الكشف عن القيم الشاذة في نموذج الانحدار الخطي متعدد المتغيرات باستخدام معاينة جبس

د. محمد نذير إسماعيل قاسم يونس حازم إسماعيل قسم الرياضيات / كلية التربية جامعة الموصل

الاستلام القبول 2011 / 06 / 01 2011 / 03 / 22

Abstract

This paper deals with finding outliers in a multivariate linear regression model after assuming a model of normal — Wishart distribution. This method is based on the estimation of probability of an outlier for each observation by mixed Bernoulli model with shifting location outlier. We show how to obtain the posterior distribution in the mixed model by Gibbs sampler algorithm. Also the determination of the number of outliers is done by criterion of marginal likelihood distribution. The theoretical results of this research are applied to real data of multivariate linear regression. The results obtained are so encouraging in determining the outliers in these data.

الملخص

يتناول هذا البحث طريقة لإيجاد القيم الشاذة في نموذج الانحدار الخطي المتعدد المتغيرات وذلك بعد افتراض نمذجة التوزيع الطبيعي – وشارت. وأساس هذه الطريقة هو تقدير احتمالية القيمة الشاذة لكل مشاهدة بواسطة نموذج برنوللي المختلط ذو القيمة الشاذة مزاحة الموقع. ونبين كيفية الحصول على التوزيع ألبعدي في النموذج المختلط بواسطة خوارزمية معاينة جبس. يتم أيضا تحديد عدد القيم الشاذة باستخدام معيار توزيع الإمكان الحدي. تم تطبيق النتائج النظرية على بيانات حقيقية لنموذج انحدار خطي متعدد المتغيرات وكانت النتائج مشجعة في تحديد القيم البيانات.

1. مقدمة

يعد موضوع الانحدار متعدد المتغيرات من الأدوات المهمة في تحليل بيانات متعددة المتغيرات في حالة وجود أكثر من متغير استجابة واحد وأكثر من متغير توضيحي واحد. ويشمل تحليل بيانات في العلوم الاقتصادية والاجتماعية والنفسية والهندسية والزراعية والطبيعية والى أخره من التطبيقات.

إن دراسة وكشف القيم الشاذة في بيانات متعددة المتغيرات تعتبر من البحوث النادرة في الدبيات القيم الشاذة. انظر (2002) Ekiz . Ekiz . Ekiz (2002) أدبيات القيم الشاذة. انظر (2002) . Ekiz . (2002) . وكلف المتخدام الأسلوب الجديد لطريقة سلسلة ماركوف مونتي كارلو . (Monte Carlo Markov الشافب الجديد لطريقة سلسلة ماركوف مونتي كارلو . ومن الرواد . ومن الرواد . ومن الرواد . Wasserman (1991) . Wasserman (1991)

سوف نستخدم طريقة معاينة جبس لهذين الباحثين في نماذج متعددة المتغيرات ذات القيمة الشاذة مزاحة الموقع Multivariate Location Shift Outlier Models لنبين كيفية الحصول على دوال الإمكان الحدية من خوارزمية معاينة جبس. إن دوال الإمكان الحدية تستخدم في الاختبار البيزي بشكل واسع وكذلك في تشخيص النموذج كما في معيار بيز في السلاسل الزمنية Bayes Factor). كما أن نسبة دوال الإمكان الحدية تعرف عامل بيز Bayes Factor والذي سيأتي ذكره لاحقا.

وكما هو معلوم أن نموذج الانحدار الخطي البسيط ممكن أن يكون ملوثا Contaminated لنموذج فيه قيمة شاذة مزاحة الموقع حيث أن احتمالية الاختلاط تتبع توزيع برنولي بمعلمة غير معروفة (انظر إسماعيل (2005)).

إن نمذجـــة القــيم الــشاذة ممكــن إيجادهــا فــي أدبيــات البحــوث مــثلا (Kitagawa and Akaike(1992), Barnett and Lewis(1984), Pettit and Smith (1985))

والأعمال التي تضمنت دراسة ونمذجة القيم الشاذة في البيانات متعددة المتغيرات مثلا .Ekiz (2006), Ekiz (2002), Chaloner and Brant (1988)

يشمل هذا البحث خمسة مباحث الأول هي المقدمة والثاني تحليل الانحدار متعدد المتغيرات ومحتوي على قيمة شاذة والذي يتضمن نموذج الانحدار البيزي بقيم شاذة مضافة. والمبحث الثالث تحديد التوزيعات الاحتمالية بطريقة معاينة جبس لنموذج الانحدار متعدد المتغيرات بقيم شاذة والمبحث الرابع خصص لدالة الإمكان الحدية لنموذج الانحدار متعدد المتغيرات والخامس عامل بيز وأخيرا الجانب التطبيقي.

2. تحليل الانحدار متعدد المتغيرات مع قيم شاذة

لنفرض أن لدينا نموذج تحليل الانحدار متعدد المتغيرات الأتي

$$Y = XB + E \qquad \dots (1)$$

حيث أن $(Y_1, Y_2, \dots, Y_n') = Y' = (Y_1, Y_2, \dots, Y_n')$ وكل صف $Y' = (Y_1, Y_2, \dots, Y_n')$ وكل متغيرات الاستجابة مقاس على N = N وكل متغيرات الاستجابة مقاس على N = N وكل متغيرات الاستجابة مقاس على N = N بحيث أن:

تمثــل $X_i'=(x_1', x_2', \dots, x_k')$ ، X مـــف مــن $X_i'=(X_1', X_2', \dots, X_n')$ المشاهدة $X_i'=(X_1', X_2', \dots, X_n')$ و Explanatory Variables و $X_i'=(X_1', X_2', \dots, X_n')$ و المشاهدة $X_i'=(X_1', X_1', \dots, X_n')$

$$E_i \sim (0, \Psi)$$
 , $Y_i \sim N_n[x_i B, \Psi]$, $i = 1, ..., n$

حيث Ψ مصفوفة التغاير ذات بعد $p \times p$. سوف نبين الآن كيفية توسيع هذا النموذج إلى نموذج متعدد المتغيرات بقيمة شاذة.

و أول نموذج من هذا النوع هو نموذج القيم الشاذة مزاحة الموقع أحادي المتغير Univariate Location Shift Outlier Model والذي تم تحليله بواسطة معاينة جبس من قبل (2005).

إن نموذج متعدد المتغيرات بقيمة شاذة مزاحة الموقع مع n من متغيرات الدليل z_1, z_2, \dots, z_n Indicator Variables الأتى:

$$f(Y_i) = P(z_i = 0)f(Y_i|z_i = 0) + P(z_i = 1)f(Y_i|z_i = 1)$$
(2)

$$= (1 - \varepsilon) N_p [Y_i | x_i B, \Psi] + \varepsilon N_p [Y_i | a_i + x_i B, \Psi] \qquad(3)$$

حيث أن X_i هي الصف X_i المصفوفة البيانات X_i و تمثل الصف X_i ذات بعد X_i من حيث أن X_i هي الصف X_i المصفوفة إزاحة الموقع X_i والتي بعدها X_i بحيث X_i بحيث X_i وسيكون القيم الشاذة متعددة المتغيرات المصفوفة X_i والتي تمثل احتمالية كون المشاهدة X_i شاذة.

وإذا كانت المشاهدات مستقلة فان النموذج الاحتمالي للانحدار متعدد المتغيرات بقيم شاذة له الشكل الأتي:

$$Y \sim N_{n \times p} [XB + D_z A, \Psi \otimes I_n]$$
(4)

$$Y' \sim N_{P \times n} \left[B' X' + A' D_z, I_n \otimes \Psi \right]$$
(5) يكون Y

حيث أن ⊗ تمثل حاصل ضرب كرونيكر. انظر (Rao and Kleffe (1988) .

المعلومات الأولية (أو السابقة) عن مصفوفات المعلمة ممكن صياغتها بشكل محكم وكالأتي:

$$B \sim N_{k \times p} [B_*, G_* \otimes H_*]$$
 $\Psi^{-1} \sim W_p [\Psi_*, n_*]$
 $A \sim N_{n \times p} [A_*, P_* \otimes I_n]$
 $z_i \sim Ber[\varepsilon_i]$
 $i = 1, 2, \dots, n$

حيث أن $0 < \varepsilon_i < 1$ الاحتمالية الأولية لمعلمة توزيع برنولي. المصفوفة A_* ذات بعد $p \times p$ تمثلان المعلومات الأولية لمعلمتي الموقع (الوسط) والتباين للقيم الشاذة. و E_i احتمالية النجاح للمشاهدة E_i أن تكون شاذة وعادة تكون ثابتة نقيمة E_i .

3. معاينة جبس في نموذج الانحدار الاحتمالي متعدد المتغيرات بقيم شاذة

في البداية لابد من الإشارة إلى مقدمة تعريفية إلى معاينة جبس حيث تعود معاينة جبس إلى عمليات جبس Gibbs Processes المأخوذة من علم الفيزياء والتي تمثل طريقة لبناء عمليات رياضية جديدة من العمليات الرياضية القديمة. إذ أننا نبدأ بعملية ابتدائية $Y^{(0)}$ ثم نعرف عملية جديدة $Y^{(0)}$ وذلك بأخذ الدالة الاحتمالية $Y^{(0)}$ لا Y بالنسبة إلى $Y^{(0)}$ وللتعرف على كيفية عمل طريقة معاينة جبس افرض انه لدينا $Y^{(0)}$ من المتغيرات العشوائية والتي نرمز لها بالرمز

$$Y = (Y_1, Y_2, \dots, Y_k)$$

ونفرض لدينا التوزيعات الشرطية التامة بالشكل:

$$p(y_i|y_j)$$
, $i \neq j$, $i, j = 1,2,...,k$

وفي هذه الخوارزمية لا نحتاج الشكل التام للتوزيعات الشرطية

$$p(y_i|y_j)$$
, $i \neq j$

لكن نحتاج فقط كتابتهم بشكل يتناسب وتحت شروط معتدلة Casella and George لكن نحتاج فقط كتابتهم بشكل يتناسب وتحت شروط معتدلة (1992).

خطوات خوارزمية معاينة جبس

وفيما يلي توضيح موجز لخوارزمية معاينة جبس:

1. افرض لدينا المتغيرات العشواية

$$\{Y_1, Y_2, ..., Y_k\}$$

ولها القيم الابتدائية

$$\left\{Y_1^{(0)}, Y_2^{(0)}, \dots, Y_k^{(0)}\right\}$$

 $Y_1^{(1)}$ من التوزيع الشرطي $Y_1^{(2)}$

 $\left\{Y_1 \mid Y_2^{(0)}, Y_3^{(0)}, \dots, Y_k^{(0)}\right\}$

3. اسحب $Y_2^{(1)}$ من التوزيع الشرطي

 $\left\{Y_{2} \mid Y_{1}^{(1)}, Y_{3}^{(0)}, \dots, Y_{k}^{(0)}\right\}$

.

الم التوزيع الشرطي $Y_k^{(1)}$ من التوزيع الشرطي .k .4

 $\left\{Y_{k} \mid Y_{1}^{(1)}, Y_{2}^{(1)}, \dots, Y_{k-1}^{(1)}\right\}$

وهذا یکمل تکرار واحد لمعاینـة جبس وعلیـه بعد تکرار واحد نکون قد حصلنا علی $\left\{Y_1^{(1)},Y_2^{(1)},\dots,Y_k^{(1)}\right\}$

 $\left\{ Y_{1}^{(k)}, Y_{2}^{(k)}, \dots, Y_{k}^{(k)} \right\}$ وبعد t مثل هذه التكرارات نكون قد حصلنا

وبتكرار هذه العملية m من المرات فان معاينة جبس ممكن أن تولد مشاهدات مستقلة وبتكرار هذه العملية $p(y_1)$ بشرط أن تكون عددها $p(y_1)$ بشرط أن تكون عددها $p(y_1)$ بالتكرارات مستقلة وباستخدام قيم ابتدائية $\{Y_1^{(0)}, Y_2^{(0)}, \dots, Y_k^{(0)}\}$

وبعد الحصول على الخصائص الإحصائية لهذا المتغير مثلا الوسط والتباين أو أية خواص إحصائية أخرى لهذا المتغير Y_1 وهكذا بالنسبة للمتغير Y_2 و Y_3 والى Y_4 انظر إسماعيل (2005).

وألان نعود إلى النموذج في (1) والذي يمكن كتابته بالشكل:

$$\mathbf{Y}_{(n \times p)} = \mathbf{X}_{(n \times K)} \underbrace{B}_{(n \times p)} + \underbrace{E}_{(n \times p)} \tag{6}$$

افرض أن A و D كما معرفة سابقا في (2). ممكن صياغة النموذج بعد افتراض التوزيع الطبيعي بالشكل:

 $Y \sim N_{n \times p} [XB + D_z A, \Psi \otimes I_n]$

 $.B, \Psi, A, D_{\bullet}$ نفرض أن ' $\theta = (B, \Psi, A, D_{\bullet})$ نفرض

التوزيع المشترك للمتغير Y والمعلمة θ يكون

 $p(\mathbf{Y}, \boldsymbol{\theta}) = N_{n \times p} \left[\mathbf{Y} | \mathbf{X} \boldsymbol{B} + D_z \boldsymbol{A}, \boldsymbol{\Psi} \otimes \boldsymbol{I}_n \right] N_{k \times p} \left[\boldsymbol{B} | \boldsymbol{B}_*, \boldsymbol{G}_* \otimes \boldsymbol{H}_* \right] W_p \left[\boldsymbol{\Psi}^{-1} | \boldsymbol{\Psi}_*, \boldsymbol{n}_* \right]$

 $N_{n \times p}[A|A_*, P_* \otimes I_n] \cdot \sum_{i=1}^n Ber(Z|\varepsilon_i)$

نستخدم في الأتي θ^c لتمثل الحجة (Argument) الشرطية في التوزيع الشرطي التام لمجموعة المعلمة θ بدون المعلمات الحالية.

يمكن إيجاد التوزيعات الشرطية التامة لمعاينة جبس كما يأتى:

أ) لمصفوفة معاملات الانحدار B

$$p(B|Y,\theta^c) = N_{K\times p}[B_{**},C_{**}]$$

حيث أن:

$$C_{**}^{-1} = G_{*}^{-1} \otimes H_{*}^{-1} + \Psi^{-1} \otimes X X,$$

$$vec B_{**} = C_{**} \left[vec \left(G_{*}^{-1} B_{*} H_{*}^{-1} + X (Y - D_{z} A) \Psi^{-1} \right) \right]$$

حيث أن vec عملية رياضية تحول المصفوفة ذات السعة $n \times p$ إلى متجه ذو بعد $n \times p$ عنصر بشكل عمود واحد تحت الأخر ابتداء من الصف الأول.

ب) لمصفوفة التغاير Ψ:

$$p(\Psi^{-1}|Y,\theta^{c})=W_{p}[\Psi_{**},n_{**}=n_{*}+n]$$

يمثل توزيع وشارت ذات بعد p وبمعلمة قياس

$$\Psi_{**} = \Psi_* + (Y - XB - D_*A)(Y - XB - D_*A)$$

تعریف: إذا کانت المصفوفة $M(p \times p)$ ممکن کتابتها بالشکل M = x'x حیث أن x مصفوفة ذات بعد $M(p \times p)$ وتمثل مصفوفة بیانات من توزیع $M(p \times p)$ فان $M \times p$ وتمثل مصفوفة بیانات من توزیع $M \times W_p[0,\Sigma]$ فان $M \sim W_p[\Sigma,m]$ ومعلمة درجات الحریة $M \sim W_p[\Sigma,m]$ فان $M \sim W_p[\Sigma,m]$

Mardia et al(1997), Johnson and Wichern (2007) انظر

A لمصفوفة إزاحة الموقع

$$p(A|Y,\theta^{c}) = N_{n \times n}[A_{**},G_{**}]$$

توزيع طبيعي متعدد المتغيرات بالمعلمات

$$\begin{split} C_{**}^{-1} &= P_{*}^{-1} \otimes I_{n} + \Psi^{-1} \otimes D_{z}^{'} D_{z} ,\\ vec \, A_{**} &= G_{**} \Big[vec \big(A_{*} \, P_{*}^{-1} + D_{z} \big(Y - XB \big) \Psi^{-1} \big) \Big]. \end{split}$$

ولكل مشاهدة فان الوسط ألبعدي Posterior Mean يمكن حسابه وذلك بتجزئة نظام المعادلات إلى معادلات بمتغير واحد:

$$C_{**_{i}}^{-1} = P_{*}^{-1} + z_{i}^{2} \Psi^{-1}, i = 1,...., n,$$

$$a_{**_{i}} = G_{**_{i}} \Big[P_{*}^{-1} a_{*_{i}} + z_{i} \Psi^{-1} (y_{i} - Bx_{i}) \Big].$$

د) لمتغيرات المؤشر :z,

$$p(z_i|Y,\theta^c) = Ber\left[\varepsilon_{i^*} = \frac{c_i}{c_i + d_i}\right], \quad i = 1,...,n,$$

توزيع برنولي بالمركبات التي نحصل عليها بطريقة تطبيق نظرية ببز ، أي بمعنى $c_i = N_p \big[y_i \big| x_i B + a_i \; , \Psi \big]. \varepsilon_i \; ,$ $d_i = N_p \big[y_i \big| x_i B \; , \Psi \big]. \big(1 - \varepsilon_i \big), \quad i = 1, \ldots, n,$

A عيث أن a_i هو الصف i هي المصفوفة X و X عيث الصف X

4. حساب الإمكان الحدى لنموذج الانحدار متعدد المتغيرات بالقيم الشاذة

باستخدام طريقة (1995) من الممكن إيجاد الإمكان الحدي عند النقطة $\hat{\theta}_0$ والتي هي:

$$\hat{\theta}_1 = (\hat{\theta}_0, \hat{A} = 0, \hat{D}_z = 0)$$
(7)

حيث $\hat{\theta}_0 = (\hat{B}, \hat{\Psi})$ نفس النقطة لنموذج الانحدار متعدد المتغيرات بدون قيم شاذة. وألان يكون لدينا التجزئة الآتية:

$$p(\hat{\theta} | Y) = p(\hat{D}_z | Y) \cdot p(\hat{A} | \hat{D}_z, Y) \cdot p(\hat{\theta}_0 | \hat{A}, \hat{D}_z, Y) \qquad \dots (8)$$

9

$$p(\hat{\theta}) = p(\hat{\theta}_0) N[A_*, I_n \otimes P_*] \prod_{l=1}^n Ber(\varepsilon).$$
 (9)

انظر (1999) Polasek.

1. نستخدم دورة جبس من J من نقاط العينة لنموذج الانحدار متعدد المتغيرات بالقيم الشاذة لحساب المقدار

حيث أن معلمة الكثافة البعدية في الموقع i لتوزيع برنولي معطاة من:

$$\begin{split} \varepsilon_{i^{**}}^{(j)} &= \frac{c_{i}^{(j)}}{c_{i}^{(j)} + d_{i}^{(j)}} \\ c_{i}^{(j)} &= N_{p} \Big[y_{i} \Big| \varepsilon^{(j)} z_{i}^{(j)} + B_{j} x_{i}^{(j)} \Big], \\ d_{i}^{(j)} &= N_{p} \Big[y_{i} \Big| B^{(j)} x_{i}^{(j)} \Big] \end{split}$$

2. المركبة الثانية ممكن الحصول عليها بدون معاينة جبس:

$$p(\hat{A} \mid \hat{D}_z = 0, Y) = N_{n \times p} \left[\hat{A} \mid A_{**}, I_n \otimes G_{**} \right] = \prod_{i=1}^n N[\hat{a}_i \mid a_{i^{**}}, G_{**}]. \tag{11}$$

ومن الممكن ملاحظة من التوزيع الشرطي التام لمعلمة إزاحة الموقع A عند وضع $\hat{D}_z=0$ فان التوزيع الشرطي يساوي التوزيع الأولى:

$$G_{**} = I_n \otimes P_n$$

و كذلك

 $vec A_{**} = vec A_*$

لمزيد من التفاصيل انظر (1995) Chib

3. ممكن الحصول على الدالة الحد في المعادلة (6) والذي هو $p(\hat{\theta}_0 \mid \hat{A}, \hat{D}_z, Y)$ بواسطة احتساب الإمكان الحدي (الهامش) لنموذج الانحدار المتعدد المتغيرات بدون القيم الشاذة.وبعد أخذ الوغاريتم نحصل على النتيجة التالية دالة الإمكان ينتج:

$$\log p(Y) = \log p(Y|\hat{\theta}_1) + \log p(\hat{\theta}_1) - \log p(\hat{\theta}|Y) \qquad \dots (12)$$

حيث أن جزء الإمكان معطى بواسطة

 $\hat{D}_g = 0$ والذي له نفس قيمة الانحدار متعدد المتغيرات بدون قيم شاذة بسبب

ولاحظ إن الصيغة (10) هي تبسيط بسبب أن مركبات إزاحات الموقع تختصر.

ولمزيد من التفاصيل التي تحتاج الى عمليات جبرية طويلة في الحصول على المعادلة (12) انظر ملحق بحث (1998).

5. اختيار النموذج مع عوامل بيز

المتراجحات البعدية تستخدم في التحليل البيزي لاختبار نموذجين أو أكثر لنفس مجموعة البيانات. الصيغة الأساسية للاختبار النموذجين M_2 , M_1 هي

المتراجحات البعدية = عامل بيز × المتراجحات الأولية

Posterior odds = Bayes factor. Prior odds

أو

$$\frac{p(M_1|Y)}{p(M_2|Y)} = BF \cdot \frac{p(M_1)}{p(M_2)}$$

حيث أن $p(M_1|Y)$, $p(M_1|Y)$ هما الاحتمالات البعدية للنموذجين M_1 , $p(M_1|Y)$, على التوالي. وفي ابسط $p(M_2|Y)$, $p(M_1)$, $p(M_1)$, $p(M_1)$, وفي ابسط حالة (اقل المعلومات الأولية) يكونان متساويين. وفي هذه الحالات تكون المتراجحات البعدية تساوي عامل بيز والذي يمكن تعريفه على انه النسبة لدالة الإمكان ألبعدي

$$BF = \frac{p(M_1|Y)}{p(M_2|Y)} = \frac{\int p(Y,\theta_1)d\theta_1}{\int p(Y,\theta_2)d\theta_2}$$
 (14)

BF>1 حيث أن θ_2 , θ_1 تمثلان معلمتي النموذجين M_2 , M_1 على الترتيب. فإذا كان BF<1 نختار النموذج M_1 مثلا عامل بيز لاختيار نموذج يحتوى على قيم شاذة مقابل نموذج بدون قيم شاذة سيعرف على انه:

$$BF = \frac{p(Y \mid outliers)}{p(Y \mid no \ outliers)}$$

إن عامل بيز يمكن حسابه من الصيغة المعرفة

$$BF = \frac{P_i F_p(a)}{(1 - P_i)(1 - F_p(a))}$$
 (15)

حيث $P_{p}(k)$ تمثل الاحتمالية البعدية للصيغة التربيعية للخطأ المشاهد وان $P_{p}(k)$ هي احتمال أن تكون القيمة المشاهدة شاذة و k تمثل قيمة حرجة تؤخذ من جدول توزيع مربع كاي استنادا إلى عدد المتغيرات المعتمدة وحجم العينة المدروسة. عبد الغنى (2010).

من الواضح أن عامل بيز يعتمد بشكل أساسي على الاحتمالية البعدية للصيغة التربيعية للخطأ المشاهد والمعرفة بالشكل:

$$P_{i} = E_{\Sigma/Y} \left\{ p_{r} \left(T_{i} > \frac{a}{\sigma_{(i)}} | \Sigma, Y \right) \right\}$$
 (16)

انظر عبد الغنى (2010).

المتراجحات البعدية لنموذج القيمة الشاذة نحصل عليه من عامل بيز بواسطة:

$$R = BF \times r \tag{17}$$

حيث أن R المتراجحات البعدية و r المتراجحات الأولية.

من المعادلة (17) وبأخذ اللوغارتم تصبح:

$$\log R = \log BF + \log r$$

حيث

$$r = \frac{p(outlier)}{1 - p(outlier)}$$

$$\Rightarrow p(outlier) = \frac{r}{1 + r}$$
 (18)

وهذه المعادلة تمثل احتمالية القيمة شاذة.

الجانب التطبيقي: تم تطبيق الجانب النظري على بيانات حقيقية لنموذج انحدار متعدد المتغيرات والتي تمثل قياسات نسبة السكر في الدم y_1 والمقاس بوحدة الملي مول (m.mol)، وضغط الدم العالي y_2 والواطي y_3 والمقاس بوحدة الملي لتر زئبقي كمتغيرات استجابة (معتمدة)، والعمر العالي y_3 مقاس بالسنة، والوزن y_3 مقاس بالكيلو غرام، والطول y_3 مقاس بالسنتمتر كمتغيرات توضيحية (مستقلة) لعينة من المرضى. إن هذه البيانات أخذت من عبد الغني (2010)، والبيانات موضحة في الجدول الأتي:

الجدول (1): قياسات نسبة السكر في الدم وضغط الدم العالي والواطي مع المتغيرات التوضيحية (العمر والوزن والطول) لعينة من المرضى.

الطول	الوزن	العمر	الضغط		السكر	التسلسل
			الواطي	العالي	استدر	العقيقين
164	66	49	8	12	5.0	.1
160	107	45	8	11	4.7	.2
158	80	18	8	16	1.9	.3
162	87	50	8	12	6.0	.4

الكشف عن القيم الشاذة في نموذج الانحدار الخطي متعدد المتغيرات باستخدام معاينة جبس.

147	95	69	8	16	23.3	.5
155	80	60	8	13	7.2	.6
162	80	72	8	13	5.5	.7
155	92	59	8	12	10.1	.8
166	85	53	7	14	6.2	.9
166	68	55	8	10	5.3	.10
154	77	50	6	12	3.8	.11
163	80	50	6	12	6.8	.12
160	90	55	8	12	8.1	.13
161	80	50	8	12	13.6	.14
156	80	54	8	15	15.7	.15
171	72	40	8	12	6.1	.16
167	75	59	8	12	5.4	.17
160	70	42	7	9	5.0	.18
170	90	47	9	12	26.7	.19
147	80	64	7	15	28.1	.20

لتكوين نموذج الانحدار المتعدد المتغيرات تم اعتبار أن متغيرات الاستجابة الثلاثة لتكوين نموذج الانحدار المتعدد المتغيرات تم اعتبار أن متغيرات الاستجابة الثلاثة $Y_i = \{Y_{i1}, Y_{i2}, Y_{i3}\}$ التي تمثل المتغيرات والتي تعتمد على المتغيرات التوضيحية الثلاثة $X_i = \{X_{i1}, X_{i2}, X_{i3}\}$ التي تمثل (الوزن، العمر، الطول). أي أن:

$$Y \sim N_{n \times p} [X B + D_z A, \Psi]$$

$$A$$
 مصفوفة الموقع المزاح وان التوزيع الأولي لـ A هو

$$A \sim N_{20\times 3} (A_*, P_* \otimes I_{20})$$
(19)

وان

$$v(a_i) = P_* = dig(var(Y))$$
, $i = 1, 2, 3$

 $.P_* = dig[60.0167, 3.3053, 0.5368]$ حيث أن

وان D_z مصفوفة المؤشر للقيم الشاذة وان التوزيع الأولي لـ D_z هو D_z

$$D_z \sim ber(\varepsilon_i)$$
(20)

حيث تم فرض $\varepsilon_i = 0.1$ قيمة ثابتة لكل القيم للسهولة.

وان التوزيع الأولي له B هو:

$$B \sim N_{3\times 3}(B_*, G_* \otimes H_*)$$
(21)

 $G_* \otimes H_* = \Gamma_3 \otimes I_3$ افرض أن

حيث أن $\Gamma_3 = dig(\lambda_i) = dig(0.0000225, 0.0006639, 0.0067794)$ وأن λ_i تمثل القيم الذاتية للمصفوفة YY. والوسط الأولي ببساطة يوضع ليمثل النقطة الوسطى لمتغير القياس والتي تساوي $T_3 = dig(\lambda_i) = dig(0.0000225, 0.0006639, 0.0067794)$

وبواسطة استخدام البرنامج الجاهز WinBUGS14 إذ تم تكرار هذا النموذج 1000 مرة والحصول على التكرارات لكل من المعلمات A,B,D_z , Ψ وكذلك الحصول على الاحتمال اللاحق. كذلك تم كتابة برنامج بنظام Matlab6.5 لحساب عامل بيز من المعادلة (15) وكانت النتائج كما موضح في الجدول الأتي:

جدول (2): عامل بيز واحتمالية المشاهدة i بأنها شاذة

p(outlier)	عامل بیز BF	t t ti
1 ()	22 7, 0	التسلسل
0.4347	0.9754	.1
0.3478	1.4069	.2
0.0783	8.7191	.3
0.2944	2.1536	.4
0.0427	16.2256	.5
0.1033	7.6636	.6
0.3817	1.4854	.7
0.0724	11.2842	.8
0.2633	2.5477	.9
0.3618	1.3677	.10
0.1254	5.8098	.11
0.2181	3.1008	.12
0.0811	10.354	.13
0.0650	12.4165	.14
0.0600	13.5809	.15
0.2535	2.3111	.16
0.4033	1.25	.17
0.4222	1.0542	.18
0.0468	18.7768	.19
0.0417	20.0644	.20

نلاحظ من خلال الجدول أعلاه إن المشاهدات ذات التسلسل (3، 5، 8، 13، 14، 15، 19، 19، 20) تملك احتمالية صغيرة جدا مقارنة مع بقية احتمالات المشاهدات الأخرى وعليه يمكن اعتبارها ممكن اعتبار هذه المشاهدات حالات متقدمة من المرض يجب الوقوف عندها من ناحية وصف الدواء واتخاذ الحمية.

الاستنتاجات

- 1. استخدام طريقة معاينة جبس يعطى نتائج دقيقة في حالة الكشف عن القيم الشاذة.
 - 2. تحديد دالة الإمكان الحدية بشكل كامل باستخدام معاينة جبس.
 - 3. إيجاد دالة الكثافة الحدية البعدية التامة.

التوصيات

- 1. يمكن استخدام أسلوب معاينة جبس في الكشف عن القيم الشاذة في بيانات تحليل التباين متعدد المتغيرات والتحليل ألعاملي.
- 2. تطبيق هذا التحليل على بيانات في مجالات العلوم الطبيعية والاقتصادية والنفسية والتربوية والاجتماعية.

المصادر

- 1) عبد الغني، تميم معاذ (2010): أسلوب بيز في الكشف عن القيم الشاذة في الانحدار الخطى متعدد المتغيرات، رسالة ماجستير غير منشورة، كلية التربية، جامعة الموصل.
- (2 إسماعيل، يونس حازم (2005): الكشف عن القيم الشاذة باستخدام التوزيع المختلط بالاعتماد على معاينة جبس وأسلوب بيز. رسالة ماجستير غير منشورة، كلية التربية، جامعة الموصل.
- 3) Barnett, V. and Lewis, T. (1984): Outliers in Statistical Data, Wiley, Chichester.
- 4) Chib, S.(1995): Marginal Likelihood from the Gibbs Output. JASA, 90, 1313-1321.
- 5) Casella, G. and George, E. (1992): "Explaining the Gibbs Sampler". The American Statistician, 46, 167-174.
- 6) Chaloner, K. and Brant, R. (1988): A Bayesian Approach to Outlier Detection and Desidual Analysis, Biometrika 75, 9 651.

- 7) Ekiz, U. (2002): A Bayesian Methods to Detect Outlier in Multivariate Linear Legression. Gazi university Faculty of science, Turkey.
- 8) Ekiz, U. (2006): Bayesian Approach to Identity Masking and Swamping Problems in Multivariate Analysis. Gazi university Faculty of Science. Turkey.
- 9) Kitagawa, G. and Akaike, H. (1992): A Quasi Bayesian Approach to Outlier Detection, *Ann. Inst. Stat. Mathematics* 34, 95–104.
- **10)** Johnson, A.R. and Wichern, D.W. (2007): Applied Multivariate Statistical Analysis. Pearson Prentice Hall. United States of America.
- 11) Mardia, K.V. and Kent, J.T. and Bibby, J.M.(1997): Multivariate Analysis. Academic press London.
- **12)** Pettit, L.I. and Smith, A.F.M. (1985): Outliers and Influential Observations in Linear Models. In Bayesian Statistics, 2, Ed. J.M. Bernardo, M.H. De Groot, D.V. Lindley and A.F.M. Smith, 473-494. Amsterdam North Holland.
- Polasek, W.(1999): Gibbs Sampling in B-VAR Models with Latent Variables. in: Innovations in Multivariate Statistical Analysis –A Festschrift for Heinz Neudecker, R.J. Heijmans, D. S. G. Pollock and A. Satorra, (eds), Kluwer Academic Publishers, Dordrecht, 137-16-56.
- 14) Polasek, W. and Ren, L (1998): Structureal Breaks and Model Selection with Marginal Likelihoods, in: W. Racugno (ed). Proceedings of the Workshop on Model Selection, Pitagora Ed. Bologna, 223-273, with discussion.
- 15) Rao, C.R. and Kleffe, J. (1988): Estimation of Variance Components and application North Holland, Amsterdam.
- Verdinelli, I. and Wasserman, L.(1991): Bayesian Analysis of Outlier Problems Using the Gibbs Sampler, Statistics and Computing 1991–1, 105–117.