



massachusetts institute of technology — computer science and artificial intelligence laboratory

---

# Regularization Through Feature Knock Out

Lior Wolf and Ian Martin

AI Memo 2004-025  
CBCL Memo 242

November 2004

# Report Documentation Page

Form Approved  
OMB No. 0704-0188

Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

1. REPORT DATE <b>NOV 2004</b>		2. REPORT TYPE		3. DATES COVERED <b>00-11-2004 to 00-11-2004</b>	
4. TITLE AND SUBTITLE <b>Regularization Through Feature Knock Out</b>				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) <b>Massachusetts Institute of Technology, Artificial Intelligence Laboratory, 77 Massachusetts Avenue, Cambridge, MA, 02139</b>				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT <b>Approved for public release; distribution unlimited</b>					
13. SUPPLEMENTARY NOTES <b>The original document contains color images.</b>					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES <b>11</b>	19a. NAME OF RESPONSIBLE PERSON
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE <b>unclassified</b>			

### Abstract

In this paper, we present and analyze a novel regularization technique based on enhancing our dataset with corrupted copies of the original data. The motivation is that since the learning algorithm lacks information about which parts of the data are reliable, it has to produce more robust classification functions. We then demonstrate how this regularization leads to redundancy in the resulting classifiers, which is somewhat in contrast to the common interpretations of the Occam's razor principle. Using this framework, we propose a simple addition to the gentle boosting algorithm which enables it to work with only a few examples. We test this new algorithm on a variety of datasets and show convincing results.

Copyright ©Massachusetts Institute of Technology, 2004

This report describes research done at the Center for Biological & Computational Learning, which is in the McGovern Institute for Brain Research at MIT, as well as in the Dept. of Brain & Cognitive Sciences, and which is affiliated with the Computer Sciences & Artificial Intelligence Laboratory (CSAIL).

This research was sponsored by grants from: Office of Naval Research (DARPA) Contract No. MDA972-04-1-0037, Office of Naval Research (DARPA) Contract No. N00014-02-1-0915, National Science Foundation (ITR/IM) Contract No. IIS-0085836, National Science Foundation (ITR/SYS) Contract No. IIS-0112991, National Science Foundation (ITR) Contract No. IIS-0209289, National Science Foundation-NIH (CRCNS) Contract No. EIA-0218693, National Science Foundation-NIH (CRCNS) Contract No. EIA-0218506, and National Institutes of Health (Conte) Contract No. 1 P20 MH66239-01A1. Additional support was provided by: Central Research Institute of Electric Power Industry, Center for e-Business (MIT), Daimler-Chrysler AG, Compaq/Digital Equipment Corporation, Eastman Kodak Company, Honda R&D Co., Ltd., ITRI, Komatsu Ltd., Eugene McDermott Foundation, Merrill-Lynch, Mitsubishi Corporation, NEC Fund, Nippon Telegraph & Telephone, Oxygen, Siemens Corporate Research, Inc., Sony MOU, Sumitomo Metal Industries, Toyota Motor Corporation, and WatchVision Co., Ltd.

# 1. Introduction

Boosting - the iterative combination of classifiers to build a strong classifier - is a popular learning technique. The algorithms based on it, such as AdaBoost and GentleBoost, are easy to implement, work reasonably fast, and in general produce classifiers with good generalization properties for large enough datasets. If the dataset is not large enough and there are many features, these algorithms tend to overfit and perform much worse than the popular Support Vector Machine (SVM) algorithm. SVM can be successfully applied to all datasets, from small to very large.

The major drawback of SVM is that it uses, at run time, when classifying a new example  $x$ , all the measurements (features) of  $x$ . This poses a problem because, while we would like to cover many promising features during training, computing all the features at run-time might be too costly. This is especially true for object detection problems in vision, where we often need to search the whole image in several scales over thousands of possible locations, each location producing one such vector  $x$ . While several approaches for combining feature selection with SVM have been suggested in the past (e.g., [17]), they are rarely used.

The complexity of the feature vector can be controlled more easily by using boosting techniques over weak classifiers based on single features (e.g., regression stumps), such as in the highly successful system of [16]. In this case, the number of features used is bounded by the number of iterations in the boosting process. However, since boosting tends to overfit on small datasets, there is a bit of a dilemma here. An ideal algorithm would enable good control over the total number of features, while being able to learn from only a few examples. Such an algorithm is presented in Sec. 5.

This new algorithm is based on GentleBoost. Within it we implemented a regularization technique based on a simple idea: add corrupted copies of your training dataset to the original one, and the algorithm will not be able to overfit. A general background on fitting and regularization is given in the next section.

## 2 Background

We are given a set of  $n$  examples  $z_i = \{(x_i, y_i)\}_{i=1}^n$ ,  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$  drawn from a joint distribution  $\mathcal{P}$  on  $\mathcal{X} \times \mathcal{Y}$ . The ultimate goal of the learning algorithm is to produce a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  such that the *expected error* of  $f$  given by the expression  $E_{(x,y) \sim \mathcal{P}}(f(x) \neq y)$  is minimized. The boolean expression inside the parentheses evaluates to one if it holds, zero otherwise.

Since we do not know the distribution  $\mathcal{P}$  we are tempted to minimize the *empirical error* given by  $\sum_{i=1}^n (f(x_i) \neq y_i)$ . The problem is that if the space of functions from which the learning algorithm selects  $f$  is too large, we are at risk

of *overfitting* (learning to deal *only* with the training error). Therefore, while the empirical error is small, the expected error is large. In other words, the *generalization error* (the difference of empirical error from expected error) is large. Overfitting can be avoided by using any one of several *regularization* techniques.

Overfitting is usually the result of allowing too much freedom in the selection of the function  $f$ . Thus, the most basic regularization technique is to limit the number of free parameters we use while fitting the function  $f$ . For example, in binary classification we may limit ourselves to learning functions of the form  $f(x) = (h^\top x > 0)$  (we assume  $\mathcal{X} = \mathbb{R}^n$ .  $h$  is a vector of free parameters). Using such functions, we reduce the risk of overfitting, but may never optimally learn the *target function* (i.e., the “true function”  $f(x) = y$  that is behind the distribution  $\mathcal{P}$ ) of other forms, e.g., we will not be able to learn  $f(x) = (x(1)^2 - x(2) > 0)$ .

Another regularization technique is to minimize the empirical error subject to constraints on the learned functions. For example, we can require that the norm of the vector of free parameters  $h$  be less than one. A related but different regularization technique is to minimize the empirical error together with a penalty term on the complexity of the function we fit. The most popular penalty term – *Tikhonov regularization* – has a quadratic form. Using the linear model above, an appropriate penalty function would be  $\|h\|_2^2$ , and we would minimize  $\sum_{i=1}^n ((h^\top x_i > 0) \neq y_i) + \|h\|_2^2$ .

Sometimes, adding a regularization term to the optimization problem solved by the algorithm is not trivial. In the most extreme case, the algorithm is a black box we cannot alter at all. Still, a simple form of regularization called noise injection can be employed. In the noise injection technique, the training dataset is enriched by multiple copies of each training data point  $x_i$ . A zero-mean, low-variance Gaussian noise (independent for each coordinate) is added to each copy, and the original label  $y_i$  is preserved. The motivation is that if two data points  $x, x'$  are close (i.e.,  $\|x - x'\|$  is small), we would like  $f(x)$  and  $f(x')$  to have similar values. By introducing many examples with similar  $x$  values, and identical  $y$  values we teach the classifier to have this stability property. Hence, the learned function is encouraged to be smooth (at least around the training points).

The study of the noise injection technique, which blossomed in the mid 90’s, established the following results on noise injection: (1) It is an effective way to reduce generalization error. (2) It has a similar effect on shrinkage (the statistical term for regularization) of the parameters in some simple models (e.g., [3]). (3) It is equivalent to Tikhonov regularization [1]. Note that this does not mean that we can always use Tikhonov regularization instead of noise injection, as for some learning algorithms it is not possible to create a regularized version.

The technique we introduce next is similar in spirit to

---

**Input:**  $(x_1, y_1), \dots, (x_m, y_m)$  where  $x_i \in \mathcal{R}^n$ ,  $y_i \in Y$ .  
**Output:** one synthesized pair  $(\hat{x}, \hat{y})$ .

---

1. Select two examples  $x_a, x_b$  at random.
  2. Select a random feature  $k \in [1..n]$ .
  3. Set  $\hat{x} \leftarrow x_a$  and  $\hat{y} \leftarrow y_a$ .
  4. Replace feature  $k$  of  $\hat{x}$ :  $\hat{x}(k) \leftarrow x_b(k)$ .
- 

Figure 1: The Feature Knockout Procedure

noise injection. However, it is different enough that the results obtained for noise injection will not hold for it. For example, the results of [1] use a Taylor expansion around the original data points. Such an approximation will not hold for our new technique, since the “noise” is too large (i.e., the new datapoint is too different). Other important properties that might not hold are the independence of noise across coordinates, and the zero mean of the noise.

Our regularization technique is based on creating corrupted copies of the dataset. Each new data point is a copy of one original training point, picked at random, where one random coordinate (feature) is replaced with a different value—usually the value of the same coordinate in another random training example. The basic procedure used to generate the new example is illustrated in Fig. 1. We call it the Feature Knockout (KO) procedure, since one feature value is being altered dramatically. It is repeated many times to create new examples. It can be used with any learning algorithm, and we use it in the analysis presented in Sec. 4. However, as we focus our application emphasis on boosting, we use the specialized version in Fig. 2.

The KO regularization technique is especially suited for use when learning from only a few examples. The robustness we demand from the selected classification function is much more than local smoothness around the classification points (c.f. noise injection). This kind of smoothness is easy to achieve when example points are far from one another. Our regularization, however, is less restrictive than demanding uniform smoothness (Tikhonov) or requiring the reduction of as many parameters as possible. Both of these approaches might not be ideal when only a few examples are available because there is nothing to balance a large amount of uniform smoothness, and it is easy to fit a model that uses very few parameters. Instead, we encourage redundancy in the classifier since, in contrast to the shortage of training examples, there is an abundance of features.

### 3 Motivation

**Intuition.** It is a common belief that simpler is better. *Ockham’s razor* - “entities should not be multiplied beyond ne-

cessity” - is often understood as suggesting the selection of the simplest possible model that fits the data well. Thus, given a dataset, we are encouraged to prefer classifier  $A$  that uses 70 features over classifier  $B$  that uses 85 if they have the same expected error, because the “simpler” classifier is believed to generalize better.

In an apparent contrast to this belief, we know from our daily experience that simpler is not always better. When a good teacher explains something to a class, he will use a lot of repetition. The teacher ensures that the students will understand the idea, even if some of the explanations were unclear. Since the idea could be expressed in a simple form without repetition, the explanation is more complex.

It is also generally accepted that biological systems use redundancy in their computations. Thus, even if several computational units break down (e.g., when a neurons die) the result of the computation is largely unaffected. We expect the learned model to be interpreted with some random error. In these cases, we should not train a single model, but instead train to optimize a distribution of models.

**Redundancy in boosted classifiers.** Boosting over a weak classifier increases the weights of those examples it does not classify well. The inclusion of a weak classifier in the strong classifier therefore inhibits future use of similar weak classifiers. In boosting over regression stumps, each weak classifier is based on one feature. Hence, from a group of similar features, we may expect to see only one participating in the strong classifier. Our boosting over regression stumps algorithm, presented in Sec. 5, modulates this effect, by creating a new example for which relying on the selected feature might lead to a mistake. However, it does not change the values of the other features, making the similar features suitable for classifying the new example.

Such a process yields a larger classifier, which uses more features. This effect is clear in our experiments, and might also be interpreted as “a more complex model is needed to deal with more complex data.” However, using more features does not necessarily mean that the classifier will lead to a worse generalization error<sup>1</sup>.

Koltchinskii and Panchenko have derived bounds on the generalization error of ensemble (voting) classifiers, such as boosting, which take redundancy into consideration [10]. A precise description of their results would require the introduction of more notation, and will not be presented here. Informally speaking, they show that one can refine the measures of complexity used for voting classifiers such that it would encourage ensembles that can be grouped into a small number of compact clusters, each including base (“weak”) classifiers that are similar to one another.

---

<sup>1</sup>The terms “simple” and “complex” are not trivial to define, and their definition usually depends on what one tries to claim. In the next section we will use more rigorous definitions for the cases we analyze.

## 4 Analysis

The effect of adding noise to the training data depends on the learning algorithm used, and is highly complex. Even for the case of adding a zero-mean, low-variance Gaussian noise (noise injection) this effect was studied only for simple algorithms (e.g. [3]) or the square loss function [1].

In Sec. 4.1 we study the effect of Feature Knockout on the well known linear least square regression problem. We show that it leads to a scaled version of Tikhonov regularization. Compare this to Bishop’s result (using a Taylor expansion) that noise injection is equivalent to Tikhonov regularization. Following in Sec. 4.2, we will try to analyze how feature KO affects the variance of the learned classifier.

### 4.1 Effect of feature KO on linear regression

One of the most basic models we can apply to the data is the linear model. In this model, the input examples  $x_i \in \mathfrak{R}^n, i = 1..m$  are organized as the columns of the matrix  $A \in \mathfrak{R}^{n \times m}$ ; the corresponding  $y_i$  values are stacked in one vector  $y \in \mathfrak{R}^m$ . The prediction made by the model is given by  $A^\top h$ , where  $h$  is the vector of free parameters we have to fit to the data. In the common least squares case,  $\|y - A^\top h\|^2$  is minimized.

In the case that the matrix  $A$  is full rank and overdetermined, it is well known that the optimal solution is  $h = A^+ y$ , where  $A^+ = (AA^\top)^{-1}A$  is known as the pseudo inverse of the transpose of  $A$  (our definition of  $A$  is the transpose of the common text book definition). If  $A$  is not full rank, the matrix inverse  $(AA^\top)^{-1}$  is not well defined. However, as an operator in the range of  $A$  it is well defined, and the above expression still holds, i.e., even if there is an ambiguity in selecting the inverse matrix, there is no ambiguity in the operation of all possible matrices on the range of the columns of  $A$ , which is what we care about.

Even so, if the covariance matrix  $(AA^\top)$  has a large condition number (i.e., it is close to being singular), small perturbations of the data result in large changes to  $h$ , and the system is unstable. The solution fits the data  $A$  well, but does not fit data which is very close to  $A$ , hence there is overfitting. To stabilize the system, we apply regularization.

Tikhonov regularization is based on minimizing  $\|y - A^\top h\|^2 + \lambda\|h\|^2$ . This is equivalent to using a regularized pseudo inverse:  $A_\lambda^+ = (AA^\top + \lambda I)^{-1}A$ , where  $I$  is the identity  $n \times n$  matrix, and  $\lambda$  is the regularization parameter.

In many applications, the linear system we need to solve is badly scaled, e.g., one variable is much larger in magnitude than the other variables. In order to rectify this, we may apply a transformation to the data that weights each variable differently, or equivalently weight the vector  $h$  by applying a diagonal matrix  $D$ , such that  $h$  becomes  $\hat{h} = Dh$ .

Instead of solving the original system  $Ah = y$ , we now solve the system  $\hat{A}\hat{h} = y$ , where  $\hat{A} = D^{-1}A$ . Solving

this system using Tikhonov regularization is termed “scaled Tikhonov regularization.” If  $D$  is unknown, a natural choice is the diagonal matrix with the entries  $D_{kk} = \sqrt{(AA^\top)_{kk}}$  [13]. We will now show that using the knockout procedure to add many new examples is equivalent to scaled Tikhonov regularization, using the weight matrix above.

**Lemma 1** *When using the linear model with a least squares fit, applying the knockout procedure in Fig. 1 to generate many examples is equivalent to applying scaled Tikhonov regularization where  $D_{kk} = \sqrt{(AA^\top)_{kk}}$ .*

**Proof** For simplicity of the proof, we will make the following two assumptions: (1) All features have an expectation of 0. Without this assumption, our derivations below will be cluttered with many more elements. (2) Instead of sampling using the knockout procedure, we will just use an augmented data matrix containing all of the points that the knockout procedure can create. This is equivalent to studying the limit of infinite new examples, and allows us to ignore all the elements that go to zero as the number of virtual examples increases.

The vector  $\hat{h}$  which is fitted by means of a scaled Tikhonov regularization technique, with a parameter  $\lambda$  is given by:  $\hat{h} = \hat{A}^\top (\hat{A}\hat{A}^\top + \lambda I)^{-1} \hat{A}y = (D^{-1}AA^\top D^{-1} + \lambda D^{-1}D^2 D^{-1})^{-1} D^{-1}Ay = D(AA^\top + \lambda D^2)^{-1} D D^{-1}Ay = D(AA^\top + \lambda D^2)^{-1} Ay$ . Therefore,  $h = D^{-1}\hat{h} = (AA^\top + \lambda D^2)^{-1} Ay$ .

Now consider  $\tilde{A}$ , the matrix whose columns contain all the possible KO examples, together with the original data points. Let  $\tilde{y}$  be the corresponding labels. By applying a least square linear fit to these inputs, we find:  $\tilde{h} = (\tilde{A}\tilde{A}^\top)^{-1} \tilde{A}\tilde{y}$ . There are  $nm^2$  total examples created. Since we assumed all features have a mean of zero, all the knockout values of each feature cancel out. What remains is  $m(n-1)$  exact copies of each variable and  $\tilde{A}\tilde{y} = m(n-1)Ay$ .

We now examine the elements of the matrix  $\tilde{A}\tilde{A}^\top$ . The off diagonal elements represent the dot product between two different variables. Each variable holds either its original value, or a different value, but it may never happen that both contain the knockout values at once. It happens  $m(n-2)$  times for each input example that both features hold the original data sets. The rest of the cases average out to zero, because while one value is fixed the other value goes over all values which have a mean of zero. For the diagonal case, because of symmetry, each value appears  $nm$  times, making the diagonal  $nm$  times the diagonal of the original matrix.

Putting it all together:  $\tilde{h} = (\tilde{A}\tilde{A}^\top)^{-1} \tilde{A}\tilde{y} = (m(n-2)AA^\top + 2mD^2)^{-1} m(n-1)Ay = \frac{n-1}{n-2} (AA^\top + \lambda D^2)^{-1} Ay$ , where  $\lambda = \frac{2}{n-2}$ . The leading fraction does not change the sign of the results, and is close to one. Ignoring it, we obtain the results of a scaled Tikhonov regularization. The param-

eter  $\lambda$  can be controlled by allowing the knockout procedure to perturb more than one element. ■

To get a better understanding of the way Feature Knockout works, we study the behavior of scaled Tikhonov regularization. As mentioned in Sec. 3, in the boosting case, the knockout procedure is expected to produce solutions which make use of more features. Are these models more complex? This is hard to define in the general case, but easy to answer in the linear least square case study.

In linear models, the predictions  $\tilde{y}$  on the training data take the form:  $\tilde{y} = Py$ . For example, in the unregularized pseudo inverse case we have  $\tilde{y} = A^\top h = A^\top (AA^\top)^{-1} Ay$ , and therefore  $P = A^\top (AA^\top)^{-1} A$ . There is a simple measure of complexity called *the effective degrees of freedom* [8], which is just  $Tr(P)$  for linear models. A model with  $P = I$  (the identity matrix) has zero training error, but may overfit. In the full rank case, it has as many effective degrees of freedom as the number of features ( $Tr(P) = n$ ).

**Lemma 2** *The linear model obtained using scaled Tikhonov regularization has a lower effective degree of freedom than the linear model obtained using unregularized least squares.*

**Proof** This claim is very standard for Tikhonov regularization. Here we present a slightly more elaborate proof. Using the same rules we can prove other claims, such as the claim that the condition number of the matrix we invert using scaled Tikhonov regularization is lower than the one achieved without regularization. A lower condition number is known to lead to better generalization.

For scaled Tikhonov regularization we have  $\tilde{y} = Py$  where  $P = A^\top (AA^\top + \lambda D^2)^{-1} A$ . Therefore,  $Tr(P) = Tr(A^\top (AA^\top + \lambda D^2)^{-1} A) = Tr((DD^{-1} AA^\top D^{-1} D + \lambda D^2)^{-1} AA^\top) = Tr((D(D^{-1} AA^\top D^{-1} + \lambda I)D)^{-1} AA^\top) = Tr(D^{-1} (D^{-1} AA^\top D^{-1} + \lambda I)^{-1} D^{-1} AA^\top) = Tr((D^{-1} AA^\top D^{-1} + \lambda I)^{-1} D^{-1} AA^\top D^{-1}) = Tr((EE^\top + \lambda I)^{-1} EE^\top)$ , where  $E = D^{-1} A$ .

Let  $USV^\top = E$  be the Singular Value Decomposition of  $E$ , where  $S$  is a diagonal matrix, and  $U$  and  $V$  are orthonormal matrices. The above trace is exactly  $Tr((US^2U^\top + \lambda I)^{-1} US^2U^\top)$ . Let  $S^*$  be the diagonal matrix with elements  $S_{kk}^* = S_{kk}^2 + \lambda$ , then  $(US^2U^\top + \lambda I)^{-1} = (US^*U^\top)^{-1} = US^{*-1}U^\top$ . The above trace becomes  $Tr(US^{*-1}U^\top US^2U^\top) = Tr(S^{*-1}S^2U^\top U) = Tr(S^{*-1}S^2) = \sum_k \frac{S_{kk}^2}{S_{kk}^2 + \lambda} < rank(E) = rank(A)$ . Compare this value with the effective degrees of freedom of the unregularized least square solution:  $Tr(A^\top (AA^\top)^{-1} A) = Tr((AA^\top)^{-1} AA^\top) = rank(A)$ . The last equality also holds in the case where  $A$  is not full rank, in which case  $(AA^\top)^{-1}$  is only defined on the range of  $AA^\top$ . ■

Similar to the work done on noise injection, we examined the effect of our procedure on a simple regression technique. We saw that feature knockout resembles the effect of scaled Tikhonov regularization, i.e., high norm features are penalized by the knockout procedure. However, boosting over regressions stumps seems to be scale invariant. Multiplying all the values of a feature by some constant does not change the resulting classifier, since the process that fits the regression stumps (see Sec. 5) uses the values of each feature to determine the thresholds that it uses. However, a closer look reveals the connection between scaling and the effect of the knockout procedure on boosting.

*Boosting over stumps* (e.g., [16]) chooses at each round one out of  $n$  features, and one threshold for this feature. The thresholds are picked from the  $m$  possible values that exist in between every two sorted feature values. The feature and the threshold define a “weak classifier” (the basic building blocks of the ensemble classifier built by the boosting procedure [14]), which predicts -1 or +1 according to the threshold. Equivalently, we can say that boosting over stumps chooses from a set of  $nm$  binary features – these features are exactly the values returned by the weak classifiers. These  $nm$  features have different norms, and are not scale invariant. Let us call each such feature an  $nm$ -feature.

Using the intuitions of the linear least squares case, we would like to inhibit features of high magnitude. All  $nm$ -features have the same norm ( $\sqrt{(m)}$ ), but different entropies (a measure which is highly related to norm). These entropies depend only on the ratio of positive values in each  $nm$ -feature - call this ratio  $p$ .

Creating new examples using the Feature Knockout procedure does not change the number of possible thresholds, and therefore the number of features remains the same. The values of the new example in the  $nm$  feature space will be the same for all features originating from the  $n - 1$  features that were not changed in the knockout procedure. The value for a knocked-out feature (feature  $k$  in Fig. 1), will change if the new value is on the other side of the threshold as compared to the old value. This will happen with probability  $2p(1 - p)$ . If this sign flip happens then the feature is inhibited because it gives two different classifications to two examples with the same label (KO leaves labels unchanged). Note that the entropy of a feature with a positive ratio of  $p$  and the probability  $2p(1 - p)$  behave similarly: both rise monotonically for  $0 \leq p \leq 1/2$  and then drop symmetrically. Hence, We obtain the following result:

**Lemma 3** *Let  $t$  be a single  $nm$ -feature created by combining a single input feature with a threshold. The amount of inhibition  $t$  undergoes, as the result of applying feature knockout, grows monotonically with the entropy of  $p$ .*

Hence, similarly to the scaling in the linear case, the knockout procedure inhibits high magnitude features (here

the magnitude is measured by the entropy). Note that in the algorithm presented in Sec. 5, a feature is used for knockout only after it was selected to be a part of the output classifier. Still, KO inhibits more weak classifiers based on these features with higher entropies, making them less likely to get picked again. It is possible to perform this higher-entropy preferential inhibition directly on all features, therefore simulating the full knockout procedure. The implementation of this is left for future experiments.

## 4.2 Bias/variance decompositions

Many training algorithms can be interpreted as trying to minimize a cost function of the form  $\sum_{i=1}^n L(f(x_i), y_i)$ , where  $L$  is a loss function. For example, in the 0/1 loss function  $L(f(x), y) = (f(x) \neq y)$ , we pay 1 if the labels are different, 0 otherwise. By applying the knockout procedure to generate more training data, an algorithm that minimizes such a cost function will actually minimize:  $\sum_{i=1}^n \mathbb{E}_{\hat{x} \sim C_X(x_i)} L(f(\hat{x}), y_i)$ , where  $C_X(x)$  represents the distribution of all knocked-out examples created from  $x$ .

In the case of the square loss function  $L(f(x), y) = (y - f(x))^2$ , the cost function can be decomposed (similarly to [9]) into bias and variance, respectively:  $\sum_{i=1}^n \mathbb{E}_{\hat{x} \sim C_X(x_i)} L(f(\hat{x}), y_i) = \sum_{i=1}^n L(\mathbb{E}_{\hat{x} \sim C_X(x_i)} f(\hat{x}), y_i) + \sum_{i=1}^n \mathbb{E}_{\hat{w} \sim C_X(x_i)} L(\hat{w}, \mathbb{E}_{\hat{x} \sim C_X(x_i)} f(\hat{x}))$ .

Consider the variance term. Since Feature Knockout changes dramatically the value of one of the features in the vector  $x_i$ , one can show that if the term  $\sum_{i=1}^n \mathbb{E}_{\hat{w} \sim C_X(x_i)} L(\hat{w}, \mathbb{E}_{\hat{x} \sim C_X(x_i)} f(\hat{x}))$  is bounded, then our learning algorithm has a *bounded differences property* [12] which is equivalent to saying that by removing one of the features, the value of the learned function  $f$  will not change by more than a bounded amount. Consider a situation (which exists in our object recognition experiments) where our features are pulled independently from a pool of many possible features. The bounded difference property guarantees that with high probability the testing error we get is close to the expectation of the testing error with regards to selecting another set of random features of the same size. The formal discussion of these results, which follows the lines of [2] is omitted.

Let us now turn our attention to another “bias-variance” decomposition. Consider one based on the 0/1 loss function, as analyzed in [4]. We follow the terminology of [4] with a somewhat different derivation, and for the presentation below we include a simplified version. Assume for simplicity that each training example occurs in our dataset with only one label, i.e., if  $x_i = x_j$  then  $y_i = y_j$ . Define the *optimal prediction*  $f_*$  to be the “true” label  $f_*(x_i) = y_i$ . Define the *main prediction* of a function  $f$  to be just the prediction  $f(x)$ . The *bias* is defined to be the loss between the

optimal and main predictions:  $B(x) = (f(x) \neq f_*(x))$ . The *variance*  $V(x)$  is defined to be the expected loss of the prediction with regard to the main prediction:  $V(x) = \mathbb{E}_{\hat{x} \sim C_X(x)} (f(x) \neq f(\hat{x}))$ . These definitions allow us to present the following observation:

**Observation 1** *Let  $B_0$  be the set of all training- example- indices for which the bias  $B(x_i)$  is zero (the unbiased set). Let  $B_1$  be the set for which  $B(x_i) = 1$  (the biased set). Then,  $\sum_{i=1}^n \mathbb{E}_{\hat{x} \sim C_X(x_i)} (f(\hat{x}) \neq y_i) = \sum_{i=1}^n B(x_i) + \sum_{i \in B_0} V(x_i) - \sum_{i \in B_1} V(x_i)$*

In the unbiased case ( $B(x) = 0$ ), the variance ( $V(x)$ ) increases the training error. In the biased case ( $B(x) = 1$ ), the variance at point  $x$  decreases the error. A function  $f$ , which minimizes the training cost function that was obtained using Feature Knockout, has to deal with these two types of variance directly while training. Define the net variance to be the difference of the biased variance from the unbiased; a function trained using the Feature Knockout procedure is then expected to have a higher net variance than a function trained without this procedure. If we assume our corruption process  $C_X$  is a reasonable model of the robustness expected from our classifier, a good classifier would have a high net variance on the **testing** data<sup>2</sup>. The net variance measured in our experiments shows the effect of the Feature Knockout approach. An interesting application that is not explored in this paper is the exploitation of net variance to derive confidence bars **at a point** (i.e., a measure of certainty in our prediction). Since Feature Knockout emphasizes these differences, it yields narrower confidence bars.

## 5 The GentleBoostKO algorithm

While our regularization procedure can be applied, in principle, to any learning algorithm, using it directly when the number of features  $n$  is high might be computationally demanding. This is because for each one of the  $m$  training examples, as many as  $n(m-1)$  new examples can be created. Covering even a small portion of this space might require the creation of many synthesized examples.

The randomized procedure in Fig. 1 samples this large space of synthesized training examples. Still, if there are many features, the sampling would probably be too sparse.

However, for some algorithms our regularization technique can be applied with very little overhead. For boosting over regression stumps, it is sufficient to modify those features that participate in the trained ensemble (i.e., those features that actually participate in the classification).

The basic algorithm used in our experiments is specified in Fig. 2. It is a modified version of the GentleBoost algo-

<sup>2</sup>We omit the formal discussion on the relation between variance on training examples, and variance on testing examples.



---

**Input:**  $(x_1, y_1), \dots, (x_m, y_m)$  where  $x_i \in \mathbb{R}^n$ ,  $y_i \in Y = \pm 1$ .  
**Output:** Composite classifier  $H(x)$ .

---

1. Initialize weights  $w_i \leftarrow 1/m$ .
  2. for  $t = 1, 2, 3, \dots, T$ .
    - (a) For each feature  $k$ , fit a regression function  $f_t^{(k)}(x)$  by weighted least squares on  $y_i$  to  $x_i$  with weights  $w_i$ ,  $i = 1..m + t - 1$ .
    - (b) Let  $k_{min}$  be the index of the feature with the minimal associated weighted least square error.
    - (c) Update the classifier  $H(x) \leftarrow H(x) + f_t^{(k_{min})}$
    - (d) Use Feature KO to create a new example  $x_{m+t}$ :
 
$$\begin{aligned} \text{Select two random indices } 1 \leq a, b \leq m \\ x_{m+t} &\leftarrow x_a \\ x_{m+t}(k_{min}) &\leftarrow x_b(k_{min}) \\ y_{m+t} &\leftarrow y_a \end{aligned}$$
    - (e) Set new example weight  $w_{m+t}$  to that of its source:
 
$$w_{m+t} \leftarrow w_a$$
    - (f) Update the weights and normalize:
 
$$\begin{aligned} w_i &\leftarrow w_i e^{-y_i f_t^{(k_{min})}(x_i)}, \quad i = 1..m + t \\ w_i &\leftarrow w_i / \sum_{i=1}^{m+t} w_i \end{aligned}$$
  3. Output the final classifier  $H(x)$
- 

Figure 2: The GentleBoostKO Algorithm. Steps d and e constitute the differences from the original GentleBoost.

rithm [7]. GentleBoost seems to converge faster than AdaBoost, and performs better for object detection problems [15]. At each boosting round, a regression function is fitted (by weighted least-squared error) to each feature in the training set. We used linear regression for our experiments, fitting parameters  $a, b$  and  $th$  so that our regression functions are of the form  $f(x) = a(x > th) + b$ . The regression function with the least weighted squared error is added to the total classifier  $H(x)$  and its associated feature ( $k_{min}$ ) is used for Feature Knockout (step d).

In the Feature Knockout step, a new example is created using the class of a randomly selected example  $x_a$  and all of its feature values except for the value at  $k_{min}$ . The value for this feature is taken from a second randomly-selected example  $x_b$ . The new example  $x_{m+t}$  is then appended to the training set. In order to quantify the importance of the new example in the boosting process, a weight has to be assigned to it. The weight  $w_{m+t}$  of the new example is estimated by copying the weight of the example from which most of the features are taken ( $x_a$ ). Alternatively, a more precise weight can be determined by applying the total classifier  $H(x)$  to the new example.

As with any boosting procedure, each iteration ends with

the update of the weights of all examples (including the new one), and a new round of boosting begins. This iterative process finishes when the weights of the examples converge, or after a fixed number of iterations. In our experiments, we stopped the boosting after 100 rounds—enough to ensure convergence in all cases.

## 6 Experiments

**UCI repository experiments.** We evaluated our methods on 10 UCI repository datasets that are suitable for binary classification, either by thresholding the value of the target function (e.g. the price in the housing dataset) at its median, or by picking a label to be the positive class. These 10 datasets were: arrhythmia, dermatology, e-coli, glass, heart, housing, letters, segmentation, wine, and yeast.

We have split each dataset randomly into 10% training, 90% testing, and ran each of the following classifiers: GentleBoost, GentleBoostKO (Fig. 2), AdaBoost, AdaBoost with a knockout procedure (similar to GentleBoostKO), and linear SVM. We also ran linear SVM on a dataset that contained 100 examples generated in accordance with the knockout procedure of Fig. 1. SVMs with different nonlinear kernels produced either similar or worse results. In addition, we report results for GentleBoost combined with noise injection (the algorithm that adds gaussian noise to the examples), with the best noise variance we found. By selecting results according to the performance on the testing data, the noise injection results were biased, and should only be taken as an upper bound for the performance of noise injection. Tab. 1 shows the results, averaged over 10 independent runs.

In this table we measured the mean error, the standard deviation of the error, and the number of features used by the classifiers (SVM always uses the maximal number of features). We also measure the variance over a distribution of knockout examples for correct classifications (unbiased variance), and incorrect classifications (biased variance) (Sec. 4.2). This variance was computed in the following way: for each **testing** example we generated 50 knockout examples (Fig. 1), and computed the variance over these 50 examples. We averaged the variance over all biased and unbiased testing examples. A good classifier produces more variance for incorrectly classified examples, and only a little variance for correctly classified ones.

It is apparent from the results that: (1) In general the knockout procedure (KO) helps GentleBoost, raising it to the same level of performance as SVM. (2) KO seems to help AdaBoost as well, but not always. It is not clear whether knockout helps SVM. (3) KO seems to help increase the net variance (which is good, see Sec. 4.2). (4) As expected KO produces classifiers that tend to use more features. (5) KO shows different, mostly better, performance

than noise injection (GentleBoostNI).

**Visual recognition using the Caltech datasets.** We tested our GentleBoostKO algorithm on several Caltech object recognition datasets that were presented in [6]. In each experiment we had to distinguish between images containing an object and background images that do not contain the object. The datasets: Airplanes, Cars, Faces, Leafs and Motorbikes, as well as the background images were downloaded from <http://www.vision.caltech.edu/>. For the experiments we used the predefined splits (available to all the datasets but the Leafs dataset). For leafs, we used a random split of 50% training and 50% testing. Note that since our methods are discriminative, we needed a negative training set. For this end, we removed 30 random examples from the negative testing set, and used them for training.

To turn each image into feature-vectors we used 500 C2 features. These extremely successful features allow us to learn to recognize objects using few training images, and the results seem to be comparable or better than the results reported in [5]. The results are shown in Fig. 3. To compare with previous work, we used the error at the equilibrium-point between false and true positives as our error-measure. It is clear that for a few dozen examples, SVM, GentleBoost and GentleBoostKO have the same performance level. However, for only a few training examples, GentleBoost does not perform as well as SVM, while GentleBoostKO achieves the same level of performance.

We also tried to apply Lowe’s SIFT features [11] to the same datasets, although these features were designed for a different task. For each image, we used Lowe’s binaries to compute the SIFT description of each key point. We then sampled from the training set 1000 random keypoints  $k_1, \dots, k_{1000}$ . Let  $\{k_i^I\}$  be the set of all keypoints associated with image  $I$ . We represented each training and testing image  $I$  by a vector of 1000 elements:  $[v^I(1) \dots v^I(1000)]$ , such that  $v^I(j) = \min_i \|k_j - k_i^I\|$ . Note that in [11] the use of the ratio of distances between the closest and the next closest points were encouraged (and not just the minimum distance). For our application, which disregards all geometric information, we found that using the minimum gives much better results. For the testing and training splits reported in [6] we got the following results (ME=mean error, EqE=error at equilibrium):

Algorithm	Planes	Cars	Faces	Leaves	Motor.
Lin. SVM ME	0.104	0.019	0.107	0.118	0.033
gentleB ME	0.118	0.036	0.168	0.137	0.026
gentleBKO ME	0.100	0.033	0.119	0.114	0.023
Lin. SVM EqE	0.108	0.018	0.111	0.126	0.007
gentleB EqE	0.120	0.037	0.166	0.132	0.003
gentleBKO EqE	0.111	0.030	0.136	0.120	0.008

**Car type identification.** This dataset consists of 480 images of private cars, and 248 images of mid sized vehicles (such as SUV’s). All images are  $20 \times 20$  pixels, and

Dataset	Algorithm	Mean Error	Unbias Var.	Biased Var.	Net Var.	Feat. Used
ARRHY.	AdaB	47.0%±4.0	0.011	0.010	-0.001	6.2
	AdaBKO	41.4%±4.3	0.022	0.026	0.004	45.0
	GentleB	40.0%±7.0	0.100	0.130	0.030	43.8
	GentleBKO	37.3%±3.7	0.042	0.068	0.026	55.2
	GentleBNI	35.2%±1.2	0.093	0.137	0.044	45.2
	LinSVM	38.5%±3.5	0.035	0.046	0.011	279.0
DERM	AdaB	4.7%±2.3	0.042	0.040	-0.002	1.6
	AdaBKO	21.8%±43.8	0.015	0.165	0.187	17.4
	GentleB	3.6%±2.5	0.651	0.712	0.061	3.8
	GentleBKO	1.7%±1.5	0.018	0.135	0.117	18.9
	GentleBNI	2.4%±0.4	0.534	0.607	0.073	6.2
	LinSVM	0.8%±0.7	0.013	0.128	0.115	34.0
ECOLI	AdaB	10.4%±5.8	0.291	0.451	0.161	4.0
	AdaBKO	11.7%±6.1	0.243	0.511	0.269	6.2
	GentleB	10.9%±5.9	0.546	0.625	0.079	3.2
	GentleBKO	5.8%±1.4	0.236	0.540	0.304	6.2
	GentleBNI	7.0%±0.9	0.485	0.714	0.229	3.6
	LinSVM	8.3%±3.3	0.235	0.483	0.248	7.0
GLASS	AdaB	37.6%±7.4	0.272	0.296	0.024	5.4
	AdaBKO	33.3%±7.0	0.302	0.337	0.035	8.0
	GentleB	34.8%±5.2	0.420	0.431	0.011	6.0
	GentleBKO	30.6%±6.9	0.287	0.320	0.033	8.0
	GentleBNI	33.6%±3.3	0.320	0.343	0.023	7.8
	LinSVM	40.8%±2.8	0.238	0.250	0.013	8.0
HEART	AdaB	37.9%±6.9	0.262	0.286	0.024	8.0
	AdaB	24.5%±3.6	0.139	0.217	0.078	6.2
	AdaBKO	21.9%±4.4	0.148	0.319	0.171	10.4
	GentleB	26.7%±5.2	0.223	0.342	0.119	9.8
	GentleBKO	23.9%±3.4	0.191	0.351	0.160	12.6
	GentleBNI	26.7%±2.4	0.215	0.334	0.119	11.4
HOUSING	LinSVM	23.8%±5.6	0.205	0.344	0.138	13.0
	LinSVM KO	24.2%±4.5	0.193	0.344	0.151	13.0
	AdaB	16.9%±1.8	0.138	0.192	0.054	4.6
	AdaBKO	16.7%±1.4	0.156	0.302	0.146	10.6
	GentleB	20.0%±3.9	0.243	0.358	0.115	10.0
	GentleBKO	17.6%±0.9	0.175	0.380	0.205	12.7
LETTERS	GentleBNI	20.0%±1.1	0.204	0.361	0.157	10.6
	LinSVM	21.2%±4.6	0.297	0.454	0.157	13.0
	LinSVM KO	17.1%±1.8	0.254	0.430	0.176	13.0
	AdaB	14.0%±0.8	0.018	0.062	0.044	4.6
	AdaBKO	11.4%±2.8	0.040	0.147	0.107	9.4
	GentleB	4.6%±0.9	0.117	0.577	0.460	15.0
SEGM.	GentleBKO	3.6%±0.6	0.094	0.551	0.457	15.5
	GentleBNI	2.8%±0.1	0.096	0.510	0.414	15.6
	LinSVM	4.2%±0.4	0.186	0.751	0.565	16.0
	LinSVM KO	4.1%±0.5	0.172	0.749	0.578	16.0
	AdaB	7.0%±0.9	0.054	0.121	0.068	3.4
	AdaBKO	9.4%±3.0	0.041	0.219	0.178	11.4
WINE	GentleB	6.7%±1.4	0.251	0.399	0.148	4.8
	GentleBKO	7.8%±2.7	0.031	0.302	0.271	14.0
	GentleBNI	6.7%±1.1	0.268	0.445	0.177	4.0
	LinSVM	2.9%±1.6	0.146	0.409	0.263	19.0
	LinSVM KO	3.4%±2.3	0.181	0.546	0.366	19.0
	AdaB	15.9%±8.2	0.128	0.155	0.027	3.8
YEAST	AdaBKO	12.9%±6.1	0.124	0.318	0.194	11.2
	GentleB	17.1%±6.9	0.605	0.723	0.119	4.2
	GentleBKO	12.2%±4.9	0.117	0.324	0.207	11.9
	GentleBNI	17.1%±2.6	0.216	0.288	0.072	9.0
	LinSVM	12.5%±9.8	0.149	0.411	0.261	13.0
	LinSVM KO	15.7%±8.5	0.166	0.379	0.213	13.0
YEAST	AdaB	31.1%±1.0	0.020	0.025	0.006	2.4
	AdaBKO	31.3%±1.1	0.023	0.034	0.011	4.6
	GentleB	33.9%±5.7	0.329	0.443	0.114	6.4
	GentleBKO	32.7%±2.6	0.280	0.438	0.158	8.0
	GentleBNI	32.4%±0.8	0.326	0.495	0.169	7.2
	LinSVM	31.2%±0.3	0.000	0.001	0.001	8.0
LinSVM KO	31.2%±0.3	0.000	0.000	0.000	8.0	

Table 1: Results for datasets from the UCI repository. Each data set was split to 10% training, 90% testing. The mean error, its standard deviation, the mean biased, unbiased and net variance as well as the mean number of features used are shown for 10 independent experiments.

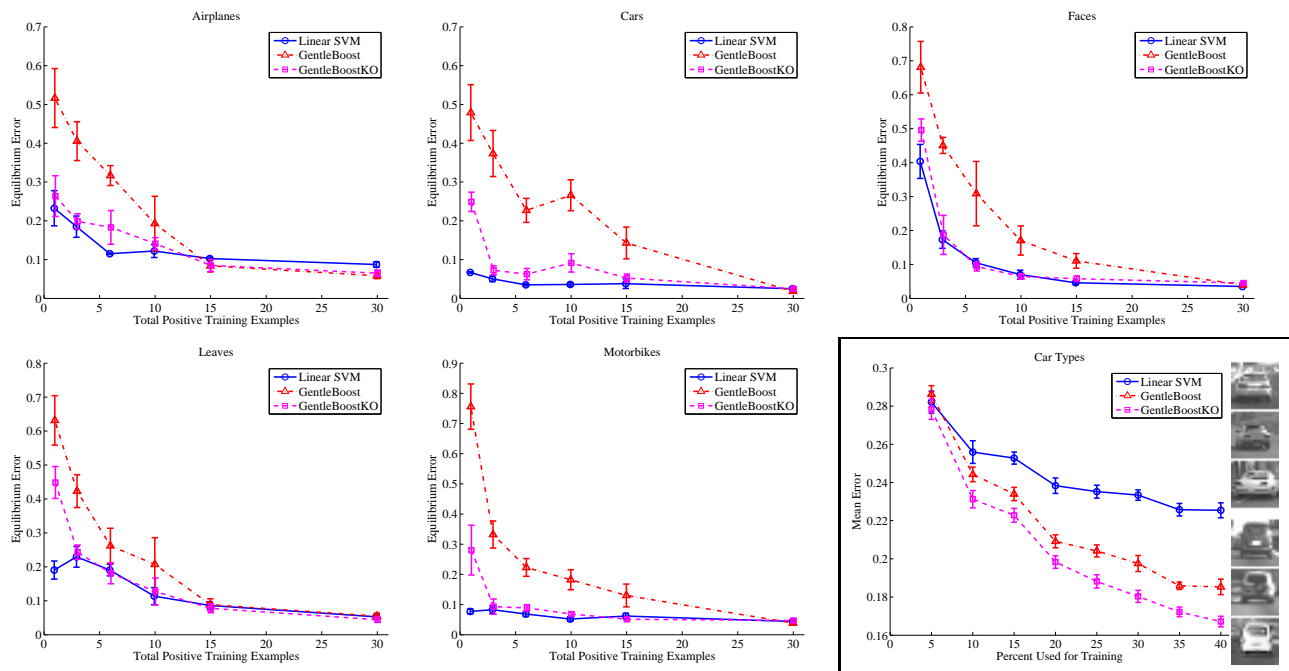


Figure 3: A comparison using the C2 features between GentleBoost, GentleBoostKO (Fig. 2), and linear SVM on the five Caltech datasets: Airplanes, Cars, Faces, Leaves and Motorbikes. The graphs show the the equilibrium error rate vs. the number of training examples used from the class we want to detect. In each experiment, the test set was fixed to be the same as those described in [6]. **Lower right corner:** The results of applying the three algorithms to the car types dataset, together with example images. The results shown are mean and standard error of 30 independent experiments versus percentile of training images.

were collected using the car detector of Mobileye corp., on a video stream taken from the front window of a moving car. The task is to learn to identify private cars from mid sized vehicles, which has some safety applications. Taking into account the low resolution and the variability in the two classes, this is a difficult task. The results are shown on the bottom right corner of Fig. 3. Each point of the graph shows the mean error when applying the algorithms to training sets of different size (between 5 and 40 percent of the data). The rest of the examples were used for testing. It is evident that for this dataset GentleBoost outperforms SVM. Still, GentleBoostKO does even better.

## 7. Summary and Conclusions

Boosting algorithms, and especially the use of GentleBoost over regression stumps, are gaining a lot of popularity in the computer vision community. However, GentleBoost does not show good performance, compared to SVM, when learning from a small dataset. In this work we propose an enhancement to the GentleBoost algorithm, that brings its level of performance to be the same of SVM.

## References

- [1] C.M Bishop. Training with Noise is Equivalent to Tikhonov Regularization. *Neural Computation*, 1995.
- [2] O. Bousquet, A. Elisseeff. Stability and Generalization. *JMLR*, 2002.
- [3] L. Breiman. Heuristics of Instability and Stabilization in Model Selection. *Ann. Statist.*, 1996.
- [4] P. Domingos. A Unifies Bias-Variance Decomposition for Zero-One and Squared Loss. *Proc. Int. Conf. AI*, 2000.
- [5] L. Fei-Fei, R. Fergus, & P. Perona. A Bayesian approach to unsupervised 1-Shot learning of Object categories. *ICCV03*.
- [6] R. Fergus, P. Perona, and A. Zisserman. Object Class Recognition by Unsupervised Scale-Invariant Learning. *CVPR 2003*
- [7] J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: a statistical view of boosting. *Ann. Statist.*, 2000.
- [8] T. Hastie, R. Tibshirani, J. H. Friedman. The Elements of Statistical Learning Springer, 2001.
- [9] S. Geman, E. Bienenstock, R. Doursat. Neural networks and the bias/variance dilemma. *Neural Computations*, 1992.

- [10] V. Koltchinskii and D. Panchenko. Complexities of convex combinations and bounding the generalization error in classification. to appear in *Ann. Statist.*
- [11] D.G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 2004.
- [12] G. Lugosi. Concentration-of-measure inequalities
- [13] A. Neumaier. Solving ill-conditioned and singular linear systems: A tutorial on regularization. *SIAM Review* 1998.
- [14] R.E. Schapire. A brief introduction to boosting. Proc. Int. Joint Conf. AI, 1999.
- [15] A. Torralba, K.P. Murphy and W.T. Freeman. Sharing features: efficient boosting procedures for multiclass object detection. *CVPR*, 2004.
- [16] P. Viola and M. Jones. Rapid object detection using a boosted cascade of simple features. *CVPR*, 2001.
- [17] J. Weston, S. Mukherjee, O. Chapelle, M. Pontil, T. Poggio and V. Vapnik. Feature Selection for SVMs. *NIPS*, 2001.