

**Blatt 2**

Abgabe **verlängert bis zum 15. Mai 2007** in der Vorlesung, oder bis 23:59:59 Uhr bei  
mikio@cs.tu-berlin.de

Die Aufgaben können auch in Gruppen bearbeitet werden, allerdings sollte die Gruppenzusammensetzung über des Semester stabil bleiben. Für praktische Aufgaben bitte ebenfalls Code und Ausgabe (d.h. Ergebnisse, Plots) abgeben. Verwende Matlab oder Octave. Lehrversionen von Matlab sind im cs-Netz (IRB) installiert. Zum Start `ml/bin/matlab` aufrufen. Octave ist ein freie Matlab-clone, der unter [www.octave.org](http://www.octave.org) verfügbar ist. Für die Abgabe von praktischen Aufgaben bitte die Coding-Richtlinien beachten, die für das Praktikum gelten (siehe hierzu [http://ml.cs.tu-berlin.de/de/ss07\\_ml\\_praktikum.html](http://ml.cs.tu-berlin.de/de/ss07_ml_praktikum.html), Hinweise auf dem ersten Übungsblatt).

**Aufgaben**

In diesem Übungsblatt wollen wir die Verbindungen zwischen Boosting und Gradientenabstieg ausnutzen, um Boosting-Methoden auf Regressionsprobleme zu erweitern, d.h. für die Beobachtungen  $(\mathbf{x}_i, y_i)$  gilt  $\mathbf{x}_i \in \mathbb{R}^d$  und  $y_i \in \mathbb{R}$ . In diesem Fall verwendet man gewöhnlich die quadratische Verlustfunktion

$$L_2(y, y') = \frac{1}{2} (y - y')^2.$$

(Der Faktor 1/2 ist irrelevant für die Minimierung des quadratischen Verlusts, vereinfacht aber die Notation im Gradientenverfahren.)

Im Gegensatz zur Klassifikation gibt es in Regressionsproblemen nur eine informelle Definition von schwachen Lernern. Ein schwacher Lerner ist ein Verfahren, das „underfittet“. Schwache Lerner sind beispielsweise Ridge-Regression oder *penalized least squares* (vgl. Blatt 9 des Wintersemesters) mit hohen Werten für den Regularisierungsparameter  $\lambda$ . Einem weiteren schwachen Lerner werden wir in Aufgabe 3 begegnen.

1.  **$L_2$ -Gradientenabstieg (5 Punkte)** Zeige, dass der Gradientenabstieg für die quadratische Verlustfunktion äquivalent zu einem iterativen Lernen auf den Residuen ist, d.h. der negative Gradient  $u_i$  aus Gleichung (2) von Blatt 1 entspricht den Residuen

$$r_i = y_i - f_m(\mathbf{x}_i).$$

*Bemerkung:* Für die meisten schwachen Lerner kann man zeigen, dass die optimale Schrittweite  $\alpha_m = 1$  ist. Eine Variante von  $L_2$ -Boosting besteht darin, die Schrittweite vorher festzulegen (s. Aufgabe 3).

2. **Linearität von  $L_2$ -Boosting (10 Punkte)** Im letzten Wintersemester (Blatt 9) haben wir den Begriff „linearer Lerner“ kennengelernt. „Linear“ bedeutet in diesem Fall, dass man den Lerner in der Form

$$\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$$

schreiben kann, wobei  $\mathbf{y}$  der Vektor ist, der alle Label des Trainingsdatensatzes enthält,  $\hat{\mathbf{y}}$  die Vorhersage auf dem Trainingsdatensatz ist, und  $\mathbf{H}$  eine Matrix ist, die von  $\mathbf{X}$ , aber nicht von  $\mathbf{y}$  abhängt.

Zeige, dass  $L_2$ -Boosting ein lineares Lernverfahren ist, falls die einzelnen schwachen Lerner linear sind. Zeige genauer: Falls in jeder Iteration  $m$  ein linearer (schwacher) Lerner  $\mathbf{H}_m$  angewandt wird, so wird  $L_2$ -Boosting durch die Matrix

$$\mathbf{B}_M = \sum_{m=1}^M \mathbf{H}_m \prod_{m'=1}^{m-1} (\mathbf{I} - \mathbf{H}_{m'})$$

repräsentiert.

Gehe hierbei wie folgt vor:

- (a) Zeige, daß der Residuenvektor im  $m$ ten Schritt  $\mathbf{r}_m$  (d.h. der Vektor  $(r_1, \dots, r_n)$  im  $m$ ten Schritt) folgende Rekursionsformel erfüllt:

$$\mathbf{r}_m = \mathbf{r}_{m-1} - \mathbf{H}_{m-1}\mathbf{r}_{m-1} = (\mathbf{I} - \mathbf{H}_{m-1})\mathbf{r}_{m-1}$$

Verwende hierzu, daß

$$\mathbf{r}_m = \mathbf{y} - \mathbf{f}_{m-1}, \quad \mathbf{f}_m = \mathbf{f}_{m-1} + \mathbf{h}_m, \quad \mathbf{h}_m = \mathbf{H}_m\mathbf{r}_m, \quad \mathbf{f}_0 = \mathbf{0}.$$

wobei  $\mathbf{h}_m = (h_m(\mathbf{x}_1), \dots, h_m(\mathbf{x}_n))$ , und entsprechend  $\mathbf{f}_m$  der Vektor aller Vorhersagen auf dem Trainingsdatensatz ist.

- (b) Zeige nun, daß

$$\mathbf{f}_M = \mathbf{B}_M\mathbf{y} = \sum_{m=1}^M \mathbf{H}_m \prod_{m'=1}^{m-1} (\mathbf{I} - \mathbf{H}_{m'})\mathbf{y}. \quad (1)$$

- (c) (*Zusatzaufgabe*) Zeige, daß sich die Matrix  $\mathbf{B}_M$  auch schreiben läßt als

$$\mathbf{B}_M = \mathbf{I}_n - (\mathbf{I}_n - \mathbf{H}_M)(\mathbf{I}_n - \mathbf{H}_{M-1}) \dots (\mathbf{I}_n - \mathbf{H}_1)$$

*Bemerkung:* Es ist also möglich, die Freiheitsgrade (vgl. Blatt 9 des Wintersemesters) von  $L_2$ -Boosting zu definieren.

3. **Implementierung (10 Punkte)** Implementiere den  $L_2$ -Algorithmus. Verwende als schwachen Lerner *komponentenweise Kleinst-Quadrate-Regression*: Wähle in jeder Boosting-Iteration die optimale Variable  $X_{j^*}$  aus. Für diese Variable ist  $h_m$  das Ergebnis der Kleinst-Quadrate-Regression

$$\gamma_0 + \gamma_1 X_{j^*}$$

angewandt auf  $\mathbf{x}_i$  und die Residuen  $r_i$  (siehe unten). Dabei ist die Variable  $X_{j^*}$  optimal, falls sie von allen  $d$  Variablen den Trainingsfehler am meisten reduziert. Wähle als konstante Schrittlänge  $\alpha_m = 0.1$ .

*Bemerkung:* Der Vorteil dieses schwachen Lernalgorithmus ist, dass i.A. nur eine kleine Anzahl der  $d$  Variablen in der endgültigen Lösung auftauchen,  $L_2$ -Boosting liefert also i.A. dünn besetzte Lösungen.

4. **Anwendung (5 Punkte)** Wende  $L_2$ -Boosting auf den `prostate`-Datensatz an, den man unter

<http://www-stat.stanford.edu/~tibs/ElemStatLearn/>

findet. Verwende dabei nur die Trainingsdaten. Plote die Anzahl der Variablen, die in der Boosting-Lösung auftauchen, als Funktion der Anzahl der Boosting-Iterationen. Bestimme die optimale Anzahl der Boosting-Iterationen, indem Du den Testfehler für alle Boostingschritte bestimmst. Plote den Trainings- und Testfehler als Funktion der Anzahl der Boosting-Iterationen. Gib die ausgewählten Variablen für die optimale Anzahl von Boosting-Iterationen aus.

## Hinweise

Wir betrachten ein eindimensionales Regressionsproblem

$$Y = \gamma_0 + \gamma_1 X + \text{Fehler}$$

d.h.  $\mathbf{x}_i \in \mathbb{R}$ . Der Kleinst-Quadrate-Schätzer ist die Lösung des Minimierungsproblems

$$\min_{\gamma_0, \gamma_1} \sum_{i=1}^n (y_i - \gamma_0 - \gamma_1 \mathbf{x}_i)^2.$$

Die Lösung lautet

$$\begin{aligned} \hat{\gamma}_1 &= \frac{\text{cov}(\mathbf{X}, \mathbf{y})}{\text{var}(\mathbf{X})} \\ \hat{\gamma}_0 &= \bar{y} - \hat{\gamma}_1 \bar{\mathbf{X}}. \end{aligned}$$