

**A PROBABILISTIC APPROACH TO FINDING
GEOMETRIC OBJECTS IN SPATIAL DATASETS OF
THE MILKY WAY**

By

Jon Purnell

A Thesis Submitted to the Graduate
Faculty of Rensselaer Polytechnic Institute
in Partial Fulfillment of the
Requirements for the Degree of
MASTER OF COMPUTER SCIENCE

Approved:

Malik Magdon-Ismail
Thesis Adviser

Rensselaer Polytechnic Institute
Troy, New York

April 2005
(For Graduation May 2005)

CONTENTS

LIST OF FIGURES

ACKNOWLEDGMENT

Initial work on this topic was done in May 2003 by Fred Liu. Also, various parts of the mathematics behind the algorithm were implemented by undergraduates at Rensselaer Polytechnic Institute. We would like to acknowledge them here: Adam Gennett, Warren Hayashi, Jim Wisniewski. Final results were obtained with the help of Nate Cole.

ABSTRACT

Data from the Sloan Digital Sky Survey has given evidence of structures within the Milky Way halo from other nearby galaxies. Both the halo and these structures are approximated by densities based on geometric objects. A model of the data is formed by a mixture of geometric densities. By using an EM-style algorithm, we optimize the parameters of our model in order to separate out these structures from the data and thus obtain a more accurate dataset for the Milky Way.

CHAPTER 1

Introduction

Recent surveys of the Milky Way halo have given astronomers a better idea of the distribution of stars in the galaxy and the location of structures that come from other nearby galaxies. In Newberg & Yanny [?, ?], the data from the Sloan Digital Sky Survey (SDSS) has a distribution inconsistent with the power-law distribution that is commonly used and also indicates the presence of a tidal stream from the Sagittarius galaxy.

A great deal of effort is spent fitting various models to this large dataset. In the case of SDSS, the dataset consists of five million stars which prohibits fitting complex distributions simultaneously by hand. The need arises for a tool that can automatically extract the background distribution of the galaxy stars as well as find certain structures of astronomical importance. (such as globular clusters, tidal streams,...). In addition to finding new structures in the Milky way, Such a tool can also yield more insight to how well a given distribution fits the Milky Way halo, which is currently a topic of controversy.

We formulate the problem of simultaneously identifying the galaxy stellar distribution and finding structures of astronomical interest as a *mixture density estimation* problem. Galaxy structure is represented as a geometric object, which is “smeared out” to obtain a parameterized probability density. The observed stellar density (from which the stars are “sampled”) is a mixture of the parameterized densities representing the geometric objects and the background stars. We assume that the observed stellar distribution forms an *i.i.d* random sample from the mixture, and the task is to extract each of the mixture components.

Our Contributions We focus on a single structure (a tidal stream) in a background, i.e. a mixture of two densities. A tidal stream is convenient because, to reasonable approximation, it can be represented as an ellipse. The main idea is that we smear out the geometric object by effectively assuming that we sample from it with noise. In the particular case we consider, the noise is cross sectional to the

ellipse. We use an EM-style algorithm to optimize for the model parameters to determine the position, orientation and size (number of stars) of the stream, in addition to the parameters describing the distribution of the background stars. In the past, in order to obtain the background distribution, one had to be careful to “look” in a direction that would avoid the stream. By *simultaneously* obtaining both distributions, we avoid this complication, and can use all the data. As a result, our estimates should be more accurate.

Our approach can be easily extended to any number of geometric objects simultaneously existing in a background, provided some parametric representation of the geometric objects exist. In particular, our approach could be used to find any geometrical object in any spatial data set, and as such could find application to other areas, for example vision. Our setup allows for additional distributions to be added to the model and for the distributions themselves to be easily changed. We give experimental results on synthetic as well as real data, which indicate that this approach performs well at automatically and simultaneously extracting both the background and structure.

Paper Organization. Next, We briefly describe the geometric tidal stream distribution and the background distribution. Following that, we describe the Maximum Likelihood framework for estimating the parameters and the resulting optimization problem. Finally, we show the results of our algorithm when applied to both synthetic data and data from the SDSS. We conclude by discussing the potential for this algorithm and future plans. These following chapters of the thesis are an extension of the work presented in Purnell, Magdon-Ismail & Newberg [?].

CHAPTER 2

Mixture Model for the Galaxy

The first step is to define the probability distributions of stars in the galaxy and in the tidal stream (the geometric object of interest). The following distributions are chosen both for their close approximation to the true distribution of stars as well as for their analytic simplicity (integrability and invertibility).

2.1 Galaxy (Background) Stellar Density

The formula for the galactic, or background, distribution, P_b , is a generalized version of the Hernquist equation:

$$P_b(x, y, z) = \frac{1}{r^\alpha (r + r_o)^{3-\alpha+\delta}} \quad (2.1)$$

where $r = \sqrt{x^2 + y^2 + (z/q)^2}$

When $\alpha = 1$ and $\delta = 1$, our formula becomes the standard Hernquist equation. q controls the scaling of the galaxy model along the z -axis and r_o controls the density of stars near the origin.

2.2 Tidal Stream Density

To represent the stars in the tidal stream, we use a longitudinal elliptical density with a two-dimensional Gaussian cross-sectional density. This has the effect of “smearing” the stars along the ellipse. An ellipse is defined by three vectors: \mathbf{a} and \mathbf{b} , the major and minor axes; and \mathbf{c} , displacement from the center of the galaxy to the center of the ellipse. Let \mathbf{x} be a generic point on the ellipse. The parametric equation for the ellipse is then

$$\mathbf{x} = \mathbf{c} + \mathbf{a} \cos t + \mathbf{b} \sin t \quad t \in [0, 2\pi] \quad (2.2)$$

In order to describe the probability density for the stream, consider $P_s(\mathbf{z})$, for a generic point \mathbf{z} . Let \mathbf{x}^* be the closest point on the stream to \mathbf{z} , and suppose that $\mathbf{x}^* = \mathbf{c} + \mathbf{a} \cos t^* + \mathbf{b} \sin t^*$. We define a cross sectional basis $\{\mathbf{E}_1, \mathbf{E}_2\}$ for the stream at the point \mathbf{x}^* by the two unit vectors

$$\mathbf{E}_1 = \frac{\mathbf{a} \times \mathbf{b}}{\|\mathbf{a} \times \mathbf{b}\|}, \quad \mathbf{E}_2 = \frac{\mathbf{E}_1 \times \dot{\mathbf{x}}^*}{\|\mathbf{E}_1 \times \dot{\mathbf{x}}^*\|} = \frac{\mathbf{E}_1 \times (-\mathbf{a} \sin t^* + \mathbf{b} \cos t^*)}{\|\mathbf{E}_1 \times (-\mathbf{a} \sin t^* + \mathbf{b} \cos t^*)\|} \quad (2.3)$$

where $\dot{\mathbf{x}}^* = -\mathbf{a} \sin t^* + \mathbf{b} \cos t^*$ is the tangent vector at \mathbf{x}^* . We can then write

$$\mathbf{z} = \mathbf{c} + \mathbf{a} \cos t^* + \mathbf{b} \sin t^* + x\mathbf{E}_1 + y\mathbf{E}_2, \quad (2.4)$$

where

$$x = (\mathbf{z} - \mathbf{c} - \mathbf{a} \cos t^* - \mathbf{b} \sin t^*) \cdot \vec{\mathbf{E}}_1 \quad (2.5)$$

$$y = (\mathbf{z} - \mathbf{c} - \mathbf{a} \cos t^* - \mathbf{b} \sin t^*) \cdot \vec{\mathbf{E}}_2 \quad (2.6)$$

We then have

$$P_s(\mathbf{z}) = \frac{1}{2\pi} \cdot \frac{1}{2\pi\sqrt{\det \Sigma}} e^{-\frac{1}{2}\mathbf{y}^T \Sigma^{-1} \mathbf{y}},$$

where $\mathbf{y} = \begin{bmatrix} x \\ y \end{bmatrix}$. This density corresponds to choosing t uniformly in $[0, 2\pi]$ along the ellipse, and then adding Gaussian cross-sectional noise with variance covariance matrix Σ . We will simplify this general model to assume that the cross-section is spherically symmetric, in which case we get the simpler density

$$P_s(\mathbf{z}) = \frac{1}{2\pi} \cdot \frac{1}{2\pi\sigma^2} e^{-\frac{1}{2\sigma^2}(x^2+y^2)}.$$

σ corresponds to the ‘‘thickness’’ of the tidal stream.

2.3 Mixture Model Density

The mixture density is obtained by combining the galaxy and tidal stream densities. To combine the distributions, they must be both normalized within the region of interest and weighted by a mixing parameter ϵ . Normalization is achieved

by dividing the distributions by their integral over the region of interest. These integrals can be calculated analytical for simple densities and by numerical integration for more complex densities. The final mixture density is given by:

$$P_m(\mathbf{x}) = \epsilon \frac{P'_b(\mathbf{x})}{\int P'_b} + (1 - \epsilon) \frac{P'_s(\mathbf{x})}{\int P'_s} \quad (2.7)$$

where $P'_b(\mathbf{x}) = \rho(\mathbf{x})P_b(\mathbf{x})$ and $P'_s(\mathbf{x}) = \rho(\mathbf{x})P_s(\mathbf{x})$. $\rho(\mathbf{x})$ is an efficiency which corresponds to the probability that a star is detected given that it is at position \mathbf{x} . Thus the presence of stars at large distances in our dataset indicates that the density there must be higher because fewer of these stars are detected.

2.4 Normalization

In normalizing our probabilities, we calculate the integral of each distribution over the area of our dataset. In our physical application, the dataset is in a wedge shape that pivots at our Sun given by:

$$307^\circ \leq \theta < 436^\circ \quad -1.25^\circ \leq \phi < 1.25^\circ \quad 1.4 \leq r < 57.5 \text{ kparsecs} \quad (2.8)$$

The analytical solution to the normalization integrals is not easy. Also, if we were to change the distributions, we would have to re-calculate the solution. For these reasons, we choose a numerical integration approach.

The integrals are obtained by dividing the wedge into several smaller volume elements and computing the Rieman sum approximation to the integrals. The small volume elements we used are defined by:

$$\Delta\theta = \frac{436 - 307}{490} \approx 0.26 \quad (2.9)$$

$$\Delta\phi = \frac{1.25 + 1.25}{10} = 0.25 \quad (2.10)$$

$$\Delta r = \frac{57.5 - 1.4}{180} \approx 0.31 \quad (2.11)$$

Let A represent a generic small volume element, V_A its volume and \mathbf{O}_A its

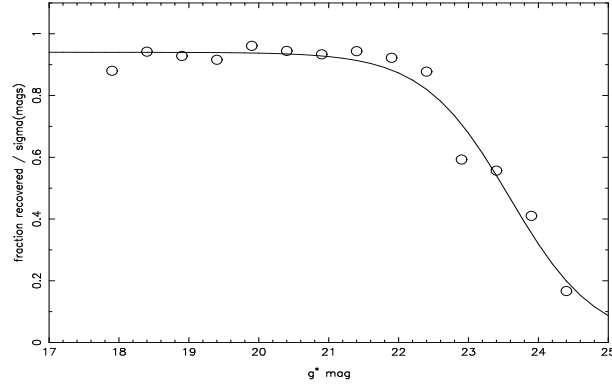


Figure 2.1: Completeness of the dataset where the color is in the range of (0.1,0.3). The circles indicate the completeness estimated by comparing two overlapping runs. The solid line is the function fitted to these estimates

center. Then the Rieman approximation to the integral is given by:

$$\int P'_b = \sum_A P'_b(\mathbf{O}_A)V_A \quad \int P'_s = \sum_A P'_s(\mathbf{O}_A)V_A \quad (2.12)$$

2.5 Data Efficiency

The ratio of the number of stars in our dataset to the number of actual stars in the galaxy is called the efficiency of our dataset and is a function of distance from the sun, r_S , which is measured in parsecs. By comparing overlapping runs, and the number of matched stars, we can estimate the efficiency function. Figure 1 shows a plot of the efficiency versus g^* where:

$$g^* = 5 \log_{10}(r_S/10) + 4.2 \quad (2.13)$$

The solid line in Figure 1 is the result of fitting a function to these data points. This derived efficiency formula is given by:

$$\rho(g^*) = \frac{0.9402}{e^{1.6171(g^*-23.5877)} + 1} \quad (2.14)$$

CHAPTER 3

Maximum Likelihood Estimation

We use maximum likelihood density estimation to estimate the parameters in our model. The dataset is composed of the stars $\{\mathbf{x}_i\}_{i=1}^N$. The likelihood is given by $P[\{\mathbf{x}_i\}_{i=1}^N | Parameters] = \prod_i P_m(\mathbf{x}_i)$. Taking the logarithm and dividing by N, we then maximize $\mathcal{E}(Parameters) = \frac{1}{N} \sum_{i=1}^N \log P_m(\mathbf{x}_i)$.

3.1 Parameter Optimization

To maximize \mathcal{E} , we use a standard conjugate gradient optimization algorithm, described in [?], together with the line search algorithm in [?].

Each step taken in the conjugate gradient algorithm involves calculating a direction vector based on the gradient at the current parameter set, finding the maximum probability that can be obtained along this direction and then testing whether the stopping conditions have been met. We use a simple finite difference scheme to obtain the gradient.

$$\left(\frac{\partial \mathcal{E}}{\partial \mathbf{p}_i}\right)_{\pm} = \pm \left(\frac{\mathcal{E}(\mathbf{p} \pm h_i \mathbf{e}_i) - \mathcal{E}(\mathbf{p})}{h_i}\right) \quad (3.1)$$

where \mathbf{e}_i is the standard unit vector.

Since \mathcal{E} has varying sensitivities to the different parameters there is no single value for h that we can use. In fact, due to the normalization constants in our probability distributions, we may find ourselves near a non-differentiable cusp on the error surface. To get around this, we start with a sufficiently large value for h , for each parameter \mathbf{p}_i . We find \mathcal{E} with that parameter at its current value, \mathbf{p}_i , and at the values of $\mathbf{p}_i \pm h$. If the signs of $\left(\frac{\partial \mathcal{E}}{\partial \mathbf{p}_i}\right)_{+}$ and $\left(\frac{\partial \mathcal{E}}{\partial \mathbf{p}_i}\right)_{-}$ are different, we know that we are near a local minimum or maximum and that our value for h is too large. In this case we decrease h by half and check the signs of the gradients again. This process continues until either the signs of the gradients agrees or until the value of h falls below a precision tolerance. If the signs agree, we use either gradient as shown

in Equation (??). If h falls below the precision tolerances, we must be sufficiently close to a maximum or minimum and so we consider the gradient to be 0 with respect to \mathbf{p}_i .

3.1.1 Parameter Representation

There are two aspects of our parameters that can help us increase the efficiency of our algorithm. First is that some parameters are redundant and can be eliminated. This will decrease the running time since there will be fewer parameters and fewer dimensions in our parameter space to search through. The second aspect is that some parameters can only take on a particular range of values. Converting these parameters into variables without constrained ranges allows us to use more efficient unconstrained optimization techniques.

3.1.1.1 Reducing the Number of Parameters

Our first reduction takes advantage of the orthogonality condition between \mathbf{a} and \mathbf{b} . Since $\mathbf{a} \cdot \mathbf{b} = 0$ we can define \mathbf{b} with a magnitude $d_{\mathbf{b}}$ and an angle $\theta_{\mathbf{b}}$ and remove one parameter.

Another reduction we make arises from approximating the stream width. As noted earlier, we made the assumption that the cross-section of a stream is circular with equal variations along both axes. So we are able to replace both widths by just one parameter, σ .

3.1.1.2 Removing Parameter Constraints

Since the constraints on our parameters are relatively simple bound constraints, we can optimize with respect to unconstrained parameters by explicitly incorporating the constraints in the objective function as follows. Suppose α is a parameter in \mathcal{E} , i.e. $\mathcal{E} = \mathcal{E}(\alpha)$ (we only show the α dependence). Suppose that α has a bound constraint $\alpha \in [A, B]$. We can write $\mathcal{E}(\alpha)$ in terms of an unconstrained parameter β by $\mathcal{E}(\alpha) \rightarrow \mathcal{E}(\alpha(\beta))$ where $\alpha(\beta) = A + (B - A)e^{-\beta^2}$. β now becomes an unconstrained parameter in the optimization and α can easily be obtained from β . Such bound constraints apply to $\epsilon \in [0, 1]$ and $q \in [0, 1]$. An unbound constraint

of the form $\alpha \in [A, \infty)$ can also be incorporated using the mapping $\alpha(\beta) = A + \beta^2$. Such an unbound constraint applies to $\delta \in [0, \infty)$.

Stopping Condition: The stopping condition is important because if we end the search too early, the resulting parameter set will not be optimal. However, if the stopping condition allows the search to continue for too long, the algorithm will be inefficient, which is disastrous on datasets of size hundreds of millions.

CHAPTER 4

Results

We applied our algorithm to two datasets, one synthetic and one from a 2.5 degree thick wedge along the Celestial equator of the SDSS dataset. For each dataset we initialized the algorithm with a parameter set that was close to the optimal value but randomized. The algorithm then ran until the maximum component of the gradient was less than 0.002.

4.1 Synthetic Data

The synthetic data was generated using the exact model mixture density for a particular setting of the parameters. For each star generated, a random number determined whether it was to be a stream star or a background star. If the star is a stream star, three random numbers are generated to determine t^* , x , and y in (??). an angle along the stream's ellipse, and the cross sectional coordinates. If the star is a background star, three random numbers are generated to determine the coordinates of the star. The z-axis component of this coordinate was then multiplied by q to take the squashness into account.

The synthetic data was generated with the following parameters:

$$\mathbf{c} = (6.9, 10.23, 0.166) \quad \mathbf{a} = (19.4, 9.8, 35.5) \quad \mathbf{b} = (18.5, -2.45, -9.43)$$
$$\sigma = 5.0 \quad q = 0.65 \quad r_o = 13.5 \quad \epsilon = -2.197$$

The algorithm ran for 19 iterations and ended with the following parameters and probability:

$$\mathbf{c} = (7.32, 6.43, -4.72) \quad \mathbf{a} = (22.26, 9.99, 33.08) \quad \mathbf{b} = (15.40, -9.64, -7.45)$$
$$\sigma = 3.32 \quad q = 0.78 \quad r_o = 13.65 \quad \epsilon = -2.78 \quad \mathcal{E} = -3.38358$$

Figure ?? shows the generated synthetic data. Figure ?? shows the separation.

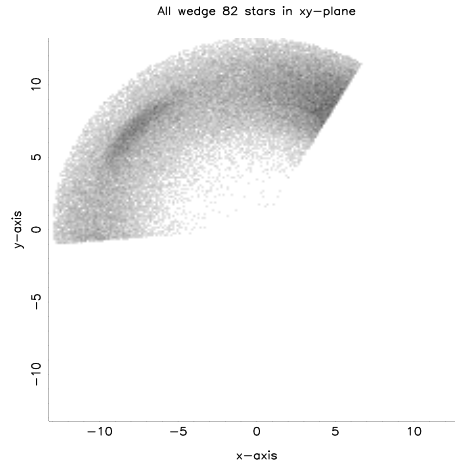


Figure 4.1: Density plot of synthetic data in log space

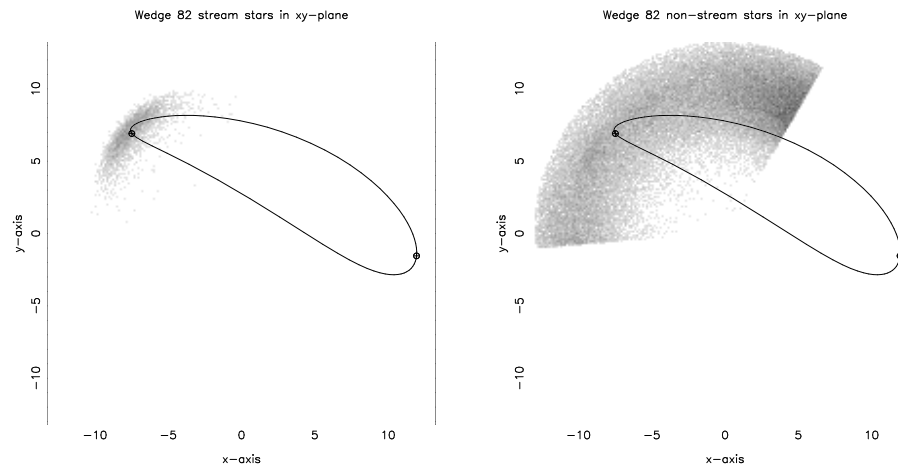


Figure 4.2: Separation plot of the synthetic data. The left plot shows the stars labeled as stream stars and the right plot shows the stars labeled as background, or galactic, stars. The line indicates the ellipse of the stream and the circles indicate the point where the stream intersects the plane of the data

4.2 Real Data

With the real data, the algorithm ran for 10 iterations and ended with the following parameters and probability:

$$\mathbf{c} = (6.06, 12.85, -0.039) \quad \mathbf{a} = (19.49, 13.34, 35.70) \quad \mathbf{b} = (19.76, -4.26, -9.20)$$

$$\sigma = 6.21 \quad q = 0.71 \quad r_o = 14.43 \quad \epsilon = -2.38 \quad \mathcal{E} = -3.39434$$

Figure ?? shows a plot of the real data. Figure ?? shows the separation.

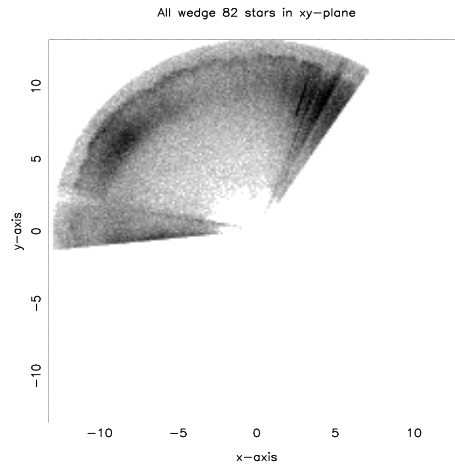


Figure 4.3: Density plot of real data from SDSS

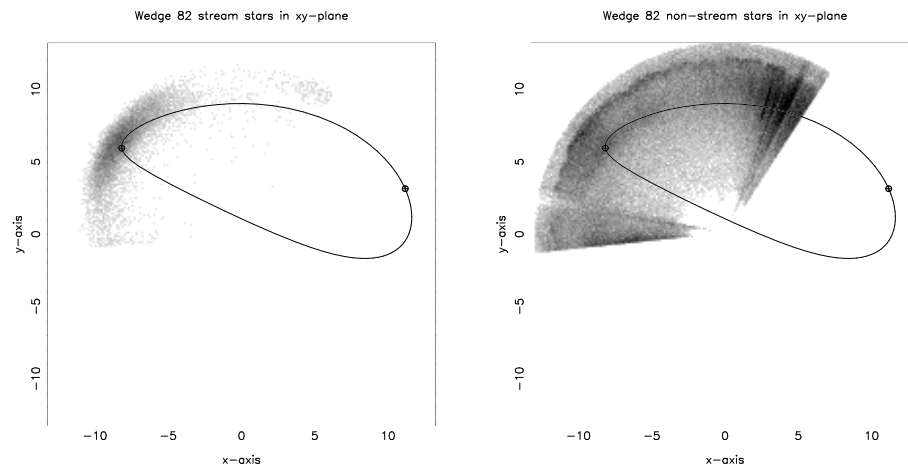


Figure 4.4: Separation plot of real data. The left plot shows the stars labeled as stream stars and the right plot shows the stars labeled as background, or galactic, stars. The line indicates the ellipse of the stream and the circles indicate the point where the stream intersects the plane of the real data

CHAPTER 5

Conclusion

We have presented a probabilistic approach to finding geometric objects in spatial databases. From the plots, the algorithm performs well with synthetic data and equally well with the real data. From these results we can get a clear idea of the direction and size of the tidal stream as well as the distribution of stars in the halo of the Milky Way. Future plans for this work are to increase the number and types of structures that are included in the mixture model, and to scale it up to handle millions of stars. A further direction is to investigate non-parametric estimates of the background densities and resulting geometric structures.

For this particular application, on going work includes different parameterizations for the tidal streams (such as piecewise linear) as well as incorporation of multiple wedges from the SDSS dataset simultaneously. This will allow astronomers a glimpse at a more global view of the tidal stream for the first time. In particular, a pressing question is whether the Sagittarius tidal stream is planar. A further application of the probabilistic viewpoint is that it allows one to probabilistically separate the geometric structure from the background. With this approach, astronomers, now, have a tool for extracting the stream from the entire collection of stars, and hence the ability to study separately the stream structure (without being obscured by galaxy structure. Such a tool can give insight into galaxy dynamics, in particular the dynamics of galaxy collisions, since it is believed that the Sagittarius tidal stream is the result of a smaller galaxy colliding with the Milky Way.

Related Work Previous research has used mixture models and EM algorithms for clustering in large databases [?, ?]. Our techniques similarly uses mixture models for extracting structure. Similar problems arise in vision where one tries to identify objects in a scene (see for example [?]). Edge detection, clustering and feature selection are often used to solve these problems.

LITERATURE CITED

- [1] Brian Yanny Heidi Jo Newberg. The ghost of sagittarius and lumps in the halo of the milky way. *The Astrophysical Journal*, 2002.
- [2] Brian Yanny Heidi Jo Newberg. Sagittarius tidal debris 90 kiloparsecs from the galactic center. *The Astrophysical Journal*, 2003.
- [3] Jon Purnell Malik Magdon-Ismail Heidi Jo Newberg. A Probabilistic Approach to Finding Geometric Objects in Spatial Datasets of the Milky way. *International Symposium on Methodologies for Intelligent Systems*, 2005.
- [4] R.P. Brent. *Algorithms for Minimization without Derivatives*. Prentice-Hall, 1973.
- [5] Christopher M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
- [6] M. Jorgensen L. Hunt. Mixture model clustering for mixed data with missing information. *Computational Statistics and Data Analysis*, 2003.
- [7] C. Reina P. Bradley, U. Fayyad. Clustering very large databases using EM mixture models. *Proc. 15th International Conference on Pattern Recognition*, 2000.
- [8] Shimon Ullman. *High-Level Vision*. MIT Press, 1996.