

استخدام طرق معايير المعلومات وطرق تشخيص النموذج لاختيار أفضل نموذج انحدار خطي متعدد مع تطبيق على أطفال مرضى الثلاسيميا بالموصل

إيمان طارق فتحي

قسم الرياضيات / كلية التربية

جامعة الموصل

القبول

٢٠١٢ / ٠٣ / ٠٧

الاستلام

٢٠١١ / ١٠ / ١٧

ABSTRACT

In this paper we compute Akaike's Information Criteria (AIC), Bayesian Information Criteria (BIC) and Schwarz Bayesian Criteria (SBC) for all possible $(2^p - 1)$ for $p \leq 10$ independent variables.

The three criteria are measurement of the difference between a given model and the real model, and the model with smallest (AIC), (BIC) and (SBC) compared with other models is the best one.

This study is applied to data gathered from children infected with Thelasyimia disease, and the data are analyzed by using SAS to evaluate which of the models $(2^5 - 1 = 31)$ is the best model to determine the best subset of variables that minimize the information criteria among all the variables in the study. We compare the three criteria with model diagnostic like root mean squared error (RMSE), Mallow's C_p and adjusted R^2 .

الملخص:

تم في هذا البحث حساب معيار اكاكي للمعلومات (AIC)، معيار بيز للمعلومات (BIC) ومعيار شوارز - بيز للمعلومات (SBC) لجميع النماذج الممكنة $(2^p - 1)$ حيث $P \leq 10$ من المتغيرات المستقلة.

ان المعايير الثلاثة أعلاه هي مقياس للفرق بين النموذج المعطى والنموذج الحقيقي، وان النموذج الذي تكون فيه قيم المعايير الثلاثة صغيرة مقارنة بالنماذج الأخرى يعتبر هو النموذج الأفضل.

لقد قمنا بتطبيق الدراسة على بيانات جمعت لأطفال مصابين بمرض الثلاسيميا في مستشفى ابن الاثير بالموصل وتم تحليلها باستخدام البرنامج الإحصائي SAS لتقييم أي من النماذج (31 = 2⁵ - 1) هو النموذج الافضل وذلك بتحديد أفضل مجموعة من المتغيرات والتي تعمل على تصغير معايير المعلومات من بين جميع المتغيرات قيد الدراسة . لقد تم مقارنة المعايير الثلاثة مع معايير تشخيص النموذج مثل ل الجذر التربيعي لمتوسط مربعات الأخطاء (RMSE)، Mallow's C_p و R² المصححة.

١ - المقدمة :

ان النموذج الحقيقي نادرا ما يكون معلوما، لذلك فان طرق اختيار النموذج تكون ذات فائدة كبيرة جدا في تحليل الانحدار الخطي ، انظر Linhart , Draper and Smith (1981) and Zucchini (1986).

ان الانحدار الخطي متعدد المتغيرات هو احد الاساليب الاحصائية المستخدمة في كشف العلاقة بين المتغيرات ، اذ انه يستخدم ليجاد النموذج الخطي الافضل للمتغير المعتمد من مجموعة من المتغيرات المستقلة.

هناك عدد كبير من البحوث تناولت مسألة اختيار المتغيرات في نموذج الانحدار الخطي (انظر Hocking(1976) and Thompson(1978 a,b)).

ان البيانات المستخدمة في هذا البحث هي بيانات لخمس متغيرات مستقلة (p = 5) ، لذا فان هناك 31 مجموعة جزئية من النماذج (31 = 2⁵ - 1) وبالتالي فان هذه النماذج الجزئية سيكون هو النموذج الافضل للبيانات المدروسة والذي يتصف بان قيمة ال RMSE فيه قليلة وقيمة R² فيه عالية وتقترب من ١ .

واخيرا فان هناك طرقا اخرى لاختيار النموذج وذلك بتصغير مقياس المعلومات لمجموعة من النماذج البديلة (المدروسة) (انظر Sparks, Zucchini (1973), Mallows (1973), and Coutsourides (1985)).

الجانب النظري

٢ - طرق تشخيص النموذج :

قمنا في هذا البحث بدراسة ثلاثة طرق إحصائية لتحديد أفضل نموذج خطي وهي A₁ - الجذر التربيعي لمتوسط مربعات الأخطاء (RMSE).

B- معامل التحديد المعدل (\bar{R}^2).

C- احصاءة Mallow's C_p .

وسوف نقوم بذكر هذه الطرق بشئ من الإيجاز :

A- الجذر التربيعي لمتوسط مربعات الأخطاء (RMSE)

ان الجذر التربيعي لمتوسط مربعات الاخطاء هي دالة لمجموع مربعات الاخطاء (SSE) وعدد المشاهدات (n) وعدد المتغيرات المستقلة ($k \leq p+1$) حيث k تشمل على المعلمة الثابتة في النموذج (Intercept) ويمكن كتابتها كما يلي:

$$RMSE = \sqrt{\frac{SSE}{n-k}} \quad \dots\dots(1)$$

حيث يتم حساب RMSE لكل النماذج ، والنموذج الذي تكون فيه قيمة RMSE قليلة هو النموذج الأفضل . ان هذه الطريقة تعتمد على عدد المعلمات في النموذج (k) لذلك فان إضافة معلمات أخرى للنموذج سوف يقلل كل من البسط والمقام انظر (Beal (2007).

B- معامل التحديد المعدل (\bar{R}^2)

ان معامل التحديد R^2 ليست دائما هي الاحصاءة المقبولة في حساب ملائمة النماذج الجزئية للنموذج الأفضل، اذ ان إضافة متغيرات مستقلة جديدة الى معادلة الانحدار يؤدي الى رفع قيمة R^2 وذلك لثبات قيمة المقام وتغير قيمة البسط غير ان الاستمرار بإضافة المتغيرات المستقلة سيؤدي الى انخفاض درجات الحرية ($n-k$)، مما يتطلب إيجاد معامل التحديد المعدل (\bar{R}^2) والذي يمثل نسبة التغيرات للمتغير المعتمد والتي تفسر التغيرات الحاصلة في المتغيرات المستقلة، ويمكن كتابته على النحو التالي:

$$\text{Adjusted } R^2 = \bar{R}^2 = 1 - \frac{n-i}{n-k} \cdot (1 - R^2) \quad \dots\dots(2)$$

حيث ان n تمثل حجم العينة، ($i=1$) اذا كان هناك β_0 في النموذج و ($i=0$) اذا لم يكن هناك β_0 في النموذج، k تمثل عدد المعلمات في النموذج . ان النموذج الذي تكون فيه قيمة R^2 المعدلة عالية هي افضل النماذج المدروسة، مع ملاحظة ان المتغيرات المستقلة الجديدة المضافة الى النموذج غالبا ما تجعل ($\bar{R}^2 = 1$) في عدة نماذج وبالتالي فان تحديد النموذج الافضل في هذه الحالة سيكون معضلة يصعب حلها، انظر (Stauffer (2008).

C- احصاءة Mallow's C_p

وضع Mallows عام ١٩٧٣ احصاءة عرفت باحصاءة C_p Mallow's واساسها تقييم ملائمة المربعات الصغرى للنماذج بأخطاء طبيعية وذات تباين ثابت وكما يلي:

استخدام طرق معايير المعلومات وطرق تشخيص النموذج لاختيار أفضل نموذج انحدار خطي متعدد مع ...

$$C_p = P + (n - P) \cdot \frac{(\hat{\sigma}^2 - \hat{\sigma}_{full}^2)}{\hat{\sigma}^2} \quad \dots\dots(3)$$

اذ ان p تمثل عدد المتغيرات في النموذج ، وان النموذج الكامل هو النموذج الذي يأخذ جميع المتغيرات المستقلة بنظر الاعتبار.

$$\hat{\sigma}^2 = \frac{\sum (y_i - \hat{y}_i)^2}{n} = \frac{n - P - 1}{n} \cdot S_{y|x}^2$$

فالنموذج الذي تكون فيه قيمة C_p Mallow's صغيرة هو النموذج الافضل او النموذج الاكثر دقة. لاستخدام احصاءة C_p Mallow's سوف نقوم بحساب قيمة C_p لكل مجموعة جزئية من المتغيرات المستقلة ثم نقوم برسم النقاط (C_p, P) لكل نموذج مدروس، والنموذج الذي تكون فيه $C_p \cong P$ يعتبر هو افضل النماذج. انظر (Staufer (2008).

٣- معايير المعلومات Information Criteria

معيار المعلومات هو معيار لحسن المطابقة او عدم التأكدية لمدى قيم مجموعة من البيانات، وفي سياق الكلام عن الانحدار الخطي المتعدد فان معيار المعلومات يقيس الفرق بين النموذج المعطى والنموذج الحقيقي . سوف نقوم بالتعرف على ثلاثة أنواع من طرق معايير المعلومات وهي:

١- معيار اكاكي للمعلومات Akaike's Information Criteria

ان احصاءة نظرية المعلومات المهمة هي اكثر فاعلية في تقييم المقارنات بين النماذج ، ان معيار اكاكي للمعلومات قدم لاول مرة عام (١٩٧٣) من قبل العالم الرياضي الياباني هيتروتكو اكاكي وتعرف بانها مقياس لنظرية المعلومات لـ (كولبك-ليبيلر) بين النموذج المعطى والنموذج الحقيقي . ان معيار اكاكي للمعلومات (AIC) ومعيار اكاكي للمعلومات المصحح (AIC_c) للانحدار الخطي المتعدد يعطيان بالصيغتين التاليتين:

$$AIC = n \ln(S_{y|x}^2 \cdot \frac{n - p - 1}{n}) + 2K \quad \dots\dots(4)$$

$$= n \ln(S_{y|x}^2 \cdot \frac{n - K + 1}{n}) + 2K$$

$$AIC_c = n \ln(S_{y|x}^2 \cdot \frac{n - p - 1}{n}) + 2K + \frac{2K(K + 1)}{(n - K - 1)} \quad \dots\dots(5)$$

$$= n \ln(S_{y|x}^2 \cdot \frac{n - K + 1}{n}) + 2K + \frac{2K(K + 1)}{(n - K - 1)}$$

حيث p تمثل عدد المتغيرات المستقلة و $(K = p + 2)$ تمثل عدد المعلمات (من ضمنها β_0, σ). ان معيار اكاكي للمعلومات (AIC) هو تقريب لسلسلة تايلر الخطية لمعلومات كولباك

- لييلر ، بينما معيار اكاكي للمعلومات المصحح (AIC_C) هو تقريب لسلسلة تايلر غير الخطية لذلك فان (AIC_C) هي اكثر دقة من (AIC). ومن الجدير بالذكر انه اذا كان حجم العينة صغير وكانت ($\frac{n}{K} < 40$) فان النموذج الافضل هو النموذج الذي يملك اقل قيمة لـ (AIC_C) (انظر (Burnham and Anderson (2002)، كما نلاحظ با نه كلما ازداد عدد المتغيرات المستقلة فان النماذج تؤول الى ان تكون ملائمة للبيانات. وبصورة عامة فانه لاي نموذج احصائي احتمالي لمجموعة من بيانات عينة ما بدالة الامكان L فان (AIC) و (AIC_C) يعرفان باستخدام الانحراف $D = -2 \cdot \ln L$ وكما يلي:

$$AIC = D + 2 \cdot K$$

$$= -2 \cdot \ln L + 2 \cdot K$$

$$AIC_C = D + 2 \cdot K + 2 \cdot \frac{K(K+1)}{n-K-1}$$

$$AIC_C = -2 \cdot \ln L + 2 \cdot K + 2 \cdot \frac{K(K+1)}{n-K-1}$$

انظر (Stauffer (2008).

٢- معيار المعلومات البيزي Bayesian Information Criterion

ان (Sawa (1978) قدم معيار لاختيار النموذج والذي اشتق من تعديل بيز لمعيار (AIC). ان معيار بيز للمعلومات (BIC) هو دالة لعدد المشاهدات (n)، SSE وتباين الخطأ (σ^2) وعدد المتغيرات المستقلة $p \leq K+1$ وكما مبين بالمعادلة التالية:

$$BIC = n \ln\left(\frac{SSE}{n}\right) + 2(K+2)\left(\frac{n}{\sigma^2}\right) - 2\left(\frac{n\sigma^2}{SSE}\right)^2 \quad \dots\dots(6)$$

حيث ان $K = p+2$ وهي عدد المعلمات المقدرة في النموذج ، والنموذج الذي تكون فيه قيمة BIC هي الأصغر هو النموذج الأفضل من بين النماذج قيد الدراسة، انظر Beal (2007).

٣- معيار شوارز البيزي Schwarz's Bayesian Criterion

في اطار الكلام عن بيز، فان (Schwarz (1978) وصف معيار اخر سمي باسمه وهو معيار شوارز البيزي للمعلومات لاختيار النموذج الأفضل والذي يرمز له بـ SBC ويمكن كتابته كما يلي:

$$SBC = n \ln\left(\frac{SSE}{n}\right) + K \ln n \quad \dots\dots(7)$$

ان البسط في الحد الاول لهذا المعيار هو $(n \ln SSE)$ والذي يقل كلما ازداد عدد المتغيرات المستقلة، اما مقام الحد الأول فهو ثابت (بحجم عينة n) والحد الثاني يزداد مع ازدياد عدد المعلمات K في النموذج. اذا كانت $(n \geq 8)$ فان جزء المعيار SBC $(K \ln n)$ هي اكبر من جزء المعيار AIC $(2K)$ لذلك فان معيار SBC هو الافضل لمعظم النماذج قيد الدراسة، انظر (Kutner et al. (2004).

الجانب التطبيقي

تم اخذ البيانات من رسالة ماجستير للطالبة اسوان محمد طيب بعنوان (اختيار المتغيرات في انحدار الحرف) حيث قامت بجمع البيانات من مستشفى ابن الاثير التعليمي للولادة والاطفال في الموصل من مرضى مصابين بفقر دم البحر الأبيض المتوسط من نوع بيتا او ما يسمى بالثلاسيميا لـ (110) مشاهدات حيث ان المتغير المعتمد او متغير الاستجابة في البحث هو العمر من العظم مقاسا بالشهر كما واختير عدد من المتغيرات التي يعتقد انها تؤثر فيه وهي:

العمر الحقيقي (X_1): مقاسا بالشهر، ان الأشخاص عرضة للإصابة بالمرض منذ مقتبل العمر ويصل اعلى نسبة للإصابة في عمر (36 الى 80) شهرا وتتركز في عمر (80 الى 85) شهرا.

تضخم الكبد (X_2): مقاسا بالسنتيمتر، اذ يعاني المصاب بالثلاسيميا من تضخم الكبد نتيجة لترسب الحديد بكميات كبيرة في الكبد وحدوث خلل في وظائفه.

الخلايا الشبكية (X_3): في دم المصاب يلاحظ ان نسبة الخلايا الشبكية تكون عالية وعندها يعتبر المريض في وضع صحي سيئ، ويمكن اعتبار الخلايا الشبكية مقياس لفاعلية نخاع العظم على تكوين كريات دم جديدة.

الارومة الحمراء (X_4): وهي عبارة عن كريات دم حمراء جديدة التكوين وتكون غير ناضجة حاوية على نواة وتكون في نخاع العظم وتنتج لكثرة الحاجة اليها، فتكون خارج نخاع العظم في الدم للمصابين بالثلاسيميا وهي مؤشر على وجود المرض.

بداية نقل الدم للمريض (X_5): حسب العمر مقاسا بالشهر اذ يلاحظ ان معظم المصابين يبدأ إعطاؤهم وحدات الدم في الأشهر الأولى من العمر وذلك لكون ظهور المرض في مقتبل العمر والكشف عنه منذ بداية ظهوره بالفحص والتحليل لدم المصاب سهلا.

اما المتغير المعتمد (Y): هو العمر من العظم مقاسا بالشهر للمرضى المصابين بالثلاسيميا ،
ويعتمد مقياس العمر من العظم على الصور الشعاعية للمرضى المصابين اذ تبدو الصور
الشعاعية للمرضى مقارنة باقرانهم الاصحاء اقل من العمر الحقيقي.
وقد تم عرض البيانات في الجدول التالي [جدول (١)]:

جدول (١): البيانات قيد الدراسة

ت	العمر الحقيقي بالعظم	العمر الحقيقي بالشهر	تضخم الكبد	الخلايا الشبكية	ارومة حمراء	ت	العمر الحقيقي بالعظم	العمر الحقيقي بالشهر	تضخم الكبد	الخلايا الشبكية	ارومة حمراء	ت	العمر الحقيقي بالعظم	العمر الحقيقي بالشهر
ت	Y	X ₁	X ₂	X ₃	X ₄	X ₅	ت	Y	X ₁	X ₂	X ₃	X ₄	X ₅	ت
١	132	147	8	6.4	120	6	٥٦	12	31	0	1	5	10	١
٢	108	153	6	0.1	17	6	٥٧	15	26	4	5.4	2	10	٢
٣	48	114	4	7.6	1	12	٥٨	72	68	6	7.6	440	30	٣
٤	72	84	7	2	0	36	٥٩	112	176	4	9.2	240	12	٤
٥	72	96	6	1.1	1	7	٦٠	24	34	1	12.8	25	8	٥
٦	108	82	4	2.7	200	12	٦١	60	44	7	2.3	0	12	٦
٧	36	68	3	0.2	2	6	٦٢	30	54	2	0.4	0	7	٧
٨	7	24	0	7.4	5	7	٦٣	9	37	3	1	14	9	٨
٩	24	65	4.5	0.2	0	4	٦٤	120	129	6	13.3	72	3	٩
١٠	68	83	10.5	1.8	18	4	٦٥	96	103	5	0.1	500	36	١٠
١١	48	51	6	2.6	4	4	٦٦	96	99	6	13.6	4	6	١١
١٢	132	145	8	3	115	9	٦٧	72	79	10	6	4	4	١٢
١٣	60	89	4	2.2	115	12	٦٨	72	85	8	0.7	100	5	١٣
١٤	96	95	6.5	2.8	126	18	٦٩	60	76	3	0.9	2	4	١٤
١٥	48	70	5	0	0	3	٧٠	8	23	2	1.8	5	7	١٥
١٦	30	44	8	2.4	3	3	٧١	96	112	8	5	2	3	١٦
١٧	24	39	2	6.4	2	3	٧٢	24	46	2	2	0	3	١٧
١٨	60	75	5	0.8	0	11	٧٣	108	117	7	0	82	7	١٨
١٩	30	56	0	2.2	1	6	٧٤	36	73	5	5	3	12	١٩
٢٠	7	24	0	6.8	10	7	٧٥	76	44	6	2.2	3	6	٢٠
٢١	64	88	7	3	390	2	٧٦	36	40	3	11.6	3	5	٢١
٢٢	108	97	7	3.4	210	2	٧٧	24	45	0	10.2	0	9	٢٢
٢٣	9	20	3	17.8	35	9	٧٨	36	26	5	11	4	48	٢٣
٢٤	132	109	6	4.7	100	12	٧٩	12	17	5.5	0	1	4	٢٤
٢٥	48	45	5	0.1	18	4	٨٠	11	21	3	0	0	7	٢٥

استخدام طرق معايير المعلومات وطرق تشخيص النموذج لاختيار أفضل نموذج انحدار خطي متعدد مع...

3	440	0.2	4	148	168	٨١		1	0	2.4	6	18	33	٢٦
18	0	0.3	3	41	36	٨٢		6	170	0.3	5	109	84	٢٧
9	5	0.2	4	38	24	٨٣		6	3	0.2	2	77	72	٢٨
6	2	0	4	69	60	٨٤		6.5	432	0	7	60	72	٢٩
4	0	0.1	5	79	66	٨٥		12	8	2	6	46	30	٣٠
5	0	4.5	5	30	7	٨٦		8	14	3.6	5	74	84	٣١
6	22	0.2	0	72	60	٨٧		7	0	4	3.5	113	96	٣٢
14	3	0	6	76	78	٨٨		6	10	9.5	3	50	30	٣٣
6	2	0.4	6	65	60	٨٩		18	6	0	2.5	82	48	٣٤
7	12	3	4	62	61	٩٠		12	20	3.6	4	90	60	٣٥
24	40	3	5	74	63	٩١		48	10	4	10	50	48	٣٦
6	0	3.4	3	63	36	٩٢		40	38	3	5	42	36	٣٧
2	10	2.6	8	71	60	٩٣		12	10	1.7	5.5	54	30	٣٨
9	10	14	2.5	37	36	٩٤		7	300	3.3	3.5	157	144	٣٩
6	0	0.2	1	111	84	٩٥		7	350	0.4	5.5	175	180	٤٠
6	0	1	2	110	144	٩٦		18	10	2	4	67	60	٤١
3	36	1.3	8	89	48	٩٧		18	0	0.3	4	45	48	٤٢
3	3	0.4	0	40	30	٩٨		7	43	0.1	8	109	48	٤٣
24	0	0	6	161	120	٩٩		48	210	4.3	8	58	66	٤٤
6	1	7.5	0	26	9	١٠٠		14	120	19.2	3	147	132	٤٥
4	1	2.9	4	31.4593	18	١٠١		18	600	10.6	3.5	137	96	٤٦
12	16	0.3	3	77.1907	60	١٠٢		6	18	3.6	4.5	40	24	٤٧
5	5	7.7	8	116.1028	120	١٠٣		30	2	8.6	6	46	48	٤٨
5	112	8.8	4	12.3102	57	١٠٤		6	0	2.4	9	88	42	٤٩
5	9	2	6	42.7778	66	١٠٥		6	3	7.4	0	65	18	٥٠
7	0	1	4	70.2398	48	١٠٦		3	432	0	7	78	72	٥١
6	25	1	8.5	117.9954	120	١٠٧		6	8	2	0	159	156	٥٢
24	1	6.6	4.5	105.6019	78	١٠٨		5	14	3.6	11	75	72	٥٣
6	1	1	4.5	166.1943	156	١٠٩		48	10	9.5	4	94	96	٥٤
4	7	6	2.5	77.3648	78	١١٠		8	172	11.4	3	33	36	٥٥

ان نموذج الانحدار الخطي المتعدد الذي يعتمد على البيانات في جدول
(١) أعلاه يمكن التعبير عنه بالشكل التالي:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \beta_4 X_{i4} + \beta_5 X_{i5} + \varepsilon_i$$

حيث ان:

β_j : تمثل معاملات غير معلومة (معاملات الانحدار)، $j=1,2,3,4,5$

وكانت نتائج التحليل الإحصائي (باستخدام البرنامج الإحصائي SAS) لهذه البيانات كما يلي:

جدول (٢) : نتائج التحليل الإحصائي للبيانات لاختيار أفضل نموذج انحدار خطي متعدد

MODEL	Dependent Variable	Intercept	x1	x2	x3	x4	x5	Independent Variables in Model	Parameters in Model
MODEL1	y	-7.22171	0.833829	1.304645		0.032337		3	4
MODEL1	y	-7.92596	0.835799	1.265563		0.031453	0.075453	4	5
MODEL1	y	-7.99977	0.834684	1.342419	0.149718	0.03194		4	5
MODEL1	y	-8.48646	0.836288	1.300602	0.122331	0.031222	0.067393	5	6
MODEL1	y	-2.61931	0.851666			0.034198		2	3
MODEL1	y	-3.87832	0.853814			0.032792	0.112806	3	4
MODEL1	y	-2.75507	0.851955		0.031521	0.034126		3	4
MODEL1	y	-8.35286	0.868901	1.408538				2	3
MODEL1	y	-3.84779	0.853746		-0.00833	0.032805	0.113285	4	5
MODEL1	y	-9.44756	0.870545	1.340372			0.122679	3	4
MODEL1	y	-9.46328	0.869517	1.461572	0.217565			3	4
MODEL1	y	-10.2107	0.870869	1.388287	0.169946		0.111001	4	5
MODEL1	y	-3.43169	0.890437					1	2
MODEL1	y	-5.21866	0.891244				0.164476	2	3
MODEL1	y	-3.8279	0.891049		0.093178			2	3
MODEL1	y	-5.33601	0.891448		0.032588		0.162523	3	4
MODEL1	y	38.61291		3.79419		0.120022		2	3
MODEL1	y	40.17681		3.881315		0.1218	-0.19801	3	4
MODEL1	y	39.69234		3.734643	-0.22107	0.120477		3	4
MODEL1	y	40.79698		3.838718	-0.1425	0.122007	-0.18844	4	5
MODEL1	y	55.54029				0.131352		1	2
MODEL1	y	57.69239			-0.58465	0.132084		2	3
MODEL1	y	56.49715				0.132351	-0.09704	2	3
MODEL1	y	58.23513			-0.56153	0.132711	-0.06367	3	4
MODEL1	y	41.95079		4.668127				1	2
MODEL1	y	42.27935		4.688507			-0.04032	2	3
MODEL1	y	41.93931		4.668718	0.002324			2	3
MODEL1	y	42.19026		4.694394	0.020352		-0.04173	3	4
MODEL1	y	65.27434			-0.43584			1	2
MODEL1	y	62.57953					0.100695	1	2
MODEL1	y	64.09329			-0.48448		0.129944	2	3

Degree of Freedom	SSE	R - Square	RMSE	Adjusted R - Square	CP	AIC	BIC	SBC
106	34755.75	0.798435	18.10758	0.79273	2.271609	641.1182	643.5516	651.9201
105	34691.76	0.798806	18.17684	0.791142	4.079626	642.9155	645.4778	656.4179
105	34714.9	0.798672	18.18291	0.791002	4.149057	642.9888	645.5443	656.4912
104	34665.21	0.79896	18.25704	0.789295	6	644.8313	647.5169	661.0342
107	35862.82	0.792015	18.30756	0.788127	3.592982	642.5674	644.7006	650.6688
106	35717.57	0.792857	18.35642	0.786994	5.157207	644.1209	646.3318	654.9229
106	35860.97	0.792025	18.39323	0.786139	5.587426	644.5617	646.7401	655.3636
107	36287.58	0.789551	18.41565	0.785618	4.867287	643.8625	645.9247	651.964
105	35717.45	0.792858	18.44359	0.784966	7.156827	646.1206	648.3845	659.623

استخدام طرق معايير المعلومات وطرق تشخيص النموذج لاختيار أفضل نموذج انحدار خطي متعدد مع...

106	36115.38	0.79055	18.45836	0.784622	6.350677	645.3393	647.4606	656.1412
106	36200.83	0.790054	18.48018	0.784112	6.607027	645.5993	647.7015	656.4012
105	36064.01	0.790848	18.53286	0.78288	8.196558	647.1827	649.3489	660.6851
108	37583.94	0.782033	18.65474	0.780015	6.756529	645.7237	647.6247	651.1246
107	37268.96	0.78386	18.66301	0.77982	7.811564	646.7979	648.7002	654.8994
107	37567.69	0.782127	18.73766	0.778055	8.707781	647.6761	649.5308	655.7776
106	37267.02	0.783871	18.75035	0.777754	9.805734	648.7922	650.6625	659.5941
107	133505.8	0.225737	35.32305	0.211265	296.5341	787.1562	783.7517	795.2576
106	133061.3	0.228315	35.43014	0.206475	297.2005	788.7893	783.9441	799.5912
106	133416.5	0.226255	35.4774	0.204356	298.2662	789.0826	784.2293	799.8845
105	133025.2	0.228524	35.59363	0.199134	299.0923	790.7595	784.4663	804.2619
108	143267.6	0.169123	36.42186	0.16143	323.8209	792.9188	790.8352	798.3198
107	142626.3	0.172842	36.50967	0.157382	323.8969	794.4253	790.8639	802.5268
107	143159.6	0.169749	36.57787	0.154231	325.4969	794.8359	791.2658	802.9373
106	142580.8	0.173106	36.67563	0.149704	325.7604	796.3902	791.3438	807.1922
108	157324.2	0.087602	38.16682	0.079154	365.9925	803.2142	800.97	808.6152
107	157305.6	0.08771	38.34248	0.070658	367.9366	805.2012	801.4233	813.3026
107	157324.2	0.087602	38.34475	0.070548	367.9925	805.2142	801.4361	813.3157
106	157304.9	0.087715	38.52283	0.061895	369.9343	807.2007	801.889	818.0026
108	172072.6	0.00207	39.91572	-0.00717	410.2393	813.071	810.6849	818.472
108	172311.4	0.000685	39.9434	-0.00857	410.9557	813.2236	810.8353	818.6245
107	171880.3	0.003185	40.0794	-0.01545	411.6626	814.9481	810.9902	823.0495

جدول (٣) : نسبة اختيار النموذج الحقيقي من بين جميع نماذج الانحدار الممكنة لمعيار AIC، BIC

ومعيار C_p

The real model	AIC	BIC	C_p
Model(1) X_1, X_2, X_4	13 %	13 %	13 %
Model(2) X_1, X_4	26 %	26 %	26 %
Model(3) X_1, X_2, X_4, X_5	6 %	6 %	6 %
AVERAGE	15 %	15 %	15 %

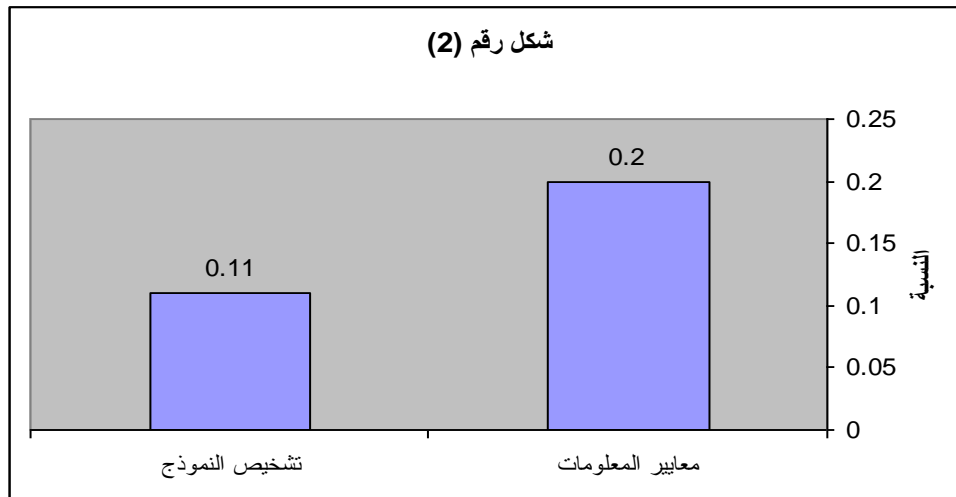
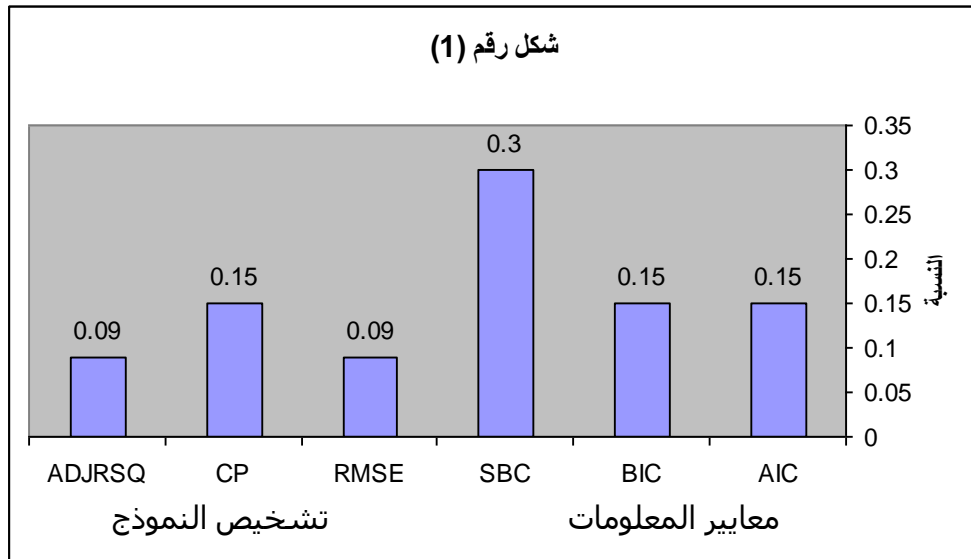
جدول (٤) : نسبة اختيار النموذج الحقيقي من بين جميع نماذج الانحدار الممكنة لمعيار \bar{R}^2 و RMSE

The real model	RMSE	\bar{R}^2
Model(1) X_1, X_2, X_4	13 %	13 %
Model(2) X_1, X_2, X_4, X_5	6 %	6 %
Model(3) X_1, X_2, X_3, X_4	6 %	6 %
AVERAGE	8.33 %	8.33 %

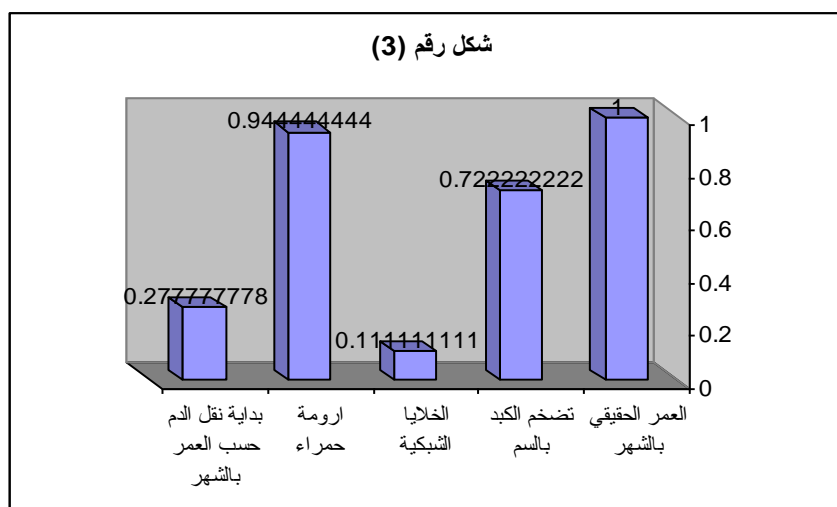
جدول (٥) : نسبة اختيار النموذج الحقيقي من بين جميع نماذج الانحدار الممكنة لمعيار SBC

The real model	SBC
Model(1) X_1, X_4	26 %
Model(2) X_1	51 %
Model(3) X_1, X_2, X_4	13 %
AVERAGE	30 %

نلاحظ من الجداول (٣)، (٤)، و (٥) أعلاه بان اعلى نسبة لاختيار النموذج الأفضل هو استخدام معيار SBC وبالتالي فهو يعتبر من أفضل المعايير المدروسة في هذا البحث ، تليه معايير AIC ، BIC ومعيار C_p .
وكما موضح بالاشكال التالية:



شكل رقم (١) و (٢) هما مدرجان يمثلان نسبة كل من طرق معايير المعلومات وطرق تشخيص النموذج لاختيار النموذج الحقيقي (الأفضل) من بين ٣١ نموذج



شكل رقم (3): يمثل نسبة كل متغير في تأثيره على العمر الحقيقي بالعظم ولجميع المعايير المدروسة وفقا للنموذج الأفضل

الأ

١. من خلال التحليل الإحصائي للبيانات نلاحظ بان هناك ثلاثة عوامل تؤثر بشكل مباشر على العمر الحقيقي بالعظم والاصابة بالمرض وهي (العمر الحقيقي بالشهر ، تضخم الكبد والارومة حمراء) مع التركيز على عاملين مهمين نلاحظ ظهورهما في جميع المعايير المحسوبة وهما (العمر الحقيقي بالشهر والارومة الحمراء).
٢. نلاحظ من الأشكال (١) و (٢) بان استخدام طرق معايير المعلومات (AIC,BIC,SBC) لاختيار أفضل نموذج انحدار خطي متعدد تعتبر أفضل من طرق تشخيص النموذج (ADJRSQ, RMSE, C_p).
٣. إن معيار (SBC) يعتبر أفضل المعايير المدروسة حيث حصل على نسبة (0.30) في اختيار النموذج الأفضل للبيانات يليه ثلاثة معايير (AIC,BIC, C_p) بنسب متساوية وهي (١٥ %) في اختيار النموذج الأفضل.
٤. نوصي باستخدام معايير أخرى مثل معيار (New Information Criteria) NIC ومقارنته مع المعايير التي تطرقنا إليها في بحثنا هذا.

المصادر :

- 1) Beal, D.J., (2007). Information Criteria Methods in SAS for Multiple Linear Regression Models. Proceedings of the Fifteenth Annual Conference of the South East SAS Users Group, Hilton Head, SC.
- 2) Burnham, K.P., and Anderson, D.R. (2002). Model selection and multimodel inference: A practical information-theoretic approach, 2nd-ed. Springer-Verlag, New York.
- 3) Draper, N.R., and Smith, H. (1981). Applied Regression Analysis. 2nd-ed. Wiley, New York.
- 4) Hocking, R.R., (1976). The analysis and selection of variables in linear regression. *Biometrics* 32, pp. 1-49.
- 5) Kundu, D., and Murali, G., (1996). Model selection in linear regression. *Computational Statist. & Data Analysis* 22, pp. 461-469.
- 6) Kutner, M.H., Nachtsheim, C.J., Neter, J., and Li, W. (2004). Applied linear statistical models. 4th-ed. McGraw Hill, Chicago, IL.
- 7) Linhart, H., and Zucchini, W. (1986). Model Selection. Wiley, New York.
- 8) Mallows, C. L. (1973). "Some Comments on Cp", *Technometric*, Vol. 15, No.4, pp 661-675.
- 9) Sparks, R. S.; Zucchini, W. and Coutsourides, D. (1985). "On Variable Selection in Multivariate Regression", *Commun Statistis-Theor. Meth.*, 14(7), 1569-1587.
- 10) Stauffer, H.B., (2008). Contemporary Bayesian and Frequentist Statistical research methods for natural resource scientists. Wiley, New York.
- 11) Thompson, M.L., (1978a). Selection of variables in multiple regression. Part I, *Int. Statist. Rev.* 46, pp. 1-19.
- 12) Thompson, M.L., (1978b). Selection of variables in multiple regression. Part II, *Int. Statist. Rev.* 46, pp. 126-146.