

بسم الله الرحمن الرحيم

"قُلْ لَوْ كَانَ الْبَحْرُ مَدَادًا لِكَلْمَاتِ رَبِّي لَنَفَدَ الْبَحْرُ قَبْلَ أَنْ تَفْدَ كَلْمَاتُ رَبِّي وَلَوْ جَئْنَا بِمِثْلِهِ مَدَادًا"

صدق الله العظيم

## ملخص

بدأت فكرة استرجاع المعلومات (Information Retrieval) خلال التحول الذي شهدناه في طرق تخزين المعلومات الثقافية والاجتماعية والعلمية من الأوراق والكتب إلى مكتبات رقمية منتشرة على الشبكات العالمية، وتأتي هذه المعلومات على شكل نصوص أو وثائق مصورة (كالصور والخرائط والرسومات التقنية وغيرها) أو صوت أو صور متحركة (الفيديو). في محاولة لجعل استرجاع هذه الكمية الهائلة من المعلومات بفاعلية كان لا بد من طريقة لاستبطاط الكلمات المفتاحية آلياً (Automatic Term Extraction). هذه الطريقة الآلية تشكل حجر الأساس للعديد من التطبيقات ومنها محركات البحث، لأنها عملية شاقة ومكلفة أن يتم استبطاط الكلمات المفتاحية يدوياً. وللتتأكد من دقة المعلومات المسترجعة فإن الكلمات المفتاحية المستتبطة يجب أن تصف محتوى الوثائق المستتبطة منها وصفاً دقيقاً.

تقدم الباحثون في هذا الموضوع بالعديد من المقترنات والطرق المختلفة لاستبطاط الكلمات المفتاحية، بعضها يعتمد على الإحصاء والاحتمالات وبعضها الآخر يعتمد على التحليل اللغوي وغيرها، تأتي هذه الأطروحة كمساهمة في كشف أوجه الاختلاف والتشابه وكذلك محاسن ومساوئ الطرق الإحصائية المستخدمة وهي: تكرار الكلمات في الوثيقة الواحدة (Term Frequency TF)، تكرار Inverse Document Frequency (IDF)، الدمج بين تكرار الكلمات وتكرار الوثائق (TFxIDF)، والطريقة الرابعة هي نموذج قيم التمييز للكلمات (Term Discrimination Value Model).

الإحصائية وتطوير أداة محوسبة لاستباط الكلمات المفتاحية آلياً (ATEWB) تستخدم للمقارنة من خلال ثلاثة تجارب تقييمية تبرز العوامل المؤثرة على دقة النتائج وذلك لتحديد الظروف المناسبة لاستخدام تلك الطرق.

من ناحية أخرى، تهدف الدراسة إلى زيادة فاعلية الطرق الإحصائية وتسريعها من خلال استخدام محركات قواعد البيانات (Database Engines) والتي تختصر العمليات الحسابية في خوارزميات هذه الطرق بشكل ملحوظ، كما استخدمت في تسريع استرجاع الوثائق عن طريق تخزين نسخة منها في قاعدة البيانات.

الوثائق المستخدمة في التجربة الأولى هي مجموعة من الملخصات (Abstracts) لأبحاث في مجال الاستباط الآلي للكلمات المفتاحية، والكلمات المفتاحية لها مستبطة يدوياً لمقارنتها بالكلمات المفتاحية المستبطة آلياً، أما تلك المستخدمة في التجربة الثانية والثالثة فهي وثائق أعدت خصيصاً لغرض البحث، وهي متوفرة على بعض موقع الإنترن特 المهمة بموضوع الاستباط الآلي للكلمات المفتاحية.