

## Tentamen - Bayesiansk statistik, 732G43, 7.5 hp

**Betygsgränser: Tentamen omfattar totalt 25 poäng. Slutbetyget på kursen baseras på en sammanvägning av tentamen (50%) och inlämningsuppgifter (50%).**

**För godkänt betyg på tentamen krävs minst 15 poäng.**

**Redovisa och motivera tydligt alla dina svar!**

1. (9 poäng)

- (a) Förklara på vilket sätt Bayes sats ligger till grund för Bayesiansk statistik genom att skriva om Bayes sats för händelserna  $A$  och  $B$  nedan till Bayes sats för en parameter  $\theta$  givet data  $x_1, \dots, x_n$  och förklara vad uttrycken i denna omskrivning av Bayes sats innebär i ord.

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}.$$

- (b) Redogör för den huvudsakliga skillnaden i inferens mellan frekventistisk och Bayesiansk inferens.
- (c) Förklara varför likelihoodfunktionen inte är en täthetsfunktion för en parameter  $\theta$  och vad likelihoodfunktionen behöver kombineras med i Bayesiansk statistik för att erhålla en giltig täthetsfunktion för  $\theta$ .
- (d) Motivera vad som menas med en informativ priorfördelning samt ge **ett** exempel på en informativ priorfördelning för ett medelvärde  $\theta$ .
- (e) Motivera vad som menas med en konjugerad priorfördelning samt ange **en** fördel med en konjugerad priorfördelning för en parameter  $\theta$  i en valfri modell för data.
- (f) Beskriv vad priorelicitering innebär för en parameter  $\theta$  samt ge **ett** exempel på hur priorelicitering kan användas utifrån experthjälp om en andel  $\theta$ .
- (g) Motivera när kvadratisk approximation av aposteriorifördelningen för  $\theta_1, \theta_2, \dots, \theta_{20}$  kan vara ett bättre samt ett sämre alternativ än MCMC skattning av aposteriorifördelningen.
- (h) Beskriv vad ett 95 % kredibilitetsintervall innebär för en parameter  $\theta$  och redogör för hur detta intervall skiljer sig från ett 95 % konfidensintervall för  $\theta$ .
- (i) Redogör för hur man kan skapa ett 95 % kredibilitetsintervall för medelvärdet  $\frac{\alpha}{\beta}$  i en gammafördelning om man känner till aposteriorifördelningen för  $\alpha$  och  $\beta$ .

2. (5 poäng) Antag följande modell för Bayesiansk linjär regression:

$$Y_1, \dots, Y_n | \mu, \sigma^2, \mathbf{x} \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

$$\mu = \beta \mathbf{x}',$$

där variansen  $\sigma^2$  är okänd och med vektorn av förklaringsvariabler  $\mathbf{x} = (1 \ x_1 \ \dots \ x_k)$  och vektorn med parametrar  $\beta = (\beta_0 \ \beta_1 \ \dots \ \beta_k)$ .

- (a) Antag en uniform prior för parametrarna  $(\boldsymbol{\beta}, \ln \sigma)$ :

$$p(\boldsymbol{\beta}, \sigma^2 | \mathbf{x}) \propto \frac{1}{\sigma^2}.$$

Vad är det för fördel med att använda sig av denna uniforma prior för  $(\boldsymbol{\beta}, \ln \sigma)$ ? Motivera.

- (b) Oavsett prior, beskriv i ord hur man kan beräkna ett 95.2 % kredibilitetsintervall för  $\mu$  utifrån samplade dragningar från posteriorfördelningen av  $\boldsymbol{\beta}$ .
- (c) Oavsett prior, beskriv i ord hur man kan skapa 90.9 % kredibilitetsintervall för  $y$  som funktion av en dummyvariabel  $x_j$  utifrån samplade dragningar från posteriorfördelningen av  $(\boldsymbol{\beta}, \sigma)$ .
- (d) Beskriv vad som skiljer en Bayesianisk förklaringsgrad från en frekventistisk förklaringsgrad för multipel linjär regression samt ange **en** fördel med den Bayesianiska förklaringsgraden i Gelman et al (2017) jämfört med andra alternativ för Bayesianisk förklaringsgrad.

3. (7 poäng)

- (a) Förklara vad det innebär att en modell överanpassar respektive underanpassar data.
- (b) Motivera hur regulariserande priors kan användas för att minska överanpassning av data. Motivera även hur regulariserande priors kan ge en underanpassning av data.
- (c) Förklara vad Akaikevikter kan användas till samt ange två informationskriterier som kan användas för att beräkna Akaikevikter och förklara vad som skiljer dessa informationskriterier åt.
- (d) Vad kan traceplottar över parameterdragningarna från olika MCMC kedjor användas till? Hur kan en plott för ackumulerade posterior medelvärden av parametrarna över MCMC dragningar användas i MCMC diagnostik? Motivera.
- (e) Förklara vad det är för skillnad mellan Bayesianisk logistisk regression och Bayesianisk binomialregression och hur dessa regressionsmodeller relaterar till varandra.

4. (4 poäng) Antag följande Bayesianiska multilevelmodell med 2 nivåer för individ  $i$  som tillhör en skolklass  $j = 1, \dots, J$ :

$$\begin{aligned} y_i &\stackrel{iid}{\sim} N(\mu, \sigma_y) \\ \mu &= \alpha_j + \boldsymbol{\beta} \mathbf{x}', \\ \alpha_j &\sim N(\alpha, \sigma_\alpha), \\ \alpha &\sim N(0, 10), \\ \sigma_\alpha &\sim \text{halfcauchy}(0, 2), \\ \beta_k &\sim N(0, 10), k = 1, \dots, p \\ \sigma_y &\sim \text{halfcauchy}(0, 2), \end{aligned}$$

där  $\mathbf{x} = (x_1 \dots x_p)$  är en radvektor med förklaringsvariabler och  $\boldsymbol{\beta} = (\beta_1 \dots \beta_p)$  är en radvektor med parametrar.

- (a) Beskriv i ord hur man kan beräkna ett 95 % kredibilitetsintervall för  $\mu$  för en individ  $i$  som tillhör en ny skolklass  $J+1$  (som inte fanns med då man anpassade modellen) utifrån samplade dragningar från posteriorfördelningen av parametrarna i modellen.
- (b) Antag att vi förändrar modellen ovan genom att sätta  $\alpha = 0$  och  $\sigma_\alpha = 10$ . Skriv upp den förändrade modellen i sin helhet och redogör för skillnaden mellan denna modell och den ursprungliga modellen ovan med avseende på poolning av information mellan grupper.
- (c) Förklara vad som menas med fullständig poolning av intercept i en regressionsmodell och hur modellen ovan hade förändrats om fullständig poolning av intercept hade varit gällande.