

Abstract

In the Arabic studies there is a theory claim that ancient authors of Arabic books more eloquent than modern authors and that lead to the hardness of reading ancient books , this project tries to proof that theory by analyze both ancient and modern books using stylometric features which is a set of methods used to analyze the texts and get meta-data about it , there is a lot of features some of them contain significant changes with 0.26 p-value and other features didn't change through the decades of Arabic writings ,also some features kept increasing until the last decade it decreased for example the letter 'هـ' (h) at the end of the words 'هـ' pronoun in Arabic- like in 'له' (lh) which mean "for him". Those features must be apply on a very large texts (corpus) to test the significant but the Arabic corpora are very rare beside the published ones doesn't contain any old texts all of them is very modern and from blogs or the websites. So there was a necessary to collect and clean new Arabic corpus with texts from the year of 100 Hijri (718) to 1439 Hijri (2017) and publish it on the Internet. The stylometric analysis which applied on the new corpus was fed to Naive Bayes classifier and when try to test a new document the classifier will predict the writing year for the document with a small error (example: the document may be written from 500 - 600 Hijri) and the result was very good with Precision of 83.7% .

ملخص

هناك فرضية تدعي أن الكتاب العرب القدماء أكثر فصاحة من الحديثين, وهذا أدى إلى الصعوبة التي نواجهها في قراءة النصوص القديمة. هذا المشروع يحاول إثبات هذه الفرضية من خلال تحليل النصوص للحصول على معلومات تصفها, هناك الكثير من الخصائص التي يمكن فحصها لعمل مقارنة بين هذه المعلومات, منها ما أظهر النتائج أنه يحتوي على فروق واضحة ومنها ما هو بقي متشابهاً مع مرور الزمن على الكتابة العربية. وهناك ما زاد استخدامه مع الزمن ومن ثم بدأ بالتراجع مثل استخدام الهاء "هـ" في نهاية الكلمة, مثل "له". هذه الخصائص تم تطبيقها على عدد كبير من النصوص "مجمّع للنصوص" لفحص مدى قدرتها على تزويدنا بالفروق التي تمكّننا من إصدار حكم بقدّم أو حداثة النص. ولإنشاء هذه المجمّع قمنا بتجميع النصوص القديمة من 100 وحتى 1439 هجري. وبعد تقسيم هذه النصوص إلى فترات وإدخالها على خوارزمية "تعلم الآلة" (Naive Bayes classifier). وعند إدخال أي نص جديد على هذه الخوارزمية سوف تقوم بتصنيفها بناء على ما تعلمته من المعلومات الواصفة للنصوص السابقة بدقة عالية.