

التعرف على المحارف العربية الأياً ضمن نصوص صور الوثائق متعددة اللغات

د. فوزي عبد الواحد العضاض
الجامعة المستنصرية
بغداد - العراق

أ.د. هلال محمد يوسف
جامعة العلوم التطبيقية
المنامة - البحرين

أ.د. عبد المنعم صالح رحمة*
جامعة التكنولوجيا
بغداد - العراق

ملخص

المعالجة الآلية للوثائق المتعددة اللغات تحتوي على عدد من المشاكل الصعبة. لأن تقنيات التعرف على الأحرف ليست متوفرة لكل أصناف اللغات. الوثيقة المصورة يجب أن تصنف طبقاً لنوع محارف اللغة أولاً، في هذا البحث تم التعرف على المحارف الرومانية (الانكليزية) أو العربية. تم احتساب ستة ميزات لتمثيل محارف الوثيقة ولكل ميزه تم احتساب ثلاثة معايير احصائية هي المتوسط، الانحراف المعياري والانحراف. استخدم البحث مصنف الذي يعتمد على مسافة "موهالمبس (Mahalanobis)" لتصنيف الوثيقة إلى الرومانية (الانكليزية) أو العربية. تم الحصول على 100 % دقة تصنيف لوثائق الاختبار..

1. مقدمة

تصنيف المحارف الى اللغة التي تنتمي اليها اصبحت عملية اساسية في المعالجة الآلية لصور الوثائق متعددة اللغات. هذه العملية تساعد على الاختيار الأمثل لنظام الألي للتعرف على الاحرف. تعتبر هذه العملية مرحلة اساسية قبل عملية الترجمة خصوصاً للأشخاص غير ملمين باللغة العربية، حيث قد تختلط لديهم شكل الاحرف العربية مع اشكال الاحرف في لغات أخرى. تصنيف محارف اللغة في اللغات التي تكون فيها الحروف مستخدمة فقط من لغة واحدة تعتبر عملية ضمن عملية التعرف الآلي على الاحرف. لكن بالنسبة للاحرف العربية والرومانية (الانكليزية) حيث تستخدم هذه الحروف في اكثر من لغة تعتبر عملية تصنيف الاحرف الى اللغة التي تنتمي اليها عملية اساسية اولية قبل عملية التعرف الآلي على الاحرف. في مجال تصنيف المحارف العربية هنالك عدة معوقات اساسية منها تعقيد النص العربي بسبب الانحناءات المتكررة للاحرف العربية، الاشكال المختلفة للاحرف العربية بحسب مواقعها داخل الكلمة، هذا بالإضافة الى تراكب وترابط الاحرف العربية في الكلمات بحسب موقعها على خط الكتابة (1).

2. الاعمال البحثية السابقة

ان اهمية هذا الموضوع قد استوجبت الكثير من البحث في هذا المجال. مرشح كابور قد اثبت فعاليته في تصنيف الحرف (2)(3)(4). وقد اعتمدت في هذه الطريقة حيث تعتمد هذه الطريقة على استخدام العديد من الميزات للتمييز على أساس الإسقاط الأفقي وتوزيع ارتفاع المكونات المتصلة لتصنيف النص العربي والروماني. لكن هذه الطريقة قد اثبتت كفاءتها على النصوص المطبوعة الأياً ولم تقيس جيداً على النص المكتوب باليد.

علي سلامات (5) اقترح تصنيف الاحرف العربية بالاعتماد على تكرار الاحرف، حيث اعتمدت خوارزمية البحث على اسلوب الرنين المتكيف في المنطق المضرب، التي تنتمي الى اسلوب تعليم الانظمة بالاعتماد على الشبكات العصبية.

بن سالم (6) عرض في بحثه طريقة لتصنيف الحروف العربية والرومانية في الوثائق متعددة اللغات. الطريقة تعتمد اولا على التحويل التشكلي لصورة خط النص. ثم تحليل خصائص النص المصور باستخراج 12 صفة للنص.

3. تعريف ميزات النصوص والعمليات الحسابية المطبقة عليها

يتم تصنيف النصوص الى اللغة التي تنتمي اليها عن بطريقتين اما عن طريق الصفات العامة أوالخاصة. حيث تستند الميزات العامة على قوام النص ويعتمد على مصفوفات لتحليل المويجات المستخرجة من النص الكتابي وتنفذ ايضا تحليل على مستوى الكلمة باستخدام مرشحات خاصة كمرشح كابور.

اما المجموعة الثانية من الميزات المحلية معضمها يعمل على سمات تستخلص على مستوى الكلمة في النص والمكونات المتصلة. الاساليب هنا كثير وتقسم الى اربع فئات:سمات هيكلية أو هندسية ، الخصائص التشكيلية ، الميزات الإحصائية و انتشار الميزات المكانية (7). على المستوى الميزات المحلية عند النظر الى صورة الوثائق تحول الصورة الى مجموعة من الكائنات المترابطة لكل كائن مجموعة من النقاط تكون هذه المجموعة من 8 نقاط مترابطة حيث تم تعريف مجموعة من الصفات المحلية التي استخدمت في هذا البحث.

- **متجه النقاط الوسطى (Centroid)** :متجه يحتوي على الاحداثي السيني و الصادي لمركز الكتلة في المنطقة.
- **المنطقة (Area)** : العدد الفعلي للنقاط في المنطقة.
- **الصندوق المحيط (Bounding box)** : هو أصغر مستطيل يمكن أن يحتوي على الكائن. و يمثله متجه يحتوي على إحداثيات الزاوية العلوية اليسرى من الصندوق ، و العرض والارتفاع للصندوق.
- **مساحة الصندوق المحيط (Area of Bounding box)** = العرض * الارتفاع من مربع الاحاطة.
- **المدى (Extent)**: نسبة النقاط في مربع الاحاطة والصندوق.
- **المحور الأكبر (Major Axis length)**: طول (ويقاس بعدد النقاط) المحور الأكبر للشكل البيضاوي للكائن
- **المحور الأصغر (Minor axis length)**: طول (ويقاس بعدد النقاط) المحور الأصغر للشكل البيضاوي للكائن
- **نسبة الارتفاع الى العرض (Aspect ratio)**: هي نسبة ارتفاع Bounding-box الى عرضه.
- **عدد اولير (Euler_number)** والفجوات (holes): عدد اولير هو عدد الكائنات في المنطقة مطروح منه عدد الفجوات ،ففي منطقة box_Bounding التي تحتوي على كائن واحد من النقاط المترابطة الفجوات تكون مساوية ل holes = 1- Euler_number
- **الانحراف (Eccentricity)**: انحراف الشكل البيضاوي للكائن ، الانحراف هو نسبة المسافة بين بؤر الشكل البيضاوي و طول المحور الرئيسي له . القيمة تكون بين 0 و 1 . (تكون القيمة 0 كلما كان الشكل البيضاوي اقرب الي الشكل الدائري، في حين تكون القيمة 1 عندما يكون الشكل البيضاوي قطعة مستقيمة) . يمكن حساب استدارة الكائن بالعلاقات الرياضية التالية (8):

$$eccentricity(ecc) = 2 \times \frac{\sqrt{\left(\frac{major}{2}\right)^2 - \left(\frac{minor}{2}\right)^2}}{major} \dots\dots\dots(1)$$

حيث ان:

$$major = 2 \times \sqrt{2 \times (uxx + uyy + common)} \dots\dots\dots(2)$$

$$minor = 2 \times \sqrt{2 \times (uxx + uyy - common)} \dots\dots\dots(3)$$

$$common = \sqrt{(uxx - uyy)^2 + 4 \times uyy^2} \dots\dots\dots(4)$$

$$uxx = \frac{1}{N} \sum_{i=1}^N x_i^2 + \frac{1}{12} \dots\dots\dots(5)$$

$$uyy = \frac{1}{N} \sum_{i=1}^N y_i^2 + \frac{1}{12} \dots\dots\dots(6)$$

علما ان قيمة (N) تمثل عدد النقاط في الكائن , و (x) (y) هي القيم المحسوبة ومخزونة في متجه النقاط الوسطى.

لو افترضنا مثلا ان صورة ثنائية (binary image) التالية (اي ان قيم النقط في الصورة اما 0 أو 1).

$$\text{binary image} = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

وكانت قيم:

$$[0.5000 \quad 0.5000 \quad 4.0000 \quad 4.0000] = (\text{Bounding_box})$$

$$4.6188 = \text{Major axis}$$

$$4.6188 = \text{Minor axis}$$

فإن قيمة eccentricity ستكون 0.

4. عملية التصنيف (classification) و خطأ التصنيف (classification error)

ان خطأ التصنيف يعتمد على تداخل (overlap) بين الاصناف. ان الابتعاد بين مراكز الاصناف مع تصغير التداخل بين الاصناف ينتج عنه امكانية تصنيف بجودة عالية. ان التصنيف الجيد يعتمد بالاساس على المتوسط الحسابي (mean) للصنف و كذلك على معامل اخر وهو معمال التشتت (variance) حيث انه كلما كان معامل التشتت قليل قلت نسبة الخطاء في التصنيف. وبالاعتماد على العلاقة السابقة فإن وجود معامل تشتت محدد فان ابعاد المتوسطات الحسابية للأصناف سوف يقلل من التداخل ويجعل التصنيف اكثر دقة. اما في الحالات التي يكون فيها معامل التشتت عالي وأن كان المتوسط للأصناف بعيد فإن التداخل موجود و يكون هنا التصنيف قليل الدقة. وبذلك يكون التصنيف الدقيق محسوب بالاعتماد على مقياس مركب يعتمد على معامل التشتت والمتوسط الحسابي في نفس الوقت.

ان المقياس الذي يحتوي على علاقة بين كل من (mean) و (variance) هو معامل (Mahalanobis distance classifier) ويحسب هذا المعامل بالاعتماد على العلاقة الرياضية التالية:

$$MD = \text{transpose} (x - u) * \Sigma^{-1} * (x - u) \dots\dots\dots(7)$$

حيث ان...

Σ هو المصفوفة التغاير (covariance matrix) لفئتين تصنيف.
 Σ^{-1} هو معكوس المصفوفة السابقة.

u هو متجه المتوسط الحسابي $u = [u_1 \ u_2 \ \dots \ u_d]'$, حيث ان d هو بعد المصفوفة.

استنادا إلى المعادلات الواردة أعلاه يمكن الاستفادة منها من ناحية زيادة المسافة بين المصنفات، وبذلك يمكن القول ان مصنف الأمثل يمكن أن يبنى بأن يكون قادرا على التمييز بين فئتين من الكائنات. أن الأمثل هنا لا يعني أن العملية ستكون خالية من الأخطاء، إلا أن المصنف قادر على تقليل الأخطاء الى الحد الأدنى (9).

5. تصنيف النصوص

الهدف الاساسي من هذه البحث هو تصنيف النصوص في صور الوثائق متعددة اللغات الى عربي وروماني (انكليزي). هناك خاصيتين اساسية يمكن تقود عملية التحليل لأجراء التصنيف و من بعده اختيار الصفات لغرض التقسيم اولا ان النصوص العربية تكتب من اليمين الى اليسار لذلك اختيار الصفات الخاصة بالتصنيف يجب ان تعكس هذه الحقيقة. ثانيا الكتابة باللغة العربية هي بشكل منحنيات وهذا يعني (Bounding_box) للمكونات المترابطة سيكون اكبر افقيا مقارنة بالكتابة الانكليزية التي ممكن ان تكون الكتابة بها بشكل منحنيات او من الممكن ان لا تكون ايضا.
 على سبيل المثال (centroid) للمكونات المترابطة في اللغة العربية يبدي ميل الى يمين مركز (Bounding_box). و (aspect ratio) للغة العربية يتكون اصغر من اللغة الانكليزية.
 هذا بالإضافة الى ان الكائنات في اللغة العربية ستميل الى الاستطالة اكثر من اللغة الانكليزية. أن قيمة (Eccentricity) تبين مقدار استدارة الكائنات. عموماً فإن اللغة الانكليزية تميل للاستدارة اكثر من اللغة العربية بسبب الطبيعة الانحنائية لأشكال اللغة العربية. وهذا السبب ايضا يجعل حجم الكائن (Extent) في اللغة العربية اصغر من الكائنات في اللغة الانكليزية.

أن المقارنات والتحليلات السابقة قد اعتمدت في البحث حيث تم تطبيق تحليلات احصائية على عينات عشوائية من الوثائق من كلا اللغتين , وكما مبين في الجدول (1).

جدول رقم (1) المتوسط الحسابي لبعض الميزات

| | Distance | Extent | Holes | eccentricity | Aspect |
|----------------|----------|--------|-------|--------------|--------|
| Arabic Script | .3 | .6 | .2 | .9 | 1.2 |
| English Script | .2 | .4 | .7 | .7 | 1.5 |

بالاعتماد على التنتائج السابقة، فإن مجموعة من الصفات قد حددت لغرض تصنيف اللغة داخل الوثائق وهذه الصفات هي:

- .centroid distance
- .angle of centroid distance vector
- .number of holes
- .eccentricity
- .aspect ration
- .extent

كل الصفات السابقة هي نسب لذلك لا حاجة لتقييسهم بحدود قيم معينة، فقط (centroid distance) يجب ان تقيس بقيم بين [0,1] . لكل وثيقة يرغب في تصنيفها يجب توليد ملخص احصائي لكل صفة من الصفات السابقة. القيم الاحصائية هي المعدل، الانحراف المعياري و الالتواء. حيث يولد متجه من 18 عنصر لغرض تصنيف كل وثيقة للاستخدام في عملية التصنيف. كما واستخدم (Mahalanobis distance classifier) كعلاقة بين (mean) و (variance).

6. تحليل النظام ونتائج اختبار نظام التصنيف

لاغراض الاختبار تم الاعتماد على من الوثائق التي تم اختيارها عشوائيا. فصلت هذه الوثائق الى قسمين القسم الاول لتدريب النظام والقسم الثاني لاغراض فحص النظام شكل (1) يبين مجموعة من الوثائق المستخدمة لتدريب النظام. الوثائق التي استخدمت للتدريب تتكون من:

- 16 وثيقة مطبوعة باللغة الانكليزية
- 16 وثيقة مكتوبة باليد انكليزية
- 16 وثيقة مطبوعة باللغة العربية
- 16 وثيقة مكتوبة باليد باللغة العربية

References in this publication to IBM products, programs or services do not imply that IBM intends to make these available in all countries in which IBM operates. Any reference to an IBM licensed program in this publication is not intended to state or imply that only IBM's licensed

Without the World Wide Web, there would probably be very few intranets. There are many forces driving Corporation to set up an intranet, but the main one is the dominating presence of the World Wide Web.

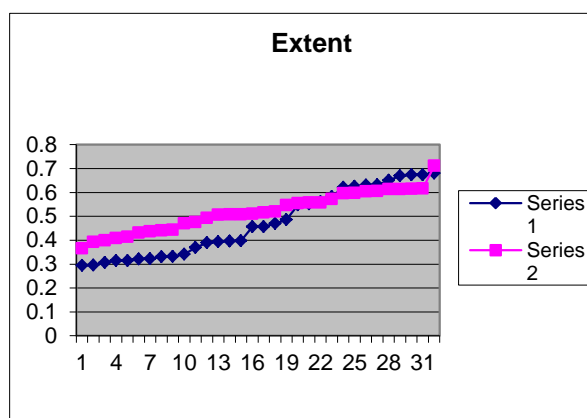
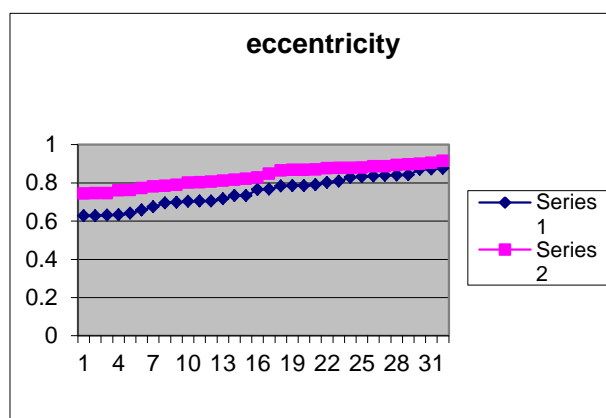
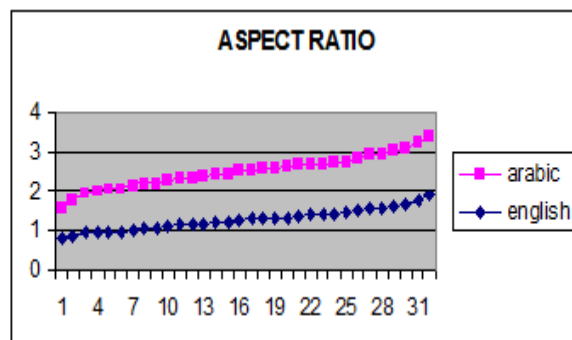
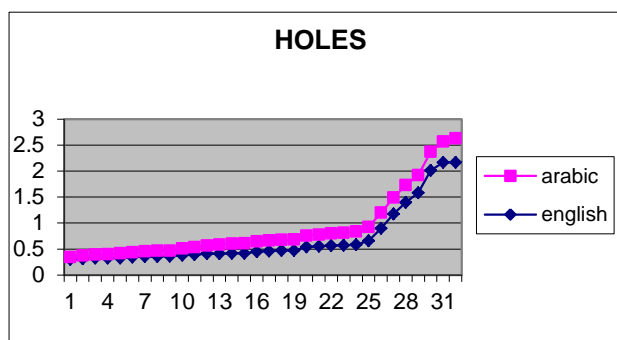
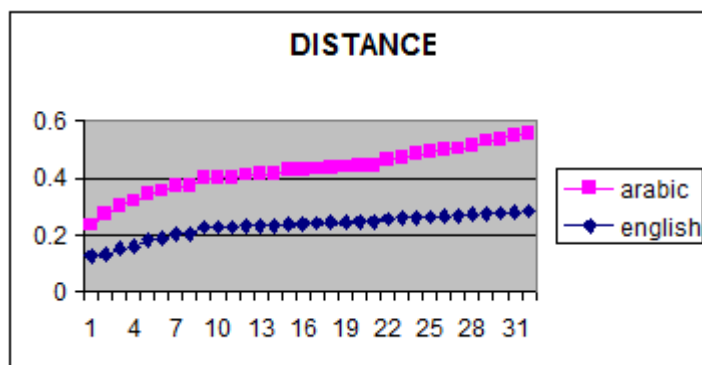
- استعمار جزر و المبالغ المخصصة من ميزانية الدراسات العليا لآراء الحاسبات والاهيئه الحكومه رسمياً من اجل تطوير استخدام الحاسب في مؤسسه معيئه بتزود الحاسبات الذي ينبغي كل رفاق فيه يصور نموذج من الانظمة ... لذا نفتح جلاز هاسوب لاستخدام قسم الحاسبات ووضع نظام مخزني سحابي يسهل عملية البحث في توفر المواد المطبوعه والمطبوعه والناتجه .

فقد ابتكر الانسان في العراق لأول مرة في التاريخ الحروف الابجدية والكتابة والارقام الحسابية وطرق القياس وجرى تدوين المعاملات والمعلومات والاحداث .. ووضع القواعد والاسس لبناء الحياة الجماعية المستقرة والمدن والقلاع الحصينة .. ووضع الانظمة والتشريعات التي تنظم حركة المجتمع .. وجرى تحديد مبدا الحقوق والواجبات ووضع اول نظام متكامل لتشكيل الجيوش النظامية الموحدة للدولة .. وكذلك وضع اول نظام

شكل (1) مجموعة نماذج من الوثائق المستخدمة لتدريب النظام

الوثائق المكتوبة باليد اخذت من اشخاص مختلفين بخلفيات مختلفة حتى لا يؤثر هذا على علمية التصنيف. متجه من 64 قيمة قد ولد لاغراض التصنيف وقد وجد ان هذا المتجه هو الحد الادنى المناسب طولاً للابتعاد عن نسبة الخطاء بالنظام.

مخطط انحناء لصفات متجه التدريب مبينة بالشكل (2). الامر الاساسي الذي يجب ملاحظته بالمخطط هو ليس القيم او شكل منحنى المخطط وانما الانفصال بين المنحنيين لكل لغة. لكل صفة بالمنحنى نجد انفصالها عن الصفة المقابلة بالمنحنى للغة الثانية.



شكل (2) منحنيات خواص النصوص.

أن بعض الصفات تبيّن انفصال أكبر من الصفات الأخرى. الجدول (2) يبيّن نتائج فحص النظام.

جدول رقم (2) نتائج تصنيف الوثائق

| نوع الوثيقة | عدد الوثائق | عدد الكتل (blocks) | النسبة المئوية (%) لمقدار صحة التصنيف |
|---|-------------|--------------------|--|
| انكليزي-عربي طباعة | 40 | 6210 | 100 |
| انكليزي-عربي خط يد | 40 | 7349 | 100 |
| انكليزي عربي خط يد و طباعة | 80 | 13559 | 100 |
| انكليزي-عربي خط يد و طباعة تدوير 90 درجة | 5 | 1053 | 100 |
| انكليزي-عربي خط يد و طباعة تدوير 180 درجة | 10 | 1752 | 100 |

شكل (3) تبيّن نماذج من الوثائق التي اخلت الى النظام وتم تصنيفها.

قاعدة البيانات بالإنجليزية (Database): هي مجموعة من عناصر البيانات المنطقية المرتبطة مع بعضها البعض بعلاقة رياضية، وتتكون قاعدة البيانات من جدول واحد أو أكثر. ويتكون الجدول من سجل (Record) أو أكثر ويتكون السجل من حقل (Field) أو أكثر.

مثال عليه السجل الخاص بموظف معين يتكون من عدة حقول مثل رقم الموظف - اسم الموظف - درجة الموظف - تاريخ التعيين - الراتب - والقسم التابع له، وغير ذلك من بيانات الموظف تخزن في جهاز الحاسوب على نحو منظم، حيث يقوم برنامج يسمى محرك قاعدة البيانات (Database Engine) بتسهيل التعامل معها والبحث ضمن هذه البيانات، وتمكين المستخدم من الإضافة والتعديل عليها.

توجد طرق متعددة لتنظيم الـ Hard Disk منها نظام FAT 32 الذي ظهر مع نظام تشغيل Windows 98 والذي يمكن استخدامه مع نظام Windows XP وهناك طريقة أخرى لتنظيم الـ Hard Disk لتنظيم البيانات هي NTFS بأنه يتبع تأمين أكبر للبيانات المخزنة سواء بأصنافها كلمة سر (Password) أو صلاحيات استخدام لزيادة مستوى الامان للملفات.

شكل (3) نماذج من الوثائق التي استخدمت للتصنيف بواسطة النظام

الخوارزمية في حالة عامة

- المدخلات: صورة ثنائية (binary image) للوثيقة المراد تصنيفها.
- المخرجات: تصنيف مجموعة احرف الموجودة في الوثيقة المدخلة الى نوع اللغة التي تنتمي اليها.

- في المرحلة الأولى: تحديد العناصر المترابطة في الصورة الثنائية (binary image) للوثيقة.
- في المرحلة الثانية: تنفيذ بالتتابع ما يلي لجميع العناصر المترابطة:
 - المرحلة الثانية -1: حساب الميزات (Centroid, Bounding_box, Eccentricity, Euler_number, Extent) لكل كائن مترابط.
 - المرحلة الثانية -2: استخلاص القيم التالية لكل كائن:

holes=1-Euler number

h=height of Bounding_box

W=width of Bounding_box

Aspect ratio=h/w

cx=x coordinate of centroid

bx=x coordinate of the center of the Bounding_box

by=y coordinate of the center of the Bounding_box

disx=x distance from the center of box to the centroid

disy=y distance from the center of box to the centroid

$$vdistx = bx + \frac{w}{2} - cx$$

$$vdisty = by + \frac{h}{2} - cy$$

$$dto = \sqrt{(vdistx^2 + vdisty^2)}$$

$$dtot(normal) = dtot - \min(dtot) / (\max(dtot) - \min(dtot))$$

theta

= angle distance (dtot) vector makes with horizontal direction

$$= \arctangent(vdisty/vdistx)$$

- في المرحلة الثالثة: احتساب المتوسط الحسابي , معامل التشتت و الألتواء للقيم المحسوبة في المرحلة الثانية.
- في المرحلة الرابعة: ينظم متجه الميزات (feature vector).
- في المرحلة الخامسة: يرسل متجه الميزات الى معادلة المصنف (classifier).
- في المرحلة السادسة: بالاعتماد على نتائج المصنف (classifier) نحصل على التصنيف النهائي لاحرف الوثيقة المدخلة للنظام.

8. الاستنتاجات

هذا البحث يعرض طريقة لتصنيف النصوص الى (العربية و الرومانية) دون أي محاولة لاكتشاف الحروف او تشكيل الكلمات. الفكرة الأساسية هي لاستخراج معلومات شكل الكلمة على مستوى الصفحة ، للحد من حساسية للضوضاء على المستندات الممسوحة ضوئياً . حيث يستند أساسا على معالجة تحليل المكونة الموصولة ، و المكونات المرتبطة تشكيمياً. وقد تم استخدام هذه الميزات الأساسية و استغلالها في مجموعة متنوعة من الطرق. هذه الطريقة سريعة وغير مكلفة حسابيا ، ليست حساسة لاختلاف الخطوط أو إلى الطابع الخط . تصنف هذه خوارزمية المتقدمة النص من ناحية اللغة وكذلك طريقة كتابة النص ان كان يدوية او طباعة.

تم التصنيف النصي بواسطة خوارزمية تصنيف الى نوعين من النصوص ، وهذه النصوص هي الرومانية والعربية. تم الحصول على دقة التصنيف من 100 % . بعد هذه النتائج، تم التأكد بأن المصنف يمكن استخدامه لتصنيف الوثائق متعددة اللغة. تم استخراج ست ناقلات الميزات من كتلة النص و يتم تغذية هذه إلى مسافة المصنف للتمييز بين النصوص.

المراجع

- [1] Dr. Firoj Parwej, “ *The State of the Art Recognize in Arabic Script through Combination of Online and Offline*”, International Journal of Computer Science and Telecommunications ,Volume 4, Issue 3, March 2013.
- [2] R.M. Haralick, “*Statistical and structural approach to texture*”. Proceeding of IEEE, Vol 67, 1997 pp. 786-804.
- [3] J. Ding, L. Lam and C.Y. Suen. “*Classification of Oriental and European scripts by using characteristic features*”. Fourth International Conference on Document Analysis and Recognition, ICDAR97, Ulm, Germany, 1997, Vol.2, pp. 1023-1027.
- [4] B Waked, S. Bergler, C.Y. Suen and S. Khoury. Skew Detection, “*Page Segmentation, and Script Classification of Printed Document Images*”. In Proceeding of the 1998 IEEE International Conference on System, Man, and Cybernetics, San Diego, CA, October 1998.
- [5] Ali Selamat and Ng Choon Ching, “*Arabic Script Documents Language Identifications Using Fuzzy ART*”, Second Asia International Conference on Modelling & Simulation, 2008 IEEE.
- [6] S. Ben Moussa, A. Zahour, A. Benabdelhafid, A.M. Alimi, “*Fractal-Based System for Arabic/Latin, Printed/Handwritten Script Identification*”, IEEE Pattern Recognition, 2008. ICPR 2008. 19th International Conference on 8-11 Dec. 2008.

- [7] Alireza Behrad, Malike Khoddami and Mehdi Salehpour, “*A Novel Framework for Farsi and Latin Script Identification and Farsi Handwritten Digit Recognition*”, JOURNAL OF AUTOMATIC CONTROL, UNIVERSITY OF BELGRADE, VOL. 20:17-25, 2010.
- [8] Haralick and Shapiro,”*Computer and Robot Vision*”, vol I, Addison-Wesley 1992, Appendix A.
- [9] C. Jose & R. Neil & W. Curt (2000) “*Numerical and Adaptive Systems: Fundamentals Through Simulations*”. John Wiley & Sons, Inc. 2000.