

تحديد القيم الشاذة باستخدام الطرق الاستكشافية

ومقارنتها مع الطرق المعلمية

الباحث

دلير صليو دغا

أ. م. د. محمود مهدي حسن البياتي

كلية الادارة والاقتصاد/جامعة بغداد/ قسم الاحصاء

الخلاصة

تلعب البيانات الاحصائية دوراً مهماً في عملية التخطيط المركزي والدراسات العلمية ، وتأتي أهمية هذه الدراسة لتعاملها حول كيفية حفظ البيانات الاحصائية من الاخطاء المتوقعة والقيم الشاذة. الهدف من البحث هو تحديد القيم الشاذة باستخدام الطرق الاستكشافية الحديثة (Exploratory Data Analysis) ومقارنتها مع الطرق المعلمية بعد تحديد القيم الشاذة باستخدام الاسلوبين. وجد هناك اختلاف بين تحديد القيم السادة بستخدام الطرق الاستكشافية وطرق المعلمية وكذلك اختلاف بين الطرق المعلمية نفسها. تبين من خلال الدراسة أنه أفضل طريقة لتحديد القيم الشاذة هي طريقة الرسم الصندوفي (Box plot) .

Abstract

The availability of statistical data plays an important role in planning process. The importance of this research which deals with safety of statistical data from errors and outliers values. The Objective of this study is to determine the outlier values in statistical data by using modern exploratory data methods and comparing them with parametric methods. The research has been divided into four chapters ,the main important conclusions reached are:1-The exploratory methods and the parametric methods showed variation between them in determining the outlier values in the data.

2-The study showed that the box plot method was the best method used in determining outlier values in data.

* بحث مستقل من أطروحة الماجستير/قسم الاحصاء، 1996 الموسمة "تحديد القيم الشاذة باستخدام الطرق الاستكشافية ومقارنتها مع الطرق المعلمية".

1- المقدمة

التخطيط الشامل هو المفتاح المركزي لعملية التنمية ، وتلعب مسألة توفر البيانات الاحصائية دوراً مهماً في عملية التخطيط اذ ان عملية التخطيط المستقبلية تحتاج الى تحليل الواقع الحالي الذي يرتكز على دقة البيانات الاحصائية للمتغيرات المراد دراستها وما لها من اثر واضح على الخطة المرسومة من هنا تبرز اهمية هذا البحث الذي يتناول التأكيد من سلامة البيانات الاحصائية المستخدمة في المتغيرات قيد الدراسة ومحاولة تصفيه هذه البيانات من الاخطاء والقيم الشاذة Outlier values) } القيم الشاذة هي القيم الغير متجانسة مع بقية البيانات، ومنفصلة عن جسم البيانات {، والتي قد تؤدي في حالة تجاهلها الى تغيير نتائج التحليل الاحصائي الى نتائج غير دقيقة او انحراف النتائج عن واقعها الذي يجب ان تكون عليه فيما لو اجري التحليل الصحيح على البيانات الصحيحة او توجيهها الى اتجاه اخر بعيد عن الواقع، ومن ثم اعتماد اساليب وهمية لمعالجة متغيرات الظاهرة المدروسة.

ان الغرض من دراسة القيم الشاذة في البيانات هو محاولة تنقية البيانات من الشوائب والقيم الغريبة بأساليب وطرق علمية، لأن وجود القيم الشاذة سوف يؤثر بشكل سلبي على نتائج التحليل الاحصائي مما يقود الى عدم الاستفادة من نتائج التحليل كما هو مخطط له. أن دراسة القيم الشاذة Outlier Values) ربما لا تقتصر على اكتشاف القيم الشاذة، وإنما تتعدى ذلك لتشمل معالجة القيم الشاذة في البيانات وقد يعزى شذوذ البيانات كما اورتها (Barnett 1978⁽⁵⁾ الى الاسباب التالية:

1- أخطاء الحساب (Calculation Error)

2- أخطاء القراءة (Reading Errors) 3- أخطاء التسجيل (Recording Errors).
وإذا تم تحديد السبب او الاسباب في شذوذ القيم المذكورة اعلاه او غيرها سهلت عملية المعالجة.

وقد يتم تصحيح القيم الشاذة او يتم حذف البيانات المسيبة للاخطاء اما في حالة الشك بوجود قيم معينة شاذة دون دليل واضح عن سبب حدوثها فلابد من اعتماد الاساليب الاحصائية المناسبة لتحديد تباعدها (أبعاد القيم عن جسم البيانات)، وأن عملية حذفها في هذه الحالة سيؤدي الى نقص في المعلومات المتوفرة للدراسة المعينة، لأنها قد تمثل حالة خاصة او ظرفًا طارئاً وأنها قيم صحيحة وأن هناك عدة احتمالات لمعالجة تلك القيم الشاذة كما جاء في دراسة المختار، سليمان⁽⁴⁾ عام 1980 وهي:

1- رفض القيمة (القيم) الشاذة ثم اجراء التحليل الاحصائي على بقية البيانات او اللجوء الى تقدير تلك القيمة (القيم) ومن ثم اجراء التحليل الاحصائي المناسب على المجموعة الجديدة من البيانات.
2- ان وجود القيم الشاذة لا يعني في جميع الاحوال بأنه حالة سلبية فقد تعطي القيم الشاذة اهتماماً خاصاً لأنها يمكن ان تمثل حالة خاصة او ربما يكون هذا الوجود مؤشرًا جديداً لعامل جديد او دليلاً على تفسير البيانات وفق سياقات معينة قد تكون ذات اهمية بالغة وقد تسهم اسهاماً بالغاً في تفسير الظاهرة المدروسة.

3- معالجة البيانات كما هي دون استثناء أي من قيمها بواسطه الاساليب الاحصائية المناسبة كأن يوضع لها نموذجاً معدلاً يأخذ بنظر الاعتبار حالة القيم الشاذة ومدى تاثيرها على طبيعة مثل هذا النموذج. وتظهر القيم الشاذة بشكل نتيجة لبعض العوامل، بعضها يتم السيطرة عليه من قبل الباحث والبعض الآخر لا يمكن السيطرة عليه وقد صنفتها Barnett 1978 كما يلي:

1- أخطاء المعاينة (Sampling Errors) وهي الأخطاء الناتجة عن اختيار العينة.
 2- الاختلاف الأصلي (Inherent Variability) لا يمكن السيطرة على هذا النوع من الاختلاف لكونه من الخصائص الطبيعية الملزمة للمجتمع الأصلي.3- الخطأ الناتج من سبب طبيعي تكون القيم الشاذة في هذه الحالة نتيجة طبيعية للأختلافات الموجودة في قيم المتغيرات الأصلية.

وهناك عدة تعاريف للقيم الشاذة حيث عرفها (Bross,1961)⁽⁹⁾ على أنها المشاهدات التي تظهر منحرفة بشكل كبير عن جميع قيم العينة. كذلك عرفها (Kishpaugh,1972)⁽¹⁸⁾ بأنها المشاهدة التي يكون ابتعادها عن باقي قيم العينة يتحقق الاهتمام، وبذلك يعني أن لأي شخصي مهم في تحديد القيمة الشاذة. كما عرفها (AL-Jobouri,1976)⁽¹⁶⁾ بأنها تلك القيمة التي تكون غير متجانسة أو منسجمة مع بقية بيانات المجموعة لمتغير من المتغيرات لظاهرة معينة او مجموعة من الظواهر، او ان القيم الشاذة هي القيم التي تأتي من مجتمع يختلف عن مجتمع العينة قيد الدراسة. أما (Barnett,1978)⁽⁵⁾ فقد عرفت القيمة الشاذة في العينة، بأنها تلك القيمة التي تبدو غير منطقية اذا ما قورنت بباقي القيم.

أن معظم الدراسات التي تناولت القيم الشاذة كانت تعامل مع القيم الشاذة كحالة (Parametric Case) وتحت افتراض ان البيانات تتبع التوزيع الطبيعي (Normal Distribution) والمشكلة تظهر عندما لا تتوزع البيانات توزيعاً طبيعياً، ولأجل الابتعاد عن مشكلة تحديد التوزيع الاحتمالي الملائم للبيانات وللتطور الذي حصل باستخدام الحسابات والذي ادى الى دفع الاحصاء مئات السنين الى الامام، واكتشاف طرق حديثة لتحليل البيانات الاحصائية فأن الكثير من الدراسات في الفترة الاخيرة بدأت تأخذ منحنى اخر وهو التعامل مع الحالة الامثلية (Non-Parametric Case) والطرق الاستكشافية الحديثة (Exploratory Data Analysis).

1.2. الهدف من البحث

قبل البدء بعملية أي تحليل للبيانات الاحصائية لابد من التأكد من صحة وسلامة البيانات وخلوها من الأخطاء والقيم الشاذة، حيث ان نتائج البحث العلمي المبنية على احصاءات غير صحيحة لا يمكن الاعتماد عليها بأي شكل من الاشكال لأنها سوف تكون بعيدة عن الواقع. أن هدف هذا البحث هو استخدام طرق استكشافية جديدة (Exploratory Data Analysis) لتحديد القيم الشاذة في البيانات الاحصائية باستخدام الرسم وهي (Box plot, Stem and Leaf, and Rangefinder Box plot) وأجراء مقارنة بين الطرق الاستكشافية وبين الطرق المعملية لتحديد القيم الشاذة وايهما افضل بالتطبيق العملي.

1.3. الدراسات السابقة عن القيم الشاذة

بداءة فكرة وجود القيم الشاذة منذ منتصف القرن الثامن عشر عندما حاول (Boscovich,1755)⁽⁵⁾ تحديد أهلية الأرض معتقداً على معدل القياسات التي سيحصل عليها وقد استطاع الحصول على عشرة قياسات، أستبعد اثنين منها لنظرهما الشديد ثم وجد المعدل للثمانية الأخرى.

يعتبر (Peirce, 1852)⁽⁴⁾ اول من اعتمد اسلوب الاختبار للقيم الشاذة فقد نص اختباره على رفض (K) من المشاهدات الشاذة في عينة حجمها (n) اذا كان احتمال منظومة الاخطاء (System of Error) الناتج عن الاحتفاظ بـ (K) من المشاهدات الشاذة اقل من احتمال رفضها مضروباً بأحتمال الحصول على هذا العدد وليس اكثراً من المشاهدات وقد حدد (Peirce) الأحتمال الاخير بالصيغة التالية $\theta^{n-k} (1-\theta)^k$ اذا عرف (θ) بأنه احتمال وجود مثل هذه المشاهدات الشاذة المرفوعة بالنظر لكبر قيمتها ومن ثم حدد له القيمة n/k .

وأقترح (Wright, 1913)⁽³⁾ رفض أي مشاهدة شاذة تتحرف عن المتوسط بأكثر من ثلاثة أضعاف الانحراف المعياري.

بعد ذلك اقترح (Goodwin, 1913)⁽¹⁶⁾ رفض المشاهدات في عينة حجمها (n) اذا زاد انحرافها عن متوسط باقي المشاهدات ($n-1$) مشاهدة بأربعة اضعاف معدل انحراف ($n-1$) مشاهدة. ولم تكن جميع الاحصاءات المقترحة قبل عام 1925 تميز بين تباين المجتمع وتباين العينة. وكان (Irwin, 1925)⁽¹⁵⁾ اول من استطاع ان يفرق بين هذين النوعين من التباينات، اذا شدد على ضرورة استخدام الانحراف المعياري للعينة (s) كتقدير الى الانحراف المعياري للمجتمع عندما يكون الانحراف المعياري للمجتمع (σ) مجهولاً... وقد اقترح استخدام مؤشرات احصائية (عندما تكون σ معلومة). وأقترح جداول بالقيم الحرجية لهذه المؤشرات الاحصائية لتحديد تباعد القيم الشاذة وللتفصيل (أرجع دليل Thompson 1996). واقترح (Thompson 1935)⁽¹⁶⁾ في عام 1935 مؤشراً احصائياً يعتمد على اختبار T.

ثم تطورت دراسة القيم الشاذة لتشمل بالإضافة الى دراستها كقيم شاذة ضمن العينة المفردة دراستها في تصميم التجارب والانحدار والبيانات متعددة المتغيرات.....الخ.

وفي عام 1950 اقترح Grubbs⁽¹³⁾ بعض المعايير لاختبار معنوية الشذوذية للمشاهدة الكبرى في عينة بحجم n مسحوبة من مجتمع طبيعي (Normal Population)، كذلك اقترح معايير اخرى لاختبار معنوية الشذوذية لاختبار ما اذا كانت المشاهدات الكبرى الاولى والثانية في العينة متطرفتين بالكبر الى حد بعيد، او ان المشاهدتين الصغرى الاولى والثانية في العينة متطرفتين في الصغر الى حد بعيد...اما التوزيعات الاحصائية التي اقترحها Grubbs⁽¹³⁾ لاحصاءاته فقد اعتمد اشتقاها الى حالة المجتمع الطبيعي.

كما اقترح جداول بالنسبة المئوية (Percentage Points) التقديرية.

واقتراح Mosteller⁽²¹⁾ عام 1948 اختبار انزلاق (Silppage) K من العينات ثم حدد النسب المؤدية لحالات معينة وعندما تكون حجوم العينات متساوية الا ان العينات قد لا تكون متساوية الحجوم...لقد ادرك Mosteller & Tukey⁽²²⁾ في عام 1950 هذه النقطة واستطاعوا تحديد مستويات معنوية دقيقة لعينات مختلفة الحجوم باقتراهما صيغة تأخذ بنظر الاعتبار تلك الاختلافات في الحجوم.

درس Zinger⁽³⁰⁾ 1961 ظاهرة الشذوذية في عدة مجتمعات طبيعية (لها تباينات معلومة)... ثم اقترح مؤشراً احصائياً لاختبار شذوذية التباعد لمجتمعات يصل عددها الى سبع مجتمعات... فقد اعتمد هذا المؤشر على الفرق المعياري بين اكبر متوسطي عينتين، وقد اعطى جداول خاصة بالقيم الحرجية.

وقدم Mcmillan⁽²⁰⁾ عام 1971 اسلوباً جديداً في اختباره لشذوذية مشاهدة واحدة او مشاهدتين في عينات طبيعية عندما يكون تباين المجتمع مجهولاً ويتضمن اسلوب Mcmillan ثلاثة طرق...

الطريقة الاولى وتتضمن التنفيذ المتسلسل لاحصاءة الباقى الاعظم (Maximum Residual). اما الطريقة الثانية فهي تعتمد على اختبار ما اذا كانت المشاهداتان الكبرى الاولى والثانية في العينة شاذتين.

اما الطريقة الثالثة والاخيرة فقد اعتبرت المشاهدين الكبرى الاولى والثانية شاذتين ولمزيد من التفاصيل ارجع (دلير 1996). وكذلك اقترح احصاء لاختبار التباعد لمشاهدة واحدة او مشاهدين في عينات طبيعية على أن يكون تباين المجتمع معلوماً.

وقد لاحظ كل من Gnanadeskan and Kettenring⁽¹²⁾ في عام 1972 ان الاعتبارات المأخوذة للقيم الشاذة في عينة ذات المتغيرات المتعددة تكون اكثر تعقيداً منها في حالة المتغير الواحد.

ان اهداف Gnanadeskan and Kettenring الاولية هي تكوين تقديرات قوية لموقع المتغيرات المتعددة (Multivariate Locations) وتحديد الاستجابات المتعددة لقيم الشاذة. ولقد تم اقتراح تصميم لرسم نقاط البيانات في الواجهات الذاتية (Eigen Vectors) لمتوسطات البيانات المعروضة لاكتشاف المشاهدات الشاذة المحتمل وجودها.

وقد عدل Kishpaugh⁽¹⁸⁾ في عام 1972 على هذا الاقتراح وسميت بتصاميم المركبات الأساسية. واقتراح كل من Tietjen, Moore and Beckman⁽²⁷⁾ في عام 1973 اسلوباً لتحويل الاخطاء الى اخطاء معيارية (Standard Residual) لاختبار التباعد لنموذج الانحدار الخطي البسيط وقد حددوا القيم الحرجية بالاعتماد على دراسة المحاكاة، وقد وجدوا ان هذه القيم مقارنة للقيم التي حصل عليها Grubbs.

اقتراح Rohlff⁽²⁴⁾ في عام 1975 صيغة لاكتشاف القيم الشاذة المتعددة (Multivariate Outliers) وسميت باختبار الفجوة العام (Generalized Gap Test)، حيث ذكر ان المشاهدات في حالة متعدد المتغيرات تتخذ شكل شجرة لها فروع متجمعة تقريباً وهناك مجموعة تكون خارجة عنها.

وقد اقترح كل من John and Prescott⁽²⁴⁾ في عام 1975 عدداً من الاحصاءات لاختبار التباعد في تصاميم التجارب ثم حدد القيم الحرجية لتلك الاحصاءات باستخدام اسلوب المحاكاة. ناقش AL-Jobouri⁽¹⁶⁾ في عام 1976 الطرق المعلمية لتحديد القيم الشاذة في حالة متغير متغيرين وعدة متغيرات في تصاميم التجارب وقد اعتمد Barnett and Lewis⁽⁵⁾ في عام 1978 دراسة القيم الشاذة بفرض ان البيانات مأخوذة من توزيع طبيعي متعدد.

استعرض المختار، سليمان في عام 1980 مختلف الطرق التي عالجت القيم الشاذة، في التجارب المصممة ونماذج الانحدار وكذلك متعدد المتغيرات.

اقتراح (الجبوري، 1988)⁽¹⁾ طريقة لاكتشاف المشاهدات الشاذة في حالة متعدد المتغيرات بأعتماد طريقة الرسم الصندوقى Box plot.

وقد اقترحت (الجبوري، منى 1990)⁽²⁾ طريقة لاكتشافالجزئي للمشاهدات الشاذة وطرق التقدير في حالة متعدد المتغيرات وبأعتماد طريقة الرسم الصندوقى Box plot .

واخيراً اقترحت (المشنو، 1993)⁽³⁾ طريقة لاكتشاف المشاهدات الشاذة في تحليل تصاميم التجارب غير المتزنة وبأعتماد طريقة الرسم الصندوقى Box plot .

نجد مما تقدم هناك العديد من الدراسات تناولت القيم الشاذة وسوف نتطرق في البند اللاحق بأختصار الى مفهوم الطرق الاستكشافية .

1.4 الطرق الاستكشافية الحديثة لتحليل البيانات الاحصائية

Exploratory data analysis statistical data analysis

تعبر هذا الطرق والاساليب التي اغلبها حديثة وتستخدم لعرض وتحليل البيانات. العالم كان له Tukey

الدور الكبير باستكشاف وتطوير هذه الاساليب واعطى الفكرة الاولية عليها عام 1970 وهذه، الاساليب يمكن ان تعطينا فكرة واضحة عن توزيع واتجاه البيانات، فهي تفصل وتشخص عناصر مكونات البيانات الاحصائية المهمة الى المحل. ان المحللين الاحصائيين الجيدين في بادئ الامر يعرضون البيانات بشكل تفصيلي قبل البدء بتحليلها واختيارها للاطلاع على مكوناتها وعلى اتجاه توزيعها، وهذه الطرق تتعامل مع البيانات بمرونة عالية في التحليل ويمكن ان تخدم هذه الطرق في المرحلة الاولى من التحليل لتحديد المقاييس المناسبة والكافحة لوصف البيانات، ويمكن ان تستخدم هذه الطرق لتحويل البيانات لاستخدام المقاييس المناسب بعد ان تقترب او تتوزع توزيع طبيعي. وقد صممت هذه الطرق لاجراء مقارنات بسيطة او بشكل تفصيلي بين توزيعات البيانات او اجزائها الرئيسية المهمة.

اما الحالات التي لا تحتاج في التحليل للبيانات عرض كل الاجزاء، يمكن ان تأخذ هذه الطرق ملخصاً لعرض توزيع البيانات وهناك قسم من هذه الاساليب يمكن ان تجز هذه المهمة بسهولة او بشكل جيد جداً مثل القيم الحرفية (latter Values) يمكن ان تعرض لنا خمس قيم من البيانات وهو الوسيط (Median) والرابع الاول (First Quartile) والرابع الثالث (Third Quartile) واكبر واصغر قيمة من البيانات وهذا الاسلوب بسيط جداً بحيث يمكن عمله بسهولة. وتوجد طريقة اخري مهمة جداً لتلخيص البيانات تعرف بالرسم الصندوفى (Box plot) تعطي هذه الطريقة انطباعاً بشكل سريع لاجزاء محدودة ومهمة من التوزيع ولشكل انتشار البيانات وعلى شكل رسم بياني. وهذا الاسلوب يعتمد على خمس قيم من البيانات ويمكن استخدام الرسم الصندوفى لاجراء مقارنات متعددة لعدة مجاميع من البيانات.

1.4.1 الطرق الحسابية والطرق الاستكشافية لتحليل البيانات⁽¹⁴⁾

ان الطرق الاستكشافية هي اكثر فائدة من الطرق الحسابية لتحليل البيانات وتحديد اجزاء محددة من البيانات، اما الطرق الحسابية فهي اكثر فائدة من الطرق الاستكشافية في حالة صنع قرارات عامة تعتمد على نتائج البيانات، على سبيل المثال بناء نموذج من العينة واستخدامه بالمجتمع.

والطرق الاستكشافية تعتمد نتائجها على القرارات الشخصية للأشخاص وتفسيراتهم وهذا يمكن ان يكون هناك فروقات بسيطة بين تفسير النتائج من شخص الى اخر بالاعتماد على هذه الطرق، ولكن الطرق الحسابية تعتمد على ملخص التحليل للبيانات الاحصائية المحددة على سبيل المثال الوسط الحسابي والانحراف المعياري لایتأثران بتصورات الاشخاص او تفسيراتهم.

1.4.2 طرق الرسم (Graphical Methods) (14))

كانت الحاجة لترتيب البيانات على شكل جداول او على شكل رسوم بيانية مع بداية الثورة الصناعية في اوربا لكثرة وجود البيانات وكان اكتشاف هالي Halliy طريقة الرسم لتحليل ضغط الباراميتر وهذا الاكتشاف قاد الى استخدام طريقة الرسم لعرض وتحليل البيانات الاحصائية، بعد ذلك تم اكتشاف طرق الرسم واحدة بعد الاخرى ولحد الان وبعد عام 1820 أصبح استخدام طرق الرسم معروفاً وشائعاً بأغلب المقالات والمجلات وكان لها وجود مثلاً في المؤتمر العلمي لللاحصاء الذي انعقد في فينا 1852 (International statistical Congress) وفي عام 1872 خصص الكونغرس الامريكي مبلغ من المال لأول مرة لتطوير طريقة الرسم في تحليل البيانات. العقدین الاخيرین شهد اهتمام واسع جداً بطرق الرسم لتحليل البيانات الاحصائية وكان Tukey الاول في هذا المضمار حيث اكتشف عدة طرق لتحليل وتفسير البيانات الاحصائية، وتوجد الان حاجة ماسةً لتحليل البيانات الاحصائية بالرسم لأن الرسم يساعد على تحليل البيانات وعرض العلاقات النظرية لها. وان الرسم هو الصورة للبيانات التي تعطي فكرة دقيقة عن المعلومات لهذه المجموعة من البيانات ولعبت الحاسيبات دوراً كبيراً في تطوير واكتشاف طرق الرسم، ولهذا تعتبر طرق الرسم الان مهمة جداً في تحليل البيانات الاحصائية وجزءاً اساسياً يربط بين العلم والتكنولوجيا وسوف يتم التطرق في الفصل الثاني الى اهم طرق الرسم.

2. الجانب النظري

2.1 المقدمة :

تم التطرق في الفقرة (1) الى بعض المفاهيم والدراسات السابقة عن القيم الشاذة ومفهوم الطرق الاستكشافية، وسوف يتم التطرق في هذه الفقرة الى الجانب النظري للبحث ويكون على مبحثين، المبحث الاول يتناول عرض بعض الطرق المعلمية (Parametric) المستخدمة لتحديد القيم الشاذة في البيانات، اما المبحث الثاني فيتناول عرض الطرق الاستكشافية الحديثة (Exploratory Data Analysis) التي تستخدم لتحديد القيم الشاذة في البيانات.

2.2 الطرق المعلمية

تم التطرق في هذا المبحث الى الطرق المعلمية التي تستخدم لتحديد القيم الشاذة في البيانات الاحصائية في حالة المتغير الواحد (Univariate case) والمتغيرين (Two Dimensions) ونمذج الانحدار (Regression models) التي تتطلب معرفة التوزيع الاحتمالي للبيانات الاحصائية.

2.2.1 الطرق العلمية المستخدمة لتحديد القيم الشاذة في حالة المتغير

الواحد (Vnivariate)

اولاً : طريقة (Grubbs) (13).

نشر Grubbs عام 1950 بحثاً حول مقياس العينة لاختبار المشاهدات الشاذة حيث اقترح فيه بعض المقاييس لاختبار دلالة المشاهدة الشاذة في عينة بحجم (n) تتوزع توزيعاً طبيعياً وقد استخدم عدة مقاييس لاختبار القيم الشاذة بعد ترتيب المشاهدات في العينة تصاعدياً حيث تم استخدام المقياس في الصيغة رقم (2-1) لاختبار دلالة اكبر مشاهدة في العينة بحجم (n) فيما اذا كانت شاذة ام لا. اما بالنسبة للمشاهدة الصغيرة فقد استخدم الاحصاء او الاختبار في الصيغة رقم (2-2)، ولاختبار فيما اذا كانت اكبر مشاهدتين شاذتين ام لا فقد استخدم الاحصاء او الاختبار في الصيغة رقم (2-3)، وبالنسبة الى اصغر مشاهدتين فقد استخدم الاحصاء او الاختبار في الصيغة رقم (2-4) وهذه الصيغ مدرجة أدناه. ونقارن هذه الاحصاءات او قيم الاختبارات بما يقابلها من القيم الحرجية في جداول خاصة، ولمزيد من التفاصيل ارجع الى (ديلر 1996).

$$\frac{S_n^2}{S^2} = \frac{\sum_{i=1}^{n-1} (X_i - X_n)^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \quad \dots \dots \dots (2-1)$$

$$\bar{X}_n = \frac{1}{n-1} \sum_{i=1}^{n-1} X_i \quad \text{عندما}$$

$$\frac{S_1^2}{S_2^2} = \frac{\sum_{i=1}^n (X_i - \bar{X}_1)^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \quad \dots \dots \dots (2-2)$$

$$\bar{X}_1 = \frac{1}{n-1} \sum_{i=2}^n X_i \quad \text{عندما}$$

$$\frac{S_{n,n-1}^2}{S^2} = \frac{\sum_{i=1}^{n-2} (X_i - X_{n,n-1})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \quad \dots \dots \dots (2-3)$$

عندما

$$\bar{X}_{n,n-1} = \frac{1}{n-2} \sum_{i=1}^{n-2} X_i$$

$$\frac{S_{1,2}^2}{S^2} = \frac{\sum_{i=3}^n (X_i - \bar{X}_{1,2})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \quad \dots \dots \dots (2-4)$$

عندما

$$\bar{X}_{1,2} = \frac{1}{n-2} \sum_{i=3}^n X_i$$

⁽²³⁾ (Quesenberry and David) طريقة ثانياً :

انصب عمل كل من (Quesenberry and David) في عام 1961 على تحديد القيم الشاذة لـ k أكثر من مجموعة من المشاهدات (k من المجتمعات) لها توزيع طبيعي شريطة أن يكون التباين σ^2 للمجتمعات غير معروف وقد استخدم المقاييس التالية في الاختبارات والتي ارقام صيغها كما يلى.

(2-5)، (2-6)، (2-7) حيث ان S^2 تساوي مجموع مربعات الانحرافات لـ K من المجتمعات. استخدم (Quesenberry and David) الاحصائين (2-6),(2-7) اثناء معالجتها لمسألة الانزلاق (التفويت) للمجتمعات الطبيعية. وبعد ان يتم اختيار احدى القيم (المشاهدات المشكوك فيها شاذة) لأي مجتمع كان واجراء الاختبار عليها لتحديد شاذة ام لا، فإن الاجراء الاتي سيكون باستبعاد المشاهدة الشاذة من أي مجموعة كانت ومن ثم الاستمرار بعملية الاختبار للقيم الأخرى ، حسب الصيغ التالية ارقامها كما يلي: (8-2),(9-2),(10-2) عندما S_1 تساوي مجموع مربعات الانحرافات لـ K من المجتمعات محسوب بدون المشاهدات الشاذة، وهكذا يستمر الاختبار وتقارن هذه الاحصاءات بما يقابلها من القيم الحرجية في جداول خاصة⁽¹⁶⁾ وهذه الصيغ مدرجة أدناه. وللاطلاع على كل التفاصيل ارجع الى (ديبلر 1996).

$$b = \text{Max}(b_i) = \frac{X_{\max} - \bar{X}}{S^*} \quad \dots \dots \dots (2-6)$$

$$b_i = \left(\sqrt{n_i \bar{X}_i} - X_w \right) / S_1^* \quad (2-8)$$

$$b = \text{Max}(b_i) = \frac{\text{Max} \sqrt{n_i \bar{X}_i} - \bar{X}_w}{S_i^*} \quad \dots \dots \dots (2-9)$$

$$b^* = \text{Max} |b_i| = \text{Max} \frac{\sqrt{n_i X_i} - \bar{X}_W}{S_i^*} \quad \dots \dots \dots \quad (2-10)$$

عندما

و

$$\overline{X}_W = \frac{1}{k} \sum_{i=1}^k \sqrt{n_i \overline{X}_i}$$

$$N = \sum_{i=1}^k n_i$$

ثالثاً: طريقة (Mcmillan) ⁽²⁰⁾

قدم (Mcmillan 1970) طريقة لاختبار واحدة او اكثر من المشاهدات الشاذة لمجتمع يتوزع توزيعاً طبيعاً بمتوسط μ وتبالين σ^2 فقد اعتمد Mcmillan ثلاثة اجراءات لمعالجة البيانات المشكوك فيها والتي تحتوي على مشاهدات شاذة متعددة، الاجراء الاول عبارة عن تطبيق متسلسل لاختبار الباقي الاعظم (Maximum Residual Test) فاذا كانت القيمة المحسوبة في الصيغة (2.11) اكبر من الجدولية كما موضح في الصيغة فأن X_n يمكن اعتبارها مشاهدة شاذة ثم يتكرر الاختبار على بقية المشاهدات عندما $V_\alpha^{(n,v)}$ عبارة عن قيم حرجة يمكن الحصول عليها من جداول خاصة و S تمثل الانحراف المعياري للعينة. فاذا كانت ايضاً القيمة كما في الصيغة (2-12) اكبر من القيمة الجدولية فأن X_{n-1} يمكن اعتبارها مشاهدة شاذة ايضاً.

الاجراء الثاني لـ Mcmillan يتمثل في اعتبار القيمتين X_{n-1}, X_n شاذتين اذا كانت القيمة المحسوبة في الصيغة (2-13) اكبر من القيمة المحسوبة $(C_\alpha^n.S)$ عندما C_α^n تمثل قيمة حرجة معطاة في جداول خاصة.

اما الاجراء الثالث مشابه لما اقترحه Grubbs وهو كما في الصيغة (2-14) وتقارن هذه الاحصاءة مع القيم الحرجة في جداول خاصة وهذه الصيغ مدرجة أدناه.

$$X_n - \overline{X} > V_\alpha^{(n,V)} . S \quad \dots \dots \dots \quad (2-11)$$

$$X_{n-1} - \overline{X} > V_\alpha^{(n-1,V)} . S_n \quad \dots \dots \dots \quad (2-12)$$

عندما

$$S_n^2 = \sum_{i=1}^{n-1} (X_i - \overline{X}_n)^2 / (n-2)$$

$$\overline{X}_n = \sum_{i=1}^{n-1} X_i / (n-1) \quad \text{و} \quad$$

$$X_n + X_{n-1} - 2\overline{X} > C_\alpha^{(n)} . S \quad \dots \dots \dots \quad (2-13)$$

$$\frac{S_{n,n-1}^2}{S^2} = \frac{\sum_{i=1}^{n-2} (X_i - \overline{X}_{n,n-1})^2}{\sum_{i=1}^n (X_i - \overline{X})^2} \quad \dots \dots \dots \quad (2-14)$$

رابعاً : طريقة (Tietjen and Moore) ⁽²⁷⁾

اقترح كل من (Tietjen and Moore 1973) أحصاءتين L_k التي تعتمد على k من القيم (المشاهدات) الكبيرة المشكوك فيها في العينة (n) ، E_k التي تعتمد على الباقي العظمى (Largest Residuals) بقيمة مطلقة، وبعد ترتيب المشاهدات للعينة الطبيعية تصاعدياً وأختيار k من المشاهدات الكبيرة المشكوك فيها نستخدم الاحصاء الآتية في الصيغة (2-15): واستخدام الاحصاء L_k^* ايضاً في اختبار k من المشاهدات الصغيرة المشكوك فيها وبالصيغة المرقمة (2-16).

اما الاحصاء E_k الصيغة (2-17) فقد اقترحت من قبل (Tietjen, Moore) لدعم وتأكيد على صحة ما تفرضه الاحصائيتان L_k , L_k^* . والاحصائيتان L_k , L_k^* تستخدمان على النحو الآتي اذا قررنا في عينة من

حجم n ان نختبر فيما اذا كانت k من القيم الكبيرة او الصغيرة المشكوك فيها شاذة ام لا فاننا نحسب L_k^* , L_k فإذا كانت قيمة هاتين الاحصائيتين اصغر من القيمة الحرجية المرغوبة فاننا نستنتج ان k من هذه القيم الكبيرة او الصغيرة هي بالفعل شاذة، واذا حسبنا E_k ثم قارناها بالقيمة الحرجية، فإذا كانت اصغر من القيم الحرجية المختارة فمن المشاهدات المشكوك فيها شاذة، وتقارن الاحصاءات E_k , L_k^* , L_k بما يقابلها من القيم الحرجية في جداول خاصة وهذه الصيغ مدرجة أدناه.

$$L_k = \frac{\sum_{i=1}^{n-k} (X_i - \bar{X}_k)^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \quad \dots \dots \dots \quad (2-15)$$

عندما

$$\bar{X}_k = \sum_{i=1}^{n-k} X_i / (n - k)$$

$$L_k^* = \frac{\sum_{i=k+1}^n (X_i - \bar{X}_k)^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \quad \dots \dots \dots \quad (2-16)$$

عندما

$$\bar{X}_k^* = \sum_{i=k+1}^n X_i / (n - k)$$

$$E_k = \frac{\sum_{i=1}^{n-k} (Z_i - \bar{Z}_k)^2}{\sum_{i=1}^n (Z_i - \bar{Z})^2} \quad \dots \dots \dots \quad (2-17)$$

$$\bar{Z}_k = \sum_{i=1}^{n-k} Z_i / (n - k) \quad \text{عندما}$$

$$Z_i = |di| = |X_i - \bar{X}| \quad \text{و}$$

خامساً : طريقة (Rosner) ⁽²⁵⁾
 اقترح (Rosner,1975) لتحديد اختبار اكثراً من مشاهدة شاذة المشكوك فيها الصغيرة او الكبيرة او الاثنان معاً كعينة من التوزيع الطبيعي بحجم (n) والاختبار كما في الصيغة (2-18).
 ولاختبار معنوية الاحصاءات المحسوبة -- $R_1, R_2 \dots R_K$ تقارن مع القيم الحرجة في جداول خاصة ⁽¹⁶⁾.

$$R - statistic = |Max(X_i) - a| / b^2 \quad \dots \dots \dots (2-18)$$

عندما

$$a = \sum_{i=k+1}^{n-k} X_i / (n - 2k)$$

$$b^2 = \sum_{i=k+1}^{n-k} (X_i - a)^2 / (n - 2k - 1) \quad \text{و}$$

$$I_o = (X_1, \dots, X_n)$$

$$R_1 = |X^{(o)} - a| / b, X^{(o)} \in I, \quad \dots \dots \dots (2-19)$$

$$I_1 = I_0 - X^{(o)}$$

$$R_k = |X^{(1)} - a| / b \quad \dots \dots \dots (2-20)$$

$$R_i = |X^{(k-1)} - a| / b \quad X^{(k-1)} \in I_{k-1} \quad \dots \dots \dots (2-21)$$

2.2.2 الطرق العلمية المستخدمة لتحديد القيم الشاذة في حالة المتغيرين (Two Dimensions) ونماذج الانحدار.

اولاً : اختبار T^2 . Hotelling

اقترح (AL-Bayati ⁽⁶⁾) في عام 1970 هذا الاختبار لمعرفة فيما اذا كان الزوج (X_o, Y_o) المشكوك فيه من المشاهدات شاذًا أم لا على النحو التالي:-

لتكن n عينة من ازواج المشاهدات (X_i, Y_i) حيث ان $(i = 1, \dots, n)$ ، مأخوذة من مجتمع طبيعي ثانى

الاتية (Bivariate Normal Distribution) له معالم محددة. اقترح AL-Bayati اختبار الفرضية

$$H_0 : P = P'$$

$$H_1 : P \neq P'$$

لاختبار الزوج المشاهد (X_o, Y_o) المشكوك فيه شاذة ام لا حيث أن ' P' يمثل معامل الارتباط محسوب بدون القيمة المشكوك فيها (X_o, Y_o) ، وباستخدام احصاءة القطع الناقص (E) كما في الصيغة (2-22). وتقديرها للعينة (\hat{E}) كما في الصيغة (2-23). والاحصاءة \hat{E} لها توزيع T^2 المستخدم عادة لاختبار متجه من المتوسطات L من المتغيرات المرتبطة هنا في هذه الحالة $P=2$ ، وباستخدام العلاقة بين الـ T^2 Hotelling وتوزيع F عند مستوى α ، والتي هي كما في الصيغة (2-24) وأن قيمة F_α يمكن الحصول عليها من جداول F بمستوى معنوية α ودرجة حرية $(n-2)$ (و عند اختبار ازواج القيم مثلاً (X_o, Y_o) فإن الاحصاءة (2-23) تصبح كما في

الصيغة (25-2)، ونقارن على النحو التالي، اذا كانت $T_\alpha^2 > T^2$ نرفض (X_o, Y_o) على انها مشاهدة شاذة واذا كانت $T_\alpha^2 \leq T^2$ تقبل (X_o, Y_o) على انها واحدة من المشاهدات الشاذة وهذه الصيغ مدرجة أدناه.

$$E = \frac{1}{1-p^2} \left[\frac{(X-u_x)^2}{\sigma_x^2} + \frac{(Y-u_y)^2}{\sigma_y^2} - \frac{2p(X-u_x)(Y-u_y)}{\sigma_x \sigma_y} \right] \dots \dots \dots (2-22)$$

وتقديرها للعينة

$$\hat{E} = \frac{1}{1-r^2} \left[\frac{(X-\bar{X})^2}{S_x^2} + \frac{(Y-\bar{Y})^2}{S_y^2} - \frac{2r(X-\bar{X})(Y-\bar{Y})}{S_x S_y} \right] \dots \dots \dots (2-23)$$

$$T^2 = \frac{P(n-p)}{n-(p+1)} F_\alpha \dots \dots \dots (2-24)$$

وبوضع $p=2$ تصبح

$$T^2 = \frac{2(n-2)}{n-3} F_\alpha$$

$$T_\alpha^2 = \frac{1}{1-r^2} \left[\frac{(X_o-\bar{X})^2}{S_x^2} + \frac{(Y_o-\bar{Y})^2}{S_y^2} - \frac{2r(X_o-\bar{X})(Y_o-\bar{Y})}{S_x S_y} \right] \dots \dots \dots (2-25)$$

ثانياً : القيم الشاذة في نماذج الانحدار (Outliers in a Regression Models) a- القيم الشاذة في النموذج الخطى البسيط (Outliers in a Regression Models)⁽⁴⁾ تصاغ معادلة النموذج الخطى البسيط كما يلى

$$Y_i = B_o + B_1 X_i + U_i \dots \dots \dots (2-26)$$

حيث Y_i يمثل المتغير المعتمد (Dependent Variable) في النموذج و χ_i يمثل المتغير المستقل (Independent Variable) في النموذج، B_o, B_1 تمثل معالم (Parameters) النموذج و U_i يمثل الخطأ العشوائي (Random Error) وان

$$U_i \sim N(O, \sigma^2), \forall i = 1, \dots, n \quad \text{cov}(\mu_i, \mu_j) = o \quad \forall i \neq j$$

فإذا قررنا القيم B_1, B_o باحدى طرق التقدير ولكن طريقة المرربعات الصغرى (Least Square Estimations) ستكون صيغة المعادلة (2-26) على النحو الآتى:

$$e_i = y_i - b_o - b_1 X_i \dots \dots \dots (2-27)$$

للبحث عن القيم الشاذة في الانحدار الخطى البسيط فإنه من المناسب دراسة اختبار حجم e_i وخواصهما والصيغة (2-28) لحساب التباين (e_i) V الى الخطأ. ونستنتج ان القيم الشاذة لها تأثير كبير على الانحراف المعياري للاخطاء ولهذا فإن اختبار تباعد القيم يعتمد على الاخطاء وكما في الصيغة (2-29) وهذه الصيغة مدرجة أدناه .

$$V(e_i) = \sigma^2 \left[\frac{n-1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right] \quad \dots \dots \dots (2-28)$$

$$\frac{e_i}{s_i} = e_i / s \sqrt{\frac{n-1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}} \quad \dots \dots \dots (2-29)$$

وقد اقترح Daniel⁽¹¹⁾ في عام 1960 استخدم الصيغة $|Max ei| / s$ اساساً لتحديد واختبار تباعد قيمة شاذة واحدة ... في حين اتخاذ كل من Behnken and Draper⁽⁸⁾ في عام 1972 الاحصاء في المعادلة المرفقة (2-30) لاختبار تباعد قيمة شاذة واحدة في النموذج البسيط، فإذا كانت قيمة الاحصاء في المعادلة (t) كبيرة بشكل كاف، عند ذلك سنحكم على المشاهدة انها قيمة شاذة متباعدة (Discordant) Outlier في حين ان الاستخدام العلمي لهذه الاحصاء يقضي بمعرفة توزيعها وهذا امر يصعب تحقيقه، لذلك اعتمد كل من Moore and Beckman⁽¹²⁾ اسلوب المحاكاة حيث قاما بعمل الالاف من التجارب لعينات يصل حجمها $n=100$ وقد استخدما مستويات المعنوية ($\alpha = 0.1, 0.05, 0.01$) وفي النهاية استطاعا تحديد القيم الحرجية للاحصاء t واستنتجوا ان القيم التي سنحصل عليها عند افتراض تساوي تباينات البوافي سوف لا تختلف كثيراً عن القيم التي ستحصل عليها بعد ان تؤخذ حالة عدم تساوي تباينات الاخطراء بنظر الاعتبار كما في الصيغة (2-30).

$$t = Max|e_i|/S_i \quad \dots \dots \dots (2-30)$$

ثم اقترح Prescott⁽¹⁷⁾ في عام 1975 أهمال التباينات المختلفة للاخطاء المقدرة e_i واحلال \bar{S}_i محل S_i في الاحصاء للمعادلة (2-30) حيث \bar{S} تمثل معدل التباين.

وبين كل من Behnken and Draper⁽⁸⁾ في عام 1972 ان معدل التباين للاخطاء غير المتحيز هو $(n-q)\sigma^2/n$ حيث ان (q) تمثل عدد المعالم (Parameters) في النموذج (وهنا $2 = q$) وعليه فأن تقدير معدل التباين سيكون كما في المعادلة (2-31) وبناءً على ذلك فأن (t^*) بعد الاشتغال سوف تكون كما في الصيغة المرفقة (2-32). اما بالنسبة لقيم الحرجية للاحصاءة تبين t^* فيمكن الحصول عليها من جداول خاصة والصيغ كما مدرجة أدناه (وللمزيد من التفاصيل انظر ديلر 1996).

$$t^* = \text{Max} |e_i| \Big/ \sqrt{(n-q)S^2/n}$$

$$= \text{Max} |e_i| \Big/ \sqrt{\sum_{i=1}^n e_i^2 / n}$$

$$= \sqrt{n} \text{ Max} |e_i| \Big/ \sqrt{\sum_{i=1}^n e_i^2}$$

b- القيم الشاذة في التموذج الخطى العام⁽⁴⁾

Outliers in a General Linear Models.

تعرف معادلة النموذج الخطى العام بصفة المصفوفات كالاتى:-

$$Y = XB + U \quad \dots\dots\dots(2-33)$$

حیث ان

Y : يمثل متجه المتغير المعتمد

$n \times 1$ بدرجة (Vector of Dependent Random variables)

X : تمثل مصفوفة المتغيرات المستقلة

$[n \times (k+1)]$ بدرجة (Matrix of Independent Random variables)

B : يمثل متجه معلمات النموذج (Vector of Parameters) بدرجة $(k+1) \times 1$

U : يمثل متوجه الأخطاء العشوائية (Vector of Random Errors) بدرجة $(n \times 1)$.

أن تقدير معلم النموذج B للمعادلة (2-33) بطريقة المربعات الصغرى هو كما في المعادلة (2-34). وتقدير تباينه كما في صيغة المعادلة (2-35). أما تقدير الخطأ (e_i) فهو كما في المعادلة (2-36) وعليه فإن تقدير تباين الخطأ يكون كما في صيغة المعادلة (2-37) ومن المعادلة (2-38) نلاحظ أن الأخطاء المقدرة (e_i) ستمتلك تباينات مختلفة مع وجود ارتباطات فيما بينها ويمكن توضيحها كما في الصيغة (2-39) وبناء على ذلك يمكن تقدير مصفوفة التباين والتباين المشترك للخطأ ($V-\text{COV}(e)$) كما في المعادلة (2-40) وسيكون التباين المقدر إلى الخطأ لـ (e_i) كما في المعادلة (2-41) والمعادلات مدرجة كما في الصيغ أدناه. ولتفاصيل ارجع إلى ديلر (1996).

$$\hat{B} = (X'X)^{-1} X' Y \quad \dots \dots \dots (Q-34)$$

$$V = \text{Cor}(\hat{B}) = (X'X)^{-1}\sigma^2 \quad \dots \dots \dots \quad (35)$$

وعلیه فأن تقدیر تباین الخطأ يكون

$$V - Cor(\underline{e}) = (I - R)\sigma^2 \quad \dots \dots \dots \quad (2-37)$$

من المعادلة (37-2) نلاحظ ان الاخطاء المقدرة e_i ستمتلك تباينات مختلفة مع وجود ارتباطات فيما بينها ويمكن توضيحها كالتالي

$$\begin{aligned}Var(e_i) &= \left[1 - X_i' (X'X)^{-1} X_i\right] \sigma^2 \\&= (1 - r_{ii}) \sigma^2\end{aligned}\dots\dots\dots(2-38)$$

حيث يمثل X_i الصف i من المصفوفة X وان

$$r_{ii} = X_i' (X X')^{-1} X_i$$

و σ^2 مجهولة القيمة بصورة عامة وان التقدير غير المتحيز لها هو

$$S_e^2 = \underline{e}' \underline{e} / (n - q) \\ = \underline{U}' [\underline{I} - \underline{R}] \underline{U} / (n - q) \quad \dots \dots \dots (2-39)$$

حيث ان : الكمية $[I - R]$ تمثل مصفوفة صماء [Idempotency Matrix] وبناء على ذلك يمكن تقدير (e) $V\text{-Cov}$ كالاتي:

$$S_e^2 = (I_n - r)\hat{\sigma}^2 \quad \dots \dots \dots \quad (2-40)$$

وبما أنه الاخطاء المقدرة (e_i) تمتلك تباينات مختلفة فيجب استخدام صيغة أخرى تأخذ بنظر الاعتبار عدم مساواة التباين، والصيغة المناسبة هي $M_i = e_i / s_i$ وتعتمد اغلب الاحصاءات المقترحة على (Maximum studentized Residual) حيث تعتبر المشاهدة التي تعطينا اكبر قيمة من قيم (m_i) مشاهدة شاذة متباعدة.

أحصاءة Strikantant ⁽²⁶⁾ عام 1961 احدى الاحصاءات المقترحة لهذا الغرض ووصيغتها كما في المعادلة (42-2). وتقارن الاحصاءة T لهذه المعادلة مع القيم الحرجية في جداول خاصة اعدها Lund ⁽¹⁹⁾ عام 1975 لهذا الغرض.

$$T = \text{Max} |e_i| S_i > n_\alpha \quad \dots \dots \dots \quad (Q-42)$$

2.3 الطرق الاستكشافية Exploratory Data Analysis

يتم التطرق في هذا المبحث الى ثلات طرق من الطرق الاستكشافية الحديثة (Exploratory Data Analysis) (التي تستخدم لتحديد القيم الشاذة في البيانات وهي طريقة العرض والورقة (steam and leaf)، الرسم الصندوفى (Box plot)، وطريقة الرسم الصندوفى المزدوج للمتغيرين (Rangefinder Box plot) (χ , y).

٢.٣.١ الغصن والورقة (Steam and Leaf)

من الطرق الحديثة لعرض البيانات الاحصائية هي طريقة الغصن والورقة، وهذا الاسلوب من العرض اسهل عند تكوينه من جداول التوزيع التكراري وكذلك من المدرج التكراري وبصورة عامة فهو يعرض معلومات اكثـر.

فهو يعرض نفس معلومات المدرج التكراري (histogram) وكذلك يعرض معلومات جداول التوزيع التكراري بالإضافة الى ذلك فهو يعرض الارقام بشكلها الاعتيادي عند ربط كل قيمة بين الغصن والورقة الخاص بها لهذا يسمى بالشكل الهجين (hybrid) لأن معلوماته تمثل الرسم وكذلك الارقام في الجدول في أن واحد كما هو موضح لاحقاً. لتوضيح الغصن والورقة لاحظ المثال التالي:-

مثال ١ : اوجد الغصن والورقة للمشاهدات التالية

((100 98 97 96 95 94 93 92 90 89 88 87 86 85 84 83 82 80 78 77 76 75 74 73 72 71 70 69 68 67 66))

الحل : بناء او تكوين الغصن والورقة والذي هو في نفس الوقت يوضح البيانات على شكل مجاميع تشبه جدول التوزيع التكراري وكذلك يعرضها على شكل رسم يشبه المدرج التكراري تقوم بالخطوات التالية:-

١. نختار ارقام من البيانات على انها الارقام التي تكون امام البيانات او تقود البيانات والتي تشمل الجزء الاول من الارقام والتي تمثل العشرات وهذا ينتج الارقام التالية (10..... 6..... 5..... 4..... 3..... 2..... 1..... 0).

نضعهم على شكل عمود كما في الشكل اللاحق (الغصن steam).

٢. بعدها نبدأ بالمرور على جميع الارقام لكتابـة الجزء الثاني من كل رقم مقابل الجزء الاول منه (Final digit) ونضعـه الى اليمـين من الجزء الاول وهو الاحد.

الرقم الاول في المثال هو (6)، ولهذا نحتاج ان نضع الرقم (6) على يمين الرقم (0)، وبعدـها نقرأ الارقام في المثال وللهذا الرقم الثاني هو (10) وللهذا نحتاج لوضع (0) الى يمين الرقم (1) ونسـتمر على هذه الطريقة نضع (0) الى يمين الرقم (5) وهـذا نستـمر.

وللهـذا سوف يكون شـكل الغـصن والـورقة لـبيانـات المـثال (1) كما موجود في الشـكل التـالـي. وكـما مـوضح في الشـكل الـارـقام التي تـتـبعـها الـارـقام الـاخـرى عـلـى الـيمـين تـدعـى الـارـقام الـبـداـية الغـصن (stem) والـارـقام النـهـائـية او الـمـكـملـة تـسـمى الـاوـاقـ (Leaves) والـشـكل هو (stem and Leaf) او الغـصن والـورـقة.

الشكل رقم (1)		N = 16
Stem-and-leaf of C1		
Leaf Unit = 1.0		
1	0 6	
2	1 0	
		2 2
4	3 29	
6	4 04	
7	5 0	
(2)	6 09	
7	7 08	
5	8 9	
4	9 08	
2	10 0	
HI	200;	

ويمكن ملاحظة اكبر رقم في الشكل رقم(1) الرقم (200) منفصل هذه القيمة يحددها الغصن والورقة على انها قيمة شاذة وبمزيد من التفاصيل أرجع ديلر 1996.

2.3.2 الرسم الصندوقى (Box plot)

الرسم الصندوقى box plot يعتبر من افضل الرسوم الاحصائية لعرض البيانات الاحصائية ولاجراء المقارنات بين عدة مجاميع من البيانات والذي اكتشف من قبل العالم (Tukey 1977). ان الفكرة الاساسية للرسم الصندوقى بسيطة وهي عرض بالرسم لخمس قيم تؤخذ من البيانات تسمى المخلصات الخمسة (5-Number Summaries) انظر الشكل رقم (2-6) وهي كما يلى (Lower Extreme) وهي اصغر قيمة ترتبط بل الصندوق غير شاذة، قيمة الربيع الاول او الربيع الادنى (Q_1) = Lower Quartile، قيمة الوسيط (Median) (Q_2)، قيمة الربيع الثالث او الربيع الاعلى (Q_3)، وكذلك اكبر قيمة غير شاذة مرتبطة مع الصندوق (Upper Extreme)، القيم التي تكون اكبر من اكبر قيمة تعتبر قيم شاذة وكذلك القيم التي هي اصغر من اصغر قيمة تعتبر قيم شاذة ولهذا كل القيم الشاذة تؤشر بشكل منفصل اذا كانت كبيرة او صغيرة.

اما كيفية بناء الرسم الصندوقى او ابعاد الصندوق فهي كما يلى، بعد ترتيب البيانات تصاعدياً.

1. طول الصندوق او المستطيل فهو الفرق بين الربيع الثالث Q_3 والربيع الاول Q_1 أي $(Q_3 - Q_1)$ ويمثل المدى الرباعي (IQR= Interquartile Range) والذي يمثل 50% من البيانات.

2. الخط داخل المستطيل او الصندوق يمثل الوسيط (Median) الى البيانات.

3. الخط الذي يوصل الصندوق با (Upper Extreme ، Lower Extreme) (Whisker Lengths) وطوله بالنسبة الى الادنى يمثل ($Q_1 - 1.5 * IQR$) وطوله بالنسبة الى الاعلى فهو ($Q_3 + 1.5 * IQR$) القيم التي تقع خارج هذين الامتدادين من الادنى والاعلى تؤشر بشكل منفصل وتسمى قيم شاذة (Outlier values).

4. عرض الصندوق (box width) فقد استخدم Tukey قاعدة تناسب مع الجذر التربيعي لحجم العينة (\sqrt{n}) .

5. حدود الشقة حول الوسيط (notch) ويعادل 5% تحسب وفقاً للقانون التالي.

$$M_e \pm \frac{1.58(IQR)}{\sqrt{n}}$$

حيث ان:

n : يمثل حجم العينة.

M_e : الوسيط للبيانات.

ويستخدم الرسم الصندوقى لاجراء المقارنات بين عدة مجاميع في ان واحد وكذلك لتحديد القيم الشاذة في البيانات فكل قيم تقع خارج (Upper extreme and lower extreme) تعتبر قيم شاذة،وكما هو واضح من الرسم الصندوقى (Box plot,Example no.2) في نهاية البحث شكل رقم (1AA) وللاطلاع على التفاصيل (ارجع الى ديلر 1996).

4-3-2 الرسم الصندوقى المزدوج (Range Finder Box plot)

يعتبر الرسم الصندوقى المزدوج من احدث الطرق الاحصائية لتحليل العلاقة بين توزيعي المتغيرين لشكل الانتشار(Scatter plot) وتحديد الاجزاء المهمة لها وكذلك القيم الشاذة لكل منها في نفس الوقت. وهو مفيد جداً حيث يربط بين شكل الانتشار الاعتيادي والرسم الصندوقى الاعتيادي ويجمع المعلومات الخاصة بالشكلين، والرسم الصندوقى المزدوج (Rangefinder Box plot) يمثل رسمين صندوقيين متلقعين في نقطة الوسيط ولكن رسم صندوقى هناك ثالث خطوط تمثل الرسم الصندوقى ثلاثة للمتغير المعتمد (Y) وثلاثة للمتغير المستقل (X).

ويمكن توضيح فكرة الرسم هو هناك خطين وسطيين (عمودي للكل متغير) طول كل منها يمثل طول الصندوق في الرسم الصندوقى الاعتيادي ويتقاطعان في موقع الوسيط وهناك فراغ يمثل طول الامتداد الى اصغر قيمة وابكر قيمة (whisker lengths) وهناك خطان افقيان لكل صندوق موقعهما يمثل اكبر قيمة واصغر قيمة (Extreme values) لكل صندوق وطولهما يمثل عرض الصندوق لكل متغير وهذا الرسمان الصندوقيان العمودي للمتغير (y) والافقى للمتغير (x). أي قيمة تقع خارج (Extreme values) لكل متغير او رسم صندوقى تؤشر على انها قيمة شاذة وتكون بشكل منفرد، وسوف يكون هناك رسم صندوقى مزدوج في الجانب العملى في الفصل التالي وللتتفاصيل رجع الى (ديлер 1996).

3. الجانب التطبيقي

3.1 المقدمة:

سوف نتطرق في هذا الجانب الى استخدام الطرق السابقة لتحديد القيم الشاذة (Outlier values) في البيانات واجراء المقارنة بين هذه الطرق جميعها طرق الرسم والطرق المعلمية (parametric methods) المستخدمة لتحديد القيم الشاذة في البيانات، الجانب العلمي يكون باستخدام الامثلة المستخدمة في الدراسات السابقة، ثم تطبيق طرق الرسم الحديثة الاستكشافية (Exploratory Data Analysis) على تلك الامثلة لتحديد القيم الشاذة فيها وتحديد الفرق بينهما ومقارنته النتائج.

3.2 القيم الشاذة في حالة المتغير الواحد

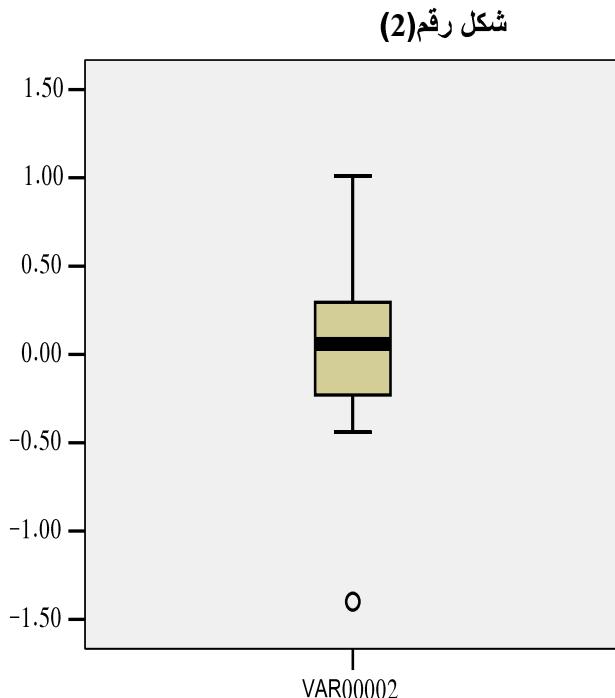
هناك ست امثلة في الاطروحة تمثل البيانات في حالة المتغير الواحد سوف نتناول ثلاثة منهم فقط لتحديد المجال الى البحث عند النشر، ومن ثم تحديد القيم الشاذة لكل مثال اولاً باستخدام الطرق المعلمية، وثانياً باستخدام الطرق الاستكشافية.

مثال رقم (2) 0.39 -0.05 -0.13 -0.24 -0.18 -0.22 0.44 0.63 0.48 0.3 -0.22 0.18 0.63 0.48 -0.22 0.18 0.44 -0.24 -0.13 0.05 -0.3 1.01 0.06 -1.4 0.20 0.10

a- تحديد القيم الشاذة باستخدام الطرق المعلمية.

استخدام AL-Jobouri¹⁶ 1976 طريقتين لاختبار وتحديد القيم الشاذة في المثال رقم (2) الطريقة الاولى طريقة Grubbs وقد وجد من نتيجة الاختبار بأن القيمة (-1.40) شاذة، اما الطريقة الثانية فهي طريقة Tietjen and moore وباستخدام الاحصاءة (E_k) وقد وجد بأن القيم (-1.40) و (1.01) هي قيم شاذة. كما استخدم الباحث طريقة Rosner وهي طريقة معلمية أخرى لتحديد القيم الشاذة لنفس المثال رقم (2) وقد طبق على القيمتين (-1.40) و (1.01) وقد وجد بأن ($R_1=2.573$) وان ($R_2=2.218$) والقيمة الحرجية للاحصاءتين R_2, R_1 على التوالي هي 2.72*, 2.68* وبالمقارنة مع القيم المحسوبة لاعتبر ايًّا من القيمتين معنوية وهذا يعني لا توجد قيم شاذة في هذا المثال حسب اختبار Rosner.

b- تحديد القيم الشاذة باستخدام الطرق الاستكشافية.
 اولاً : بطريقة الرسم الصندوقي (Box plot)
 من خلال النظر الى الشكل رقم (2) التالي وهو الرسم الصندوقي لهذا المثال نجد ان هناك قيمة
 شاذة واحدة وهي (-1.40). شكل رقم (2) حيث تم رسمها بشكل منفصل خارج الرسم الصندوقي.



ثانياً: باستخدام طريقة الغصن والورقة وهي موضحة في الشكل رقم (3)
 شكل رقم(3)

Stem-and-leaf of C1

N

= 15

Leaf Unit = 0.10

LO -14;

2 -0 4

5 -0 322

7 -0 10

(3) 0 011

5 0 23

3 0 4

2 0 6

1 0

1 1 0

ايضاً باستخدام هذه الطريقة نلاحظ ان هناك قيمة شاذة واحدة فقط وهي (1.40) حيث تم فصلها عن باقي البيانات وقد اشرت تحت حقل (LO-1.40) وهذا يعني وجود تطابق بين الطرق الاستكشافية في تحديد القيم الشاذة وتخالف عن الطرق المعلمية كما هناك اختلاف ايضاً بين الطرق المعلمية في تحديد القيم الشاذة فحسب طريقة Grubbs هناك قيمة شاذة واحدة وهي (-1.40) وقيمتان شاذة وهما (1.40) و(1.01) حسب طريقة Tietjen and Moore ولا توجد قيم شاذة حسب طريقة Rosner وهذا يعني ان الاختلافات واضحة بين الطرق المعلمية بعكس الطرق الاستكشافية بالرسم فهي واضحة جداً ومتطابقة في هذا المثال.

بيانات المثال رقم (3)

-1.056	-1.008	-0.34	0.533	0.109	0.661	1.638	-0.413	-0.667	-0.57
1.207	-0.550	2.290	0.504	-2.215	2.139	-0.048	-0.909	0.967	-0.143

a- تحديد القيم الشاذة باستخدام الطرق المعلمية.

استخدام AL-Jobouri⁽¹⁶⁾ عام 1976 طريقة Rosner لتحديد القيم الشاذة في المثال رقم 3 اعلاه، وقد وجد بأنه لا توجد اية قيمة من القيم شاذة باستخدام هذه الطريقة. في حين استخدام الباحث اضافة الى ذلك ثلاثة طرق معلمية اخرى لتحديد القيم الشاذة في المثال رقم 3 وعلى النحو التالي:-

1. الطريقة الاولى:

باستخدام الاحصاءة التي افترضها Mcmillan لاختبار اعلى قيمتين في البيانات وهي:

$$X_n + X_{n-1} - 2\bar{X} > C_{\alpha}^n S$$

$$(2.29 + 2.139) - 2(0.1064) > (1.145)(0.637)^*$$

$$4.2101 > 0.729$$

وبذلك تعتبر القيمتان 2.139 ، 2.2906 شاذتان حسب اختبار Mcmillan .

2-الطريقة الثانية :- باستخدام اختبار Grubbs لاختبار اعلى قيمتين 2.290 ، 2.139

$$\frac{S_{n,n-1}^2}{S^2} = \frac{15.0129}{24.899} = 0.6029$$

نستخدم الاحصاءة

ولاختبار ادنى قيمتين -1.056 ، -2.215

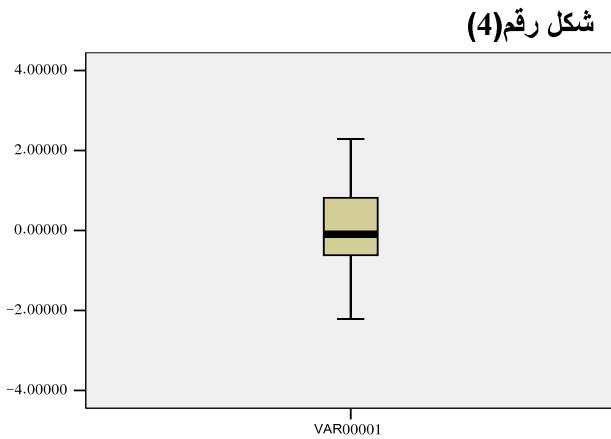
$$\frac{S_{1,2}^2}{S^2} = \frac{21.896}{24.896} = 0.879$$

وبالمقارنة مع القيمة الحرجة للاحصاءتين اعلاه عند مستوى $\alpha = 0.05$ والتي تساوي 0.5269 ، نجد بأن القيم اعلاه ليست شاذة. وهذا يعني عدم وجود قيم شاذة في المثال رقم (3) باستخدام طريقة Grubbs المعلمية.

3-الطريقة الثالثة باستخدام اختبار Tietjen and Moore

وقد وجد بأن الاحصاءة ($E_2 = 0.537$) وبالمقارنة مع القيمة الحرجية والتي تساوي $E^*_{2,16} = 0.416$ ^{١٦} ، أي ان القيمة غيرمعنوية وهذا يعني لاتوجد قيم شاذة للمثال رقم(3) اعلاه باستخدام هذه الطريقة.

b- تحديد القيم الشاذة للمثال رقم (3) باستخدام الطرق الاستكشافية.
اولاً: بطريقة الرسم الصندوقى الشكل رقم (4) لا توجد قيم شاذة لهذه البيانات للمثال رقم (3).



ثانياً: بطريقة الغصن والورقة (Stem and Leaf) وهي موضحة في الشكل رقم (5) للغصن والورقة ونلاحظ من خلال الشكل لاتوجد قيم شاذة أيضاً.

وهنا توجد فروقات بين الطرق المعلمية وطرق الرسم في تحديد القيم الشاذة، حيث نرى بوضوح عدم وجود قيم شاذة للمثال رقم (3) اعلاه باستخدام ثلاثة طرق معلمية مختلفة في حين هناك قيمتان شاذتان باستخدام اختبار Mcmillan اما بالنسبة الى طرق الرسم فنلاحظ هناك تطابق بين طريقة (stem and leaf) و طريقة الرسم الصندوقى (Box plot) حيث لاتوجد قيم شاذة.

شكل رقم (5)

	Stem-and-leaf	of	C1
N = 20	Leaf	Unit	=
0.10			
2	1	-2	
	1	-1	
	3	-1	
00	7	-0	
9655	(4)	-0	
4310	9	0	
1	8	0	
5569	4	1	
2	3	1	
6	2	2	
12			

بيانات المثال رقم (4)

4.57 5.62 4.12 5.29 4.64 4.31 4.30 4.39 4.45 5.67 4.39 4.52 4.26

4.26 4.40 5.78 4.73 4.56 5.08 4.41 4.12 5.51 4.82 4.63 4.29 4.60

a- تحديد القيم الشاذة باستخدام الطرق المعلمية.

تم استخدام ثلاثة طرق معلمية لتحديد القيم الشاذة في المثال اعلاه وعلى النحو التالي:-

اوألاً: باستخدام طريقة (Grubbs)
لاختبار القيمة العليا : تأخذ أعلى قيمة 5.78

$$\frac{S_n^2}{S^2} = \frac{4.91}{6.165} = 0.79$$

وباستخدام الاحصاءة

ولاختبار القيم الصغرى: تأخذ اصغر قيمة 4.12 وحساب

$$\frac{S_1^2}{S^2} = \frac{5.837}{6.165} = 0.94$$

وبالمقارنة مع القيمة الحرجة 0.71 نرى بأن القيمتين 5.78، 4.12 ليست شاذتين وهذا يعني عدم وجود قيم شاذة في مثال رقم 4 اعلاه باستخدام اختبار Grubbs.

ثانياً: باستخدام اختبار (Tietjen and Morre) وحساب

$$R_1 = |5.78 - 4.681| / 0.49$$

$$R_1 = 2.24$$

وبالمقارنة مع الحد الاعلى للاحصاء وهي $R_1^* = 2.93$ وهذا معناه عدم وجود قيم شاذة في هذا المثال حسب الاختبار (Rosner).

b- تحديد القيم الشاذة باستخدام الطرق الاستكشافية

أولاً: باستخدام طريقة الرسم الصندوفى.

ابعاد الرسم الصندوفى هي:-

الوسيط (Median) 4.54 –

الربع الاول (Q_1) – 4.31

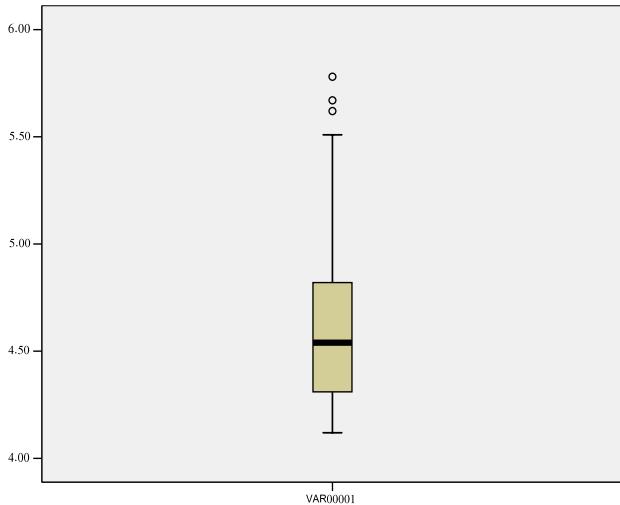
الربع الثالث (Q_3) – 4.82

ادنى قيمة غير شاذة (L) – 4.12

اعلى قيمة غير شاذة (U) – 5.51

والشكل رقم (6) يمثل الرسم الصندوفي لبيانات المثال رقم (4).

شكل رقم (6)



وكما مبين من الشكل اعلاه وجود ثلاث قيم شاذة يمكن ملاحظتها بوضوح تم رسمها بشكل منفصل خارج الرسم الصندوفي وهي (5.78, 5.67, 5.62).

ثانياً: باستخدام طريقة الغصن والورقة وهي موضحة بالشكل رقم (7)
 شكل رقم (7)

Stem-and-leaf of C1

N =

26

Leaf Unit = 0.10

2	4	11
9	4	2223333
(6)	4	444555
11	4	6667
7	4	8
6	5	0
5	5	2
4	5	5

HI 56;

56; 57;

ايضاً باستخدام هذه الطريقة هناك ثلات قيم شاذة حيث تم فصلها عن باقي البيانات وقد حدة تحت حق (HI)، وهي نفس القيم ظهرت باستخدام الرسم الصندوقى ، وهذا يعني هناك تطابق بين الطرق الاستكشافية لتحديد القيم الشاذة وهناك تطابق أيضاً بين الطرق المعلمية لتحديد القيم الشاذة حيث لم تؤشر ايً منها هناك قيم شاذة وهذا يوضح لنا هناك فروقات كبيرة بين الطرق المعلمية والطرق الاستكشافية لتحديد القيم الشاذة.

3.3 : القيم الشاذة في حالة المتغيرين ونماذج الانحدار
 سنتناول في هذا المبحث مثالين من مجموعة اربع امثلة تمثل البيانات في حالة المتغيرين، ثم تحديد القيم الشاذة لكل مثال باستخدام الطرق المعلمية اولاً والطرق الاستكشافية ثانياً.

بيانات المثال رقم (5)

X: 3.25 3.79 3.4 3.68 3.76 3.3 2.4 2.94 2.68 3.27 2.46 3.46 2.14 3.22
 2.91 1.98

Y: 2.85 3.91 2.33 3.29 3.71 3.52 4.07 2.98 3.18 2.72 2.41 3.01 2.56 2.89
 2.64 2.22

a- تحدد القيم الشاذة لبيانات المثال رقم(5) باستخدام الطرق المعلمية.

استخدم AL-Jobouri ¹⁶ عام 1976 اختبار البياتي (Hotelling T² test) لتحديد القيم الشاذة في بيانات المثال (رقم 5) فقد وجد من نتيجة الاختبار وجود زوج من قيم المشاهدات شاذة وهي (2.4,4.07) عن بقية قيم المشاهدات كما استخدم الباحث طريقة معلمية اخرى لاختبار وجود القيم الشاذة باستخدام الاحصاءة $t=Max |ei/s_i|$ وعلى النحو التالي:-

لحساب قيم e $t=Max |ei/s_i|$ على النحو التالي:-

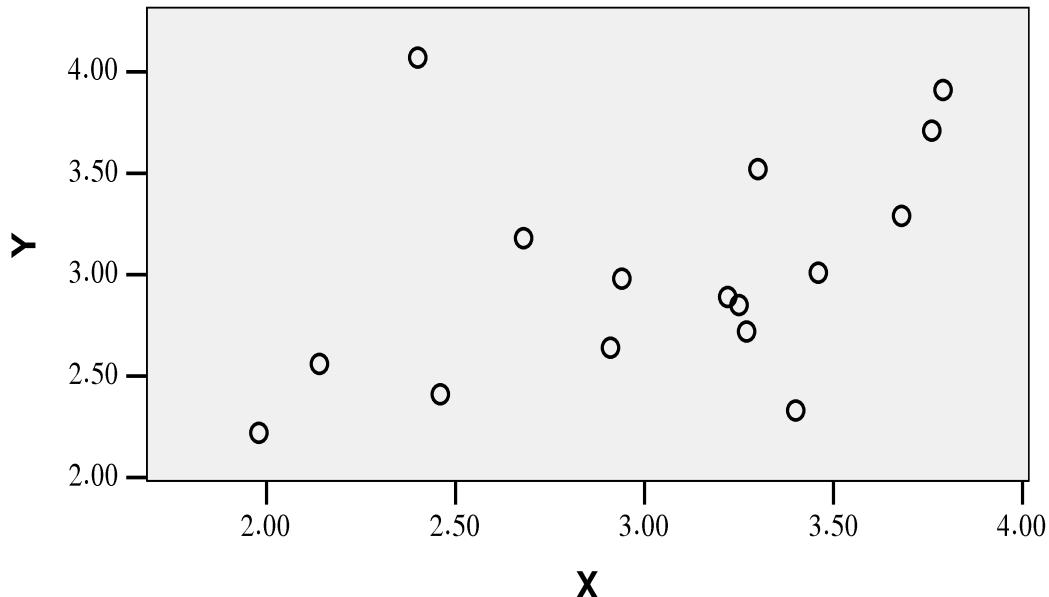
$$\hat{y}_i = 1.9 + 0.36 X_i$$

وبعد حساب قيم $|e_i / s_i|$ وجد الباحث بأن أكبر قيمة هي 2.61 وهي القيمة التي تقابل زوج المشاهدات ($x = 2.4, y = 4.07$) وبالمقارنة مع القيمة الحرجة للاحصاء اعلاه والتي تساوي $t^* = 2.6^{16}$ وبناءً على ذلك يعتبر زوج القيم (2.4, 4.07) شاذًا عن باقي القيم حسب هذا الاختبار وهي نفس النتيجة التي حصل عليها الجبوري في اختباره.

b- تحديد القيم الشاذة باستخدام طريقة الرسم الصندوقي المزدوج للمتغيرين (x,y)

(بيانات المثال رقم (5) اعلاه وهي موضحة بالشكل رقم (8) التالي).

شكل رقم (8)



ونلاحظ من الشكل اعلاه الرسم الصندوقي المزدوج للمتغيرين (x, y) عدم وجود ايَّة قيمة شاذة للمتغيرين.

وهذه النتيجة تختلف عن نتيجة الطرق المعلمية حيث لاحظنا ان الطريقتين المعلمتين اعطت نفس النتائج وهي وجود قيمتين شاذة ، وهذا مؤشر قوي على انه هناك فروقات معنوية بين الطرق المعلمية والطرق الاستكشافية في تحديد القيم الشاذة.

بيانات المثال رقم (6)

X : 51.3 49.9 50 44.2 48.5 47.8 47.3 45.1 46.3 42.1 44.2 43.5 42.3 40.2
31.8 34.0

Y : 102.5 104.5 100.4 95.9 87.0 95 88.6 29.2 78.9 84.6 81.7 72.2 65.1 68.1
67.3 52.5

a- تحديد القيم الشاذة لبيانات مثال رقم (6) باستخدام الطرق المعلمية:-

استخدم الباحث طريقتين لاختبار وجود القيم الشاذة، في هذا المثال وهي كما يلي: الطريقة الاولى باستخدام اختبار البياني وخطواته على النحو التالي:

$$\text{اولاً: ايجاد قيمة } \hat{E} \text{ وعند اختبار القيمة } (X_0 = 31.8, Y_0 = 67.3)$$

$$\text{فإن } \hat{E} = 8.86$$

$$T^2 = \frac{2(n-2)}{n-3} F_\alpha \quad \text{وباستخدام العلاقة}$$

$$T^2 = 8.05$$

أي ان $\hat{E} > T^2$

وهذا يعني ان زوج القيم اعلاه $(x_0 = 31.8, Y_0 = 67.3)$ يعتبر شاذًا باستخدام طريقة البياني.

والطريقة الثانية باستخدام الاحصاءة $t = \text{Max}|e_i|/s_i|$ على النحو التالي:

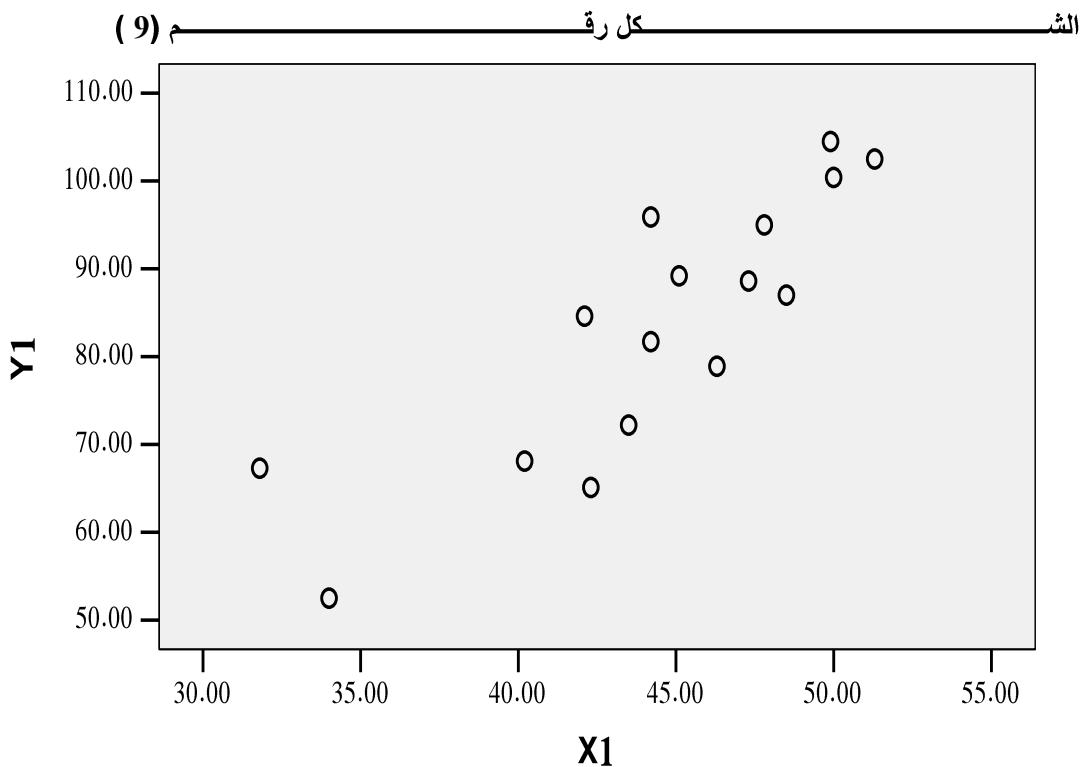
ان تقديرات معالم معادلة خط الانحدار بطريقة المربعات الصغرى هي كما يلي:

$$\hat{Y} = -21.8 + 2.4 X_i$$

ولحساب المعلومات المطلوبة لحساب الاحصاءة t اعلاه.

وجد في جدول المعلومات في العمود $|e_i|/s_i|$ بأن أعلى قيمة له هي 2.00 وهي القيمة التي تقابل كل من $x = 42.3$, $y = 65.1$ وبالمقارنة مع القيمة الحرجية $(t^* = 2.9)^{16}$ وببناءً على ذلك لا توجد قيم شاذة في هذا المثال حسب هذا الاختبار ولمزيد المعلومات يمكن الرجوع الى المصدر (ديلر 1996).

b- تحديد القيم الشاذة باستخدام طريقة الرسم الصندوقى المزدوج للمثال رقم (6) اعلاه وهي موضحة بالشكل رقم (9) التالي:



ونلاحظ من الشكل اعلاه بوضوح وجود قيمة شاذة واحدة من ضمن القيم التي يمكن رؤيتها اذا رسمت بشكل منفصل خارج الرسم الصندوقى المزدوج وهي ($x = 31.8$, $y = 67.3$) وهذه النتيجة مطابقة الى النتيجة التي تم التوصل اليها باستخدام احدى الطريقتين المعلمتين ومخالفة لنتيجة الطريقة المعلميمية الاخرى، ومرة اخرى ظهرت فروقات في النتائج بين الطرق المعلميمية والطرق الاستكشافية.

4- الاستنتاجات والتوصيات

4.1- المقدمة:

سنتطرق في هذا المبحث اولاً الى اهم النتائج التي تم التوصل اليها في الجانب العملي:-
في المثال الاول والذي يمثل البيانات الخاصة بمتغير واحد كانت الطرق الاستكشافية متطابقة فيما بينها في تحديد القيم الشاذة حيث اظهرت كل من طريقتين الرسم الصندوقى والغصن والورقة بوجود قيمة شاذة واحدة، في حين ان الطرق المعلمية اختلفت فيما بينهما في تحديد القيم الشاذة، فحسب طريقة Grubbs هناك قيمة شاذة واحدة وقيمتان شاذتان حسب طريقة Tietjen and Moore.

في المثال الثاني في الاطروحة والذي لم تطرق اليه في الجانب العملي تطابقت الطرق الاستكشافية ايضاً في تحديد القيم الشاذة، حيث اظهرت كل من طريقتين الرسم الصندوقى والغصن والورقة عدم وجود قيم شاذة اما الطرق المعلمية فكانت الاختلافات بينها واضحة، وكانت النتائج وجود قيمتين شاذتين حسب طريقة Grubbs وعدم وجود اي قيمة شاذة باستخدام طريقة Tietjen and Moore.

في المثال الثالث في الاطروحة الذي تم التطرق اليه في الجانب العملي مثال رقم 2، اختلفت الطرق الاستكشافية في تحديد القيم الشاذة حيث اظهرت طريقة الرسم الصندوقى ببعد وجود قيم شاذة، في حين عند استخدام طريقة الغصن والورقة كانت هناك ثلاثة قيم شاذة، وكذلك الطرق المعلمية اختلفت ايضاً في تحديد القيم الشاذة وكانت نتيجة ثلاثة طرق منها وهي طريقة Ronser وطريقة Grubbs وطريقة Tietjen and Moore متباينة وهي عدم وجود قيم شاذة ووجود قيمتين شاذتين باستخدام طريقة McMillan. اما في المثال الرابع في الاطروحة والذي لم تطرق اليه في الجانب العملي وموضح بشكل مفصل في المصدر (ديلر 1996) فقد استخدم طريقة معلمية واحدة لتحديد القيم الشاذة نظراً لطبيعة البيانات والتي كانت تتكون من ثلاثة متغيرات ومجموعات وكانت النتائج وجود قيمة شاذة واحدة في كل متغير او مجموعة بينما في الطرق الاستكشافية كانت نتائجها مختلفة بالنسبة للمجموعة الاولى باستخدام الرسم الصندوقى لم تكن هناك قيمة شاذة وجود قيمة شاذة واحدة باستخدام الغصن والورقة وللمجموعة الثانية وجود قيمة شاذة واحدة باستخدام طريقة الغصن والورقة واخيراً للمجموعة الثالثة وجود قيمة شاذة واحدة باستخدام الرسم الصندوقى وقيمتين شاذتين باستخدام طريقة الغصن والورقة.

وهذا ايضاً مؤثر جديد لهذه الحالة من البيانات، حيث كانت هناك اختلافات كبيرة في النتائج باستخدام الطرق المعلمية والطرق الاستكشافية.

اما في المثال الخامس في الاطروحة والذي تم التطرق الى نتائجه في المثال رقم (3) في الجانب العملي في هذا المبحث وكانت نتائج الطرق الاستكشافية متباينة نوعاً ما فيما بينها في تحديد القيم الشاذة حيث اظهرت ثلاثة قيم شاذة باستخدام طريقة الرسم الصندوقى واربع قيم شاذة باستخدام طريقة الغصن والورقة، اما الطرق المعلمية فكانت نتائجها في هذا المثال متطابقة حيث لم تظهر اي من الطرق الثلاث المستخدمة طريقة Rosner وطريقة Tietjen and Moore.

وطريقة Grubbs وجود اي قيمة شاذة.
في المثال السادس في الاطروحة والذي لم تطرق اليه في هذا البحث في الجانب العملي والذي يمثل بيانات متغير واحد فقد كانت النتائج متباينة نوعاً ما بين الطرق المعلمية والطرق الاستكشافية حيث اظهرت جميع الطرق نفس النتائج وهي عدم وجود قيمة شاذة فقط الاختلاف كان في طريقة الغصن والورقة حيث هناك قيمة شاذة واحدة فقط.

في المثال السابع في الاطروحة والذي يمثل الرابع في هذا المبحث في الجانب العملي والذي يمثل البيانات الخاصة بالمتغيرين (y, x) ثم استخدام طريقتين معلمتين طريقة اختبار T^2 وطريقة الاحصاء t لتحديد القيم الشاذة وكانت النتيجة متطابقة في الطريقتين وهي وجود زوج من

المشاهدات شاداً اما باستخدام الرسم الصندوقى المزدوج للمتغيرين (x, y) فلم يكن هناك أي زوج من القيم الشاذة.

اما المثال الثامن في الاطروحة والذي لم ننطرق اليه في البحث في الجانب العملي وهو ايضاً يمثل حالة العلاقة بين المتغيرين (x, y) كانت النتيجة وجود زوج من قيم المشاهدات شاداً باستخدام أحد الطرقين المعلميتين المستخدمتين وعدم وجود قيمة شادة باستخدام الطريقة الاخرى وعدم وجود أي زوج من قيم المشاهدات شاداً باستخدام الرسم الصندوقى المزدوج للمتغيرين (x, y). في المثال التاسع في الاطروحة والذي لم ننطرق اليه في هذا البحث في الجانب العملي والذي يمثل حالة العلاقة بين المتغيرات (x_1, x_2, y) كانت النتيجة متطابقة بين الطريقة المعلمية المستخدمة وطريقة الرسم الصندوقى المزدوج في تحديد القيم الشاذة وهي وجود قيمة شادة بالنسبة لكل من المتغيرات (x_1, x_2, y).

اما في المثال العاشر والذي تم التطرق اليه في هذا البحث والذي يمثل العلاقة بين (x, y) فقد اختلفت الطرقان المعلميتان المستخدمتان في تحديد القيم الشاذة فقد اظهرت قيمة من المشاهدات شادة باستخدام احدى الطرقين ولم يظهر اي من القيم شاداً باستخدام الطريقة المعلمية الاخرى، اما باستخدام الرسم الصندوقى المزدوج فكانت هناك قيمة من قيم المشاهدات شادة وهذا يعني وجود اختلافات بين الطرق المعلمية والاستكشافية وكذلك بين الطرق المعلمية نفسها.

4.2 : الاستنتاجات

1. اظهرت الدراسة تفاوت الطرق المعلمية في تحديد القيم الشاذة في حالة المتغير الواحد وهي وطريقة Grubbs وطريقة Tietjen وطريقة Rosner فيما بينهما بشكل واضح جداً حيث كانت طريقة Grubbs اكثر حساسية لتحديد القيم الشاذة من غيرها من الطرق المعلمية وبعكس طريقة Rosner التي كانت اقل الطرق المعلمية حساسية لتحديد القيم الشاذة.
2. بينت الدراسة وجود تطابق ملحوظ في النتائج بين الطرقين المعلميتين المستخدمتين طريقة اختبار T2 وطريقة اختبار الاحصاء χ^2 وبين طريقة الرسم الصندوقى المزدوج Rangefinder Box plot في تحديد القيم الشاذة في حالة المتغيرين في اغلب الامثلة المستخدمة في البحث.
3. اظهرت الدراسة اختلاف بين الرسم الصندوقى والغضن والورقة في تحديد القيم الشاذة وان الغصن والورقة اكثر حساسية من الرسم الصندوقى.
4. بينت الدراسة هناك اختلافات واضحة بين الطرق المعلمية والطرق الاستكشافية في تحديد القيم الشاذة.
5. اظهرت الدراسة وجود خلاف كبير في النتائج بين الطرق المعلمية وطريقة الغصن والورقة (stem and Leaf) في اغلب الامثلة المستخدمة في البحث.
6. وجد من خلال النتائج للدراسة هناك تشابه في النتائج بين الطرق المعلمية (Tietjen and Moore).
7. الطرق الاستكشافية تحدد القيم الشاذة مباشرة بينما الطرق المعلمية تتطلب تحديد قيم مشكوك فيها شاذة كاجراء اولى، ومن ثم الاعتماد على الاختبارات الاحصائية الخاصة بتحديد فيما اذا كانت تلك القيم شاذة ام لا.
8. توصل الباحث الى ان الرسم الصندوقى (Box plot) هو افضل الطرق الاستكشافية وان طريقة (Tietjen and Moore) افضل الطرق المعلمية وهناك تشابه بين نتائج الطرقين كما مبين في الفقرة 6.
9. توصل الباحث بأن طريقة الرسم الصندوقى (Box plot) هي افضل طريقة لتحديد القيم الشاذة بالنسبة للطرق الاستكشافية والمعلمية.

10. اثبتت الدراسة كفاءة ونجاح طريقة الرسم الصندوقى المزدوج (Rangefinder box plot) في تحديد القيم الشاذة في حالة المتغيرين، حيث انها تقوم بتحديد القيم الشاذة في كل من المتغيرين (x, y) على انفراد والقيم الشاذة في المتغيرين (x, y) معاً وتجمع المعلومات الخاصة بالمتغيرين (x, y) في رسم يمثل رسمنين صندوقيين معاً.

4.3 : التوصيات

على ضوء الاستنتاجات السابقة يوصي الباحث بما يلي:

1. ضرورة اجراء عرض للبيانات قبل القيام بأية عملية تحليل احصائى حيث ان مثل هذا الاجراء يؤدي الى الحصول على عوامل ومتغيرات جديدة ستسهم في تطوير الظاهرة تحت الدراسة، كذلك الحصول على مؤشر جديد لمدى دقة العمل، وكلا الامرین سيسهمان في دقة النتائج.
2. استخدام الطرق الاستكشافية Box plot, stem and leaf, and Rangefinder Box plot لتحديد القيم الشاذة في البيانات في حالة المتغير الواحد والمتغيرين وتفضيلهما على الطرق المعلمية الأخرى.
3. اذا كان لابد من استخدام الطرق المعلمية لتحديد القيم الشاذة يوصي الباحث باستخدام طريقة Tetjen and Moore في حالة المتغير الواحد وطريقة اختبار الاحصاء ($t = \frac{Max|e_i|}{|s_i|}$) في حالة المتغيرين (x, y) وتفضيلهما على الطرق المعلمية الأخرى.
4. اجراء دراسات على بيانات مختلفة ومتنوعة يستخدمها الباحث باستخدام الطرق المعلمية والاستكشافية لتحديد بالضبط أي من الطرق سوف يكون اكثر كفاءة.

المصادر

المصادر العربية

1. الجبوري، شلال حبيب (1988) "اسلوب جديد لاكتشاف وتقدير المشاهدات الشاذة في حالة متعدد المتغيرات"، بحث القى في المؤتمر الثاني للجمعية العراقية للعلوم.
2. الجبوري، منى حسين(1990) "الاكتشاف الجزئي للمشاهدات الشاذة وطرق التقدير في حالة متعدد المتغيرات"، رسالة ماجستير في الاحصاء، الجامعة المستنصرية.
3. المتنو، نعم مسلم(1993) "تقييم البيانات المفقودة والشاذة في تحليل تصميم التجارب غير المتزنة، رسالة ماجستير في الاحصاء، جامعة بغداد.
4. المختار، سليمان محمد امين (1980) "القيم الشاذة واثرها في تحليل البيانات الاحصائية، رسالة ماجستير في الاحصاء، جامعة بغداد.

المصادر الاجنبية

5. Barnett, V. & Lewis, T.(1978) "Outlier in Statistical Data", John Wiley and Sons, New York.
6. Al-Bayati, H.A.(1973)"Procedure for Detecting Observation in Samples two related Variables", The Proceeding of the Ninth Conference of Statistics and Computation to Multivariate Statistical Analysis, John Wiley and Sons,New York.
7. Becketti, S. & Gould, W.(1987) "Rangefinder Box plot Anote", The American Statistician, V.41,No.2,p.(149).
8. Behken, D. W. & Draper, N. R.(1972) "Residuals and their Variance Patterns", Technometrics, V.17,pp(127-128).
9. Bross, L. D. J.(1961) "Outliers in Patternend Exprements a Strategic reappraisal", Technometrics, V.3,pp(91-102).
10. Cleveland, W.S.(1985) "The Elements of Graphing Data", Monterey, CA; Wads Worth.
11. Daniel, C. (1960) "Location Outliers in Factorial Experiments", Technometrics, V.2,pp(140-150).
12. Gnanadesikan, R.& kettenring, J.R.(1972) "Robust Estimates, Residuals and Outlier Detection with Multiresponse Data", Biometrics,V.28,pp.(81-124).
13. Grubbs,F.E.(1950) "Sample Cirteria for Testing Outling Observation", Ann. Math. Stat., V.s1,pp(27-58).
14. Hussin, M. M. (1989) "Some studies of Graphical Methods in Statistical Data Analysis ; subjective Judgements in the Interpretation of Box plot", Unpublished Ph. D. Thesis, Keele University, Uk.
15. Irwin, J.O. (1925) "On a Criterion for the Rejection of Outlying Observation", Biometrika, V.17,pp.(238-250).
16. Al-Jobouri, S.(1976) "Test of Outliers", Unpublished M.Sc Thesis, University of Baghdad.

17. John & Prescott, P. (1975) "Critical Values of a Test to detect Outliers in Factorial Experiments", *Appl. Sta.*, V.24,pp.(56-59).
18. Kishpaugh, J.R.L.(1972) "Experiments in Outliers Detection in Multivariate Data", M. Sc. Thesis, State University of New York, Stony Brook, NY.
19. Lund,R.E.(1975) "Table for an Approximate Test for Outliers in Linear Models", *Technometrics*, V.17,pp.(473-476).
20. McMillan, R.G. (1971) "Tests for one or Two Outliers in Normal Samples with unknown variance", *Technometrics*, V.13,No.1,pp.(87-100).
21. Mosteller, F.(1948) "AK-Sample Slippage Test for an Extreme Population", *Ann. Math.Stat.*,Vol.19,pp.(58-65).
22. -----& Tukey,J.W.(1950) "Significance Levels for a k-Sample Slippage Test", *Ann.Math.Stat.*,Vol.21,pp.(120-123).
23. Quesenberry, C. P. & David(1961) "Some Tests for Outliers", *Biometrika*,V.48,pp.(379-387).
24. Rohlff, F.J.(1975) "Generalization) of Gap Test for the Detection of Multivariate Outliers", *Biometrics*,V.31,pp.(93-101).
25. Rosner,B.(1975)"On the Detection of many Outliers",*Technometrics*,Vol.17,No.2,pp.(120-135).
26. Strikantan, K.S.(1961) "Testing for Single Outliers in a Regression Model", *Sankhya*,A.Vol.23,pp.(251-260).
27. Tietjen, G.L., Moore, R.H.& Beckman, R.J.(1973) "Testing for a Single Outlier in a Simple Linear Regression", *Technometrics*,Vol.15,pp.(717-721).
28. Tukey, J.W.(1977) "Exploratory Data Analysis", Addison-Wesley Publishing Company.
29. Velleman,P.F.& Hoaglin, D.C.(1981) "Applications, Basics and Computing Exploratory Data Analysis", Boston, M. Sc., Duxbury Press.
30. Zinger, A.(1961) "Detection of Best and Outline Normal Distributions with known Variances", *Biometrika*,vol.48,p.(457).

شكل رقم (1AA)

