

Efficient multiplicative updates for SVM

Sergey Plis

Computer Science Department
University of New Mexico

NIMH grant number 1R01MH076282-01

joint work with Vamsi Potluru, Morten Mørup,
Vince Calhoun, and Terran Lane

Outline

- 1 Introduction
 - Support Vector Machines
 - Non-negative Matrix Factorization
 - NQP
- 2 SVM as NMF
- 3 Experiments

Outline

- 1 Introduction
 - Support Vector Machines
 - Non-negative Matrix Factorization
 - NQP

- 2 SVM as NMF

- 3 Experiments

Maximum Margin Classifiers

- Given two classes A and B
Data is a set of labeled examples $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$

- Bayesian Decision
Boundary: *Distributions are known*

$$g(\mathbf{x}) = P(A|\mathbf{x}) - P(B|\mathbf{x}) = 0$$

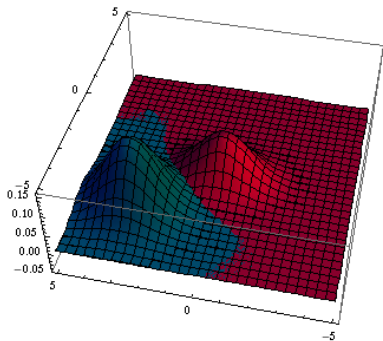
- Maximum Margin
Hyperplane: *Only data is given*

Maximum Margin Classifiers

- Given two classes A and B
Data is a set of labeled examples $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$
- Bayesian Decision
Boundary: *Distributions are known*

$$g(\mathbf{x}) = P(A|\mathbf{x}) - P(B|\mathbf{x}) = 0$$

- Maximum Margin
Hyperplane: *Only data is given*

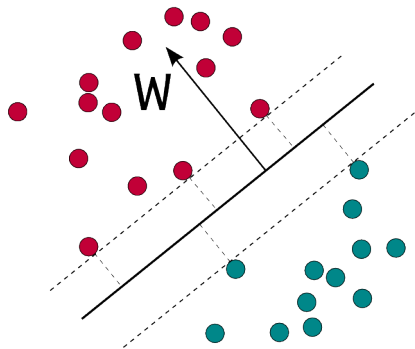


Maximum Margin Classifiers

- Given two classes A and B
Data is a set of labeled examples $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$
- Bayesian Decision
Boundary: *Distributions are known*

$$g(\mathbf{x}) = P(A|\mathbf{x}) - P(B|\mathbf{x}) = 0$$

- Maximum Margin
Hyperplane: *Only data is given*



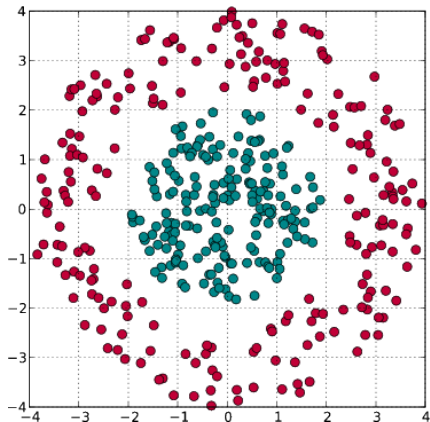
Mapping to Higher Dimensions

- A two class problem: **not separable by a line**
- Map each point into a higher dimensional space:

$$\eta_i = \Phi(\mathbf{x}_i)$$

- Choose Φ so the data becomes linearly separable

$$\Phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



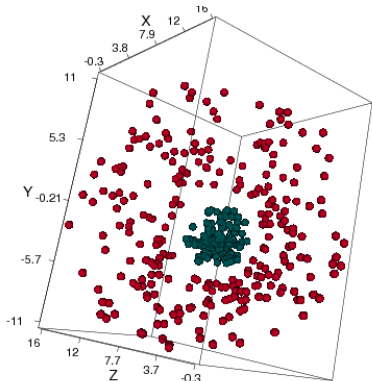
Mapping to Higher Dimensions

- A two class problem: **not separable by a line**
- Map each point into a higher dimensional space:

$$\boldsymbol{\eta}_i = \Phi(\mathbf{x}_i)$$

- Choose Φ so the data becomes linearly separable

$$\Phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



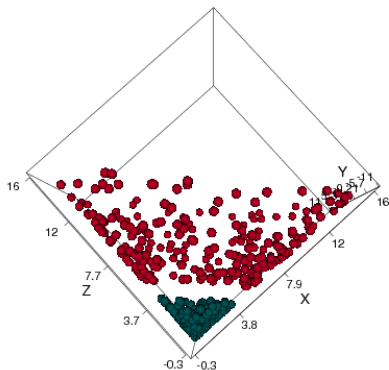
Mapping to Higher Dimensions

- A two class problem: **not separable by a line**
- Map each point into a higher dimensional space:

$$\eta_j = \Phi(\mathbf{x}_j)$$

- Choose Φ so the data becomes linearly separable

$$\Phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$



Kernel trick

- Dimensions can be many (even ∞)!. Have to compute very expensive inner product?
- Avoid it by defining Hilbert spaces with kernels

$$\begin{aligned} \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle &= k(\mathbf{x}, \mathbf{y}) \\ (x_1^2, \sqrt{2}x_1x_2, x_2^2)(y_1^2, \sqrt{2}y_1y_2, y_2^2)^T &= (\mathbf{x} \cdot \mathbf{y})^2 \\ x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 &= (x_1y_1 + x_2y_2)^2 \end{aligned}$$

- Take a linear algorithm, replace inner products with kernels and get a nonlinear algorithm as a result!

Kernel trick

- Dimensions can be many (even ∞)!. Have to compute very expensive inner product?
- Avoid it by defining Hilbert spaces with kernels

$$\begin{aligned} \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle &= \mathbf{k}(\mathbf{x}, \mathbf{y}) \\ (x_1^2, \sqrt{2}x_1x_2, x_2^2)(y_1^2, \sqrt{2}y_1y_2, y_2^2)^T &= (\mathbf{x} \cdot \mathbf{y})^2 \\ x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 &= (x_1y_1 + x_2y_2)^2 \end{aligned}$$

- Take a linear algorithm, replace inner products with kernels and get a nonlinear algorithm as a result!

Kernel trick

- Dimensions can be many (even ∞)!. Have to compute very expensive inner product?
- Avoid it by defining Hilbert spaces with kernels

$$\begin{aligned}\langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle &= k(\mathbf{x}, \mathbf{y}) \\ (x_1^2, \sqrt{2}x_1x_2, x_2^2)(y_1^2, \sqrt{2}y_1y_2, y_2^2)^T &= (\mathbf{x} \cdot \mathbf{y})^2 \\ x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 &= (x_1y_1 + x_2y_2)^2\end{aligned}$$

- Take a linear algorithm, replace inner products with kernels and get a nonlinear algorithm as a result!

Primal form

- Separating hyperplane

$$y = \text{sign}((\mathbf{w} \cdot \mathbf{x}) + b)$$

- Classification error

$$y_i((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b)$$

- Quadratic optimization problem

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $y_i((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b) \geq 1, i \in \{1..n\}$

Dual form

- Introduce Lagrange multipliers:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i ((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b) - 1)$$

- The dual quadratic optimization problem for SVM Schölkopf and Smola [2001] is given by minimizing the following loss function:

$$S(\alpha) = \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i$$

subject to $\alpha_i \geq 0, i \in \{1..n\}$,

Dual form

- Introduce Lagrange multipliers:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i ((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) + b) - 1)$$

- The dual quadratic optimization problem for SVM Schölkopf and Smola [2001] is given by minimizing the following loss function:

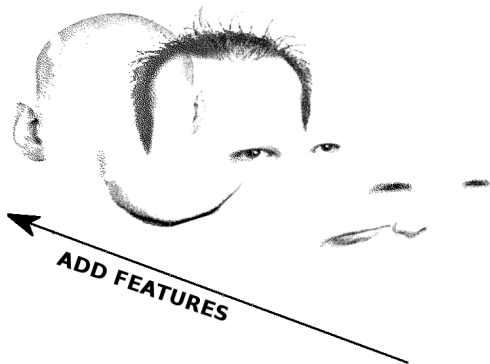
$$S(\alpha) = \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i$$

subject to $\alpha_i \geq 0, i \in \{1..n\},$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

Additive Features

- Features are non-negative and only add up
- Features are unknown: data comes as their combination



Additive Features

- Features are non-negative and only add up
- Features are unknown: data comes as their combination



Mathematical Formulation

- Given data \mathbf{X} find its factorization:

$$\mathbf{X} \approx \mathbf{WH}$$
$$\mathbf{X}_{ij} \geq 0 \quad \mathbf{W}_{ij} \geq 0 \quad \mathbf{H}_{ij} \geq 0$$

- Minimize the objective function:

$$E = \frac{1}{2} \|\mathbf{X} - \mathbf{WH}\|_F^2$$

- Ignore other possible objectives

Mathematical Formulation

- Given data \mathbf{X} find its factorization:

$$\mathbf{X} \approx \mathbf{WH}$$
$$\mathbf{X}_{ij} \geq 0 \quad \mathbf{W}_{ij} \geq 0 \quad \mathbf{H}_{ij} \geq 0$$

- Minimize the objective function:

$$E = \frac{1}{2} \|\mathbf{X} - \mathbf{WH}\|_F^2$$

- Ignore other possible objectives

Mathematical Formulation

- Given data \mathbf{X} find its factorization:

$$\mathbf{X} \approx \mathbf{WH}$$
$$\mathbf{X}_{ij} \geq 0 \quad \mathbf{W}_{ij} \geq 0 \quad \mathbf{H}_{ij} \geq 0$$

- Minimize the objective function:

$$E = \frac{1}{2} \|\mathbf{X} - \mathbf{WH}\|_F^2$$

- Ignore other possible objectives

Gradient Descent

- Compute the derivative and find its zero

$$\frac{\partial E}{\partial \mathbf{W}} = \mathbf{W}\mathbf{H}\mathbf{H}^T - \mathbf{X}\mathbf{H}^T$$

$$\frac{\partial E}{\partial \mathbf{H}} = \mathbf{W}^T\mathbf{W}\mathbf{H} - \mathbf{W}^T\mathbf{X}$$

- Classical solution

$$\mathbf{H} = \mathbf{H} + \eta \odot (\mathbf{W}^T\mathbf{X} - \mathbf{W}^T\mathbf{W}\mathbf{H})$$

- Exponentiated gradient

$$\mathbf{H} = \mathbf{H} \odot e^{\eta \odot (\mathbf{W}^T\mathbf{X} - \mathbf{W}^T\mathbf{W}\mathbf{H})}$$

Gradient Descent

- Compute the derivative and find its zero

$$\frac{\partial E}{\partial \mathbf{W}} = \mathbf{W}\mathbf{H}\mathbf{H}^T - \mathbf{X}\mathbf{H}^T$$

$$\frac{\partial E}{\partial \mathbf{H}} = \mathbf{W}^T\mathbf{W}\mathbf{H} - \mathbf{W}^T\mathbf{X}$$

- Classical solution

$$\mathbf{H} = \mathbf{H} + \eta \odot (\mathbf{W}^T\mathbf{X} - \mathbf{W}^T\mathbf{W}\mathbf{H})$$

- Exponentiated gradient

$$\mathbf{H} = \mathbf{H} \odot e^{\eta \odot (\mathbf{W}^T\mathbf{X} - \mathbf{W}^T\mathbf{W}\mathbf{H})}$$

Gradient Descent

- Compute the derivative and find its zero

$$\frac{\partial E}{\partial \mathbf{W}} = \mathbf{W}\mathbf{H}\mathbf{H}^T - \mathbf{X}\mathbf{H}^T$$

$$\frac{\partial E}{\partial \mathbf{H}} = \mathbf{W}^T\mathbf{W}\mathbf{H} - \mathbf{W}^T\mathbf{X}$$

- Classical solution

$$\mathbf{H} = \mathbf{H} + \eta \odot (\mathbf{W}^T\mathbf{X} - \mathbf{W}^T\mathbf{W}\mathbf{H})$$

- Exponentiated gradient

$$\mathbf{H} = \mathbf{H} \odot e^{\eta \odot (\mathbf{W}^T\mathbf{X} - \mathbf{W}^T\mathbf{W}\mathbf{H})}$$

Multiplicative Updates

- Setting the learning rate:

$$\eta = \frac{H}{W^T W H}$$

- Results in updates:

$$W = W \odot \frac{X H^T}{W H H^T},$$

$$H = H \odot \frac{W^T X}{W^T W H},$$

- Advantages:

- automatic non-negativity constraint satisfaction
- adaptive learning rate
- no parameter setting

Multiplicative Updates

- Setting the learning rate:

$$\eta = \frac{H}{W^T W H}$$

- Results in updates:

$$W = W \odot \frac{X H^T}{W H H^T},$$

$$H = H \odot \frac{W^T X}{W^T W H},$$

- Advantages:

- automatic non-negativity constraint satisfaction
- adaptive learning rate
- no parameter setting

Multiplicative Updates

- Setting the learning rate:

$$\eta = \frac{H}{W^T W H}$$

- Results in updates:

$$W = W \odot \frac{X H^T}{W H H^T},$$

$$H = H \odot \frac{W^T X}{W^T W H},$$

- Advantages:

- automatic non-negativity constraint satisfaction
- adaptive learning rate
- no parameter setting

Non-negative Quadratic Programming

$$F(\alpha) = \frac{1}{2}\alpha^T \mathbf{A}\alpha - \mathbf{1}^T \alpha,$$

subject to $\alpha_i \geq 0, i \in \{1..n\}$

Outline

- 1 Introduction
 - Support Vector Machines
 - Non-negative Matrix Factorization
 - NQP

- 2 SVM as NMF

- 3 Experiments

SVM as NMF-type problem

■ SVM dual formulation

$$S(\alpha) = \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i$$

subject to $\alpha_i \geq 0, i \in \{1..n\}$,

■ Looks like NMF:

$$\min_{\alpha} \frac{1}{2} \|\Phi(\mathbf{X}_A)\alpha_A - \Phi(\mathbf{X}_B)\alpha_B\|_2^2 - \sum_{i \in \{A,B\}} \alpha_i$$

subject to $\alpha_i \geq 0$,

■ NMF objective function:

$$E = \frac{1}{2} \|\mathbf{X} - \mathbf{WH}\|_F^2$$

SVM as NMF-type problem

- SVM dual formulation

$$S(\alpha) = \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i$$

subject to $\alpha_i \geq 0, i \in \{1..n\}$,

- Rewrite the square

$$\min_{\alpha} \frac{1}{2} \left(\sum_{ij \in A} \alpha_i \alpha_j k_{ij} - 2 \sum_{\substack{i \in B \\ j \in A}} \alpha_i \alpha_j k_{ij} + \sum_{ij \in B} \alpha_i \alpha_j k_{ij} \right) - \sum_{i=1}^n \alpha_i$$

- Looks like NMF:

$$\min_{\alpha} \frac{1}{2} \|\Phi(\mathbf{X}_A)\alpha_A - \Phi(\mathbf{X}_B)\alpha_B\|_2^2 - \sum_{i \in \{A,B\}} \alpha_i$$

subject to $\alpha_i \geq 0$,

- NMF objective function:

SVM as NMF-type problem

- Looks like NMF:

$$\min_{\alpha} \frac{1}{2} \|\Phi(\mathbf{X}_A)\alpha_A - \Phi(\mathbf{X}_B)\alpha_B\|_2^2 - \sum_{i \in \{A, B\}} \alpha_i$$

subject to $\alpha_i \geq 0$,

- NMF objective function:

$$E = \frac{1}{2} \|\mathbf{X} - \mathbf{WH}\|_F^2$$

SVM as NMF-type problem

- Looks like NMF:

$$\min_{\alpha} \frac{1}{2} \|\Phi(\mathbf{X}_A)\alpha_A - \Phi(\mathbf{X}_B)\alpha_B\|_2^2 - \sum_{i \in \{A, B\}} \alpha_i$$

subject to $\alpha_i \geq 0$,

- NMF objective function:

$$E = \frac{1}{2} \|\mathbf{X} - \mathbf{WH}\|_F^2$$

Multiplicative updates

- Differentiate the objective

$$\begin{aligned} \frac{\partial \mathcal{S}}{\partial \alpha_A} &= \langle \Phi(\mathbf{X}_A), \Phi(\mathbf{X}_A) \rangle \alpha_A - \langle \Phi(\mathbf{X}_A), \Phi(\mathbf{X}_B) \rangle \alpha_B - \mathbf{1} \\ &= K(\mathbf{X}_A, \mathbf{X}_A) \alpha_A - (K(\mathbf{X}_A, \mathbf{X}_B) \alpha_B + \mathbf{1}) \end{aligned}$$

- Simple multiplicative updates for SVM

$$\begin{aligned} \alpha_A &= \alpha_A \odot \frac{K(\mathbf{X}_A, \mathbf{X}_B) \alpha_B + \mathbf{1}}{K(\mathbf{X}_A, \mathbf{X}_A) \alpha_A} \\ \alpha_B &= \alpha_B \odot \frac{K(\mathbf{X}_B, \mathbf{X}_A) \alpha_A + \mathbf{1}}{K(\mathbf{X}_B, \mathbf{X}_B) \alpha_B}, \end{aligned}$$

Multiplicative updates (cont.)

■ Multiplicative Updates for Sign Insensitive Kernel SVM

$$\alpha_A = \alpha_A \odot \frac{\mathbf{K}_{AB}^+ \alpha_B + \mathbf{K}_A^- \alpha_A + \mathbf{1} + \mathbf{D}_A \alpha_A}{\mathbf{K}_A^+ \alpha_A + \mathbf{K}_{AB}^- \alpha_B + \mathbf{D}_A \alpha_A}$$

$$\alpha_B = \alpha_B \odot \frac{\mathbf{K}_{BA}^+ \alpha_A + \mathbf{K}_B^- \alpha_B + \mathbf{1} + \mathbf{D}_B \alpha_B}{\mathbf{K}_B^+ \alpha_B + \mathbf{K}_{BA}^- \alpha_A + \mathbf{D}_B \alpha_B}$$

■ semiNMF-type SVM updates

$$\alpha_A = \alpha_A \odot \sqrt{\frac{\mathbf{K}_{AB}^+ \alpha_B + \mathbf{K}_A^- \alpha_A + \mathbf{1}}{\mathbf{K}_A^+ \alpha_A + \mathbf{K}_{AB}^- \alpha_B}}$$

$$\alpha_B = \alpha_B \odot \sqrt{\frac{\mathbf{K}_{BA}^+ \alpha_A + \mathbf{K}_B^- \alpha_B + \mathbf{1}}{\mathbf{K}_B^+ \alpha_B + \mathbf{K}_{BA}^- \alpha_A}}$$

Sum Constraint and Box Constraint

- Soft Margin SVM (the box constraint)

$$\min_{\alpha} \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i$$

subject to $0 \leq \alpha_i \leq l, i \in \{1..n\}$.

- Bias

- Introduce λ and rewrite $\sum_i y_i \alpha_i = 0$ as:

$$\sum_{i \in A} \alpha_i = \lambda, \sum_{i \in B} \alpha_i = \lambda$$

- Redefine $\beta_k = \alpha_k / \lambda$ and solve resulting SVM for λ and β using multiplicative updates

Outline

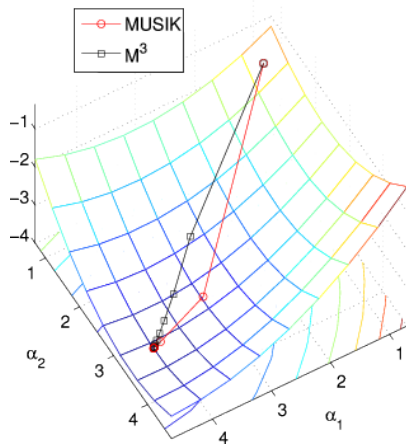
- 1 Introduction
 - Support Vector Machines
 - Non-negative Matrix Factorization
 - NQP

- 2 SVM as NMF

- 3 Experiments**

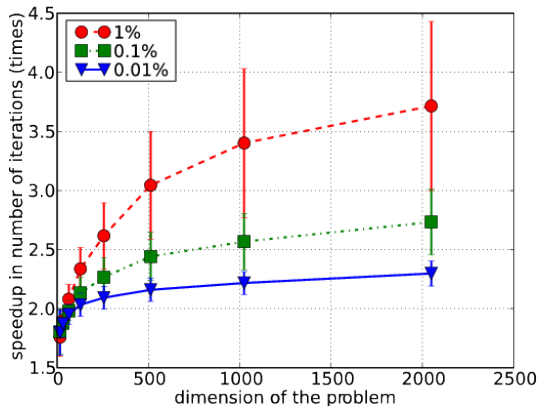
Simulations

- MUSIK takes bigger steps and follows a different path
- Converges faster within a given tolerance










Simulations

- MUSIK takes bigger steps and follows a different path
- Converges faster within a given tolerance



it works correctly

- Check on UCI dataset Newman and Merz [1998]
- Sonar and Breast cancer data
- Convergence is fast (breast cancer dataset)

i	support vectors	$\epsilon_t(\%)$	$\epsilon_g(\%)$
0		3.8	0.0
1		2.5	3.0
2		1.5	1.5
4		0.5	1.5
8		0.2	2.3
16		0.0	2.3
64		0.0	2.3

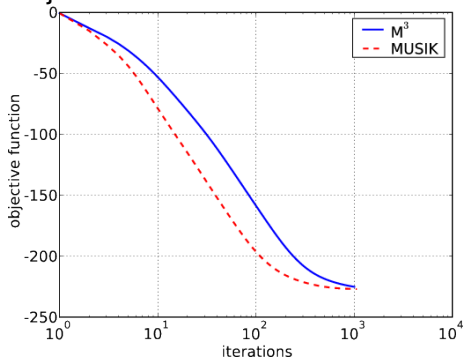
it works correctly (cont.)

Kernel		Breast			Sonar		
		M ³	M	KA	M ³	M	KA
Poly	4	2.26	2.26	2.26	9.62	9.62	9.62
	6	3.76	3.76	3.76	10.58	10.58	10.58
Gaussian	3	2.26	2.26	2.26	11.53	11.53	11.53
	1	0.75	0.75	0.75	7.69	7.69	7.69

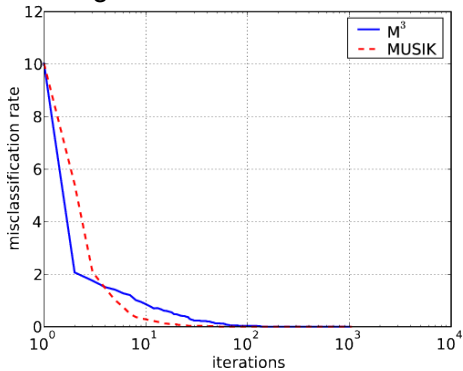
Converges to the exact same global solution

it works fast

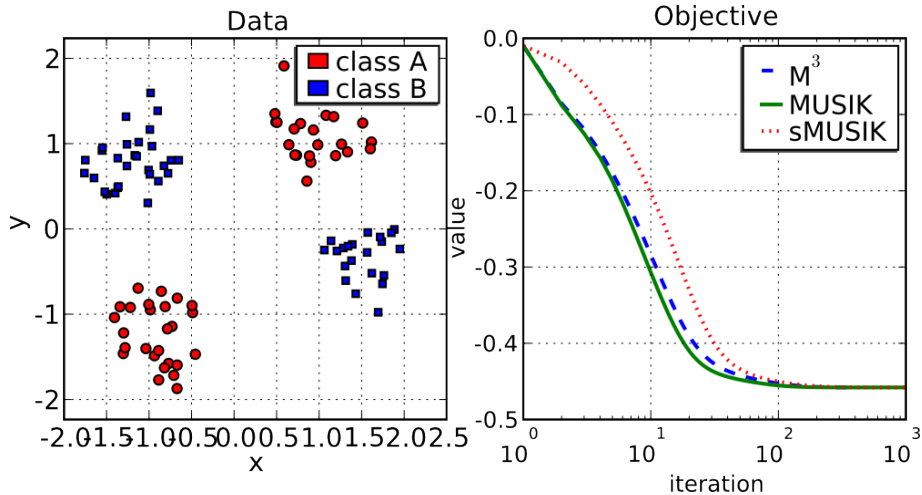
Objective



Training error



... and it is general



Conclusions

- Clean connection between SVM and NMF
- Fully multiplicative algorithm for SVM
- Simple to code algorithm: about 5 Matlab lines
- Speed Improvements
Theoretically (asymptotic convergence rates) and practically faster
- Possibility for algorithm reuse
 - SVM for NMF
 - NMF for SVM
- Further details in SDM 2009 paper

Thank you!

Bibliography

C. L. Blake D. J. Newman and C. J. Merz. UCI repository of machine learning databases, 1998. URL <http://www.ics.uci.edu/~lms/mllearn/MLRepository.html>.

<http://www.ics.uci.edu/~lms/mllearn/MLRepository.html>.

Bernhard Schölkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond (Adaptive Computation and Machine Learning)*. The MIT Press, 2001. ISBN 0262194759. URL

<http://www.amazon.ca/exec/obidos/redirect?tag=citeulike09-20&path=ASIN/0262194759>.