

Support vector machines unit:

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

Today:

- ▶ Hard-margin
- ▶ Kernelized hard-margin
- ▶ Kernelized soft-margin

The most important source for all of this was Stephen Marsland, *Machine Learning: An Algorithmic Perspective*, 2015, pg 179–183. Notation adjusted to agree better with Deisenroth et al.

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

A *quadratic programming* problem can be stated as, minimize

$$\frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x}$$

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 ℓ be the number of equality constraints.

(Deisenroth et al combine the constraints into a single set of inequalities, $\mathbf{A} \mathbf{x} \leq \mathbf{b}$.)

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} .

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

With constraint

$$\forall n \in [0, N), \quad y_n \cdot (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1$$

maximize the margin between the hyperplane and the closest point by finding

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \right\} \\ &= \operatorname{argmin}_{\mathbf{w}, b} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 \right\} \\ &= \operatorname{argmin}_{\mathbf{w}, b} \left\{ \frac{1}{2} \mathbf{w}^T \mathbf{w} \right\} \end{aligned}$$

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

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

where $\boldsymbol{\alpha}$ is the vector of *Lagrangian multipliers*. This function has the *dual representation*

$$\mathcal{D}(\boldsymbol{\alpha}) = \sum_{n=0}^{N-1} \alpha_n - \frac{1}{2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} \alpha_i \alpha_j y_i y_j \cdot k(\mathbf{x}_i, \mathbf{x}_j)$$

which we want to maximize subject to constraints

$$0 \leq \alpha_n \leq C, \quad \forall n \in [0, N-1)$$

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

where $k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$

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

$$\mathcal{D}(\boldsymbol{\alpha}) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \cdot k(\mathbf{x}_i, \mathbf{x}_j)$$

becomes

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

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

In the formula

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

the Hadamard product \circ gives us

$$\mathbf{y} \mathbf{y}^T \circ \mathbf{K} = \begin{pmatrix} y_1 \cdot y_1 \cdot k(\mathbf{x}_1, \mathbf{x}_1) & y_2 \cdot y_1 \cdot k(\mathbf{x}_2, \mathbf{x}_1) & \cdots & y_N \cdot y_1 \cdot k(\mathbf{x}_N, \mathbf{x}_1) \\ y_1 \cdot y_2 \cdot k(\mathbf{x}_1, \mathbf{x}_2) & y_2 \cdot y_2 \cdot k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & y_N \cdot y_2 \cdot k(\mathbf{x}_N, \mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ y_1 \cdot y_N \cdot k(\mathbf{x}_1, \mathbf{x}_N) & y_2 \cdot y_N \cdot k(\mathbf{x}_2, \mathbf{x}_N) & \cdots & y_N \cdot y_N \cdot k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

A *quadratic programming* problem can be stated as, minimize

$$\frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x}$$

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 ℓ be the number of equality constraints.

(Deisenroth et al combine the constraints into a single set of inequalities, $\mathbf{A} \mathbf{x} \leq \mathbf{b}$.)

Quadratic programming problem:

$$\min \quad \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x}$$

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

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

Our problem:

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

$$0 \leq \alpha_i \leq C$$

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

We want to find $\boldsymbol{\alpha}$. Let $\mathbf{Q} = \mathbf{y} \mathbf{y}^T \circ \mathbf{K}$, $\mathbf{c} = \begin{pmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{pmatrix}$, $\mathbf{A} = (y_1 \ y_2 \ \cdots \ y_n)$, and

$$\mathbf{b} = (0).$$

Quadratic programming problem:

$$\min \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x}$$

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

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

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

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

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

Quadratic programming problem:

$$\min \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x}$$

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

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

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

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

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

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

Weights are

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

Intercept is

$$b = \frac{1}{|\{\alpha_n \mid \alpha_n \neq 0\}|} \sum_{j=1 \mid \alpha_j \neq 0}^N \left(y_j - \sum_{i=1 \mid \alpha_i \neq 0}^N \alpha_i y_i \cdot k(\mathbf{x}_i, \mathbf{x}_j) \right)$$

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 $\mathbf{Q} = \mathbf{y}\mathbf{y}^T \circ \mathbf{K}$

Assemble \mathbf{c} 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 0s

Compute α vector by feeding \mathbf{Q} , \mathbf{c} , \mathbf{G} , \mathbf{h} , \mathbf{A} , and $\mathbf{b}=\mathbf{0}$ into QP solver

Select support vectors from α that are not zero

Compute b

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

To classify new data point \mathbf{x} , compute $\sum_{n=1}^N \alpha_n y_n k(\mathbf{x}_n, \mathbf{x}) + b$