

---

# Discovering Natural Kinds of Robot Sensory Experiences in Unstructured Environments

---

Daniel H Grollman  
Odest Chadwicke Jenkins  
Frank Wood

DANG@CS.BROWN.EDU  
CJENKINS@CS.BROWN.EDU  
FWOOD@CS.BROWN.EDU

Brown University Department of Computer Science, Providence, RI 02906

## Abstract

We address the symbol grounding problem for robot perception through a data-driven approach to deriving sensor categories. Unlike model-based approaches, our method learns intrinsic categories (or natural kinds) from the raw data itself. We approximate a manifold underlying sensor data using Isomap nonlinear dimension reduction and apply Bayesian clustering (Gaussian mixture models) to discover categories. We demonstrate our method through the learning of sensory kinds from trials in various indoor environments with multiple sensor modalities. Learned kinds are used to classify new sensor data (out-of-sample readings). We present results indicating greater consistency in classifying sensor data employing mixture models in low-dimensional embeddings.

## 1. Introduction

The symbol grounding problem in robotics deals with connecting arbitrary symbols with entities in the robot's world. Names such as 'door', 'hallway', and 'tree' must be associated with sensor readings so that an autonomous robot can reason about them at a higher level. Traditionally, a human programmer is relied upon to provide these connections by identifying categories in the world and building *models* of how they would appear to the robot. However, actual sensory information is dictated by the robot's embodiment and may not accord with models of sensor function. Consequently, our understanding of a robot's perception of the world is often biased and heuristic.

A data-driven approach to sensor analysis could discover a more appropriate interpretation of sensor readings. Sensor data collected during robot operation are observations of the underlying sensory process and, if teleoperation is involved, the control policy of the operator. We posit that the intrinsic structure underlying

robot sensor data can be uncovered using recent techniques from manifold learning. Once uncovered, sensory structures can provide a solid foundation for autonomous sensory understanding as a robot's perceptual system is allowed to develop classes of sensor data based on its own, unique, experiences.

We present a data-driven method for classifying robot sensor input via unsupervised dimension reduction and Bayesian clustering. We view the input of the system as a high dimensional space where each dimension corresponds to a reading from one of the robot's sensors. This sensory space is likely to be sparse and described by a lower dimensional subspace. Our approach is to embed sensor data into a lower-dimensional manifold that condenses this space and captures latent structure. By clustering in this embedded space we generate simpler probability densities while grouping together areas that appear similar to the robot. We take each cluster of sensor readings in the reduced-dimensional space as a kind<sup>1</sup> of entity as viewed by the robot.

Once classes are learned, new sensor readings can be quickly classified with an out-of-sample (OOS) classification procedure: After projecting new samples into the embedding space they are classified with a Gaussian mixture model (GMM). When revisiting a location this procedure should embed the new readings near the old ones, consistently classifying the space.

We use consistency as an evaluation metric because ground truth is unknown and often subjective. The clusters developed by this technique reflect areas that are perceived similarly by the robot and may not reflect any categories we would develop ourselves. It is however important that the found kinds be consistent, by which we mean that similar inputs should belong to the same kind and be classified similarly.

---

<sup>1</sup>Philosophically, a natural kind is a collection of objects that all share salient features. For instance, the 'Green Kind' includes all green objects. We use the terms 'kind', 'class' and 'category' interchangeably.

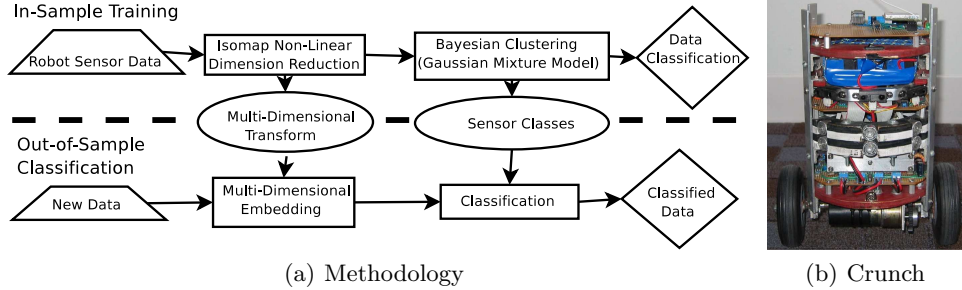


Figure 1. Our method in flowchart form. Data from robot sensors are analyzed with Isomap to obtain a low-dimensional embedding. The embedded data is then clustered to develop sensor classes. Out-of-sample data can be quickly projected and classified using models learned during the in-sample training.

## 2. Related Work

Topological mapping depends on the ability to discover regions in an explored area (Thrun, 1998). This process is usually done by extracting features from sensor data that indicate the robot’s current location. When a human decides which region types exist in the robot’s world and which features are important (Tomatis et al., 2003), biases from models of sensor operation are introduced. We attempt to remove these biases by deriving classes directly from sensor data.

Localization techniques also depend on region identification. Landmarking, or the identification of unique places, is commonly used to let a robot know when it has returned to a previously visited location on a map (i.e., revisiting, loop closure). The revisiting problem is key when it comes to map-making because it allows a robot to discover loops in the world (Howard, 2004) or, in the case of multiple exploration robots, it allows one robot to discover when it has entered space explored by another (Stewart et al., 2003). Often, landmarking is accomplished by modifying the environment to disambiguate similar places. We hypothesize that with a data-driven classification technique, it will become clearer which areas of the world look similar to the robot and require disambiguation.

A semi-supervised approach to discovering clusters in vision data is introduced in (Grudic & Mulligan, 2005). By allowing each cluster to self-optimize its parameters, they are able to discover clusters that more accurately correspond to the predefined ones, as well as detect outlying points that do not belong to any cluster. However, the original clusters must be decided upon by human operators and exemplar photographs of each cluster are provided to the algorithm. In contrast, our approach is completely unsupervised and allows for the discovery of space classes and outliers that are potentially non-obvious to humans.

In order to tie sensing and action together, (Klingspor

et al., 1996) learn sensory and action concepts directly from the sonar data of a robot, after the data is segmented and categorized by hand. By utilizing sensor information related to actions (such as teleoperation data), we can determine the usual action performed in each space class in an unsupervised way and use these actions as a first-attempt control policy.

## 3. Methodology

Our method, outlined in figure 1(a), views  $d$ -dimensional robot sensor data as lying on a manifold in  $\mathbb{R}^d$ . We model each sensor datum  $\vec{x}$  as having been generated by a mixture model on this manifold, where each mixture density corresponds to a natural kind. Here we closely follow the methods and notation of (Bengio et al., 2004).

**For training**, let  $D = \{\vec{x}_1, \dots, \vec{x}_N\}$  be the collection of readings from  $S$  sensors at  $N$  time instances. We compute an affinity matrix  $M$  by approximating the geodesic distance between points on the sensor data manifold. As in (Tenenbaum et al., 2000), we define the geodesic distance between points  $a$  and  $b$  to be:

$$\tilde{D}(a, b) = \min_p \sum_i d(p_i, p_{i+1})$$

where  $p$  is a sequence of points of length  $l \geq 2$  with  $p_1 = a$ ,  $p_l = b$ , and  $p_i \in D \forall i \in \{2, \dots, l-1\}$  and  $(p_i, p_{i+1})$  are neighbors as determined by a  $k$ -nearest neighbors algorithm. We compute  $\tilde{D}$  by applying Dijkstra’s algorithm (Cormen et al., 1990) to the graph  $V = D, E = \{p_i, p_{i+1}\}$  where edge length is the Euclidean distance between neighbors.

$M$  is formed with elements  $M_{ij} = \tilde{D}^2(x_i, x_j)$  and then converted to equivalent dot products using the “double-centering” formula to obtain  $\tilde{M}$ .

$$\tilde{M}_{ij} = -\frac{1}{2}(M_{ij} - \frac{1}{2}S_i - \frac{1}{N}S_j + \frac{1}{N^2} \sum_k S_k)$$

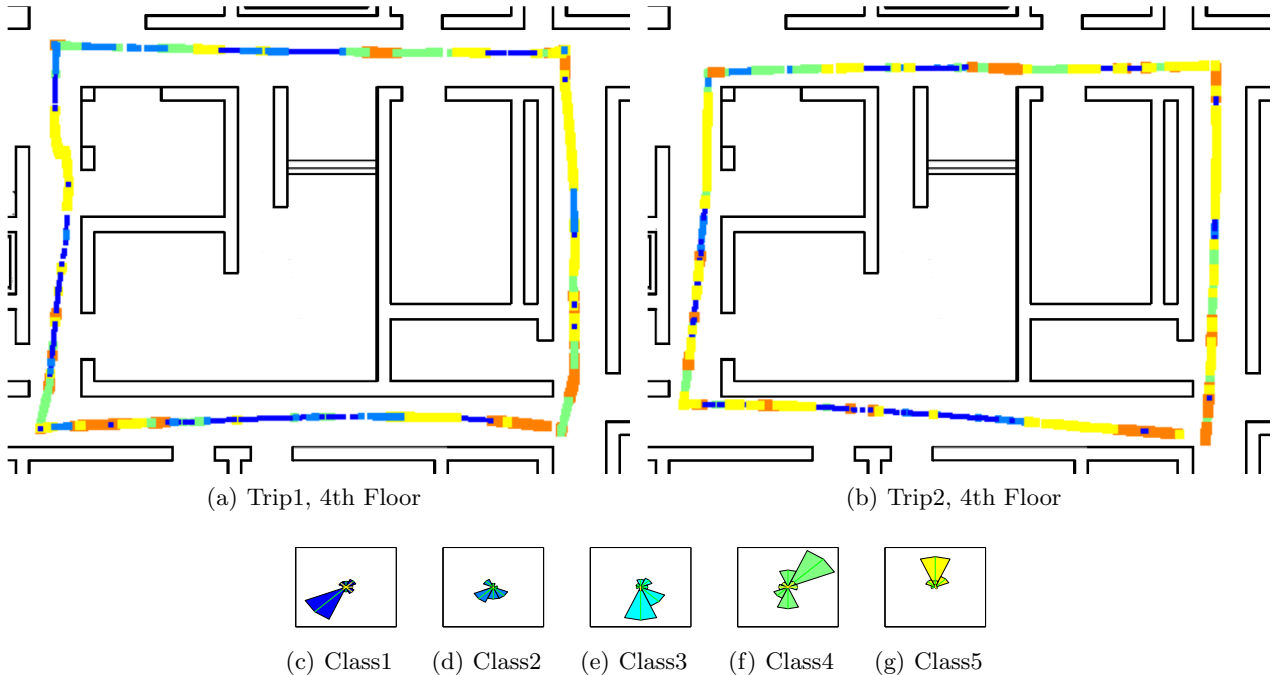


Figure 2. Results from Crunch on the fourth floor of the Brown University CIT Building. ???: Sensor data from trip 1 has been clustered into 5 classes. A unique width and color value for each class is overlaid on registered odometry to show the classification of regions of space. ???: The learned classes were used to classify data from a second trip. 2(c)-2(g): The mean-centered expected sensor readings for each class under the standard ray model.

where  $S_i = \sum_j M_{ij}$ . In practice, this grows as  $N^2$  and is thus currently infeasible to calculate for more than a few thousand points. For larger datasets, only a subset of the data can be fully processed (landmarks).

The  $k$  dimensional embedding  $\vec{e}_i$  of each sensor output  $\vec{x}_i$  on the sensor data manifold is obtained via Multi-Dimensional Scaling (MDS). Here the embedding is approximated by the vector  $\vec{e}_i = [\sqrt{\lambda_1}v_{1i}, \sqrt{\lambda_2}v_{2i}, \dots, \sqrt{\lambda_k}v_{ki}]$  where  $\lambda_k$  is the  $k^{th}$  largest eigenvalue of  $M$  and  $v_{ki}$  is the  $i^{th}$  element of the corresponding eigenvector. We reduce the dimensionality of the sensor data by setting  $k < d$ , thus removing many of the low eigenvalue coordinates of the embedding.  $k$  is selected by comparing the error between distances in the input and reduced spaces for different dimensionalities. In particular, we look for an ‘elbow’, a point after which increasing dimensionality does not lead to a significant decrease in residual variance. In practice we take  $k$  to be few dimensions higher than the elbow. We then define  $E = \vec{e}_1, \dots, \vec{e}_N$  to be the reduced dimensionality embedding of the training sensor data  $D$  in  $k$  dimensions, henceforth referred to as the “sensor embedding.”

Initially, we assume that the sensor embeddings were generated by exactly  $J$  statistically distinct intrinsic classes of sensor readings. We assume that the distri-

bution of each of these classes is Gaussian and fit  $E$  with a mixture model with  $J$  components.

The probability that  $\vec{e}_i$  was output by the robot’s sensors while it was in a physical space corresponding to sensor class  $j$ ,  $1 < j < J$  given these assumptions is:

$$P(\vec{e}_i|j) = \frac{1}{(2\pi)^{\frac{k}{2}} \sqrt{\det(\Sigma_j)}} \exp\left(-\frac{1}{2}(\vec{e}_i - \mu_j)^T \Sigma_j^{-1} (\vec{e}_i - \mu_j)\right)$$

where  $\mu_j$  and  $\Sigma_j$  are the mean and covariance of the sensor output while in class  $j$ .

Assuming that each sensor datum is independent, then the probability of  $E$  according to the mixture model is:

$$P(E) = \prod_{i=1}^N \sum_{j=1}^J \alpha_j P(\vec{e}_i|j)$$

where the  $\alpha_j > 0$  are mixing coefficients and  $\sum_{j=1}^J \alpha_j = 1$ .

The EM algorithm (McLachlan & Basford, 1988) is used to maximize  $P(E)$  by solving for optimal distribution parameters and membership weights.

Model selection is a central issue in clustering and corresponds to determining the number of clusters (intrinsic classes) in the data. We employ two existing empirical criteria for model selection, Bayesian Information Criteria (BIC) and cross-validation. The BIC

penalizes likelihood as a function of the complexity of the model. If  $\kappa$  is the number of free parameters in the model, then we calculate the BIC as:

$$-2\log(\mathcal{L}(\Theta|E, \mathcal{Y})) - \kappa(\log(N) + 1)$$

Since in practice the BIC often doesn't sufficiently penalize complex models, we additionally use cross-validation on held-out data to check for overfitting: We train our model on half the training data and then compute the unpenalized likelihood of the remainder. When too many classes are posited, i.e. the model may be over-fit, the likelihood of the held-out data may decrease relative to simpler models. These two techniques guide us in selecting  $J$ .

**Online classification** of a new point  $\vec{p}$  is simple and rapid. We refer the reader to (Bengio et al., 2004) for full details. The embedding is given by:

$$e_k(\vec{p}) = \frac{1}{2\sqrt{\lambda_k}} \sum_i v_{ki} (E_{\vec{x}}[\tilde{D}^2(\vec{x}, \vec{x}_i)] + E_{\vec{x}'}[\tilde{D}^2(\vec{p}, \vec{x}')] - E_{\vec{x}, \vec{x}'}[\tilde{D}^2(\vec{x}, \vec{x}')] - \tilde{D}^2(\vec{x}_i, \vec{p}))$$

where  $E$  is an average over the training data set. Using the GMM from the training stage, we determine the probability of this newly embedded point belonging to each cluster.

## 4. Experiments

To test our algorithms, we collected robotic sensory data as a robot was teleoperated through an environment several times. Data from one trip was analyzed using our training algorithm to learn embedding and clustering parameters. Then, data from other trips were run through our out-of-sample algorithm. Categorizations from multiple trips in the same environment were then examined for consistency. We compared the results of our Isomap based algorithm to one based on Principle Component Analysis (PCA). For PCA we used the first  $k$  principle components of the data as the embedding space, where  $k$  is the number of intrinsic dimensions discovered by Isomap.

Data was collected in an office environment, using **Crunch**, the small, cylindrical, inverted pendulum robot pictured in figure 1(b). It has eight sonar and eight IR sensors arranged in dual rings around its body as well as wheel encoders that record wheel rotation. During operation, these sensors are sampled and transmitted back to a base laptop where they are logged at around 10Hz.

For the training phase, one set of data was used to discover embedding and clustering parameters. Af-

ter computing the geodesic distance and MDS embedding of the training data, we examined the residual variances and retained 8 of the resulting dimensions for future processing. Based on the BIC and hold-out calculations, we judged that there were 5 classes in the data. The resulting mixture model was used to assign each datapoint to a class. For display purposes, we manually registered the odometry with the underlying floor plan and overlaid these classes on the path that the robot followed. This assignment is illustrated in figure ???. Figures 2(c)-2(g) show expected readings from each of the 5 classes discovered by our method. These images were generated using a "ray model" of Crunch's IR and sonar sensors and the values were computed from a weighted average of the mean-centered datapoints. Under this model, many of these shapes are hard to interpret as corresponding to a hallway, doorway, corner, etc, but these are the sensor readings that are most distinguishable to the robot.

### 4.1. Consistent Sensory Classification

Our first experiment was designed to test the consistency of our classification when a location is revisited. We used the parameters learned from the training stage to classify data from a second trip in the same environment. As the robot followed the same general path as it did in the first trip, we expected the sensory readings along the path to be classified similarly across trips. The results from the out-of-sample classification of Crunch's second trip are shown in figure ??.

We used the registered odometry to compare classifications of the same physical space across trips. Given a position  $(x, y)$  in the first trip that has been classified as generating sensor readings of kind  $k$ , we compare it to all points from the second trip within a certain radius,  $r$ . If more than a given percentage  $z$  of these points have also been classified as kind  $k$ , we declare a match. Figure 3 shows a plot of how measured consistency varies with these two parameters for both our approach and the PCA based one applied to the Crunch data set. As you can see, manifold learning greatly improves the consistency of classification. For example, if we require 50% of points within a 1 foot radius to be classified the same, our approach achieves 40% consistency, while PCA performs at less than 10%. The consistency between random assignments over these paths is  $\sim 0.5\%$ .

### 4.2. Consistency in New Spaces

To evaluate the applicability of our approach to new, but similar, spaces, we ran Crunch on a different floor

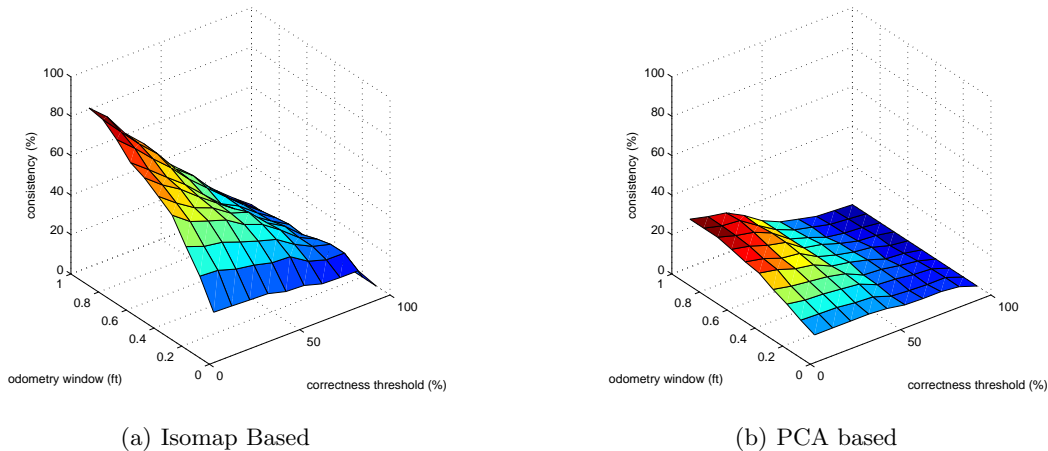


Figure 3. The consistency metric is highly sensitive to constant selection and registration errors. Here we show the measured consistency of our Isomap based technique as the constants  $(r, z)$  are varied, and compare it to a PCA-based version 3(b).

of our CIT building. Data from two trials in this new space were collected and separately classified using our out-of-sample technique. Results are shown in Figure 4. If the learned classes were non-applicable to the space, that is, if areas that looked similar to the robot were not assigned to the same cluster, we would expect to see successive data points assigned to different classes. Instead, there are several large contiguous sections of points that are all assigned to the same class. Furthermore, by repeating the consistency test from above, and classifying data from a second trip on the fifth floor using the same classes, we see that these classifications are usable in this area, even though they were not learned here.

## 5. Discussion and Conclusion

We attempt to remove human bias from the analysis of robotic sensor data by identifying latent structure in the sensor readings themselves. Currently, we empirically determine the neighborhood function and size, the number of embedding coordinates to retain, and the number of intrinsic sensor classes. In theory, each of these can be determined automatically, and perhaps even adaptively, from the data. In particular, model selection is very difficult. There are techniques to alleviate this issue, such as infinite mixture models (Blei et al., 2003) that can be incorporated in future work. The main contribution of the work presented here is in demonstrating that intrinsic sensor classes may form a better foundation for applications that require classifying sensor data. In addition, we treat each sensor reading as independent and compute affinities based on an approximation of geodesic distance. Better per-

formance may result from modeling spatial and temporal correlations as in ST-Isomap (Jenkins & Mataric, 2004), or by using different affinity propagation techniques such as SDE and HLLE (Donoho & Grimes, 2003).

Because our technique operates in a space defined by robot sensors, the results are sometimes difficult to reconcile with human intuition. In particular, when the “canonical” sensor reading for a Crunch class is examined, it does not correspond to any class that we, as humans, would have developed for the robot. In fact, even the *number* of classes in the space differs. However, as Crunch is a small wheeled robot equipped with sonar and IR and we are tall humans with eyes, it makes sense that our world views, and our divisions of that world into categories, would be different. Our intuition is further bolstered by noting that armed with the kinds discovered by our system, a human crawling on his hands and knees through the area explored by Crunch can see how they match up.

Alas, there is no “ground truth” we may use to evaluate our model. By design we cannot determine the “correct” classification of each point in robot sensor space. At most, we can use an ad-hoc metric to test classifications for consistency. The metric described herein is highly sensitive to registration errors and constant selection. It served only to help us intuit that our classification scheme is consistent and reapplicable.

### 5.1. Mapping and Control

One use of our system would be the creation of topological maps of the robot’s environment. Such “robot-centric” maps (Grudic & Mulligan, 2005) re-

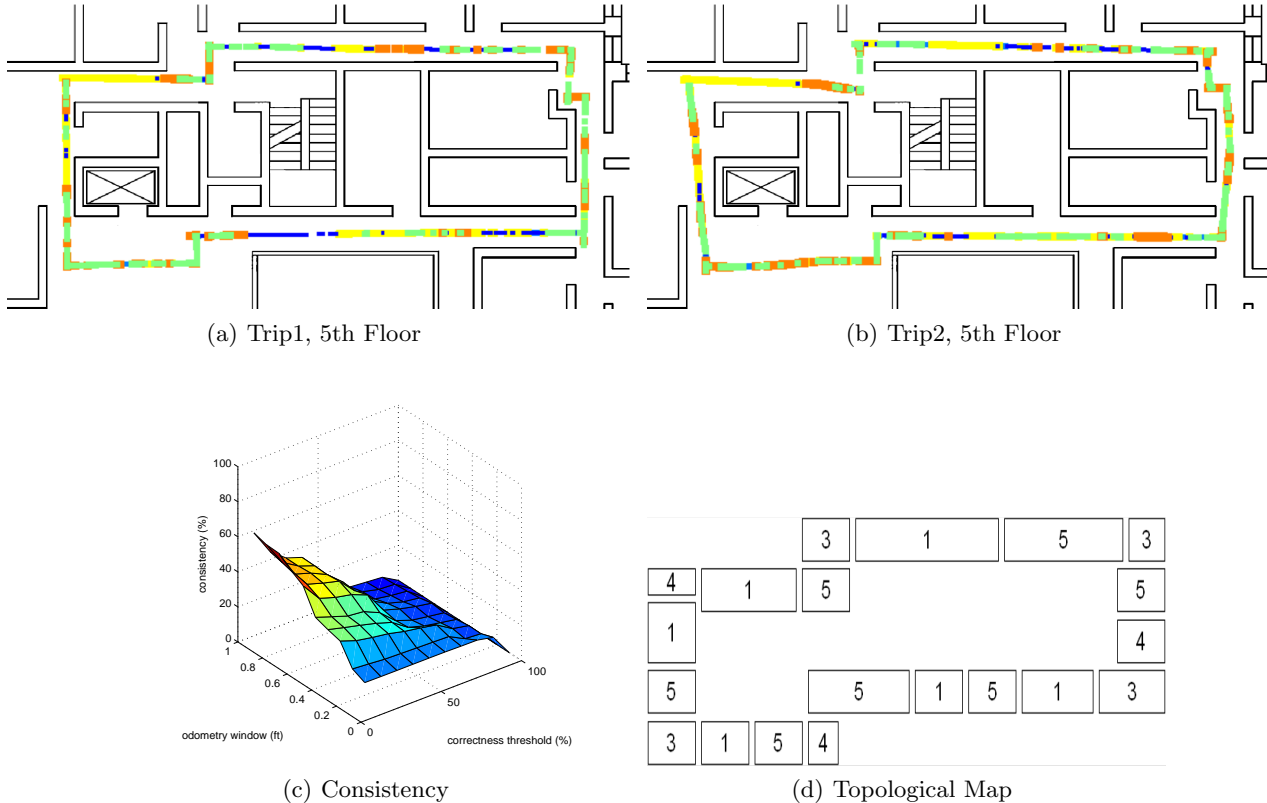


Figure 4. Using the sensor classes discovered on the fourth floor, Crunch took two trips around the fifth floor of our building and classified each datapoint. As before, the consistency metric, 4(c) shows that the two trips are classified similarly. Thus the learned classes are applicable in other (although similar) locations. 4(d): a topological map derived from our method, see text.

quire that the robot accurately recognize when it is in certain types of space. By combining our classification with odometric data, rough topological maps can be derived. Figure 4(d) shows a topological map derived from 4(a) by dividing the space into regions based on classification. Further processing with loop-closure algorithms and landmark identification techniques (Howard, 2004) can refine these maps into useful tools for autonomous robot navigation.

In addition, control algorithms can be derived from the motor data associated with each class. Firstly, we can use the average movement of the robot in each space class as a first-pass control policy for what the robot should do if it finds itself in that class. Furthermore, we can include the teleoperation data in the training process, so areas that are clustered together not only look similar, but are areas where the robot should behave similarly as well (at least according to the teleoperator). We plan to use this ability to perform robotic learning by demonstration (Nicolescu & Matarić, 2003). After being led through a task by a human teleoperator, a robot can segment the task and associate actions with each segment in an unsupervised

manner. As the task is repeated, more data become available to fine tune the robot’s actions. Standard reinforcement learning techniques can also be applied to allow a human trainer better control.

## 6. Conclusion

This paper presents an extensible method for data-driven discovery of intrinsic classes in robot sensor data. We demonstrate that classes discovered with manifold-learning techniques are more consistently recognizable than those found using PCA. We also show that these classes are reapplicable to new data using out-of-sample techniques. We believe that this technique can provide a basis for future work in autonomous robot operation.

## Acknowledgements

This work was supported in part by the NSF. The authors would also like to thank E. Chris Kern for his support and assistance.

## References

- Bengio, Y., Paiement, J., Vincent, P., Delalleau, O., Roux, N. L., & Ouimet, M. (2004). Out-of-sample extensions for lle, isomap, mds, eigenmaps, and spectral clustering. In *Advances in neural information processing systems 16*. MIT Press.
- Blei, D., Ng, A., & Jordan, M. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*.
- Cormen, T. H., Leiserson, C. E., & Rivest, R. L. (1990). *Introduction to algorithms*. MIT Press/McGraw-Hill.
- Donoho, D., & Grimes, C. (2003). *Hessian eigenmaps: new locally linear embedding techniques for highdimensional data* (Technical Report TR2003-08). Stanford.
- Grudic, G., & Mulligan, J. (2005). Topological mapping with multiple visual manifolds. *Robotics: Science and Systems (R:SS 2005)*.
- Howard, A. (2004). Multi-robot mapping using manifold representations. *IEEE International Conference on Robotics and Automation* (pp. 4198–4203).
- Jenkins, O. C., & Matarić, M. J. (2004). A spatio-temporal extension to isomap nonlinear dimension reduction. *The International Conference on Machine Learning (ICML 2004)* (pp. 441–448).
- Klingspor, V., Morik, K. J., & Rieger, A. D. (1996). Learning concepts from sensor data of a mobile robot. *Machine Learning*, 23, 305–332.
- McLachlan, G. J., & Basford, K. E. (1988). *Mixture models: Inference and applications to clustering*. Marcel Dekker.
- Nicolescu, M. N., & Matarić, M. J. (2003). Natural methods for robot task learning: Instructive demonstrations, generalization and practice. *Joint Conference on Autonomous Agents and Multi-Agent Systems*.
- Stewart, B., Ko, J., Fox, D., & Konolige, K. (2003). The revisiting problem in mobile robot map building: A hierarchical bayesian approach. *The Conference on Uncertainty in Artificial Intelligence*.
- Tenenbaum, J. B., de Silva, V., & Langford, J. C. (2000). A global geometric framework for nonlinear dimensionality reduction. *Science*, 290.
- Thrun, S. (1998). Learning maps for indoor mobile robot navigation. *Artificial Intelligence*, 99, 21–71.
- Tomatis, N., Nourbakhsh, I., & Siegwart, R. (2003). Hybrid simultaneous localization and map building: a natural integration of topological and metric. *Robotics and Autonomous Systems*, 44, 3–14.