

Heterogeneous Data Fusion to Type Brain Tumor Biopsies

Vangelis Metsis¹, Heng Huang¹, Fillia Makedon¹, and Aria Tzika²

¹ University of Texas at Arlington, USA, vangelj.meci@mavs.uta.edu, heng@uta.edu, makedon@uta.edu

² NMR Surgical Laboratory, Department of Surgery, Harvard Medical School and Massachusetts General Hospital, Boston, USA, tzika@hms.harvard.edu

Abstract Current research in biomedical informatics involves analysis of multiple heterogeneous data sets. This includes patient demographics, clinical and pathology data, treatment history, patient outcomes as well as gene expression, DNA sequences and other information sources such as gene ontology. Analysis of these data sets could lead to better disease diagnosis, prognosis, treatment and drug discovery. In this paper, we use machine learning algorithms to create a novel framework to perform the heterogeneous data fusion on both metabolic and molecular datasets, including state-of-the-art high-resolution magic angle spinning (HRMAS) proton (¹H) Magnetic Resonance Spectroscopy and gene transcriptome profiling, to intact brain tumor biopsies and to identify different profiles of brain tumors. Our experimental results show our novel framework outperforms any analysis using individual dataset.

1 Introduction

Brain tumors are the second most common cancer of childhood, and comprise approximately 25% of all pediatric cancers. Over 3,400 children are diagnosed in the U.S. each year; of that, about 2,600 will be under the age of 15. Brain tumors are the leading cause of solid tumor cancer death in children; they are the third leading cause of cancer death in young adults ages 20-39. Many researchers are looking for efficient and reliable ways to early diagnose brain tumor types and detect related biomarkers through different biomedical images or biological data. The machine learning algorithms have been playing the most important role during those heterogamous biomedical/biological datasets analysis to classify different brain tumor types and detect biomarkers.

Magnetic resonance spectroscopic (MRS) studies of brain biomarkers can provide statistically significant biomarkers for tumor grade differentiation and im-

proved predictors of cancer patient survival [1]. Instead of selecting biomarkers based on microscopic histology and tumor morphology, the introduction of microarray technology improves the discovery rates of different types of cancers through monitoring thousands of gene expressions in a parallel, in a rapid and efficient manner [23][8]. Because the genes are aberrantly expressed in tumor cells, researchers can use their aberrant expression as biomarkers that correspond to and facilitate precise diagnoses and/or therapy outcomes of malignant transformation.

Different data sources are likely to contain different and partly independent information about the brain tumor. Combining those complementary pieces of information can be expected to enhance the brain tumor diagnosis and biomarkers detection. Recently, several studies have attempted to correlate imaging findings with molecular markers, but no consistent associations have emerged and many of the imaging features that characterize tumors currently lack biological or molecular correlates [7][6]. Much of the information encoded within neuroimaging studies therefore remains unaccounted for and incompletely characterized at the molecular level [4]. This paper presents a computational and machine learning based framework for integrating heterogeneous genome-scale gene expression and MRS data to classify the different brain tumor types and detect biomarkers. We employ wrapper method to integrate the feature selection process of both gene expression and MRS. Three popular feature selection methods, Relief-F (RF), Information Gain (IG) and χ^2 -statistic (χ^2), are performed to filter out the redundant features in both datasets. The experimental results show our framework using the combination of two datasets outperforms any individual dataset on sample classification accuracy that is the standard validation criterion in cancer classification and biomarker detection. Our data fusion framework exhibits great potential on heterogeneous data fusion between biomedical image and biological datasets and it could be extended to another cancer diseases study.

2 Methodology

Advancements in the diagnosis and prognosis of brain tumor patients, and thus in their survival and quality of life, can be achieved using biomarkers that facilitate improved tumor typing. In our research, we apply state-of-the-art, high-resolution magic angle spinning (HRMAS) proton (1H) MRS and gene transcriptome profiling to intact brain tumor biopsies, to evaluate the discrimination accuracy for tumor typing of each of the above methods separately and in combination. We used 46 samples of normal (control) and brain tumor biopsies from which we obtained ex vivo HRMAS 1H MRS and gene expression data respectively. The samples came from tissue biopsies taken from 16 different people. Out of the forty-six biopsies that were analyzed, 9 of them were control biopsies from epileptic surgeries and the rest 37 were brain tumor biopsies. The tumor biopsies belonged to 5 different categories: 11 glioblastoma multiforme (GBM); 8 anaplastic astrocytoma (AA); 7 meningioma; 7 schwannoma; and 5 from adenocarcinoma.

HRMAS 1H MRS. Magnetic resonance spectroscopic (MRS) studies of brain biomarkers can provide statistically significant biomarkers for tumor grade differentiation and improved predictors of cancer patient survival [1]. Ex vivo high-resolution magic angle spinning (HRMAS) proton (^1H) MRS of unprocessed tissue samples can help interpret in vivo ^1H MRS results, to improve the analysis of micro-heterogeneity in high-grade tumors [3]. Furthermore, two-dimensional HRMAS ^1H MRS enables more detailed and unequivocal assignments of biologically important metabolites in intact tissue samples [16]. In Fig.1, an ex vivo HRMAS ^1H MR spectrum of a 1.9 mg anaplastic ganglioglioma tissue biopsy is shown together with metabolites values that correspond to each frequency of the spectrum. Please see more detailed information in [29].

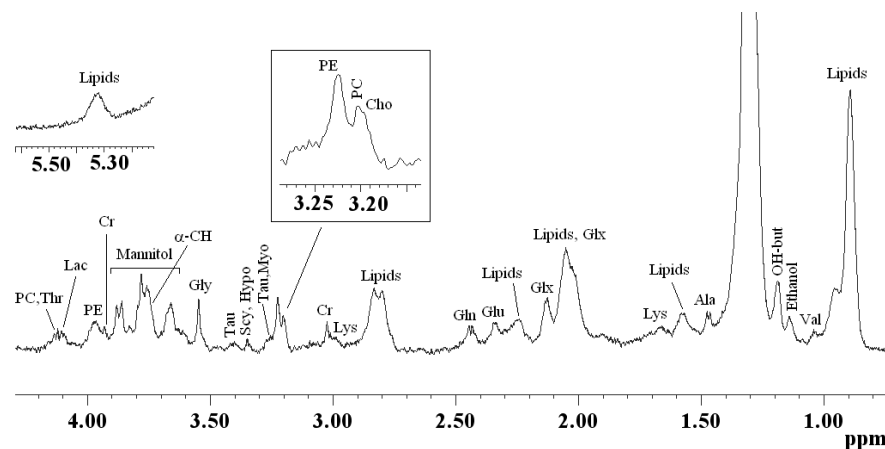


Fig. 1: Ex vivo HRMAS ^1H MR spectrum of a 5.8 mg glioblastoma multiforme (GBM) tissue biopsy. Val, Valine; OH-but, OH-butyrate; Lac, Lactate; Ala, Alanine; Lys, Lysine; Glx, \square -CH₂ of Glutamine and Glutamate; Glu, Glutamate; Gln, Glutamine; Cr, Creatine; Tau, Taurine; Myo, Myo-inositol; Hypo, Hypotaurine; Scy, Scyllo-inositol; Gly, Glycine; \square -CH of aliphatic amino-acids; PE, PhosphoEtanolamine; Thr, Threonine; PC, PhoshoCholine; Cho, Choline. The insert shows the choline containing compounds region.

Microscale genomics. A major focus in cancer research is to identify genes, using DNA-microarrays that are aberrantly expressed in tumor cells, and to use their aberrant expression as biomarkers that correspond to and facilitate precise diagnoses and/or therapy outcomes of malignant transformation [17]. In our study, the Affymetrix gene-chip U133Plus[®] DNA microarray of the complete human genome was used to perform transcriptome profiling on each specimen for two different experimental conditions, minus or plus previous HRMAS NMR analysis. The raw expression data were analyzed for probe intensities using the Affymetrix GeneChip expression analysis manual procedures; and the data were normalized using current R implementations of RMA algorithms [10].

Combining MRS and genomic data. While several studies have utilized MRS data or genomic data to promote cancer classification, to date these two methods have

not been combined and cross-validated to analyze the same cancer samples. Herein, we implement a combined quantitative biochemical and molecular approach to identify diagnostic biomarker profiles for tumor fingerprinting that can facilitate the efficient monitoring of anticancer therapies and improve the survival and quality of life of cancer patients. The MRS and genomic data strongly correlate, to further demonstrate the biological relevance of MRS for tumor typing [21]. Also, the levels of specific metabolites, such as choline containing metabolites, are altered in tumor tissue, and these changes correspond to the differential expression of Kennedy cycle genes responsible for the biosynthesis of choline phospholipids (such as phosphatidylcholine) and suggested to be altered with malignant transformation [18]. These data demonstrate the validity of our combined approach to produce and utilize MRS/genomic biomarker profiles to type brain tumor tissue.

2.1 Classification and feature selection methods

Classification aims to build an efficient and effective model for predicting class labels of unknown data. In our case the aim is to build a model that will be able to discriminate between different tumor types given a set of gene expression values or MRS metabolite values or a combination of them. Classification techniques have been widely used in microarray analysis to predict sample phenotypes based on gene expression patterns. Li et al. have performed a comparative study of multiclass classification methods for tissue classification based on gene expression [12]. They have conducted comprehensive experiments using various classification methods including SVM [22] with different multiclass decomposition techniques, Naïve Bayes [14], K-nearest neighbor and decision trees [20].

Since the main purpose of this study is not to assess the classification performance of different classification algorithms but to evaluate the potential gain of combining more than one type of data for tumor typing, we only experimented with Naïve Bayes (NB) and Support Vector Machines (SVM) with RBF kernel.

Another related task is *feature selection* that selects a small subset of discriminative features. Feature selection has several advantages, especially for the gene expression data. First, it reduces the risk of over fitting by removing noisy features thereby improving the predictive accuracy. Second, the important features found can potentially reveal that specific chromosomal regions are consistently aberrant for particular cancers. There is biological support that a few key genetic alterations correspond to the malignant transformation of a cell [19]. Determination of these regions from gene expression datasets can allow for high-resolution global gene expression analysis to genes in these regions and thereby can help in focusing investigative efforts for understanding cancer on them.

Existing feature selection methods broadly fall into two categories, wrapper and filter methods. Wrapper methods use the predictive accuracy of predetermined classification algorithms, such as SVM, as the criteria to determine the goodness of a subset of features [9]. Filter methods select features based on discriminant cri-

teria that rely on the characteristics of data, independent of any classification algorithm [5]. Filter methods are limited in scoring the predictive power of combined features, and thus have shown to be less powerful in predictive accuracy as compared to wrapper methods [2]. In our experiments we used feature selection method from both major categories. We experimented with Relief-F (RF), Information Gain (IG), and χ^2 -statistic (χ^2), filter methods and we also used wrapper feature selection for each of the two types of classification algorithms.

The basic idea of *Relief-F* [11] is to draw instances at random, compute their nearest neighbors, and adjust a feature weighting vector to give more weight to features that discriminate the instance from neighbors of different classes. Specifically, it tries to find a good estimate of the following probability to assign as the weight for each feature f .

$$w_f = P(\text{different value of } f | \text{different class}) - P(\text{different value of } f | \text{same class})$$

Information Gain (IG) [15] measures the number of bits of information obtained for class prediction by knowing the value of a feature. Let $\{c_i\}_{i=1}^m$ denote the set of classes. Let V be the set of possible values for feature f . The information gain of a feature f is defined to be:

$$G(f) = -\sum_{i=1}^m P(c_i) \log P(c_i) + \sum_{v \in V} \sum_{i=1}^m P(f=v) P(c_i | f=v) \log P(c_i | f=v)$$

The χ^2 -statistic (χ^2) [13] measures the lack of independence between f and c . It is defined as follows:

$$\chi^2(f) = \sum_{v \in V} \sum_{i=1}^m \frac{(A_i(f=v) - E_i(f=v))^2}{E_i(f=v)}$$

where V is the set of possible values for feature f , $A_i(f=v)$ is the number of instances in class c_i with $f=v$, $E_i(f=v)$ is the expected value of $A_i(f=v)$. $E_i(f=v)$ is computed with $E_i(f=v) = P(f=v)P(c_i)N$, where N is the total number of instances.

3 Experimental Results

Initially we aimed at evaluating how well the classifiers would perform when applying them to each of our datasets separately. For that purpose we performed 10-fold cross validation over our 46 samples by using a combination of feature selection and classification methods.

Table 1 shows the classification accuracy of Naïve Bayes (NB) and SVM classifiers when using all 16 *metabolites* and when using a feature selection method. Clearly the wrapper feature selection method gives the better accuracy across all classifiers, followed by the case where we use all metabolites for classification. The SVM classifier using RBF kernel consistently shows the best performance in this type of data. The decision of keeping the top 6 metabolites when using the filter feature selection methods was based on the fact that that was the best number

of features that were selected by using the wrapper feature selection method for each classification algorithm.

Table 1: Classification accuracy for the 6-class problem using MRS data only.

	NB	SVM
All metabolites	70.21 %	72.34 %
χ^2 (top 6)	46.81 %	51.06 %
IG (top 6)	46.81 %	51.06 %
RF (top 6)	63.83 %	68.09 %
Wrapper	72.34 %	78.72 %

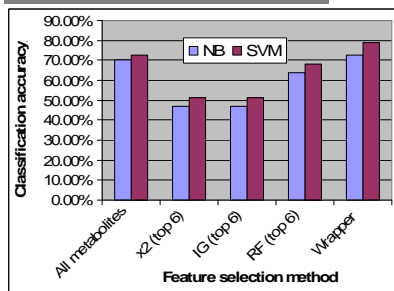
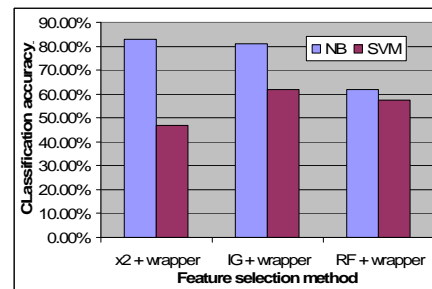


Table 2: Classification accuracy using gene expression data only.

	NB	SVM
χ^2 + wrapper	82.98 %	46.81 %
IG + wrapper	80.85 %	61.70 %
RF + wrapper	61.70 %	57.44 %



For the problem of the multiclass classification using gene expression data only, we followed a hybrid feature selection method combining filter and wrapper approaches. Using wrapper approach to select a few top genes starting from an initial number of thousands of genes is computationally prohibiting, and using filter approach to select less than 100 genes does not give good classification accuracy because the final set of selected genes contains genes that are highly correlated to each other, thus giving a redundant set of genes. In our approach, first we selected the top 100 genes using filter feature selection and then we used wrapper feature selection to further reduce the number of genes to be used resulting usually in a number between 5 and 15 genes.

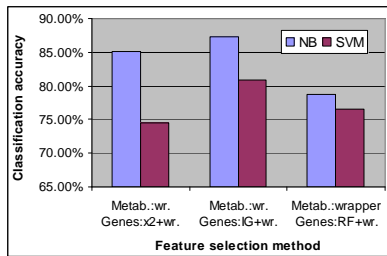
The experimental results (Table 2) show that in this type of data the Naïve Bayes was by far the best classification algorithm obtaining a maximum accuracy of 82.98% accuracy when combined with χ^2 and wrapper feature selection.

Finally, we tested the classification accuracy of our methods by using a combination of features from both gene expression and MRS data. For the MRS data we tested the wrapper feature selection method which performed best in our previous experiments. For the gene expression data we used the feature selection method that we described above, i.e. combination of filter and wrapper feature selection. After completing the feature selection stage separately for each of the datasets we combined the selected features by putting them in the same feature vector space and using that space for classification. Table 3 shows the classification accuracy results of our experiments. In most cases, the combination of features sets from the two datasets yield significantly better accuracy than each of them separately.

In general Naïve Bayes gives the best performance with a maximum accuracy of 87.23% when using wrapper feature selection for metabolites and a combination of Information Gain and wrapper feature selection for genes.

Table 3: Classification accuracy using a combination of features from gene NB SVM expression and MRS datasets.

	NB	SVM
Metabolite selection: wrapper	85.11%	74.46%
Gene selection: χ^2 + wrapper		
Metabolite selection: wrapper	87.23%	80.85%
Gene selection: IG + wrapper		
Metabolite selection: wrapper	78.72%	76.59%
Gene selection: RF + wrapper		



4 Conclusion

In this paper, we propose a machine learning based data fusion framework which integrates heterogeneous data sources to type different brain tumors. Our method employs real biomedical/biological MRS and genomic data and applies a combination of popular feature selection and classification methods to evaluate the tumor type discrimination capabilities of the two datasets separately and together. The feature selection process identifies a number of biomarkers from each dataset which are subsequently used as features for the classification process. The experimental results show that our data fusion framework outperforms each individual dataset in the brain tumor multi-class classification problem. Since our framework is a general method, it can also be applied to any other biomedical and biological data fusion for sample classification and biomarker detection.

References

1. LG Astrakas, D Zurakowski, and AA Tzika et al. Noninvasive magnetic resonance spectroscopic imaging biomarkers to predict the clinical grade of pediatric brain tumors. Clin Cancer Res, 10:8220–8228, 2004.

2. H. Chai and C. Domeniconi. An Evaluation of Gene Selection Methods for Multi-class Microarray Data Classification. ECML/PKDD 2004.
3. L.L. Cheng, D.C. Anthony, A.R. Comite, P.M. Black, A.A. Tzika, and R.G. Gonzalez. Quantification of microheterogeneity in glioblastoma multiforme with ex vivo high-resolution magic-angle spinning (HRMAS) proton magnetic resonance spectroscopy. *Neuro-Oncology*, 2(2):87–95, 2000.
4. M Diehn, C Nardini, and et al. Identification of noninvasive imaging surrogates for brain tumor gene-expression modules. *Proc Natl Acad Sci*, 105(13):5213–5218, 2008.
5. C. Ding and H. Peng. Minimum redundancy feature selection from microarray gene expression data. *Journal of Bioinformatics and Computational Biology*, 3(2):185–205, 2005.
6. Carlson MR et al. Relationship between survival and edema in malignant gliomas: role of vascular endothelial growth factor and neuronal pentraxin. *Clin Cancer Res.*, 13(9):2592–2598, 2007.
7. Hobbs SK et al. Magnetic resonance image-guided proteomics of human glioblastoma multiforme. *Magn. Reson. Imaging*, 18(5):530–536, 2003.
8. J.N. Rich et al. Gene expression profiling and genetic markers in glioblastoma survival. *Cancer Research*, 65:4051–4058, 2005.
9. Guyon and A. Elisseeff. An introduction to variable and feature selection. *The Journal of Machine Learning Research*, 3:1157–1182, 2003.
10. R.A. Irizarry, B. Hobbs, F. Collin, Y.D. Beazer-Barclay, K.J. Antonellis, Uwe Scherf, and T.P. Speed. Exploration, normalization, and summaries of high density oligonucleotide array probe level data. *Biostatistics*, 4(2):249, 2003.
11. Kononenko. Estimating Attributes: Analysis and Extensions of Relief. *Lecture Notes in Computer Science*, pages 171–171, 1994.
12. T. Li, C. Zhang, and M. Ogihara. A comparative study of feature selection and multiclass classification methods for tissue classification based on gene expression, 2004.
13. H. Liu and R. Setiono. Chi2: Feature selection and discretization of numeric attributes. In *Proceedings 7th International Conference on Tools with Artificial Intelligence*, page 88. IEEE Computer Society Washington, DC, 1995.
14. V. Metsis, I. Androutsopoulos, and G. Paliouras. Spam filtering with naive bayes – which naive bayes. In *Third Conference on Email and Anti-Spam (CEAS)*, 2006.
15. T.M. Mitchell. *Machine Learning*. 1997. Burr Ridge, IL: McGraw Hill.
16. D. Morvan, A. Demidem, J. Papon, M. De Latour, and J.C. Madelmont. Melanoma Tumors Acquire a New Phospholipid Metabolism Phenotype under Cystemustine As Revealed by High-Resolution Magic Angle Spinning Proton Nuclear Magnetic Resonance Spectroscopy of Intact Tumor Samples 1, 2002.
17. C.L. Nutt, DR Mani, R.A. Betensky, P. Tamayo, J.G. Cairncross, C. Ladd, U. Pohl, C. Hartmann, M.E. McLaughlin, T.T. Batchelor, et al. Gene Expression-based Classification of Malignant Gliomas Correlates Better with Survival than Histological Classification 1, 2003.
18. F. Podo. Tumour phospholipid metabolism. *NMR in Biomedicine*, 12(7):413–439, 1999.
19. MJ Renan. How many mutations are required for tumorigenesis? Implications from human cancer data. *Mol Carcinog*, 7(3):139–46, 1993.
20. PN Tan, M. Steinbach, and V. Kumar. *Introduction to Data Mining*. 2005.
21. A.A. Tzika, L. Astrakas, H. Cao, D. Mintzopoulos, O.C. Andronesi, M. Mindrinos, J. Zhang, L.G. Rahme, K.D. Blekas, A.C. Likas, et al. Combination of high-resolution magic angle spinning proton magnetic resonance spectroscopy and microscale genomics to type brain tumor biopsies. *International Journal of Molecular Medicine*, 20(2):199, 2007.
22. V. Vapnik. *Statistical Learning Theory*. 1998. NY Wiley.
23. S.J. Watson, F. Meng, R.C. Thompson, and H. Akil. The chip as a specific genetic tool. *Biol Psychiatry*, 48:1147–1156, 2000.