

## تحليل الارتباط القوي اللخطي باستخدام بعض دوال كيرنل

علي جواد كاظم\*

### الخلاصة :

يعد تحليل الارتباط القوي من تحليلات متعددة المتغيرات التقليدية حيث يعمل على ايجاد تحويل خطى لزوج من المتغيرات المتعددة بحيث يحقق هذا التحويل تعظيم معامل الارتباط و يجعله اعظم ما يمكن ومن وجہة النظر النظرية فان هذا التحويل يعمل على تعظيم المعلومات المستخلصه عن المعالم .

اذا كانت العلاقة بين المتغيرات غير خطيه فان تحليل الارتباط القوي لا يستطيع استخلاص معلومات مفيدة عن تلك المعالم . في هذا البحث تم استخدام ما يدعى بخدعة كيرنل (Kernel trick) والتي تعمل على جعل البيانات ملائمه للتحليل وتم استخدام دوال مختلفة من نوع كيرنل منها دالة (Uniform Kernel) و دالة (Gaussian Kernel) .

### المقدمة :

يفترض هذا التحليل وجود زوج من المتغيرات المتعدده  $X \in R^{n_x}, Y \in R^{n_y}$  ويعلم على ايجاد زوج من التحويلات الخطيه التي تجعل معامل الارتباط بين المعالم المستخلصه اعظم ما يمكن وبافتراض ان  $X, Y$  لكل منهما متوسط مقداره صفر وتباعي مقداره واحد فان التحويلات تكون كما يأتي<sup>(3)</sup>:

$$U = \langle a, x \rangle, V = \langle b, y \rangle$$

حيث ان :

$\langle a, x \rangle, \langle b, y \rangle$  تمثل حاصل الضرب الداخلي (inner product) بين  $a, x$  و  $b, y$  على التوالي .

والهدف هنا هو ايجاد التحويل المناسب الى  $a, b$  والذي يجعل معامل الارتباط اعظم ما يمكن . وبمعنى اخر اذا كان لدينا متغير عشوائي يتكون من  $p$  من الأبعاد هو  $X$  ومتغير يتكون من  $q$  من الأبعاد هو  $Y$  ، هدفا هو الحصول على تركيبات خطيه  $a^T X$  ،  $b^T Y$  من المتغيرات الأصلية تمتلك اعظم ارتباط ، ويمكن التعبير عن ذلك بصيغة رياضية بالشكل التالي<sup>(8)</sup> :

\* مدرس مساعد قسم الاحصاء/جامعة القادسية

$$\rho = \text{MAX} |\text{corr}(a^T X, b^T Y)| \dots \dots (1)$$

ان المتغيرات الناتجة  $X = a^T U$  ،  $Y = b^T V$  تدعى المتغيرات القوية وان الارتباط القوي الأول  $\rho$  يعرف كقيمه مطلقه لارتباط بين مجموعتين من المتغيرات القوية وكما في العلاقة رقم (1) ، ان المتغير القوي ذو الرتبة  $k < k < \text{Min}(p, q)$  سيكون غير مرتبط مع كل المتغيرات القوية ذات الرتب الدنيا .

ويمكن ايجاد  $a^T, b^T$  كمتجه مميز مرفق لقيم مميزه عظمى <sup>(3)</sup> .

افرض ان  $\Sigma$  تمثل مصفوفة التباين المشترك للمجتمع للمتغير العشوائي  $Z$  حيث ان  $Z = (X^T, Y^T)^T$  ويمكن وضع  $\Sigma$  بالشكل التالي :

$$\Sigma = \begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix}$$

وقد نستخدم مصفوفة الارتباطات  $R$

$$R = \begin{bmatrix} R_{XX} & R_{XY} \\ R_{YX} & R_{YY} \end{bmatrix}$$

حيث يتم حساب المتجه المميز  $b^T$  من العلاقة التالية <sup>(7)</sup> :

$$(M - \lambda I)b = 0 \dots \dots (2)$$

بفرض ان

$$M = \Sigma_{YY}^{-1} \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$$

عند استخدام مصفوفة التباين المشترك .

او

$$M = R_{YY}^{-1} R_{YX} R_{XX}^{-1} R_{XY}$$

عند استخدام مصفوفة الارتباطات .

ان قيم المتجه المميز  $b^T$  والذي يرتبط بكل قيمة من القيم المميزه ( $\lambda$ ) الناتجه من المعادله (2) تمثل اووزانا للمجموعه  $Y$  من المتغيرات .

وبعد حساب  $b^T$  يمكن ان نحدد المتجه  $a^T$  المقابل للمجموعه  $X$  من المتغيرات من خلال العلاقة التالية <sup>(7)</sup>:

$$a^T = \frac{1}{\sqrt{\lambda}} R_{XX}^{-1} R_{XY} d^T$$

يستخدم تحليل الارتباط القوي عندما تكون  $\mathbf{X}$ ,  $\mathbf{Y}$  متغيرات تتبع التوزيع الطبيعي المشتركة وحتى عندما يكون هذا الافتراض غير متحقق يمكن استخدامه في بعض الحالات وعندما يكون هدفنا هو الحصول على انحدار تلك المتغيرات فاننا نرغب في الحصول على قيم كبيرة لمعاملات الارتباط . ولكن احيانا قد نحصل على قيم صغيرة لتلك المعاملات وقد يكون ذلك عائدا الى احد السببين التاليين<sup>(3)</sup> :

1- عدم وجود أي علاقة بين  $\mathbf{X}$ ,  $\mathbf{Y}$  .

2- هنالك علاقة غير خطية قوية بين  $\mathbf{X}$ ,  $\mathbf{Y}$  .

في الحال الاولى لا يمكن الحصول على افضل مما تم الحصول عليه ولا يمكن اجراء أي تحسينات تذكر .اما في الحاله الثانية فيمكن الحصول على العلاقة باستخدام بعض الطرق منها طريقة كيرنل (Kernel Method) في تحليل الارتباط القوي .

### طريقة كيرنل في تحليل الارتباط القوي Kernel CCA<sup>(3)</sup>

في هذه الطريقة يتم تحويل  $\mathbf{Y}$ ,  $\mathbf{X}$  في فضاء هيلبرت<sup>\*</sup> (Hilbert Space) حيث :

$$\phi_x(x) \in H_x, \phi_y(y) \in H_y$$

وبأخذ حاصل الضرب الداخلي (inner product) للمعلم في فضاء هيلبرت

$$a \in H_x, b \in H_y$$

نحصل على ما ياتي :

$$U = \langle a, \phi(x) \rangle, V = \langle b, \phi(y) \rangle$$

ويفرض وجود عينات  $\{(x_i, y_i)\}_{i=1}^N$  يمكن ايجاد  $a, b$  وذلك عن طريق حل علاقة لاكرانج التالية :

$$L_0 = E[(U - E[U])(V - E[V])] - \frac{\lambda_1}{2} E[(U - E[U])^2] - \frac{\lambda_2}{2} E[(V - E[V])^2]$$

ان علاقه لاكرانج لا تعمل بشكل جيد عندما تكون ابعاد فضاء هيلبرت كبيره ولهذا يجب اضافة حد انتظام تربيعي (Quadratic Regularization Term) وذلك من اجل تحسين عمل علاقه لاكرانج وكما ياتي :

$$L = L_0 + \frac{\eta}{2} (\|a\|^2 + \|b\|^2)$$

حيث ان  $\eta$  تمثل ثابت انتظام (Regularization Constant) .  
ان معدل  $\eta$  يمكن حسابه من الصيغه التالية :

$$E[U] = \frac{1}{N} \sum_i \langle a, \phi_x(x_i) \rangle$$

• ان فضاء هيلبرت الابتدائي الكامل يدعى فضاء هيلبرت. بعده اخرى اذا كان  $X$  فضاء

متوجهات على الحقل  $F$  مع الضرب الداخلي  $\langle , \rangle$  فان  $X$  يكون فضاء هيلبرت اذا

كان الفضاء المترى المتولد بواسطة المعيار  $\|x\|^2 = \langle x, x \rangle$  فضاء كاملا<sup>(1)</sup>.

ويمكن توضيح فضاء هيلبرت الابتدائي او ما يسمى احيانا بفضاء الضرب الداخلي من خلال التعريف التالي:

تعريف: ليكن  $X$  فضاء متوجهات على الحقل  $F$ . الدالة

بدالة الضرب الداخلي على  $X$  ، اذا تحقق البديهيات التالية<sup>(1)</sup>:

$$\text{. } x \in X \text{ لكل } \langle x, x \rangle \geq 0 \quad -1$$

$$\text{. } x \in X \text{ اذا وفقط اذا كان } \langle x, x \rangle = 0 \quad -2$$

$$\text{. } \langle x, y \rangle \text{ حيث } \overline{\langle x, y \rangle} = \langle y, x \rangle \quad -3$$

$$\text{. } x, y, z \in X \text{ لكل } \langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle \quad -4$$

$$\text{. } \alpha, \beta \in F$$

فضاء هيلبرت الابتدائي ( pre-Hilbert Space) او فضاء الضرب الداخلي هو

الثاني  $(X, \langle \rangle, F)$  حيث  $X$  فضاء متوجهات على الحقل  $F$  ضرب داخلي على

$$\text{. } X$$

وكذلك معلم UV يمكن حسابه بالشكل التالي:

$$E[UV] = \frac{1}{N} \sum_{i,j} \langle a, \phi_x(x_i) \rangle \langle b, \phi_y(y_j) \rangle$$

وباشتقاق  $L$  بالنسبة الى  $a$  ومساواة المشتقه الى الصفر نحصل على ما يأتي:

$$a = \sum_i \alpha_i \phi_x(x_i)$$

حيث ان  $\alpha_i$  تمثل ثابت (Scalar).

ونتيجه للخطوه السابقه فان :

$$U = \sum_i \alpha_i \langle \phi_x(x_i), \phi_x(x) \rangle$$

والتي يمكن حسابها عن طريق حاصل الضرب الداخلي . ان ما يسمى بخدعة كيرنل (Kernel Trick) يتلخص باستخدام دالة كيرنل  $K_x(X_1, X_2)$  بدلا من حاصل الضرب الداخلي بين  $\phi_x(X_1), \phi_x(X_2)$  وبذلك لسنا بحاجه الى صيغة  $\phi_x(X)$  نحتاج فقط الى تحديد  $K_x$  والتي تكون معرفه موجبه متماثله <sup>(2)</sup>. ويمكن اعادة كتابة L من حسب وجهة نظر طريقة كيرنل كما يأتي :

افرض ان

$$\beta = (\beta_1, \beta_2, \dots, \beta_N)^T, \quad \alpha = (\alpha_1, \alpha_2, \dots, \alpha_N)^T$$

$$(K_x)_{ij} = K_x(x_i, x_j), \quad (K_y)_{ij} = K_y(y_i, y_j)$$

: فان

$$L = \alpha^T M \beta - \frac{\lambda_1}{2} \alpha^T L \alpha - \frac{\lambda_2}{2} \beta^T N \beta$$

حيث ان :

$$M = \frac{1}{N} K_x^T J K_y$$

$$L = \frac{1}{N} K_x^T J K_x + \eta_1 K_x$$

$$N = \frac{1}{N} K_y^T J K_y + \eta_2 K_y$$

$$J = I - \frac{1}{N} 1 1^T$$

$$1 = (1, 1, \dots, 1)^T$$

$$\eta_1 = \frac{\eta}{\lambda_1}, \quad \eta_2 = \frac{\eta}{\lambda_2}$$

وإذا كانت  $\eta > 0$  فان كل من  $N, L$  سيكون مصفوفة معرفة موجبه (Positive definite matrix) ويمكن اثبات ان  $\lambda_1 = \lambda_2 = \lambda$  ونتيجة لذلك يكون لدينا ما يدعى بمشكلة قيمة مميزة عامة (Generalized eigenvalue problem) لـ  $\alpha, \beta$  وكما ياتي:

$$M\beta = \lambda L\alpha$$

$$M^T\alpha = \lambda N\beta$$

والتي يمكن حلها كمشكلة قيمة مميزة عامة<sup>(2)</sup>.

#### المحاكاة:

لقد تم استخدام المحاكاة لغرض إجراء المقارنة بين مقدرات الارتباط القوي المحسوبة وفق الطريقة التقليدية ومقدرات الارتباط القوي المحسوبة وفق طريقة كيرنل باستخدام ثلاث دوال من نوع كيرنل هي دالة (Gaussian Kernel) ودالة (Uniform Kernel) ودالة (Triangle Kernel).

Kernel وذلك بهدف معرفة أفضلية المقدرات وذلك بالاعتماد على كون الارتباط القوي يمثل أكبر الارتباطات للحكم على المقدر الأفضل وتم كتابة برنامج باستخدام لغة (V.B) خاص بالتجربة لتحقيق هذا الهدف حيث تم تكرارها (1000) مرة لغرض الوصول إلى نتائج مُقنعة ويمكن أدراج خطوات أجراء المحاكاة وكذلك :

1— نولد قيم للمتغيرات ( $x_1, x_2$ ) لتكوين تركيبة خطية هي  $UC_i$  والتي يتم استخدامها في حساب الارتباط القوي وفق الطريقة التقليدية وتركيبه خطية أخرى هي  $UK_i$  والتي يتم استخدامها في حساب الارتباط القوي وفق طريقة كيرنل وبأحجام عينات ( $n=10, 20, 30, 40, 50, 60, 70, 80, 90$ ) أي لكل متغير من متغيرات  $x$  نولد  $n$  من المشاهدات المذكورة حيث إن :

$$UC_i = a_1x_1 + a_2x_2$$

$$UK_i = a_1K(x_1) + a_2K(x_2)$$

وان صيغ دوال كيرنل المستخدمة هي الصيغ التالية :

$$a) K(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

وهذه الدالة تدعى دالة (Gaussian Kernel).

b)  $K(x) = 1/2, |x| \leq 1$

و هذه الدالة تدعى دالة (Uniform Kernel).

c)  $K(x) = (1 - |x|), |x| \leq 1$

و هذه الدالة تدعى دالة (Triangle Kernel).

وان :

$$X = (\theta) \sin 3\theta + \varepsilon_1$$

حيث يتم توليد قيم  $\theta$  على اساس انها تتبع التوزيع المنتظم المستمر على الفترة  $[-\pi, \pi]$ .  
 $\varepsilon_1$  : يمثل حد الخطأ العشوائي ويتم توليده على اساس انه يتبع التوزيع الطبيعي القياسي في  
 الحالة الأولى ويتبع التوزيع الاسي بالمعلمة  $\lambda=1/3$  في الحالة الثانية أما في الحالة الثالثة فيتم  
 توليده على اساس انه يتبع توزيع مربع كاي بالمعلمة  $n$ .

2- يتم توليد متغيرات  $(y_1, y_2)$  لتكوين تركيبة خطية هي  $VC_i$  والتي يتم استخدامها في  
 حساب الارتباط القوي وفق الطريقة التقليدية وتركيبه خطية أخرى هي  $VK_i$  والتي يتم  
 استخدامها في حساب الارتباط القوي وفق طريقة كيرنل وب أحجام عينات  
 $(n=10, 20, 30, 40, 50, 60, 70, 80, 90)$  أي لكل متغير من متغيرات  $y$  نولد  $n$  من  
 المشاهدات المذكورة حيث إن :

$$VC_i = b_1 y_1 + b_2 y_2$$

$$VK_i = b_1 K(y_1) + b_2 K(y_2)$$

وان  $K(y)$  تأخذ نفس الصيغ السابقة وان الصيغة التي يتم الحصول على قيم  $y$  منها هي  
 الصيغة التالية:

$$Y = e^{\theta/4} (\cos 2\theta) (\sin 2\theta) + \varepsilon_2$$

حيث يتم توليد قيم  $\theta$  على اساس انها تتبع التوزيع المنتظم المستمر على الفترة  $[-\pi, \pi]$ .

ونولد قيم للمتغير  $\epsilon_2$  مرة بالتوزيع الطبيعي القياسي ومرة أخرى بالتوزيع الأسوي بالمعلمة  $\lambda = 1/3$  ومره ثالثه بتوزيع مربع كاي بدرجة حرية (n).

3- يتم حساب معاملات الارتباط بين  $VC_i$  و  $UC_i$  حيث ان اكبرها يمثل معامل الارتباط

القويم المحسوب وفق الطريقة التقليدية ويتم حساب معاملات الارتباط بين  $VK_i$  و  $UK_i$  حيث ان اكبرها يمثل معامل الارتباط القوي المحسوب وفق طريقة كيرنل حسب الدالة المستخدمة أي يتم حساب معامل الارتباط القوي نوع كيرنل لدوال كيرنل الثلاث.

4- يتم المقارنة بين معاملات الارتباط القوي المحسوبه وفق طريقة كيرنل وبالدوال الثلاث ومعامل الارتباط القوي المحسوب وفق الطريقة التقليدية.

### تحليل النتائج :

بعد اجراء تجربة المحاكاة تم الحصول على النتائج التالية:

1- من الجدول رقم (1) والذي يمثل قيم معاملات الارتباط القوي المحسوبه باستخدام طريقة كيرنل وبالدوال الثلاث وكذلك معاملات الارتباط القوي المحسوبه وفق الطريقة التقليدية عندما يكون  $\epsilon_1, \epsilon_2$  يتبع كل منهما التوزيع الطبيعي القياسي نلاحظ ان معاملات الارتباط القوي المحسوبه وفق طريقة كيرنل كانت اكبر من تلك المحسوبه وفق الطريقة التقليدية ولجميع حجوم العينات وكانت اكبر تلك الارتباطات باستخدام دالة (Gaussian Kernel) ومن ثم باستخدام دالة (Uniform Kernel) وبعد ذلك باستخدام دالة (Triangle Kernel). ونلاحظ بشكل عام ان الارتباطات المحسوبه وفق الطريقة التقليدية كانت صغيره جدا.

2- نلاحظ من الجدول رقم (2) والذي يمثل قيم معاملات الارتباط القوي المحسوبه باستخدام طريقة كيرنل وبالدوال الثلاث وكذلك معاملات الارتباط القوي المحسوبه وفق الطريقة التقليدية عندما يكون  $\epsilon_1, \epsilon_2$  يتبع كل منهما التوزيع الأسوي بالمعلمه  $\lambda = 1/3$  ان معاملات الارتباط القوي المحسوبه وفق طريقة كيرنل كانت اكبر من تلك المحسوبه وفق الطريقة التقليدية عند استخدام دالة (Gaussian Kernel) وكذلك عند استخدام دالة (Uniform Kernel) ولجميع حجوم العينات بينما كانت تلك الارتباطات اقل مما كانت عليه ارتباطات الطريقة التقليدية عند استخدام دالة (Triangle Kernel). ونلاحظ بشكل عام ان الارتباطات المحسوبه وفق الطريقة التقليدية كانت اكبر مما كانت عليه في الجدول رقم (1).

3- من الجدول رقم (3) والذي يمثل قيم معاملات الارتباط القوي المحسوبه باستخدام طريقة كيرنل وبالدوال الثالث وكذلك معاملات الارتباط القوي المحسوبه وفق الطريقة التقليديه عندما يكون  $\epsilon_1, \epsilon_2$  يتبع كل منها توزيع مربع كاي بالمعلمee  $n$  نلاحظ ان معاملات الارتباط القوي المحسوبه وفق طريقة كيرنل كانت اكبر من تلك المحسوبه وفق الطريقة التقليديه ولجميع حجوم العينات وكانت اكبر تلك الارتباطات باستخدام دالة (Gaussian Kernel) ومن ثم باستخدام دالة (Uniform Kernel) وبعد ذلك باستخدام دالة (Triangle Kernel). . ونلاحظ بشكل عام ان الارتباطات المحسوبه وفق الطريقة التقليديه كانت اكبر مما كانت عليه في الجدول رقم (1) ، ولكنها كانت اقل مما كانت عليه في الجدول رقم (2) عند حجوم العينات المتشابهه.

#### الاستنتاجات:

- 1- عندما يكون  $\epsilon_1, \epsilon_2$  يتبع كل منها التوزيع الطبيعي القياسي او التوزيع الاسي بالمعلمee  $\lambda=1/3$  او توزيع مربع كاي بالمعلمee  $n$  نلاحظ ان معاملات الارتباط القوي المحسوبه وفق طريقة كيرنل كانت اكبر من تلك المحسوبه وفق الطريقة التقليديه ولجميع حجوم العينات فيما عدا في حالة التوزيع الاسي حيث كانت الارتباطات القويه المحسوبه باستخدام طريقة كيرنل اقل من الطريقة التقليديه عند استخدام دالة (Triangle Kernel) .
- 2- كانت الارتباطات القويه المحسوبه وفق طريقة كيرنل باستخدام دالة (Gaussian Kernel) تمثل اكبر الارتباطات في جميع الحالات وعند جميع حجوم العينات.

#### المصادر:

- 1- المياحي ، نوري فرحان وعلى حسين بتور (2005)."مقدمه في التحليل الدالي ".مؤسسة النبراس للطباعه والنشر والتوزيع ، النجف الاشرف،العراق.
- 2-Anderson, T.W (1984)."An Introduction to Multivariate statistical analysis".Jonwiley and sons.
- 3-Akaho,S.(2001)"A kernel method for canonical correlation analysis". International meeting of psychometric society, Osaka.
- 4- Donald F. Morrison (1988)."Multivariate statistical methods". Second edition, McGraw. Hill series in probability and statistics.
- 5- Florian,M.(2003) "Canonical Correlation Analysis With Kernels". Computational Diagnostics Group Seminar .Berlin.
- 6-Lai,p.1 and C.fyfe(2000)."Kernel and nonlinear canonical correlation analysis " . International Journal Of Neural Systems 10(5) , 365-377.

- 7- Lindeman,R,H(1980)."Introduction to bivariate and multivariate analysis".Scott-Foresman company,Illinois,U.S.A.
- 8- Romanazzi,m.(1992)"Influence in canonical correlation analysis", Psychometrika,57,237-259.

### الملحق

جدول رقم (1)

يمثل الارتباطات القوية المحسوبة وفق الطريقتين وبالدوال المختلفة عندما تتبع

$\epsilon_1, \epsilon_2$  الاخطاء

### التوزيع الطبيعي القياسي

الطريقه	الدالة المستخدمة	n=10	n=20	n=30	n=40	n=50	n=60	n=70	n=80	n=90
Kernel CCA	Gaussian	0.82	0.80	0.79	0.78	0.78	0.77	0.77	0.77	0.77
	Uniform	0.63	0.62	0.61	0.59	0.58	0.57	0.57	0.57	0.58
	Triangle	0.55	0.52	0.50	0.47	0.46	0.45	0.45	0.46	0.46
CCA	Classical	0.10	0.03	0.02	0.002	0.02	0.04	0.04	0.04	0.03

جدول رقم (2)

يمثل الارتباطات القوية المحسوبة وفق الطريقتين وبالدوال المختلفة عندما تتبع

$\epsilon_1, \epsilon_2$  الاخطاء

### التوزيع الاسي بالمعلمه $\lambda=1/3$

الطريقه	الدالة المستخدمة	n=10	n=20	n=30	n=40	n=50	n=60	n=70	n=80	n=90
Kernel CCA	Gaussian	0.73	0.72	0.71	0.71	0.71	0.71	0.72	0.72	0.72
	Uniform	0.56	0.54	0.54	0.53	0.52	0.52	0.53	0.54	0.54
	Triangle	0.47	0.44	0.43	0.42	0.41	0.42	0.42	0.42	0.42
CCA	Classical	0.50	0.47	0.45	0.44	0.44	0.44	0.44	0.44	0.44

**جدول رقم (3)**  
**يمثل الارتباطات القوية المحسوبة وفق الطريقتين وبالدوال المختلفة عندما تتبع**  
**الخطاء  $\epsilon_1, \epsilon_2$**

**توزيع مربع كاي بالمعلمه n**

الطريقة	الدالة المستخدمة	n=10	n=20	n=30	n=40	n=50	n=60	n=70	n=80	n=90
Kernel CCA	Gaussian	0.75	0.74	0.72	0.72	0.71	0.70	0.71	0.71	0.71
	Uniform	0.62	0.59	0.57	0.56	0.55	0.54	0.55	0.55	0.55
	Triangle	0.52	0.51	0.48	0.46	0.45	0.44	0.45	0.45	0.45
CCA	Classical	0.43	0.40	0.39	0.36	0.35	0.34	0.34	0.34	0.34