



Computer Science and Artificial Intelligence Laboratory
Technical Report

MIT-CSAIL-TR-2010-052

October 13, 2010

**One-Shot Learning with a Hierarchical
Nonparametric Bayesian Model**

Ruslan Salakhutdinov, Josh Tenenbaum, and
Antonio Torralba

One-Shot Learning with a Hierarchical Nonparametric Bayesian Model

Ruslan Salakhutdinov

*Brain and Cognitive Sciences and CSAIL,
Massachusetts Institute of Technology
Cambridge, MA, USA*

RSALAKHU@MIT.EDU

Josh Tenenbaum

*Brain and Cognitive Sciences and CSAIL,
Massachusetts Institute of Technology
Cambridge, MA, USA*

JBT@MIT.EDU

Antonio Torralba

*CSAIL,
Massachusetts Institute of Technology
Cambridge, MA, USA*

TORRALBA@MIT.EDU

Abstract

We develop a hierarchical Bayesian model that learns to learn categories from single training examples. The model transfers acquired knowledge from previously learned categories to a novel category, in the form of a prior over category means and variances. The model discovers how to group categories into meaningful super-categories that express different priors for new classes. Given a single example of a novel category, we can efficiently infer which super-category the novel category belongs to, and thereby estimate not only the new category's mean but also an appropriate similarity metric based on parameters inherited from the super-category. On MNIST and MSR Cambridge image datasets the model learns useful representations of novel categories based on just a single training example, and performs significantly better than simpler hierarchical Bayesian approaches. It can also discover new categories in a completely unsupervised fashion, given just one or a few examples.

1. Introduction

In typical applications of machine classification algorithms, learning curves are measured in tens, hundreds or thousands of training examples. For humans learners, however, the most interesting regime occurs when the training data are very sparse. Just a single example is often sufficient for people to grasp a new category and make meaningful generalizations to novel instances, if not to classify perfectly (Pinker (1999)). Human categorization often asymptotes after just three or four examples (Xu and Tenenbaum (2007); Smith et al. (2002); Kemp et al. (2006); Perfors and Tenenbaum (2009)). Here we present a nonparametric hierarchical Bayesian model that aims to capture this human-like pattern of one-shot learning, and test its performance against several alternatives on two standard benchmark datasets of visual categories.

At a minimum, categorizing an object requires information about the category's mean and variance along each dimension in an appropriate feature space. This is a similarity-based approach,

where the mean represents the category prototype, and the inverse variances (or precisions) correspond to the dimensional weights in a category-specific similarity metric. One-shot learning may seem impossible because a single example provides information about the mean or prototype of the category, but not about the variances or the similarity metric. Giving equal weight to every dimension in a large a priori-defined feature space, or using the wrong similarity metric, is likely to be disastrous.

Our model leverages higher-order knowledge abstracted from previously learned categories to estimate the new category’s prototype as well as an appropriate similarity metric from just one example. These estimates are also improved as more examples are observed. To illustrate, consider how human learners seeing one example of an unfamiliar animal, such as a wildebeest (or gnu), can draw on experience with many examples of ‘horse’, ‘cows’, ‘sheep’, and more familiar related categories. These similar categories have similar prototypes – horses, cows, and sheep look more like each other than like furniture or vehicles – but they also have similar variability in their feature-space representations, or similar similarity metrics: The ways in which horses vary from the ‘horse’ prototype are similar to the ways in which sheep vary from the ‘sheep’ prototype. We may group these similar basic-level categories into an ‘animal’ super-category, which captures these classes’ similar prototypes as well as their similar modes of variation about their respective prototypes. If we can identify the new example of ‘wildebeest’ as belonging to this ‘animal’ super-category, we can transfer an appropriate similarity metric and thereby generalize informatively even from a single example.

For many real-world applications, we must be able to learn tens of thousands of different categories, and to learn new categories building on (and not disrupting) representations of old ones (Bart and Ullman (2005); Biederman (1995)). In these settings, learning from one or a few labeled examples and performing efficient inference will be crucial, and our method is designed to scale up in precisely these ways. A nonparametric prior allows new categories to be formed at any time in either supervised or unsupervised modes, and conjugate distributions allow most parameters to be integrated out analytically for very fast inference.

2. Related Prior Work

Hierarchical Bayesian models have previously been proposed (Kemp et al. (2006); Perfors and Tenenbaum (2009); Heller et al. (2009)) to describe how people learn to learn categories from one or a few examples, or learn similarity metrics, but these approaches were not focused on machine learning settings – large-scale problems with many categories and high-dimensional natural image data. Most similar to our work is Heller et al. (2009)’s account of how people learn dimensional biases in categorization tasks, but in their model, the analog of our super-categories capture only shared covariance of basic-level categories, rather than both means and variances as we do. As we show in our experimental results, this prevents their model from generalizing any better than baseline when given just one or two examples of a novel category.

A large class of models based on hierarchical Dirichlet processes (Teh et al. (2006)) have also been used for transfer learning (Sudderth et al. (2008); Canini and Griffiths (2009)). There are two key difference between our approach and previous applications of HDPs to cross-task transfer or multi-task learning. First, HDPs typically assume a fixed hierarchy of classes for sharing parameters, while we learn the hierarchy in an unsupervised fashion. Second, HDPs are typically given many examples for each category rather than the one-shot learning cases we consider here, and

it is not clear how well they would work for our problems. Recently introduced nested Dirichlet processes can also be used for transfer learning tasks (Rodriguez and Vuppala (2009); Rodriguez et al. (2008)). However, this work assumes a fixed number of classes (or groups) and did not attempt to address one-shot learning problem: their motivation was to use a multilevel nonparametric mixture to capture more complex within-class structure (Rodriguez and Vuppala (2009)). Our approach allows for new categories to be formed at multiple levels of the hierarchy in order to support generalization from few examples.

The multi-level structure of our model is similar to the recently introduced nested Dirichlet process of (Rodriguez and Vuppala (2009); Rodriguez et al. (2008)). However, there are crucial differences both in our goals and our mathematical formulation: Rodriguez and Vuppala (2009) attach class labels to the top level of their hierarchy, while we attach class labels to the bottom level. This allows the upper-level classes in our model to capture super-categories (e.g., ‘animal’ as a superclass for ‘dog’, ‘horse’, ‘cow’) with learned prior knowledge about means and covariances for categories of that kind. This learned prior can be transferred to new categories in order to support generalization from very few examples. In contrast, Rodriguez and Vuppala (2009) do not learn super-categories of labeled classes. Their approach aims to capture more complex within-class structure using their lower-level mixture components which allow learning of nonlinear decision boundaries from large sets of examples (hence their focus on a 2-class, 2-dimensional ‘spiral’ problem or Fisher’s 3-class, 4-dimensional Iris task). A recent hierarchical model of Adams et al. (2011) could also be used for transfer learning tasks. However, this model does not learn hierarchical priors over covariances, which is crucial for transferring an appropriate similarity metric to new basic-level categories. These recently introduced models are complementary to our approach, and can be combined productively, although we leave that as a subject for future work.

There are several related approaches in the computer vision community. A hierarchical topic model for image features (Bart et al. (2008); Sivic et al. (2008)) can discover visual taxonomies in an unsupervised fashion from large datasets but was not designed for one-shot learning of new categories. Congealing methods (Miller et al. (2000)) support one-shot category learning with a hierarchical probabilistic model, but they are designed primarily for black-and-white images using special purpose image representations. Perhaps closest to our work, Fei-Fei et al. (2006) also gave a hierarchical Bayesian model for visual categories, with a prior on the parameters of new categories that was induced from other categories. However their approach is not well-suited as a generic approach to one-shot learning. They learned a single prior shared across all categories and the prior was learned only from three categories, chosen by hand.

3. Hierarchical Bayesian Model

Consider observing a set of N *i.i.d* input feature vectors $\{\mathbf{x}^1, \dots, \mathbf{x}^N\}$, $\mathbf{x}^n \in R^D$. In general, features will be derived from high-dimensional, highly structured data, such as images of natural scenes, in which case the feature dimensionality D can be quite large (e.g. 50,000). For clarity of presentation, let us first assume that our model is presented with a fixed two-level category hierarchy. In particular, suppose that N objects are partitioned into C basic-level (or level-1) categories. We represent such partition by a vector \mathbf{z}^b of length N , each entry of which is $z_n^b \in \{1, \dots, C\}$. We also assume that our C basic-level categories are partitioned into K super-categories (level-2 categories), which we represent by \mathbf{z}^s of length C , with $z_c^s \in \{1, \dots, K\}$. We will relax these assumption later by placing a hierarchical nonparametric prior over the category assignments.

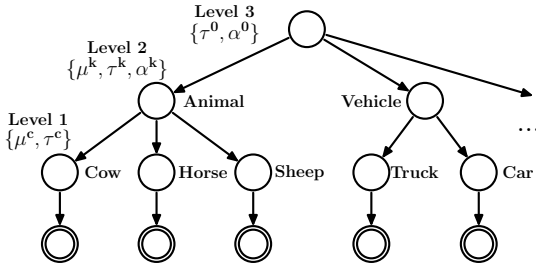


Figure 1: **Left:** Hierarchical Bayesian model that assumes a fixed tree hierarchy for sharing parameters. **Right:** Generative process of the corresponding nonparametric model.

A schematic representation of the overall model is shown in Fig. 1, left panel, and we now formalize it more precisely. For any basic-level category c , the distribution over the observed feature vectors is assumed to be Gaussian with a category-specific mean μ^c and a category-specific *diagonal* precision matrix, whose entries are $\{\tau_d^c\}_{d=1}^D$. The distribution takes the following product form:

$$P(\mathbf{x}^n | z_n^b = c, \theta^1) = \prod_{d=1}^D \mathcal{N}(x_d^n | \mu_d^c, 1/\tau_d^c), \quad (1)$$

where $\mathcal{N}(x | \mu, 1/\tau)$ denotes a Gaussian distribution with mean μ and precision τ and $\theta^1 = \{\mu^c, \tau^c\}_{c=1}^C$ denotes the level-1 category parameters. We next place a conjugate Normal-Gamma prior over $\{\mu^c, \tau^c\}$. Let $k = z_c^s$, i.e. let the level-1 category c belong to level-2 category k , where $\theta^2 = \{\mu^k, \tau^k, \alpha^k\}_{k=1}^K$ denote the level-2 parameters. Then: $P(\mu^c, \tau^c | \theta^2, \mathbf{z}^s) = \prod_{d=1}^D P(\mu_d^c, \tau_d^c | \theta^2, \mathbf{z}^s)$, where for each dimension d we have:

$$\begin{aligned} P(\mu_d^c, \tau_d^c | \theta^2) &= P(\mu_d^c | \tau_d^c, \theta^2) P(\tau_d^c | \theta^2) = \\ &= \mathcal{N}(\mu_d^c | \mu_d^k, 1/(\nu \tau_d^c)) \Gamma(\tau_d^c | \alpha_d^k, \alpha_d^k / \tau_d^k). \end{aligned} \quad (2)$$

Note that our parameterization of the Gamma density is in terms of its shape α^k and mean τ^k parameters:

$$\Gamma(\tau | \alpha^k, \alpha^k / \tau^k) = \frac{(\alpha^k / \tau^k)^{\alpha^k}}{\Gamma(\alpha^k)} \tau^{\alpha^k - 1} \exp\left(-\tau \frac{\alpha^k}{\tau^k}\right). \quad (3)$$

Such a parameterization is more interpretable and is much easier to work with, since $E[\tau] = \tau^k$. In particular, from Eq. 2, we can easily derive that $E[\mu^c] = \mu^k$ and $E[\tau^c] = \tau^k$. This gives our model a very intuitive interpretation: the expected values of the basic level-1 parameters θ^1 are given by the corresponding level-2 parameters θ^2 . The parameter α^k further controls the variability of τ^c around its mean, i.e. $\text{Var}[\tau^c] = (\tau^k)^2 / \alpha^k$. For the level-2 parameters θ^2 , we shall assume the following conjugate priors:

$$P(\mu_d^k) = \mathcal{N}(\mu_d^k | 0, 1/\tau^0), \quad P(\alpha_d^k | \alpha^0) = \text{Exp}(\alpha_d^k | \alpha^0), \quad P(\tau_d^k | \theta^0) = \text{IG}(\tau_d^k | a^0, b^0),$$

where $\text{Exp}(x | \alpha)$ denotes an exponential distribution with rate parameter α , and $\text{IG}(x | \alpha, \beta)$ denotes an inverse-gamma distribution with shape parameter α and scale parameter β . We further place a diffuse Gamma prior $\Gamma(1, 1)$ over hyperparameters α^0 and τ^0 . Throughout our experimental results, we also set $a^0 = 1$ and $b^0 = 1$.

3.1 Modelling the number of super-categories

So far we have assumed that our model is presented with a two-level partition $\mathbf{z} = \{\mathbf{z}^s, \mathbf{z}^b\}$. This model corresponds to a standard hierarchical Bayesian model that assumes a fixed hierarchy for sharing parameters. If, however, we are not given any level-1 or level-2 category labels, we need to infer the distribution over the possible category structures. We place a nonparametric two-level nested Chinese Restaurant Prior (CRP) (Blei et al. (2003, 2010)) over \mathbf{z} , which defines a prior over tree structures and is flexible enough to learn arbitrary hierarchies. The main building block of the nested CRP is the Chinese restaurant process, a distribution on partition of integers. Imagine a process by which customers enter a restaurant with an unbounded number of tables, where the n^{th} customer occupies a table k drawn from:

$$P(z_n = k | z_1, \dots, z_{n-1}) = \begin{cases} \frac{n^k}{n-1+\gamma} & n^k > 0 \\ \frac{\gamma}{n-1+\gamma} & k \text{ is new} \end{cases}, \quad (4)$$

where n^k is the number of previous customers at table k and γ is the concentration parameter.

The Nested CRP, nCRP(γ), extends CRP to nested sequence of partitions, one for each level of the tree. In this case each observation n is first assigned to the super-category z_n^s using Eq. 4. Its assignment to the basic-level category z_n^b , that is placed under a super-category z_n^s , is again recursively drawn from Eq. 4 (for details see Blei et al. (2010)). For our model, a two-level nested CRP allows flexibility of having a potentially unbounded number of super-categories as well as an unbounded of basic-level categories placed under each super-category. Finally, we also place a Gamma prior $\Gamma(1, 1)$ over γ . The full generative model is given in Fig. 1, right panel. Unlike in many conventional hierarchical Bayesian models, here we infer both the model parameters as well as the hierarchy for sharing those parameters.

Our model can be readily used in unsupervised or semi-supervised modes, with varying amounts of label information. Here we focus on two settings. First, we assume basic-level category labels have been given for all examples in a training set, but no super-category labels are available. We must infer how to cluster basic categories into super-categories at the same time as we infer parameter values at all levels of the category hierarchy. The training set includes many examples of familiar basic categories but only one (or few) example for a novel class. The challenge is to generalize the new class intelligently from this one example by inferring which super-category the new class comes from and exploiting that super-category’s implied priors to estimate the new class’s prototype and similarity metric most accurately. Second, we consider a similar labeled training set but now the test set consists of many unlabeled examples from an unknown number of basic-level classes – including both familiar and novel classes. This reflects the problem of “unsupervised category learning”: How to discover when the model has encountered novel categories, and how to break up new instances into categories in an intelligent way that exploits knowledge abstracted from a hierarchy of more familiar categories.

4. Inference

Inferences about model parameters at all levels of hierarchy can be made by running a Markov chain whose stationary distribution is the posterior distribution over the model parameters. When the tree structure \mathbf{z} of the model is not given, the inference process will alternate between fixing \mathbf{z} while sampling the space of model parameters θ and fixing θ while sampling category assignments. The

use of conjugate priors allows for an efficient Gibbs sampler.

Sampling level-1 parameters: Given level-2 parameters θ^2 and \mathbf{z} , the conditional distribution $P(\mu^c, \tau^c | \theta^2, \mathbf{z}, \mathbf{x})$ is Normal-Gamma (Eq. 2), which allows us to easily sample level-1 parameters $\{\mu^c, \tau^c\}$. Making inferences about precision terms in our model can be thought of as learning a category-specific similarity metric. Note that the conditional distribution over θ^1 factorizes into the product of conditional distributions over the parameters of individual categories:

$$P(\{\mu^c, \tau^c\}_{c=1}^C | \theta^2, \mathbf{z}) = \prod_{c=1}^C \prod_{d=1}^D P(\mu_d^c, \tau_d^c | \theta^2, \mathbf{z}).$$

We can therefore easily speed up our inference process by sampling from these conditional distributions in parallel. The speedup could be substantial as the number of the basic-level categories becomes large.

Sampling level-2 parameters: Given \mathbf{z} , θ^1 , and θ^3 , the conditional distributions over the mean μ^k and precision τ^k take Gaussian and Inverse-Gamma forms. The only complicated step involves sampling α^k that control the variation of the precision term τ^c around its mean (Eq. 3). The conditional distribution over α^k cannot be computed in closed form and is proportional to:

$$p(\alpha^k) \propto \frac{(\alpha^k / \tau^k)^{\alpha^k n_k}}{\Gamma(\alpha^k)^{n_k}} \exp\left(-\alpha^k \left(\alpha^0 + S^k / \tau^k - T^k\right)\right),$$

where $S^k = \sum_{c:z(c)=k} \tau^c$ and $T^k = \sum_{c:z(c)=k} \log(\tau^c)$. For large values of α^k the density, specified by Eq. 5, is similar to a Gamma density (Wiper et al. (2001)). We therefore use Metropolis-Hastings with a proposal distribution given by the Gamma density. In particular, we generate a new candidate

$$\alpha^* \sim Q(\alpha^* | \alpha^k) \quad \text{with} \quad Q(\alpha^* | \alpha^k) = \Gamma(\alpha^* | t, t / \alpha^k)$$

and accept it with M-H rule. In all of our experiments we use $t = 3$, which gave an acceptance probability of about 0.6. Finally, sampling level-3 parameters is similar to sampling level-2 parameters.

Sampling assignments \mathbf{z} : Given model parameters $\theta = \{\theta^1, \theta^2\}$, combining the likelihood term with the nCRP(γ) prior, the posterior over the assignment \mathbf{z}_n can be calculated as follows:

$$p(\mathbf{z}_n | \theta, \mathbf{z}_{-n}, \mathbf{x}^n) \propto p(\mathbf{x}^n | \theta, \mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{-n}), \quad (5)$$

where \mathbf{z}_{-n} denotes variables \mathbf{z} for all observations other than n . We can further exploit the conjugacy in our hierarchical model when computing the probability of creating a new basic-level category. Using the fact the Normal-Gamma prior $p(\mu^c, \tau^c)$ is the conjugate prior of a normal distribution, we can easily compute the following marginal likelihood:

$$p(\mathbf{x}^n | \theta^2, \mathbf{z}_n) = \int_{\mu^c, \tau^c} p(\mathbf{x}^n, \mu^c, \tau^c | \theta^2, \mathbf{z}_n) = \int_{\mu^c, \tau^c} p(\mathbf{x}^n | \mu^c, \tau^c) p(\mu^c, \tau^c | \theta^2, \mathbf{z}_n).$$

Integrating out basic-level parameters θ^1 lets us more efficiently sample over the tree structures¹. When computing the probability of placing \mathbf{x}^n under a newly created super-category, its parameters are sampled from the prior.

5. One-shot Learning

Consider observing a single new instance \mathbf{x}^* of a *novel category* c^* ². Conditioned on the current setting of the level-2 parameters θ^2 and our current tree structure \mathbf{z} , we can first infer which super-category the novel category should belong to, i.e. we can compute the posterior distribution over the assignments \mathbf{z}_c^* using Eq. 5. We note that our new category can either be placed under one of the existing super-categories, or create its own super-category, if it is sufficiently different from all of the existing super-categories.

Given an inferred assignment³ \mathbf{z}_c^* and using Eq. 2, we can infer the posterior mean and precision terms (or similarity metric) $\{\mu^*, \tau^*\}$ for our novel category. We can now test the ability of the HB model to generalize to new instances of a novel category by computing the conditional probability that a new test input \mathbf{x}^t belongs to a novel category c^* :

$$p(c^*|\mathbf{x}^t) = \frac{p(\mathbf{x}^t|\mathbf{z}_c^*)p(\mathbf{z}_c^*)}{\sum_{\mathbf{z}} p(\mathbf{x}^t|\mathbf{z})p(\mathbf{z})},$$

where the prior is given by the nCRP(γ) and the likelihood takes form:

$$\log p(\mathbf{x}^t|c^*) = \frac{1}{2} \sum_d \log(\tau_d^*) - \frac{1}{2} \sum_d \tau_d^* (x_d^t - \mu_d^*)^2 + C,$$

where C is a constant that does not depend on the parameters. Observe that the relative importance of each feature in determining the similarity is proportional to the category-specific precision of that feature. Features that are salient, or have higher precision, within the corresponding category contribute more to the overall similarity of an input.

It is informative to better understand what kind of similarity metric transfer our model is performing based on a single example of a novel category. Let us examine the posterior mean and precision of the d^{th} feature. The inferred mean is given by:

$$\mu_d^* = \frac{\nu \mu_d^k + x_d^*}{\nu + 1}.$$

The parameter ν controls the blend between an observation and the mean of the global super-category. In all of our experiments we set $\nu = 0.1$. Inferred precision (or similarity metric) is given by the Gamma density, whose expected value is equal to α_d^*/β_d^* . Provided α_d^* is large (in our experiments α_d^* is typically much larger than 1), so that $\alpha_d^*/(\alpha_d^* + 0.5) \approx 1$, the expected value of the precision parameter takes the following form:

$$\mathbb{E}[\tau_d^*] = \frac{\tau_d^k}{\frac{\alpha_d^k}{\alpha_d^k + 0.5} \left(1 + \frac{\nu}{1+\nu} \frac{\tau_d^k}{\alpha_d^k + 0.5} (x_d^* - \mu_d^k)^2\right)} \approx \frac{\tau_d^k}{1 + \frac{1}{\alpha_d^k + 0.5} \frac{\nu}{1+\nu} (\tau_d^k (x_d^* - \mu_d^k)^2)}. \quad (6)$$

-
1. In the supervised case, inference is simplified by only considering which super-category each basic-level category is assigned to.
 2. Observing several examples of a new category is treated similarly.
 3. In our experiments, for faster inference, we simply compute the most probable assignment $\mathbf{z}_c^* = \operatorname{argmax} p(\mathbf{z}_c^*|\theta^2, \mathbf{z}_{-c}^*, \mathbf{x}^*)$ with parameters θ^1 integrated out.

The above formula has a very intuitive interpretation. If a new observation x_d^* is relatively close to the mean μ_d^k of the global super-category (with distances scaled by the precision terms), then the expected similarity metric will be close to the similarity metric defined by the global super-category. Otherwise, the model will become more uncertain about the value of the observed feature, and this feature will contribute less to the overall similarity of an input.

6. Experimental results

We now present experimental results on the MNIST handwritten digit and MSR Cambridge object recognition image datasets. During the inference step, we run our hierarchical Bayesian (HB) model for 200 full Gibbs sweeps, which was sufficient to reach convergence and obtain good performance. We normalize input vectors to zero mean and scale the entire input by a single number to make the average feature variance be one.

In all of our experiments, we compare performance of the HB model to the following four baseline models. The first model, called “Euclidean”, uses a Euclidean metric, i.e. all precision terms are set to one and are never updated, hence all dimensions are equally important for all categories. The second model, that we call “HB-Flat”, always uses a single super-category. When presented with a single example of a new category, HB-Flat will inherit a similarity metric that is shared by all existing categories. This approach, similar in spirit to Fei-Fei et al. (2006), could potentially identify a set of useful features common to all categories and learn to ignore irrelevant features. Our third baseline model, which we refer to as “HB-Var”, is similar in spirit to the approach of Heller et al. (2009) and is based on clustering only covariance matrices without taking into account the means of the super-categories. Our final baseline model, called “MLE” ignores hierarchical Bayes altogether and estimates a category-specific mean and precision from sample averages. If a category contains only one example, the model resorts to using the Euclidean metric. Finally, we also compare to the the “Oracle” model that always uses the correct, instead of inferred, similarity metric.

6.1 MNIST dataset

The MNIST dataset contains 60,000 training and 10,000 test images of ten handwritten digits (zero to nine), with 28×28 pixels. For our experiments, we randomly choose 1000 training and 1000 test images (100 images per class). We work directly in the pixel space because all handwritten digits were already properly aligned. Fig. 2 shows a typical partition over the basic level categories, along with corresponding mean and similarity metrics, that our model discovers.

We first study the ability of the HB model to generalize from a single training example of handwritten digit ‘nine’. To this end, we trained the HB model on 900 images (100 images of each of zero-to-eight categories), while withholding all images that belong to category ‘nine’. Given a single new instance of a novel ‘nine’ category our model is able to discover that the new category is more like categories that contain images of seven and four, and hence this novel category can inherit the mean and the similarity metric, shared by categories ‘seven’ and ‘four’. Fig. 2 precisely illustrates the kind of transfer our model is performing. The transferred similarity metric allows HB model to generalize much better to new instances of a novel category.

Figure 3 and Table 1 further quantifies performance using the area under the ROC curve (AU-ROC) for classifying 1000 test images as belonging to the ‘nine’ vs. all other categories (an area of 0.5 corresponds to the classifier that makes random predictions). The HB model achieves an AUROC of 0.81, considerably outperforming HB-Flat, HB-Var, Euclidean, and MLE that achieve

Table 1: Performance results using the area under the ROC curve (AUROC) on the MNIST dataset. The rightmost Average panel shows results averaged over all 10 categories, using leave-one-out test format.

Model	Category: Digit 9				Category: Digit 6				Average			
	1 ex	2 ex	4 ex	20 ex	1 ex	2 ex	4 ex	20 ex	1 ex	2 ex	4 ex	20 ex
HB	0.81	0.85	0.88	0.90	0.85	0.89	0.92	0.97	0.85	0.88	0.90	0.93
HB-Flat	0.71	0.77	0.84	0.90	0.73	0.79	0.88	0.97	0.74	0.79	0.86	0.93
HB-Var	0.72	0.81	0.86	0.90	0.72	0.83	0.90	0.97	0.75	0.82	0.89	0.93
Euclidean	0.70	0.73	0.76	0.80	0.74	0.77	0.82	0.86	0.72	0.76	0.80	0.83
Oracle	0.87	0.89	0.90	0.90	0.95	0.96	0.96	0.97	0.90	0.92	0.92	0.93
MLE	0.69	0.75	0.83	0.90	0.72	0.78	0.87	0.97	0.71	0.77	0.84	0.93

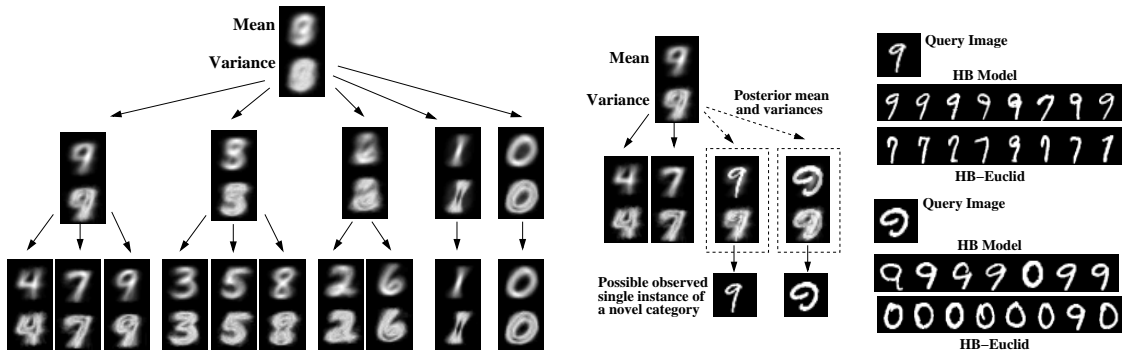


Figure 2: MNIST dataset. **Left:** A typical partition over the 10 basic-level categories discovered by the HB model. Top panels display means and bottom panels display variances (white encodes larger values). **Middle:** Transfer of similarity metric based on a single example of a novel 'nine' category. **Right:** Retrieval results: Top eight most similar images retrieved from the test set of 1000 images corresponding to 10 categories. Note that due to metric transfer, the HB model is able to avoid mistakes made by the Euclidean model.

an AUROC of 0.71, 0.72, 0.70, and 0.69 respectively. Table 1 further reveals with just a single example, the HB model performs comparable to both HB-Flat and MLE that use 4 examples. This result clearly demonstrates that the HB model is able to successfully transfer appropriate metrics from previously learned categories. Moreover, with just four examples, the HB model is able to achieve performance close to that of the Oracle model. This is in sharp contrast to HB-Flat, MLE and Euclidean models, that even with four examples perform far worse.

Finally, Fig. 4, left panel, shows that using the wrong 'two' similarity metric for the novel 'nine' category can significantly deteriorate model's prediction accuracy. Indeed, the model performs worse than the Euclidean model which does not learn a similarity metric at all. This example clearly demonstrates that our model learns meaningful super-categories and is indeed able to transfer good similarity metric.

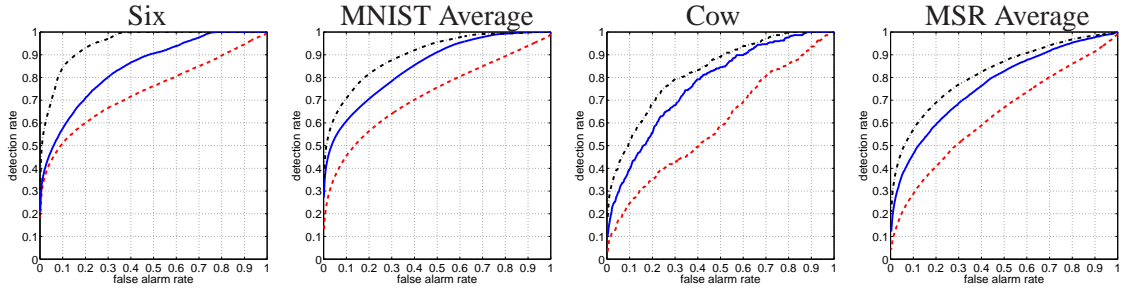


Figure 3: ROC curves for classifying test images belonging to a novel category vs. the rest based on observing a *single* instance of a new category. Three curves represent Euclidean (lower red), HB (middle blue), and Oracle (upper red) models. All curves are averaged over 100 (MNIST) or 25 (MSR) possible examples corresponding to a novel category. The ‘Average’ represents results averaged over all 10 (MNIST) or 24 (MSR) categories, using leave-one-out test format.

6.2 MSR Cambridge Dataset

We now present results on a considerably more difficult MSR Cambridge dataset, that contains MSR Cambridge dataset⁴ contains images of 24 different categories. Figure 4, right panel, shows 24 basic-level categories along with a typical partition that our model discovers, where many super-categories contain semantically similar basic-level categories. For all experiments we use 15 and 25 images per category for testing and training.

6.2.1 DETAILS OF IMAGE REPRESENTATION

We use a simple “texture-of-textures” framework for constructing image features. In particular, we use the algorithm of DeBonet and Viola (1997) that extracts 46,875 very specific features that respond to edge orientation, color, texture, and many local properties at multiple scales. Each image is convolved with a set of 25 local linear filters including bars and oriented edges. Filter response is then rectified by squaring and further downsampled by a factor of two. Convolution, rectification and downsampling is repeated two more times, producing a vector of size $25^3 = 15,625$. The same operation is then applied to each of the three RGB channels, yielding a total of 46,875 features. We emphasize that presented model is not restricted to using this type of features and we expect that performance could potentially be improved by using more advanced features. For simple comparison, we also present results of the HDP model, where each image was represented as a bag of 2000 visual words derived from texture-of-textures features.

6.2.2 RESULTS

We first tested the ability of our model to generalize from a single image of a cow. Similar to the experiments on the MNIST dataset, we first train the HB model on images corresponding to 23 categories, while withholding all images of cows. In general, our model is able to discover that the new ‘cow’ category is more like the ‘sheep’ category, as opposed to categories that contain images of cars, or forks, or buildings. This allows the new ‘cow’ category to inherit sheep’s similarity metric.

4. Available at <http://research.microsoft.com/en-us/projects/objectclassrecognition/>

Table 2: Performance results categories using the area under the ROC curve (AUROC) on the MSR dataset. The rightmost Average panel shows results averaged over all 24 categories, using leave-one-out test format.

Model	Category: Cow				Category: Flower				Average			
	1 ex	2 ex	4 ex	20 ex	1 ex	2 ex	4 ex	20 ex	1 ex	2 ex	4 ex	20 ex
HB	0.77	0.81	0.84	0.89	0.71	0.75	0.78	0.81	0.76	0.80	0.84	0.87
HB-Flat	0.62	0.69	0.80	0.89	0.59	0.64	0.75	0.81	0.65	0.71	0.78	0.87
HB-Var	0.61	0.73	0.83	0.89	0.60	0.68	0.77	0.81	0.64	0.74	0.81	0.87
Euclidean	0.59	0.61	0.63	0.66	0.55	0.59	0.61	0.64	0.63	0.66	0.69	0.71
Oracle	0.83	0.84	0.87	0.89	0.77	0.79	0.80	0.81	0.82	0.84	0.86	0.87
MLE	0.58	0.64	0.78	0.89	0.55	0.62	0.72	0.81	0.62	0.67	0.77	0.87
HDP	0.64	0.71	0.82	0.90	0.61	0.67	0.77	0.83	0.67	0.72	0.79	0.89

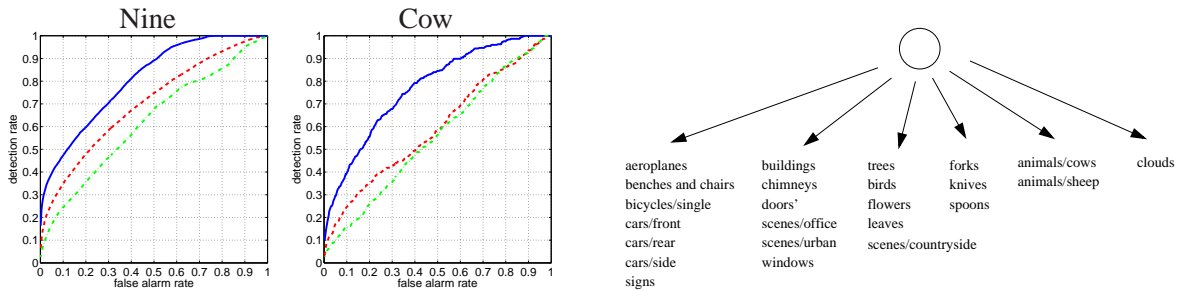


Figure 4: Results when using the wrong similarity metric. Three curves represent HB (top blue), Euclidean (middle red), and the model that uses the wrong similarity metric (lower green). **Left:** Using the ‘two’ similarity for the novel ‘nine’ category. **Middle:** Using the ‘fork’ similarity for the novel ‘cow’ category. **Right:** MSR Cambridge dataset: A typical partition over the 24 basic-level categories discovered by the HB model.

Figure 3 and Table 2 show that the HB model, based on a single example of cow, achieves an AUROC of 0.77. This is compared to an AUROC of only 0.64, 0.62, 0.61, 0.59, and 0.58 achieved by the HDP, HB-Flat, HB-Var, Euclidean, and MLE models. As the number of training examples increases, the HB model still consistently outperforms all the other methods. Similar to the results on the MNIST dataset, the HB model with just one example performs comparably with the HB-Flat and MLE models that make use of four examples. This clearly demonstrates that the HB model is able to successfully transfer metric from similar categories. In particular, the improvement over HDP, Euclidean, HB-Flat HB-Var, and MLE models is particularly striking when learning with only one example. With 20 examples, however, the part-based HDP model slightly outperforms our HB model.

Fig. 5 further displays retrieval results based on a single image of a cow. As expected, the HB model performs much better compared to the simple Euclidean model that does not learn a similarity metric. Fig. 5, right panel, further shows an example where the HB model fails, since it retrieves many images of the wrong ‘sheep’ category. This is in sharp contrast to the Euclidean model, that tends to retrieve images from very unrelated categories.

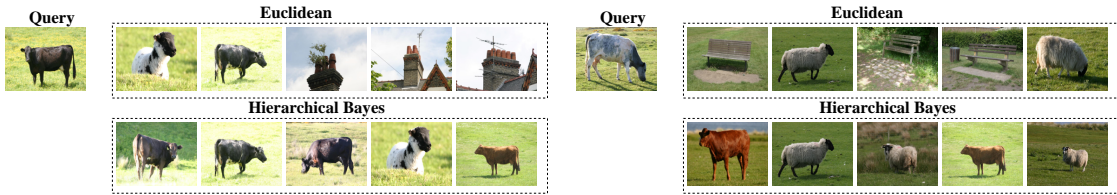


Figure 5: Retrieval results based on observing a single example of cow. Top five most similar images were retrieved from the test set, containing 360 images corresponding to 24 categories.

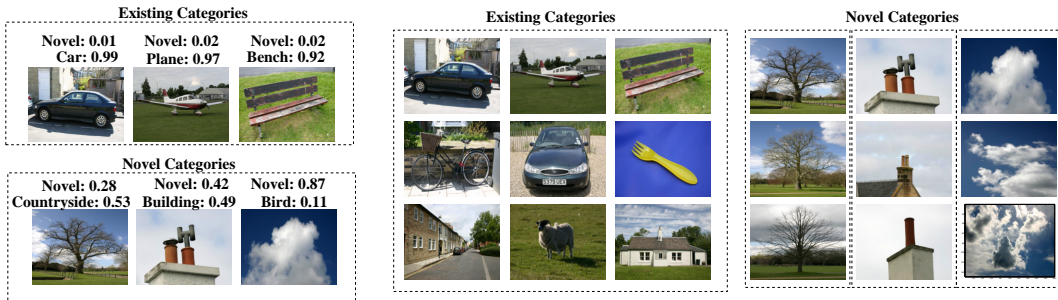


Figure 6: Unsupervised category discovery. **Left:** Six representative test images, sorted by the posterior probability of forming a novel category. **Right:** The model is presented with 18 unlabeled test images. After running a Gibbs sampler for 100 steps, the model correctly places nine ‘familiar’ images in nine different basic-level categories, while also correctly forming three novel basic-level categories with three examples each.

6.3 Unsupervised Category Discovery

Another key advantage of the hierarchical nonparametric Bayesian model is its ability to infer category structure in an unsupervised fashion, discovering novel categories at both levels 1 and 2 of the hierarchy. We explored the HB model’s category discovery ability by training on labeled examples of 21 basic-level MSR categories, leaving out clouds, trees, and chimneys. We then provided six test images: one in each of the three unseen categories and one in each of three familiar basic-level categories (car, airplane, bench). For each test image, using Eq. 6, we can easily compute the posterior probability of forming a new basic-level category. Figure 6, left panel, shows six representative test images, sorted by the posterior probability of forming a novel category. The model correctly identifies the car, the airplane and the bench as belonging to familiar categories, and places much higher probability on forming novel categories for the other images. With only one unlabeled example of these novel classes, the model still prefers two of them in familiar categories: the ‘tree’ is interpreted as an atypical example of ‘countryside’ while the ‘chimney’ is classified as an atypical ‘building’. However, the model can correctly discover novel categories given only a little more unlabeled data.

With 18 unlabeled test images (see Fig. 6, right panel), after running a Gibbs sampler for 100 steps, the model correctly places nine ‘familiar’ images in nine different basic-level categories, while also correctly forming three novel basic-level categories with three examples each. Most interestingly, these new basic-level categories are placed at the appropriate level of the category hierarchy: the novel ‘tree’ category is correctly placed under the super-category containing ‘leaves’ and ‘countrysides’; the novel ‘chimney’ category is placed together with ‘buildings’ and ‘doors’;

while ‘clouds’ category is placed in its own super-category – all consistent with the hierarchy we originally found from a fully labeled training set with many examples in each of these three classes (see Fig. 4). Other models we tried for this unsupervised task, including HB-Euclid and HB-Flat, perform much worse; they confuse ‘chimneys’ with ‘cows’ (see Fig. 5) and ‘trees’ with ‘country-sides’.

7. Conclusions

In this paper we developed a hierarchical nonparametric Bayesian model for learning a novel category based on a single training example. Our experimental results, conducted on realistic datasets, further demonstrate that our model is able to effectively transfer appropriate similarity metric from the previously learned categories to a novel category based on just observing a single example.

There are several key advantages to our model. First, due to efficient Gibbs moves that can exploit conjugacy, the model can be efficiently trained. Many of the Gibbs updates can be run in parallel, which will allow our model to potentially handle a large number of basic-level categories. Second, the model is able to discover meaningful super-categories and be able to form coherent novel categories. Finally, given a single example of a novel category, the model is able to quickly infer which super-category the new basic-level category should belong to. This in turns allows us to efficiently infer the appropriate similarity metric for this novel category.

Acknowledgments

We acknowledge the financial support from NSERC, Shell, and NTT Communication Sciences Laboratory.

References

- R. Adams, Z. Ghahramani, and M. Jordan. Tree-structured stick breaking processes for hierarchical data. In *To appear in NIPS*. MIT Press, 2011.
- E. Bart and S. Ullman. Cross-generalization: Learning novel classes from a single example by feature replacement. In *CVPR*, pages 672–679, 2005.
- E. Bart, I. Porteous, P. Perona, and M. Welling. Unsupervised learning of visual taxonomies. In *CVPR*, pages 1–8, 2008.
- I. Biederman. Visual object recognition. *An Invitation to Cognitive Science*, 2:152–165, 1995.
- David M. Blei, Thomas L. Griffiths, Michael I. Jordan, and Joshua B. Tenenbaum. Hierarchical topic models and the nested chinese restaurant process. In *NIPS*. MIT Press, 2003.
- David M. Blei, Thomas L. Griffiths, and Michael I. Jordan. The nested chinese restaurant process and bayesian nonparametric inference of topic hierarchies. *J. ACM*, 57(2), 2010.
- Kevin R. Canini and Thomas L. Griffiths. Modeling human transfer learning with the hierarchical dirichlet process. In *NIPS 2009 workshop: Nonparametric Bayes*, 2009.
- Jeremy S. DeBonet and Paul A. Viola. Structure driven image database retrieval. In Michael I. Jordan, Michael J. Kearns, and Sara A. Solla, editors, *NIPS*. The MIT Press, 1997.

- Li Fei-Fei, R. Fergus, and P. Perona. One-shot learning of object categories. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 28(4):594–611, April 2006.
- K. Heller, A. Sanborn, and N. Chater. Hierarchical learning of dimensional biases in human categorization. In *NIPS*, 2009.
- C. Kemp, A. Perfors, and J. Tenenbaum. Learning overhypotheses with hierarchical Bayesian models. *Developmental Science*, 10(3):307–321, 2006.
- E. G. Miller, N. E. Matsakis, and P. A. Viola. Learning from one example through shared densities on transforms. In *CVPR*, pages 464–471, 2000.
- A. Perfors and J.B. Tenenbaum. Learning to learn categories. In *31st Annual Conference of the Cognitive Science Society*, pages 136–141, 2009.
- S. Pinker. *How the Mind Works*. W.W. Norton, 1999.
- A. Rodriguez and R. Vuppala. Probabilistic classification using Bayesian nonparametric mixture models. Technical Report, 2009.
- A. Rodriguez, D. Dunson, and A. Gelfand. The nested Dirichlet process. *Journal of the American Statistical Association*, 103:1131–1144, 2008.
- J. Sivic, B. C. Russell, A. Zisserman, W. T. Freeman, and A. A. Efros. Unsupervised discovery of visual object class hierarchies. In *CVPR*, pages 1–8, 2008.
- L.B. Smith, S.S. Jones, B. Landau, L. Gershkoff-Stowe, and L. Samuelson. Object name learning provides on-the-job training for attention. *Psychological Science*, pages 13–19, 2002.
- E. B. Sudderth, A. Torralba, W. T. Freeman, and A. S. Willsky. Describing visual scenes using transformed objects and parts. *International Journal of Computer Vision*, 77(1-3):291–330, 2008.
- Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- Michael Wiper, David Rios Insua, and Fabrizio Ruggeri. Mixtures of gamma distributions with applications. *Journal of Computational and Graphical Statistics*, 10(3), September 2001.
- Fei Xu and Joshua B. Tenenbaum. Word learning as bayesian inference. *Psychological Review*, 114(2), 2007.

