

I Want To Know More — Efficient Multi-Class Incremental Learning Using Gaussian Processes

Alexander Lütz, Erik Rodner and Joachim Denzler

Computer Vision Group, Friedrich Schiller University Jena
{alexander.freytag, erik.rodner, joachim.denzler}@uni-jena.de
<http://www.inf-cv.uni-jena.de>

One of the main assumptions in machine learning is that sufficient training data is available in advance and batch learning can be applied. However, because of the dynamics in a lot of applications, this assumption will break down in almost all cases over time. Therefore, classifiers have to be able to adapt themselves when new training data from existing or new classes becomes available, training data is changed or should be even removed. In this paper, we present a method allowing for efficient incremental learning of a Gaussian process classifier. Experimental results show the benefits in terms of needed computation times compared to building the classifier from the scratch. In addition we highlight the general benefits of incremental learning.

1 Introduction

In the last decade, research in visual object recognition has focused mostly on batch learning, *i.e.*, a classifier is build using a pre-defined and fixed set of training examples. However, if we think of the main goal of closing the gap between human and computer vision, which still exists when focusing on the recognition performance in real-world scenarios, batch learning is not appropriate. Even by using the large number of training examples already present in current image databases [1], a fixed learning system would not be able to adapt to new tasks and object classes. Due to this reason, it is important to incrementally learn a classifier in an efficient manner, *i.e.*, by utilizing previously calculated model parameters. Such a procedure is often referred to as *online or incremental learning*.

Incremental learning has been studied for various types of classifiers. The work of Cauwenberghs and Poggio [2] presents how to carefully update the set of support vectors and their weights when incrementally adding new training data to existing SVM classifiers. Building on the algorithmic ideas of [2], Tax and Laskov [3] focus on incrementally learning a support vector data description classifier, which is a one-class classification method. For on-

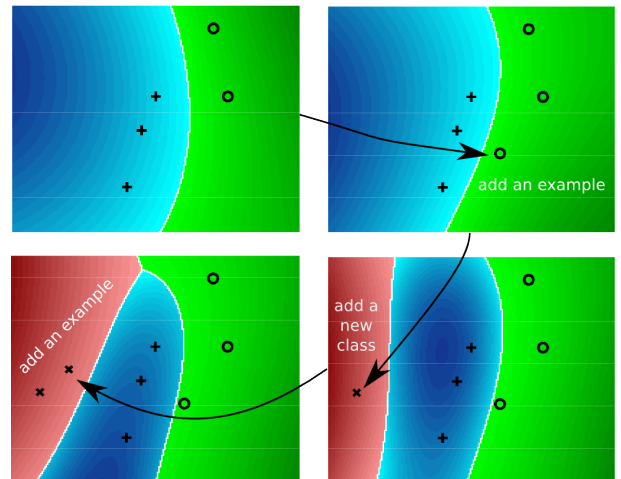


Fig. 1. Multi-class incremental learning: adding examples to already existing or new classes. Color brightness shows the classifier class score.

line learning on a class level, *i.e.*, when examples of new classes become available, the methods proposed by Yeh *et al.* [4] can be used to perform online learning of one-vs-all SVMs. Csató and Opper [5] show how to apply the idea of Bayesian online learning to sparse Gaussian process (GP) models.

This paper focuses on exploring the advantages of a one-vs-all Gaussian process classifier in incremental learning scenarios. In contrast to previous work, we do not utilize a sparse GP

representation or approximate inference techniques. We briefly review the GP framework and present the update formulas necessary to add new examples. The presented approach is applied to a 2d toy example in Fig. 1.

2 Regression with Gaussian Processes

In the following, we briefly describe the basic principles of Gaussian process regression. We are given a set of input examples $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \Omega^n$, where Ω is the space of possible inputs. Let us further assume a vector $\mathbf{t}_L \in \mathbb{R}^n$ containing the corresponding function values is given. For the regression case, we are interested in estimating the relationship $f : \Omega \rightarrow \mathbb{R}$ between inputs and outputs. Because of the finite size of \mathbf{X} , we are not able to determine a single unique solution for f . Therefore, we model f to be drawn from a distribution over functions, *i.e.*, a random process with index set \mathbb{R} . One possible choice for an underlying random process is a Gaussian process $\mathcal{GP}(\mu, \kappa)$ [6]. The mean function $\mu : \Omega \rightarrow \mathbb{R}$ with

$$\mu(\mathbf{x}) = \mathbb{E}_f[f(\mathbf{x})]$$

is used to specify our prior assumptions about the mean value of the function at certain positions. The covariance function $\kappa : \Omega^2 \rightarrow \mathbb{R}$:

$$\kappa(\mathbf{x}, \mathbf{x}') = \mathbb{E}_f[(f(\mathbf{x}) - \mu(\mathbf{x}))(f(\mathbf{x}') - \mu(\mathbf{x}'))]$$

models the expected covariance of the function values, which is directly related to the expected local variability of functions drawn from the process. Learning from output and input values and estimating the output t_* of a new example \mathbf{x}_* can be done with Bayesian inference. In the following, we assume a zero mean function and that observed values \mathbf{t}_L of f are corrupted with independent Gaussian noise with variance σ_n^2 . In this special case, Bayesian inference leads to a closed-form solution for the posterior [6]:

$$p(t_* | \mathbf{x}_*, \mathbf{X}, \mathbf{t}_L) = \mathcal{N}(t_* | \mu_*, \sigma_*^2) . \quad (1)$$

This means that the posterior is Gaussian with the following mean and variance:

$$\mu_* = \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{t}_L \quad (2)$$

$$\sigma_*^2 = k_{**} - \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}_* + \sigma_n^2 . \quad (3)$$

With a slight abuse of notation, we use $k_{**} = \kappa(\mathbf{x}_*, \mathbf{x}_*)$, $\mathbf{k}_* = \kappa(\mathbf{X}, \mathbf{x}_*)$ and $\mathbf{K} = \kappa(\mathbf{X}, \mathbf{X})$

as abbreviations for the kernel values of the training set and the new example. It should be noted that the possibility to estimate the variance of the classifier output is an important benefit for advanced incremental learning, because it allows for active learning [7] and out-of-vocabulary detection [8].

3 GP Regression for Multi-class Classification

As we have seen in the previous section, the GP model and a Gaussian noise assumption leads to simple inference equations. If we consider non-Gaussian noise for classification problems, there are no closed-form solutions available and approximation methods are necessary [6]. However, it has been demonstrated in several papers [7, 8, 9] that applying GP regression directly to binary classification problems, without considering the discrete nature of the labels $t \in \{-1, 1\}$, gives comparable results. Therefore, we assume the discrete labels to be generated by a real-valued function which we estimate using GP regression.

For multi-class classification problems with $t \in \{1, \dots, M\}$, the one-vs-all technique can be applied, which leads to the following scores for each class $c \in \{1, \dots, M\}$:

$$\mu_*^{(c)} = \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{t}_L^{(c)} , \quad (4)$$

with $\mathbf{t}_L^{(c)} \in \{-1, 1\}^n$ representing the vector of binary training labels for class c [7]. The final score of the multi-class classifier is achieved by taking the maximum of the scores of all classes:

$$\mu_*^{mc} = \max_{c=1 \dots M} \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 \cdot \mathbf{I})^{-1} \mathbf{t}_L^{(c)} , \quad (5)$$

and returning the corresponding class c .

Since the predictive variance σ_*^2 for a new test example \mathbf{x}_* is not affected by the training labels (see eq. (3)) the resulting variance is the same for all classes.

4 Efficient Updates of the Kernel Matrix

As shown in equations (2), (3), (4), we only need \mathbf{t}_L , \mathbf{k}_* , and $(\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1}$ for computing the predictive mean and predictive variance of a new input \mathbf{x}_* , respectively. Therefore, if training data is changed or new training data becomes available, we have to consider methods for updating the corresponding variables.

In the following, we show how to update the inverse of the regularized kernel matrix ($\mathbf{K} + \sigma_n^2 \mathbf{I}$) in an efficient manner. Another option is to update the Cholesky factorization, as described in [10, 11].

Let \mathbf{K} be a $n \times n$ -Matrix, \mathbf{U} a $n \times p$ -Matrix, \mathbf{C} a $p \times p$ -Matrix, and \mathbf{V} a $p \times n$ -Matrix. Let further \mathbf{K} and \mathbf{C} be invertible. Woodbury's Formula [12] states that

$$\mathbf{K}'^{-1} = (\mathbf{K} + \mathbf{UCV})^{-1} = \mathbf{K}^{-1} - \mathbf{K}^{-1} \mathbf{U} (\mathbf{C}^{-1} + \mathbf{VK}^{-1} \mathbf{U})^{-1} \mathbf{VK}^{-1}. \quad (6)$$

If we already precomputed \mathbf{K}^{-1} and \mathbf{K} is only changed by adding \mathbf{UCV} , we are able to compute \mathbf{K}'^{-1} in a more efficient manner compared to standard approaches if $p \ll n$. For the following derivations we set \mathbf{C} to the identity matrix.

Update of a single example Changes in only one specific training example \mathbf{x}_i to \mathbf{x}'_i can be modeled by

$$\mathbf{K}' = (\mathbf{K} + \mathbf{UV}) \quad (7)$$

with $\mathbf{U} = [\mathbf{a}^{(i)} \ \mathbf{e}_i]$, $\mathbf{V} = [\mathbf{e}_i^T \ \mathbf{a}^{(i)T}]^T$, \mathbf{e}_i being the i th unity vector, and

$$\mathbf{a}_q^{(i)} = \begin{cases} \kappa(\mathbf{x}'_i, \mathbf{x}_q) - \mathbf{K}_{i,q} & \text{if } i \neq q \\ \frac{1}{2} (\kappa(\mathbf{x}'_i, \mathbf{x}_q) - \mathbf{K}_{i,q}) & \text{if } i = q \end{cases} \quad (8)$$

as the vector of differences between old and new kernel values for \mathbf{x}_i . Because of $\text{rk}(\mathbf{UV}) = 2$, we call \mathbf{K}' a rank-2-update of \mathbf{K} . Thereby, we can directly apply (6) for computing the new inverse of the kernel matrix efficiently. The resulting asymptotic runtimes are derived in section 5.

Update of multiple examples If not one but a set $\mathcal{S} = \{\mathbf{x}_{s_1}, \dots, \mathbf{x}_{s_k}\}$ of training examples changes, the update can be done either by *iteratively* changing one example or *directly* updating all examples simultaneously. For the direct update, the presented ideas for the rank-2-update can be generalized in the following way. We build the vector of differences between old and new kernel values

$$\mathbf{a}_q^{(i)} = \begin{cases} \kappa(\mathbf{x}'_i, \mathbf{x}_q) - \mathbf{K}_{i,q} & \text{if } \mathbf{x}_q \notin \mathcal{S} \\ \frac{1}{2} (\kappa(\mathbf{x}'_i, \mathbf{x}_q) - \mathbf{K}_{i,q}) & \text{if } \mathbf{x}_q \in \mathcal{S} \end{cases} \quad (9)$$

with $i \in \mathcal{S} = \{s_1, \dots, s_k\}$. Furthermore, we set the two matrices \mathbf{U} and \mathbf{V} with $\mathbf{U} = [\mathbf{a}^{(s_1)} \dots \mathbf{a}^{(s_k)} \ \mathbf{e}_{s_1} \dots \mathbf{e}_{s_k}]$ and $\mathbf{V} = [\mathbf{e}_{s_1}^T \dots \mathbf{e}_{s_k}^T \ \mathbf{a}^{(s_1)T} \dots \mathbf{a}^{(s_k)T}]^T$ such that \mathbf{UV} transforms \mathbf{K} to \mathbf{K}' , which means that \mathbf{K}' is a rank- $2k$ -Update of \mathbf{K} . Thereby, we again can apply (6) for efficiently computing the modified inverse \mathbf{K}'^{-1} . Resulting asymptotic runtimes for both approaches are given in section 5.

Efficient Incremental Learning The previous explanations focused on efficiently updating \mathbf{K} when existing examples are modified. If new training data becomes available, we can directly make use of these ideas. For adding one example (rank-2-update), we increase the size of \mathbf{K} by one, whereas for k new examples (rank- $2k$ -update) we attach k new rows and columns to \mathbf{K} with value one on the main diagonal and zero elsewhere. After preparing \mathbf{K} in this way, we can apply the presented update formulas and end up with a GP classifier incrementally trained in an efficient way.

Efficient Decremental Learning When thinking about a system for lifelong learning, the most prominent tasks are to efficiently deal with new or altered data. Nonetheless, such a system has also to be capable of dealing with data which becomes invalid over time. This task is known as decremental learning and can be seen as counterpart to the previously presented problem.

To enable efficient decremental learning, we apply the same approaches as already used for incremental learning, but in an inverse manner. More precisely, if the current example i becomes invalid, we use the update rules from eq. (6) such that the resulting row and column contains one on the i th position and zero elsewhere. Finally, we remove the i th row and column of the kernel matrix. The same holds for deleting several examples using either the iterative or the direct update.

Updates for multi-class classification The presented approach is independent of the number of classes taken into account. As presented in eq. (4) every one-vs-all classifier uses the same matrix \mathbf{K} and differs merely in the corresponding binary label vector $\mathbf{t}_L^{(c)}$. Therefore, only a single update of the shared inverse covariance matrix \mathbf{K}'^{-1} is required. Even for new classes

a corresponding classifier can be trained with minimal effort by computing the new binary label vector and using the efficiently updated \mathbf{K}'^{-1} .

5 Asymptotic Runtimes

Baseline Training a GP classifier with $n + k$ training examples and ignoring any previous knowledge needs $\mathcal{O}((n + k)^3)$ operations. This term is dominated by the time spent for computing the inverse of \mathbf{K} . Even with methods such as Cholesky decomposition, which takes advantage of the positive definiteness of \mathbf{K} , the main complexity stays cubically in the number of examples.

Update of a single example If k equals one, we just add a single example. If we have a closer look on eq. (6) we notice, that the most time-consuming operation is the multiplication of $\mathbf{K}^{-1}\mathbf{U}$ with time requirement quadratic in the number of examples. Since we update a single example, the matrix $(\mathbf{C}^{-1} + \mathbf{V}\mathbf{K}^{-1}\mathbf{U})$ is of size 2×2 and can therefore be inverted with a constant number of operations. Consequently, the computations necessary for updating the inverse can be performed within $\mathcal{O}((n + 1)^2)$ operations. Since the kernel matrix is of size $(n + 1)^2$ and all entries of its inverse have to be updated, this asymptotic performance is the best achievable one.

Update of multiple examples Using the presented *iterative* IL approach, *i.e.*, performing k updates with a single new example, the asymptotic bound is $\mathcal{O}(k(n + k)^2)$. With the *direct* approach, we can perform the computation in $\mathcal{O}(k^3 + (n + k)^2)$ operations. Note, that for this approach the involved constants are relatively high, which is caused by the fact that $(\mathbf{C}^{-1} + \mathbf{V}\mathbf{K}^{-1}\mathbf{U})$ (see eq. (6)) is not symmetric and can therefore not be inverted using Cholesky decomposition.

As a rule of thumb we can observe, that if $n \in \mathcal{O}(k)$, the iterative update should be preferred, whereas for large datasets with only small changes ($k \ll n$) the direct approach should be the method of choice. We validate this fact in our experiments given in section 6. In all cases the update of the corresponding \mathbf{t}_L can be done in constant time. If an example of a new class becomes available the new label vec-

tor has to be computed requiring $\mathcal{O}(n + k)$ and not affecting the complexity at all.

6 Experiments on the Caltech-101 dataset

Our update rules do not use any approximations, therefore, the differences between the results of batch and incremental learning are marginal and caused by numerical differences. The recognition rates are not affected. Due to this reason, we concentrate on *evaluating the runtime* in this section only. How incrementally added examples affect the resulting performances is evaluated in section 7.

Data We follow the experimental setup of [4] and use the well-known Caltech-101 dataset [13]. For each step, we randomly sample n images and train a GP regression classifier based on these examples. Afterwards, we additionally sample k examples to add them incrementally or for batch training. Using suitable features and multi-class GP classification allows average recognition rates of up to 74% with 15 training examples on this dataset [7].

Features For realistic evaluations, we again follow the setup of [4] and choose feature vectors corresponding to the spatial pyramid match kernel [14] as a representation for each image. Given computation times include the time needed for feature calculation.

Evaluation The experimental results shown in Fig. 2 confirm a significant benefit arising from incrementally training the classifier compared to batch training. Note that for $k = 1$, iterative and direct updates result in the same computations. As expected the resulting gain depends on the ratio of n and k . In our experiments, the speed-up was up to a factor of 100. Additionally, the theoretical result that a direct IL approach is most useful when k is large but $\frac{k}{n}$ is small, can be verified.

7 Experiments on the 15 scenes dataset

Data We use the 15Scenes database [14] to experimentally analyze the benefits of incremental learning strategies. As previously suggested [15], we scale all images to a size of 256×256 pixels to avoid results biased by the specific image sizes of different classes. We sample different numbers of images from the training set with respect to the known classes in

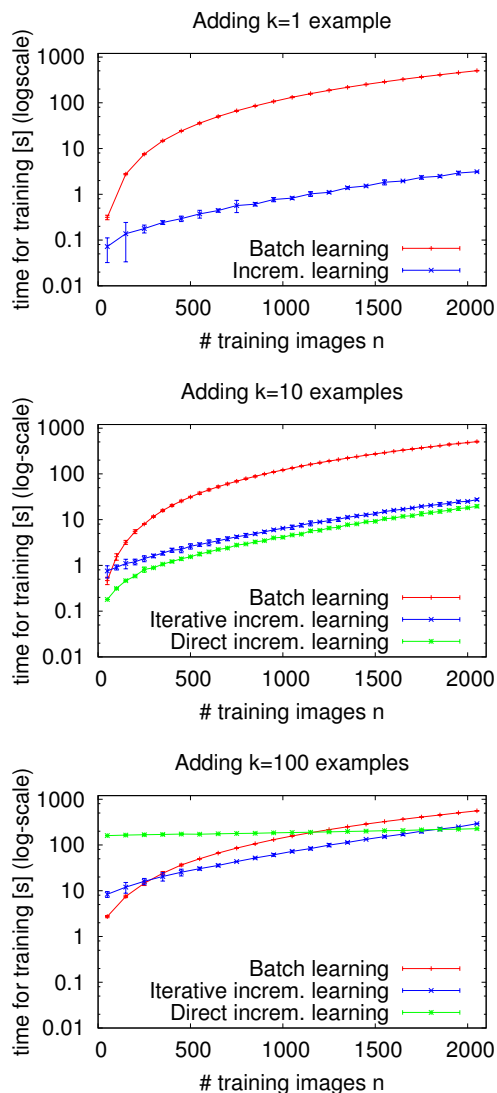


Fig. 2. Time measurements for incremental compared to batch learning with GP regression. Times are given in seconds and displayed in logarithmic scale.

each step. For testing, 100 images are randomly sampled for every class, resulting in 1,500 test images. Results are averaged over 10 runs.

Features We represent images by bag of visual words (BoV) features computed using the toolkit provided with the ILSVRC'10 database [16].

Evaluation In a first experiment, we increase the number of known classes during training over time but perform evaluations on all available classes. As can be seen in Fig. 3, the performance increases significantly with the number of categories known to the system. This behaviour is as expected, since only classes known to the system can be classified correctly. Obviously, this fact validates the necessity of

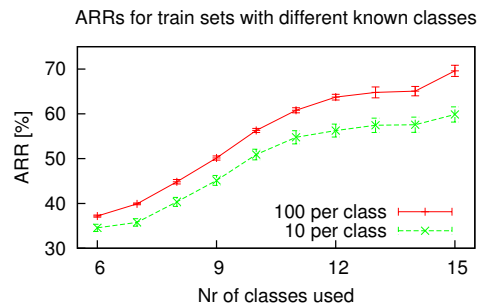


Fig. 3. Recognition rates for different numbers of classes known during training. Results are conducted on the 15 scenes dataset and averaged over 10 runs.

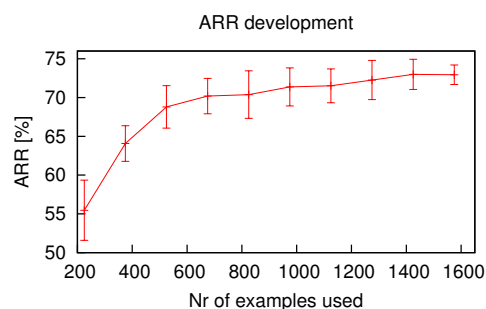


Fig. 4. Development of recognition rates while incrementally adding new examples. In each step, we add 10 randomly chosen examples per class. Results are conducted on the 15 scenes dataset and averaged over 10 runs.

incrementally increasing the data accessible by the system.

In a second experiment, we train the classifier with all available 15 classes, but use only a small fraction with 5 examples per class. After evaluating the resulting performance on 100 distinct examples per class, we make new examples available in an incremental way, *i.e.*, we incrementally add 10 examples per class in one step. Fig. 4 shows recognition rates achieved after different steps. It can clearly be seen, that the system benefits from the availability of new examples, especially when only a limited amount of training data was available before.

8 Conclusions and Future Work

The aim of this paper was to show that Gaussian Processes, which are a powerful machine learning tool, can be applied to incremental learning scenarios using efficient retraining methods. We presented in-depth how to utilize the well-known Woodbury formula to incrementally up-

date a GP-based classifier and performed experiments for the tasks of object recognition. To incrementally add several new examples, we introduced two strategies: iterative update and direct update. Both approaches reduce the necessary computations for retraining from cubic to quadratic in the number of known examples. The experimental results clearly demonstrate the large benefit of incremental learning for improved classification capabilities.

Future research will focus on incremental kernel hyperparameter optimization, which is especially important when combining several different feature types with multiple kernel learning. Another interesting direction is to perform updates of a tree-based GP classifier [17]. This would allow us to use large-scale training datasets [1] as initial knowledge base for incremental learning.

References

- [1] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009, pp. 248–255.
- [2] G. Cauwenberghs and T. Poggio, "Incremental and decremental support vector machine learning," in *NIPS*, 2001, vol. 13, pp. 409–415.
- [3] D. Tax and P. Laskov, "Online svm learning: from classification to data description and back," in *Workshop on Neural Networks for Signal Processing (NNSP'03)*, 2003, pp. 499–508.
- [4] T. Yeh and T. Darrell, "Dynamic visual category learning," in *CVPR*, 2008, pp. 1–8.
- [5] L. Csató and M. Opper, "Sparse representation for gaussian process models," in *NIPS*, 2000, pp. 444–450.
- [6] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. The MIT Press, 01 2006.
- [7] A. Kapoor, K. Grauman, R. Urtasun, and T. Darrell, "Gaussian processes for object categorization," *IJCV*, vol. 88, pp. 169–188, 2010.
- [8] M. Kemmler, E. Rodner, and J. Denzler, "One-class classification with gaussian processes," in *ACCV*, 2010, pp. 489–500.
- [9] E. Rodner, D. Hegazy, and J. Denzler, "Multiple kernel gaussian process classification for generic 3d object recognition from time-of-flight images," in *Proceedings of the IVCNZ'10*, 2010.
- [10] D. Nguyen-Tuong, M. W. Seeger, and J. Peters, "Model learning with local gaussian process regression." *Advanced Robotics*, vol. 23, no. 15, pp. 2015–2034, 2009.
- [11] M. W. Seeger, "Low rank updates for the cholesky decomposition," University of California at Berkeley, Tech. Rep., 2004.
- [12] W. W. Hager, "Updating the inverse of a matrix," *Society for Industrial and Applied Mathematics (SIAM) Review*, vol. 31, no. 2, pp. 221–239, 1989.
- [13] L. Fei-Fei, R. Fergus, and P. Perona, "Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories," in *CVPR '04: Workshop on Generative-Model Based Vision*, 2005.
- [14] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *CVPR*, 2006, pp. 2169–2178.
- [15] A. Quattoni and A. Torralba, "Recognizing indoor scenes," in *CVPR*, 2009, pp. 413–420.
- [16] A. Berg, J. Deng, and L. Fei-Fei, "Large scale visual recognition challenge," 2010, <http://www.image-net.org/challenges/LSVRC/2010/>.
- [17] B. Fröhlich, E. Rodner, M. Kemmler, and J. Denzler, "Efficient gaussian process classification using random decision forests," *Pattern Recognition and Image Analysis*, vol. 21, pp. 184–187, 2011.



Alexander Lütz Alexander Lütz, born November 20, 1988, received the Diploma degree in Computer Science with honours in 2011 from the University of Jena, Germany. Currently, he is a PhD student under supervision of Joachim Denzler at the Chair for Computer Vision, University of Jena. His research interests cover the fields of machine learning and incremental learning, computer vision, object recognition and novelty detection.



Erik Rodner Erik Rodner, born May 22, 1983 earned the Diploma degree in Computer Science with honours in 2007 from the University of Jena, Germany. He pursued and received his Ph.D. in 2011 with summa cum laude for his work on learning with few examples, which was done under supervision of Joachim Denzler at the computer vision research group of the University of Jena. Erik

is currently continuing his research as a post-doctoral researcher. His research interests include kernel methods, visual object discovery, rare animals, and scene understanding.



Joachim Denzler Joachim Denzler

zler, born April 16, 1967, earned the degrees "Diplom-Informatiker", "Dr.-Ing.", and "Habilitation" from the University of Erlangen in the years 1992, 1997, and 2003, respectively. Currently, he holds a position of full professor for computer science and is head of the Chair for Computer Vision, Faculty of Mathematics and Informatics, Friedrich-Schiller-University of Jena. His research interests comprise active computer vision, object recognition and tracking, 3D reconstruction, and plenoptic modeling, as well as computer vision for autonomous systems. He is author and coauthor of over 200 journal papers and technical articles. He is a member of the IEEE, IEEE computer society, DAGM, and GI. For his work on object tracking, plenoptic modeling, and active object recognition and state estimation, he was awarded the DAGM best paper awards in 1996, 1999, and 2001, respectively.