

# Support Vector Machines

Charlie Frogner <sup>1</sup>

MIT

2011

---

<sup>1</sup>Slides mostly stolen from Ryan Rifkin (Google).

- Regularization derivation of SVMs.
- Analyzing the SVM problem: optimization, duality.
- Geometric derivation of SVMs.
- Practical issues.

# The Regularization Setting (Again)

Given  $n$  examples  $(x_1, y_1), \dots, (x_n, y_n)$ , with  $x_i \in \mathbb{R}^n$  and  $y_i \in \{-1, 1\}$  for all  $i$ .

We can find a classification function by solving a regularized learning problem:

$$\operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n V(y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}}^2.$$

Note that in this class we are specifically considering **binary classification**.

# The Hinge Loss

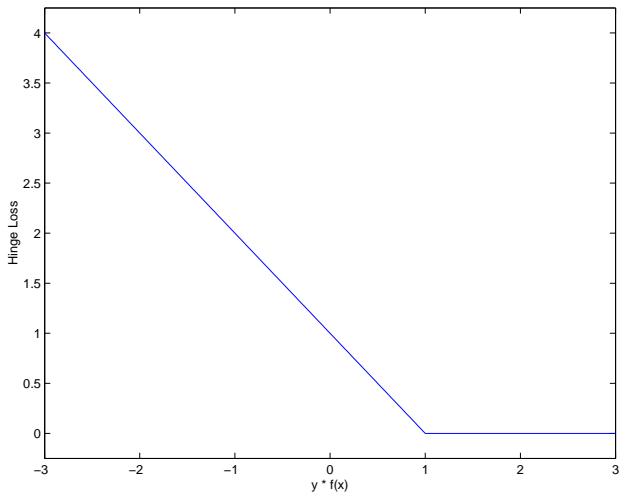
The classical SVM arises by considering the specific loss function

$$V(f(x), y) \equiv (1 - yf(x))_+,$$

where

$$(k)_+ \equiv \max(k, 0).$$

# The Hinge Loss



# Substituting In The Hinge Loss

With the hinge loss, our regularization problem becomes

$$\operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (1 - y_i f(x_i))_+ + \lambda \|f\|_{\mathcal{H}}^2.$$

Note that we don't have a  $\frac{1}{2}$  multiplier on the regularization term.

This problem is non-differentiable (because of the “kink” in  $V$ ).  
So rewrite the “max” function using slack variables  $\xi_i$ .

$$\begin{aligned} \operatorname{argmin}_{f \in \mathcal{H}} \quad & \frac{1}{n} \sum_{i=1}^n \xi_i + \lambda \|f\|_{\mathcal{H}}^2 \\ \text{subject to :} \quad & \xi_i \geq 1 - y_i f(x_i) \quad i = 1, \dots, n \\ & \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

# Applying The Representer Theorem

Substituting in:

$$f^*(x) = \sum_{i=1}^n c_i K(x, x_i),$$

we get a constrained quadratic programming problem:

$$\begin{aligned} & \underset{c \in \mathbb{R}^n, \xi \in \mathbb{R}^n}{\operatorname{argmin}} && \frac{1}{n} \sum_{i=1}^n \xi_i + \lambda c^T \mathbf{K} c \\ \text{subject to :} &&& \xi_i \geq 1 - y_i \sum_{j=1}^n c_j K(x_i, x_j) \quad i = 1, \dots, n \\ &&& \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$



# Adding A Bias Term

Adding an unregularized bias term  $b$  (which presents some theoretical difficulties) we get the “primal” SVM:

$$\begin{aligned} & \underset{\mathbf{c} \in \mathbb{R}^n, b \in \mathbb{R}, \xi \in \mathbb{R}^n}{\operatorname{argmin}} && \frac{1}{n} \sum_{i=1}^n \xi_i + \lambda \mathbf{c}^T \mathbf{K} \mathbf{c} \\ & \text{subject to :} && \xi_i \geq 1 - y_i \left( \sum_{j=1}^n c_j K(\mathbf{x}_i, \mathbf{x}_j) + b \right) \quad i = 1, \dots, n \\ & && \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

In most of the SVM literature, instead of  $\lambda$ , a parameter  $C$  is used to control regularization:

$$C = \frac{1}{2\lambda n}.$$

Using this definition (after multiplying our objective function by the constant  $\frac{1}{2\lambda}$ ), the regularization problem becomes

$$\operatorname{argmin}_{f \in \mathcal{H}} C \sum_{i=1}^n V(y_i, f(x_i)) + \frac{1}{2} \|f\|_{\mathcal{H}}^2.$$

Like  $\lambda$ , the parameter  $C$  also controls the tradeoff between classification accuracy and the norm of the function. The primal problem becomes ...

# The Reparametrized Problem

$$\begin{aligned} & \underset{c \in \mathbb{R}^n, b \in \mathbb{R}, \xi \in \mathbb{R}^n}{\operatorname{argmin}} && C \sum_{i=1}^n \xi_i + \frac{1}{2} c^T K c \\ & \text{subject to :} && \xi_i \geq 1 - y_i \left( \sum_{j=1}^n c_j K(x_i, x_j) + b \right) \quad i = 1, \dots, n \\ & && \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

# How to Solve?

$$\begin{aligned} & \underset{c \in \mathbb{R}^n, b \in \mathbb{R}, \xi \in \mathbb{R}^n}{\operatorname{argmin}} && C \sum_{i=1}^n \xi_i + \frac{1}{2} c^T K c \\ & \text{subject to :} && \xi_i \geq 1 - y_i (\sum_{j=1}^n c_j K(x_i, x_j) + b) \quad i = 1, \dots, n \\ & && \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

- This is a constrained optimization problem. The general approach:
  - Form the *primal* problem – we did this.
  - *Lagrangian* from primal – just like Lagrange multipliers.
  - *Dual* – one dual variable associated to each primal constraint in the Lagrangian.

We derive the dual from the primal using the Lagrangian:

$$\begin{aligned} L(\mathbf{c}, \xi, \mathbf{b}, \alpha, \zeta) &= \mathbf{C} \sum_{i=1}^n \xi_i + \frac{1}{2} \mathbf{c}^T \mathbf{K} \mathbf{c} \\ &\quad - \sum_{i=1}^n \alpha_i (y_i \{ \sum_{j=1}^n c_j \mathbf{K}(x_i, x_j) + \mathbf{b} \} - 1 + \xi_i) \\ &\quad - \sum_{i=1}^n \zeta_i \xi_i \end{aligned}$$

Dual problem is:

$$\operatorname{argmax}_{\alpha, \zeta \geq 0} \inf_{\mathbf{c}, \xi, b} L(\mathbf{c}, \xi, b, \alpha, \zeta)$$

First, minimize  $L$  w.r.t.  $(\mathbf{c}, \xi, b)$ :

$$(1) \quad \frac{\partial L}{\partial \mathbf{c}} = 0 \quad \Longrightarrow \quad \mathbf{c}_i = \alpha_i \mathbf{y}_i$$

$$(2) \quad \frac{\partial L}{\partial b} = 0 \quad \Longrightarrow \quad \sum_{i=1}^n \alpha_i \mathbf{y}_i = 0$$

$$(3) \quad \frac{\partial L}{\partial \xi_i} = 0 \quad \Longrightarrow \quad \mathbf{C} - \alpha_i - \zeta_i = 0$$

$$\Longrightarrow \quad 0 \leq \alpha_i \leq \mathbf{C}$$

Dual:

$$\operatorname{argmax}_{\alpha, \zeta \geq 0} \inf_{\mathbf{c}, \xi, \mathbf{b}} L(\mathbf{c}, \xi, \mathbf{b}, \alpha, \zeta)$$

Optimality conditions:

- (1)  $\mathbf{c}_i = \alpha_i \mathbf{y}_i$
- (2)  $\sum_{i=1}^n \alpha_i \mathbf{y}_i = \mathbf{0}$
- (3)  $\alpha_i \in [0, C]$

Plug in (2) and (3):

$$\operatorname{argmax}_{\alpha \geq 0} \inf_{\mathbf{c}} L(\mathbf{c}, \alpha) = \frac{1}{2} \mathbf{c}^T \mathbf{K} \mathbf{c} + \sum_{i=1}^n \alpha_i \left( 1 - y_i \sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j) \mathbf{c}_j \right)$$

Dual:

$$\operatorname{argmax}_{\alpha, \zeta \geq 0} \inf_{\mathbf{c}, \xi, \mathbf{b}} L(\mathbf{c}, \xi, \mathbf{b}, \alpha, \zeta)$$

Optimality conditions:

- (1)  $\mathbf{c}_i = \alpha_i \mathbf{y}_i$
- (2)  $\sum_{i=1}^n \alpha_i \mathbf{y}_i = 0$
- (3)  $\alpha_i \in [0, C]$

Plug in (1):

$$\begin{aligned} \operatorname{argmax}_{\alpha \geq 0} L(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \mathbf{y}_i \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) \alpha_j \mathbf{y}_j \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \alpha^T (\operatorname{diag} \mathbf{Y}) \mathbf{K} (\operatorname{diag} \mathbf{Y}) \alpha \end{aligned}$$



# The Primal and Dual Problems Again

$$\begin{aligned} & \underset{\mathbf{c} \in \mathbb{R}^n, \mathbf{b} \in \mathbb{R}, \xi \in \mathbb{R}^n}{\operatorname{argmin}} && \mathbf{C} \sum_{i=1}^n \xi_i + \frac{1}{2} \mathbf{c}^T \mathbf{K} \mathbf{c} \\ & \text{subject to :} && \xi_i \geq 1 - y_i (\sum_{j=1}^n c_j K(\mathbf{x}_i, \mathbf{x}_j) + b) \quad i = 1, \dots, n \\ & && \xi_i \geq 0 \quad i = 1, \dots, n \end{aligned}$$

$$\begin{aligned} & \underset{\alpha \in \mathbb{R}^n}{\max} && \sum_{i=1}^n \alpha_i - \frac{1}{2} \alpha^T \mathbf{Q} \alpha \\ & \text{subject to :} && \sum_{i=1}^n y_i \alpha_i = 0 \\ & && 0 \leq \alpha_i \leq C \quad i = 1, \dots, n \end{aligned}$$

- Basic idea: solve the dual problem to find the optimal  $\alpha$ 's, and use them to find  $b$  and  $c$ .
- The dual problem is easier to solve the primal problem. It has simple box constraints and a single equality constraint, and the problem can be decomposed into a sequence of smaller problems (see appendix).

# Interpreting the solution

$\alpha$  tells us:

- $c$  and  $b$ .
- The identities of the misclassified points.

How to analyze? Use the *optimality conditions*.

- Already used: derivative of  $L$  w.r.t.  $(c, \xi, b)$  is zero at optimality.
- Haven't used: complementary slackness, primal/dual constraints.

# Optimality Conditions: all of them

All optimal solutions must satisfy:

$$\sum_{j=1}^n c_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{j=1}^n y_j \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) = 0 \quad i = 1, \dots, n$$

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

$$C - \alpha_i - \zeta_i = 0 \quad i = 1, \dots, n$$

$$y_i \left( \sum_{j=1}^n y_j \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) + b \right) - 1 + \xi_i \geq 0 \quad i = 1, \dots, n$$

$$\alpha_i \left[ y_i \left( \sum_{j=1}^n y_j \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) + b \right) - 1 + \xi_i \right] = 0 \quad i = 1, \dots, n$$

$$\zeta_i \xi_i = 0 \quad i = 1, \dots, n$$

$$\xi_i, \alpha_i, \zeta_i \geq 0 \quad i = 1, \dots, n$$

These optimality conditions are both necessary and sufficient for optimality:  $(c, \xi, b, \alpha, \zeta)$  satisfy all of the conditions if and only if they are optimal for both the primal and the dual. (Also known as the Karush-Kuhn-Tucker (KKT) conditions.)

$$\frac{\partial L}{\partial \mathbf{c}} = \mathbf{0} \implies \mathbf{c}_i = \alpha_i \mathbf{y}_i, \forall i$$

Suppose we have the optimal  $\alpha_j$ 's. Also suppose that there exists an  $i$  satisfying  $0 < \alpha_i < C$ . Then

$$\alpha_i < C \implies \zeta_i > 0$$

$$\implies \xi_i = 0$$

$$\implies y_i \left( \sum_{j=1}^n y_j \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) + b \right) - 1 = 0$$

$$\implies b = y_i - \sum_{j=1}^n y_j \alpha_j K(\mathbf{x}_i, \mathbf{x}_j)$$

(Remember we defined  $f(\mathbf{x}) = \sum_{i=1}^n y_i \alpha_i K(\mathbf{x}, \mathbf{x}_i) + b$ .)

$$\begin{aligned} y_i f(\mathbf{x}_i) > 1 &\Rightarrow (1 - y_i f(\mathbf{x}_i)) < 0 \\ &\Rightarrow \xi_i \neq (1 - y_i f(\mathbf{x}_i)) \\ &\Rightarrow \alpha_j = 0 \end{aligned}$$



# Interpreting the solution — support vectors

$$\begin{aligned}y_i f(\mathbf{x}_i) < 1 &\Rightarrow (1 - y_i f(\mathbf{x}_i)) > 0 \\ &\Rightarrow \xi_i > 0 \\ &\Rightarrow \zeta_i = 0 \\ &\Rightarrow \alpha_i = C\end{aligned}$$

So

$$y_i f(x_i) < 1 \Rightarrow \alpha_i = C.$$

Conversely, suppose  $\alpha_j = C$ :

$$\begin{aligned} \alpha_j = C &\implies \xi_j = 1 - y_j f(x_j) \\ &\implies y_j f(x_j) \leq 1 \end{aligned}$$

# Interpreting the solution

Here are all of the derived conditions:

$$\alpha_j = 0 \implies y_j f(\mathbf{x}_j) \geq 1$$

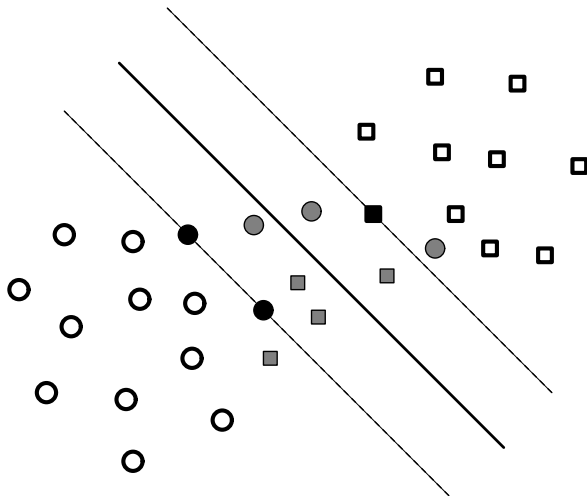
$$0 < \alpha_j < C \implies y_j f(\mathbf{x}_j) = 1$$

$$\alpha_j = C \iff y_j f(\mathbf{x}_j) < 1$$

$$\alpha_j = 0 \iff y_j f(\mathbf{x}_j) > 1$$

$$\alpha_j = C \implies y_j f(\mathbf{x}_j) \leq 1$$

# Geometric Interpretation of Reduced Optimality Conditions



# Summary so far

- The SVM is a Tikhonov regularization problem, using the hinge loss:

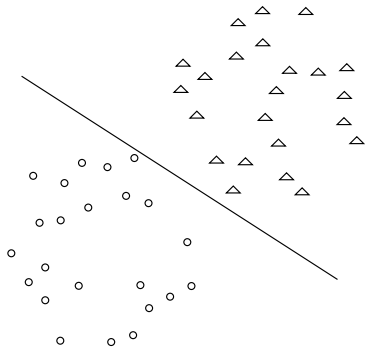
$$\operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (1 - y_i f(\mathbf{x}_i))_+ + \lambda \|f\|_{\mathcal{H}}^2.$$

- Solving the SVM means solving a constrained quadratic program.
- Solutions can be *sparse* – some coefficients are zero.
- The nonzero coefficients correspond to points that aren't classified correctly enough – this is where the “support vector” in SVM comes from.

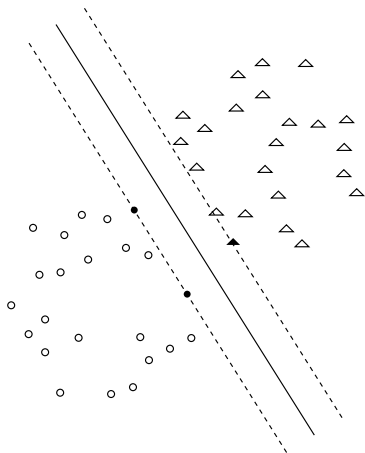
# The Geometric Approach

The “traditional” approach to developing the mathematics of SVM is to start with the concepts of *separating hyperplanes* and *margin*. The theory is usually developed in a linear space, beginning with the idea of a perceptron, a linear hyperplane that separates the positive and the negative examples. Defining the margin as the distance from the hyperplane to the nearest example, the basic observation is that intuitively, we expect a hyperplane with larger margin to generalize better than one with smaller margin.

# Large and Small Margin Hyperplanes



(a)



(b)

# Maximal Margin Classification

Classification function:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x}). \quad (1)$$

$\mathbf{w}$  is a normal vector to the hyperplane separating the classes. We define the boundaries of the margin by  $\langle \mathbf{w}, \mathbf{x} \rangle = \pm 1$ .

What happens as we change  $\|\mathbf{w}\|$ ?

We push the margin in/out by rescaling  $\mathbf{w}$  – the margin moves out with  $\frac{1}{\|\mathbf{w}\|}$ . So maximizing the margin corresponds to minimizing  $\|\mathbf{w}\|$ .



# Maximal Margin Classification

Classification function:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x}). \quad (1)$$

$w$  is a normal vector to the hyperplane separating the classes. We define the boundaries of the margin by  $\langle w, \mathbf{x} \rangle = \pm 1$ .

What happens as we change  $\|w\|$ ?

We push the margin in/out by rescaling  $w$  – the margin moves out with  $\frac{1}{\|w\|}$ . So maximizing the margin corresponds to minimizing  $\|w\|$ .

# Maximal Margin Classification, Separable case

Separable means  $\exists w$  s.t. all points are beyond the margin, i.e.

$$y_i \langle w, x_i \rangle \geq 1, \forall i.$$

So we solve:

$$\begin{aligned} \underset{w}{\operatorname{argmin}} \quad & \|w\|^2 \\ \text{s.t.} \quad & y_i \langle w, x_i \rangle \geq 1, \forall i \end{aligned}$$

# Maximal Margin Classification, Non-separable case

Non-separable means there are points on the wrong side of the margin, i.e.

$$\exists i \text{ s.t. } y_i \langle w, x_i \rangle < 1 .$$

We add slack variables to account for the wrongness:

$$\begin{aligned} \operatorname{argmin}_{\xi_i, w} \quad & \sum_{i=1}^n \xi_i + \|w\|^2 \\ \text{s.t.} \quad & y_i \langle w, x_i \rangle \geq 1 - \xi_i, \quad \forall i \end{aligned}$$

# Historical Perspective

Historically, most developments begin with the geometric form, derived a dual program which was identical to the dual we derived above, and only then observed that the dual program required only dot products and that these dot products could be replaced with a kernel function.

# More Historical Perspective

In the linearly separable case, we can also derive the separating hyperplane as a vector parallel to the vector connecting the closest two points in the positive and negative classes, passing through the perpendicular bisector of this vector. This was the “Method of Portraits”, derived by Vapnik in the 1970’s, and recently rediscovered (with non-separable extensions) by Keerthi.

- The SVM is a Tikhonov regularization problem, with the hinge loss:

$$\operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (1 - y_i f(\mathbf{x}_i))_+ + \lambda \|f\|_{\mathcal{H}}^2.$$

- Solving the SVM means solving a constrained quadratic program.
  - It's better to work with the dual program.
- Solutions can be *sparse* – few non-zero coefficients.
- The non-zero coefficients correspond to points not classified correctly enough – a.k.a. “support vectors.”
- There is alternative, geometric interpretation of the SVM, from the perspective of “maximizing the margin.”

- We can also use RLS for classification. What are the tradeoffs?
- SVM possesses sparsity: can have parameters set to zero in the solution. This enables potentially faster training and faster prediction than RLS.
- SVM QP solvers tend to have many parameters to tune.
- SVM can scale to very large datasets, unlike RLS – for the moment (active research topic!).

# Good Large-Scale SVM Solvers

- SVM Light: <http://svmlight.joachims.org>
- SVM Torch: <http://www.torch.ch>
- libSVM:  
<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>



(Follows.)

Our plan will be to solve the dual problem to find the  $\alpha$ 's, and use that to find  $b$  and our function  $f$ . The dual problem is easier to solve the primal problem. It has simple box constraints and a single inequality constraint, even better, we will see that the problem can be *decomposed* into a sequence of smaller problems.

We can solve QPs using standard software. Many codes are available. Main problem — the  $Q$  matrix is dense, and is  $n$ -by- $n$ , so we cannot write it down. Standard QP software requires the  $Q$  matrix, so is not suitable for large problems.

Partition the dataset into a *working set*  $W$  and the remaining points  $R$ . We can rewrite the dual problem as:

$$\begin{aligned} & \max_{\alpha_W \in \mathbb{R}^{|W|}, \alpha_R \in \mathbb{R}^{|R|}} && \sum_{i \in W} \alpha_i + \sum_{i \in R} \alpha_i \\ & \text{subject to :} && -\frac{1}{2} [\alpha_W \ \alpha_R] \begin{bmatrix} Q_{WW} & Q_{WR} \\ Q_{RW} & Q_{RR} \end{bmatrix} \begin{bmatrix} \alpha_W \\ \alpha_R \end{bmatrix} \\ & && \sum_{i \in W} y_i \alpha_i + \sum_{i \in R} y_i \alpha_i = 0 \\ & && 0 \leq \alpha_i \leq C, \forall i \end{aligned}$$

Suppose we have a feasible solution  $\alpha$ . We can get a better solution by treating the  $\alpha_W$  as variable and the  $\alpha_R$  as constant. We can solve the reduced dual problem:

$$\begin{aligned} \max_{\alpha_W \in \mathbb{R}^{|W|}} \quad & (\mathbf{1} - Q_{WR}\alpha_R)\alpha_W - \frac{1}{2}\alpha_W Q_{WW}\alpha_W \\ \text{subject to :} \quad & \sum_{i \in W} y_i \alpha_i = - \sum_{i \in R} y_i \alpha_i \\ & \mathbf{0} \leq \alpha_i \leq \mathbf{C}, \forall i \in W \end{aligned}$$

The reduced problems are fixed size, and can be solved using a standard QP code. Convergence proofs are difficult, but this approach seems to always converge to an optimal solution in practice.

# Selecting the Working Set

There are many different approaches. The basic idea is to examine points not in the working set, find points which violate the reduced optimality conditions, and add them to the working set. Remove points which are in the working set but are far from violating the optimality conditions.