

Support Vector Components Analysis

Michiel H. van der Ree¹, Jos B.T.M. Roerdink², Christophe Phillips³,
Gaëtan Garraux³, Eric Salmon³ and Marco A. Wiering⁴

1- Semiotic Labs B.V.
Science Park 402, Amsterdam - The Netherlands

2- Johann Bernoulli Institute for Mathematics and Computer Science
University of Groningen, Nijenborgh 9, Groningen - The Netherlands

3- Cyclotron Research Centre
University of Liège, Allée du Six Août 8 B30, Liège - Belgium

4- Institute of Artificial Intelligence and Cognitive Engineering
University of Groningen, Nijenborgh 9, Groningen - The Netherlands

Abstract. In this paper we propose a novel method for learning a distance metric in the process of training Support Vector Machines (SVMs) with the radial basis function kernel. A transformation matrix is adapted in such a way that the SVM dual objective of a classification problem is optimized. By using a wide transformation matrix the method can effectively be used as a means of supervised dimensionality reduction. We compare our method with other algorithms on a toy dataset and on PET-scans of patients with various Parkinsonisms, finding that our method either outperforms or performs on par with the other algorithms.

1 Introduction

The Support Vector Machine [1] is one of the most popular algorithms for solving both regression and classification problems in machine learning. The algorithm is robust and offers an excellent generalization performance, making it very well suited for small datasets with many features. One of the drawbacks of SVMs not using a linear kernel is that the algorithm is a *black box*: The model can't be inspected to see what features of the data are decisive for the eventual prediction. In addition, when using the radial basis function (RBF) kernel, SVMs are very sensitive to a proper scaling of the input data.

The method proposed in this paper aims to tackle both aforementioned problems. By learning a quadratic distance metric during SVM training the model becomes less sensitive to the scaling of data. By forcing the distance metric to be of rank 2 or 3, we can visualize a lower dimensional representation of the input data and make the relations learned by the SVM more intelligible.

Outline Section 2 will introduce the support vector components analysis (SVCA) algorithm. Next, we illustrate how the proposed algorithm can be used as a means of supervised dimensionality reduction (SDR). Section 4 will cover the

setup and results of experiments conducted with the SVCA and other SDR algorithms. A conclusion is presented in Section 5.

2 Support Vector Components Analysis

We have a labeled dataset consisting of P real-valued input vectors $\mathbf{x}_1, \dots, \mathbf{x}_P$ in \mathbb{R}^N and corresponding class labels c_1, \dots, c_P . A matrix $\mathbf{T} \in \mathbb{R}^{M \times N}$ defines a linear map from \mathbb{R}^N to \mathbb{R}^M by mapping \mathbf{x} to $\mathbf{T}\mathbf{x}$. We try to optimize \mathbf{T} such that when the transformation is applied to the dataset, class differences are emphasized in the transformed space. Our algorithm tries to maximize the margins of one-versus-rest Support Vector Machines in the transformed space. Prior to explaining our approach in more detail, we review the objective used in support vector classification.

2.1 Support Vector Classification

In binary soft margin linear support vector classification, training consists of solving the following constrained optimization problem:

$$\min_{\mathbf{w}, \xi, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \quad (1)$$

subject to constraints $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$ and $\xi_i \geq 0$. Here, \mathbf{w} is a weight vector, b is a bias value, (\mathbf{x}_i, y_i) is a training sample and its associated label in $\{-1, 1\}$, ξ_i is a so-called “slack variable” that measures the degree of constraint violation \mathbf{x}_i and C is a constant determining the trade-off between margin maximization and error minimization.

Introduction of Lagrange multipliers $\boldsymbol{\alpha}$ and solving for the coordinates of a saddle point allow us to reformulate the primal objective and its constraints as:

$$\max_{\boldsymbol{\alpha}} \mathcal{Q}(\boldsymbol{\alpha}) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) \quad (2)$$

subject to constraints $0 \leq \alpha_i \leq C$ and $\sum_i \alpha_i y_i = 0$. Once the $\boldsymbol{\alpha}$ maximizing (2) is found, the linear support vector classifier determines its output using:

$$f(\mathbf{x}) = \text{sgn} \left(\sum_i \alpha_i (\mathbf{x}_i \cdot \mathbf{x}) + b \right). \quad (3)$$

Since both the dual objective (2) and the model output (3) only depend on inner products between patterns, the model can be made non-linear by using a *kernel function* $K(\mathbf{x}_i, \mathbf{x})$. such as polynomial functions and the radial basis function.

2.2 The SVCA Objective

The basic idea of our algorithm is to learn a projection matrix \mathbf{T} such that the margin between classes in the transformed space is maximized. If we train a one-versus-rest SVM for each class ℓ in the transformed space, the primal objective

of each linear binary classification SVM becomes:

$$\min J^\ell(\mathbf{w}, \boldsymbol{\xi}) = \left[\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \right] \quad (4)$$

now subject to constraints $y_i^\ell(\mathbf{w} \cdot \mathbf{T}\mathbf{x}_i + b) \geq 1 - \xi_i$ and $\xi_i \geq 0$, where we use $\boldsymbol{\xi}$ to denote the vector containing all ξ_i 's. Correspondingly, the new kernelized 'dual' objective is defined as:

$$\min_{\mathbf{T}} \max_{\boldsymbol{\alpha}} \mathcal{Q}^\ell(\boldsymbol{\alpha}; \mathbf{T}) = \left[\sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i^\ell y_j^\ell K(\mathbf{T}\mathbf{x}_i, \mathbf{T}\mathbf{x}_j) \right] \quad (5)$$

subject to constraints $0 \leq \alpha_i \leq C$ and $\sum_i \alpha_i y_i^\ell = 0$. Note that the dual objective needs to be minimized w.r.t. \mathbf{T} as is also the case in other (multi)-kernel learning approaches [2].

2.3 Training Procedure

We use the following procedure to find the "support vector components": First, we solve the quadratic programming subproblem of finding the $\boldsymbol{\alpha}^\ell$ that maximizes the expression in (5) for each of the one-versus-rest support vector machines. Since that expression is equal to the objective being maximized in normal support vector machine training, we can use tried and tested optimization methods such as sequential minimal optimization (SMO) [3] to do so.

Then, for the optimized $\boldsymbol{\alpha}^\ell$'s, we can minimize the sums of all dual-objectives w.r.t. \mathbf{T} using stochastic gradient descent. In stochastic gradient descent, we minimize $\sum_\ell \mathcal{Q}^\ell$ by minimizing the gradients of single examples. Writing the expression in (5) as $\mathcal{Q}^\ell = \sum_i q_i^\ell$ with single example terms

$$q_i^\ell = \alpha_i^\ell - \frac{1}{2} \alpha_i^\ell y_i^\ell \sum_j \alpha_j^\ell y_j^\ell K(\mathbf{T}\mathbf{x}_i, \mathbf{T}\mathbf{x}_j) \quad (6)$$

we find the following derivative of q_i^ℓ w.r.t. \mathbf{T} :

$$\frac{\partial q_i^\ell}{\partial \mathbf{T}} = -\frac{1}{2} \alpha_i^\ell y_i^\ell \sum_j \alpha_j^\ell y_j^\ell \frac{\partial K(\mathbf{T}\mathbf{x}_i, \mathbf{T}\mathbf{x}_j)}{\partial \mathbf{T}}. \quad (7)$$

We alternate between optimizing all $\boldsymbol{\alpha}^\ell$'s and \mathbf{T} a preset number of times. Alternatively, we can use batch gradient descent or batch methods such as resilient backprop (RPROP, [4]) and adjust \mathbf{T} using $\sum_\ell \partial \mathcal{Q}^\ell / \partial \mathbf{T}$ instead of its single example based estimate.

3 SVCA as Supervised Dimensionality Reduction

When using a wide matrix \mathbf{T} such that $M < N$, the SVCA algorithm can be used as a means of supervised dimensionality reduction. SDR can have multiple

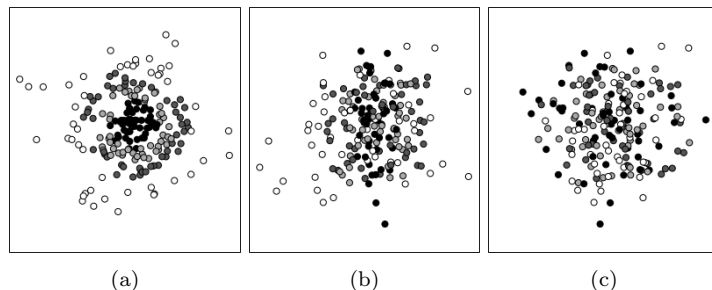


Figure 1: The artificial dataset of concentric rings. Shown are scatter plots in which both features are relevant to the labels (a), only one feature is relevant (b) and none of the features are relevant (c).

advantages. When using an RBF kernel all training patterns have to be stored in memory. With a wide \mathbf{T} , the memory cost of saving these patterns is reduced by a factor M/N . The most interesting application is when we set M to 2 or 3, so we can visualize the low-dimensional representation of the dataset. This can be useful for exploring relations between and separability of the classes. Therefore, in this paper we use $M = 2$ or 3. Similar SDR methods learning a linear transformation matrix are neighbourhood components analysis (NCA) [5], local Fisher discriminant analysis (LFDA) [6] and Limited Rank Matrix Learning Vector Quantization (LiRaM LVQ) [7].

4 Experiments and Results

We report the performance of the SVCA algorithm compared to other SDR algorithms on two datasets: a toy problem of concentric rings and FDG-PET scans of patients with various Parkinsonisms. For SVCA, we use the RPROP algorithm to optimize the transformation matrix and we use an RBF kernel.

4.1 Experiments on Artificial Data

Inspired by the concentric ring data in [5], we create an artificial dataset in the following way: First, we create patterns $\mathbf{x}_1 \dots \mathbf{x}_P$ in \mathbb{R}^8 by drawing from $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where $\boldsymbol{\mu} = \mathbf{0}$ and $\boldsymbol{\Sigma}$ is the 8×8 identity matrix. Then we assign labels based solely on the distance from the origin in the first two dimensions, i.e. $\sqrt{x_1^2 + x_2^2}$. This results in classes that take the shape of concentric rings in the class-relevant subspace. In total, we create 200 patterns belonging to four different classes. In defining class boundaries, we ensure that each class has about the same number of patterns. Figure 1 shows the dataset thus generated.

We compare the performance of SVCA with the three other SDR techniques mentioned in Section 3: NCA, LiRaM LVQ and LFDA. The algorithms are compared on 100 randomly generated datasets. We find that LiRaM LVQ and LFDA never succeed in finding the underlying structure. For SVCA and NCA, we have

		SVCA	
		incorrect	correct
NCA	incorrect	$e_{00} = 30$	$e_{01} = 20$
	correct	$e_{10} = 6$	$e_{11} = 44$

Table 1: Contingency table of SVCA and NCA error rates in the concentric ring experiment.

simply counted the number of times each algorithm finds the ‘right’ projection by learning a transformation matrix with $M = 2$, resulting in the contingency table shown in table 1. Under the null hypothesis that NCA and SVCA have the same error rate. We computed the number of correctly learned projections, and used McNemar’s test to obtain a p -value of 0.011. SVCA therefore significantly outperforms NCA and the other algorithms in this experiment.

4.2 Experiments on FDG-PET Scans

Here we apply the SDR algorithms to the same set of PET scans as used in [8]. These scans were obtained in two different locations between 1993 and 2009 and are comprised of 42 Parkinson’s disease patients, 31 multiple system atrophy patients, 26 progressive supranuclear palsy patients and 21 corticobasal syndrome patients. Each scan consists of 153,594 voxels. We preprocess the data using the Scaled Subprofile Modelling routine [9], leaving us with projections onto principal components. We retain the first n principal components that explain at least 75% of the variance in the data. This procedure is applied in a k -fold fashion, so the number of selected components will differ per fold.

We predefine 100 splits of the data. In each split, 10% of the patterns has been randomly assigned to the test set, the rest of the patterns are used for training. We report mean test accuracies and their standard deviations in Table 2. NCA and LFDA do not provide an explicit prediction for new patterns. We have chosen to assign labels according to the nearest neighbor classification in the transformed space, where the number of neighbors was determined through cross-validation. Running paired t -tests on the different fold error rates, we find no significant differences between the various algorithms. However, we do find that for $M = 3$ the performance of the SDR algorithms rival that of an RBF SVM with parameters (C, γ) optimized through cross-validation. The results for the SDR algorithms are impressive since unlike these algorithms, the RBF SVM does not act on data transformed by a matrix with limited rank.

	SVCA	LRMLVQ	NCA	LFDA	RBF SVM:
$M = 2$	0.58 ± 0.15	0.56 ± 0.13	0.59 ± 0.12	0.61 ± 0.12	0.68 ± 0.13
$M = 3$	0.68 ± 0.12	0.67 ± 0.13	0.66 ± 0.12	0.68 ± 0.13	

Table 2: Average test accuracies and standard deviations of the various algorithms on 100 test/train splits on the data from [8].

5 Conclusion

We have presented the novel learning method SVCA that can be used for both distance metric learning and dimensionality reduction. In our experiment on toy data, we found that SVCA is most likely to succeed in finding the rings hidden in the data, only having NCA as a true competitor. In [10], we explore the relation between NCA and SVCA in more detail and find that NCA can be seen as doing SVCA with fixed α values. These results suggest that adapting the alpha values as we do in SVCA helps in finding the latent structure in a noisy dataset.

The results of the experiment on FDG-PET scans do not show any significant differences between the different SDR methods, so SVCA can only be considered to perform “on-par” with the other SDR methods in this experiment. In turn, all SDR algorithms rival the performance of an optimized RBF SVM while still allowing the relations they discover in the training set to be inspected.

In future work, we will examine the use of non-linear transformation functions. Furthermore, it would be very interesting to integrate the SVCA algorithm in the multi-layer SVM architecture [11]. Finally, we would like to compare our method to other kernel or distance-function learning algorithms.

References

- [1] V.N. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, 1995.
- [2] A-D. Pietersma, L.R.B. Schomaker, and M.A. Wiering. Kernel learning in support vector machines using dual-objective optimization. In *Proceedings of the 23rd Belgian-Dutch Conference on Artificial Intelligence*, pages 167–174, 2011.
- [3] J. Platt. Sequential minimal optimisation: a fast algorithm for training support vector machines. Technical Report MSR-TR-98-14, Microsoft Research, 1998.
- [4] M. Riedmiller and H. Braun. A direct adaptive method for faster backpropagation learning: The Rprop algorithm. In *Proceedings of the IEEE International Conference on Neural Networks*, pages 586–591, 1993.
- [5] J. Goldberger, S. Roweis, G. Hinton, and R. Salakhutdinov. Neighbourhood components analysis. In *Advances in Neural Information Processing Systems 17*, pages 513–520. MIT Press, 2004.
- [6] M. Sugiyama. Dimensionality reduction of multimodal labeled data by local Fisher discriminant analysis. *Journal of Machine Learning Research*, 8:1027–1061, 2007.
- [7] K. Bunte, P. Schneider, B. Hammer, F. Schleif, T. Villmann, and M. Biehl. Limited rank matrix learning, discriminative dimension reduction and visualization. *Neural Networks*, 26:159–173, 2012.
- [8] G. Garraux, C. Phillips, J. Schrouff, A. Kreisler, C. Lemaire, Degueldre C., C. Deltour, R. Hustinx, A. Luxen, A. Desée, and E. Salmon. Multiclass classification of FDG PET scans for the distinction between Parkinson’s disease and atypical parkinsonian syndromes. *NeuroImage: Clinical*, pages 883–893, 2013.
- [9] G.E. Alexander and J.R. Moeller. Application of the scaled subprofile model to functional imaging in neuropsychiatric disorders: a principal component approach to modeling brain function in disease. *Human Brain Mapping*, 2:79–94, 1994.
- [10] M.H. van der Ree. Explorations in intelligible classification. Master’s thesis, University of Groningen, the Netherlands, 2014.
- [11] M.A. Wiering and L.R.B. Schomaker. Multi-layer support vector machines. In J.A.K. Suykens, M. Signoretto, and A. Argyriou, editors, *Regularization, Optimization, Kernels, and Support Vector Machines*, chapter 20. Chapman and Hall, 2014.