

Learning with Few Examples for Binary and Multiclass Classification Using Regularization of Randomized Trees

Erik Rodner*, Joachim Denzler

Chair for Computer Vision, Friedrich Schiller University of Jena, Germany

Abstract

The human visual system is often able to learn to recognize difficult object categories from only a single view, whereas automatic object recognition with few training examples is still a challenging task. This is mainly due to the human ability to transfer knowledge from related classes. Therefore, an extension to Randomized Decision Trees is introduced for learning with very few examples by exploiting interclass relationships. The approach consists of a maximum a posteriori estimation of classifier parameters using a prior distribution learned from similar object categories. Experiments on binary and multiclass classification tasks show significant performance gains.

Key words: object categorization, randomized trees, few examples, interclass transfer, transfer learning

1. Introduction

During the last few decades, research in machine learning and computer vision has led to many new object representations and improved algorithms for numerical classification. Despite the success of this development, there is still an unanswered question: how does one learn object models from few training examples?

On the one hand, this question is motivated from industrial demand. In many applications, gathering hundreds or thousands of training images is either expensive or nearly impossible (Platzer et al., 2008). Building robust classification systems in those settings therefore requires complex specialized methods, that indirectly incorporate human prior knowledge about the task.

On the other hand, progress on learning with few examples is an important challenge and an essential step towards closing the gap between human and computer vision abilities. The human visual recognition system is often easily able to learn a new object category, such as a new animal class, from just a single view.

At first glance, this observation seems to contradict to the classical theory. The parameters of object models often exceed the available number of training examples

in multiple dimensions. From a mathematical point of view, this results in an ill-posed optimization problem, especially in cases with only a few training examples.

Therefore the only possibility to solve this problem is to regularize the optimization using prior knowledge. In previous algorithms this prior knowledge was often derived from abstract assumptions or was manually tuned during the development. However psychological studies (Jones and Smith, 1993) suggest that a key component of the human ability to recognize a class from a limited number of examples is the concept of *interclass transfer*. This paradigm is also known as *knowledge transfer*, *learning to learn* or *transfer learning*. It states that prior knowledge from previously learned object categories is the most important additional information source when learning object models from weak representations (Fei-Fei, 2006). To give an illustrative example of this idea, consider the recognition of a new animal class such as an okapi. With the aid of our prior knowledge from related animal classes (giraffe, zebra, antelope, etc.), we are able to generalize quickly from a single view.

In this work, a concept is presented how prior knowledge of related classes (often also called support classes) can be used to increase the generalization ability of a discriminative classifier. The underlying idea is a maximum a posteriori (MAP) estimation of parameters using a prior distribution estimated from similar

*Corresponding author

Email addresses: rodner@informatik.uni-jena.de (Erik Rodner), denzler@informatik.uni-jena.de (Joachim Denzler)

Preprint submitted to Pattern Recognition Letters

51 object categories. Furthermore, the application of this
52 idea to Randomized Decision Trees, as introduced by
53 Geurts et al. (2006), is demonstrated. The paper is based
54 on our previous work in Rodner and Denzler (2008) that
55 concentrates on multiclass classification. Studies are ex-
56 tended by showing the applicability of the approach to
57 binary classification. An additional experiment also em-
58 phasizes that the information transferred is not generic
59 prior knowledge unrelated to interclass relationships.

60 The remainder of the paper is organized as follows.
61 After previous work in the field of learning with weak
62 representations is briefly reviewed, it is shown that
63 Bayesian estimation using a prior distribution is a well
64 founded possibility to transfer knowledge from related
65 classes (Bayesian Interclass Transfer). This is followed
66 by a detailed description of an extension to Randomized
67 Decision Trees in Section 4, which can be regarded as
68 an application of Bayesian Interclass Transfer. Experi-
69 ments in binary and multiclass classification settings us-
70 ing publicly available image databases demonstrate the
71 benefits of the proposed algorithm in Sections 6 to 9. A
72 summary of our findings and a discussion about further
73 research steps conclude the paper.

74 2. Related Work

75 Previous work on interclass transfer varies signifi-
76 cantly in the type of information transferred from re-
77 lated classes. An intuitive assumption is that similar
78 classes share common intraclass geometric transforma-
79 tions. The *Congeaing* approach of Miller et al. (2000)
80 therefore tries to estimate those transformations and use
81 them to increase the amount of training data of a new
82 class. For example, a single training image of a letter
83 in a text recognition setting can be transformed using
84 typical rotations estimated from other letters.

85 Another idea is to assume shared structures in fea-
86 ture space and estimate a metric or transformation from
87 support classes (Fink, 2004; Quattoni et al., 2007). Tor-
88 ralba et al. (2007) used a discriminative boosting tech-
89 nique that exploits shared class boundaries within fea-
90 ture space. In contrast, Fei-Fei et al. (2006) devel-
91 oped a generative framework with MAP estimation of
92 model parameters using a prior distribution estimated
93 from support classes. A similar idea in the context of
94 shape based image categorization is presented in Stark
95 et al. (2009). In general the concept of shared priors for
96 a set of related classification problems can be used to
97 extend several classification techniques to multi-task ap-
98 proaches, such as generalized linear models (Lee et al.,
99 2007) or Gaussian processes (Bonilla et al., 2008).

100 Our work on regularized decision trees using transfer
101 learning is related to the approach of Lee and Giraud-
102 Carrier (2007). The key idea of their method is the
103 reusability of a decision tree structure from a related
104 binary classification task. In contrast, this paper intro-
105 duces a technique that also reuses estimated class proba-
106 bilities in leaf nodes and performs a re-estimation based
107 on a Bayesian framework.

108 3. Bayesian Interclass Transfer

109 The interclass transfer paradigm leads quickly to two
110 important questions: What type of information can be
111 transferred, and how can this be done using machine
112 learning techniques? The first question is answered in
113 Section 4.2. Here we concentrate on the description of
114 how prior knowledge can be incorporated.

Let a set \mathcal{S} of support classes and a class γ with few
training examples be given. In the remainder of this pa-
per, class γ is called the new class. The overall goal
of Bayesian Interclass Transfer is to estimate an ob-
ject model $\theta(\gamma)$ (parameters of a distribution, param-
eters of a classifier, etc.) with the help of prior knowl-
edge from related object models $\theta(i)$ where $i \in \mathcal{S}$. Us-
ing the Bayesian principle, this can be formulated as the
following maximum a posteriori estimation

$$\theta^{\text{MAP}}(\gamma) = \arg \max_{\theta} p(T^{\gamma} | \theta) p(\theta | T^{\mathcal{S}}), \quad (1)$$

115 where T^{γ} denotes the training data of the new class and
116 $T^{\mathcal{S}}$ denotes the training data of all support classes. The
117 fundamental assumption is that it is possible to estimate
118 a suitable prior distribution and use it to regularize the
119 parameter estimation of a related class.

120 The application of the principle of Bayesian Inter-
121 class Transfer (or Generative Transfer Learning) was
122 limited to generative approaches (Fei-Fei et al., 2006).
123 As we show in this paper it is also possible to en-
124 hance a discriminative classifier. The key idea is the
125 re-estimation of parameters of a discriminative classi-
126 fier by MAP estimation.

127 For this reason we propose to estimate the param-
128 eters $\theta(i)$ ($i \neq \gamma$) using a state-of-the-art discriminative
129 approach and only recompute the parameters of the new
130 class $\theta(\gamma)$ with further regularization. Figure 1 gives an
131 overview of this concept.

132 4. Regularized Randomized Trees

133 This section describes how to apply the previous idea
134 of Bayesian Interclass Transfer to decision tree classi-
135 fiers. Although the approach can be easily applied to

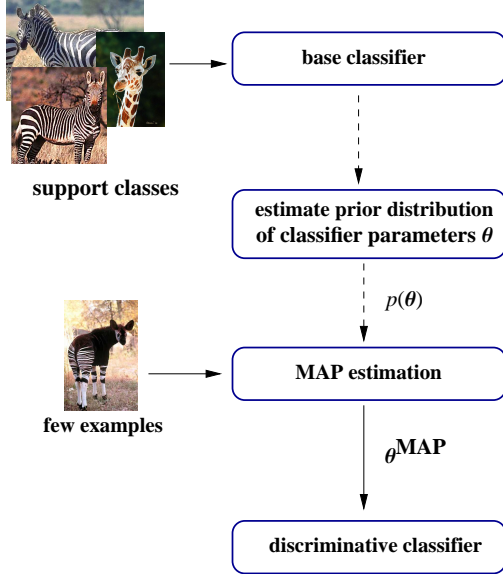


Figure 1: Overview of our approach using Bayesian Interclass Transfer for parameter estimation within a discriminative classification approach.

arbitrary decision tree approaches, the Randomized Decision Forest (RDF) approach is used, because of its superior generalization performance and its widely use in different applications (Marée et al., 2005; Shotton et al., 2008). In this section, we review RDF before providing a step-wise description of our method.

4.1. Randomized Decision Trees

Decision tree classifiers are commonly binary trees with two types of nodes. Each inner node represents a weak classifier (one-dimensional feature and threshold) which defines a hyperplane in feature space and thus determines the traversal of a new example within the tree. The traversal of the tree ends in a leaf node n .

We use n or n^l with $l = 1 \dots m$ to denote the event of an example reaching a single leaf node of a decision tree. This event also corresponds to the infinite set of all such examples (feature vectors). The total number of all leaf nodes in a single decision tree is denoted by m . Each leaf node is associated with a posterior distribution $p(\Omega_i | n)$, which is an estimate of the probability of class i given that this specific leaf is reached. We denote by Ω_i the event of an example belonging to the class i . These general principles and terms are illustrated in Figure 2.

Standard decision tree approaches suffer from two serious problems: long training time and over-fitting. The RDF approach solves both issues by random sampling.

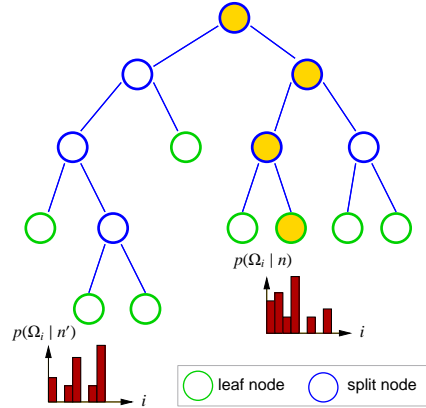


Figure 2: General principle and terms of decision trees. Diagrams illustrate the posterior distribution within each leaf node. Traversal of the tree (nodes filled with grey/yellow color) is done using features stored within each split node.

Instead of evaluating every feature and threshold, the training time is reduced by an approximate search for the most informative weak classifier in each node. The selection is made by choosing the weak classifier with the highest gain in information from a random fraction of features and thresholds.

Given enough training data for each class i , the generalization performance can be improved by learning an ensemble of M decision trees (often called a forest) using a random subset of the training data. From the final leaf nodes of the forest $\mathbf{n} = (n_1, \dots, n_M)$, the overall posterior can be obtained by voting with equal weights:

$$p(\Omega_i | \mathbf{n}) = \frac{1}{M} \sum_{s=1}^M p(\Omega_i | n_s) . \quad (2)$$

This special case of Bagging (Breiman, 2001) reduces the over-fitting effects without the need for additional tree pruning.

4.2. Transfer Learning Using RDF

The transfer learning idea can be applied to each tree of the forest individually; therefore, the details of our method are explained using only a single decision tree. Two different types of information are transferred: a discriminative tree structure and a prior distribution on leaf probabilities.

4.2.1. Recycling of Decision Trees

The selection of discriminative features in high dimensional spaces using few examples is a highly ill-posed problem. Therefore, we construct a discriminative tree structure using all the available training data of

183 all classes. This concept has also been used in Hoiem
 184 et al. (2007) and Lepetit et al. (2005), to recycle fea-
 185 tures and to reduce computation time. The assumption
 186 of shared discriminative features (or weak learners) is
 187 closely related to the use of shared features in the work
 188 of Torralba et al. (2007).

189 4.2.2. Re-estimation of leaf probabilities

Although decision tree approaches can be considered
 as discriminative, they are closely related to individual
 density estimation. The tree structure is a partitioning of
 the whole feature space into several cells n^l represented
 by leaf nodes. This corresponds to an approximation of
 a class distribution using a piecewise constant density
 or discrete probability distribution. The leaf probabilities
 $\theta_l^i = p(n^l | \Omega_i)$ are the maximum likelihood (ML)
 estimates of a multinomial distribution estimating the
 density of each cell:

$$\theta_l^{\text{ML}(i)} = \frac{|n^l \cap T^i|}{|T^i|} . \quad (3)$$

190 Note that $|n^l \cap T^i|$ is the number of examples of class
 191 i reaching a node n_l during the training step. It should
 192 be noted that with a careful implementation of decision
 193 trees, which store those unnormalized values instead of
 194 the posterior probability, a complicated recursive com-
 195 putation of leaf probabilities as presented in Rodner and
 196 Denzler (2008) is not necessary.

197 It is obvious that with only a few training examples
 $\mathbf{x} \in T^\gamma$, the vector $\theta^{\text{ML}}(\gamma)$ is sparse and is unable to pro-
 198 vide a good approximation of the underlying distribu-
 199 tion. The overall goal of our approach is to re-estimate
 200 $\theta(\gamma)$ by using maximum a posteriori estimation, which
 201 leads to a smoother solution $\theta^{\text{MAP}}(\gamma)$. Since the leaves
 202 of a decision tree induce a partitioning in disjoint sub-
 203 sets n^l , each instance of the parameter vector θ is a dis-
 204 crete multinomial distribution. For this reason any suit-
 205 able distribution of discrete distributions can be used to
 206 model the prior distribution.
 207

208 4.3. Constrained Gaussian Prior

We propose to use a constrained Gaussian distribu-
 tion (CGD), which is a simple family of parametric dis-
 tributions and can serve as an alternative to a standard
 Dirichlet distribution. For all $l : \theta_l \geq 0$ the density is
 defined as

$$p(\theta|T^S) \propto \mathcal{N}(\theta|\mu^S, \sigma^2\mathbf{I}) \delta\left(1 - \sum_l \theta_l\right) . \quad (4)$$

209 The factor of δ ($\delta(0) = 1, \forall x \neq 0 : \delta(x) = 0$) is essential
 210 to ensure that the support of the density function is the

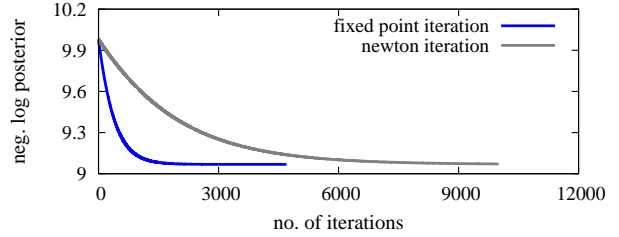


Figure 3: Comparison between the convergence of the newton method and a simple fixed point iteration.

211 simplex of all feasible discrete distributions. The use of
 $\sigma^2\mathbf{I}$ as a covariance matrix is an additional assumption
 212 that will be useful in deriving an efficient MAP estima-
 213 tion algorithm (Section 4.4).
 214

This simple model allows us to estimate hyper-
 parameters μ^S and σ in an usual way. Because the sim-
 plex is a convex set, the mean vector μ^S can be esti-
 mated analogously to a non-constrained Gaussian. In
 our application on decision trees, μ^S is estimated using
 the leaf probabilities of the support classes:

$$\mu^S = \frac{1}{|S|} \sum_{i \in S} \theta(i) . \quad (5)$$

215 Our choice to model the unknown distribution by a
 216 Gaussian parametric family is mostly due to practical
 computational considerations rather than theoretical re-
 sults. Of course, one could argue, that using a sym-
 metric Dirichlet prior leads to the same set of param-
 eters as a CGD and is additionally a conjugate prior. In
 our application for Regularized trees, we expect a sym-
 metric Dirichlet prior to yield similar results. But in
 our opinion the use of a constrained Gaussian prior is
 scientifically interesting and we will show in the fol-
 lowing that even without a conjugate prior, one can de-
 rive a simple inference method using an easy to solve
 one-dimensional optimization problem. An investiga-
 tion and analysis of other parametric distributions and
 more sophisticated priors would be an interesting topic
 for future research.
 228
 229
 230

231 4.4. MAP Estimation using a CGP

232 The process of MAP estimation using complex para-
 233 metric distribution often requires nonlinear optimiza-
 234 tion techniques. In contrast to these approaches we
 235 briefly show that by using our constrained Gaussian as
 236 a prior of a multinomial distribution, it is possible to
 237 derive a closed-form solution of the global optimum de-
 238 pending on a single Lagrange multiplier.

We start by writing the objective function of the MAP estimation as a Lagrange function of our simplex constraint and the posterior:

$$L(\theta, \lambda) = \log(p(T^\gamma|\theta) p(\theta|T^S)) + \lambda \left(\sum_l \theta_l - 1 \right). \quad (6)$$

The likelihood has a simple multinomial form and depends on a discrete histogram $\mathbf{c} = (c_l)_{l=1}^m$ representing the number of samples of each component:

$$p(T^\gamma|\theta) \propto \prod_l (\theta_l)^{c_l}. \quad (7)$$

In our application to leaf probabilities of decision trees, the absolute number of examples reaching a node $c_l = |n^l \cap T^\gamma|$ is used, where m is the number of all leaves. With the CGD prior in equation 4 we obtain the overall objective function

$$\sum_l \left(c_l \log(\theta_l) - \frac{1}{2\sigma^2} (\theta_l - \mu_l)^2 + \lambda \theta_l \right) - \lambda.$$

This objective function is convex and therefore has a unique solution. Setting the gradient $\left(\frac{\partial L}{\partial \theta_l}\right)(\theta, \lambda)$ to zero leads to the m independent equations

$$0 = \frac{c_l}{\theta_l} - \frac{1}{2\sigma^2} \cdot 2 \cdot (\theta_l - \mu_l) + \lambda. \quad (8)$$

Note that we get a non-informative prior, which reduces MAP to ML estimation as $\sigma^2 \rightarrow \infty$. With positive discrete probabilities ($\theta_l > 0$), it is possible to obtain a simple quadratic equation in θ_l :

$$0 = \theta_l^2 + \theta_l (-\mu_l - \lambda\sigma^2) - \sigma^2 c_l. \quad (9)$$

A stationary point with $\theta_l = 0$ is only possible with $c_l = 0$ or $\sigma^2 \rightarrow 0$, which is also reflected by the above equation. Therefore the optimization problem has only a single non-negative solution depending on λ :

$$\theta_l = \frac{\mu_l + \lambda\sigma^2}{2} + \sqrt{\left(\frac{\mu_l + \lambda\sigma^2}{2}\right)^2 + \sigma^2 c_l}. \quad (10)$$

This solution depends on the Lagrange multiplier, for which an optimal value can be found using a simple fixed point iteration:

$$\lambda^{j+1} = \frac{1}{m\sigma^2} \left(1 - 2 \sum_l \sqrt{\left(\frac{\mu_l + \lambda^j \sigma^2}{2}\right)^2 + \sigma^2 c_l} \right). \quad (11)$$

As an initial value, it is possible to use the optimal Lagrange multiplier in the case of no prior knowledge and maximum likelihood estimation. Figure 3 shows the convergence of our technique compared to that of a Newton iteration, which converges much slower than our simple recursion formula of Equation (11).

5. Binary and Multiclass Transfer Learning

Transfer learning for binary classification relies on a set of support tasks that try to separate a class i and a background class \mathcal{B} . Regularized trees can be applied straight-forwardly to this setting if a single support classification task is given. After building a random forest using training data from $\mathcal{S} = \{i\}$ and \mathcal{B} , we can apply the re-estimation method as explained in Section 4.4 using the mean vector $\mu^S = \theta(i)$. Finally the class probabilities of γ are substituted for all probabilities of i , so the decision tree now tries to separate between γ and \mathcal{B} .

In contrast to previous work, which often concentrate on the binary case (Fei-Fei et al., 2006), Regularized Trees are even suitable for multiclass classification problems. Given the leaf probabilities θ_l^i for each class i and leaf l and prior probabilities $p(\Omega_i)$ for each class, one can easily calculate the needed posterior probabilities for each class in the multiclass problem:

$$p(\Omega_i | n_l) = \frac{p(n_l | \Omega_i) p(\Omega_i)}{\sum_j p(n_l | \Omega_j) p(\Omega_j)}. \quad (12)$$

Reducing Confusion with Support Classes. All machine learning approaches using the interclass paradigm within a multi-classification task have to cope with a common issue: transferring knowledge from support classes can lead to confusion with the new class. For example, using prior information from camel images to support the class dromedary enables us to transfer shared features like fur color or head appearance. However, the we have to use additional features (e.g. shape information) to discriminate between both categories.

To solve this problem, we propose to build additional discriminative levels of the decision tree after MAP estimation of the leaf distributions. Starting from a leaf node n^l with non-zero posterior probability $p(\Omega_\gamma | n^l)$, the tree is further extended by the randomized training procedure described in Section 4.1. The training data in this case consists of all samples of the new class and samples of all support classes which reached the leaf n^l . All of the training examples are weighted by the values of the posterior distribution $p(\Omega_i | n^l)$ of the leaf n^l . This technique allows us to find new discriminative

277 features especially between the new class and the sup-
278 port classes. We observed that often only one additional
279 level can be build using the few examples of γ .

280 6. Experimental Setup and Overview

281 The approach presented is evaluated experimentally
282 to analyze the benefits and the limitations of all our as-
283 sumptions. Three experiments are performed to provide
284 empirical proof of the following statements:

- 285 1. Regularized Trees lead to a significant perfor-
286 mance gain for multiclass classification with few
287 training examples (Exp. 1, Sect. 7).
- 288 2. The performance of binary classification can be
289 improved by our method (Exp. 2, Sect. 8).
- 290 3. Our method uses prior knowledge that relies on vi-
291 sual similarity, and is thus not related to generic
292 prior knowledge (Exp. 3, Sect. 9).

293 For the comparative analysis, three types of public
294 datasets with different characteristics are used: a dataset
295 of handwritten Latin letters provided by Fink (2004), a
296 combination of the bird and butterfly datasets used in
297 Lazebnik et al. (2004, 2006) and a dataset for binary
298 classification using images from the database of mam-
299 mals presented in Fink and Ullman (2008).

300 The evaluation criteria are the unbiased average
301 recognition rates of the whole classification task and
302 single recognition rates of the new class. Monte Carlo
303 analysis is performed by randomly selecting f examples
304 of the new class for training and the remainder for test-
305 ing. To estimate the recognition rates for a fixed value
306 of f the results of multiple runs are averaged. This also
307 averages out the influence of our randomized classifier.

308 The experimental evaluation aims to analyze the gain
309 of our transfer learning approach compared to the RDF
310 classifier Geurts et al. (2006). We do not focus on the
311 development of new feature types that would be suitable
312 for special recognition tasks. For this reason, our choice
313 of features is not optimized. The variance σ^2 of the CGP
314 is an important parameter of our method, which we fix
315 to the value of 10^{-5} in all experiments. It controls the in-
316 fluence of the prior distribution and therefore, indirectly,
317 our assumption of how much the new class is related to
318 support classes. We decided to use a constant value for
319 this parameter, because cross-validation is impossible
320 with a single training example.

321 Furthermore we select support classes manually in all
322 the experiments. Our main assumption in Equation (1)
323 is that those categories have to share common features,
324 shape or appearance. Estimating the class similarities

325 automatically would be optimal to provide support class
326 subsets. Regarding the selection of support classes as a
327 model selection problem allows to use cross-validation
328 or leave-one-out estimates (cf. Tommasi and Caputo
329 (2009)). However, this can be rather difficult and results
330 in ill-posed problems themselves. Hence, we leave the
331 estimation of a set of similar classes as a task for future
332 research.

333 7. Experiment 1: Multiclass Classification

334 This experiment shows the benefits of our method in a
335 high-level image categorization task and a simpler letter
336 recognition task. We explain all features used and give
337 a detailed discussion of all results in section 7.3.

338 7.1. Letter Recognition

339 The database of Fink (2004) is a collection of im-
340 ages containing handwritten Latin letters resulting in 26
341 object categories. For each object class 60 images are
342 provided. For classification an ensemble of 10 decision
343 trees is used and the following classification scenario is
344 selected: new class e and support classes a, b, c, d .

345 *Features.* The images in this database are binary, so
346 a very simple feature extraction method is used. The
347 whole image is divided into an equally spaced $w_x \times w_y$
348 grid. In each cell of the grid, the ratio of black pixels to
349 all pixels within the cell is used as a single feature. This
350 leads to a feature vector with $w_x w_y$ dimensions. In all
351 experiments, the values $w_x = 8$ and $w_y = 12$ are used.

352 7.2. Image Categorization

353 To demonstrate the behavior of the method on a
354 high-level image categorization task, we combine the
355 birds (Lazebnik et al., 2006) and the butterflies dataset
356 (Lazebnik et al., 2004) into one single multiclass clas-
357 sification task. The object categories can therefore be
358 divided into two different semantic sets. The category
359 *black swallowtail* is used as a new class γ , and all the
360 other butterfly categories serve as support classes \mathcal{S} .
361 Thus, training data consists of a variable number of
362 training images for γ and 26 images for each of the re-
363 maining classes. This classification task is more diffi-
364 cult than our letter recognition setting. For this reason
365 an ensemble of 500 decision trees was used.



Figure 4: Example images of all datasets used for experimental evaluation: Top row: combined bird and butterfly dataset of Lazebnik et al. (2006, 2004). Middle row: latin letter dataset of Fink (2004). Bottom row: zebra and okapi images used for binary classification obtained from the mammals dataset of Fink and Ullman (2008) and google image search.

366 *Features.* A standard approach to image categorization 400
 367 is the bag-of-features idea. A quantization of local fea- 401
 368 tures is computed, which is often called a codebook, is 402
 369 computed at the time of training. An image can then 403
 370 be represented as a histogram of local features with re- 404
 371 spect to codebook entries. The method of Moosmann 405
 372 et al. (2006), which utilizes a random forest as a cluster- 406
 373 ing mechanism, is used to construct the codebook. This 407
 374 codebook generation procedure shows superior results 408
 375 compared to standard k -Means within all experiments. 409
 376 It also allows us to create large codebooks (a size of 410
 377 13000 used in all experiments) in a few minutes on a 411
 378 standard PC. A combined SIFT descriptor computed on 412
 379 normalized RGB channels, as described in van de Sande 413
 380 et al. (2010), is used as a local feature representation. 414

381 7.3. Evaluation

382 The results of this experiment evaluating multiclass 417
 383 classification performance can be found in Figures 5 and 418
 384 6. The plots show the average recognition rate of the 419
 385 whole task (plots on the left side) and the recognition 420
 386 rate of the new class (plots on the right side) compared 421
 387 to those of the original method of RDF.

388 It can be seen that our method improves the recog- 422
 389 nition rate of the new class and the average recognition 423
 390 rate, in the range with few training examples (1 to 8 424
 391 examples, marked with green color). The regulariza- 425
 392 tion is therefore able to transfer knowledge from sup- 426
 393 port classes without violating the separation between 427
 394 the other classes.

395 After a specific number of training examples the aver- 428
 396 age recognition rate decreases while the recognition rate 429
 397 or hit-rate of class γ (plots on the right side) still grows. 430
 398 This critical area is highlighted in yellow in Figures 5 431
 399 and 6. The effect corresponds to over-regularization. 432

The influence of the prior distribution is controlled by 400
 hyper-parameter σ^2 which is kept a fixed value inde- 401
 pendent of the training examples used. Therefore the 402
 MAP estimation of leaf probabilities leads to many leafs 403
 with non-zero posterior probabilities for the new class. 404
 This corresponds to a large variance of the distribution 405
 in feature space, which dominates the distribution for 406
 all other classes. The variance of the class distribution 407
 reaches a critical threshold, which leads to an overes- 408
 timation of the distribution corresponding to the new 409
 class. The classifier prefers the new class, which results 410
 in a worse average recognition rate (or an increasing 411
 number of false positives) on the whole classification 412
 task. It should be noted that this phenomenon is unique 413
 to our application of transfer learning in a multi-class 414
 classification task. Other transfer learning algorithms 415
 converge to the performance of independent learning af- 416
 ter a specific number of training examples, due to their 417
 treatment of a support and new class as independent bi- 418
 nary classification tasks. A similar effect has been ob- 419
 served in the context of zero-shot learning Rohrbach 420
 et al. (2010) (cf. their Fig. 3). 421

422 8. Experiment 2: Binary Classification

423 For an experimental evaluation of the method on a 424
 binary classification task, images from the animal cate- 425
 gories *zebra* and *okapi* from the mammals database of 426
 Fink and Ullman (2008) are used. In order to increase 427
 the number of test images, additional images from the 428
 category *okapi* were downloaded using Google Image 429
 Search and filtered manually to delete wrong search re- 430
 sults. The new dataset includes a total of 231 images 431
 of okapis and 200 images of zebras. The image set of 432
 the background class \mathcal{B} was generated by obtaining 300

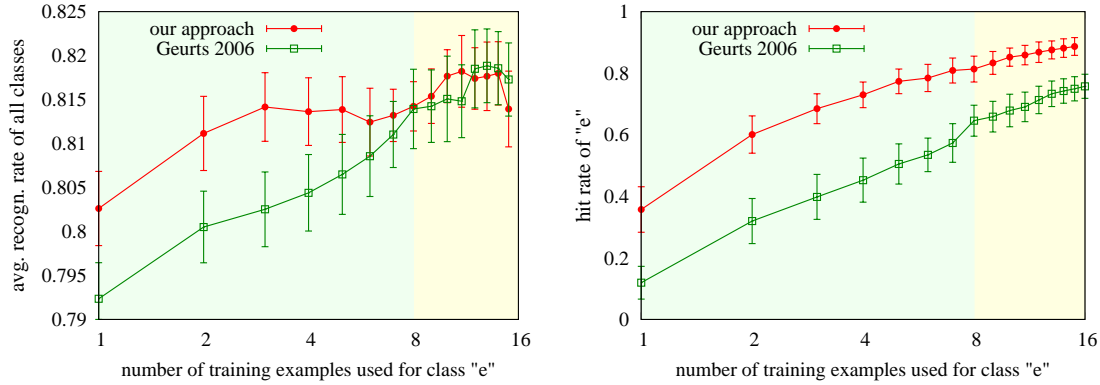


Figure 5: Comparison to Geurts et al. (2006) in a multiclass classification task using the letter recognition dataset of Fink (2004). The left plot shows the average recognition rate of the whole classification task with respect to the number of training examples (log scale) of a specific class. On the right side the single recognition rate of this class is plotted. Highlighted green area corresponds to the working range of our algorithm before over-regularization effects. False alarm rates are skipped because we concentrate on the categorization performance.

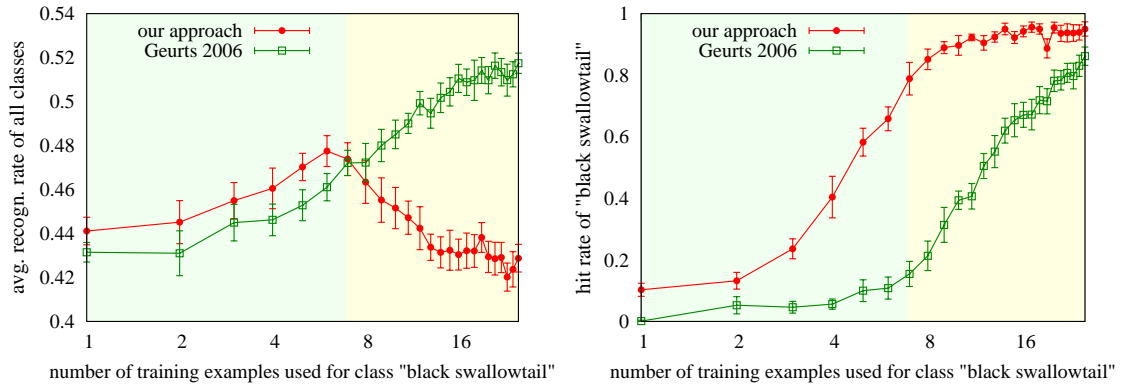


Figure 6: Comparison to Geurts et al. (2006) within a high level multiclass classification task using the bird-and-butterfly dataset as used in Lazebnik et al. (2004, 2006). Semantic of the plot is analogous to Figure 5.

433 random images from Google Image Search (using *the* 434
 434 as a search key word). Our algorithm was tested with 451
 435 two scenarios: using few training examples of the class 452
 436 okapi with the support of the class zebra and vice versa. 453
 437 Feature extraction was done as described in Section 7.2.

438 8.1. Evaluation

439 Figure 7 shows the results of our approach (red plot, 455
 440 circular dots) compared to the standard approach of 456
 441 Randomized Decision Forest (green plot, rectangular 457
 442 dots). We also tested the performance of a random forest 458
 443 built by using the supporting classification task without 459
 444 our re-estimation technique (blue plot, triangular dots). 460

445 First of all, it is apparent that our method significantly 461
 446 increases the classification performance compared to 462
 447 the standard approach in both cases. Using a random 463
 448 forest without re-estimation of leaf probabilities does 464
 449 not use training examples of the new class and is there- 465

450 fore independent of the number of training examples. 451
 452 Additionally one can see that the “okapi” task seems to 453
 453 be much harder, and benefits of knowledge transfer for 454
 454 a wider range of training examples.

454 9. Experiment 3: Similarity Assumption

455 What happens if support classes are selected that do 456
 456 not share common features with the new class?

457 As mentioned in Section 3 the concept of Bayesian 458
 458 Interclass Transfer is based on the main assumption that 459
 459 the support classes \mathcal{S} are somehow similar to the new 460
 460 class γ . Therefore, it is possible to further assume that 461
 461 those similarities can be captured in feature space by a 462
 462 distribution $p(\theta)$. The following experiment tries to un- 463
 463 cover whether the knowledge transferred is related to a 464
 464 generic prior or is more category-specific and thus trans- 465
 465 fers more detailed elements, such as object parts.

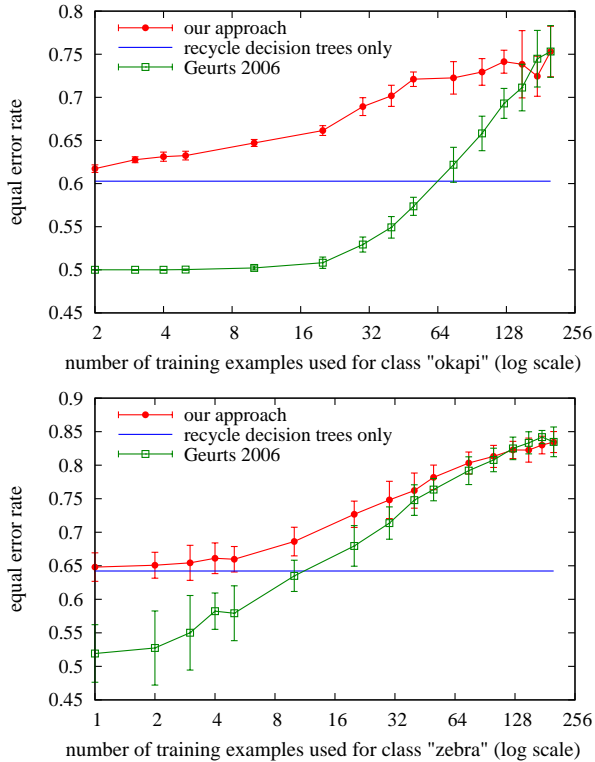


Figure 7: Results of the comparison of our method with the RDF classifier of Geurts et al. (2006) using binary classification tasks.

466 To answer this question, an experiment using the letter
 467 recognition scenario (Section 7.1) is performed. As
 468 a new class with a weak representation of 4 training ex-
 469 amples we selected the letter e and used two different
 470 sets of similar support classes (a, b, c, d) and dissimilar
 471 support classes (m, n, w, v, z). Figure 8 shows a scatter plot
 472 of several runs, where each point corresponds to the av-
 473 erage recognition rate of a Randomized Decision Forest
 474 without (ML estimation) and with our transfer learning
 475 method (MAP estimation). All points above the diag-
 476 onal therefore indicate a clear benefit from prior knowl-
 477 edge. It can be seen that visually dissimilar classes (tri-
 478 angular dots in red color) do not lead to a performance
 479 gain and can even decrease the performance.

480 9.1. Discussion of Experiment 3

481 Our results clearly show that our transfer learning
 482 method learns prior knowledge that is not related to
 483 generic prior knowledge. This is an important dif-
 484 ference to a lot other approaches which capture more
 485 generic prior knowledge. For example in Fei-Fei et al.
 486 (2006), Bayesian Interclass Transfer is applied to trans-
 487 fer knowledge between object categories such as: mo-
 488 torbikes, faces, airplanes and wild cats. Therefore, their

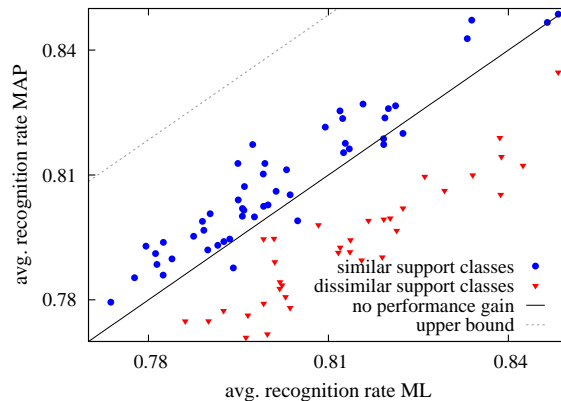


Figure 8: Average recognition rate of the ML approach in comparison to the rate after applying our MAP re-estimation technique. The regularization results in a performance gain only if support classes are (visually) similar to the new class.

489 method seems to use a generic prior of object category
 490 images (e.g. size and location of objects are not uni-
 491 formly distributed). Bart and Ullman (2005) also tested
 492 their approach with a large set of various unrelated cate-
 493 gories of the Caltech-101 database and showed that the
 494 knowledge transferred by their approach, represented by
 495 shared image fragments, helped to improve the recog-
 496 nition performance. In general the use of generic prior
 497 knowledge has its own tradition and motivation, espe-
 498 cially in the context of natural image statistics (Torralba
 499 and Oliva, 2003). In our opinion the use of category-
 500 specific in addition to generic priors is essential to cap-
 501 ture available knowledge as much as possible and thus
 502 allows efficient learning with few examples similar to
 503 the development of the human visual system.

504 10. Conclusion

505 We argue that learning with few examples can benefit
 506 from incorporating prior knowledge of related classes
 507 (interclass transfer paradigm). Therefore, we proposed
 508 to reuse (transfer) the discriminative structure of a Ran-
 509 domized Decision Forest and apply a subsequent maxi-
 510 mum a posteriori estimation of leaf probabilities in
 511 each tree. This Bayesian formulation allows us to infer
 512 knowledge as a prior distribution obtained from related
 513 classes and can be seen as a regularization technique.
 514 The method is able to exploit interclass relationships to
 515 support learning of a class with few training examples.

516 Experiments on several public datasets showed a sig-
 517 nificant performance gain in dealing with a weak train-
 518 ing representation. In contrast to other work (Fei-Fei
 519 et al., 2006), transfer learning of Randomized Decision

Trees is applicable for binary and even for multiclass classification, where information is transferred within the task. An additional experiment validated that the transferred prior information captures (visual) similarities of related classes unlike a generic prior.

11. Further Work

Regularization with a meaningful prior derived from similar object categories is an interesting research direction. Especially for learning with few training examples, transferring knowledge from similar object categories currently seems to be the only way to handle the underlying ill-posed problems.

Despite the benefits presented in this paper, the proposed method has two drawbacks: the support classes have to be selected manually and the influence of the prior has to be controlled by the variance σ^2 of the underlying distribution. The optimal parameter σ^2 could be found by a typical method for estimating the regularization parameter using the L-curve (Kilmer and O'Leary, 2001). An alternative would be to use cross validation, which is a common tool for all parameter estimation problems within a classification task. Automatically selecting the support classes is more complex. In our case it is yet unknown whether the information of few examples is sufficient to estimate the similarity to other categories that would be useful for regularization.

References

Bart, E., Ullman, S., 2005. Cross-generalization: Learning novel classes from a single example by feature replacement. In: Proceedings of the 2005 IEEE Conference on Computer Vision and Pattern Recognition (CVPR'05). pp. 672–679.

Bonilla, E., Chai, K. M., Williams, C., 2008. Multi-task gaussian process prediction. In: Advances in Neural Information Processing Systems 20. MIT Press, pp. 153–160.

Breiman, L., October 2001. Random forests. *Machine Learning* 45 (1), 5–32.

Fei-Fei, L., 2006. Knowledge transfer in learning to recognize visual objects classes. In: Proceedings of the International Conference on Development and Learning (ICDL).

Fei-Fei, L., Fergus, R., Perona, P., 2006. One-shot learning of object categories. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 28 (4), 594–611.

Fink, M., 2004. Object classification from a single example utilizing class relevance pseudo-metrics. In: Advances in Neural Information Processing Systems. Vol. 17. The MIT Press, pp. 449–456.

Fink, M., Ullman, S., 2008. From aardvark to zorro: A benchmark for mammal image classification. *Int. J. Comput. Vision* 77 (1-3), 143–156.

Geurts, P., Ernst, D., Wehenkel, L., 2006. Extremely randomized trees. *Maching Learning* 63 (1), 3–42.

Hoiem, D., Rother, C., Winn, J., 2007. 3d layoutcrf for multi-view object class recognition and segmentation. In: Proceedings of the

2007 IEEE Conference on Computer Vision and Pattern Recognition (CVPR'07). pp. 1–8.

Jones, S. S., Smith, L. B., April-June 1993. The place of perception in children's concepts. *Cognitive Development* 8, 113–139.

Kilmer, M., O'Leary, D., 2001. Choosing regularization parameters in iterative methods for ill-posed problems. *SIAM J. Matrix Anal. Appl* 22 (4), 1204–1221.

Lazebnik, S., Schmid, C., Ponce, J., 2004. Semi-local affine parts for object recognition. In: British Machine Vision Conference. Vol. 2. pp. 779–788.

Lazebnik, S., Schmid, C., Ponce, J., 2006. A discriminative framework for texture and object recognition using local image features. In: Ponce, J., Hebert, M., Schmid, C., Zisserman, A. (Eds.), *Toward Category-Level Object Recognition*. Vol. 4170 of Lecture Notes in Computer Science. Springer, pp. 423–442.

Lee, J. W., Giraud-Carrier, C., Aug. 2007. Transfer learning in decision trees. In: International Joint Conference on Neural Networks (IJCNN) 2007. pp. 726–731.

Lee, S.-I., Chatalbashev, V., Vickrey, D., Koller, D., 2007. Learning a meta-level prior for feature relevance from multiple related tasks. In: ICML '07: Proceedings of the 24th International Conference on Machine Learning. pp. 489–496.

Lepetit, V., Lagger, P., Fua, P., 2005. Randomized trees for real-time keypoint recognition. In: Proceedings of the 2005 IEEE Conference on Computer Vision and Pattern Recognition (CVPR'05). pp. 775–781.

Marée, R., Geurts, P., Piater, J., Wehenkel, L., June 2005. Random subwindows for robust image classification. In: Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR'05). Vol. 1. pp. 34–40.

Miller, E. G., Matsakis, N. E., Viola, P. A., 2000. Learning from one example through shared densities on transforms. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR'00). pp. 464–471.

Moosmann, F., Triggs, B., Jurie, F., 2006. Fast discriminative visual codebooks using randomized clustering forests. In: Advances in Neural Information Processing Systems. pp. 985–992.

Platzer, E.-S., Denzler, J., Süsse, H., Nägele, J., Wehking, K.-H., October 2008. Challenging anomaly detection in wire ropes using linear prediction combined with one-class classification. In: Proceedings of the Vision, Modelling, and Visualization Workshop. Konstanz, pp. 343–352.

Quattoni, A., Collins, M., Darrell, T., 2007. Learning visual representations using images with captions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR'07). pp. 1–8.

Rodner, E., Denzler, J., October 2008. Learning with few examples using a constrained gaussian prior on randomized trees. In: Proceedings of the Vision, Modelling, and Visualization Workshop. Konstanz, pp. 159–168.

Rohrbach, M., Stark, M., Szarvas, G., Schiele, B., Gurevych, I., 2010. What helps where - and why? semantic relatedness for knowledge transfer. In: CVPR'10: Proceedings of the Computer Vision and Pattern Recognition Conference.

Shotton, J., Johnson, M., Cipolla, R., 2008. Semantic texton forests for image categorization and segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR'08). pp. 1–8.

Stark, M., Goesele, M., Schiele, B., 2009. A shape-based object class model for knowledge transfer. In: Proceedings of the International Conference on Computer Vision (ICCV). pp. 373–380.

Tommasi, T., Caputo, B., 2009. The more you know, the less you learn: from knowledge transfer to one-shot learning of object categories. In: BMVC.

Torralba, A., Murphy, K. P., Freeman, W. T., 2007. Sharing visual fea-

637 tures for multiclass and multiview object detection. *IEEE Transactions*
638 *on Pattern Analysis and Machine Intelligence* 29 (5), 854–
639 869.

640 Torralba, A., Oliva, A., 2003. Statistics of natural image categories.
641 *Network: Computation in Neural Systems* 14 (1), 391–412.

642 van de Sande, K. E., Gevers, T., Snoek, C. G., 2010. Evaluating color
643 descriptors for object and scene recognition. *IEEE Transactions on*
644 *Pattern Analysis and Machine Intelligence* 32 (9), 1582–1596.