

## Support vector machines unit:

- ▶ Linear programming (last week Wednesday)
- ▶ SVM concepts (last week Friday)
- ▶ Lab: SVM applications (Monday)
- ▶ The math of SVMs (Wednesday)
- ▶ SVM algorithms (**today**)

## Today:

- ▶ Summary of hard-margin version, with algorithm
- ▶ Quadratic programming and QP solver library
- ▶ Soft-margin version
- ▶ Kernelized version
- ▶ Put it all together in an algorithm

The most important source for all of this was Stephen Marsland, *Machine Learning: An Algorithmic Perspective*, 2015, pg 179–183.

Given training data  $\mathbf{X}, \mathbf{y}$ , where  $y_n \in \{-1, +1\}$ , find  $\mathbf{w}$ ,  $b$ , and  $r$ , specifically

maximize  $r$

subject to the constraints  $\forall \mathbf{x}_n, y_n, y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq r$   
 $\|\mathbf{w}\| = 1$   
 $r > 0$

Or, equivalently

minimize  $\frac{1}{2} \|\mathbf{w}\|^2$

subject to the constraints  $\forall \mathbf{x}_n, y_n, y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1$

alternately written as  $\forall \mathbf{x}_n, y_n, 1 - y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 0$

Combining the margin maximization with the fact that examples need to be on the correct side of the hyperplane (based on their labels) gives us

$$\max_{w, b} \frac{1}{\|w\|} \quad (12.16)$$

$$\text{subject to } y_n((w, x_n) + b) \geq 1 \text{ for all } n = 1, \dots, N. \quad (12.17)$$

Instead of maximizing the reciprocal of the norm as in (12.16), we often minimize the squared norm. We also often include a constant  $\frac{1}{2}$  that does not affect the optimal  $w, b$  but yields a tidier form when we compute the gradient. Then our objective becomes

$$\min_{w, b} \frac{1}{2} \|w\|^2 \quad (12.18)$$

$$\text{subject to } y_n((w, x_n) + b) \geq 1 \text{ for all } n = 1, \dots, N. \quad (12.19)$$

Equation (12.18) is known as the *hard margin SVM*. The reason for the expression "hard" is because the preceding formulation does not allow for any violations of the margin condition. We will see in Section 12.2.4 that this "hard" condition can be relaxed to accommodate violations if the data is not linearly separable.

### 12.2.3 Why We Can Set the Margin to 1

In Section 12.2.1, we argued that we would like to maximize some value  $r$ , which represents the distance of the closest example to the hyperplane. In Section 12.2.2, we scaled the data such that the closest example is of distance 1 to the hyperplane. In this section, we relate the two derivations and show that they are equivalent.

**Theorem 12.1.** *Maximizing the margin  $r$ , where we consider normalized weights as in (12.10),*

$$\max_{w, b, r} \underbrace{r}_{\text{margin}} \quad (12.20)$$

$$\text{subject to } \underbrace{y_n((w, x_n) + b)}_{\text{data fitting}} \geq r, \quad \underbrace{\|w\| = 1}_{\text{normalization}}, \quad r > 0,$$

is equivalent to scaling the data, such that the margin is unity:

$$\min_{w, b} \frac{1}{2} \|w\|^2 \quad (12.21)$$

$$\text{subject to } \underbrace{y_n((w, x_n) + b)}_{\text{margin}} \geq 1, \quad \underbrace{\text{data fitting}}$$

*Proof* Consider (12.20). Since the square is a strictly monotonic transformation for nonnegative arguments, the maximum stays the same if we consider  $r^2$  in the objective. Since  $\|w\| = 1$ , we can reparametrize the equation with a new weight vector  $w'$  that is not normalized by explicitly using  $\frac{w'}{\|w'\|}$ . We obtain

### 12.2 Primal Support Vector Machine

$$\max_{w', b, r} r^2 \quad (12.22)$$

$$\text{subject to } y_n \left( \left\langle \frac{w'}{\|w'\|}, x_n \right\rangle + b \right) \geq r, \quad r > 0.$$

Equation (12.22) explicitly states that the distance  $r$  is positive. Therefore, we can divide the first constraint by  $r$ , which yields

$$\max_{w', b, r} r^2 \quad (12.23)$$

$$\text{subject to } y_n \left( \left\langle \frac{w'}{\|w'\| r}, x_n \right\rangle + \frac{b}{r} \right) \geq 1, \quad r > 0$$

renaming the parameters to  $w''$  and  $b''$ . Since  $w'' = \frac{w'}{\|w'\| r}$ , rearranging for  $r$  gives

$$\|w''\| = \left\| \frac{w'}{\|w'\| r} \right\| = \frac{1}{r} \left\| \frac{w'}{\|w'\|} \right\| = \frac{1}{r}. \quad (12.24)$$

By substituting this result into (12.23), we obtain

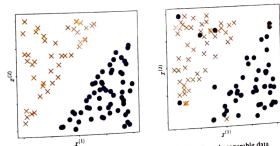
$$\max_{w'', b''} \frac{1}{\|w''\|^2} \quad (12.25)$$

$$\text{subject to } y_n((w'', x_n) + b'') \geq 1.$$

The final step is to observe that maximizing  $\frac{1}{\|w''\|^2}$  yields the same solution as minimizing  $\frac{1}{2} \|w''\|^2$ , which concludes the proof of Theorem 12.1.  $\square$

### 12.2.4 Soft Margin SVM: Geometric View

In the case where data is not linearly separable, we may wish to allow some examples to fall within the margin region, or even to be on the wrong side of the hyperplane as illustrated in Figure 12.6.



(a) Linearly separable data, with a large margin

(b) Nonlinearly separable data

Figure 12.6 (a) Linearly separable and (b) nonlinearly separable data.

The squared norm results in a convex quadratic programming problem for the SVM (Section 12.5).

hard margin SVM

Note that  $r > 0$  because we assumed linear separability, and hence there is no issue to divide by  $r$ .

rescaling the data is what we are charging, not X

The Lagrangian of the hard-margin version is

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=0}^{N-1} \lambda_n (1 - y_n(\mathbf{w}^T \mathbf{x}_n + b))$$

Take the gradient with respect to  $\mathbf{w}$  and  $b$ , set to  $\mathbf{0}$  or 0

$$\nabla_{\mathbf{w}} \mathcal{L} = \mathbf{w} - \sum_{n=0}^{N-1} \lambda_n y_n \mathbf{x}_n$$

$$\nabla_b \mathcal{L} = - \sum_{n=0}^{N-1} \lambda_n y_n$$

$$\mathbf{0} = \mathbf{w} - \sum_{n=0}^{N-1} \lambda_n y_n \mathbf{x}_n$$

$$0 = - \sum_{n=0}^{N-1} \lambda_n y_n$$

$$\mathbf{w}^* = \sum_{n=0}^{N-1} \lambda_n y_n \mathbf{x}_n$$

Wait, how does this help us find  $b^*$ ?

... where  $\mathbf{w}^*$  and  $b^*$  are the optimal weights and bias.

$$\mathbf{w}^* = \sum_{n=0}^{N-1} \lambda_n y_n \mathbf{x}_n = \left[ \dots \sum_{n=0}^{N-1} \lambda_n y_n \mathbf{x}_{n,i} \dots \right] \quad \text{and} \quad \sum_{n=0}^{N-1} \lambda_n y_n = 0.$$

Theorem 3.4.3 in Han Veiga and Ged tells us that with  $\mathbf{w}^*$ ,  $b^*$ , and  $\lambda^*$ ,

$$\forall \mathbf{x}_n, y_n, \quad \lambda_n (1 - y_n (\mathbf{w}^{*T} \mathbf{x}_n + b^*)) = 0$$

... which implies  $\lambda_n = 0$  for all *non-support vectors*  $\mathbf{x}_n$

Substitute  $\mathbf{w}^*$  and  $b^*$  into the Lagrangian to make it a function just of  $\lambda$ .

$$\underbrace{\mathcal{D}(\lambda)}_{\text{dual}} = \underbrace{\mathcal{L}(\mathbf{w}^*, b^*, \lambda)}_{\text{filled-in}} = \frac{1}{2} \|\mathbf{w}^*\|^2 + \sum_{n=0}^{N-1} \lambda_n (1 - y_n (\mathbf{w}^{*T} \mathbf{x}_n + b^*))$$

Simplify this based on results above into a quadratic program involving only  $\mathbf{X}$ ,  $\mathbf{y}$ , and  $\lambda$  with constraints  $\sum_{n=0}^{N-1} \lambda_n y_n = 0$ .

A *quadratic programming* problem (or a *quadratic program*) can be stated as, minimize

$$\frac{1}{2} \mathbf{x}^T \mathbf{P} \mathbf{x} + \mathbf{q}^T \mathbf{x}$$

subject to

$$\mathbf{G} \mathbf{x} \leq \mathbf{h}$$

$$\mathbf{A} \mathbf{x} = \mathbf{b}$$

In this formula, let  $n$  be the number of variables,  $m$  be the number of inequality constraints, and  $\ell$  be the number of equality constraints.

Soft-margin form:

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=0}^{N-1} \xi_n$$

subject to the constraints  $\forall \mathbf{x}_n, y_n, y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n$

Define a hyperplane

$$\mathbf{w}^T \phi(\mathbf{x}) + b = 0$$

such that

$$f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$$

classifies  $\mathbf{x}$ .

$\phi$  is a feature map that projects the vectors  $\mathbf{x}$  into higher dimensions. Assume  $k$  is a kernel function corresponding to  $\phi$  for efficiently computing dot products in these higher dimensions.

$$k(\mathbf{x}_a, \mathbf{x}_b) = \phi(\mathbf{x}_a)^T \phi(\mathbf{x}_b)$$



Kernelized form (hard- or soft-margin):

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 \left[ + C \sum_{n=0}^{N-1} \xi_n \right]$$

subject to the constraints  $\forall \mathbf{x}_n, y_n, y_n(\mathbf{w}^T \phi(\mathbf{x}_n + b)) \geq 1 [-\xi_n]$

(This doesn't actually use the kernel function... but when we put the problem in this form, we anticipate using the kernel.)

This *quadratic programming problem* has an equivalent *Lagrangian function*

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=0}^{N-1} \lambda_n y_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)$$

This function has the *dual representation*

$$\mathcal{D}(\boldsymbol{\lambda}) = \sum_{n=0}^{N-1} \lambda_n - \frac{1}{2} \sum_{n=0}^{N-1} \sum_{m=0}^{N-1} \lambda_n \lambda_m y_n y_m \cdot k(\mathbf{x}_n, \mathbf{x}_m)$$

which we want to maximize subject to constraints

$$0 \leq \lambda_n [\leq C],$$

$$\sum_{n=0}^{N-1} \lambda_n y_n = 0$$

where  $k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$

Let  $\mathbf{K}$  be the *kernel matrix* for data set  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{N-1}$ :

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_0, \mathbf{x}_0) & k(\mathbf{x}_1, \mathbf{x}_0) & \cdots & k(\mathbf{x}_{N-1}, \mathbf{x}_0) \\ k(\mathbf{x}_0, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_{N-1}, \mathbf{x}_1) \\ \vdots & & & \vdots \\ k(\mathbf{x}_0, \mathbf{x}_{N-1}) & k(\mathbf{x}_1, \mathbf{x}_{N-1}) & \cdots & k(\mathbf{x}_{N-1}, \mathbf{x}_{N-1}) \end{pmatrix}$$

Then

$$\mathcal{D}(\boldsymbol{\lambda}) = \sum_{n=1}^{N-1} \lambda_n - \frac{1}{2} \sum_{n=0}^{N-1} \sum_{m=0}^{N-1} \lambda_n \lambda_m y_n y_m \cdot k(\mathbf{x}_n, \mathbf{x}_m)$$

becomes

$$\mathcal{D}(\boldsymbol{\lambda}) = \mathbf{1}^T \boldsymbol{\lambda} - \frac{1}{2} \boldsymbol{\lambda}^T (\mathbf{y} \mathbf{y}^T \circ \mathbf{K}) \boldsymbol{\lambda}$$

where  $\mathbf{1} = [1 \ 1 \ 1 \ \cdots \ 1]^T$  and  $\circ$  indicates the Hadamard product.

In the formula

$$\mathcal{D}(\boldsymbol{\lambda}) = \mathbf{1}^T \boldsymbol{\lambda} - \frac{1}{2} \boldsymbol{\lambda}^T (\mathbf{y} \mathbf{y}^T \circ \mathbf{K}) \boldsymbol{\lambda}$$

the Hadamard product  $\circ$  gives us

$$\mathbf{y} \mathbf{y}^T \circ \mathbf{K} = \begin{pmatrix} y_0 \cdot y_0 \cdot k(\mathbf{x}_0, \mathbf{x}_0) & y_1 \cdot y_0 \cdot k(\mathbf{x}_1, \mathbf{x}_0) & \cdots & y_{N-1} \cdot y_0 \cdot k(\mathbf{x}_{N-1}, \mathbf{x}_0) \\ y_0 \cdot y_1 \cdot k(\mathbf{x}_0, \mathbf{x}_1) & y_1 \cdot y_1 \cdot k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & y_{N-1} \cdot y_1 \cdot k(\mathbf{x}_{N-1}, \mathbf{x}_1) \\ \vdots & \vdots & \ddots & \vdots \\ y_0 \cdot y_{N-1} \cdot k(\mathbf{x}_0, \mathbf{x}_{N-1}) & y_1 \cdot y_{N-1} \cdot k(\mathbf{x}_1, \mathbf{x}_{N-1}) & \cdots & y_{N-1} \cdot y_{N-1} \cdot k(\mathbf{x}_{N-1}, \mathbf{x}_{N-1}) \end{pmatrix}$$

Quadratic programming problem:

$$\min \quad \frac{1}{2} \mathbf{x}^T \mathbf{P} \mathbf{x} + \mathbf{q}^T \mathbf{x}$$

$$\mathbf{G} \mathbf{x} \leq \mathbf{h}$$

$$\mathbf{A} \mathbf{x} = \mathbf{b}$$

Our problem:

$$\max \quad \mathbf{1}^T \boldsymbol{\lambda} - \frac{1}{2} \boldsymbol{\lambda}^T (\mathbf{y} \mathbf{y}^T \circ \mathbf{K}) \boldsymbol{\lambda}$$

$$0 \leq \lambda_i \quad [\leq C]$$

$$\sum_{i=0}^{N-1} \lambda_i y_i = 0$$

We want to find  $\boldsymbol{\lambda}$ . Let  $\mathbf{P} = \mathbf{y} \mathbf{y}^T \circ \mathbf{K}$ ,  $\mathbf{q} = \begin{pmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{pmatrix}$ ,  $\mathbf{A} = (y_0 \quad y_1 \quad \cdots \quad y_{N-1})$ , and  $\mathbf{b} = (0)$ .

Quadratic programming problem:

$$\min \frac{1}{2} \mathbf{x}^T \mathbf{P} \mathbf{x} + \mathbf{q}^T \mathbf{x}$$

$$\mathbf{G} \mathbf{x} \leq \mathbf{h}$$

$$\mathbf{A} \mathbf{x} = \mathbf{b}$$

For *hard* margin classification ( $0 \leq \lambda_i$ ),

$$\mathbf{G} = -\mathcal{I} = \begin{pmatrix} -1 & 0 & \dots & 0 \\ 0 & -1 & \dots & 0 \\ \vdots & & & \vdots \\ 0 & 0 & \dots & -1 \end{pmatrix}$$

$$\mathbf{h} = \mathbf{0} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Quadratic programming problem:

$$\min \frac{1}{2} \mathbf{x}^T \mathbf{P} \mathbf{x} + \mathbf{q}^T \mathbf{x}$$

$$\mathbf{G} \mathbf{x} \leq \mathbf{h}$$

$$\mathbf{A} \mathbf{x} = \mathbf{b}$$

For *soft* margin classification ( $0 \leq \lambda_i \leq C$ ),

$$\mathbf{G} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & & & \vdots \\ 0 & 0 & \vdots & 1 \\ -1 & 0 & \dots & 0 \\ 0 & -1 & \dots & 0 \\ \vdots & & & \vdots \\ 0 & 0 & \dots & -1 \end{pmatrix}$$

$$\mathbf{h} = \begin{pmatrix} C \\ C \\ \vdots \\ C \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Support vectors are  $\{\mathbf{x}_i \mid \lambda_i \neq 0\}$ .

Weights are

$$\mathbf{w} = \sum_{n=0}^{N-1} \lambda_n y_n \mathbf{x}_n = \sum_{n \mid \lambda_n \neq 0} \lambda_n y_n \mathbf{x}_n$$

Intercept/bias is

$$b = \frac{1}{|\{\lambda_n \mid \lambda_n \neq 0\}|} \sum_{m \mid \lambda_j \neq 0} \left( y_j - \sum_{n \mid \lambda_n \neq 0} \lambda_n y_n \cdot k(\mathbf{x}_n, \mathbf{x}_m) \right)$$



To train a classifier for hard margin classification:

Given **data**  $\mathbf{X}$ , **targets**  $\mathbf{y}$ , and **kernel function**  $k$ ,

Compute kernel matrix  $\mathbf{K}$

Compute  $\mathbf{P} = \mathbf{y}\mathbf{y}^T \circ \mathbf{K}$

Assemble  $\mathbf{q}$  vector of  $-1$ s

Assemble  $\mathbf{A}$  matrix of  $y_i$  along the diagonal

Assemble  $\mathbf{G}$  matrix of  $-1$ s along the diagonal

Assemble  $\mathbf{h}$  vector of  $0$ s

Compute  $\boldsymbol{\lambda}$  vector by feeding  $\mathbf{P}$ ,  $\mathbf{q}$ ,  $\mathbf{G}$ ,  $\mathbf{h}$ ,  $\mathbf{A}$ , and  $\mathbf{b} = [0]$  into QP solver

Select support vectors from  $\boldsymbol{\lambda}$  that are not zero

Compute  $b$

(For soft margin, modify  $\mathbf{G}$  and  $\mathbf{h}$ .)

To classify new data point  $\mathbf{x}$ , compute  $\text{sign}(\sum_{n \mid \lambda_n \neq 0} \lambda_n y_n k(\mathbf{x}_n, \mathbf{x}) + b)$

## Coming up:

**Due Fri, Mar 7:**

*Take SVM quiz*

**Due Wed, Mar 19:**

*Implement SVM classification*

*(Midterm on Fri, Mar 21)*