

Research

## A kernel-based integration of genome-wide data for clinical decision support

Anneleen Daemen\*, Olivier Gevaert\*, Fabian Ojeda\*,  
Annelies Debucquoy<sup>†</sup>, Johan AK Suykens\*, Christine Sempoux<sup>‡</sup>,  
Jean-Pascal Machiels<sup>§</sup>, Karin Haustermans<sup>†</sup> and Bart De Moor\*

Addresses: \*Department of Electrical Engineering (ESAT-SCD), Katholieke Universiteit Leuven, Kasteelpark Arenberg, 3001 Leuven, Belgium. <sup>†</sup>Department of Experimental Radiotherapy, Katholieke Universiteit Leuven, UZ Herestraat, 3000 Leuven, Belgium. <sup>‡</sup>Department of Pathology, Université Catholique de Louvain, St Luc University Hospital, Avenue Hippocrate, 1200 Brussels, Belgium. <sup>§</sup>Department of Medical Oncology, Université Catholique de Louvain, St Luc University Hospital, Avenue Hippocrate, 1200 Brussels, Belgium.

Correspondence: Anneleen Daemen. Email: [anneleen.daemen@esat.kuleuven.be](mailto:anneleen.daemen@esat.kuleuven.be)

Published: 3 April 2009

*Genome Medicine* 2009, **1**:39 (doi:10.1186/gm39)

The electronic version of this article is the complete one and can be found online at <http://genomemedicine.com/content/1/4/39>

© 2009 Daemen *et al.*; licensee BioMed Central Ltd.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Received: 4 November 2008

Revised: 20 March 2009

Accepted: 3 April 2009

### Abstract

**Background:** Although microarray technology allows the investigation of the transcriptomic make-up of a tumor in one experiment, the transcriptome does not completely reflect the underlying biology due to alternative splicing, post-translational modifications, as well as the influence of pathological conditions (for example, cancer) on transcription and translation. This increases the importance of fusing more than one source of genome-wide data, such as the genome, transcriptome, proteome, and epigenome. The current increase in the amount of available omics data emphasizes the need for a methodological integration framework.

**Methods:** We propose a kernel-based approach for clinical decision support in which many genome-wide data sources are combined. Integration occurs within the patient domain at the level of kernel matrices before building the classifier. As supervised classification algorithm, a weighted least squares support vector machine is used. We apply this framework to two cancer cases, namely, a rectal cancer data set containing microarray and proteomics data and a prostate cancer data set containing microarray and genomics data. For both cases, multiple outcomes are predicted.

**Results:** For the rectal cancer outcomes, the highest leave-one-out (LOO) areas under the receiver operating characteristic curves (AUC) were obtained when combining microarray and proteomics data gathered during therapy and ranged from 0.927 to 0.987. For prostate cancer, all four outcomes had a better LOO AUC when combining microarray and genomics data, ranging from 0.786 for recurrence to 0.987 for metastasis.

**Conclusions:** For both cancer sites the prediction of all outcomes improved when more than one genome-wide data set was considered. This suggests that integrating multiple genome-wide data sources increases the predictive performance of clinical decision support models. This emphasizes the need for comprehensive multi-modal data. We acknowledge that, in a first phase, this will substantially increase costs; however, this is a necessary investment to ultimately obtain cost-efficient models usable in patient tailored therapy.

## Background

Kernel methods are a powerful class of methods for pattern analysis. In recent years, they have become a standard tool in data analysis, computational statistics, and machine learning applications [1]. Based on a strong theoretical framework, their rapid uptake in applications such as bioinformatics [2], chemoinformatics, and even computational linguistics is due to their reliability, accuracy, and computational efficiency. In addition, they have the capability to handle a very wide range of data types (for example, kernel methods have been used to analyze sequences, vectors, networks, phylogenetic trees, and so on). The ability of kernel methods to deal with complex structured data makes them ideally positioned for heterogeneous data integration. More specifically, in this study we used a weighted least squares support vector machine (LS-SVM), an extension of the support vector machine (SVM) for supervised classification [3-5]. Compared to the SVM, the LS-SVM is easier and faster for high dimensional data because the quadratic programming problem is converted into a linear problem. To account for the unbalancedness in many two-class problems, this linear problem is extended with weights that are different for the positive and negative classes.

The growing amount of data combined with factors such as time, cost, and personalized treatment is complicating clinical decision making. Using advanced mathematical models such as the above mentioned LS-SVM can aid clinical decision support because information arising from clinical risk factors (for example, tumor size, number of positive lymph nodes) is not accurate enough to reliably predict patient prognoses. Patients with the same clinical and pathological characteristics but different clinical outcomes can potentially be discerned with microarray technology. This technology investigates the transcriptomic make-up of a tumor in one experiment. A decade ago, it was first used in cancer studies to classify tissues as cancerous or non-cancerous [6,7]. Within the domain of cancer, microarray technology has earned a prominent place for its capacity to characterize underlying tumor behavior in detail. Although the first gene expression profile signature is being validated in clinical trials [8-10], microarray technology can not measure the complete transcription profile due to the limited number of probes per gene on a chip; nor does the transcriptome completely reflect the biology underlying a disease.

Besides transcription, pathological conditions such as cancer also influence alternative splicing, chromosomal aberrations, and methylation [11,12]. For example, chromosomal aberrations have been found in the general population as well as in all major tumor types [13,14]. These regions of increased or decreased DNA copy number can be detected using, for example, array comparative genomic hybridization (CGH) technology. This technique measures copy number variations (CNVs) within the entire genome of a disease sample compared to a normal sample [11]. Many

small aberrations have emerged as prognostic and predictive markers. Numerous aberrations, however, also affect large genomic regions, encompassing multiple genes or whole chromosome arms.

Due to differential splicing or post-translational modifications such as phosphorylation or acetylation, the proteome is many orders of magnitude bigger than the transcriptome. This makes the proteome, which reflects the functional state of the cell, a potentially richer source of data for unraveling diseases [15]. It can be measured using mass spectrometry [16], or protein or antibody microarrays [17]. Additionally, other available omics data, such as epigenomics - the study of epigenetic changes such as DNA methylation and histone modifications [12] - and single nucleotide polymorphism genotyping [18], should be considered as they promise to be useful in unraveling cancer mechanisms and the refinement of their molecular descriptions. Although the technologies are available, joint analysis of multiple hierarchical layers of biological regulation is at a preliminary stage.

In this study we investigate whether the integration of information from multiple layers of biological regulation improves the prediction of cancer outcome.

## Related work

Other research groups have already proposed the idea of data integration, but most groups have only investigated the integration of clinical and microarray data. Tibshirani and colleagues [19] proposed such a framework by reducing the microarray data to one variable, addable to models based on clinical characteristics such as age, grade, and size of the tumor. Nevins and colleagues [20] combined clinical risk factors with metagenes (that is, the weighted average expression of a group of genes) in a tree-based classification system. Wang *et al.* combined microarray data with knowledge on two clinicopathological variables by defining a gene signature only for the subset of patients for whom the clinicopathological variables were not sufficient to predict outcome [21].

A further evolution can be seen in studies in which two omics data sources are simultaneously considered, in most cases microarray data combined with proteomics or array CGH data. Much literature on such studies involving data integration already exists. However, the current definition of the integration of high-throughput data sources as it is used in the literature differs from our point of view.

In a first group of integration studies, heterogeneous data from different sources were analyzed sequentially; that is, one data source was analyzed while the second was used as confirmation of the found results or for further deepening the understanding of the results [22]. Such approaches are used for biological discovery and a better understanding of the development of a disease, but not for predictive pur-

poses. For example, Fridlyand and colleagues [23] found three breast tumor subtypes with a distinct CNV pattern based on array CGH data. Microarray data were subsequently analyzed to identify the functional categories that characterized these subtypes. Tomioka *et al.* [24] analyzed microarray and array CGH data of patients with neuroblastoma in a similar way. Genomic signatures resulted from the array CGH data, while molecular signatures were found after the microarray analysis. The authors suggested that a combination of these independent prognostic indicators would be clinically useful.

The term data integration has also been used as a synonym for data merging in which different data sets are concatenated at the database level by cross-referencing the sequence identifiers, which requires semantic compatibility among data sets [25,26]. Data merging is a complex task due to, for example, the use of different identifiers, the absence of a 'one gene-one protein' relationship, alternative splicing, and measurement of multiple signals for one gene. In most studies, the concordance between the merged data sets and their interpretation in the context of biological pathways and regulatory mechanisms are investigated. Analyses of the merged data set by clustering or correlating the protein and microarray data can help identify candidate targets when changes in expression occur at both the gene and protein levels. However, there has been only modest success from correlation studies of gene and protein expression. Bitton *et al.* [27] combined proteomics data with exon array data, which allowed a much more fine-grained analysis by assigning peptides to their originating exons instead of mapping transcripts and proteins based on their IDs.

Our definition for the combination of heterogeneous biological data is different. We integrate multiple layers of experimental data into one mathematical model for the development of more homogeneous classifiers in clinical decision support. For this purpose, we present a kernel-based integration framework. Integration occurs within the patient domain at a level not so far described in the literature. Instead of merging data sets or analyzing them in turn, the variables from different omics data are treated equally. This leads to the selection of the most relevant features from all available data sources, which are combined in a machine learning-based model. We were inspired by the idea of Lanckriet and colleagues [28]. They presented an integration framework in which each data set is transformed into a kernel matrix. Integration occurs on this kernel level without referring back to the data. They applied their framework to amino acid sequence information, expression data, protein-protein interaction data, and other types of genomic information to solve a single classification problem: the classification of transmembrane versus non-transmembrane proteins. In this study by Lanckriet and colleagues, all considered data sets were publicly available. This requires a computationally intensive framework for determining the

relevance of each data set by solving an optimization problem. Within our set-up, however, all data sources are derived from the patients themselves. This makes the gathering of these data sets highly costly and limits the number of data sets, but guarantees more relevance for the problem at hand.

We previously investigated whether the prediction of distant metastasis in breast cancer patients could be improved when considering microarray data besides clinical data [29]. In this manuscript, we consider not only microarray data but also high-throughput data from multiple biological levels. Three different strategies for clinical decision support are proposed: the use of individual data sets (referred to as step A); an integration of each data type over time by manually calculating the change in expression (step B); and an approach in which data sets are integrated over multiple layers in the genome (and over time) by treating variables from the different data sets equally (step C).

We apply our framework to two cases, summarized in Table 1. In the first case on rectal cancer, tumor regression grade, lymph node status, and circumferential margin involvement (CRM) are predicted for 36 patients based on microarray and proteomics data, gathered at two time points during therapy. The second case on prostate cancer involves microarray and copy number variation data from 55 patients. Tumor grade, stage, metastasis, and occurrence of recurrence were available for prediction [30,31].

## Materials and methods

### Data set I: rectal cancer

#### Patients and treatment

Forty patients with rectal cancer (T3-T4 and/or N+) from seven Belgian centers were enrolled in a phase I/II study investigating the combination of cetuximab, capecitabine, and external beam radiotherapy in the preoperative treatment of patients with rectal cancer [32]. These patients received preoperative radiotherapy (1.8 Gy, 5 days/week for 5 weeks) in combination with cetuximab (initial dose 400 mg/m<sup>2</sup> intravenous given 1 week before the beginning of radiation followed by 250 mg/m<sup>2</sup>/week for 5 weeks) and capecitabine for the duration of radiotherapy (first dose level, 650 mg/m<sup>2</sup> orally twice-daily; second dose level, 825 mg/m<sup>2</sup> twice-daily; including weekends). Details of the eligibility criteria, pretreatment evaluation, radiotherapy, chemotherapy and cetuximab administration, surgery, follow-up, and histopathological assessment of response to chemoradiation have been published [32].

#### Data preprocessing

Tissue and plasma samples were gathered at three time points: before treatment ( $T_0$ ); after the first loading dose of cetuximab but before the start of radiotherapy with capecitabine ( $T_1$ ); and at the moment of surgery ( $T_2$ ). All

**Table 1****Overview of the two case studies on rectal and prostate cancer**

|  | Data set I: rectal cancer  | Data set II: prostate cancer               |
|--|--|--|
| Number of samples                        | 36   | 55   |
| Data sources                             | Microarray<br>Proteomics   | Microarray<br>Genomics                     |
| Number of features (after preprocessing) | $T_0$ : 6,913 genes; 90 proteins<br>$T_1$ : 6,913 genes; 92 proteins | 6,974 genes<br>7,305 CNVs                  |
| Outcomes                                 | WHEELER<br>pN-STAGE<br>CRM   | GRADE<br>STAGE<br>METASTASIS<br>RECURRENCE |

experimental procedures were done following standard laboratory procedures, or following the manufacturers' instructions. Because of the exclusion of some patients due to a missing outcome value, death before surgery, or not having surgery, the data set ultimately contained 36 patients.

The frozen tissue samples were hybridized to Affymetrix human U133 2.0 plus gene chip arrays. The resulting data were first preprocessed for each time point separately using robust multichip analysis [33]. Secondly, the number of features was reduced from 54,613 probe sets to 27,650 genes by taking the median of all probe sets that matched on the same gene. Probe sets that matched on multiple genes were excluded because of the danger of cross-hybridization. Taking into account the low signal-to-noise ratio of microarray data, we finally filtered out genes with low variation across all samples. Only retaining the genes with a variance in the top 25% reduced the number of features to 6,913 genes.

Ninety-six proteins known to be involved in cancer were measured in the plasma samples using a Luminex 100 instrument. Proteins that had absolute values above the detection limit in less than 20% of the samples were excluded for each time point separately. This resulted in the exclusion of six proteins at  $T_0$ , four at  $T_1$ , and six at  $T_2$ . The proteomics expression values of transforming growth factor alpha, which had too many values below the detection limit, were replaced by the results of ELISA tests performed at the Department of Experimental Oncology in Leuven, Belgium. For the remaining proteins the missing values were replaced by half of the minimum detected for each protein over all samples, and values exceeding the upper limit were replaced by the upper limit value. Because most of the proteins had a positively skewed distribution, a log transformation (base 2) was performed.

In this paper, only the data sets at  $T_0$  and  $T_1$  were used because our goal is to predict the four different outcomes before therapy or early in therapy.

*Response classification*

A semiquantitative classification system has been described by Wheeler *et al.* [34] for determining histopathological tumor regression (that is, the therapy response). There are also two prognostic factors important in rectal cancer: pathologic lymph node involvement and CRM [35]. Because the completeness of tumor resection relies on the assessment of resection margins by the pathologist, knowledge of the CRM before therapy provides important prognostic information for local recurrence and for development of distant metastasis and survival [36].

These three outcomes were registered for 36 patients at the moment of surgery. For all these outcomes, 'responders' are distinguished from 'non-responders'. The grading of regression established by Wheeler and colleagues [34] (from now on referred to as WHEELER) is a modified pathological staging system for irradiated rectal cancer. It includes a measurement of tumor response after preoperative therapy: grade 1, good responsiveness (tumor is sterilized or only microscopic foci of adenocarcinoma remain); grade 2, moderate responsiveness (marked fibrosis but still with a macroscopic tumor); grade 3, poor responsiveness (little or no fibrosis with abundant macroscopic tumor). Tumors are classified as 'responder' when assigned to grade 1 (26 patients) and 'non-responder' when assigned to grade 2 or 3 (10 patients). Response can also be evaluated with the pathologic lymph node stage at surgery (pN-STAGE). The 'responder' class contains 22 patients with no lymph nodes found at surgery while the 'non-responder' class contains 14 patients with at least 1 regional lymph node. CRM was measured according to the guidelines of Quirke *et al.* [37]. CRM was considered positive when the distance between the tumor and the mesorectal fascia was  $\leq 2$  mm. Tumors with a negative CRM are classified as 'responder' (27 patients), while tumors with a positive CRM belong to the 'non-responder' class (9 patients). Thirteen patients belong to the 'responder' class for all three outcomes, while there is an overlap of two patients between the 'non-responder' classes.

## Data set II: prostate cancer

### Patients and treatment

We also applied our method to a publicly available data set of prostate cancer. Lapointe and colleagues [30] first profiled gene expression in 71 prostate tumor cases of which 62 were primary and 9 had lymph node metastases. All tumors were removed by radical prostatectomy (that is, the surgical removal of the prostate gland). A cDNA microarray was used, containing 39,711 human cDNAs representing 26,260 mapped genes. Additionally, DNA CNVs were profiled on cDNA microarrays for CGH for 64 prostate tumor cases, among which 55 were primary tumors and 9 had pelvic lymph node metastases. The arrays were obtained from the Stanford Functional Genomics Facility and included 39,632 human cDNAs corresponding to 22,279 genes [31]. Among the primary tumors, the available gene expression and genomics data were in common for 55.

### Data preprocessing

Median fluorescence ratios were calculated for genes represented by multiple arrayed cDNAs. Missing gene expression values were imputed unsupervised using the k-nearest neighbors method of Troyanskaya *et al.* [38]. The parameter k was set to 15 such that a missing value for a spot S in a sample was estimated as the weighted average of the 15 spots that are most similar to spot S in the remaining samples. The same unsupervised prefiltering as applied on the rectal cancer data set was used for both the microarray and genomics data sets. Features with a variance in the top 50% were retained, reducing the data sets to 6,974 genes and 7,305 CNVs, respectively.

### Response classification

Two pathological variables, stage and grade, metastasis of the tumor, as well as the outcome after prostatectomy defined as recurrence were considered. For grade (from now on referred to as GRADE), the Gleason Grading system was used, which is based on the most common and second most common architectural patterns of the glands of the tumor [39]. Two groups could be distinguished based on the architecture of the most common pattern: 36 tumors were well differentiated (that is, low-grade), 19 were poorly differentiated (that is, high-grade). According to the extent of the primary tumor (STAGE), 25 samples were of stage T2 (that is, the cancer is confined within one lobe of the prostate gland), while 25 samples were of advanced stage T3 (that is, the tumor has extended through the fibrous tissue surrounding the prostate gland but no other organs are affected). The stage of the remaining five patients was not known. The cancer had metastasized to distant lymph nodes in 12 tumors, while the cancer had not spread beyond the regional lymph nodes in 38 of the tumors (METASTASIS). Tumor recurrence was defined as a rise in prostate-specific antigen of at least 0.07 ng/ml or as occurrence of clinical metastasis (RECURRENCE). Seven tumors recurred while 22 tumors did not. The recurrence status of the remaining 26 patients was not available.

## Kernel methods and weighted least squares support vector machines

Kernel methods are a group of algorithms that can handle a very wide range of data types, such as vectors, sequences, networks, and so on. They map the data  $x$  from the original input space to a high dimensional feature space with the mapping function  $\Phi(x)$ . This embedding into the feature space is performed by a mathematical object  $K(x_k, x_l)$ , called a 'kernel function'. This function efficiently computes the inner product  $\langle \Phi(x_k), \Phi(x_l) \rangle$  between all pairs of data items  $x_k$  and  $x_l$  in the feature space, resulting in the kernel matrix. The size of this matrix is determined only by the number of data items, whatever the nature or the complexity of these items. For example, a set of 100 patients each characterized by 6,913 gene expression values is still represented by a  $100 \times 100$  kernel matrix [40]. The representation of all data sets by this real-valued square matrix, independent of the nature or complexity of the data to be analyzed, makes kernel methods ideally positioned for heterogeneous data integration.

Any symmetric, positive semidefinite function is a valid kernel function, resulting in many possible kernels - for example, linear, polynomial, and diffusion kernels. They all correspond to a different transformation of the data, meaning that they extract a specific type of information from the data set. In this paper, the normalized linear kernel function:

$$\tilde{K}(x_k, x_l) = K(x_k, x_l) / \sqrt{\{K(x_k, x_k) K(x_l, x_l)\}}$$

where  $K(x_k, x) = x_k^T x$  is used instead of the linear kernel function  $K(x_k, x_l) = x_k^T x_l$ . With the normalized version, the values in the kernel matrix will be bounded because the data points are projected onto the unit sphere while these elements can take very large values without normalization. Normalizing is thus required when combining multiple data sources to guarantee the same order of magnitude for the kernel matrices of the data sets.

A kernel algorithm for supervised classification is the SVM developed by Vapnik [41] and others. Contrary to most other classification methods and due to the way data are represented through kernels, SVMs can tackle high dimensional data (for example microarray data). Given a training set  $(x_k, y_k)_{k=1}^N$  of N samples with feature vectors  $x_k \in R^n$  and output labels  $y_k \in \{-1, +1\}$ , the SVM forms a linear discriminant boundary  $y(x) = \text{sign}[w^T \Phi(x) + b]$  in the feature space with maximum distance between samples of the two considered classes, with  $w$  representing the weights for the data items in the feature space and  $b$  the bias term. This corresponds to a non-linear discriminant function in the original input space. A modified version of SVM, LS-SVM, was developed by Suykens *et al.* [3,4]. On high dimensional data sets, this modified version is much faster for classification because a linear system instead of a quadratic programming problem needs to be solved.

The constrained optimization problem for an LS-SVM has the following form:

$$\min_{w,b,e} \left( \frac{1}{2} w^T w + \gamma \frac{1}{2} \sum_{k=1}^N e_k^2 \right)$$

subject to:

$$y_k [w^T \Phi(x_k) + b] = 1 - e_k, k = 1, \dots, N$$

with  $e_k$  the error variables, tolerating misclassifications in cases of overlapping distributions, and  $\gamma$  the regularization parameter, which allows tackling the problem of overfitting. It has been shown that regularization seems to be very important when applying classification methods on high dimensional data [42].

In many two-class problems, data sets are skewed in favor of one class such that the contribution of false negative and false positive errors to the performance assessment criterion are not balanced. We therefore used a weighted LS-SVM in which a different weight  $\zeta_k$  is given to positive and negative samples in order to account for the unbalancedness in the data set [5]. The objective function changes into:

$$\min_{w,b,e} \left( \frac{1}{2} w^T w + \gamma \frac{1}{2} \sum_{k=1}^N \zeta_k e_k^2 \right)$$

with

$$\zeta_k = \begin{cases} \frac{N}{2N_p} & \text{if } y_k = +1 \\ \frac{N}{2N_N} & \text{if } y_k = -1 \end{cases}$$

and  $N_p$  and  $N_N$  representing the number of positive and negative samples, respectively.

### Feature selection

Univariate feature selection techniques are computationally simple but do not incorporate feature-feature interactions. However, due to small sample size limitations, multivariate approaches are often not appropriate for discovering the underlying complex, multivariate correlations. Because it has been shown that univariate gene selection methods lead to good and stable performances across many cancer types and yield in many cases consistently better results than multivariate approaches [43], we used the method DEDS (differential expression via distance synthesis) [44]. This technique is based on the integration of different univariate test statistics via a distance synthesis scheme because features highly ranked simultaneously by multiple statistics are more likely to be differentially expressed than features highly ranked by a single test statistic. The statistical tests combined are ordinary fold changes, ordinary  $t$ -statistics, SAM (significance analysis for microarrays) statistics and moderated  $t$ -statistics. DEDS is available as a BioConductor package in R.

We applied DEDS to the microarray data sets as well as the genomics data set. From our experience, DEDS is less appropriate for data with a limited set of features (data not shown). Since the proteomics data on rectal cancer contain only 90-92 cancer-related proteins, one test statistic suffices, for which we chose the Wilcoxon rank sum test.

### Model building

To determine the optimal number of features, we use a leave-one-out (LOO) cross-validation approach in which we increase the number of included features iteratively according to the obtained feature ranking but in which we do not include more features than the number of samples in the data set on which the optimal number of features is determined, as discussed by Li and Yang [45]. Besides the number of features, the parameters of the kernel method (parameter  $\gamma$  for LS-SVM with normalized linear kernel) also need to be selected. This selection occurs on a  $k$ -dimensional grid with  $k-1$  the number of data sets included. We considered 40 possible values for  $\gamma$ , ranging from  $10^{-4}$  to  $10^6$  on a logarithmic scale. In each LOO iteration, a sample is left out, feature selection is performed on the remaining  $n-1$  samples, and models are built for all possible combinations of parameters on this grid. Each model with the instantiated parameters is evaluated on the left out sample. This whole procedure is repeated for all samples. The model parameters are chosen corresponding to the model with the highest LOO area under the receiver operating characteristic (ROC) curve (AUC). If multiple models have the same AUC, the model with the lowest balanced error rate and an as high as possible sum of sensitivity and specificity is chosen. For each considered outcome, the AUC of the best performing model is compared with the AUC of the other models using the method of Hanley and McNeil [46]. The final features are chosen as those that occurred most often in the top rankings determined in each LOO iteration.

Three kinds of model building strategies are proposed, different in the degree of integration. Figure 1 shows these strategies in more detail. The data sets are represented as matrices with rows corresponding to patients and columns corresponding to genes, proteins, or CNVs. The matrices representing microarray or genomics data are larger than those for the proteomics data to emphasize the difference in dimensionality.

All three strategies were applied to the microarray and proteomics data sets of rectal cancer. For the prostate cancer data set, however, only two strategies were applicable due to a lack of measurements repeated over time. For all models the parameters were trained according to the same approach, which makes the corresponding LOO results comparable for each outcome separately.

#### Step A models: single data set

In a first step, LS-SVM models are built on each data set separately, mimicking the results that would have been

obtained when only static data from one platform were available. For rectal cancer, the single data sets are microarray at  $T_0$ , microarray at  $T_1$ , proteomics at  $T_0$ , and proteomics at  $T_1$  for the prediction of a regression grading system and two prognostic factors (Figure 1a). For prostate cancer, LS-SVM models are built on the microarray and genomics data separately for the prediction of grade, stage, metastasis, and recurrence. Because of only one set of features, a two-dimensional grid is used for the optimization of the regularization parameter and the number of features.

#### Step B models: manual integration of data over time

When measurements are repeated at multiple time points, knowledge over time can be exploited. For rectal cancer, data were available before and early in therapy and, therefore, can be combined in the models. This is done for each data type separately by manually calculating the change in gene expression or protein abundance between the first two time points ( $T_0$ - $T_1$ ). These changes over time are used as features for the models as shown in Figure 1b. Also for these models, a two-dimensional grid suffices for the optimization of the regularization parameter and the number of features.

#### Step C models: multiple omics integration approach

The previous two types of models (steps A and B) are considered to verify whether complex integration of data over multiple layers of biological regulation is crucial. The ability of kernel methods to deal with complexly structured data makes them ideally positioned for more advanced integration of heterogeneous data sources. We will use the intermediate integration method proposed in [47] in which a kernel matrix is computed for each data source separately. Subsequently, these data sources can be integrated in a straightforward way by summing the multiple kernel matrices. Positive semidefiniteness of the linear combination of kernel matrices is guaranteed by constraining the weights of the kernels to be non-negative. A weighted LS-SVM is trained on the explicitly heterogeneous kernel matrix. The choice of the weights to give to each data set is important. A kernel framework for optimizing weights is proposed in [48]. This optimization is important when dealing with many data sets of which only several are relevant. However, when the number of data sets is limited and most of them are reliable and relevant to the problem at hand, a trade-off needs to be made between performance and computational burden (for example, extra required cross-validation loops). Due to the rather small sample size in both case studies, weights were chosen equally. Moreover, our aim is to emphasize that classification becomes more accurate when data from multiple layers in the genome are available and to offer a machine learning-based method for integrating these data sources, rather than to improve an algorithm for the optimization of weights (for example, [48]). A three-dimensional grid is used for the optimization of the parameters, that is, the regularization parameter, the number of genes selected from the microarray data sets, and

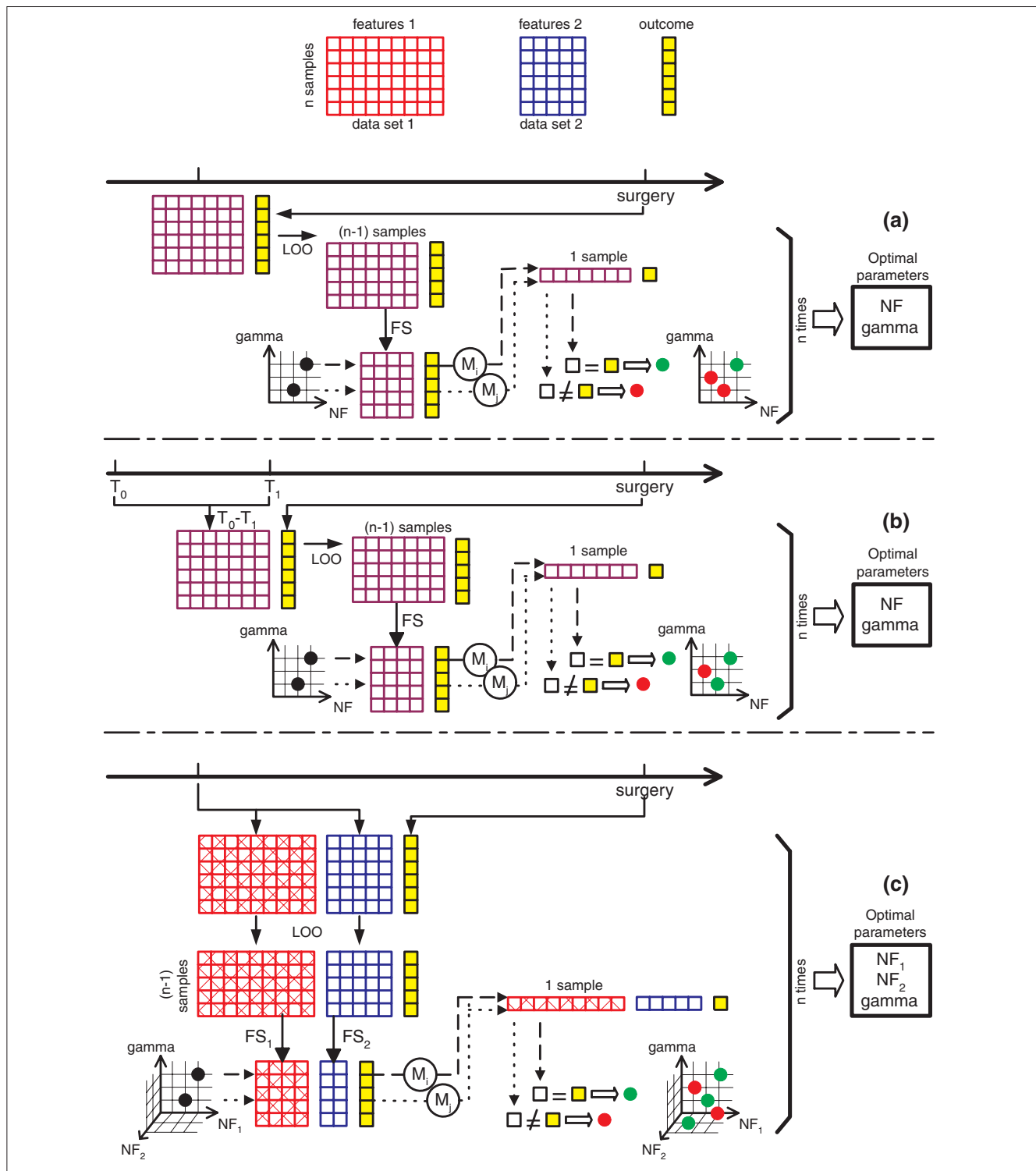
the number of proteins or CNVs obtained from the proteomics data sets or the genomics data set, respectively. For the data on rectal cancer, the number of genes/proteins selected at  $T_0$  and  $T_1$  were taken equally when data from both time points were considered. Figure 1c gives an overview of the strategy.

## Results

### Study I: rectal cancer

Using the methodologies shown in Figure 1, models were built using microarray and proteomics data of 36 rectal cancer patients at two time points during therapy for the prediction of three outcomes registered at the moment of surgery: a tumor regression grading system (WHEELER) and two prognostic factors, pathologic N stage at surgery (pN-STAGE) and the circumferential margin involvement (CRM). The models with the highest AUC, lowest balanced error rate and an as high as possible sum of sensitivity and specificity are shown in Table 2. The step A models are  $MT_0$  (model based on microarray data at  $T_0$ ),  $MT_1$  (model based on microarray data at  $T_1$ ),  $PT_0$  (model based on proteomics data at  $T_0$ ), and  $PT_1$  (model based on proteomics data at  $T_1$ ). The step B models consist of  $MT_0$ - $T_1$  (model based on change in gene expression between  $T_0$  and  $T_1$ ) and  $PT_0$ - $T_1$  (model based on change in protein abundances between  $T_0$  and  $T_1$ ). Finally, the step C models comprise  $MT_{01}$  (model based on microarray data at both time points),  $PT_{01}$  (model based on proteomics data at both time points),  $MPT_0$  (model based on microarray and proteomics data at  $T_0$ ),  $MPT_1$  (model based on microarray and proteomics data at  $T_1$ ), all possible combinations of three data sets (using the same name convention), and  $MPT_{01}$  (model based on all data (microarray and proteomics data at both time points)). The numbers of genes and proteins were chosen to optimize the LOO performance of the LS-SVM models. The features selected most often in the 36 LOO iterations are listed and discussed. For each outcome, the ROC curve of the best model was compared with the ROC curves of all other models [46]. The  $P$ -values of these significance tests are reported as well.

Table 2 shows the LS-SVM models for the considered combinations of data sets to predict WHEELER, pN-STAGE, and CRM with the optimal number of genes and proteins selected with DEDS and the Wilcoxon rank sum test, respectively. The corresponding ROC curves are shown in Additional data file 1. The performance of the models based on three data sets is given in Additional data file 2. Due to the slightly, but not significantly, better performance for each outcome of one model based on three data sets compared to models based on two data sets, we report the results for the best model combining two data sets. Such models would only require a sample to be taken at one time point ( $MPT_0$ ,  $MPT_1$ ) or one technology to be applied on two time points ( $MT_{01}$ ,  $PT_{01}$ ). For the prediction of WHEELER, the expression of 25 genes and 12 proteins at  $T_1$  was best,



**Figure 1**  
 Overview of the three applied model building strategies. **(a)** Use of a single data set; **(b)** manual integration of data over time; **(c)** a genome-wide integration approach. The data sets are represented as matrices with rows corresponding to patients and columns corresponding to genes, proteins, or CNVs. In step A, LS-SVM models are built on each data set separately. A two-dimensional grid is used for the optimization of the regularization parameter and the number of features. For step B, data sets over time are combined. By using the changes in expression or abundance as features, a two-dimensional grid is sufficient. In step C, an intermediate integration method is used for the integration of all available data sets. A k-dimensional grid is required for optimizing the regularization parameter and the number of features selected from the (k - 1) integrated data sets. FS, feature selection; M<sub>i</sub>, model for parameter combination i; NF, number of features; T, time point.



Table 2

LS-SVM models for the prediction of WHEELER, pN-STAGE and CRM in rectal cancer

| Outcome         | Model                       | NG*       | NP†       | AUC (SE)‡                          | p-value§     |
|-----------------|-----------------------------|-----------|-----------|------------------------------------|--------------|
| <b>WHEELER</b>  |                             |           |           |                                    |              |
| A               | $MT_0$                      | 4         |           | 0.7538 (0.1085)                    | 0.0987       |
|                 | $MT_1$                      | 29        |           | 0.9038 (0.0502)                    | 0.6861       |
|                 | $PT_0$                      |           | 35        | 0.7423 (0.0867)                    | 0.0540       |
|                 | $PT_1$                      |           | 11        | 0.9038 (0.0575)                    | 0.7273       |
| B               | $MT_{\sigma}T_1$            | 32        |           | 0.6846 (0.1215)                    | 0.0598       |
|                 | $PT_{\sigma}T_1$            |           | 5         | 0.8654 (0.0621)                    | 0.4135       |
| C               | $MT_{01}$                   | 3¶        |           | 0.7808 (0.0985)                    | 0.1320       |
|                 | $PT_{01}$                   |           | 21¶       | 0.7692 (0.0831)                    | 0.0831       |
|                 | $MPT_0$                     | 3         | 35        | 0.8461 (0.0718)                    | 0.2760       |
|                 | $MPT_1$                     | <b>25</b> | <b>12</b> | <b>0.9269 (0.0425)</b>             |              |
|                 | $MPT_{01}$                  | 2¶        | 31¶       | 0.8846 (0.0558)                    | 0.4858       |
|                 | $MT_0PT_1$                  | 2         | 4         | 0.9385 (0.0444)                    | 0.8101*      |
| <b>pN-STAGE</b> |                             |           |           |                                    |              |
| A               | $MT_0$                      | 25        |           | 0.6493 (0.0914)                    | 2.315e-4     |
|                 | $MT_1$                      | 22        |           | 0.8506 (0.0665)                    | 0.0362       |
|                 | $PT_0$                      |           | 2         | 0.6753 (0.0906)                    | 6.659e-4     |
|                 | $PT_1$                      |           | 12        | 0.8409 (0.0652)                    | 0.0238       |
| B               | $MT_{\sigma}T_1$            | 4         |           | 0.6071 (0.0986)                    | 1.359e-4     |
|                 | $PT_{\sigma}T_1$            |           | 9         | 0.7662 (0.0900)                    | 0.0153       |
| C               | $MT_{01}$                   | 24¶       |           | 0.9286 (0.0450)                    | 0.1998       |
|                 | $PT_{01}$                   |           | 34¶       | 0.8182 (0.0695)                    | 0.0145       |
|                 | $MPT_0$                     | 27        | 27        | 0.9188 (0.0469)                    | 0.1591       |
|                 | $MPT_1$                     | <b>21</b> | <b>14</b> | <b>0.9870 (0.0135)</b>             |              |
|                 | $MPT_{01}$                  | 23¶       | 16¶       | 0.9610 (0.0280)                    | 0.3421       |
|                 | $MT_0PT_{01}$               | 26        | 20¶       | 1 (0)                              | 0.3347*      |
| <b>CRM</b>      |                             |           |           |                                    |              |
| A               | $MT_0$                      | 33        |           | 0.6790 (0.1016)                    | 0.0072       |
|                 | $MT_1$                      | 9         |           | 0.9259 (0.0472)                    | 0.4955       |
|                 | $PT_0$                      |           | 34        | 0.8518 (0.0624)                    | 0.0935       |
|                 | $PT_1$                      |           | 34        | 0.7654 (0.0831)                    | 0.0281       |
| B               | $MT_{\sigma}T_1$            | 6         |           | 0.9136 (0.0480)                    | 0.4030       |
|                 | $PT_{\sigma}T_1$            |           | 2         | 0.8272 (0.0709)                    | 0.0849       |
| C               | $MT_{01}$                   | 16¶       |           | 0.8066 (0.0846)                    | 0.0468       |
|                 | $PT_{01}$                   |           | 3¶        | 0.7531 (0.0865)                    | 0.0227       |
|                 | $MPT_0$                     | 7         | 27        | 0.8477 (0.0688)                    | 0.1340       |
|                 | $MPT_1$                     | <b>7</b>  | <b>33</b> | <b>0.9630 (0.0344)</b>             |              |
|                 | $MPT_{01}$                  | 2¶        | 3¶        | 0.8230 (0.0771)                    | 0.0973       |
|                 | $MT_1PT_0$<br>$MT_{01}PT_1$ | 16<br>9¶  | 14<br>29  | 0.9630 (0.0376)<br>0.9876 (0.0146) | 1<br>0.4924* |

\*Number of genes selected in each LOO iteration. †Number of proteins selected in each LOO iteration. ‡Area under the ROC curve (standard error) obtained with leave-one-out. §Comparison of AUC between each model and the best model in bold [46]. ¶Number of features used at both time points. \*This model is better than the model in bold we compare with.

although not significantly, with an AUC of 0.9269. Also for pN-STAGE, combining both data sets at  $T_1$  using the expression of 21 genes and 14 proteins resulted in the best LOO AUC of 0.9870. This performance is significantly better than all step A and B models as well as  $PT_{01}$ . Finally, the

inclusion of 7 genes and 33 proteins at  $T_1$  led to an AUC of 0.9630 for the prediction of CRM. Four models based on only one data type perform significantly worse compared to  $MPT_1$ . For all outcomes, none of the selected proteins are a product of the selected genes.

**Table 3****Features for (colo)rectal cancer selected by  $MPT_1$  and known to be involved in this type of cancer**

| Outcome* | Gene/protein        | Hits† | Region               | Function                                    | Up/down‡ | Reference |
|----------|---------------------|-------|----------------------|---|----------|-----------|
| W        | Cox-2               | 36    | 1q25.2-q25.3         | Progression                                 | Up       | [50]      |
| W        | IL-1B               | 36    | 2q14                 | Inflammatory response                       | Up       | [50]      |
| W        | Ferritin            | 36    | 11q13; 19q13.3-q13.4 | Iron storage                                | Down     | [63]      |
| W        | EGF                 | 36    | 4q25                 | Cell growth/proliferation/ differentiation  | Up       | [64]      |
| W        | MMP-2               | 36    | 16q13-q21            | Invasion/metastasis                         | Up       | [65]      |
| W        | TGF $\alpha$        | 36    | 2p13                 | Angiogenesis/cell proliferation             | Down     | [51]      |
| W        | SELE                | 25    | 1q22-q25             | Progression/metastasis                      | Up       | [66]      |
| W        | GM-CSF              | 24    | 5q31.1               | Maintenance of granulocytes/macrophages     | Up       | [67]      |
| W        | MMP-1               | 15    | 11q22.3              | Tumor invasion/ metastasis/poor prognosis   | Up       | [68]      |
| N        | Reg4                | 36    | 1p13.1-p12           | Early carcinogenesis                        | Down     | [69]      |
| N        | MUC2                | 36    | 11p15.5              | Deregulated by TNF $\alpha$                 | Down     | [70]      |
| N        | CA1                 | 36    | 8q13-q22.1           | Carbonate dehydratase activity              | Down     | [71]      |
| N        | CA2                 | 36    | 8q22                 | Carbonate dehydratase activity              | Down     | [71]      |
| N        | CLDN8               | 36    | 21q22.11             | Tumorigenesis                               | Down     | [72]      |
| N        | CEA                 | 36    | 19q13.1-q13.2        | Cell adhesion; tumor marker for recurrence  | Down     | [53]      |
| N        | IL-1ra              | 36    | 2q14.2               | Carcinogenesis                              | Up       | [73]      |
| N        | CA19-9              | 36    |                      | Tumor marker for recurrence                 | Down     | [53]      |
| N        | Ferritin            | 36    | 11q13; 19q13.3-q13.4 | Iron storage                                | Down     | [63]      |
| N        | IL-1beta            | 36    | 2q14                 | Inflammatory response                       | Down     | [50]      |
| N        | beta2-microglobulin | 36    | 15q21-q22.2          | Metastasis                                  | Up       | [74]      |
| N        | RARRES1             | 31    | 3q25.32-q25.33       | Cell proliferation                          | Down     | [75]      |
| N        | IL-8                | 28    | 4q13-q21             | Progression/metastasis                      | Down     | [52]      |
| N        | TNFR11              | 24    | 1p36.3-p36.2         | Apoptosis                                   | Up       | [76]      |
| C        | ICAM-1              | 36    | 19p13.3-p13.2        | Metastasis                                  | Down     | [77]      |
| C        | CEA                 | 36    | 19q13.1-q13.2        | Cell adhesion; tumor marker for recurrence  | Down     | [53]      |
| C        | MMP-2               | 36    | 16q13-q21            | Invasion/metastasis                         | Up       | [65]      |
| C        | Adiponectin         | 36    | 3q27                 | Metabolic/hormonal processes                | Down     | [78]      |
| C        | Thrombospondin-1    | 36    | 15q15                | Angiogenesis/tumor growth                   | Up       | [79]      |
| C        | EGFR                | 36    | 7p12                 | Cell growth/ proliferation/ differentiation | Up       | [49]      |
| C        | Tissue factor       | 35    | 1p22-p21             | Angiogenesis/metastasis                     | Up       | [80]      |
| C        | CYP1B1              | 35    | 2p21                 | Drug metabolism                             | Down     | [81]      |
| C        | EGF                 | 32    | 4q25                 | Cell growth/proliferation/ differentiation  | Up       | [64]      |

\*W, WHEELER; N, pN-STAGE; C, CRM. †Number of occurrences of the gene/protein in the 36 LOO iterations. ‡Up/down-regulation in the good responders with respect to moderate or poor responders; no lymph nodes with respect to at least one regional lymph node; negative CRM with respect to positive CRM. CRC, (colo)rectal cancer.

The contribution of the genes and/or proteins in rectal or colorectal cancer that were selected most often in the LOO iterations of  $MPT_1$  and predicted most accurately WHEELER,

pN-STAGE, or CRM are shown in Table 3. A protein important for CRM, for example, is the epidermal growth factor receptor (EGFR), involved in signaling pathways

Table 4

## LS-SVM models for the prediction of GRADE, STAGE, METASTASIS and RECURRENCE in prostate cancer

| Outcome           | Model     | NG*       | NC†       | AUC (SE)‡              | p-value§ |
|-------------------|-----------|-----------|-----------|------------------------|----------|
| <b>GRADE</b>      |           |           |           |                        |          |
| A                 | <i>M</i>  | 24        |           | 0.8304 (0.0623)        | 0.2727   |
|                   | <i>G</i>  |           | 8         | 0.7822 (0.0632)        | 0.0503   |
| C                 | <i>MG</i> | <b>6</b>  | <b>8</b>  | <b>0.9006 (0.0413)</b> |          |
| <b>STAGE</b>      |           |           |           |                        |          |
| A                 | <i>M</i>  | 18        |           | 0.6576 (0.0778)        | 0.0191   |
|                   | <i>G</i>  |           | 32        | 0.7936 (0.0631)        | 0.3466   |
| C                 | <i>MG</i> | <b>42</b> | <b>22</b> | <b>0.8528 (0.0550)</b> |          |
| <b>METASTASIS</b> |           |           |           |                        |          |
| A                 | <i>M</i>  | 18        |           | 0.9759 (0.0178)        | 0.4392   |
|                   | <i>G</i>  |           | 12        | 0.8114 (0.0755)        | 0.0166   |
| C                 | <i>MG</i> | <b>18</b> | <b>3</b>  | <b>0.9868 (0.0121)</b> |          |
| <b>RECURRENCE</b> |           |           |           |                        |          |
| A                 | <i>M</i>  | 24        |           | 0.7208 (0.0936)        | 0.5392   |
|                   | <i>G</i>  |           | 26        | 0.4481 (0.1433)        | 0.0354   |
| C                 | <i>MG</i> | <b>32</b> | <b>2</b>  | <b>0.7857 (0.0934)</b> |          |

\*Number of genes selected in each LOO iteration. †Number of copy number variations selected in each LOO iteration. ‡Area under the ROC curve (standard error) obtained with leave-one-out. §Comparison of AUC between each model and the best model in bold [46].

affecting cellular growth, differentiation, and proliferation. This protein represents one of the most promising targets allowing progress in colorectal cancer treatment. It has been suggested that EGFR polymorphisms as well as polymorphisms of other genes active in the EGFR pathway may be potential indicators of radiosensitivity in patients with rectal cancer treated with chemoradiation [49]. In colorectal cancer, pro-inflammatory cytokines such as interleukin-1 beta and interleukin-6 may be accountable for the over-expression of *Cox-2*, important in the early stage and for progression [50]. Transforming growth factor alpha, down-regulated in our patients with a good responsiveness to preoperative therapy, is implicated in metastatic spread of colon cancer cells [51]. The expression of interleukin-8 is associated with induction and progression of colorectal carcinoma and the development of colorectal liver metastases [52]. In our data set, it is down-regulated in the group of patients with no lymph nodes found at surgery. Finally, elevated carcinoembryonic antigen and cancer antigen 19-9 are related to poor outcome in colorectal cancer [53]. Their levels are low in patients with no lymph nodes, while carcinoembryonic antigen is also less expressed in patients with a negative CRM, that is, belonging to the class of 'responders'. A complete list of the genes and proteins chosen by the models  $MPT_1$  are shown, for each outcome separately, in Additional data file 3. The predictions seem to depend on mainly different subsets of features. The gene encoding PAI-2 is important for both WHEELER and CRM, while the proteins important for two of the three outcomes are interleukin-4, ferritin, apolipoprotein H, epidermal growth factor, matrix metalloproteinase-2, and lymphotactin. Notably, these genes and proteins were also selected

by the other models based on microarray and/or proteomics data at  $T_1$ , although the specific feature ranking depends on the number of features included. Some of these genes and proteins were also included in the models based on data at  $T_0$ .

#### Study II: prostate cancer

The same methodology was applied to microarray and genomics data of 55 patients with prostate cancer. Table 4 shows the results for the prediction of the grade and stage of the tumor (GRADE and STAGE), as well as the tumors that metastasized to distant lymph nodes (METASTASIS) or that recurred (RECURRENCE). Because the data were gathered at one time point, only step A and C models are applicable. The step A models are represented as *M* (model based on microarray data) and *G* (model based on genomics data), and the step C model based on both microarray and genomics data as *MG*. Also, after having optimized the essential number of features to be included using a LOO cross-validation, the final genes and CNVs were selected based on their position and number of occurrences in the 55 LOO rankings.

We obtained similar results as for rectal cancer. Combining gene expression with measurements at the DNA level (*MG*) led, for all four outcomes, to an improvement in classification accuracy and was significant in some cases (Table 4). For the prediction of GRADE, six genes and eight CNVs selected with DEDS resulted in an AUC of 0.9006. For STAGE, 42 genes and 22 CNVs were needed for a performance of 0.8528. The model *MG* for the prediction of METASTASIS had an AUC of 0.9868 when fusing the expression of 18 genes with 3 CNVs. Finally, the prediction

**Table 5****Features for prostate cancer selected by MG and known to be involved in this type of cancer**

| Outcome* | Gene/CNV       | Hits† | Region        | Function   | Up/down‡ | Reference |
|----------|----------------|-------|---------------|--|----------|-----------|
| G        | <i>SFRP4</i>   | 55    | 7p14.1        | Inhibitor of PT growth/invasion  | Up       | [55]      |
| G        | <i>VCAN</i>    | 55    | 5q14.3        | Contributor to PC pathology  | Up       | [82]      |
| G        | <i>ALOX15B</i> | 36    | 17p13.1       | Suppressor of PT development   | Down     | [54]      |
| S        | <i>MAGEA4</i>  | 50    | Xq28          | Only expressed in PC (diagnosis and therapy)                                     | Down     | [83]      |
| S        | <i>ANPEP</i>   | 50    | 15q25-q26     | PT cell invasion   | Down     | [84]      |
| S        | <i>POU4F1</i>  | 50    | 13q31.1       | PC cell growth   | Down     | [85]      |
| S        | <i>CXCL14</i>  | 48    | 5q31          | Inhibitor of PT growth   | Up       | [56]      |
| S        | <i>RNASEL</i>  | 48    | 1q25          | Polymorphic changes as tumor; suppressor in hereditary PC                        | Up       | [62]      |
| S        | <i>GDEP</i>    | 41    | 4q21.1        | Prostate-specific gene   | Down     | [86]      |
| M        | <i>ERG</i>     | 50    | 21q22.3       | Proto-oncogene; early prostate carcinogenesis                                    | Up       | [57]      |
| M        | <i>AREG</i>    | 49    | 4q13-q21      | PC progression/growth via TARP   | Down     | [87]      |
| M        | <i>VAV3</i>    | 49    | 1p13.3        | Oncogene; PC development/ progression  | Up       | [59]      |
| M        | <i>ADAMTS1</i> | 26    | 21q21.2       | Negatively affected by TGFbeta1, which increases VCAN-expression                 | Down     | [82]      |
| R        | <i>AZGP1</i>   | 29    | 7q22.1        | Inversely associated to tumor stage; predictor of biochemical recurrence         | Down     | [88]      |
| R        | <i>TIAMI</i>   | 29    | 21q22.1-11    | Predictor of decreased disease-free survival/recurrence                          | Up       | [60]      |
| R        | <i>FGG</i>     | 28    | 4q28          | PC cell growth   | Down     | [89]      |
| R        | <i>ATF3</i>    | 26    | 1q32.3        | Inversely related to invasion/ angiogenesis; positively correlated to metastases | Down     | [90]      |
| R        | <i>JAG1</i>    | 26    | 20p12.1-11.23 | Cell growth/progression/metastasis   | Up       | [61]      |
| R        | <i>ERG</i>     | 14    | 21q22.3       | Proto-oncogene; early prostate carcinogenesis                                    | Up       | [57]      |
| R        | <i>ALOX15B</i> | 14    | 17p13.1       | Suppressor of PT development   | Down     | [54]      |

\*G, GRADE; S, STAGE; M, METASTASIS; R, RECURRENCE. †Number of occurrences of the gene/CNV in all LOO iterations (number of LOO iterations for G = 55, S = 50, M = 50, R = 29). ‡Up/down-regulation in high-grade with respect to low-grade; advanced stage with respect to early stage; metastasis with respect to no metastasis; recurrence with respect to no recurrence. PC, prostate cancer; PT, prostate tumor.

of RECURRENCE was most difficult, with an AUC of 0.7857 when combining 32 genes and 2 CNVs. Additional data file 1 shows the ROC curves of the models listed in Table 4.

Several genes and CNVs have been selected by MG and are known to be involved in, and important for, prostate cancer (Table 5). The gene *ALOX15B* is a suppressor of prostate tumor development [54] and in this data set is down-regulated in tumors of high-grade and in tumors that recurred. Both *SFRP4* and *CXCL14* on the other hand are inhibitors of prostate tumor growth [55,56]. *SFRP4* is up-regulated in tumors of high-grade, and *CXCL14* in tumors of advanced stage. A small deletion involving chromosomal

band 21q22.3 fuses all coding exons of *ERG* to androgen-related sequences in the promoter of the prostate-specific *TMPRSS2* gene. This chromosomal rearrangement is a highly prevalent oncogenic alteration in prostate tumor cells and leads to an aberrant expression of the *ERG* proto-oncogene, important for early prostate carcinogenesis [57]. In this data set, *ERG* is overexpressed in tumors in which the cancer metastasized to distant lymph nodes. It has been shown that this genetic biomarker is a strong prognostic factor for disease recurrence, and can be used for early detection and outcome prediction in prostate cancer [58]. *VAV3*, an oncogene involved in development and progression of prostate cancer, is up-regulated in tumors that

**Table 6****Comparison of our kernel-based integration approach with the ensemble approach**

| Outcome    | AUC (SE)*: $MPT_1/MG$ | AUC (SE)*: ensemble approach | p-value       |
|------------|-----------------------|------------------------------|---------------|
| WHEELER    | 0.9269 (0.0425)       | 0.9500 (0.0339)              | 0.6160        |
| pN-STAGE   | 0.9870 (0.0135)       | 0.9253 (0.0432)              | 0.1422        |
| CRM        | 0.9630 (0.0344)       | 0.7860 (0.0783)              | <b>0.0384</b> |
| GRADE      | 0.9006 (0.0413)       | 0.8567 (0.0521)              | 0.3745        |
| STAGE      | 0.8528 (0.0550)       | 0.8304 (0.0582)              | 0.6836        |
| METASTASIS | 0.9868 (0.0121)       | 0.9452 (0.0309)              | 0.1313        |
| RECURRENCE | 0.7857 (0.0934)       | 0.4545 (0.1352)              | <b>0.0182</b> |

\*Area under the ROC curve (standard error) obtained with leave-one-out. †Comparison in AUC between the best models obtained with our strategy ( $MPT_1$  for rectal cancer,  $MG$  for prostate cancer) and the corresponding ensemble models based on the same number of features [46].

metastasized [59]. It has previously been shown that strong overexpression of *TIAM1* is significantly associated with disease recurrence and a decreased disease-free survival [60]. Also, *JAG1* is significantly associated with recurrence [61] and plays a role in cell growth, progression, and metastasis. In this data set, both genes are up-regulated in the group of tumors that recurred. Finally, several germline mutations or variants in *RNASEL* have been observed among hereditary prostate cancer cases, indicating that polymorphic changes within the *RNASEL* gene may be associated with increased risk of familial but not sporadic prostate cancer [62]. A list of all the genes and CNVs selected by the models  $MG$  are shown in Additional data file 3. As for rectal cancer, the outcomes for prostate cancer seem to be characterized by mainly different sets of features. Five genes overlap between at least two outcomes (*ERG*, *AHSG*, *SEMA4G*, *F5*, and *ALOX15B*), while the same holds for four CNVs of the genes *GPD1L*, *KCTD12*, *SMYD5*, and *TRO*.

#### Comparison with an ensemble approach

To assess the benefit of our kernel-based integration approach over standard data fusion techniques, we implemented an ensemble approach in which each data set gives rise to a separate LS-SVM classifier. These individual LS-SVM models were built similarly to the step A models, with the same number of genes, proteins or CNVs selected as included in the best models  $MPT_1$  and  $MG$ . Subsequently, as a late integration step, the continuous outputs of these models were added.

For the study on rectal cancer, the AUC values of the ensemble models integrating the microarray and proteomics data set gathered at  $T_1$ , and the corresponding AUC values of the best model obtained with our strategy ( $MPT_1$ ) are shown in Table 6. The *P*-values of the significance tests comparing the ROC curves are reported as well [46]. For CRM, our strategy was significantly better than the ensemble approach at a significance level of 0.05. For WHEELER and pN-

STAGE, the AUC values did not differ significantly. Similarly for the study on prostate cancer, the AUC values of  $MG$  were compared with the AUC values of the ensemble models combining microarray and genomics data (Table 6). For all four outcomes, the AUC of  $MG$  was better than the AUC of the ensemble models, although being significantly better for RECURRENCE only.

#### Correlation analysis

We additionally verified whether, in both cases, data from multiple layers of molecular biology were complementary. After mapping the entities of the data sets based on their entrez gene IDs, we investigated the correlation between the microarray and proteomics data of rectal cancer on the one hand, and between the microarray and genomics data of prostate cancer on the other hand. Using the Spearman correlation coefficient, there was no significant correlation for rectal cancer between the abundances of the 90-92 proteins and their corresponding transcripts at a significance level of 0.05. The microarray and genomics data sets for prostate cancer were slightly more correlated. While for GRADE the 6 genes selected by the model  $MG$  did not correlate with their DNA expression, 2 of the 42 selected genes for STAGE were significantly correlated ( $P < 0.05$ ). For METASTASIS and RECURRENCE, there was a significant correlation for one and three genes, respectively. The regions, with involved CNVs selected from the genomics data, were also compared with the regions in which the selected genes from the microarray data were located. For the majority of regions, there was no overlap. For the other regions with the same rough chromosomal location, the genes selected by both data sets were different.

#### Discussion

The proposed integration approach has been applied to two patient data sets, each with two high-throughput data sources. Microarray and proteomics were gathered from 36 patients with rectal cancer at two time points during

preoperative treatment, while microarray and genomics were gathered from 55 patients with prostate cancer. To verify the merit of our integration approach over the use of a single omics data source, models were built for classifying cancer patients according to therapy response, prognostic factors, metastasis, or recurrence. In many studies, only single data sources are explored for the development of such profiles. However, in our opinion, a single layer of molecular information is inadequate to explain the complete network of molecules underlying a disease. In this study, LS-SVMs were first built on all data sets individually (Figure 1). Next, we manually integrated data measured at multiple time points by building LS-SVMs using the change in expression between two time points. Because the integration of data may be more complex than the change in expression over time, we subsequently applied an intermediate integration approach in which data from multiple omics were combined at the kernel level within the patient domain.

For the data on rectal cancer, all three outcomes - a tumor regression grading system and two prognostic factors - could be predicted most accurately and most cost-efficiently with an AUC ranging from 0.9269 to 0.9870 when fusing microarray and proteomics data gathered during therapy ( $MPT_0$ ; Table 2). For WHEELER, for example,  $MPT_0$  performance is better than each of the models based on data from an individual technology ( $MT_0$  and  $PT_0$ ), as is the case for  $MPT_{01}$  compared to  $MT_{01}$  and  $PT_{01}$ . This trend of increased performance when combining data from two different technologies was further confirmed by our second data set for prostate cancer patients. Best results for the prediction of grade, stage, metastasis, and recurrence were obtained when integrating microarray and genomics data ( $MG$ ). The corresponding AUC values were 0.9006, 0.8528, 0.9868, and 0.7857, respectively (Table 4). For many of the genes, proteins, and CNVs included in these models, involvement in rectal or prostate cancer has been defined, indicating the reliability of the selected features (Tables 3 and 5). These models were compared with models obtained with an ensemble approach in which classifiers are combined instead of data sets at the kernel level. Globally, our approach performed better, although not always significantly (Table 6).

By looking at the correlation between two data sets gathered from the same set of patients, we show that data from different layers are mainly complementary. For rectal cancer, there was a lack of correlation between the selected genes and their corresponding proteins. Also, the selected proteins did not significantly correlate with their transcript level, suggesting alternative splicing and post-translational modification. With newer technologies such as mass spectrometry, the whole proteome will become measurable. For prostate cancer, up to three genes included in the model  $MG$  were significantly correlated with their corresponding CNV.

More specific for the study on rectal cancer, we can conclude from Table 2 that data gathered after an initial dose of cetuximab are more informative for prediction of therapy response than data gathered before the start of the therapy. Neither microarray nor proteomics data can predict the outcomes more accurately at  $T_0$  than  $T_1$ , except for the proteomics data at  $T_0$  being more informative for the prediction of CRM. Moreover, when combining both data types at one time point ( $MPT_0$  and  $MPT_1$ ), the models applicable after the initial dose of cetuximab outperform those at  $T_0$ .

We acknowledge that the models proposed in this manuscript are quite expensive. Applying a model for rectal cancer would require microarray and/or proteomics data, gathered at one or two time points during therapy. However, we have attempted to keep the cost to a minimum. The performance difference between models combining two data sets, only requiring a sample to be taken at one time point or one technology to be applied at two time points, and models requiring a sample to be taken at both time points and both technologies to be performed was minimal and not statistically significant. We therefore chose the best model among the models based on two data sets. We admit that there may exist other, less expensive data sources that can contain complementary information as well. Firstly, clinical information is routinely gathered during therapy, such as tumor size, tumor location and number of positive lymph nodes. However, we only had access to the clinical parameter age, for which we performed an additional analysis to verify whether this parameter could be of use. A univariate analysis based on the Wilcoxon rank sum test showed no significant difference in age between the two classes of samples according to the considered outcomes. In a multivariate logistic regression model, the parameter age was not significant as well. Secondly, there is an increasing need for multi-modal studies in which, among others, clinical, genomic and genetic data are collected. Also, imaging, such as computed tomography (CT) and magnetic resonance imaging (MRI) can be a potential predictor to use in combination with high-throughput data sources. Such studies are required to determine which data sets are most relevant for the problem at hand and which data sets should be combined to become good performing, affordable models that are clinically applicable.

## Conclusions

The results suggest that the use of our integration approach on experimental data from multiple levels in the genome can improve the performance of decision support in cancer. For both data sets studied in this manuscript, combining high-throughput data sets (transcriptomics with proteomics, or genomics with transcriptomics) outperformed the models based on data from a single layer of biological information, independent of the outcome considered for prediction.

These results emphasize the need for comprehensive multi-modal data gathered with high-throughput technologies as well as imaging, because it is unknown which technologies, and thus which levels of molecular biology, are the most relevant for prognostic prediction. We acknowledge that this will substantially increase costs in a first exploratory phase. However, this is a necessary investment to ultimately obtain cost-efficient models usable in patient tailored therapy.

In the near future, we will compare our kernel-based integration method with a Bayesian network integration framework. These frameworks are complementary. We also plan to apply an ensemble approach for integrating these two frameworks because more accurate classifiers are not only obtained by combining different data types but also by combining individual decisions of multiple classifiers. In this way, the advantages of both methods can be exploited.

### Abbreviations

AUC, area under the ROC curve; CGH, comparative genomic hybridization; CNV, copy number variation; CRM, circumferential margin involvement; DEDS, differential expression via distance synthesis; EGFR, epidermal growth factor receptor; *G*, model based on genomics data; LOO, leave-one-out; LS-SVM, least squares support vector machine; *M*, model based on microarray data; *MG*, model based on both microarray and genomics data; *MPT*<sub>0</sub>, model based on microarray and proteomics data at *T*<sub>0</sub>; *MPT*<sub>1</sub>, model based on microarray and proteomics data at *T*<sub>1</sub>; *MPT*<sub>01</sub>, model based on all data (microarray and proteomics data at both timepoints); *MT*<sub>0</sub>, model based on microarray data at *T*<sub>0</sub>; *MT*<sub>1</sub>, model based on microarray data at *T*<sub>1</sub>; *MT*<sub>01</sub>, model based on microarray data at both time points; *MT*<sub>0</sub>-*T*<sub>1</sub>, model based on change in gene expression between *T*<sub>0</sub> and *T*<sub>1</sub>; *PT*<sub>0</sub>, model based on proteomics data at *T*<sub>0</sub>; *PT*<sub>1</sub>, model based on proteomics data at *T*<sub>1</sub>; *PT*<sub>01</sub>, model based on proteomics data at both time points; *PT*<sub>0</sub>-*T*<sub>1</sub>, model based on change in protein abundances between *T*<sub>0</sub> and *T*<sub>1</sub>; ROC, receiver operating characteristic; SVM, support vector machine; *T*<sub>0</sub>, time point before treatment; *T*<sub>1</sub>, time point after the first loading dose of cetuximab but before the start of radiotherapy with capecitabine; *T*<sub>2</sub>, time point at moment of surgery.

### Competing interests

The authors declare that they have no competing interests.

### Authors' contributions

ADa performed the kernel-based integration modeling and drafted the manuscript. OG, FO, and JS participated in the design and implementation of the framework. ADa and OG performed pre-processing of the data. OG, JS, and BDM helped draft the manuscript. ADe, JPM, and KH provided

clinical input, looked up patient records in the database, performed sample annotation, and gathered follow-up of patients. All authors read and approved the final manuscript.

### Additional data files

The following additional data files are available with the online version of this paper. Additional data file 1 shows the ROC curves of the optimal LS-SVM models for all considered combinations of data sets shown in Tables 2 and 4. Additional data file 2 shows the results for the prediction of WHEELER, pN-STAGE, and CRM in rectal cancer, using step C models for which a sample is required at both time points and for which both technologies need to be performed. Additional data file 3 contains additional tables 1-3 showing all genes and proteins selected by the best performing models MPT1 for the prediction of WHEELER, pN-STAGE, and CRM in rectal cancer. Additional data file 3 also contains additional tables 4-7 showing, for prostate cancer, the genes and CNVs selected by the best performing models MG for the prediction of GRADE, STAGE, METASTASIS, and RECURRENCE. All tables in additional data file 3 show the number of LOO iterations in which each gene, protein, or CNV was selected, their chromosomal region, and whether it is up- or down-regulated.

### Acknowledgements

ADa is research assistant of the Fund for Scientific Research-Flanders (FWO-Vlaanderen). BDM is a full professor at the Katholieke Universiteit Leuven, Belgium. The authors are grateful to Anja von Heydebreck, Detlef Guessow and Christopher Stroh for their contribution at Merck Serono. This work is partially supported by the following. (1) Research Council KUL: GOA AMBioRICS, CoE EF/05/007 SymBioSys, PROMETA, several PhD/postdoc and fellow grants. (2) Flemish Government: (a) FWO: PhD/postdoc grants, projects G.0241.04 (Functional Genomics), G.0499.04 (Statistics), G.0318.05 (subfunctionalization), G.0302.07 (SVM/Kernel), research communities (ICCoS, ANMMM, MLDM); (b) IWT: PhD Grants, GBOU-McKnow-E (Knowledge management algorithms), GBOU-ANA (biosensors), TAD-BioScope-IT, Silicos; SBO-BioFrame, SBO-MoKa, TBM-Endometriosis. (3) Belgian Federal Science Policy Office: IUAP P6/25 (BioMaGNet, Bioinformatics and Modeling: from Genomes to Networks, 2007-2011). (4) EU-RTD: ERNSI: European Research Network on System Identification; FP6-NoE Biopattern; FP6-IP e-Tumors, FP6-MC-EST Biop-train, FP6-STREP Strokemap.

### References

1. Shawe-Taylor J, Cristianini N: *Kernel Methods for Pattern Analysis*. Cambridge: Cambridge University Press; 2004.
2. Bhaskar H, Hoyle DC, Singh S: **Machine learning in bioinformatics: a brief survey and recommendations for practitioners**. *Comput Biol Med* 2006, **36**:1104-1125.
3. Suykens JAK, Vandewalle J: **Least squares support vector machine classifiers**. *Neural Processing Lett* 1999, **9**:293-300.
4. Suykens JAK, Van Gestel T, De Brabanter J, De Moor B, Vandewalle J: *Least Squares Support Vector Machines*. Singapore: World Scientific; 2002.
5. Cawley GC: **Leave-one-out cross-validation based model selection criteria for weighted LS-SVMs**. *Proc Int Joint Conf on Neural Networks 2006*:1661-1668.
6. Alon A, Barkai N, Notterman DA, Gish K, Ybarra S, Mack D, Levine AJ: **Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays**. *Proc Natl Acad Sci USA* 1999, **96**:6745-6750.

7. Golub TR, Slonim DK, Tamayo P, Huard C, Gaasenbeek M, Mesirov JP, Coller H, Loh ML, Downing JR, Caligiuri MA, Bloomfield CD, Lander ES: **Molecular classification of cancer: class discovery and class prediction by gene expression monitoring.** *Science* 1999, **286**:531-537.
8. Cardoso F, van't Veer L, Rutgers E, Loi S, Mook S, Piccart-Gebhart MJ: **Clinical application of the 70-gene profile: the MINDACT trial.** *J Clin Oncol* 2008, **26**:729-735.
9. Sparano JA: **TAILORx: trial assigning individualized options for treatment (Rx).** *Clin Breast Cancer* 2006, **7**:347-350.
10. Sparano JA, Paik S: **Development of the 21-gene assay and its application in clinical practice and clinical trials.** *J Clin Oncol* 2008, **26**:721-728.
11. Pinkel D, Albertson DG: **Array comparative genomic hybridization and its applications in cancer.** *Nat Genet* 2005, **37**:S11-S17.
12. Esteller M: **Epigenetics in cancer.** *N Engl J Med* 2008, **358**:1148-1159.
13. Redon R, Ishikawa S, Fitch KR, Feuk L, Perry GH, Andrews TD, Fiegler H, Shaperro MH, Carson AR, Chen W, Cho EK, Dallaire S, Freeman JL, González JR, Gratacòs M, Huang J, Kalaitzopoulos D, Komura D, MacDonald JR, Marshall CR, Mei R, Montgomery L, Nishimura K, Okamura K, Shen F, Somerville MJ, Tchinda J, Valsesia A, Woodwark C, Yang F, et al.: **Global variation in copy number in the human genome.** *Nature* 2006, **444**:444-454.
14. Frohling S, Dohner H: **Chromosomal abnormalities in cancer.** *N Engl J Med* 2008, **359**:722-734.
15. Kolch W, Mischak H, Pitt AR: **The molecular make-up of a tumor: proteomics in cancer research.** *Clin Sci* 2005, **108**:369-383.
16. Aebersold R, Mann M: **Mass spectrometry-based proteomics.** *Nature* 2003, **422**:198-207.
17. MacBeath G, Schreiber SL: **Printing proteins as microarrays for high-throughput function determination.** *Science* 2000, **289**:1760-1763.
18. Cooper GM, Zerr T, Kidd JM, Eichler EE, Nickerson DA: **Systematic assessment of copy number variant detection via genome-wide SNP genotyping.** *Nat Genet* 2008, **40**:1199-1203.
19. Tibshirani RJ, Efron B: **Pre-validation and inference in microarrays.** *Stat Appl Genet Mol Biol* 2002, **1**:Article 1.
20. Nevins JR, Huang ES, Dressman H, Pittman J, Huang AT, West M: **Towards integrated clinico-genomic models for personalized medicine: combining gene expression signatures and clinical factors in breast cancer outcomes prediction.** *Hum Mol Genet* 2003, **12**:R153-R157.
21. Wang SM, Ooi LL, Hui KM: **Identification and validation of a novel gene signature associated with the recurrence of human hepatocellular carcinoma.** *Clin Cancer Res* 2007, **13**:6275-6283.
22. Mathew JP, Taylor BS, Bader GD, Pyarajan S, Antoniotti M, Chinnaiyan AM, Sander C, Burakoff SJ, Mishra B: **From bytes to bedside: data integration and computational biology for translational cancer research.** *PLoS Comput Biol* 2007, **3**:e12.
23. Fridlyand J, Snijders AM, Ylstra B, Li H, Olshen A, Segraves R, Dairkee S, Tokuyasu T, Ljung BM, Jain AN, McLennan J, Ziegler J, Chin K, Devries S, Feiler H, Gray JW, Waldman F, Pinkel D, Albertson DG: **Breast tumor copy number aberration phenotypes and genomic instability.** *BMC Cancer* 2006, **6**:96.
24. Tomioka N, Oba S, Ohira M, Misra A, Fridlyand J, Ishii S, Nakamura Y, Isogai E, Hirata T, Yoshida Y, Todo S, Kanedo Y, Albertson DG, Pinkel D, Feuerstein BG, Nakagawara A: **Novel risk stratification of patients with neuroblastoma by genomic signature, which is independent of molecular signature.** *Oncogene* 2008, **27**:441-449.
25. Waters KM, Pounds JG, Thrall BD: **Data merging for integrated microarray and proteomic analysis.** *Brief Funct Genomic Proteomic* 2006, **5**:261-272.
26. Goble C, Stevens R: **State of the nation in data integration for bioinformatics.** *J Biomed Inform* 2008, **41**:687-693.
27. Bitton DA, Okoniewski MJ, Connolly Y, Miller CJ: **Exon level integration of proteomics and microarray data.** *BMC Bioinformatics* 2008, **9**:118.
28. Lanckriet GRG, De Bie T, Cristianini N, Jordan MI, Noble WS: **A statistical framework for genomic data fusion.** *Bioinformatics* 2004, **20**:2626-2635.
29. Daemen A, Gevaert O, Moor BD: **Integration of clinical and microarray data with kernel methods.** *Conf Proc IEEE Eng Med Biol Soc* 2007:5411-5415.
30. Lapointe J, Li C, Higgins JP, van de Rijn M, Bair E, Montgomery K, Ferrari M, Egevad L, Rayford W, Bergerheim U, Ekman P, DeMarzo AM, Tibshirani R, Botstein D, Brown PO, Brooks JD, Pollack JR: **Gene expression profiling identifies clinically relevant subtypes of prostate cancer.** *Proc Natl Acad Sci USA* 2004, **101**:811-816.
31. Lapointe J, Li C, Giacomini CP, Salari K, Huang S, Wang P, Ferrari M, Hernandez-Boussard T, Brooks JD, Pollack JR: **Genomic profiling reveals alternative genetic pathways of prostate tumorigenesis.** *Cancer Res* 2007, **67**:8504-8510.
32. Machiels JP, Sempoux C, Scalliet P, Coche JC, Humblet Y, Van Cutsem E, Kerger J, Canon JL, Peeters M, Aydin S, Laurent S, Kartheuser A, Coster B, Roels S, Daisne JF, Honhon B, Duck L, Kirkove C, Bonny MA, Haustermans K: **Phase I/II study of preoperative cetuximab, capecitabine, and external beam radiotherapy in patients with rectal cancer.** *Ann Oncol* 2007, **18**:738-744.
33. Irizarry RA, Hobbs B, Collin F, Beazer-Barclay YD, Antonellis KJ, Scherf U, Speed TP: **Exploration, normalization and summaries of high density oligonucleotide array probe level data.** *Biostatistics* 2003, **4**:249-264.
34. Wheeler JMD, Warren BF, Mortensen NJ, Ekanyaka N, Kulacoglu H, Jones AC, George BD, Kettlewell MGW: **Quantification of histologic regression of rectal cancer after irradiation.** *Dis Colon Rectum* 2002, **45**:1051-1056.
35. Machiels JP, Aydin S, Bonny MA, Hammouch F, Sempoux C: **What is the best to predict disease-free survival after preoperative radiochemotherapy for rectal cancer patients: tumor regression grading, nodal status or circumferential resection margin invasion?** *J Clin Oncol* 2006, **24**:1319-1321.
36. Adam IJ, Mohamdee MO, Martin IG, Scott N, Finan PJ, Johnston D, Dixon MF, Quirke P: **Role of circumferential margin involvement in the local recurrence of rectal cancer.** *Lancet* 1994, **344**:707-711.
37. Quirke P, Durdey P, Dixon MF, Williams NS: **Local recurrence of rectal adenocarcinoma due to inadequate surgical resection: histopathological study of lateral tumor spread and surgical excision.** *Lancet* 1986, **2**:996-999.
38. Troyanskaya O, Cantor M, Sherlock G, Brown P, Hastie T, Tibshirani R, Botstein D, Altman RB: **Missing value estimation methods for DNA microarrays.** *Bioinformatics* 2001, **17**:520-525.
39. Gleason DF: **Classification of prostatic carcinomas.** *Cancer Chemother Rep* 1966, **50**:125-128.
40. Scholkopf B, Tsuda K, Vert JP: **Kernel Methods in Computational Biology.** Cambridge, MA: MIT Press; 2004.
41. Vapnik V: **Statistical Learning Theory.** New York: Wiley; 1998.
42. Pochet N, De Smet F, Suykens J, Moor BD: **Systematic benchmarking of microarray data classification: assessing the role of nonlinearity and dimensionality reduction.** *Bioinformatics* 2004, **20**:3185-3195.
43. Lai C, Reinders MJT, van't Veer LJ, Wessels LFA: **A comparison of univariate and multivariate gene selection techniques for classification of cancer datasets.** *Bioinformatics* 2006, **7**:235-244.
44. Yang YH, Xiao Y, Segal MR: **Identifying differentially expressed genes from microarray experiments via statistic synthesis.** *Bioinformatics* 2005, **21**:1084-1093.
45. Li W, Yang Y: **How many genes are needed for a discriminant microarray data analysis.** In *Methods of Microarray Data Analysis*. Edited by Lin SM, Johnson KF. Kluwer Academic; 2002:137-150.
46. Hanley JA, McNeil BJ: **A method of comparing the areas under receiver operating characteristics curves derived from the same cases.** *Radiology* 1983, **148**:839-843.
47. Pavlidis P, Weston J, Cai J, Grundy WN: **Gene functional classification from heterogeneous data.** In *Proceedings of the Fifth Annual International Conference on Computational Biology: April 22-25, 2001; Montreal, Quebec, Canada*. New York, NY: ACM; 2001:242-252.
48. Zien A, Ong CS: **Multiclass multiple kernel learning.** In *Proceedings of the 24th International Conference on Machine Learning: June 20-24, 2007; Corvallis, Oregon*. New York, NY: ACM; 2007:1191-1198.
49. Zhang W, Park DJ, Lu B, Yang DY, Gordon M, Groshen S, Yun J, Press OA, Vallbohmer D, Rhodes K, Lenz HJ: **Epidermal growth factor receptor gene polymorphisms predict pelvic recurrence in patients with rectal cancer treated with chemoradiation.** *Clin Cancer Res* 2005, **11**:600-605.
50. Maihofner C, Charalambous MP, Bhambra U, Lightfoot T, Geisslinger G, Gooderham NJ; The Colorectal Cancer Group: **Expression of cyclooxygenase-2 parallels expression of interleukin-1beta, interleukin-6 and NF-kappaB in human colorectal cancer.** *Carcinogenesis* 2003, **24**:665-671.
51. Sawhney RS, Sharma B, Humphrey LE, Brattain MG: **Integrin alpha2 and extracellular signal-regulated kinase are functionally linked in highly malignant autocrine transforming growth factor-alpha-driven colon cancer cells.** *J Biol Chem* 2003, **278**:19861-19869.



52. Rubie C, Frick VO, Pfeil S, Wagner M, Kollmar O, Kopp B, Graber S, Rau BM, Schilling MK: **Correlation of IL-8 with induction, progression and metastatic potential of colorectal cancer.** *World J Gastroenterol* 2007, **13**:4996-5002.
53. Louhimo J, Carpelan-Holmstrom M, Alfthan H, Stenman UH, Jarvinen HJ, Haglund C: **Serum HCG beta, CA 72-4 and CEA are independent prognostic factors in colorectal cancer.** *Int J Cancer* 2002, **101**:545-548.
54. Bhatia B, Maldonado CJ, Tang S, Chandra D, Klein RD, Chopra D, Shappell SB, Yang P, Newman RA, Tang DG: **Subcellular localization and tumor-suppressive functions of 15-lipoxygenase 2 (15-LOX2) and its splice variants.** *J Biol Chem* 2003, **278**:25091-25100.
55. Horvath LG, Lelliott JE, Kench JG, Lee CS, Williams ED, Saunders DN, Grvgiel JJ, Sutherland RL, Henshall SM: **Secreted frizzled-related protein 4 inhibits proliferation and metastatic potential in prostate cancer.** *Prostate* 2007, **67**:1081-1090.
56. Schwarze SR, Luo J, Isaacs WB, Jarrard DF: **Modulation of CXCL14 (BRAK) expression in prostate cancer.** *Prostate* 2005, **64**:67-74.
57. Furusato B, Gao CL, Ravindranath L, Chen Y, Cullen J, McLeod DG, Dobi A, Srivastava S, Petrovics G, Sesterhenn IA: **Mapping of TMPRSS2-ERG fusions in the context of multi-focal prostate cancer.** *Mod Pathol* 2008, **21**:67-75.
58. Nam RK, Sugar L, Yang W, Srivastava S, Klotz LH, Yang LY, Stanimirovic A, Encioiu E, Neill M, Loblaw DA, Trachtenberg J, Narod SA, Seth A: **Expression of the TMPRSS2:ERG fusion gene predicts cancer recurrence after surgery for localised prostate cancer.** *Br J Cancer* 2007, **97**:1690-1695.
59. Dong Z, Liu Y, Lu S, Wang A, Lee K, Wang LH, Revelo M, Lu S: **Vav3 oncogene is overexpressed and regulates cell growth and androgen receptor activity in human prostate cancer.** *Mol Endocrinol* 2006, **20**:2315-2325.
60. Engers R, Mueller M, Walter A, Collard JG, Willers R, Gabbert HE: **Prognostic relevance of Tiam1 protein expression in prostate carcinomas.** *Br J Cancer* 2006, **95**:1081-1086.
61. Santagata S, Demichelis F, Riva A, Varambally S, Hofer MD, Kutok JL, Kim R, Tang J, Montie JE, Chinnaiyan AM, Rubin MA, Aster JC: **JAGGED1 expression is associated with prostate cancer metastasis and recurrence.** *Cancer Res* 2004, **64**:6854-6857.
62. Silverman RH: **Implications for RNase L in prostate cancer biology.** *Biochemistry* 2003, **42**:1805-1812.
63. Rajc D, Mukhtar H, Oshowo A, Clark CI: **What proportion of patients referred to secondary care with iron deficiency anemia have colon cancer?** *Dis Colon Rectum* 2007, **50**:1211-1214.
64. Ciardiello F, Tortora G: **Epidermal growth factor receptor (EGFR) as a target in cancer therapy: understanding the role of receptor expression and other molecular determinants that could influence the response to anti-EGFR drugs.** *Eur J Cancer* 2003, **39**:1348-1354.
65. Kim TD, Song KS, Li G, Choi H, Park HD, Lim K, Hwang BD, Yoon WH: **Activity and expression of urokinase-type plasminogen activator and matrix metalloproteinases in human colorectal cancer.** *BMC Cancer* 2006, **6**:211.
66. Uner A, Akcali Z, Unsal D: **Serum levels of soluble E-selectin in colorectal cancer.** *Neoplasma* 2004, **51**:269-274.
67. Eksioglu EA, Mahmood SS, Chang M, Reddy V: **GM-CSF promotes differentiation of human dendritic cells and T lymphocytes toward a predominantly type I proinflammatory response.** *Exp Hematol* 2007, **35**:1163-1171.
68. Zinzindohoue F, Lecomte T, Ferraz JM, Houllier AM, Cugnenc PH, Berger A, Blons H, Laurent-Puig P: **Prognostic significance of MMP-1 and MMP-3 functional promoter polymorphisms in colorectal cancer.** *Clin Cancer Res* 2005, **11**:594-599.
69. Zhang Y, Lai M, Lv B, Gu X, Wang H, Zhu Y, Zhu Y, Shao L, Wang G: **Overexpression of Reg IV in colorectal adenoma.** *Cancer Lett* 2003, **200**:69-76.
70. Ahn DH, Crawley SC, Hokari R, Kato S, Yang SC, Li JD, Kim YS: **TNF-alpha activates MUC2 transcription via NF-kappaB but inhibits via JNK activation.** *Cell Physiol Biochem* 2005, **15**:29-40.
71. Kummola L, Hala J, Kivelainen JM, Kivela AJ, Saarnio J, Karttunen T, Parkkila S: **Expression of a novel carbonic anhydrase, CA XIII, in normal and neoplastic colorectal mucosa.** *BMC Cancer* 2005, **5**:41.
72. Gropcke S, Mannone J, Weber B, Staub E, Heinze M, Klamann I, Pilarsky C, Hermann K, Castanos-Velez E, Ropcke S, Mann B, Rosenthal A, Buhr HJ: **Differential expression of genes encoding tight junction proteins in colorectal cancer: frequent dysregulation of claudin-1, -8 and -12.** *Int J Colorectal Dis* 2007, **22**:651-659.
73. Viet HT, Wagsater D, Hugander A, Dimberg J: **Interleukin-1 receptor antagonist gene polymorphism a gs in human colorectal cancer.** *Oncol Rep* 2005, **14**:915-918.
74. Kloor M, Michel S, Buckowitz B, Ruschoff J, Buttner R, Holinski-Feder E, Dippold W, Wagner R, Tariverdian M, Benner A, Schwitalle Y, Kuchenbuch B, von Knebel Doeberitz M: **Beta2-microglobulin mutations in microsatellite unstable colorectal tumors.** *Int J Cancer* 2007, **121**:454-458.
75. Youssef EM, Chen Xq, Higuchi E, Kondo Y, Garcia-Manero G, Lotan R, Issa JPJ: **Hypermethylation and silencing of the putative tumor suppressor Tazartene-induced gene 1 in human cancers.** *Cancer Res* 2004, **64**:2411-2417.
76. Muc-Wierozon M, Nowakowska-Zajdel E, Kokot T, Kozowicz A, Zubelewicz B, Klakla K, Mazurek U, Cholewa K, Wilczok T, Wierozon J, Sosada K: **Genetic dysregulation of gene coding tumor necrosis factor alpha receptors (TNFalpha Rs) in colorectal cancer cells.** *J Exp Clin Cancer Res* 2004, **23**:651-660.
77. Maeda K, Kang SM, Sawada T, Nishiguchi Y, Yashiro M, Ogawa Y, Ohira M, Ishikawa T, Hirakawa YS, Chung K: **Expression of intercellular adhesion molecule-1 and prognosis in colorectal cancer.** *Oncol Rep* 2002, **9**:511-514.
78. Ferroni P, Palmirotta R, Spila A, Martini F, Raparelli V, Fossile E, Mariotti S, Del Monte G, Buonomo O, Roselli M, Guadagni F: **Prognostic significance of adiponectin levels in non-metastatic colorectal cancer.** *Anticancer Res* 2007, **27**:483-489.
79. Miyayama K, Kato Y, Nakamura T, Matsumura M, Amaya H, Horiuchi T, Chiba Y, Tanaka K: **Expression and role of thrombospondin-1 in colorectal cancer.** *Anticancer Res* 2002, **22**:3941-3948.
80. Wan Y, Wu N, Wang Z, Ju X, Zhu J, Liu Y, Tang J, Huang Y: **Relationship between tissue factor expression and hepatic metastasis and prognosis in rectal cancer.** *Zhonghua Zhong Liu Za Zhi* 2002, **24**:378-380.
81. Bethke L, Webb E, Sellick G, Rudd M, Penegar S, Withey L, Qureshi M, Houlston R: **Polymorphisms in the cytochrome P450 genes CYP1A2, CYP1B1, CYP3A4, CYP3A5, CYP11A1, CYP17A1, CYP19A1 and colorectal cancer risk.** *BMC Cancer* 2007, **7**:123.
82. Cross NA, Chandrasekharan S, Jokonya N, Fowles A, Hamdy FC, Buttler DJ, Eaton CL: **The expression and regulation of ADAMTS-1, -4, -5, -9, and -15, and TIMP-3 by TGFbeta1 in prostate cells: relevance to the accumulation of versican.** *Prostate* 2005, **63**:269-275.
83. Hudolin T, Juretic A, Spagnoli GC, Pasini J, Bandic D, Heberer M, Kosicek M, Cacic M: **Immunohistochemical expression of tumor antigens MAGE-A1, MAGE-A3/4, and NY-ESO-1 in cancerous and benign prostatic tissue.** *Prostate* 2006, **66**:13-18.
84. Ishii K, Usui S, Sugimura Y, Yoshida S, Hioki T, Tatematsu M, Yamamoto H, Hirano K: **Aminopeptidase N regulated by zinc in human prostate participates in tumor cell invasion.** *Int J Cancer* 2001, **92**:49-54.
85. Diss JK, Faulkes DJ, Walker MM, Patel A, Foster CS, Budhrum-Mahadeo V, Djamgoz MB, Latchman DS: **Brrn-3a neuronal transcription factor functional expression in human prostate cancer.** *Prostate Cancer Prostatic Dis* 2006, **9**:83-91.
86. Cross DS, Burmester JK: **Functional characterization of the GDEP promoter and three enhancer elements in retinoblastoma and prostate cell lines.** *Med Oncol* 2008, **25**:40-49.
87. Wolfgang CD, Essand M, Lee B, Pastan I: **T-cell receptor gamma chain alternate reading frame protein (TARP) expression in prostate cancer cells leads to an increased growth rate and induction of caveolins and amphiregulin.** *Cancer Res* 2001, **61**:8122-8126.
88. Descazeaud A, de la Taille A, Allory Y, Faucon H, Salomon L, Bismar T, Kim R, Hofer MD, Chopin D, Abbou CC, Rubin MA: **Characterization of ZAG protein expression in prostate cancer using a semi-automated microscope system.** *Prostate* 2006, **66**:1037-1043.
89. Sahni A, Simpson-Haidaris PJ, Sahni SK, Vaday GG, Francis CW: **Fibronogen synthesized by cancer cells augments the proliferative effect of fibroblast growth factor-2 (FGF-2).** *J Thromb Haemost* 2008, **6**:176-183.
90. Bandyopadhyay S, Wang Y, Zhan R, Pai SK, Watabe M, Iizumi M, Furuta E, Mohinta S, Liu W, Hirota S, Hosobe S, Tsukada T, Miura K, Takano Y, Saito K, Commes T, Piquemal D, Hai T, Watabe K: **The tumor metastasis suppressor gene Drg-1 down-regulates the expression of activating transcription factor 3 in prostate cancer.** *Cancer Res* 2006, **66**:11983-11990.