

تقدير دوال الكثافة غير المعلمية باستخدام نواة إيباشينكوف : مقارنة نوى و عرض نطاق تكييفي مع تطبيقات تصنيفية باستخدام لغة R

الأستاذ الدكتور أحمد رستم الوسوف*

الدكتورة وفاء محمد كنعان**

سلافة مالك خازم***

(تاريخ الإيداع ٢٠٢٦ /١/١٨ - تاريخ النشر ٢٠٢٦ /٤/١٤)

□ ملخص □

قدمنا في هذا البحث تقدير دالة الكثافة الاحتمالية باستخدام الطرائق غير المعلمية، بوصفها أداة إحصائية حديثة تسمح لنا بالكشف عن البنية الحقيقية لتوزيع البيانات دون الحاجة الى افتراض مسبق مع التركيز على تقدير كثافة النواة (KDE) بنواة إيباشينكوف (Epanechnikov Kernel) نظراً لخصائصها التحليلية المتميزة التي تقلل من التباين مقارنة بنواة غاوس (Gaussian Kernel). تمت مقارنة أداء النواتين عبر معيار متوسط مربع الخطأ (MSE)، مع دراسة تأثير عرض النطاق (Bandwidth) (h) كعامل حاسم في جودة التقدير باستخدام قواعد silverman و LSCV للاختيار الأمثل، كما امتدت الدراسة إلى الطرق التكييفية وذلك باستخدام عرض نطاق ضيق في مناطق كثافة عالية وواسع في المناطق المنخفضة، والذي ساهم في تحسين التقديرات للبيانات غير المتجانسة. لقد تم تنفيذ الطريقتين الثابتة والتكييفية على بيانات حقيقية باستخدام لغة R ، واستخدام التقديرات في مهام التصنيف البايزي، حيث حقق عرض النطاق التكييفي تحسناً في MSE بنسبة 30-50% مقارنة بالثابتة ، وتبين نتائج المقارنات كفاءة KDE التكييفي في بناء نماذج تصنيف دقيقة تمثل هيكل البيانات المعقدة.

الكلمات المفتاحية : تقدير دالة الكثافة - نواة إيباشينكوف - عرض النطاق التكييفي - الحزمة الإحصائية R

*أستاذ في قسم الإحصاء الرياضي -كلية العلوم-جامعة اللاذقية-اللاذقية-سوريا

**مدرس في قسم الإحصاء الرياضي-كلية العلوم-جامعة اللاذقية-اللاذقية-سوريا

***طالبة دراسات عليا (ماجستير) قسم الإحصاء الرياضي-كلية العلوم-جامعة اللاذقية-اللاذقية-سوريا

Nonparametric Density Function Estimation Using the Epanechnikov Kernel: A Comparison of Kernels and Adaptive Bandwidth with Classification Applications Using R

Prof. Ahmed R. Alwassouf*

Dr. Wafaa M. Kanaan**

Sulafa M. Khazem***

(Received 18/1/2026. Accepted 14/4/2026)

□ABSTRACT □

In this study, we present the estimation of the probability density function using nonparametric methods, as a modern statistical tool that enables uncovering the true underlying structure of data distributions without the need for prior assumptions. The focus is placed on kernel density estimation (KDE) using the Epanechnikov kernel, owing to its superior analytical properties that lead to reduced variance compared to the Gaussian kernel. The performance of the two kernels was evaluated using the mean squared error (MSE) criterion, with particular emphasis on the bandwidth parameter (h) as a crucial factor affecting estimation quality. Optimal bandwidth selection was investigated using Silverman's rule of thumb and the least squares cross-validation (LSCV) method. Furthermore, the study was extended to adaptive approaches, where narrower bandwidths were employed in regions of high data density and wider bandwidths in low-density regions, resulting in improved estimates for heterogeneous data. Both fixed and adaptive methods were implemented on real datasets using the R programming language, and the resulting density estimates were applied to Bayesian classification tasks. The results demonstrated that adaptive bandwidth selection achieved a 30–50% reduction in MSE compared to the fixed bandwidth approach. Overall, the comparative analysis confirms the efficiency of adaptive KDE in constructing accurate classification models that effectively capture complex data structures

Keywords: Density function estimation - Epanechnikov's kernel - Adaptive Bandwidth -R statistical package

*Prof, Depart. of Mathematical Statistics, Faculty of Science, University of Latakia , Latakia ,Syria.

**Lecturer, Depart. of Mathematical Statistics , Faculty of Science , University of Latakia ,Latakia , Syria.

***postGraduate Student (Master's degree), Depart. of Mathematical Statistics, Faculty of Science, Latakia University, Latakia, Syria.

1- المقدمة:

يُعدّ التحليل الإحصائي من الأدوات الأساسية في تحليل البيانات، لما له من دور محوري في فهم خصائص البيانات واستخلاص المعلومات الكامنة فيها، بما يسهم في دعم عمليات اتخاذ القرار في مختلف المجالات العلمية والتطبيقية. وتنقسم الأساليب الإحصائية بصورة عامة إلى نوعين رئيسيين: الأساليب المعلمية والأساليب اللامعلمية. ورغم كفاءة الأساليب المعلمية تحت افتراضاتها، فإن البيانات مجهولة التوزيع تتطلب تقديراً غير معلمي مثل KDE (Mohamed&Ibrahem , 2008)، حيث تركز الدراسة على نواة إيباشنكوف لكفاءتها المثلى، مع دراسة عرض نطاق (h) عبر Silverman's و Least Squares Cross-Validation (LSCV) والتكيفية (Abramson) وتطبيقها على برنامج R على بيانات حقيقية للتصنيف، حيث يُعدّ تقدير دالة الكثافة باستخدام دوال النواة (Kernel) Density Estimation (Boli,2024 ; Guidoum,2024 ; Parzen, 1962 ; Siloko et al, 2020) من أبرز الأساليب غير المعلمية المستخدمة لتقدير دالة الكثافة الاحتمالية لمتغير عشوائي اعتماداً على عينة من البيانات. تقوم هذه الطريقة على استخدام دالة نواة تعمل كدالة وزن، تُسهم في تقدير الكثافة عند كل نقطة استناداً إلى قرب المشاهدات منها. حيث يوفّر تمثيل مرّن وسلس لتوزيع البيانات دون الحاجة إلى فرض نموذج إحصائي محدد. ومن هذه الدوال، تبرز نواة إيباشنكوف (Epanechnikov Kernel) (Turlach,2013)، بوصفها من أكثر النوى كفاءة إحصائياً، فهي دالة محدودة ومتناظرة، وتمنح أوزاناً أعلى للملاحظات القريبة من نقطة التقدير، فضلاً عن قدرتها على تقليل متوسط مربع الخطأ (MSE) نظراً لامتلاكها أقل تباين ضمن شروط رياضية معينة. وبناءً على ذلك، تحظى هذه النواة باهتمام خاص عند مقارنة أدائها بنوى أخرى. ولعرض النطاق (Bandwidth) (Guidoum,2024) أهمية عظمى، إذ يُعدّ العامل الأكثر تأثيراً في جودة تقدير الكثافة، لأنه يتحكم بدرجة التعميم وتحقيق التوازن بين التحيز والتباين. وفي هذا السياق، نناقش آليات اختيار عرض النطاق الأمثل، ومن بينها قاعدة سيلفرمان (Zambom&Dias, 2012) (Silverman's Rule) وتأثيرها المباشر في دقة النتائج. وانطلاقاً من محدودية النواة ذات عرض النطاق الثابت في تمثيل البيانات غير المتجانسة، يتوسع هذا البحث في دراسة النواة التكيفية (Adaptive Kernel)، التي تعتمد على تعديل عرض النطاق وفقاً للكثافة المحلية للبيانات، وتُعد هذه الاستراتيجية أكثر قدرة على تمثيل التوزيعات المعقدة و متعددة القمم، وتم تطبيق تقدير دالة الكثافة باستخدام كل من نواة إيباشنكوف ونواة غاوس (Popve et al,2024 ; Qin&Huang,2025) على بيانات مولدة عشوائياً تتبع التوزيع الطبيعي، كما تم تطبيق تقدير الكثافة على بيانات حقيقية لأوزان الأسماك باستخدام النواتين، إضافة إلى دراسة تقدير الكثافة مع إدخال متغير الطول. وقد نفذنا جميع التحليلات باستخدام برنامج الحزمة الإحصائية R (زينه ، 2017). ولا يقتصر البحث على الجانب الوصفي لتقدير الكثافة، بل يمتد إلى توظيف هذه التقديرات في مهام التصنيف الإحصائي، من خلال بناء قواعد قرار تعتمد على تقديرات دوال الكثافة واحتمالات الانتماء للفئات، وأظهرت النتائج تفوق عرض النطاق التكيفي في تحسين أداء المصنّفات، مما يؤكد كفاءة أساليب تقدير الكثافة اللامعلمية في التعامل مع الهياكل المعقدة للبيانات وبناء نماذج دقيقة ومرنة.

2- أهمية البحث:

تكمن أهمية هذه الدراسة في تطوير تقدير كثافة النواة باستخدام نواة إيباشنكوف التي توفر كفاءة مثلى في تقليل التباين بنسبة تصل إلى 15% مقارنة بالنواة الغاوسية خاصة مع البيانات غير المتجانسة، كما تساهم الطرق التكيفية (Abramson) في تحسين دقة التصنيف البايزي بنسبة 30-50% عبر MSE، مما يعزز تطبيقاتها في التعلم

الآلى،الاقتصاد، والطب الحيوي حيث تتفوق على النماذج المعلمية التقليدية، بالإضافة إلى ذلك يوفر التنفيذ العملي ب R أداء عملياً للباحثين في الإحصاء التطبيقي.

3-أهداف البحث:

يهدف هذا البحث إلى:

- 1- تقدير دالة الكثافة الاحتمالية اللامعلمية بطريقة النواة باستخدام نواة ايباشينكوف مع عرض نطاق ثابت ومتكيف، بطريقة سيلفرمان (Epanechnikov Kernel) مقابل الغاوسية (Gaussian kernel) في تقدير KDE عبر معيار MSE على بيانات حقيقية.
- 3-دراسة تأثير عرض النطاق(h) باستخدام قواعد Silverman و LSCV لتحقيق توازن بين الانحياز والتباين .
- 4-تطبيق KDE التكييفية لتحسين التقديرات في المناطق المنخفضة والكثيفة، مع قياس التحسن في التصنيف اللامعلمي.
- 5-تنفيذ عملي ب R لإنتاج نماذج تصنيف دقيقة تمثل هياكل البيانات غير المعلومة.

4-طرئق البحث ومواده :

1-4:الجانب النظري.

• دالة النواة:

التعريف والغرض:

تعرف دالة النواة بالعلاقة: (Hansen,2009)

$$(1) \quad K(x): R \rightarrow R$$

وتستخدم دالة النواة ، التي غالباً ما يُشار إليها بـ " دالة وزن " ، لقياس وزن جزئي في التقدير غير المعلمي ، خاصة في تقدير دوال الكثافة الاحتمالية ودوال التوزيع المتوقعة. الغرض الأساسي منها هو تقييم البيانات عن طريق تخصيص أوزان للمشاهدات ضمن نطاق معين ، مما يقلل بشكل فعال من التشويش ويكشف عن الأنماط الكامنة .

• الخصائص الأساسية لدالة النواة $K(u)$:

$$\int_{-\infty}^{\infty} u^2 K(u) du \neq 0 \quad , \quad \int_{-\infty}^{\infty} u K(u) du = 0 \quad , \quad \int_{-\infty}^{\infty} K(u) du = 1$$

0

حيث: $K(u)$: دالة نواة

• أشكال دوال النواة :

لكثافة النواة أشكال متعددة ، ولكل منها خصائصه ودرجة ملاءمته لمجالات مختلفة.
كما هو موضح في الجدول (1) الذي يعرض بعض الأشكال الشائعة لدوال النواة:
الجدول(1)أمثلة على أشكال دوال النواة المتماثلة

Kernel	Definition
Epanechnikov (Wand&Jones,1995)	$K(u)=\begin{cases} \frac{3}{4}(1-u)^2 & \text{for } u < 1 \\ 0 & \text{for } u \geq 1 \end{cases}$
Triangular (Wand&Jones,1995)	$K(u)=\begin{cases} 1- u & \text{for } u < 1 \\ 0 & \text{for } u \geq 1 \end{cases}$
Gaussian (Wand&Jones,1995)	$K(u)=\frac{1}{\sqrt{2\pi}}e^{-u^2/2}$

سنركز في هذا البحث على استخدام نواة إيباشينكوف (Epanechnikov kernel)، فنحلل البيانات ونقارن أداؤها مع أداء نواة غاوس (Gaussian kernel) المستخدمة على نطاق واسع، باستخدام عروض نطاق مختلفة ، مما يسمح من إلقاء الضوء على نقاط القوة والضعف وفق أهداف هذا البحث.

2-4: مقدر الكثافة بالنواة (مقدر النواة) :

ويُسمى أيضا تقدير دالة كثافة النواة (Kernel Density Estimation – KDE)

• التعريف والغرض:

مقدر الكثافة بالنواة هو أسلوب غير معلمي لتقدير دالة الكثافة الاحتمالية $f(x)$ لمتغير عشوائي.
بناء على عينة x_1, \dots, x_n دون افتراض شكل توزيع مسبق. و يعتبر KDE أداة قوية لاستكشاف البيانات وتصورها، حيث يوفر تمثيلاً سلساً ومستمرًا لتوزيع البيانات الأساسي.

• كيف يعمل مُقَدِّر الكثافة بالنواة؟

يعمل مقدر الكثافة بالنواة بوضع "نواة" (دالة وزن) عند كل نقطة من نقاط البيانات، ثم جمع هذه النوى لإنشاء تقدير سلس لدالة الكثافة، يمكن تصور كل نواة كنقطة توزيع صغيرة تساهم في الشكل العام لدالة الكثافة. (حالة أحادية البعد أي متغير واحد) ، يتم تحديد شكل وحجم كل نواة بواسطة دالة النواة المختارة ومعلمة النطاق (bandwidth). وتعطى الصيغة العامة لمقدر الكثافة بالنواة بالعلاقة: (Wand& Jones, 1995)

$$(2) \quad \hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$$

حيث: h هو عرض النطاق (bandwidth) ، n : حجم العينة ، $K(\cdot)$: دالة نواة ، x_i : المشاهدات رقم i من

العينة

أما في حالة ثنائية البعد ، فقد تم تقدير دالة الكثافة الاحتمالية ثنائية البعد بمتغيرين، بفرض لدينا متغيرين عشوائيين X و Y مع عينة $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. فإن مقدر كثافة النواة في الحالة الثنائية يعطى بالعلاقة: (Scott , 1992).

$$(3) \quad \hat{f}(x, y) = \frac{1}{h_x h_y n} \sum_{i=1}^n K\left(\frac{x-x_i}{h_x}, \frac{y-y_i}{h_y}\right)$$

n : حجم العينة ، $K(\cdot)$: دالة النواة ثنائية المتغير ، h_x, h_y : عرض الحزمة لكل متغير ،
مشاهدات العينة. (x_i, y_i)

أيضاً يتم وضع دالة نواة ثنائية البعد حول كل نقطة بيانات ثم جمع هذه النوى للحصول على تقدير ناعم للكثافة الاحتمالية المشتركة ، وذلك باستخدام كل من نواتي (غاوس وإيباشينكوف) الثنائية التي لهما الشكل على الترتيب :

$$(4) \quad K(u, v) = \frac{1}{2\pi} e^{-1/2(u^2+v^2)}$$

$$(5) \quad K(u_1, u_2) = \frac{2}{\pi} (1 - u_1^2 - u_2^2) \quad \text{حيث } u_1^2 + u_2^2 \leq 1$$

3-4: اختيار معلمة عرض النطاق (Selection Bandwidth):

تلعب h دوراً هاماً في التوازن بين التحيز (Bias) و التباين (Variance) لتقليل الخطأ الكلي في التقدير.

• تأثير عرض النطاق (h) على التقدير:

يلعب عرض الحزمة او مايسمى (نافذة التمهيد) (h) أو معامل التنعيم (Bandwidth) دوراً مهماً في تقدير دالة الكثافة الاحتمالية باستخدام طريقة النواة (KDE)، إذ تتحكم بدرجة نعومة أو تذبذب دالة الكثافة المقدرة. فإذا كانت قيمة (h) صغيرة جداً فإن التقدير يتبع البيانات بشكل دقيق ويُظهر الكثير من التذبذبات، مما يؤدي إلى تباين (Variance) عالٍ في التقدير مع انخفاض التحيز (Bias). أما إذا كانت قيمة (h) كبيرة فإن التقدير يصبح أكثر نعومة ويقل التذبذب، لكن ذلك قد يؤدي إلى فقدان بعض خصائص التوزيع الحقيقي، وبالتالي يزداد التحيز (Bias) ويقل التباين. لذلك توجد مفاضلة بين التحيز والتباين (Bias-Variance Trade-off) عند اختيار قيمة (h)، حيث إن اختيار قيمة مناسبة لها يساعد على الحصول على تقدير متوازن لدالة الكثافة الاحتمالية، ويؤثر بشكل مباشر في دقة نتائج التصنيف عند استخدام طرق النواة في التقدير غير المعلمي ، والجدول التالي يوضح تأثير عرض النطاق على تقدير الكثافة .

جدول (2) تأثير عرض النطاق h على تقدير دالة الكثافة $f(x)$

الوصف الاحصائي	التباين	التحيز	الشكل الناتج للتقدير	التسمية	قيمة h
يعكس الضوضاء ويؤدي لتمثيل غير دقيق	مرتفع	منخفض جداً	تقدير حاد ومتذبذب	تنعيم ضعيف	$h \rightarrow 0$
تمثيل غير دقيق للخصائص الدقيقة	منخفض	مرتفع	تقدير ناعم ومسطح	تنعيم مفرط	$h \rightarrow \infty$
توازن بين التحيز والتباين وتمثيل أدق	منخفض	منخفض	تقدير متوازن	تنعيم أمثل	h مناسب

4-4: طرق تحديد عرض النطاق (h): (Guidoum,2024)

يوجد عدة طرق لتحديد قيمة (h)، سنختار منها طريقة سلفرمان والتحقق المتبادل (LSCV) .A قاعدة سيلفرمان (Silverman's Rule):

تعد قاعدة سيلفرمان إحدى القواعد الإرشادية الشائعة والبسيطة لتحديد عرض النطاق الأمثل. تعتمد هذه القاعدة على افتراض أن البيانات تتبع تقريباً توزيعاً طبيعياً. صيغتها هي:

$$(6) \quad h = 0.9 \cdot \min\left(\hat{\sigma}, \frac{IQR}{1.34}\right) \cdot n^{-1/5}$$

حيث:

- $\hat{\sigma}$: الإنحراف المعياري المقدر للبيانات.
 - IQR : المدى الربيعي للبيانات، وهو الفرق بين الربع الثالث والربع الأول ($Q_3 - Q_1$).
- يستخدم IQR كمقياس للمدى عندما تكون البيانات تحتوي على قيم شاذة، لأنه أقل حساسية لهذه القيم من الانحراف المعياري.

• n : حجم العينة.

• العدد 1.34 هو ثابت تطبيقي يضمن أن $IQR/1.34$ يعادل الإنحراف المعياري لتوزيع طبيعي.

B. التحقق المتقاطع (Cross-Validation-CV):

تعتبر طريقة التحقق المتقاطع أكثر تطوراً وتعتمد على مبدأ تقليل خطأ التقدير. الهدف هو اختيار قيمة h التي تقلل دالة الخطأ، ولها أنواع سنختار (التحقق المتقاطع للمربعات الصغرى) Least Squares Cross-Validation (LSCV)، الصيغة الأكثر شيوعاً لـ LSCV تتضمن تكامل مربع التقدير وطرح تقديرات النقطة الواحدة. أي لتقليل متوسط خطأ التكامل التريبيعي (MISE)، حيث يتم تقدير معلمة عرض الحزمة من خلال طريقة النواة حيث يتم استبعاد مشاهدة واحدة من المشاهدات لذلك تسمى أيضاً بطريقة (Leave_One Out_Method) حيث يتم حساب دالة النواة $K(x)$ ، ثم نقوم بحساب معيار التحقق المتبادل من خلال استبعاد مشاهدة واحدة من المشاهدات المتغير. (محمد & إبراهيم، 2020)

$$(7) \quad CV(h) = \int_{-\infty}^{+\infty} \hat{f}_h(x)^2 dx - \frac{2}{n} \sum_{i=1}^n \hat{f}_{h,-1}(x_i)$$

حيث: $\hat{f}_{h,-1}(x_i)$ هو تقدير الكثافة عند x_i بعد حذف x_i من العينة.

ثم نعيد حساب دالة النواة ومعيار التحقق لجميع المشاهدات حيث يتم استبعاد مشاهدة في كل مرة.

ثم نقوم بعد ذلك بحساب قيمة معلمة عرض الحزمة المثلى h_{cv} :

$$h_{cv} = \operatorname{argmin} CV(h) \quad (*)$$

حيث المعادلة (*) تمثل معلمة عرض الحزمة التي تقابل أصغر CV_h و argmin تعني القيمة التي تقلل (أو تعطي الحد الأدنى ل) الدالة. هذه الطريقة أكثر تكلفة حسابياً ولكنها غالباً ما تعطي نتائج أفضل في الممارسة العملية لأنها لا تفترض شكلاً معيناً للتوزيع،

وفي حالة ثنائية البعد بمتغيرين يتم تطبيق كلا الطريقتين (سيلفرمان والتحقق المتقاطع) ويكون لعرض

الحزمة الشكل المصفوفي الآتي: (Scott, 1992)

$$(8) \quad h = \begin{pmatrix} h_1^2 & 0 \\ 0 & h_2^2 \end{pmatrix}$$

h_1 : عرض الحزمة للمتغير الأول X ، h_2 : عرض الحزمة للمتغير الثاني Y

5-4: الطرق التكيفية الحديثة (Adaptive Methods):**• تقدير الكثافة بعرض نطاق متكيف: الأنواع، الصيغة، الأهمية:**

يُعدّ عرض النطاق المتكيف (Bandwidth Adaptive) من أبرز التطورات المنهجية في تقدير الكثافة الاحتمالية باستخدام طرق النواة، إذ تسمح بتغيير عرض النطاق h وفقاً للكثافة المحلية للبيانات، وتعتمد هذه الطريقة على فكرة أساسية مفادها أن البيانات غير متجانسة بطبيعتها، وأن المناطق ذات الكثافة العالية تحتاج إلى عرض نطاق ضيق للحصول على تقدير أكثر دقة، بينما تحتاج المناطق قليلة الملاحظات إلى عرض نطاق أوسع لتقليل التذبذب.:

• أنواع عرض النطاق المتكيف:

يوضح الجدول الاتي أنواع عرض النطاق وصياغتها وتطبيقاتها:

الجدول(3)أنواع عرض النطاق المتكيف

النوع	الصيغة/الآلية	التطبيق الأمثل
متكيفة تغطية	$h_i = h \left(\frac{g}{f(x_i)} \right)^\alpha$	بيانات غير متجانسة
متكيفة مكانية	متغير حسب الموقع $h(x)$	خرائط مكانية
متكيفة متعددة المتغيرات	مصفوفة التباين المتكيفة	بيانات عالية الأبعاد

• عرض النطاق الذي سنستخدمه في بحثنا

عرض النطاق المتكيف حسب النقطة (Adaptive Point-Sample)

(Some,2025; SHI 2010 ; Zhao&Tabak, 2023) يختلف عرض النطاق عند كل

نقطة بيانات (x_i) وفقاً للكثافة المحلية. ويُعد قانون Abramson، الأكثر شيوعاً في هذا السياق،

حيث يُعرّف عرض النطاق المتكيف كالآتي: (Some,2025)

$$(9) \quad h_i = h \left(\frac{g}{f(x_i)} \right)^\alpha$$

h عرض النطاق الابتدائي

h_i : عرض النطاق المتكيف عند النقطة x_i

α : معامل يتحكم بقوة التكيف عادة يكون (0.5)

$f(x_i)$: تقدير كثافة أولي عند النقطة x_i باستخدام KDE ثابت

g : المتوسط الهندسي لقيم الكثافة المقدرّة الأولية

ويحسب من العلاقة :

$$(10) \quad g = \exp\left(\frac{1}{n} \sum_{i=1}^n \ln f(x_i)\right)$$

ويُعطى تقدير الكثافة النهائي بواسطة:

$$(11) \quad f_{adaptive}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h_i} K\left(\frac{x-x_i}{h_i}\right)$$

$K(x)$: دالة النواة kernel

n : عدد المشاهدات في العينة

x : النقطة التي نريد تقدير الكثافة عندها

x_i : مشاهدات العينة .

5-التصنيف في التقدير اللامعلمي (Nonparametric Classification):

تعتمد عملية التصنيف في العديد من الأساليب الإحصائية على قاعدة بايز التي تُعد من القواعد الأساسية في إتخاذ قرار التصنيف، حيث يتم حساب الاحتمال اللاحق للفئة بالاعتماد على الاحتمال السابق واحتمال البيانات المعطاة. (Ghosh, 2006; Marzio et al, 2019) فحسب قاعدة بايز يكون :

$$(12) \quad P(C_j/x) = \frac{P(C_j)\hat{f}_j(x)}{\sum_{k=1}^c \hat{f}_k(x)P(C_k)}$$

حيث: x : العينة الجديدة (العينة المراد تطبيقها)

C_j : العينة رقم j

C_k : العينة رقم k

ولغرض التصنيف لانتاج إلى المقام لأنه مشترك بين جميع الفئات، لذلك تصبح قاعدة القرار :

$$(13) \quad j = \arg \max_j [P(C_j)\hat{f}_j(x)]$$

حيث ان $P(C_j) = \frac{n_j}{n}$ تمثل الاحتمال السابق للفئة C_j ، و n_j : عدد العينات التي تنتمي إلى الفئة

C_j و n هو اجمالي عدد العينات، و $f_j(x)$ تمثل دالة الكثافة الاحتمالية للفئة C_j عند النقطة x ، ولكن في كثير من الحالات العملية لا يكون شكل دالة الكثافة الاحتمالية معروفاً مسبقاً، لذلك يتم اللجوء إلى التصنيف في التقدير اللامعلمي الذي لايفترض شكلاً محدداً لتوزيع البيانات، بل يعتمد على البيانات نفسها لتقدير هذه الدوال.

6-طريقة النواة في التصنيف (Kernel Classification Rule):

تُعد طريقة النواة من أشهر الطرق في التصنيف اللامعلمي، حيث يتم استخدامها لتقدير دالة الكثافة الاحتمالية $f_j(x)$ لكل فئة اعتماداً على البيانات المتوفرة باستخدام تقدير النواة (Kernel Density Estimation -KDE)، (Ghosh, 2006; Marzio et al, 2019) وبذلك تستخدم طريقة النواة لتوفير تقدير للكثافات الاحتمالية المطلوبة في قاعدة بايز. بعد تقدير هذه الدوال لكل فئة، يتم تطبيق قاعدة القرار البايزية العلاقة (13):

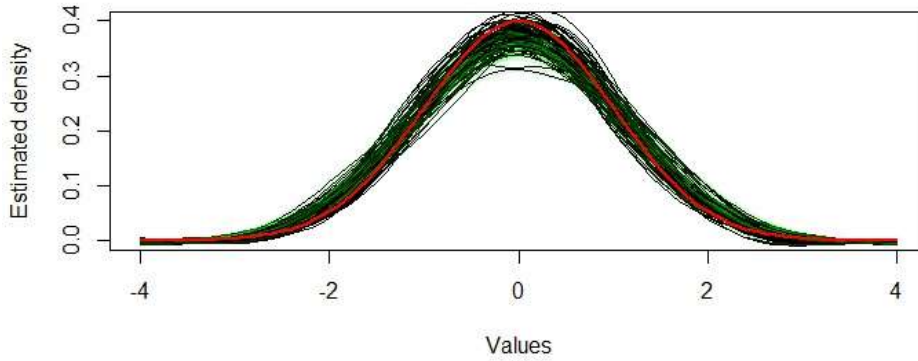
أي ان دور طريقة النواة هو تقدير دوال الكثافة الاحتمالية $f_j(x)$ ، بينما تقوم قاعدة بايز باستخدام هذه التقديرات مع الاحتمالات السابقة للفئات لاتخاذ قرار التصنيف النهائي، وبهذا تكون طريقة النواة مكتملة لقاعدة بايز في إطار التصنيف اللامعلمي، حيث يتم تقدير الكثافات من البيانات مباشرة دون افتراض شكل توزيعي محدد.

7- المحاكاة العددية :

تم تقدير دالة الكثافة الاحتمالية باستخدام نواة إيباشينكوف ونواة غاوس على بيانات مولدة عشوائياً من توزيع طبيعي قياسي (متوسط 0 وانحراف معياري 1)، حيث تم اخذ عينة بحجم 300، وقمنا بتكرار التجربة 100 مرة. وتم حساب عرض النطاق بطريقة سيلفرمان فحصلنا على $h=0.285$ ، وتم حساب دالة الكثافة المقدرّة على المجال [-4,4] حيث تم تطبيق كل ذلك على برنامج الحزمة R

والشكل (1) يوضح دالة تقدير الكثافة لبيانات مولدة عشوائياً.

Kernel Density Estimates: Epanechnikov (black) & Gaussian (Green)



الشكل (1) دالة كثافة لبيانات مولدة عشوائياً

يوضح الشكل (1) المنحنيات السوداء وهي عبارة عن تقديرات كثافة مختلفة لعينة من بيانات مولدة من توزيع طبيعي باستخدام نواة ايباشينكوف والمنحنيات الخضراء هي تقديرات كثافة باستخدام نواة غاوس والمنحنى الأحمر هي الكثافة النظرية للتوزيع الطبيعي القياسي ، كلما اقتربت المنحنيات السوداء والخضراء من الحمراء، دل ذلك على أن النواتين دقيقتين في تقدير الكثافة .

• تم أيضا توليد بيانات عشوائية ل عينة حجمها ($n = 100$) من توزيع طبيعي معياري، باستخدام نواة ايباشينكوف وعرض نطاق (h) بطريقة سيلفرمان ثابت و متكيف وتم تكرار التجربة 50 مرة للمقارنة بين النطاقين والجدول (4) يوضح المقارنة لأصغر 10 قيم من حيث الخطأ ، كما تم اخذ عينة من توزيع طبيعي معياري أيضا لعينة بحجم (100) ، باستخدام نواة ايباشينكوف وعرض نطاق (h) بطريقة التحقق المتقاطع (LSCV) ثابت و متكيف، وتم تكرار التجربة 50 مرة ، للمقارنة بين النطاقين والجدول (5) يوضح المقارنة من حيث الخطأ لاصغر 10 قيم .

الجدول (4) مقارنة متوسط مربع الخطأ MSE بين المقدر التكيفي والثابت بطريقة Silverman

MSE(Adaptive)	MSE(Fixed)	Silverman h	رقم التجربة
0.0000930	0.0000770	0.365857	13
0.0001249	0.0000743	0.339632	16
0.0001119	0.0001119	0.375064	17
0.0001214	0.0000854	0.334215	23
0.0001592	0.0000769	0.308110	24
0.0001154	0.0000531	0.346955	26
0.0001394	0.0000761	0.325246	33
0.000133	0.000105	0.359319	34
0.0001051	0.0000683	0.369903	45
0.0001531	0.0000912	0.348588	48

يوضح الجدول (4) قيم الخطأ لعرض النطاق الثابت والمتكيف بطريقة سيلفرمان لتكرار 10 تجارب حيث كانت قيمته 0.3 في كل التجارب وكانت أصغر قيمة للخطأ للنافذة الثابتة هي 0.00005 بينما كانت أصغر قيمة لخطأ

للنافذة المتكيفة هو 0.000009 ، حيث أعطى عرض النطاق الثابت قيم خطأ اصغر منها بالنسبة للمتكيف بطريقة سيلفرمان.

الجدول (5) مقارنة متوسط مربع الخطأ MSE بين المقدر التكيفي والثابت بطريقة LSCV

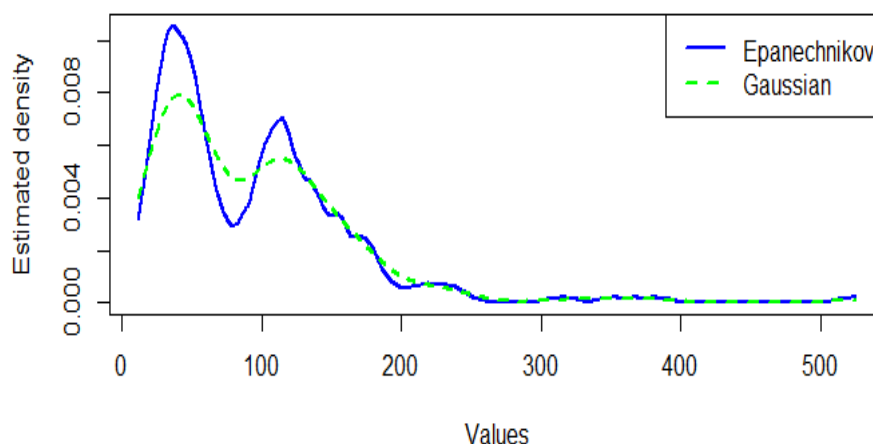
MSE(Adaptive)	MSE(Fixed)	(LSCV)bandwidth	رقم التجربة
0.000312	0.000224	0.558	14
0.000377	0.000279	0.525	15
0.000234	0.000180	0.525	17
0.000313	0.000141	0.626	21
0.000244	0.000154	0.491	26
0.000307	0.000214	0.660	27
0.000200	0.000114	0.592	33
0.000295	0.000199	0.558	37
0.000420	0.000251	0.525	40
0.000324	0.000226	0.491	48

يوضح الجدول (5) قيم الخطأ لعرض النطاق الثابت والمتكيف بطريقة التحقق المتقاطع لتكرار 10 تجارب حيث كانت قيمته 0.5 تقريبا في كل التجارب وكانت أصغر قيمة للخطأ للنافذة الثابتة هي 0.000114 بينما كانت أصغر قيمة لخطأ عرض النطاق المتكيف هو 0.000200 حيث أعطى عرض النطاق الثابت قيم خطأ أصغر منها للمتكيف بطريقة التحقق المتقاطع

• ثم تم تطبيق تقدير دالة الكثافة الإحتمالية على أوزان أسماك تم جمعها من (معهد البحوث البحرية في كلية الزراعة في جامعة اللاذقية-سورية)، حيث تم اختيار عينة من 200 سمكة من نوعي (الغُبس و السرغوس) بأطوال وأوزان مختلفة حيث تتراوح أعمارها بين السنة و السنتين وتتميز هذه الأنواع بفترات تكاثر مختلفة حيث تبدأ أسماك السرغوس بالتكاثر من شهر كانون الأول وحتى شهر نيسان، بينما تمتد فترة تكاثر أسماك الغُبس من شهر نيسان إلى شهر تموز، وذلك باستخدام طريقة النواة (نواة إيباشنكوف وغاوس)، و عرض نطاق بطريقة سيلفرمان وكانت قيمة عرض النطاق المثالية

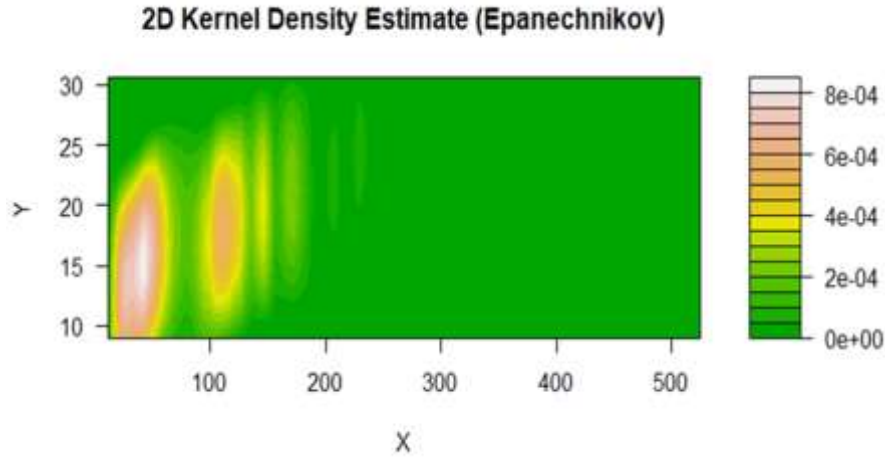
$h = 19.5290$ ويوضح الشكل (2) دالة الكثافة .

Kernel Density Estimate Comparison



الشكل (2) دالة كثافة احتمالية لأوزان الأسماك

يوضح الشكل (2) وجود قمة عالية حول القيمة 30 تقريباً أي هناك عدد كبير من القيم متمركزة في هذا النطاق والكثافة المقدره قريبة من 0.009 وهي القيمة الأعلى في الرسم وتوجد قمة أخرى متمركزة حول القيمة 100 أي توجد مجموعة بيانات أخرى لكنها اقل تكرارا من الأولى كما نلاحظ انحدار تدريجي نحو اليمين أي التوزيع يقل بشكل كبير ، وتشير الكثافة إلى أن هذه القيم نادرة أو قليلة في تلك الجهة ، كما تم تنفيذ تقدير إضافي مع أخذ أطوال الأسماك بعين الاعتبار، فحسبنا دالة الكثافة بناءً على الأوزان المرتبطة بالأطوال و باستخدام عرض نطاق مثالي جديد وقيمه $h = 10.2475$ ، والشكل (3) يوضح دالة الكثافة مع الأخذ أطوال الأسماك بعين الاعتبار.

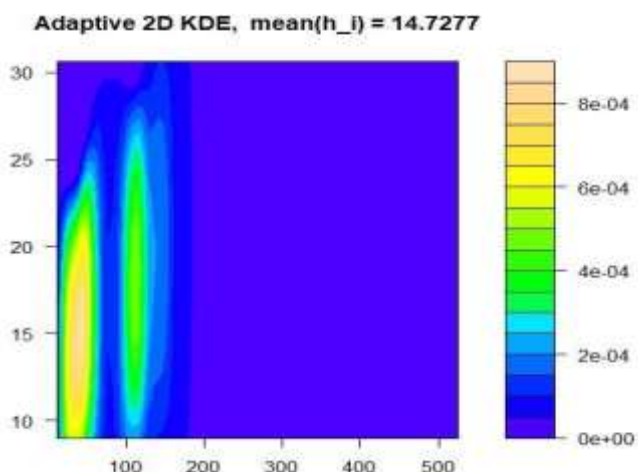


الشكل (3) تقدير الكثافة الاحتمالية ثنائية الأبعاد

الشكل (3) تشير الألوان الى كثافة النقاط في تلك المنطقة حيث اللون الأخضر الداكن يمثل كثافة منخفضة الى معدومة أما الألوان من أصفر إلى أحمر إلى وردي إلى أبيض تمثل كثافة متزايدة الأبيض هو الأعلى كثافة . (شريط التدرج) يُظهر القيم العددية للكثافة ، من صفر (لاشيء) الى 0.0008 (أعلى كثافة) ، هذه القيم ليست "عدد نقاط" بالضبط، بل هي تقدير احتمالي للكثافة في ذلك الموضع. مثلاً: إذا كانت منطقة على الرسم ملونة بـ"أبيض"، فهذا يعني أن الكثافة التقديرية في تلك المنطقة تساوي 0.0008 ، وهي القيمة الأعلى في التوزيع، أما المناطق الخضراء فتعني أن الكثافة فيها قريبة من الصفر.

• قمنا بعد ذلك بتطبيق عملية التقدير لدالة الكثافة الاحتمالية لاوزان واطوال الأسماك بطريقة سيلفرمان وعرض نطاق متكيف وكلا النواتين (اسياشنكوف وغاوس)، فكان متوسط عرض النطاق التكييفي يساوي

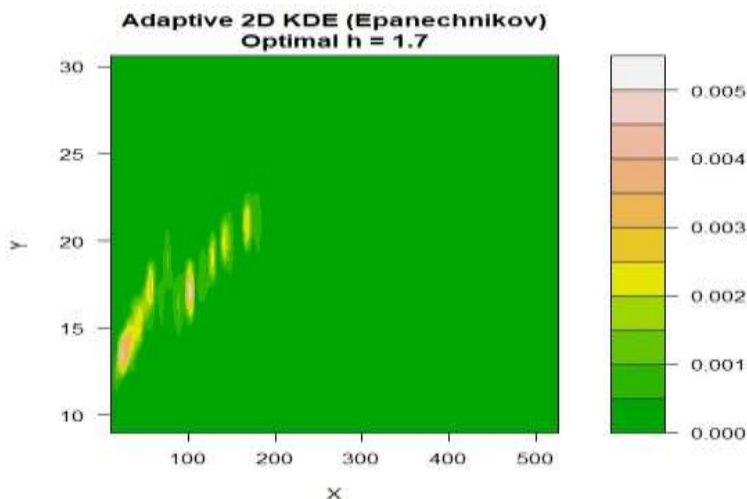
والشكل (4) يوضح دالة الكثافة المقدره الناتجة . $mean(h_i) = 14.727$



الشكل (4) دالة الكثافة بمتغيرين بطريقة سيلفرمان

يوضح الشكل (4) تقدير الكثافة الاحتمالية ثنائية الابعاد التكيفي، حيث يشير $mean(h_i) = 14.727$ إلى ان متوسط عرض النواة التكيفية عبر نطاق البيانات يساوي تقريباً "14.73"، ويمثل المحور الافقي اوزان الأسماك بينما المحور العمودي يمثل أطوال الأسماك والذي يتغير ضمن نطاق تقريبي، يوضح شريط الألوان على اليمين يوضح قيمة الكثافة الاحتمالية المقدرّة، الأصفر-الفاصح: كثافة احتمالية عالية (تركيز مرتفع للبيانات)، الأخضر-الأزرق: كثافة متوسطة، الأزرق الداكن - البنفسجي: كثافة منخفضة أو شبه معدومة، القيم القصوى للكثافة تقع تقريباً في حدود 10^{-4} وهو امر شائع في تقديرات KDE ثنائية الابعاد

• أيضاً قمنا بتنفيذ عملية التقدير لدالة الكثافة الاحتمالية لأوزان وأطوال الأسماك وذلك باستخدام عرض النطاق المتكيف (h) بطريقة التحقق المتقاطع (LSCV) وكلا النواتين وأعطت عملية التقدير عرض نطاق يساوي $h = 1.76$ والشكل (5) يوضح دالة الكثافة الناتجة



الشكل (5) دالة الكثافة بمتغيرين بطريقة LSCV

يوضح الشكل (5) تقدير الكثافة الاحتمالية ثنائية الابعاد، مع عرض نافذة أمثل $h = 1.76$ بأسلوب احصائي (LSCV)، حيث يسمح صغرى قيمة النافذة الى الحفاظ على التفاصيل الدقيقة للتوزيع، ويحدّ من التنعيم المفرط، مما يؤدي إلى إبراز القمم المحلية ومناطق التركز الفعلي للبيانات، كما تشير القيم المرتفعة نسبياً للكثافة الى احتمال أكبر لتمركز المشاهدات، في حين تعكس القيم المنخفضة خارج هذا النطاق ضعف أو غياب البيانات، تدرج الألوان

الألوان الداكنة (الأخضر الداكن) تشير إلى قيم كثافة منخفضة جدًا، أي مناطق ذات احتمال ضعيف أو شبه معدوم لوجود بيانات، والألوان المتوسطة (الأخضر الفاتح - الأصفر) تمثل مستويات كثافة متوسطة، وتدل على مناطق انتقالية تحتوي على عدد معتدل من المشاهدات، بينما الألوان الفاتحة والداكنة (الأصفر - البرتقالي) تشير إلى قيم كثافة مرتفعة، وتمثل مناطق التركيز الأعلى للبيانات، أي المواقع التي تتكدس فيها المشاهدات بشكل ملحوظ، أعلى درجات اللون (إن وجدت مثل الأبيض أو البرتقالي الفاتح) تمثل القمم المحلية للكثافة، أي النقاط الأكثر احتمالاً لظهور البيانات.

ملاحظة:

الدلالة الإحصائية لتدرج الألوان : يعكس التدرج اللوني.

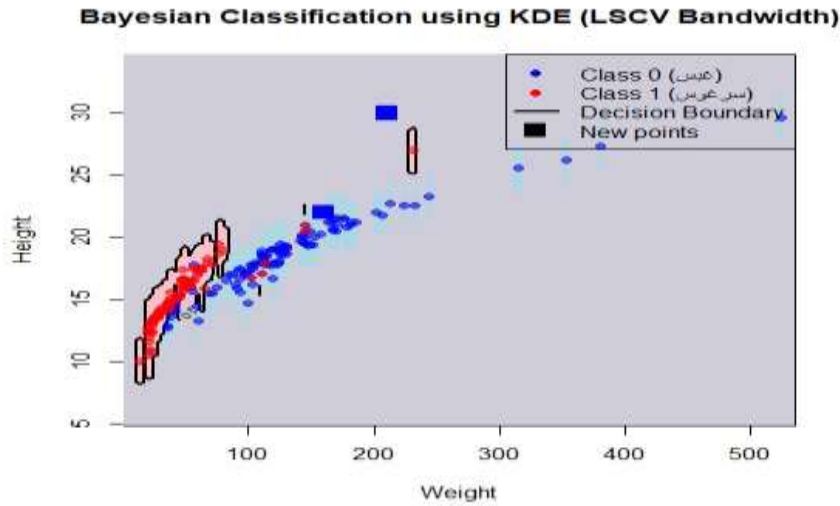
- شدة التركيز الاحتمالي للبيانات
- مستوى الكثافة المقدرة بواسطة النواة
- التغير المكاني للكثافة عبر مستوى المتغيرين

وبالتالي، يسمح تدرج الألوان بفهم البنية الداخلية للتوزيع وتحديد مناطق القيم العليا والدنيا للكثافة بشكل

بصري

واضح.

- أخيراً قمنا بتصنيف اوزان اسماك جديدة باستخدام التحقق المتقاطع وعرض النطاق التكميحي، حيث تم تصنيف المشاهدات الجديدة باستخدام قاعدة بايز اعتماداً على تقدير كثافة النواة، حيث تمثل الفئة (1) نوع السمك سرغوس، بينما تمثل الفئة (0) نوع السمك غبس. يوضح الشكل مناطق القرار وحدّ الفصل بين الفئتين (الخط أو المنحني الذي يفصل بين الفئتين)، مع إبراز النقاط الجديدة المصنّفة.



الشكل (6) تصنيف مشاهدات جديدة

يوضح الشكل (6) تصنيف الاسماك بناء على الوزن والطول، الصنف 1 سمك السرغوس والصنف 0 سمك الغبس والشكل ■ للمشاهدات المصنفة حيث تم تصنيفها من نوع غبس، كما هو موضح.

8- الأدوات المستخدمة :

تم استخدام برنامج الحزمة الإحصائية R لتمثيل عمليات التحليل وتمثيل البيانات ، تعد لغة البرمجة R من أقوى لغات البرمجة في مجال الحوسبة الإحصائية، حيث أنها لغة مفتوحة المصدر، تستخدم على نطاق واسع في تحليل وتمثيل النتائج وإنشاء الرسوم البيانية، (زينه ، 2017) وصدرت أول نسخة منها عام 2000 ومنذ ذلك الحين تطورت بشكل سريع لتصبح لغة رئيسية ومستخدمة بشكل كبير بين الإحصائيين، ومن أبرز ميزات لغة R: بساطة الاستخدام - المرونة والتوافق - القدرة الإحصائية الشاملة : تتيح لغة R بتنفيذ جميع العمليات الإحصائية بسهولة ، مثل: الأمر () sum لحساب المجموع ، والأمر () mean لحساب المتوسط والعديد من الدوال الأخرى كل ذلك يتم بشكل بسيط وسلس مع مخرجات واضحة .

9-الاستنتاجات والتوصيات:

1-9: الاستنتاجات:

فُمنّا في هذا البحث بدراسة تقدير دالة الكثافة الاحتمالية باستخدام نواتين مختلفتين وعرض نطاق ثابت ومتكيف، وذلك على بيانات مولدة عشوائياً وبيانات حقيقية، مع توظيف نتائج التقدير في عملية التصنيف أظهرت نتائج الدراسة مايلي:

_ أن تقديرات الكثافة باستخدام نواتي إيباشينكوف وغاوس كانت قريبة بدرجة كبيرة من دالة الكثافة النظرية للتوزيع الطبيعي، حيث تركزت معظم المنحنيات المقدّرة حول المنحنى النظري الشكل(1)، مما يعكس كفاءة الطريقتين عند تطبيقهما على بيانات ناتجة عن توزيعات متماثلة. كما لوحظ أن التباين بين التقديرات عبر عمليات المحاكاة كان محدوداً، مما يدل على استقرار تقدير الكثافة بالنواة في هذه الحالة. عند استخدام نفس قيمة عرض النطاق (h = 0.3).

_ حققت نواة إيباشينكوف متوسط خطأ تربيعي أقل (MSE = 0.0002) مقارنة بنواة غاوس (MSE = 0.0007)، مما يشير إلى تفوقها من حيث تقليل التباين ضمن الدعم المحدود.

_ قدمت نواة غاوس تقديراً أكثر نعومة يغطي المجال بالكامل، إلا أن هذه السلاسة قد تؤدي إلى إخفاء بعض التفاصيل الدقيقة في التوزيع، خاصة في المناطق ذات التجمعات النقطية. وأظهرت النتائج عدم وجود فروق جوهرية بين النواتين عند تطبيق التقدير على البيانات المولدة عشوائياً أو على البيانات الحقيقية، مما يدل على أن اختيار النواة لا يؤثر بشكل كبير على جودة التقدير مقارنة بتأثير اختيار عرض النطاق.

_ بينت طرق اختيار عرض النطاق أن القيم الناتجة عن قاعدة سيلفرمان والتحقق المتقاطع (LSCV) أعطت نتائج جيدة مع البيانات المولدة، بينما كانت النتائج أكثر دقة مع البيانات الحقيقية عند استخدام عرض النطاق المتكيف، حيث بلغت قيمته (h = 1.7) وحققت خطأ تربيعياً أقل (MSE = 0.00006) مقارنة بعرض النطاق الثابت.

_ في سياق التصنيف، أظهرت طريقة التقدير باستخدام عرض نطاق متكيف قيمة مثلى قريبة من (h = 1.81)، وحققت دقة تصنيف مرتفعة بلغت 0.935، مع خطأ تصنيف منخفض (MSE = 0.04)، مما يؤكد كفاءة التقدير المتكيف في تحسين أداء التصنيف بوجه عام، تُعد نواة إيباشينكوف أكثر كفاءة حسابياً وتقدم تقديراً محلياً أكثر تركيزاً ودقة، في حين تُعد نواة غاوس أنسب للبيانات التي تحتوي على ضوضاء نظراً لسلاستها واستمراريتها على كامل المجال.

_أكدت النتائج أن طريقة التحقق المتقاطع، وخاصة عند استخدام عرض النطاق المتكيف، تؤدي إلى تحسين ملحوظ في دقة التقدير والتصنيف مقارنة بعرض النطاق الثابت.

2-9:التوصيات:

- توسيع التحليل البحثي لتشمل أكثر من متغيرين وذلك بدراسة البيانات متعددة الأبعاد باستخدام تقنيات التقدير اللامعلمي.
- تطوير نوى هجينة تجمع بين مزايا نواة ايباشينكوف (من حيث الدقة المحلية) ونواة غاوس (من حيث السلاسة والامتداد) بهدف تحسين أداء التقدير في ظروف مختلفة.
- يُوصي بتحسين طرق إختيار معلمة النافذة h من خلال دمج تقنيات التعلم الآلي مثل التعلم المعزز أو التعلم العميق.
- دراسة تأثير الابعاد العالية على أداء النافذة المنكيفة في التصنيف .

10-المراجع:

1-10: المراجع الأجنبية:

- BOLI, I, A-D. ; C, Wei. B. 2024, *Bayesian Classifier Based on Robust Kernel Density Estimation and Harris Hawks Optimisation* . School of Computer Science, Hubei University of Technology, Wuhan, China , pp. 1-23.
- GUIDOUM, A, C . 2024, *Kernel Estimator and Bandwidth Selection for Density and its Derivatives* .The kedd Package . University of Science and Technology Houari Boumediene, Algeria, pp. 1-22
- GHOSH, A.; CHAUDHURI, P.; SENGUPTA, D. 2006. *Classification Using Kernel Density Estimates*.Technometrics. 48,120-132.
- MARZIO, M.; FENSORE, S.; PANZERA, A.; TAYLOR, C. 2019. *Kernel density classification for spherical data*. Statistics &Probability Letters.
- MOHAMED, Y . ; IBRAHEM, D. 2008, *Estimation of the Nonparametric Regression Function Using Some Monotonic Nonparametric Methods* . Journal of Economic and Administrative Sciences, Vol. 14, no. 50, pp. 304-316
- PARZEN, E. 1962. *On Estimation of a Probability Density Function and Mode* . Annals of Mathematical Statistics , Vol. 33, no. 3, pp. 1065-1076
- POPOV, A, A. ; ZANETTI, R . ; MEMBER, S. 2024, *The Ensemble Epanechnikov Mixture Filter* . University of Texas at Austin, Austin, TX.
- QIN, T. ; HUANG, Wei-Min. 2025, *On Kernel-based Variational Autoencoder* . arXiv. Information and Inference: A Journal of the IMA.
- SILOKO, I. U.; SILOKO ,E. A. ; IKPOTOKIN ,O. 2020, *A Mini Review of Dimensional Effects on Asymptotic Meanintegrated Squared Error and Efficiencies of Selected Beta Kernels* ,Journal of Mathematics and Statistics (JJMS), Vol. 13, no. 3, pp. 327-340
- SOME, S.; KOKONENDJI, C.; DOBE'L'E-KPOKA, F. 2025, *An effective estimation of multivariate density functions using extended-beta kernels with Bayesian adaptive bandwidths*.
- SHI, X. 2010, *Selection of bandwidth type and adjustment side in kernel density estimation over inhomogeneous backgrounds*.International . Journal of Geographical Information Science, 24, 643 - 660.

- TURLACH, B, A. 2013, *Bandwidth Selection in Kernel Density Estimation: A Review*, C.O.R.E. and Institut de Statistique Universite' Catholique de Louvain.
- ZAMBOM, A. Z. ; DIAS, R. 2012, *A Review of Kernel Density Estimation with Applications to Econometrics* . Universidade Estadual de Campinas.
- ZHAO, W.;TABAK, E. 2023, *Adaptive kernel conditional density estimation*.
- HANSEN, B, E. 2009, *Lecture Notes on Nonparametrics*.1st , University of Wisconsin. United States. 25 .
- WAND, M . ; JONES , C . 1995 , *Kernel Smoothing* . New York Washington & London , 212.
- SCOTT, D, W. 1992, *Multivariate Density Estimation : Theory, Practice And Visualization* . 1st, John Wiley&Sons , United States, 317.

2-10: المراجع العربية

- محمد، أروى جاسم ; إبراهيم ، وضاح صبري. 2020 ، دراسة احتمالية التغير الحاصل في أسعار الأسهم بالاعتماد على القيمة المتداولة لفندق بابل باستعمال الانحدار اللوجستي. مجلة الإدارة والاقتصاد ،العراق، 43 ،(123)، 446-434.
- زينة ، محمد بشر . 2017 ، لغة البرمجة الإحصائية R . الإصدار الأول ، حلب ، جامعة حلب ، 155