

Accuracy Estimation of Detection of Casting Defects in X-Ray Images Using Some Statistical Techniques

Romeu Ricardo da Silva and Domingo Mery

Departamento de Ciencia de la Computación,
Pontificia Universidad Católica de Chile, Vicuña Mackenna 4860 (143)
romeu@romeu.eng.br, dmey@ing.puc.cl
www.romeu.eng.br
http://dmery.puc.cl

Abstract. Casting is one of the most important processes in the manufacture of parts for various kinds of industries, among which the automotive industry stands out. Like every manufacturing process, there is the possibility of the occurrence of defects in the materials from which the parts are made, as well as of the appearance of faults during their operation. One of the most important tools for verifying the integrity of cast parts is radioscopy. This paper presents pattern recognition methodologies in radioscopic images of cast automotive parts for the detection of defects. Image processing techniques were applied to extract features to be used as input of the pattern classifiers developed by artificial neural networks. To estimate the accuracy of the classifiers, use was made of random selection techniques with sample reposition (Bootstrap technique) and without sample reposition. This work can be considered innovative in that field of research, and the results obtained motivate this paper.

Keywords: Casting Defects, Radioscopy, Image Processing, Accuracy Estimation, Bootstrap.

1 Introduction

Shrinkage as molten metal cools during the manufacture of die castings can cause defect regions within the workpiece. These are manifested, for example, by bubble-shaped voids, cracks, slag formation, or inclusions. Light-alloy castings for the automotive industry, such as wheel rims, steering knuckles, and steering gear boxes are considered important components for overall roadworthiness. To ensure the safety of construction, it is necessary to check every part thoroughly. Radioscopy rapidly became the accepted way for controlling the quality of die castings through computer-aided analysis of X-ray images [1]. The purpose of this nondestructive testing method is to identify casting defects, which may be located within the piece and thus are undetectable to the naked eye.

Two classes of regions are possible in a digital X-ray image of an aluminium casting: regions belonging to regular structures (RS) of the specimen, and those relating to defects (D). In an X-ray image we can see that the defects, such as voids, cracks and bubbles (or inclusions and slag), show up as bright (or dark) features. The reason is that X-ray attenuation in these areas is lower (or higher). Since contrast in

the X-ray image between a flaw and a defect-free neighbourhood of the specimen is distinctive, the detection is usually performed by analysing this feature (see details in [2] and [3]). In order to detect the defects automatically, a pattern recognition methodology consisting of five steps was developed [1]: a) Image formation, in which an X-ray image of the casting that is being tested is taken and stored in the computer. b) Image pre-processing, where the quality of the X-ray image is improved in order to enhance its details. c) Image segmentation, in which each potential flaw of the X-ray image is found and isolated from the rest of the scene. d) Feature extraction, where the potential flaws are measured and some significant features are quantified. e) Classification, where the extracted features of each potential flaw are analysed and assigned to one of the classes (regular structure or defect).

Although several approaches have been published in this field (see for example a review in [1]), the performance of the classification is usually measured without statistical validation. This paper attempts to make an estimation of the true accuracy of a classifier using the Bootstrap technique [4] and random selection without repositioning applied to the automated detection of casting defects. The true accuracy of a classifier is usually defined as the degree of correctness of data classification not used in its development. The great advantage of this technique is that the estimation is made by sampling the observed detection distribution, with or without repositioning, to generate sets of observations that may be used to correct for bias. The technique provides nonparametric estimates of the bias and variance of a classifier, and as a method of error rate estimation it is better than many other techniques [5].

The rest of the paper is organised as follows: Section 2 outlines the methodology used in the investigation. Section 3 shows the results obtained recently on real data. Finally, Section 4 gives concluding remarks.

2 Methodologies

2.1 Processing of the Casting Images

The X-ray image taken with an image intensifier and a CCD camera (or a flat panel detector), must be pre-processed to improve the quality of the image. In our approach, the pre-processing techniques are used to remove noise, enhance contrast, correct the shading effect, and restore blur deformation [1].

The segmentation of potential flaws identifies regions in radioscopic images that may correspond to real defects. Two general features of the defects are used to identify them: a) a flaw can be considered as a connected subset of the image, and b) the grey level difference between a flaw and its neighbourhood is significant. According to these features, a simple automated segmentation approach was suggested in [6] (see Fig. 1). First, a Laplacian of Gaussian (LoG) kernel and a zero crossing algorithm [7] are used to detect the edges of the X-ray images. The LoG-operator involves a Gaussian lowpass filter which is a good choice for pre-smoothing our noisy images that are obtained without frame averaging. The resulting binary edge image should produce closed and connected contours at real flaws which demarcate regions. However, a flaw may not be perfectly enclosed if it is located at an edge of a regular structure as shown in Fig. 1c. In

order to complete the remaining edges of these flaws, a thickening of the edges of the regular structure is performed as follows: a) the gradient of the original image is calculated (see Fig. 1d); b) by thresholding the gradient image at a high grey level a new binary image is obtained; and c) the resulting image is added to the zero crossing image (see Fig. 1e). Afterwards, each closed region is segmented as a potential flaw. For details see a description of the method in [6].

All regions enclosed by edges in the binary image are considered 'hypothetical defects' (see example in Fig. (1e)). During the feature extraction process the properties of each of the segmented regions are measured. The idea is to use the measured features to decide whether the hypothetical defect corresponds to a flaw or a regular structure.

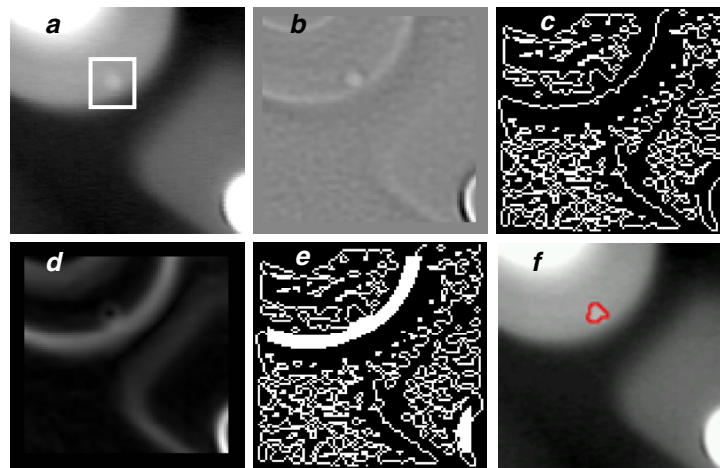


Fig. 1. Detection of flaws: a) radioscopic image with a small flaw at an edge of a regular structure, b) Laplacian-filtered image with $\sigma = 1.25$ pixels (kernel size = 11×11), c) zero crossing image, d) gradient image, e) edge detection after adding high gradient pixels, and f) detected flaw using feature $F1$ extracted from a crossing line profile [2]

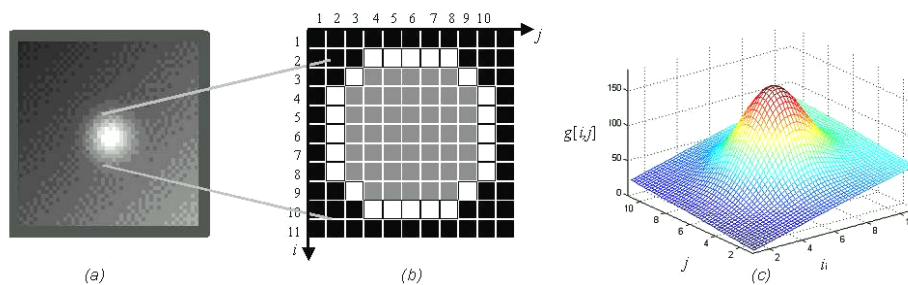


Fig. 2. Example of a region. (a) X-Ray image, (b) segmented region, (c) 3D representation of the intensity (grey value) of the region and its surroundings [8].

Table 1. Descriptions of the features extracted

f1 and f2	Height (f1) and width (f2): height (h) and width (w) of the region [9].
f3	Area (A): number of pixels that belong to the region [9].
f4	Mean grey value (G): mean of the grey values that belong to the region [9].
f5	Mean second derivative (D): mean of the second derivative values of the pixels that belong to the boundary of the region [9].
f6	Crossing Line Profile (F_1): Crossing line profiles are the grey level profiles along straight lines crossing each segmented potential flaw in the middle. The profile that contains the most similar grey levels in the extremes is defined as the best crossing line profile (BCLP). Feature F_1 corresponds to the first harmonic of the fast Fourier transformation of BCLP [2].
f7	Contrast $K\sigma$: standard deviation of the vertical and horizontal profiles without offset [9].
f8	High contrast pixels ratio (r): ratio of number of high contrast pixels to area [3].

The features extracted in this investigation are described below (Table 1), and they provide information about the segmented regions and their surroundings.

The total number of features extracted is 8 divided into 3 geometric features and 5 intensity features. In our work we present results obtained on 72 radioscopic images of aluminium die castings. The size of the images is 572×768 pixels. About 25% of the defects of the images were existing blow holes (with $\varnothing = 2.0 - 7.5$ mm). They were initially detected by visual (human) inspection. The remaining 75% were produced by drilling small holes (with $\varnothing = 2.0 - 4.0$ mm) in positions of the casting which were known to be difficult to detect. In these experiments, 424 potential defects were segmented, 214 of them correspond to real defects, while the others are regular structures (210).

2.2 Development of the Nonlinear Classifiers

The non-linear classifiers were implemented using a two-layer neural network with training by error backpropagation. The first step taken in the development of a non-linear classifier was to optimize the number of neurons used in the intermediate layer in order to obtain the best accuracy possible for the test sets. Some tests were carried out in terms of training parameters of the network, and the best result (fastest convergence) was found when the moment ($\beta=0.9$) and α (training rate) variables were used [10, 11]. The initialization of the synapses and bias used the Widrow [12] method. All these training variations resulted in a convergence for the same range of error.

2.3 Accuracy Estimation

There are various techniques to estimate the *true accuracy* of a classifier, which is usually defined as being the degree of correctness of classification of data not used in its development. The three that are most commonly used are: simple random selection of data, cross validation that really presents diverse implementations [13], and the bootstrap technique [4, 14]. It is not really possible to confirm whether one method is better than the other for any specific pattern classification system. The choice of one of these techniques will depend on the quantity of data available and the specific classification to be made.

As described in [4], two properties are important when evaluating the efficiency of an estimator $\hat{\theta}$, its bias and its variation, that are defined by the equations below:

$$Bias = E[\hat{\theta}] - \theta \tag{1}$$

$$Var(\hat{\theta}) = E\left[(\hat{\theta} - E[\hat{\theta}])^2 \right] \tag{2}$$

where,

$E[\hat{\theta}]$: expected value of estimator $\hat{\theta}$.

$Var(\hat{\theta})$: variation of estimator.

An estimator is said to be reliable if it contains low values of bias (trend) and variation. However, in practice an appropriate relation between both is desirable when looking for a more realistic objective [4, 14]. When dealing with the accuracy of a classifier, bias and variation of the estimated accuracy are going to vary as a function of the number of data and the accuracy estimation technique used.

In this work, to calculate the classification accuracy of casting defects we first carried out the bootstrap technique as follows:

A set of bootstrap data (size n), following Efron's definition [4], is made up of $x_1^*, x_2^*, \dots, x_n^*$ data, obtained in a random way and with repositioning, from an original set of data x_1, x_2, \dots, x_n (also size n). In this way it is possible for some data to appear 1, 2, 3 or n times or no times [4]. With this technique the classifier implemented using the i^{th} training set is tested with data that were not used in the make up of this set, resulting in an accuracy estimator of $\hat{\theta}_i$ (for test data). This is repeated b times. The model of bootstrap accuracy estimation $\hat{\theta}_B$ of frequently used pattern classifiers is defined by

$$\hat{\theta}_B = \frac{1}{b} \sum_{i=1}^b (\hat{\omega} \hat{\theta}_i + (1 - \hat{\omega}) \hat{\theta}_c) \tag{3}$$

where $\hat{\theta}_c$ is the apparent accuracy (calculated with the training set data only) and the weight $\hat{\omega}$ varies between 0.632 and 1, which is normally taken as being equal to 0.632 [4, 14].

As a second way of estimating the accuracy of the developed classifiers, the form of random selection without data reposition was used for the formation of the training and testing sets, different from the Bootstrap technique [15]. In addition to that, ROC curves were drawn to verify the reliability of the results achieved with this technique [11].

3 Results

3.1 Features Selection

An optimized way of representing the domains of the classes of patterns of multivariate system in a two-dimensional space is by obtaining the two main discrimination components. It is known that the main linear discrimination address is called Fisher's Discriminator [11], and it maximizes the interclass covariance matrix and minimizes the intraclass covariance matrix [11, 16]. In this case, the first linear discrimination address of classes RS and D can be obtained going over a supervised neural network of the backpropagation type with only one neuron [10]. Then it is possible to obtain a second main linear discrimination address, also with a neural network with only one neuron, using for the training of the network the residual information of the projection of the original information in the first discrimination address, what is called independent components (orthogonals). A detailed description of this technique is found in [17].

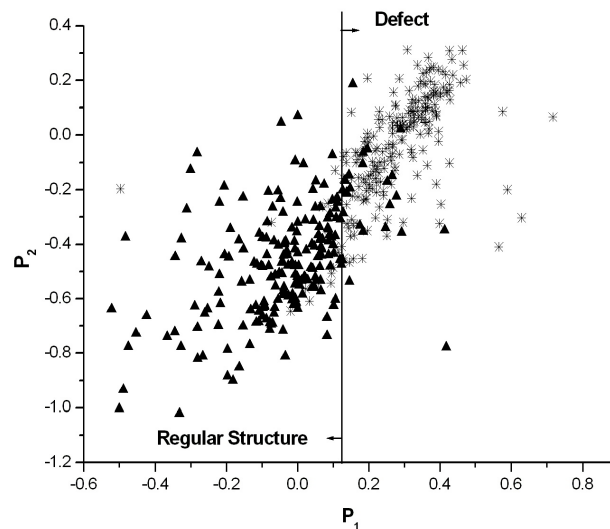


Fig. 3. Graphs made with the two principal linear discrimination components

In this way the two main components of the linear discrimination of classes RS e D with a neural network of only one neuron which was trained through the error backpropagation algorithm using batch training (3000 periods), parameter $\beta=0.9$ and α variable, were obtained. Figure 3 shows the graph obtained with those two main linear discrimination addresses. It is evident that the separation of classes RS and D is more efficient in that representation space, because a visual analysis will make it possible to identify that there are few false positive (RS inputs in the domain space of D) and false negative (D inputs in the domain space of RS) errors. The projection of the data on the x axis (p1) represents what would be the best discrimination of these classes, and a projection on y (p2), the second best discrimination. From this graph it is concluded that the separation between RS and D can achieve good indices of success with well developed pattern classifiers.

3.2 Study of Neuron Number in the Intermediate Layer

The graph of Figures 3 showed the problem of classification of classes RS and D only from the two principal linear discrimination components. However, it is well known that the linear pattern classifiers solve well very easy class separation problems [11]. To optimize the separation between the classes of patterns RS and D, non-linear pattern classifiers will be developed through supervised neural networks with two layers of neurons and error backpropagation training [10].

Since non-linear classifiers can have network overtraining problems, whose probability increases with increasing number of neurons in the second layer, thereby losing the capacity to generalize [10], to decrease the probability of the existence of *overfitting* the parameters of the non-linear classifier, a study was made of the optimum number of neurons in the intermediate layer of the classifier that would make possible the best result with test sets. For that purpose, from the initial set of data with the eight features, a training set was chosen with 75% of the data chosen randomly and without reposition, and a test set with the remaining 25%, keeping the proportion between the classes. In this way the training set contained 158 samples of RS and 160 of D, and the test set had 52 of RS and 54 of D. The number of neurons in the intermediate layer of the network was varied one at a time up to 20 neurons, and the indices of success in classification and testing were recorded. It should be noted that, since we are dealing with only two classes of patterns, the last layer of the classifier can contain only one neuron.

The results obtained from the study of the number of neurons are shown in Table 2. In the table it is seen that the smallest difference between the results of the training and the tests, which theoretically can indicate a good generalization capacity of the classifier, occurs for two neurons in the intermediate layer. However, if we analyse the increase of the performance of the classifier, which occurs significantly with the increase in the number of neurons, which is expected, a second lowest difference occurs for 11 neurons, achieving 94.34% of success with the test set. For that reason, 11 neurons were used in the intermediate layer of the neural network for the development of all the classifiers of this work having in view the estimation of the accuracy of the classification.

Table 2. Optimization of the number of neurons in the intermediate layer

Number of Neurons	Training Performance (%)	Test Performance (%)
1	90.57	86.80
2	90.25	89.63
3	94.66	89.63
4	97.50	91.51
5	96.90	91.51
6	97.80	89.63
7	96.90	91.51
8	99.06	89.63
9	98.75	93.40
10	98.43	92.46
11	98.43	94.34
12	99.06	92.46
13	99.38	92.46
14	99.38	93.40
15	99.38	92.46
16	99.06	92.46
17	99.38	90.57
18	99.38	89.63
19	99.38	92.46
20	99.70	94.34

3.3 Accuracy Estimation by the Bootstrap Technique

To estimate the accuracy of the non-linear classifiers, the first technique used was a random selection with data repositioning. Once more, as an example of the operation of this technique, one can imagine a “bag” with all the original data, then we choose data randomly from this bag to form the training sets, but every piece of data that goes to these sets returns to the “bag” and can be chosen again several times. The data that are not chosen for training are used in the formation of the test sets.

In this way, 10 pairs of training and test sets were formed. It should be noted that, with this selection technique the training sets always have the same number of data as the original set (in this work, 424 data). In this way, the test sets had a number of data between 150 ($\approx 35\%$) and 164 ($\approx 38\%$). For the development of the classifiers, with the aim of decreasing the possibility of overtraining of its parameters (synapses and bias), use was made of a validation set formed by samples selected randomly from the bootstrap training sets in a 10% proportion. This technique is well known as cross validation (Haykin, 1994), and the end of the training was set to when the validation error increases or remains stable for 100 epochs or a maximum of 3000 epochs, obviously choosing the values of the network parameters in the situation of the least validation error. The results are presented in Table 3.

Analysing these results, it is seen that the training indices were quite high, with a mean of 98.46%, however the success indices of the test were significantly lower, with a mean of 55.61%. Calculating the accuracy estimator according to the

weighting factor of 0.632 for the test set estimator and 0.368 for the training estimator [4, 13], the estimated accuracy result is 71.40%, which can be considered unsatisfactory for the classification of patterns for this problem of fault detection in automobile rims.

The great problem for the classification of patterns, which is common to almost all work in that relation, is the lack of data to estimate with precision the true classification accuracy, so that it can be trusted that all the success indices will always be similar when the classifier is tested with a new set of data. The main objective of the use of the bootstrap technique was to try to reproduce several sets for training and testing the classifiers as well as for estimating the accuracy expected for classes RS and D. One justification that can be thought of for the low success indices with that technique is the fact that the test sets have a large number of data in relation to the number of data used for training. Normally, in terms of pattern classification, the test or validation sets contain between 20 and 30% of data, and by the bootstrap technique, in this paper, some test sets get to contain almost 40% of the data, and this can in fact affect the correct training of the network parameters, even using a cross validation technique to interrupt the trainings. This is even more feasible if we think that the original data did not contain a large number of samples. To expect a success index of only about 55%, or even 71.40%, for this classification problem is too pessimistic having in mind the efficiency of the image processing techniques used and the relevance of the extracted features.

Table 3. Result of classification with the bootstrap input sets (%)

Input Sets	Training (%)	Test (%)
1	418/98.60	75/50.00
2	422/98.60	88/53.66
3	410/96.70	92/56.10
4	421/99.30	94/57.31
5	405/95.52	94/57.32
6	421/99.30	88/53.66
7	416/98.12	86/52.45
8	424/100	100/61.00
9	420/99.05	86/52.45
10	422/99.53	102/62.20
Mean	98.47	55.61
Bootstrap accuracy estimation $\hat{\theta}_b = \frac{1}{b} \sum_{i=1}^b (0,632\hat{\theta}_i + 0,368\hat{\theta}_c)$	71.40	

3.4 Accuracy Estimation by Random Selection Without Repositioning

In the simple method of evaluation with random sampling, the original data set (with n data) is partitioned randomly into two sets: a training set containing $p \times n$ data,

and a test set containing $(1-p) \times n$ data (the values of p are chosen in a variable way case by case). This process is repeated a number of times, and the mean value is the accuracy estimator [13]. This technique was used for the first selection and formation of sets with the purpose of choosing the number of neurons of the classifier's intermediate layer. Using that simple yet very efficient technique, 10 pairs of data sets for training and testing of the classifier were chosen, and the percent proportion chosen (based on experience from other work) was 75% for training (318) and 25% for testing (106).

Table 4 contains the results achieved successfully with these sets. The fourth and fifth columns of the table refer to the number of data of each class contained in the corresponding sets. The mean was approximately 53 data of each class in each set, that is, in general there was not a significant disproportion between the number of data of each class that would affect the trainings and tests of the classifiers. The training column contains not only the percentages of success, but also the number of data classified correctly, which were as high as those obtained with the bootstrap sets. However, it is seen that the test results were considerably higher than those achieved with the bootstrap technique, with a mean estimated accuracy of 90.30% for the 10 test sets selected, a very satisfactory index close to the mean of 97.52% obtained for the training sets. That small difference of about 7% is perfectly acceptable, and it shows the generalization of the classifiers (confirmed also by the low values found for standard deviation). It should be noted that with these sets cross validation was also used for interrupting the training in a manner similar to that used for the bootstrap sets.

Table 4 also contains the false negative (FN) indices, real defects classified as regular structures as well as the false positive (FP) indices, regular structures classified as defects. The mean values achieved of 7.69% and 11.64%, respectively, can be considered satisfactory, especially if we consider that the most critical situation

Table 4. Results of classification with the input sets of the random selection without repositioning (%)

Input Sets	Training (%)	Test (%)	RS	D	FN (%)	FP (%)
1	314/98.75	95/89.63	57	49	3.51	18.37
2	311/97.80	98/92.46	52	54	11.54	3.70
3	315/99.06	101/95.30	50	56	2.00	7.14
4	312/98.12	95/89.63	55	51	18.18	1.96
5	314/98.75	93/87.74	45	61	4.44	18.03
6	307/96.55	93/87.74	60	46	6.67	19.57
7	299/94.03	96/90.57	53	53	5.66	13.21
8	314/98.75	94/88.68	46	60	6.52	15.00
9	311/97.80	96/90.57	55	51	7.27	11.76
10	304/95.60	96/90.57	54	52	11.11	7.69
Mean (%)	97.52	90.30	≈53	≈53	7.69	11.64
Standard Deviation(%)	1.28	1.61			13.03	12.53

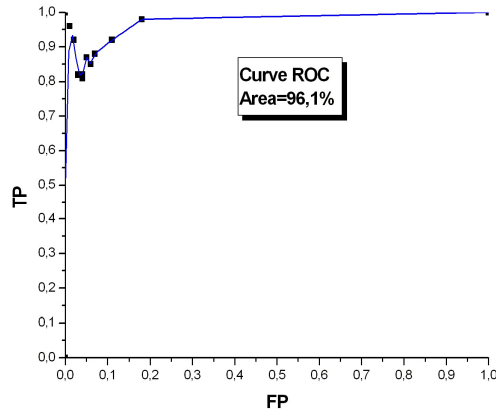


Fig. 4. ROC resultant curve of the randomly selected sets without data repositioning (sixth and seventh columns of Table 5)

is always that of false negative, and less than 8% of errors in the classification of real defects is an index that cannot be considered high for a fault detection situation in these kinds of images.

Figure 4 shows the ROC (Receiver Operating Characteristic) curve obtained from the interpolation of true positive (TP), 1- FN, and false positive points of Table 4. The area over the curve, calculated by simple integration of the interpolated curve, represents the efficiency of the system used for the detection of the real defects in the acquired images (probability of detection, PoD). In this case the value found for the area was 96.1%, which can be considered an optimum index of the efficiency and reliability of the system, higher than the 90.30% estimated accuracy value of Table 4.

4 Conclusions

As to the bootstrap technique, the accuracy results were well on this side of acceptable, and that can be explained by the small amount of data available in the training sets.

The estimation of the accuracy of classification with the random selection technique without data repositioning, with fixed values of 25% of data for the test sets, had high indices of correctness, showing the efficiency of the system developed for the detection of defects, which was also evident from the drawing of the ROC curve for the system.

It must be pointed out that this work does not exhaust the research in this field, and that much can still be done to increase the reliability of the results obtained as well as to increase the number of features to be extracted to increase the degree of success in the detection of faults. However, this paper can be considered pioneering dealing with defects in automobile wheels, and there are no results on estimated accuracy in other papers that could be used for comparison with these results.

Acknowledgment. This work was supported in part by FONDECYT – Chile (International Cooperation), under grant no. 7060170. This work has been partially supported by a grant from the School of Engineering at Pontificia Universidad Católica de Chile. We acknowledge the permission granted for publication of this article by Insight, the Journal of the British Institute of Non-Destructive Testing.

References

1. Mery, D.: Automated Radioscopic Testing of Aluminium die Castings. *Materials Evaluation* 64, 135–143 (2006)
2. Mery, D.: Crossing line profile: a new approach to detecting defects in aluminium castings. In: Bigun, J., Gustavsson, T. (eds.) SCIA 2003. LNCS, vol. 2749, pp. 725–732. Springer, Heidelberg (2003)
3. Mery, D.: High contrast pixels: a new feature for defect detection in X-ray testing. *Insight* 46, 751–753 (2006)
4. Efron, B., Tibshirani, R.J.: *An Introduction to the Bootstrap*. Chapman & Hall/CRC, New York (1993)
5. Webb, A.: *Statistical Pattern Recognition*, 2nd edn. John Wiley & Sons Inc, Chichester (2002)
6. Mery, D., Filbert, D.: Automated flaw detection in aluminum castings based on the tracking of potential defects in a radioscopic image sequence. *IEEE Trans. Robotics and Automation* 18, 890–901 (2002)
7. Castleman, K.: *Digital Image Processing*. Prentice-Hall, Englewood Cliffs, New Jersey (1996)
8. Mery, D., Silva, R.R., Caloba, L.P., Rebello, J.M.A.: Pattern Recognition in the Automatic Inspection of Aluminium Castings. *Insight* 45, 431–439 (2003)
9. Mery, D., Filbert, D.: Classification of Potential Defects in Automated Inspection of Aluminium Castings Using Statistical Pattern Recognition. In: 8th European Conference on Non-Destructive Testing (ECNDT 2002), Barcelona (June 17–21, 2002)
10. Haykin, S.: *Neural Networks - A Comprehensive Foundation*. Macmillan College Publishing, Inc, USA (1994)
11. Duda, R.O., Hart, P.E., Stork, D.G.: *Pattern Classification*, 2nd edn. John Wiley & Sons, U.S.A (2001)
12. Beale, M.: *Neural Network Toolbox for Use with Matlab User's Guide Version 4*. USA. The MathWorks (2001)
13. Diamantidis, N.A., Karlis, D., Giakoumakis, E.A.: Unsupervised Stratification of Cross-Validation for Accuracy Estimation. *Artificial Intelligence* 2000 116, 1–16 (2002)
14. Efron, B., Tibshirani, R.J.: Cross-Validation and the Bootstrap: Estimating the Error Rate of the Prediction Rule. Technical Report 477, Stanford University (1995), <http://utstat.toronto.edu/tibs/research.html>
15. Silva, R.R., Siqueira, M.H.S., Souza, M.P.V., Rebello, J.M.A., Calôba, L.P.: Estimated accuracy of classification of defects detected in welded joints by radiographic tests. *NDT & E International* UK 38, 335–343 (2005)
16. Silva, R.R., Soares, S.D., Calôba, L.P., Siqueira, M.H.S., Rebello, J.M.A.: Detection of the propagation of defects in pressurized pipes by means of the acoustic emission technique using artificial neural networks. *Insight* 48, 45–51 (2006)
17. Silva, R.R., Calôba, L.P., Siqueira, M.H.S., Rebello, J.M.A.: Pattern recognition of weld defects detected by radiographic test. *NDT&E International* 37, 461–470 (2006)