

القسم الثالث

المعالجة الأولية للبيانات

يتناول هذا القسم عرضاً للخطوة الثالثة من خطوات تكوين المؤشر المركب، وهى خطوة المعالجة الأولية للبيانات، حيث يعرض الجزء الأول من هذا القسم تقدير القيم المفقودة بالبيانات، بينما يعرض الجزء الثانى التأكد من مدى جودة البيانات الأساسية من خلال دراسة مدى اتفاق هذه البيانات مع معايير جودة البيانات العالمية.

١.٣ تقدير البيانات المفقودة (Imputation):

تعتبر البيانات المفقودة من أكثر المعوقات التى تحد من متانة وجوده المؤشر المركب، لذلك فقد كان من الضرورى استعراض بعض الطرق التى تُستخدم لمعالجة وتقدير تلك البيانات، إلا أنه قبل استعراض هذه الطرق كان لابد من سرد أنواع البيانات المفقودة، حيث تختلف طرق التقدير تبعاً لاختلاف نوع البيانات المفقودة.

١.١.٣ أنواع البيانات المفقودة:

يمكن تقسيم البيانات المفقودة إلى ثلاثة أنواع أساسية: بيانات مفقودة بشكل عشوائى تماماً (Missing completely at random)، بيانات مفقودة بشكل عشوائى (Missing at random)، وبيانات مفقودة بشكل غير عشوائى (Not missing at random)، وفيما يلى عرض لتلك الأنواع بشيء من التفصيل:

(أ) بيانات مفقودة بطريقة عشوائية تماماً (MCAR) Missing completely at random:

فى هذا النوع من فقد البيانات نجد أن البيانات المفقودة لا تعتمد على المتغير نفسه - الذى يحتوى على هذه القيم المفقودة - أو أى متغير آخر فى قاعدة البيانات، فعلى سبيل المثال تعتبر القيم المفقودة من متغير الدخل من نوعية البيانات المفقودة بطريقة عشوائية تماماً فى حالة كون القيم المفقودة غير مرتبطة بالدخل أو بأى من المتغيرات الأخرى فى قاعدة البيانات. بمعنى أنه إذا ما كان الأشخاص الذين لم يقرؤوا بقيم دخولهم لهم - فى المتوسط - نفس قيم دخول الأشخاص الذين أقرؤوا بالدخل، أو إذا

كان كل متغير من المتغيرات الأخرى الموجودة في قاعدة البيانات هو نفسه - في المتوسط - لكل من الأشخاص الذين أقرؤا بقيمة الدخل والأشخاص الذين لم يقرؤا به.

(ب) بيانات مفقودة بطريقة عشوائية (MAR) Missing at random:

نجد أن البيانات المفقودة في هذا النوع لا تعتمد على المتغير محل الدراسة، إلا أنها ترتبط بمتغيرات أخرى في قاعدة البيانات، فعلى سبيل المثال قد تعتمد نسبة البيانات المفقودة للدخل على الحالة الاجتماعية، إلا أنه في داخل كل فئة من فئات الحالة الاجتماعية تكون نسبة البيانات المفقودة غير مرتبطة بقيمة الدخل، وهذا ما يعنى أن القيم المفقودة لا تعتمد على متغير الدخل إلا أنها تعتمد على الحالة الاجتماعية للأفراد، وكذلك من حالات فقد البيانات بطريقة عشوائية هي حالة البيانات المفقودة نتيجة لتصميم الإستمارة (Missing data by design، مثل حالة عدم انطباق السؤال على المجيب (not applicable) فمثلاً إذا كانت إجابة سؤال معين بنعم فهذا يتبعه عدم الإجابة عن سؤال آخر بالإستمارة.

(ج) بيانات مفقودة بطريقة غير عشوائية (NMAR) Not missing at random:

في هذه الحالة تعتمد البيانات المفقودة على القيم نفسها، فعلى سبيل المثال، الأفراد ذوى الدخل المرتفعة غالباً لا يذكرون دخولهم. وغالباً لا يوجد أساس معين للحكم على ما إذا كانت البيانات مفقودة بطريقة عشوائية أم بطريقة غير عشوائية، إلا أن معظم طرق تقدير البيانات المفقودة تتطلب أن يكون السبب عشوائياً (MCAR, MAR)، وإذا كان هناك اتجاهات تشير إلى وجود نمط غير عشوائى لفقد البيانات، فإنه يجب توضيح هذا النمط ونمذجته بدقة، إلا أن هذه العملية تتسم بالصعوبة والتعقيد.

٣.١.٢ أساليب تقدير البيانات المفقودة:

بعد أن تم التعرف على الأنواع الثلاثة لفقد البيانات، يتم الانتقال لعرض الطرق الأساسية التى يمكن استخدامها لمعالجة البيانات المفقودة، حيث تتمثل هذه الطرق فى الآتى:

(أ) حذف المفردة (Case deletion):

تعتمد هذه الطريقة على حذف المشاهدة - التى تحتوى على قيمة مفقودة لأحد المتغيرات - نهائياً من التحليلات، ولذلك فهى تسمى أيضاً تحليل الحالة الكاملة (complete case analysis) حيث أنها تعتمد ببساطة على حذف الحالات غير الكاملة من التحليل، إلا أنه من عيوب هذه الطريقة أنها

تتجاهل احتمال وجود اختلافات بين العينات الكاملة وغير الكاملة، كما أنها تتطلب فقد البيانات بطريقة عشوائية تماماً، حيث أنها تؤدي إلى تقديرات غير متحيزة فقط إذا كانت المشاهدات المحذوفة تعتبر عينة جزئية عشوائية من العينة الأصلية. هذا بالإضافة إلى أن الأخطاء المعيارية سوف تكون أكثر في العينة غير الكاملة نظراً لاستخدام معلومات أقل، ووفقاً للتجربة يمكن القول أن القاعدة العامة التي يمكن الاعتماد عليها في هذه الطريقة هي عدم حذف أي من الحالات إذا كان المتغير لديه أكثر من ٥٪ من القيم المفقودة.

(ب) التقدير الفردي (Single imputation):

تعتمد هذه الطريقة على تقدير البيانات المفقودة من خلال عمل توزيع تنبؤي (Predictive distribution) لهذه البيانات المفقودة، ويجب أن يُولد هذا التوزيع التنبؤي عن طريق توظيف البيانات من خلال تصميم إما نماذج ظاهرة (Explicit modeling)، أو نماذج ضمنية (Implicit modeling).

تعتمد النماذج الضمنية على نظام حسابي ذو افتراضات ضمنية يجب اختبار كونها مناسبة وجيدة للغرض محل الاهتمام أم لا، ومن أمثلتها:

● تقدير هوت دك (Hot deck imputation):

حيث يتم تقدير البيانات المفقودة ببيانات مسحوبة من وحدات مشابهة، على سبيل المثال: القيم المفقودة للدخل لأحد الأفراد يمكن التعويض عنها بقيم الدخل لأفراد لهم نفس الخصائص.

● التعويض (Substitution):

تعتمد هذه الطريقة على إحلال الوحدات غير المستجيبة بوحدات لم يتم اختيارها في العينة، فمثلاً إذا لم يتاح الاتصال بأحد الأسر، فإننا نقوم باختيار أسرة تنتمي إلى نفس الكتلة السكنية لم تكن مختارة في العينة.

● تقدير كولد دك (Cold deck imputation):

حيث يتم استبدال البيان الناقص ببيان من مصدر آخر، كأخذ البيان من إصدار سابق لنفس الاستقصاء.

أما بالنسبة للنماذج الظاهرة، فيعتمد التقدير فيها على نماذج إحصائية حيث تكون الافتراضات واضحة وظاهرة، ومن أمثلة هذه النماذج:

● **التقدير باستخدام المتوسط الحسابي / الوسيط / المنوال الغير شرطي (Unconditional mean/median/mode imputation):**

وتعتمد هذه الطريقة على استبدال البيانات المفقودة بقيمة المتوسط (mean) – أو الوسيط (median) أو المنوال (mode) – المقدر باستخدام المشاهدات الغير مفقودة.

● **تقدير البيانات عن طريق تحليل الانحدار (Regression imputation):**

تعتمد هذه الطريقة على التعويض عن القيم المفقودة بالقيم المتنبأ بها من تحليل الانحدار، على أن يكون المتغير التابع في تحليل الانحدار هو المؤشر الفرعي الذي لديه قيم مفقودة، ويكون المتغير المستقل (المتغيرات المستقلة) هو المؤشر الفرعي (المؤشرات الفرعية) الذي يعكس علاقة قوية مع المتغير التابع، أى توجد درجة عالية من الارتباط بين المتغير المستقل (المتغيرات المستقلة) والمتغير التابع، ومن الجدير بالذكر هنا أنه يجب أن يتم عمل الانحدار باستخدام المشاهدات الكاملة، أى التى لا يوجد بها قيم مفقودة للمتغير محل الاهتمام.

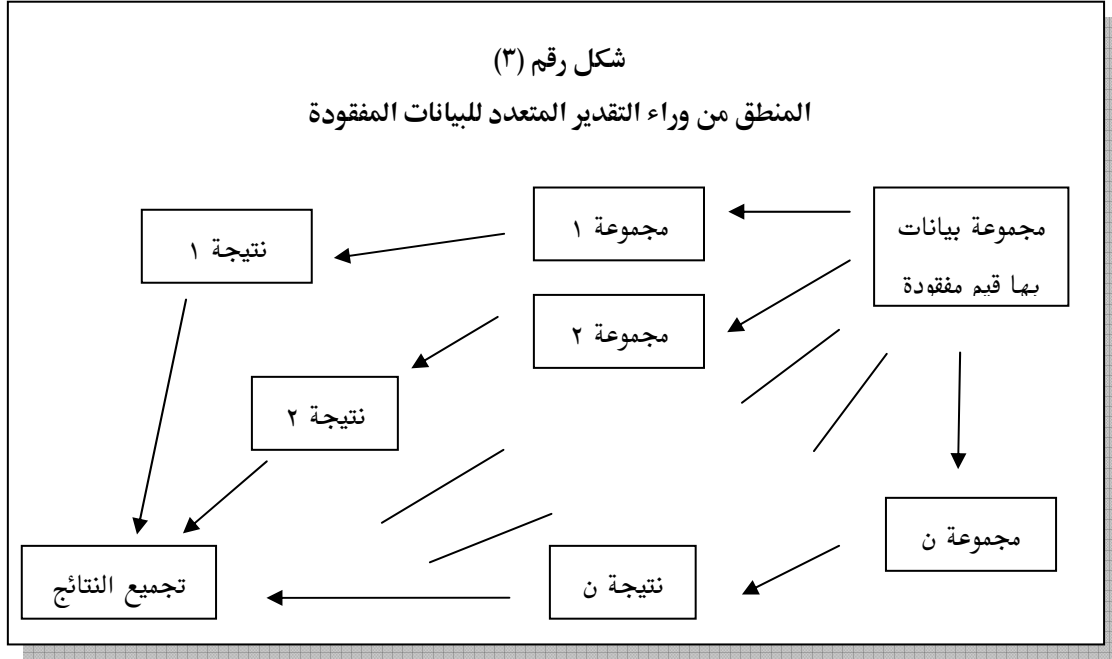
● **التقدير باستخدام القيمة العظمى المتوقعة (Expected Maximization Imputation):**

تعتمد هذه الطريقة على التعويض عن القيم المفقودة من خلال مقدرات يتم الحصول عليها من خلال عملية تقدير متكررة (Iterative process)، حيث يتم أولاً التنبؤ بالقيم المفقودة بناءً على مقدرات أولية (Initial estimates) لقيم معالم النموذج، ثم تستخدم هذه التنبؤات بعد ذلك لتعديل قيم المعالم، ويتم تكرار عملية التنبؤ باستخدام معالم النموذج المعدلة حتى تقترب المعالم من تقديرات الإمكان الأكبر (Maximum-likelihood estimates)، وبذلك فإن تقدير البيانات المفقودة باستخدام هذه الطريقة يتكون من جزأين: الأول هو تقدير معالم النموذج عن طريق الإمكان الأكبر (Maximum Likelihood) بفرض أنه لا توجد بيانات ناقصة، والثانى هو حساب القيمة المتوقعة للبيان المفقود بواسطة المعالم المقدرة.

(ج) التقدير المتعدد للبيانات المفقودة (Multiple imputation):

تعتمد هذه الطريقة على تقدير البيانات المفقودة (ن) مرة، بهدف الحصول على (ن) مجموعة كاملة من البيانات، بحيث أنه فى كل مجموعة من البيانات يتم تقدير المعلمات محل الاهتمام، وكذلك تقدير أخطائها المعيارية، ثم يتم حساب متوسط التقديرات (وسط حسابى أو وسيط) لجميع المجموعات

(ن مجموعة)، وكذلك التباين داخل المجموعات، والتباين بين المجموعات. والشكل التالي يوضح المنطق من وراء التقدير المتعدد للبيانات المفقودة.



ومن الجدير بالذكر هنا أنه يمكن استخدام أى طريقة من طرق التقدير السابقة لتقدير البيانات في كل مجموعة من المجموعات.

٢.٣ جودة البيانات الأساسية:

يعتبر الأساس الذي يُبنى عليه المؤشر المركب هو البيانات الأساسية للمؤشرات الفرعية، لذلك فيجب أن يكون الأساس جيد حتى يكون البناء جيد، وحتى نستطيع الاعتماد على المؤشر المركب فيما بعد. وحيث إن جودة المؤشر المركب الكلية تعتمد على جودة البيانات المبدئية المستخدمة في بناءه، فإن هذا يستدعي التأكد من جودة البيانات الأساسية حتى يتم اختيارها بشكل يعظم الجودة الكلية للنتائج النهائية.

١.٢.٣ معايير الجودة:

البيانات الأساسية المستخدمة في تكوين المؤشر المركب لا بد أن تُحقق معايير جودة البيانات الأساسية التي تتمثل في:

● نفعية البيانات (Relevance):

يتم تقييم نفعية البيانات الأساسية المكونة للمؤشر المركب مع الأخذ في الاعتبار الهدف النهائي من المؤشر، بحيث يتم تقييم واختيار البيانات الأساسية بحذر لضمان أن هذا المعيار تم تغطيته بشكل مناسب. فمثلاً بافتراض عدم التوافر الفعلي للبيانات، ففي أغلب الأحيان يتم استخدام سلاسل بديلة، ولكن في هذه الحالة لا بد أن تتوفر بعض الأدلة عن علاقة هذه السلاسل البديلة بسلاسل الهدف.

● الدقة (Accuracy):

يعكس هذا المعيار مدى دقة وصف البيانات للظاهرة التي صممت من أجل قياسها، كما يوضح مدى القرب بين القيمة المقدرة والقيمة المجهولة، وتعتبر دقة البيانات الأساسية من المعايير الهامة جداً للجودة، وهذا ما يستدعي أن يتم الحصول على البيانات الأساسية للمؤشر المركب بمستوى عالٍ من الدقة، إلا أنه من الناحية العملية لا يوجد مقياس للدقة، ولكنه قد يمكن قياس مستوى الدقة بالبيانات الأساسية، من خلال قياس الخطأ في المقدرات التي تعتمد على عينات المسوح، ويعتبر أهم مصادر الأخطاء في هذه المقدرات هو خطأ التغطية، وخطأ المعاينة، وخطأ عدم الإستجابة.

● حداثة البيانات (Timeliness):

يتمثل هذا المعيار في الحصول على البيانات في الوقت المناسب ووفقاً لتواريخ مسبقة، كما يعكس الفترة ما بين المرحلة البحثية أو نهاية الفترة البحثية والوقت الذي تكون فيه البيانات متاحة، ويعتبر هذا المعيار من الأبعاد الهامة جداً لتقليل الاحتياج إلى تقدير البيانات المفقودة ومراجعة البيانات المنشورة مسبقاً لتنقيحها، ومن الجدير بالذكر في هذا المعيار هو أن البيانات التي تغطي مجالات مختلفة تكون متاحة في أوقات مختلفة خلال الزمن، بمعنى أنه قد تكون بيانات كل مؤشر من المؤشرات الفردية متاحة في وقت مختلف عن المؤشرات الأخرى، ولذلك فيجب أن يكون هناك تناسب بين توقيتات الحصول على البيانات المكونة للمؤشر المركب.

● إمكانية الوصول إليها (Accessibility):

يشير هذا المعيار إلى مدى سهولة الحصول على البيانات، وهذا يتضمن سهولة الوصول إلى البيانات، بالإضافة إلى وضوح طريقة عرض تلك البيانات، ونلاحظ أن سهولة أو صعوبة الوصول

إلى البيانات يؤثر على التكلفة الإجمالية لإنتاج وتجديد المؤشر المركب خلال الزمن، كذلك يستطيع هذا المعيار التأثير على مصداقية المؤشر المركب إذا كان هناك صعوبة فى الوصول إلى البيانات، حيث أنه قد يؤثر على المعايير الأخرى، فقد تؤثر صعوبة الوصول إلى البيانات على حداتها حيث يتم الوصول إلى البيانات فى توقيت غير مناسب، وقد تؤثر صعوبة الوصول للبيان على دقة البيانات التى يتم الحصول عليها، ونلاحظ أنه فى ظل التطورات التكنولوجية والقدرة على الوصول إلى قواعد البيانات المصدرة من جهات مختلفة، أصبح من السهل الوصول إلى الكثير من البيانات، وكذلك إيجاد العلاقات بين المصادر المختلفة للبيانات، وبالتالي الحصول على أكبر قدر ممكن من المنفعة للبيانات، ولكن من الجدير بالذكر هنا أنه لا بد عند اختيار مصدر البيانات أن لا يتم تفضيل المصدر السهل الوصول إليه فقط، ولكن يجب النظر فى معايير الجودة الأخرى.

● إمكانية التفسير (Interpretability):

يعكس هذا المعيار مدى السهولة التى يستطيع بها المستخدم تفسير البيانات، وذلك ما يتم عن طريق توافر المعلومات الإضافية لتفسير واستخدام البيانات بشكل ملائم، وتشتمل هذه المعلومات عادةً على المفاهيم الأساسية والمتغيرات والتصنيفات المستخدمة ومنهجية جمع البيانات، ونظراً لاتساع مدى البيانات المستخدمة لبناء المؤشر المركب، وكذلك الصعوبات المتعلقة بالإجراءات التجميعية المطلوبة، فإن ذلك يحتاج إلى تفسير كامل، ويعتبر هذا المعيار من المعايير الهامة للجودة، حيث أن توافر التعريفات والتصنيفات المستخدمة لإنتاج البيانات يكون ضرورى لتقييم المقارنات بين البيانات عبر الزمن وبين الدول المختلفة.

● التماسك (Coherence):

يقيس هذا المعيار درجة سهولة مقارنة البيانات باختلاف الزمان أو المكان، أى تماسك البيانات عبر الزمن وبين الدول المختلفة، ويشير التماسك عبر الزمن ضمناً إلى أن البيانات مستندة على مفاهيم وتعريفات ومنهجيات عامة ومشاركة عبر الزمن، حيث يشير عدم التماسك عبر الزمن إلى التغيرات التى تحدث فى السلاسل الزمنية نتيجة للتغير فى المفاهيم أو التعريفات أو المنهجيات المستخدمة. أما بالنسبة للوجه الثانى للتماسك، فهو التماسك عبر الدول البلدان، وهو ما يشير ضمناً إلى أنه من بلد إلى بلد يظل هناك اشتراك فى المفاهيم والتعريفات والتصنيفات والمنهجيات. وتظهر أهمية الارتباط والاتساق المنطقى فى أن المسوح الميدانية الإحصائية قد يتم تطبيقها على

مجتمعات مختلفة وفي أوقات مختلفة، لذلك فيجب أن تكون الإحصاءات الناتجة عنها مترابطة، ويمكن المقارنة بينها حتى تكون بيانات المؤشرات الفرعية متماسكة للدول المختلفة والأزمنة المختلفة.