

# Pattern Detection Using a Maximal Rejection Classifier

Michael Elad \*

Yacov Hel-Or †

Renato Keshet ‡

## Abstract

In this paper we propose a new classifier - the Maximal Rejection Classifier (MRC) - for target detection. Unlike pattern recognition, pattern detection problems require a separation between two classes, *Target* and *Clutter*, where the probability of the former is substantially smaller, compared to that of the latter. The MRC is a linear classifier, based on successive rejection operations. Each rejection is performed using a projection followed by thresholding. The projection vector is designed to minimize the expected number of operations until detection. In the case where the probabilities of target and clutter signals are equal, it is shown that the Fisher linear discriminant is optimal in the above sense. However, in more common cases where the probabilities are quite different, a new optimal classifier is suggested. An application of detecting frontal faces in images is demonstrated using the MRC with encouraging results.

## 1 Introduction

In target detection applications, the aim is to detect occurrences of a specific *Target* in a given signal. In general, the target is subjected to some particular type of transformation, hence we have a set of target signals to be detected. In this context, the set of non-*Target* samples are referred to as *Clutter*. In practice, the target detection problem can be characterized as designing a classifier  $C(z)$ , which, given an input vector  $z$ , has to decide whether  $z$  belongs to the *Target* class  $\mathbf{X}$  or the *Clutter* class  $\mathbf{Y}$ . In example based classification, this classifier is designed using two training sets -  $\hat{\mathbf{X}} = \{x_i\}_{i=1..L_x}$  (*Target* samples) and  $\hat{\mathbf{Y}} = \{y_i\}_{i=1..L_y}$  (*Clutter* samples), drawn from the above two classes.

Various types of example-based classifiers are suggested in the literature [2, 3, 1]. The most simple and fast are the linear classifiers, where  $C(z)$  is based on a projection operation followed by a thresholding. The projection of  $z$  is performed onto a projection vector  $u$ , thus,  $C(z) = f(u^t z)$  where  $f(*)$  is a thresholding operation (or some other decision rule). The Support Vector Machine (SVM) [3] and the Fisher Linear Discriminant (FLD) [2] are two examples of linear classifiers. In both cases the kernel  $u$  is chosen in some optimal manner. In the FLD,  $u$  is chosen such that the Mahalanobis distance of the two classes after projection will be maximized. In the SVM approach the motive is similar, but the vector  $u$  is chosen such that it maximizes the margin between the two sets. In both these classifiers, it is assumed that the two classes have equal importance. In typical target detection applications, however, the above assumption is not valid since the probability of  $z$  belonging to  $\mathbf{X}$  is substantially smaller, compared to that of belonging to  $\mathbf{Y}$ . Both the FLD and the SVM do not exploit this property. Moreover, in both of these methods, it is assumed that the classes

---

\*HP Israel Science Center, Technion City, Haifa 32000 Israel, elad@hp.technion.ac.il.

†The Interdisciplinary Center, Herzliya, toky@idc.ac.il.

‡HP Israel Science Center, Technion City, Haifa 32000 Israel, renato@hp.technion.ac.il.

are linearly separable. This requirement is not valid in typical pattern detection problems where the *Target* class is surrounded by the *Clutter* class. In order to be able to treat more complex, and unfortunately, more common scenarios, non-linear extensions of these algorithms are required [2, 3]. Such extensions are typically at the expense of much more computationally intensive algorithms.

In this paper we propose a new classifier - the Maximal Rejection Classifier (MRC) which is appropriate for the pattern detection problems. The MRC is a linear classifier that overcomes the above two drawbacks. While maintaining the simplicity of a linear classifier, it can also deal with non linearly separable cases. The only requirement is that the *Clutter* class and the convex hull of the *Target* class are disjoint. We define this property as two convexly-separable classes, which is a much weaker condition compared to linear-separability. In addition, this classifier exploits the property of high *Clutter* probability. Hence, it attempts to give very fast *Clutter* labeling, even if at the expense of slow *Target* labeling. Thus, the entire input signal is classified very fast.

The MRC is an iterative rejection based classification algorithm. The main idea is to apply at each iteration a linear projection followed by a thresholding, similar to the SVM and the FLD. However, as opposed to these two methods, the projection vector and the corresponding thresholds are chosen such that at each iteration the MRC attempts to maximize the number of rejected *Clutter* samples. This means that following the first classification iteration, many of the *Clutter* samples are already classified as such, and discarded from further consideration. The process is continued with the remaining *Clutter* samples, again searching for a linear projection vector and thresholds that maximizes the rejection of *Clutter* points from the remaining set. This process is repeated iteratively until a small number or non of the *Clutter* points remain. The remaining samples at the final stage are considered as *Targets*. The idea of rejection-based classifier was already introduced by [1]. However, in this work we extend the idea by using maximal rejection.

In order to demonstrate the behavior of the MRC, this algorithm is applied to the problem of detecting frontal and vertical faces in images. It is demonstrated that the MRC is a very efficient algorithm, requiring an effective computation of close to two convolutions of the input image per each resolution layer in order to reliably detect faces at all scales and all spatial positions.

## 2 The MRC in Theory

Assume two classes are given in  $\mathfrak{R}^n$ ,  $\mathbf{X}$  (the *Target* class) and  $\mathbf{Y}$  (the *Clutter* class). It is required to discriminate between these two classes, i.e., given a point  $z$  drawn from one of these classes, we would like to be able to label it correctly as either *Target* or *Clutter*. We permit also to label an input as *Unknown* as an intermediate label. One important point, however, is that we know a-priori that for a typical input stream the vast majority of the inputs are *Clutters*, i.e.:

$$P\{\mathbf{X}\} \ll P\{\mathbf{Y}\}. \quad (1)$$

where  $P\{\mathbf{X}\}$  is the a-priori probability that an input signal will be a *Target*, and  $P\{\mathbf{Y}\}$  is defined similarly. Based on this knowledge, we would like the classifier to give a decision as fast as possible (i.e., with as few operations as possible). Thus, *Clutter* labeling should be performed fast, even if at the expense of slow *Target* labeling.

Similar to other linear classifiers [2, 3], we suggest to first project the sample  $z$  onto a vector  $u$ , and label it based on the projected value  $\alpha = u^T z$ . Projecting the *Target* class onto  $u$  results with a Probability Density Function (PDF)  $P\{\alpha|\mathbf{X}\}$ . Similarly  $P\{\alpha|\mathbf{Y}\}$  represents the PDF of the *Clutter* class after projection. Our aim is to minimize the overlap between these two PDF's. Note, that this can be achieved not necessarily by keeping the PDFs far away from each other. What is

really needed is that each member from one class will be as distant as possible from members of the other class. We will define this requirement using the following expected distance between a point  $\alpha_0$  and a distribution  $P\{\alpha\}$ :

$$D(\alpha_0 || P\{\alpha\}) = \int_{\alpha} \frac{(\alpha_0 - \alpha)^2 P\{\alpha\}}{\sigma^2} d\alpha = \frac{(\alpha_0 - \mu)^2 + \sigma^2}{\sigma^2}$$

where  $\sigma$  is the variance of  $P\{\alpha\}$  and  $\mu$  is the mean of  $P\{\alpha\}$ . The division by  $\sigma$  is performed in order to make the distance scale-invariant (or unit-invariant). Using this distance definition, the distance of  $P\{\alpha|\mathbf{Y}\}$  from  $P\{\alpha|\mathbf{X}\}$  can be defined as the expected distance of  $(\alpha|\mathbf{Y})$  from  $P\{\alpha|\mathbf{X}\}$ :

$$\begin{aligned} D(P\{\alpha|\mathbf{Y}\} || P\{\alpha|\mathbf{X}\}) &= \int_{\alpha} D(\alpha || P\{\alpha|\mathbf{X}\}) P\{\alpha|\mathbf{Y}\} d\alpha = \\ &= \int_{\alpha} \frac{(\alpha - \mu_x)^2 + \sigma_x^2}{\sigma_x^2} P\{\alpha|\mathbf{Y}\} d\alpha = \frac{(\mu_y - \mu_x)^2 + \sigma_x^2 + \sigma_y^2}{\sigma_x^2} \end{aligned}$$

where  $[\mu_x, \sigma_x]$  and  $[\mu_y, \sigma_y]$  are the mean-variance pairs of  $P\{\alpha|\mathbf{X}\}$  and  $P\{\alpha|\mathbf{Y}\}$ , respectively. Since we want the two distributions to have as small an overlap as possible, we would like to maximize this distance or minimize the *proximity* between  $P\{\alpha|\mathbf{Y}\}$  and  $P\{\alpha|\mathbf{X}\}$ , which can be defined as the inverse of their mutual distance. Note, that this measure is asymmetric with respect to the two distributions, i.e the proximity defines the closeness of  $P\{\alpha|\mathbf{Y}\}$  to  $P\{\alpha|\mathbf{X}\}$ , but not vice versa. Therefore, we define the overall proximity between the two distributions as follows:

$$Prox(P\{\alpha|\mathbf{Y}\}, P\{\alpha|\mathbf{X}\}) = P\{\mathbf{X}\} \frac{\sigma_y^2}{\sigma_x^2 + \sigma_y^2 + (\mu_y - \mu_x)^2} + P\{\mathbf{Y}\} \frac{\sigma_x^2}{\sigma_x^2 + \sigma_y^2 + (\mu_y - \mu_x)^2}.$$

Compared to the original expression in Equation ??, the minimization of this term with respect to  $u$  is easier. If  $P\{\mathbf{X}\} = P\{\mathbf{Y}\}$ , i.e. if there is an even chance to obtain *Target* or *Clutter* inputs, the proximity becomes similar the cost function minimized by the Fisher Linear Discriminant (FLD)[2]. In our case  $P\{\mathbf{X}\} \ll P\{\mathbf{Y}\}$  (Equation 1), thus, the first term is negligible in Equation 2 and can be omitted. Therefore, the optimal  $u$  should minimize the resulting term:

$$d(u) = \frac{\sigma_x^2}{\sigma_x^2 + \sigma_y^2 + (\mu_y - \mu_x)^2} \quad (2)$$

where  $\sigma_y^2, \sigma_x^2, \mu_y$  and  $\mu_x$  are all a function of the projection vector  $u$ .

Minimization of this expression can be achieved in two ways: Maximizing the distance between the means of the two PDFs, or maximizing  $\sigma_y$  while minimizing  $\sigma_x$ . The second alternative is more common in pattern detection. This means that the projection of *Target* inputs tend to concentrate near a constant value, whereas the *Clutter* inputs will spread with a large variance (see e.g. Fig 2). According to this, for a classification decision, we set two threshold values defining an interval  $C_u = [d_1, d_2]$ , where any value  $\alpha = z^T u$  is compared to. If  $\alpha$  is outside  $C_u$  (which we denote as  $C_c$ ), then the input  $z$  is labeled as Clutter. In the other case where  $\alpha \in C_u$ , the input  $z$  is considered Unknown and will be left for next steps. For the optimal  $u$ , most of the *Clutter* inputs will be projected in  $C_c$ . Subsequently, after projection, many of the *Clutter* inputs are usually classified, whereas *Target* labeling may not be immediately possible. This serves our purpose because for a *Clutter* input, there is a high probability that a decision will be made. Since these inputs are more frequent, this means a faster decision for the vast majority of the inputs.

The method which we suggest follows this scheme: The classifier works in an *iterative* manner, projecting and thresholding with different parameters at each iteration sequentially. Since the classifier is asymmetric, the classification is based on *rejections*; *Clutter* inputs are classified and removed from further consideration while the remaining inputs are kept as suspected *Targets*. The iterations and the *rejection* approaches are both key concepts of the proposed scheme.

### 3 The MRC in Practice

Let us return to Equation 2 and find the optimal projection vector  $u$ . In order to do so, we have to express  $\sigma_y^2$ ,  $\sigma_x^2$ ,  $\mu_y$  and  $\mu_x$  as functions of  $u$ . It is easy to see that for the Target class:

$$\mu_x = u^T \mathbf{M}_x \quad \text{and} \quad \sigma_x^2 = u^T \mathbf{R}_{xx} u \quad (3)$$

where we define:

$$\mathbf{M}_x = \int_z z P\{z|\mathbf{X}\} dz \quad \text{and} \quad \mathbf{R}_{xx} = \int_z (z - \mathbf{M}_x)(z - \mathbf{M}_x)^T P\{z|\mathbf{X}\} dz \quad (4)$$

Similar terms can be calculated for the Clutter class. As can be seen, only the first and second moments of the classes play a role in the choice of the projection vector  $u$ .

In practice we usually do not have the the probabilities  $P\{z|\mathbf{X}\}$ ,  $P\{z|\mathbf{Y}\}$ , and inference on the *Target* or *Clutter* class is achieved through examples. For the Target example-sets  $\hat{\mathbf{X}} = \{x_k\}_{k=1}^{L_x}$  the mean-covariance pair  $(\mathbf{M}_x, \mathbf{R}_{xx})$  are replaced with empirical approximations:

$$\hat{\mathbf{M}}_x = \frac{1}{L_x} \sum_{k=1}^{L_x} x_k \quad \text{and} \quad \hat{\mathbf{R}}_{xx} = \frac{1}{L_x} \sum_{k=1}^{L_x} (x_k - \hat{\mathbf{M}}_x)(x_k - \hat{\mathbf{M}}_x)^T \quad (5)$$

Similar terms are valid for the Clutter class as well. The function we aim to minimize is therefore:

$$d(u) = \frac{u^T \hat{\mathbf{R}}_{xx} u}{u^T \left[ \hat{\mathbf{R}}_{xx} + \hat{\mathbf{R}}_{yy} + (\hat{\mathbf{M}}_y - \hat{\mathbf{M}}_x) (\hat{\mathbf{M}}_y - \hat{\mathbf{M}}_x)^T \right] u} = \frac{u^T A u}{u^T B u} \quad (6)$$

Similarly to [2, ?, ?], it is easy to show that  $u$  that minimizes the above expression satisfies the generalized eigenvalue equation:  $Au = \lambda Bu$ . The optimal  $u$  correspond to the smallest possible  $\lambda$ . Notice that given any solution  $u$  for this equation,  $\beta u$  is also a solution with the same  $\lambda$ . Therefore, without loss of generality, the normalized solution  $\|u\| = 1$  is used.

Input vectors whose projected values are in  $C_u$  are not labeled. For these inputs we apply another steps of classification, where the design of the optimal projection vectors in the next steps are performed according to the remaining examples which was not rejected in previous steps.

### 4 Face Detection Using the MRC

The face detection problem can be specified as the need to detect all instances of faces in a given image, at all spatial positions, all scales, all facial expressions, all poses, of all people, and under all lighting conditions. All these requirements should be met, while having few or no false alarms and miss-detections, and with as fast an algorithm as possible. This description reveals the complexity

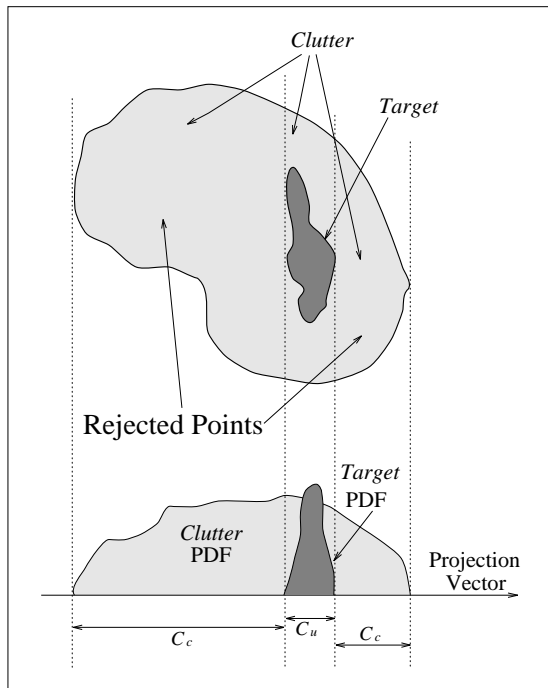


Figure 1: The first rejection stage for a 2D example.

of the detection problem at hand. As opposed to other pattern detection problems, faces are expected to appear with considerable variations, even for detecting only frontal and vertical faces. Variations are expected because of changes in skin color, facial hair, glasses, face shape, and more.

Several papers already addressed the face detection problem using various methods, such as SVM [3, ?], Neural Networks [?, ?, ?], and other methods [?, ?, ?, ?]. In all of these studies, the above complete list of requirements is relaxed in order to obtain practical detection algorithms. Following [?, ?, ?, ?, ?], we deal with the detection of frontal and vertical faces only.

In all these algorithms, spatial position and scale are treated through the same method, in which the given image is decomposed into a Gaussian pyramid with near-unity (e.g., 1.2) resolution ratio. The search for faces is performed in each resolution layer independently, thus enabling the treatment of different scales. In order to be able to detect faces at all spatial positions, fixed sized blocks of pixels are extracted from the image at all positions (with full or partial overlap) for testing. In addition to the pyramid part, which treats varying scales and spatial positions, the core part of the detection algorithm is essentially a classifier which provides a *Face/Non-Face* decision for each input block.

In this paper we propose the application of the MRC for this task. In the face-detection application, *Faces* take the role of targets, and *Non-Faces* are the clutter. In a typical image having millions of pixels, it is expected to detect a few dozens of faces at the most, which means that picking a *Non-Face* block from the image is much more probable. This property is exploited by the MRC in order to obtain an efficient face-detection classifier. The MRC produces very fast *Non-Face* labeling (i.e., with a low computational cost), at the expense of slow *Face* labeling. Thus, on the average, it has a short decision time per input block.

The first stage in the MRC is to gather two example-sets, *Faces* and *Non-Faces*. As mentioned earlier, large enough sets are needed in order to guarantee good generalization for the faces and



the non-faces that may be encountered in images. As to the *Face* set, the ORL data-base<sup>1</sup> was used. This database contains 400 frontal and vertical face images of 40 different individuals. By extracting the face portion from each of these images and scaling to  $15 \times 15$  pixels, we obtained the set  $\hat{\mathbf{X}} = \{x_k\}_{k=1}^{L_x}$  (with  $L_x = 400$ ).

The *Non-Face* set is required to be much larger, in order to represent the variability of *Non-Face* patterns in images. We took 54 arbitrary images containing various textures, natural scenes, graphic images, etc. Common to all these images is that they contain no faces. Each of the 54 images was decomposed into a Gaussian pyramid with a 1.2 resolution ratio, thus creating 1290 images. Using the pyramids is beneficial both for enriching the *Non-Face* set, and for including multi-resolution versions of the same patterns in the data-base. Each possible block of  $15 \times 15$  pixels in these 1290 images is considered as a candidate example of *Non-Face*. Thus, we have effectively collected more than 20 million *Non-Face* examples.

## 5 Results

We trained the MRC for detecting faces by computing 50 sets of kernels  $\{u_k\}_{k=1}^{50}$  and associated thresholds  $\{[T_1^k, T_2^k]\}_{k=1}^{50}$ , using the above described databases of *Faces* and *Non-Faces*. The following figures show the obtained results for several images.

res11.00Face detection with the MRC - Example 1 res20.60Face detection with the MRC - Example 2 res30.70Face detection with the MRC - Example 3 res40.70Face detection with the MRC - Example 4 res50.70Face detection with the MRC - Example 5 res60.70Face detection with the MRC - Example 6

In all these examples, the first stage rejected close to 90% of the candidates. This stage is merely a convolution of the input image (at every scale) with the first kernel,  $u_1$ , followed by thresholding. For these examples, the complete MRC classification required an effective number

<sup>1</sup><http://www.cam-orl.co.uk/facedatabase.html>: ORL database web-site

of close to two convolutions per each pixel in each resolution layer. As can be seen, there are few false alarms, which typically correspond to blocks of pixels having a pattern which may resemble a face. Generally speaking, the algorithm performs very well in terms of detection rate, false alarm rate, and most important of all, computational complexity.

## 6 Conclusion

In this paper we presented a new classifier for target detection, which discriminates between *Target* and *Clutter* classes. The proposed classifier exploits the fact that the probability of a given input to belong to the *Target* class is much lower, as compared to its probability to belong to the *Clutter* class. This assumption, which is valid in many pattern detection applications, is exploited in designing an optimal classifier that detects *Target* signals as fast as possible. Moreover, exact classification is possible when the *Target* and the *Clutter* classes are convexly separable. The Fisher Linear Discriminant (FLD) is a special case of the proposed framework when the *Target* and *Clutter* probabilities are equal. In addition, the proposed scheme overcomes the instabilities arising in the FLD in cases where the mean of the two classes are close to each other. An improvement of the proposed technique is possible by rejecting *Target* patterns instead of *Clutter* patterns in advanced stages, when the probability of *Clutter* is not larger anymore.

The performance of the MRC is demonstrated in the face detection problem. The obtained face detection algorithm is shown to be both computationally very efficient and accurate.

Further details on the theory of the MRC and its application to face detection can be found in [?, ?].

## Acknowledgments

The authors would like to thank Olliveti Research Laboratory (ORL) for making their database available through the net with no restrictions. Thanks are also in order to Dr. HonjZiang Zhang from Microsoft Research Center in China, for his helpful suggestions and fruitful discussions.

## References

- [1] S. Baker and S.K. Nayar. Pattern rejection. In *Proceedings 1996 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - San Francisco, CA, USA, 18-20 June, 1996*.
- [2] Richard O. Duda and Peter E. Hart. *Pattern Classification And Scene Analysis*. Wiley-Interscience Publication, 1973. 1st Edition.
- [3] Valdimir N. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, 1995. 1st Edition.