

# Automated Prostate Cancer Diagnosis via Pattern Recognition Approach

Anthony Karali<sup>1</sup>, Miguel García-Torres<sup>2</sup>, Federico Divina<sup>2</sup>, Alcides Chaux<sup>3</sup>,  
Anahí Chaux<sup>3</sup>, and George J. Netto<sup>4</sup>

<sup>1</sup> Faculté dinformatique, Université de Namur, Namur, Belgique  
Anthony.karali@gmail.com

<sup>2</sup> Computer Science, Universidad Pablo de Olavide, ES-41013, Seville, Spain  
{mgarciaf,fdivina}@upo.es

<sup>3</sup> Universidad del Norte, Paraguay {Acahux,Ana.i.marvil}@gmail.com

<sup>4</sup> Departments of Pathology, Urology and Oncology, The Johns Hopkins Medical  
Institutions (Baltimore, MD)

**Abstract.** Traditionally, pathologists make diagnostic assessment based on cell morphology and tissue distribution. However, this diagnosis depends on the experience of the pathologists and, so, it leads to high variability. Image analysis approach enables to perform an objective judgment by characterizing the images extracting quantitative measures. In this work, we develop a pipeline that, using image analysis tools and machine learning techniques, can produce a diagnosis. Moreover, we apply feature selection strategies to analyze the contribution of each extracted feature to the predictive model.

**Keywords:** Automated cancer diagnosis, biomedical image analysis, segmentation, feature extraction, classification, feature selection

## 1 Introduction

According to the World Health Organization, cancer figures among the leading causes of death worldwide with over 14 million of cases and over 8 million cancer related deaths in 2012. Prostate cancer is the most common noncutaneous cancer among males and the sixth leading cause of death for men worldwide [8]. Although the causes of prostate cancer are not yet fully understood, it is known that the chances of developing it increases with age.

The treatment of this kind of cancer depends on its malignancy level. The chances of survival are generally high if diagnosed at an early stage, and decrease at more advanced stages. Nowadays the diagnosis depends on the pathologists personal experience and, therefore, this subjective judgment often leads to considerable variability [7]. Therefore, to improve the reliability of the diagnosis, it is necessary to develop a mathematical model to characterize the prostate cancer features.

In order to characterize numerically the prostate cancer from Tissue Microarray (TMA) images, it is necessary to extract features from each single image.

Most cancer image analysis systems have been developed from cytological specimens, which only capture cells and, so, do not use any information at the tissue level [14,13]. However, the analysis at tissue level provides information about the structure of different pathological elements that are more important for the diagnosis than the appearance of individual level. Despite the important role at tissue level, research on image analysis targeting cancer tissue is not widely available due to the difficulty and complexity involved in performing quantitative tissue image analysis.

Previous works in automatic cancer diagnosis have focused on detecting among several Gleason grades [3,11], while others on detecting cancer on such images [9,2]. All these works base their prediction on features that are not normally used by pathologists, e.g., the entropy and energy of the multiwavelet coefficients.

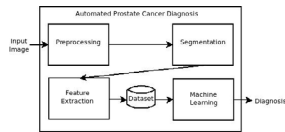


Fig. 1: Computational steps in automatic cancer diagnosis.

Figure 1 shows the computational steps in automatic cancer diagnosis: pre-processing, segmentation, feature extraction and machine learning. However, in our approach, we consider the segmentation step as part of the preprocessing step.

The goal of this work is to develop a prostate cancer diagnosis system for the automatic detection of the presence of cancer given an image of a tissue. This is done by examining histological properties of the tissue and cell nucleoids. We also analyze three different preprocessing approaches and study the predictive power of the features extracted. The software developed is available upon request.

## 2 Material and methods

In this section we describe the data and the methods used in this paper.

The study proposed in this paper was carried out on tissue samples obtained from 50 patients with localized prostate cancer who underwent radical retropubic prostatectomy at The Johns Hopkins Medical Institutions (Baltimore, MD) between 1993 and 2001. These 50 patients were randomly selected from a previously published cohort that included 524 matched cases and controls [1]. A total of 400 tissue cores were obtained, of which 196 cores of tumor and 102 cores of paired normal tissue from the patients with prostate cancer, plus 102 cores of nonprostate tissue.

Images have to be preprocessed in order to remove noise and useless information from them. In this work, we analyze three different preprocessing methods that explore different approaches to characterize the cells. Below we explain them in more detail.

**Morphological image filtering** This preprocessing pipeline is identified, from now on, as  $P_1$ . As we can see in the Figure 2,  $P_1$  involves six steps in order to isolate cells from the tissue.

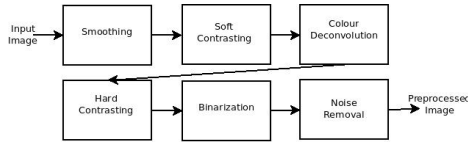


Fig. 2: Steps of the second preprocessing pipeline used in this study.

The first step, smoothing, consists of reducing the amount of intensity variation between one pixel and the next. We use a  $3 \times 3$  mean filtering, which replaces each pixel value with the mean value of its neighbours, including itself. In the next step, contrast enhancement (soft contrasting in the pipeline), is performed by histogram stretching keeping 0.5% saturated pixels. The values of the saturated pixels determines the number of pixels in the image that are allowed to become saturated. A colour deconvolution phase follows, in order to separate the image into three channels, corresponding to the actual colours of the stains used. Since hematoxylin mainly stains the cell nuclei, we can remove, from the image, most of the tissue and artifacts. However, it is still hard to distinguish the cells via image through analysis due to noise. Therefore, we perform a new contrast enhancement method that uses anisotropic diffusion [12] to highlight the cells and be able to remove the noise surrounding the cell nuclei. The remaining artifacts will be removed in the next step by applying a binary filter, which transforms the current image into a binary one. To set the threshold, the method uses an iterative procedure based on the isodata algorithm. This step results in a series of cells nuclei and artifacts that, in some cases, are empty. We use a four-way Flood fill algorithm to fill such regions followed of the watershed segmentation method. Finally, we remove such isolated regions that are smaller than a given threshold. In this phase all object with a number of pixels lower than 10 are removed.



Fig. 3: Steps of the second preprocessing pipeline used in this study.

**H&E clustering based pipeline** This preprocessing pipeline, identified as  $P_2$ , consists of three steps (see Figure 3). The first one consists of removing the background pixels of the image by using the clustering technique  $K$ -means ( $K = 2$ ). As seeds, we used the vectors associated to black ( $rgb = (0, 0, 0)$ ) and white ( $rgb = (255, 255, 255)$ ) colours. Those pixels associated with the later cluster will be considered background pixels and, therefore, set to white pixels. The next step is to perform clustering taking into account the H&E-stain and the white colour to avoid the contamination of background pixels on any of the other clusters. Therefore, we set  $K = 3$  with the following seeds: 1) white ( $rgb = (255, 255, 255)$ ); 2) haematoxylin-stained pixels ( $rgb = (0.490157, 0.768971, 0.410402)$ ); and 3) eosin-stained pixels ( $rgb = (0.046153, 0.842068, 0.537393)$ ). This step will generate, at the cell-level, many isolated objects. Then, after a visual analysis, we remove those objects with a number of pixels smaller than 10.

**H&E based rule** This pipeline, identified as  $P_3$ , only requires the use of a single filter. We use the property of the H&E stain, which allows to distinguish between cells and tissue or artifacts. Hematoxylin is a dark blue or violet stain that reacts with cells while eosin is red or pink and stains the non-cell objects of the image.

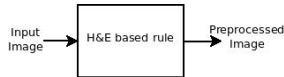


Fig. 4: Steps of the third preprocessing pipeline used in this study.

Therefore, in order to filter out non background pixels, we impose the following constrains on pixels in:

- The difference between the green and blue values must be greater than a given threshold  $\mu_{gb}$ .
- The difference between the red and blue values must be greater than a given threshold  $\mu_{rb}$ .
- The difference between the highest total intensity (sum of the red, green and blue colours of a white pixel) and the total intensity must be lower than a given threshold  $\lambda_t$ .

The last constrain defines how close the total intensity of a cell pixel should be to the white one. We tested several values and finally we set  $\mu_{gb} = \mu_{rb} = -5$  and  $\lambda_i = 420$

## 2.1 Feature extraction

In this work, we aim at extracting features to quantify the properties of cell structures and changes in the cell distribution across the tissue. Therefore, prostate cancer images are characterized by extracting features at cell-level. In order to do so, we measure morphological features (which provides information about the size of a nucleus or a cell), color-intensity (measured on the intensity histogram of the pixels located in a nucleus or a cell) and topological features (which inform on the cellular structure of a tissue by quatifying the spatial distribution of its cells).

In this work we consider the following features.

**Morphological features** Given  $S = \{s_1, \dots, s_n\}$  a set of boundary pixels of a cell, and  $C$  its centroid, we extract the following morphological features defined on  $S$ , where  $\mathcal{M}^{(\mu)}$  and  $\mathcal{M}^{(\sigma)}$  denotes the mean and standrard deviation of a measure  $\mathcal{M}$ :

- *Area*  $\mathcal{M}_a$  is the number of pixels within the boundary. We consider  $\mathcal{M}_a^{(\mu)}$  and  $\mathcal{M}_a^{(\sigma)}$ .
- *Perimeter*  $\mathcal{M}_p$  is the number of pixels in the boundary, and is measured as the sum of the distances between every consecutive boundary pixels:  $\mathcal{M}_p = |s_n s_1| + \sum_{i=1}^{n-1} |s_i s_{i+1}|$
- *Compactness*  $\mathcal{M}_{Co} = \frac{\mathcal{M}_p^2}{\mathcal{M}_a}$ , associates the perimeter and the area of  $S$
- *Centroid Gravity*  $\mathcal{M}_G$  determines the average coordinates of  $S$ . We extract a feature for each coordinate, and is computed with the following formula:

$$\mathcal{M}_G = \begin{cases} \mathcal{M}_{G_x} = \sum_{i=1}^n x_i, \\ \mathcal{M}_{G_y} = \sum_{i=1}^n y_i. \end{cases}$$

**Colour-intensity features** Colour-intensity features are potentially useful in prostate cancer images since our sample is composed of H&E stained prostate cancer images. As colour space, we use the RGB model, and, for each colour channel, we extract the first two statistical moments and the best range of pixels from histogram. For a given colour channel  $C$  and the set of pixel  $p_C = \{p_1, \dots, p_n\}_C$ ,  $C = \{r, g, b\}$ , the features extracted are shown below:

- *Mean* ( $I_C^{(\mu)} = \frac{1}{n} \sum_{j=1}^n p_{Cj}$ ), obtained by averaging the pixel values for each colour channel.
- *Standard deviation* ( $I_C^{(\sigma)} = \left( \frac{1}{n} \sum_{j=1}^n (p_{Cj} - I_C^{(\mu)})^2 \right)^{1/2}$ ), calculated over the pixel values on the given colour channel.
- *Interval*  $[I_C^l, I_C^h]$  is obtained from the histogram and it has been defined so that it is centered at  $I_C^{(\mu)}$  and it includes 90% of the pixel values.

**Topological features** These features measure the structure of a tissue by quantifying the spatial distribution of its cells. For that, it is necessary to encode the spatial interdependency of the cells prior to the feature extraction. In this work we encode this dependency between adjacent cells by using the Voronoi diagram, the Delaunay triangulation and the graph approach.

- *Voronoi Diagrams* represent a partitioning of the image into a set of non-overlapping regions that constitutes convex polygons. Each polygon contains a cell and every point in its region is closer to this cells than to another one in the tissue. When creating the Voronoi diagram on tissue image, border polygons tend to be larger and so we discarded them to avoid any bias due to border effect. We remove the 4% larger polygons. The list of features extracted is, then:

- *Area*,  $\mathcal{V}_a$  of the polygons. We extract  $\mathcal{V}_a^{(\mu)}$ ,  $\mathcal{V}_a^{(\sigma)}$  and the median  $\mathcal{V}_a^{(m)}$ .
- *Area disorder* [10],  $\mathcal{V}_d = 1 - \frac{1}{1 + \frac{\mathcal{V}_a^{(\sigma)}}{\mathcal{V}_a^{(\mu)}}}$  reflects the variation in the area of

the Voronoi polygons.

- *Average Roundness Factor*  $\mathcal{V}_r = \frac{4\pi}{p^2}$  [10] calculates the average of the roundness factor of the polygons. with p the perimeter of a polygon.
- *Roundness Factor disorder*  $\mathcal{V}_{rd} = 1 - \frac{1}{1 + \frac{\mathcal{V}_r^{(\sigma)}}{\mathcal{V}_r^{(\mu)}}}$  [10] measures the variation

of  $\mathcal{V}_r$  for all polygons. A value of 1 means that all  $\mathcal{V}_r$  are equal while 0 otherwise.

- *Density*  $\mathcal{V}_\rho = \frac{\#polygons}{\sum \mathcal{V}_a}$  [10] measures feature the density of the voronoi polygons.

- *Delaunay triangulation* is the dual graph of the Voronoi diagram. It constructs a triangular graph that cover the area of the image. In tissue images, such network consist of non-overlapping triangles so that each vertices correspond to nuclear centroids. The extracted features are:

- *Length of segments*  $\mathcal{D}_\ell$  generated in the Delaunay graph.
- *Delaunay segment length disorder*  $\mathcal{D}_d = \frac{1}{1 + \frac{\mathcal{D}_\ell^{(\sigma)}}{\mathcal{D}_\ell^{(\mu)}}}$  measures the variation

of  $\mathcal{D}_\ell$  for all segments.

- *Graph-based features* In this approach, the dependency between every pair of cells is encoded with a graph. Vertices correspond to nuclear centroids while edges are probabilistically assigned between the vertices; the probability of the existance of an edge between a pair of vertices decays with the increasing Euclidean distance between them [5]. To construct the graph we use the Waxman model [16]. The list of features extracted are:

- *Degree of nodes*  $\mathcal{G}_d$  of the graph.
- *Weighted degree of nodes*  $\mathcal{G}_w$  of the graph.
- *Clustering coefficient*  $\mathcal{G}_c$  [15] measures, for a given vertex  $v$ , the fraction between the number of connections  $k_v$  of such vertex, and the maximum number of connections it can have  $\left(\frac{k_v(k_v-1)}{2}\right)$ . The feature is define as the average over all vertices of the graph.

## 2.2 Machine Learning methods

It is well known that, in general, not all features contribute equally to the classification. Therefore, removing irrelevant feature may yield better predictive models. In this context the aim of the feature selection is to find the optimal feature subset, from the original feature set. The goodness of a particular feature subset is evaluated using an objective function,  $J(S)$ , where  $S$  is a feature subset of size  $|S|$ .

In this work we use three popular and widely used classifiers due to their good performance in general: Bayesian Network Classifier (BNC), the open source Java implementation of the C4.5 algorithm termed J48 and the Support Vector Machine (SVM). As feature selection algorithms, we selected the Fast Correlation Based Filter [17] (FCBF) and the Scatter Search (SS) metaheuristic [4]. FCBF is a popular and efficient while SS is an evolutionary strategy that achieves very competitive results. In this work, SS measures the quality of the subsets by means of the Correlation Feature Selection (CFS) function [6].

## 3 Results

In this section we present the results of experiments performed in order to assess the effectiveness of the proposed pipelines and the classification performance of the features extracted from the images. Moreover, we performed a feature selection analysis to select the most informative features. K-fold cross validation was used, where  $k$  was set to 5.

As performance measures, we use the classification error, the sensitivity and specificity averaged over the folds. In our data positive examples refer to prostatic carcinoma images while negatives to control cases. Therefore, sensitivity, also called *true positive rate* or recall, measures the proportion of actual positives which are correctly identified as such. Higher values mean that more cases of carcinoma are detected. Specificity is the proportion of actual negatives which are identified as such. Higher values correspond to lower probability of false positives, i.e., that a control case be classified as carcinoma case. Finally, we also report the average number of features selected by each strategy.

### 3.1 Baseline classification results

The performances of the baseline classifiers are shown in Table 1. The first column shows the classifier used. Then, for each pipeline developed, the table presents the sensitivity, the specificity and the accuracy, respectively.

The highest performance is achieved using  $P_2$  with BNC. With the same classifier,  $P_1$  obtains results slightly lower. However, the capacity for detecting carcinoma cases (sensitivity) is higher for  $P_1$  than for  $P_2$ , although it also presents a higher false positive rate. Therefore,  $P_2$  presents a more balanced performance. With J48,  $P_3$  obtains the highest accuracy but in all cases, the specificity is lower than 0.7. Finally, linear SVM is the classifier that obtains the

Table 1: Baseline classification results achieved with the different pipelines for BNC, J48 and linear SVM.

pipeline	P <sub>1</sub>			P <sub>2</sub>			P <sub>3</sub>		
	clf	sens.	spec.	acc.	sens.	spec.	acc.	sens.	spec.
BNC	0.827	0.637	76.18	0.791	0.716	76.51	0.765	0.657	72.82
J48	0.750	0.598	69.80	0.791	0.618	73.15	0.791	0.667	74.83
SVM	0.832	0.343	66.44	0.454	0.627	51.34	0.668	0.696	67.80

lowest results. With P<sub>1</sub>, the sensitivity is high but the specificity is too low to be considered a good predictive model. Only using the third pipe line does SVM achieves acceptable results. Pipeline P<sub>3</sub> seems to be the most robust, since all the classifiers obtains relatively good results with it.

### 3.2 Feature selection analysis

Results obtained when applying feature selection are shown in Table 2. The first column refers to the feature selection algorithm. Then, the classifier used or the number of features is presented. The information about the number of features is given in the last row of each algorithm. Finally, for each pipeline, the table shows the sensitivity, specificity and accuracy respectively.

Table 2: Classification results achieved after applying feature selection with the different pipelines for BNC, J48 and linear SVM.

A	pipeline	P <sub>1</sub>			P <sub>2</sub>			P <sub>3</sub>		
		clf	sens.	spec.	acc.	sens.	spec.	acc.	sens.	spec.
SS	BNC	0.842	0.647	77.52	0.801	0.637	74.50	0.760	0.696	73.83
	J48	0.750	0.637	71.14	0.760	0.647	72.15	0.776	0.657	73.50
	SVM	0.270	0.735	42.95	0.847	0.471	71.82	0.531	0.578	54.70
	#feats	7			11			8		
FCBF	BNC	0.786	0.637	73.49	0.867	0.392	70.47	0.745	0.598	69.46
	J48	0.714	0.598	67.45	0.791	0.559	71.14	0.781	0.510	68.80
	SVM	0.827	0.431	69.13	0.724	0.510	65.10	0.526	0.480	51.01
	#feats	2			2			2		



The first thing we notice, is that FCBF only selects two features.  $\mathcal{M}_a^{(\sigma)}$  is selected in all the cases,  $\mathcal{V}_d$  in pipeline 1 and 3 and  $I_b^{(\sigma)}$  in pipeline 2. SS is the strategy that achieves the highest accuracy using  $P_1$  with BNC. In this case, the model outperforms the baseline classifier in all performance measures. In  $P_2$  the application of SS degrades the model with BNC and J48, while it improves the results achieved by SVM. Finally, in  $P_3$ , only results with BNC outperforms the baseline results. As a conclusion, we can state that FCBF finds, in all cases, smaller subsets of features than SS at the expense of degrading the classifier performance. In all cases, except in  $P_1$  with SVM, SS outperforms FCBF. The list of features selected in each case can be found in Table 3.

Table 3: List of features selected by each SS and FCBF.

A	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>
SS	$\mathcal{M}_a^{(\mu)}, \mathcal{M}_a^{(\sigma)}, \mathcal{V}_d^{(\mu)}$ $\mathcal{V}_a^{(\sigma)}, \mathcal{V}_d, \mathcal{D}_d, \mathcal{G}_w^{(\sigma)}$	$\mathcal{V}_a^{(\mu)}, \mathcal{V}_a^{(\sigma)}, \mathcal{V}_d, \mathcal{D}_d$ $\mathcal{D}_d^{(\mu)}, \mathcal{D}_d^{(\sigma)}, I_g^{(\sigma)}, I_b^{(\sigma)}$ $I_g^f, I_b^f$	$\mathcal{V}_a^{(\sigma)}, \mathcal{V}_a^{(m)}, \mathcal{V}_p$ $\mathcal{V}_d, \mathcal{D}_d, \mathcal{D}_d, I_b^f$
FCBF	$\mathcal{M}_a^{(\sigma)}, \mathcal{V}_d$	$\mathcal{M}_a^{(\sigma)}, I_b^{(\sigma)}$	$\mathcal{M}_a^{(\sigma)}, \mathcal{V}_d$

## 4 Conclusion

In this work we have proposed three different pipelines to characterize H&E stained prostate cancer images. Among these pipelines, the morphological-based pipeline is the one that achieves the highest results after applying feature selection. In this case, the image is characterized by the mean and standard deviation of the cell area, the mean and standard deviation of the voronoi diagrams, the area disorder and the standard deviation of the weighted degree of the nodes.

The results achieved in this paper are promising. However more development are necessary in order to improve the results obtained. With this in mind, we are planning to study other data representations and to apply other feature selection mechanisms. Moreover, we believe that the results achieved could be improved by considering a wider range of unsupervised and supervised strategies. This would most likely also contribute to increase the knowledge in this domain. We also intend to apply the strategy proposed in this paper to other images datasets relative to other type of cancer.

## References

1. Chaux, A., Peskoe, S.B., Gonzalez-Roibon, N., Schultz, L., Albadine, R., Hicks, J., Marzo, A.M.D., Platz, E.A., Netto, G.J.: Loss of PTEN expression is associated with increased risk of recurrence after prostatectomy for clinically localized prostate cancer. *Mod Pathol* 25(11), 1543–1549 (jun 2012)
2. DiFranco, M.D., Reynolds, H.M., Mitchell, C., Williams, S., Allan, P., Haworth, A.: Performance assessment of automated tissue characterization for prostate h and e stained histopathology. In: *SPIE Medical Imaging*. vol. 9420, pp. 94200M–94200M–9 (2015)
3. Doyle, S., Feldman, M.D., Shih, N., Tomaszewski, J., Madabhushi, A.: Cascaded discrimination of normal, abnormal, and confounder classes in histopathology: Gleason grading of prostate cancer. *BMC Bioinformatics* 13, 282 (2012)
4. García-López, F., García-Torres, M., Melián-Batista, B., Moreno-Pérez, J., Moreno-Vega, J.: Solving feature subset selection problem by a parallel scatter search. *European Journal of Operational Research* 169(2), 477–489 (2006), cited By 69
5. Gunduz, C., Yener, B., Gultekin, S.: The cell graphs of cancer. *Bioinformatics* 20, i145–i151 (2004)
6. Hall, M.A.: Correlation-based Feature Subset Selection for Machine Learning. Ph.D. thesis, University of Waikato (1998)
7. Ismail, S., Colclough, A.B., Dinnen, J.S., Eakins, D., Evans, D.M., Gradwell, E., OSullivan, J.P., Summerell, J., Newcombe, R.G.: Observer variation in histopathological diagnosis and grading of cervical intraepithelial neoplasia. *British Medical Journal* 298(6675), 707–710 (1989)
8. Krnjacki, L., Baade, P., Youlten, D.: International epidemiology of prostate cancer: Geographical distribution and secular trends. *Molecular Nutrition and Food Research* 53(2), 171–184 (2009)
9. Litjens, G., Bejnordi, B.E., Timofeeva, N., Swadi, G., Kovacs, I., van de Kaa, C.H., van der Laak, J.: Automated detection of prostate cancer in digitized whole-slide images of h and e-stained biopsy specimens. In: *SPIE Medical Imaging*. vol. 9420, pp. 94200B–94200B–6 (2015)
10. Marcelpoi, R., Sudbo, J., Reith, A.: New algorithms based on the voronoi diagram applied in a pilot study on normal mucosa and carcinomas. *Analytical Cellular Pathology* 21(2), 71–86 (2000)
11. Nguyen, K.: Contributions to Computer-Aided Diagnosis of Prostate Cancer in Histopathology. Ph.D. thesis, Michigan State University (2013)
12. Perona, P., Malik, J.: Scale-space and edge detection using anisotropic diffusion. *IEEE Transactions On Pattern Analysis and Machine Intelligence* 12(7), 629–639 (1990)
13. Swindle, P.W., Kattan, M.W., Scardino, P.T.: Markers and meaning of primary treatment failure. *Urologic Clinics of North America* 30(2), 377–401 (2003)
14. Wahlbj, C., Lindblad, J., Vondrus, M., Bengtsson, E., Björkstén, L.: Algorithms for cytoplasm segmentation of fluorescence labeled cells. *Analytical Cellular Pathology* 24, 101–111 (2002)
15. Watts, D., Strogatz, S.: Collective dynamics of 'small-world' networks. *Nature* 393, 440–442 (1998)
16. Waxman, B.: Routing of multipoint connections. *IEEE Journal on Selected Areas in Communications* 6(9), 1617–1622 (1988)
17. Yu, L., Liu, H.: Efficient feature selection via analysis of relevance and redundancy. *Journal of Machine Learning Research* 5, 1205–24 (2004)