

# BAYESIÁNSKÁ ANALÝZA – CVIČENÍ 1

## 10.10.2011 – 14.10.2011

Toto cvičení je založeno na znalosti prvních tří kapitol z učebnice Koop (2003): *Bayesian econometrics*, případně na odpovídajících kapitolách podkladového učebního textu *Bayesiánská analýza*.

### Co bude náplní cvičení?

- ✎ Tvorba umělého souboru dat a jeho využití při zkoumání vlastností odhadových metod a postupů.
- ✎ Odhad a posteriorní analýza normálního lineárního regresního modelu s přirozeně konjugovanou apriorní hustotou (jedna i více vysvětlujících proměnných).
- ✎ Citlivostní analýza volby apriorní hustoty pravděpodobnosti.
- ✎ Odhad a posteriorní analýza na příkladech s využitím reálných dat.

### Zadání příkladů

K řešení příkladů využijte již hotové funkce, případně si vytvořte své vlastní. Poslední příklady jsou z knížky Hill, Griffiths, Lim (2008): *Principles of Econometrics* a mohou sloužit jako inspirace pro problémy řešené ve vašich semestrálních projektech. Hovoří-li se zde o testování hypotéz, má se za to, že tento test provedeme za pomoci porovnávání modelů.

1. *Empirická ilustrace z druhé kapitoly Koop (2003)*. Projděte si řešení příkladu a diskutujte nejasnosti.
2. *Empirická ilustrace z třetí kapitoly Koop (2003)*, data o prodeji domů (soubor `HPRICE.TXT`). Projděte si řešení příkladu a podpůrné funkce a diskutujte nejasnosti. Rozšířte analýzu o další vysvětlující proměnné a řešte příklad s případnými vlastními apriorními představami o hodnotách apriorních hyperparametrů.
3. (*Tvorba umělých datových souborů*) Vytvořit si vlastní datový soubor je užitečné pro pochopení vlastností modelu a pro analýzu kvality počítačových algoritmů. Pokud si totiž zvolíme hodnoty parametrů, víme jaké hodnoty by nám měly přibližně vycházet, použijeme-li tu či onu ekonometrickou metodu.
  - (a) Vygenerujte umělý datový soubor v rámci normálního lineárního regresního modelu:
    - Zvolte si hodnoty  $\beta$ ,  $h$  a  $N$  (např.  $\beta = 2$ ,  $h = 1$  a  $N = 100$ ).
    - Vygenerujte  $N$  hodnot pro vysvětlující proměnnou z rozdělení dle vašeho výběru (např. proveďte  $N = 100$  náhodných výběrů z uniformního rozdělení  $U(0, 1)$ ).

- Vygenerujte  $N$  chybových členů pomocí  $N$  i.i.d. výběrů z rozdělení  $N(0, h^{-1})$ .
  - Vytvořte vektor vysvětlované proměnné  $y$  v rámci NLRM (tj.  $y_i = \beta x_i + \epsilon_i$  pro  $i = 1, \dots, N$ ).
- (b) Vytvořte bodový graf (graf XY) každého datového souboru pro zobrazení toho, jak se váš výběr  $\beta$ ,  $h$  a  $N$  promítne do generovaných dat.

4. (Bayesiánská analýza v NLRM – citlivostní analýza prioru)

- (a) Vytvořte umělý datový soubor pro  $\beta = 2$ ,  $h = 1$  a  $N = 100$  použitím uniformního rozdělení pro generování vysvětlující proměnné.
- (b) Předpokládejme apriorní hustotu v podobě  $\beta, h \sim NG(\underline{\beta}, \underline{V}, \underline{s}^{-2}, \underline{\nu})$ , kde  $\underline{\beta} = 2$ ,  $\underline{V} = 1$ ,  $\underline{s}^{-2} = 1$ ,  $\underline{\nu} = 10$ , a spočítejte posteriorní střední hodnoty a směrodatné odchylky pro  $\beta$  a  $h$ . Spočítejte Bayesův faktor porovnávající model, kde  $\beta = 0$ , s modelem, kde  $\beta \neq 0$ . Spočítejte střední hodnotu a směrodatnou odchylku prediktivní posteriorní hustoty pro  $x^* = 0.5$ .
- (c) Jak se změní odpovědi na otázku (4b) pokud  $\underline{V} = 0.01$ ? Jak v případě  $\underline{V} = 0.1$ ,  $\underline{V} = 10$ ,  $\underline{V} = 100$  a  $\underline{V} = 1000000$ ?
- (d) Jak se změní odpovědi na otázku (4b) pokud  $\underline{\nu} = 5$ ? Jak v případě  $\underline{\nu} = 100$ ,  $\underline{\nu} = 1000$  a  $\underline{\nu} = 1000000$ ?
- (e) Nastavte apriorní střední hodnotu  $\beta$  na hodnotu odlišnou od své skutečné hodnoty (např.  $\underline{\beta} = 0$ ) a zopakujte část (4c).
- (f) Nastavte apriorní střední hodnotu  $h$  na hodnotu odlišnou od své skutečné hodnoty (např.  $\underline{s}^{-2} = 100$ ) a zopakujte část (4d).
- (g) S ohledem na výsledky dosažené v částech (4b)–(4f) diskutujte citlivost posteriorní střední hodnoty, směrodatné odchylky a Bayesova faktoru na volbu apriorní hustoty.
- (h) Zopakujte části (4a)–(4g) za použití více informativní datové sady (např.  $N = 1000$ ) a méně informativní sady (např.  $N = 10$ ).
- (i) Zopakujte části (4a)–(4h) za použití jiných hodnot  $\beta$  a  $h$  pro generování umělého souboru dat.

5. Odhad a Monte Carlo integrace v modelu vícenásobné regrese.

- (a) Vytvořte umělý datový soubor pro velikosti  $N = 100$  pro normální lineární regresní model s úroňovou konstantou a jednou vysvětlující proměnnou. Úroňovou konstantu položte rovnu 0 a koeficient sklonu regresní přímky položte roven jedné a  $h = 1$ . Vysvětlující proměnnou vezměte z uniformního rozdělení  $U(0, 1)$ .
- (b) Spočítejte posteriorní střední hodnotu a směrodatnou odchylku pro tato data při použití přirozeně kojugované normální-gama apriorní hustoty s  $\underline{\beta} = (0, 1)'$ ,  $\underline{V} = I_2$ ,  $\underline{s}^{-2} = 1$  a  $\underline{\nu} = 10$ .

- (c) Vykreslete posteriorní hustotu pro  $\beta_2$ , a to jak z definice její posteriorní marginální hustoty, tak pomocí Monte Carlo integrace (užijte histogram nebo funkci `ksdensity` pro vykreslení jádrové hustoty vašeho výběru). Pro různě velké velikosti výběru spočítejte numerickou standardní chybu aproximace střední hodnoty parametrů.
- (d) Spočítejte Bayesův faktor porovnávající model  $M_1 : \beta_2 = 0$  s  $M_2 : \beta_2 \neq 0$ .
- (e) Vykreslete predikční hustotu pro pozorování s hodnotou  $x_2^* = 0.5$ .
- (f) Proveďte citlivostní analýzu apriorní hustoty nastavením  $V = cI_2$  a opakujte kroky (5b), (5d) a (5e) pro hodnoty  $c = 0.01, 1.0, 100.0, 1 \times 10^6$ . Diskutujte citlivost posteriorní hustoty, Bayesova faktoru a predikční hustoty rozdělení.
- (g) Spočítejte posteriorní střední hodnotu a směrodatnou odchylku vektoru parametrů  $\beta$  za použití neinformativního prioru.
- (h) Spočítejte 99% HPDI pro  $\beta_2$  užitím neinformativního prioru a užit jej pro ověření hypotézy, že  $\beta_2 = 0$ . Porovnejte své výsledky s výsledky dosaženými v části (5d).

**Poznámka:** Můžete samozřejmě rozšířit model o další proměnné (při generování umělých dat) a měnit jejich nastavení včetně volby priorů..

6. Soubor `cocaine.m` obsahuje 56 pozorování proměnných vztahujících se k prodeji kokainu v severovýchodní Kalifornii v období 1984-1991. Data jsou podmnožinou dat použitých ve studii Culkins, J.P. a Padman, R. (1993): „Quantity Discounts and Quality Premia for Illicit Drugs,“ *Journal of the American Statistical Association*, 88, 748-757. Proměnné jsou

- *price* = cena za gram kokainu v rámci dané transakce;
- *quant* = počet gramů kokainu prodaných v dané transakci;
- *qual* = kvalita kokainu vyjádřená jako procento čistoty;
- *trend* = časová proměnná s hodnotami od 1984=1 až po 1991=8.

Předpokládejme regresní model

$$price = \beta_0 + \beta_1 quant + \beta_2 qual + \beta_3 trend + \epsilon.$$

- (a) Jaká znaménka koeficientů byste očekávali u parametrů  $\beta_1$ ,  $\beta_2$  a  $\beta_3$ ?
- (b) Odhadněte daný model (předpokládáme, že se jedná o NLRM s přirozeně konjugovanou apriorní hustotou). Zvolte si vhodné hyperparametry dle vašich zkušeností. Jsou znaménka parametrů v souladu s vašim očekáváním?
- (c) Říká se, že čím větší objem obchodů, tím větší riziko, že vás dostihne ruka zákona. Prodejci tak jsou ochotni akceptovat nižší cenu, pokud prodávají větší množství. Pokuste se testovat tuto hypotézu.
- (d) Ověřte hypotézu, že kvalita kokainu nemá vliv na jeho cenu.

- (e) Jaká je průměrná roční změna ceny kokainu? Zamyslete se nad tím, proč by se měla cena takto měnit.
7. Každé ráno mezi 6:30 a 8:00 opouští Bill Melbournské předměstí Carnegie, aby se dostal do práce na University of Melbourne. Čas, který Bill stráví cestou do práce,  $time$ , závisí na času odjezdu,  $depart$ , počtu červených světel na semaforech,  $reds$  a počtu vlaků, kvůli kterým musí čekat na Murrumbeenském přejezdu,  $trains$ . Pozorování těchto proměných je celkem získáno za 231 pracovních dní v roce 2006 a jsou obsahem souboru `commute.m`. Proměnná  $time$  je měřena v minutách,  $depart$  je počet minut po 6:30, které uplynou než Bill vyrazí z domu.
- (a) Odhadněte rovnici

$$time = \beta_0 + \beta_1 depart + \beta_2 reds + \beta_3 trains + \epsilon.$$

- (b) Jaká znaménka koeficientů byste očekávali u parametrů  $\beta_1$ ,  $\beta_2$  a  $\beta_3$ ?
- (c) Otestujte hypotézu, že každé červené světlo zpozdí Billa nejméně o 2 minuty.
- (d) Testujte hypotézu, že čas odjezdu nemá vliv na čas strávený cestováním.
- (e) Otestujte hypotézu, čas cestování navíc díky čekání na jednom semaforu je stejný jako čas čekání průjezdu jednoho vlaku.