

تحليل الاقتران اللامعلمي في ظل زمن الفشل لبيانات المخاطر المنافسة متعددة المتغيرات

أ.م.د. لقاء علي محمد م.م. امال صادق الموسوي (*)

كلية الادارة والاقتصاد / جامعة بغداد

المستخلص :

ان التحليل اللامعلمي لبيانات المخاطر المنافسة الثنائية ذو اهمية كبيرة في الدراسات الاجتماعية والطبية لاسيما علم الامراض الجينية اذ ان اسباب الفشل المتعددة تتفاعل مع بعضها مما يبطل التحليل اللامعلمي بسبب انتهاك فرضية استقلالية بيانات المخاطر المنافسة وهو شائع في المجتمعات التي تعتمد على بداية مؤكدة للفشل كالامراض المزمنة والاضطرابات النفسيةالخ اذ تكون المراقبة مرتبطة مع الموت في هذا البحث قمنا بتطوير مقدرات لامعلمية لدالة خطورة سبب- معين الثنائية ودالة الحادثة التجميعية الثنائية دون الحاجة الى وجود فرضيات حول استقلالية بيانات المخاطر المنافسة ان هذه الطريقة لم تطرح سابقا لاستعمالها بيانات المخاطرة المنافسة بدلا من بيانات اوقات الفشل الاعتيادية وان مقدرات هذه الطريقة تتصف بكونها متسقة ،موحدة وتقترب بضعف الى عمليات Gaussain

Abstract:

The non parametric analysis of bivariate competing risks data is too importance in social studies and medical studies ,especially pathologist genetic where the multi causes of failure interaction ,with the others ,violating independence hypothesis of censoring data which is common in communities that depend on the beginning of the proven failure as in chronic disease estimators where the censoring are associated with death developed non parametric estimators for the bivariate cause-specific hazard function and the bivariate cumulative incidence function without make any assumptions about the independence of the risks data. This method has not been touched upon previously used data competition rather than risk the usual times of failure data .

(*) جزء مسئل من اطروحة دكتوراه للباحثة الثانية .

And estimators of this method are characterized by being consistent ,uniform and approaching weakly to Gaussain process.

الهدف من البحث The Aim of Re search :

ان الهدف الرئيسي من هذا البحث هو تقديم صياغة لبيانات ثنائية وتوسيعها الى بيانات متعددة المتغيرات في ظل تحليل لامعلمي للاقتران باستخدام طرق تجريبية في تقدير الدوال الامعلمية الاساسية في صياغة المخاطر المنافسة دون الحاجة الى وجود افتراضات لاجراء الاختبار .

مشكلة البحث The problem of Re search :

ان تقنيات تحليل البقاء الكلاسيكية هي غير ملائمة الاستخدام في بيانات المخاطر المنافسة ((compent risk data)) لما لهذه البيانات من خاصية انما يختار حدثا واحدا فقط لتحليل الاسباب المنافسة للفشل ،اما بقية الاسباب فهي تهمل وتعتبر مشاهدات مراقبة يمنية مما يؤدي الى تحيز في اهم طرق التقدير الا وهي طريقة كابلان – مير ((kaplan-meier)) وطريقة نيلسن – الن Nelson-Alen بسبب انتهاك احدى الفرضيات الاساسية للمخاطر المنافسة الا وهي فرضية استقلال اوقات الفشل والاستقلالية بين الاسباب المختلفة للفشل مما يجعل استخدام هذه الطرق مستحيلا .

الفكرة المقترحة proposed idea -:

الحصول على مقياس استدلالي كمي قادر على تمثيل احتمالات الفشل للإصابة بمرض السكري لبيانات المخاطر المنافسة وفي نفس الوقت يوصف الاقترانات الثنائية في حالة الإصابة او الموت عند اعمار معينة للاباء والابناء ضمن العائلة الواحدة التي لها استعداد جيني للإصابة باي مرض وراثي (مرض السكري) من خلال تطوير مقدرين لامعلميين لبيانات المخاطر المنافسة احدهما هو مقدر دالة الخطورة الثنائية التجميعية بسبب الإصابة بالمرض والثاني هو مقدر دالة الحادثة (المرض) الثنائية التراكمية من دون ان يكون هناك فرضيات حول استقلالية بيانات المخاطر المنافسة .

عينة البحث sample of re search -:

لقد تم استخدام عينة تتالف من 500 فردا بشكل ازواج (125) زوجا ذكورا واناثا ولمختلف الاعمار ممن لهم سابقة ظهور مرض السكري الوراثي في العائلة.وقد وجدت ان هذه البيانات تتبع توزيع كاما ذو معلمين حسب اختبار Goodness of fit و استخدمت بيانات مراقبة مستقلة من النوع الثاني يمين censoring type tow right data.

ان عينة البحث قد شملت جميع الاعداد ذكوراً واناثاً وحسب الجدول الاتي :-

العمر الحالة	0-10	10-20	20-30	30-40	40-50	50-60	60-70	70-	المجموع	
مصاب	ذ	28	16	33	45	41	23	20	6	216
	أ	16	15	15	21	24	10	7	-	108
غير مصاب	ذ	14	19	8	8	7	6	4	1	67
	أ	11	27	35	13	10	8	7	2	113
المجموع		77	91	87	87	52	38	38	9	500

مقدمة البحث :

تظهر بيانات المخاطر المتنافسة competing risks عادة في الدراسات التي يكون فيها فشل الفرد يصنف الى واحد من k حيث $(k \geq 1)$ من الاسباب المختلفة للفشل ، ولعل اهم مشكلة تواجه علم الامراض الوراثية هي في تحديد العوامل الوراثية التي تتفاوت من مريض لآخر فيحاول البحث دوماً عن الاقتران بين الجينات والامراض والصفات الوراثية من جيل الى اخر ولعل اهمها تلك الامراض الوراثية المرتبطة بالموت .

وفي هذا البحث ندرس مدى اقتران مثل تلك الامراض بين افراد العائلة الواحدة ضمن جيلين هما (الاباء و الابناء) ولماذا يوجد في اشخاص وينعدم في اشخاص اخرين لهم نفس الظروف ، وما هي العوامل المسببة في تعجيل زمن الفشل (ظهور المرض) عند اشخاص في اعمار مبكرة او مبكرة جداً وباعمار متاخرة عند اشخاص اخرين ضمن نفس نطاق العائلة واجراء اختبار حول اعتمادية هذا الزمن .

تحت فرضية الاستقلال يكون توزيع متعدد المتغيرات قابل للتشخيص وبانعدام وجود هذه الفرضية تصبح بيانات مخاطر المنافسة متعدد المتغيرات غير قابل للتشخيص وغير قابلة ايضاً للتقدير لامعلما بسبب وجود عدة اسباب مستترة أهمها عدم وجود الاستقلالية بين اوقات الفشل عندها يصبح التطبيق غير صحيح .

لذلك كان علينا صياغة بيانات ثنائية تتوسع الى بيانات متعددة المتغيرات باستخدام طرق وشكل دوال لامعلمية اساسية في صياغة بيانات المخاطر المنافسة وقادرة على تحليل الاقتران العائلي اللامعلمي لاي مرض جيني .

- وصف البيانات Data description : البيانات تتكون من n من الأزواج أباء و أبناء

في i^{th} عنقود لعينة تتوزع مشاهدتها (iid) بالشكل الاتي :

$$(Y_{i1}, \dots, Y_{in}, \eta_{i1}, \dots, \eta_{in}), \quad i = 1, \dots, n$$

مع وجود بيانات مراقبة مستقلة هي :

$$(C_1, \dots, C_n), \quad (T_1, \dots, T_n, \epsilon_1, \dots, \epsilon_n)$$

$$Y_j = T_j^{(1)}, \eta_j = 1 \quad T_1^{(1)} \quad \text{إذا كان الفرد مريضا}$$

$$Y_j = T_j^{(2)}, \eta_j = 2 \quad T_2^{(j)} \quad \text{إذا مات الفرد دون مرض}$$

$$Y_j = c_j, \eta_j = 0 \quad \text{إذا كان الفرد على قيد الحياة وخال من المرض}$$

علما ان : $T_j^{(1)}$ يمثل الاب او الام و $T_j^{(2)}$ يمثل الابن او البنت

هؤلاء الافراد هم ضمن عنقود خاضعين لنوعين من انواع الفشل ز حيث ان $J = 1, 2$ هما :

وقت الفشل المستتر المهم (وهو يمثل الاصابة بالمرض) $T_1^{(J)}$

وقت الفشل المستتر للخطورة المنافسة (ويمثل الموت) $T_2^{(J)}$

أن صيغة المشاهدة الواحدة تكون كالآتي : $TJ = \min(T_1^{(J)}, T_2^{(J)}) = T_1^{(J)} \wedge T_2^{(J)}$

$$\epsilon_j = I(T_j = T_1^{(j)}) + 2 I(T_j = T_2^{(j)}) \quad \text{for } J = 1, 2$$

حيث ان I هو مؤشر الدالة .

وبوجود بيانات مراقبة $C^{(1)}$ و $C^{(2)}$ صيغة البيانات المشاهدة ستكون بحسب الآتي :

$$sJ = I(T_j \leq C^{(j)})$$

$$Y_j = T_1^{(j)} \wedge T_2^{(j)} \wedge C^{(j)}$$

الشكل النهائي للمشاهدة في العنقود الواحد ستكون :

$$\eta_j = \epsilon_j \quad s_j = I(T_j = T_1^{(j)}) + 2 I(T_j = T_2^{(j)}) \quad I(T_j \leq C^{(j)})$$

وكما ذكرنا سابقا بسبب وجود ارتباط اوقات الفشل للأفراد ستكون غير قابلة للاختبار كما

ان توزيعاتها الحدية هي غير مشخصة لا معلما ، على هذا الاساس انصب التركيز حول

أي من الكميات يمكن ان تكون قابلة للتشخيص لا معلما وايهما يمكن الاستفادة منها و

تكيفها والاستفادة منها بشكل كبير في تطبيق الاختبار .

يعتمد تكوين بيانات مخاطر المنافسة على دالتين اساسيتين هما :

١- دالة خطورة سبب – معين التجميعية

٢- دالة الحادثة التجميعية

والسبب في ذلك يرجع الى ان هاتان الدالتان تصفان الزمن الاول للفشل T والحالة المسببة للفشل

وان صيغة دالة خطورة سبب – معين التجميعية هي :

$$\Lambda^{(j)}_k(t) = \int_0^t \lambda_k^{(j)}(s) ds \quad k, j = 1, 2$$

اما دالة الحادثة التجميعية F فتستخرج بالاعتماد على دالة الباقيين على قيد الحياة اللامعلمية وهي دالة بسيدو للبقاء :

$$\text{Exp} \{-\Lambda^{(j)}_k(t)\}, \quad k, j = 1, 2$$

عندئذ يمكن استخراج $Fk^{(j)}(s)$ وكالاتي :

$$Fk^{(j)}(s) = p(T_j \leq s, E_j = k) \quad k, j = 1, 2$$

وبتوسيع الدوال السابقة الاحاديتان الى دوال ثنائية فان دالة خطورة سبب – معين التجميعية الثنائية ستكون كالاتي :

$$\Lambda_{kl}(s, t) = \int_0^s \int_0^t \lambda_{kl}(u, v) du dv \quad k, j = 1, 2$$

و دالة الحادثة الثنائية التجميعية فهي :

$$Fkl(s, t) = p(T_1 \leq s, E_1 = k, T_2 \leq t, E_2 = L)$$

$$= \int_0^s \int_0^t \lambda_{kl}(u, v) S(\bar{U}, \bar{V}) du dv$$

اذ ان

$$S(u, v) = p(T_1 \geq u, T_2 \geq v)$$

تمثل دالة الباقيين على قيد الحياة الثنائية الكلية .

الا ان الذي يهمنا هنا هو λ_{11} و Λ_{11} و F_{11} لأنها تمثل الكميات المشخصة لاقتترانات سبب 1- ضمن العناقيد وان دالة الخطورة λ_{11} ودالة الخطورة التجميعية Λ_{11} هي كمية العنقدة في خطوات سبب 1 – الآنية عند النقطة (s, t) حيث s يمثل عمر عند الاصابة للاباء و t عمر الاصابة للابناء بشرط ان يكون الشخص على قيد الحياة عند (s, t) ، لذلك نجد ان اغلب اختبارات الاقتران يكون عن طريق F أكثر مما هي عن طريق Λ بسبب كونها الأكثر اظهارة وتشخيصا لأسباب ذلك المرض وبالوقت نفسه هي قابلة للمقارنة. الا ان هذه الكميات نفسها مع البيانات الثنائية تكون غير قابلة للتبديل اذ لا يمكن افتراض ان

$$\Lambda_1^{(1)} = \lambda_1^{(2)} \quad \text{or} \quad F_1^{(1)} = F_1^{(2)}$$

فالتحليلات اللامعلمية هي معقدة حتى مع السبب الفردي للفشل ولحسم هذه التحديات اوجدنا اسلوبا اخر في اجراء اختبار الاقتران أولا نسبة الخطورة البديلة وهي :

$$\Phi(s, t) = \Lambda_{11}(s, t) \{ \Lambda_1^{(1)}(s) \cdot \Lambda_1^{(2)}(t) \}^{-1}$$

وهي بالحقيقة تمثل المساهمة النسبية للمخاطر التجميعية لسبب — 1 للآباء والابناء مقسوما على المخاطر التجميعية للسبب — 1 الاحادية للآباء والابناء
وهي معرفة جيدا عند النقطة (s,t) عمر الاصابة بالمرض للآباء والابناء بسبب إن الكميات في Φ هي عبارة عن مجموع للخطورات الاحادية والثنائية .
الاجراء الثاني : اعتمدنا على نسبة دالة الحادثة الثنائية التجميعية

$$\Psi(s,t) = F11(s, t) \{ F1^{(1)}(s) \cdot F1^{(2)}(t) \}^{-1}$$

وهي معرفة ايضا عند نقاط الزمن لـ Φ نفسها .
اصبح الان اختبار استقلالية بيانات المخاطر المنافسة ممكنا وبالأخص استقلالية سبب — 1 (مرض السكري الوراثي) اذ يمكننا التعرف على وجوده عن طريق 11λ عندما تكون :

$$\lambda 11(s,t) = \lambda 1^{(1)}(s) \cdot \lambda 1^{(2)}(s) \quad \text{for all } s, t$$

$$\Phi(s, t) = 1 \quad \text{for all } s, t$$

وكذلك نفس الشيء فيما يخص F11 هو :

$$F11(s,t) = F1^{(1)}(s) \cdot F1^{(2)}(t) \quad \text{for all } s, t$$

$$\Psi(s, t) = 1 \quad \text{for all } s, t$$

إن كل من $\Phi(s,t)$ و $\Psi(s,t)$ مؤشران على نسبة الزيادة والنقصان للاحتمالات المشتركة للإصابة بالمرض للآباء والابناء عند زمن الاصابة بالمرض (s,t) فعندما يكون كل من Φ و Ψ مساويان إلى 1 الصحيح فهذا يعني ان كل من $(t1, \in 1)$ و $(t2, \in 2)$ هما مستقلين في النقطة (s,t) .

ولكن وبصورة عامة إن تكون $\Phi(s,t)=1$ هذا لا يعني بالضرورة إن $\Psi(s,t) = 1$ والعكس صحيح ويرجع هذا الاختلاف بسبب الاقتران في F11 أو بسبب تأثير بقية الاقترانات القوية في 12λ ، 22λ ، وهي عكس التي في 11λ .

بعد ذلك وللحصول على مقدرات لـ Φ و Ψ نستخدم طريقة التمهيد smoothing من خلال الممهد الخطي LLS الذي يعتبر افضل ممهد خطي لأنه يتكيف مع كل التصاميم الثابتة والعشوائية وعن طريق الملء (plug in) لـ Φ و Ψ من خلال برنامج bootstrap نحصل على :

$$\hat{\Phi} = \hat{\Lambda} 11, \hat{\Lambda} 1^{(1)} \hat{\Lambda} 1^{(2)} \}^{-1} \quad \& \quad \hat{\Psi} = \hat{F} 11 \{ \hat{F} 1^{(1)} \hat{F} 1^{(2)} \}^{-1}$$

ان المقدرين اعلاه هما معلمتي التمهيد band width وهما يؤثران تأثيرا كبيرا على تمهيد المنحنى المقدر وأقترابه الى المنحنى الحقيقي وأن اي زيادة لهاتين المعلمتين يؤديان الى

تعظيم التحيز وتصغير التباين والعكس صحيح وهنا معلمتي التمهيد Ψ^{\wedge} و Φ^{\wedge} تمثلان معلمتي التباين.

ولتوحيد المعلومات عبر الزمن لتكون اكثر اختصارا نكامل المعدلات الموزونة الى كل من Ψ^{\wedge} و Φ^{\wedge} :

$$\Phi^{\wedge*} = \frac{\int (t1,t2)w^{\wedge}(s,t)\Phi^{\wedge}(s,t)dsdt}{\int [t1,t2]w^{\wedge}(s,t)dsdt} = \int (t1,t2) w^{\wedge}(s,t) \Phi^{\wedge}(s,t) dsdt$$

$$\Psi^{\wedge*} = \frac{\int (T1,T2)w^{\wedge}(s,t)\Psi^{\wedge}(s,t)dsdt}{\int [T1,T2]w^{\wedge}(s,t)dsdt} = \int (t1,t2) w^{\wedge}(s,t) \Psi^{\wedge}(s,t) dsdt$$

حيث ان

$$\tilde{w}(s,t) = \frac{w^{\wedge}(s,t)}{\int [T1,T2]w^{\wedge}(s,t)dsdt} \quad \text{as } n \rightarrow \infty$$

وان $w^{\wedge}(s,t)$ تمثل دالة المعدل الموزون للاستجابات اللذان يساويان الواحد الصحيح تحت الاستقلالية لكل من النقطتين s, t ويمكن استخراج حدود الثقة للمقدر اللامعلمي F^{11} (وهو الاكثر قوة) عند مستوى دلالة $\alpha = 0.05$ من خلال حساب ولكل مجموعة بيانات $F^{11}(s,t)$ عند فترة الثقة 0.95 حسب الاتي :

$$\Phi^{\wedge}(s,t) = \text{logit} \{ F^{11}(s,t) \}$$

$$\Phi^{\wedge}(s,t) = \log \left\{ \frac{F^{11}(s,t)}{1 - F^{11}(s,t)} \right\}$$

حدود الثقة لكل مجموعة ستكون

$$\Phi^{\wedge}(s,t) \pm 1.96 \text{ Se } \{ \Phi^{\wedge}(s,t) \}$$

اذ ان $\{ \Phi^{\wedge}(s,t) \}$ هو الخطا المعياري المستخرج من bootstrap لعينة تتكون من n من الأزواج مع ارجاع وان مقدرات طريقة البوتستراب تتصف بانها متسقة ، موحدة وتقرب من عمليات (Gaussian) .

المحاكاة : simulation

- توليد ارقام عشوائية عن طريق RND (I) بين الصفر والواحد تتوزع توزيعا منتظما $U(0,1)$.
- توليد بيانات الالباء $T_2^{(1)}$, $T_1^{(1)}$ والابناء $T_2^{(2)}$, $T_1^{(2)}$ تتوزع توزيع كاما الضعيف بالمعلمتين $\alpha = \beta = \frac{1}{\sigma^2}$.

- ترتيب البيانات بشكل عناقيد وحسب طريقة الربط المفرد .
- توليد بيانات مراقبة مستقلة $C^{(1)}$ و $C^{(2)}$ من النوع الثاني يمين censor type two data تتوزع توزيع منتظم $U(0,2)$ بعدد $n = 250$.
- استخراج قيم الدوال الاحادية والثنائية لكل من λ و Λ و F .
- تقدير الدوال الثنائية لـ S و Λ_{11} و F_{11} عن طريق برنامج بوتستراب للحصول على المقدرات $S^{\Lambda\#}$ و $F^{\Lambda\#}_{11}$ و $\Lambda^{\Lambda\#}_{11}$.
- استخراج مقدر التباين للـ Φ^{\wedge} و Ψ^{\wedge} باستخدام دالة انفلونس التجريبية $I_i(s, t)$.

استخراج الخطأ القياسي لبوتستراب وحسب الصيغة :-

$$BMSE = \frac{\sum_{i=1}^N ei^2}{N-p}$$

- مقارنة الأخطاء القياسية لبوتستراب مع الأخطاء التجريبية للمقدر المقترح .

التوصيات :-

يفضل ان يكون اختبار استقلالية الإصابة بمرض السكري الوراثي عن طريق استخدام $\Psi(s, t)$ بدلا من $\Phi(s, t)$ لأنها تزيد من قوة الاختبار عند نقاط زمن الإصابة بالمرض عند اعمار معينة اي (s, t) .

المصادر :

- ١- علوان ، اسراء سعدون " مقارنة بين طريقتي simex و Lls لتقدير دالة الانحدار اللامعلمية باستخدام المحاكاة " أطروحة ماجستير جامعة بغداد / كلية الادارة والاقتصاد ٢٠٠٣ م .
- ٢- النزال ، رافد اسماعيل " طرق دراسة الـ jack knife والـ bootstrap اطروحة ماجستير جامعة بغداد / كلية الادارة والاقتصاد ١٩٨٩ م .
- ٣- *Bandeem – roche , k . and liang , k (2002) modeling multivariate failure time association in the presence of a competing risk . biometrika , 89, 299-314 .*
- ٤- *Dabrowska , D.M (1989) Kaplan Meier estimate on the plane : work convergence ,LIL and the bootstrap , J. multiv. Anal,29,308-25.*
- ٥- *Kosorok , M.R.(2006,to appear) introduction to empirical processes and semi parametric inference New York ;Springer – verlag .*

