

Visual features selection using Naïve Bayes Models

Xavier Baró and Jordi Vitrià

Centre de Visió per Computador, Edifici O
Departament de ciències de la Computació, Edifici Q
Campus UAB, Bellaterra, Catalonia (Spain)
E-mail: {xbaro,jordi}@cvc.uab.cat

Abstract This paper describes an object detection approach based on the use of Evolutive Algorithms based on Probability Models (EAPM), a new paradigm in evolutionary computation consisting on iteratively building a probability model to describe the space of promising solutions to a given problem. After the parametrization of the detection problem by means of the weighted dissociated dipoles, a new EAPM based on Naïve Bayes Models estimation is used to find good features.

Keywords: Object Detection, Evolutionary Computation, Probability estimation

1 Introduction

The detection and classification of objects in images that have been acquired in unconstrained environments is a challenging problem because objects can occur under different poses, lighting conditions, backgrounds and clutter. This variation in the object appearance makes unfeasible the design of hand-crafted methods for object detection. Although this problem has been the subject of research from the early beginning of the computer vision field, it has not been until the recent past years that researchers have developed generic object recognition systems for a broad class of real world objects. The key point for this achievement has been the use of a machine learning framework that makes use of very large sets of sample images to learn robust models: Given a training set of n pairs (\mathbf{x}_i, y_i) , where \mathbf{x}_i is the i th image and y_i is the category of the object present in \mathbf{x}_i , we would like to learn a model, $f(\mathbf{x}_i) = y_i$ that maps images to object categories.

One of the most extended approaches to object detection is based on *local methods*, which model an

object as a collection of local visual features or "patches". Thus an image \mathbf{x}_i can be considered to be a vector $(\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,m})$ of m patches. Each patch $\mathbf{x}_{i,j}$ has a feature-vector representation $F(\mathbf{x}_{i,j}) \in \mathbb{R}^d$; this vector might represent various features of the appearance of a patch, as well as features of its relative location and scale. We can choose from a wide variety of features, such as the fragments-based representation approach of Ullman [1], the gradient orientation-based SIFT [2], or some forms of geometric invariant descriptors. One of the most successfully used *a priori* image feature, at least for a broad class of visual objects, is known as Haar feature. These features, which are related to the wavelet decomposition, were originally proposed in the framework of object detection by Viola and Jones [3] in their face detection algorithm.

The two-dimensional Haar decomposition of a square image with n^2 pixels consists of n^2 wavelet coefficients, each of which corresponds to a distinct Haar wavelet. The first such wavelet is the mean pixel intensity value of the whole image; the rest of the wavelets are computed as the difference in mean intensity values of horizontally, vertically, or diagonally adjacent squares. Haar features extend this definition to all possible adjacent rectangles in the window.

Haar features present an interesting property in the context of object recognition: they can be computed very fast using the *integral image*. Integral image at location of x, y contains the sum of the pixel values above and left of x, y inclusive. From that, any Haar feature can be computed from a few values of the integral image.

Viola and Jones [3] used the Adaboost algorithm[4] to learn a real time face detector with a very high classification performance. In their proposal, weak

classifiers are threshold-based classification rules on the values of several Haar features of the image. In order to select the best features, they performed an exhaustive search on the whole set of Haar features for a window. Finally, they built a cascade of strong classifiers to achieve a very low level of false positives.

It is important to note that Haar features constitute an overcomplete dictionary of the image and that there are more than 2^{18} different features for a small image window of 576 pixels (24x24 pixels). This fact imposes a high computational cost on the learning step of the Adaboost algorithm, which involves several rounds of exhaustive searches. From a practical point of view, the development of a high performance object detector represents, when using conventional hardware, a learning time of the order of several hundred hours.

The work of Viola and Jones was extended by Lienhart and Maydt [5], who showed that the use of a larger feature set may improve the convergence and the performance of the final classifier. The extension of the feature set was done by adding rotated versions of original Haar-like features, and thus adding a factor to the exponential relation between the size of the feature set and the training time. Another natural extension of the Haar features is the dissociated dipoles proposed by Balas and Sinha in [6] in the context of computational neuroscience. As in the case of Haar-like features, they are region based comparisons, but in this case we always have only two regions, and these regions do not have to be adjacent regions. The main reason to consider these features is because while Haar-like features are local descriptors, the dissociated dipoles maintain this ability but also includes non-local descriptors of the images. In addition, they Sinha demonstrated that human visual system can perform that type of non-local comparisons, and that these comparisons benefits recognition systems.

The proposal of Balas and Sinha, in spite of the fact that can be easily adapted to Adaboost-based detectors, has not been used in any real world object detector due to computational limitations: there are more than 2^{28} different dissociated dipoles in standard image window of 24×24 pixels, what makes the use of the Adaboost algorithm unfeasible.

Using a parameterizable feature set as Haar-like features or dissociated dipoles, the object detec-

tion problem can be reformulated as an optimization problem, where we must find the best parameters to minimize an error function (i.e. the weighted classification error in the Adaboost algorithm, the AUC or the number of miss-classified samples). The new problem corresponds to search in a large and sparse solution space in order to deal with the best parameters. Nowadays, an emerging field is growing up to deal with this type of problems, the Evolutive Algorithms Based on Probabilistic Models (EAPM).

EAPMs are a new paradigm in the evolutionary computing field. This paradigm was started with the publication of simple algorithm named *Population Based Incremental Learning* (PBIL)[7], based on a simple univariate model, where all the variables are assumed to be independent. The best individuals of each generation are used to update these variables, and finally the model is sampled to obtain a new generation. Although its simplicity, these algorithms demonstrated to converge to good solutions for several problems. Few years later, Schmidh et al. [8] re-introduced the genetic operators to the PBIL algorithm improving significantly its performance. In spite of this return to the origins, the PBIL algorithm introduced an interesting view on evolutionary computation: the extraction of a statistical description of the promising solutions, in terms of a probability distribution is the base of EAPMs and the new systematic way to solve hard search and optimization problems that they represent.

In the literature we can find a wide variety of EAPMs, in which the most important difference is the used probability model. Taking into account the considered interactions between variables, we can classify the models within three main types: Univariate models where no interactions are considered, bivariate models with only pair-wise interactions and finally the models that allow multiple interactions. Once the most convenient probability model is selected, different estimation and sample strategies can be used, thus, we can find different algorithms that share the same type of model. The most known and used algorithms are the UMDA, PBIL and cGA for univariate models, MIMIC, COMIT and BMMA in the case of bivariate models and finally, considering models with multiple interactions the FDA, BOA and EBNA [17].

The use of a complex model allows to better represent the features space, but it adds complexity to the

estimation and sampling stages. In [17] a comparison between different EAPMs over several optimization problems suggests that for simple functions, where there are no interaction between the variables, the performance of the univariate and bivariate models perform as well as more complex models, but when we face more complex problems, a more sophisticated probability model is required. As a general rule, more complex models are more reliable but at the expense of bigger execution times.

2 Object detection formulation

To formulate the problem, in this section we first define the object detection strategy and the used features. At the end, a parametric definition which can be used in optimization approaches is obtained. The object detection used in this paper is based on the Adaboost algorithm and Haar-like features of Viola and Jones in [3], but in our case we do not use combinations of features and threshold to obtain the simple classifiers, only the sign is used instead. This type of features are called ordinal features, which have demonstrated to be robust against various intra-class variations and illumination changes [18]. In addition, we experiment with two larger feature sets than Haar-like features, which cannot be used with the classical Adaboost approach, but can be used with the Adaboost based on EAPMs. These features set are the *Dissociated Dipoles* introduced by Sinha in [6] and the *Weighted Dissociated Dipoles*, an extension of the dissociated dipoles that allows to represent most part of the Haar-like features (see Figure 1).

The main idea of Adaboost is to associate to each

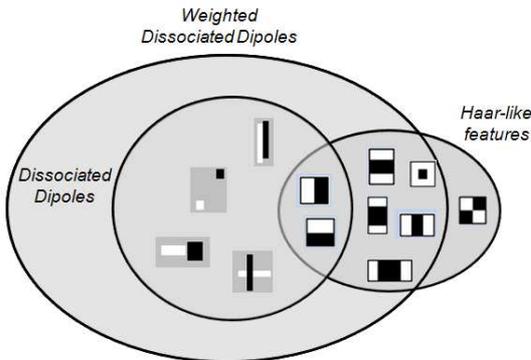


Figure 1: Graphical comparison between the different feature sets.

sample a measure of how many times each sample has been missclassified by the previous added weak learners, and therefore, allowing to add new simple classifiers that concentrate on those samples which are systematically miss classified. At the end, the final classifier is a weighted combination of several simple classifiers. In Adaboost jargon, the simple classifier is named *Weak Classifier* and the final classifier is named *Strong Classifier*. The process of learning the best weak classifier using the weights distribution is named *Weak Learner* and it is the central point in the section.

Given a training set $\langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_M, y_M) \rangle$, where $y_i \in \{-1, +1\}$ is the target value for sample \mathbf{x}_i , the goal of an object detection learning algorithm is to deal with the strong classifier $H(\mathbf{x}_i) = y_i$. In the boosting framework, we define a distribution $W = \{w_1, \dots, w_M\}$ over the training set, where each w_i is the weight associated to the sample \mathbf{x}_i , and $H(\mathbf{x})$ corresponds to an additive model $H(\mathbf{x}) = \sum_t \alpha_t h_t(\mathbf{x})$ where the final decision is a combination of the decisions of several *weak classifiers* $h(\mathbf{x}) \in \{-1, +1\}$. In contrast to the strong classifier $H(\mathbf{x})$ where we expect a good performance for any sample \mathbf{x}_i in the training set, in the case of weak classifier we only expect they are better than a random decision.

Given \mathcal{H} the set of all possible weak classifiers, $h^s \in \mathcal{H}$ a certain weak classifier defined by parameters \mathbf{s} , W the weights distribution of the Adaboost and $\epsilon(h^s) = Pr_{i \sim W}[h^s(\mathbf{x}_i) \neq y_i]$ the error function, the regression step consists on finding \mathbf{s}^* that $\epsilon(h^{\mathbf{s}^*}) \leq \epsilon(h^s) \forall h^s, h^s \in \mathcal{H}$, where the complexity of finding \mathbf{s}^* depends on the size of \mathcal{H} .

Using ordinal measures, Haar-like features can be parameterized by the upper-left position of one of the regions (X, Y) and their size (W, H) , because the size and position of all the regions in Haar-like features is predefined for each type T . Therefore, a Weak Classifier with ordinal Haar-like[5] features can be defined as:

$$h^{\mathbf{s}}(\mathbf{x}) \mapsto \{-1, +1\}$$

$$\text{where } \mathbf{s} = (X, Y, W, H, T) \quad (1)$$

To evaluate this feature, the mean intensity value of negative region (black one) is subtracted from mean value of positive region. The sign of this subtraction is used as the final class value. In the case of dissociated dipoles we can analogously parameterize the two regions as:

$$h^s(\mathbf{x}) \mapsto \{-1, +1\}$$

$$\mathbf{s} = (X_e, Y_e, W_e, H_e, X_i, Y_i, W_i, H_i) \quad (2)$$

To deal with the weighted dissociated dipoles, we just add a weight parameter $RW \in \{1, 2\}$ to each region:

$$h^s(\mathbf{x}) \mapsto \{-1, +1\}$$

$$\mathbf{s} = (X_e, Y_e, W_e, H_e, RW_e, X_i, Y_i, W_i, H_i, RW_i) \quad (3)$$

In addition to the mandatory parameters defined above, we can add extra parameters to improve the results. For instance, we add a polarity parameter which inverts the classification value. This parameter allows a fast step on the search process, inverting the regions. Other parameters can be added to extend the features to multichannel images, but the use of color images is out of the focus of this work.

Once the differences between the feature sets are described, we summarize the problem as finding the best instance of a set of random variables $\mathbf{s} = \{X_1, \dots, X_K\}$, to minimize $\epsilon(h^s) = Pr_{i \sim W}[h^s(\mathbf{x}_i) \neq y_i]$, where the only difference between each feature set is how to evaluate h^s and the dimension of \mathbf{s} .

If we analyze the random variables, it is easy to discover that there are multiple dependencies between them, for instance, regions with positions near the right or bottom sides of the training window cannot have large sizes. The rest of the paper is concentrated on defining an evolutionary approach based on a novel EAPM which allows to solve this optimization problem taking into account these dependencies. From this point, we can forget about the fact we are in an object detection problem and the fact that we use Adaboost, and the problem now is how to find the best parameters of a feature given a weight distribution over the samples.

3 Evolutive Algorithm based on Naïve Bayes Estimation

As we show before, the optimization formulation derived from the object detection problem has multiple dependencies between the variables, and thus, we need to use a model that allows to represent all these interactions. In addition, as the optimization

step must be repeated several times in order to obtain a good detector, we need a fast method. In this section we present a novel EAPM algorithm which accomplish both requirements. First we introduce the probability model, and lastly the evolutive algorithm. The most extended probabilistic model to capture multiple dependencies is the Bayesian network (or a belief network), represented as a probabilistic graphical model, specifically an acyclic directed graph, where each node corresponds to a variable (measured parameter, latent variable or hypothesis) and whose arcs encode the dependence between variables.

Learning a Bayesian network from data is a two-fold problem: Structure learning and parameter estimation. Although there exist good methods to estimate the structure and parameters of a Bayesian network, because exact inference is $\#P$ -complete and thus the existent methods are often too costly, approximate methods like Markov Chain Monte Carlo [19] and loopy belief propagation [20] must be used. The applicability of the Bayesian networks is limited by the fact that these methods have an unpredictable inference time and its convergence is difficult to diagnose. An alternative to the Bayesian networks are the Naïve Bayes models, where with the "naive" assumption that all variables are mutually independent given a "special" variable C the model is simplified. Although the resulting model is very simple (see Fig. 2), in [21] Lowd and Domingos demonstrated from an empirical point of view that this simpler model has the same representation power than a Bayesian network. In addition to several experiments comparing both models, in this work they also propose the Naïve Bayes Estimation (NBE) algorithm to efficiently estimate Naïve Bayes Models from data. This algorithm consists on an Expectation Maximization (EM) wrapped in an outer loop that progressively adds and prunes mixture components.

Once the probability model is defined, and we have

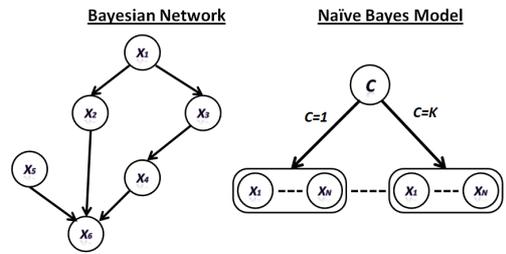


Figure 2: Probabilistic models.

the NBE algorithm to estimate this model from a given data, the final algorithm is presented in Alg. 1. The algorithm is based on iteratively estimating and sampling the model to find promising regions on the features space.

Algorithm 1 EANBE learning algorithm

Require: Initial population P_0 , evaluation function $\mathcal{E} : \mathbb{R}^D \mapsto \mathbb{R}$ (where D is the dimension of each individual $I \in P_0$ and L the number of individuals in P_0), the selection percentage P_S and the hold-out percentage P_H .
Ensure: The best individual I_{best} among all evaluated individuals.

$L_S \leftarrow \text{round}(P_S \times L)$

$L_H \leftarrow \text{round}(P_H \times L_S)$

repeat

Evaluate the population P_i using the evaluation function \mathcal{E}

Select the best individual I_i in the population P_i

Select the L_S best individuals using the evaluation value

if $\mathcal{E}(I_i)$ is better than $\mathcal{E}(I_{Best})$ or it is the first iteration **then**

$I_{Best} \leftarrow I_i$

end if

Divide the selected individuals into two random subsets T and H with $L_S - L_H$ and L_H individuals respectively.

Use the NBE algorithm to learn a model M using T as the training set and H as the hold-out set.

Sample the model M to obtain a new population P_{i+1} with L individuals.

until The maximum number of generations, the optimum value or some other stopping criteria are achieved

return I_{Best} and $\mathcal{E}(I_{Best})$.

4 Results

To validate the proposed object detection method, we have developed the *caplib*, a C++ based library that implements a general framework to work with EAPMs. The learning process is done by a Gentle Adaboost version in Matlab, using MEX files to communicate with the *caplib*. All the source codes for the library, MEX files and Matlab files and the 10-fold sets are available in an online appendix at www.cvc.uab.cat/~xbaro/eanbe/.

In the first part of this section we describe the data sets used to perform the experiments, and after we describe the experiments. Finally, we analyze the obtained results.

4.1 Data Sets

In order to verify the usefulness of the presented methodology we have selected different real world object detection problems, using for each one a public database to facilitate the reproduction and comparison of the results. In some cases, the databases contain both, the object image and background images to use as negative samples, and when the background images are not given, we create the negative samples from the *Corel Photo Libraries*, cutting random regions at different position and scales. A detailed description of each data set is given in the following.

- *Faces:* We use the MIT-CBCL face database with a random selection of 1.000 face images and 3.000 non-face images. All the images correspond to frontal faces with several illumination changes.
- *Text:* We use the text location dataset from the *7th International Conference on Document Analysis and Recognition (ICDAR03)*. The images correspond to text regions over a wide set of surfaces and with several illumination changes. To obtain the training data, each text region has been split into overlapped subregions with the same size.
- *Cars:* We use the UIUC cars database, with a total of 1.050 images containing 550 instances of lateral views of different cars in urban scenes and 500 images of background.
- *Pedestrians:* We use the INRIA Person Dataset, with 2.924 images divided into 924 pedestrian instances and 2.000 background images.
- *Traffic Signs:* In this case we use real images acquired in the context of a mobile mapping project provided by the Institut Cartogràfic de Catalunya. The database consists on 1.000 images containing a traffic sign and 3.000 background images.

4.2 Measures

To compare the performance of each configuration of the system, the classification Balanced Error (BER) is used. This measure takes into account the different number of examples in each of the classes in order

to face unbalanced problems. Given a samples set (\mathbf{x}_i, y_i) with $y_i \in \{-1, 1\}$ and a certain hypothesis $h(\mathbf{x}) \mapsto \{-1, 1\}$ the balanced error is defined as:

$$BER = \frac{FN}{N^+} + \frac{FP}{N^-} \quad (4)$$

where N^+ and N^- are the number of positive and negative samples respectively, and FP and FN the number of false positive and false negative samples.

4.3 Experiments

To validate the proposed schema and compare each selection of feature set and evolutive algorithm, we perform a stratified 10-fold cross-validation with a confidence interval at 95%. The population size for genetic algorithm and EANBE is set to 100 individuals, using $T = 200$ as the number of iterations in the Gentle Adaboost algorithm. In the case of genetic algorithm, the individuals are encoded using binary gray code, while in the case of EANBE, the decimal code is used directly. The BER mean values and confidence intervals obtained in the experiments are shown in the following table:

Table 1: Results obtained in the experiments.

Data Set	Dissociated Dipoles	Haar-like features	Weighted Diss. Dip.
Cars	78.59% \pm 7.50	70.62% \pm 4.98	94.64% \pm 4.81
Faces	69.40% \pm 5.09	55.43% \pm 4.93	88.27% \pm 5.25
Text	50.48% \pm 5.02	47.92% \pm 4.85	73.27% \pm 3.36
Pedestrians	70.37% \pm 3.72	55.91% \pm 8.24	90.01% \pm 1.89
Traffic Signs	63.06% \pm 3.91	65.37% \pm 3.92	89.67% \pm 2.51

5 Conclusions and Future Work

Although the tests must be extended to a larger set of databases, the results obtained with weighted dissociated dipoles in combination to the new evolutive algorithm encourage to continue with this line. As a future work we want to extend this approach to multiclass learning and use a priori knowledge from the learning data in order to speed-up the selection process and obtain not only good features but meaningful features.

Acknowledgments

This work has been developed in a project in collaboration with the *Institut Cartogràfic de Catalunya* under the supervision of Maria Pla. This work is supported by MEC grant TIN2006-15308-C02-01, Ministerio de Ciencia y Tecnología, Spain.

References

- [1] S. Ullman and E. Sali, "Object classification using a fragment-based representation," in *BMVC '00: Proceedings of the First IEEE International Workshop on Biologically Motivated Computer Vision*. London, UK: Springer-Verlag, 2000, pp. 73–87.
- [2] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proc. of the International Conference on Computer Vision ICCV, Corfu, 1999*, pp. 1150–1157. [Online]. Available: citeseer.ist.psu.edu/lowe99object.html
- [3] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 1, 2001, pp. I–511–I–518.
- [4] Y. Freund and R. E. Schapire, "Experiments with a new boosting algorithm," in *International Conference on Machine Learning*, 1996, pp. 148–156.
- [5] R. Lienhart and J. Maydt, "An extended set of haar-like features for rapid object detection," in *Proceedings of the International Conference on Image Processing*. Rochester, USA: IEEE, September 2002, pp. 900–903.
- [6] B. Balas and P. Sinha, "Dissociated dipoles: Image representation via non-local comparisons," Annual meeting of the Vision Sciences Society, Sarasota, FL., 2003.
- [7] S. Baluja and R. Caruana, "Removing the genetics from the standard genetic algorithm," in *The Int. Conf. on Machine Learning 1995*, A. Prieditis and S. Russel, Eds. San Mateo, CA: Morgan Kaufmann Publishers, 1995, pp. 38–46.
- [8] M. Schmidt, K. Kristensen, and T. R. Jensen, "Adding genetics to the standard PBIL algorithm," in *Proceedings of the 1999 Congress on Evolutionary Computation*, vol. 2, 1999, pp. 1527–1534.
- [9] H. Muhlenbein, "The equation for response to selection and its use for prediction," *Evolutionary Computation*, vol. 5, no. 3, pp. 303–346, 1997.
- [10] G. R. Harik, F. G. Lobo, and D. E. Goldberg, "The compact genetic algorithm," *IEEE-EC*, vol. 3, no. 4, p. 287, November 1999.
- [11] J. S. de Bonet, C. L. Isbell, Jr., and P. Viola, "MIMIC: Finding optima by estimating probability densities," in *Advances in Neural Information Processing Systems*, M. C. Mozer, M. I. Jordan, and T. Petsche, Eds., vol. 9. The MIT Press, 1997, p. 424.

- [12] S. Baluja and S. Davies, "Using optimal dependency-trees for combinatorial optimization," in *ICML '97: Proceedings of the Fourteenth International Conference on Machine Learning*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1997, pp. 30–38.
- [13] M. Pelikan and H. Mühlenbein, "The bivariate marginal distribution algorithm," in *Advances in Soft Computing - Engineering Design and Manufacturing*, R. Roy, T. Furuhashi, and P. K. Chawdhry, Eds. London: Springer-Verlag, 1999, pp. 521–535.
- [14] H. Mühlenbein and T. Mahning, "The factorized distribution algorithm for additively decomposed functions," in *Second Symposium on Artificial Intelligence. Adaptive Systems, CIMAFA 99. La Habana*, 1999, pp. 301–313.
- [15] M. Pelikan, D. E. Goldberg, and E. Cantú-Paz, "BOA: The Bayesian optimization algorithm," in *Proceedings of the Genetic and Evolutionary Computation Conference GECCO-99*, W. Banzhaf, J. Daida, A. E. Eiben, M. H. Garzon, V. Honavar, M. Jakiela, and R. E. Smith, Eds., vol. I. Orlando, FL: Morgan Kaufmann Publishers, San Francisco, CA, 13-17 1999, pp. 525–532.
- [16] R. Etxeberria and P. L. naga, "Global optimization with bayesian networks." in *In Second Symposium on Artificial Intelligence(CIMAFA-99), Cuba*, 1999, pp. 332–339.
- [17] R. Blanco and J. Lozano, *An Empirical Comparison of Discrete Estimation of Distribution Algorithms*. Kluwer Academic Publishers, 2001, ch. chapter 7, pp. 167–180.
- [18] K. Thoresz and P. Sinha, "Qualitative representations for recognition," *Journal of Vision*, vol. 1, no. 3, pp. 298–298, 12 2001. [Online]. Available: <http://journalofvision.org/1/3/298/>
- [19] W. Gilks, S. Richardson, and D. J. Spiegelhalter, *Markov Chain Monte Carlo in Practice*, ser. Interdisciplinary Statistics Series. CRC Press, 1996.
- [20] J. S. Yedidia, W. T. Freeman, and Y. Weiss, "Generalized belief propagation," in *NIPS*, pp. 689–695.
- [21] D. Lowd and P. Domingos, "Naive bayes models for probability estimation," in *ICML '05: Proceedings of the 22nd international conference on Machine learning*. New York, NY, USA: ACM Press, 2005, pp. 529–536.
- [22] "Corel stock photo library." [CD-ROM], 1996.
- [23] "MIT-CBCL face database." [Online]. Available: cbcl.mit.edu/projects/cbcl/software-datasets/FaceData1Readme.html
- [24] S. Agarwal, A. Awan, and D. Roth, "UIUC cars database." [Online]. Available: l2r.cs.uiuc.edu/~cogcomp/Data/Car/