

# Machine Learning Problemset 4

Aimee Barciauskas

12 March 2016

## Problem 17

**Let**  $(x_1, y_1), \dots, (x_n, y_n)$  **be data in**  $\mathbb{R}^d \times \{-1, 1\}$ . Suppose the data are *linearly seperable*...

Since the data are linearly seperable, we can define a separating hyperplane:

$$\left\{x : f(x) = w^T x = 0\right\} \text{ where } |w| = 1$$

The seperating hyperplane returns a signed distance to the plane for each  $x$ , and classification is done according to the sign:

$$f(x) : \text{sgn}[w^T x]$$

The margin of this classifier is as stated in the problem:

$$\gamma(w) = \min_i \frac{y_i w^T x_i}{|w|}$$

An optimal  $f(x)$  is one which maximizes the margin  $\gamma$ , i.e.:

$$\gamma^* = \max_{\gamma, |w|=1} \text{s.t. } y_i w^T x_i \geq \gamma$$

The constraint on the norm of  $w$ ,  $|w| = 1$  can be removed when:

$$\frac{y_i w^T x_i}{|w|} \geq \gamma$$

and by setting  $|w| = \frac{1}{\gamma}$  the maximization optimization problem becomes a convex minimization problem of the form:

$$\min_w |w| \text{ s.t. } y_i w^T x_i \geq 1$$

Part 2

HELP

$w$  is a linear combination of all the points that are at the margin distance so  $w$  is a linear combination of these 2.

$$w = \sum_i \gamma_i w$$

# Problem 18

Part 1

$$y_n \in \{-1, 1\}$$

$$x_n \in \{0, 1\}$$

Data are linearly separable by the weight vector:

$$w = (w_0, \dots, w_d)$$

Being linearly separable means, by definition, there exists such a  $w^T x_i > 0$  whenever  $y_i = 1$  and  $w^T x_i \leq 0$  whenever  $y_i = -1$

By the statement of the problem,  $w^T x_i > 0$  whenever at least one  $x_{i,j} = 1$ . The  $w$  satisfying this is the  $w$  where  $w_0$  is  $(-1, 0)$  and all  $w_1 = \dots = w_d = 1$ :

$$\left\{ w : -1 < w_0 < 0 \text{ and } w_1 = \dots = w_d = 1 \right\}$$

So  $w^T x$  will be greater than 0 where at least one  $x_i$  is 1 and  $w_0$  otherwise. The  $x_i$ 's nearest the  $w$  vector are the  $x_i$ 's where all elements are zero and their counterparts: those with only one dimension being equal non-zero, i.e. the  $x_i$ 's satisfying:

$$\sum_{m=1}^d x_{m,i} = 1 \text{ } d \text{ is the dimension of } x$$

These two types of points define the location of the hyperplane. The margin maximizing the distance between between such  $x_i$ 's is where the margin equivalent for both these points:

$$\frac{y_j w^T x_j}{|w|} = \frac{y_i w^T x_i}{|w|}$$

Where  $y_j = -1$ ,  $y_i = 1$  and  $\sum_{m=1}^d x_{m,j} = 0$ ,  $\sum_{m=1}^d x_{m,i} = 1$

To solve for the optimal  $w^*$ , where we know all  $w_1 = \dots = w_d = 1$  so we are just solving for  $w_0$ :

$$w_0 = -\frac{1}{2}$$

Plugging this back into the equation for  $\gamma$ :

$$\gamma^*(w) = \frac{1}{2\sqrt{\frac{1}{4} + d}}$$

Part 2

$$w^T x > 0 \text{ when } \sum_i x_i > \frac{d}{2}$$

$w^T x = w_0^* + kw_{1,\dots,d}^*$  where  $k$  is the number of non-zero elements of  $x$  and  $w^*$  is the optimal weight vector

$$w^T x > 0 \text{ when } k \geq \frac{d+1}{2} \text{ and}$$

$$w^T x < 0 \text{ when } k < \frac{d-1}{2}$$

Say  $-w_0^* = kw_{1,\dots,d}^*$  and we can set  $w_{1,\dots,d}^* = 1$  and there are two scenarios:

1.  $w_0^* + \frac{d+1}{2} > 0$  when  $k \geq \frac{d+1}{2}$
2.  $w_0^* + \frac{d-1}{2} < 0$  when  $k < \frac{d-1}{2}$

So the possible values of  $w_0$  are:

$$\frac{d-1}{2} < -w_0^* < \frac{d+1}{2}$$

$x_i$ 's nearest the classifying hyperplane, are those where:

$$\sum_{m=1}^d x_{ij} = \frac{d-1}{2} \text{ or}$$

$$\sum_{m=1}^d x_{ij} = \frac{d+1}{2}$$

I.e.:

$$y_i = 1, \sum_i x_i = \frac{d+1}{2}$$

$$y_j = -1, \sum_j x_j = \frac{d-1}{2}$$

$\gamma(w)$  is the same for all  $x_i$  so:

$$y_i w^T x_i = y_j w^T x_j$$

Substituting  $y_i, y_j$  with  $-1, 1$ :

$$w^T x_i = -w^T x_j$$

$$w_0 + \sum_{m=1}^d w_m x_{i,m} = -w_0 - \sum_{m=1}^d w_m x_{j,m}$$

$$w_0 + \frac{d+1}{2} = -w_0 - \frac{d-1}{2}$$

$$w_0 = -\frac{d}{2}$$

Plugging this into the equation for  $\gamma$ :

$$\gamma(w) = \frac{w_0 + \frac{d+1}{2}}{\sqrt{w_0^2 + d}}$$

$$\gamma(w) = \frac{1}{2\sqrt{\frac{d^2}{4} + d}}$$

## Problem 19

*Part 1*

$$K(x, y) = \langle \phi(x), \phi(y) \rangle$$

where:

$$\phi(x)_n = \frac{1}{\sqrt{n!}} x^n e^{-\frac{x^2}{2}}$$

$$\langle \phi(x), \phi(y) \rangle = \sum_{n=0}^{\infty} \frac{1}{\sqrt{n!}} x^n e^{-\frac{x^2}{2}} \frac{1}{\sqrt{n!}} y^n e^{-\frac{y^2}{2}}$$

Simplifying and using that  $\sum_{n=0}^{\infty} \frac{x^n}{n!} = e^x$ :

$$= e^{xy - \frac{1}{2}(x^2 + y^2)}$$

Multiply the exponent by  $\frac{2}{2}$ :

$$= e^{-\frac{(x-y)^2}{2}}$$

*Is this the gaussian kernel?*

*Part 2*

Generalize to

$$\mathbb{R}^d$$

$$= \sum_{i=1}^{\infty} \frac{1}{n!} (\sqrt{(x^T x)} \sqrt{y^T y})^n \exp(-\frac{1}{2}(x^T x) - \frac{1}{2}(y^T y))$$

$$= \exp(\text{sqrt}(x^T x y^T y) - \frac{1}{2}x^T x - \frac{1}{2}y^T y)$$

$$= \exp(x^T y - \frac{1}{2}x^T x - \frac{1}{2}y^T y)$$

$$\begin{aligned}
&= \exp\left(\frac{2x^T y - x^T x - y^T y}{2}\right) \\
&= \exp\left(-\frac{(\text{norm}(x - y))^2}{2}\right)
\end{aligned}$$

## Problem 20

*Part 1*

$$\begin{aligned}
K(x, y) &= \langle \phi(x), \phi(y) \rangle \\
&= \left\langle \begin{pmatrix} \phi_1 \\ \phi_2 \end{pmatrix}_1(x), \begin{pmatrix} \phi_1 \\ \phi_2 \end{pmatrix}_1(y) \right\rangle \\
&= \langle \phi_1(x), \phi_1(y) \rangle + \langle \phi_2(x), \phi_2(y) \rangle \\
&= K_1(x, y) + K_2(x, y)
\end{aligned}$$

*Part 2*

$$\begin{aligned}
K(x, y) &= \phi(x)^T \phi(y) \\
&= \phi_1(x)^T \phi_1(y) \phi_2(x)^T \phi_2(y) \\
&= (\phi_1(x) \phi_2(x))^T (\phi_1(y) \phi_2(y)) \\
&= K_1(x, y) K_2(x, y)
\end{aligned}$$

## Problem 21

*Part 1*

There are  $2^m$  possible substrings  $s$ , so we define the dimension to be  $\mu = 2^m$

$$\phi(x)_\mu = \sum_{i=1}^{\mu} \mathbb{I}_{s_i \text{ substring of } x}$$

$$\langle \phi(x)_\mu, \phi_\mu(y) \rangle = K(x, y)$$

$$\begin{aligned}
&= \sum_{i=1}^{2^m} \mathbb{I}_{s_i \in x} \mathbb{I}_{s_i \in y} \\
&= \sum_{i=1}^{2^m} \mathbb{I}_{s_i \in x, y} \\
&= \sum_{s_i \in 0, 1^m} \mathbb{I}_{s_i \in x, y}
\end{aligned}$$

*Part 2*

Let  $J$  be the set of all possible  $s$ , the magnitude of  $J$  is the dimension of the  $\mathcal{H}$  the hilbert space of this kernel function, that is the dimension of the Hilbert space of the string kernel is  $2^m$ , i.e.:

$$\phi(x) = \mathbb{R}^n \rightarrow \mathbb{R}^\mu$$