

Semi-Supervised Support Vector Machines and Application to Spam Filtering

Alexander Zien

Empirical Inference Department, Bernhard Schölkopf
Max Planck Institute for Biological Cybernetics

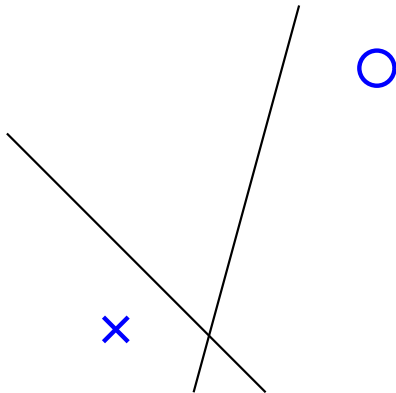
ECML 2006 – Discovery Challenge

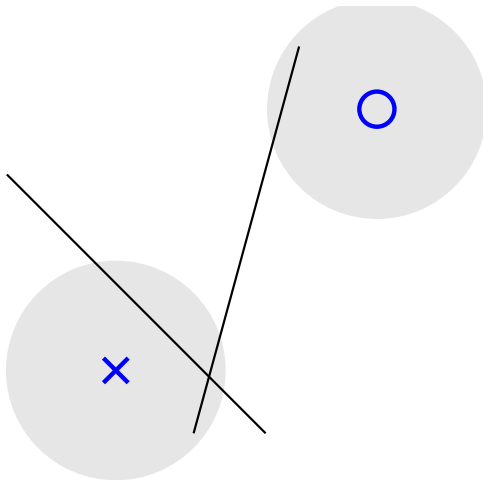


- 1 Introduction
- 2 Training a S^3VM
 - Why It Matters
 - Some S^3VM Training Methods
 - Gradient-based Optimization
 - The Continuation S^3VM
- 3 Overview of SSL
 - Assumptions of SSL
 - A Crude Overview of SSL
 - Combining Methods
- 4 Application to Spam Filtering
 - Naive Application
 - Proper Model Selection
- 5 Conclusions



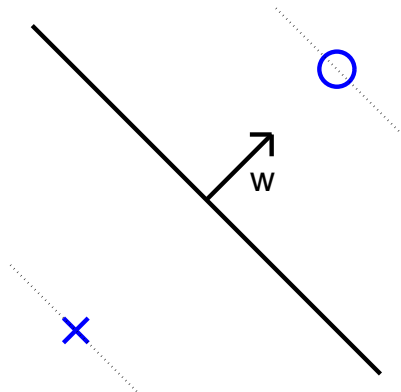
find a linear classification boundary





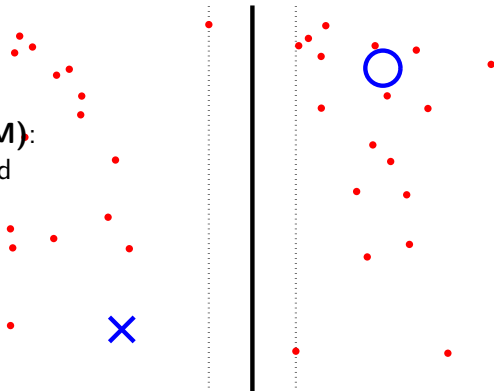
not robust wrt input noise!

SVM:
maximum margin
classifier



$$\min_{\mathbf{w}, b} \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w}}_{\text{regularizer}} \quad \text{s.t.} \quad y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

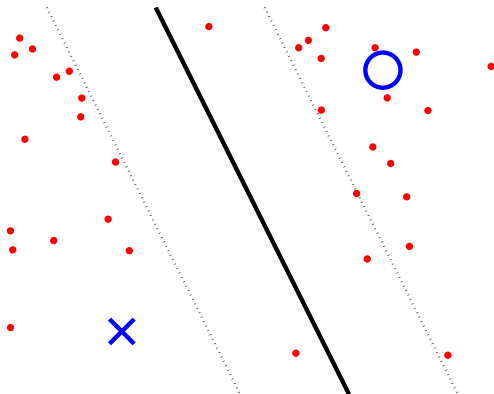
S^3VM (TSVM):
 semi-supervised
 (transductive)
 SVM



$$\min_{\mathbf{w}, b, (y_j)} \underbrace{\frac{1}{2} \mathbf{w}^\top \mathbf{w}}_{\text{regularizer}}$$

$$\text{s.t.} \quad \begin{aligned} y_i (\mathbf{w}^\top \mathbf{x}_i + b) &\geq 1 \\ y_j (\mathbf{w}^\top \mathbf{x}_j + b) &\geq 1 \end{aligned}$$

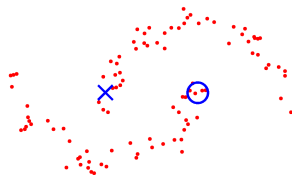
soft margin
 S^3VM



$$\begin{aligned}
 \min_{\mathbf{w}, b, (y_j), (\xi_k)} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\
 & + C \sum_i \xi_i \\
 & + C^* \sum_j \xi_j \quad \text{s.t.} \quad \xi_i \geq 0 \quad \xi_j \geq 0 \\
 & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \\
 & y_j (\mathbf{w}^T \mathbf{x}_j + b) \geq 1 - \xi_j
 \end{aligned}$$

“Two Moons” toy data

- easy for human (0% error)
- hard for S^3VM s!



S ³ VM optimization method		test error	objective value	
<i>global min.</i> {Branch & Bound		0.0%	7.81	
<i>find local minima</i>	{	CCCP	64.0%	39.55
		S ³ VM ^{light}	66.2%	20.94
		∇S ³ VM	59.3%	13.64
		cS ³ VM	45.7%	13.25

- objective function is good for SSL
- ⇒ **try to find better local minima!**

$$\min_{\mathbf{w}, b, (y_j), (\xi_k)} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j$$

$$\text{s.t.} \quad \begin{aligned} y_i (\mathbf{w}^\top \mathbf{x}_i + b) &\geq 1 - \xi_i & \xi_i &\geq 0 \\ y_j (\mathbf{w}^\top \mathbf{x}_j + b) &\geq 1 - \xi_j & \xi_j &\geq 0 \end{aligned}$$

Mixed Integer Programming [Bennett, Demiriz; NIPS 1998]

- global optimum found by standard optimization packages (eg CPLEX)
- **combinatorial & NP-hard !**
- only works for small sized problems

$$\begin{aligned} \min_{\mathbf{w}, b, (y_j), (\xi_k)} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j \\ \text{s.t.} \quad & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \quad \xi_i \geq 0 \\ & y_j (\mathbf{w}^\top \mathbf{x}_j + b) \geq 1 - \xi_j \quad \xi_j \geq 0 \end{aligned}$$

S^3VM^{light} [T. Joachims; ICML 1999]

- train SVM on labeled points, predict y_j 's
- in prediction, always make sure that

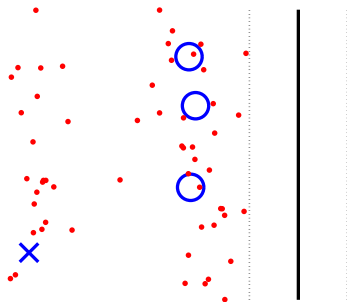
$$\frac{\#\{y_j = +1\}}{\#\text{ unlabeled points}} = \frac{\#\{y_i = +1\}}{\#\text{ labeled points}} \quad (1)$$

- with stepwise increasing C^* do
 - ① train SVM on all points, using labels (y_i) , (y_j)
 - ② predict new y_j 's s.t. "balancing constraint" (*)

$$\min_{\mathbf{w}, b, (y_j), (\xi_k)} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j$$

$$s.t. \quad \begin{aligned} y_i(\mathbf{w}^\top \mathbf{x}_i + b) &\geq 1 - \xi_i & \xi_i &\geq 0 \\ y_j(\mathbf{w}^\top \mathbf{x}_j + b) &\geq 1 - \xi_j & \xi_j &\geq 0 \end{aligned}$$

Balancing constraint required to avoid **degenerate solutions!**



$$\min_{\mathbf{w}, b, (y_j), (\xi_k)} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j$$

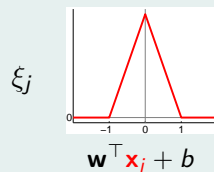
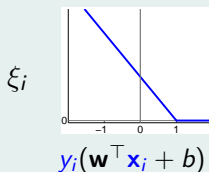
$$\text{s.t.} \quad \begin{aligned} y_i(\mathbf{w}^\top \mathbf{x}_i + b) &\geq 1 - \xi_i & \xi_i &\geq 0 \\ y_j(\mathbf{w}^\top \mathbf{x}_j + b) &\geq 1 - \xi_j & \xi_j &\geq 0 \end{aligned}$$

Effective Loss Functions

$$\xi_i = \min \left\{ 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b), 0 \right\}$$

$$\xi_j = \min_{y_j \in \{+1, -1\}} \left\{ 1 - y_j(\mathbf{w}^\top \mathbf{x}_j + b), 0 \right\}$$

loss
functions



$$\min_{\mathbf{w}, b, (y_j), (\xi_k)} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j$$

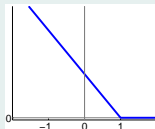
$$\text{s.t.} \quad \begin{aligned} y_i(\mathbf{w}^\top \mathbf{x}_i + b) &\geq 1 - \xi_i & \xi_i &\geq 0 \\ y_j(\mathbf{w}^\top \mathbf{x}_j + b) &\geq 1 - \xi_j & \xi_j &\geq 0 \end{aligned}$$

Resolving the Constraints

$$\frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i l_1 \left(y_i(\mathbf{w}^\top \mathbf{x}_i + b) \right) + C^* \sum_j l_u \left(\mathbf{w}^\top \mathbf{x}_j + b \right)$$

loss
functions

l_1



l_u



$$\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \ell_l \left(y_i (\mathbf{w}^T \mathbf{x}_i + b) \right) + C^* \sum_j \ell_u \left(\mathbf{w}^T \mathbf{x}_j + b \right)$$

CCCP-S³VM [R. Collobert et al.; ICML 2006]

- CCCP: “Concave Convex Procedure”
- objective = convex function + concave function
- starting from SVM solution, iterate:
 - ① approximate concave part by linear function at given point
 - ② solve resulting convex problem

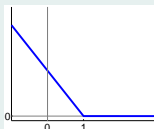
[Fung, Mangasarian; 1999]

- similar approach
- restricted to linear S³VMs

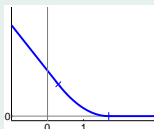
$$\frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i l_l \left(y_i (\mathbf{w}^\top \mathbf{x}_i + b) \right) + C^* \sum_j l_u \left(\mathbf{w}^\top \mathbf{x}_j + b \right)$$

S³VM as Unconstrained Differentiable Optimization Problem

original
loss
functions

 l_l

 l_u


smooth
loss
functions

 l_l

 l_u


$$\frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \ell_l \left(y_i (\mathbf{w}^\top \mathbf{x}_i + b) \right) + C^* \sum_j \ell_u \left(\mathbf{w}^\top \mathbf{x}_j + b \right)$$

∇ S³VM [Chapelle, Zien; AISTATS 2005]

- simply do gradient descent!
- thereby stepwise increase C^*

contS³VM [Chapelle et al.; ICML 2006]

... in more detail on next slides!

$$\frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \ell_l \left(y_i (\mathbf{w}^\top \mathbf{x}_i + b) \right) + C^* \sum_j \ell_u \left(\mathbf{w}^\top \mathbf{x}_j + b \right)$$

Hard Balancing Constraint

S^3VM^{light}
constraint

$$\frac{\#\{y_j = +1\}}{\#\text{unlabeled points}} = \frac{\#\{y_i = +1\}}{\#\text{labeled points}}$$

equivalent
constraint

$$\underbrace{\frac{1}{m} \sum_j \text{sign}(\mathbf{w}^\top \mathbf{x}_j + b)}_{\text{average prediction}} = \underbrace{\frac{1}{n} \sum_i y_i}_{\text{average label}}$$

Making the Balancing Constraint Linear

hard / non-linear	$\underbrace{\frac{1}{m} \sum_j \text{sign}(\mathbf{w}^\top \mathbf{x}_j + b)}_{\text{average prediction}} = \underbrace{\frac{1}{n} \sum_i y_i}_{\text{average label}}$
soft / linear	$\underbrace{\frac{1}{m} \sum_j \mathbf{w}^\top \mathbf{x}_j + b}_{\text{mean output on unlabeled points}} = \underbrace{\frac{1}{n} \sum_i y_i}_{\text{average label}}$

Implementing the linear soft balancing:

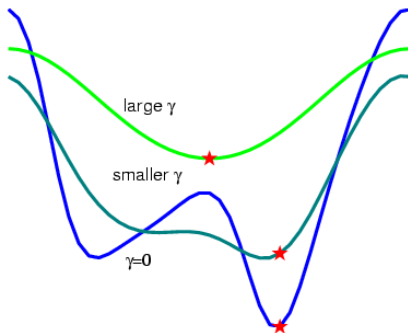
- center the unlabeled data: $\sum_j \mathbf{x}_j = \mathbf{0}$
- \Rightarrow just fix b ; unconstrained optimization over \mathbf{w} !

The Continuation Method in a Nutshell

Procedure

- 1 smooth function until convex
- 2 find minimum
- 3 track minimum while decreasing amount of smoothing

Illustration



Smoothing the S^3VM Objective $f(\cdot)$

Convolution of $f(\cdot)$ with Gaussian of width $\sqrt{\gamma/2}$:

$$f_{\gamma}(\mathbf{w}) = (\pi\gamma)^{-d/2} \int f(\mathbf{w} - \mathbf{t}) \exp(-\|\mathbf{t}\|^2/\gamma) d\mathbf{t}$$

Closed form solution!

Smoothing Sequence

choose $\gamma_0 > \gamma_1 > \dots > \gamma_{p-1} > \gamma_p = 0$

- choose γ_0 such that $f_{\gamma_0}(\cdot)$ is convex
- choose γ_{p-1} such that $f_{\gamma_{p-1}}(\cdot) \approx f_{\gamma_p}(\cdot) = f(\cdot)$
- $p = 10$ steps (equidistant on log scale) sufficient

Handling Non-Linearity

Consider non-linear map $\Phi(\mathbf{x})$, kernel $k(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^\top \Phi(\mathbf{x}_j)$.

Representer Theorem: S^3VM solution is in span E of data points

$$E := \text{span}\{\Phi(\mathbf{x}_i)\} \stackrel{\Delta}{=} \mathbb{R}^{n+m}$$

Implementation

- ① expand basis vectors \mathbf{v}_i of E :

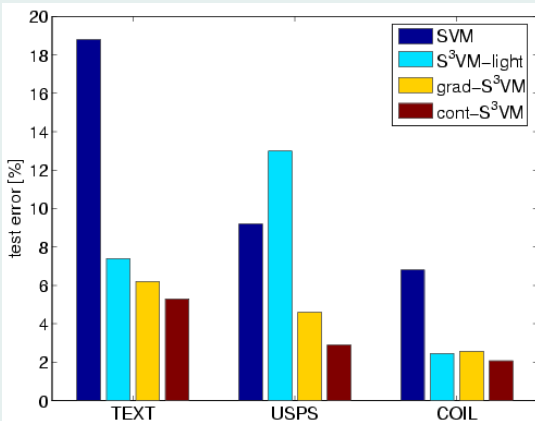
$$\mathbf{v}_i = \sum_k A_{ik} \Phi(\mathbf{x}_k)$$

- ② orthonormality gives:
solve for A , eg by KPCA or Choleski

$$(A^\top A)^{-1} = K$$

- ③ project data $\Phi(\mathbf{x}_i)$ on basis $V = (\mathbf{v}_j)_j$: $\tilde{\mathbf{x}}_i = V^\top \Phi(\mathbf{x}_i) = (A)_i$

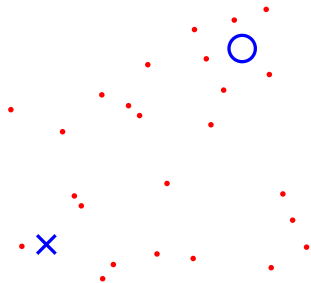
Comparison of S^3VM Optimization Methods



- averaged over splits (and pairs of classes)
- fixed hyperparams (close to hard margin)
- similar results for other hyperparameter settings

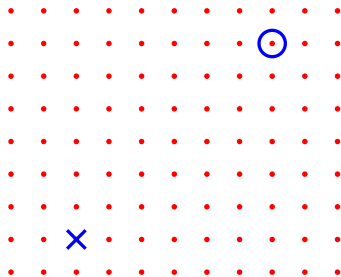
[Chapelle, Chi, Zien; ICML 2006]

Why would unlabeled data be useful at all?



Uniform data do not help.

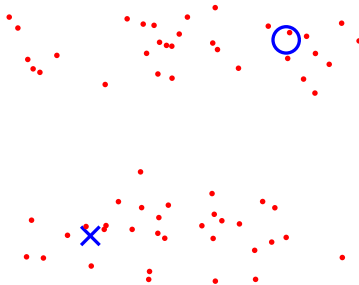
Why would unlabeled data be useful at all?



Uniform data do not help.

Cluster Assumption

Points in the **same cluster** are likely to be of the **same class**.



Algorithmic idea: **Low Density Separation**

Manifold Assumption

The data lie on (close to) a low-dimensional manifold.



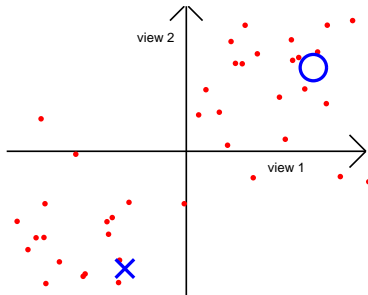
[images from "The Geometric Basis of Semi-Supervised Learning", Sindhwani, Belkin, Niyogi
in "Semi-Supervised Learning" Chapelle, Schölkopf, Zien]

Algorithmic idea: use **Nearest-Neighbor Graph**

Assumption: Independent Views Exist

There exist **subsets of features, called views**, each of which

- is **independent** of the others given the class;
- is **sufficient** for classification.



Algorithmic idea: **Co-Training**

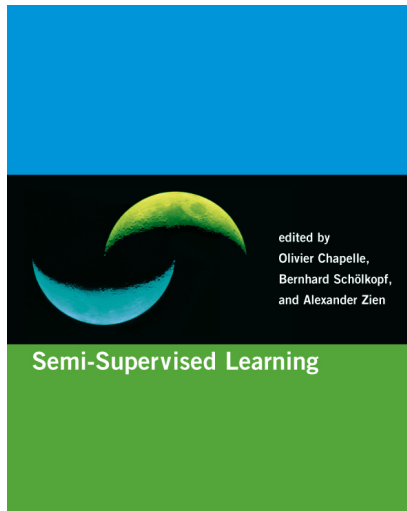
Assumption	Approach	Example Algorithm
Cluster Assumption	Low Density Separation	S^3VM ; Entropy Regularization; Data-Dependent Regularization; ...
Manifold Assumption	Graph-based Methods	<ul style="list-style-type: none"> • build weighted graph (w_{kl}) • $\min_{(y_j)} \sum_k \sum_l w_{kl} (y_k - y_l)^2$ • relax y_j to be real \Rightarrow QP
Independent Views	Co-Training	<ul style="list-style-type: none"> • train two predictors $y_j^{(1)}, y_j^{(2)}$ • couple objectives by adding $\sum_j (y_j^{(1)} - y_j^{(2)})^2$

Discriminative Learning (Diagnostic Paradigm)

- **model** $p(y|\mathbf{x})$ (or just boundary: $\{\mathbf{x} \mid p(y|\mathbf{x}) = \frac{1}{2}\}$)
- examples: S³VM, graph-based methods

Generative Learning (Sampling Paradigm)

- **model** $p(\mathbf{x}|y)$
- predict via Bayes: $p(y|\mathbf{x}) = \frac{p(y)p(\mathbf{x}|y)}{\sum_{y'} p(y')p(\mathbf{x}|y')}$
- \Rightarrow missing data problem
- EM algorithm (expectation-maximization) is a natural tool
- successful for text data [Nigam et al.; Machine Learning, 2000]



SSL Book

- MIT Press, Sept. 2006
- edited by B. Schölkopf, O. Chapelle, A. Zien
- contains many state-of-art algorithms by top researchers
- extensive SSL benchmark
- online material:
 - sample chapters
 - benchmark data
 - more information

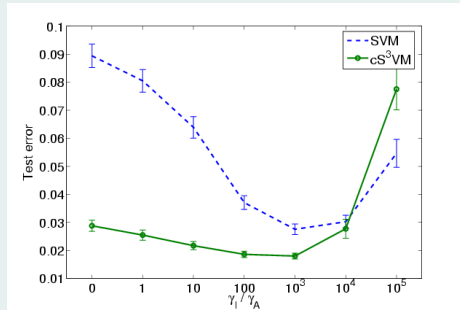
<http://www.kyb.tuebingen.mpg.de/ssl-book/>

SSL Book – Text Benchmark

	error [%]		AUC [%]	
	l=10	l=100	l=10	l=100
1-NN	38.12	30.11	–	–
SVM	45.37	26.45	67.97	84.26
MVU + 1-NN	45.32	32.83	–	–
LEM + 1-NN	39.44	30.77	–	–
QC + CMN	40.79	25.71	70.71	84.62
Discrete Reg.	40.37	24.00	53.79	71.53
TSVM	31.21	24.52	73.42	80.96
SGT	29.02	23.09	80.09	85.22
Cluster-Kernel	42.72	24.38	73.09	85.90
LDS	27.15	23.15	80.68	84.77
Laplacian RLS	33.68	23.57	76.55	85.05

Combining S^3VM with Graph-based Regularizer

- LapSVM [1]: modify kernel using graph, then train SVM
- combination with S^3VM even better [2]
- MNIST, “3” vs “5”



[1] “Beyond the Point Cloud”; Sindhwani, Niyogi, Belkin; ICML 2005

[2] “A Continuation Method for S^3VM ”; Chapelle, Chi, Zien; ICML 2006

Combining S^3VM with Co-Training

“SSL for Structured Output Variables”; Brefeld, Scheffer; ICML 2006

$$\min_{\mathbf{w}, b, (y_j), (\xi_k)} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j$$

$$s.t. \quad \begin{aligned} y_i (\mathbf{w}^\top \mathbf{x}_i + b) &\geq 1 - \xi_i & \xi_i &\geq 0 \\ y_j (\mathbf{w}^\top \mathbf{x}_j + b) &\geq 1 - \xi_j & \xi_j &\geq 0 \end{aligned}$$

How to set C ?

- data fitting, $y_i \mathbf{w}^\top \mathbf{x}_i \geq 1$, and regularization, $\min \|\mathbf{w}\|^2$:

$$|\mathbf{w}^\top \mathbf{x}_i| = \mathcal{O}(1) \quad \Rightarrow \quad \|\mathbf{w}\|^2 \approx \text{Var}[\mathbf{x}]^{-1}$$

- balance influence: $\|\mathbf{w}\|^2 \approx C \xi_i \quad \Rightarrow \quad C \approx \text{Var}[\mathbf{x}]^{-1}$

How to set C^* ?

- $C^* = C$
- $C^* = \lambda \frac{\# \text{ unlabeled points}}{\# \text{ labeled points}} C$

Naive Application:

- **Transductive setting** on each user/inbox:
 - use inbox of given user as unlabeled data
 - test data = unlabeled data
- **Guess the model:**
 - $Var[\mathbf{x}] \approx 1$, so set $C = 1$
 - $C^* = C$
 - linear kernel

Results: AUC (rank) [rank in unofficial list]

	task A	task B
S^3VM^{light}	94.53% (4) [6]	92.34% (2) [4]
∇S^3VM	96.72% (1) [3]	93.74% (2) [4]
cont S^3VM	96.01% (1) [3]	93.56% (2) [4]

Model selection:

- $C \in \{10^{-2}, 10^{-1}, 10^0, 10^{+1}, 10^{+2}\}$
- $C^* \in \{10^{-2}, 10^{-1}, 10^0, 10^{+1}, 10^{+2}\} \cdot C$
- cross-validation (3-fold for task A; 5-fold for task B)

Results: AUC for cont S^3VM

	task A	task B
$C = C^* = 1$ (guessed model)	96.01%	93.56%
model selection	89.31%	90.09%

- **significant drop in accuracy!**
- CV relies on **iid assumption**:
that the data are independent identically distributed

Take Home Messages

- S^3VM implements “**low density separation**” (margin maximization)
- optimization technique matters (non-convex objective)
- **works well for text** classification (texts form clusters)
- S^3VM -based **hybrids** may be even better
- for **spam filtering**, further methods needed to cope with **non-iid** situation (mail inboxes)!

Thank you!