

# Oppgaver, STK4040, uke 44

30. oktober 2005

## Oppgave 1: MSE

La  $f$  være en regresjon trent på et datasett  $\mathbf{X}$  og  $\mathbf{y}$ , og la  $\hat{y} = f(\mathbf{x})$  være predikert respons for en gitt ny observasjon  $\mathbf{x}$ . Vis at

$$\text{MSE}(\hat{y}) = \text{Bias}^2(\hat{y}) + \text{Var}(\hat{y} - Y), \quad (1)$$

der  $Y$  er (den ukjente) responsverdien tilhørende  $\mathbf{x}$ . (Vi regner  $\mathbf{X}$  og  $\mathbf{x}$  som gitte.) Merk at dette gjelder generelt for regresjoner; vi har ikke antatt noe om formen til  $f$  (f.eks. at det er en lineær regresjon).

## Oppgave 2: leverage

Anta modellen

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (2)$$

hvor  $\mathbf{X}$  er kjent og  $\boldsymbol{\varepsilon} \sim N_n(0, \sigma^2\mathbf{I})$ .

Gitt et datasett  $\mathbf{X}$  og  $\mathbf{y}$ . La  $e_i = \hat{y}_i - y_i$  være residualene fra en lineær regresjon av  $\mathbf{y}$  på  $\mathbf{X}$ . La  $h_i$  være diagonalelementene til hattmatrisen  $\mathbf{H} = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t$ .  $h_i$  kalles «leverage»n til  $\mathbf{x}_i$ , og er et mål på hvor mye  $\mathbf{x}_i$  potensielt kan påvirke regresjonen. Merk at  $h_i = \mathbf{x}_i^t(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{x}_i$ .

a)

La  $\mathbf{X}_{(i)}$  og  $\mathbf{y}_{(i)}$  være  $\mathbf{X}$  og  $\mathbf{y}$  uten den  $i$ te observasjonen, og la  $\mathbf{b}_{(i)}$  være estimatet av  $\boldsymbol{\beta}$  fra  $\mathbf{X}_{(i)}$  og  $\mathbf{y}_{(i)}$ . Vis at

$$e_{(i)} = y_i - \mathbf{x}_i^t\mathbf{b}_{(i)} = e_i/(1 - h_i). \quad (3)$$

(Med andre ord: man kan beregne residualene fra en full kryssvalidering («leave one out cross-validation») av en lineærregresjon uten å måtte gjøre selve regresjonen om igjen for hver observasjon.)

Tips: for å spare litt regning, kan dere bruke at

$$(\mathbf{X}_{(i)}^t\mathbf{X}_{(i)})^{-1} = (\mathbf{X}^t\mathbf{X})^{-1} + \frac{(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{x}_i\mathbf{x}_i^t(\mathbf{X}^t\mathbf{X})^{-1}}{1 - h_i}. \quad (4)$$

(Ekstraoppgave: Bruk at  $\mathbf{X}_{(i)}^t\mathbf{X}_{(i)} = \mathbf{X}^t\mathbf{X} - \mathbf{x}_i\mathbf{x}_i^t$  til å vise (4).)

b)

Anta at  $X = (\mathbf{1} \ X_1)$ . La  $X_c$  være  $X_1$  sentrert, og la  $\mathbf{y} = \beta_0 \mathbf{1} + X_c \boldsymbol{\beta}_1 + \boldsymbol{\varepsilon}$  være den ekvivalente modellen. Kall leverage'ne til denne regresjonen for  $g_i$ , og vis at

$$h_i = g_i + 1/n, \quad (5)$$

hvor  $n$  er antall observasjoner.