اتساق تابع قرار الجوارات الـ k الأكثر قرباً باستخدام الطوريات العشوائية α

د. زباد قناية *

د. أحمد يونسو **

نور أزهرى * * *

(تاريخ الإيداع 30 / 2 / 2021. قُبل للنشر في 29 / 4 /2021

□ ملخّص □

يعد التصنيف الإحصائي من المواضيع المتقدمة في الإحصاء بطرائقه المختلفة ومن هذه الطرق طريقة النواة, المدرج التكراري و الجوارات الـk الأكثر قرباً المستخدمة في هذا البحث, ويعد موضوع دراسة اتساق توابع القرار من المواضيع التي تشغل العديد من الباحثين,حيث تم في دراسات سابقة دراسة اتساق قاعدة اتساق قاعدة الجوارات الـ k الأكثر قرباً في الحالة المستقلة, وعندما تكون العينة المدروسة مرتبطة يصبح من الضروري استخدام مفهوم المزج باستخدام معاملات المزج المختلفة, لذلك نهدف في هذا البحث إلى اثبات اتساق تابع الجوارات الـ k الأكثر قرباً في حالة الارتباط الضعيف أي سنثبت أن تابع الجوارات الـ k الأكثر قرباً متسق عندما تكون العينة التدريبية مشاهدات لطورية عشوائية α مزوجة أو (مزوجة بقوة) ,كما توصلنا من خلال محاكاة عينة تجريبية لدوال مختلفة أن قيمة k الأفضل (عدد الجوارات الأفضل) هو ثلاث جوارات , أوصت الدراسة بإجراء أبحاث لاحقة للحصول على أقوى تقارب لقوابع القرار لذلك كدراسة التقارب شبه الأكيد لتابع الجوارات الـ k الأكثر قرباً تحت شروط المزج الأخرى.

الكلمات المفتاحية: الاتساق, الجوارات الـ k الأكثر قرباً, المزج, التصنيف الموجه, تابع القرار,طورية عشوائية.

^{*} أستاذ مساعد, قسم الرباضيات, كلية العلوم, جامعة تشربن , اللاذقية , سوربة.

^{**} أستاذ مساعد, قسم الإحصاء الرباضي, كلية العلوم, جامعة دمشق , دمشق , سوربة.

^{** *} طالب دراسات عليا دكتوراه ,قسم الرياضيات, كلية العلوم , جامعة تشرين , اللاذقية , سورية.

Consistency the Decision Function of K-nearest Neighbors by Using α-mixing Random Stochastic

Dr. Ziyad Kanaya* Dr. Ahmad Younso** Nour Azhari***

(Received 30 / 2 / 2021. Accepted 29 / 4 /2021)

\square ABSTRACT \square

The statistical classification is one of the advanced topics in the statistics in its many ways, including the Kernel method, the Histogram and the k-nearest neighbors used in this research, and the subject of the study of the consistency of decision functions is one of the topics that concern many researchers, in previous studies has been studied the consistency of the k-nearest neighbors in the independent case, and when the sample studied is dependent it becomes necessary to use the concept of mixing using different mixing coefficients Therefore, in this research, we aim to prove the consistency of the k-nearest neighbors in the case of weak correlation, i.e. we will prove that the follow-up of the k-nearest neighbors is consistent when the training sample is α -mixing random stochastic, or strong mixing, we came through a simulation of an experimental sample of different functions that the best k value (the number of best neighborhoods) is three neighborhoods, the study recommended subsequent research to obtain stronger convergence of the decision function as a study of the almost certain convergence of the k-nearest neighbors under other mixing conditions.

Keywords: consistency, k-nearest neighbors, mixing, supervised classification, decision functions, random stochastic .

^{*}Associate Prof, Depart. Of Mathematics, Faculty of Science, Tishreen University, Lattakia, Syria. **Associate Prof, Depart. Of Mathematical Statistics, Faculty of Science, Damascus University, Damascus, Syria.

^{***}Postgraduate Student(Ph.D.), Depart. Of Mathematical Statistics, Tishreen University, Lattakia, Syria.

مقدمة:

التصنيف الإحصائي عملية إحصائية يتم فيها توزيع البيانات لمجتمع إحصائي إلى مجموعات مختلفة بناء على معلومات كمية تستند إلى واحدة او اكثر من الخواص الأساسية لهذه البيانات أو أعضاء المجتمع الإحصائي . تستند على عملية التصنيف هذه على خاصيات أصيلة في العناصر (التي قد تكون : رموزا أو متغيرات) و تستند على مجموعة تدريب, يستخدم التصنيف بشكل واسع في حل الكثير من المشكلات خاصة تلك التي تتعلق بالأعمال والطب والجرائم وغيرها من خلال تحليل مجموعة من البيانات ووضعها على شكل أصناف أو أقسام يمكن استخدامها فيما بعد لتصنيف البيانات مستقبلاً , وهناك عدد من الطرق التي يمكن استخدامها في تصنيف البيانات المستخدام خوارزميات مختلفة مثل الخوارزميات الإحصائية Statistical Algorithms والشبكات العصبية المستخدام والخوارزميات الجينية Rearest Neighbor Method وطريقة الجار الأقرب Network وبعد خوارزميات التعلم الآلي والتي تعمل بمشرف (موجه) وتعد خوارزمية الجار الأقرب من خوارزميات التصنيف التنبؤبة والوصفية.

وللتصنيف نوعان رئيسيان إما تصنيف موجه (Training Samples) وفي هذا النوع يتم مطابقة كل وحدة من الوحدات التجريبية مع العينات التدريبية باستخدام (Training Samples) وفي هذا النوع العديد من الطرق منها طريقة التصنيف الأكثر احتمالية (Likelihood Classification الحسابات الإحصائية ولهذا النوع العديد من الطرق منها طريقة التصنيف الأكثر احتمالية بعضها مع بعض وعلى أساسها يتم توزيع الوحدات التجريبية غير المعروفة إلى تلك الفئات المعروفة[1], أو تصنيف غير موجه أساسها يتم توزيع الوحدات التجريبية في الطريقة في حالة عدم توفر العينات التدريبية واساس عملها يعتمد على كون اي نوع من الاصناف الموجودة في الدراسة متكوناً من وحدات ذات قيم متقاربة مع بعضها، وتتضمن هذه الطريقة حسابات رياضية تختبر عددا كبيرا من الوحدات المجهولة وتقسيمها الى مجاميع معتمدة على القيمة الطيفية لكل وحدة من هذه الوحدات، وهناك تقنيات احصائية عديدة متوفرة بشكل برامج جاهزة بالإمكان استخدامها مع الحاسبات الالكترونية ومن ابرز هذه التقنيات عملية التحليل العنقودي Cluster Analysis. [2]

من أجل إجراء التصنيف الموجه نحتاج إلى بناء تابع قرار تعتمد كفاءته على قربه من دالة قرار أمثلية تسمى دالة قرار بايز (أي أن يكون الخطأ المرتكب أقرب ما يمكن إلى خطأ بايز) وسندرس في هذا البحث اتساق تابع قرار الجوارات الكثر قرباً عندما تكون العينة التدريبية مشاهدات لطورية عشوائية مزوجة.

مشكلة البحث:

تم بناء معظم نظريات الاحصاء التقليدي اعتماداً على متحولات عشوائية مستقلة وذلك بالاستفادة من مبرهنات النهاية المركزية وقانون الأعداد الكبيرة, إلا أنه في بعض الدراسات كالسلاسل الزمنية وبعض المقدرات التابعية يكون شرط الاستقلال غير محقق فقد تكون المتحولات العشوائية المدروسة مرتبطة وفق صيغة معينة , لذلك سندرس أسرة من الطوريات العشوائية المؤلفة من متحولات عشوائية مرتبطة احتمالياً وفق مفهوم معين يسمى المزج mixing ووفق هذا المفهوم تميل متحولات الطورية الى الاستقلال مع تباعد الزمن والتكرارات .ففي السلاسل الزمنية على سبيل المثال (دراسة تغيرات درجة الحرارة,دراسات الطقس...) يلاحظ كلما زادت الفروق الزمنية أو التكرارات بين المشاهدات قل تأثير المشاهدات السابقة على المشاهدات الحالية أو المستقبلية حتى أنه من أجل قفزات كبيرة للزمن أو المشاهدات

تميل الى أن تصبح مستقلة عن المشاهدات السابقة ومن هنا انطلقت العديد من الدراسات حول مفهوم المزج الذي له أشكال عديدة تتمثل مشكلة البحث في دراسة اتساق تابع قرار الجوارات الـk – الأكثر قرباً من تابع قرار بايز وذلك في حالة الارتباط أي عندما تكون العينة التدريبية مشاهدات لطورية عشوائية مزوجة.

أهمية البحث و أهدافه:

يهدف البحث الى دراسة اتساق قاعدة تابع قرار الجوارات الله - الأكثر قرباً من تابع قرار بايز وذلك في حالة الارتباط أي عندما تكون العينة التدريبية مشاهدات لطورية عشوائية مزوجة وذلك كتعميم لحالة الاتساق المثبت في الحالة المستقلة [3] ثم إجراء محاكاة باستخدام البرنامج الإحصائي R.

الإطار النظري:

الطوريات العشوائية المزوجة بالمفهوم $(\alpha - mixing)$ (4,5,6]:

لتكن $(Z_i, i \geq 1)$ طورية عشوائية معرفة على فضاء احتمالي (Ω, \mathcal{F}, P) وتأخذ قيمها في فضاء قيوس $(Z_i, i \geq 1)$, عندها نقول عن الطورية $(Z_i, i \geq 1)$ أنها $(Z_i, i \geq 1)$ أنها عندها نقول عن الطورية (أو مزوجة بقوة) إذا تحقق الشرط التالي:

على الترتيب يكون: $(Z_t, t=1,\dots,k), (Z_t, t=\ell+n,\dots)$ المولدان ب $\mathcal{F}_1^\ell, \mathcal{F}_{\ell+n}^{+\infty}$ على الترتيب يكون: $\alpha(n)=\sup_{\ell\geq 1}\sup_{A\in\mathcal{F}_1^\ell,B\in\mathcal{F}_{\ell+n}^{+\infty}}|P(A\cap B)-P(A)P(B)|\underset{n\to\infty}{\longrightarrow}0$ (1)

 $\alpha(n)=0, n\geq 0$ من الصفر كلما نزعت الطورية نحو الاستقلال وخصوصاً عندما $\alpha(n)$ من الصفر كلما نزعت الطورية نحو الاستقلال وخصوصاً عندما $\alpha(n)$ من العديد من نماذج $\alpha(n)$ تحقق شرط المزج (1) تحت شروط معينة [7] فمثلاً $\alpha(n)$ عرف بالشكل $\alpha(n)$ المديد من نماذج $\alpha(n)$ تحقق شرط المزج (1) تحت شروط معينة $\alpha(n)$ فمثلاً α

قاعدة الجوارات الـ k – الأكثر قرباً [3]:

لتكن $\{(X_i,Y_i),i\geq 1\}$ طورية عشوائية مستقرة بقوة (بمعنى أن الخصائص الاحتمالية للمتجهات تبقى ثابتة عند تغير الزمن بمقدار ثابت) معرفة على فضاء احتمالي (Ω,\mathcal{F},P) وتأخذ قيمها في الفضاء $\{0,1\}$, في التصنيف الموجه يسمى X_i متجه السمات (الخصائص) ويسمى X_i الصف الموافق لا X_i ونسعى من خلال هذا التصنيف إلى التنبؤ بالصنف X_i لمتجه السمات X_i الموافق لمشاهدة جديدة ليست من ضمن العينة.

بما أن الطورية مستقرة بقوة فرضاً يمكن اعتبار التوزيع الاحتمالي لأي شعاع (X_i,Y_i) مطابق لتوزع الشعاع (X,Y) بما أن الطورية مستقرة بقوة فرضاً يمكن اعتبار التوزيع الاحتمالي لل μ μ μ μ μ μ μ μ و μ تابع والذي يعرف جيداً من خلال μ و μ القيام μ أي من خلال μ عندما μ عندما μ يأخذ القيمة μ أي من خلال μ μ أي من خلال μ عندما μ عندما μ عندما μ يأخذ القيمة μ أي من خلال μ عندما μ عندما μ عندما μ أي من خلال μ عندما μ أي من خلال μ عندما μ عندما μ عندما μ أي من خلال μ عندما μ أي من خلال μ عندما μ أي من خلال μ عندما μ عندما μ عندما μ أي من خلال μ عندما μ أي من خلال μ عندما μ عندما μ عندما μ عندما μ أي من خلال μ أي من خلال أي من أي من خلال أي من أي م

من أجل إجراء التصنيف الموجه نحتاج إلى بناء تابع قرار $g:\mathbb{R}^d \to \{0,1\}$ بحيث g(X) تقابل صنف X وعندها نرتكب خطأ عندما $Y \neq g(X)$ حيث $Y \neq g(X)$ هو الصنف الفعلى لـ X.

نرمز لاحتمال الخطأ لتابع القرار g بg بg بg بg وإن احتمال هذا الخطأ يكون أصغرياً من أجل تابع القرار g^* المعرف بالشكل [3]:

$$g^*(x) = \begin{cases} 0 & if \ P(Y = 0/X = x) \ge P(Y = 1/X = x) \\ 1 & otherwise \end{cases}$$

يسمى $L^* = L(g^*)$ تابع قرار بايز ونرمز لاحتمال الخطأ الموافق بالرمز $g^*(x)$ ويسمى خطأ بايز .

لسوء الحظ تابع قرار بايز غير قابل للاستخدام مباشرة في التصنيف لأنه يعتمد على التوزيع الاحتمالي لـ $D_n = \{(X_i, Y_i), i = g^*(x) \text{ من خلال عينة عشوائية عشوائية على معلوم لذلك نسعى إلى تقدير <math>g^*(x)$ من خلال عينة عشوائية وبفضل هذه العينة التدريبية نعرف تابع قرار الجوارات الـ k - l الأكثر قرباً $g_n(x)$ كما يلى حيث k عدد صحيح موجب تماماً أي $l \geq 1$:

$$g_n(x) = \begin{cases} 0 & \text{if } \sum_{i=1}^n w_{n_i} Y_i \le \frac{1}{2} \\ 1 & \text{otherwise} \end{cases}$$

حيث $W_{n_i} = w_{n_i}(x, D_n)$ يساوي $w_{n_i} = w_{n_i}(x, D_n)$ و يساوي الصفر خلاف ذلك $k \to \infty$ يساوي $k \to \infty$ يساوي $k \to \infty$ الشرطين: $k \to \infty$ الشرطين: $k \to \infty$ الشرطين: $k \to \infty$ (فرضيات تقليدية في الحالة المستقلة[3]).

لنفرض أن X له كثافة احتمالية f حيث يمكن من خلال ذلك تجنب وقوع ربطات ناتجة عن تساوي بعد مجاورتين عن x , لنأخذ :

$$\eta_n(x) = \sum_{i=1}^n w_{n_i} Y_i$$

يسمى $\eta_n(x)$ التقدير الk مجاورة أكثر قرباً لدالة الانحدار $\eta(x)$ ولنرمز ب $\eta_n(x)$ بتابع القرار التجريبي واحتمال خطأ التصنيف له $L_n = L(g_n)$ ويسمى الخطأ التجريبي حيث يمكن أن نكتب:

$$g_n(x) = \begin{cases} 0 & if \ \eta_n(x) \le \frac{1}{2} \\ 1 & otherwise \end{cases}$$

أفضل ما نتوقعه من $g_n(x)$ هو أن يكون الخطأ أقرب ما يكون إلى خطأ بايز L^* وهذا يعبر عنه من خلال دراسة أشكال الاتساق المختلفة لـ $g_n(x)$ نحو تابع قرار بايز $g^*(x)$

تعريف الاتساق: نقول عن تابع القرار $g_n(x)$ أنه متسق إذا تحقق الشرط:

$$E(L_n) = \xrightarrow[n \to \infty]{} L^*$$

. L^* متسق عندما يكون متوسط الخطأ التجريبي $E(L_n)$ يتقارب عددياً نحو $g_n(x)$

سنثبت في هذه الورقة أن تابع الجوارات الـ k – الأكثر قرباً متسق عندما تكون العينة التدريبية مشاهدات لطورية عشوائية مزوجة .

قبل أن نتناول النتيجة الرئيسية لنقدم التمهيديات التالية:

تمهیدیة (1) یکن z_1, z_2 متحولین عشوائیین مستمرین بقیم حقیقیة ولنفرض أنهما محدودان عندئذ: $|cov(z_1, z_2)| \leq 4 \|z_1\|_{\infty} \|z_2\|_{\infty} \propto \{\sigma(z_1), \sigma(z_2)\}$

حبث:

$$\propto \{\sigma(z_1), \sigma(z_2)\} = \sup_{\substack{A \in \sigma(z_1) \\ B \in \sigma(z_2)}} |P(A \cap B) - P(A)P(B)|$$

معامل المزج بين الجبرين التامين $\sigma(z_1), \sigma(z_2)$ المولدين بـ z_1, z_2 على الترتيب.

تمهيدية (2) [9]: لتكن المجموعة:

$$B_a(x^{\setminus}) = \left\{ x \in \mathbb{R}^d; \mu(S_{x, \|x - x^{\setminus}\|}) \le a \right\}$$

 $x \in \mathbb{R}^d$ عندها من أحل كل

$$\mu(B_a(x^{\setminus})) \le a\gamma_a$$

حيث γ_a العدد الأصغري للمخاريط المتمركزة في المبدأ وبزاوية $\frac{\pi}{6}$ وتغطي γ_a .

: مبرهنة: بفرض أن العينة التدريبية D_n مشاهدات لطورية عشوائية α -مزوجة تحقق الشرط

: عيث $n o \infty$ الشرطين من أجل $au o \alpha(t) = O(t^{- heta})$

$$\frac{k}{\sqrt{n}} \to \infty$$

$$k \to \infty$$

حيث هنا الشروط على k أضعف من الشروط المقترحة في [8] عندئذ:

$$E(L_n) = \underset{n \to \infty}{\longrightarrow} L^*$$

قبل البدء بإثبات المبرهنة لنعرف بعض الرموز المهمة في عملية البرهان:

لنرمز بـ $ho_n=
ho_n(x)$ لحل المعادلة $ho_n=\mu(S_{x,
ho_n})$ (*) لحل المعادلة الم ولنعرف أيضاً:

$$\widehat{\eta_n}(x) = \frac{1}{k} \sum_{i=1}^n Y_i \mathbb{I}_{(X_i \in S_{x,\rho_n})}$$

$$i$$
=1 حيث A التابع المميز للمجموعة A حيث: \mathbb{I}_A حيث $\mathbb{I}_A(w)=\{1\ ; w\in A\ 0\ ; w\notin A$

بحسب المبرهنة في [3] يكفى أن نثبت أن:

$$E\int_{\mathbb{R}^d} |\eta(x) - \eta_n(x)| \mu(dx) \xrightarrow[n \to \infty]{} 0$$

لكن:

$$|\eta(x) - \eta_n(x)| \le |\eta(x) - E\widehat{\eta_n}(x)| + |E\widehat{\eta_n}(x) - \eta_n(x)| \quad (i)$$

بما أن $n o\infty$ فرضاً فإنه وباستخدام (*) يتحقق $ho_n o0$ عندما $ho_n o0$ وبحسب مبرهنة لوبيغ نحصل من أجل

$$E\widehat{\eta_n}(x) = \frac{1}{\mu(S_{x,\rho_n})} \int_{S_{x,\rho_n}} E(Y/X = x) \mu(dx)$$

$$\to E(Y/X = x) = \eta(x) \ \forall x \bmod \mathcal{M}$$

 $n \to \infty$ نحصل حسب مبرهنة التقارب الراجح عندما $|Y| \le 1$ بما أن

$$\int_{\mathbb{R}^d} |\eta(x) - E\widehat{\eta_n}(x)| \mu(dx) \xrightarrow[n \to \infty]{} 0 \qquad (ii)$$

 $n o \infty$ بحسب (ii) و (ii) بكفى أن نثبت أنه عندما

$$E\int_{\mathbb{R}^d}|E\widehat{\eta_n}(x)-\eta_n(x)|\mu(dx)\underset{n\to\infty}{\longrightarrow}0$$

ولدينا المتراجحة التالية:

$$E\int_{\mathbb{R}^d}|E\widehat{\eta_n}(x)-\eta_n(x)|\mu(dx)$$

$$\leq E \int_{\mathbb{R}^d} |E\widehat{\eta_n}(x) - \widehat{\eta_n}(x)| \mu(dx) + E \int_{\mathbb{R}^d} |\widehat{\eta_n}(x) - \eta_n(x)| \mu(dx)$$

 $n o \infty$ الحدين في يمين المتراجحة يتقاربان نحو الصفر عندما

من أجل الحد الأول وحسب متراجحة كوشي شوارتز:

$$E \int_{\mathbb{R}^d} |E\widehat{\eta_n}(x) - \widehat{\eta_n}(x)| \mu(dx) \le \int_{\mathbb{R}^d} \sqrt{E(E\widehat{\eta_n}(x) - \widehat{\eta_n}(x))^2} \mu(dx)$$

$$= \int_{\mathbb{R}^d} \sqrt{Var(\widehat{\eta_n}(x))} \mu(dx)$$

$$\le \int_{\mathbb{R}^d} \sqrt{\frac{n}{k^2} Var(Y\mathbb{I}_{(X \in S_{x,\rho_n})}) + C_n(x)} \mu(dx)$$

حيث:

$$C_n(x) = \frac{1}{k^2} \sum_{i \neq j} \left| cov\left(Y_i \mathbb{I}_{\left(X_i \in S_{x,\rho_n}\right)}, Y_j \mathbb{I}_{\left(X_j \in S_{x,\rho_n}\right)}\right) \right|$$

من جهة أخري لدينا:

$$\frac{n}{k^2} Var\left(Y \mathbb{I}_{\left(X \in S_{x,\rho_n}\right)}\right) \leq \frac{n}{k^2} E\left(\mathbb{I}_{\left(X \in S_{x,\rho_n}\right)}\right) = \frac{n}{k^2} \mu\left(S_{x,\rho_n}\right) = \frac{1}{k}$$

من جهة أخرى لدينا بحسب التمهيدية 1:

$$C_n(x) \le \frac{4}{k^2} \sum_{i \ne j} \alpha(|i - j|) \le \frac{4n}{k^2} \sum_{i=1}^{\infty} \alpha(i) \le \frac{4n}{k^2} C \sum_{l=1}^{\infty} i^{-\theta}, \theta > 1$$
$$\le C \frac{4n}{k^2} \underset{n \to \infty}{\longrightarrow} 0$$

 $\frac{k}{\sqrt{n}} \to \infty$ لأن $\infty \to \frac{k}{\sqrt{n}}$ فرضاً.

$$\Rightarrow E \int_{\mathbb{R}^d} |E\widehat{\eta_n}(x) - \widehat{\eta_n}(x)| \mu(dx) \underset{n \to \infty}{\longrightarrow} 0$$

$$\begin{aligned} |\widehat{\eta_n}(x) - \eta_n(x)| &= \left| \frac{1}{k} \sum_{i=1}^n Y_i \mathbb{I}_{(X_i \in S_{x,\rho_n})} - \frac{1}{k} \sum_{i=1}^n Y_i \mathbb{I}_{(X_i \in S_{x,r_n})} \right| \\ &\leq \frac{1}{k} \sum_{i=1}^n \left| \mathbb{I}_{(X_i \in S_{x,\rho_n})} - \mathbb{I}_{(X_i \in S_{x,r_n})} \right| \leq \left| \frac{1}{k} \sum_{i=1}^n \mathbb{I}_{(X_i \in S_{x,\rho_n})} - 1 \right| \\ &= |\widehat{\eta_n}(x) - E\widehat{\eta_n}(x)| \end{aligned}$$

 $\widehat{\eta_n}(x) = \frac{1}{k} \sum_{i=1}^n \mathbb{I}_{\left(X_i \in S_{x,\rho_n}\right)}$

وعليه يكفى أن نثبت أنه عندما $\infty \to \infty$ فإن:

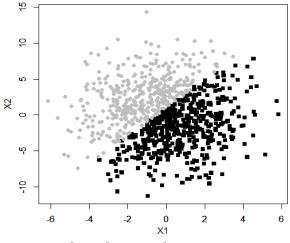
$$E\int_{\mathbb{R}^d} |\widetilde{\eta_n}(x) - E\eta_n(x)| \mu(dx) \underset{n\to\infty}{\longrightarrow} 0$$

لاحظ أن $\widetilde{\eta_n}(x)=\widehat{\eta_n}(x)$ عندما $Y_i=1$ من أجل $Y_i=1$ من أجل عندما والما قمنا به $\widetilde{\eta_n}(x)=\widehat{\eta_n}(x)$ من أجل من أجل من أجل الإثبات مشابه تماماً لما قمنا به سابقاً حول $\eta_n(x)$

محاكاة النتائج:

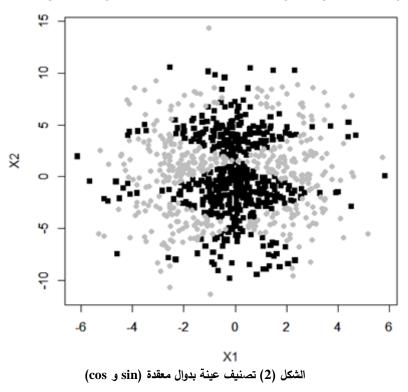
تم استخدام أسلوب المحاكاة بواسطة برنامج تم كتابته بلغة R [10,11] حيث نبين من خلاله آلية التصنيف وفق الخطوات التالية:

- 1. تابع لحساب مصفوفة المسافات بين الأدلة.
 - 2. تابع لحساب مصفوفة التباين والتغاير.
 - 3. محاكاة عينة حجمها ألف مرتبطة.
- 4. تقسيم العينة إلى عينة تدرببية حجمها 900 وعينة اختبار حجمه 100.
 - 5. تابع لحساب قاعدة التصنيف لقيمة واحدة.
 - 6. تابع لحساب قاعدة التصنيف لمجموعة قيم.
 - 7. رسم العينة التدريبية لدوال مختلفة.
- 8. حساب مجموع مربعات الخطأ لعينة الاختبار من أجل قيم مختلفة لـ k.
 - -من أجل دوال بسيطة (مثلاً كمعادلة مستقيم) لاحظ دقة التصنيف:

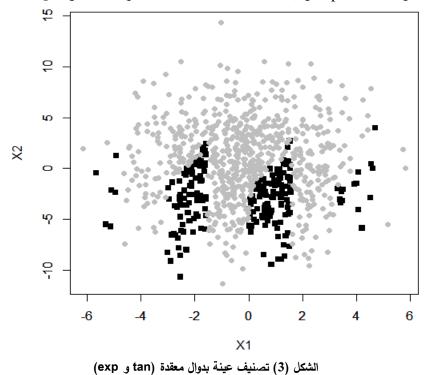


الشكل (1) تصنيف عينة بدوال بسيطة (معادلة مستقيم)

- من أجل دوال أكثر تعقيداً sin و cos وعينة حجمها 1000 مقسمة كما ذكرنا سابقاً يكون من أجل القرارين الشكل:



- من أجل دوال افتراضية tan و exp وعينة حجمها 1000 مقسمة كما ذكرنا سابقاً يكون من أجل القرارين الشكل:



وتكون قيم مجموع مربعات الخطأ من أجل قيم مختلفة له k والدوال السابقة كما هو موضح بالجدول(1) التالي:

البدول (1) سيم مجموع مربات العلق عليم المداد							
معادلة	k	1	3	5	7	9	33
مستقيم	MSE	1	3	2	2	2	1
sin و	k	1	3	5	7	9	33
cos	MSE	6	5	4	5	6	10
tan و	k	1	3	5	7	9	33
exp	MSF	7	5	7	8	8	16

الحدول(1) قيم محموع مربعات الخطأ لقيم مختلفة لـk

حيث نلاحظ أن أفضل قيمة لـ k في الحالة البسيطة 1 أو 33 وفي الحالة الأكثر تعقيداً (cos و sin) كانت عندما

```
. 3 هي k-3 وعندما (exp و k=3 كانت عندما k=3 وبالتالي يمكن اعتبار القيمة المثلي k=5
                                           وقمنا بكتابة تابع لتصنيف قيمة جديدة:
                                        CLASS = function(x, X1, X2, Y, k){
                                                       XX = cbind(X1, X2)
                                                     d=rep(0,length(X1))
                                                  for (i in 1:length(X1)){
                                           d[i]=sqrt(sum((x-XX[i,])^2))
                                                               ds = sort(d)
                                                               D=ds[1:k]
                                                              m=max(D)
                                                              dc=(d \le m)
                                                                gn=Y[dc]
                                                           gn=sum(gn)/k
                                                          Class=(gn>0.5)
                                                                  Class;}
                                 كما قمنا بكتابة تابع لتصنيف مجموعة قيم جديدة:
                                        CLASS = function(x, X1, X2, Y, k)
                                                      XX = cbind(X1, X2)
                                                     d=rep(0,length(X1))
                                                  for (i in 1:length(X1)){
                                           d[i]=sqrt(sum((x-XX[i,])^2))
                                                               ds = sort(d)
                                                               D=ds[1:k]
                                                              m=max(D)
                                                              dc=(d \le m)
                                                                gn=Y[dc]
                                                           gn=sum(gn)/k
                                                          Class=(gn>0.5)
                                                                    Class;
                                      MCLASS = function(x, X1, X2, Y, k)
```

n=nrow(x)h=rep(0,n)

for (i in 1:n){
h[i]=CLASS(x[i,],X1,X2,Y,k)}
h:}

وبالتالي من أجل أي قيمة جديدة ولتكن مثلاً 0.5 وعند k=12 مثلاً نستخدم التابع:

CLASS(0.5,X1,X2,Y,12)

تكون النتيجة true

و من أجل أي قيمة جديدة ولتكن مثلاً 6 وعند k=12 مثلاً نستخدم التابع:

CLASS(6,X1,X2,Y,12)

تكون النتيجة false

الاستنتاجات والتوصيات:

توصلنا في هذا البحث إلى اثبات اتساق تابع الجوارات الـ k – الأكثر قرباً في حالة الارتباط الضعيف أي أن تابع الجوارات الـ k – الأكثر قرباً متسق عندما تكون العينة التدريبية مشاهدات لطورية عشوائية مزوجة كما توصلنا من خلال محاكاة عينة تجريبية لدوال مختلفة أن قيمة k الأفضل (عدد الجوارات الأفضل) هو ثلاث جوارات , ونسعى في دراسات لاحقة للحصول على أقوى تقارب لتوابع القرار لذلك يمكن دراسة التقارب شبه الأكيد لتابع الجوارات الـ k – الأكثر قرباً تحت شروط المزج الأخرى.

References:

- [1]Li, N., Martin, A., & Estival, R. Heterogeneous information fusion: combination of multiple supervised and unsupervised classification methods based on belief functions(2020).
- [2] Ueda, R. M., Souza, A. M., & Menezes, R. M. C. P. How macroeconomic variables affect admission and dismissal in the Brazilian electro-electronic sector: A VAR-based model and cluster analysis. Physica A: Statistical Mechanics and Its Applications, (2020) 124872.
- [3] Luc Devroye ,Laszlo Gyorfi,Gabor Lugosi, *A Probabilistic Theoryof Pattern Recognition*, (1996) Springer-Verlag New York, Inc,P26-27.P78-79.
- [4]Rosenblatt, M. Remarks on some non-parametric estimates of the density function. Annals Math. Statist. 27, (1956) 832-837.
- [5]E. Rio, Th' eorie asymptotique des processus al' eatoires faiblement d' ependants. Springer Verlag, Berlin Heidelberg (2000).p40.
- [6] M. Rosenblatt, A central limit theorem and a strong mixing condition. Proc. Nat. Acad. Sci., USA 42 (1956) 43–47.
- [7] Bandyopadhyay, Soutir; *A NOTE ON STRONG MIXING*, Department of Statistics, Iowa State University, April 21, 2006.
- [8]H.C.P. Berbee, Random walks with stationary increments and renewal theory. Math. Cent. Tract. Amsterdam (1979).
- [9]Devroye, L. and Gyorfi, L. *Nonparametric Density Estimation: The L1 View.* John Wiley, New York(1985).
- [10]CRAWLEY J. M. The R book. 2nd. ed., John Wiley & Sons, Ltd., (2013) 1060.
- [11] COHEN, Y.; COHEN, J.Y. Statistics and Data with R: An Applied Approach Through Examples. A John Wiley & Sons, Ltd. (2008).