

STATISTIKA NONPARAMETRIK: PENGGABUNGAN DIAGRAM POHON POLYA YANG BERHINGGA

Asri Ode Samura

Tadris Matematika, IAIN Ternate, Asri22samura@gmail.com

Abstrak

Sebuah pendekatan tentang nonparametrik bayes, penggabungan diagram pohon polya (MPT). MPT menggunakan sebuah partisi dari dukungan kepadatan distribusi aslinya. Kepadatan umumnya mempertahankan bentuk distribusi asli pada tiap-tiap rangkaian partisi dan menambahkan parameter baru dari peluang-peluang bersyarat. Model diagram pohon polya dapat diaplikasikan secara luas dan mudah diprogramkan dengan diberikannya skema MCMC untuk penyesuaian model parametrik aslinya.

Kata Kunci. Bayes, Model gabungan linier tergeneralisasi, MGLT

1. Pendahuluan

Analisis Bayes membutuhkan probabilitas prior untuk parameter-parameter distribusi sampling. Akan bisa sulit memilih prior yang masuk akal dalam masalah-masalah dengan banyak parameter, sehingga parameter-parameter tersebut disini dipecah menjadi dua kelompok yang tepat: (1) yang berhubungan dengan keluarga parametrik asli, dan (2) yang berhubungan dengan generalisasi keluarga asli tersebut. Metode standard dapat digunakan untuk membentuk sebuah prior bagi parameter asli; lihat Bedrick, Christensen, dan Johnson (1996) untuk sebuah pendekatan dalam model linier tergeneralisasi. Prior acuan yang tepat ada untuk parameter-parameter dari suatu generalisasi. Diagram pohon Polya adalah sebuah distribusi probabilitas random. Bersyarat pada sebuah anggota himpunan parametrik asli yang ditetapkan (contohnya distribusi normal dengan rata-rata μ dan ragam σ^2), himpunan distribusi tergeneralisasi bersama dengan prior acuan merupakan diagram pohon Polya yang berhingga. Integrasi prior terhadap parameter himpunan asli (misalnya μ dan σ^2) membentuk sebuah gabungan diagram pohon Polya.

MPT sudah sangat luas digunakan akan tetapi pada contoh kita akan melibatkan penggunaan model gabungan linier tergeneralisasi (GLMM), lihat Breslow dan Clayton (1993), dan keduanya melibatkan generalisasi distribusi normal yang biasanya diasumsikan bersama dengan keduanya. GLMM memberikan sebuah

kerangka kerja yang sangat dikenal untuk analisis ukuran-ukuran longitudinal dan data berkelompok yang muncul dalam banyak bidang seperti pertanian, biologi, epidemiologi, ekonomi, dan geofisika. Model-model tersebut menghitung korelasi antar pengamatan berkelompok dengan memasukkan efek random ke dalam komponen prediktor linier dari model tersebut. Meskipun penyesuaian GLMM secara khusus bersifat kompleks, intersep random standar dan model intersep/slop random dengan efek random berdistribusi normal sekarang dapat secara rutin disesuaikan dalam paket-paket software komersil seperti SAS dan Stata. Model-model tersebut cukup fleksibel dalam mengakomodasi perilaku heterogen, akan tetapi mereka memiliki kekurangan yang sama yaitu kurangnya ketahanan terhadap tujuan dari asumsi distribusi seperti model statistika lain yang berdasarkan distribusi Gauss.

Suatu strategi untuk menjaga terhadap asumsi normalitas yang tidak sesuai adalah dengan memasukkan asumsi distribusi yang lebih fleksibel untuk efek random ke dalam model. Sehingga tampaknya diperlukan pengembangan nonparametrik dari GLMM parametrik. Kita menggambarkan perluasan tersebut beserta akibat dari penggunaan yang salah untuk asumsi model tradisional dalam GLMM dengan menggunakan pada contoh kehidupan nyata kita. Demi tujuan perbandingan, kita menyesuaikan semua model berdasarkan kedua asumsi efek random berdistribusi normal dan generalisasi.

2. Landasan Teori

2.1. Theorema Bayes

Misal $y' = (y_1, \dots, y_n)$ adalah sebuah vektor dari n pengamatan yang distribusi peluangnya $p(y|\theta)$ bergantung pada nilai k parameter $\theta' = (\theta_1, \dots, \theta_k)$. Dimisalkan bahwa θ sendiri memiliki distribusi peluang $p(\theta)$. Maka

$$P(y|\theta)P(\theta) = P(y, \theta) = P(\theta|y)P(y) \quad (2.1)$$

Dengan diberikan data pengamatan y , distribusi bersyarat dari θ menjadi

$$P(\theta|y) = \frac{P(y|\theta)P(\theta)}{P(y)}. \quad (2.2)$$

Maka

$$P(y) = EP(y|\theta) = C^{-1} = \begin{cases} \int P(y|\theta)P(\theta)d\theta, & \text{jika } \theta \text{ kontinu} \\ \sum P(y|\theta)P(\theta), & \text{jika } \theta \text{ diskrit} \end{cases} \quad (2.3)$$

dimana sigma atau integralnya diambil pada rentang dari θ yang dapat diterima, dan dimana $E[f(\theta)]$ adalah ekspektasi matematika dari $f(\theta)$ sehubungan dengan distribusi $P(\theta)$. Oleh karena itu kita dapat menuliskan sebagai berikut

$$P(\theta|y) = cP(y|\theta)P(\theta). \quad (2.4)$$

2.2. Distribusi Posterior

Distribusi posterior diberikan x (atau disingkat posterior) akan dinotasikan $\delta(\theta|x)$, dan didefinisikan sebagai distribusi bersyarat dari θ , diberikan x sampel pengamatan. Dengan memperhatikan bahwa θ dan X memiliki (subjektif) kerapatan

$$h(x|\theta) = \delta(\theta)f(x|\theta), \tag{2.5}$$

dan bahwa X telah marjinal (tanpa syarat) kerapatan

$$m(x) = \begin{cases} \int f(x|\theta)\delta(\theta)d\theta, & \text{jika } \theta \text{ kontinu} \\ \sum_{\theta} f(x|\theta)\delta(\theta), & \text{jika } \theta \text{ diskrit} \end{cases} \tag{2.6}$$

jelas bahwa (dengan memberikan $m(x) \neq 0$)

$$\delta(\theta|x) = \frac{h(x|\theta)}{m(x)}. \tag{2.7}$$

dengan menentukan estimator bayes, maka persamaan yang dipakai

$$E(\delta(\theta|x)) = \begin{cases} \int \theta f(x|\theta)\delta(\theta)d\theta & \text{jika } \theta \text{ kontinu} \\ \sum_{\theta} \theta f(x|\theta)\delta(\theta) & \text{jika } \theta \text{ diskrit} \end{cases} \tag{2.8}$$

Adapun dapat menentukan estimator dengan cara Box-Tiao. Dari persamaan (2.1) dapat tulis dalam bentuk

$$\delta(\theta|x) \propto f(x|\theta)\delta(\theta) \tag{2.9}$$

Konstanta kesebandingan dipilih sedemikian hingga $\delta(\theta|x)$ merupakan densitas.

Teorema : 2.1 : *Jika x_1, x_2, \dots, x_n adalah sampel random dari distribusi $N(\mu, \sigma^{-1})$ dengan presisi σ , ($\sigma > 0$) diketahui, misalkan distribusi prior untuk rata-rata μ adalah distribusi $N(\mu_0, (n_0\sigma^{-1}))$, maka distribusi posterior μ bersyarat bilamana diberikan nilai $X_i = x_i$, untuk $i = 1, 2, \dots, n$ adalah berdistribusi $N(\mu_n, \sigma_n^{-1})$ dimana*

$$\sigma_n = (n_0 + n)\sigma \quad \text{dan} \quad \mu_n = \frac{n\bar{x} + n_0\mu_0}{n_0 + n} \tag{2.10}$$

Bukti:

Diberikan $X_i = x_i$, untuk $i = 1, 2, \dots, n$ karena σ diketahui, dengan mensubstitusi fungsi likelihood $L(\mu|x) \propto \text{Exp}\{-n\sigma(\mu - \bar{x})^2/2\}$ dan fungsi densitas prior μ pada $\mu(\mu|\sigma) \propto \text{Exp}\{-n_0\sigma(\mu - \mu_0)^2/2\}$ kedalam persamaan

$$\delta(\mu|x) \propto L(\mu|x)\delta(\mu) \tag{2.11}$$

Sehingga diperoleh

$$\begin{aligned} \delta(\mu|x) &\propto L(\mu|x)\delta(\mu) \\ &\propto \text{Exp}\left\{-\frac{1}{2}[n\sigma(\mu - \bar{x})^2 + n_0\sigma(\mu - \mu_0)^2]\right\} \end{aligned} \tag{2.12}$$

Ekspresikan $n\sigma(\mu - \bar{x})^2 + n_0\sigma(\mu - \mu_0)^2$ pada ruas kanan (2.12) diuraikan menjadi

$$n\ddot{e}(\hat{i} - \bar{x})^2 + n_0\ddot{e}(\hat{i} - \hat{i}_0)^2 = (n_0 + n)\ddot{e}(\hat{i} - \hat{i}_n) + \frac{\ddot{e}nn_0(\bar{x}-\hat{i}_0)^2}{n_0+n} \quad (2.13)$$

Dimana \hat{i}_n adalah sebagaimana pada persamaan (2.10), karena suku terakhir pada ruas kanan persamaan (2.13) adalah konstan terhadap \hat{i} , maka suku tersebut dapat diabaikan, dan setelah itu hasilnya disubstitusikan kembali kedalam persamaan (2.12) sehingga diperoleh fungsi densitas posterior

$$\delta(\hat{i}|\underline{x}) \propto \text{Exp}\left\{-\left(\frac{1}{2}\right)n\ddot{e}(\hat{i} - \hat{i}_n)^2\right\} \quad (2.14)$$

Yang merupakan fungsi densitas dari distribusi normal $N(\hat{i}_n, \ddot{e}_n^{-1})$ dimana \hat{i}_n dan \ddot{e}_n adalah seperti pada persamaan (2.10) jadi teorema terbukti. ■

Teorema: 2.2 :Jika x_1, x_2, \dots, x_n adalah sampel random dari distribusi $N(\hat{i}, \ddot{e}^{-1})$; untuk \hat{i} dan \ddot{e} kedua-duanya tidak diketahui. Jika distribusi prior bersama dari (\hat{i}, \ddot{e}) adalah distribusi $NG(\hat{i}, \ddot{e}; \hat{i}_0, \ddot{e}_0, \acute{a}_0, \hat{a}_0)$; $-\infty < \hat{i}_0 < \infty, n_0 > 0, \ddot{e} > 0, \acute{a} > 0, \hat{a} > 0$, maka distribusi posterior bersama $(\hat{i}, \ddot{e}|\underline{x})$ adalah distribusi $NG(\hat{i}, \ddot{e}|\hat{i}_n, \ddot{e}_n, \acute{a}_n, \hat{a}_n)$ dimana

$$\hat{i}_n = (n\bar{x} + n_0\hat{i}_0)/(n_0 + n) \text{ dan } \ddot{e}_n = (n_0 + n)\ddot{e} \quad (2.15)$$

Dan

$$\acute{a}_n = \acute{a}_0 + n/2 \text{ dan } \hat{a}_n = \hat{a}_0 + \frac{1}{2}\sum_{i=1}^n(x_i - \bar{x})^2 + \frac{nn_0(\bar{x}-\hat{i}_0)^2}{2(n_0+n)} \quad (2.16)$$

Bukti:

Diberikan $X_i = x_i; i = 1, 2, 3, \dots, n$, karena \hat{i} dan \ddot{e} kedua-duanya tidak diketahui, disubstitusi fungsi likelihood

$$L(\hat{i}, \ddot{e}|\underline{x}) \propto \ddot{e}^{n/2} \text{Exp}\left\{-\left(\ddot{e}/2\right)\sum_{i=1}^n(x_i - \bar{x})^2\right\} \text{exp}\left\{-n\ddot{e}(\hat{i} - \bar{x})^2/2\right\} \quad (2.17)$$

Dan fungsi densitas prior normal gamma pada

$$\begin{aligned} \delta(\hat{i}, \ddot{e}|n_0, n_1, \dots, n_n) &= \delta(\hat{i}, \ddot{e}|n_0, \hat{i}_0, \acute{a}_0, \hat{a}_0) \\ &\propto \left\{\ddot{e}^{1/2} \text{Exp}\left\{-\frac{n_0\ddot{e}}{2}(\hat{i} - \hat{i}_0)^2\right\}\right\} \left(\ddot{e}^{\acute{a}_0-1} \text{Exp}\{-\hat{a}_0\ddot{e}\}\right) \end{aligned} \quad (2.18)$$

Kedalam persamaan $\delta(\ddot{e}|\underline{x}) \propto L(\ddot{e}|\underline{x})\delta(\ddot{e})$ dimana $\ddot{e} = (\ddot{e}_1, \ddot{e}_2) = (\hat{i}, \ddot{e})$, sehingga diperoleh

$$\begin{aligned} \delta(\hat{i}, \ddot{e}|\underline{x}) &\propto \ddot{e}^{(n/2)+\acute{a}_0-1} \ddot{e}^{1/2} \text{Exp}(-\hat{a}_0\ddot{e}) \\ &\cdot \text{Exp}\left\{-\left(\frac{\ddot{e}}{2}\right)\left[\sum_{i=1}^n(x_i - \bar{x})^2 + n(\hat{i} - \bar{x})^2 + n_0(\hat{i} - \hat{i}_0)^2\right]\right\} \end{aligned} \quad (2.19)$$

Karena $\sum_{i=1}^n(x_i - \bar{x})^2 + n(\hat{i} - \bar{x})^2 + n_0(\hat{i} - \hat{i}_0)^2$ dapat diuraikan menjadi

$$(n_0 + n)(\hat{i} - \hat{i}_n)^2 + \sum_{i=1}^n(x_i - \bar{x})^2 + \frac{nn_0(\bar{x}-\hat{i}_0)^2}{n_0+n} \quad (2.20)$$

Dengan \hat{i}_n seperti pada persamaan (2.15) dan persamaan (2.19) sehingga menjadi

$$\delta(\hat{i}, \ddot{e}|\underline{x}) \propto \ddot{e}^{(n/2)+\acute{a}_0-1} \ddot{e}^{1/2}$$

$$\begin{aligned} & \cdot \text{Exp} \left\{ -\frac{\ddot{e}}{2} (n_0 + n) (\hat{i} - \hat{i}_n)^2 + \sum_{i=1}^n (x_i - \bar{x})^2 + \frac{nn_0(\bar{x} - \hat{i}_0)^2}{n_0 + n} \right\} \\ & \cdot \text{Exp}(-\hat{a}_0 \ddot{e}) \propto \ddot{e} \text{Exp} \left\{ -\ddot{e}^{1/2} (n_0 + n) (\hat{i} - \hat{i}_n)^2 / 2 \right\} \\ & \cdot \ddot{e}^{(n/2) + \hat{a}_0 - 1} \text{Exp} \left\{ -\left[\hat{a}_0 + \frac{1}{2} \sum_{i=1}^n (x_i - \bar{x})^2 + \frac{1}{2} \frac{nn_0(\bar{x} + \hat{i}_0)^2}{n_0 + n} \right] \ddot{e} \right\} \end{aligned}$$

Atau setara dengan

$$\delta(\hat{i}, \ddot{e} | \underline{x}) \propto \left\{ \ddot{e}^{1/2} \text{Exp}[-\ddot{e}_n (\hat{i} - \hat{i}_n)^2 / 2] \right\} \left(\ddot{e}^{\hat{a}_n - 1} \text{Exp}(-\hat{a}_n \ddot{e}) \right) \quad (2.21)$$

Ini adalah fungsi densitas dari distribusi normal gamma $NG(\hat{i}, \ddot{e}; \hat{i}_n, \ddot{e}_n, \hat{a}_n, \hat{a}_n)$ dengan $\hat{i}_n, \ddot{e}_n, \hat{a}_n, \hat{a}_n$ sebagaimana di dalam persamaan (2.15) dan (2.16) ■

Teorema; 2.5.3: Diberikan x_1, x_2, \dots, x_n dari distribusi $N(\hat{i}, \ddot{e}^{-1})$; dengan \hat{i} dan \ddot{e} kedua-duanya tidak diketahui. Jika digunakan distribusi prior-informatif $\delta(\hat{i}, \ddot{e}) = \ddot{e}^{-1}$ maka distribusi posterior bersama $(\hat{i}, \ddot{e} | \underline{x})$ adalah distribusi $NG(\hat{i}'_n, \ddot{e}'_n, \hat{a}'_n, \hat{a}'_n)$, dimana

$$\hat{i}'_n = \bar{x} \quad \text{dan} \quad \ddot{e}'_n = n\ddot{e} \quad (2.22)$$

Dan

$$\hat{a}'_n = (n - 1)/2 \quad \text{dan} \quad \hat{a}'_n = \sum_{i=1}^n (x_i - \bar{x})^2 / 2 \quad (2.23)$$

Bukti:

Diberikan $X_i = x_i$ dimana $i = 1, 2, \dots, n$, dengan \hat{i} dan \ddot{e} kedua-duanya tidak diketahui maka dapat kita pakai persamaan $\delta(\ddot{e} | \underline{x}) \propto L(\ddot{e} | \underline{x}) \delta(\ddot{e})$ dapat digunakan untuk $\underline{\hat{e}} = (\hat{e}_1, \hat{e}_2) = (\hat{i}, \ddot{e})$ kedalam persamaan $\delta(\hat{e} | \underline{x}) \propto L(\hat{e} | \underline{x}) \delta(\hat{e})$ ini disebut fungsi likelihood pada persamaan $L(\hat{i} | \underline{x}) \propto \text{Exp}\{-n\ddot{e}(\hat{i} - \bar{x})^2 / 2\}$, dan fungsi diensitas prior tak informatif pada persamaan $\delta(\hat{i}, \ddot{e}) = \ddot{e}^{-1}$ maka diperoleh

$$\begin{aligned} & \delta(\hat{i}, \ddot{e} | \underline{x}) \propto L(\hat{i}, \ddot{e} | \underline{x}) \delta(\hat{i}, \ddot{e}) \\ & \propto \ddot{e}^{-1} \ddot{e}^{n/2} \text{Exp} \left\{ -\left(\frac{\ddot{e}}{2}\right) \sum_{i=1}^n (x_i - \bar{x})^2 \right\} \frac{\text{Exp}\{-n\ddot{e}(\hat{i} - \bar{x})^2\}}{2} \\ & \propto n^{1/2} \ddot{e}^{1/2} \text{Exp}\{-n\ddot{e}(\hat{i} - \bar{x})^2 / 2\} \cdot ((n - 1)/2)^{-1} \\ & \cdot \text{Exp} \left\{ -(\ddot{e}/2) \sum_{i=1}^n (x_i - \bar{x})^2 \right\} \end{aligned}$$

$$\delta(\hat{i}, \ddot{e} | \underline{x}) \propto (n\ddot{e})^{1/2} \text{Exp}\{-n\ddot{e}(\hat{i} - \bar{x})^2 / 2\} \cdot \ddot{e}^{\hat{a}'_n - 1} \cdot \text{Exp}(-(\hat{a}'_n / 2) \ddot{e}) \quad (2.24)$$

Ini menunjukkan distribusi normal-Gamma $NG(\hat{i}, \ddot{e} | \hat{i}'_n, \ddot{e}'_n, \hat{a}'_n, \hat{a}'_n)$, dengan $\hat{i}'_n, \ddot{e}'_n, \hat{a}'_n, \hat{a}'_n$ sebagaimana pada persamaan (2.22) dan (2.23) ■

2.3. Pohon Polya

Definisi. Probabilitas random P pada R mempunyai distribusi pohon Polya dengan parameter (Π, \mathcal{A}) terdapat variabel random $Y = \{Y_0, Y_{00}, Y_{10}, \dots\}$ jika memenuhi:

1. Variabel-variabel random pada Y independen;
2. Untuk setiap $\epsilon \in E^*$, $Y_{\epsilon_0} \sim \text{Beta}(\acute{a}_{\epsilon_0}, \acute{a}_{\epsilon_1})$;
3. Untuk setiap $m = 1, 2, \dots$ dan untuk setiap $\epsilon \in E^m$, maka:

$$P(B_\epsilon) = \left[\prod_{\{\epsilon_j=0\}}^m Y_{\epsilon_1, \dots, \epsilon_j 0} \right] \left[\prod_{\{\epsilon_j=1\}}^m 1 - Y_{\epsilon_1, \dots, \epsilon_{j-1} 0} \right] \quad (2.25)$$

Yang dinotasikan

$$P \sim PT(\Pi, \mathcal{A}). \quad (2.26)$$

Varibel random W_1, \dots, W_n dikatakan sampel berukuran n dari P jika:

$$\begin{aligned} P[W_i \in B_\epsilon] &= EP[W_i \in B_\epsilon | P] \\ &= E[P(B_\epsilon)] \\ &= E[P(B_{\epsilon_1})P(B_{\epsilon_1 \epsilon_2} | B_{\epsilon_1}) \dots P(B_\epsilon | B_{\epsilon_1 \dots \epsilon_{m-1}})] \\ &= \frac{\acute{a}_{\epsilon_1}}{\acute{a}_0 + \acute{a}_1} \dots \frac{\acute{a}_{\epsilon_1 \dots \epsilon_{m-1}}}{\acute{a}_{\epsilon_1 \dots \epsilon_{m-1} 0} + \acute{a}_{\epsilon_1 \dots \epsilon_{m-1} 1}} \end{aligned} \quad (2.27)$$

2.4. Rantai Markov Monte Carlo (MCMC).

Rantai Markov dapat digunakan untuk berbagai keperluan dalam perhitungan dan optimasi. MCMC algoritma yang paling dikenal adalah Gibbs sampler dan algoritma Metropolis-Hastings (M-H).

Rantai Markov

- a. Rantai markov waktu diskrit adalah proses stokastik dengan keadaan diskrit yang memenuhi:

$$P\{X_{n+1} = j | X_n = i, X_{n-1} = i_{n-1}, \dots, X_0 = i_0\} = P\{X_{n+1} = j | X_n = i\} \quad (2.28)$$

Sehingga

$$P_{ij} = P\{X_{n+1} = j | X_n = i\} \quad (2.29)$$

- b. Rantai markov waktu kontinu adalah proses stokastik yang memiliki distribusi bersyarat dari masa yang akan datang dengan $X(t+s)$, masa sekarang $X(s)$, dan masa lalu $X(u)$, $0 \leq u < s$, tergantung pada masa sekarang dan independen dari masa lalu, Yaitu:

$$\begin{aligned} P\{X(t+s) = j | X(s) = i, X(u) = x(u), 0 \leq u < s\} \\ = P\{X(t+s) = j | X(s) = i\} \end{aligned} \quad (2.30)$$

Dimana $s, t \geq 0$

Proses $\{X(t), t \geq 0\}$ memiliki probabilitas transisi stasioner atau homogen jika

$$P[X(t + s) = j | X(s) = i] = P_{ij}(t) \tag{2.31}$$

2.5. Metode Simulasi Monte Carlo

Simulasi Monte Carlo dapat diartikan sebagai metode simulasi statistik, di mana simulasi statistik didefinisikan sebagai istilah umum untuk semua metode simulasi yang menggunakan rangkaian bilangan acak. Simulasi Monte Carlo adalah suatu metode untuk mengevaluasi suatu model yang melibatkan bilangan acak sebagai salah satu input.

Simulasi Monte Carlo dapat diartikan juga sebagai metode untuk menganalisa perambatan ketidakpastian, di mana tujuannya adalah untuk menentukan bagaimana variasi acak atau error mempengaruhi sensitivitas, dari sistem yang sedang dimodelkan. Simulasi Monte Carlo digolongkan sebagai metode sampling karena input yang dibangkitkan secara acak dari suatu distribusi probabilitas untuk proses sampling dari suatu populasi nyata.

2.6. Gibbs Sampling

Gibbs sampling adalah algoritma yang didasarkan pada generasi yang berurutan dari kepadatan bersyarat penuh $p(\theta_i | \theta_{\neq i}, y)$, yaitu, kepadatan posterior pada elemen ke- i dari $\theta \doteq (\theta_1, \dots, \theta_d)'$, diberikan semua elemen lainnya, di mana elemen θ dapat berupa skalar atau sub-vektor.

Menghitung vektor dalam beberapa urutan sebagai vektor $1, 2, \dots, s$ dan mengidentifikasi vektor j dengan keadaan j pada rantai Markov dengan probabilitas transisi. Jika vektor-vektor i dan j berbeda lebih dari satu dalam komponen, dimana $p_{ij} = 0$. Jika mereka berbeda dalam maksimal satu komponen, misalkan, menjadi konkrit, berbeda dalam komponen pertama ($i \neq j$). Dengan vektor i adalah (y_1, y_2, \dots, y_k) dan vektor j sebagai (y_1^*, y_2, \dots, y_k) . Maka

$$\begin{aligned} p_{ij} &= \text{prob}(Y_1 = y_1^* | Y_2 = y_2, Y_3 = y_3, \dots, Y_k = y_k) \\ &= \frac{\text{Prob}(Y_1 = y_1^*, Y_2 = y_2, Y_3 = y_3, \dots, Y_k = y_k)}{\text{Prob}(Y_2 = y_2, Y_3 = y_3, \dots, Y_k = y_k)} \end{aligned} \tag{2.32}$$

Probabilities pada pembilang dan penyebut dihitung dengan menggunakan $P_Y(y)$, kami mengklaim bahwa rantai Markov tereduksi dan aperiodik, dan selanjutnya memiliki distribusi stasioner $P_Y(y)$.

3. Pembahasan

3.1. Penggabungan Diagram Pohon Polya yang Berhingga

Dengan menggunakan definisi yang relatif sederhana dari diagram pohon Polya, kita dapat menggambarkan proses generalisasi sebuah keluarga distribusi parametrik dengan menggunakan keluarga $N(\mu, \sigma^2)$. Keluarga parametrik lain digeneralisasi dengan cara yang sama.

Generalisasi tersebut melalui sejumlah tahapan, sebut saja J. Pada tiap tahapan kita memperkenalkan parameter baru untuk menggeneralisasi tahapan sebelumnya. Pada tahap pertama, kita membagi garis bilangan riil, yaitu pendukung distribusi normal, menjadi dua interval yang terbagi oleh median μ . Selanjutnya kita membiarkan perubahan yang terjadi dalam probabilitas menjadi di bawah atau di atas μ , tetapi kita mempertahankan bentuk kepadatan normal baik di bawah μ dan di atas μ .

Parameter baru yang ada pada tahap pertama adalah θ_{11} probabilitas yang tidak lebih besar dari μ , dan θ_{12} probabilitas yang berada di atas μ . Secara formal, diketahui X_1 memiliki distribusi tahap pertama, maka

$$\theta_{11} \equiv P[X_1 \leq i], \tag{3.1}$$

dan

$$\theta_{12} \equiv P[X_1 > i] = 1 - \theta_{11}. \tag{3.2}$$

Karena kita mempertahankan bentuk normal pada kedua himpunan, jika $a \leq \mu$ dan $Y \sim N(\mu, \sigma^2)$, secara kondisional kita akan mendapatkan

$$P[X_1 \leq a | X_1 \leq i] = \frac{P[Y \leq a]}{0.5}, \tag{3.3}$$

Dimana

$$\begin{aligned} P[Y \leq a] &= P\left(\frac{y - \mu}{\sigma} < \frac{a - \mu}{\sigma}\right) \\ &= P\left(Z < \left(\frac{a - \mu}{\sigma}\right)\right) \\ &= \Phi\left(\frac{a - \mu}{\sigma}\right) \end{aligned}$$

Sehingga

$$\begin{aligned} P[X_1 \leq a | X_1 \leq i] &= \frac{\Phi\left(\frac{a - \mu}{\sigma}\right)}{1/2} \\ &= 2\Phi[(a - i)/\sigma] \end{aligned} \tag{3.4}$$

dimana $\Phi(\cdot)$ merupakan fungsi kepadatan bersyarat (cdf) dari suatu normal standar. Begitu pula, jika $b > \mu$,

$$\begin{aligned} P[X_1 > b | X_1 > i] &= \frac{P[Y > b]}{1/2} \\ &= 2P[Y > b], \end{aligned} \tag{3.5}$$

Dengan

$$\begin{aligned} P[Y > b] &= P\left(\frac{y + i}{\sigma} > \frac{b - i}{\sigma}\right) \\ &= P\left((1 - Z)\left(\frac{b - i}{\sigma}\right)\right) \\ &= 1 - \Phi[(b - i)/\sigma] \end{aligned}$$

Maka

$$P[X_1 > b | X_1 > i] = 2P[Y > b] = 2\{1 - \Phi[(b - \mu)/\sigma]\} \quad (3.6)$$

Dengan cara lain kita bisa menuliskan

$$P[X_1 \leq a] = P[Y \leq a]2\epsilon_{11} \quad (3.7)$$

$$P[X_1 > b] = P[Y > b]2\epsilon_{12}. \quad (3.8)$$

Dengan $I_A(x)$ fungsi indikator dari A, kepadatan tahap distribusi 1 adalah

$$f(x_1 | i, \sigma^2, \epsilon_{11}, \epsilon_{12}) = \frac{1}{\sqrt{2\sigma^6}} e^{-(x_1 - i)^2 / 2\sigma^2} \times 2^1 [\epsilon_{11} I_{(-\infty, i]}(x_1) + \epsilon_{12} I_{(i, \infty)}(x_1)]. \quad (3.9)$$

Ditentukan X_2 memiliki distribusi tahap kedua, kita memperkenalkan parameter baru, $\theta_{21}, \theta_{22}, \theta_{23}, \theta_{24}$, yang didefinisikan sebagai peluang bersyarat relatif terhadap himpunan yang digunakan pada tahap 1:

$$\epsilon_{21} = P[X_2 \leq q_1 | X_2 \leq i] \quad (3.10)$$

$$\epsilon_{22} = P[q_1 < X_2 \leq i | X_2 \leq i] \quad (3.11)$$

$$\epsilon_{23} = P[i < X_2 \leq q_3 | X_2 > i] \quad (3.12)$$

$$\epsilon_{24} = P[q_3 < X_2 | X_2 > i]. \quad (3.13)$$

Catat bahwa $\theta_{21} = 1 - \theta_{22}$ dan $\theta_{23} = 1 - \theta_{24}$. Dengan tanpa syarat, keempat himpunan memiliki peluang

$$P[X_2 \leq q_1] = \epsilon_{11} \epsilon_{21} \quad (3.14)$$

$$P[q_1 < X_2 \leq i] = \epsilon_{11} \epsilon_{22} \quad (3.15)$$

$$P[i < X_2 \leq q_3] = \epsilon_{12} \epsilon_{23} \quad (3.16)$$

$$P[q_3 < X_2] = \epsilon_{12} \epsilon_{24}. \quad (3.17)$$

Dalam tiap himpunan, kita menggunakan bentuk kepadatan normal asli sehingga, sebagai contoh, jika $\mu < a < b \leq q_3$ dan $Y \sim N(\mu, \sigma^2)$,

$$\begin{aligned} & P[a < X_2 \leq b] \\ &= P[a < X_2 \leq b | i < X_2 \leq q_3] P[i < X_2 \leq q_3] \\ &= P[a < Y \leq b | i < Y \leq q_3] P[i < X_2 \leq q_3] \\ &= \frac{P[a < Y \leq b]}{P[i < Y \leq q_3]} P[i < X_2 \leq q_3] \\ &= P[a < Y \leq b] \frac{\epsilon_{12} \epsilon_{23}}{0.25}. \end{aligned} \quad (3.18)$$

Pada umumnya, kepadatan pada distribusi tingkat 2 adalah

$$f(x_2 | i, \sigma^2, \epsilon_{11}, \epsilon_{12}, \epsilon_{21}, \epsilon_{22}, \epsilon_{23}, \epsilon_{24}) = \frac{1}{\sqrt{2\sigma^6}} e^{-(x_2 - i)^2 / 2\sigma^2} 2^2 \left[\epsilon_{11} \epsilon_{21} I_{(-\infty, q_1]}(x_2) + \epsilon_{11} \epsilon_{22} I_{(q_1, i]}(x_2) + \epsilon_{12} \epsilon_{23} I_{(i, q_3]}(x_2) + \epsilon_{12} \epsilon_{24} I_{(q_3, \infty)}(x_2) \right] \quad (3.19)$$

Untuk menampilkan analisis Bayes dengan distribusi sampling ini, kita perlu sebuah distribusi prior gabungan pada parameter θ_{js} . θ_{js} sangat mudah

diinterpretasikan, informasi prior yang sangat berarti bisa muncul pada mereka. Sebuah kasus ekstrim memilih $\theta_{11} = \theta_{12} = 0,5$ dengan peluang prior 1. Hal ini memberikan peluang prior satu untuk median menjadi μ untuk suatu J . Namun, terdapat terlalu banyak sekali parameter untuk memilih sebuah distribusi yang menggambarkan informasi prior yang berarti pada semua θ_{js} , sehingga prior acuan juga dimasukkan. Secara khusus, informasi prior yang berarti akan dibatasi untuk parameter dari beberapa tahap j yang pertama.

Dengan selalu berhubungan, distribusi prior dan posterior yang fokus pada probabilitas yang tinggi di area sekitar $\theta_{js} = 0,5$ untuk semua js , akan berperilaku seperti distribusi normal. Hal ini terjadi ketika c besar dalam $\alpha_{js} = cp(j)$. Dengan $q(j)$ yang meningkat, nilai j yang besar mengimplikasikan bahwa probabilitas prior yang tinggi diletakkan pada θ_{js} di dekat 0,5. Hal ini adalah sifat yang telah disebutkan dalam pendahuluan yang memperbolehkan jumlah parameter yang besar sehubungan dengan jumlah pengamatan independennya.

Di lain pihak, ketika c kecil, distribusi tersebut lebih “nonparametrik”. Ditentukan A_j merupakan sebuah himpunan dalam pembagian level ke- J . Ketika c kecil, sebuah pengamatan dalam A_j memiliki pengaruh yang besar pada semua distribusi beta posterior dari θ_{js} sehubungan dengan A_j , sehingga menyebabkan probabilitas yang tinggi untuk A_j dalam distribusi posterior. Karena A_j adalah sebuah himpunan dalam pembagian terbaik yang ditentukan, sehingga menyebabkan perilaku yang menonjol yang diperkirakan diskrit dalam posterior.

Distribusi sampling tergeneralisasi tahap ke- J , sebut G , tergantung pada θ_{js} . G bersama dengan prior pada θ_{js} , menentukan diagram pohon Polya yang berhingga, dimana kita tulis

$$G \sim PT_j(c, \tilde{n}, N(\hat{i}, \hat{\sigma}^2)). \tag{3.20}$$

Sebuah prior pada (μ, σ) mengimplikasikan bahwa median μ , kuartil, oktil, dan sebagainya, adalah random. Hal ini memiliki pengaruh pada proses smoothing.

Kepadatan sampling random yang diperoleh dengan membangkitkan himpunan $\{\theta_{js}\}$ berdasarkan prior mereka, tetapi dirata-ratakan atas sebuah prior pada (μ, σ) diistilahkan sebagai sebuah gabungan diagram pohon Polya, dan ditulis

$$G \sim \int PT_j(c, \tilde{n}, N(\hat{i}, \hat{\sigma}^2))p(d\hat{i}, d\hat{\sigma}^2). \tag{3.21}$$

Hanson (2006) menunjukkan bahwa untuk prior khusus, kepadatan MPT randomnya halus. Diagram pohon Polya dan versi lain dari penggabungan diagram pohon Polya tidak perlu memiliki sifat ini; lihat Barron et al. (1999), Paddock (1999), dan Berger dan Guglielmi (2001).

3.2. Perhitungan Posterior

Sangat penting dalam menggunakan metode bayes nonparametrik dibandingkan dengan analisis parametrik karena kemampuan untuk memasukkan lebih banyak ketidakpastian tentang distribusi sampling. Akan tetapi, fleksibilitas ini meningkatkan kompleksitas perhitungan analisis. Banyak perkembangan model Bayes nonparametrik yang telah menjadi hasil langsung dari bantuan-bantuan dalam simulasi yang didasarkan pada metode perhitungan, khususnya metode MCMC. Pengenalan metode MCMC dalam bidang tersebut dimulai dengan penelitian Escobar (1994) untuk penggabungan proses Dirichlet. Dalam bagian ini kita membahas beberapa aspek perhitungan dari sampling posterior. Modelnya diberikan sebagai

$$\begin{aligned} X_1, \dots, X_n | G \sim^{iid} G, \\ G | c, \hat{\mu}, \hat{\sigma}^2 \sim PT_j(c, \hat{\mu}, N(\hat{\mu}, \hat{\sigma}^2)), \end{aligned} \quad (3.22)$$

dan

$$(\hat{\mu}, \hat{\sigma}^2) \sim p(\hat{\mu})p(\hat{\sigma}^2).. \quad (3.23)$$

Disini c ditentukan tetapi dapat juga diperlakukan sebagai random.

4. Simulasi

Pendekatan yang digunakan untuk menganalisa data yang bersifat ordinal dengan menyesuaikan serangkaian model logistik seperti rasio kontinuitas atau logit kumulatif. Dengan secara khusus, ditentukan $Y_{ij} = 1$ jika individu i mengalami kesulitan pada waktu j , dengan $Y_{ij} = 0$ jika tidak mengalami kesulitan atau ringan dan dapat menentukan model,

$$\begin{aligned} \text{logit}\{P(Y_{ij} = 1 | \hat{\mu}_i, \hat{\sigma}_i)\} \\ = \hat{\mu}_i + \hat{\alpha}_1 Trt_i + \hat{\alpha}_2 Time_{ij} + \hat{\alpha}_3 Trt_i \times Time_{ij} \end{aligned} \quad (4.1)$$

Dimana $i = 1, \dots, n, j = 1, \dots, N_i$. Disini y_i merupakan pengaruh random untuk masing-masing subyek, dan $\beta = (\beta_1, \beta_2, \beta_3)'$ adalah parameter regresi yang berhubungan dengan Trt , indikator perlakuan biner, waktu dalam periode, dan interaksi $Trt \times Waktu$. Secara khusus, y_i akan diasumsikan independen $N(\mu, \sigma^2)$. Kita mengganti asumsi normalitas dengan prior MPT,

$$\begin{aligned} \tilde{\mu}_1, \dots, \tilde{\mu}_n | G \sim^{iid} G, \\ G | \hat{\mu}, \hat{\sigma}^2 \sim PT_4(c, j^2, N(\hat{\mu}, \hat{\sigma}^2)). \end{aligned} \quad (4.2)$$

Metodologi dalam tulisan ini memberikan asumsi kenormalan yang pengaruh random dan implikasi kesalahan spesifikasi model normal yang mungkin. Prior acuan yang berbeda untuk θ_j s ditentukan dengan menggunakan tiga nilai c , yaitu $c = 0,1, c = 1$, dan $c = 10$, untuk menggambarkan peningkatan tingkat kepercayaan dalam normalitas untuk pengaruh random. Analisis yang lebih awal didasarkan pada asumsi normal menampilkan ketidakkonsistenan antara marginal

dan pengaruh khusus subyek yang membuat dalam menduga validitas analisis berdasarkan asumsi kenormalan.

Analisis Bayes membutuhkan prior baik untuk parameter regresi $\beta = (\beta_1, \beta_2, \beta_3)'$ maupun parameter μ dan σ^2 dari keluarga normal asli yang digeneralisasi oleh MPT. Kita menggunakan normal konjugat independen – prior gamma invers yang menampilkan informasi prior yang lemah.

Adapun simulasi dijalankan dengan menggunakan softwer dan program R, dimana data idnr, trt, dan time digeneret, dan ditentukan model dengan menggunakan persamaan (4.1) dan (4.2). parameter diasumsikan bahwa $\hat{\alpha}_0 = \text{trt}$, $\hat{\alpha}_1 = \text{time}$, dan $\hat{\alpha}_3 = \text{trt} \times \text{time}$ dengan parameter lainnya yaitu $\hat{\alpha}_2$ dan $\hat{\alpha}_4$. Untuk menentukan model mana yang terbaik, cukup kita melihat nilai (Deviance Information Criterion) DIC dan (Log Pseudo Marginal Likelihood) LPML yang terkecil.

Untuk menentukan model masing-masing dihitung dengan nilai c. Hasil simulasi untuk tiap-tiap model yang diberikan dalam tabel sebagai berikut:

Tabel 1. Hasil simulasi untuk posterior padamodelMPT dengan ukuran $c = 0.1$

Parameter	95%HPD				
	Mean	Median	Std. Dev	Lower	Upper
(Intercept)	-1.8328	-1.8228	0.4552	-2.6822	-0.9701
Trt	0.2922	0.2412	0.4913	-0.6572	1.2069
Time	-0.3544	-0.3524	0.0411	-0.4291	-0.2756
Trt*time	-0.1268	-0.1253	0.0637	-0.2404	0.0017
$\hat{\alpha}_2$	-0.7987	-0.8020	0.0589	-0.9292	0.6969
$\hat{\alpha}_4$	8.4236	8.5307	1.0476	6.2112	10.1113

Untuk hasil simulasi pada posterior model MPT dengan $c = 0,1$ memperoleh (Deviance Information Criterion) DIC = 930.4 dan (Log Pseudo Marginal Likelihood) LPML = - 477.6

Tabel 2. Hasil simulasi untuk posterior padamodelMPT dengan ukuran $c = 1$

Parameter	95%HPD				
	Mean	Median	Std. Dev	Lower	Upper
Intercept	-0.462823	-0.490515	0.530045	-1.382388	0.670885
Trt	0.136144	0.128177	0.717310	-1.250292	1.472916
Time	-0.177666	-0.166746	0.096298	-0.368031	0.006302
Trtxtime	0.029648	0.030614	0.121340	-0.191784	0.276181
$\hat{\alpha}_2$	-0.52133	-0.55759	0.68342	-1.89266	0.79928
$\hat{\alpha}_4$	0.68683	0.42639	0.85987	0.06532	2.17872

Untuk hasil simulasi pada posterior model MPT dengan $c = 1$ memperoleh (Deviance Information Criterion) DIC = 924.2 dan (Log Pseudo Marginal Likelihood) LPML = -468.0.

Tabel 3. Hasil simulasi untuk posterior pada model MPT dengan ukuran $c = 10$

Parameter	95%HPD				
	Mean	Median	Std. Dev	Lower	Upper
Intercept	-1.7524	-1.7491	0.4429	-2.5020	-0.8428
Trt	-0.1527	-0.1743	0.5711	-1.2850	0.9228
Time	-0.3931	-0.3924	0.0448	-0.4898	-0.3142
Trt*time	-0.1349	-0.1357	0.0705	-0.2730	-0.0076
$\hat{\mu}$	-1.9454	-1.9418	0.6870	-3.3357	-0.7314
$\hat{\sigma}^2$	16.8547	16.2695	3.9546	10.4916	24.8731

Untuk hasil simulasi pada posterior model MPT dengan $c = 10$ memperoleh (Deviance Information Criterion) DIC = 955.3 dan (Log Pseudo Marginal Likelihood) LPML = -481.3.

Tabel 4. Cara posterior, kriteria perbandingan model, dan interval 95% HPD untuk model gabungan linier tergeneralisasi (GLMM) parameter

Parameter	MPT		
	$\hat{\mu} = 10$	$\hat{\mu} = 1$	$\hat{\mu} = 0,1$
Intercept	-1.7524	-2.2108	-1.8328
Trt	-0.1527	0.2310	0.2922
Time	-0.3931	-0.3841	-0.3544
Trt*time	-0.1349	-0.1301	-0.1268
$\hat{\mu}$	-1.9454	-1.2101	-0.7987
$\hat{\sigma}^2$	16.8547	22.7622	8.4236
DIC	955.3	924.2	930.4
LPML	-481.3	-468.0	-477.6
Posterior Interval			
Intercept	-2.5020, -0.8428	-3.3212, -1.2262	-2.6822, -0.9701
Trt	-1.2850, 0.9229	-0.7307, 1.0223	-0.6572, 1.2069
Time	-0.4898, 0.3142	-0.4735, -0.3076	-0.4291, -0.2756
Trt*time	-0.2730, 0.0076	-0.2634, -0.0020	-0.2404, 0.0017
$\hat{\mu}$	-3.3357, 0.7314	-3.0559, 0.0173	-0.9292, 0.6969
$\hat{\sigma}^2$	10.4916, 24.8731	8.8015, 44.5144	6.2112, 10.1113

Hasil simulasi di atas yang ditunjukkan pada tabel 4.4 dapat kita simpulkan bahwa model yang dipakai atau yang terbaik yaitu pada c sama dengan 1 dengan nilai DIC dan LPML nya adalah model penyesuaian terbaik. Model ini lebih baik dibandingkan dengan model $c = 0.1$ dan $c = 10$.

5. Kesimpulan dan Saran

5.1. Kesimpulan

Kemampuan untuk menyelesaikan model yang rumit selalu digunakan analisis bayes. Hal ini bergantung pada kemampuan untuk menggunakan metode Rantai Markof Monte Carlo (MCMC) yang dapat memprediksi distribusi posterior. Dengan kata lain analisis bayes mampu memberikan fleksibilitas yang besar akan tetapi tidak menggunakan fleksibilitas itu jika data tersebut benar-benar membutuhkannya.

Diagram pohon polya merupakan sebuah distribusi probabilitas random, dengan mempelajari MPT kita dapat menggambarkan sebuah proses generalisasi distribusi parametrik yang normal. Generalisasi tersebut dengan beberapa tahapan dapat kita memperkenalkan parameter baru untuk menggeneralisasi tahapan sebelumnya, parameter yang dipakai berupa α atau β

Berdasarkan simulasi dengan menggunakan tiga nilai acuan $c = 0,1$, $c = 1$ serta $c = 10$, diperoleh hasil model penyesuaian terbaik pada $c = 1$ didasarkan pada nilai DIC dan LPMLnya.

5.2. Saran

Sangat menarik untuk penulisan lebih lanjut merupakan konstruksi diagram pohon polya untuk ruang multidimensional, dan memperluas diagram pohon polya untuk menentukan model probabilitas atas distribusi yang berhubungan, yaitu diagram pohon polya dependen dalam waktu, ruang, atau nilai kovariant silang. Diagram pohon polya dapat juga diadaptasi untuk menghasilkan realisasi dependen yang sesuai untuk memodelkan data kurva pertumbuhan. Merupakan pengaruh random yang mudah dirubah, atau aplikasi lain yang membutuhkan evolusi kompleks sebuah kepadatan dalam waktu, ruang atau dengan kovariatnya.

Daftar Pustaka

- Bain, L.J., and Max Engelhardt (1992) Introduction to Probability and Mathematical Statistics, 2nd edn, Duxbury Press, California.
- DeGroot, M.H. (1970) Optimal Statistical Decision, McGraw-Hill.
- Rohatgi, V.K. (1976) An Introduction to Probability Theory and Mathematical Statistics, John Wiley-Sons.

- W.J.,Conover. (1971) *Practical Nonparametric Statistics*, John Wiley-Sons, New York-London-Sydney-Toronto.
- George E.P.Box and George C.Tiao. (1973) *Bayesian Inference In Statistical Analysis*, Addison-Wesley Publishing Company.
- James O. Berger (1980) *Statistical Decision Theory Foundations, and Methods*, Springer-Verlag, New York Heidelberg Berlin.
- Mauldin,R.D.,Sudderth, W.D., and Williams, S.C. (1992). "Polya Trees and Random Distributions," *The Annals of Statistics*, 20, 1203-1221.
- Reich, B.J., and Fuentes, M. (2007), "Multivariate Semiparametric Bayesian Spatial Modeling Framework for Hurricane Surface Wind Fields," *Annals of Applied Statistics*, 1, 249-264.
- Dey,D., Muller, P., and Sinha, D. (1998), *Partical Nonparametric and Semiparametric Bayesian Statistics*, New York: Springer.
- Paddock, S.M.(1999), "Randomized Polya Trees: Bayesian Nonparametrics for Multivariate Data Analysis," Unpublished doctoral thesis, Institute of Statistics and Decision Sciences, Duke University.
- Spiegelhalter, D.J., Best, N.G., Carlin, B.P., Van der Line, A. (2000),"Bayesian Measures of Model Complexity and Fit," *Journal of the Royal Statistical Society, Series B*, 64, 583-639.
- Paddock, S.M., Ruggeria, F., Lavine M., and West, M. (2003), "Randomized Polya Tree Models for Nonparametric Bayesian inference," *Statistica Sinica*, 13, 443-460.
- Walker, S.G., and Mallick, B.K. (1999), "Semiparametric Accelerated Life Time Model," *Biometrics*, 55, 477-483.
- Warren, J., Ewens, and Gregory, R., Grant. (2005), "Statistical Methods in Bioinformatics: An Introduction, Springer.