

Invariance in Kernel Methods - Distance and Integration Kernels

PASCAL Workshop on Machine Learning,
SVM and Large Scale Optimization 2005

Dipl.-Math. Bernard Haasdonk



University of Freiburg, Germany
Computer Science Department
Chair of Pattern Recognition and Image Processing

Abstract

In pattern analysis with kernel methods, it is widely accepted, that the kernel function is the main ingredient to represent any kind of prior knowledge. A frequently observed kind of prior knowledge is invariance with respect to transformations of single patterns. For this case we present two generic methods of incorporating such knowledge into quite arbitrary kernels. The first is based on invariant distances, the second on invariant integration. As the former method can not guarantee to result in positive definite kernels, we discuss the use of indefinite kernels in machine learning. Application on several classification problems with SVMs demonstrates the competitiveness in terms of recognition accuracy and computational complexity with existing methods.

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 2

Overview

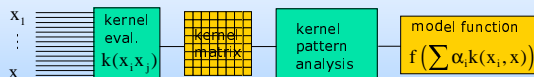
- Motivation and Notions
 - Kernel Methods and Invariance
- Invariant **Distance** Substitution Kernels
 - Distance Substitution Kernels
 - Tangent Distance Kernels
 - Application Raman-Microspectroscopy
- Transformation **Integration** Kernels
 - Application OCR
- Indefinite Kernels in SVM
 - Interpretation in pE spaces
- Summary and Perspectives

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 3

Motivation and Notions

Kernel Methods [SS02,SC04]

- Typical tasks:
 - Classification, Regression, Clustering, Novelty Detection,...
- Multitude of kernel methods: SVM, SVR, KPCA,...
- Analysis chain:

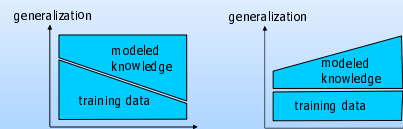


- Multitude of kernels for various datatypes
 - Vectorial, sequences, graphs, finite state machines...
- Kernel matrix is information „bottleneck”
=> importance of kernel choice!

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 5

Importance of Prior Knowledge

- Without prior assumptions no generalization possible
There is no best classification method („No Free Lunch Thm.”)
There is no best feature representation („Ugly Duckling Thm.”)
- Equivalence of training data and model complexity



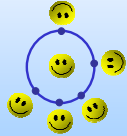
- Consequences of modelling more prior knowledge
 - less training samples required
 - better generalization in case of equal number of samples
 - feature representation or classification can be simplified

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 6

Invariance as Prior Knowledge

- transformation knowledge
i.e. samples maintain inherent meaning under certain transformations

- continuous



- discrete

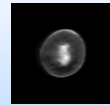


- continuous and discrete



Transformation Examples

- complete group



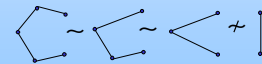
- subset of group

6 ~ 6 ~ 6 + 9
N ~ N ~ N + Z
M ~ M ~ M + W

- reversible



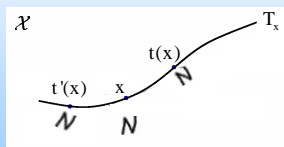
- irreversible



Transformation Formalization

- Patterns x stem from pattern space \mathcal{X}
- Transformations $t: \mathcal{X} \rightarrow \mathcal{X}$
- Set of transformations T
- Set of transformed patterns of x

$T_x := \{t(x) | t \in T\}$
= „similar“
meaning as x



- Total invariance: $k(x, x') = k(t(x), t'(x'))$

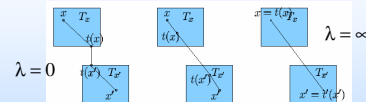
Specific Goals

- Wanted properties for general approach:
 - support for various kernel methods, not only SVM
 - support for many kernels, not only Gaussian
 - support for infinite set of transformations, exponentially many transformation combinations
 - support for discrete, continuous, group and non-group transformations
 - Applicability, good memory, time and model complexities
 - Adjustability of degree of invariance

Invariant Distance Substitution Kernels

Invariant Distances

- Distances of $T_x, T_{x'}$ better than $d(x, x')$



- Invariant Distances

- two-sided distance

$$d_{2S}(x, x') := \min_{t, t'} (d(t(x), t'(x')) + \lambda \Omega(t, t'))$$

- cost function

$$\Omega(t, t') \geq 0, \quad \Omega(t, t') = 0 \Leftrightarrow t = t' = \text{id}$$

- λ is regularization parameter

- similar: one-sided distance + symmetrization

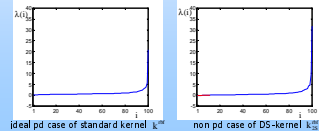
Distance Substitution Kernels [H804]

- Distance d : symmetric, nonnegative, zero-diagonal
- Distance-based kernels: $k(\|x - x'\|) \Rightarrow k_d(x, x') := k(d(x, x'))$
- Inner-product case: $k(\langle x, x' \rangle) \Rightarrow k_d(x, x') := k(\langle x, x' \rangle_d^O)$ where O is an arbitrary origin and $\langle x, x' \rangle_d^O := -\frac{1}{2}(d(x, x')^2 - d(x, O)^2 - d(x', O)^2)$
- Examples:
 - $k_d^{\text{lin}}(x, x') := \langle x, x' \rangle_d^O$ $k_d^{\text{nd}}(x, x') := -d(x, x')^\beta, \beta \in [0, 2]$
 - $k_d^{\text{pol}}(x, x') := (1 + \gamma \langle x, x' \rangle_d^O)^p$ $k_d^{\text{rbf}}(x, x') := e^{-\gamma d(x, x')^2}, \gamma \in \mathbb{N}, \gamma \geq 0$
- Expectation: Similar behaviour as standard kernels

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 13

Positive Definiteness

- Equivalent conditions:
 - d is Hilbertian metric
 - k_d^{nd} is cpd for all $\beta \in [0, 2]$
 - k_d^{lin} is pd
 - k_d^{rbf} is pd for all $\gamma \geq 0$
 - k_d^{pol} is pd for all $\gamma \geq 0, p \in \mathbb{N}$
- DS-kernels mostly non-(c)pd: counterexample by invalid Δ -inequality
- Empirical observations:
 - only weak pd-ness violation on real world data,

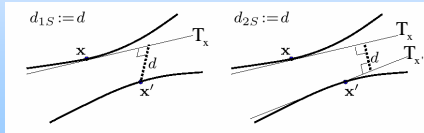


17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 14

Tangent Distance Kernels

- Assumption: differentiable transformations
- Tangent Distance: local linear approximation [SLD93]
 - distance of tangent spaces instead of manifolds
 - symmetric two-sided distance $d_{2S}(x, x')$
 - nonsymmetric one-sided distance $d_{1S}(x, x')$

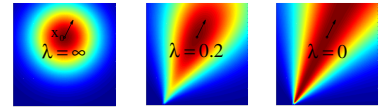
$$d_{1S}(x, x') = \min_p \|x + \sum_i p_i t_i - x'\|^2 + \lambda \|p\|^2$$



17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 15

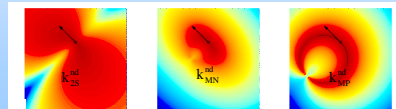
Invariance in 2D

- Scaling: $k_{\text{MP}}^{\text{rbf}}(x_i, x)$ varying x



- invariance can be adjusted by regularization λ

- Rotation:



- More than linear invariances can be captured

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 16

Application Raman Spectra

- Project OMIB in BMBF framework „Biophotonics“
 - detection of clean room contaminations by Raman microscopy
- Raman spectroscopy



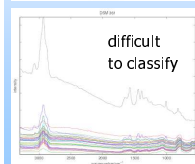
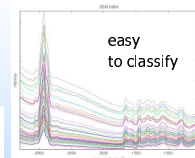
- Use of the „Raman effect“, inelastic scattering of light
- Distribution of energy differences => Raman spectrum
- non-destructive method, fingerprint of chemical bindings
- Pattern variations due to
 - measurement duration, background radiation
 - thickness of sample, heterogeneity, growing time, nutrition condition, temperature, photo bleaching

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 17

Application Raman Spectra

- Details on dataset [R&al05]:
 - 2545 spectra, 1833 features
 - 20 classes,
 - unequally distributed

Name	Code	Class Number	Number of Samples
B. parvulus	DSM 27	1	57
B. parvulus	DSM 361	5	43
B. sphaerulus	DSM 28	2	53
B. sphaerulus	DSM 396	6	42
B. subtilis subsp. subtilis	DSM 10	0	206
B. subtilis subsp. spizizenii	DSM 147	3	42
E. coli	DSM 423	7	51
E. coli	DSM 898	8	24
E. coli	DSM 499	9	20
M. luteus	DSM 348	4	619
M. luteus	DSM 20030	13	48
M. luteus	DSM 20315	16	20
M. luteus	DSM 20318	18	20
S. colibaci subsp. colibaci	DSM 6669	10	67
S. colibaci subsp. colibaci	DSM 20260	15	65
S. colibaci subsp. urealyticum	DSM 6718	11	65
S. colibaci subsp. urealyticum	DSM 6719	12	65
S. Epidermidis	DSM RFP 62A	19	805
S. Warneri	DSM 20636	14	67
S. Warneri	DSM 20316	17	74



17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 18

Application Raman Spectra

- Details on tangents:

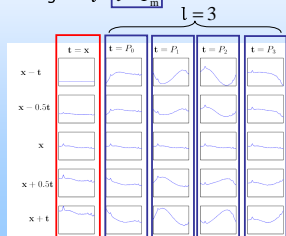
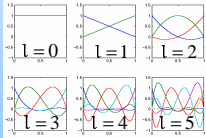
- intensity scale: scale tangent $t = x$

- baseline shift:

Lagrange polynomials of degree l : $t = P_m$

$$P_m(x) = \prod_{i \neq m} \frac{(x - t_i)}{(t_m - t_i)}$$

$m = 0, \dots, l$



Application Raman Spectra

- Recognition results

- LIBSVM, one-vs-one SVM, k_{25}^{rbf}

- Simultaneous scaling and Lagrange tangents

- 8x8 grid-search for (C, γ)

- LOO and class-wise average LOO minimization

reg.pur λ	LOO-error [%], varying degree l					av-LOO-error [%], varying degree l				
	0	1	2	3	4	0	1	2	3	4
1.000000	4.72	4.56	4.56	4.44	4.28	10.71	10.33	10.20	9.83	9.54
0.100000	4.48	4.44	4.56	4.72	4.48	10.16	9.63	9.85	10.49	9.76
0.010000	4.05	4.09	3.85	3.93	3.97	9.06	9.18	8.53	8.71	9.49
0.001000	3.65	3.69	3.77	3.77	3.73	8.61	8.46	8.55	8.65	8.50
0.000100	3.85	3.03	2.99	3.18	3.18	8.87	6.87	7.26	7.48	7.32
0.000010	5.78	3.65	3.34	3.38	3.65	12.70	8.09	7.46	7.68	8.26
0.000001	11.43	7.27	6.05	5.38	5.38	26.96	16.68	14.17	12.86	12.51
base SVM	4.20					9.58				

- LOO and av-LOO clearly improved by tangents

- Regularization required due to scaling tangent

Transformation Integration Kernels

Transformation Integration Kernels

- Motivation:

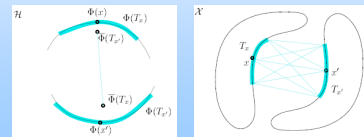
- Maintain pd-ness by avoiding min/max-operations

- Extend Haar-integration framework [594] to kernels

- Adjustability of degree of invariance

- Definition of TI-Kernels [HWB05]:

$$k_{TI}(x, x') := \left(\int_T \Phi(t(x)) dt, \int_T \Phi(t'(x')) dt' \right) = \iint_{T \times T} k(t(x), t'(x')) dt dt'$$

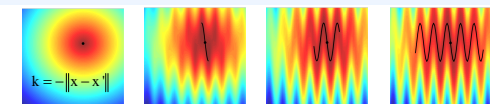


- (c)pd-ness transferred from base-kernel

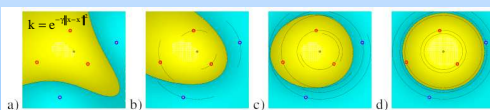
Invariance in 2D

- Invariance adjustability:

highly nonlinear sinus shifts



- Effect in SVM: rotation



- Similar: Total invariance of linear, polynomial kernels

Acceleration

- Single Kernel Evaluation: Integral Reduction (IR)

If T are invertible and compatible with k

$$\iint_{T \times T} k(t(x), t'(x')) dt dt' = \iint_{T \times T} k(x, t^{-1} \circ t'(x')) dt dt' = \int_{T^{-1} \circ T} k(x, \bar{t}(x')) d\bar{t} \Rightarrow \text{squareroot complexity reduction}$$

- SVM-Training: SV-extraction (SV)

- Perform initial ordinary SVM-training

- Extract SVs

- Train invariant SVM on the SVs

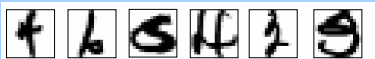
\Rightarrow linear complexity reduction

Application USPS-Digits

- Details on dataset:
 - standard benchmark dataset
 - 7291 training-, 2007 testpatterns of handwritten digits given by 16x16 grayvalue bitmaps
 - samples: easy to classify



difficult to classify



Application USPS-Digits

- SVM Recognition Results:

increasing rotation integration range increasing x-y translation integration range

ϕ -range [rad]	k^{TM}	test error [%]
0		4.5
$\pm 0.04\pi$		4.1
$\pm 0.08\pi$		4.2
$\pm 0.12\pi$		3.9
$\pm 0.16\pi$		4.2

x-y-range [pixels]	k^{TM}	test error [%]
0		4.5
± 1		3.7
± 2		3.2
± 3		3.3
± 4		3.2

- Both TI-kernels superior over base kernel
- x-y-integration better than rotations
- TI-kernels yield state-of-the-art result as VSV [SBV96]

Application USPS-Digits

- Complexity comparison to VSV:
 - shared non-optimized parameters, x-y-translation, + -2 pixels, 3x3 shifts per sample

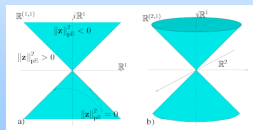
Method	test-error [%]	train-time [sec]	test-time [sec]	average #SV
TI-SVM	3.6	1771	479	412
TI-SVM, IR	3.6	810	176	412
TI-SVM, SV	3.6	113 + 130 + 297	466	410
TI-SVM, SV + IR	3.6	113 + 130 + 91	172	410
VSV-SVM	3.5	113 + 864 + 1925	177	4240

- TI-kernels produce small models
- No recognition degradation by SV-acceleration
- Acceleration techniques clearly successful
- Accelerated TI-SVM consistently faster than VSV

Indefinite Kernels in SVM

Pseudo-Euclidean Spaces [G85]

- Real finite dimensional vector spaces $\mathbb{R}^{(p,q)} := \mathbb{R}^p \times i\mathbb{R}^q$ of signature (p,q)
- symmetric (indefinite) inner-product $\langle z, z' \rangle_{pE} := z_p^T z'_p - z_q^T z'_q = z^T M z'$ $M = \text{diag}(I_p, -I_q)$
- squared norm $\|z\|_{pE}^2 := \langle z, z \rangle_{pE} = z^T M z$ can be negative:
- squared distance $\|z - z'\|_{pE}^2 := \langle z - z', z - z' \rangle_{pE}$
- orthogonality $\langle z, z' \rangle_{pE} = z^T M z' = 0$
- hyperplanes $\mathcal{H}: \langle z, w \rangle_{pE} + b = 0$



pE Feature Space Embedding

- Data dependent pE-embedding:
 - Given data $\{x_i\}_{i=1}^n \Rightarrow$ Existence of pE space $\mathbb{R}^{(p,q)}$ + sym. kernel $k \Rightarrow$ + embedding $\Phi: \mathcal{X} \rightarrow \mathbb{R}^{(p,q)}$
 - Representation of kernel $k(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle_{pE}$
 - Construction by Eigendecomposition [LM04, GHB09, PPD01]:
 - $K = U \Lambda U^T$ $p := \dim(\lambda^+), q := \dim(\lambda^-)$
 - $\Lambda = \text{diag}(\lambda^+, \lambda^-)$ $\Phi(x_i) := (\sqrt{|\Lambda|} U^T)_i$

Geometrical Interpretation of Kernel Operations

- Norm of Feature vector
 $k(x_i, x_i) = \|\Phi(x_i)\|_{pE}^2$
- Kernel induced distance
 $d^2(x_i, x_j) := k(x_i, x_i) - 2k(x_i, x_j) + k(x_j, x_j)$
 $d^2(x_i, x_j) = \|\Phi(x_i) - \Phi(x_j)\|_{pE}^2$
- Linear combinations
 $\sum \alpha_i k(x_i, x) = \langle \sum \alpha_i \Phi(x_i), x \rangle_{pE}$
- Centering of Kernel Matrix
 $\mathbf{K}_C = \mathbf{J} \cdot \mathbf{K} \cdot \mathbf{J}$
 $\mathbf{J} := \mathbf{I} - \mathbf{1} \cdot \mathbf{1}^T / n \Rightarrow \sum \Phi_C(x_i) = \mathbf{0}$
- Projections, variance analysis, Eigendecompositions, Optimization problems ... general kernel methods?

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 31

Capacity Estimate

- VC-bound for pE hyperplanes in $\mathbb{R}^{(p,q)}$
 - data embedded s.th. $r^2 \|\Phi(x_i)\|_E^2 \leq \|\Phi(x_i)\|_{pE}^2 \leq R^2$
 - $\mathcal{F} := \{\text{sgn} \langle \mathbf{w}, \mathbf{z} \rangle\}$ set of hyperplanes, which are canonical wrt the data and satisfy
 $\lambda^2 \|\mathbf{w}\|_E^2 \leq \|\mathbf{w}\|_{pE}^2 \leq \Lambda^2$
- then holds

$$h(\mathcal{F}) \leq \left(\frac{R\Lambda}{r\lambda} \right)^2$$
- Proof: Vapnik [v95] + modif. Cauchy Schwartz.
- \Rightarrow SRM: minimize $\|\mathbf{w}\|_{pE}^2 = \mathbf{w}^T \mathbf{M} \mathbf{w}$ while maintaining strict positivity

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 32

Optimal Separation of Convex Hulls

[H05, extension of B800, C800]

- Fomalization
 train: $\min_{z^+, z^-} \|\mathbf{z}^+ - \mathbf{z}^-\|_{pE}^2$
 s.t. $\mathbf{z}^\pm \in \text{conv}_n \{\Phi(x_i) | y_i = \pm 1\}$
 classify: sign of $f(x) = \|\Phi(x) - \mathbf{z}^+\|_{pE}^2 - \|\Phi(x) - \mathbf{z}^-\|_{pE}^2$
 - Separable case Non-separable case
-
- Training and test avoid explicit pE-embedding
 - Reasonable solution: $\mathbf{w}^T \mathbf{M} \mathbf{w} \geq 0$
 - Optimum exists and obtained in $\text{span} \{\Phi(x_i)\}$

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 33

Indefinite SVM in pE Space

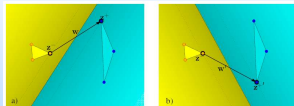
- SVM-"primal": $\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{M} \mathbf{w}$
 s.t. $y_i (\mathbf{w}^T \mathbf{M} \Phi(x_i) + b) \geq 1$
- Usual interpretation of SV + coefficients
- Non-SV are correctly classified,
- Only bounded SV can be misclassified
- SVM is CH-classifier:
 - Feasible, stationary points + local optima transfer between SVM and CH, if $\mathbf{w}^T \mathbf{M} \mathbf{w} \geq 0$
 - E.g Nonzero SVM-solution $\alpha \Rightarrow$ CH-solution $\bar{\alpha}$

$$\bar{\alpha} := 2\alpha \sum_i \alpha_i \quad \mu := 2C / \sum_i \alpha_i$$
- CH and SVM hyperplanes are parallel, even identical if coefficients are not bounded

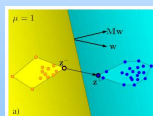
17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 34

Numerics

- Convergence to stationary point, libsvm [LL03]
- Multiple solutions



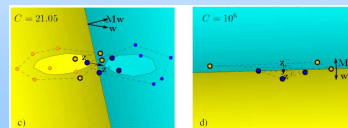
- Uniqueness in extreme indefinite cases



17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 35

Practical Criteria for Indefinite SVM

- Criterion for suitability: #bSV
 - No (few) bounded $\alpha_i \Rightarrow$ no (few) training errors
- Criterion for unsuitability: $\mathbf{w}^T \mathbf{M} \mathbf{w} \leq 0$
 - after training: $\sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j)$
 - before training: signature (p,q)
- Criterion for suitability: Distance of Class Means
 - If DCM is positive, sufficiently low C yields solution



$$DCM^2 = \sum_{i,j} c_i c_j k(x_i, x_j)$$

$$c_i = \begin{cases} 1/n^+ & \text{for } y_i = +1 \\ -1/n^- & \text{for } y_i = -1 \end{cases}$$

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 36

Summary and Perspectives

Summary

- Principle „Incorporation of transformation knowledge in kernels“ for improving kernel methods
- General framework I: IDS-Kernels
 - Distance substitution kernels with invariant distances
 - Application: Tangent Distance Kernels on Spectra
 - Efficient for multiple invariances, but indefiniteness
- General framework II: TI-Kernels
 - Integration over transformations
 - Application: Offline HWR (USPS digits)
 - Efficient for few invariances, but positive definiteness
- Indefinite SVM
 - Constructive geometric interpretation: pE CH-separation
 - Convergent implementation required
 - Practical criteria for suitability

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 38

Perspectives

- Invariance
 - combination with invariant representations
 - Learning of transformation directions
- Kernel design
 - Relaxing the pd-ness of kernels yields much wider flexibility
 - Applications with *proximity data*: nonlinear analysis!
- Kernel methods
 - Acceptance/robustness against indefinite kernels?
 - Interpretation in pE-spaces as basis for geometrical/numerical/statistical analysis and new methods
- Non-convex optimization
 - Large scale implementations with convergence statements
 - number and quality of solutions.

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 39

References

- [BB00] K.P. Bennett, E.J. Breitenstein, „Duality and geometry in SVM classifiers“, Proc. 17th ICML, pp. 57-64, 2000.
- [CB00] D.J. Crijp, C.J.C. Burges, „A geometric interpretation on nu-svm classifiers“, NIPS 12, pp. 233-239, MIT-Press, 2000.
- [G85] L. Goldfarb, „A new approach to pattern recognition“, Progress in Pattern Recognition 2, pp. 291-402, Elsevier, 1985.
- [GH099] T. Graedel, R. Herbrich, P. Bollmann-Sdorra, K. Obermayer, „Classification on pairwise proximity data“, NIPS 11, pp. 438-444, MIT Press, 1999.
- [HB04] B. Haasdonk, C. Bahlmann, „Learning with Distance Substitution Kernels“, Proc. 26th DAGM, pp. 220-227, Springer, 2004.
- [H05] B. Haasdonk, „Feature space interpretation of SVMs with indefinite kernels“, IEEE TPAMI, 27(4):482-492, april 2005.
- [HV05] B. Haasdonk, A. Vossen, H. Burkhardt, „Invariance in Kernel Methods by Haar-Integration Kernels“, Proc. 13th ICCV, 2005, submitted.
- [LM04] J. Laub, K.-R. Müller, „Feature Discovery in non-metric pairwise data“ Journal of Machine Learning Research, 5:810-818, 2004.
- [L03] H.-T. Lin, C.-J. Lin, „A study on sigmoid kernels for SVM and the training of non-PSD kernels by SMO-type methods.“ Technical Report, National Taiwan University, March 2003.
- [PPD01] E. Pekala, P. Paclik, R. Dujin, „A generalized kernel approach to dissimilarity based classification“ Journal of Machine Learning Research, 2:175-211, 2001.
- [R&05] P. Rösch, M. Schmitt, K.-D. Peschke, O. Ronneberger, H. Burkhardt, H.-W. Moitzkus, M. Lankers, S. Hofer, H. Thiele, J. Popp, „Chemosensory identification of single bacteria relevant for clean room production by micro-Raman spectroscopy“, Applied and Environmental Microbiology, 2005, to appear.
- [S02] B. Schölkopf, A. Smola, „Learning with Kernels“, MIT-Press, 2002.
- [S94] H. Schölkopf, „Constructing invariant features by averaging techniques“, Proc. 12th ICPR, IEEE, 1994.
- [S04] J. Shawe-Taylor, N. Cristianini, „Kernel Methods for Pattern Analysis“, Cambridge University Press, 2004.
- [SLD93] P. Y. Simard, Y. A. LeCun, J.S. Denker, „Efficient Pattern Recognition using a new transformation distance“, NIPS 5, pp. 50-58, Morgan Kaufman 1993.
- [V95] V. Vapnik, „The nature of statistical learning theory“, Springer, New York, 1995.

17.3.2005 B. Haasdonk, Computer Science Department, University of Freiburg, Germany 40

Thank You! 😊

Any questions?