

تقدير القيم المفقودة لمتغير الاستجابة في نموذج الانحدار الخطي المتعدد
Estimating of lost value of Responding variable in the Multi Regression Model


م.ع.ع. علي ناصر حسين
Assist. L. A. N. Hussein
قسم الإحصاء / كلية الإدارة والاقتصاد
جامعة البصرة

المستخلص

EM

Algorithm

(-) ()



١- المقدمة

$$\underline{Y} = X\underline{B} + e \quad \dots\dots\dots (1)$$

2-أنماط البيانات المفقودة

١-٢ Univariate missing data

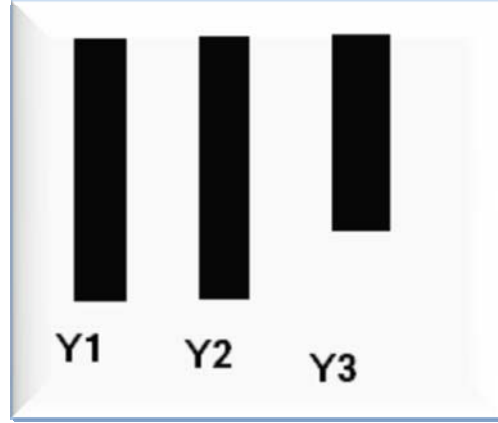
k

k-1



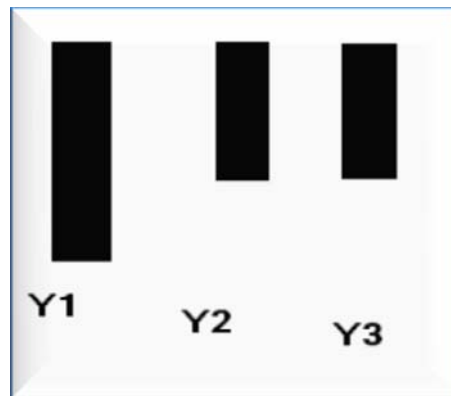


شكل (١) النمط الأول من البيانات غير التامة



Multivariate two patterns ٢-2

شكل (٢) النمط الثاني من البيانات غير التامة



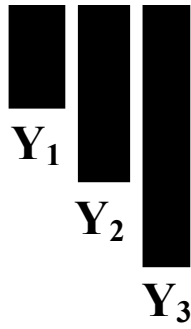
(-)

()

.....

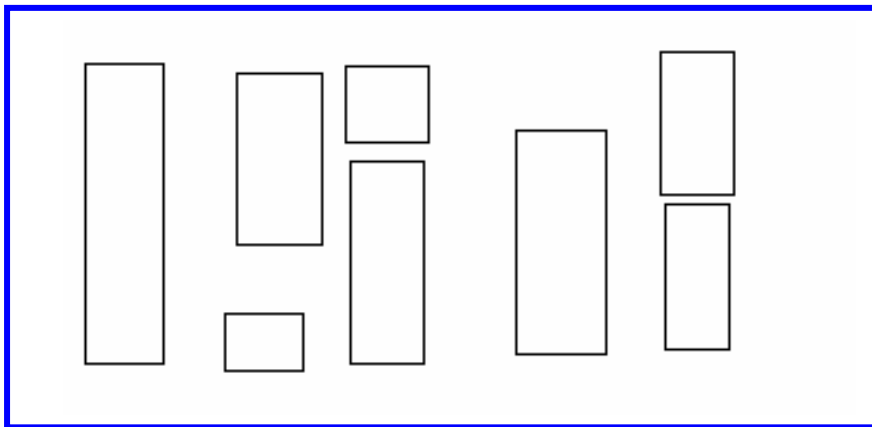
Monotone ٣-٢

شكل رقم (٣) النمط الثالث من البيانات غير التامة



General ٤-٢

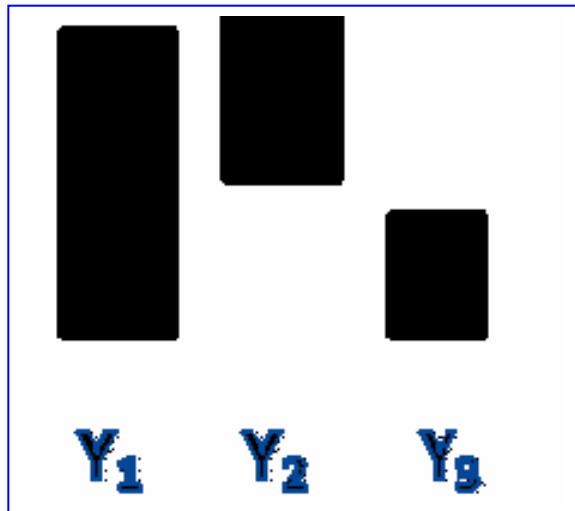
شكل رقم (٤) النمط الرابع من البيانات غير التامة





File matching ٥-٢

الشكل (٥) النمط الثالث من البيانات غير التامة



3-آليات الفقدان (1,3,4,6)

()

(Missing Data Indicator

(n*p)

:R

matrix)

m

R

M

$$m_{ij} = \begin{cases} 1 & \text{if } r_{ij} \text{ is observed} \\ 0 & \text{if } r_{ij} \text{ is missing} \end{cases}$$

:

M

R

(-) () ..

١-٣ آلية الفقدان بشكل عشوائي تام (MCAR) Missing Completely at Random

$$p_r(M_m/R, \emptyset) = p_r(M_m/\emptyset) \quad \forall M_m \dots\dots (2)$$

: M_m
 \emptyset

٢-٣ آلية الفقدان بشكل عشوائي (MAR) Missing at Random

$$p_r(M_m/R, \emptyset) = p_r(M_m/R_o, \emptyset) \quad \forall M_m \dots\dots\dots (3)$$

٣-٣ آلية الفقدان بشكل غير عشوائي (NMAR) Missing not at Random

: MNAR

$$p_r(M_m/M_o R, \emptyset) = p_r(M_m/R_m, \emptyset) \quad \forall M_m \dots\dots\dots(4)$$

R_o :

4- طريقة المربعات الصغرى بالتعويض عن البيانات المفقودة:



1-4- طريقة المتوسط غير الشرطي:

$$\hat{y}_i = \sum_{n_j} \frac{y_{obs}}{n_j} \dots\dots\dots (5)$$

.(y) :n_j

$$n_j - 1/n - 1 \dots\dots\dots (6)$$

$$n_{jk} - 1/n - 1 \dots\dots\dots (7)$$

(x)

$$E(y_i/x_i) = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} \dots\dots\dots (8)$$

$$E(y_i/x_{i2} \dots x_{ip}) = \beta_0 + \beta_1 x_{i1}^* + \sum_{j=1}^p \beta_j x_{ij} \dots\dots\dots (9)$$

Where $\beta_1 x_{i1}^* = E(x_{i1}/x_{i2} \dots x_{ip})$

2-4- طريقة التعويض بالانحدار:



(-) ()

$$y = x\beta + \epsilon \dots \dots \dots (10)$$

	(n*1)		y
p-1	(n*P)		x
		(p*1)	β
			ε
	n		m

(10)

$$\begin{bmatrix} y_{obs} \\ y_{mis} \end{bmatrix} = \begin{bmatrix} x_{obs} \\ x_{mis} \end{bmatrix} \beta + \begin{bmatrix} \epsilon_{obs} \\ \epsilon_{mis} \end{bmatrix} \dots \dots \dots (11)$$

	(r=n-m,1)	: y _{obs}
.	(m,1)	y _{mis} :
.	(r,p)	x _{obs} :
	(m,p)	x _{mis}
		ε _{obs} :
		ε _{mis}

$$y_{mis} = x_{mis} \hat{\beta} \dots \dots \dots (12)$$





$$\hat{\beta} = (X_{obs} X_{obs}')^{-1} X_{obs} Y_{obs} \dots \dots \dots (13)$$

E M Algorithm -5

(EM)

(Maximum Likelihood ML)

:

E-)

(step

$$E(\sum_{i=1}^n y_{ij} / Y_{obs} \theta^t) = \sum_{i=1}^n y_{ij} \dots \dots \dots (14)$$

j=1,2,...k

$$E(\sum_{i=1}^n y_{ij} y_{ik} / Y_{obs}, \theta^t) = \sum_{i=1}^n (y_{ij}^t y_{ik}^t + c_{jki}^t) \dots \dots \dots (15)$$

Where I,j=1,2 ...k , i≠k

$$y_{ij}^t = \begin{cases} y_{ij} & \text{if } y_{ij} \text{ is observed} \\ E(y_{ij} / Y_{obs}, \theta^t) & \text{if } y_{ij} \text{ is missing} \end{cases} \dots \dots \dots (16)$$

$$c_{jki}^t = \begin{cases} 0 & \text{if } y_{ij} \text{ or } y_{ik} \text{ is observed} \\ cov(y_{ij}, y_{ik} | Y_{obs}, \theta^t) & \text{if } y_{ij} \text{ and } y_{ik} \text{ are missing} \end{cases} \dots \dots \dots (17)$$

Θ : M

: E θ^{t+1}

$$\hat{\mu}_j = \frac{\sum_{i=1}^n y_{ij}^t}{n} \dots \dots \dots (18)$$

$$\sigma_{ij}^{t+1} = \frac{1}{n} \sum_{i=1}^n [(x_{ij}^t - \mu_j^t)(x_{ik}^t - \mu_k^t) + c_{jki}^t] \dots \dots \dots (19)$$

-6 الطريقة المقترحة :

E step (E M Algorithm)

:

$$\hat{y}_{ij} = [1 - \lambda \sigma^2 tr(w) / (\hat{\beta}' x_m' \Sigma^{-1} x_m \hat{\beta})] x_m \hat{\beta} \dots \dots \dots (20)$$



(-) () ..
 ..

$$\lambda = (n - p)(m - 2) / ((n - p + 2)tr(w)) \dots \dots \dots (21)$$

$$\Sigma = x_{mis}(x'_{obs}x_{obs})^{-1}x'_{mis} \dots \dots \dots (22)$$

$$w = m'\Sigma^{-1}m \dots \dots \dots (23)$$

$$m = x_{mis}(x'_{obs}x_{obs})^{-1} \dots \dots \dots (24)$$

□

M

M E

$$|\hat{\mu}^{t+1} - \hat{\mu}^t| < \delta I \dots \dots \dots (25)$$

$$|\hat{\Sigma}^{t+1} - \hat{\Sigma}^t| < \delta J \dots \dots \dots (26)$$

7- الجانب التطبيقي:

matlab

(μ, σ^2)

(Σ)

$$\mu = (\mu_1, \mu_2, \dots, \mu_3)'$$

x₁, x₂, ..., x_d as iid N(0, 1) .I

$$x_i = \mu_i + \sum_{j=1}^i c_{ij} z_j \dots \dots \dots .II$$

i=1, 2 ... d

(cholesky

(Σ)

: *d*d* c

decomposition)

(c)

(Σ)

$$cc' = \Sigma \dots \dots \dots (27)$$





: c

$$c_{ii} = \left(\sum_{k=1}^{i-1} c_{ik}^2 \right)^{\frac{1}{2}} \dots \dots \dots (28)$$

$$c_{jt} = \frac{1}{c_{ii}} \left(\sum_{k=1}^{i-1} c_{ik} c_{jk} \right) \dots \dots \dots (29)$$

:

$$\underline{y} = \underline{x}\underline{\beta} + \underline{e} \dots \dots \dots (30)$$

8-العينة:

(0.5) ()

(5%,10%,15%) (MCAR)

MCAR

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-p} \dots \dots \dots (31)$$

9-تحليل النتائج:

$$y_i = 2 + 1.5x_{i1} + .5x_{i2} + .2x_{i3} + e_i$$

MCAR y

:

EM % -

MSE

.(0.92) (0.872) (0.821)



(-) () ..

EM % -

MSE

(2.491) (2.492) (2.451)

EM % -

MSE

(0.952) (0.973) (0.695)

جدول (١) نتائج تقدير البيانات المفقودة لطرائق التقدير

حجم العينة	نسبة الفقدان	طريقة التقدير		B0	B1	B2	B3	MSE
50	5%	طريقة EM المقترحة	قيمة المعلمة	0.123	1.032	0.584	0.316	0.821
		الانحراف المعياري	0,094	0.102	0.092	0.102		
		طريقة التعويض بالانحدار	قيمة المعلمة	0.127	1.05	0.591	1.314	0.872
		الانحراف المعياري	1	0.11	0.097	0.107		
		طريقة المتوسط غير الشرطي	قيمة المعلمة	0.15	0.957	0.549	0.317	0.92
		الانحراف المعياري	0.099	0.108	0.097	0.108		
	10%	طريقة EM المقترحة	قيمة المعلمة	0.2204	0.3846	0.0602	0.2337	2.4 ^٥ 1
			الانحراف المعياري	0.1608	0.1384	0.1614	0.1605	
		طريقة التعويض بالانحدار	قيمة المعلمة	0.2462	0.4173	0.0610	0.2879	2.4 ^٩ 2
			الانحراف المعياري	0.1596	0.1373	0.1601	0.1593	
		طريقة المتوسط غير الشرطي	قيمة المعلمة	0.2204	0.3846	0.0602	0.2337	2.491
			الانحراف المعياري	0.1608	0.1384	0.1614	0.1605	
15%	طريقة EM المقترحة	قيمة المعلمة	0.248	0.982	0.503	0.1806	0.695	
		الانحراف المعياري	0.084	0.092	0.074	0.0819		
	طريقة التعويض بالانحدار	قيمة المعلمة	0.1832	0.7658	0.4376	0.1589	0.973	
		الانحراف المعياري	0.099	0.109	0.088	0.0969		
	طريقة المتوسط غير الشرطي	قيمة المعلمة	0.221	0.776	0.436	0.1566	0.952	
		الانحراف المعياري	0.098	0.108	0.087	0.0958		



ثانيا : في حالة حجم العينة 100 نلاحظ الاتي

EM % -

MSE

(1.026) (0.995) (0.967)

EM % -

MSE

(1.042) (0.168) (0.158)

EM % -

MSE

(1.576) (0.961) (0.942)

(-)

()

..... ..

جدول (٢) نتائج تقدير البيانات المفقودة لطرق التقدير عند حجم عينة (١٠٠)

	نسبة الفقدان	طريقة التقدير	طريقة التقدير					
1٠0	5%	طريقة EM المقترحة	قيمة المعلمة	0.302	1.037	0.61	0.209	0.967
			الانحراف المعياري	0.082	0.083	0.081	0.077	
		طريقة التعويض بالانحدار	قيمة المعلمة	0.309	1.007	0.601	0.189	0.995
			الانحراف المعياري	0.083	0.084	0.082	0.078	
		طريقة المتوسط غير الشرطي	قيمة المعلمة	0.317	0.994	0.594	0.177	1.026
			الانحراف المعياري	0.084	0.086	0.083	0.079	
	10%	طريقة EM المقترحة	قيمة المعلمة	0.6250	0.0730	0.3320	0.1070	0.158
			الانحراف المعياري	0.0220	0.0720	0.0750	0.0770	
		طريقة التعويض بالانحدار	قيمة المعلمة	0.6360	0.0810	0.3410	0.1160	0.168
			الانحراف المعياري	0.0250	0.0780	0.0760	0.0775	
		طريقة المتوسط غير الشرطي	قيمة المعلمة	0.6720	0.0830	0.3250	0.1240	1.042
			الانحراف المعياري	0.0290	0.0820	0.0861	0.0768	
15%	طريقة EM المقترحة	قيمة المعلمة	0.523	0.568	0.246	0.127	0.942	
		الانحراف المعياري	0.622	0.006	0.047	0.0757		
	طريقة التعويض بالانحدار	قيمة المعلمة	0.628	0.582	0.285	0.088	0.961	
		الانحراف المعياري	0.733	0.014	0.053	0.078		
	طريقة المتوسط غير الشرطي	قيمة المعلمة	0.736	0.812	0.254	0.092	1.576	
		الانحراف المعياري	0.963	0.123	0.099	0.096		

10-الاستنتاجات :

:

MSE

-



المصادر :

أولاً : المصادر العربية

" -

" ECME, ECM , EM

(James – stein) " -

" -

" -

" -

" -

ثانياً : المصادر الانكليزية .

3- C.Yang Yuan “ Multiple Imputation for Missing Data : concepts and
New Development

<http://support.sas.com/rnd/app/papers/multipleimputation.pdf>

4- Dempster .A.p & N.M Laird “Maximum Likelihood from Incomplete
Data Via the EM Algorithm . (1977) “ Journal of the Royal Statistical
Society

Vol.39.No.1<http://people.csail.mit.edu/jrennie/trg/papers/rubin-missing-76.pdf>



(-) () ..

-
- 5- D.paul .Allison (2001) “Missing Data “
<http://www.research.ibm.com/dar/papers/pdf/rc21783.pdf>
 - 6- E.Trivellore .Raghunathan & James M. Lepkowski “A Multivariate Technique for Multiply Imputing Missing Values Using a Sequence of Regression Models <http://www.statcan.gc.ca/ads-annonces/12-001-x/5857-eng.pdf>
 - 7- J.Skyler.Cranmer (2005)” Methods Exam Review “
<http://jmlr.csail.mit.edu/papers/volume8/saar-tsechansky07a/saar-tsechansky07a.pdf>
 - 8- Law.M Averill & W. David Kelton (2000) “ Simulation Modeling And Analysis “ third edition ,New York.
 - 9- Nittner “ The Additive Model with missing Values in the independent variable – theory and simulation”; http://epub.ub.uni-muenchen.de/1653/1/paper_272.pdf
 - 10- Roche .Alexis EM algorithm and variants : an information tutorial
http://webcache.googleusercontent.com/search?hl=ar&lr=&spell=1&q=cache:0tlu1vBml4gJ:ftp://ftp.cea.fr/pub/dsv/madic/publis/Roch_e_em.pdf+e+m+algorithm+and+variants&ct=clnk
 - 11- Sundberg. Rolf (1974) “ Maximum Likelihood Theory for Incomplete Data from an Exponential Family
<http://www2.math.su.se/~rolfs/Publikationer/SJS1974.pdf>

