

## Bayes-päätely, 6. harjoitukset (1.–2.3.2017)

1. Lineaarinen regressio. Tarkastellaan aineistoa<sup>1</sup> `diabetes.csv`, joka sisältää 442 diabetespotilaan iän, sukupuolen, painoindeksin, verenpaineen ja mittaukset kuudesta eri veren (ehkä kolesteroli?)-arvosta (S1-S6) ja selitettävän muuttujan  $y$ , joka kuvaa taudin etenemistä vuosi mitausten tekemisen jälkeen. Tarkoituksena on siis selittää, mitkä riskitekijät vaikuttavat diabeteksen pahenemiseen. Käytä tämän ja seuraavien tehtävien analyyseihin tiedostoa `diabetes_sd.csv`, joka sisältää saman aineiston, mutta normalisoituna siten, että jokaisen muuttujan keskiarvo on 0, ja kaikkien selittävien muuttujien euklidinen normi on 1.

- (a) Suorita tavallinen lineaarinen regressio R:n `lm`-funktiolla, siten että selität muuttujaa  $y$  kaikilla aineiston muilla muuttujilla. Mitkä selittäjistä ovat tilastollisesti merkitseviä merkitsevyytasolla 0.05?
- (b) Suorita vastaava bayesiläinen lineaarinen regressio stanilla. Malli on tarkemmin

$$Y_i | \alpha, \beta, \sigma \sim N(\alpha + \beta \mathbf{x}_i, \sigma^2) \quad \text{kaikille } i = 1, \dots, N,$$

tai matriisimuodossa

$$\mathbf{Y} | \alpha, \beta, \sigma \sim N_n(\alpha \mathbf{1}_n + \beta \mathbf{X}, \sigma^2 \mathbf{I}),$$

missä

$$\mathbf{X} = \begin{bmatrix} x_{11} & \dots & x_{1K} \\ \vdots & & \vdots \\ x_{N1} & \dots & x_{NK} \end{bmatrix}$$

on matriisi, joka sisältää  $K = 10$  selittävän muuttujan arvot  $N = 442$  havainnolle,

$$\beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_K \end{bmatrix}$$

on vektori joka sisältää regressiokerrointen arvot,  $\alpha$  on regression vakiotermin, ja  $\sigma$  virhetermien keskihajonta. Voit käyttää epäoleellista prioria

$$p(\alpha, \beta, \sigma) \propto 1, \quad \sigma > 0,$$

jolloin prioria ei tarvitse määrittää stanille erikseen (mutta muista rajoittaa  $\sigma$  positiiviseksi!).

Jos määrität selittävät muuttujat matriisina: `matrix[N,K] x`; ja selitettävän muuttujan (ja regressiokertoimet) vektorina: `vector[N] y`; voit määrittää mallin kätevästi matriisimuodossa: `y ~ normal(x * beta + alpha, sigma)`;

2. Jatkoa edelliseen tehtävään.

- (a) Piirrä kuva regressiokerrointen  $\beta_1, \dots, \beta_K$  mediaaneista ja 50% ja 95% uskottavuusväleistä (Tämän pitäisi onnistua automaattisesti `plot(stan_fit)`-komennolla. Piirrettäviä parametreja voi säätää `par`-argumentilla, ja ulompaa uskottavuusväliä `outer_level`-argumentilla), missä `stan_fit` on `stan`-funktion palauttama olio. Uskottavuusvälit (credible interval), eli ”Bayesiläiset luottamusvälit” ovat Bayes-päätelyn vastine frekventistisille luottamusväleille. Tulkinallisesti ne ovat huomattavasti miellyttävämpiä, sillä nyt voit sanoa, että parametrin todellinen arvo on 95% todennäköisyydellä sen 95% uskottavuusvälillä!

---

<sup>1</sup>Tätä on käytetty esimerkkiaineistona artikkeleissa Least Angle Regression [1] ja Bayesian Lasso [2]

- (b) Piirrä vastaava kuva, johon merkitset (a)-kohdassa estimoimasi regressiokerrointen  $\beta_1, \dots, \beta_K$  suurimman uskottavuuden estimaatit  $\hat{\beta}_1, \dots, \hat{\beta}_K$ , ja regressiokerrointen luottamustasojen  $\alpha = 0.5$  ja  $\alpha = 0.05$  (frekventistiset) luottamusvälit. Luottamustason  $\alpha$  frekventistinen luottamusväli regressiokertoimelle  $\beta_k$  ovat muotoa

$$\left( \hat{\beta}_k - t_{N-K-1}(\alpha/2) \cdot \text{se}(\hat{\beta}_k), \hat{\beta}_k + t_{N-K-1}(\alpha/2) \cdot \text{se}(\hat{\beta}_k) \right),$$

missä  $t_\nu(u)$  on  $t$ :n jakauman vapausasteella  $\nu$  yläkvantiili, eli piste, jonka oikealle puolelle jää osuus  $u$  jakauman todenäköisyysmassasta,<sup>2</sup> ja  $\text{se}(\hat{\beta}_j)$  on estimaattorin  $\beta_j$  keskivirhe<sup>3</sup>. Lisäksi voit hyödyntää `segments`-funktia.

Vertaa kuvia. Millä jakaumatuloksilla selittäisit havaintojasi?

- 3.** Jatkoa edellisille tehtäville. Tehdään samalle aineistolle Lasso, eli L1-regularisoitu lineaarinen regressio. Tämä on toteutettu valmiiksi R:n `glmnet`-kirjastossa.

- (a) Lue "Quick start"-osio osoitteesta [https://web.stanford.edu/~hastie/glmnet/glmnet\\_alpha.html](https://web.stanford.edu/~hastie/glmnet/glmnet_alpha.html), ja tee Lasso diabetes-aineistolle siten, että selität jälleen  $y$ :tä kaikilla muilla muuttujilla. (Vihje: Tämän pitäisi onnistua suoraan komennolla `glmnet(x, y)`, missä  $y$  on selitettävä muuttuja, ja  $x$  matriisi, joka sisältää selittäjät.) Piirrä lasso-polut (tämän pitäisi onnistua suoraan komennolla `plot(lasso_fit)`, missä `lasso_fit` on `glmnet`:in palauttama olio), eli regressiokerrointen arvot niiden yhteenlaskettujen pituuksien funktiona. Katso saadko samanlaiset lasso-polut, kuin artikkelin [1] kuvan 1 vasemmanpuolimmaisessa paneelissa!
- (b) Etsi 'rankaisu-parametrin'  $\lambda$  optimaalinen arvo ristiinvalidoinnin avulla (Vihje: voit käyttää valmista funktiota `cv.glmnet`. Tarkemmat ohjeet löytyvät yllä linkitetystä tutoriaalista). Mitkä ovat regressiokerrointen arvot tällä  $\lambda$ :n arvolla?

- 4.** Jatkoa edellisille tehtäville. Tehtävän 2 Bayesiläinen lineaarinen regressio, mutta tällä kertaa epäinformatiivisen priorin sijasta oletetaan, että osa regressiokertoimista on lähellä nollaa, eli asetetaan regressiokertoimille informatiivinen Laplace-priori, joka kutistaa regressiokertoimia kohti nollaa Lasso-tyylisesti. Tehdään täysi Bayesiläinen malli antamalla myös parametrin  $\lambda$  neliölle Gamma-priori. Malli on siis muuten sama kuin tehtävässä 2, mutta

$$\begin{aligned} \beta_k | \sigma, \lambda &\sim \text{Laplace}(0, \sigma/\lambda) \quad \text{kaikille } k = 1, \dots, K, \\ \lambda^2 &\sim \text{Gamma}(1, 1.78). \end{aligned}$$

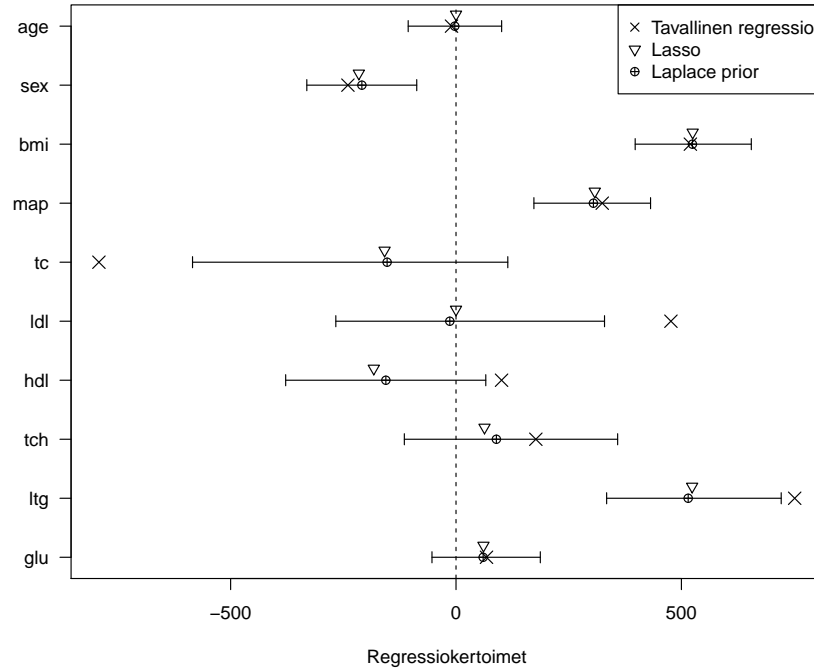
Keskihajonnalle voit olettaa edelleen epäinformatiivisen priorin  $p(\sigma) \propto 1, \sigma > 0$ .

- (a) Sovita malli stanilla (Vihje Laplace-jakauman saat komennolla `double_exponential(0, sigma / lambda)`). Piirrä histogrammi parametrin  $\lambda$  posteriorijakaumasta. Vertaa tulosta tehtävän 3 ristiinvalidoinnilla saatuun parametrin  $\lambda$  arvoon.
- (b) Koetetaan replikoida artikkelin [1] kuvaa 2! Piirrä parametrien  $\beta_1, \dots, \beta_{10}$  posteriorijakaumien mediaanit ja 95% luottamusvälit. Lisää kuvaan tehtävässä 1 laskemasi suurimman uskottavuuden estimaatit  $\hat{\beta}_{\text{MLE}}$ , ja tehtävässä 3 laskemasi lasso-estimaatit  $\hat{\beta}_{\text{Lasso}}$  ristiinvalidoimalla valitsemallasi parametrin arvolla  $\lambda$ . Lopputulos voisi muistuttaa kuvaa 1.
- (c) Vertaa tehtävän 1 tuloksiin. Mille parametreille 95%:n uskottavuusväli ei sisällä nollaa?

<sup>2</sup>saadaan R:ssä funktiolla `qt(u, nu, lower = FALSE)`.

<sup>3</sup>Löytyvät esimerkiksi komennon `summary(lm_fit)$coefficients`, missä `lm_fit` on `lm`-funktion palauttama olio, palauttamasta matriisista.

Kuva 1: Regressiokerrointen  $\beta$  95% mediaanit ja uskottavuusvälit. Lisäksi myös suurimman uskottavuuden estimaatit  $\hat{\beta}_{MLE}$  ja Lasso-estimaatit  $\hat{\beta}_{Lasso}$ .



## Viitteet

- [1] Bradley Efron, Trevor Hastie, Iain Johnstone, Robert Tibshirani, et al. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- [2] Trevor Park and George Casella. The bayesian lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.