

Learning to Classify Galaxy Shapes using the EM Algorithm

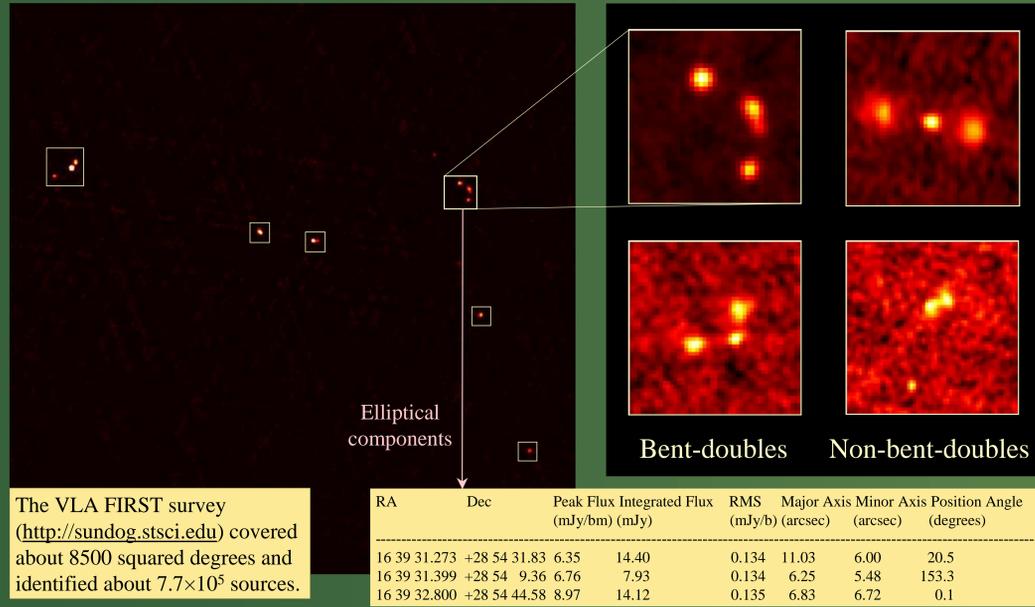
Sergey Kirshner¹, Igor V. Cadez², Padhraic Smyth¹, Chandrika Kamath³

¹ Information and Computer Science,
University of California,
Irvine, CA 92697-3425
skirshne@ics.uci.edu, smyth@ics.uci.edu

² Sparta Inc.,
23382 Mill Creek Drive #100,
Laguna Hills, CA 92653
icadez@ics.uci.edu

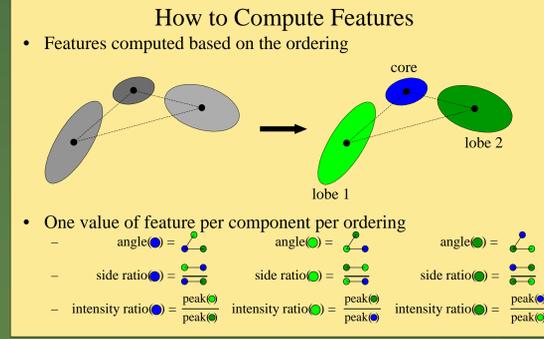
³ Center for Applied Scientific Computing,
Lawrence Livermore National Laboratory,
Livermore, CA 94551
kamath2@llnl.gov

We describe an application of probabilistic learning to the problem of automatic identification of astronomical radio sources with a bent-double morphology. Calculation of object features requires identifying an ordering on the “corners” of the object of interest. We describe an approach for automatic ordering and calculation of features and demonstrate how classes of objects can be learned in both a supervised and unsupervised manner.



- Set labeled by scientists
 - 128 bent-doubles
 - 22 non-bent-doubles
- Objective: to classify new configurations

- ### Orderings and Features
- Build probabilistic model for bent-doubles
 - Features based on symmetry
 - Need to order components to calculate features
 - Compute features given the ordering:
 - angles
 - side ratios
 - peak intensity ratios
 - Assume independence of features given the ordering
 - Assume pairwise independence for each feature
 - Incorporate constraints for features
 - angles add up to π
 - product of ratios is 1



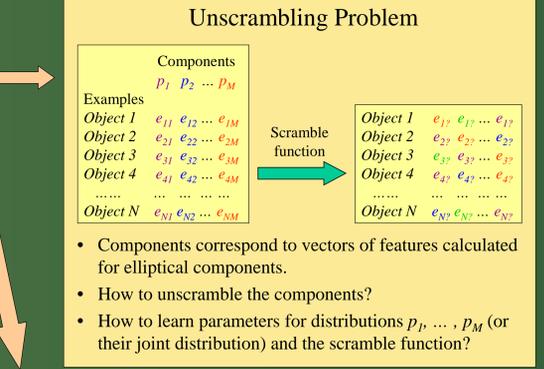
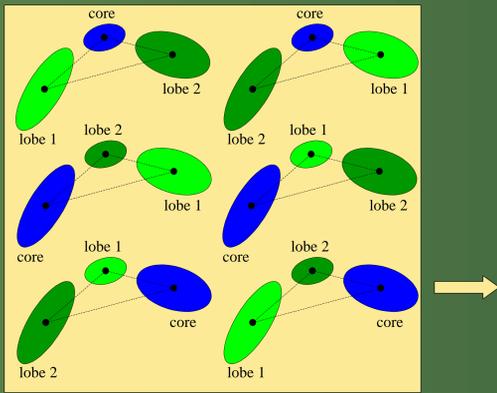
$$P(\text{features}) = P(\text{angles}) \times P(\text{side ratios}) \times P(\text{intensity ratios})$$

$$P(\text{angles}) = P(\text{angle}_{\text{core}}, \text{angle}_{\text{lobe 1}}, \text{angle}_{\text{lobe 2}}) = P(\text{angle}_{\text{core}} | \text{angle}_{\text{lobe 1}}, \text{angle}_{\text{lobe 2}}) \times P(\text{angle}_{\text{lobe 1}} | \text{angle}_{\text{lobe 2}}) \times P(\text{angle}_{\text{lobe 2}}) = P(\text{angle}_{\text{core}}) \times P(\text{angle}_{\text{lobe 1}}) \times P(\text{angle}_{\text{lobe 2}})$$

$$P(\text{side ratios}) = P(\text{side ratio}_{\text{core}}, \text{side ratio}_{\text{lobe 1}}, \text{side ratio}_{\text{lobe 2}}) = P(\text{side ratio}_{\text{core}}) \times P(\text{side ratio}_{\text{lobe 1}}) \times P(\text{side ratio}_{\text{lobe 2}})$$

$$P(\text{intensity ratios}) = P(\text{intensity ratio}_{\text{core}}, \text{intensity ratio}_{\text{lobe 1}}, \text{intensity ratio}_{\text{lobe 2}}) = P(\text{intensity ratio}_{\text{core}}) \times P(\text{intensity ratio}_{\text{lobe 1}}) \times P(\text{intensity ratio}_{\text{lobe 2}})$$

- ### Two Problems
- There are six possible orderings – which one is best? Proper orderings for bent-doubles in the data may not be known.
 - Non-bent-doubles do not have a correct ordering since they may not exhibit any symmetry that the features are designed to capture. Need to come up with a scheme to be able to assign a probability score to all configurations.

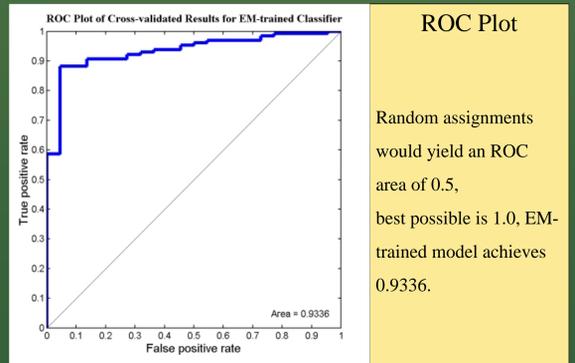
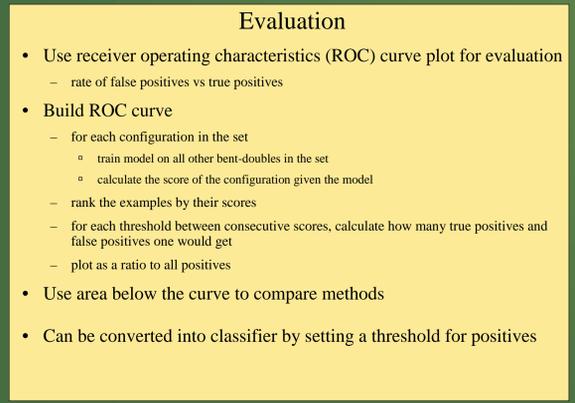


Mixture of Orderings

- Reformulate probability of a configuration as a mixture over possible orderings

$$P(\text{features}) = P(\text{features} | \text{ordering 1}) \times P(\text{ordering 1}) + \dots + P(\text{features} | \text{ordering N}) \times P(\text{ordering N})$$

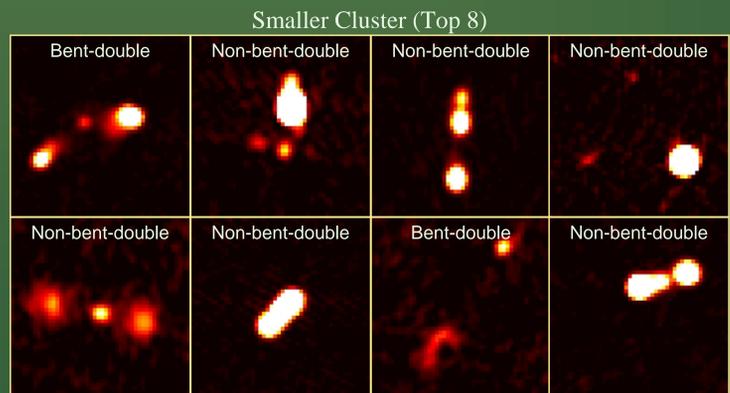
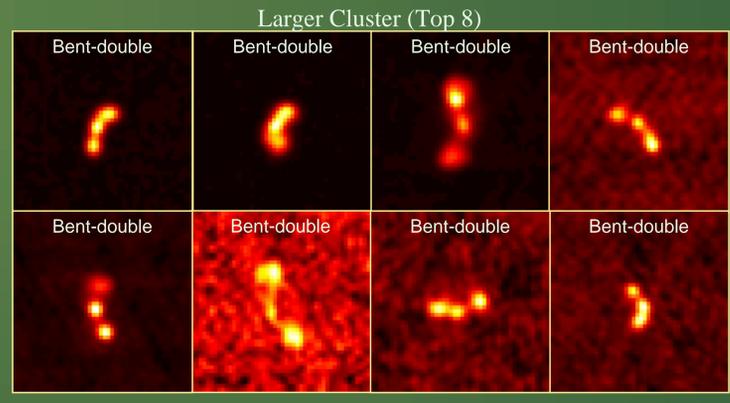
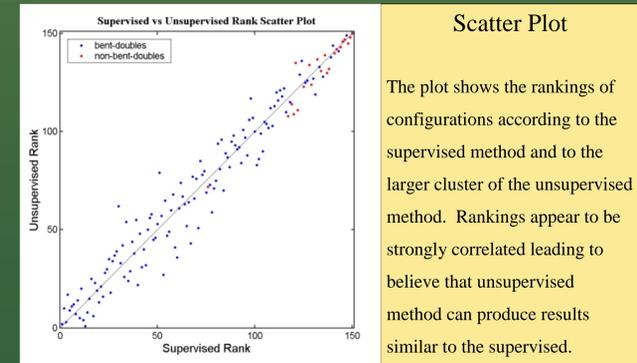
- ### Parameterization
- Parametrize P_e and P_o for all features as Gaussians
 - use transformed features
 - Want to learn the parameters of the model for bent-doubles to fit the available data best
 - maximize the likelihood of the data:
$$L = \prod_{e \in \text{bent-doubles}} P(e) = \prod_{e \in \text{bent-doubles}} \sum_{o \in \text{orderings}} P(\text{features}(e) | o) \times P(o)$$
 - sum defined in Mixture of Orderings
 - consider distribution for orderings uniform
 - learn the parameters of the Gaussians (means and variances)
 - Use Expectation-Maximization algorithm to learn the parameters
 - Once the parameters known, can calculate probability given the model for any configuration



- ### Expectation Maximization (EM)
- Iterative algorithm changing the parameters to increase the log-likelihood of the data in two steps per iteration
 - Initialize – can select parameters at random
 - E-step – calculate probability distribution
 - for each example in the training set, calculate the probability of each possible ordering given the current model
 - M-step – maximize the expected log-likelihood
 - select parameters for the new model to maximize the expected log-likelihood of data and orderings according to the distribution from E-step
 - Repeat until convergence reached
 - Convergence guaranteed; possibly, local maximum

- ### Motivation for Unsupervised Learning
- Only a very small fraction of the large data set is labeled
 - Indicates how much class information is contained in the features
 - Can be used for automated discovery of sub-classes of galaxies

- ### Unsupervised Approach
- Disregard the labels and probabilistically cluster the data using EM algorithm
 - Parameters:
 - probability for each class
 - probabilities of each orientation per class – uniform
 - parameters for Gaussians for P_e and P_o per orientation and class
 - Use all available examples (all 150)
 - Ideally, for 2 clusters, want one cluster to correspond to bent-double, another to non-bent-doubles



- ### 2-Cluster Results
- 89 in the larger cluster; 61 in the smaller cluster
 - 88 out of 89 examples in the larger cluster are bent-doubles
 - Bent-double cluster as hoped
 - Top 8 configurations for each cluster
 - larger cluster (left, top)
 - smaller cluster (left, bottom)

- ### Conclusions and Future Directions
- Method potentially accurate enough to be used by astronomers in screening candidate galaxies from large catalogs
 - Strong structure in the data as indicated by the unsupervised algorithm
 - Tests on larger set (underway)
 - Compare with other methods (underway)
 - Generalize to different number of components (2, 4+) and to missing components
 - Generalize to weighted misclassification function
 - Worse to classify a bent-double as a non-bent-double than vice versa
 - Consider other features and parametrizations