

CS725: Tutorial 6

1 Detecting spam mails

One of the fundamental tasks of machine learning is to detect spam e-mails. You are given some words and a label of +1 if it is spam or -1 if it is not. Here **1** indicates the presence¹ of word and **0** the absence of word. Assume the learning rate η is $\frac{1}{2}$. Find the separating hyperplane using perceptron training algorithm

	area	click	your	in	singles	y
a	1	1	0	1	1	+1
b	0	0	1	1	0	-1
c	0	1	1	0	0	+1
d	1	0	0	1	0	-1
e	1	0	1	0	1	+1
f	1	0	1	1	0	-1

Solution: Since this is a programming exercise. I would like you to share and discuss solutions to this and evaluate others based on your personal solutions.

One possible solution

$$w_{click} = 1, w_{in} = -1, w_{singles} = 1$$

2 Computing power of perceptrons

Perceptrons can only separate linearly separable data as discussed in class. Given n variables we can have 2^{2^n} boolean functions, but not all of these can be represented by a perceptron. For example when $n=2$ the XOR and XNOR cannot be represented by a perceptron. Given n boolean variables how many of 2^{2^n} boolean functions can be represented by a perceptron?

Solution: <http://unbc.arcabc.ca/islandora/object/unbc\%3A6871/datastream/PDF/view>

¹<https://preview.overleaf.com/public/vgbycngdqhgc/images/a9c18fe31ba566c1dc8ecd306bd0463d880f856b.jpeg>

3 Kernel Perceptron

Recall the proof for convergence of the perceptron update algorithm. Now can this proof be extended to the kernel perceptron?

Recall that Kernelized perceptron² is specified as:

$$f(x) = \text{sign} \left(\sum_i \alpha_i^* y_i K(x, x_i) + b^* \right)$$

The perceptron update algorithm for the Kernelized version is:

- INITIALIZE: $\alpha = \text{zeros}()$
- REPEAT: for $\langle x_i, y_i \rangle$
 - If $\text{sign} \left(\sum_j \alpha_j y_j K(x_j, x_j) + b \right) \neq y_i$
 - then, $\alpha_j = \alpha_j + 1$
 - endif

Solution: Yes, in fact kernel perceptron can be derived from the perceptron update rule as follows:

$$f(x) = \text{sign} \left((w^*)^T \phi(x) \right) = \text{sign} \left(\sum_i \alpha_i^* y_i K(x, x_i) + b^* \right)$$

- INITIALIZE: $w = [0, 0, \dots, 0, 1] \Rightarrow f(x) = \text{sign} \left((w)^T \phi(x) \right) = \text{sign} \left(\sum_i \alpha_i y_i K(x, x_i) + b \right)$

with $\alpha_i = 0$ and $b = 1$

Note: $\phi^T(\hat{x})\phi(x)\hat{y} = \hat{y}K(\hat{x}, x) + \hat{y}$

- REPEAT: for each $\langle \hat{x}, \hat{y} \rangle$
 - If $\hat{y}w^T\phi(\hat{x}) < 0$
 - $\Rightarrow f(\hat{x}) = \text{sign} \left((w)^T \phi(\hat{x}) \right) = \text{sign} \left(\sum_i \alpha_i y_i K(\hat{x}, x_i) + b \right) \neq \hat{y}$
 - then, $w' = w + \Phi(\hat{x})\hat{y}$
 - $\Rightarrow f(x) = \text{sign} \left((w')^T \phi(x) \right) = \text{sign} \left(\sum_i (\alpha_i y_i K(x, x_i) + \phi^T(\hat{x})\phi(x)\hat{y}) + b \right)$
 - $= \text{sign} \left(\sum_i \alpha'_i y_i K(x, x_i) + b' \right)$ where $\alpha'_i = \alpha_i$ for all i except that $\alpha'_x = \alpha_x + 1$ and $b' = b + \hat{y}$
 - endif

$$\text{Thus, } f(x) = \text{sign} \left((w^*)^T \phi(x) \right) = \text{sign} \left(\sum_i^* \alpha_i y_i K(x, x_i) \right)$$

²In the original tutorial problem, b was missing. Re-introducing b helps state the equivalence of kernel perceptron to regular perceptron more easily.

4 Number of iterations for convergence of perceptron update

Prove the following:

If $\|\mathbf{w}^*\| = 1$ and if there exists $\theta > 0$ such that for all $i = 1, \dots, n$, $y_i(\mathbf{w}^*)^T \phi(\mathbf{x}_i) \geq \theta$ and $\|\phi(\mathbf{x}_i)\|^2 \leq \Gamma^2$ then the perceptron algorithm will make at most $\frac{\Gamma^2}{\theta^2}$ errors (that is take at most $\frac{\Gamma^2}{\theta^2}$ iterations to converge)

Solution:

We know that $\|\mathbf{w}^*\|_2^2 = 1$ and $y_i \phi(\mathbf{x}_i) \mathbf{w}^* \geq \theta$ for all i . We assume that $\mathbf{w}^{(0)} = 0$

Now consider $(\mathbf{w}^*)^T \mathbf{w}^{(k)} = (\mathbf{w}^*)^T (\mathbf{w}^{(k-1)} + y_i \phi(\mathbf{x}_i)) = (\mathbf{w}^*)^T \mathbf{w}^{(k-1)} + y_i (\mathbf{w}^*)^T \phi(\mathbf{x}_i) \geq (\mathbf{w}^*)^T \mathbf{w}^{(k-1)} + \theta \geq (\mathbf{w}^*)^T \mathbf{w}^{(k-2)} + 2\theta \geq (\mathbf{w}^*)^T \mathbf{w}^{(0)} + k\theta = k\theta$

Thus,

$$(\mathbf{w}^*)^T \mathbf{w}^{(k)} \geq k\theta$$

and because

$$\|\mathbf{w}^*\| \|\mathbf{w}^{(k)}\| = \|\mathbf{w}^{(k)}\| \geq |(\mathbf{w}^*)^T \mathbf{w}^{(k)}|$$

we must have

$$\|\mathbf{w}^{(k)}\| \geq k\theta$$

Similarly,

$$\|\mathbf{w}^{(k)}\|_2^2 = \|\mathbf{w}^{(k-1)} + y_i \phi(\mathbf{x}_i)\|_2^2 = \|\mathbf{w}^{(k-1)}\|_2^2 + y_i^2 \|\phi(\mathbf{x}_i)\|_2^2 + 2y_i (\mathbf{w}^{(k-1)})^T \phi(\mathbf{x}_i) < \|\mathbf{w}^{(k-1)}\|_2^2 + \Gamma^2 < \|\mathbf{w}^{(k-2)}\|_2^2 + 2\Gamma^2 < \|\mathbf{w}^{(0)}\|_2^2 + k\Gamma^2 = k\Gamma^2$$

since $y_i^2 = 1$ and it must have been that (as per perceptron update rule)

$$y_i (\mathbf{w}^{(k-1)})^T \phi(\mathbf{x}_i) < 0$$

Thus,

$$\|\mathbf{w}^{(k)}\|_2^2 < k\Gamma^2$$

and

$$\|\mathbf{w}^{(k)}\|_2^2 \geq k^2 \theta^2$$

which implies

$$k^2 \theta^2 < k\Gamma$$

that is,

$$k < \frac{\Gamma}{\theta^2}$$

which proves our claim.

<http://www.cs.columbia.edu/~mcollins/courses/6998-2012/notes/perc.converge.pdf>