

العنوان:	استخدام خوارزميات تعلم الآلة لتصنيف همزتي الوصل والقطع
المصدر:	مجلة كلية الآداب والعلوم الإنسانية
الناشر:	جامعة قناة السويس - كلية الآداب والعلوم الإنسانية
المؤلف الرئيسي:	الثقفي، طلال أحمد شداد
مؤلفين آخرين:	علي، ياسر نصر الدين السيد، عبدالمعطي، محمد فتحي عبدالفتاح، بشير، طلال الطاهر قطيبي(م. مشارك)
المجلد/العدد:	ع35
محكمة:	نعم
التاريخ الميلادي:	2020
الشهر:	ديسمبر
الصفحات:	11 - 47
رقم MD:	1167410
نوع المحتوى:	بحوث ومقالات
اللغة:	Arabic
قواعد المعلومات:	AraBase, HumanIndex
مواضيع:	اللغة العربية، القواعد النحوية، البرمجة اللغوية، لغات الآلة، الذكاء الاصطناعي
رابط:	http://search.mandumah.com/Record/1167410

استخدام خوارزميات تعلم الآلة لتصنيف همزتي الوصل والقطع

- د. طلال أحمد شداد الثقفي
الأستاذ المساعد بقسم اللغة العربية،
الكلية الجامعية بترية، جامعة الطائف
- د. ياسر نصر الدين السيد علي
الأستاذ المساعد بقسم الرياضيات، برنامج علوم
الحاسب بالكلية الجامعية بترية، جامعة الطائف
البريد الجامعي/ ynali@tu.edu.sa
- د. محمد فتحي عبد الفتاح عبد المعطي
الأستاذ المشارك بقسم اللغة العربية،
الكلية الجامعية بترية، جامعة الطائف
البريد الجامعي/ m.fathe@tu.edu.sa
- د. طلال الطاهر قطبي بشير
الأستاذ المشارك بقسم اللغة العربية،
الكلية الجامعية بترية، جامعة الطائف

تم تمويل هذا البحث برعاية ودعم عمادة البحث العلمي،
جامعة الطائف - المملكة العربية السعودية [مجموعة بحثية
رقم ١٠١-١٤٤١-١]

الملخص:

تعتبر اللغة العربية كائن حي ينمو ويتطور بالممارسة والتطبيق الصحيح لجميع أديانها وفروعها النحوية والصرفية والدلالية والمعجمية، وفي هذه الدراسة البحثية نستعرض مساهمة التقنية في تطور اللغة العربية لاسيما الرسم الصحيح لهمزتي الوصل والقطع.

يهدف هذا البحث بشكل خاص لبناء نموذج مصنف ذكي يعمل على تصنيف الكلمات العربية المبتدئة بحرف الألف، وتصنيف همزتها إلى همزة وصل أو همزة قطع من خلال استخدام تقنيات الذكاء الاصطناعي بصورة عامة، وخوارزميات تعلم الآلة بصورة خاصة من أجل تأسيس معايير دقيقة وصحيحة في رسم همزتي الوصل والقطع بالصورة الصحيحة، وبذلك نكون قد طوعنا التقنية لتسهم في خدمة اللغة العربية.

اعتمد هذا البحث على تجميع الكلمات العربية المبتدئة بالهمزة وذلك عن طريق تصميم استبانة رقمية مهمتها تجميع أكبر قدر من الكلمات المبتدئة بالهمزة وتصنيفها لهمزتي وصل أو قطع وفقاً للقواعد النحوية المتبعة في ذلك، تم نشر الاستبانة على الويب وتمت تعبئتها بواسطة ٥٠ متخصص في النحو بدرجات علمية متفاوتة، حيث بلغ عدد الكلمات الكلية المصنفة ٤٠٠ كلمة،

وبعد معالجة واستبعاد الكلمات المكررة وعددها ١٠١ كلمة، حصلنا على عدد ٢٩٩ كلمة صالحة لتطبيقها على نموذج المصنف، وبناءً على حجم ونوع البيانات المجمعة وآلية التصنيف المتبعة تم تطبيق خوارزميات تصنيف تتناسب مع العينة المجمعة مثل خوارزمية آلة دعم المتجه (SVM) وخوارزمية نايف بيز (NB) وخوارزمية الجار الأقرب (KNN) وذلك من خلال استخدام لغة Python ومكتبة (sk-learn).

بعد تدريب نماذج المصنفات المستخدمة وقياس دقة الخوارزميات تبين أن خوارزمية آلة دعم المتجه (SVM) قد حصلت على أعلى دقة للنموذج بنسبة ٩٢% وهي نسبة مرتفعة وكافية لحل مشكلة البحث.

الكلمات المفتاحية:

الذكاء الاصطناعي، خوارزميات تعلم الآلة، همزة الوصل، همزة القطع، خوارزميات التصنيف، اللغة العربية.

Abstract

The Arabic language is considered a living organism that grows and develops through the correct practice and application of all its literature and its syntactic, morphological, semantic, and lexical branches. In this research study, we review the contribution of technology to the development of Arabic language, especially the correct writing of Hamzat Alwasl and Hamzat Alqatae.

This research aims to build a smart Classifier model to classify Arabic words beginning with the letter alif and their hamza into conjunctive and disjunctive through the using of artificial intelligence techniques in general and machine learning algorithms to establish accurate and correct criteria in writing the conjunctive and disjunctive hamza correctly.

Consequently, technology would be adapted to contribute to the service of Arabic language.

This research relied on the compilation of Arabic words beginning with hamza by designing a digital questionnaire. The task of this questionnaire is to collect the largest number of words beginning with hamza and to classify them as Alwasl and Hamzat Alqatae according to the grammatical rules followed in this process. The questionnaire was circulated in the internet web and was filled out by fifty specialists in syntax with different academic ranks. The total number of classified words reached ٤٠٠ words, and after processing and excluding the repeated words, ١٠١ words, we obtained ٢٩٩ valid words to be applied to the Classifier model, and based on the size and type of the collected data and the classification mechanism followed, classification algorithms were applied that fit the collected sample such as the vector support machine algorithm, Naif Biz algorithm, and Nearest neighbor algorithm by using Python language and sk-learn library.

After training the used classificatory models and measuring the accuracy of the algorithms, it was quite apparent that the vector support machine (SVM) algorithm had obtained the highest accuracy of the model (٩٢%), a high and enough percent to solve the research problem.

key words:

Artificial Intelligent, Machine Learning Algorithms, Hamza, Wasl, Gtaa, Classification Algorithms, Arabic Language.

تمهيد:

يمكننا تعريف التعلم الآلي على أنه يتبع لعلم الذكاء الاصطناعي الذي يهتم بتصميم خوارزميات وتقنيات تسمح لأجهزة الحاسب الآلي بعملية التعلم الذاتي مع إمكانية تطوير هذه الخوارزميات.

يمكننا استعراض نوعين من أنواع التعلم الآلي هما: التعلم الاستقرائي والتعلم الاستنتاجي، يعرف التعلم الاستقرائي أيضاً بالتعلم الاستكشافي ويعمل على مبدأ استنتاج قواعد عامة من البيانات، وهذا يختلف عن التعلم الاستنتاجي حيث يتم إعطاء المتعلم القواعد التي يحتاج لتطبيقها.

يعتمد التعلم الآلي على مبدأ تعلم الأنظمة من البيانات المتاحة وتحديد الأنماط المناسبة واتخاذ القرارات بدون تدخل بشري، وحيث أن المهمة الأساسية للتعلم الآلي هي استخراج المعلومات من البيانات، فإن التعلم الآلي هو طريقة من طرق تحليل البيانات تعمل على أتمتة بناء نماذج تحليل بيانات تعرف بنماذج التحليل التنبؤية، يُعرف هذا النمط باسم نموذج التحليل التنبؤي أو التصنيفي، تتيح هذه النماذج للباحثين وعلماء البيانات والمحللين اتخاذ قرارات تمكن من الوصول لنتائج موثوقة.

يخضع بناء نموذج تحليل البيانات لعدة آليات أو طرق للوصول لنتائج دقيقة، حيث أن مرحلة بناء النموذج تستلزم تنفيذ عدد من الخطوات الرياضية أو المنطقية المرتبة والمتسلسلة للوصول لنموذج تحليلي مدرب يستطيع التنبؤ بقيم أقرب للحقيقة أو تصنيف البيانات بشكل دقيق، يطلق على هذه المجموعة من الخطوات اسم الخوارزمية، كما يطلق على نموذج تصنيف البيانات اسم المصنف.

وفي هذا البحث سنتناول استخدام نماذج خوارزميات تعلم الآلة في تصنيف همزتي الوصل والقطع وسيتم تطبيق هذه الخوارزميات على الكلمات العربية نسبة لأن اللغة العربية لم تنل حظاً من الدراسة والبحث كغيرها من اللغات الأخرى.

أهداف البحث:

- حل مشكلة الاستخدام الخاطيء لرسم همزتي الوصل والقطع في النصوص العربية.
- التعرف على أنسب خوارزميات التصنيف ضمن خوارزميات تعلم الآلة لتحديد المواضع الصحيحة لاستخدام همزتي الوصل والقطع في بداية الكلمة.

- قياس جودة أشهر خوارزميات التصنيف في التمييز بين همزتي الوصل والقطع في بداية الكلمة.

مشكلة البحث وأسئلته:

تكمن مشكلة البحث الأساسية في الاستخدام الخاطئ لرسم همزتي الوصل أو القطع في غير موضعها الصحيح، وذلك عند كتابة الكلمات المبدوءة بحرف الألف بصورة خاصة، حيث يخطئ الكثيرين في رسم الهمزة أو إسقاطها الشيء الذي يضعف قوة اللغة.

يسهم هذا البحث، في تطوير اللغة العربية، وذلك عن طريق بناء نموذج تصنيف ذكي يمكن المحررين والممارسين، وكل مستخدمي اللغة من التصنيف الصحيح لهمزتي الوصل والقطع.

يسعى هذا البحث للإجابة عن الأسئلة التالية في ضوء الاستفادة من خوارزميات تعلم الآلة المستخدمة في التصنيف والتنبؤ:

- كيف يمكننا التمييز بين همزتي الوصل والقطع في النصوص العربية؟

- ما الفوائد التي تجنيها عند تصنيف همزتي الوصل والقطع في النصوص العربية؟

- ما خوارزميات تعلم الآلة المناسبة لتصنيف همزتي الوصل والقطع في النصوص العربية؟
الدراسات، السابقة:

تعد هذه الدراسة من الدراسات الجديدة في حقل الدراسات النحوية والتقنية التي تتناول استخدام خوارزميات تعلم الآلة لتصنيف همزتي الوصل والقطع في النصوص العربية، وعلى حد علمي، ومن خلال بحثي، لم أجد أي دراسات تعرضت لهذا الموضوع من قبل نسبة لصعوبة استخدام المكتبات البرمجية في التعامل مع النصوص العربية بمقدارة، ولكن هناك وجه شبه كبير وتداخل لدراسات أخرى ارتبطت بالبحث لاستخدام هذه الخوارزميات في تطبيقات مماثلة ومشابهة، وقد استفدت منها، وساعدتني كثيراً في الولوج لبحثي، ومنها:

١. دراسة بعنوان: تنقيب الآراء في جمل المقارنة العربية^١

تناولت هذه الدراسة مشكلة التعرف على جمل المقارنة في تنقيب الآراء المستخدمة في النص العربي، وقد ركزت الدراسة على استخلاص الآراء من جمل المقارنة وذلك بمعرفة المنتج الذي يفضله كاتب الرأي مقارنة مع منتج آخر أو أكثر.

وقد ذكر الباحث أن هناك بعض الأبحاث في هذا المجال بالنسبة لجمل اللغة الإنجليزية وغيرها من اللغات، ولكن بالنسبة للجمل العربية فهذه أول دراسة، كما أن الدراسة استخدمت تقنية تعتمد على التصنيف اللغوي وتقنية أخرى تعتمد على تعلم الآلة.

٢. دراسة بعنوان: دراسة مقارنة لخوارزميات التنقيب في الآراء وتحليل العواطف

وتطبيقاتها^٢

تناولت هذه الدراسة مشكلة تعدد وجهات نظر الزبائن المخزونة في مستودعات البيانات الخاصة بالإنترنت، الشيء الذي أعطى للتنقيب في البيانات وتحليل المشاعر اهتماماً في السنوات الأخيرة، وقد ذكر الباحث أن الناس قد اعتمدوا على الآلة في تصنيف البيانات ومعالجتها، إذ أن توافر كميات هائلة من وجهات النظر حول منتج واحد يساعد على التنبؤ بمشاعر الزبون وذلك عن طريق تحليل وجهات النظر التي تساعد ليس فقط في زيادة الأرباح ولكن أيضاً في تحسين المنتج، وقد قارنت هذه الدراسة التقانات المتوفرة حالياً والمستخدمة في التطبيقات المتعددة في مجال التنقيب في الآراء.

بعد استعراض الدراسات أعلاه نبعت فكرة تصنيف الهمة انطلاقاً من استخدام الباحثين في الدراساتين أعلاه مفهوم التنقيب في النصوص العربية باستخدام طريق التنقيب في النصوص العربية مع اختلاف خوارزميات التصنيف، المستخدمة.

منهجية البحث:

- ١- استخدمنا في هذا البحث أسلوب الاستبانة الرقمية، واقتصرت مهمة الاستبانة في تجميع أكبر قدر من الكلمات المبتدئة بالهمزة وتصنيفها (بواسطة مختصين في اللغة العربية) لهمزتي وصل أو قطع، وفقاً للقواعد النحوية المتبعة في تصنيف همزتي الوصل والقطع.
- ٢- استخدمنا في هذا البحث أيضاً النهج الإحصائي الذي يعتمد على الاستبانة المصممة وعلى (مجموعة تدريب الكلمات) باستخدام نموذج التصنيف، ومن أجل هذا البحث سنستخدم خوارزميات تعلم الآلة وذلك بعد النظر في تكرار الشروط والخصائص (القواعد النحوية المرتبطة بالكلمات المبتدئة بالهمزة والمستخلصة عبر الاستبانة) لاستنتاج العناصر الأساسية في نموذج التصنيف.

هيكل البحث:

يتكون البحث من ثلاثة مباحث:

المبحث الأول (الإطار النظري للبحث): يتناول مفهوم الذكاء الاصطناعي، مفهوم تعلم الآلة، خوارزميات تصنيف البيانات وهمزتي الوصل والقطع.

المبحث الثاني (المعالجة والتطبيق): يتناول مراحل تطبيق نماذج المصنفات التي تشمل: جمع ووصف وتنظيف ومعالجة وترميز البيانات وتحديد المتغيرات المستقلة والتابعة في عينة البيانات كما يتناول أيضاً بناء نماذج المصنفات.

المبحث الثالث (النتائج مخرجات البحث): يتناول اختبار نماذج المصنفات، قياس دقة النماذج، النتائج، ومخرجات البحث والفوائد المرجوة منه.

المبحث الأول: الإطار النظري للمبحث

أولاً: مفهوم الذكاء الاصطناعي:

يُعتبر علم الذكاء الاصطناعي (Artificial Intelligence) هو أحد فروع علوم الحاسب، وأحد ركائز صناعة التقنية في عصرنا الحديث، ويشار له بالاختصار (AI) ويمكن تعريف علم الذكاء الاصطناعي على أنه قدرة الآلات وأجهزة الحاسب الآلي على القيام بمهام تحاكي إلى حد كبير ما يقوم به العقل البشري والذي يتميز بالذكاء، ويمكن تلخيص هذه المهام في القدرة على التفكير أو التعلم من تجاربه السابقة، إذن يمكننا القول بأن الذكاء الاصطناعي يهدف إلى الوصول إلى أنظمة تتصرف وتتعلم وتفهم كما يتصرف ويتعلم ويفهم البشر كما تمتلك خاصية الذكاء^٢.

أنواع الذكاء الاصطناعي:

يمكن تقسيم أنواع الذكاء الاصطناعي حسب قدراته إلى ثلاثة أنواع كالآتي:

١- الذكاء الاصطناعي المحدود:

وهو من الأنواع التي تستطيع القيام بمهام محددة وواضحة، مثل تطبيقات السيارات ذاتية القيادة أو برامج التعرف على الكلام أو الصور، أو لعبة الشطرنج، ويُعتبر هذا النوع من أنواع الذكاء الاصطناعي أكثرها شيوعاً.

٢- الذكاء الاصطناعي العام:

وهو من الأنواع التي لها قدرات تفكيرية تُشابه قدرة الإنسان، إذ أنه يجعل الآلة تكون قادرة على التفكير من تلقاء نفسها وبشابه بشكل كبير التفكير البشري، وفي الحقيقة لا توجد أي تطبيقات عملية لهذا النوع، بل توجد دراسات بحثية تحتاج للكثير من الجهد لتحويلها إلى واقع، وتعد طريقة الشبكات العصبية أحد نماذج الذكاء الاصطناعي العام، إذ أنها تُعنى بإنتاج نظام شبكات عصبية للآلة تكون مُشابهة لتلك التي يحتويها العقل البشري^٣.

٣- الذكاء الاصطناعي غير المحدود:

يعتبر الذكاء الاصطناعي غير المحدود من النوع الذي قد يتجاوز مستوى الذكاء البشري، وقدرته على أداء المهام بصورة يمكن أن تكون أفضل من قدرة البشر المتخصصين ذوي المعرفة، وهذا النوع له العديد من الخصائص الضرورية، مثل: التعلم والتخطيط التلقائي والقدرة على التواصل وإصدار القرار المناسب، لكن مفهوم الذكاء الاصطناعي الفائق يعتبر مفهوماً افتراضياً غير موجود في عصرنا. يمكن أيضاً تصنيف الذكاء الاصطناعي وفقاً للوظائف التي يمكن أن يؤديها إلى الأنواع الأربعة المختلفة الآتية:

١- الآلات التفاعلية:

هي أبسط أنواع الذكاء الاصطناعي نظراً لأنها تفتقر إلى القدرة على التعلم من الخبرة السابقة لتطوير الأعمال المستقبلية، لذا فأما سوف تتفاعل مع الخبرة الحالية لإنتاج أفضل طريقة ممكنة، من أمثلة هذا النوع معدات Deep Blue التي طورها شركة IBM ونظام AlphaGo من شركة جوجل.

٢- الذاكرة المحدودة:

يستطيع الذكاء الاصطناعي من فئة الذاكرة المحدودة (Limited Memory) تخزين بيانات تاريخية سابقة عن النظام الحالي لمدة زمنية مقيّدة، ويُعدّ نَحج القيادة الذاتي من أجدد الأمثلة على ذلك النمط، إذ يقوم بحفظ السرعة الأخيرة للسيارات الأخرى، ومعدل المسافة الفاصلة بين تلك السيارات، والحد الأقصى للسرعة، وغيرها من المعلومات الأخرى الأساسية للقيادة عبر طرق النقل.^٦

٣- نظرية العقل

يُعنى ذلك النوع من أنواع الذكاء الاصطناعي باستيعاب الآلة للمشاعر البشرية، والتفاعل مع البشر والتواصل معهم، وتُعدّ الإشارة أنه لا نجد أية تطبيقات عملية حتى هذه اللحظة على ذلك النمط من الذكاء الاصطناعي.

٤- الإدراك الذاتي:

يُعتبر فئة الإدراك الذاتي (Self-Awareness) من التنبؤات المستقبلية التي يصبو إليها الذكاء الاصطناعي، ويعمل على مبدأ تقني وحسي (إدراكي) حديث للغاية حيث يمكن أن تولد عند الآلة معرفة ذاتية وأحاسيس خاصة بها، الشأن الذي سيجعلها أكثر ذكاءً من الكائن البشري، وما يزال ذلك المفهوم غير حاضر في الواقع.^٧

الحقول الفرعية للذكاء الاصطناعي:

يحتوي علم الذكاء الاصطناعي على العديد من المجالات الفرعية، مثل: التعلم الآلي، والذي يتضمن تمكين أجهزة الكمبيوتر من التعلم بشكل مستقل عن أي تجربة سابقة، حتى تتمكن أجهزة الكمبيوتر من التنبؤ لاتخاذ القرار المناسب بسرعة، يتم من خلال تطوير خوارزمية تسمح بهذا الموقف. وتجدر الإشارة إلى أن هذا المصطلح اقترحه آرثر صموئيل لأول مرة في عام ١٩٥٩. سنشير أدناه إلى بعضاً من أشهر المجالات الفرعية للذكاء الاصطناعي على النحو الآتي:

- تنقيب البيانات:

يُقصد به البحث والتنقيب عن بيانات مُحددة وأنماط مُعينة ضمن مجموعة كبيرة من البيانات بواسطة برامج حاسوبية، إذ تستطيع الشركات الاستفادة من تنقيب البيانات، في تطوير أدائها وزيادة مبيعاتها وتقليص تكاليف الإنتاج^١.

- استرجاع المعلومات والويب الدلالي:

يشير مفهوم استرجاع المعلومات إلى عملية البحث عن أي نوع من البيانات والمستندات التي قد تكون موجودة على الإنترنت من خلال مفهوم الويب الدلالي. تقوم خدمة الويب الدلالية بتحويل البيانات الموجودة على شبكة الويب إلى قاعدة بيانات علمية للمعلومات المترابطة، بحيث يمكن للآلات أن تفهمها، ولا يقتصر استخدامها على البشر. بهذه الطريقة، فيمكن الآلة حجز التذاكر عبر الإنترنت، أو استخدام قاموس على الإنترنت، أو أشياء أخرى تتطلب في البداية استخداماً يدوياً لإكمالها^٢.

- تمثيل المعرفة:

يُعتبر تمثيل المعرفة أحد مجالات الذكاء الاصطناعي الذي يهتم بجعل الآلات تفكر وتتخذ القرار المناسب، حيث أنه يتم تجميع وتخزين المعارف المكتسبة بواسطة الآلة في قاعدة بيانات يتم استخدامها لتبادل المعرفة وإدارة مكوناتها، ونكون مرجعاً لاتخاذ أية قرارات ذكية تصدر عن الآلة في المستقبل.

- التفكير المنطقي والتفكير الاحتمالي

التفكير المنطقي في الذكاء الاصطناعي هو أحد أشكال التفكير المختلفة، لأن الحقائق يتم استنتاجها بناءً على البيانات المتاحة. يتوافق التفكير المنطقي مع ما يسمى بالتفكير الاحتمالي، والذي يستخدم مفاهيم الاحتمالية وعدم اليقين في المعرفة للتعامل مع جميع أوجه عدم اليقين في المستقبل لجميع الأحداث التي قد يشتهب في حدوثها^٣.

- تعلم الآلة:

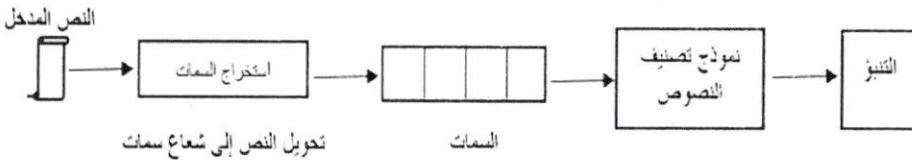
يعد التعلم الآلي أحد فروع الذكاء الاصطناعي، والذي يتضمن تصميم وتطوير خوارزميات وتقنيات تسمح لأجهزة الكمبيوتر بامتلاك خصائص "التعلم". وبشكل عام، ينقسم التعلم إلى مستويين هما: الاستقرائي والاستنباطي، حيث يقوم المنهج الاستنباطي باستنتاج القواعد والأحكام العامة من البيانات الضخمة.

ثانياً: مفهوم تعلم الآلة (Machine Learning):

تتمثل المهمة الرئيسية للتعلم الآلي في استخلاص معلومات قيمة من البيانات، لذا فهي قريبة جداً من استخراج البيانات، يستخدم التعلم الآلي في مجال تحليل البيانات وهو طريقة لتطوير النماذج المعقدة والخوارزميات المناسبة لاستخلاص البيانات باستخدام عمليات تنبؤية، يسمى هذا التحليل بالتحليل التنبؤي. تتيح هذه النماذج التحليلية للباحثين والمحللين البيانات تعلم قرارات ونتائج موثوقة وتستطيع إدراك البيانات المخزنة وعلاقتها.

كما يمكن تعريف نظم تعلم الآلة على أنها أنظمة تقوم بتنبؤات بناءً على ما تعلمته من المعطيات السابقة وتحتاج هذه الأنظمة إلى التدريب على العديد من أمثلة النصوص والتنبؤات (العلامات) المتوقعة لكل منها، وتسمى المعطيات المستخدمة للتدريب بمجموعة بيانات التدريب وتكون هذه المعطيات مصنفة مسبقاً بميزات وكلما كانت مجموعة التدريب أكثر دقة والسمات المختارة مناسبة كانت تنبؤات المصنف أفضل، فعندما يتم تدريب مصنف بطريقة التعلم الآلي يجب تحويل بيانات التدريب إلى شيء يمكن أن تفهمه الآلة، حيث يتم استخراج السمات وتحويلها إلى أشعة (تمثيل النصوص بواسطة الأرقام) والتي سوف تساعد على التعلم من البيانات الموجودة ووضع التنبؤات حول النصوص القادمة^(١)

يمكن للنموذج المدرب استخراج السمات من النص الجديد والتنبؤ أو تصنيف النصوص حسب خصائص محددة باستخدام خوارزميات تصنيف البيانات كما في الشكل (١-١) أدناه:



شكل (٩-١): يوضح الية تصنيف النصوص باستخدام خوارزميات تصنيف البيانات

ثالثاً: خوارزميات تصنيف البيانات:

هناك عدة خوارزميات لتصنيف البيانات تناسب تطبيقات التقيب في البيانات النصية بكل سهولة بعد معالجتها، كما أنها سهلة التدريب سواءً مع الكميات الكبيرة أو الصغيرة من البيانات المقدمة، وفيما يلي سنستعرض أشهر خوارزميات تعلم الآلة الموجه لتصنيف البيانات النصية والتي تم استخدامها في هذا البحث:

١- خوارزمية آلة دعم المتجه SVM-Support Vector Machines:

تدرج هذه الخوارزمية المعروفة اختصاراً بـ (SVM) تحت خوارزميات التعلم الموجه لئلا الذي يعتمد على مجموعة بيانات معروفة النتائج مسبقاً (بمجموعة التدريب) في تدريب الخوارزمية حتى تتمكن من تحليل أي مجموعة جديدة من البيانات وتصنيفها أو تحديد ميولها، وضعت هذه الخوارزمية من قبل العالمان فلاديمير فابنك وأليكسي شيرفونينكييز عام ١٩٦٣م ومن ثم طورها كل من كورينا كورتز وفابنك عام ١٩٩٣م ونشرت عام ١٩٩٥م^{١١}.

تعد خوارزمية (SVM) من أشهر طرق التصنيف الآلي التي تعتمد على إيجاد منحنى أو مستوى فاصل، بفصل العينات التي تم إدخالها عن بعضها البعض، وتتميز الخوارزمية باستخدامها في تصنيف المسائل ذات الفئات الثنائية حصراً.

بصورة عامة يمكننا القول إن خوارزمية آلة دعم المتجه تستخدم للتصنيف وتميز الأنماط. والهدف منها هو إيجاد أفضل دالة تصنيف وأيضاً تحدف إلى التمييز بين أعضاء فئتين من بيانات التدريب، وكما ذكرنا مسبقاً فإن الفكرة من الخوارزمية هي إيجاد مستوى مثالي يفصل بين الفئتين والذي يستخدم للتصنيف وتحديد كل نمط، ومن مميزات الدقة العالية في التصنيف وتطبيق في مجالات واسعة منها تحديد فئات النص حسب تصنيف الصورة^{١٢}.

في حال استطاعت الخوارزمية إيجاد مستوى فاصل يبعد أقل بواحد من بعد متجهات النقاط يكون التصنيف خطياً، وإلا فإن التصنيف يعتبر غير خطي. وفي حالة وجود أكثر من فاصل خطي، يتم اختيار الفاصل الذي يحقق هامش أوسع بين أقرب نقطتين من نوعين مختلفين لبعضهم، وهو ما يسمى بالمستوى ذو الهامش الأكبر^{١٣}.

تحدد دقة الخوارزمية بقدرتها على فصل نوعين بحيث تكون أقرب نقطتين لبعضهما البعض أبعد ما يكون، ويمكننا تسمية المستوى الفاصل بالحافة أو هامش الفصل. وبصفة عامة كلما ازدادت الحافة أو هامش الفصل، كلما قل الخطأ عند تعميم على مجموعة بيانات جديدة.

٢- خوارزمية نايف بيز : (Naive Bayes)

تعتبر خوارزمية نايف بيز (NB) من خوارزميات التعلم الموجه للآلة أيضاً، وتعتمد على قواعد الاحتمال الشرطي التي صاغها العالم توماس بيز^{١٥}، حيث تحسب الاحتمال باستخدام عدد التكرارات للقيم وتكرارات وتركيبات القيم في البيانات المعروفة النتائج مسبقاً (بيانات التدريب). يُعرف مصنف نايف بيز كمجموعة من المصنّفات الاحتمالية البسيطة القائم على فرضية عامة مفادها أن جميع السمات مستقلة عن بعضها البعض وفقاً للمصنف المحدد، ولسهولة تطبيق هذا المصنف وسرعته فهو يعتبر خط الأساس في تصنيف النصوص ويعتبر فعالاً في العديد من المجالات بالرغم من وجود عدد من المصنّفات الأخرى بدقة، أعلى مثل نموذج SVM، حيث أن نموذج Naive Bayes يقوم بتوزيع النصوص لكل صنف باستخدام نموذج احتمالي مع افتراضات مستقلة، هذه الطريقة شائعة جداً في مجال تصنيف النصوص، حيث إن المصنف الثنائي واحداً من أفضل الطرق المعروفة لنموذج Naive Bayes الذي استخدم تمثيل شعاعي ثنائي القيمة للنصوص.

تم إجراء العديد من التحسينات لمصنف بيز منها تعديل حساب الاحتمالات وتقليل السمات وتقليل من الخصائص الأخرى، وحيث إن نظرية بيز تبحث عن احتمال وقوع حدث ما علماً بأن حدث آخر قد وقع مسبقاً.

كما تعتبر خوارزميه نايف بيز واحدة من أهم خوارزميات تعلم الآلة الموجه لعدة أسباب منها: سهولة بناء المصنف، كما أنها لا تحتاج لاستخدام ما يسمى بمخططات التخمين estimation schemes لأي متغيرات تكرارية معقدة، وقد تطبق بسهولة على مجموعة بيانات ضخمة والهدف من الخوارزمية هو بناء قاعدة تسمح بتخصيص هياكل مستقبلية إلى صنف معين وذلك بإعطاء متجهات للمتغيرات التي تصف ذلك الهيكل، ويمكن بواسطتها إجراء المستخدم للعديد من الإحصاءات لسهولة استخدامها^{١٦}.

٣- خوارزمية الجار الأقرب K-Nearest Neighbor

يمكن استخدام خوارزمية الجار الأقرب (KNN) على أنها مصنف بسيط وفعال لتصنيف النصوص، يحتوي المصنف KNN على عيينين هما: التعقيد الحسابي في حال تشابه العينات كبير جداً، وأنه يتأثر أداءها بسهولة في حال كانت العينات التدريبية فردية.

يمكن تقليل تعقيد KNN من خلال استخدام ثلاث طرق: إما بالحد من أبعاد المتجه الممثل للنص، أو بالحد من كمية العينات التدريبية أو بالحد من إيجاد أقرب الجيران أي قيمة الـ K.

تستخدم KNN في تصنيف النصوص عن طريق حساب المسافة بين النص وكل النصوص في مجموعة البيانات التدريبية باستخدام مقياس للاختلاف أو التشابه فيما بينها، ثم العثور على أقرب K مجاورة بين جميع نصوص التدريب ويتم تحديد صنف النص إلى الصنف الذي يضم أكبر عدد من النصوص الموجودة في أقرب الجيران من النصوص وكمباقي الخوارزميات فقد تم التحسين عليها بأكثر من طريقة^{١٧}.

رابعاً: همزتا الوصل والقطع

في هذا الجزء من الجانب النظري سنتناول همزتي الوصل والقطع في أول الكلمة تحت المسميات الآتية:

• معنى الهمزة (لغة):

ورد في كثير من معاجم اللغة ومنها لسان العرب لابن منظور إن الهمزة في اللغة بمعنى الغمز واللمز والضغط^{١٨}.

• معنى الهمزة (اصطلاحاً):

الهمزة في اصطلاح اللغويين وردت لها عدة تعريفات نورد منها تعريف الأزهري في كتابه تهذيب اللغة حيث يقول: "اعلم أن الهمزة لا هجاء لها، وإنما تُكتب مرةً ألفاً، ومرةً واواً، ومرةً ياءً، والألف اللينة لا حرف لها، وإنما هي جزء من ما في بعد فتحة، والحروف ثمانية وعشرون حرفاً مع الواو والألف والياء، وتتم بالهمزة تسعة وعشرين حرفاً، والهمزة كالحرف الصحيح، غير أن لها حالات، من التلين والحذف والإبدال والتخفيف... وليست من حروف الجوف، إنما هي حلقية من أقصى الفم".

• أقسام الهمزة في بداية الكلمة:

تنقسم الهمزة في بداية الكلمة إلى قسمين:

أ/ همزة الوصل:

• تسميتها:

سميت همزة الوصل بهذا الاسم، لأنها يتوصل بها إلى النطق بالحرف الساكن الواقع في ابتداء الكلام، حيث أن القاعدة المهمة في ذلك أنه لا يبتدأ بساكن ولا يوقف بحركة.

• تعريفها:

هي الهمزة التي تثبت في الابتداء، وتستقط في حالة الوصل، وذكر ابن مالك: "إذا سمي بما أوله همزة وصل قطعت الهمزة إن كانت في منقول من فعل مثل: (رأيت اعلم)" يعني إذا سمي به شخص^{١٩}.

• حركة همزة الوصل:

الأصل في همزة الوصل أن تكون مكسورةً مثل: انطلق، ماعدا ألف التعريف تكون مفتوحة، وكذلك ألف أيمن مثل قولك: "أيمن الله بفتح الألف"^{٢٠}.

• مواضع همزة الوصل:

لهمزة الوصل مواضع ثلاث كما يلي:

أ/ الأسماء:

تقع همزة الوصل في أول الكلمة في عشرة أسماء في اللغة هي (اسم، ابن، ابنة، أثنان، اثنتان، امرؤ، امرأة، ابنم، اسم الله، است) ولقد ورد عدد من هذه الأسماء في القرآن الكريم.

ب/ الأفعال:

همزة الوصل في الأفعال المبدوءة بما هي فعلان:

- الفعل الماضي: وتكون همزة الوصل في ماضي الخماسي مثل: انفتح، وكذلك ماضي السداسي مثل: استغفر.

- فعل الأمر: وتكون همزة الوصل في فعل الأمر في المواضع التالية:

○ أمر الفعل الثلاثي مثل: نصر -> انصر.

○ أمر الفعل الخماسي مثل: انتصر للحق.

○ أمر الفعل السداسي مثل: استغفر لذنبك.

ج/ الحروف:

لا توجد همزة الوصل في الحروف إلا في حرف واحد وهو (ال) التعريفية أو ما يعرف باللام

الشمسية مثل: الشمس، وكذلك في اللام القمرية مثل: القمر^{٢١}

أ/ همزة القطع:

• تسميتها:

سميت همزة القطع بهذا الاسم لأنها تقطع وتحجز ما قبلها عمّا بعدها من الحروف مثل قولك: درس - أدرس، فالهمزة قطعت الدال عن الهمزة التي قبلها وكذلك نصر - أنصر، فالهمزة حجزت ما بعدها وهو حرف النون عن الهمزة التي قبلها.

• مواضع همزة القطع:

لهمزة القطع مواضع ثلاث كما يلي:

أ/ الأسماء:

تقع همزة القطع في كل الأسماء في اللغة العربية مثل: (أحمد - أجد - أكرم - أسعد) وكذلك الأسماء الستة فكل ما يبدأ فيها بهمزة فهي همزة قطع مثل: (أخوك - أبوك). ماعدا الأسماء العشرة التي تكون همزتها همزة وصل.

ب/ الأفعال:

تقع همزة القطع في الأفعال كما يلي:

○ ماضي الفعل الثلاثي مثل: أخذ - أكل.

○ ماضي الفعل الرباعي مثل: أكرم - أنجز.

○ أمر الفعل الرباعي مثل: أكرم.

○ المصدر الرباعي مثل: إنذار من أنذر.

ج/ الحروف:

كل الحروف في اللغة العربية همزتها همزة قطع ماعدا (ال) التعريفية مثل: (إن - أن - إذا -

إلى) ^{٢٢}.

• الفرق بين همزتي الوصل والقطع:

لكي نميز ونفرض بين همزة الوصل وهمزة القطع لابد من معرفة الآتي:

- تأتي همزة القطع ساكنة أو متحركة.

- همزة الوصل لا تأتي إلا متحركة.

- همزة القطع تثبت في بداية الكلام وفي أثنائه وفي طرفه مثل: سأل في وسطها وقرأ في طرفها.

- همزة الوصل تسقط أثناء الكلام بينما تثبت في البدء فقط.

- همزة الوصل لا تكون إلا زائدة بينما همزة القطع تكون أصلية أو زائدة ^{٢٣}.

المبحث الثاني: التطبيق والتنفيذ

تمهيد:

إن البيانات في يومنا هذا تتضخم يوماً بعد يوم وتتعدد مصادرها، وهذا يقود إلى تعرض هذه البيانات إلى الكثير من المشاكل التي تقلل من جودة البيانات مثل كثرة البيانات المفقودة وعدم تناسق البيانات، لذلك قمنا في هذا البحث بتقسيم مراحل تطبيق نماذج مصنفات البيانات إلى ستة مراحل بدءاً بمرحلة تصميم الاستبانة وجمع البيانات وانتهاءً بمرحلة قياس دقة نماذج المصنفات، كما في الشكل (٢-١) التالي:



شكل رقم (٢-١): مراحل تطبيق نماذج المصنفات

في هذا المبحث سنتناول الأربعة مراحل الأول من مراحل تطبيق نماذج المصنفات وسيتم تناول المرحلتين المتبقيتين في المبحث الثالث.

- مرحلة تصميم الاستبانة ووصف البيانات

تم تجميع الكلمات العربية المبتدئة بالهمزة وذلك عن طريق تصميم استبانة رقمية انظر (ملحق أ) وقد تم تصميمها باستخدام نماذج جوجل وتم تسميتها (استبانة تصنيف الكلمات المبتدئة بـهمزة وصل أو قطع)، وقد اقتصرتم مهمة الاستبانة في تجميع أكبر قدر من الكلمات المبتدئة بالهمزة بواسطة متخصصين في اللغة العربية وقد قام المتخصصون بتصنيفها لهمزتي وصل أو قطع وفقاً للقواعد النحوية المعروفة، كما تم نشر الاستبانة على عنوان الويب التالي:

<https://docs.google.com/forms/d/e/1FAIpQLSdtwnfWQvhoTdVbKfZUHDVixξfmXPErgrXMiQGIDem9Dsc.Q/viewform>

تمت تعبئة بواسطة ٥٠ متخصص في النحو بدرجات علمية متفاوتة بهدف الحصول على عينة متجانسة، وقد بلغ عدد الكلمات الكلية المصنفة ٤٠٠ كلمة.

تم التصنيف للكلمات المستخلصة من نصوص العينة إلى كلمات تبدأ بهمزة وصل وأشير إليها في هذا البحث بـ (Wasl) وكلمات تبدأ بهمزة قطع أشير إليها في هذا البحث بـ (Gtaa) وقد تم تحديد المتغير التابع (Val) وفقاً للقيمتين (Wasl/Gtaa) أما المتغيرات المستقلة (الخصائص) فقد تم تقسيمها لثلاث خصائص طبقاً للقواعد النحوية المشار إليها لاحقاً في الجدول (١-٢).

تم اختيار صيغة ملف مجموعة البيانات (Dataset) من نوع (csv) وقد احتوى على الخصائص (المتغيرات المستقلة) وهي:

- المتغير (diacritic): متغير رقمي بطول (١) رقم يشير إلى حركة حرف الألف في بداية الكلمة (فتحة، كسرة، ضمة) تم تمثيلها بالقيم الرقمية (١، ٢، ٣) على التوالي.
- المتغير (Count): متغير رقمي يشير لعدد حروف الكلمة حيث إن الكلمة التي تحتوي على حرفين، ثلاثة أحرف أو أربعة... إلخ، يتم تمثيلها رقمياً بالقيم (٢، ٣ أو ٤... ١٠) على التوالي.
- المتغير (morphological): متغير رقمي يشير للوزن الصرفي للأفعال والمصادر (فعل، أفعال، أنفعال، فعال) يتم تمثيلها رقمياً بالقيم (٠، ١، ٢، ٣... ١١) على التوالي.
- المتغير (noun): متغير رقمي يشير إلى أن الكلمة (ليست اسم، اسم عادي، من الأسماء العشرة، من الأسماء الستة، اسم موصول، اسم فاعل) وتم تمثيله رقمياً بالأرقام (٠، ١، ٢، ٣، ... ٦) على التوالي.
- المتغير (verb): متغير رقمي يشير إلى أن الكلمة (ليست فعلاً، فعل ماضي، فعل مضارع، فعل أمر) وتم تمثيله رقمياً بالأرقام (٠، ١، ٢، ٣) على التوالي.

- المتغير (adjective): متغير رقمي يشير إلى أن الكلمة (ليست صفة، صفة) وتم تمثيله رقمياً بالأرقام (٠، ١) على التوالي.
- المتغير (letter): متغير رقمي يشير إلى أن الكلمة (ليست حرفاً، حرف) وتم تمثيله رقمياً بالأرقام (٠، ١) على التوالي.
- المتغير (the): متغير رقمي يشير إلى أن الكلمة (غير معرفة بـ آل، معرفة بـ آل) وتم تمثيله رقمياً بالأرقام (٠، ١) على التوالي.
- تم تسمية ملف مجموعة البيانات (Dataset) بالاسم (Arabic Words)، الجدول (١-٢) يوضح جانب من توصيف بيانات الملف (Arabic_Words.csv) كما يلي:

191	افتعل	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
192	افتقد	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
193	اكتفى	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
194	امتحن	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
195	انتصر	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
196	انقذ	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
197	انتكس	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
198	انتمى	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
199	انتهى	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
200	اكتشف	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
201	التف	5	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افتعل	غير معرفة بال	كسرة	Wasl
202	أحمد	4	اسم عادي	ليس فعلاً	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa
203	أيمن	4	من الأسماء العشرة	ليس فعلاً	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa
204	أثنى	4	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa
205	أبعد	4	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa
206	أبلغ	4	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa
207	أثرى	4	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa
208	أحرم	4	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa
209	أدرج	4	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa
210	أسرى	4	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa
211	أسلم	4	ليس اسماً	فعل ماضي	ليس صفة	ليس حرفاً	افعل	غير معرفة بال	فتحة	Gtaa

وبالتأكيد عندما تكون جودة البيانات منخفضة فهذا سيؤثر حتماً على نتائج التحليل. في هذا البحث استخدمنا عدة أساليب لتنظيف البيانات (Cleaning Data) على النصوص المجمعة، وقد اشتملت مرحلة تنظيف البيانات على المراحل الآتية:

-التعامل مع البيانات المفقودة

- حذف البيانات المكررة.

بعد إجراء عمليات تنظيف البيانات (Cleaning Data) على النصوص المجمعة، وبعد معالجة واستبعاد الكلمات المكررة وكذلك إكمال البيانات الناقصة أصبح عدد الكلمات المستبعدة ١٠١ كلمة، حصلنا على عدد ٢٩٩ كلمة صالحة لتطبيقها على نموذج المصنف.

• مرحلة ترميز وتمثيل البيانات:

بعد دراسة القواعد النحوية التي توضح مواضع همزتي الوصل والقطع في أول الكلمة تم تحديد الأسس والخصائص (Features) التي يمكن الاعتماد عليها في تحديد قيمة المتغير التابع (Outcome)، الجدول (٢-٢) يوضح القواعد المتبعة لتحديد مواضع همزتي الوصل والقطع.

ال تعريف	عدد الحروف	الميزان الصرفي		حركة الهمزة		صفة		حرف		فعل		اسم	
		ال	Count	morphology	diacritic	Subjective	Letter	verb	Noun				
ال تعريف	ال	Count	morphology	diacritic	Subjective	Letter	verb	Noun					
ترويز	خانات	ترويز	خانات	ترويز	خانات	ترويز	خانات	ترويز	خانات	ترويز	خانات	ترويز	خانات
٠	غير مضافة	٢	حرفين	١	فعل	١	فتحة	٠	ليس حرفاً	٠	ليس فعلاً	٠	ليس اسماً
١	مضافة	٣	ثلاثة	٢	افعل	٢	كسرة	١	حرف	١	ماضي	١	اسم عادي
		٤	اربعة	٣	انفعل	٣	ضمة			٢	مضارع	٢	الاسماء العشرة
		٥	خمس	٤	افتعال					٣	امر	٣	الاسماء الستة
		٦	سنة	٥	افاعيل							٤	اسماء الموصول
		٧	سبعة	٦	افعال							٥	اسماء الافعال
		٨	ثمانية	٧	افتعل								
		٩	تسعة	٨	استفعل								
		١٠	عشرة	٩	استفعال								
				١٠	انفعال								
				١١	فعال								

الجدول (٢-٢) يوضح القواعد المتبعة لتحديد مواضع همزتي الوصل والقطع

تم ترميز البيانات الموجودة في جدول (٢-١) وتحويل القيم الحرفية إلى قيم رقمية حتى تتمكن الخوارزميات من التعامل معها وفقاً للجدول (٢-٢) الذي يوضح القواعد المتبعة لتحديد مواضع همزتي الوصل والقطع وتصيح البيانات بعد ترميزها كما في الجدول (٢-٣) الآتي:

Outcome	diacritic	the	morphological	letter	adjective	verb	noun	Count	Word
Gtaa	1	0	2	0	1	0	0	4	أصفر
Gtaa	1	0	2	0	1	0	0	4	أعذب
Gtaa	1	0	2	0	1	0	0	4	أعزب
Gtaa	1	0	2	0	1	0	0	4	أعور
Gtaa	1	0	2	0	1	0	0	4	أسود
Gtaa	1	0	2	0	1	0	0	4	أعزل
Gtaa	1	0	2	0	0	0	1	4	أحمد
Gtaa	1	0	2	0	0	1	0	4	أثنى
Gtaa	1	0	2	0	0	0	2	4	أبين
Wasl	2	0	3	0	0	1	0	5	انحصر
Wasl	2	0	3	0	0	1	0	5	انحصر
Wasl	2	0	3	0	0	1	0	5	انطلق
Wasl	2	0	3	0	0	1	0	5	انفعل
Wasl	2	0	4	0	0	0	1	6	اجتهد
Wasl	2	0	4	0	0	0	1	6	اجتار
Wasl	2	0	4	0	0	0	1	6	ابنأه
Wasl	2	0	4	0	0	0	1	6	ابنأه
Wasl	2	0	4	0	0	0	1	6	ابنأه
Wasl	2	0	4	0	0	0	1	6	ابنأه
Wasl	2	0	4	0	0	0	1	6	اجتهد

جدول (٢-٣) ترميز البيانات في ملف مجموعة البيانات (arabic.csv)

بعد اكتمال عملية ترميز البيانات يتم تمثيل البيانات باستخدام لغة بايثون Python

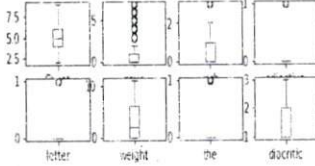
واستيراد مكتبات (numpy- scipy- sklearn) باستخدام محرر (Jupyter) كما في

الشكل (٢-٢) الآتي:

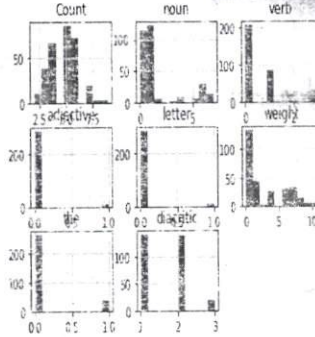
4. Data Visualization

4.1 Univariate Plots

```
In [20]: # box and whisker plots
dataset.plot(kind='box', subplots=True, layout=(4,4), sharex=False, sharey=False)
pyplot.show()
```



```
In [21]: # histograms
dataset.hist()
pyplot.show()
```



شكل (٢-٢): تمثيل البيانات باستخدام محرر (Jupyter Notebook)

• مرحلة بناء وتدريب نماذج المصنفات:

تم تقسيم البيانات المثلة إلى بيانات تدريبية بنسبة (٦٧٪) وبيانات تجريبية بنسبة (٣٣٪) من أصل البيانات الكلية (٢٩٩) سجل تمهيداً لبناء نموذج تدريب البيانات باستخدام خوارزميات التصنيف المختارة (خوارزمية الة دعم المتجه SVM - خوارزمية نايف بيز NB - خوارزمية الجار الأقرب KNN) وقد تم اختيارها نسبة لتناسب حجم العينة وقيم البيانات مع الخوارزميات أعلاه وقد تم بناء النموذج باستخدام دوال مكتبة (Sk-Learn) كما في الشكل (٢-٣) أدناه:

5.2 Build Models

```
In [35]: # Spot Check Algorithms
models = []

models.append(('KNN', KNeighborsClassifier()))
models.append(('NB', GaussianNB()))
models.append(('SVM', SVC(gamma='auto')))
# evaluate each model in turn
results = []
names = []
for name, model in models:
    # Test options and evaluation metric
    kfold = StratifiedKFold(n_splits=10, random_state=None)
    cv_results = cross_val_score(model, X_train, Y_train, cv=kfold, scoring='accuracy')
    results.append(cv_results)
    names.append(name)
    print('%s: %f (%f)' % (name, cv_results.mean(), cv_results.std()))

KNN: 0.940000 (0.037417)
NB: 0.540000 (0.053852)
SVM: 0.955000 (0.041533)
```

شكل (٢-٣): بناء نماذج التصنيف باستخدام محرر (Jupyter Notebook)

المبحث الثالث: النتائج ومخرجات البحث

أولاً: اختبار نماذج المصنفات

بعد اكتمال مرحلة بناء نماذج المصنفات باستخدام خوارزميات (SVM, NB, KNN) انتقلنا لمرحلة اختبار هذه النماذج حتى تتمكن من قياس جودتها فيما بعد، وقد تم اختبار النماذج بإجراء تنبؤات على النماذج بعد تدريبها باستخدام دوال مكتبة (sklearn)، الشكل (١-٣) يوضح اختبار نماذج المصنفات باستخدام دالة التنبؤ (predict) كما يلي:

6. Make Predictions

6.1 Make Predictions

```
In [47]: # Make predictions on validation dataset For SVM
model1 = SVC(gamma='auto')
model1.fit(X_train, Y_train)
predictions1 = model1.predict(X_validation)
```

```
In [51]: # Make predictions on validation dataset FOR Naive Bayes
from sklearn.naive_bayes import MultinomialNB
model2 = MultinomialNB()
model2.fit(X_train, Y_train)
predictions2 = model2.predict(X_validation)
```

```
In [76]: # Make predictions on validation dataset FOR Naive Bayes
from sklearn.naive_bayes import MultinomialNB
model3 = KNeighborsClassifier(n_neighbors = 6)
model3.fit(X_train, Y_train)
predictions3 = model3.predict(X_validation)
```

شكل (٣-١) يوضح اختبار نماذج المصنفات باستخدام دالة التنبؤ (predict)

ثانياً: قياس دقة نماذج المصنفات

تعتبر مرحلة تقييم نتائج نماذج التنقيب في البيانات من المراحل المهمة التي يمكننا من تعريف النموذج الأكثر فعالية، وتقاس فعالية النموذج من خلال دقة الخطة المعمول بها وتلعب طبيعة البيانات المستخدمة في بناء النماذج دوراً أساسياً في فعاليتها، ويوجد العديد من الطرق الإحصائية التي تختبر نماذج التصنيف نبين أهمها فيما يلي:

- حساب متوسط الدقة Average Accuracy

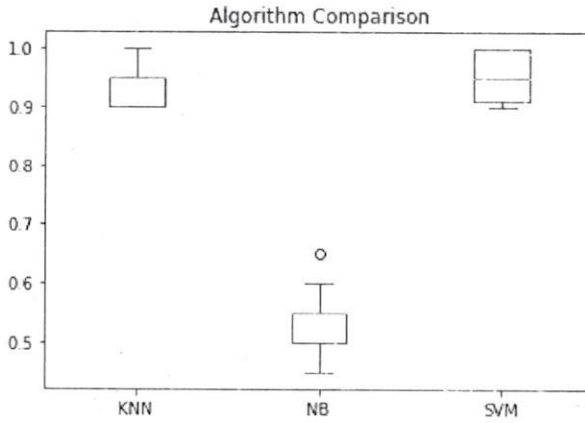
هو المتوسط الحسابي لنسب دقة التوقعات الصحيحة لكل فئة يقدمها النموذج إلى عدد التصنيفات الفعلية لهذه الفئة في مجموعة بيانات الاختبار، وقد بلغ متوسط دقة التوقعات للخوارزميات المستخدمة كما هو مبين في الجدول (٣-١) الذي يبين أن خوارزمية (SVM) قد حصلت على أعلى نسبة عند حساب متوسط الدقة حيث بلغت دقتها (٩٤٪) بينما حصلت خوارزمتي (NB) و (KNN) على متوسط دقة (٥٤٪) و (٨٥٪) على الترتيب كما يلي:

متوسط الدقة	الخوارزمية
٪٩٤	SVM
٪٥٤	NB
٪٨٥	KNN

جدول (٣-١): متوسط دقة الخوارزميات

الشكل (٣-٢) يوضح رسم بياني لمتوسط دقة خوارزميات (SVM) و (NB) و

(KNN)



شكل (٣-٢): رسم بياني لمتوسط دقة الخوارزميات

- حساب الدقة الإجمالية Total Accuracy

هي نسبة مجموع التوقعات الصحيحة المقدمة من النموذج إلى مجموع التصنيفات الفعلية في

مجموعة بيانات الاختبار، وقد بلغت الدقة الإجمالية للخوارزميات المستخدمة كما هو مبين في

الجدول (٣-٢):

الدقة الإجمالية	الخوارزمية
٪٩٢	SVM
٪٧٨	NB
٪٩١	KNN

جدول (٢-٣): الدقة الاجمالية للخوارزميات

من الجدول أعلاه يتضح أن خوارزمية KNN حصلت على دقة إجمالية بلغت ٩١٪. مقارنة بخوارزمية SVM التي حصلت على أعلى معدل دقة إجمالية بلغت ٩٢٪. بينما حصلت خوارزمية NB على أقل معدل دقة إجمالية بنسبة ٧٨٪.

- إيجاد مصفوفة الارتباك Confusion Matrix

تظهر مصفوفة الارتباك عدد الحالات المتوقعة بشكل صحيح وعدد الحالات المتوقعة بشكل خاطئ في مجموعة الاختبار لكل صنف من الأصناف مقارنة مع عدد الحالات الفعلية لتلك الأصناف. رتبة هذه المصفوفة هي $n \times n$ حيث n هي عدد الأصناف في عمود المتغير الهدف (المتغير التابع)، وقد تم إيجاد مصفوفة الارتباك للخوارزميات المستخدمة لتصنيف المتغير التابع (Outcome) إلى صنفين (همزة وصل Wasl وهمزة قطع Gtaa) كما في الشكل (٣-٣) التالي:

Confusion matrix for SVM:

: Gtaa pred: Wasl

true: Gtaa	٤٩	٥
true: Wasl	٣	٤٢

confusion matrix for NB:

pred: Gtaa pred: Wasl

true: Gtaa	٤١	١٣
true: Wasl	٩	٣٦

confusion matrix for KNN:

pred: Gtaa pred: Wasl

true: Gtaa	٤٩	٥
true: Wasl	٤	٤١

شكل (٣-٣): مصفوفة الارتباك لخوارزميات (SVM, NB, KNN)

المستخدمة في تصنيف همزتي الوصل والقطع

ثالثاً: النتائج

وقد قمنا باستخدام خوارزمية آلة دعم المتجه (SVM) وخوارزمية نايف بيز (NB) وخوارزمية الجار الأقرب (KNN) لتصنيف همزتي الوصل والقطع (في بداية الكلمة) وقد تم اختيار هذه الخوارزميات بناء على عوامل عدة أبرزها تناسب حجم مجموعة البيانات مع هذه الخوارزميات وكثرة الخصائص المعتمد عليها في عملية التصنيف، تم عمل نموذج تدريب لجميع هذه الخوارزميات وقد اتضح من خلال حساب دقة النموذج لكل خوارزمية تفوق خوارزمية آلة دعم المتجه (SVM) حيث حصلت على أعلى نسبة دقة للنموذج بلغت ٩٢٪، وهي نسبة عالية وكافية لحل مشكلة البحث الأساسية، وقد كان السبب الأساسي في حصول الخوارزمية على هذه النسبة العالية هو كثرة عدد الخصائص المستخدمة في التصنيف (تم استخدام ٨ متغيرات مستقلة) مما أسهم في جعل نموذج التصنيف أكثر واقعية.

رابعاً: مخرجات البحث والفوائد المرجوة منه:

كان الهدف من هذا البحث المنجز هو بناء نموذج ذكي يعمل على تصنيف همزتي الوصل والقطع (في بداية الكلمة) من خلال استخدام خوارزميات تصنيف البيانات من أجل تأسيس معايير دقيقة وصحيحة في كتابة النصوص العربية بصورة دقيقة لتسهل وتساعدهم في تطوير التقنية لخدمة اللغة العربية، كما هدف البحث أيضاً لقياس جودة أشهر خوارزميات التصنيف في التمييز بين همزتي الوصل والقطع في بداية الكلمة، ويمكن لهذا النموذج الذكي الذي تم تصميمه أن يساهم في تطوير اللغة العربية كما يلي:

- النموذج الجديد يمكن استخدامه في مراجعة البحوث العلمية والتأكد من رسم الهمزة في الكلمات بصورة صحيحة مما يساهم في اكتمال عناصر البحث العلمي.
- النموذج الجديد يمكن استخدامه في مراجعة المقالات الإخبارية المنشورة في وسائل التواصل الاجتماعي المختلفة من خلال تثبيته في جزء الإضافات الخاصة بالمستعرض للتأكد من رسم الهمزة بشكل صحيح في جميع كلمات المقالات المنشورة مما يساهم في تطوير الإعلام العربي الجديد.

- النموذج الجديد يمكن استخدامه مع تطبيقات الجوال و اعتباره كجزء أصيل من أنظمة تشغيل الجوال للتأكد من كتابة الأوامر المترجمة للغة العربية والمبتدئة بهمزة بشكل صحيح.

- النموذج الجديد يمكن دمج مع أنظمة الترجمة المستخدمة في المؤتمرات والطائرات والقطارات حيث يساعد في عرض الهمزة بصورة صحيحة.

خلاصة القول يمكننا استخدام النموذج الجديد كأداة برمجية يمكن أن يتم دمجها مع كافة الأنظمة البرمجية والتقنية التي تعرض نصوص عربية مقروءة أو مكتوبة باستخدام أنظمة برمجية لتصنيف الهمزة في الكلمات المبدوءة بهمزة وصل أو قطع.

التوصيات

١- تطبيق خوارزميات تصلح للتعامل مع الأغراض ذات الفئات التصنيف المتعددة باستخدام خوارزميات فهرسة أخرى غير التي تعالج الأغراض ذات الفئات الثنائية فقط.

٢- تطوير مكتبات sklearn لتدعم اللغة العربية بشكل كامل.

٣- تطوير النموذج ليعمل كمصنف للنصوص العربية المسموعة.

الهوامش:

- ١ - د. علاء مصطفى الهليس: تنقيب الآراء في جمل المقارنة العربية، المحلة العربية الدولية للمعلوماتية، المجلد الثاني، العدد الرابع، ٢٠١٣م.
- ٢ د. رنا زهير عبد الغني العبيدي، د. غيداء عبد العزيز الطالب: دراسة مقارنة لخوارزميات التنقيب في الآراء وتحليل العواطف وتطبيقاتها، مجلة الرافدين لعلوم الحاسوب والرياضيات المجلد (١٢)، العدد الثاني، ٢٠١٨م.
- ٣ امل كاظم ميرة &، تحرير جاسم كاطع. (٢٠١٩). تطبيقات الذكاء الاصطناعي في التعليم من وجهة نظر تدريسي الجامعة.
- ٤ امل كاظم ميرة &، تحرير جاسم كاطع. (٢٠١٩). تطبيقات الذكاء الاصطناعي في التعليم من وجهة نظر تدريسي الجامعة.
- ٥ د. م. مصطفى عبيد، (٢٠١٨). التحليل المتقدم وتنقيب البيانات، دار الفكر العربي، الطبعة الأولى، المجلد الأول، القاهرة.

٦ Burkov, A. (٢٠١٩). The hundred-page machine learning book (Vol. ١). Canada: Andriy Burkov.

٧ Florin, G. (٢٠١١). Data mining: concepts, models and techniques. Springer-Verlag. Berlin Heidelberg.

٨ Ian, H. W and Eibe, F. (٢٠٠٥). Data Mining: practical machine learning tools and techniques. Second Edition. Elsevier Inc. San Francisco: USA.

٩ Jiawei, H., Micheline, K. and Jian P. (٢٠١٢). Data mining: concepts and techniques, Third edition. Elsevier Inc: USA.

١٠ Kalyani, G. and Jaya, A. Lakshmi Performance assessment of different classification techniques for intrusion detection. Journal of Computer Engineering (IOSRJCE).

¹¹ Kalyani, G. and Jaya, A. Lakshmi Performance assessment of different classification techniques for intrusion detection. Journal of Computer Engineering (IOSRJCE).

¹² Michael, J.A. B. and Gordon, S. L. (٢٠٠٤). Data mining techniques for marketing, sales, and customer relationship management, Second edition. Wiley Publishing, Inc. Indianapolis, Indiana: USA.

¹³ امل كاظم ميرة &، تحرير جاسم كاطع. (٢٠١٩). تطبيقات الذكاء الاصطناعي في التعليم من وجهة نظر تدريسي الجامعة.

¹⁴ د. م. مصطفى عبيد، (٢٠١٨). التحليل المتقدم وتنقيب البيانات، دار الفكر العربي، الطبعة الأولى، المجلد الأول، القاهرة.

¹⁵ العالم الإحصائي الإنجليزي توماس بيز (Thomas Bayes) عاش خلال الفترة (١٧٠١-١٧٦١م). هو من قام بصياغة حالة خاصة من النظرية المشهورة والتي تحمل اسمه وهي نظرية بيز (Bayes' theorem) رغم أنها لم تنشر في حياته وإنما نشرت بعد وفاته بواسطة ريتشارد برايس (Richard Price).

¹⁶ د. م. مصطفى عبيد، (٢٠١٨). التحليل المتقدم وتنقيب البيانات، دار الفكر العربي، الطبعة الأولى، المجلد الأول، القاهرة.

¹⁷ امل كاظم ميرة &، تحرير جاسم كاطع. (٢٠١٩). تطبيقات الذكاء الاصطناعي في التعليم من وجهة نظر تدريسي الجامعة.

¹⁸ محمد بن مكرم بن علي، أبو الفضل، جمال الدين ابن منظور الأنصاري الرويفعي الإفريقي. (١٩٩٣). لسان العرب. دار صادر - بيروت الطبعة: الثالثة (١٧/١).

¹⁹ محمد بن عبد الله، ابن مالك الطائي الجباني، أبو عبد الله، جمال الدين. تحقيق: عبد المنعم أحمد هريدي. (د.ت). شرح الكافية الشافية. جامعة أم القرى مركز البحث العلمي وإحياء التراث الإسلامي كلية الشريعة والدراسات الإسلامية مكة المكرمة الطبعة: الأولى (١٤٦٦/٣).

- ٢٠ أبو الفتح عثمان بن جني الموصلي ، تحقيق: فائز فارس. (د.ت) اللمع في العربية. دار الكتب الثقافية - الكويت (٢٢٦-٢٢٥).
- ٢١ محمد رفيق مؤمن الشويكي. (٢٠١٥) اللآلئ الذهبية في شرح المقدمة الجزرية. غزة- فلسطين الطبعة: الأولى (٧٧).
- ٢٢ عبد الله محمد النقراط ، تحقيق: محمد خليل هراس. (٢٠٠٣) الشامل في اللغة العربية. دار الكتب الوطنية - لبنان الطبعة: الأولى (١٦٢-١٧٨).
- ٢٣ الطاهر بن محمد زواوي البيروني الجزائري. الشيخ. (٢٠١٩). الميسر المفيد في فن التلاوة والتجويد من قراءة نافع المدني وعاصم الكوفي ومن رواية ورش وقالون. دار الكتب العلمية - بيروت - لبنان.

المصادر والمراجع

أولاً: المصادر والمراجع العربية

- [١] محمد بن مكرم بن علي، أبو الفضل، جمال الدين ابن منظور الأنصاري الرويعي الإفريقي. (١٩٩٣). لسان العرب. دار صادر - بيروت الطبعة: الثالثة (١٧/١).
- [٢] محمد بن أحمد بن الأزهرى الطروي، أبو منصور. (٢٠٠١) تهذيب اللغة. دار إحياء التراث العربي - بيروت الطبعة: الأولى، (٣٣/١).
- [٣] محمد بن عبد الله، ابن مالك الطائي الجبائي، أبو عبد الله، جمال الدين. تحقيق: عبد المنعم أحمد هريدي. (د.ت). شرح الكافية الشافية. جامعة أم القرى مركز البحث العلمي وإحياء التراث الإسلامي كلية الشريعة والدراسات الإسلامية مكة المكرمة الطبعة: الأولى (١٤٦٦/٣).
- [٤] أبو الفتح عثمان بن جني الموصلي، تحقيق: فائز فارس. (د.ت) اللمع في العربية. دار الكتب الثقافية - الكويت (٢٢٦-٢٢٥).
- [٥] محمد رفيق مؤمن الشوبكي. (٢٠١٥) اللآلئ الذهبية في شرح المقدمة الجزرية. غزة- فلسطين الطبعة: الأولى (٧٧).
- [٦] عبد الله محمد النفراط، تحقيق: محمد خليل هراس. (٢٠٠٣) الشامل في اللغة العربية. دار الكتب الوطنية - لبنان الطبعة: الأولى (١٧٨-١٦٢).
- [٧] الطاهر بن محمد زاوي البيريني الجزائري، الشيخ. (٢٠١٩). المسر المفيد في فن التلاوة والتجويد من قراءة نافع المدني وعاصم الكوفي ومن رواية ورش وقالون. دار الكتب العلمية - بيروت - لبنان.

[٨] م. د. نعيم سلمان البديري. (٢٠٠٨) همزة الوصل في اللغة العربية. جامعة واسط - مجلة كلية التربية - لبنان المجلد: الأول (٢٩-١٥).

[٩] الهليس ، علاء مصطفى. (٢٠١٣). تنقيب الآراء في جمل المقارنة العربية Arab

International Informatics Journal, ١٧٨(٢١١٧), ١-١٤

[١٠] د. رنا زهير عبد الغني العبيدي، د. غيداء عبد العزيز الطالب: دراسة مقارنة لخوارزميات التنقيب في الآراء وتحليل العواطف وتطبيقاتها، مجلة الرادين لعلوم الحاسوب والرياضيات المجلد (١٢)، العدد الثاني، ٢٠١٨ م.

[١١] امل كاظم ميرة & تحرير جاسم كاطع. (٢٠١٩). تطبيقات الذكاء الاصطناعي في التعليم من وجهه نظر تدريسي الجامعة. (٢٢). Psychological Science,

[١٢] د. م. مصطفى عبيد، (٢٠١٨). التحليل المتقدم وتنقيب البيانات، دار الفكر العربي، الطبعة الأولى، المجلد الأول، القاهرة.

ثانياً: المصادر والمراجع الأجنبية:

[١٣] Burkov, A. (٢٠١٩). The hundred-page machine learning book (Vol. ١). Canada: Andriy Burkov.

[١٤] Florin, G. (٢٠١١). Data mining: concepts, models and techniques. Springer-Verlag. Berlin Heidelberg.

[١٥] Ian, H. W and Eibe, F. (٢٠٠٥). Data Mining: practical machine learning tools and techniques, Second Edition. Elsevier Inc. San Francisco: USA.

[١٦] Jiawei, H., Micheline, K. and Jian P. (٢٠١٢). Data mining: concepts and techniques, Third edition. Elsevier Inc: USA.

[١٧] Kalyani, G. and Jaya, A. Lakshmi Performance assessment of different classification techniques for intrusion detection. Journal of Computer Engineering (IOSRJCE).

[١٨] Michael, J.A. B. and Gordon, S. L. (٢٠٠٤). Data mining techniques for marketing, sales, and customer relationship management, Second edition. Wiley Publishing, Inc. Indianapolis, Indiana: USA.

ملحق (أ)

نموذج استبانة تصنيف الكلمات المبتدئة بهمزة وصل أو قطع

استبانة تصنيف الكلمات المبتدئة بهمزة لهزمة وصل أو قطع

*مضروب

اختر التخصص *

لغة عربية أخرى:

الترجمة العلمية *

معيد محاضر أستاذ مساعد أستاذ مشارك أستاذ

صفحة 1 من 2

التالي

ذلك كلمة كلمة تكون متبينة بهمزة وصل أو قطع

السؤال متعلقة بالكلمة التي اخترتها

اكتب الكلمة *

حافظك

عدد حروف الكلمة *

جاءك

هل الكلمة اسم *

اسم عادي

من الأسماء السكينة

من الأسماء العشرة

اسم موصول

اسم اعجمي

ليس اسماً

هل الكلمة فعل *

ليس فعلاً

فعل منطسي

فعل مضارع

فعل أمر

هل الكلمة حرفاً* ؟

ليس حرفاً

حرف

هل الكلمة صفة* ؟

ليست صفة

صفة

حدد حركة الهمزة* ؟

فتحة

كسرة

ضمة

هل الكلمة معرفة بـ (ال)* ؟

تتبع معرفة بـ (ال)

معرفة بـ (ال)

