

ji.

# Probablistic Model-Based Detection of Bent-Double Radio Galaxies

S. Kirshner, I. Cadez, P. Smyth, C. Kamath, and E. Cantu-Paz

This article was submitted to 16<sup>th</sup> International Conference on Pattern Recognition, Quebec City, Canada, August 11-15, 2002

### U.S. Department of Energy



# April 23, 2001

Approved for public release; further dissemination unlimited

#### DISCLAIMER

This document was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor the University of California nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or the University of California, and shall not be used for advertising or product endorsement purposes.

This is a preprint of a paper intended for publication in a journal or proceedings. Since changes may be made before publication, this preprint is made available with the understanding that it will not be cited or reproduced without the permission of the author.

This report has been reproduced directly from the best available copy.

Available electronically at http://www.doc.gov/bridge

Available for a processing fee to U.S. Department of Energy And its contractors in paper from U.S. Department of Energy Office of Scientific and Technical Information P.O. Box 62 Oak Ridge, TN 37831-0062 Telephone: (865) 576-8401 Facsimile: (865) 576-5728 E-mail: reports@adonis.osti.gov

Available for the sale to the public from U.S. Department of Commerce National Technical Information Service 5285 Port Royal Road Springfield, VA 22161 Telephone: (800) 553-6847 Facsimile: (703) 605-6900 E-mail: <u>orders@ntis.fedworld.gov</u> Online ordering: <u>http://www.ntis.gov/ordering.htm</u>

OR

Lawrence Livermore National Laboratory Technical Information Department's Digital Library http://www.llnl.gov/tid/Library.html

### **Probabilistic Model-Based Detection of Bent-Double Radio Galaxies**

Sergey Kirshner\* skirshne@ics.uci.edu Igor V. Cadez\* icadez@ics.uci.edu Padhraic Smyth\* smyth@ics.uci.edu Chandrika Kamath<sup>‡</sup> kamath2@llnl.gov

Erick Cantú-Paz<sup>‡</sup> cantupaz@llnl.gov

### Abstract

## 2. Data

We describe an application of probabilistic modeling to the problem of recognizing radio galaxies with a bentdouble morphology. The type of galaxies in question contain distinctive signatures of geometric shape and flux density that can be used to be build a probabilistic model that is then used to score potential galaxy configurations. The experimental results suggest that even relatively simple probabilistic models can be useful in identifying galaxies of interest in an automatic manner.

### **1. Introduction**

In this paper we investigate the problem of identifying bent-double radio galaxies in the FIRST (Faint Images of the Radio Sky at Twenty-cm) Survey data set [1]. FIRST produces large numbers of radio images of the deep sky using the Very Large Array at the National Radio Astronomy Observatory. It is scheduled to cover more that 10,000 square degrees of the northern and southern caps (skies). Of particular scientific interest to astronomers is the identification and cataloging of sky objects with a "bent-double" morphology, indicating clusters of galaxies. (For an example, see Figure 1.) Due to the very large number of observed deep-sky radio sources, it is infeasible for the astronomers to manually label all of them by hand.

The data from the FIRST Survey is available in In the "raw image" format, two different formats. image cut-outs are available from the FIRST website (http://sundog.stsci.edu/). The second data format is in the form of a catalog of features that have been automatically derived from the raw images by an image analysis program [5]. Each entry corresponds to a single detectable "blob" of bright intensity (a sky object) relative to the sky background: these entries are called components. The "blob" of intensities for each component is fitted with an ellipse (details in [5]). The ellipses and intensities for each ellipse are described by a set of estimated features such as sky position of the centers, (RA (right ascension) and Dec (declination)), peak density flux and integrated flux, RMS noise, lengths of the major and minor axes, and the position angle of the major axis of the ellipse counterclockwise from the north. The goal is to find sets of components that are spatially close and that resemble a bent-double. In the results in this paper we focus on the classification of candidate sets of components that have been detected by an existing spatial clustering algorithm [3] where each set consists of three components from the catalog (three ellipses). As of 2000, the catalog contained over 15,000 three-component configurations and over 600,000 configurations total. Three-component bentdouble configurations typically consist of a center or "core" component and two other side components called "lobes".

The set which we use to build and evaluate our models consists of a total of 128 examples of bent-double galaxies and 22 examples of non-bent-double configurations. A configuration is labeled as a bent-double if at least two astronomers labeled it as such. Note that the visual identification process is the bottleneck in the process since it requires significant time and effort from the scientists, and is subjective and error-prone. This motivates the creation of automated methods for identifying bent-doubles. This data set is also considerably biased towards the bent-double class (i.e.,

<sup>\*</sup>Department of Information and Computer Science, University of California, Irvine, CA 92697-3425, U.S.A

<sup>&</sup>lt;sup>†</sup>Center for Applied Scientific Computing, Lawrence Livermore National Laboratory, Livermore, CA 94551, U.S.A

<sup>&</sup>lt;sup>‡</sup>The work of Chandrika Kamath and Erick Cantú-Paz was performed under the auspices of the U.S. Department of Energy by University of California Lawrence Livermore National Laboratory under contract No. W-7405-Eng-48.



Figure 1. An example of a bent-double (left) and non-bent-double (right) galaxies. Notice that the configuration on the right does not have enough "bend".

bent-doubles are far more prevalent in this training data set than they are in the catalog in general). This is an artifact of the manner in which scientists generated a labeled data set. However, since we use a likelihood-based approach for ranking candidate objects, where a model is built only on positive examples (bent-doubles), the training methodology presented below is not sensitive to such an imbalance in the training data.

Previous work on automated classification of threecomponent candidate sets has focused on the use of decision-tree classifiers using a variety of geometric and image intensity features [2] [3]. A limitation of the decisiontree approach is its relative inflexibility in handling uncertainty about the object being classified, e.g., the identification of which of the three components should be treated as the core of a candidate object. A primary motivation for the development of a probabilistic approach is to provide a framework that can handle such uncertainties in a coherent manner. In particular, in this paper, we focus on a probabilistic mixture model that treats the identification of the center component as a hidden variable, providing a natural framework for handling this uncertainty both in the modelbuilding phase (on training data) and in the detection phase (on test data).

# 3. Probabilistic Modeling of Bent-Double Galaxies

We denote a three-component **configuration** by  $C = (\mathbf{c}_1, \mathbf{c}_2, \mathbf{c}_3)$ , where the  $\mathbf{c}_i$ 's are the elliptical components as described in the previous section. Each component  $\mathbf{c}_x$  is represented as a feature vector, where the specific features will be defined later. Our approach focuses on building a probabilistic model for bent-doubles: p(C) =

 $p(\mathbf{c}_1, \mathbf{c}_2, \mathbf{c}_3)$ , the likelihood of the observed  $\mathbf{c}_i$  under a bent-double model where we implicitely condition on "bent-double". Our general approach is to define this likelihood, then estimate its parameters from training data, and use it to rank candidate configurations.

#### 3.1. Modeling Orientation

By looking at examples of bent-double galaxies and by talking to the scientists studying them, we have been able to establish a number of potentially useful characteristics of the components, the primary one being geometric symmetry. In bent-doubles, two of the components will look close to being mirror images of one another with respect to a line through the third component. We will call mirrorimage components lobes components, and the other one the core component. It also appears that non-bent-doubles either don't exhibit such symmetry, or the angle formed at the core component is too straight-the configuration is not "bent" enough. Once the core component is identified, we can calculate symmetry-based features. However, identifying the most plausible core component requires either an additional algorithm or human expertise. In our approach we use a probabilistic framework that averages over different possible orientations weighted by their likelihood.

To formalize the estimation of the core and the lobes, consider the following. Without loss of generality assign the numbers 1, 2, 3 to the components. In general we do not know which of 1, 2, or 3 is the core (under a bent-double assumption). By an **orientation** we mean a mapping of vertices to a set of labels  $\{a, b, c\}$  which preserves the neighbor relation in a cyclical order. Figure 2 shows an example of elliptical representation with possible orientations. For the set of three vertices, all 6 mappings preserve the neighbor relation. (In general, for configurations of *n* components, there will be 2n such mappings.)

The mapping from components 1, 2, 3 to a, b, c is defined by orientation  $\theta_i$ . We can then write

$$p(\mathbf{c}_a, \mathbf{c}_b, \mathbf{c}_c) = \sum_{i=1}^{6} p(\mathbf{c}_a, \mathbf{c}_b, \mathbf{c}_c | \theta_i) p(\theta_i), \quad (1)$$

i.e., a mixture over all possible orientations. Each orientation is assumed a priori to be equally likely, i.e.,  $p(\theta_i) = \frac{1}{6}$ . Intuitively for a configuration that clearly looks like a bent-double, the terms in the mixture corresponding to the correct core component would dominate, while the other core interpretations would have much lower likelihood.

We represent each component  $c_x$  by three features. Note that the features can only be calculated conditioned on



component 1

Figure 2. Elliptical components of a hypothetical bent-double. Assuming that label *a* would correspond to a core component, a good choice of orientations would be  $\{1 \rightarrow b, 2 \rightarrow a, 3 \rightarrow c\}$  or  $\{1 \rightarrow c, 2 \rightarrow a, 3 \rightarrow b\}$ .

a particular mapping since they rely on properties of the (assumed) core and lobe components. Thus, conditioned on a particular mapping or orientation  $\theta$ , assuming label  $x \in \{a, b, c\}$  where a, b, c are defined in a cyclical order, the features are defined as:

Angle α<sub>x,θ</sub>—the angle formed at the center of the component (a vertex of the configuration) mapped to label x:

• Side ratios 
$$sr_{x,\theta} = \frac{|\text{center of } x \text{ to center of } next(x)|}{|\text{center of } x \text{ to center of } prev(x)|}$$
;

• Intensity ratios 
$$ir_{x,\theta} = \frac{\text{peak flux of } next(x)}{\text{peak flux of } prev(x)}$$

and so  $\mathbf{c}_x|\theta = (\alpha_{x,\theta}, sr_{x,\theta}, ir_{x,\theta})$ . Other features are also possible. Nonetheless this particular set appears to capture the more obvious visual properties of bent-doubles.

Rather than modeling the full joint distribution of all features, we make some approximating conditional independence assumptions (motivated by the relatively small amount of training data). In particular, we assume that

$$P((\mathbf{c}_{a}, \mathbf{c}_{b}, \mathbf{c}_{c}) | \theta) = P(\alpha_{a,\theta}, \alpha_{b,\theta}, \alpha_{c,\theta}) P(sr_{a,\theta}, sr_{b,\theta}, sr_{c,\theta}) \times P(ir_{a,\theta}, ir_{b,\theta}, ir_{c,\theta}).$$

For all ratio features r (either of sr, ir),  $r_{a,\theta} \cdot r_{b,\theta} \cdot r_{c,\theta} = 1$ . For the angle features,  $\alpha_{a,\theta} + \alpha_{b,\theta} + \alpha_{c,\theta} = \pi$ . Assume that label a corresponds to the choice of the core component. If we further assume a conditional independence for the features of any two components we can obtain further simplifications:

$$P(\alpha_{a,\theta}, \alpha_{b,\theta}, \alpha_{c,\theta}) = P(\alpha_{a,\theta}) P(\alpha_{b,\theta} | \alpha_{a,\theta}) P(\alpha_{c,\theta} | \alpha_{a,\theta}, \alpha_{b,\theta})$$
  

$$= P(\alpha_{a,\theta}) P(\alpha_{b,\theta}) ;$$
  

$$P(r_{a,\theta}, r_{b,\theta}, r_{c,\theta})$$
  

$$= P(r_{a,\theta}) P(r_{b,\theta} | r_{a,\theta}) P(r_{c,\theta} | r_{a,\theta}, r_{b,\theta})$$
  

$$= P(r_{a,\theta}) P(r_{b,\theta}) .$$

Given  $\theta$ , let  $P_a(\alpha) = P(\alpha_{a,\theta})$ ,  $P_a(r) = P(r_{a,\theta})$ , and let  $P_b(\alpha) = P(\alpha_{b,\theta})$ ,  $P_b(r) = P(r_{b,\theta})$ . If we know for every training example which component is the core (and is mapped to label *a*) we can then unambiguously estimate each of these distributions, e.g., by using kernel-density estimators. In practice, however, the identity of the core component is unknown.

We use our model to estimate which components are likely to be cores, using the following iterative scheme. Initially, core components for the bent-double examples in the training set are chosen at random. At each step of the iteration, we build the corresponding  $P_a$  and  $P_b$  distributions (using kernel density estimators) from the training set using the currently estimated orientations (and labels a). The estimated  $P_a$  and  $P_b$  distributions are then used on all of the examples in the training set to calculate the probability of each component being a core. This is done by summing  $P(\mathbf{c}_a, \mathbf{c}_b, \mathbf{c}_c | \theta_i)$  in Equation 1 over the 2 (out of 6 possible) orientations  $\theta_i$  that map that component to label a. The most likely core components for each example are chosen to be the cores for the next iteration (in effect this is an approximation to a full expectation-maximization procedure, where the most likely core component is chosen rather than averaging over core components). The likelihood (probability of the training set under the currently estimated distributions) is recorded at each iteration. The algorithm stops either after a prespecified maximum number of iterations or when there are no changes from one iteration to the next.

This procedure yields estimates of the  $P_a$  and  $P_b$ distributions for each feature, allowing calculation of  $P(\mathbf{c}_a, \mathbf{c}_b, \mathbf{c}_c | \theta_i)$  for any particular orientation  $\theta_i$ . Thus, for a new unlabeled example we can now calculate a full likelihood  $P(\mathbf{c}_a, \mathbf{c}_b, \mathbf{c}_c)$  (Equation 1), i.e. we average over all 6 possible orientations. For a set of unlabeled examples this yields a set of likelihood scores under the bentdouble model, which can be sorted and thresholded to yield a receiver-operating characteristic. If the likelihood of the data under a non-bent-double model is assumed to be roughly uniform in feature-space, then these likelihoods will be roughly monotonically proportional to the posterior probability of a bent-double given the observed data. Here we choose not to build an explicit model of non-bentdoubles given that they can exhibit considerable variation, and instead rely on a model only of the positive examples for detection.

### **3.2. Kernel Densities**

The choice of kernels is discussed in detail in [4].

### 4. Experimental Results

For our experiments, to score a positive example, we use all positive examples to build the model except for the one being scored-leave-one-out cross-validation. All negative examples are scored by the model trained on the set of all 128 positive examples. The examples are then sorted in decreasing order by their likelihood score and analyzed using receiver operating characteristics (ROC curves). If the two classes can be perfectly separated by these scores, i.e. scores of all negative examples would appear after scores of all positive examples, then the curve would coincide with the left and upper sides of the  $[0,1] \times [0,1]$  square. If the scores are chosen at random, then ROC curve would approach y = x line. We use  $A_{ROC}$ , the area above the curve, as the measure of goodness of the model. A random score assignment would yield  $A_{ROC} = 0.5$  while perfect assignment would have  $A_{ROC} = 0$ .

We have tried a number of different choices of bandwidth for kernel density estimators for the features, and the results appear relatively insensitive to the particular bandwidths chosen. One set of bandwidths resulted in the plot shown in Figure 3. From the plot we can infer, among other things, that the highest score for a negative example appears after scores of 74% or 95 out of 128 positive examples. We can also find that 55% or 12 out of 22 negative examples are among the lowest 14 scores. Thus, the model appears to be quite useful at detecting bent-double galaxies.

### 5. Conclusions

A probabilistic model for the identification of bentdouble galaxies appears quite useful for detection, based on cross-validation results. A general mixture model framework allows for a principled and effective approach to orientation estimation. The experimental results are accurate enough to suggest that the technique may be quite useful for automated identification of likely bent-double candidates from very large astronomy catalogs. Future work includes a comparison of this method with decision trees.



Figure 3. ROC curve plot for model using angle, ratio of sides, and ratio of intensities, as features,  $A_{ROC} = 0.069602$ 

### 6. Acknowledgements

This work was performed under a sub-contract from the ASCI Scientific Data Management Project of the Lawrence Livermore National Laboratory. We gratefully acknowledge our FIRST collaborators, in particular, Robert H. Becker for sharing his expertise on the subject.

### References

- R. H. Becker, R. L. White, and D. J. Helfand. The FIRST Survey: Faint Images of the Radio Sky at Twenty-cm. Astrophysical Journal, 450:559, 1997.
- [2] E. Cantú-Paz and C. Kamath. Combining evolutionary algorithms with oblique decision trees to detect bent-double galaxies. In *Proceedings of International Symposium on Optical Science and Technology (SPIE Annual meeting)*, San Diego, July 30-August 4 2000.
- [3] I. K. Fodor, E. Cantú-Paz, C. Kamath, and N. A. Tang. Finding bent-double radio galaxies: A case study in data mining. In *Proceedings of the Interface: Computer Science and Statistics Symposium*, volume 33, 2000.
- [4] S. Kirshner, I. V. Cadez, P. Smyth, C. Kamath, and E. Cantú-Paz. Probabilistic model-based detection of bent-double radio galaxies. Technical Report 02-14, Department of Information and Computer Science, University of California, Irvine.
- [5] R. L. White, R. H. Becker, D. J. Helfand, and M. D. Gregg. A catalog of 1.4 GHz radio sources from the FIRST Survey. *Astrophysical Journal*, 475:479, 1997.