

استخدام الأساليب اللامعلمية لتقدير الدالة التمييزية

م. طاهر ريسان دخيل قسم الإحصاء/ جامعة القادسية/كلية الإدارة والاقتصاد
م. غسق فاضل عبد الرزاق /الجامعة المستنصرية / قسم الإحصاء / كلية
الإدارة والاقتصاد م. زينب
يوسف داود/ الجامعة المستنصرية /كلية الإدارة والاقتصاد / قسم الإحصاء

الخلاصة

هناك بعض الافتراضات التي يجب أن تتوفر في المتغيرات العشوائية بهدف استخدام الطرائق المعلمية في تقدير الدالة التمييزية ومن هذه الافتراضات أن يكون توزيع هذه المتغيرات هو توزيعاً طبيعياً ولكن هذا الافتراض قد لا يتحقق مما يتطلب البحث عن طرائق بديلة لتقدير الدالة ومن هذه الطرائق أسلوب التقدير اللامعلمي المسمى بتقدير Kernel أو طريقة Kernel والتي تعتمد على البيانات في الحصول على تقدير للدالة مما يجنبنا المرور بتقدير المعلمات للدالة التمييزية، ففي هذا البحث تم استخدام أسلوب Kernel ولثلاث دوال هي دالة Epanchinicove و Gaussian و Sin ومقارنة النتائج مع الطريقة المعلمية التقليدية وذلك بالاعتماد على مصادر خطأ التصنيف في الدالة التمييزية .

Abstract

Some assumptions must available in random variables to allow us to use the parametric methods in discriminate function estimation. One of these assumptions is that the distribution of the variables are normal distribution ,but sometimes this is not achievement so we must seek for an alternative methods to estimate the discriminant function ,one of these methods is the nonparametric estimation procedure which called Kernel estimation or Kernel method which depends on data to get the estimation. In this paper we use Kernel procedure with three functions, Epanchinikov, Gaussian & Sine functions and compare results with the parametric method by using misclassification as criterion of comparing.

المقدمة [5][1]

أن الهدف من الدالة التمييزية هو تصنيف المشاهدات إلى مجموعتين أو أكثر بحيث يكون التباين بين تلك المجاميع أعظم ما يمكن بينما يكون التباين داخل كل مجموعة اقل ما يمكن وذلك من خلال بناء معادلة ذات تركيبية خطية وبالتالي يمكن استخراج الطرائق المعلمية في تقدير تلك الدالة وهذا يمكن أن يتم في حالة وجود الفروض الأساسية لهذا التحليل ولكن عندما لا يتوفر ذلك يمكن الاتجاه إلى الأساليب اللامعلمية، حيث اتجهت البحوث حديثاً إلى دراسة طرائق التقدير اللامعلمي والتقدير اللامعلمي لا يعني أن النموذج لا يحوي على معلمات وإنما هو طريقة لا تعتمد على المعلمات في التقدير وإنما تعتمد على البيانات مباشرة، ومن تلك الأساليب أسلوب Kernel في التقدير والذي يستخدم آلية خاصة في عملية التقدير .

لقد اقترحت طريقة Kernel في التقدير ابتداءً من قبل الباحث Osenblatt [5] عام 1956 وطورت هذه الطريقة من قبل الباحث Parzen عام 1962 وهناك بعض البحوث التي قدمت ضمن التقدير اللامعلمي والتي تضمنت تقديرات Kernel ومنه البحث الذي قدم من قبل الباحث Silverman عام 1986 والذي تضمن استخدام دوال Kernel في تقدير الدوال التمييزية .

هدف البحث

يهدف هذا البحث إلى المقارنة بين الطريقة المعلمية وطريقة Kernel اللامعلمية في استخدام عدة دوال وذلك من خلال استخدام المحاكاة في حالة عدم تحقق الفروض الأساسية للدالة التمييزية وقد استخدم معيار خطأ التصنيف كأساس للمقارنة .

الدالة الخطية للتصنيف [1][4]

يعتبر أسلوب التصنيف احد الأساليب المهمة ضمن متعدد المتغيرات وذلك بهدف معرفة إلى أي مجتمع من المجتمعات تعود مشاهدة ما والمبدأ الأساس الذي يستند عليه هذا التصنيف هو بتقليص خطأ التصنيف أو سوء التصنيف وهو أن نضيف مشاهدة معينة بأنها تعود للمجموعة الأولى مثلاً بينما هي في حقيقة الأمر تعود إلى المجموعة الثانية وهذا المعيار هو معيار مهم فمن خلاله يمكن معرفة قوة الدالة التمييزية أو التصنيفية فكلما كان عدد المشاهدات المصنفة خطأ قليل فإن تلك الدالة جيدة وقادرة على فرز المشاهدات إلى مجتمعاتها الأصلية بصورة صحيحة والعكس بالعكس.

ولغرض شرح الأساس النظري لأسلوب التمييز لنفرض أن هناك عينتين عشوائيتين هما n_1 و n_2 للمجموعة الأولى والثانية على التوالي والتي تتكون من p من المتغيرات العشوائية ويفترض أن هاتين العينتان سحبنا من مجتمعين يتوزعان توزيعاً طبيعياً

بمتوسط μ_1 و μ_2 للمجموعة الأولى والثانية على التوالي وبتباين Σ للمجموعتين ومن هاتين العينتين يتم حساب \bar{Y}_1 و \bar{Y}_2 واللذان يمثلان متجها المتوسطات المقدرة من العينة للمجموعة الأولى والثانية على التوالي وكذلك حساب S_{pl} والتي تمثل مصفوفة التباين المشترك المقدر، وبعدها يمكن حساب الدالة التمييزية المقدرة وبالشكل الآتي

$$Z = a'Y = (\bar{Y}_1 - \bar{Y}_2)'(S_{pl})^{-1}Y \dots\dots\dots 1$$

حيث إن Y هو متجه المشاهدات الجديدة والتي نرغب في تصنيفها إما إلى المجموعة الأولى أو الثانية وحسب المعلومات التي تحملها.

ولغرض تحديد إن Y تعود لإحدى المجموعتين فأنا نعتد بذلك على أن Z في المعادلة ١ هل هي قريبة إلى المتوسط \bar{Z}_1 أو \bar{Z}_2 حيث إننا نستخدم بيانات المجموعة الأولى Y_{1i} لتحديد القيم $Z_{11}, Z_{12}, Z_{13}, \dots, Z_{1n_1}$ ومن ثم حساب \bar{Z}_1 والذي يمكن أيضا أن يحسب وفق الصيغة الآتية

$$\bar{Z}_1 = a'\bar{Y}_1 = (\bar{Y}_1 - \bar{Y}_2)'(S_{pl})^{-1}\bar{Y}_1 \dots\dots\dots 2$$

وبشكل مشابه يمكن حساب قيمة \bar{Z}_2 والتي تساوي $\bar{Z}_2 = a'\bar{Y}_2$ وبالتالي فإن المشاهدة تعود للمجموعة الأولى إذا كانت $\bar{Z} = a'Y$ قريبة إلى المجموعة \bar{Z}_1 أكثر من \bar{Z}_2 وتعود إلى المجموعة الثانية إذا كانت قريبة من \bar{Z}_2 أكثر من \bar{Z}_1 أو بشكل آخر فإن المشاهدة تعود إلى المجموعة الأولى إذا كانت $Z > \frac{1}{2}(\bar{Z}_1 + \bar{Z}_2)$ وحيث أن المسافة بين \bar{Z}_1 و \bar{Z}_2 هي نفس المسافة بين \bar{Y}_1 و \bar{Y}_2 إذ يمكن كتابة قاعدة التصنيف وذلك بالاعتماد على قيم Y بالشكل الآتي

$$\frac{1}{2}(\bar{Z}_1 + \bar{Z}_2) = \frac{1}{2}(\bar{Y}_1 - \bar{Y}_2)'(S_{pl})^{-1}(\bar{Y}_1 + \bar{Y}_2) \dots\dots\dots 3$$

وبالتالي فإن قاعدة التصنيف تصبح بالشكل الآتي

$$a'Ya = (\bar{Y}_1 + \bar{Y}_2)'S_{pl}^{-1}Y > \frac{1}{2}(\bar{Y}_1 - \bar{Y}_2)'S_{pl}^{-1}(\bar{Y}_1 + \bar{Y}_2) \dots\dots\dots 4$$

وفي هذه الحالة فإن المشاهدة تصنف إلى المجموعة الأولى بينما تصنف المشاهدة إلى المجموعة الثانية إذا كانت

$$a'Ya = (\bar{Y}_1 - \bar{Y}_2)'S_{pl}^{-1}Y < \frac{1}{2}(\bar{Y}_1 - \bar{Y}_2)'S_{pl}^{-1}(\bar{Y}_1 + \bar{Y}_2) \dots\dots\dots 5$$

خطا التصنيف Miss Classification [1][5]

أن خطأ التصنيف يكمن في ان الدالة تصنف مشاهدة معينة في المجموعة الأولى وهي في حقيقة الأمر تعود الى المجموعة الثانية او العكس أي أنها تصنف مشاهدة معينة الى المجموعة الأولى وهي تعود الى المجموعة الثانية ويمكن توضيح ذلك بالجدول الاتي

المجموعة الحقيقية		عدد المشاهدات	المجموعة المتنبأ بها	
1	2		1	2
n_{11}	n_{12}	n_1	n_{21}	n_{22}

من خلال الجدول أعلاه يمكن إيجاد نسبة خطأ التصنيف للدالة التمييزية المقدرة ولنرمز لها بالرمز E

$$E = \frac{n_{12} + n_{21}}{n_1 + n_2} \dots\dots\dots 6$$

وكذلك يمكن حساب نسبة التصنيف الصحيح بالشكل الاتي

$$C = \frac{n_{11} + n_{22}}{n_1 + n_2} \dots\dots\dots 7$$

حيث ان C يمثل نسبة التصنيف الصحيح

التصنيف بالاعتماد على دوال Kernel [4][5][6]

إن التصنيف يعتمد على الفروض الأساسية لقواعد التصنيف الخطية والتربيعية ومنها أن المتغيرات تتبع التوزيع الطبيعي وفي حالة عدم تحقق هذا الشرط فإنه يمكن الانتقال إلى الطرائق اللامعلمية وأحدى هذه الطرائق تسمى طريقة Kernel وهنا يمكن إعطاء شرحاً عن طريقة Kernel، فعلى افتراض وجود متغير Y والذي يمثل الدالة f(y) والتي يطلب تقديرها باستخدام العينة العشوائية $y_1, y_2, y_3, \dots, y_n$ إن التقدير ل f(y) لأي نقطة (y_0) يمكن أن تعتمد على النسبة للنقاط ضمن الفترة $(y_0 - h, y_0 + h)$ فإذا رمزنا لعدد

النقاط ضمن الفترة المذكورة بالرمز $N(y_0)$ فستكون النسبة $N(\frac{y_0}{h})$ تقديرا إلى $p(y_0 - h < y_0 < y_0 + h)$ والتي تكون مساوية بشكل تقريبي للمقدار $2h(y_0)$ ولذلك يمكن تقدير $f(y_0)$ باستخدام المعادلة الآتية

$$f(y_0) = \frac{N(y_0)}{2hn} \dots\dots\dots 8$$

ويمكن وضع المقدار $\hat{f}(y_0)$ كدالة لكل قيم y_i في العينة وذلك من خلال الدالة الآتية

$$K(u) = \begin{cases} \frac{1}{2} & \text{for } |u| \leq 1 \\ 0 & \text{for } |u| > 1 \end{cases} \dots\dots\dots 9$$

حيث أن الدالة أعلاه تسمى بالدالة المنتظمة Uniform function وبعد التعويض عن قيمة u بالمقدار $\frac{y_0 - y_i}{h}$ نحصل على معادلة

$$Ny_0 = 2 \sum_{i=1}^n k(\frac{y_0 - y_i}{h}) \dots\dots\dots 10$$

لتصبح المعادلة رقم ٨ بالشكل الآتي

$$f(y_0) = \frac{1}{hn} \sum_{i=1}^n k(\frac{y_0 - y_i}{h}) \dots\dots\dots 11$$

وان المقدار $k(\frac{y_0 - y_i}{h})$ يعني أنها تساوي نصف لأي نقطة y_i ضمن الفترة $(y_0 - h, y_0 + h)$ وتساوي صفر لأي نقطة خارج هذه الفترة . ويمكن كذلك استخدام دوال تابعة لأسلوب Kernel غير الدالة المنتظمة uniform function ومن تلك الدوال

$$k(u) = \frac{1}{\pi} \frac{\sin^2(u)}{u^2} \dots\dots\dots 12$$

حيث تسمى هذه الدالة بالدالة الجيبية sine function .
والدالة التالية تسمى بالدالة الكاوسية (Gaussian function)

$$k(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2} \dots\dots\dots 13$$

اما الدالة التالية فتسمى بدالة Epanechnikov

$$k(u) = \begin{cases} \frac{3}{4}(1-u^2) & \text{if } (|u| \leq 1) \\ 0 & \text{otherwise} \dots\dots\dots 14 \end{cases}$$

حيث ان دوال kernel هذه لا تتضمن أي افتراضات حول دالة الكثافة f(y) أن عملية اختيار قيمة h التي تستخدم مع دوال kernel يجب أن تكون بحذر شديد، إذ أن إذا كانت قيمة h صغيرة فان دالة $\hat{f}(y_o)$ ستكون لها قيمة عالية عند كل y_i بينما إذا كانت h كبيرة فان $\hat{f}(y_o)$ تكون غالباً منتظمة لذلك فان اختيار قيمة h يجب ان تعتمد على حجم العينة وذلك لتجنب عملية التمهيد الكبيرة او الصغيرة بحيث لو كان حجم العينة كبيراً فيجب اختيار قيمة h صغيرة والعكس بالعكس او يمكن اختيار عدة قيم لـ h والمقارنة بينها من خلال قيمة الخطأ .

ولغرض استخدام طريقة Kernel في تقدير الدالة التمييزية فانه يمكن تطبيقها على كل مجموعة وذلك لغرض تحديد $\hat{f}(y_o|g_1), \hat{f}(y_o|g_2) \dots \hat{f}(y_o|g_k)$ ان $g_i, i=1, 2, \dots, k$ تمثل المجاميع الغير معلومة ، لتصبح قاعدة التصنيف هي تعظيم للاحتمال الاتي

$$p_i(\hat{f}(y_o|g_i))$$

إذ أن g_i تمثل المجموعة i

بناء تجربة المحاكاة

لقد تم بناء تجربة المحاكاة باستخدام لغة البيسك بحيث يتم خرق الافتراض القائل بان المتغيرات الداخلة في عملية تقدير الدالة التمييزية ذات توزيع طبيعي وذلك من خلال توليد مشاهدات تتبع التوزيع الاسي لأحد المتغيرات ومن ثم تقدير الدالة التمييزية بالطريقة المعلمية التقليدية أولاً وبعدها استخدام طريقة Kernel اللامعلمية وبثلاث دوال وهي دالة

Epanchinicove ودالة Gaussian والدالة الجيبية Sine function وذلك عند أحجام عينات مختلفة هي 250,200,150,100,50,30,10 وتكرار التجربة بواقع ١٠٠٠ مرة وملاحظة سلوك خطأ التصنيف لنتم المقارنة على أساسه. ويمكن تلخيص خطوات المحاكاة بالشكل الاتي تحديد حجم العينة وكما مذكور سابقاً. توليد مشاهدات لثلاث متغيرات ولمجموعتين بحيث يكون توزيع المتغير الثاني ذو توزيع أسّي بالمعلمة $\lambda = 0.3$. حساب تقديراً لدالة التمييزية على وفق الطريقة المعلمية التقليدية وطريقة Kernel اللامعلمية وحسب الدوال المستخدمة والتي تم شرحها سابقاً. تكرار التجربة ١٠٠٠ مرة. حساب خطأ التصنيف وذلك بتحديد المشاهدات التي صنفت إلى المجموعة الخاطئة وحساب شبه الخطأ.

تحليل النتائج

١- عند استخدام دالة Epanchinicove نلاحظ بصورة عامة وعند استخدام هذه الدالة أن طريقة Kernel اللامعلمية تكون أفضل من الطريقة المعلمية التقليدية لتقدير دالة التصنيف فمثلاً عند حجم العينة $n=10$ نلاحظ أن الطريقة المعلمية التقليدية صنفت بعض المشاهدات خطأ وبنسبة 45% بينما صنفت طريقة Kernel اللامعلمية 40% من المشاهدات خطأ. ونلاحظ أيضاً أن خطأ التصنيف يتناقص عند زيادة حجم العينة ولكلنا الطريقتين مع بقاء الأفضلية لطريقة Kernel اللامعلمية، هذا ما يوضحه الجدول رقم ١ في الملاحق فعند حجم العينة $n=30$ كانت نسبة خطأ التصنيف 16%، 28% بالنسبة للطريقة المعلمية التقليدية وطريقة Kernel اللامعلمية على التوالي وهكذا لجميع أحجام العينات المستخدمة حيث نلاحظ التقارب بين قيم خطأ التصنيف للطريقة المعلمية 7% بينما الطريقة اللامعلمية 6%.

٢- عند استخدام الدالة الكاوسية Gaussian function عند استخدام هذه الدالة فإننا نلاحظ أن نفس السلوك الذي تسلكه الدالة السابقة يمكن أن يسري على هذه الدالة إن 50% من المشاهدات صنفت خطأ عندما يتم استخدام الطريقة المعلمية التقليدية بينما 16% من المشاهدات صنفت خطأ استخدام طريقة Kernel اللامعلمية وذلك حجم العينة

$n=10$ وهكذا بالنسبة لبقية أحجام العينات فعند حجم العينة $n=250$ نلاحظ ان خطأ التصنيف بلغت نسبته 6% لكلتا الطريقتين وبقية نسب خطأ التصنيف موضحة في جدول رقم ٢ وعموماً فان هناك تقارب بين قيم نسبة خطأ التصنيف عند استخدام هذه الدالة والدالة السابقة.

٣- عند استخدام الدالة الجيبية Sine function توضح قيم خطأ التصنيف في الجدول رقم ٣ أن 20% من المشاهدات صنفت خطأ عند استخدام الطريقة المعلمية بينما 32% من المشاهدات صنفت خطأ عند استخدام طريقة Kernel اللامعلمية وذلك عند حجم العينة $n=10$ وهكذا بالنسبة لبقية أحجام العينات المستخدمة ، حيث نلاحظ ان هذه الدالة مع الدالة السابقة (الدالة الكاوسية Gaussian function) افضل من دالة Epanchinicove وذلك بموجب النتائج التي أعطتها.

الاستنتاجات

- ١- كانت طريقة Kernel اللامعلمية افضل من الطريقة المعلمية التقليدية .
- ٢- ان دالتي Sine function و Gaussian function تعطي نتائج افضل من دالة Epanchinicove عند استخدام طريقة Kernel
- ٣- نلاحظ ان عند زيادة حجم العينة تقل نسبة خطأ التصنيف لكلتا الطريقتين .

المصادر

الجبوري ، شلال وعبد ، صلاح حمزة (٢٠٠٠) "تحليل متعدد المتغيرات "، الجامعة المستنصرية

Gavin, C (2005) "Efficient Cross -Validation of Kernel Fisher discriminant classifiers "Elsevier Scic,uk.

Hyunsoo Kim and Hyson park (2005) ,,Relations lips between support vector classifiers and Generalized linen discriminant analysis on support vectors " Technical Report university of Minnesota

Rencher, A (2002)" Methods of Multivariate Analysis " John Wiley & Sons ,Inc.

Shan ,L.and Tad E.C (2005) ' Nonlinear Kernel MSE. Methods for Cancer classification " Springer verlag Berlin Heidelberg p .p 975 – 984

Tao,X. , Y. ,Jieping , Qi, L. , C. , Vladimir &J. ,Ravi (2005)" Efficient Kernel Discriminate Analysis via QR decomposition " The Army high per furnace computing research center.

الملاحق

جدول رقم (١)

ويمثل قيم خطأ التصنيف عندما يتم استخدام دالة Epanchinicove

حجم العينة	الطريقة التقليدية	طريقة كيرنل
١٠	٠,٤٥	٠,٤٠
٣٠	٠,١٦	٠,٢٨
٥٠	٠,٢١	٠,١٦
١٠٠	٠,٢٣	٠,١٧
١٥٠	٠,١٠	٠,٠٩
٢٠٠	٠,٠٨	٠,٠٨
٢٥٠	٠,٠٧	٠,٠٦

جدول رقم (٢)

ويمثل قيم خطأ التصنيف عندما يتم استخدام دالة Gaussian

حجم العينة	الطريقة التقليدية	طريقة كيرنل
١٠	٠,٥٠	٠,١٦
٣٠	٠,٢١	٠,١٨
٥٠	٠,١٨	٠,١٥
١٠٠	٠,١٤	٠,١٤
١٥٠	٠,٠٩	٠,٠٦
٢٠٠	٠,٠٨	٠,٠٦
٢٥٠	٠,٠٦	٠,٠٦

جدول رقم (٣)

ويمثل قيم خطأ التصنيف عندما يتم استخدام الدالة الجيبية Sine

حجم العينة	الطريقة التقليدية	طريقة كيرنل
١٠	٠,٢٠	٠,٣٢
٣٠	٠,١٦	٠,١٣

المحور الاحصائي _____ استخدام الأساليب اللامعلمية لتقدير الدالة التمييزية

٥٠	٠,١٢	٠,٢١
١٠٠	٠,٠٩	٠,٠٨
١٥٠	٠,٠٤	٠,٠٣
٢٠٠	٠,٠٣	٠,٠٥
٢٥٠	٠,٠٣	٠,٠٣

دورية فصلية علمية محكمة تصدر عن كلية الإدارة والاقتصاد

This document was created with Win2PDF available at <http://www.daneprairie.com>.
The unregistered version of Win2PDF is for evaluation or non-commercial use only.