

Define a hyperplane

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

such that

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

classifies \mathbf{x} .

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

With constraint

$$\forall i \in [1, N], \quad t_i \left(\mathbf{w}^T \phi(\mathbf{x}_i) + b \right) \geq 1$$

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

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_i [t_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b)] \right\} \\ &= \operatorname{argmax}_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \right\} \\ &= \operatorname{argmin}_{\mathbf{w}, b} \left\{ \frac{1}{2} \|\mathbf{w}\| \right\} \\ &= \operatorname{argmin}_{\mathbf{w}, b} \left\{ \frac{1}{2} \mathbf{w}^T \mathbf{w} \right\} \end{aligned}$$

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

$$\mathcal{L}(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N a_i \left(t_i \left(\mathbf{w}^T \phi(\mathbf{x}_i) + b \right) \right)$$

where \mathbf{a} is the vector of *Lagrangian multipliers*. This function has the *dual representation*

$$\tilde{\mathcal{L}}(\mathbf{a}) = \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N a_i a_j t_i t_j k(\mathbf{x}_i, \mathbf{x}_j)$$

which we want to maximize subject to constraints

$$0 \leq a_i \leq C \quad \forall i \in [1, N]$$

$$\sum_{i=1}^N a_i t_i = 0$$

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 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$:

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_1) \\ k(\mathbf{x}_1, \mathbf{x}_2) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_2) \\ \vdots & & & \vdots \\ k(\mathbf{x}_1, \mathbf{x}_N) & k(\mathbf{x}_2, \mathbf{x}_N) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

Then

$$\tilde{\mathcal{L}}(\mathbf{a}) = \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N a_i a_j t_i t_j k(x_i, x_j)$$

becomes

$$\tilde{\mathcal{L}}(\mathbf{a}) = \mathbf{i}^T \mathbf{a} - \frac{1}{2} \mathbf{a}^T \mathbf{t} \mathbf{t}^T \mathbf{K} \mathbf{a}$$

where \mathbf{i} is the identity vector.

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

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

subject to

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

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

Quadratic programming problem:

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

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

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

Our problem:

$$\max \quad \mathbf{i}^T \mathbf{a} - \frac{1}{2} \mathbf{a}^T \mathbf{t} \mathbf{t}^T \mathbf{K} \mathbf{a}$$

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

$$\sum_{i=1}^N a_i t_i = 0$$

We want to find \mathbf{a} . Let $\mathbf{P} = \mathbf{t} \mathbf{t}^T \mathbf{K}$, $q = \begin{pmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{pmatrix}$, $\mathbf{A} = \begin{pmatrix} t_1 & 0 & \dots & 0 \\ 0 & t_2 & \dots & 0 \\ \vdots & & & \vdots \\ 0 & 0 & \dots & t_n \end{pmatrix}$, and $b = 0$.

Quadratic programming problem:

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

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

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

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

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

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

Quadratic programming problem:

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

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

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

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

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

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

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

Weights are

$$\mathbf{w} = \sum_{i=1}^N a_i t_i \mathbf{x}_i = \sum_{i=1 \mid a_i \neq 0} a_i t_i \mathbf{x}_i$$

Intercept is

$$b = \frac{1}{|\{a_i \mid a_i \neq 0\}|} \sum_{j=1 \mid a_j \neq 0}^N \left(t_j - \sum_{i=1 \mid a_i \neq 0}^N a_i t_i \mathbf{x}_i^T \mathbf{x}_j \right)$$

To train a classifier for hard margin classification:

Given **data**, **targets**, and k ,

Compute kernel matrix **K**

Compute $\mathbf{P} = \mathbf{t}\mathbf{t}^T \mathbf{K}$

Assemble **q** vector of -1 s

Assemble **A** matrix of t_i along the diagonal

Assemble **G** matrix of -1 s along the diagonal

Assemble **h** vector of 0s

Compute **a** vector by feeding **P**, **q**, **G**, **h**, **A**, and $\mathbf{b}=\mathbf{0}$ into QP solver

Select support vectors from **a** that are not zero

Compute b

(For soft margin, modify **G** and **h**.)

To classify new data point \mathbf{x} , compute $\sum_{i=1|a_i \neq 0}^N a_i t_i k(\mathbf{x}_i, \mathbf{x}) + b$