

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

د. زياد قناية*

د. أحمد يونسو**

نور أزهي***

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

□ ملخص □

يعد التصنيف الإحصائي من المواضيع المتقدمة في الإحصاء بطرائقه المختلفة ومن هذه الطرق طريقة النواة، المدرج التكراري و الجوارات الـ k 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 متنسق عندما تكون العينة التدريبية مشاهدات لطوريات 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، وأوصت الدراسة بإجراء أبحاث لاحقة للحصول على أقوى تقارب لتوابع 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، طوريات عشوائية.

* أستاذ مساعد، قسم الرياضيات، كلية العلوم، جامعة تشرين، اللاذقية، سورية.

** أستاذ مساعد، قسم الإحصاء الرياضي، كلية العلوم، جامعة دمشق، دمشق، سورية.

***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)

□ ABSTRACT □

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 أو أكثر من الخواص الأساسية لهذه البيانات أو أعضاء المجتمع الإحصائي . تستند عملية التصنيف هذه على خاصيات أصيلة في العناصر (التي قد تكون : رموزاً 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 باستخدام خوارزميات مختلفة مثل الخوارزميات الإحصائية Statistical Algorithms والشبكات العصبية Neural Network والخوارزميات الجينية Genetic Algorithms وطريقة الجار الأقرب Nearest Neighbor Method المستخدمة في هذا البحث وهي إحدى خوارزميات التعلم الآلي والتي تعمل بمشرف (موجه) وتعد خوارزمية الجار الأقرب من خوارزميات التصنيف التنبؤية والوصفية.

وللتصنيف نوعان رئيسيان إما تصنيف موجه Supervised Classification والذي يحتاج لوجود عينات تدريبية (Training Samples) وفي هذا النوع يتم مطابقة كل وحدة من الوحدات التجريبية مع العينات التدريبية باستخدام الحسابات الإحصائية ولهذا النوع العديد من الطرق منها طريقة التصنيف الأكثر احتمالية (Maximum Likelihood Classification) والتي تحسب مقدار الارتباط والاختلاف للعينات التدريبية بعضها مع بعض وعلى أساسها يتم توزيع الوحدات التجريبية غير المعروفة إلى تلك الفئات المعروفة [1]، أو تصنيف غير موجه Unsupervised Classification وتستخدم هذه الطريقة في حالة عدم توفر العينات التدريبية وأساس عملها يعتمد على كون أي نوع من الأصناف الموجودة في الدراسة متكوناً من وحدات ذات قيم متقاربة مع بعضها، وتتضمن هذه الطريقة حسابات رياضية تختبر عدداً كبيراً من الوحدات المجهولة وتقسيمها إلى مجاميع معتمدة على القيمة الطيفية لكل وحدة من هذه الوحدات، وهناك تقنيات إحصائية عديدة متوفرة بشكل برامج جاهزة بالإمكان استخدامها مع الحاسبات الإلكترونية ومن أبرز هذه التقنيات عملية التحليل العنقودي Cluster Analysis [2].

من أجل إجراء التصنيف الموجه نحتاج إلى بناء تابع قرار تعتمد كفاءته على قرينه 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 ال k- الأكثر قريباً عندما تكون العينة التدريبية مشاهدات لطوريات 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 التابعة 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 إلى 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 الـ k - الأكثر قرأً 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 الـ k - الأكثر قرأً 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 α (α - mixing) [4,5,6]:

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

من أجل الجبرين التامين $\mathcal{F}_1^\ell, \mathcal{F}_{\ell+n}^{+\infty}$ المولدان بـ $(Z_t, t = \ell + n, \dots), (Z_t, t = 1, \dots, k)$ على الترتيب يكون:

$$\alpha(n) = \sup_{\ell \geq 1} \sup_{A \in \mathcal{F}_1^\ell, B \in \mathcal{F}_{\ell+n}^{+\infty}} |P(A \cap B) - P(A)P(B)| \xrightarrow{n \rightarrow \infty} 0 \quad (1)$$

لاحظ أنه كلما اقتربت $\alpha(n)$ من الصفر كلما نزعنا الطورية نحو الاستقلال وخصوصاً عندما $\alpha(n) = 0, n \geq 0$ ويمكن التحقق من أن العديد من نماذج $ARMA$ تحقق شرط المزج (1) تحت شروط معينة [7] فمثلاً $AR(1)$ تعرف بالشكل $X_t = \rho X_{t-1} + \varepsilon_t ; |\rho| < 1$

قاعدة الجوارات الـ k - الأكثر قرأً [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 $\mathbb{R}^d \times \{0,1\}$ ، في التصنيف الموجه يسمى X_i متجه السمات (الخصائص) ويسمى Y_i الصف الموافق لـ X_i ونسعى من خلال هذا التصنيف إلى التنبؤ بالصف Y_j لمتجه السمات X_j الموافق لمشاهدة جديدة ليست من ضمن العينة.

بما أن الطورية مستقرة بقوة فرضاً يمكن اعتبار التوزيع الاحتمالي لأي شعاع (X_i, Y_i) مطابق لتوزيع الشعاع (X, Y) والذي يعرف جيداً من خلال μ و η حيث μ القياس الاحتمالي لـ X حيث $\mu(A) = P(x \in A)$ و η تابع انحدار Y على X عندما X يأخذ القيمة x أي من خلال $\eta(x) = E(Y/X = x)$.

من أجل إجراء التصنيف الموجه نحتاج إلى بناء تابع قرار $g: \mathbb{R}^d \rightarrow \{0,1\}$ بحيث $g(X)$ تقابل صف X وعندها نرتكب خطأً عندما $Y \neq g(X)$ حيث Y هو الصف الفعلي لـ X .

نرمز لاحتمال الخطأ لتابع القرار g بـ $L(g) = \mathbb{P}(Y \neq g(X))$ وإن احتمال هذا الخطأ يكون أصغرياً من أجل تابع القرار g^* المعروف بالشكل [3]:

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

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

لسوء الحظ تابع قرار بايز غير قابل للاستخدام مباشرة في التصنيف لأنه يعتمد على التوزيع الاحتمالي لـ (X, Y) والذي هو في الغالب غير معلوم لذلك نسعى إلى تقدير $g^*(x)$ من خلال عينة عشوائية $D_n = \{(X_i, Y_i), i = 1, \dots, n\}$ تسمى عينة تدريبية وبفضل هذه العينة التدريبية نعرف تابع قرار الجوارات ال k -الأكثر قريباً $g_n(x)$ كما يلي حيث k عدد صحيح موجب تماماً أي $k \geq 1$: [3]

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

حيث $w_{n_i} = w_{n_i}(x, D_n)$ يساوي $\frac{1}{k}$ عندما X_i واحدة من ال k الأقرب لـ x في D_n و يساوي الصفر خلاف ذلك حيث $k = k(n)$ متتالية من الأعداد الصحيحة الموجبة تماماً التي تحقق من أجل $n \rightarrow \infty$ الشرطين: $k \rightarrow 0, \frac{k}{n} \rightarrow 0$ (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 ، لتأخذ:

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

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

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

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:

$$E(L_n) \xrightarrow{n \rightarrow \infty} L^*$$

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

سنثبت في هذه الورقة أن تابع قرار الجوارات ال k -الأكثر قريباً متسق عندما تكون العينة التدريبية مشاهدات لطورية عشوائية مزوجة .

قبل أن نتناول النتيجة الرئيسية لنقدم التمهيديات التالية:

تمهيدية (1) [8]: ليكن Z_1, Z_2 منحولين عشوائيين مستمرين بقيم حقيقية ولنفرض أنهما محدودان عندئذ:

$$|cov(Z_1, Z_2)| \leq 4 \|Z_1\|_{\infty} \|Z_2\|_{\infty} \propto \{\sigma(Z_1), \sigma(Z_2)\}$$

حيث:

$$\propto \{\sigma(Z_1), \sigma(Z_2)\} = \sup_{\substack{A \in \sigma(Z_1) \\ B \in \sigma(Z_2)}} |P(A \cap B) - P(A)P(B)|$$

معامل المزج بين الجبرين التامين $\sigma(Z_1), \sigma(Z_2)$ المولدين بـ Z_1, Z_2 على الترتيب.

تمهيدية (2) [9]: لتكن المجموعة:

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

عندها من أجل كل $x^\setminus \in \mathbb{R}^d$:

$$\mu(B_a(x^\setminus)) \leq a\gamma_d$$

حيث γ_d العدد الأصغري للمخاريط المتمركزة في المبدأ وبزاوية $\frac{\pi}{6}$ وتغطي \mathbb{R}^d .

النتيجة الرئيسية:

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

$$\alpha(t) = O(t^{-\theta}) \quad \text{حيث } \theta > 1 \text{ إذا تحقق من أجل } n \rightarrow \infty \text{ الشرطين:}$$

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

$$k \rightarrow \infty$$

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

$$E(L_n) \xrightarrow{n \rightarrow \infty} L^*$$

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

لنرمز بـ $\rho_n = \rho_n(x)$ لحل المعادلة (*) $\frac{k}{n} = \mu(S_{x, \rho_n})$ حيث إن الحل موجود لأن X يملك تابع كثافة ولنعرف أيضاً:

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

حيث $\mathbb{1}_A$ التابع المميز للمجموعة A حيث:

$$\mathbb{1}_A(w) = \begin{cases} 1 & ; w \in A \\ 0 & ; w \notin A \end{cases}$$

اثبات المبرهنة:

بحسب المبرهنة في [3] يكفي أن نثبت أن:

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

لكن:

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

بما أن $\frac{k}{n} \rightarrow 0$ فرضاً فإنه وباستخدام (*) يتحقق $\rho_n \rightarrow 0$ عندما $n \rightarrow \infty$ وبحسب مبرهنة لوبيغ نحصل من أجل $n \rightarrow \infty$:

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

$$\rightarrow E(Y/X = x) = \eta(x) \quad \forall x \text{ mod } \mathcal{M}$$

بما أن $|Y| \leq 1$ نحصل حسب مبرهنة التقارب الراجح عندما $n \rightarrow \infty$:

$$\int_{\mathbb{R}^d} |\eta(x) - E\widehat{\eta}_n(x)| \mu(dx) \xrightarrow{n \rightarrow \infty} 0 \quad (ii)$$

بحسب (i) و (ii) يكفي أن نثبت أنه عندما $n \rightarrow \infty$:

$$E \int_{\mathbb{R}^d} |E\widehat{\eta}_n(x) - \eta_n(x)| \mu(dx) \xrightarrow{n \rightarrow \infty} 0$$

ولدينا المتراحة التالية:

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

لنثبت أن الحدين في يمين المتراحة يتقاربان نحو الصفر عندما $n \rightarrow \infty$.

من أجل الحد الأول وحسب متراحة كوشي شوارتز :

$$\begin{aligned} E \int_{\mathbb{R}^d} |E\widehat{\eta}_n(x) - \widehat{\eta}_n(x)| \mu(dx) &\leq \int_{\mathbb{R}^d} \sqrt{E(E\widehat{\eta}_n(x) - \widehat{\eta}_n(x))^2} \mu(dx) \\ &= \int_{\mathbb{R}^d} \sqrt{\text{Var}(\widehat{\eta}_n(x))} \mu(dx) \end{aligned}$$

$$\leq \int_{\mathbb{R}^d} \sqrt{\frac{n}{k^2} \text{Var}(Y \mathbb{1}_{(X \in S_{x, \rho_n})})} + C_n(x) \mu(dx)$$

حيث :

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

من جهة أخرى لدينا:

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

من جهة أخرى لدينا بحسب التمهيدية 1:

$$\begin{aligned} C_n(x) &\leq \frac{4}{k^2} \sum_{i \neq j} \alpha(|i - j|) \leq \frac{4n}{k^2} \sum_{i=1}^{\infty} \alpha(i) \leq \frac{4n}{k^2} C \sum_{l=1}^{\infty} i^{-\theta}, \theta > 1 \\ &\leq C \frac{4n}{k^2} \xrightarrow{n \rightarrow \infty} 0 \end{aligned}$$

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

$$\Rightarrow E \int_{\mathbb{R}^d} |E\widehat{\eta}_n(x) - \widehat{\eta}_n(x)| \mu(dx) \xrightarrow{n \rightarrow \infty} 0$$

الآن سنبرهن أن الحد الثاني في المتراحة المذكورة يسعى للصفر ولأجل ذلك لنرمز $X_{(k)}$ للمجاورة الك الأقرب لـ x

وليكن $r_n = r_n(x) = \|X_{(k)}(x) - x\|$ عندها واضح أن :

$$|\widehat{\eta}_n(x) - \eta_n(x)| = \left| \frac{1}{k} \sum_{i=1}^n Y_i \mathbb{1}_{(X_i \in S_{x,\rho_n})} - \frac{1}{k} \sum_{i=1}^n Y_i \mathbb{1}_{(X_i \in S_{x,r_n})} \right|$$

$$\leq \frac{1}{k} \sum_{i=1}^n \left| \mathbb{1}_{(X_i \in S_{x,\rho_n})} - \mathbb{1}_{(X_i \in S_{x,r_n})} \right| \leq \left| \frac{1}{k} \sum_{i=1}^n \mathbb{1}_{(X_i \in S_{x,\rho_n})} - 1 \right|$$

$$= |\widehat{\eta}_n(x) - E\widehat{\eta}_n(x)|$$

حيث $\widehat{\eta}_n(x) = \frac{1}{k} \sum_{i=1}^n \mathbb{1}_{(X_i \in S_{x,\rho_n})}$ وعليه يكفي أن نثبت أنه عندما $n \rightarrow \infty$ فإن:

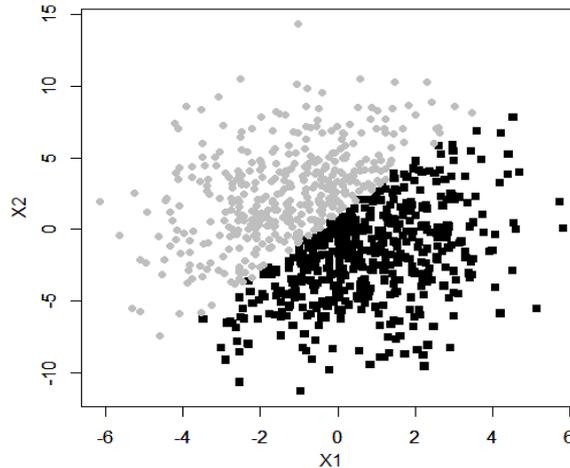
$$E \int_{\mathbb{R}^d} |\widehat{\eta}_n(x) - E\eta_n(x)| \mu(dx) \xrightarrow{n \rightarrow \infty} 0$$

لاحظ أن $\widehat{\eta}_n(x) = \eta_n(x)$ عندما $Y_i = 1$ من أجل $i = 1, \dots, n$ لذلك الإثبات مشابه تماماً لما قمنا به سابقاً حول $\eta_n(x)$.

محاكاة النتائج:

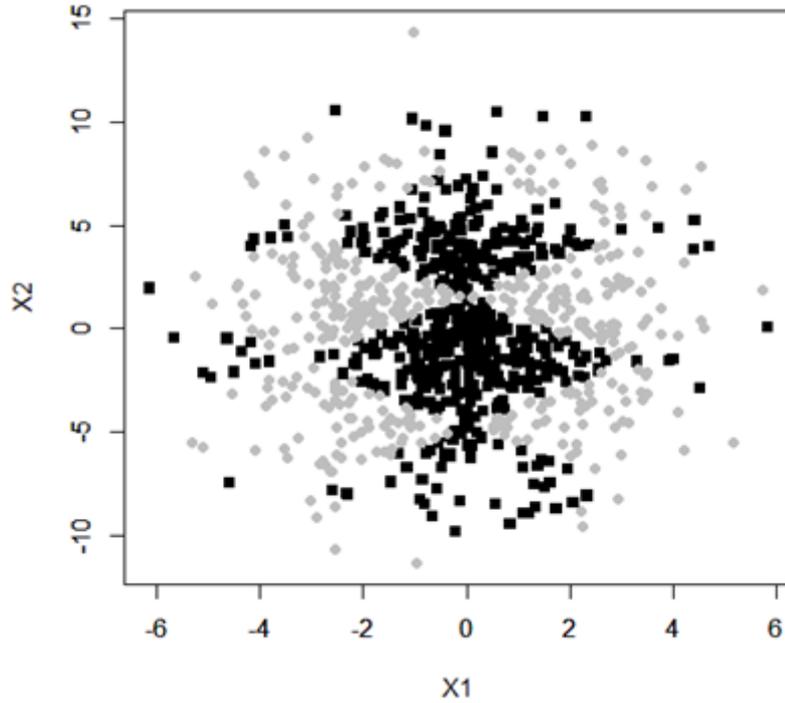
تم استخدام أسلوب المحاكاة بواسطة برنامج تم كتابته بلغة R [10,11] حيث نبين من خلاله آلية التصنيف وفق الخطوات التالية:

1. تابع لحساب مصفوفة المسافات بين الأدلة.
 2. تابع لحساب مصفوفة التباين والتغاير.
 3. محاكاة عينة حجمها ألف مرتبطة.
 4. تقسيم العينة إلى عينة تدريبية 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:



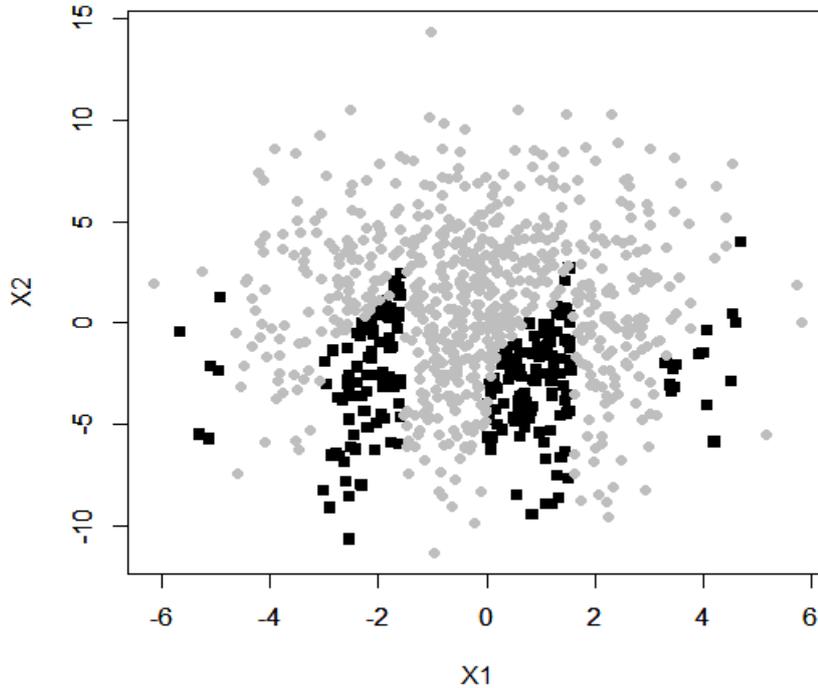
الشكل (1) تصنيف عينة بدوال بسيطة (معادلة مستقيم)

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



الشكل (2) تصنيف عينة بدوال معقدة (\cos و \sin)

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



الشكل (3) تصنيف عينة بدوال معقدة (\tan و \exp)

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

الجدول (1) قيم مجموع مربعات الخطأ لقيم مختلفة لـ k

33	9	7	5	3	1	k	معادلة
1	2	2	2	3	1	MSE	مستقيم
33	9	7	5	3	1	k	sin و
10	6	5	4	5	6	MSE	cos
33	9	7	5	3	1	k	tan و
16	8	8	7	5	7	MSE	exp

حيث نلاحظ أن أفضل قيمة لـ k في الحالة البسيطة 1 أو 33 وفي الحالة الأكثر تعقيداً (sin و cos) كانت عندما k=5 وعندما (tan و exp) كانت عندما k=3 وبالتالي يمكن اعتبار القيمة المثلى لـ k هي 3 .

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<=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<=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 ال k - الأكثر قرباً في حالة الارتباط الضعيف أي أن تابع الجوارات ال k - الأكثر قرباً متسق عندما تكون العينة التدريبية مشاهدات لطوريات 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 لتتابع القرار لذلك يمكن دراسة التقارب شبه الأكيد لتابع الجوارات ال 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 Gyorfı,Gabor Lugosi, *A Probabilistic Theory of 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 Gyorfı, 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).