

Analisis Perbandingan Berbagai Uji Pencilan Pada Analisis Regresi

Admi Nazra

Jurusan Matematika FMIPA Universitas Andalas

Abstrak

Dalam tulisan ini disimpulkan bahwa jika suatu data terdeteksi sebagai pencilan dengan uji sisa ter-*student*, maka data tersebut belum tentu pencilan jika diuji dengan simpangan baku. Secara teoritis Uji AIC dan BIC tidak mungkin dibandingkan dengan kedua metode terdahulu mengingat prosedurnya yang berbeda. Ketiga jenis uji ini belum tentu memberikan kasimpulan yang sama untuk menentukan data mana yang dikategorikan sebagai data pencilan pada suatu masalah regresi linier. Disamping itu uji pencilan dengan simpangan baku dan uji pencilan dengan sisa ter-*student* lebih sederhana untuk digunakan dibandingkan dengan uji pencilan kriteria informasi. Uji pencilan dengan sisa ter-*student* memberikan hasil yang relatif akurat dari dua uji lainnya.

Pendahuluan

Analisis regresi merupakan metoda statistika yang amat banyak digunakan, tidak hanya di lingkungan para peneliti di bidang matematika atau statistika namun juga di bidang-bidang lain seperti biologi, kimia, pertanian, ekonomi dan lain-lain. Umumnya analisis regresi digunakan dalam rangka mengolah data untuk menentukan hubungan antara dua atau lebih peubah sehingga diperoleh model atau hubungan fungsional antara peubah tersebut. Sehingga dengan model tersebut para peneliti dapat berusaha memahami, menerangkan, mengendalikan dan kemudian memprediksikan kelakuan sistem yang mereka teliti. Secara umum model merupakan penyederhanaan abstraksi dari keadaan alam yang sesungguhnya. Keadaan alam yang ingin diteliti biasanya amat rumit dan kemampuan kita menelitinya secara keseluruhan amat terbatas karena itu kita perlu menyederhanakannya sesuai dengan kemampuan

akal kita menghadapinya. Model yang dimaksud disini akan selalu berbentuk fungsi, dan regresi merupakan alat yang ampuh dalam pembentukannya.

Model regresi secara umum dapat ditulis sebagai $Y=X\beta + e$ dengan Y merupakan vektor respon $n \times 1$, X menyatakan matriks peubah bebas $n \times (k+1)$ (tuliskan $k+1=p$), β vektor parameter $(k+1) \times 1$ dan e vektor galat $n \times 1$. Metoda yang digunakan untuk menaksir β adalah metoda kuadrat terkecil.

Kebaikan model dapat dilihat dari nilai R^2 (koefisien determinasi) dan pengujian hipotesis terhadap parameter. Sedangkan kecocokan model dengan data dilihat dengan mengamati sisa (residual). Salah satu yang dapat dilihat dari sisa adalah pencilan (*outlier*). Pencilan adalah data yang tidak mengikuti pola umum model. Pada beberapa literatur telah dikemukakan beberapa uji untuk memeriksa apakah suatu data dapat dikategorikan sebagai pencilan. Secara kasar dapat diambil patokan yaitu yang sisanya berjarak dua simpangan baku atau lebih dari rata-ratanya dapat dikategorikan sebagai data pencilan (Sembiring, 1995). Namun (Weisberg, 1985) dan (Gentlemen & Wilk, 1975a-b) mengemukakan suatu cara yang sederhana ialah dengan menggunakan sisa *terstudent* dengan derajat bebas $dk=n-p-1$. Dimana bila sisa untuk data yang bersangkutan lebih besar dari nilai $t(n-p-1, \alpha)$ dari table-t maka dianggap data tersebut terpencil. Pynnönen (1992) menggunakan kriteria informasi untuk mendeteksi pencilan. Ada dua kriteria yaitu $AIC=n \log(1-R^2)-2 \log(n-k)!+2m$ dan $BIC=n \log(1-R^2)-2 \log(n-k)!+m \log n$.

Dari beberapa uji pencilan yang disebutkan di atas, maka muncul permasalahan, apakah data yang dikategorikan sebagai pencilan oleh suatu uji tertentu dapat juga terdeteksi sebagai pencilan oleh uji yang lain. Artinya apakah setiap uji uji tersebut memberikan hasil yang sama dalam menentukan suatu data apakah pencilan atau tidak. Jika memberikan hasil yang sama, uji manakah yang lebih sederhana atau lebih efektif dan efisien untuk digunakan. Hal ini sangat penting, sebab jika dalam penelitian ini nantinya ditemukan bahwa uji pencilan ini tidak memberikan informasi pencilan yang sama untuk berbagai metoda atau uji yang disebutkan di atas, maka tentunya seorang peneliti harus berhati-hati dalam menggunakan atau memilih uji uji tersebut.

Pendeteksian pencilan ini tidak hanya penting dalam rangka untuk memperbaiki model yang kita cari, namun juga dengan diketahuinya suatu data adalah pencilan maka seorang peneliti dapat menelusuri data tersebut untuk mengetahui dan mempelajari lebih mendalam mengenai data pencilan yang bersangkutan.

Tujuan Dan Manfaat Penelitian

Tujuan dari penelitian ini adalah:

1. Untuk mencoba berbagai uji pencilan yaitu uji simpangan baku, uji sisa *ter-student*, dan uji kriteria informasi untuk kasus yang sama dengan menggunakan data simulasi.
2. Untuk menganalisa dan membandingkan secara teoritis matematis berbagai uji pencilan diatas.
3. Untuk mengetahui apakah data yang dikategorikan sebagai data pencilan oleh uji tertentu, dapat juga terdeteksi sebagai pencilan oleh uji-uji yang lain.
4. Untuk mengetahui uji manakah yang lebih sederhana atau lebih efektif dan efisien untuk digunakan.
5. Untuk mengetahui data-data yang berpola atau berperilaku seperti apa yang cocok untuk uji-uji pencilan tertentu.

Manfaat Penelitian.

Dari hasil penelitian ini seorang peneliti khususnya peneliti yang akan menggunakan regresi dalam pengolahan datanya, dapat mengetahui dan menggunakan uji-uji pencilan yang lebih cocok dan tepat. Jadi hasil ini akan memberikan kontribusi penting bukan hanya untuk para peneliti namun juga bagi pengembangan konsep analisis regresi khususnya analisa data pencilan.

Metode penelitian

Metoda yang digunakan dalam penelitian ini adalah studi kepustakaan dan metoda simulasi dengan menggunakan program komputer (SAS, MINITAB, MATLAB), dan komputer yang digunakan adalah yang berkecepatan tinggi.

Langkah-langkah yang dilakukan dalam penelitian ini adalah:

1. Membangkitkan data dari populasi yang berdistribusi normal dan beberapa data yang bukan berdistribusi normal (yang diduga sebagai pencilan).
2. Melakukan berbagai uji pencilan seperti yang disebutkan di atas dengan menggunakan data-data simulasi pada langkah 1.
3. Membandingkan dan menganalisa hasil yang diperoleh dari langkah 2.
4. Gambaran yang diperoleh dari hasil langkah 3, dijadikan bahan perbandingan dan pemandu dalam menganalisa dan membandingkan secara teoritis matematis berbagai uji pencilan tersebut.
5. Menyimpulkan hasil yang diperoleh dari langkah 4.
6. Mencobakan uji-uji ini pada beberapa kasus yang diambil dari beberapa literatur serta dibandingkan dengan hasil simulasi dan hasil analisa teoritis

Hasil dan pembahasan

Formula untuk metoda uji simpangan baku dan uji sisa *ter-student* dapat dibandingkan secara teoritis dimana jika suatu data terdeteksi sebagai pencilan dengan uji sisa *ter-student* maka data tersebut belum tentu pencilan jika diuji dengan simpangan baku. Berdasarkan hasil ini tentu uji sisa *ter-student* mempunyai kelebihan karena dengan terdeteksinya beberapa data sebagai pencilan, maka hal ini dapat membantu kita untuk menyelidiki data-data tersebut lebih jauh, sesuai dengan salah satu tujuan mendeteksi pencilan.

Uji AIC dan BIC didapat dengan cara menghitung nilai R^2 dengan terlebih dahulu membuang data-data yang diasumsikan sebagai pencilan. Selanjutnya membandingkan nilai AIC dan BIC yang terkecil dari semua kemungkinan model yang dihasilkan. Secara teoritis ini tentu tidak mungkin dibandingkan dengan kedua metode terdahulu mengingat prosedurnya yang berbeda. Saat ini belum ada software khusus untuk mendapatkan nilai AIC dan BIC ini. Kita harus melakukan regresi yang berulang-ulang untuk mendapatkan nilai R^2 setelah membuang data-data yang diasumsikan sebagai pencilan. Jadi dari segi perhitungan, kedua metode terdahulu lebih sederhana. Yang mungkin

kita lakukan untuk melihat perbedaan hasil dari metode ini dengan dua metode sebelumnya adalah dengan melihat hasil beberapa ilustrasi berikut.

Untuk melihat perbedaan hasil uji pencilan untuk berbagai uji pencilan yaitu uji simpangan baku, uji sisa terstudent dan uji kriteria informasi maka digunakan tiga ilustrasi berikut ini.

Ilustrasi 1 diambil dari Barnett (1983) dengan model regresi $y = \beta_0 + \beta_1 x + \varepsilon$.

Datanya sebagai berikut:

(x)	4	5	7	9	11	14	17	20	23	26	30	35
(z)	110	81	90	74	20	30	37	22	38	25	18	9
(y = logz)	4.7	4.4	4.5	4.3	3.0	3.4	3.6	3.1	3.6	3.2	2.9	2.2

Dengan data ini diperoleh model dan hasil-hasil lainnya sebagai berikut:

$$Y = 4.66 - 0.0645x \quad S = 0.4103 \quad R-Sq = 73.8\%$$

Dengan data ini kita peroleh nilai AIC dan BIC yang terkecil adalah data $(x,y)=(11;3)$. Dari tabel sisa, terlihat bahwa harga mutlak sisa yang lebih besar dari $2s=0.8206$ adalah sisa untuk data $(x,y)=(11;3)$ yaitu -0.953487 . Untuk uji ter-student dimana data dengan harga mutlak sisa lebih besar dari nilai tabel $t(9;0,05)=1,833$ adalah data $(x,y)=(11;3)$ dengan sisa -2.46645 .

Ketiga jenis uji memberikan hasil yang sama dimana data pencilan adalah data $(x,y)=(11;3)$. Grafik 1 pada lampiran memperlihatkan juga bahwa data $(x,y)=(11;3)$ terlihat memencil dari data-data lainnya.

Sebagai ilustrasi 2 untuk membandingkan ketiga jenis uji pencilan digunakan data dari Chatterjee S. and B. Price (1977)

Dengan menggunakan model $y = \beta_0 + \beta_1 x + \varepsilon$ diperoleh hasil-hasil sebagai berikut.

$$Y = 1.71 + 0.665 X \quad S = 1.402 \quad R-Sq = 39.6\%$$

Nilai AIC dan BIC yang diperoleh dari data ini yang terkecil terjadi jika data pencilannya adalah $(x,y)=(7,3;9,5)$ yaitu observasi 29. Hasil ini juga bersesuaian dengan uji simpangan baku dimana data dengan sisa lebih besar dari $2s=2,804$ adalah data ke-29 juga. Namun kedua uji diatas memberikan hasil yang berbeda dengan uji sisa ter-student. Dimana dengan uji ini data pencilan adalah data dengan sisanya lebih besar dari $t(27;0,05)=1,703$ yaitu data observasi ke-27, 28, 29 dan 30.

Berikut ditampilkan data simulasi $y=5+2x$ dengan data x adalah 2.0 6.0 3.0 8.0 5.0 2.0 4.0 7.0 5.0 2.0 1.0 6.0 8.0 4.0 dan 9.0. Untuk memunculkan data yang diasumsikan sebagai data pencilan maka ditambah pasangan data $(x;y)$ $(5,5;17)$ dan $(2,5;8)$.

Dengan model $y=\beta_0+\beta_1x+\varepsilon$ diperoleh hasil-hasil sebagai berikut.

$$y = 4.68 + 2.06 x \quad S = 0.5570 \quad R-Sq = 98.8\%$$

Berdasarkan hasil perhitungan AIC dan BIC diatas diperoleh kesimpulan bahwa data pencilan adalah $(2,5;8)$. Hasil ini sama dengan uji simpangan baku. Sedangkan dengan uji sisa ter-student diperoleh bahwa data pencilan adalah $(2,5;8)$ dan $(5,5;17)$.

Melihat kepada hasil ketiga ilustrasi diatas dapat diperoleh hasil bahwa ketiga uji pencilan belum tentu memberikan kesimpulan yang sama. Jadi ada data yang oleh suatu uji dapat dikategorikan sebagai data pencilan namun dengan uji lain data tersebut tidak dapat dikategorikan sebagai data pencilan. Kalau dibandingkan uji-uji yang dipakai dengan grafik data aslinya, terlihat dari ilustrasi ini bahwa data-data yang dikategorikan sebagai data pencilan menurut uji sis ter-student, hampir bersesuaian juga dengan posisi data tersebut dilihat dari grafik, yang memang cukup memencil dibandingkan dengan data-data lainnya. Melihat hasil ini sepertinya uji sisa ter-student memberikan hasil yang relatif akurat dibandingkan dua uji lainnya walaupun dari ilustrasi ini uji simpangan baku memberikan hasil yang sama dengan uji AIC dan BIC.

Dilihat dari kesederhanaan dan keefektifannya, uji sisa ter-student ini lebih mudah untuk digunakan dibandingkan dengan uji AIC dan BIC. Karena saat ini perangkat lunak statistika selalu menyediakan menu untuk menghitung sisa ter-student ini. Sedangkan uji AIC dan BIC belum tersedia fasilitas untuk itu. Disamping itu dengan uji AIC dan BIC ini kita harus memprediksi dulu data-data yang mungkin dikategorikan sebagai data pencilan dan kita harus melakukan regresi beberapa kali tergantung kapada berapa macam kelompok data yang diasumsikan sebagai data pencilan. Penggunaan ketiga uji pencilan ini tidak tergantung kepada pola-pola datanya.

KESIMPULAN DAN SARAN

Dari pembahasan dan penjelasan diatas dapat diambil beberapa kesimpulan sebagai berikut:

1. Jika suatu data terdeteksi sebagai pencilan dengan uji sisa ter-*student*, maka data tersebut belum tentu pencilan jika diuji dengan simpangan baku.
2. Uji pencilan dengan simpangan baku, uji pencilan dengan sisa ter-*student* dan uji pencilan dengan kriteria informasi, belum tentu memberikan kesimpulan yang sama untuk menentukan data mana yang dikategorikan sebagai data pencilan pada suatu masalah regresi linier.
3. Uji pencilan dengan simpangan baku dan uji pencilan dengan sisa ter-*student* lebih sederhana untuk digunakan dibandingkan dengan uji pencilan kriteria informasi.
4. Uji pencilan dengan sisa ter-*student* memberikan hasil yang relatif akurat dari dua uji lainnya.

Daftar Pustaka

- Barnett, V. Principles and methods for handling outliers in data sets. *Statistical Methods and The Improvement of Data Quality*, pp. 131-166, 1983
- Chatterjee S. and B. Price, *Regression Analysis by Example*, John Wiley & Sons, New York, 1977.
- Gentleman, J.F. dan Wilk, M.B., Detecting Outliers in a two-way table I, *Technometrics*, 17,h.1-14, 1975a.
- Gentleman, J.F. dan Wilk, M.B., Detecting Outliers II , *Biometrics* , 31,h.387-400, 1975b.
- Pynnönen, S. Detection of Outlier in Regression Analysis by Information Criteria, www.twas.fi/~sjp/, 1992.
- Ryan, T.P., *Modern Regression Methods*, John Wiley & Sons, New York, 1997
- Sembiring, R. K., *Analisis Regresi*, ITB, 1995.
- Weisberg, S., *Applied Linear Regression* ,ed-2. John Wiley & Sons, New York, 1985