

# Automatic recognition of complete palynomorphs in digital images

J. J. Charles

Received: date / Accepted: date

**Abstract** Images of dispersed kerogen preparation are analysed in order to detect palynomorphs of elliptical/spherical shape. This process consists of three automatic stages. Firstly, the background of the image is segmented from the foreground. Secondly the foreground particles are segmented into individual regions. Finally a trained classifier is used to label a region as either containing a palynomorph or containing other material. Ten classifiers were trained and then tested using a 10 times 10-fold cross validation. Typically the number of regions in the image containing other material exceeds by far the number of regions with palynomorphs. Hence the problem of imbalanced classes was addressed. Training data was sampled 10 different times maintaining a balanced class distribution. Thus the accuracy for each classifier was assessed on 1000 testing sets. The logistic classifier was chosen and a certainty threshold was selected by ROC curve analysis. The final automatic recognition has accuracy of 88%, sensitivity of 87% and specificity of 88%.

**Keywords** classification · microfossils · image analysis · segmentation · palynomorph

## 1 Introduction

Microfossils are used extensively by the petroleum industry when exploring for oil and hydrocarbon palaeontologists consider them to be one of their main tools. When drilling for oil a fluid lubricates the drill bit and

helps flush small pieces of rock from the bottom of the drill hole, these small pieces are known as cuttings. Because microfossils are small ( $\leq 1\text{mm}$ ) they are mostly undamaged by the drilling process. A sample of cuttings contains three groups of microfossil, in this study we are interested in the group known as palynomorphs.

Palynomorphs are important in hydrocarbon exploration to construct biostratigraphies, chronostratigraphies, palaeoenvironmental determinations and maturity assessments (1). These studies require a specialist to examine a slide containing a sample of cuttings. The task of automating quantitative palynofacies studies has been of interest for over 20 years. In 1988 an attempt to extract the outline of microfossils from reflected light images was made by (20). Later in 1989 a classification system to assist with identifying fossils was implemented by (23). An expert system for visually identifying microfossils was constructed by (24) in 1992 and one of the most recent studies in automating identification of palynofacies was conducted by (25) in 2005. Although much work has been done such tasks are still considered a challenging problem. It would be advantageous not only to an automatic system but also a specialist, to devise a method which assists in locating regions on the slide that contain more easily identifiable palynomorphs.

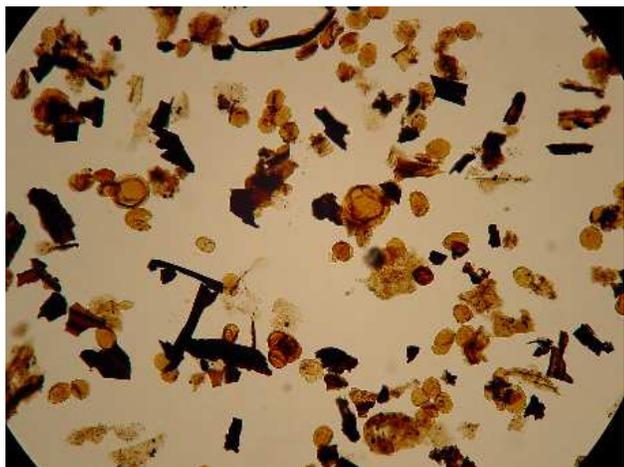
An image of sedimentary organic matter is captured at a resolution of 2272x1704 pixels through a digital camera mounted on a microscope, an example of a typical image is shown in Figure 1. Microfossils on the slide can be broadly classified into 3 groups: Kerogen<sup>1</sup>, palynomorphs and amorphous material. It has been shown that images of single particles can be recognised auto-

---

J. J. Charles  
School of Computer Science, Bangor University, Bangor, LL57  
1UT, UNITED KINGDOM  
Tel.: +44 1248 383661  
Fax: +44 1248 361429  
E-mail: j.j.charles@bangor.ac.uk

---

<sup>1</sup> An organic matter that can yield hydrocarbons upon heating



**Fig. 1** Typical microscopy image of sedimentary organic matter containing kerogen, palynomorphs and amorphous material

matically with an accuracy of 87% (25). However for a slide containing many particles it is necessary to first locate and extract them individually. Previous work has demonstrated that a stable segmentation and detection system can be applied to the kerogen material with an accuracy of 91% (21). Consequently our attention now focuses on the detection of palynomorphs.

Palynomorphs are difficult to extract automatically, firstly due to their haphazard arrangement on the slide, secondly their structure is deformable and thirdly palynomorphs can be semi-transparent. When there are multiple overlapping pieces it can be extremely difficult to distinguish between them, even by human eye. Therefore an unsupervised segmentation procedure is unlikely to be perfect and under/over-segmentation of some particles will occur. By first identifying the regions containing a palynomorph, a classification system need only be concerned with individual particles. These regions will contain a single, complete, elliptic palynomorph that has not been folded, torn, squashed etc. The main goal of this system is to extract complete palynomorphs from a slide containing overlapping microfossils. To the best of the author's knowledge, there is no other system for accomplishing this task. Complete palynomorphs are desirable because they can be further classified by existing systems. However, such systems find it very difficult to classify palynomorphs if they are visually distorted. Distortions occur when the microfossils are heavily overlapping and have folded over each other or occluded one another. An expert system for recognising microfossils under these types of conditions does not yet exist.

It is our goal to disregard regions that contain kerogen and amorphous material. The type of palynomorphs we are interested in are those with an elliptic/spherical morphology. These types of palynomorph can come from a variety of classes but particularly acritarchs, spores and pollen. Acritarchs was a term first introduced by (14). There are a large amount of these types of palynomorph preserved in the geological record. Hence they are extremely usefull for quantitative biostratigraphic and palaeobiological studies.

We have chosen to use the Centre Supported Segmentation (CSS) algorithhm (9) as our method of segmentation. This technique works well when segmenting kerogen pieces and can be applied to any binary image separated into foreground and background. It was shown by (8; 9) that this method is simple to implement, insensitive to small changes in its parameters and relatively quick to run compared to other automatic methods such as Randomized Hough Transform for ellipse detection as proposed by (22). Ellipse detection could be applied to an edge image of the slide however the performance of such an approach was found to be slow and inaccurate. This was due to a) most of the palynomorphs on a slide are not perfect ellipses and b) noise in the edge image. On the other hand the CSS algorithm is robust against noise and changes in object boundaries. The CSS algorithm can be controlled to only segment regions which overlap up to a certain limit. Heavily overlapping regions will not be segmented and this reduces the number of segmented regions containing visually distorted microfossils. Subsequent to segmentation, the types of region fall into three categories: complete palynomorphs, non-palynomorphs or clumps of heavily overlapping microfossils.

We propose to segment the foreground particles in the image and identify those that contain a single complete elliptic palynomorph. This is accomplished by training a classifier to distinguish between a region containing a single palynomorph and one containing othe material (kerogen and heavily overlapping microfossils). The novelty here is in automatically removing from the image complete palynomorphs which can be further classified by machine, leaving recognition of more complex regions to human expert. Furthermore such a system can be used to improve segmentation results by filtering out regions that need to be merged or further segmented.

The rest of the paper is split into four sections. The first section describes the methods and analysis used for image pre-processing and classification. The second section compares ten state of the art classifiers and the third section demonstrates the results of the logistic

classifier on our images. In the fourth section we summarise and discuss future work.

## 2 Methods and analysis

The microscopy image of a sieved rock sample contains a dispersed arrangement of kerogen and palynomorph material. It is unlikely that each palynomorph will appear disconnected and unobstructed. In most cases palynomorphs will be touched and occluded by other material, even folded, squashed or torn. A segmentation technique can be applied to split the image into regions with the hope that each region contains a separate palynomorph. A classifier can be applied to each region to determine whether it contains a single complete palynomorph or something in which we have no interest. This section will describe the techniques used to accomplish segmentation and classification. Methods used for classification analysis are also discussed.

### 2.1 Pre-processing

#### 2.1.1 Background segmentation

The first step is to correct the uneven light intensity across the image which ensures a consistency in colour for all objects on the slide. In microscopy images the main cause of uneven lighting will be an inherent problem i.e. the lens allows more light to hit the centre of the image than its edges, this is known as vignetting. Also in transmitted microscopy a bulb is placed behind the slide, this will amplify the light intensity at the centre of the image. The overall effect is a drop off in light intensity from the centre of the image towards the edges. The corrective procedure we chose (10) fits horizontal and vertical parabolas to the background intensities of the greyscale image shown in Figure 2 (a). By combining these parabolas a model of the background is formed. The image is corrected by this model to even out the background.

In the next step we segment all foreground objects into individual regions. First the background is segmented from the foreground by thresholding the corrected greyscale image. The image histogram typically has two pronounced peaks corresponding to background and foreground. The threshold value is found by locating the global minimum between the two peaks. The image is converted into black and white. The black regions are the foreground containing all palynomorphs and the white regions are the background (Figure 2 (b)).

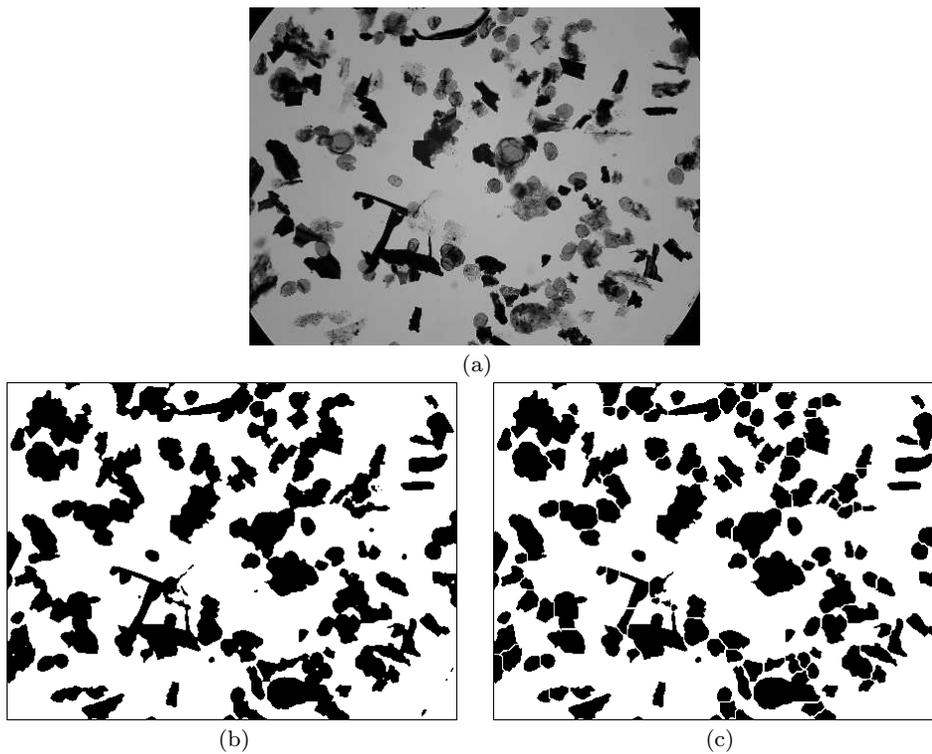
#### 2.1.2 Microfossil segmentation

A recently proposed algorithm known as Centre Supported Segmentation (CSS) can segment individual microfossils based upon this binary image. The foreground regions will either contain a single object or a collection of touching and overlapping pieces. CSS finds the centres of all objects which could have formed a region. For each centre an overlap value  $d \in [0, 1]$  is found. This expresses the amount of overlapping caused by an object at that centre. The degree of overlap increases as  $d$  increases, for example a separate object has  $d = 0$  and a completely occluded object has  $d = 1$ . Centres are disregarded if they have overlap greater than 0.5, this reduces the amount of over-segmentation. Classification of kerogen material with respect to the overlap value of the CSS algorithm was found to be stable even with dramatic changes in  $d$  (11). We recommend using a value of 0.5. Particles that are too small can be eliminated, (25) suggests removing particles with a diameter less than  $14\mu\text{m}$ . The result of this can be seen in Figure 2 (c).

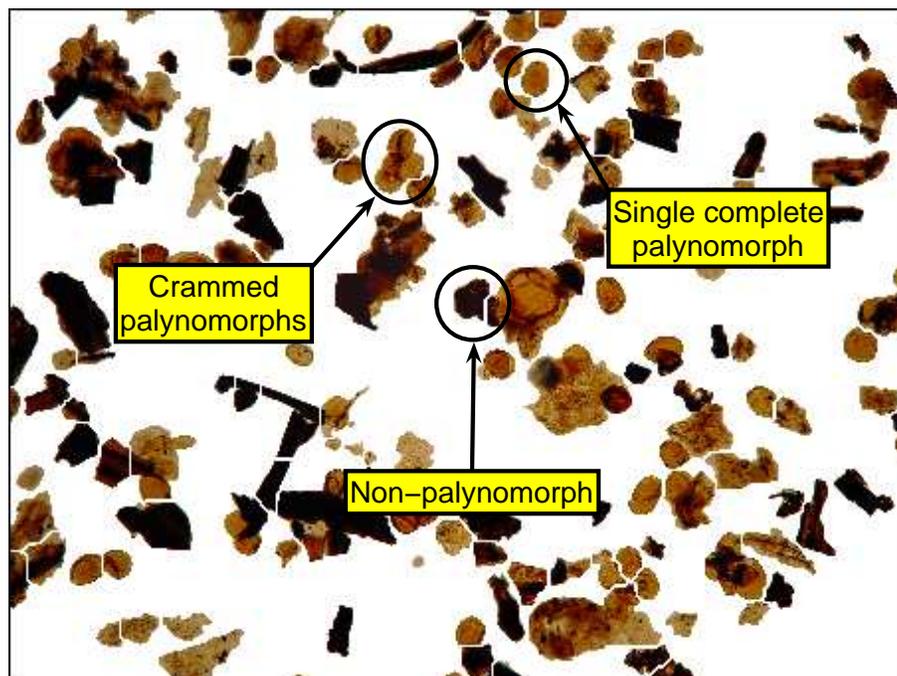
The CSS algorithm was originally designed to segment touching kerogen pieces. Under transmitted light microscopy kerogen will appear relatively dark. The only visual cue to determine the position of individual pieces is the silhouette of the material. Here we are applying CSS to the silhouette of all foreground material including kerogen and amorphous matter. A palynomorph region will contain a single complete elliptical palynomorph. An “other” region will contain a non-palynomorph such as kerogen, amorphous material or more than one palynomorph crammed together. Examples of these types of regions are shown in Figure 3. Therefore we expect the set of segmented regions found using CSS to contain both palynomorph and “other” regions.

### 2.2 Classification

Classifiers are a mapping from an input space consisting of object features to an output space of discrete labels. In our case the classification is learnt from training data consisting of segmented objects which have been hand labelled. An object is represented by a vector of feature values extracted from each region. We have used 32 features from four groups describing colour, shape, size and texture. The features are general descriptions for any type of region and are not specific to palynomorphs. Each feature is explained in Table 1.



**Fig. 2** Demonstrating the pre-processing steps before classification. (a) Image converted to grayscale. (b) Foreground/background segmentation subsequent to image normalisation using the background estimate. (c) Segmentation of microfossils using CSS algorithm.



**Fig. 3** Final segmentation annotated to show the types of segmented regions

**Table 1** Groups of features

Type	Feature	Notation	Explanation
Colour	mean red		
	mean blue		
	mean green		
	mean gray		
Size	inner radius	$r_i$	Radius of the largest circle contained entirely within the object
	outer radius	$r_o$	Radius of the smallest circle which contains the entire object
	diameter	$d$	The maximum distance between 2 contour pixels
	perimeter	$p$	Number of pixels on the border between object and background
	circle difference	$r_o - r_i$	Difference between the outer and inner radii
	area	$a$	Total number of pixels comprising the object
	distance	$\delta$	Mean distance from centre of gravity to all contour pixels
Texture	entropy		Entropy of the grey-level histogram taken as a pdf
	anisotropy		Symmetry of the grey-level histogram about its median
	correlation		Correlation between grey level intensity of neighbouring pixels
	homogeneity		Homogeneity of neighbouring pixels in the grey level image
	contrast		Contrast of neighbouring pixels in the grey level image
	energy		Energy of neighbouring pixels in the grey level image
	rim variability		Variance of the gray level intensity in a "rim" of width $r_i/5$
Shape	anisometry	$e_+/e_-$	Ratio of the lengths of the major and minor elliptic semi-axes
	eccentricity	$d_-/d$	Ratio of the length of the minor axis of the object to $d$
	rectangularity	$a/a_b$	Ratio of object area to the area of smallest bounding rectangle
	bulkiness	$\frac{\pi(e_+)(e_-)}{a}$	Ratio of the areas of a corresponding ellipse and the object
	convexity	$a/a_c$	Ratio of the object's area to its convex area
	variance x		Variance across x-axis with respect to centre of gravity
	variance y		Variance across y-axis with respect to centre of gravity
	covariance		
	compactness	$4\pi a/p^2$	Ratio of the area to that of a circle with the same perimeter
	sigma	$\sigma$	Standard deviation of distances from centre of gravity to contour
	roundness	$1 - \sigma/\delta$	
	sides	$1.41 \left(\frac{\delta}{\sigma}\right)^{0.4724}$	Number of pieces of a regular polygon
	equant/lath	$r_i/d$	Equant/lath ratio
	structure factor	$\frac{\pi(e_+)^2}{a} - 1$	anisometry $\times$ bulkiness - 1

### 2.2.1 Data classifiers

We have chosen to test ten well known classifiers, all have been shown to achieve a good standard across a diverse range of datasets. The classifiers chosen are logistic (19), bagging (6), support vector machines (SVM) (12), multilayer perceptron (3), random forest (7), logitboost (17), adaboost (16), decision tree (5), nearest neighbour (13) and naive bayes (18). All classifiers are implemented in Weka (27) using their default parameter settings.<sup>2</sup>

### 2.2.2 Cross-validation

Cross validation is used to determine the accuracy of a classifier. The dataset is split into two sections called the training and testing set. We train the classifiers on the training set and then calculate its accuracy on the testing set.  $K$ -fold cross-validation partitions the dataset into  $K$  subsamples. One subsample is used for testing and the other  $K - 1$  are used for training. This process is carried out  $K$  times where each of the  $K$  subsamples are used exactly once for testing. At this point the  $K$  accuracies can then be averaged. However, to remove a bias towards the initial partitioning

of the dataset the whole process is repeated  $N$  times each time using a different partitioning. The final accuracy obtained is found by averaging the  $N \times K$  results. In this study we choose to use 10-fold cross-validation 10 times.

### 2.2.3 Feature selection

A subset of the 32 features can be chosen to train the classifier and will hopefully improve classification accuracy. A subset of features is selected using a greedy step-wise approach within a 10-fold cross-validation. The method begins by using the full set of 32 features and removing a feature if it reduces the classification accuracy. Each feature is then ranked depending on when it was removed. For example, if feature A is removed first then it will be ranked 1 and if feature B is removed last it will be ranked 32. An average of the 10 ranks for each feature, obtained from the 10 folds, is used as a measure of feature importance. Feature selection also serves an additional purpose, which is to guard against possible overtraining.

### 2.2.4 Class imbalance

In some cases the dataset contains many more samples of one class than the other. The trivial (largest prior)

<sup>2</sup> Weka is a free software environment for machine learning and data mining. <http://www.cs.waikato.ac.nz/ml/weka/>

classifier labels all samples according to the most popular class in the training set. Although reasonable accuracies can be achieved, such a classifier would be useless. Problems like this are termed *imbalanced*. Three main approaches have been employed to solve imbalanced problems (2): i) we can assign a cost to classification errors. ii) the discrimination process can be internally biased to account for the class imbalance. iii) we can sample from the training set to balance the class distribution by either over-sampling from the minor class or under sampling from the major class. Another alternative to is to sample from both classes with replacement maintaining a balanced class distribution. In this study, we choose to use the alternative approach.

### 2.2.5 ROC analysis

The posterior probability of the palynomorph class can be estimated from the output of the classifier. We use the inbuilt functions in Weka to estimate the posterior probabilities. The classification output can be decided by thresholding the posterior probability that a given fossil is a true palynomorph, we call this the *certainty threshold* (CT). A ROC space is formed by two axes. The  $x$  axis is (1-specificity) and the  $y$  axis is the sensitivity. A perfect classifier will be at point (0,1) in this space. Random classification lies on the diagonal line running from the point (0,0) to the point (1,1). Points above this line are better than a random guess. For each CT we can calculate the sensitivity and specificity of the results and plot a point in this space. In this way a curve is formed known as a ROC curve (15). Classifier performance can be measured as the total area under this curve (AUC), the closer this measure is to 1 the better the classifier. The best CT is found by locating the point on the ROC curve closest to (0,1).

## 3 Experimental results

The pre-processing step is applied to seven microscopy images containing palynofacies. In total 1139 objects are extracted. Each region is hand labelled as a palynomorph if it consists of a single palynomorph with an elliptic shape otherwise it is labelled as “other”. To achieve manual classification the palynofacies are initially segmented using the CSS algorithm. A human expert is then able to sit in front of a monitor and select the regions segmented by the CSS algorithm which they deem to be complete elliptic palynomorphs. This data is stored and a ground truth is obtained. For each region the 32 features from Table 1 are extracted and a dataset is formed. We will use this to train and compare classifier accuracy and AUC.

The dataset contained 142 palynomorph objects and 997 classed as “other”. Due to class imbalance a bootstrap sample of the training data is drawn, so that the classes have approximately 50/50 representation. This type of sampling was also done using Weka. Accuracy is calculated by performing a 10-fold cross validation 10 times. To remove a bias towards the training sample we repeated the bootstrap sampling 10 times. Classification accuracy and AUC is therefore found as an average of 1000 testing sets of size 113 objects each. A 95% confidence interval (CI) is retrieved by finding the 26th and 975th largest accuracy. We trialled four bootstrap sample sizes: 30%, 50%, 70% and 90% of the training data. It was found that higher accuracies and AUC’s were achieved when a sample size of 90% is used.

For a sample size of 90%, the accuracies of each classifier together with their CI’s are shown in Table 2, AUC’s are shown in Table 3. The CI’s are heavily overlapping indicating no clear winner between the various classifiers. AUC values indicate the logisitc classifier to be the best. We have chosen to use the logistic classifier due to its simplicity, high classification accuracy, high AUC and speed of training/testing.

**Table 2** Accuracy and 95% CI’s of the ten classifiers using a 10 times 10-fold cross-validation.

Classifier	Accuracy (%)	95% CI (%)
Logisitic	88.31	(82.45, 92.11)
Bagging	90.77	(85.09, 94.74)
SVM	85.33	(78.94, 92.11)
MLP	88.24	(81.58, 93.86)
Random Forest	92.89	(88.59, 93.86)
LogitBoost	88.04	(81.58, 93.86)
AdaBoost	85.29	(77.19, 92.98)
Decision Tree	88.24	(81.58, 92.11)
Nearest Neighbour	88.49	(83.33, 92.11)
Naive Bayes	75.74	(67.54, 84.96)

**Table 3** AUC’s of the ten classifiers using a 10 times 10-fold cross-validation.

Classifier	AUC
Logisitic	0.946
Bagging	0.934
SVM	0.938
MLP	0.873
Random Forest	0.933
LogitBoost	0.937
AdaBoost	0.926
Decision Tree	0.822
Nearest Neighbour	0.845
Naive Bayes	0.824

## 4 Classification using the Logistic Classifier

Let  $\mathbf{x}$  be the region to be classified, where  $\mathbf{x} = [x_1, \dots, x_n]^T \in \mathbb{R}^n$ . Let  $w_1, \dots, w_c$  be the class labels,  $P(w_i)$  be the prior probabilities and  $P(w_i|\mathbf{x})$  be the posterior probabilities for the classes,  $i = 1, \dots, c$ . The logistic classifier relies on the assumption that the log-odds of the posterior probabilities for any two classes can be approximated as a linear function. Without loss of generality, we can pick class  $w_c$  and fix its discriminant function to be  $g_c(\mathbf{x}) = 0$  for any  $\mathbf{x}$ . The remaining  $c - 1$  discriminant functions are calculated as

$$g_i(\mathbf{x}) = \log \frac{P(w_i|\mathbf{x})}{P(w_c|\mathbf{x})} \\ = \beta_{i0} + \sum_{j=1}^n \beta_{ij} x_j, \quad i = 1, \dots, c - 1,$$

where  $\beta_{ij}$ , are the coefficients obtained through training the classifier. The training is done by the *Iterative Reweighted Least Squares (IRLS)* method using the *Newton-Raphson* updates (4). Let  $g(\mathbf{x})$  be the output of the discriminant function produced by the logistic classifier. An estimate of the posterior probability  $P(\text{palynomorph}|\mathbf{x})$  is  $\frac{1}{1 + \exp(-g(\mathbf{x}))}$ .

The greedy stepwise feature selection method is performed and the top 10 ranked features are selected. These features are listed in table 4 along with their average ranks. Classification using the logistic classifier will be conducted using only these 10 features.

Because of the class imbalance, the trivial (largest prior) classifier has specificity 100% and sensitivity 0%, while its accuracy is a reasonable 87%. By applying a threshold to the posterior probability we will be able to adjust the sensitivity and specificity until a good compromise is found. We evaluate the performance using ROC analysis.

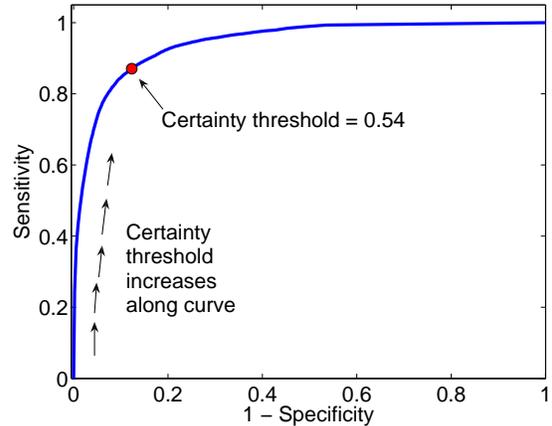
In this study the CT was increased from 0 to 1 in steps of 0.001. The specificity and sensitivity of the logistic classifier was found at each step as an average of 1000 testing sets (in the same way we calculated the accuracy). This yields 1001 points in ROC space, joining the points forms a ROC curve shown in Figure 4. The best CT was found at 0.54 with a specificity of 88%, sensitivity of 87% and an accuracy of 88%.

### 4.1 Example

For demonstration purposes we will train the classifier on six of our images. The seventh image will be used to test the classifier and display the objects labelled as a palynomorph. The training set is re-sampled to attain a balanced class distribution. The logistic classifier is

**Table 4** Top 10 features ranked by the greedy stepwise selection procedure using the logistic classifier (average of 10-fold cross-validation).

Feature	Rank $\pm$ std
distance	30.20 $\pm$ 0.87
covariance	24.70 $\pm$ 6.34
eccentricity	24.40 $\pm$ 5.62
outer radius	24.20 $\pm$ 7.59
mean red	24.00 $\pm$ 6.28
mean blue	23.60 $\pm$ 4.72
mean green	22.00 $\pm$ 7.60
rectangularity	21.60 $\pm$ 6.55
anisotropy	21.50 $\pm$ 8.64
variance y	21.40 $\pm$ 10.50



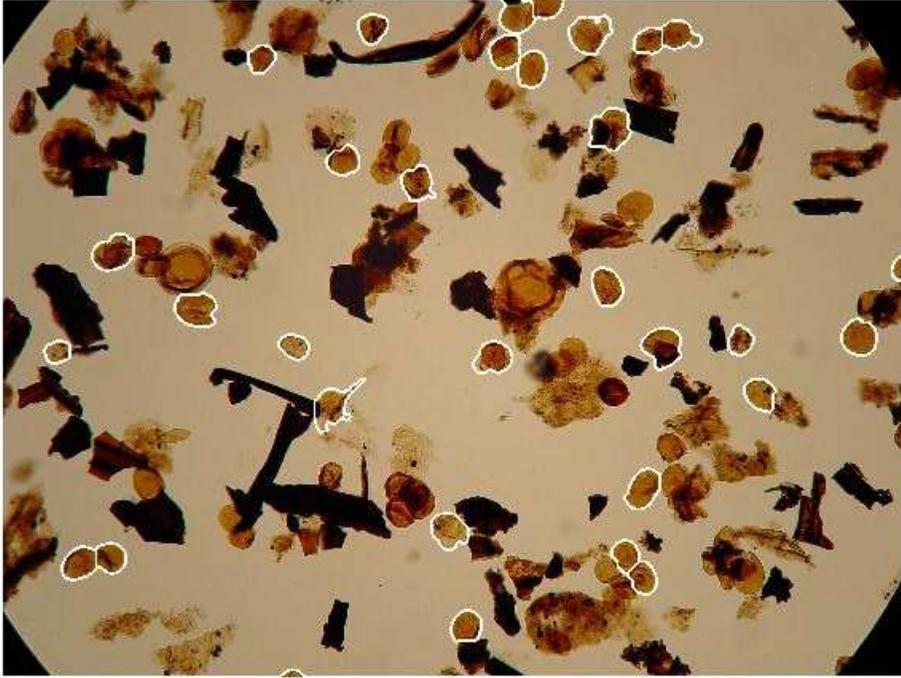
**Fig. 4** ROC curve of logistic classifier with respect to changes in the CT

trained and then tested on the segmented image shown in Figure 3 using the 10 features listed in table 4. The ROC analysis is performed and the best CT was found to be 0.53. The accuracy is 91% with a sensitivity of 94% and specificity of 90%. The regions hand labelled as palynomorph by human expert are shown in Figure 5 highlighted with a white border. The results of the logistic classifier are illustrated in Figure 6 and can be visually compared with the human expert.

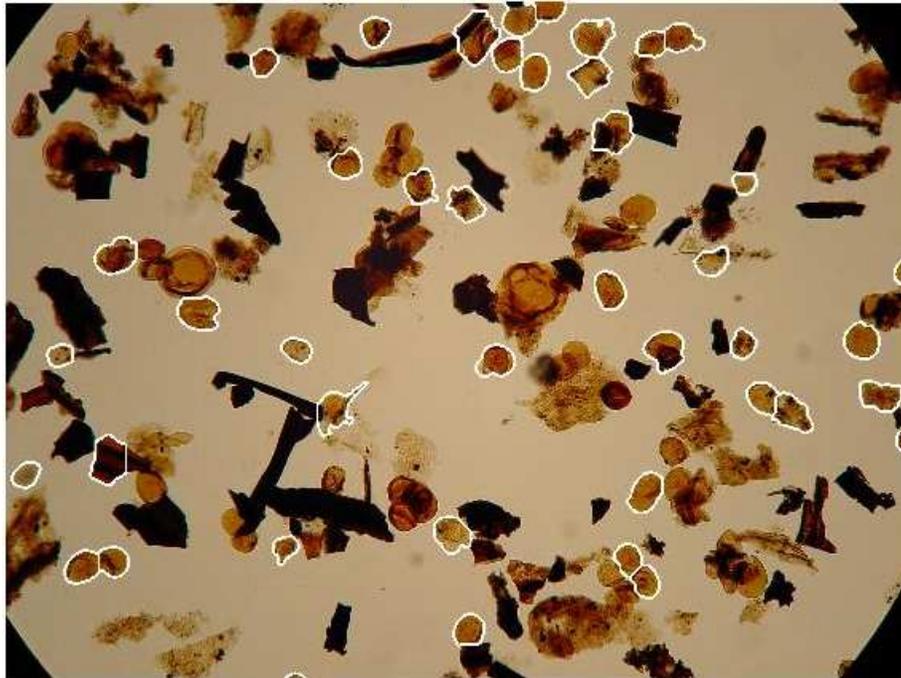
The 11 non-palynomorph regions that were labelled as palynomorph by the logistic classifier contained mainly amorphous material. Because the amorphous material has a rough appearance the number of false positives could be reduced by adding more specific texture features. Only two regions hand labelled as palynomorph were undetected by the classifier.

## 5 Conclusion

Automatic classification of single complete palynomorphs can be achieved with high accuracy. However to obtain a single palynomorph from the image of a slide



**Fig. 5** Original image displaying regions labelled by *human expert* as single elliptic palynomorphs. Regions are highlighted with a white border.



**Fig. 6** Original image displaying regions labelled by *logistic classifier* as single elliptic palynomorphs. Regions are highlighted with a white border.

containing highly irregular, overlapping kerogen pieces and other debris requires an advanced automatic segmentation technique. Such techniques are not 100% accurate and some regions will contain more than one palynomorph or a non-palynomorph. The experiments conducted here have shown that it is possible to automatically locate regions in an image that contain only palynomorphs with an accuracy of 88%. The system consists of three steps: 1) a pre-processing of the original image to correct background intensity levels and segment the background from the foreground. 2) Image segmentation using the CSS algorithm. 3) Classification of segmented regions using the logistic classifier.

Having identified regions containing single complete palynomorphs, a second classifier can be trained to identify their subgroup, species etc. For example the neural network constructed by (26) can now be applied to an image of a microscope slide containing many microfossils rather than images of a single complete microfossil.

Avenues for improvement consist of increasing classification accuracy and performance by tuning the parameters of the classification algorithms.

Future work will focus firstly on the idea of improving the segmentation of palynomorphs using a trained classifier. For instance when automatically segmenting a set of images the usual approach is to fix the parameters of the segmentation algorithm and then apply it to all images. The preferred approach would be to adjust these parameters for each image until the best segmentation is found. The best segmentation in our case would be one that maximises the total number of segmented palynomorphs. By using the trained logistic classifier to count the number of palynomorph regions it may be possible to automatically tune the parameters of the CSS algorithm until a maximum number of palynomorphs have been detected. Secondly, we would also like to develop a system for automatically classifying visually distorted microfossils and segmenting them from the heavily overlapping regions.

## References

- Athersuch, J., Banner, F., Higgins, A., Howarth, R., Swaby, P.: The application of expert systems to the identification and use of microfossils in the petroleum industry **26**(4), 483–489 (1994)
- Barandela, R., Sanchez, J., Garcia, V., Rangel, E.: Strategies for learning in class imbalance problems. *Pattern Recognition* **36**, 849–851 (2003)
- Bishop, C.: *Neural Networks for Pattern Recognition*. Clarendon Press, Oxford (1995)
- Bishop, C.: *Pattern Recognition and Machine Learning*. Springer, New York (2006)
- Breiman, L.: *Classification and Regression Trees*. Wadsworth International, Belmont (1984)
- Breiman, L.: Bagging predictors. *Machine Learning* **26**(2), 123–140 (1996)
- Breiman, L.: Random forests. *Machine Learning* **45**, 5–32 (2001)
- Charles, J., Kuncheva, L., Wells, B., Lim, I.: An evaluation measure of image segmentation based on object centres. *LNCS Image analysis and recognition* **4141**, 283–294 (2006)
- Charles, J., Kuncheva, L., Wells, B., Lim, I.: Object segmentation within microscope images of palynofacies. *Computers & Geosciences* (2008a). (accepted, <http://dx.doi.org/10.1016/j.cageo.2007.09.014>)
- Charles, J., Kuncheva, L., Wells, B., Lim, I.: Background segmentation in microscope images. In *Proc 3rd International Conference on Computer Vision Theory and Applications VISAPP08* (2008b)
- Charles, J., Kuncheva, L., Wells, B., Lim, I.: Stability of kerogen classification with regard to image segmentation. *Mathematical Geology* (2008c). (submitted)
- Cristianini, N., Shawe-Taylor, J.: *An Introduction to Support Vector Machines*. Cambridge University Press, Cambridge, UK. (2000)
- Duda, R., Hart, P., Stork, D.: *Pattern Classification*. John Wiley & Sons, Ny, second edition (2001)
- Evitt, W.: A discussion and proposals concerning fossil dinoflagellates, hystrichospheres and acritarchs. *Proceedings of National Academy of Sciences* **49**, 158–164 (1963)
- Fawcett, T.: ROC Graphs: Notes and Practical Considerations for Researchers. *Machine Learning* **31** (2004)
- Freund, Y., Schapire, R.: A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Science* **55**(1), 119–139 (1997)
- Friedman, J., Hastie, T., Tibsharani, R.: Additive logistic regression: a statistical view of boosting. *Annals of Statistics* **28**(2), 337–374 (2000)
- Hand, D., Yu, K.: Idiot's Bayes - not so stupid after all? *International Statistical Review* **69**, 385–398 (2001)
- Hastie, T., Tibsharani, R., Friedman, J.: *Elements of Statistical Learning*. Springer, New York (2001)
- Hill, S.: Outline extraction of microfossils in reflected light images. *Computers & Geosciences* **14**(4), 481–488 (1988)
- Kuncheva, L., Charles, J., Wells, B., Lim, I.: Automated kerogen classification in microscope images of dispersed kerogen preparation. *Mathematical*

- Geology (2008). (accepted)
22. McLaughlin, R.: Randomized Hough Transform: Improved ellipse detection with comparison. *Pattern Recognition Letters* **19**(3-4), 299–305 (1998)
  23. Riedel, W.: A prolog program to help identify fossils. *Computers & Geosciences* **15**(5), 809–823 (1989)
  24. Swaby, P.: Vides: An expert system for visually identifying microfossils. *IEEE Expert* pp. 36–42 (1992)
  25. Weller, A., Corcoran, J., Harris, A., Ware, J.: The semi-automated classification of sedimentary organic matter in palynological preparations. *Computers & Geosciences* **31**(10), 1213–1223 (2005)
  26. Weller, A., Harris, A., Ware, J., Jarvis, P.: Determining the saliency of feature measurements obtained from images of sedimentary organic matter for use in its classification. *Computers & Geosciences* **32**(9), 1357–1367 (2006)
  27. Witten, H., Frank, E.: *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, 2nd edition (2005)