

Blatt 6Abgabe bis Montag, 1. Juni 2008, 13 Uhr bei Mikio Braun (mikio@cs.tu-berlin.de).

Ausarbeitung schriftlich abgeben, Programm per Email.

Aufgaben

1. (10 Punkte) Zeige, dass die duale Darstellung der One-Class-SVM mit Soft Margin (vgl. Seite 26 der Folien) tatsächlich die Gestalt hat

$$\max_{\alpha} \sum_{i=1}^n \alpha_i k(x_i, x_i) - \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j)$$

so dass $\sum_{i=1}^n \alpha_i = 1$ and $0 \leq \alpha_i \leq C$ for $i = 1, \dots, n$,

wobei $\{x_1, \dots, x_n\}$ Trainingsdaten sind und k einer beliebigen Kernfunktion entspricht.

2. (5 Punkte) Gib die allgemeine duale quadratische Darstellung der One-Class SVM in Matrix-/Vektornotation in der Form

$$\max_{\alpha} c^T \alpha + \alpha^T B \alpha$$

so dass $d^T \alpha = e$,

$$l_i \leq \alpha_i \leq u_i, \text{ für } 1 \leq i \leq n.$$

3. (7 Punkte) Ergänze im Programmskelett die Funktion `oneclass`, die gegeben die Kernmatrix $K \in \mathbb{R}^{n \times n}$ die Kernmatrix und die Regularisierungskonstante $C \geq 0$, den Vektor $\alpha \in \mathbb{R}^n$ von Kernkoeffizienten aus der dualen Darstellung berechnet. Verwende zur Optimierung die Funktion `pr_loqo2`.
4. (8 Punkte). Finde mit der One-Class SVM Hacker-Angriffe in einem vorverarbeiteten Datensatz von Netzwerkverkehr. Du findest den Datensatz auf der Veranstaltungsseite. Der Datensatz ist in eine Trainings- und Testpartition unterteilt. Wie in einem realistischen Einsatz enthalten beide Partitionen Angriffe.

Verwende einen linearen Kern. Die Regularisierungskonstante C kann durch Raten, Probieren oder durch Kreuzvalidierung ermittelt werden. Gib das von Dir gewählte C explizit an, und dokumentiere, wie Du C gewählt hast. Übergib eine Datei mit den Positionen der Angriffe des Testdatensatzes.

Hinweise: Der Datensatz wurde aus echten Verbindungen und Angriffen des HTTP-Protokolls generiert. Jede Verbindung x wurde über 3-Gramme in einen Vektor $\phi(x)$ abgebildet. Die Daten sind normalisiert, d.h. $\|\phi(x)\|_1 = 1$.

Das Programmskelett enthält keine Kreuzvalidierung oder ähnliches. Dies müsste ggf. im oberen Teil des Programmskelettes eingefügt werden.

Wie in der Vorlesung besprochen, wird die Abweichung eines Punktes z vom erlernten Model α ermittelt durch

$$a(z) = k(z, z) - 2 \sum_{i=1}^n \alpha_i k(x_i, z) + \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j).$$

Die Funktion `compute_scores` berechnet diese Größe.


```
% ...  
  
% 4. Compute the outlier scores given  
%  
%   KR: kernel matrix on training data  
%   KS: kernel matrix on test data  
%   KSR: kernel matrix on test data / training data  
%   alpha: learned kernel coefficients  
function scores = compute_scores(KS, KSR, KR, alpha)  
% ...
```
