



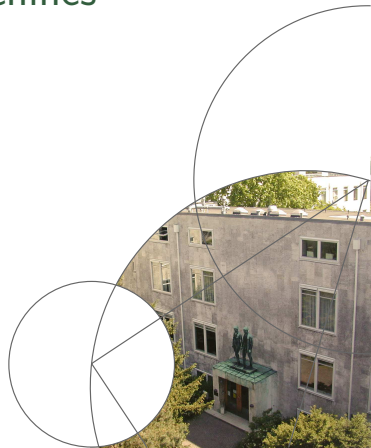
Faculty of Science



Training Support Vector Machines

Advanced Topics in Machine Learning

Christian Igel
Department of Computer Science



What's up?

- We consider training of *non-linear* support vector machines.
- We focus on batch learning using sequential minimal optimization (SMO), which can be adapted for online and active learning.
- For training linear SVMs, other methods – working directly on the primal optimization problem – are preferable.

Why?

- ...to deepen our understanding of SVMs, one of the most important classification methods, which helps using them successfully in practice.
- ...to learn general techniques from optimization that are widely applicable.



Outline

- ➊ Recall: Linear Classification
- ➋ Recall: Support Vector Machines
- ➌ Karush-Kuhn-Tucker Theorem
- ➍ SVM Dual Optimization Problems
- ➎ SMO for SVM Training
- ➏ Working Set Selection
- ➐ Second-order SMO for Online Learning



Outline

- 1 Recall: Linear Classification
- 2 Recall: Support Vector Machines
- 3 Karush-Kuhn-Tucker Theorem
- 4 SVM Dual Optimization Problems
- 5 SMO for SVM Training
- 6 Working Set Selection
- 7 Second-order SMO for Online Learning



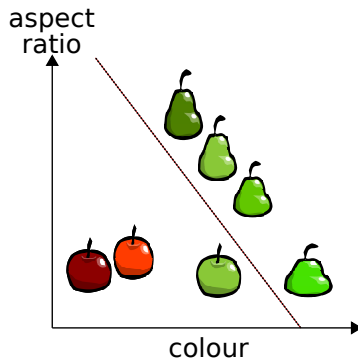
Affine linear decision functions I

Affine linear decision functions

$$f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$$

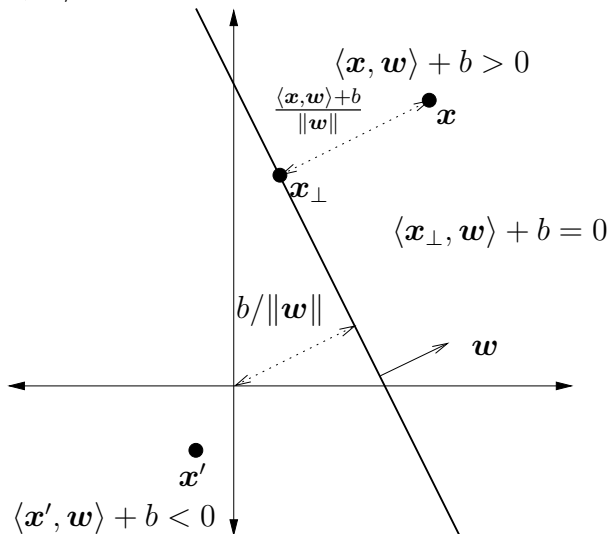
with $\mathcal{X} = \mathbb{R}^n$, $\mathbf{w}, \mathbf{x} \in \mathbb{R}^n$, $b \in \mathbb{R}$
lead to hypothesis with (affine)
linear decision boundaries:

$$\mathbf{x} \mapsto \begin{cases} 1 & \text{if } \langle \mathbf{w}, \mathbf{x} \rangle + b > 0 \\ -1 & \text{otherwise} \end{cases}$$



Linear decision functions II

$$f(\mathbf{x}) = \langle \mathbf{x}, \mathbf{w} \rangle + b$$



Margins

The *functional margin* of an example (\mathbf{x}_i, y_i) with respect to a hyperplane (\mathbf{w}, b) is

$$\gamma_i := y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \quad .$$

and its *geometric margin* is

$$\rho_i := y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) / \|\mathbf{w}\| = \gamma_i / \|\mathbf{w}\| \quad .$$

A positive margin implies correct classification.

The functional margin γ_S of a hyperplane (\mathbf{w}, b) with respect to a training set S is $\min_i \gamma_i$.



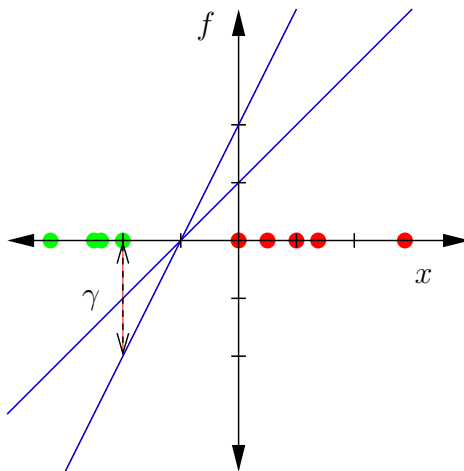
Outline

- 1 Recall: Linear Classification
- 2 Recall: Support Vector Machines
- 3 Karush-Kuhn-Tucker Theorem
- 4 SVM Dual Optimization Problems
- 5 SMO for SVM Training
- 6 Working Set Selection
- 7 Second-order SMO for Online Learning



“Inherent degree of freedom”

Inherent degree of freedom: (cw, cb) leads to same decision boundary for all $c \in \mathbb{R}^+$



Large margin classifier for separable data

Given linearly separable training data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell)\}$, we get rid of the inherent degree of freedom in

$$\begin{array}{ll} \text{maximize}_{\mathbf{w}, b} & \rho = \gamma / \|\mathbf{w}\| \\ \text{subject to} & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \gamma, \quad i = 1, \dots, \ell \end{array}$$

by fixing $\gamma = 1$

$$\begin{array}{ll} \text{maximize}_{\mathbf{w}, b} & \rho = 1 / \|\mathbf{w}\| \\ \text{subject to} & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \quad i = 1, \dots, \ell \end{array}$$

is equal to:

$$\begin{array}{ll} \text{minimize}_{\mathbf{w}, b} & \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle \\ \text{subject to} & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \quad i = 1, \dots, \ell \end{array}$$



Linear hard margin SVM primal

Given linearly separable data $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell)\}$ the hyperplane (\mathbf{w}, b) solving

$$\begin{aligned} \text{minimize}_{\mathbf{w}, b} \quad & \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle \\ \text{subject to} \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \quad i = 1, \dots, \ell \end{aligned}$$

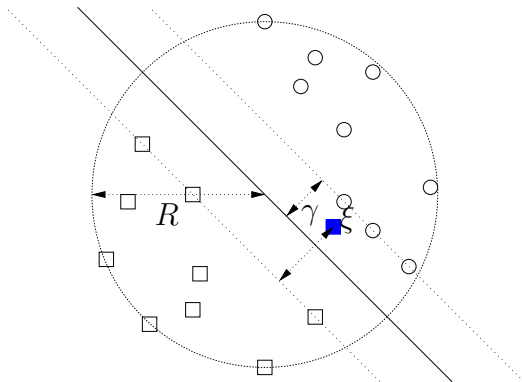
realizes the maximal margin hyperplane with margin $\rho = 1/\|\mathbf{w}\|$.



Tolerating margin violations: Slack variables

For fixed target margin $\gamma > 0$, the margin *slack variable* ξ_i of (\mathbf{x}_i, y_i) with respect to the hyperplane (\mathbf{w}, b) is

$$\xi((\mathbf{x}_i, y_i), (\mathbf{w}, b), \gamma) = \xi_i := \max(0, \gamma - y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b)) \ .$$



1-norm linear soft margin SVM primal

Penalizing the absolute values of the slack variables gives:

Given $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell)\} \in (\mathbb{R}^d \times \{-1, 1\})^\ell$ and a regularization parameter $C \geq 0$, a 1-norm linear soft margin SVM computes an affine linear decision function

$f(\mathbf{x}) = \langle \mathbf{x}, \mathbf{w} \rangle + b$ by solving:

$$\begin{aligned} \text{minimize}_{\xi, \mathbf{w}, b} \quad & \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle + C \sum_{i=1}^{\ell} \xi_i \\ \text{subject to} \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \quad i = 1, \dots, \ell \\ & \xi_i \geq 0, \quad i = 1, \dots, \ell \end{aligned}$$



Non-linear SVM

- In SVMs elements from \mathcal{X} occur only in scalar products – we can apply the “kernel trick”!
- Consider an arbitrary input space \mathcal{X} and a feature map $\Phi : \mathcal{X} \rightarrow \mathcal{H}_k$, where

$$\mathcal{H}_k = \text{span}\{k(x, \cdot) \mid x \in \mathcal{X}\}$$

is the RKHS induced by the (symmetric, positive-definite) kernel function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$.

- The non-linear SVM learns decision functions of the form

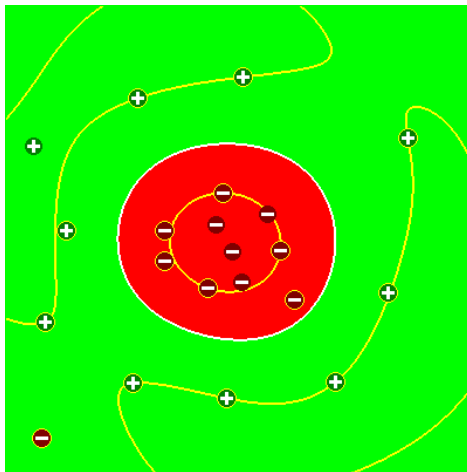
$$f(x) = \langle \Phi(x), \mathbf{w} \rangle + b .$$

Here $\mathbf{w} \in \mathcal{H}_k$ and $\Phi : x \mapsto k(x, \cdot)$.

- The scalar product $\langle \Phi(x), \mathbf{w} \rangle$ will be computed using the kernel trick.



Regularization and kernel representation



Kernel k : Represent data for linear classification (ideally, $h^{\text{Bayes}} \in \mathcal{H}_k^b$)
Slack variables: Deal with noise and outliers (i.e., $\mathcal{R}_p^{\text{Bayes}} > 0$)



1-norm non-linear soft margin SVM primal

Given $S = \{(x_1, y_1), \dots, (x_\ell, y_\ell)\} \in (\mathcal{X} \times \{-1, 1\})^\ell$ a regularization parameter $C \geq 0$, and a kernel k on \mathcal{X} , a 1-norm soft margin SVM computes a linear decision function $f(x) = \langle \Phi(x), \mathbf{w} \rangle + b$ by solving:

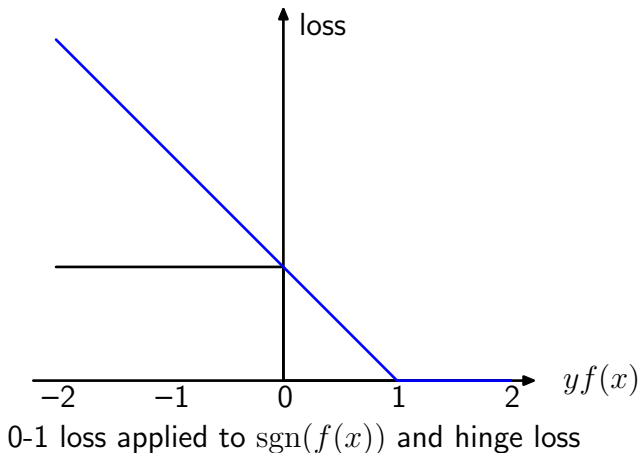
$$\begin{aligned} \text{minimize}_{\xi, \mathbf{w}, b} \quad & \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle + C \sum_{i=1}^{\ell} \xi_i \\ \text{subject to} \quad & y_i(\langle \mathbf{w}, \Phi(x_i) \rangle + b) \geq 1 - \xi_i, \quad i = 1, \dots, \ell \\ & \xi_i \geq 0, \quad i = 1, \dots, \ell, \end{aligned}$$

where $\Phi(x) = k(x, \cdot)$.



Hinge loss as convex surrogate for 0-1 loss

$$L_{\text{hinge}}(y, f(x)) = [1 - yf(x)]_+ = \max(0, 1 - yf(x))$$



1-norm soft margin SVM and regularization I

- 1-norm soft margin SVM, primal

$$\begin{aligned} \text{minimize}_{\boldsymbol{\xi}, \boldsymbol{w}, b} \quad & \frac{1}{2} \langle \boldsymbol{w}, \boldsymbol{w} \rangle + C \sum_{i=1}^{\ell} \xi_i \\ \text{subject to} \quad & y_i(\langle \boldsymbol{w}, \Phi(\boldsymbol{x}_i) \rangle + b) \geq 1 - \xi_i, \quad i = 1, \dots, \ell \\ & \xi_i \geq 0, \quad i = 1, \dots, \ell \end{aligned}$$

- For fixed \boldsymbol{w} optimal slack variables are:

$$\xi_i = \max(0, 1 - y_i(\langle \boldsymbol{w}, \Phi(\boldsymbol{x}_i) \rangle + b)) = L_{\text{hinge}}(y_i, f(\boldsymbol{x}_i))$$

- Hypothesis classes
 - \mathcal{H}_k : RKHS induced by k
 - $\mathcal{H}_k^b = \{f(x) = g(x) + b \mid g \in \mathcal{H}_k, b \in \mathbb{R}\}$



1-norm soft margin SVM and regularization II

- $L_{\text{hinge}}(y, \hat{y}) = \max(0, 1 - y\hat{y})$
- 1-norm soft margin SVM

$$\begin{aligned} \text{minimize}_{\xi, \mathbf{w}, b} \quad & \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle + C \sum_{i=1}^{\ell} \xi_i \\ \text{subject to} \quad & y_i(\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b) \geq 1 - \xi_i, \quad i = 1, \dots, \ell \\ & \xi_i \geq 0, \quad i = 1, \dots, \ell \end{aligned}$$

corresponds to

$$\text{minimize}_{f \in \mathcal{H}_k^b} \quad \frac{1}{\ell} \sum_{i=1}^{\ell} L_{\text{hinge}}(y_i, f(x_i)) + \gamma_{\ell} \|f\|_k^2$$

where $\gamma_{\ell} = (2\ell C)^{-1}$ and $\|\cdot\|_k$ inherited from \mathcal{H}_k to \mathcal{H}_k^b is only a semi-norm



Representer theorem

Let $\Omega : [0, \infty[\rightarrow \mathbb{R}$ be a strictly monotonic increasing function, \mathcal{H} a RKHS with kernel k on \mathcal{X} and L a loss function. Given $S = \{(x_1, y_1), \dots, (x_\ell, y_\ell)\} \subset (\mathcal{X} \times \mathbb{R})^\ell$, each minimizer $f \in \mathcal{H}^b$ of the regularized empirical risk

$$\sum_{i=1}^{\ell} L(y_i, f(x_i)) + \Omega(\|f\|_k^2)$$

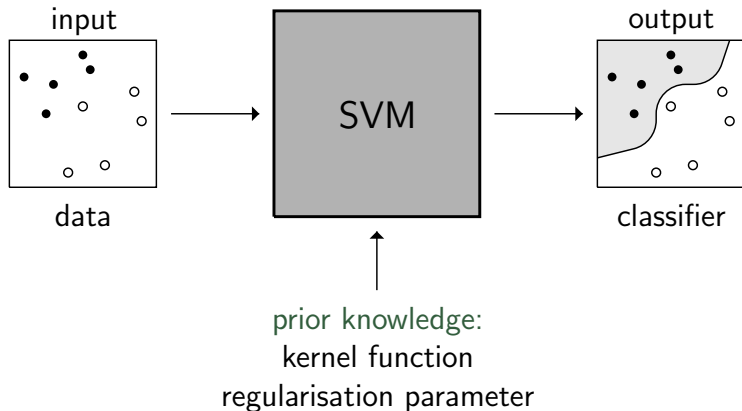
admits a representation of the form

$$f(x) = \sum_{i=1}^{\ell} \beta_i k(x_i, x) + b$$

with $\alpha_1, \dots, \alpha_\ell, b \in \mathbb{R}$.



Binary SVMs



Cortes, Vapnik: Support-Vector Networks, *Machine Learning* 20(3):273–297, 1995



Outline

- 1 Recall: Linear Classification
- 2 Recall: Support Vector Machines
- 3 Karush-Kuhn-Tucker Theorem**
- 4 SVM Dual Optimization Problems
- 5 SMO for SVM Training
- 6 Working Set Selection
- 7 Second-order SMO for Online Learning



Lagrangian Function

Given real-valued functions f, g_i, h_j ($i = 1, \dots, k$ and $j = 1, \dots, m$) defined on $M \subset \mathbb{R}^n$, we consider (*primal*) *optimization problems* of the form

$$\begin{array}{ll} \text{minimize} & f(\mathbf{w}) \quad \mathbf{w} \in M \\ \text{subject to} & g_i(\mathbf{w}) \leq 0 \quad i = 1, \dots, k \\ & h_j(\mathbf{w}) = 0 \quad j = 1, \dots, m \end{array} .$$

For a given primal optimization problem, the *Lagrangian function* is defined as

$$L(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{w}) + \sum_{i=1}^m \lambda_i h_i(\mathbf{w}) + \sum_{i=1}^k \mu_i g_i(\mathbf{w}) .$$



Karush-Kuhn-Tucker (KKT) theorem

Given a minimization problem with differentiable convex functions f, g_i , and affine functions h_j ($i = 1, \dots, k$ and $j = 1, \dots, m$) on a convex set $M \subset \mathbb{R}^n$. Suppose the differentials $dh_i(\mathbf{w}^*)$ and $dg_i(\mathbf{w}^*)$ are linearly independent at $\mathbf{w}^* \in M$. Then \mathbf{w}^* is a global minimum iff there exist $\boldsymbol{\lambda} \in \mathbb{R}^m$ and $\boldsymbol{\mu} \in \mathbb{R}^k$ with

$$\frac{\partial L(\mathbf{w}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)}{\partial \mathbf{w}} = \mathbf{0} ,$$

$$\frac{\partial L(\mathbf{w}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)}{\partial \boldsymbol{\lambda}} = \mathbf{0} ,$$

$$\mu_i^* g_i(\mathbf{w}^*) = 0 \text{ for } i = 1, \dots, k ,$$

$$g_i(\mathbf{w}^*) \leq 0 \text{ for } i = 1, \dots, k ,$$

$$\mu_i^* \geq 0 \text{ for } i = 1, \dots, k .$$



KKT Complementarity Condition

The conditions

$$\mu_i^* g_i(\mathbf{w}^*) = 0 \text{ for } i = 1, \dots, k ,$$

$$g_i(\mathbf{w}^*) \leq 0 \text{ for } i = 1, \dots, k ,$$

$$\mu_i^* \geq 0 \text{ for } i = 1, \dots, k$$

are called Karush-Kuhn-Tucker complementarity conditions.



Dual optimization problem

Given a primal optimization problem with Lagrangian function L , the *dual optimization problem* is defined by

$$\begin{array}{ll} \text{maximize} & \vartheta(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\boldsymbol{w} \in M} L(\boldsymbol{w}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \quad \boldsymbol{\lambda} \in \mathbb{R}^m, \boldsymbol{\mu} \in \mathbb{R}^k \\ \text{subject to} & \mu_i \geq 0 \quad i = 1, \dots, k. \end{array}$$

Under the assumptions of the KKT theorem, the *duality gap*

$$f(\boldsymbol{w}^*) - \sup_{\boldsymbol{\lambda} \in \mathbb{R}^m, \mu_1, \dots, \mu_k \in \mathbb{R}_0^+} \left[\inf_{\boldsymbol{w} \in M} L(\boldsymbol{w}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \right]$$

is zero.



Solution strategy

Given a constraint minimization problem, we maximize the dual, where we get rid of the inner minimization task using the following procedure:

- 1 Set derivatives of Lagrangian w.r.t. primal variables to zero;
- 2 Solve analytically w.r.t. primal variables;
- 3 Substitute primal variables into Lagrangian;
- 4 Maximize Lagrangian with w.r.t. dual variables.



Outline

- 1 Recall: Linear Classification
- 2 Recall: Support Vector Machines
- 3 Karush-Kuhn-Tucker Theorem
- 4 SVM Dual Optimization Problems**
- 5 SMO for SVM Training
- 6 Working Set Selection
- 7 Second-order SMO for Online Learning



Linear hard margin SVM: Primal to dual

Primal form:

$$\begin{aligned} \text{minimize}_{\mathbf{w}, b} \quad & \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle \\ \text{subject to} \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \quad i = 1, \dots, \ell \end{aligned}$$

Dual form:

$$\begin{aligned} \text{maximize}_{\boldsymbol{\alpha}} \quad & \inf_{\mathbf{w}, b} L(\mathbf{w}, b, \boldsymbol{\alpha}) \\ \text{subject to} \quad & \alpha_i \geq 0, \quad i = 1, \dots, \ell \end{aligned}$$

with Lagrangian:

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^{\ell} \alpha_i [y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1]$$



Eliminating primal variables I

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^{\ell} \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1]$$

Karush-Kuhn-Tucker (KKT) theorem requires

$$\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \boldsymbol{\alpha}) = 0 \quad \frac{\partial}{\partial b} L(\mathbf{w}, b, \boldsymbol{\alpha}) = 0$$

yielding

$$\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \boldsymbol{\alpha}) = \mathbf{w} - \sum_{i=1}^{\ell} \alpha_i y_i \mathbf{x}_i \quad \text{and} \quad \frac{\partial}{\partial b} L(\mathbf{w}, b, \boldsymbol{\alpha}) = \sum_{i=1}^{\ell} \alpha_i y_i$$

implying $\mathbf{w} = \sum_{i=1}^{\ell} \alpha_i y_i \mathbf{x}_i$ and $0 = \sum_{i=1}^{\ell} \alpha_i y_i$.



Eliminating primal variables II

Using $\mathbf{w} = \sum_{i=1}^{\ell} \alpha_i y_i \mathbf{x}_i$ gives

$$\begin{aligned} L(\mathbf{w}, b, \boldsymbol{\alpha}) &= \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle - \sum_{i=1}^{\ell} \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1] \\ &= \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^{\ell} \alpha_i \\ &= \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle . \end{aligned}$$



Linear hard margin SVM dual

For linearly separable $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell)\}$ solving

$$\text{maximize}_{\alpha} \quad \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{subject to} \quad \sum_{i=1}^{\ell} \alpha_i y_i = 0$$

$$\alpha_i \geq 0, \quad i = 1, \dots, \ell$$

leads to the maximal margin hyperplane with margin $\rho = 1/\|\mathbf{w}^*\|$ with

$$\mathbf{w}^* = \sum_{i=1}^{\ell} \alpha_i y_i \mathbf{x}_i.$$



KKT complementarity condition I

- KKT complementarity condition requires

$$\alpha_i^*[y_i(\langle \mathbf{w}^*, \mathbf{x}_i \rangle + b^*) - 1] = 0$$

for $i = 1, \dots, \ell$.

- KKT condition can be used to compute b^* , e.g.:

$$b^* = -\frac{\max_{y_i=-1}(\langle \mathbf{w}^*, \mathbf{x}_i \rangle) + \min_{y_i=1}(\langle \mathbf{w}^*, \mathbf{x}_i \rangle)}{2}$$

- Solution is sparse:

$$SV = \{\mathbf{x}_i \mid \alpha_i^* \neq 0\}$$

$$f(\mathbf{x}, \boldsymbol{\alpha}^*, b^*) = \sum_{i=1}^{\ell} y_i \alpha_i^* \langle \mathbf{x}_i, \mathbf{x} \rangle + b^* = \sum_{\mathbf{x}_i \in SV} y_i \alpha_i^* \langle \mathbf{x}_i, \mathbf{x} \rangle + b^*$$



KKT complementarity condition II

For $\mathbf{x}_j \in \text{SV}$

$$y_j f(\mathbf{x}_j, \boldsymbol{\alpha}^*, b^*) = y_j \left(\sum_{\mathbf{x}_i \in \text{SV}} y_i \alpha_i^* \langle \mathbf{x}_i, \mathbf{x}_j \rangle + b^* \right) = 1$$

and therefore

$$\begin{aligned} \langle \mathbf{w}^*, \mathbf{w}^* \rangle &= \sum_{i,j=1}^{\ell} y_i y_j \alpha_i^* \alpha_j^* \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ &= \sum_{\mathbf{x}_j \in \text{SV}} \alpha_j^* y_j \sum_{\mathbf{x}_i \in \text{SV}} \alpha_i^* y_i \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ &= \sum_{\mathbf{x}_j \in \text{SV}} \alpha_j^* (1 - y_j b^*) = \sum_{\mathbf{x}_j \in \text{SV}} \alpha_j^* . \end{aligned}$$



Hard margin SVM dual

For data $\{(x_1, y_1), \dots, (x_\ell, y_\ell)\}$ linearly separable in the feature space \mathcal{H}_k^b defined by the kernel k the solution α^*, b^* of

$$\text{maximize}_{\alpha} \quad \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

$$\text{subject to} \quad \sum_{i=1}^{\ell} \alpha_i y_i = 0, \quad \alpha_i \geq 0, \quad i = 1, \dots, \ell$$

leads to the decision rule $\text{sgn} \left(\sum_{i=1}^{\ell} y_i \alpha_i^* k(x_i, x) + b^* \right)$ corresponding to the maximal margin hyperplane in \mathcal{H}_k^b with margin $\rho = 1/\|\mathbf{w}^*\| = 1/\sqrt{\sum_{x_j \in \text{SV}} \alpha_j^*}$.



1-norm soft margin SVM: Primal to dual

Primal (linear case for ease of notation)

$$\begin{aligned} \text{minimize}_{\boldsymbol{\xi}, \boldsymbol{w}, b} \quad & \frac{1}{2} \langle \boldsymbol{w}, \boldsymbol{w} \rangle + C \sum_{i=1}^{\ell} \xi_i \\ \text{subject to} \quad & y_i(\langle \boldsymbol{w}, \boldsymbol{x}_i \rangle + b) \geq 1 - \xi_i, \quad i = 1, \dots, \ell \\ & \xi_i \geq 0, \quad i = 1, \dots, \ell \end{aligned}$$

has Lagrangian under constraints $\alpha_i, \beta_i \geq 0, i = 1, \dots, \ell$:

$$\begin{aligned} L(\boldsymbol{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \\ \frac{1}{2} \|\boldsymbol{w}\|^2 + C \sum_{i=1}^{\ell} \xi_i - \sum_{i=1}^{\ell} \alpha_i [y_i(\langle \boldsymbol{w}, \boldsymbol{x}_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^{\ell} \beta_i \xi_i \end{aligned}$$



Eliminating primal variables III

Lagrangian under constraints $\alpha_i, \beta_i \geq 0, i = 1, \dots, \ell$:

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^{\ell} \xi_i - \sum_{i=1}^{\ell} \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^{\ell} \beta_i \xi_i$$

KKT:

$$\frac{\partial L}{\partial \mathbf{w}}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \mathbf{w} - \sum_{i=1}^{\ell} \alpha_i y_i \mathbf{x}_i = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^{\ell} \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L}{\partial \xi_i}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = C - \alpha_i - \beta_i = 0, \quad i = 1, \dots, \ell$$

$$\frac{\partial L}{\partial b}(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \sum_{i=1}^{\ell} \alpha_i y_i = 0$$



Eliminating primal variables IV

$$\begin{aligned}
 & L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \\
 &= \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^{\ell} \xi_i - \sum_{i=1}^{\ell} \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^{\ell} \beta_i \xi_i \\
 &= \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^{\ell} [C \xi_i - \alpha_i \xi_i - \beta_i \xi_i] - \sum_{i=1}^{\ell} \alpha_i y_i \langle \mathbf{w}, \mathbf{x}_i \rangle - \sum_{i=1}^{\ell} \alpha_i y_i b + \sum_{i=1}^{\ell} \alpha_i \\
 &= \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^{\ell} [C - \alpha_i - \beta_i] \xi_i - \|\mathbf{w}\|^2 - b \sum_{i=1}^{\ell} \alpha_i y_i + \sum_{i=1}^{\ell} \alpha_i \\
 &= \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle
 \end{aligned}$$



Box constraint

KKT complementarity conditions:

$$\alpha_i[y_i(\langle \mathbf{w}, x_i \rangle + b) - 1 + \xi_i] = 0, \quad i = 1, \dots, \ell$$

$$\beta_i \xi_i = 0, \quad i = 1, \dots, \ell$$

$$\xi_i(\alpha_i - C) = 0, \quad i = 1, \dots, \ell \quad (\text{using } \beta_i = C - \alpha_i \text{ and } \beta_i \xi_i = 0)$$

- Bounded SVs: $\xi \neq 0 \Rightarrow \alpha_i = C$
- Box constraint: $\alpha_i, \beta_i \geq 0$ and $\beta_i = C - \alpha_i$ imply $0 \leq \alpha_i \leq C$
- Free SVs: $0 < \alpha_i < C \Rightarrow \xi_i = 0$
- How to get b : $0 < \alpha_i < C \wedge \xi_i = 0 \Rightarrow y_i(\langle \mathbf{w}, x_i \rangle + b) = 1$



1-norm soft margin SVM

For $\{(x_1, y_1), \dots, (x_\ell, y_\ell)\}$ and kernel k solving

$$\text{maximize}_{\alpha} \quad \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

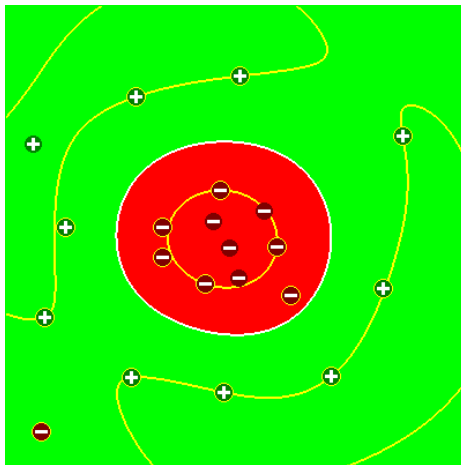
$$\text{subject to} \quad \sum_{i=1}^{\ell} \alpha_i y_i = 0, \quad C \geq \alpha_i \geq 0, \quad i = 1, \dots, \ell$$

leads to the decision rule $h(x) = \text{sgn}(f(x))$ with

$f(x) = \sum_{i=1}^{\ell} y_i \alpha_i^* k(x_i, x) + b^*$, where b^* is chosen so that $y_i f(x_i) = 1$ for any i with $C > \alpha_i > 0$ and the slack variables of the “corresponding hyperplane” in \mathcal{H}_k^b are defined relative to the margin $\rho = 1/\|\mathbf{w}^*\| = 1/\sqrt{\sum_{x_i, x_j \in \text{SV}} y_i y_j \alpha_i^* \alpha_j^* k(x_i, x_j)}$.



Inspecting the solution



Bounded SV: $\alpha_i = C$, $\xi_i \geq 0$, $y_i f(x_i) \leq 1$

Free SV: $0 < \alpha_i < C$, $\xi_i = 0$, $y_i f(x_i) = 1$

Non-SV: $\alpha_i = 0$, $\xi_i = 0$, $y_i f(x_i) > 1$



Outline

- 1 Recall: Linear Classification
- 2 Recall: Support Vector Machines
- 3 Karush-Kuhn-Tucker Theorem
- 4 SVM Dual Optimization Problems
- 5 SMO for SVM Training**
- 6 Working Set Selection
- 7 Second-order SMO for Online Learning



1-norm soft margin SVM: Dual

1-norm Soft Margin SVM Dual optimization problem:

$$\text{maximize}_{\alpha} \quad \mathcal{D}(\alpha) = \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

$$\text{subject to} \quad \sum_{i=1}^{\ell} \alpha_i y_i = 0 \quad \text{and} \quad \forall i \in \{1, \dots, \ell\} :$$

$$C \geq \alpha_i \geq 0, y_i \alpha_i \in [a_i, b_i] = \begin{cases} [0, C] & \text{if } y_i = +1 \\ [-C, 0] & \text{if } y_i = -1 \end{cases}$$

Decision rule $\text{sgn}(d(x))$ with $d(x) = \sum_{i=1}^{\ell} y_i \alpha_i^* k(x_i, x) + b^*$,
where b^* is chosen so that $y_i d(x_i) = 1$ for any i with
 $C > \alpha_i > 0$.



Restricting the quadratic program

1-norm Soft Margin SVM dual optimization problem restricted to *working set* B :

$$\text{maximize}_{\hat{\alpha}} \quad \mathcal{D}(\hat{\alpha}) = \sum_{i=1}^{\ell} \hat{\alpha}_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \hat{\alpha}_i \hat{\alpha}_j y_i y_j k(x_i, x_j)$$

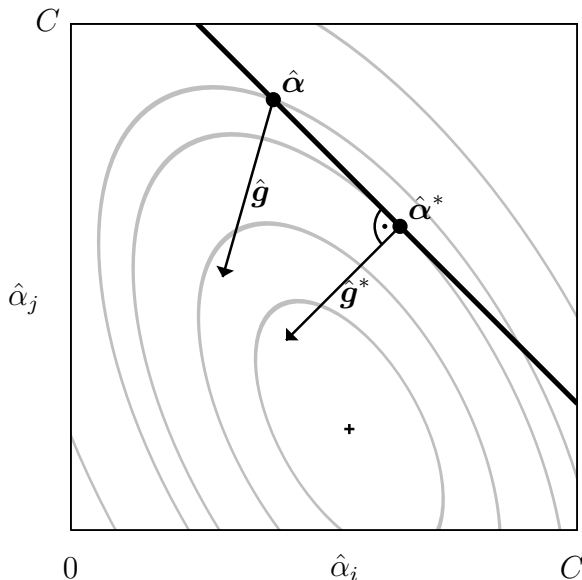
$$\text{subject to} \quad \sum_{i=1}^{\ell} \hat{\alpha}_i y_i = 0$$

$$\forall i \in \{1, \dots, \ell\} : C \geq \hat{\alpha}_i \geq 0$$

$$\forall i \notin B : \hat{\alpha}_i = 0$$



Two-dimensional subproblem



Some definitions

Let's define Gram matrix entry

$$K_{ij} = k(x_i, x_j)$$

gradient of dual

$$g_i = \frac{\partial \mathcal{D}(\boldsymbol{\alpha})}{\partial \alpha_i} = 1 - y_i \sum_{j=1}^{\ell} y_j \alpha_j K_{ij}$$

and the index sets:

$$I_{\text{up}} = \{i \mid y_i \alpha_i < b_i\} \quad (y_i \alpha_i \text{ may increase})$$

$$I_{\text{down}} = \{i \mid y_i \alpha_i > a_i\} \quad (y_i \alpha_i \text{ may decrease})$$



Decomposition algorithms

Strategy: Iteratively solve dual optimization problem

Decomposition Algorithm

- 1 $\alpha \leftarrow$ feasible starting point
 - 2 **repeat**
 - 3 select working set B
 - 4 solve QP restricted to B resulting in $\hat{\alpha}$
 - 5 $\alpha \leftarrow \hat{\alpha}$
 - 6 **until** *stopping criterion is met*
-

$B = \{i, j\}$, $i < j$: Sequential minimal optimization (SMO)
search directions are just

$$\pm(0, \dots, y_i, 0, \dots, 0, -y_j, 0, \dots, 0) = \pm \mathbf{u}_{ij}$$

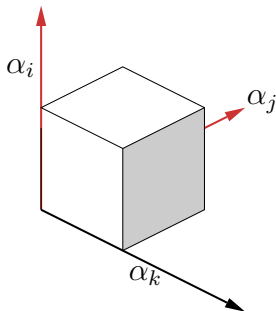


Sequential minimal optimization

repeat until target accuracy reached

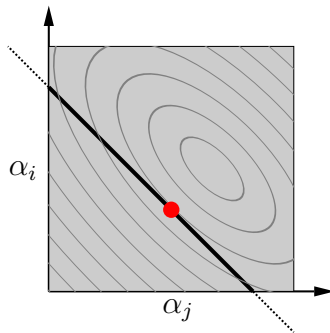
{

select working set



}

solve optimally



Solving the two-dimensional subproblem

Hessian of dual problem has elements:

$$\frac{\partial^2 \mathcal{D}(\boldsymbol{\alpha})}{\partial \alpha_i \partial \alpha_j} = y_i y_j K_{ij}$$

Maximizing w.r.t. λ (ignoring box constraints)

$$\mathcal{D}(\boldsymbol{\alpha} + \lambda \mathbf{u}_{ij}) - \mathcal{D}(\boldsymbol{\alpha}) = \lambda(y_i g_i - y_j g_j) - \frac{\lambda^2}{2}(K_{ii} + K_{jj} - 2K_{ij})$$

by Newton step gives optimal λ^* :

$$\lambda^* = \frac{y_i g_i - y_j g_j}{K_{ii} + K_{jj} - 2K_{ij}}$$



Recomputing gradient & stopping criterion

Gradient of full problem can be adjusted after optimizing on B by:

$$\forall k \in \{1, \dots, \ell\} : g_k \leftarrow g_k - y_k \sum_{i \in B} y_i (\hat{\alpha}_i - \alpha_i) K_{ik}$$

Stopping criterion (which needs gradient)

$$\max_{i \in I_{\text{up}}} y_i g_i - \min_{j \in I_{\text{down}}} y_j g_j \leq \epsilon$$

for $\epsilon > 0$.



Sequential minimal optimization

Sequential minimal optimization

- 1 $\alpha \leftarrow 0, g \leftarrow 1$
 - 2 **repeat**
 - 3 select indices $i \in I_{\text{up}}$ and $j \in I_{\text{down}}$
 - 4 $\lambda = \min \left\{ b_i - y_i \alpha_i, y_j \alpha_j - a_j, \frac{y_i g_i - y_j g_j}{K_{ii} + K_{jj} - 2K_{ij}} \right\}$
 - 5 $\forall k \in \{1, \dots, \ell\} : g_k \leftarrow g_k - \lambda y_k K_{ik} + \lambda y_k K_{jk}$
 - 6 $\alpha_i \leftarrow \alpha_i + y_i \lambda$
 - 7 $\alpha_j \leftarrow \alpha_j - y_j \lambda$
 - 8 **until** $\max_{i \in I_{\text{up}}} y_i g_i - \min_{j \in I_{\text{down}}} y_j g_j \leq \epsilon$
-



How long does training an SVM take?

Intuitive bounds:

- Assume an oracle tells us the unbounded SVs $F = \{x_i \mid 0 < \alpha_i < C\}$ and bounded SVs, then computing the α s takes $\mathcal{O}(|F|^3)$.
- Checking the optimality condition by computing the gradient from scratch takes $\mathcal{O}(\ell \cdot \#SV)$.

Batch SVM training scales between quadratically and cubically in the number of training points.



Outline

- 1 Recall: Linear Classification
- 2 Recall: Support Vector Machines
- 3 Karush-Kuhn-Tucker Theorem
- 4 SVM Dual Optimization Problems
- 5 SMO for SVM Training
- 6 Working Set Selection**
- 7 Second-order SMO for Online Learning



Working set selection: Most violating pair

Problem: How to select the working set B such that

- 1 much progress is made / only few iterations are needed, and
- 2 few kernel evaluations are required?

We ignore Gram matrix caching / “chunking” / “shrinking” in this course and just consider the selection of B .

Standard algorithm: *Most violating pair* working set selection.

- 1 first index $i = \operatorname{argmax}_{k \in I_{\text{up}}} y_k g_k$
 - 2 second index $j = \operatorname{argmin}_{k \in I_{\text{down}}} y_k g_k$
- *first order working set selection*, recall that for $\lambda \rightarrow 0$

$$\mathcal{D}(\boldsymbol{\alpha} + \lambda \mathbf{u}_{ij}) - \mathcal{D}(\boldsymbol{\alpha}) = \lambda(y_i g_i - y_j g_j) + \mathcal{O}(\lambda^2)$$

- requires just $\mathcal{O}(\ell)$ computations



Maximum gain

Maximizing gain of subproblem in search direction \mathbf{u}_{ij} ignoring box constraints corresponds to maximizing w.r.t. λ :

$$\mathcal{D}(\boldsymbol{\alpha} + \lambda \mathbf{u}_{ij}) - \mathcal{D}(\boldsymbol{\alpha}) = \lambda(y_i g_i - y_j g_j) - \frac{\lambda^2}{2}(K_{ii} + K_{jj} - 2K_{ij})$$

Newton step gives optimal λ^*

$$\lambda^* = \frac{y_i g_i - y_j g_j}{K_{ii} + K_{jj} - 2K_{ij}}$$

yielding a gain $\mathcal{D}(\boldsymbol{\alpha} + \lambda^* \mathbf{u}_{ij}) - \mathcal{D}(\boldsymbol{\alpha})$ of

$$\frac{(y_i g_i - y_j g_j)^2}{2(K_{ii} + K_{jj} - 2K_{ij})} \ .$$



Maximum gain working set selection

Idea: Select i and j such that gain

$$\mathcal{D}(\alpha + \lambda^* \mathbf{u}_{ij}) - \mathcal{D}(\alpha) = \frac{(y_i g_i - y_j g_j)^2}{2(K_{ii} + K_{jj} - 2K_{ij})}$$

is maximized (ignoring box constraints).

Problem: Checking all $\ell(\ell - 1)/2$ index pairs is not feasible.

Solution:

- ① First index i is picked according to most violating pair heuristic.
- ② Second index j is selected to maximize gain.
 - *Second order working set selection*
 - Requires just $\mathcal{O}(\ell)$ computations (given reasonable caching strategy).



Outline

- 1 Recall: Linear Classification
- 2 Recall: Support Vector Machines
- 3 Karush-Kuhn-Tucker Theorem
- 4 SVM Dual Optimization Problems
- 5 SMO for SVM Training
- 6 Working Set Selection
- 7 Second-order SMO for Online Learning



Online learning

- We consider an ordered series of patterns $(x_1, y_1), (x_2, y_2), \dots$ presented in a serial fashion.
- Online learning algorithms adapt hypotheses by successively processing new training examples.

At step $t > 0$, the goal of online learning is to infer a new hypothesis h_t given

- ① the current sample (x_t, y_t) ,
- ② the previous hypothesis h_{t-1} (where induction starts from some a priori hypothesis h_0), and
- ③ some finite memory $M \subseteq \{(x_1, y_1), (x_2, y_2), \dots, (x_{t-1}, y_{t-1})\}$ of previous data points.

At each step t , h_t should minimize some loss function.



Online learning algorithms. . .

- provide usable solution without processing complete training set,
- allow to deploy an adaptive system that further improves over time (life-long learning),
- are ideal for streaming data,
- can be adopted for *novelty detection*,
- allow to track changes in the environment (i.e., can be applied if i.i.d. assumption is violated),
- need not be worse than batch learning algorithms even in batch scenario.



Active learning has successfully been applied to address sample selection bias.

If

- there is a wealth of input data from \mathcal{X} , but the corresponding targets are difficult or expensive to determine, or
- the data set is so large that processing every element is not feasible

then we need systems which autonomously choose the input data points from which they expect to learn most – that learn *actively*!

Active learning usually generates a new hypothesis based on an existing hypothesis and new sample data, that is, it requires online learning.



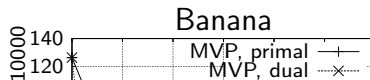
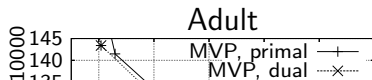
Online learning using LASVM

LASVM

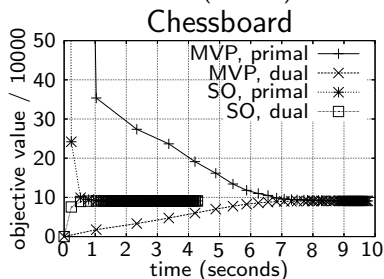
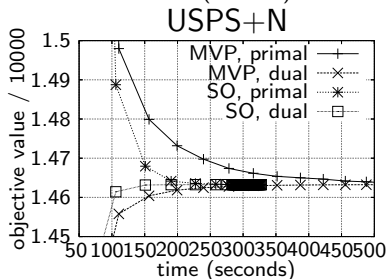
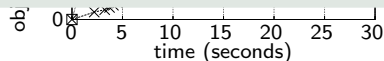
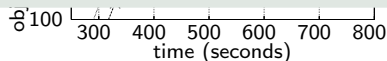
- 1 SVM with SVs M
 - 2 **repeat**
 - 3 observe pattern (x_t, y_t)
 // Step 1
 - 4 select working set $B = \{t, i\}, i \in M$
 - 5 $M \leftarrow M \cup \{t\}$
 - 6 solve QP restricted to B
 - 7 remove non-SVs from M
 // Step 2
 - 8 select working set $B = \{i, j\}, i, j \in M$
 - 9 solve QP restricted to B
 - 10 remove non-SVs from M
 - 11 **until** *stopping criterion is met*
-



Online SVMs results I



- For difficult data sets (i.e., that cannot be learnt in one pass), second-order working set selection is clearly superior.



Online SVMs results II

dataset	algorithm	#SV	cl. rate	primal	dual
---------	-----------	-----	----------	--------	------

Second-order working set selection makes better decisions, which lead to more efficient solutions in terms of

- classification error, and
- number of support vectors.

USPS+N	MVP	3870	99.09%	17,214	14,207
$\ell = 7,329$	SO	2501	99.36%	15,867	14,454
10 partitions	LIBSVM	2747	99.51%	14,635	14,633
Chessboard	MVP	3395	97.86%	782,969	60,771
$\ell = 5000$	SO	944	98.91%	143,614	83,686
1 partition	LIBSVM	873	99.40%	90,959	90,946



Summary

- Machine learning requires solving optimization problems, and knowledge in mathematical optimization helps to scale-up learning algorithms.
- Decomposition algorithms are state-of-the-art for training *non-linear* SVMs.
- Decomposition algorithms can be used for online and active learning.
- Second-order working set selection is in general better than first order methods.



References I

Overview

J. Shawe-Taylor and S. Sun. A review of optimization methodologies in support vector machines. *Neurocomputing* 74(17): 3609–3618, 2011.

L. Bottou and C.-J. Lin. Support Vector Machine Solvers. In L. Bottou, O. Chapelle, D. DeCoste, and J. Weston, eds.: *Large Scale Kernel Machines*, MIT Press, 2007.

SMO classics

T. Joachims. Making Large-Scale SVM Learning Practical. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, eds.: *Advances in Kernel Methods – Support Vector Learning*, pp. 185–208, MIT Press, 1999.

J. Platt. Fast training of support vector machines using sequential minimal optimization. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, eds.: *Advances in Kernel Methods – Support Vector Learning*, pp. 185–208, MIT Press, 1999.

H.-U. Simon, N. List. SVM-Optimization and Steepest-Descent Line Search. *Proceedings of the 22nd Conference on Learning Theory (COLT)*, 2009.



References II

Second-order working set selection

R.-E. Fan, P.-H. Chen, and C.-J. Lin. Working set selection using the second order information for training SVM. *Journal of Machine Learning Research* 6:1889–1918, 2005.

T. Glasmachers and C. Igel. Maximum-Gain Working Set Selection for SVMs. *Journal of Machine Learning Research* 7: 1437–1466, 2006.

SVM online and active learning

A. Bordes, S. Ertekin, J. Weston, and L. Bottou. Fast Kernel Classifiers with Online and Active Learning. *Journal of Machine Learning Research* 6:1579-1619, 2005

T. Glasmachers and C. Igel. Second-Order SMO Improves SVM Online and Active Learning, *Neural Computation* 20(2):374–382, 2008

Active learning and selection bias

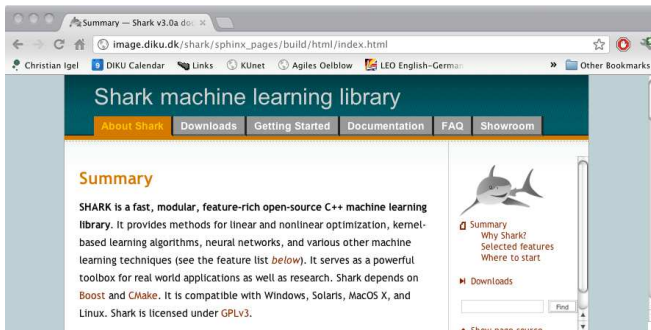
J. Kremer, F. Gieseke, K. Steenstrup Pedersen, and C. Igel. Nearest Neighbor Density Ratio Estimation for Large-Scale Applications in Astronomy. *Astronomy and Computing* 12:67–72, 2015



Shark

`image.diku.dk/shark`

Igel, Glasmachers, Heidrich-Meisner: Shark, *Journal of Machine Learning Research* 9:993–996, 2008



Gold Prize at Open Source Software World Challenge 2011

