مقارنة بين المكونات الرئيسة وطريقة المربعات الصغرى الجزئية لمعالجة مشكلة التعدد الخطي مع التطبيق على معمل السمنت

فرح عبد الغني يونس الصالح

قسم الاحصاء ، كلية علوم الحاسوب والرياضيات ، جامعة الموصل ، الموصل ، العراق (تاريخ الاستلام: ١٩ / ١٠١) ٢٠١٠)

الملخص:

تعتبر مشكلة التعدد الخطي مشكلة خطيرة في مسائل تحليل الانحدار لذلك تم تسليط الضوء على هذه المشكلة بشكل اكثر دقة وايجاد الحلول لها ، وفي هذا البحث تم مقارنة المكونات الرئيسية pc ومقدرات المربعات الصغرى الجزئية pls والتي تعتبر احدى طرق معالجة مشكلة التعدد الخطي، لاختبار افضل طريقة بالاعتماد على معيار اقل مربعات الخطأ MSE كمعيار للمقارنة بين المقدرات

الجانب النظري

ان تحليل الأنحدار من الطرق الاحصائية الواسعة الأستخدام وان الهدف الأساسي منها هو السعى الدائم لمعرفة وتفسير العلاقات المختلفة بين الظواهر وذلك بتحديد العلاقات بين المتغيرات ، والذي يوضح العلاقة بين متغير واحد يدعى بالمتغير التابع (متغير الاستجابة (Response variable) ومتغير او متغيرات أخرى يعتقد أنها تسبب التباين في المتغير التابع وتسمى بالمتغيرات التوضيحية (التنبؤية) (predicted variables) كذلك يصف العلاقة بين المتغير المعتمد والمتغيرات التوضحية على هيئة نموذج رياضي وان دقة هذا النموذج الرياضي تعتمد على فروض التحليل حيث ان بعض هذه الفروض مرتبط بالعلاقة الدالية والبعض الاخر مرتبط بالمتغير العشوائي ، وهذه الفروض عدة واهمها مشكلة التعدد الخطى حيث أكددت العديد من البحوث والدراسات الاهتمام الكبير بمشكلة التعدد الخطى وايجاد الحل المناسب لها واهمها المكونات الرئيسية التي تعود الى العام ١٩٠١ عندما اقترحها karl person بوصفها وسيلة لحل مشكلة تعدد العلاقة الخطية واعتمدها طريقة استكشافية يمكن الاستفادة منها للتوصيل الى تفسير وفهم العلاقة المتداخلة بين المتغيرات. ويعرف تحليل المكونات الرئيسية بانه اسلوب يهدف الى ايجاد عوامل factors اوتوليفات خطية تسمى بالمكونات الرئيسية قليلة مشتقة من المتغيرات الاصلية لتحل محلها بحيث تكون مؤهلة لتفسير معظم التباين الكلى للقيم الاصلية وتكون هذه المكونات الرئيسية متعامدة أي لايوجد ارتباط فيما بينها ويمكن كتابة المكونات الرئيسية كالاتى:

$$pc_i = a_{1i} x_1 + a_{2i} x_2 + \dots + a_{mn} x_m$$

$$pc_i = \sum_{j=1}^{m} a_{ji} x_j$$
 $(i, j = 1, 2...m)$

حيث ان

Pc_i: تمثل المكون الرئيسي

يمثل معامل j في المكون الرئيسي i الذي يمثل قيم المتجهات المميزة a_{ij} المرافقة المستخدمة وباستخدام الملوب المصفوفات فان i i ول المصفوفة i في الصيغة اعلاه اعمدتها تمثل المتجهات المميزة المرافقة المصفوفة المستخدمة والمرتبقة وفقاً للمصفوفة المستخدمة والمرتبقة وفقاً للمصفوفة المستخدمة والمرتبقة وفقاً المقال المتجهات المقال المتحدمة والمرتبقة وفقاً المتحدمة والمتحدمة والمرتبقة وفقاً المتحدمة والمتحدمة والمتحدمة والمتحدمة والمتحدمة والمتحدمة وفقاً وفقاً المتحدمة والمتحدمة والمتحد

المميزة $A_1 > A_2 > \dots$ المحونة المصغوفة $A_1 > A_2 > \dots$ المكونات الرئيسية وتعتمد فكرة احتساب المكونات الرئيسية على خصائص الجذور المميزة ومايرافقها من المتجهات المميزة لمصفوفة الارتباط او مصفوفة التباين والتباين المشترك المتغيرات الاصلية $X_1 > X_2 > X_3$ وحسب وحدات قياسها من حيث كونها متشابهة او مختلفة . حيث ان (s) تمثل جميع المتغيرات المدروسة من $X_1 > X_2 > X_3 > X_4 > X_4 > X_5 > X_5 > X_5 > X_5 > X_6 > X$

$$v = \begin{bmatrix} v_{11} & v_{12} & \cdots & v_{1m} \\ v_{21} & v_{22} & \cdots & v_{2m} \\ \vdots & \vdots & & \vdots \\ v_{m1} & v_{m2} & \cdots & v_{mm} \end{bmatrix}$$

ولايجاد الجذور المميزة λ نقوم بطرح من القيم القطرية للمصفوفة ν ثم نجعل قيمة محدثها تساوي صفر فنحصل على المعادلة المميزة للمصفوفة ν أي

$$|v - \lambda I| = 0$$

 λ وشكل المعادلة المميزة هو متعدد الحدود (polynomial) في λ من الدرجة m أي

 $\lambda^m + C_{m-1}\lambda^{m-1} + \dots + C_1\lambda + C_0 = 0$ وعند حـل هـنه المعادلـة سنحصـل علــى مـن الجـنور وعنـد حـل هـنه المعادلـة سنحصـل علــى $\lambda_1\lambda_2\dots \lambda_m$ هـ و اكبر قيمـة ويليـه λ_2 وهكذا $\lambda_1\lambda_2\dots \lambda_m$ وان لكـل جـنر مميز خـاص يقابلـه ولايجـاد هذه المتجهات تستخدم المعادلـة

$$(v - \lambda I)a = 0$$

ونختار a_{ji} بحیث تساوی قیم المتجة الممیز القیاسی أی تکون a'a=1 وعندئذ تدخل a_{ji} کمعادلة المتغیرات a'a=1 الرئیسیة اما عندما یکون المتغیرات x'_{s} وحدات قیاس مختلفة یستحسن تحویل x الی متغیرات قیاسیة لها وسط حسابی = صفر

وتباین = قبل إیجاد a_{ii} یتم اختیار المکونات الرئیسیة المؤثرة تأثير معنويا باختبار النسبة التجميعية للتباين المفسر لكل مكون كما ذكر (الراوي -١٩٨٧) ان عدد المكونات الرئيسية المختارة يكون بعدد الجذور المميزة الاكبر من واحد $(1 \ \langle \ \mathcal{R})$ اما بالنسبة لطريقة المربعات الصغرى الجزئية (pls) فقد ظهرت في العديد من الدراسات فيها الكيمياء والعلوم الاجتماعية ويعد الباحث (wold;1966) اول من وجد هذه الطريقة واصبحت شائعة الاستخدام في الطب وخاصة في المعالجات السريرية حيث قلة عدد المشاهدات (عدد المرضى) مع كثرة عدد المتغيرات (x'_s) المتمثلة بأغراض المرض مع عدد المتغيرات المعتمدة (y'_{c}) والمتمثلة بالمستوى الصحى للمريض مع تحسن حالته الصحية ونتائج اخرى كما طورت هذه الطريقة من قبل (Friedman;1993) . تعتمد طريقة انحدار المربعات الصغرى الجزئية على مصفوفة التباين والتباين المشترك (cov(x,y) اذ ان هذه الطريقة تسمح بتحديد العوامل بالمتغيرات الصماء (latent (variables والتي بدورها تكون افضل نموذج للمتغير المعتمد ويمكن وصفها بالشكل التالى : ليكن لدينا مدور مصفوفة $({oldsymbol y}_{{oldsymbol s}}')$ المتغيرات التنبؤية (ر) مضرورية بمتجه عشوائي من U أي

$$X'U = \begin{bmatrix} \sum x_{i1} U_i \\ \sum x_{i2} U_i \\ \vdots \\ \sum x_{ip} U_i \end{bmatrix}$$

ومن خلالها يمكن الحصول على مصفوفة اله β والتي توضح نوع وشكل العلاقة بين كل متغير معتمد مع كل متغير تنبؤكا لاتي

	y_1 y_2 y_j
x ₁ x ₂	$oldsymbol{eta}_{11}$ $oldsymbol{eta}_{12}$ $oldsymbol{eta}_{1j}$
:	eta_{21} eta_{22} eta_{2j}
x _k	: $oldsymbol{eta}_{k1}$ $oldsymbol{eta}_{k2}$ $oldsymbol{eta}_{kj}$

ان طريقة المربعات الصغرى الجزئية عبارة عن تركيبة خطية للمربعات الصغرى لمصفوفة الارتباط والتباين المشترك بين المتغيرات التتبؤية والمتغيرات المعتمدة ، والجزء الذي تعتمد عليه طريقة المربعات الصغرى الجزئية في مصفوفة الارتباط والتباين هو جزء التقاطع الحريقة المرتباط والتباين هو جزء التقاطع Cross Block أي الارتباطات بين المتغيرات التتبؤية والمتغيرات المعتمدة . كما تقدم هذه الطريقة عوامل Factor scores على شكل مجموعات خطية بين المتغيرات التبؤية الاصلية ، بذلك سوف لايكون هناك ارتباطات بين عوامل المتغيرات المستخدمة في النموذج الانحدار التنبؤي يمكن توضيح عمل طريقة المربعات الصغرى الجزئية

T من حيث حساب معلمات الانحدار ومصفوفة المتغيرات الصماء P وكخطوة اولية يتم تحويل مصفوفة المتغيرات التنبؤية X ومصفوفة المتغيرات المعتمدة Y الى الصيغة القياسية بحيث تصبح E,F على التوالي وخطوات هذه الخوارزمية كما يلي (Abdi,2003) .

۱ - اخذ قيم ابتدائية عشوائية للمتجه U

W=E'U ايجاد المتجه W من خلال العلاقة W

T ايجاد المتجة t وهو احد الاعمدة المصفوفة t للمتغيرات المستخلصة من المصفوفة t أي ان t_{old} =EW

 $\frac{1}{\|t\|}$ بضربه ب normalized خحویل المتجه $\frac{1}{\|t\|}$

$$t_{new} = \frac{1}{\|t\|} t_{old}$$

$$\|t\| = \sqrt{t_1^2 + t_2^2 + t_3^2 + \dots + t_n^2}$$

 $\mathbf{c} = \mathbf{F't}$. Y المصفوفة \mathbf{c} المتجه الموزون \mathbf{c} للمتغيرات \mathbf{U} المتغيرات \mathbf{U} المتخلصة للمصفوفة \mathbf{U} المتخلصة للمصفوفة \mathbf{v} المتخلصة المصفوفة \mathbf{v} المتخلصة المصفوفة \mathbf{v} المتخلصة المصفوفة \mathbf{v}

$$U_{old} = FC$$

 $\frac{1}{\parallel U \parallel}$ بضریه ب normalized الی صیغهٔ U المتجه ۷–تحویل المتجه ال

$$egin{aligned} U_{NEW} &= rac{1}{\|U\|} \cdot U_{old} \ &\|U\| &= \sqrt{U_1^2 + U_2^2 + U_3^2 + \dots + U_n^2} \end{aligned}$$

 $V_{\rm new}$ في الخطوة $V_{\rm new}$ في الخطوة $V_{\rm new}$ في الخطوة (7) فإذا كانت القيم متقاربة يتم حساب قيم $V_{\rm new}$ في حالة عدم التقارب يتم حساب قيمة المحمل $V_{\rm new}$ المصفوفة $V_{\rm new}$ من خلال المعادلة التالية .

$$P = E't$$

t بعد ذلك يتم حساب المصفوفتين الجديدين E_1,F_1 من طرح المتجه E_1,F_1 من كل من المصفوفتين E_1,F_2 وكما يلى :-

$$E_1=E-tp'$$
 $F_1=F-btc'$
 $P=E't$ آبالمعادلة C قيمة P بالمعادلة $C=F't$ قيمة $E_1=E-t(E't)'$
 $E=F-ttE$
 $E=F-bt(F't)'$
 $E=F-btt'F$

بعد ذلك تعاد الخوارزمية من جديد بنفس الخطوات ولكن باستخدام المصغوفتين الجديدتين F_{1},E_{1} مع العلم ان في كل عملية تكرار يتم مقارنة قيم المتجه (t) للدورة الاخيرة مع قيم المتجه (t) للدورة السابقة وتستمر عملية التكرار الى ان يتم الحصول على جميع اعمدة المصغوفة T (والتي عددها مساو الى رتبة المصغوفة T) ومن ثم يتم ايجاد مصغوفة المعلمات \mathcal{A} وكما يلى

 $\beta = E'U(TEEU)^{-1}T'F$

الجانب التطبيقي:

ان البيانات في هذا البحث تم الحصول عليها من كتاب DRAPER) (البيانات في هذا البحث تم الحصول عليها من كتاب ASMTTH, 1966) أولية لانتاج مادة السمنت والتي تمثل متغيرات مستقلة يحتوي كل متغير على (١٣) مشاهدة ومتغير معتمد y كمايلى:

$$y_i = B_0 + B_1 x_1 + B_2 x_2 + B_3 x_3 + B_4 x_4 + ei$$

: وكما موضحة بالجدول الآتى

جدول (١)

X ₁ : 3 Ca O . AL ₂ O ₃	الومينات الكالسيوم الثلاثي
X ₂ : 3 Ca O . S _i O ₂	سليكات الكالسيوم الثلاثي
X ₃ : 4 Ca O . AL ₂ O ₃ .Fe ₂ O ₃	الوهينات الكالسيوم الرباعي الحديد يكسي
X ₄ : 2 Ca O . S _i O ₂	سليكات الكالسيوم الثنائي

Y : الحرارة المنبعثة (سعرة / غم) بالنسبة المئوية لوزن مادة الذي يصنع منها مادة السمنت e_i : الخطأ العشوائي ويكون ذي توزيع طبيعي بوسط (٥) وتباين 2 اجري تحليل الانحدار بصورة عامة بين المتغير المعتمد y مع بقية المتغيرات المؤثرة الاربعة بأستخدام البرمجة الجاهزة v_{15} المتائج التحليل بالشكل التالي : وكانت نتائج التحليل بالشكل التالي :

$$y_i = 74.0 + 1.41x_1 + 0.390x_2 - 0.031x_3 - 0.266x_4$$

T والمختبر الاحصائي VIF والمختبر الاحصائي ودول (٢) والمختبر

predictor	coef	sEcoef	T	P	VIF
constant	73.97	72.39	1.02	0.337	-
X_1	1.4111	0.7695	1.83	0.104	38.5
X_2	0.3896	0.7478	0.52	0.616	254.4
X ₃	-0.0314	0.7797	-0.04	0.969	46.9
X_4	-0.2656	0.7326	-0.36	0.726	282.5

S=2.527 R-sq=98.1 % R-sq (adj)=97.2%

وللتأكيد على وجود تداخل خطي بين المتغيرات التنبؤية فقد تم استخدام بعض اساليب للكشف عنها منها:

۱-۱ استخدام عامل تضخم التباین Variance Inflation (VIF)

من ملاحظة قيم عوامل تضخم التباين في الجدول اعلاه نجد انها مرتفعة جدا ولقد (marquardt , 1970) على وجد تعدد العلاقة الخطية بين المتغيرات التنبؤية في حالة كون قيم العوامل تضم التباين VIF اكبر من (٤) او (١٠) ومن خلال التحليل السابق يلاحظ ان هذه القيم جداً مرتفعة وهذا يدل على وجدود تداخل خطي كبير بين المتغيرات التنبؤية .

|x'x| ايجاد محدد المصفوفة ا

اقترح هذا المعيار من قبل (Mason & Webster , 1975) اذ ينص على انه اذا كان $|x'|_{|x'|} = |x'|_{|x'|}$ فهذا يدل على وجود تداخل تام بين المتغيرات التتبؤية اما اذا كان $|x'|_{|x'|} = |x'|_{|x'|}$ بين الصفر المتغيرات التتبؤية فيما بينها . اما اذ كانت قيمة $|x'|_{|x'|} = |x'|_{|x'|}$ بين الصفر والواحد فهذا يدل على وجود تداخل خطي شبه تام ومن خلال القيم الذاتية للمصفوفة (x'x) تم حساب محدد المصفوفة (x'x) والذي كان يساوى

$$|x'x| = \prod_{i=1}^{p} Lj = 0.000014$$

وهذه القيم صغيرة جداً وقريبة من الصفر وهذا دليل اخر على وجود تداخل خطى كبير بين المتغيرات التنبؤية .

٣-١ استخدام عناصر مصفوفة الارتباط الواقعة خارج القطر:

اقترح هذا الاسلوب من قبل (Gunst & Mason,1980) حيثُ من خلال مصفوفة الارتباط (x'x) والموضحة ادناه

$$(x'x) = \begin{bmatrix} 1.000 & 0.229 & -0.824 & -0.245 \\ 0.229 & 1.000 & -0.139 & -0.973 \\ -0.824 & -0.139 & 1.000 & 0.030 \\ -0.245 & -0.973 & 0.030 & 1.000 \end{bmatrix}$$

نلاحظ ان هنالك ارتباط تام بين جميع المتغيرات مما يدل وجود تداخل خطي كبير بين هذه المتغيرات والذي يؤدي الى عدم دقة نتائج التحليل . كذلك نلاحظ دليل اخر على التداخل بين المتغيرات التنبؤية هو قيمة \mathbb{R}^2 أي تمثل (معامل التحديد) حيث ظهرت كبيرة جداً

٢ - طريقة المكونات الرئيسة

لأجل معالجة مشكلة التعدد الخطي تم استخدام طريقة المكونات الرئيسية بالاعتماد على مصفوفة التباين والتباين المشترك لحساب معاملات المكونات الرئيسية aij كالاتى:

 $v = \begin{bmatrix} 34.603 & 20.923 & -31.051 & -24.167 \\ 20.923 & 242.141 & -13.878 & -253.417 \\ -31.051 & -13.878 & 41.026 & 3.167 \\ -241.167 & -253.417 & 3.167 & 280.167 \end{bmatrix}$

ثم نجد الجذور المميزة وهي:

$$\lambda_1 = 517.80$$
 $\lambda_2 = 67.50$
 $\lambda_3 = 12.41$
 $\lambda_4 = 0.24$

ومن الجذور المميزة نحسب مصفوفة المتجهات المميزة:

	-0.068	0.646	-0.567	0.506	
	-0.679	0.020	0.544	0.493	
	0.029	-0.755	-0.404	0.516	
	0.731	0.108	0.468	0.484	
1 =	= -0.068x	$\frac{1}{1} - 0.679$	$x_2 + 0.021$	$x_3 + 0.73$	$1x_4$
2 =	$= 0.646x_1$	$+0.020x_2$	-0.755x	$\frac{1}{3} + 0.108$	x_4

 $pc_3 = -0.567x_1 + 0.544x_2 - 0.404x_3 + 0.468x_4$ $pc_4 = 0.506x_1 + 0.493x_2 + 0.516x_3 + 0.484x_4$

وتكون المكونات الرئيسية هي:

اذ تبين من نتائج جدول (*) رفض فرضية العدم وقبول الفرضية البديلة عند مستوى معنوية (α) * وهذا يعني ان هناك اختلافات معنوية بين المتغيرات التوضيحية على المتغير المعتمد ولايوجد مشكلة تعدد العلاقة الخطية وكذلك الحالة بالنسبة الى نتائج جدول (*) رفض فرضية العدم وقبول الفرضية البديلة عند مستوى معنوية (*) * الله لايوجد مشكلة تعدد العلاقة الخطية .

الاستنناجات

لقد تبين اهمية دراسة مشكلة تعدد العلاقة الخطية وتأثيرها على نتائج تحليل الانحدار

١ – تـم اختبار فـروض التحليـل للبيانـات واستتجنا مـن خـلال الاختبارات بأن البيانات تعاني من مشكلة التعدد الخطي حيث تم كشفه وفق اختبار تصخم التباين (VIF) ومصفوفة الارتباط حيث ظهرت ارتباط تام بين المتغيرات وكذلك عن طريق معامل الارتباط حيث ظهرت قيمة كبيرة جداً

 $Y - \dot{e}_{\omega}$ بحثنا هذا تم معالجة مشكلة تعدد العلاقة الخطية بطريقتين المكونات الرئيسية PC المربعات الصغرى الجزئية (PLS) حيث تبين من خلال التحليل الاحصائي Mse متوسط مربعات الخطأ ان طريقة المكونات الرئيسية اكثر كفاءة من طريقة المربعات الصغرى الجزئية تمثلك اقل قيمة لـ Mse.

F - كما تبين من خلال التحليلات الاحصائية والمقارنة بين قيمة F المحسوية مع F الجدولية ان هنالك اختلافات معنوية بين المتغيرات التوضيحية على المتغير المعتمد وكما اظهرت الدراسة انه لايوجد مشكلة تعدد خطى .

Nonlinear Estimation)), Technometrics, Vol. 12, pp (591-612).

- 5- Mason ,R.L.,Gunst,R.F. and Webster, I.T., (1975), "Regression Analysis and problems of Multicollinearity", comm. . In statistics .VOL.(4) No.(3) pp(277-292).
- 6-Saikat Maitra and Jun yan , (2008) , "Principle component Analysis and Partial Least squares: Two Dimension Reduction Techiques for Regression ," casualty Actuarial society , pp (79-90).
- Y- Wold , H . (1966) ⁽⁽⁾ La Regression PLS . Paris : Technip Iterative Least squares ⁽⁾⁾ , In P . R . Kishnaiaah (Ed) . Multivairte Analysis , New york, Academic Press . PP (391 -420).

حيث نختار الجذور المميزة الاكبر من واحد وهي الجذور (الاول والثاني والثانث) ثم نحسب معادلة الانحدار المتعدد وهي كالاتي : $y = 1.02557 + 0.01393 \, pc1 - 0.42407 \, pc2 - 0.63449 \, pc3$ والان نعوض عن كل من Pc_1 و Pc_2 في معادلة الانحدار المتعدد اعلاه فنحصل على المعادلة الانحدار الاتبة :

 $y = 1.02557 + 2.15701x_1 + 1.15367x_2 + 0.73568x_3 + 0.4823x_4$ وان جدول تحليل التباين لـ y کالاتي

جدول (٣) تحليل التباين لانتاج السمنت بطريقة (pc)

souree	DF	SS	MS	F	P
Regession Residnal Error Total	3 9 12	2664.18 52.74 2716.92	888.06 5.86	151.54	0.000

٣- طريقة المربعات الصغرى الجزئية (Pls)

ان الطريقة الثانية لاجل معالجة مشكلة التعدد الخطي تم استخدام طريقة المربعات الصغرى الجزئية (pls) التي تم توضيحها بالجانب النظري وكانت النتائج كما يلى:-

جدول (٤) تحليل التباين لانتاج السمنت بطريقة (pls)

souree	DF	SS	MS	F	P
Regession	4	26665.83 51.09	666 159	10/13	
Residual Error	8	51.09	6.386	104.3	0.000
Total		2716.92	0.380	6	

وللتعرف على مدى معنوية العلاقة الخطية المقترحة واختبار مدى تأثير المتغيرات التوضيحية على المتغير المعتمد . اختبرت الفرضية الخاصة نموذج الانحدار كالاتى :

 $H_0: B_1 = B_2 = B_3 = B_4$ $H_1: B_1 \neq B_2 \neq B_3 \neq B_4$

المصادر

١.الراوي ، خاشع محمود ، ١٩٨٧ ، " المدخل الى تحليل الانحدار "
 ، دار الكتب للطباعة والنشر جامعة الموصل .

- 3- Frank , I . E . and friedman , J . H . (1993) . ((Astatistical view of chemometrics Regression Tools (Technometrics , 35 , pp (109-148)
- 4- Marquardt, D. W. (1970), "Generalized Inverse, Ridge Regression, Biased Linear Estimation and

Comparision between principal components and partial least square method in multicollinrarity

Farah Abdul – Ghana y. Al – Saleh

Dept. of statistic, Coll. Of comp.& math. Science, Mosul University, Mosul, Iraq (Received: 19/10/2010 ---- Accepted: 16/3/2011)

Abstract: -

Multicollinearity is considered as one of the serions problem in regression analysis , so in this paper we give a big accuracy intention to this problem and to find its solutions .We compare between pc and pls , which is considered as one of important methods of multicollinraity MSE criterion is used for choosing the best estimation Method