pooled controls showed a %P equal to zero. The controls had a much lower %P than in the previous material, but the number is small. The %P of the cases was lower than in the case-control study, which can be due to the small tumours found in the screening.

## 5 Conclusion

The analyses has demonstrated a high %P for controls taken from the NOWAC postgenome biobank leaving the test difficult to use. However, the stress test %P could indicate that in a clinical situation the test might prove to work fairly well.

The prediction results (Section 7.3) demonstrates that the BLOBREC primarily is a test for metastatic clinical cancer.

## 6 References

[1] Vanessa Dumeaux, Josie Ursini-Siegel, Arnar Flatberg, Hans E. Fjosne, Jan-Ole Frantzen, Marit Muri Holmen, Enno Rodegerdts, Ellen Schlichting and Eiliv Lund. Peripheral blood cells inform on the presence of breast cancer: A population-based case-control study. *Int. J. Cancer*: 136, 656–667 (2015).

[2] Marit Holden, Clara-Cecilie Günther and Lars Holden. Verification of a blood-based test for breast cancer (BLOBREC): Distinguishing breast-cancer patients from population-based controls. NR note SAMBA/33/15, 2015.

## 7 Appendix

### 7.1 The set of 345 differentially expressed gene

ABHD10 ABI3 ABR ACTB ACTG1 AGPAT3 AIFM1 AIMP2 ALG8 ALKBH5 ANXA1 ANXA2 ANXA5  
APBB3 APEX2 APOBEC3C APOL3 APP AQP9 ARCN1 ARF3 ARFIP1 ARHGAP1 ARHGAP17  
ARHGDI1 ARPC5L ASPHD2 ATF5 ATG12 ATP1A1 ATP2B4 ATP5B AXIN1 BHLHE40 BRE C11orf57  
C12orf47 C14orf2 C16orf72 C17orf63 C18orf8 C20orf4 C21orf33 C4orf3 CALHM2 CALM1  
CAMLG CAPNS1 CAPZA2 CASP4 CBARA1 CCDC86 CCDC92 CCDC97 CCT7 CD74 CDK19 CDKN1C  
CECR1 CKAP5 CLN5 CLPTM1L CLSTN1 CNDP2 COBRA1 COMT COPB2 COPS7B CPD CPEB3 CRKL  
CS CSK CSTF2 CTBP1 CTCF CTNNA1 CTNNA2 CTNNA3 CTNNA4 CTNNA5 CTNNA6 CTNNA7  
CTNNA8 CTNNA9 CTNNA10 CTNNA11 CTNNA12 CTNNA13 CTNNA14 CTNNA15 CTNNA16  
CTNNA17 CTNNA18 CTNNA19 CTNNA20 CTNNA21 CTNNA22 CTNNA23 CTNNA24  
CTNNA25 CTNNA26 CTNNA27 CTNNA28 CTNNA29 CTNNA30 CTNNA31 CTNNA32  
CTNNA33 CTNNA34 CTNNA35 CTNNA36 CTNNA37 CTNNA38 CTNNA39 CTNNA40  
CTNNA41 CTNNA42 CTNNA43 CTNNA44 CTNNA45 CTNNA46 CTNNA47 CTNNA48  
CTNNA49 CTNNA50 CTNNA51 CTNNA52 CTNNA53 CTNNA54 CTNNA55 CTNNA56  
CTNNA57 CTNNA58 CTNNA59 CTNNA60 CTNNA61 CTNNA62 CTNNA63 CTNNA64  
CTNNA65 CTNNA66 CTNNA67 CTNNA68 CTNNA69 CTNNA70 CTNNA71 CTNNA72  
CTNNA73 CTNNA74 CTNNA75 CTNNA76 CTNNA77 CTNNA78 CTNNA79 CTNNA80  
CTNNA81 CTNNA82 CTNNA83 CTNNA84 CTNNA85 CTNNA86 CTNNA87 CTNNA88  
CTNNA89 CTNNA90 CTNNA91 CTNNA92 CTNNA93 CTNNA94 CTNNA95 CTNNA96  
CTNNA97 CTNNA98 CTNNA99 CTNNA100 CTNNA101 CTNNA102 CTNNA103  
CTNNA104 CTNNA105 CTNNA106 CTNNA107 CTNNA108 CTNNA109 CTNNA110  
CTNNA111 CTNNA112 CTNNA113 CTNNA114 CTNNA115 CTNNA116 CTNNA117  
CTNNA118 CTNNA119 CTNNA120 CTNNA121 CTNNA122 CTNNA123 CTNNA124  
CTNNA125 CTNNA126 CTNNA127 CTNNA128 CTNNA129 CTNNA130 CTNNA131  
CTNNA132 CTNNA133 CTNNA134 CTNNA135 CTNNA136 CTNNA137 CTNNA138  
CTNNA139 CTNNA140 CTNNA141 CTNNA142 CTNNA143 CTNNA144 CTNNA145  
CTNNA146 CTNNA147 CTNNA148 CTNNA149 CTNNA150 CTNNA151 CTNNA152  
CTNNA153 CTNNA154 CTNNA155 CTNNA156 CTNNA157 CTNNA158 CTNNA159  
CTNNA160 CTNNA161 CTNNA162 CTNNA163 CTNNA164 CTNNA165 CTNNA166  
CTNNA167 CTNNA168 CTNNA169 CTNNA170 CTNNA171 CTNNA172 CTNNA173  
CTNNA174 CTNNA175 CTNNA176 CTNNA177 CTNNA178 CTNNA179 CTNNA180  
CTNNA181 CTNNA182 CTNNA183 CTNNA184 CTNNA185 CTNNA186 CTNNA187  
CTNNA188 CTNNA189 CTNNA190 CTNNA191 CTNNA192 CTNNA193 CTNNA194  
CTNNA195 CTNNA196 CTNNA197 CTNNA198 CTNNA199 CTNNA200 CTNNA201  
CTNNA202 CTNNA203 CTNNA204 CTNNA205 CTNNA206 CTNNA207 CTNNA208  
CTNNA209 CTNNA210 CTNNA211 CTNNA212 CTNNA213 CTNNA214 CTNNA215  
CTNNA216 CTNNA217 CTNNA218 CTNNA219 CTNNA220 CTNNA221 CTNNA222  
CTNNA223 CTNNA224 CTNNA225 CTNNA226 CTNNA227 CTNNA228 CTNNA229  
CTNNA230 CTNNA231 CTNNA232 CTNNA233 CTNNA234 CTNNA235 CTNNA236  
CTNNA237 CTNNA238 CTNNA239 CTNNA240 CTNNA241 CTNNA242 CTNNA243  
CTNNA244 CTNNA245 CTNNA246 CTNNA247 CTNNA248 CTNNA249 CTNNA250  
CTNNA251 CTNNA252 CTNNA253 CTNNA254 CTNNA255 CTNNA256 CTNNA257  
CTNNA258 CTNNA259 CTNNA260 CTNNA261 CTNNA262 CTNNA263 CTNNA264  
CTNNA265 CTNNA266 CTNNA267 CTNNA268 CTNNA269 CTNNA270 CTNNA271  
CTNNA272 CTNNA273 CTNNA274 CTNNA275 CTNNA276 CTNNA277 CTNNA278  
CTNNA279 CTNNA280 CTNNA281 CTNNA282 CTNNA283 CTNNA284 CTNNA285  
CTNNA286 CTNNA287 CTNNA288 CTNNA289 CTNNA290 CTNNA291 CTNNA292  
CTNNA293 CTNNA294 CTNNA295 CTNNA296 CTNNA297 CTNNA298 CTNNA299  
CTNNA300 CTNNA301 CTNNA302 CTNNA303 CTNNA304 CTNNA305 CTNNA306  
CTNNA307 CTNNA308 CTNNA309 CTNNA310 CTNNA311 CTNNA312 CTNNA313  
CTNNA314 CTNNA315 CTNNA316 CTNNA317 CTNNA318 CTNNA319 CTNNA320  
CTNNA321 CTNNA322 CTNNA323 CTNNA324 CTNNA325 CTNNA326 CTNNA327  
CTNNA328 CTNNA329 CTNNA330 CTNNA331 CTNNA332 CTNNA333 CTNNA334  
CTNNA335 CTNNA336 CTNNA337 CTNNA338 CTNNA339 CTNNA340 CTNNA341  
CTNNA342 CTNNA343 CTNNA344 CTNNA345 CTNNA346 CTNNA347 CTNNA348  
CTNNA349 CTNNA350 CTNNA351 CTNNA352 CTNNA353 CTNNA354 CTNNA355  
CTNNA356 CTNNA357 CTNNA358 CTNNA359 CTNNA360 CTNNA361 CTNNA362  
CTNNA363 CTNNA364 CTNNA365 CTNNA366 CTNNA367 CTNNA368 CTNNA369  
CTNNA370 CTNNA371 CTNNA372 CTNNA373 CTNNA374 CTNNA375 CTNNA376  
CTNNA377 CTNNA378 CTNNA379 CTNNA380 CTNNA381 CTNNA382 CTNNA383  
CTNNA384 CTNNA385 CTNNA386 CTNNA387 CTNNA388 CTNNA389 CTNNA390  
CTNNA391 CTNNA392 CTNNA393 CTNNA394 CTNNA395 CTNNA396 CTNNA397  
CTNNA398 CTNNA399 CTNNA400 CTNNA401 CTNNA402 CTNNA403 CTNNA404  
CTNNA405 CTNNA406 CTNNA407 CTNNA408 CTNNA409 CTNNA410 CTNNA411  
CTNNA412 CTNNA413 CTNNA414 CTNNA415 CTNNA416 CTNNA417 CTNNA418  
CTNNA419 CTNNA420 CTNNA421 CTNNA422 CTNNA423 CTNNA424 CTNNA425  
CTNNA426 CTNNA427 CTNNA428 CTNNA429 CTNNA430 CTNNA431 CTNNA432  
CTNNA433 CTNNA434 CTNNA435 CTNNA436 CTNNA437 CTNNA438 CTNNA439  
CTNNA440 CTNNA441 CTNNA442 CTNNA443 CTNNA444 CTNNA445 CTNNA446  
CTNNA447 CTNNA448 CTNNA449 CTNNA450 CTNNA451 CTNNA452 CTNNA453  
CTNNA454 CTNNA455 CTNNA456 CTNNA457 CTNNA458 CTNNA459 CTNNA460  
CTNNA461 CTNNA462 CTNNA463 CTNNA464 CTNNA465 CTNNA466 CTNNA467  
CTNNA468 CTNNA469 CTNNA470 CTNNA471 CTNNA472 CTNNA473 CTNNA474  
CTNNA475 CTNNA476 CTNNA477 CTNNA478 CTNNA479 CTNNA480 CTNNA481  
CTNNA482 CTNNA483 CTNNA484 CTNNA485 CTNNA486 CTNNA487 CTNNA488  
CTNNA489 CTNNA490 CTNNA491 CTNNA492 CTNNA493 CTNNA494 CTNNA495  
CTNNA496 CTNNA497 CTNNA498 CTNNA499 CTNNA500

### 7.2 The 50 genes included in the BLOBREC test

ABHD10 AIMP2 APOL3 ARHGAP1 ARHGAP17 C14orf2 C16orf72 C18orf8 CAMLG CD74 CECR1  
COPS7B DCP1A DENR DHX40 DYNLRB1 EP300 FIP1L1 FRYL GLRX GMFG GSTP1 HNRNPAB HPS6  
IK JAK1 KARS KIAA0930 KIAA1310 LOC100290936 LRFN3 PAPOLA PHF5A PPP2R5A RASSF5  
RPL21 RPL5 RPS3A S100A8 SDHA SEC31A SMARCAL1 SP2 TMEM39B TUBB UBAC2 WBP11  
YWHAB ZNF266 ZNF319

### 7.3 Prediction results stratified on screening status

Table 6 Prediction results for the 539 cases and 539 controls in the prospective dataset. "N", "P and "%P" are defined as in Table 1.

Controls	N	P	%P	The BLOBREC test (41 genes included in the test)														
	338	201	37	Screening									Clinical					
Cases	In-situ			Without Met.			With Met.			In-situ			Without Met.			With Met.		
	N	P	%P	N	P	%P	N	P	%P	N	P	%P	N	P	%P	N	P	%P
Year 1	12	4	25	33	19	37	8	4	33	5	0	0	8	4	33	2	8	80
Year 2	4	5	56	40	16	29	13	6	32	2	0	0	5	6	55	6	2	25
Year 3	10	5	33	33	25	43	16	8	33	2	1	33	6	2	25	1	1	50
Year 4	10	9	47	40	13	25	24	6	20	1	0	0	12	2	14	3	2	40
Year 5	3	4	57	23	12	34	10	5	33	1	0	0	4	2	33	4	0	0
Year 6	0	0	0	4	6	60	4	2	33	1	0	0	1	2	67	1	0	0
Year 7	0	0	0	1	2	67	1	0	0	0	0	0	1	0	0	0	0	0
Year 8	0	0	0	0	1	100	0	0	0	0	0	0	0	0	0	0	0	0

Controls	N	P	%P	The extended BLOBREC test (300 genes included in the test)														
	342	197	37	Screening									Clinical					
Cases	In-situ			Without Met.			With Met.			In-situ			Without Met.			With Met.		
	N	P	%P	N	P	%P	N	P	%P	N	P	%P	N	P	%P	N	P	%P
Year 1	12	4	25	32	20	38	8	4	33	5	0	0	8	4	33	2	8	80
Year 2	5	4	44	43	13	23	13	6	32	2	0	0	4	7	64	6	2	25
Year 3	10	5	33	33	25	43	18	6	25	2	1	33	6	2	25	1	1	50
Year 4	10	9	47	38	15	28	21	9	30	1	0	0	12	2	14	3	2	40
Year 5	3	4	57	23	12	34	11	4	27	1	0	0	4	2	33	4	0	0
Year 6	0	0	0	4	6	60	3	3	50	1	0	0	1	2	67	1	0	0
Year 7	0	0	0	1	2	67	1	0	0	0	0	0	1	0	0	0	0	0
Year 8	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0

## 7.4 Selecting 10 of 50 genes for the predictor

The CC3 dataset is preprocessed as described in [1]. In [2] this preprocessing method is denoted preprocessing method A and differs only slightly from the preprocessing procedure described in Section 3.1.

Table 7 Prediction results for the 118 individuals in the CC3 dataset.

	FN	FP	TN	TP	P-value
50-gene best predictor in [1] (result reported previously in Table 6 a) in [2])	10	22	37	49	2.98e-07
10-gene best predictor (genes selected from the 50-gene best predictor using the same procedure as was used in [1] for the 50-gene best predictor)	9	21	38	50	3.38e-08

Table 8 Prediction results for the 118 individuals in the CC3 dataset using leave-one-out prediction. Results are shown for five different simulations.

	FN	FP	TN	TP	P-value
Naïve Bayes (selected 50 of 345 genes, simulation 1) Result reported previously in Table 7 in [2]	14	23	36	45	3.78e-05
Naïve Bayes (selected 50 of 345 genes, simulation 2) Result reported previously in Table 7 in [2]	18	21	38	41	2.06e-04
Naïve Bayes (selected 10 of 50 genes, simulation 3)	17	27	32	42	4.30e-03
Naïve Bayes (selected 10 of 50 genes, simulation 4)	19	23	36	40	1.49e-03
Naïve Bayes (selected 10 of 50 genes, simulation 5)	19	21	38	40	4.19e-04

Note that the 10 genes that are selected for the predictor (Table 3) define one of many possible «10-gene best predictors» as this predictor is selected as the best of 100 000 randomly chosen predictors (we choose 10 out of 50 genes 100 000 times so that we obtain 100 000 different predictors). Choosing one predictor from 100 000 other randomly chosen predictors, will result in another «10-gene best predictors» (i.e. we cannot expect that the same 10 genes as in Table 3 are selected for the predictor).

Table 9 This table is a simplified version of Additional file 3 from [1] where only rows with the 50 genes of the 50-gene best predictor are included. The 10 genes selected from these 50 genes are highlighted in yellow.

Gene Symbol	Gene Name	Keywords	logFC CC1	logFC CC2	logFC CC3	FDR CC3
DYNLRB1	dynein, light chain, roadblock-type 1	ER vesicles transport, tumorigenesis	-0.18	-0.12	-0.12	0.01
TUBB	tubulin beta	cytoskeleton, cell shape, cell cycle	-0.20	-0.14	-0.12	0.01
DHX40	DEAH (Asp-Glu-Ala-His) box polypeptide 40	regulation transcription, helicase	0.13	0.08	0.07	0.03
PAPOLA	poly(A) polymerase alpha	RNA processing, polyA	0.42	0.21	0.17	0.03
KARS	lysyl-tRNA synthetase	regulation transcription, immune, monocyte/macrophage	-0.09	-0.12	-0.08	0.03
CAMLG	calcium modulating ligand	immune, T cell, signalling	0.26	0.20	0.10	0.03
DCP1A	DCP1 decapping enzyme homolog A ( <i>S. cerevisiae</i> )	RNA metabolism, decay, tgfb	-0.22	-0.16	-0.08	0.07
GSTP1	glutathione S-transferase pi 1	metabolism, xenobiotic, tumorigenesis	-0.27	-0.18	-0.10	0.08
ZNF266	zinc finger protein 266	regulation transcription, zinc-finger	0.16	0.20	0.12	0.08
RPS3A	ribosomal protein S3A	translation, erythropoiesis, ribosome	0.87	0.64	0.30	0.08
RPL5	ribosomal protein L5	translation, ribosome	0.37	0.29	0.16	0.09
DENR	density-regulated protein	translation, tumorigenesis	0.11	0.12	0.07	0.10
GLRX	glutaredoxin (thioltransferase)	Metabolism	0.48	0.26	0.17	0.13
LRFN3	leucine rich repeat and fibronectin type III domain containing 3	cell adhesion	-0.23	-0.14	-0.07	0.14
ARHGAP1	Rho GTPase activating protein 1	cell cycle, ras signalling	-0.40	-0.24	-0.10	0.14
APOL3	apolipoprotein L3	intracellular lipid	-0.17	-0.19	-0.08	0.15
SMARCAL1	SWI/SNF related, matrix associated, actin dependent regulator of chromatin, subfamily a-like 1	regulation transcription, helicase	-0.20	-0.11	-0.06	0.22
FRYL	FRY-like		0.21	0.22	0.05	0.25
C18orf8	chromosome 18 open reading frame 8		-0.12	-0.11	-0.04	0.25
RPL21	ribosomal protein L21	translation, ribosome	0.65	0.33	0.11	0.29
PHF5A	PHD finger protein 5A	RNA metabolism, splicing	0.19	0.18	0.05	0.29
HPS6	Hermansky-Pudlak syndrome 6	secretory pathway, platelet	-0.23	-0.18	-0.05	0.29
KIAA1310			-0.18	-0.16	-0.05	0.31
IK	IK cytokine, down-regulator of HLA II		-0.16	-0.18	0.05	0.32
SDHA	succinate dehydrogenase complex, subunit A, flavoprotein (Fp)	energy, mitochondrion	-0.21	-0.20	-0.04	0.35
PPP2R5A	protein phosphatase 2, regulatory subunit B', alpha	cell growth	0.31	0.26	0.06	0.36



CECR1	cat eye syndrome chromosome region, candidate 1	cell proliferation, cell differentiation	-0.28	-0.47	-0.09	0.37
WBP11	WW domain binding protein 11	RNA processing, splicing	-0.11	-0.19	0.05	0.37
ABHD10	abhydrolase domain containing 10		0.28	0.21	0.05	0.38
COPS7B	COP9 constitutive photomorphogenic homolog subunit 7B (Arabidopsis)	protein modification, ubiquitination, pluripotent	-0.14	-0.09	-0.04	0.45
GMFG	glia maturation factor, gamma		0.34	0.33	0.07	0.46
KIAA0930			-0.14	-0.09	-0.03	0.48
S100A8	S100 calcium binding protein A8	pleiotropic, immune, inflammation, cell cycle, cell differentiation, apoptosis	0.76	0.44	0.12	0.49
SP2	Sp2 transcription factor	regulation transcription	-0.32	-0.25	-0.05	0.51
LOC100290936	phosphoglycerate mutase 1-like		-0.27	-0.27	-0.05	0.52
RASSF5	Ras association (RalGDS/AF-6) domain family member 5		-0.17	-0.19	-0.05	0.52
UBAC2	UBA domain containing 2		0.10	0.09	0.03	0.53
ARHGAP17	Rho GTPase activating protein 17	cytoskeleton, Rho signalling	-0.17	-0.13	-0.03	0.53
ZNF319	zinc finger protein 319	regulation transcription, zinc-finger	-0.23	-0.16	-0.03	0.58
TMEM39B	transmembrane protein 39B		-0.18	-0.12	-0.03	0.59
C14orf2	chromosome 14 open reading frame 2		0.34	0.26	0.04	0.61
CD74	CD74 molecule, major histocompatibility complex, class II invariant chain	immune, MHC II, ER transport	-0.37	-0.33	-0.06	0.62
EP300	E1A binding protein p300	transcription regulation, chromatin, cell proliferation, cell differentiation	-0.22	-0.21	0.03	0.68
AIMP2	aminoacyl tRNA synthetase complex-interacting multifunctional protein 2	regulation transcription, RNA processing, tRNA synthetase	-0.13	-0.12	-0.02	0.73
FIP1L1	FIP1 like 1 (S. cerevisiae)	RNA metabolism, polyA	-0.11	-0.11	-0.01	0.77
SEC31A	SEC31 homolog A (S. cerevisiae)	ER vesicles transport	-0.18	-0.11	-0.01	0.86
YWHAB	tyrosine 3-monooxygenase/tryptophan 5-monooxygenase activation protein, beta polypeptide	signaling, cell cycle, ras signalling	-0.26	-0.21	-0.01	0.90
C16orf72	chromosome 16 open reading frame 72		0.25	0.22	-0.01	0.94
JAK1	Janus kinase 1	immune, ifn	-0.16	-0.15	0.01	0.96
HNRNPAB	heterogeneous nuclear ribonucleoprotein A/B	RNA processing, RNA metabolism, RBP, RRM	-0.16	-0.17	0.00	0.98