

التعويض الجزئي لقيم المفقودة

* م. م. هيثم عبد الأمير الوردي

المُسْنَدُ

إذا كانت مشكلة فقدان القيم، في مرحلة جمع البيانات، قد تمت بشكل عشوائي تام، فيمكننا معالجتها بالحذف بسهولة، أما إذا تم فقدان بشكل عشوائي، فنستطيع غالباً جعله عشوائي تام بالتعويض الجزئي، الذي يستند إلى إحلال قيم مقدرة محل مجموعة من القيم المفقودة وليس جميعها، وباحتى طرق التعويض الأحادي للسهولة، المتوفرة في معظم الحزم الإحصائية، ثم استخدام الحذف للمعالجة. إن بساطة هذه الخطوات تمكن الباحث غير المتخصص، وهو غالباً ما يكون كذلك، من تطبيقها بمفرده، دون طلب مساعدة، يصعب غالباً الحصول عليها.

Abstract

If the problem of missing values, during the data collection stage, happened completely at random (MCAR), we can simply use the deletion method to solve it. But if the missing happened at random (MAR), we can often make it as MCAR by partial imputation. That is based on replacing estimated values with some of the missing values instead of all of them, using one of the single imputation methods for simplicity, which are available in most of the statistical packages. Then we can use the deletion method as a solution. The simplicity of these steps enables unspecialized researcher to imply them alone, without asking help, which is usually hard to obtain.

1- المقدمة

غالباً ما تتضمن البيانات المجمعة على قيم مفقودة، تحتاج إلى معالجة دقيقة للحصول على أفضل النتائج. لقد تعددت طرق المعالجة وفقاً لآليات فقدان، وأبسطها هي حذف حالات القيم المفقودة، ولكن الحذف يتطلب فرضية فقدان العشوائي التام لها، فعندما لا تكون البيانات كذلك، يلجأ الباحث إلى طرق التعويض الكلي، ومعظمها تحتاج إلى خبرة لا يستطيع غير المتخصص تطبيقها، مثل طريقتي التعويض المتعدد والإمكان الأعظم، وهذا البحث هو من أجل تقديم طريقة جديدة، بسيطة وفعالة، تمكن الباحثين في مختلف التخصصات من تطبيقها، وتتلخص بتتويج جزئي عن القيم المفقودة بحيث نحصل على حالة فقدان عشوائي تام، وبالتالي المعالجة بالحذف حيث نحصل على تقديرات غير متحيزة للمعلمات. لقد قسمت البحث على أربعة أقسام، تناولت في مقدمته ما يتعلق بالمشكلة والهدف والفرضية وعينة البحث والدراسات السابقة، وشرحت في الإطار النظري منه آليات فقدان القيم وطرق المعالجة وصيغ الاختبارات، أما القسم الثالث فهو الجانب التطبيقي، إذ استخدمت فيه بيانات ملفين جاهزين ضمن حزمة SPSS v18، واستخدمت هذه الحزمة في استخراج النتائج، وأخيراً أدرجت الاستنتاجات وما أخطط له من عمل مستقبلي في القسم الرابع.

1-1 مشكلة البحث وأهميته

غالباً ما تحتوي البيانات المجمعة على قيم مفقودة، ومعالجة الشكل العشوائي التام (Missing Completely at Random) لها يمكن أن يتم ببساطة بإحدى طرقتي الحذف (Deletion)، إذ تتوافق طرفيته، القائمة (Pairwise) والمزدوج (Pairwise)، في معظم الحزم الإحصائية، وستكون المعلمات المقدرة غير متحيزة (Unbiased). [7] إن تطبيق الحذف مرتبط بفرضية MCAR للبيانات، ولكن عند عدم توافر هذا الفرض، وعندما تكون البيانات قد فقدت بشكل عشوائي (Missing at Random)، وليس عشوائياً تماماً (MCAR)، يلجأ الباحث إلى استخدام إحدى طرق التعويض (Imputation)، التي تستند إلى إحلال قيم تقديرية محل جميع القيم المفقودة، ومشكلة التقدير تتطلب استخدام الأفضل من بين طرق التعويض للحصول على تقديرات غير متحيزة (Unbiased) وفعالة (Efficient). وبالفعل نحصل باستخدام طريقة التعويض المتعدد (Multiple imputation)، أو الإمكان الأعظم (Maximum likelihood) [1] على مثل هذه التقديرات، ولكن تطبق أي منها يحتاج إلى خبرة عالية في هذا المجال، نادراً ما تتوافق لدى الباحث، فضلاً عن عدم توافر طريقة الإمكان الأعظم في معظم الحزم الإحصائية. من هنا تأتي أهمية هذا البحث في تقديم طريقة التعويض الجزئي (Partial imputation) للقيم المفقودة، التي تستند إلى خطوات بسيطة، تمكن الباحث غير المتخصص من تطبيقها.

2-1 هدف البحث

يهدف البحث إلى تقديم طريقة جديدة في معالجة القيم المفقودة، تستند إلى تحديد المتغير المسبب للمشكلة أولاً، ثم التعويض الجزئي بإحدى طرق التعويض الأحادي (Single imputation) البسيطة لقيمة المفقودة ثانياً، بحيث نحصل على حالة MCAR لفقدان القيمة، وأخيراً استخدام الحذف لمعالجة المشكلة.

3-1 فرضية البحث

يمكن الحصول على حالة MCAR لفقدان القيم، الضرورية لتطبيق الحذف، بالتعويض الجزئي لقيم المفقودة من المتغير المسبب للمشكلة.

4-1 عينة البحث

عند تثبيت حزمة SPSS (مثلاً إصدار 18 الذي استخدمه) على القرص الصلب للحاسوب، تنسخ معها تلقائياً مجموعة ملفات بيانات جاهزة، منها الملفان telco.sav و telco_missing.sav، اللذان سأستخدمهما في هذا البحث. يحتوي الأول على بيانات كاملة (Complete data) عن أشخاص، ولعينة عشوائية حجمها 1000، أما الثاني فهو نسخة ناقصة من الأول؛ لاحتوائه على قيم مفقودة (Missing values) في عدة متغيرات. البيانات خاصة اتصالات هاتفية، هدفها تحسين خدمات الشركة. يمثل كل سطر بيانات عن شخص، أي حالة (Case)، أما بيانات كل عمود فهي قيم متغير (Variable).

5-1 الدراسات السابقة

تستند معظم الدراسات، على حد ما هو متوفّر على شبكة المعلومات العالمية (الإنترنت)، وفيما يخص آلية MAR لفقدان القيم، إلى استخدام طرق التعويض الكلي عن القيم المفقودة، وهي تقارن بين تلك الطرق، أما بأسلوب المحاكاة (Simulation) أو بالصيغ الرياضية (Mathematical formulas)، وبالرغم من أن تلك الدراسات لا تستخدم تعبير "التعويض الكلي" (Complete imputation)، إلا أنه الأسلوب المفترض، أما الدراسات عن التعويض الجزئي، فقد اقتصرت على الأسلوب المسمى من قبل Yang and Shoptaw [12]، وبالطبع يتم فصل القيم المفقودة على مجموعتين، الأولى خاصة بالقيم المفقودة بشكل متقطع (Intermittent Missing Values)، والثانية خاصة بالقيم المتروكة (Dropouts)، ويتم التعويض عن الأولى فقط عدة مرات، فنحصل على عدة مجموعات من قيم التعويض، وقد طبق أسلوب التعويض الجزئي المتعدد بدراسة طبية عام 2011 على تجارب سريرية [11]، استخدمت فيها الحزمة الإحصائية MPI 2.0 الخاصة بذلك.

2- الأطر النظرية لفقدان القيمة

1-2 آليات فقدان القيمة

هناك ثلاثة آليات لفقدان القيم، هي:

1-1-1 فقدان عشوائي تام Missing Completely at Random (MCAR)

لفترض أن Y متغيراً يحتوي على قيم مفقودة، وأن X هو متجه (Vector) من متغيرات لا تحتوي على قيم مفقودة. هنا ، تعد القيم (أو البيانات) مفقودة بشكل عشوائي تام (MCAR) إذا كان احتمال فقدان في Y لا يعتمد على X أو على Y نفسها [5]. وقد صاغ Allison [1] ذلك بالشكل التالي:

$$\Pr(R=1|X,Y)=\Pr(R=1) \quad (1)$$

حيث R مؤشر استجابة (Response indicator)، يأخذ القيمة 1 في حالة فقدان في Y ، ويأخذ القيمة 0 في حالة المشاهدة في Y ؛ فالقيم المفقودة في Y لا علاقة لها بالمشاهدات في X ولا بالقيم المفقودة في Y نفسها. ومثل، إذا كان احتمال فقدان قيم في الدخل لا يرتبط بالعمر ولا بالدخل نفسه، فلدينا حالة MCAR.

يمكن تعليم الحالة السابقة عند وجود قيم مفقودة في عدة متغيرات في التحليل، فاحتمال فقدان $\Pr(R=1)$ لأي متغير لا يعتمد على أي متغير آخر في التحليل ولا على القيم المفقودة نفسها [8].
بقي لنا أن نذكر خاصية مهمة لحالة فقدان العشوائي التام للقيم، وهي أن المعلومات المقدرة غير متحيزة (Unbiased)，أي أن عدم التحيز في التحليل الإحصائي يبقى كذلك بالرغم من فقدان في القيم، كما تشير إلى ذلك مصادر عديدة، منها [5] و [9].

1-1-2 فقدان عشوائي Missing at Random (MAR)

تعد القيم المفقودة في Y مفقودة عشوائياً إذا كان احتمال فقدان لا يعتمد على Y نفسها، ولكنه يمكن أن يعتمد على متغيرات أخرى في التحليل. ومثل، إذا كان احتمال فقدان قيم في الدخل لا يعتمد على الدخل نفسه، ولكن يمكن أن يعتمد على العمر. وباستخدام نفس الرموز السابقة، صاغ Allison [1] فرضية MAR ، كالتالي :

$$\Pr(R=1|X,Y)=\Pr(R=1|X) \quad (2)$$

إن MAR يمكن أن تصبح MCAR عندما لا يعتمد احتمال فقدان في Y على أي من المتغيرات الأخرى في التحليل فضلاً عن عدم اعتماده على Y نفسها. وبتعبير آخر، إذا كانت البيانات MCAR فهي أيضاً .MAR

1-1-3 فقدان غير عشوائي Not Missing at Random (NMAR)

بساطة، إذا لم تكن القيم المفقودة في Y مفقودة عشوائياً (MAR)، فإنها مفقودة بشكل غير عشوائي (NMAR). هنا، يرتبط احتمال فقدان في Y على Y نفسها، كما هو موضح في الصيغة الآتية (لم ترد في المصدر السابق):

$$\Pr(R=1|X,Y)=\Pr(R=1|Y) \quad (3)$$

ومثال، الأشخاص ذوو الدخل العالي لا يميلون إلى كتابة دخلكم في استماراة جمع البيانات؛ فقدان القيم في متغير الدخل هنا مرتبط بالدخل نفسه، وبالتالي يكون تقدير متوسط الدخل مثلاً، من بيانات تحتوي على قيم مفقودة، متحيزاً (Biased) بالتأكيد؛ فهناك بيانات مهمة مفقودة بشكل غير عشوائي. وللأسف، لا تتوافر لغاية الآن طريقة فعالة منقولة عليها في معالجة هذه المشكلة، ولكن توجد طرق مقترنة، معقدة جداً، منها تلك المقدمة من قبل Heckman [1] و Dunning and Freedman [2]. ولهذا السبب، تستند معظم الحزم الإحصائية، التي تعالج مشكلة فقدان القيم، على فرضية فقدان العشوائي .

2- الوقاية والمعالجة للقيم المفقودة

الوقاية خير من العلاج، وهذه حقيقة وليس مجرد مثل، فالباحث يريد بيانات دقيقة وكاملة من المشمولين بالاستبيان، وهناك عوامل يجب أن تأخذ بالاعتبار قبل البدء بعملية جمع البيانات، منها يتعلق بالاستبيان نفسها من حيث التصميم وعدد الأسئلة وصياغتها، ومنها يتعلق بزمان ومكان جمع البيانات، ومنها يتعلق بآليات جمع البيانات، وغير ذلك من الأمور المهمة، التي يعتبر الالتزام بها في غاية الأهمية.

بعد جمع البيانات، نلاحظ غالباً، رغم اتباعنا طرق الوقاية، نقصاً في البيانات المجمعة، متمثلة بقيم مفقودة في متغير أو أكثر، وقد يكون الفقدان هذا ذي تأثير غير فعال على النتائج، لأن نحصل على تقديرات غير متحيزة (Unbiased) للمعلمات، ولها أقل تباين (فعالة Efficient)، وقد يكون مؤثراً إلى درجة كبيرة ويطلب علاجاً.

سأعرض باختصار طرق معالجة القيم المفقودة، وتتوفر في [5] مقارنة تفصيلية بينها.
هناك نوعان من المعالجة، هما:

1-2-2-1 الحذف (Deletion).

تتوفر طريقتان للحذف، تستند الأولى على حذف كامل الحالات التي تحتوي على قيم مفقودة، والبقاء على الحالات التامة فقط، وتسمى بطريقة حذف القائمة (Listwise deletion)، أو طريقة الحالة التامة، وتستند الثانية على حذف القيم المفقودة فقط من الحالات، والتعامل مع ما متوفّر من البيانات، وتسمى بطريقة الحذف المزدوج (Pairwise deletion)، أو طريقة الحالة المتوفّرة، وتفترض أي من الطريقتين الآلية لفقدان القيم MCAR.

1-2-2-2 التعويض (Imputation).

ويقصد به إحلال قيم تدريبية محل القيم المفقودة. وهو أما أن يكون أحادياً (Single)， أو متعددًا (Multiple). ومن طرق التعويض الأحادي هي التعويض بالمتوسط، والتعويض بالحالة الأقرب Hot(-) Deck Expectation-Maximization، وتعويض الانحدار، وتعويض EM الذي يستخدم خوارزمية (Markov) MCMC لتقدير القيمة المفقودة. أما أكثر الطرق استخداماً في التعويض المتعدد فهي خوارزمية MAR (Chain Monte Carlo) المستندة إلى الانحدار الخطي. وتستند طرق التعويض إلى فرضية لفقدان القيم.

إن التعويض المتعدد أفضل من التعويض الأحادي [8,10]، وطريقة الإمكان الأعظم (Maximum likelihood) أفضل من التعويض المتعدد (Multiple imputation) [1] ، وأي منها تحتاج إلى خبرة عالية في التطبيق.

3- صيغ الاختبارات

سأستخدم اختبار مربع كاي لمعرفة ما إذا كانت البيانات مفقودة بشكل عشوائي تام (MCAR)، وهي فرضية عدم له، وهو متوفّر ضمن الحزمة الاحصائية SPSS، ويستند إلى الصيغة الآتية [6]:

$$\chi^2_{\text{MCAR}} = \sum_{\text{each unique pattern}} (\text{no. of cases in pattern}) * (\text{MD}) \quad (4)$$

بدرجات حرية محددة بالصيغة الآتية :

$$DF_{\text{MCAR}} = \sum_{\text{each unique pattern}} (\text{no. of nonmissing variables}) - v$$

عندما :

.Mahalanobis D² of pattern mean: MD

v : عدد المتغيرات.

وأسأستخدم اختبار t في حالتي الفقدان والمشاهدة للقيم، وباستخدام متغير تأشير (Indicator variable) يقوم بتحديد ما إذا كان متغير في حالة فقدان أو مشاهدة لحالة (Case)، وهو أيضاً متوفّر ضمن حزمة SPSS، وفيما يلي صيغته [8]:

$$t_{jk} = \frac{\bar{x}_{jk} - \bar{x}_{kj} | \text{variable } j \text{ is missing}}{\left(\frac{\hat{\sigma}_{jk}^P}{n_{jk}} + \frac{\hat{\sigma}_{kj}^P | \text{variable } j \text{ is missing}}{n_{kk} - n_{jk}} \right)^{1/2}} \quad (5)$$

عندما :

k المتغيرات الكمية.

z جميع المتغيرات.

n عدد الحالات.

$$\bar{\mathbf{x}}^P = \left[\bar{x}_{ik}^P \right] = \left[\sum_i x_{ik} / n_{ik}; \quad i \in I(I, k) \right]$$

$$\hat{\sigma}^P = \left[\hat{\sigma}_{ik}^P \right] = \left[\left(\sum_i (x_{ik} - \bar{x}_{ik}^P)^2 / (n_{ik} - 1) \right)^{1/2}; \quad i \in I(I, k) \right]$$

3- الأطر الاجرامي

سأستخدم بيانات الملفين telco_missing.sav و telco.sav، المذكورين في الفقرة 1-4 من هذا البحث، وهي خاصة بشركة اتصالات هانفية، والبيانات مجتمعة من عينة عشوائية حجمها 1000. يحتوي الملف الأول على بيانات كاملة عن أشخاص، أما الثاني فهو نسخة ناقصة من الأول؛ لاحتوائه على قيم مفقودة في عدة متغيرات. كتبت بيانات كل شخص في سطر، وهو يمثل حالة (Case)، أما بيانات الأعمدة فهي قيم متغيرات، ويستدل بحثنا إلى ستة متغيرات كمية، كالتالي :

المعنى	اسم المتغير
عدد الأشهر بالخدمة	Tenure
العمر بالسنوات	Age
عدد السنوات في العنوان الحالي	Address
دخل الأسرة بالألاف	Income
عدد السنوات مع رب العمل الحالي	Employ
عدد أفراد العائلة	Reside

يبين الجدول (1) قيمتي الوسط الحسابي والانحراف المعياري لكل من المتغيرات، وللحالات الكاملة للبيانات، من الملف telco.sav ، فضلاً عن عدد القيم الكاملة لكل متغير (1000)، وعدد القيم المتطرفة، ونلاحظ عدم وجود قيم مفقودة؛ لأن هذا الملف يحتوي على بيانات كاملة، وسأستفيد منه لأغراض المقارنة فيما بعد.

جدول (1)
إحصائيات عن كامل البيانات من الملف .telco.sav

	N	Mean	Std. Deviation	Missing		No. of Extremes	
				Count	Percent	Low	High
tenure	1000	35.53	21.360	0	.0	0	0
age	1000	41.68	12.559	0	.0	0	0
address	1000	11.55	10.087	0	.0	0	13
income	1000	77.5350	107.04416	0	.0	0	93
employ	1000	10.99	10.082	0	.0	0	17
reside	1000	2.33	1.436	0	.0	0	35

سأستخدم بيانات الملف telco_missing.sav لاستخراج الإحصائيات ولنفس المتغيرات، وكما هو موضح في الجدول (2).

جدول (2)
إحصائيات عن البيانات غير الكاملة من الملف .telco_missing.sav

	N	Mean	Std. Deviation	Missing		No. of Extremes	
				Count	Percent	Low	High
tenure	968	35.56	21.268	32	3.2	0	0
age	975	41.75	12.573	25	2.5	0	0
address	850	11.47	9.965	150	15.0	0	9
income	821	71.1462	83.14424	179	17.9	0	71
employ	904	11.00	10.113	96	9.6	0	15
reside	966	2.32	1.431	34	3.4	0	33

السؤال المهم هنا هو: هل أن البيانات مفقودة بشكل عشوائي تام (MCAR)؟

يمكن استخدام اختبار مربع كاي [6] مباشرة لهذا الغرض (وهو متوفّر في معظم الحزم الإحصائية)، وهو الصيغة (4) في البند 3-2، وتستند فرضية عدم له على أن البيانات مفقودة بشكل عشوائي تام، ونتيجة الاختبار مبينة في الجدول رقم (3):

جدول (3)
اختبار فقدان العشوائي التام للبيانات غير الكاملة

EM Means^a

Tenure	age	address	income	employ	reside
35.62	41.72	11.44	76.6841	11.02	2.32

a. Little's MCAR test: Chi-Square = 187.284, DF = 107, Sig. = .000

نلاحظ أن البيانات ليست مفقودة بشكل عشوائي تام؛ لأن القيمة المعنوية (Significant value) للاختبار هي أقل من 0.05، وهنا يلجأ الباحث ذو الخبرة إلى استخدام التعويض المتعدد أو الإمكان الأعظم لمعالجة المشكلة، ولكن ما الذي يمكن أن يفعله من لا يمتلك تلك الخبرة في هذه الحالة، ولدينا نسب فقدان 15.0 % ، 17.9 % ، 9.6 % ، في المتغيرات employ ، income ، address على التوالي (أكثر من 10 % في متغيرين)؟ إن التعويض الكلي للبيانات بالطرق الأحادية يؤدي إلى تقديرات متحيزه (Biased) للمعلمات، وكذلك إلى مشكلة في الأخطاء المعيارية (Standard errors) للتقديرات [1]؛ لأننا نتعامل مع قيم تقديرية مفروضة بفعل التعويض.

سنرى فيما بعد كيف أن طريقة التعويضالجزئي للقيم المفقودة، وبماحدى طرق التعويض الأحادي، البسيطة التطبيق (مثل التعويض بالمتوسط)، تجعل البيانات في حالة MCAR، ومن ثم جواز استخدام الحذف، وسنرى أيضاً أن سبب مشكلة عدم فقدان العشوائي التام للبيانات هو ارتباط فقدان القيم في الدخل (income) بمتغيرات أخرى قيد الدراسة.

سنطبق اختبار t في حالتي فقدان والمشاهدة للقيم، وهو الصيغة رقم (5) في البند 3-2، إذ بين الجدول (4) قيم الاختبار باستخدام متغير تأشير (Indicator variable) يقوم بتحديد ما إذا كان متغير في حالة فقدان أو مشاهدة لحالة (Case)، وتم أيضاً حساب متوسطات المجموعات الجزئية لمتغير التأشير.

جدول (4)
قيم اختبار t ومتوسطات المجموعات الجزئية في حالتي فقدان والمشاهدة.

Separate Variance t Tests^a

	tenure	age	address	income	employ	reside
address						
t	.4	.3				
# Present	819	832	850	693	766	824
# Missing	149	143	0	128	138	142
Mean(Present)	35.68	41.79	11.47	74.08	11.20	2.34
Mean(Missing)	34.91	41.49	.	55.27	9.86	2.21
income	-5.0	-8.3	-3.9	.	-5.9	3.6
t	793	801	693	821	741	792
# Present	175	174	157	0	163	174
# Missing	33.93	40.01	10.67	71.15	9.91	2.39
Mean(Present)	42.97	49.73	14.97	.	15.93	2.02
employ	-1.0	-.4	-.7	.5	.	-.3
t	877	881	766	741	.	874
# Present	91	94	84	80	904	92
# Missing	35.34	41.69	11.37	71.50	0	2.31
Mean(Present)	37.70	42.27	12.32	67.91	11.00	2.37

a. Indicator variables with less than 5% missing are not displayed.

نلاحظ أن متوسط العمر في حالة فقدان الدخل هو 49.73 بينما في حالة عدم فقدان (المشاهدة) هو 40.01، ونلاحظ أيضاً أن فقدان في الدخل أثر على متوسط متغيرات أخرى، وكل ذلك مؤشر على أن البيانات ليست مفقودة بشكل عشوائي تام. إنها ليست مفقودة بشكل كامل بسبب فقدان في متغير الدخل، فهو

استبعينا هذا المتغير من الدراسة، سلاحظ أن البيانات أصبحت MCAR، كما هو موضع في الاختبار الوارد في جدول (5)، حيث نلاحظ أن القيمة المعنوية 0.708، أي لا نرفض فرضية عدم في كون البيانات MCAR، عند مستوى معنوية 0.05.

جدول (5)

اختبار فقدان العشوائي التام للبيانات بعد استبعاد متغير الدخل.

EM Means^a

Tenure	age	address	employ	Reside
35.61	41.72	11.45	11.05	2.32

a. Little's MCAR test: Chi-Square = 45.049, DF = 51, Sig. = .708

ولكننا في الواقع لا نريد استبعاد متغير الدخل من دراستنا. ومن جهة أخرى، لو حصلنا على بيانات كاملة للدخل؛ فنتوقع أن نحصل على نتيجة معنوية لاختبار MCAR. وللتتأكد من ذلك، سأقوم بإدخال قيم متغير الدخل (income) من الملف الكامل للبيانات telco.sav وإحلالها محل البيانات الناقصة للمتغير في الملف telco_missing.sav، وكانتا قمنا بجمع البيانات الناقصة عن متغير الدخل، ويوضح هذا الإجراء من الجدول رقم (6)؛ إذ لا توجد قيم مفقودة لمتغير الدخل.

جدول (6)

إحصائيات عن البيانات غير الكاملة من الملف telco_missing.sav باستثناء متغير الدخل

Univariate Statistics

	N	Mean	Std. Deviation	Missing		No. of Extremes	
				Count	Percent	Low	High
tenure	968	35.56	21.268	32	3.2	0	0
age	975	41.75	12.573	25	2.5	0	0
address	850	11.47	9.965	150	15.0	0	9
income	1000	77.5350	107.04416	0	0	0	93
employ	904	11.00	10.113	96	9.6	0	15
reside	966	2.32	1.431	34	3.4	0	33

وبإجراء اختبار MCAR نحصل على النتيجة المتوقعة، وهي أن البيانات MCAR، وكما هو موضع من القيمة المعنوية للاختبار في الجدول (7).

جدول (7)

اختبار فقدان العشوائي التام للبيانات بعد إحلال بيانات تامة لمتغير الدخل فقط

EM Means^a

tenure	age	address	income	Employ	reside
35.62	41.72	11.45	77.5350	11.02	2.32

a. Little's MCAR test: Chi-Square = 58.411, DF = 67, Sig. = .764

لقد قمنا في الواقع بالتعويضالجزئي عن القيم المفقودة؛ لأننا أحللنا قيمًا محل القيم المفقودة لمتغير الدخل فقط، بعد تشخيصنا له بكونه متغير المشكلة في دراستنا، وتركتنا القيم المفقودة لبقية المتغيرات كما هي، وتحصلنا نتيجة لذلك على نتيجة معنوية لاختبار MCAR. لقد كان مجموع القيم المفقودة للملف tele_missing.sav وللمتغيرات الستة قيد الدراسة هو 516 (34+96+179+150+25+32)، وبعد الإحلال أصبح المجموع 337، فالإحلال شمل 179 قيمة مفقودة خاصة بمتغير الدخل حصرًا.

نستطيع أن نذهب أبعد من ذلك، ونقوم بالتعويضالجزئي للقيم المفقودة لمتغير الدخل نفسه؛ فنستطيع مثلاً تعويض مجموعة من القيم المفقودة فقط بطريقة عشوائية وباستخدام إحدى طرق التعويض الأحادي البسيطة، مثل طريقة المتوسط أو طريقة الحالة الأقرب أو غيرهما، وسنحصل على نتيجة معنوية لاختبار MCAR ، وكما هو موضع في الجدول (8) الخاص بالإحصائيات، والجدول (9) الخاص باختبار MCAR.

جدول (8)

إحصائيات عن البيانات من الملف telco_missing.sav بعد التعويض الجزئي لقيم المفقودة
لمتغير الدخل

Univariate Statistics

	N	Mean	Std. Deviation	Missing		No. of Extremes ^a	
				Count	Percent	Low	High
tenure	968	35.56	21.268	32	3.2	0	0
age	975	41.75	12.573	25	2.5	0	0
address	850	11.47	9.965	150	15.0	0	9
income	867	72.1326	81.84437	133	13.3	0	64
employ	904	11.00	10.113	96	9.6	0	15
reside	966	2.32	1.431	34	3.4	0	33

جدول رقم (9)

اختبار الفقدان العشوائي التام للبيانات بعد التعويض الجزئي لقيم المفقودة لمتغير الدخل.

EM Means^a

tenure	age	address	income	employ	reside
35.62	41.72	11.44	74.3177	11.03	2.32

a. Little's MCAR test: Chi-Square = 120.836, DF = 107, Sig. = .170

نلاحظ من جدول رقم (8) أن التعويض شمل 46 ($46 = 133 - 179$) قيمة مفقودة فقط من مجموع القيم المفقودة لمتغير الدخل (179). وبتعبير أعم، إن التعويض شمل 46 قيمة مفقودة فقط من المجموع الكلي لقيم المفقودة للمتغيرات الستة والبالغ 516 قيمة، أي أن نسبة تعويض بمقدار 9% تقريباً كانت كافية لجعل البيانات MCAR.

في الواقع العملي، غالباً ما نحاول جمع بيانات حقيقة من المستجيبين تحل محل البيانات الناقصة لمتغير الدخل؛ فمن 179 شخص قد تحصل على 46 (%) قيمة فعلية، تجعل البيانات MCAR، وهذا بالتأكيد أفضل من التعويض بقيم تقديرية.

4 الاستنتاجات وعمل مستقبلي (Conclusions and Future Work)

4 - 1 الاستنتاجات (Conclusions)

عند ملاحظة وجود قيم مفقودة بعد جمع البيانات، علينا أولاً محاولة جمع بيانات حقيقة من المستجيبين تحل محل البيانات الناقصة؛ لأن ذلك أفضل من التعويض بقيم تقديرية، وقد يكفي ما نحصل عليه لجعل البيانات MCAR، إن لم تكن كذلك، وبالتالي نستطيع استخدام الحذف للمعالجة، وهي طريقة في متداول الجميع، ومتوفرة في الحزم الإحصائية، أما عند تعذر ذلك، وعندما لا تكون البيانات MCAR، وفقاً للاختبار، فيمكننا استخدام التعويض الجزئي لجعلها كذلك، وذلك بإدخال قيمة مقدرة محل مجموعة من القيم المفقودة وليس جميعها، وبأحدى طرق التعويض الأحادي، السهلة التطبيق، والمتوافرة في معظم الحزم الإحصائية، ثم استخدام الحذف للمعالجة.

إجراء التعويض الجزئي بسيط، واختبار MCAR أيضاً، وكذلك استخدام الحذف للمعالجة، وكل ذلك ممكن التطبيق من قبل الباحث غير المتخصص بمفردته، دون مساعدة خبير في الإحصاء، أما طرق التعويض المتعدد والإمكان الأعظم، فتحتاج إلى خبرة في التطبيق، بدونها لن يحصل الباحث إلا على نتائج غير دقيقة، وإن توفرت هذه الطرق المتقدمة في الحزمة الإحصائية.

4 - 2 عمل مستقبلي (Future Work)

لاحظنا من جدول رقم (8) أن التعويض شمل 46 قيمة مفقودة فقط من مجموع القيم المفقودة لمتغير الدخل (179)، ومن المجموع الكلي لقيم المفقودة للمتغيرات الستة والبالغ 516 قيمة. وبحساب النسبة فإنها تشكل نسبة تعويض 9% تقريباً، وكانت كافية لجعل البيانات MCAR. إن هذه النسبة لا تشكل الحد الأدنى المطلوب من القيم التعويضية، فقد استندت إلى التجربة كما لاحظنا، وليس إلى صيغة إحصائية محددة؛ لعدم توفر مثل هذه الصيغة لغاية الآن.

إن العمل على إيجاد صيغة أو معيار نستند إليه كدليل لمعرفة الحد الأدنى لنسبة القيم التعويضية، التي تجعل البيانات MCAR، والذي سيأخذ بالاعتبار طبيعة البيانات الخاصة للتحليل، يحتاج بالتأكيد إلى مجهود كبير ووقت طويل، وهو ما أنوي القيام به مستقبلاً، فضلاً عن تطبيق طريقة التعويض الجزئي لقيم المفقودة، المقترنة في هذا البحث، على أحجام صغيرة ومتوسطة من البيانات؛ للتأكد من كفافتها.

المصادر (References)

- 1- Allison, P.D. (2012)'Handling Missing Data by Maximum Likelihood', SAS Institute, Statistical Horizons, Haverford, PA, USA.
 - 2- Dunning, T. and Freedman, D.A. (2008)'Modeling section effects', London, UK.
 - 3- Haitovsky, Y. (1968) 'Missing data in regression analysis', *Journal of the Royal Statistical Society, Series B* 30: 67–82.
 - 4- Heckman, J. J. (1979) 'Sample selection bias as a specification error', *Econometric* 47: 153–161.
 - 5- Little, R. J. A. and Rubin, D. B. (1987) 'Statistical Analysis with Missing Data', New York, Wiley.
 - 6- Little, R.J.A. (1988) 'A Test of Missing Completely at Random for Multivariate Data with Missing Values', *Journal of the American Statistical Association* 83.
 - 7- Little, R.J.A. (1992)'Regression with missing X's: a review', *Journal of the American Statistical Association* 87.
 - 8- Rubin, D. B. (1987) 'Multiple Imputation for Non-response in Surveys', New York, Wiley.
 - 9- Schafer, J. L. (1997) 'Analysis of Incomplete Multivariate Data', London, Chapman and Hall.
 - 10- Scheffer, J. (2002) 'Dealing with Missing Data', Available online at <http://www.massey.ac.nz/~wwiims/research/letters/>.
 - 11- Xiaowei Yang, Jinhui Li, and Steven Shoptaw (2011)' Multiple Partial Imputation for Longitudinal Data with Missing Values in Clinical Trials', University of California, <https://escholarship.org/uc/item/9733x421>.
 - 12- Yang X, Shoptaw S. (2005) 'Assessing missing data assumptions in longitudinal studies: an example using a smoking cessation trial', *Drug and Alcohol Dependence*; 77: 213-225
-
.....
.....