

# Gaussian Processes for Object Categorization

Ashish Kapoor · Kristen Grauman · Raquel Urtasun · Trevor Darrell

Received: 22 July 2008 / Accepted: 1 July 2009

© The Author(s) 2009. This article is published with open access at Springerlink.com

**Abstract** Discriminative methods for visual object category recognition are typically non-probabilistic, predicting class labels but not directly providing an estimate of uncertainty. Gaussian Processes (GPs) provide a framework for deriving regression techniques with explicit uncertainty models; we show here how Gaussian Processes with covariance functions defined based on a Pyramid Match Kernel (PMK) can be used for probabilistic object category recognition. Our probabilistic formulation provides a principled way to learn hyperparameters, which we utilize to learn an optimal combination of multiple covariance functions. It also offers confidence estimates at test points, and naturally allows for an active learning paradigm in which points are optimally selected for interactive labeling. We show that with an appropriate combination of kernels a significant boost in classification performance is possible. Further, our experiments indicate the utility of active learning with probabilistic predictive models, especially when the amount of training data labels that may be sought for a category is ultimately very small.

**Keywords** Object recognition · Gaussian process · Kernel combination · Active learning

## 1 Introduction

Object categorization is a fundamental problem in image understanding. It remains a challenging learning task given both the variability of images that objects from the same class can produce, as well as the substantial expense of providing high quality image annotations needed to train accurate models. Discriminative methods for visual category learning have yielded promising results in recent years, including various approaches based on support vector machines or nearest neighbor classification (Grauman and Darrell 2005; Zhang et al. 2006; Wallraven et al. 2003; Nister and Stewenius 2006; Lazebnik et al. 2006; Varma and Ray 2007; Bosch et al. 2007; Frome et al. 2007; Kumar and Sminchisescu 2007). However, such methods typically are not explicitly probabilistic, which makes them inadequate when estimates of uncertainty are required. At the same time, probabilistic generative methods that attempt to directly model the joint distribution of object classes and their features—though appealing for their ability to estimate uncertainty during inference—can be impractical for image recognition applications due to the complexity of representing the data's underlying density.

In this work we provide a probabilistic discriminative approach to object categorization, with the goal of exercising the advantages of both types of methods. We introduce a new Gaussian Process (GP) regression method for object category recognition using a local feature correspondence kernel. Local feature-based object recognition has several important advantages, including invariance to various translational, rotational, affine and photometric transformations,

---

A. Kapoor (✉)

Microsoft Research, Redmond, WA 98052, USA  
e-mail: [akapoor@microsoft.com](mailto:akapoor@microsoft.com)

K. Grauman

University of Texas at Austin, Austin, TX 78712, USA  
e-mail: [grauman@cs.utexas.edu](mailto:grauman@cs.utexas.edu)

R. Urtasun · T. Darrell

UC Berkeley EECS & ICSI, Berkeley, CA 94720, USA

R. Urtasun

e-mail: [rurtasun@csail.mit.edu](mailto:rurtasun@csail.mit.edu)

T. Darrell

e-mail: [trevor@eecs.berkeley.edu](mailto:trevor@eecs.berkeley.edu)

and robustness to partial occlusions. Our method is based on a GP with a covariance function derived from a Pyramid Match Kernel (Grauman and Darrell 2005), which offers an efficient approximation to a partial-match distance function and can therefore handle outliers and occlusions. Our model offers some of the known benefits of probabilistic techniques, while still maintaining the power of a discriminative learner. In particular, we show how it enables both *active* visual category learning, as well as learning from multiple image feature sources with an optimal combination of covariance functions.

Collecting training data for large-scale image category models is a potentially expensive process. While certain categories may have a large number of training images available, many more will have relatively few. A number of ingenious schemes have been developed to obtain labeled data from people performing other tasks (e.g., von Ahn et al. 2006; von Ahn and Dabbish 2004), or directly labeling objects in images (<http://labelme.csail.mit.edu/>). To make the most of scarce human labeling resources it is imperative to carefully select points for user labeling. The paradigm of active learning has been introduced in the machine learning community to address this issue (Freund et al. 1997; Tong and Koller 2000; McCallum and Nigam 1998; Muslea et al. 2002; Zhu et al. 2003); with an active learning method, generally new test points are selected so as to minimize the model entropy.

GPs have received limited attention in the computer vision literature to date perhaps due to the fact that they are conventionally limited to modest amounts of training data: the learning complexity is  $O(n^3)$ , cubic in the number of training examples. While recent advances in sparse GPs are promising (e.g., Lawrence et al. 2002; Shen et al. 2006; Snelson and Ghahramani 2006; Urtasun and Darrell 2008), we focus here on the case of active learning with relatively small numbers of labeled examples (10–100), which is feasible with existing implementations. In this realm, we show that active learning provides significantly more accurate estimates per labeled point than does a conventional random selection of training points.

Specific choices made regarding image representations and kernel parameters can greatly influence a classifier's potential. Even within the domain of local image features and matching kernels, a variety of alternative interest point detectors, descriptors, match criteria, and feature space quantization strategies are available. Rather than require a user to decide *a priori* which particular items will define the GP's covariance function, we show how to automatically optimize the combination of kernels for the recognition task using the GP marginal likelihood function. As a result, one can compute a set of potential kernels using a variety of local feature types, and then directly learn a weight for each such that the final combination is highly discriminative. While

recent work has considered multiple kernel learning (Varma and Ray 2007; Kumar and Sminchisescu 2007) and cross-validation approaches (Bosch et al. 2007) to combine image feature types within SVM classifiers, to our knowledge our approach is the first to consider kernel combinations in a probabilistic setting.

The three main contributions of this paper are (1) a probabilistic discriminative category recognition scheme based on a Gaussian Process prior with a covariance function defined using the Pyramid Match Kernel, (2) the introduction of an active learning paradigm for object category learning which optimally selects unlabeled test points for interactive labeling, and (3) a probabilistic approach to learn discriminative kernel combinations for multiple local feature types within a GP framework. We show that with active learning small amounts of interactively labeled data can provide very accurate category recognition performance, while with covariance functions that optimally combine multiple matching kernels our method obtains state-of-the-art results with benchmark datasets.

## 2 Previous Work

Object category recognition has been a topic of active interest in the computer vision literature. Methods based on local feature descriptors (cf. Lowe 2004; Mikolajczyk and Schmid 2001) have been shown to offer invariance across a range of geometric and photometric conditions. Early models captured appearance and shape variation in a generative probabilistic framework (Fergus et al. 2003), but more recent techniques have typically exploited methods based on SVMs or nearest neighbors in a bag-of-visual-words feature space (Sivic and Zisserman 2003; Nister and Stewenius 2006; Zhang et al. 2006; Moosmann and Jurie 2007).

Several authors have explored correspondence-based kernels (Zhang et al. 2006; Wallraven et al. 2003), where the distance between a set of local feature descriptors—potentially including appearance and shape/position—is computed based on associating pairs of descriptors. However, the polynomial-time computational cost of correspondence-based distance measures makes them unsuitable for domains where there are large databases or large numbers of features per image. In Grauman and Darrell (2005) the authors introduced the Pyramid Match Kernel (PMK), an efficient linear-time approximation to a partial match correspondence, and in Lazebnik et al. (2006) it was demonstrated that a spatial variant—which efficiently represents the distinction between appearance and image location features—outperformed many competing methods.

Semi-supervised or unsupervised visual category learning methods are related to active learning, in that they also leverage unlabeled examples to learn more accurately when

limited labeled examples are available. Generative models which model visual words as arising from a set of underlying objects or “topics” based on recently introduced methods for Latent Dirichlet Allocation have been developed (Sivic et al. 2005; Sudderth et al. 2005) but as yet have not been applied to active learning nor evaluated on purely supervised tasks. A semi-supervised method using normalized cuts to cluster a graph defined by Pyramid Match distances between examples was presented in Grauman and Darrell (2006b), but this method is not probabilistic nor does it provide for an active learning formalism.

In the machine learning literature active learning has been a topic of recent interest, and numerous schemes have been proposed for choosing unlabeled points for tagging. For example, in Freund et al. (1997) the authors propose using the disagreement among the committee of classifiers as a criterion for active learning, and show an application to image classification (Abramson and Freund 2004). In Tong and Koller (2000), unlabeled examples to query are selected based on minimizing the version space within the SVM formulation, while in Chang et al. (2005) an SVM-based active learner is applied for image retrieval using color and texture features.

Within the Gaussian Process framework, the method of choice has been to look at the expected informativeness of an unlabeled data point (Lawrence et al. 2002; MacKay 1992). Specifically, the idea is to choose to query cases that are expected to maximally influence the posterior distribution over the set of possible classifiers. Additional studies have sought to combine active learning with semi-supervised learning (McCallum and Nigam 1998; Muslea et al. 2002; Zhu et al. 2003). Our work is significantly different as we focus on local feature approaches for the task of object categorization. We explore the GP models, which provide estimates for uncertainty in prediction and can be easily extended to active learning.

Recent work has shown the value of combining multiple local image feature types into a single kernel matrix, either by using cross-validation with a held-out set of labeled images to adjust the weight attached to each (Bosch et al. 2007), or by optimizing the weights to align the combined kernel with the ideal kernel matrix reflecting the labels on the training data (Kumar and Sminchisescu 2007; Lin et al. 2007; Varma and Ray 2007). Both tactics have yielded impressive results in practice. Our proposed method to optimize kernel weights fits directly within our GP learning framework, and is distinct in that rather than target the labels of training examples, it maximizes the evidence of the probabilistic model.

Gaussian Processes have been recently introduced to the computer vision literature. While they have been used in Urtasun et al. (2005, 2006) for human motion modeling, gender classification (Kim et al. 2006) and in Williams (2006)

for stereo segmentation, we are unaware of any prior work on visual object recognition in a Gaussian Process framework.<sup>1</sup>

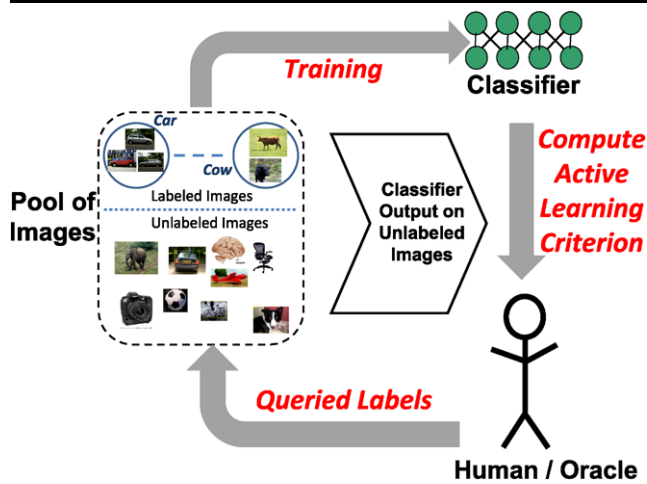
### 3 Approach Overview

Gaussian Processes provide an appealing probabilistic framework where the uncertainty is modeled conditioned on the observations. This circumvents the need to explicitly model the probability distribution of observations, which are often high-dimensional. Some popular GP models include GP classification and regression (Rasmussen and Williams 2006), non-linear dimensionality reduction using the Gaussian Process Latent Variable Model (GPLVM) (Lawrence 2004), and Mixture of GPs (Tresp 2000). In the context of this work, the Gaussian Process classification and regression framework is especially appealing as it can be considered a probabilistic counterpart to Support Vector Machines (SVM) and its variants, which have already been thoroughly used with much success for object recognition tasks. Beyond providing good accuracy, GP models for classification and regression are probabilistic, and also enjoy the benefits common to all kernel-based models.

The main idea of our approach is to construct probabilistic discriminative classifiers for object recognition using Gaussian Process priors, with covariance functions defined by the Pyramid Match Kernel (GP-PMK). In addition to offering a novel approach to supervised visual category learning, we show how this framework also allows both the learning of optimal combinations of covariance functions, as well as an active learning strategy—which is especially preferable when minimal labeling effort is available. Figure 1 shows the proposed framework for active image categorization. Given a pool of images of which few are labeled, the system aims to actively seek labels for unlabeled images by considering information from both the labeled and unlabeled sets of images. With the uncertainty estimates the GP classifier provides, we are able to designate an active learning criterion to focus labeling efforts on the most ambiguous unlabeled examples.

In the next section we review classification using GP priors and discuss the distributions and parameters we employ for our model. Then in Sect. 5 we present our GP-PMK model, which is directly suitable for supervised learning with or without active learning. Then in Sect. 6 we describe how to optimize the weights to combine multiple matching kernels computed from different feature sets. Finally, we

<sup>1</sup> This paper expands on our previous conference publication (Kapoor et al. 2007); here we provide further explanation of our Gaussian Process model, extend it to allow combinations of multiple kernel functions, and report and discuss a number of additional experiments.



**Fig. 1** The active learning framework. The goal of the system is to query labels for images that are most useful in training

derive an active learning variant that can optimally select points for interactive labeling in Sect. 7.

Note that throughout we assume that there is one primary object of interest in an image. Handling multiple objects in the same image is also an interesting and challenging problem, and will be the focus of future work.

#### 4 Categorization with Gaussian Processes

Gaussian Process (GP) classification is related to kernel machines such as Support Vector Machines (SVMs) (Evgeniou et al. 2000) and Regularized Least Square Classification (RLSC) and has been well-explored in machine learning. In contrast to these methods, GPs provide probabilistic prediction estimates and thus are well-suited for active learning. In this section we briefly review regression and classification with Gaussian Process priors and describe our model choices.

Given a set of labeled data points  $\mathbf{X}_L = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , with class labels  $\mathbf{t}_L = \{t_1, \dots, t_n\}$ , we are interested in classifying the unlabeled data  $\mathbf{x}_u$ . Under the Bayesian paradigm, we are interested in the distribution  $p(t_u | \mathbf{X}, \mathbf{t}_L)$ . Here  $\mathbf{X} = \{\mathbf{X}_L, \mathbf{x}_u\}$ , and  $t_u$  is the random variable denoting the class label for the unlabeled point  $\mathbf{x}_u$ . For sake of simplicity in discussion we limit ourselves to two-way classification, hence, the labels are  $t_i \in \{-1, 1\}$ , but this can be extended to multi-label classification; see Rasmussen and Williams (2006) for a detailed discussion.

With GP models, a discrete label  $t$  for a data point  $\mathbf{x}$  can be considered to be generated via a continuous hidden random variable  $y$ . The soft-hidden label arises due to a Gaussian Process, which in turn imposes a smoothness constraint on the possible solutions. A likelihood model  $p(t|y)$  characterizes the relationship between the soft label  $y$  and

the observed annotation  $t$ . Thus, when we infer the label  $t_u$  for the unlabeled data point  $\mathbf{x}_u$ , we probabilistically combine the smoothness constraint and the information obtained by observing the annotations  $\mathbf{t}_L$ .

##### 4.1 Smoothness Constraints via the GP Prior

There exist two different perspectives for regression and classification with Gaussian Process: the process perspective and the weight perspective. We overview both in the following in order to provide background on the basic concepts underlying the GP model.

*The process perspective:* The smoothness constraint is imposed using a Gaussian Process prior that defines the probabilistic relationship between the images  $\mathbf{X}$  and the soft labels  $\mathbf{Y}$ . The distribution  $p(\mathbf{Y}|\mathbf{X})$  gives higher probability to the labelings that respect the similarity between the data points. Intuitively, the assumption is that similar data points should have the same class assignments/regression values; the similarity between two points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  is defined via a kernel  $k(\mathbf{x}_i, \mathbf{x}_j)$ . Probabilistic constraints are imposed on the collection of soft labels  $\mathbf{Y} = \{y_1, \dots, y_n, y_u\}$ . In particular, the soft labels are assumed to be jointly Gaussian and the covariance between two outputs  $y_i$  and  $y_j$  is typically specified using a kernel function<sup>2</sup> applied to  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . Formally,  $p(\mathbf{Y}|\mathbf{X}) \sim \mathcal{N}(0, \mathbf{K})$  where  $\mathbf{K}$  is a  $(n+1)$ -by- $(n+1)$  kernel matrix with  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ , and  $n+1$  reflects the  $n$  labeled examples and one unlabeled example.

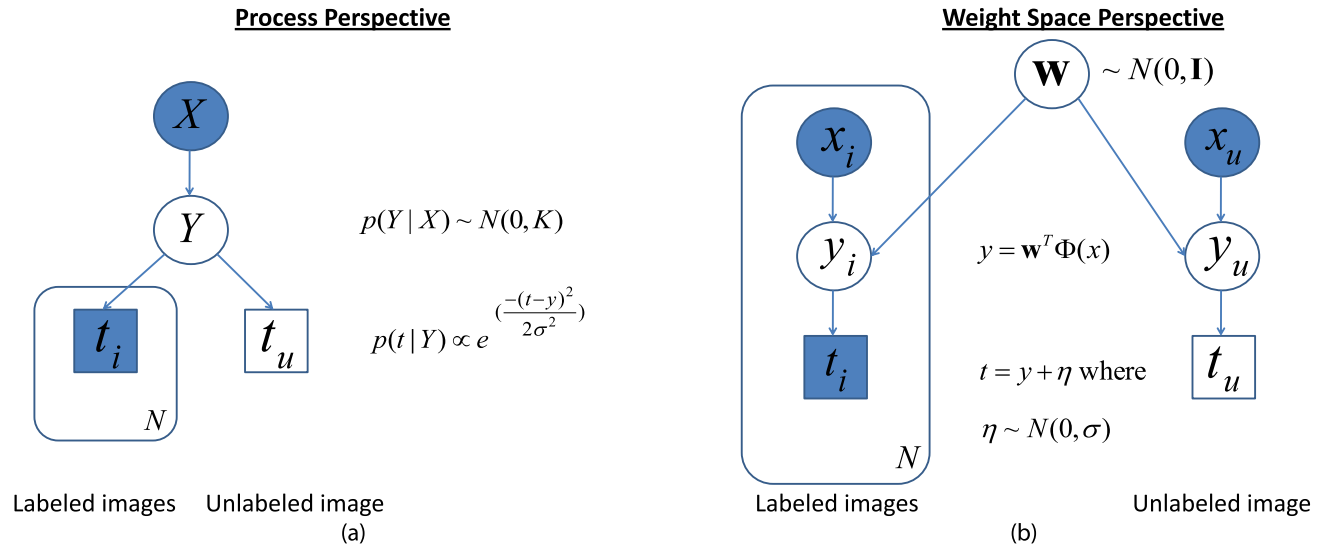
*The weight perspective:* What we have described above is the process perspective for regression and classification with the GP priors. An alternate but mathematically equivalent interpretation is based on the weight perspective. In this perspective the hidden soft-label  $y$  arises due to application of a function  $f(\cdot)$  directly on the input data point (i.e.  $y = f(\mathbf{x})$ ), which takes the form of a linear combination of orthonormal basis functions:

$$f(\mathbf{x}) = \sum_k w_k v_k^{1/2} \phi_k(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}), \quad (1)$$

where  $\phi_k$  are the eigenfunctions of the operator induced by  $k$  in the Reproducing Kernel Hilbert Space (Evgeniou et al. 2000),  $v_k$  are the corresponding eigenvalues,  $w_k$  are the weights, and  $\Phi(\mathbf{x}) = [v_1^{1/2} \phi_1(\mathbf{x}), v_2^{1/2} \phi_2(\mathbf{x}), \dots]^T$ . Note that the dimensionality of the basis can be infinite. Assuming a spherical Gaussian prior over the weights, that is  $\mathbf{w} = [w_1, w_2, \dots]^T \sim \mathcal{N}(0, \mathbf{I})$ , it can be shown that the hidden soft labels  $\mathbf{Y}$  (which result from evaluation of the function  $f(\cdot)$  on the input data points  $\mathbf{X}$ ) are jointly Gaussian

<sup>2</sup>One can use a non-parametric covariance function, but the number of parameters to estimate grows exponentially with the amount of training data.





**Fig. 2** Graphical models in plate notation for classification via Gaussian Processes. The *rounds* and *squares* represent continuous and discrete random variables, respectively. A *filled* (*unfilled*) *round/square* denotes that the random variable is fully observed (unobserved).  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_u\}$  is the set of all images and is ob-

served for both labeled and unlabeled data points. The corresponding  $\mathbf{Y} = \{y_1, \dots, y_n, y_u\}$  is completely unobserved and the labels  $\{t_1, \dots, t_n\}$  are observed only for the training images  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and unobserved for the unlabeled image  $\mathbf{x}_u$

with zero mean and with the covariance given by the kernel matrix  $\mathbf{K}$ .

These two different but equivalent perspectives for regression and classification with the GP priors are illustrated in Fig. 2. Both views lead to different implementations, but are conceptually equivalent. The process perspective is easier to implement for the cases when  $\mathbf{K}$  is non-parametric or when it is difficult to determine  $\Phi(\mathbf{x})$ . Similarly, difficulty arises in using the weights perspective when  $\Phi(\mathbf{x})$  is high dimensional (possibly infinite dimensional). In such cases it is easier to follow the process perspective as inference over the labels can be done by computing the kernel matrix, which circumvents the need to explicitly evaluate the eigenfunctions. In our work we determine similarities between images using primarily non-parametric techniques for which the explicit form of the eigenfunctions is unknown. Thus, in this work, we follow the process perspective. For details, please see Seeger (2004).

#### 4.2 The Likelihood Model

The likelihood models the probabilistic relationship between the observed label  $t$  and the hidden label  $y$ . The majority of the likelihood models proposed for GP classification use additional latent “squashing” variables that transform unconstrained variables into labels. A wide range of squashing functions have been developed in the literature (Rasmussen and Williams 2006) and examples include the logistic and the probit functions. To make predictions based on the training set for a test set in these models (i.e., probit and logit) one has to integrate out the prediction over the

posterior. Since the likelihood is not Gaussian, neither the posterior, the marginal likelihood, nor the predictions can be computed analytically. Instead, one has to rely on numerical methods, such as MCMC (Williams and Barber 1998), or approximations of the posterior, e.g. Laplace and Expectation Propagation (Minka 2001).

In contrast to GP classification, GP regression leads to efficient analytic solutions for prediction. For Gaussian Process regression using a Gaussian noise model, the relation between  $t$  and  $y$  is given by

$$p(t|y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(t-y)^2}{2\sigma^2}}, \quad (2)$$

where  $\sigma$  is the noise model variance. Since this likelihood model is Gaussian, it leads to a closed form solution for inference. Although originally developed for regression, the Gaussian noise model has also proven effective for classification,<sup>3</sup> and its performance typically matches the more complex probit and logit likelihood models noted above. Due to its simplicity and good performance, in our experiments we use regression (i.e., the Gaussian noise model) to label variables. Non-Gaussian noise models could also be applied within the proposed framework, and exploring them is a topic of interest for future work.

<sup>3</sup>This method is referred to as least-squares classification in the literature (see Sect. 6.5 of Rasmussen and Williams 2006) and often demonstrates performance competitive with more expensive Gaussian Process classification methods that require approximate inference.

### 4.3 Inference

Given the labeled and unlabeled data points, our goal is then to infer  $p(t_u|\mathbf{X}, \mathbf{t}_L)$ . Specifically:

$$p(t_u|\mathbf{X}, \mathbf{t}_L) \propto \int_{\mathbf{Y}} p(t_u|\mathbf{Y}) p(\mathbf{Y}|\mathbf{X}, \mathbf{t}_L). \quad (3)$$

For a Gaussian noise model we can compute this integral using closed form expressions. Note that the key quantity to compute is the posterior  $p(\mathbf{Y}|\mathbf{X}, \mathbf{t}_L)$ , which can be written as:

$$p(\mathbf{Y}|\mathbf{X}, \mathbf{t}_L) \propto p(\mathbf{Y}|\mathbf{X}) p(\mathbf{t}_L|\mathbf{Y}) = p(\mathbf{Y}|\mathbf{X}) \prod_{i=1}^n p(t_i|y_i). \quad (4)$$

This equation probabilistically combines the smoothness constraints  $p(\mathbf{Y}|\mathbf{X})$  imposed via the GP prior and the information provided in the labels ( $p(\mathbf{t}_L|\mathbf{Y})$ ). The posterior as shown in (4) is simply a product of Gaussians, and the posterior over the soft label  $y_u$  has a particularly simple form. Specifically,  $p(y_u|\mathbf{X}, \mathbf{t}_L) \sim \mathcal{N}(\bar{y}_u, \sigma_u^2)$ , where:

$$\bar{y}_u = \mathbf{k}(\mathbf{x}_u)^T (\sigma^2 \mathbf{I} + \mathbf{K}_{LL})^{-1} \mathbf{t}_L, \quad (5)$$

$$\Sigma_u = k(\mathbf{x}_u, \mathbf{x}_u) - \mathbf{k}(\mathbf{x}_u)^T (\sigma^2 \mathbf{I} + \mathbf{K}_{LL})^{-1} \mathbf{k}(\mathbf{x}_u). \quad (6)$$

Here,  $\mathbf{k}(\mathbf{x}_u)$  is the vector of kernel function evaluations with  $n$  training points, and  $\mathbf{K}_{LL} = \{k(\mathbf{x}_i, \mathbf{x}_j)\}$ , is the training covariance, with  $\mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}_u$ . Further, due to the Gaussian noise model that links  $t_u$  to  $y_u$ , the predictive distribution over the unknown label  $t_u$  is also a Gaussian:  $p(t_u|\mathbf{X}, \mathbf{t}_L) \sim \mathcal{N}(\bar{y}_u, \Sigma_u + \sigma^2)$ .

Note that the posterior mean for both  $t_u$  and  $y_u$  is the same; thus, the unlabeled point  $\mathbf{x}_u$  can be classified according to the sign of  $\bar{y}_u$ . Note that despite the fact that we only consider sign of the posterior mean for classification, the uncertainties modeled using the GP prove very useful in performing two key tasks. First, GPs inherently provide a principled way to do kernel combination as the probabilistic framework allows us to determine how well a particular combination of kernels explains the data well. Secondly, unlike the Regularized Least Square Classification (RLSC) methods we also get the variance in prediction. As we will show in Sect. 7, we can exploit these measures of uncertainty to guide an active learning procedure.

### 4.4 Training with the Gaussian Process Models

The performance of Gaussian Process-based classification depends upon the chosen kernel used to capture the similarity between examples, as well as the kernel's hyperparameters, such as the length-scale, the noise variance, and other parameters determining local feature-based image similarity. Finding the right set of all these parameters can be a

challenge. Many discriminative models (including SVMs) often use cross-validation, which is a robust measure but can be prohibitively expensive and problematic when we have few labeled data points. Learning in a Gaussian Process framework is equivalent to choosing the kernel hyperparameters of the covariance function. Ideally we would like to marginalize over these hyperparameters. While approaches based on Hybrid Monte Carlo have been explored to perform this marginalization (Williams and Barber 1998), such techniques are relatively expensive.

Empirical Bayes is a more computationally efficient alternative where the idea is to maximize the marginal likelihood or the evidence, which is nothing but the constant  $p(\mathbf{t}_L|\mathbf{X})$  that normalizes the posterior. This methodology of tuning the hyperparameter is often called *evidence maximization*, and has been one of the favorite tools for performing model selection. Evidence is a numerical quantity and signifies how well a model fits the given data. By comparing the evidence corresponding to the different models (or hyperparameters that determine the model), we can choose the model and the hyperparameters suitable for the task.

The idea is to choose a set of hyperparameters  $\Theta$  that maximize the evidence:  $\hat{\Theta} = \arg \max_{\Theta} \log[p(\mathbf{t}_L|\mathbf{X}, \Theta)]$ . Note that the log evidence  $\log(p(\mathbf{t}_L|\mathbf{X}, \Theta))$  can be written as a closed form equation for the Gaussian noise model (GP regression):

$$\begin{aligned} \log p(\mathbf{t}_L|\mathbf{X}, \Theta) = & -\frac{1}{2} \mathbf{t}_L^T (\sigma^2 \mathbf{I} + \mathbf{K}_{LL})^{-1} \mathbf{t}_L \\ & -\frac{1}{2} \log |\sigma^2 \mathbf{I} + \mathbf{K}_{LL}| - Const. \end{aligned}$$

This objective can be maximized using non-linear optimization techniques, such as gradient descent. In this work, we use gradient-descent to maximize evidence. The optimization procedure can perform multiple searches with different initializations to deal with the fact that the evidence will have multiple local optima.

This scheme of learning hyperparameters by maximizing evidence lets us find the correct parameters without the need of cross-validation. Further, this procedure can also be used to learn an ideal linear combination of covariance functions, which is a useful tool in practice to combine various local feature object categorization schemes. We show this combination strategy in Sect. 6.

## 5 Pyramid Match Kernel Gaussian Processes (GP-PMK)

To use GPs for object categorization, we need to define a suitable covariance function. We would like to exploit local feature methods for object and image representations. However, GP priors require covariance functions which are

positive semi-definite (a Mercer kernel) and traditional covariance functions (e.g., RBF) are not suitable for representations that are comprised of sets of features.

We wish to define a GP with a covariance function based on a partial match distance function. The idea is to first represent an image as an unordered set of local features, and then use a matching over these sets of features to compute a smoothness prior between images. The optimal least-cost partial matching takes two sets of features, possibly of varying sizes, and pairs each point in the smaller set to a unique point in the larger one, such that the sum of the distances between the matched points is minimized. The cubic cost of the optimal matching makes it prohibitive for recognition with a large number of local image features, yet rich image descriptions comprised of densely sampled local features are known to often yield better recognition accuracy (Nowak et al. 2006).

Therefore, rather than adopt a full partial match kernel for the GP prior, we use the Pyramid Match (Grauman and Darrell 2005). The Pyramid Match is a linear-time kernel function over unordered feature sets that approximates the similarity measured by the optimal partial matching, and it forms a Mercer kernel. A multi-resolution partition (pyramid) carves the feature space into increasingly larger regions. At the finest resolution level in the pyramid, the partitions are very small; at successive levels they continue to grow in size until a single partition encompasses the entire feature space. The insight of the Pyramid Match algorithm is to treat points which share a bin in this pyramid as being matched, and to use the histograms to read off the number of possible matches without explicitly searching for correspondences. Histogram intersection (the sum of the minimum number of points in a given histogram bin) is used to count the number of new matches that occur at each resolution level.

The input space  $S$  contains sets of feature vectors drawn from feature space  $\mathcal{F}$ :  $S = \{\mathbf{F} | \mathbf{F} = \{\mathbf{f}_1, \dots, \mathbf{f}_m\}\}$ , where each feature  $\mathbf{f}_i \in \mathcal{F} \subseteq \mathbb{R}^d$ , and  $m = |\mathbf{F}|$ . For example,  $\mathcal{F}$  might be the space of SIFT (Lowe 2004) descriptors ( $d = 128$ ), or image coordinate positions ( $d = 2$ ), etc.; a set  $\mathbf{F}$  contains a collection of these descriptors extracted from a single image or object. An  $L$ -level histogram pyramid for input example  $\mathbf{F} \in S$  is defined as:  $\Psi(\mathbf{F}) = [H_0(\mathbf{F}), \dots, H_{L-1}(\mathbf{F})]$ , where  $H_i(\mathbf{F})$  is a histogram vector formed over points in  $\mathbf{F}$  using multi-dimensional bins.

The Pyramid Match Kernel (PMK) value between two input sets  $\mathbf{F}_1, \mathbf{F}_2 \in S$  is defined as the weighted sum of the number of feature matches found at each level of their pyramids (Grauman and Darrell 2005):

$$\mathbf{K}_\Delta(\Psi(\mathbf{F}_1), \Psi(\mathbf{F}_2))$$

$$= \sum_{i=0}^{L-1} w_i (\mathcal{I}(H_i(\mathbf{F}_1), H_i(\mathbf{F}_2)) - \mathcal{I}(H_{i-1}(\mathbf{F}_1), H_{i-1}(\mathbf{F}_2))),$$

where  $\mathcal{I}$  denotes histogram intersection, and the difference in intersections across levels serves to count the number of new matches formed at level  $i$  that were not already counted at any finer resolution level. Note that  $H_i(\cdot)$  corresponds to histogram at level  $i$ , where  $H_{-1}(\cdot)$  is always zero and bins at level  $i$  are always larger than those at level  $i - 1$ . The weights are set to be inversely proportional to the size of the bins, in order to reflect the maximal distance two matched points could be from one another. As long as  $w_i \geq w_{i+1}$ , the kernel is Mercer.

We thus define a Pyramid Match Gaussian Process model (GP-PMK) using the prior

$$p(\mathbf{Y}|\mathbf{X}) \sim \mathcal{N}(0, \mathbf{K}_\Delta). \quad (7)$$

In contrast to previous GP priors, this prior is well-suited for visual category recognition as it naturally handles representations based on sets of local image features.

A variety of Pyramid Match Kernels (and thus GP priors) are possible, given that we have flexibility in choosing the interest operator used to sample local image regions, the type of descriptor used to describe each region, and the partitioning strategy used to form the pyramid histogram bins. To extract local features, we can exploit a wealth of interest operators designed to detect a sparse set of salient regions (e.g., Lowe 2004; Mikolajczyk and Schmid 2004; Kadir and Brady 2003), or simply sample densely at regular intervals and at multiple scales. To describe each region or patch, we can choose from an array of descriptors designed to capture local texture while maintaining some invariance to small shifts and rotations, such as SIFT (Lowe 2004), shape context (Belongie et al. 2001), or geometric blur (Berg and Malik 2001).

For low-dimensional feature spaces, the partitions within each histogram  $H_i$  may be placed at uniform intervals to divide the feature space into equally sized grid cells, as in Grauman and Darrell (2005) and Lazebnik et al. (2006). For higher-dimensional feature spaces it is better to place the partitions non-uniformly in a data-dependent manner, as described in Grauman and Darrell (2006a). To encode spatial position together with the region appearance, each feature  $\mathbf{f}_i$  can be expanded to include both the image descriptor concatenated with the normalized image coordinate at which it occurred; however, doing so requires standardizing the dimensions carefully. An efficient way to incorporate both feature channels is to use the spatial pyramid match (Lazebnik et al. 2006), a variant of the PMK that first quantizes the appearance feature descriptors to form a bag-of-words representation, and then sums over the PMK values for each word in the space of image coordinates. Depending on the image data, such choices are likely to influence the accuracy of the GP-PMK model. In the next section, we describe a technically sound strategy to combine all of these different kernels such that the resulting kernel is highly discriminatory.

## 6 Combining Multiple Covariance Functions

Given multiple kernels  $\mathbf{K}^{(1)}, \dots, \mathbf{K}^{(k)}$  we seek a linear combination of the base kernels such that the resulting kernel  $\mathbf{K}$  has good discriminatory power. Formally, we have

$$\mathbf{K} = \sum_{i=1}^k \alpha_i \mathbf{K}^{(i)}, \quad (8)$$

where  $\alpha = \{\alpha_1, \dots, \alpha_k\}$  are the weight parameters that we wish to optimize. We can take an evidence maximization approach as described in Sect. 4.4 to solve for these weights. However, note that the procedure of evidence maximization is a type-2 maximum likelihood estimation technique and akin to the principle of MAP estimation we can additionally regularize the objective function. In particular instead of finding the hyperparameters by maximum likelihood, we assume a prior distribution over the hyperparameters,  $p(\alpha)$ , and choose the maximum-a-posteriori (MAP) estimate. It can be easily shown that various choices of priors lead to different choices of regularization. For instance assuming a Gaussian and a Laplacian prior on  $\alpha$  leads to an  $L2$  and  $L1$  regularized formulation respectively; the latter is well known to enforce a degree of sparsity on the kernel weights. Formally the objective we minimize is:

$$\arg \min_{\alpha} -\log p(\mathbf{t}_L | \mathbf{X}, \alpha) + \gamma_1 \|\alpha\|_1 + \gamma_2 \|\alpha\|_2$$

subject to:  $\alpha_i \geq 0$  for  $i \in \{0, \dots, k\}$ .

Here,  $\gamma_1$  and  $\gamma_2$  are regularization constants for  $L1$  and  $L2$  norms respectively and often boost recognition performance when the amount of labeled data is low. The non-negativity constraints on  $\alpha$  ensure that the resulting  $\mathbf{K}$  is positive-semidefinite and can be used in a GP formulation (or other kernel-based methods).

The proposed objective is a non-linear program and can be solved using any gradient-descent based procedure. In our implementation we use a gradient descent procedure where we limit based on the projected BFGS method using a simple line search. The gradients of the objective are efficient to compute and can be written as:

$$\begin{aligned} \frac{\delta \mathcal{L}(\alpha)}{\delta \alpha_i} = & -\frac{1}{2} \mathbf{t}_L^T \mathbf{A}^{-1} \mathbf{K}_{LL}^{(i)} \mathbf{A}^{-1} \mathbf{t}_L + \frac{1}{2} \text{Tr}(\mathbf{A}^{-1} \mathbf{K}_{LL}^{(i)}) \\ & + \gamma_1 + 2\gamma_2 \alpha_i, \end{aligned}$$

where  $\mathbf{A} = \sigma^2 \mathbf{I} + \mathbf{K}_{LL}$  and  $\alpha_i \geq 0$  for all  $i$ . In our implementation, the non-negativity constraints for  $\alpha_i$  are enforced by considering their form to be an exponential  $\alpha_i = e^{\beta_i}$  and then performing an unconstrained optimization to determine optimal values for each  $\beta_i$ . Once the parameters  $\alpha$  are found, then the resulting linear combination of kernels ( $\mathbf{K}$ ) can be used for classification. By selecting the kernel

weights within the GP framework, we allow a user to provide several feature choices and PMK kernel variants that seem plausible, with the system itself selecting the most discriminative combination.

Note that learning the kernel combination as described in this section is a particular parameterization of the GP kernel being learned, and the method to optimize the particular objective is principally an instantiation of the general method described in Sect. 4.4. Also, here we assume that Pyramid Match kernels (or other similarity measures) are given, and that aside from the values of  $\alpha$ , there are no additional parameters to be optimized in the individual kernels. However, should the individual kernels also have parameters to be set, then it would be straightforward to extend the optimization scheme to include those parameters as well.

### 6.1 Extension to Multi-Class Problems

Object categorization is typically a multiclass problem and consequently requires a multiclass extension of the kernel learning framework. Popular techniques include 1-vs-all or 1-vs-1 formulations, where outputs from multiple binary classifiers trained on 1-vs-rest and pairwise classification problems are combined respectively. Learning a kernel introduces additional complexity as the optimization procedure for kernel combination should consider all the labels and result in a single set of global parameters that are informative about the entire classification task. Learning a global set of parameters is in general non-trivial, and learning separate kernels for each binary subproblem has been proposed (Varma and Ray 2007). Despite the fact that such *classwise* parameterizations offer flexibility in modeling each individual class, these strategies are more prone to overfitting than global ones when dealing with small number of examples. As shown below, global optimization of the parameters consistently outperforms classwise optimization in our experiments.

Furthermore, classwise techniques require solving as many classification tasks as the number of classes; with large datasets such as Caltech-101 this means that learning has to be repeated 101 times. While it is unclear how to overcome these issues in non-probabilistic approaches such as Varma and Ray (2007), GPs provide a principled and computationally efficient scheme of finding globally optimal parameters.

Lets consider a 1-vs-all formulation of GP classifiers, where multiple binary classifiers correspond to each individual class. Similar to binary classifiers we optimize kernels weights by considering the log evidence, however, for the multiclass case we consider a joint log-likelihood over all the classifiers:

$$\mathcal{L}(\alpha) = -\sum_i \log p(\mathbf{Y}^{(i)} | \mathbf{X}, \alpha).$$



Here the sum is taken over all the class labels, and  $\mathbf{Y}^{(i)}$  are the labels for  $i$ th 1-vs-all problem. This joint likelihood corresponds to a probabilistic model that assumes that given the input images the binary outputs of 1-vs-all problems are independent. Note that, this assumption is well justified as given an image its class label is determined by the image content only. Further, this model allows us to optimize for a global set of kernel parameters that maximize the joint likelihood over all the class labels. Thus, instead of learning a kernel for every individual class, we can learn an optimal parameterization that is globally discriminative.

There are additional computational benefits of the above scheme. Note that in the proposed GP framework, given a test observation  $\mathbf{x}_*$ , the mean prediction for a binary classifier can be computed as:

$$\bar{y}_* = \mathbf{k}(\mathbf{x}_*)^T \mathbf{A}^{-1} \mathbf{t}_L, \quad (9)$$

where  $\mathbf{k}(\mathbf{x}_*)$  is the kernel computed between the training and test data and  $\mathbf{A} = \sigma^2 \mathbf{I} + \mathbf{K}_{LL}$ . The most expensive operation in such computation is the matrix inversion which has a time complexity of  $O(n^3)$  for  $n$  training examples and is independent of the training labels  $\mathbf{t}_L$ . Consequently, once the inverse is computed, estimating predictions for 1-vs-all models corresponds to a multiplication with the relevant label vectors.

This is a significant advantage since the cost of training all the classifiers in a 1-vs-all formulation is the same as the cost of training a single classifier. This is especially beneficial in cases with a large number of classes, and provides a significant advantage over other methods which separately need to train different classifiers per class. This observation readily extends to the kernel learning scenario with multiple classes. As before, the primary operation is a matrix inversion (computing  $\mathbf{A}^{-1}$ ) that is independent of the labels. Thus, learning kernels for multiple class problems using the joint likelihood has similar cost as that of learning a kernel in a binary problem.<sup>4</sup> There are various ways to make this computation even more efficient. Specifically, a lot of research has gone into *sparsifying* GP (Lawrence et al. 2002; Shen et al. 2006; Snelson and Ghahramani 2006; Urtasun and Darrell 2008) where the aim is to select a subset of points that are informative or important with respect to the classification task. Also note that in addition to reducing manual labeling effort, an active learning formulation does help us reduce the computational overhead in inference by reducing the number of needed training points.

Competing approaches for kernel combination (Varma and Ray 2007) are based on support vector machines. In contrast to our approach, the method proposed by Varma and

Ray (2007) is non-probabilistic and is based on second order cone programming (SOCP), which has similar or worse time complexity (Tsang and Kwok 2006). It might be possible to further improve the performance of SVM-based kernel combination by using cross validation. We can do a local search for kernel combination parameters around an initial solution found by the method of Varma and Ray (2007). The time required to run such an approach is bounded from below by the time required to first optimize the SVM parameters. Further, a grid search even within a limited range quickly becomes infeasible as we increase the number of kernels. For example, for eight kernels on the Caltech-101 dataset, grid search over 21 possible values of each parameter requires training of  $101 \cdot 21^8 \approx 3.8$  trillions SVMs.

## 7 Active Learning for Object Categorization

In this section we consider the scenario where our visual category learner has access to a pool of unlabeled data  $\mathbf{X}_U = \{\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+m}\}$ . The task in active learning is to seek the label for one of these examples and then update the classification model by incorporating it into the existing training set. The goal is to select the sample that would maximize the benefit in terms of the discriminatory capability of the system.

With non-probabilistic classification schemes, a popular heuristic for establishing the confidence of estimates and identifying points for active learning is to simply use the distance from the classification boundary (margin). This approach can also be used with GP classification models, by inspecting the magnitude of the posterior mean  $|\bar{y}_u|$ : one would then choose the next point  $\mathbf{x}^*$  as  $\arg \min_{\mathbf{x}_u \in \mathbf{X}_U} |\bar{y}_u|$ .

However, GP classification provides us with both the posterior mean as well as the posterior variance for the unknown label  $t_u$ . An alternative criteria could be to look at the variances and select the point that has the maximum variance, i.e.  $\mathbf{x}^* = \arg \max_{\mathbf{x}_u \in \mathbf{X}_U} \Sigma_u$ . However such an approach does not consider the mean  $\bar{y}_u$  at all! Further, the expression for  $\Sigma_u$  does not consider labels from the annotated training data; this scheme uses only a very limited amount of information.

We therefore propose an approach which considers both the posterior mean as well as the posterior variance. Specifically, we select the next point according to:

$$\mathbf{x}^* = \arg \min_{\mathbf{x}_u \in \mathbf{X}_U} \frac{|\bar{y}_u|}{\sqrt{\Sigma_u + \sigma^2}}, \quad (10)$$

where  $\sigma^2$  is the noise model variance. This formulation considers uncertainty in the labeling  $\mathbf{x}_u$  as  $\pm 1$ . Note that the predictive distribution for  $t_u$  is a Gaussian; however, we are interested in the binary label decided according to the

<sup>4</sup>The computational cost is dominated by the  $O(n^3)$  cost of inverting  $\mathbf{A}$ , with  $n$  the number of examples.

**Table 1** Active learning criteria

Method	Criteria
Distance from Boundary (SVM)	$\mathbf{x}^* = \arg \min_{\mathbf{x}_u \in \mathbf{X}_U}  \bar{y}_u $
Variance	$\mathbf{x}^* = \arg \max_{\mathbf{x}_u \in \mathbf{X}_U} \Sigma_u$
Uncertainty (GP)	$\mathbf{x}^* = \arg \min_{\mathbf{x}_u \in \mathbf{X}_U} \frac{ \bar{y}_u }{\sqrt{\Sigma_u + \sigma^2}}$

sign of  $t_u$ . To this end we should consider the value  $p(t_u \geq 0) = \phi(\frac{\bar{y}_u}{\sqrt{\Sigma_u + \sigma^2}})$ , where  $\phi(\cdot)$  denotes the cdf of a standard normal distribution, to provide the hard label  $\pm 1$ . Further, we are interested in selecting those samples where the uncertainty is maximum. The points where the classification model is most uncertain should have a value for  $p(t_u \geq 0)$  that is close to 0.5—equivalently, a value of  $\frac{|\bar{y}_u|}{\sqrt{\Sigma_u + \sigma^2}}$  that is very close to zero. Thus, the criterion in (10) chooses the unlabeled point where the classification is the most uncertain.

We summarize the methods for identifying points to be labeled in Table 1, with our strategy given in the third row. Our active learning approach looks at all the points before choosing the points to actively label; thus it considers the whole dataset instead of just looking at individual points. Further, this scheme considers both the distance from the boundary as well as the variance in selecting the points; this is only possible due to the availability of the predictive distribution in GP regression. In results below we show that in practice we can effectively choose useful examples to label, allowing our active GP approach to fare much better with minimal labeled data than a “passive” random selection scheme.

## 8 Experiments and Results

In this section we report results from experiments to demonstrate (1) the effectiveness of the GP-PMK classification framework, (2) the ability of the proposed framework to identify good kernel combinations, and (3) how active learning can guide the learning procedure to select critical examples to be labeled. We show how kernel combination and active learning with Gaussian Process priors yield classifiers which can learn object categories from relatively few examples.

### Datasets and Implementation Details

We performed supervised and active learning experiments on two different datasets that are considered standards for the object categorization task: the Caltech-4 dataset and the Caltech-101 dataset (which is a superset of Caltech-4). We

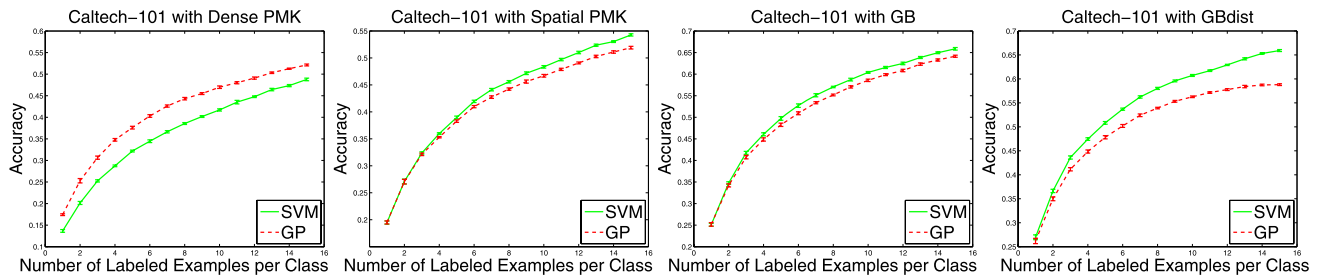
compute the similarity between all pairs of images in each database using the PMK. We use LIBSVM (Chang and Lin 2001) for the SVM baseline tests. In our experiments we set the noise model variance  $\sigma = 10^{-5}$  for the Gaussian Process models and fix  $C = 10000$  for SVM models. These parameter values worked well; we experimented with other values but found that both SVM and GP classification schemes were fairly insensitive to the choice of these parameters. Further, we initialize all the weights as  $\alpha_i = 1$  for the kernel learning procedure. We found the GP-based kernel learning to be extremely stable with respect to the initializations.

The object categorization task is a multi-class problem ( $n_{class} = 101$  and  $n_{class} = 4$  for the Caltech-101 and the Caltech-4, respectively). To handle multiple classes we use the one-vs-all formulation, where we choose the label corresponding to the class with maximum value of the soft label  $y$ . For kernel combination under the one-vs-all classification scheme we assume a joint model by summing the log evidence over all the binary classification problems. For multi-class active learning in every round we select one example from each of the one-vs-all classifiers, thus adding  $n_{class}$  examples every time.

The Caltech-4 database contains 3188 images with four object classes. There are 1155 rear views of cars, 800 images of airplanes, 435 images of frontal faces, and 798 images of motorcycles. The second database is the Caltech-101 database of 101 object categories (Fei-Fei et al. 2006); there are 8677 images in this data set, with between 31 to 800 images for each of the 101 categories. Our experiments for kernel combination use 30 images per class (3030 images in total), and are exactly the same as the ones used in Varma and Ray (2007). We perform active learning experiments using the complete Caltech-101 dataset.

We consider various shape and appearance features and sampling strategies, which are useful to capture the intra-class variation present in the Caltech-101 images. Specifically, we look at the following four combinations of matching kernels and features:

- **Dense PMK:** the PMK with uniformly shaped pyramid bins, using SIFT descriptors extracted densely from the images at every 8th pixel in the image from a region of 16 pixels in diameter, with each SIFT descriptor concatenated with its normalized image position. We use PCA to reduce the dimensionality of the SIFT descriptors to 10 before adding the position, yielding features having a total of 12 dimensions. See Grauman and Darrell (2005) for details.
- **Spatial PMK:** The spatial variant of Dense-PMK. We take the same raw SIFT features, but quantize them into visual words, and then build one pyramid per word, each with uniform bins in the space of image coordinates. See Lazebnik et al. (2006) for details.



**Fig. 3** Performance comparison of GP and SVM classification on each of the four different kernels used in this work. The figures show that using the GP framework with supervised learning achieves com-

- **GB:** The Geometric Blur feature of Berg and Malik (2001) is extracted at sampled edge points. For the kernel values, the exact correspondences are computed based on the average minimum distance between points in the two input sets of features, as in Zhang et al. (2006).
- **GBdist:** Same as GB, except the feature representation has an additional geometric distortion term.

The kernel matrices for this dataset using each of these variants were provided directly by the authors. The first two kernels are variations of the PMK with different feature spaces, while the last two kernels use an explicit (non-approximate) correspondences between geometric blur features.

In all our experiments in which comparisons are made against other methods, we follow the standard testing protocol, where a given number of training images (say 15) are taken from each class at random, and the rest of the data is used for testing. The mean recognition rate per class is used as a metric of performance. Note that this metric ensures that the recognition accuracies are such that classes with large numbers of examples are not favored. This process is repeated 10 times and the average correctness rate is reported.

#### Classification and Kernel Combination

First, we explore classification performance on individual kernels using different classification strategies. Figure 3 graphically shows the performance of classification with Gaussian Process as compared to an SVM classifier. From the eight graphs we can observe that overall the performance obtained using either the GP or the SVM is very similar. However, we note some deviations in performance: for example GP is significantly better on the Dense-PMK, whereas SVM performs very well with GBdist. We also compare performance of these methods with 1-Nearest Neighbor as a baseline; Table 2 summarizes the accuracy obtained with 15 labeled points per class. Both GP and SVM perform significantly better than the 1-Nearest Neighbor classifier. We find these experiments encouraging, since they indicate that we

parable performance to that of SVMs, but with the additional benefit of retaining a probabilistic formulation

**Table 2** Recognition accuracy on the Caltech-101 using 15 labeled points per class

Method	GP	SVM	1-NN
	1-vs-All	1-vs-All	
Dense PMK	52.13 $\pm$ 0.69	48.77 $\pm$ 0.95	24.20 $\pm$ 0.48
Spatial PMK	51.90 $\pm$ 0.78	54.26 $\pm$ 0.65	41.10 $\pm$ 0.78
GB	64.15 $\pm$ 0.76	65.87 $\pm$ 0.92	45.58 $\pm$ 0.79
GBDist	58.81 $\pm$ 0.58	65.91 $\pm$ 0.66	50.23 $\pm$ 0.66
Combination	73.95 $\pm$ 1.13	-NA-	-NA-

need not give up the accuracy of other well-used discriminative methods (like the SVM) in order to gain the other benefits of having the probabilistic GP model.

In general, the superior performance of a particular classification algorithm with a specific kernel might be due to several reasons. For any classification strategy to work well, the underlying data must support the assumptions made by the model; whenever the data is favorable to the assumptions of a method, then we can hope that the algorithm will perform well. The point we wish to make here is that GP classification can often provide comparable or slightly improved classification performance when compared to SVMs; we do not have to lose accuracy to gain the predictive uncertainty offered by probabilistic recognition models.

**Effect of Regularization:** We also studied how regularization of the log evidence affects classification performance. Figure 4 demonstrates recognition performance of different regularized and unregularized probabilistic schemes averaged over 10 different splits. As shown in the figure regularization results in a significant improvement in performance, especially  $L_2$ ; the average gain is very high for small number of examples ( $\approx 12\%$  for 2 examples per class). With two examples per class, the average accuracy, 47%, is higher than many prior methods that used a much larger number of label examples.

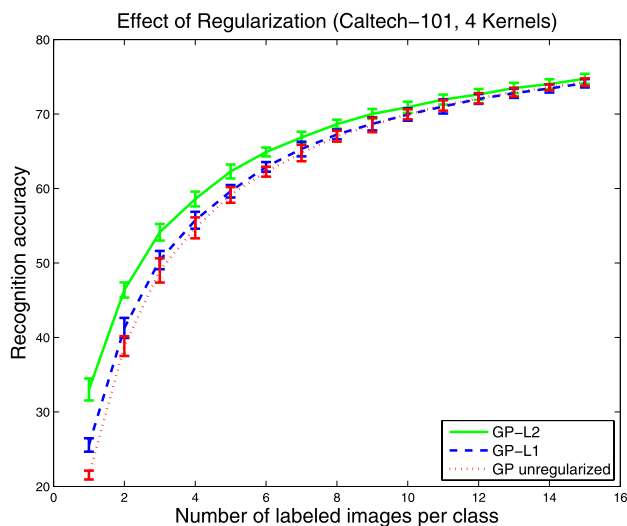
**Global versus Classwise Learning:** Next, we explore the performance difference when parameters are trained globally versus trained separately for each class. Figure 5 (left) shows a scatter plot where each point represents test accuracy obtained on a single train-test split of Caltech-101 data with 5 labeled examples per class. The figure illustrates performance on 35 different train-test splits when combinations of the four kernels are learned. Most of the points lie above the diagonal, which suggests that training parameters globally is better than classwise training. To judge the significance of the results we performed a paired-t test and found the performance difference to be significant at  $p = 10^{-3}$  level. By jointly maximizing parameters for all classes the classifier learns representations that are maximally informative with respect to all the classes simultaneously, as is evident from superior performance across all three choice of ker-

nel combinations. Classwise training has been the method of choice as non-probabilistic alternatives such as SVMs do not have straightforward formulations to optimize the weights globally—however, as described above, GPs avoid this problem by forming the joint likelihood of all the labels given the training data.

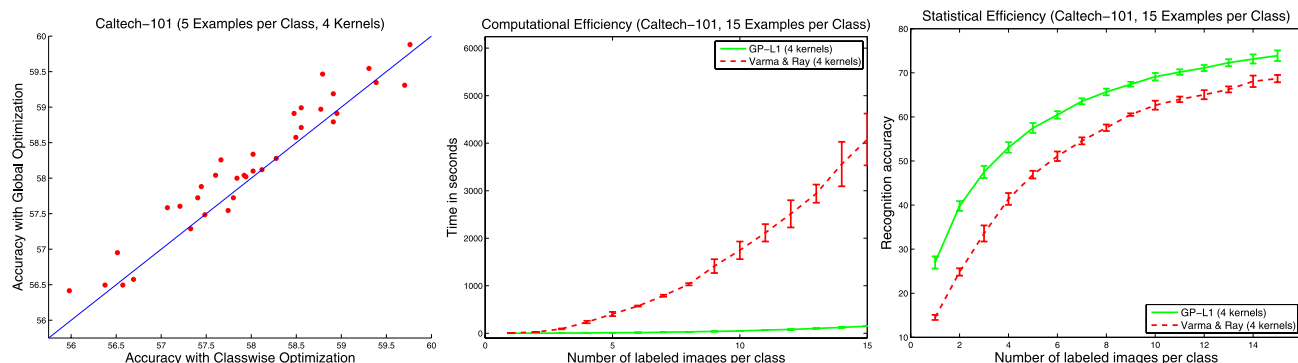
We also compare the computational efficiency of our approach with the scheme of Varma and Ray (2007). Figure 5 (center) shows the time required to learn the kernel combination with four kernels. The method of Varma and Ray (2007) learns kernels separately for each class using the one-vs-all formulation of binary SVMs and takes significantly longer time than the probabilistic combination based on GP. The GP-based approach learns the kernels simultaneously for all the classes and has clear computational advantages vs. training one-vs-all classifiers. In terms of accuracy GP-based formulation significantly outperforms the SVM formulation (see Fig. 5 (right)).

Finally, we compare GP-PMK classification with state-of-the-art supervised visual category learning methods that have been tested on the Caltech-101. Figure 6 shows the performance of an SVM and the classification with GP priors using the PMK along with other recent methods using the same evaluation methodology. The PMK was also earlier used in Grauman and Darrell (2005) with SVMs. Similar to our findings in Fig. 3, we show in Fig. 6 that classification with GP-PMK and an identical kernel (Dense-PMK) actually slightly outperforms the SVM, thus demonstrating the value in the proposed approach.

We also plot the recognition results obtained with the combination of kernels learned via evidence maximization (magenta curve). At 15 points per class we achieve 73.95% accuracy, which shows the power of combining different correspondence kernels in the probabilistic framework. In fact with just eight labeled examples per class the combined kernel achieves an accuracy of 65.68%, beating recognition performance by most of the other methods trained with an individual kernel with *any* size of the training set. Ta-

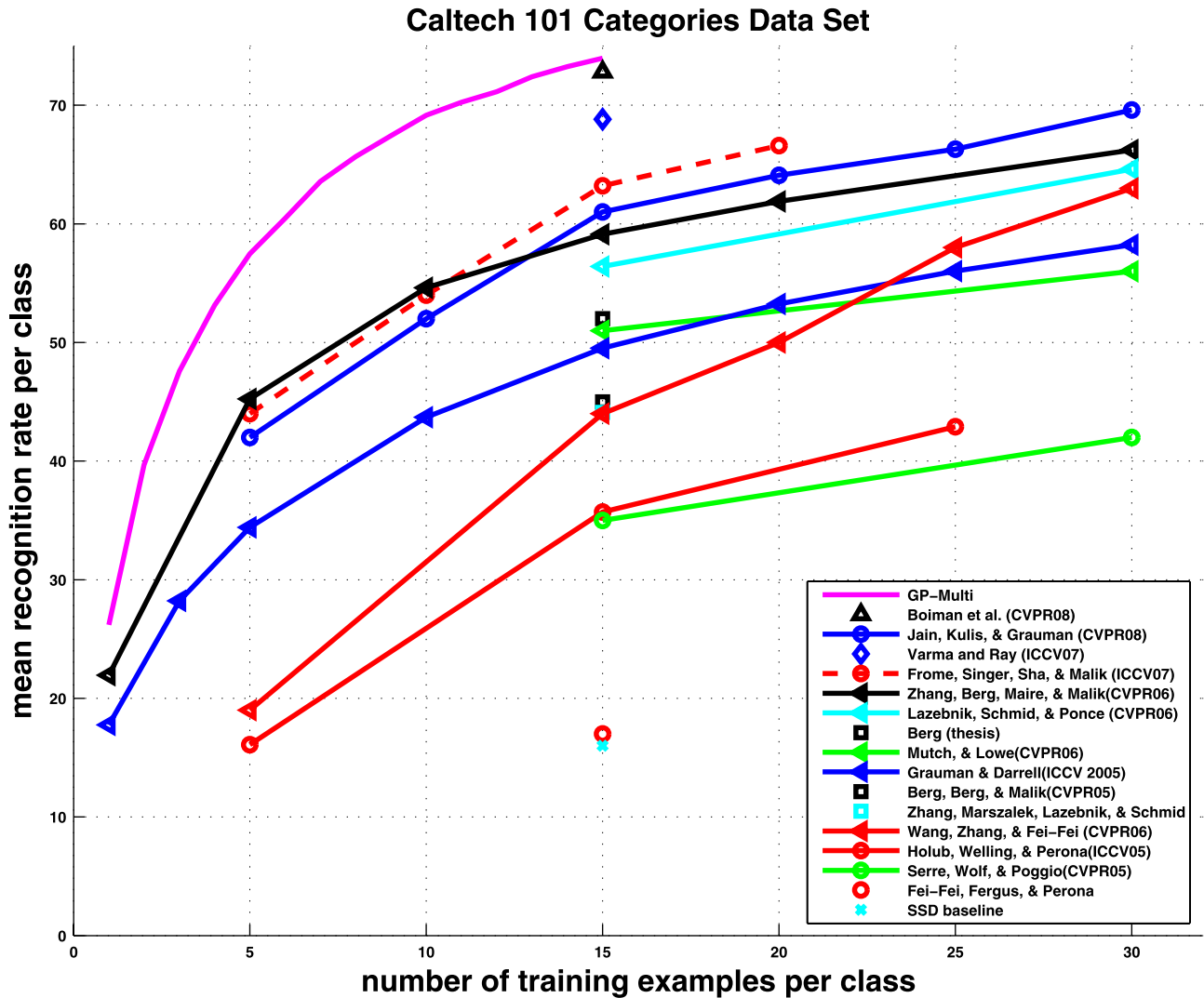


**Fig. 4** Performance comparison of different regularized and unregularized version of the probabilistic kernel combination scheme. A strong regularization results in significantly higher gains specially when the amount of labeled data is sparse



**Fig. 5** (Left) Comparison of accuracy obtained when optimizing kernels globally versus class wise. Comparison of (center) computational efficiency and (right) statistical efficiency for learning using GPs and Varma and Ray (2007) that uses class wise optimization with SVMs





**Fig. 6** Performance comparison of GP-PMK and GP-Multi-Kernel classification with reported results from the literature. Using the same PMK kernel and features, our GP-PMK approach outperforms earlier SVM-PMK results (Grauman and Darrell 2005). Furthermore, with an appropriate combination of various kernels (GP-Multi-Kernel), we

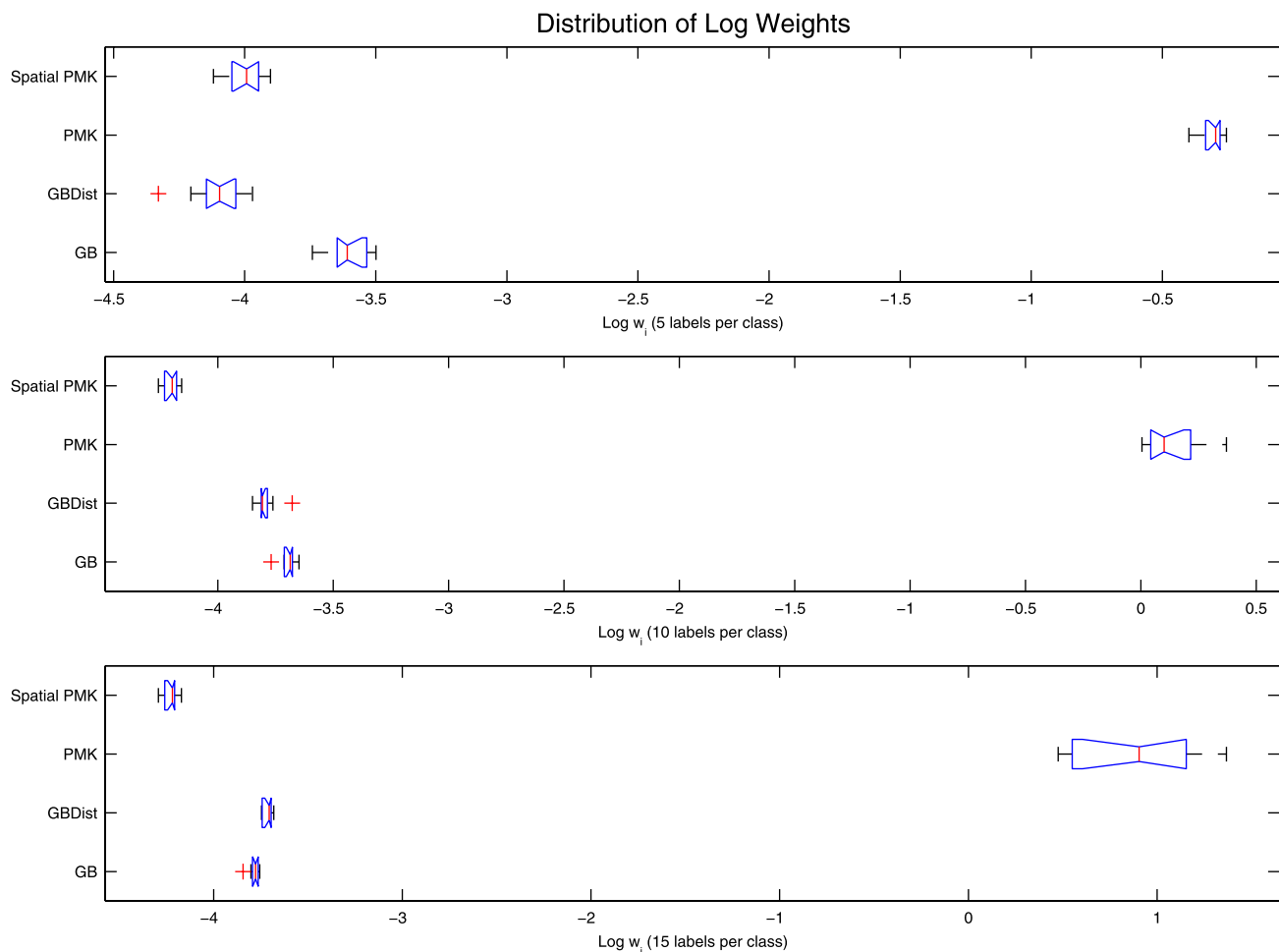
obtain recognition performance very competitive with the state-of-the-art. In fact, we believe our approach is yielding the highest accuracy to-date on this dataset when learning with few training examples (1 to 10 labeled examples per class)

ble 3 provides a direct comparison with results previously reported for other approaches that include kernel combination. Also we plot the accuracy obtained by the method of Varma and Ray (2007) (68.82%) on the Caltech-101 data with the same four kernels, and from the results we can conclude that GPs provide an effective unifying framework for classification as well as kernel combination.

We chose to work with this dataset due to the variety of categories it contains, as well as the large number of existing results published using it. The accuracy of our GP framework is a compelling result, with what appear to be the best performance numbers obtained to-date. We would like to point out that we have certainly benefited from progress in recent years due to other work that has, for example,

determined good features that are applicable for this data. Nonetheless, it is an encouraging result that our method can automatically learn a good combination from a variety of existing matching kernels, and greatly improve the state-of-the-art with quite small labeled training sets.

We would like to point out that the recognition performance can be further improved by using additional kernels, as well as region of interest (ROI) segmentation masks and hierarchical classification. Also, Varma and Ray (2007) in their original paper present a higher recognition accuracy of 87.82% using additional kernels, however, there were doubts about reproducibility of a subset of kernels used in that work. Consequently, we limit our detailed evaluation to the four kernels that are known to be correct and were re-



**Fig. 7** Distribution of weights plotted with MATLAB's *boxplot* command, showing the kernel combinations for various sizes of training set. All the different kernels are adding discriminatory power to the classification task. From these plots we see that the weights can be

learned effectively even with a small number of training examples, as the relative weights are fairly consistent across the plots from top (5 training examples per class) to bottom (15 examples per class)

**Table 3** Accuracy reported with kernel combination on the Caltech-101

Method	Accuracy
<b>Ours (4 Kernels)</b>	<b><math>73.95 \pm 1.13</math></b>
Boiman et al. (2008)	$72.80 \pm 0.39$
Varma and Ray (2007) (4 Kernels)	$68.82 \pm 1.00$
Frome et al. (2007)	$60.30 \pm 0.70$
Lin et al. (2007)	$59.80 \pm \text{NA}$
Zhang et al. (2006)	$59.08 \pm 0.37$
Kumar and Sminchisescu (2007)	$57.83 \pm \text{NA}$

produced by other researchers independently. Our method using the same set of six kernels as used in Varma and Ray (2007) achieves an accuracy of 88.15%.

Also as described above, we compared the kernel combination with GPs and SVMs and our results indicate that

GPs provide a powerful discriminative probabilistic framework for the purpose of object categorization by learning appropriate descriptors. We also note that the segmentation implicit in the ROI kernels as shown in Chum and Zisserman (2007) is quite powerful and should yield higher overall performance than with whole-image kernels. GPs can be extended easily to the hierarchical combination, which together with ROIs, should provide a significant performance boost.

As a final experiment with supervised GP-PMK, we investigate the distribution of learned weights attributed to the different matching kernels. Figure 7 displays the range of weights over the 10 different runs of the algorithm using MATLAB's *boxplot*, which is a graphical representation of the statistics of the log weights.<sup>5</sup> Each row corresponds to

<sup>5</sup>We use log weights for clarity in plotting.

a kernel, and the red line in each row denotes the median over the 10 different runs. The end lines are at the lower and upper quartile values, and the outliers are data with values beyond the ends of the whiskers. We show boxplots for runs with five, 10 and 15 training examples per class. From the figure we see that this distribution is fairly similar for different sizes of training sets. This highlights that we can hope to learn the kernel weights even with very few data points. Further, we also notice that the weight corresponding to the spatial PMK is often fairly high. However, we cannot interpret the weights directly as a measure of importance. This is due to the fact that the scale of each kernel is different; thus, the weight encompasses both the discriminatory power as well as automatic scale adjustment.

Finally, we would like to point out that the mean accuracy with a learned kernel combination using 15 training examples per class over 10 random train-test splits was 73.95% with a standard deviation of 1.13 (see Table 2). The low standard deviation value highlights the stability of the method's classification accuracy with respect to the data used to learn the weights. Moreover, in our experiments we found that GP-based kernel combination (sum of weighted kernels) was extremely stable with respect to the initialization. In fact, in all our experiments we perform the optimization only once (as opposed to optimizing with multiple initializations). As expected, the classification accuracy steadily increases as we increase the number of labeled points; however, the performance is very competitive even with small number of examples per class (for example GP-Multi-Kernel is better than most of the methods for 5 examples per class).

### Active Learning for Object Categorization

In this section, we show the value of active learning in selecting examples to annotate. For these experiments, we test the classification performance on a validation set that includes 10 examples from each class. We first consider the *binary* problem of detecting an object class. Starting with one labeled example per class, the procedure chooses the next image to query from the set of images not in the validation set. We compare the active version of the GP classification with a version that selects the points to query randomly. We again use the mean classification rate per class to compare the methods. We repeat this procedure for 100 different validation sets.

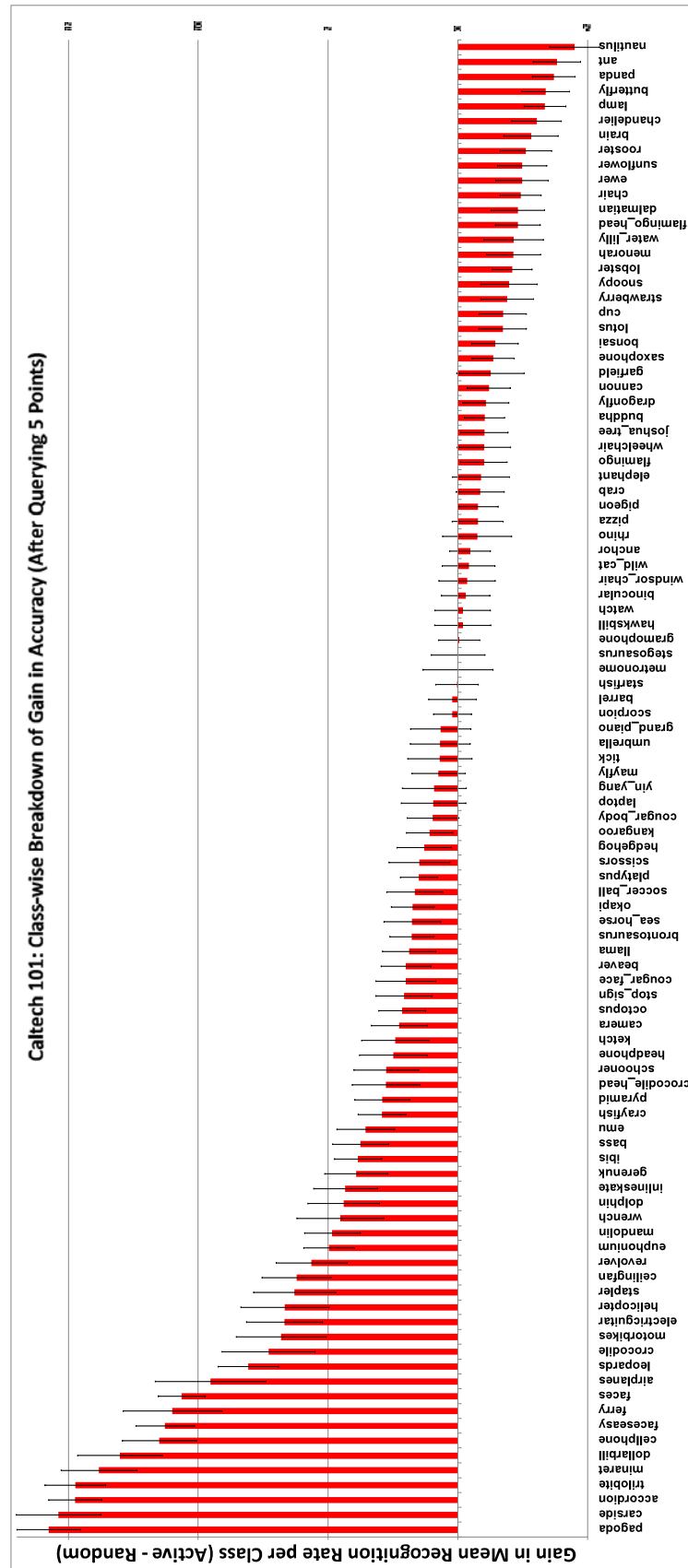
Figure 8 shows the gain in performance on all the 101 binary problems, averaged over the 100 runs, made by the active learning scheme on the validation set after 5 examples are chosen. We can clearly see that for most of the categories there is a significant positive gain showing the benefit of the active learning scheme. Further, Fig. 9 shows the performance on various binary problems as we increase the size

of the training set. The figure depicts that the active learning scheme quickly exploits the uncertainty in its estimates to select appropriate examples to seek the annotation for. The random policy on the other hand performs poorly. The fact that the Caltech-101 dataset has unbalanced numbers of examples per category affects the random sampling policy; it does not work well in these unbalanced scenarios because the training set will usually be skewed towards one class, resulting in poor accuracy. However, selecting points via active learning focuses on points with maximum uncertainty, irrespective of their label, making the procedure highly effective.

Next we describe active learning experiments with the Caltech-4 dataset using multiple feature sampling and pyramid partitioning strategies. The goal here was to investigate the benefits of the proposed scheme across the spectrum of kernels available for the task of object categorization. For this experiment we again experimented with three different flavors of the Pyramid Match Kernel. Besides Dense PMK, we also used PMK computed using only sparse interest points where salient points in the images are detected with a Harris-Affine interest operator (Harris PMK). The third PMK variant was vocabulary guided (Vocabulary Guided PMK) where the features were binned non-uniformly in a data-dependent manner, as in Grauman and Darrell (2006a).

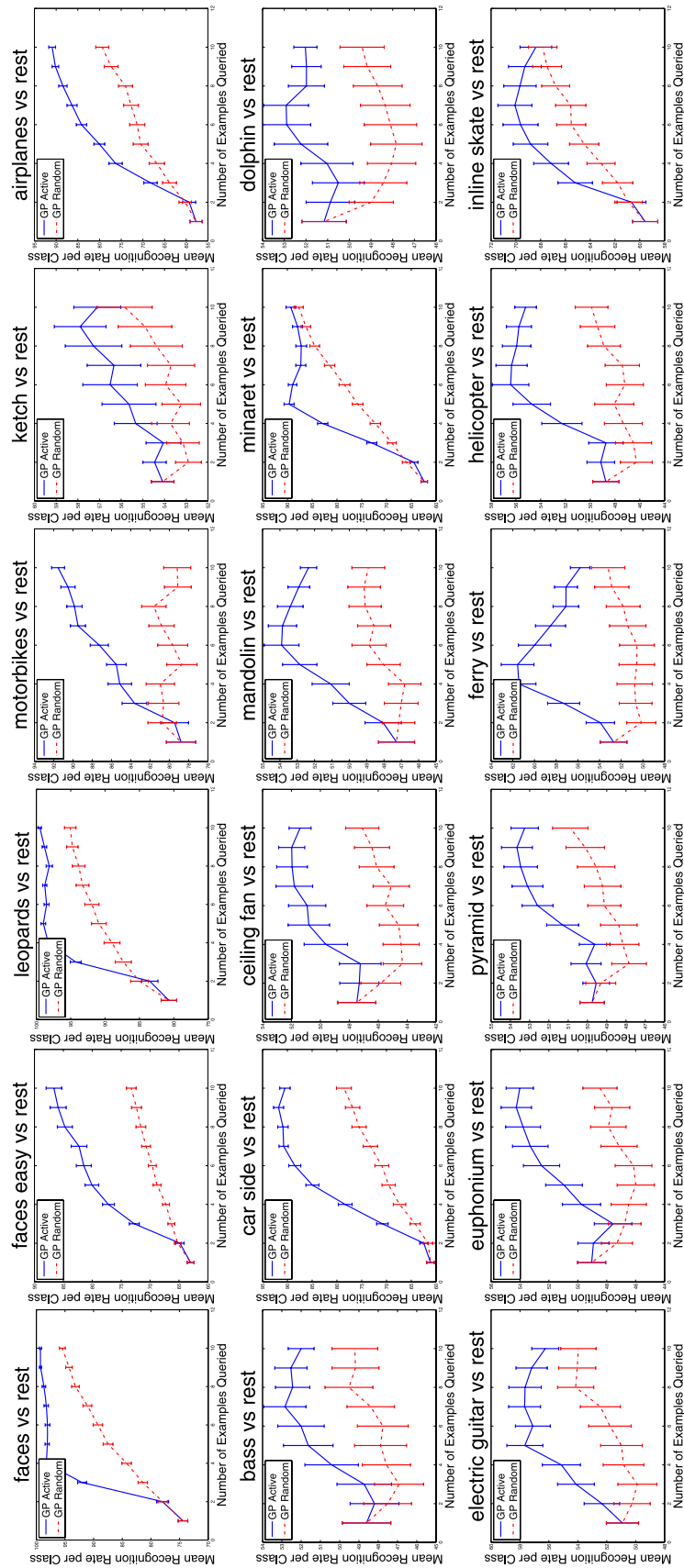
Figure 10 compares different classification approaches on the Caltech-4 database for different kinds of kernels. Essentially, the plot shows mean classification accuracy per class as we vary the total number of examples in the training data. The images not in the training set are considered as the test set to compute the classification performance. We plot the performance of the SVM and the GP classification with and without active learning. We start with one labeled point per class. For the SVM and supervised GP without active learning, we randomly select points as we increase the size of the training set, whereas for the active learning with GP classification and SVM we used all the criteria mentioned in Table 1. This process was repeated 40 times. Figure 10 shows the mean performance together with error bars denoting the standard error.

From Fig. 10 we observe that GP classification again performs competitively with SVM, and using active learning further improves the performance. Note that active learning using uncertainty is superior to active learning using just variance or just the classifier output (margin). The active learning results that use only the variance are far inferior to the other methods, mostly due to the fact that this variant of the criterion does not consider the class label of the data points already in the training set. A similar observation has also been made by Krause et al. (2008). On the other hand, the active learning criterion that uses uncertainty is effective. In fact we can see that a mean accuracy per class close to 90% can be obtained with just 20 labeled examples, whereas

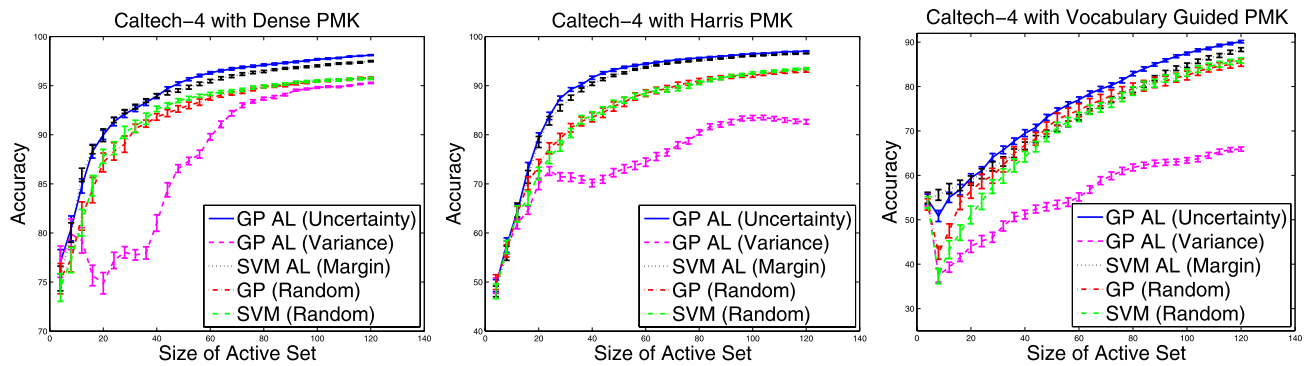


**Fig. 8** Average gain in performance over random selection when choosing five points using active learning. The graph shows the mean and the standard error for 100 runs for all the object classes in Caltech-101 database using GP-PMK





**Fig. 9** Performance comparison of GP classification with active learning and GP with random supervision for various object detection problems (binary) in the Caltech-101 database



**Fig. 10** Active learning on the Caltech-4 database using different kinds of Pyramid Match Kernels and feature types. In each case, our active learning (uncertainty) approach provides significant gains over

traditional passive approaches that label points at random, while the GP classification even shows some gains over the SVM

the non-active learners achieve around 85% accuracy for the same amount of labeled data. This demonstrates that active learning can provide a significant boost in accuracy across different flavors of kernels and feature types used in object categorization. Further, the scheme also makes it possible for the learning algorithm to learn the object classes even with very few labeled examples.

Table 4 shows the confusion matrix resulting after incorporating only 120 examples in the training set using the active learning methodology with Dense PMK. We obtain an overall accuracy of 98.48%, which demonstrates the effectiveness of the framework. The completely supervised GP classification and SVM achieved a mean classification accuracy per class of 95.6% and 95.19%, respectively. This shows that our active learning strategy allows us to learn object categories much more effectively than traditional supervised classification.

## 9 Discussion

The experiments in this paper indicate that classification using GP priors performs competitively with SVM on the object categorization task with Caltech-101 and Caltech-4 data. Of course, these experiments cannot serve as conclusive proof that classification using a GP prior is inherently superior than other classification techniques or vice-versa. Yet for this object categorization task and data, the underlying data density is favorable to the assumptions of the classification model we are using. The experiments in this paper strongly suggest that there is a value in looking at GP classification models for object categorization.

Another important aspect of our framework lies in its seamless extension to kernel combination and active learning. The probabilistic paradigm allows us to exploit the evidence maximization framework to principally combine different correspondence kernels. Furthermore, the Bayesian

**Table 4** Confusion matrix obtained for Caltech-4 database using active learning with the Pyramid Match Kernel (Dense PMK). (120 labeled images, mean accuracy over all the classes = 98.48%)

True Class	Recognized Class			
	Cars	Faces	Airplanes	Motorbikes
<b>Cars</b>	1121	0	0	1
<b>Faces</b>	0	416	0	2
<b>Airplanes</b>	0	2	753	20
<b>Motorbikes</b>	10	0	10	733

formulation lets us incorporate measures such as uncertainty, variance, and expected information gain that could be highly valuable in guiding a supervised learning procedure. One of the challenges in computer vision is the ability to learn object categories with a low number of examples. Humans are able to learn object categories and generalize from a very small number of examples. However, current machine vision systems are far from achieving performance akin to humans. One of the principal differences among humans and existing object classification systems is that humans have the ability to actively seek supervision from the environment and other sources of information. We believe that active learning might enable us to move towards vision systems that require few examples to learn successfully.

## 10 Conclusion and Future Work

We have presented a discriminative probabilistic framework based on Gaussian Process priors and the local feature-based correspondence kernels, and have shown its utility for visual category recognition. Gaussian Process regression provides a principled framework to combine different correspondence kernels, which results in performance superior to individual kernels. Further, the modeling with

Gaussian Process priors provides direct estimates of prediction uncertainty using a smoothness prior that captures a correspondence-based notion of similarity between sets of local image features. We introduced an active learning method for visual category recognition based on the GP-PMK uncertainty estimates, and showed empirically that active learning can be used to achieve very good recognition results using far fewer training images than standard supervised learning approaches.

We plan to extend the framework to adopt non-Gaussian noise models, and investigate other active learning formulations such as value of information and/or criteria previously developed for sparsifying GPs (Lawrence et al. 2002). By incorporating decision-theoretic formulations we should be able to learn object categories within a given budget. We also plan to extend the model to handle multiple objects in the same image, incorporate semi-supervised learning, and explore sparse GP techniques for large training sets.

**Acknowledgements** We thank the following for providing kernel matrices: Alex Berg, Anna Bosch, Jitendra Malik and Andrew Zisserman. We also wish to thank Manik Varma for his help in obtaining the required data and for many helpful discussions. Kristen Grauman is supported in part by NSF CAREER award #0747356, a Microsoft Research New Faculty Fellowship, Texas Higher Education Coordinating Board award #003658-01-40-2007, DARPA VIRAT, and the Henry Luce Foundation.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution Noncommercial License which permits any noncommercial use, distribution, and reproduction in any medium, provided the original author(s) and source are credited.

## References

- Abramson, Y., & Freund, Y. (2004). *Active learning for visual object recognition* (Technical report). UCSD.
- Belongie, S., Malik, J., & Puzicha, J. (2001). Matching shapes. In *ICCV*.
- Berg, A., & Malik, J. (2001). Geometric blur for template matching. In *CVPR*.
- Boiman, O., Shechtman, E., & Irani, M. (2008). In defense of nearest-neighbor based image classification. In *CVPR*.
- Bosch, A., Zisserman, A., & Muñoz, X. (2007). Representing shape with a spatial pyramid kernel. In *CIVR*.
- Chang, C., & Lin, C. (2001). *LIBSVM: a library for SVMs*.
- Chang, E. Y., Tong, S., Goh, K., & Chang, C. (2005). Support vector machine concept-dependent active learning for image retrieval. *IEEE Transactions on Multimedia*.
- Chum, O., & Zisserman, A. (2007). An exemplar model for learning object classes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Evgeniou, T., Pontil, M., & Poggio, T. (2000). Regularization networks and support vector machines. *Advances in Computational Mathematics*, 13(1).
- Fei-Fei, L., Fergus, R., & Perona, P. (2006). One-shot learning of object categories. *IEEE Transaction on Pattern Recognition and Machine Intelligence*.
- Fergus, R., Perona, P., & Zisserman, A. (2003). Object class recognition by unsupervised scale-invariant learning. In *CVPR*.
- Freund, Y., Seung, H. S., Shamir, E., & Tishby, N. (1997). Selective sampling using the query by committee algorithm. *Machine Learning*, 28(2–3).
- Frome, A., Singer, Y., Sha, F., & Malik, J. (2007). Learning globally-consistent local distance functions for shape-based image retrieval and classification. In *ICCV*.
- Grauman, K., & Darrell, T. (2005). The pyramid match kernel: Discriminative classification with sets of image features. In *ICCV*.
- Grauman, K., & Darrell, T. (2006a). Approximate correspondences in high dimensions. In *NIPS*.
- Grauman, K., & Darrell, T. (2006b). Unsupervised learning of categories from sets of partially matching image features. In *CVPR*.
- Kadir, T., & Brady, M. (2003). Scale saliency: A novel approach to salient feature and scale selection. In *International conference visual information engineering*.
- Kapoor, A., Grauman, K., Urtasun, R., & Darrell, T. (2007). Active learning with Gaussian processes for object categorization. In *ICCV*.
- Kim, H. C., Kim, D., Ghahramani, Z., & Bang, S. Y. (2006). Appearance-based gender classification with Gaussian processes. *Pattern Recognition Letters*.
- Krause, A., Singh, A., & Guestrin, C. (2008). Near-optimal sensor placements in Gaussian processes: Theory, efficient algorithms and empirical studies. In *JMLR*.
- Kumar, A., & Smorchil, C. (2007). Support kernel machines for object recognition. In *ICCV*.
- Lawrence, N. (2004). Gaussian process latent variable models for visualisation of high dimensional data. In *NIPS*.
- Lawrence, N., Seeger, M., & Herbrich, R. (2002). Fast sparse Gaussian process method: Informative vector machines. In *NIPS*.
- Lazebnik, S., Schmid, C., & Ponce, J. (2006). Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR*.
- Lin, Y. Y., Liu, T. Y., & Fuh, C. S. (2007). Local ensemble kernel learning for object category recognition. In *CVPR*.
- Lowe, D. (2004). Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2).
- MacKay, D. (1992) Information-based objective functions for active data selection. *Neural Computation*, 4(4).
- McCallum, A. K., & Nigam, K. (1998). Employing EM in pool-based active learning for text classification. In *ICML*.
- Mikolajczyk, K., & Schmid, C. (2001). Indexing based on scale invariant interest points. In *ICCV*.
- Mikolajczyk, K., & Schmid, C. (2004). Scale and affine invariant interest point detectors. *IJCV*, 1(60), 63–86.
- Minka, T. P. (2001). *A family of algorithms for approximate Bayesian inference*. PhD thesis, MIT.
- Moosmann, B. T. F., & Jurie, F. (2007). Fast discriminative visual codebooks using randomized clustering forests. In *NIPS*.
- Muslea, I., Minton, S., & Knoblock, C. A. (2002). Active + semi-supervised learning = robust multi-view learning. In *ICML*.
- Nister, D., & Stewenius, H. (2006). Scalable recognition with a vocabulary tree. In *CVPR*.
- Nowak, E., Jurie, F., & Triggs, B. (2006). Sampling strategies for bag-of-features image classification. In *ECCV*.
- Rasmussen, C. E., & Williams, C. (2006). *Gaussian processes for machine learning*. Cambridge: MIT Press.
- Seeger, M. (2004). Gaussian processes for machine learning. *International Journal of Neural Systems*, 14(2).
- Shen, Y., Ng, A., & Seeger, M. (2006). Fast Gaussian process regression using kd-trees. In *NIPS*.
- Sivic, J., & Zisserman, A. (2003). Video Google: a text retrieval approach to object matching in videos. In *ICCV*.
- Sivic, J., Russell, B., Efros, A., Zisserman, A., & Freeman, W. (2005). Discovering object categories in image collections. In *ICCV*.
- Snelson, E., & Ghahramani, Z. (2006). Sparse Gaussian processes using pseudo-inputs. In *NIPS*.

- Sudderth, E., Torralba, A., Freeman, W., & Willsky, A. (2005). Describing visual scenes using transformed Dirichlet processes. In *NIPS*.
- Tong, S., & Koller, D. (2000). Support vector machine active learning with applications to text classification. In *ICML*.
- Tresp, V. (2000). Mixtures of Gaussian processes. In *NIPS*.
- Tsang, I. W.-H., & Kwok, J. T.-Y. (2006). Efficient hyperkernel learning using second-order cone programming. *IEEE Transactions on Neural Networks*.
- Urtasun, R., & Darrell, T. (2008). Local probabilistic regression for activity-independent human pose inference. In *CVPR*.
- Urtasun, R., Fleet, D. J., Hertzman, A., & Fua, P. (2005). Priors for people tracking from small training sets. In *ICCV*.
- Urtasun, R., Fleet, D. J., & Fua, P. (2006). Gaussian process dynamical models for 3d people tracking. In *CVPR*.
- Varma, M., & Ray, D. (2007). Learning the discriminative power-invariance trade-off. In *ICCV*.
- von Ahn, L., & Dabbish, L. (2004). Labeling images with a computer game. In *ACM CHI*.
- von Ahn, L., Liu, R., & Blum, M. (2006). Peekaboom: A game for locating objects in images. In *ACM CHI*.
- Wallraven, C., Caputo, B., & Graf, A. (2003). Recognition with local features: the kernel recipe. In *ICCV*.
- Williams, C., & Barber, D. (1998). Bayesian classification with Gaussian processes. *IEEE Transaction on Pattern Recognition and Machine Intelligence*, 20(12), 1342–1351.
- Williams, O. (2006). A switched Gaussian process for estimating disparity and segmentation in binocular stereo. In *NIPS*.
- Zhang, H., Berg, A., Maire, M., & Malik, J. (2006). SVM-KNN: Discriminative nearest neighbor classification for visual category recognition. In *CVPR*.
- Zhu, X., Lafferty, J., & Ghahramani, Z. (2003). Combining active learning and semi-supervised learning using Gaussian fields and harmonic functions. In *Workshop on the continuum from labeled to unlabeled data in machine learning and data mining at ICML*.