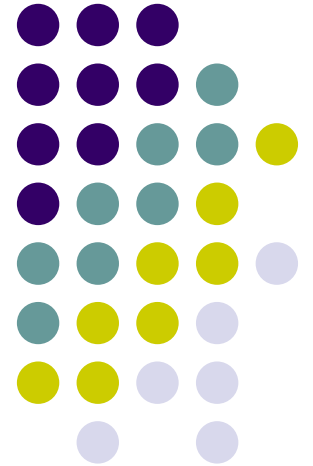


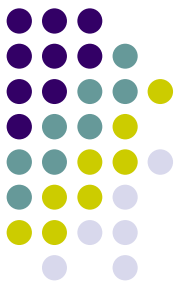
استخدام ذخائر النصوص لاستخلاص المصطلحات المتخصصة

عبدالمحسن بن عبيد الثبيتي

مركز المعلومات

الرياض – المملكة العربية السعودية

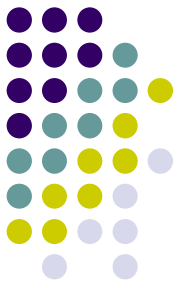




ما لذي تقدمه هذه الورقه

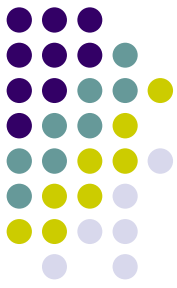
- طريقة إحصائية بسيطة لاستخلاص المصطلحات المتخصصة ذات المفردة الواحدة اعتمادا على ذخائر النصوص.

- تقيس الطريقة مدى ابتعاد استخدام المفردة في مدونة المجالات المتخصصة عنه في مدونة المجالات العامة أو اللغة اليومية.



مالفائدة؟

- يمثل استخلاص المصطلحات المتخصصة آليا الخطوة الأولى نحو التعامل مع النتاج المعرفي المتخصص كقاعدة وأساس لعمليات أخرى مهمة مثل
- بناء المعاجم والقواميس المتخصصة.
- المساعدة في عمليات الترجمة الآلية للنصوص العلمية والفنية.
- استرجاع وفهرسة النصوص العلمية والفنية.
- اختيار الخصائص اللغوية الملائمة لعمليات التصنيف الآلي للنصوص.
- بناء شبكة العلاقات بين المفاهيم في المجالات المتخصصة.



الأساليب المتبعة لاستخلاص المصطلحات

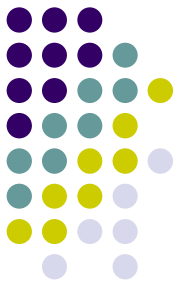
- الاكتشاف الآلي للتركيب النحوي للمصطلحات.
تعاني أدوات معالجة اللغة العربية آليا من شح وندرة في مثل هذا الدراسات وتوافر هذه الأدوات.
- الطرق الإحصائية لاستخلاص المصطلحات في النصوص المتخصصة
- في كلتا الطريقتين يبرز استخدام ذخائر النصوص (المدونات) كطريقة تثبت صلاحية الطرق المستخدمة لاستخلاص المصطلح العلمي آليا كون المدونات تعطي الدليل على التراكم النحوية الأكثر استخداما التي يظهر فيها المصطلح أو حتى التوزيع الإحصائي لمفرداته.

الفرضية



- تفترض هذه الدراسة اختلاف التوزيع الإحصائي للمفردات العربية في المجالات المتخصصة عنه في المجالات العامة أو اللغة العربية اليومية وذلك اعتماداً على الدراسات السابقة الخاصة بمدونات المجالات المتخصصة في اللغة الانجليزية

توصيف المدونات (مدونة المجالات العامة)



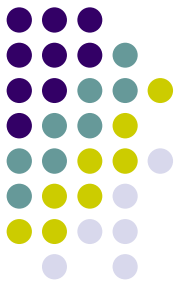
- تمثل اللغة اليومية وهي مكونة من حوالي ١,١ مليون كلمة مختلفة بهيكله مشابهة لبنية مدونة (Collins-Birmingham) للنصوص الانجليزية.
- تم جمع نصوص المدونة من شبكة الانترنت من مواقع سعوديه ولكتاب سعوديين في الأغلب موزعة كالتالي:
 - الصحف اليومية (٥٧%).
 - المجلات الأسبوعية (٢١%).
 - الكتب (١٥%).
 - المطويات التعريفية (٦%).
 - المواقع الشخصية (١%).

توصيف المدونات (مدونة المجالات المتخصصة)



- تمثل المجالات المتخصصة يقارب حجمها ٥٨ ألف كلمة تتكون مدونة المجالات المتخصصة من مدونتين فرعيتين.
- مدونة الذكاء الصناعي (٤٢%) وقد غطت مواضيع متعددة وخصوصا موضوعي الشبكات العصبية والأنظمة الخبيرة.
- مدونة الفيزياء (٥٨%). غطت مواضيع فيزياء الكم والفيزياء الذرية والجزيئية.

الكلمات المائة الأكثر تكرارا في مدونة اللغة اليومية.



| النسبة | المفردة | |
|--------|--|----|
| ١٩,١ | و، في، من، على، أن، إلى، التي، لا، ما، عن | ١ |
| ٣,٤ | هذا، <u>الله</u> ، الذي، هذه، أو، مع، ان، ذلك، هو، كان | ٢ |
| ١,٩ | قد، لم، كل، بين، بن، بعد، هي، إن، كما، الى | ٣ |
| ١,٣ | خلال، <u>العربية</u> ، حتى، <u>عبد</u> ، له، <u>قال</u> ، حيث، <u>المملكة</u> ، كانت، غير | ٤ |
| ١,١ | لكن، بعض، <u>محمد</u> ، عليه، <u>السعودية</u> ، إلا، قبل، به، بها، فيها | ٥ |
| ٠,٩ | يكون، ثم، هناك، لها، أي، تلك، عام، <u>العمل</u> ، فيه، إذا | ٦ |
| ٠,٨ | <u>مجلس</u> ، أنه، العام، ليس، <u>الدكتور</u> ، منها، <u>العالم</u> ، رئيس، بل، الذين | ٧ |
| ٠,٧ | عند، مثل، أكثر، يمكن، اليوم، <u>الرياض</u> ، تكون، فقد، <u>الملك</u> ، <u>الناس</u> | ٨ |
| ٠,٦ | <u>العزیز</u> ، عدد، أما، عدم، <u>السعودي</u> ، <u>التعليم</u> ، <u>عبدالله</u> ، عليها، تم، بما | ٩ |
| ٠,٥ | <u>عبدالعزیز</u> ، أخرى، جميع، عندما، فإن، دون، <u>المجتمع</u> ، بأن، مما، يجب | ١٠ |
| ٣٠,٣ | المجموع | |

الكلمات المائة الأكثر تكرارا في مدونة الذكاء الصناعي

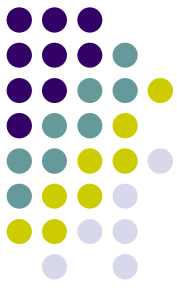


| النسبة | المفردة | |
|--------|---|----|
| ١٩,٥ | و، في، من، على، التي، أن، إلى، أو، هذه، <u>الذكاء</u> | ١ |
| ٤,٨ | عن، هي، هذا، هو، ما، لا، مع، <u>الشبكة</u> ، كل، الذي | ٢ |
| ٢,٨ | <u>الصناعي</u> ، ذلك، <u>الإنسان</u> ، حيث، <u>النظام</u> ، قد، <u>الاصطناعي</u> ، <u>الخبير</u> ، <u>المعرفة</u> ، بين | ٣ |
| ٢,٢ | أي، يمكن، <u>أنظمة</u> ، لكن، <u>عملية</u> ، كما، يتم، <u>الشبكات</u> ، بها، كان | ٤ |
| ١,٨ | بعد، غير، تكون، <u>نظام</u> ، <u>التعرف</u> ، ثم، <u>المعلومات</u> ، يكون، <u>الأنظمة</u> ، مثل | ٥ |
| ١,٥ | <u>طريق</u> ، <u>البحث</u> ، <u>طبقة</u> ، <u>مجموعة</u> ، <u>الخبيرة</u> ، <u>علم</u> ، <u>قاعدة</u> ، إذا، <u>تقوم</u> ، لها | ٦ |
| ١,٤ | <u>الأوزان</u> ، <u>الحاسوب</u> ، فيها، <u>الخبراء</u> ، بعض، <u>معالجة</u> ، <u>الطريقة</u> ، <u>العصبية</u> ، <u>مجال</u> ، <u>يقوم</u> | ٧ |
| ١,٢ | فإن، يجب، أنه، حتى، <u>الصوتية</u> ، <u>العصبونية</u> ، <u>الكلام</u> ، <u>المشاكل</u> ، ذات، لم | ٨ |
| ١,٢ | <u>التعليم</u> ، <u>الحاسب</u> ، <u>العمليات</u> ، <u>المشكلة</u> ، منها، أكثر، إن، بواسطة، <u>حل</u> ، <u>مرحلة</u> | ٩ |
| ١,١ | هناك، <u>المعالجة</u> ، كانت، إلا، <u>الآلية</u> ، <u>البشري</u> ، <u>الكمبيوتر</u> ، <u>المجال</u> ، <u>المنطق</u> ، <u>النظم</u> | ١٠ |
| ٣٧,٥ | المجموع | |

الكلمات المائة الأكثر تكرارا في مدونة الفيزياء.



| النسبة | المفردة | |
|--------|--|----|
| ١٦,٣ | و، في، من، على، أن، هذه، إلى، هذا، التي، عن | ١ |
| ٤,٥ | ان، لا، <u>النظرية</u> ، يمكن، مع، هو، <u>الضوء</u> ، هي، <u>النسبية</u> ، <u>نظرية</u> | ٢ |
| ٣,٠ | ما، الذي، <u>الكم</u> ، بين، <u>سرعة</u> ، كل، <u>بالنسبة</u> ، أو، <u>الزمن</u> ، حيث | ٣ |
| ٢,٤ | ذلك، كان، يكون، <u>الطاقة</u> ، لكن، <u>الكون</u> ، كما، بعد، أي، <u>الكمومية</u> | ٤ |
| ١,٩ | الى، فإن، <u>عام</u> ، <u>النواة</u> ، حول، تكون، <u>الإلكترون</u> ، <u>طاقة</u> ، غير، <u>الجسيمات</u> | ٥ |
| ١,٦ | لم، عند، <u>الجسيم</u> ، <u>جسيمات</u> ، <u>حالة</u> ، كانت، <u>جسيم</u> ، خلال، سوف، قد | ٦ |
| ١,٤ | <u>الأرض</u> ، <u>المراقب</u> ، <u>الذرة</u> ، <u>ذرة</u> ، <u>بسرعة</u> ، مثل، <u>العالم</u> ، <u>الفيزياء</u> ، عندما، <u>الهيدروجين</u> | ٧ |
| ١,٢ | <u>نموذج</u> ، <u>العلماء</u> ، <u>القوة</u> ، <u>الموجة</u> ، <u>الفأ</u> ، <u>الخاصة</u> ، <u>العامة</u> ، هنا، أنه، <u>تفسير</u> | ٨ |
| ١,١ | <u>الموجية</u> ، <u>حركة</u> ، <u>مبدأ</u> ، <u>معادلة</u> ، نفس، <u>الحالة</u> ، حتى، هناك، <u>المادة</u> ، <u>بشكل</u> | ٩ |
| ١,٠ | ذات، <u>ميكانيكا</u> ، احد، <u>الثانية</u> ، <u>الجسم</u> ، <u>الطبيعة</u> ، <u>كتلة</u> ، <u>الأبعاد</u> ، <u>الذرية</u> ، <u>الفضاء</u> | ١٠ |
| ٣٤,٤ | المجموع | |



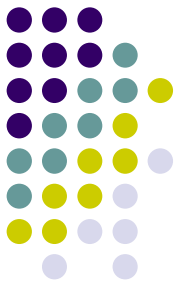
معامل الغرابة

- يقيس مدى ابتعاد معدل استخدام المفردة في المجال المتخصص عن استخدامها في المجالات العامة أو اللغة اليومية

معدل استخدام المفردة في المجال المتخصص

= معامل الغرابة

معدل استخدام المفردة في المجالات العامة



استخلاص المصطلحات

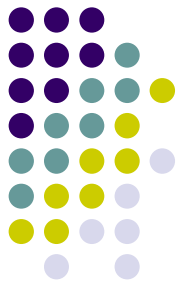
- يخضع للاعتبارات التالية:
- أن الكلمات التي يطبق عليها معامل الغرابة تنتمي للمجموعة المفتوحة من الكلمات.
- أن تكرار الكلمات مهم جدا في اعتبار الكلمة مصطلحا مقبولا أو متعارفا عليه في المجال المتخصص فهناك علاقة مباشرة بين التكرار وقبول الكلمة في المجال (Quirk, 1985).
- أن تكون قيمة معامل الغرابة عالية للكلمات التي تراعي الفقرتين السابقتين أعلاه.

الكلمات العشرين ذات معامل الغرابة الأعلى في مدونة الذكاء الصناعي.



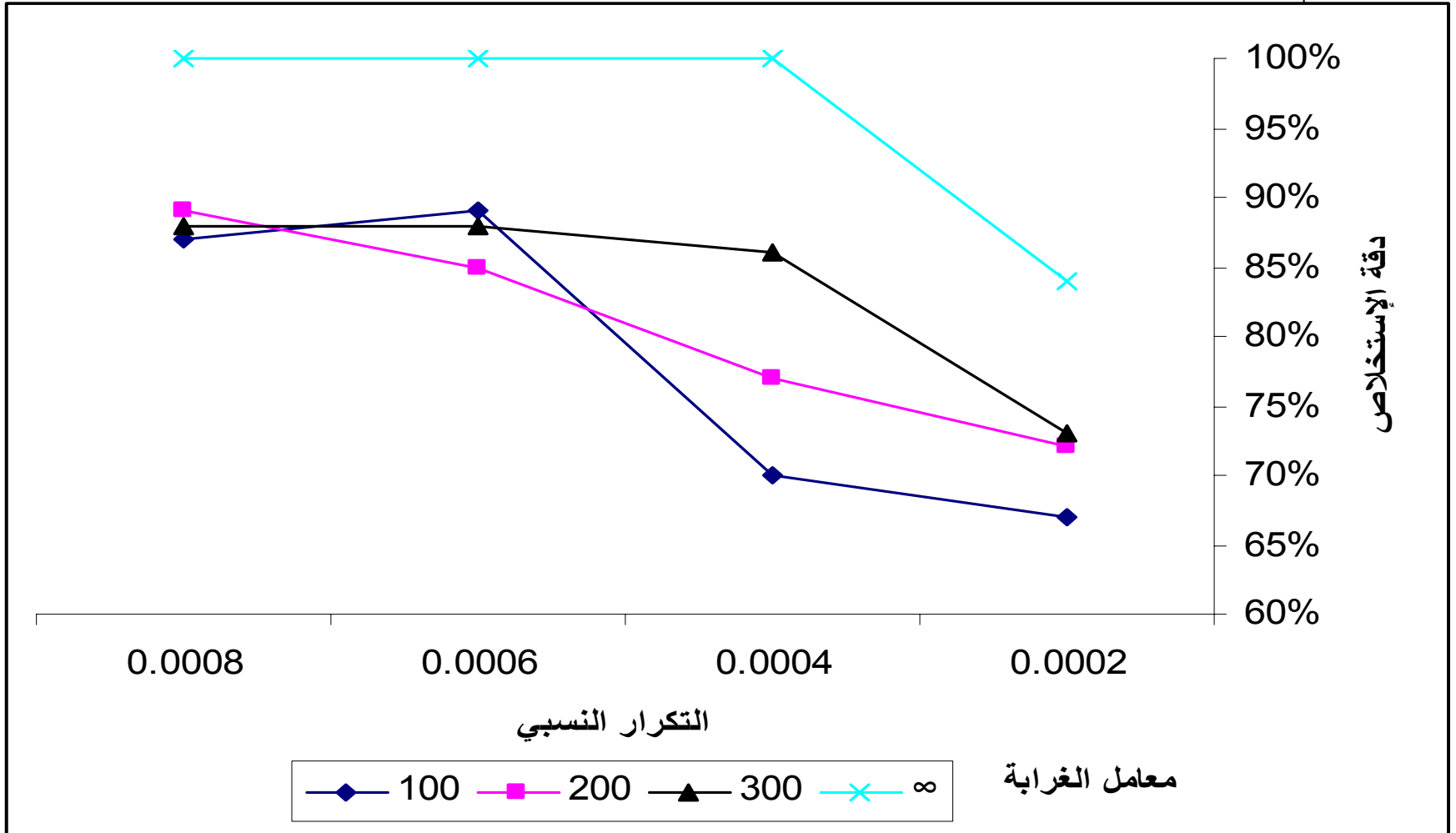
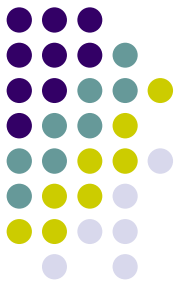
| معامل الغرابة | الكلمة | معامل الغرابة | الكلمة |
|---------------|---------|---------------|-------------------|
| ١٠٣٨ | درجة | ∞ | العصبونية/العصبية |
| ٥٠٧ | الاوزان | ∞ | الاصطناعي |
| ٤١٥ | محاكاة | ∞ | الضبابي |
| ٣١١ | الأذرع | ∞ | المدخلات |
| ٢٩٤ | البرمجة | ∞ | النظم |
| ٢٧٧ | مسموعة | ∞ | استنتاجات |
| ٢٧٧ | شبكة | ∞ | الانحيازات |
| ٢١٦ | صورة | ∞ | الخوارزميات |
| ٢٠٦ | طبقة | ∞ | مصفوفة |
| ٢٠٦ | الأنماط | ١٦٢٦ | الخبيرة |

الكلمات العشرين ذات معامل الغرابة الأعلى في مدونة الفيزياء.

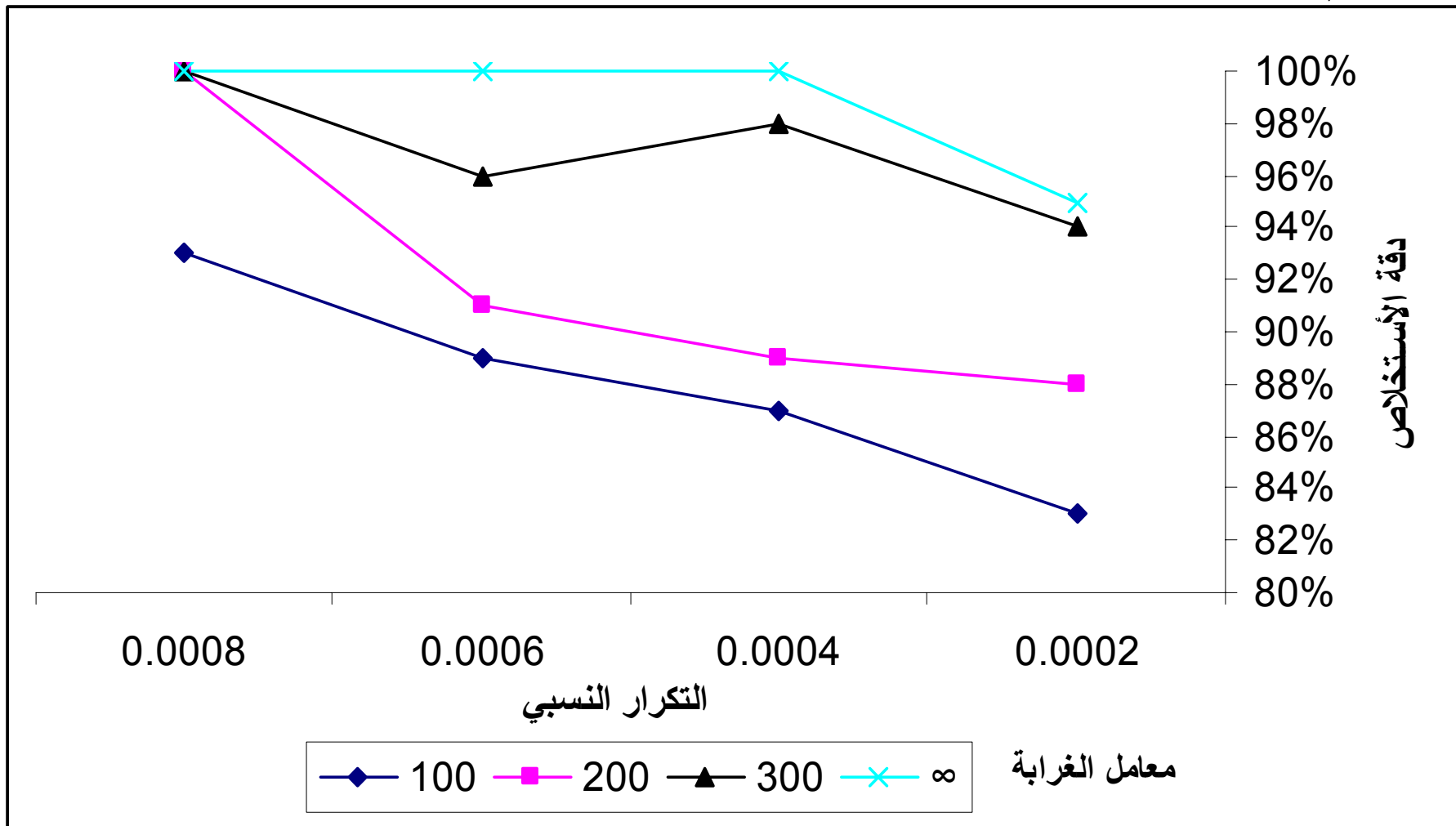


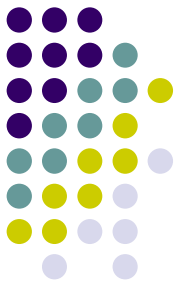
| معامل الغرابة | الكلمة | معامل الغرابة | الكلمة |
|---------------|-----------|---------------|-----------|
| ∞ | بوهر | ∞ | الكمومية |
| ∞ | ديراك | ∞ | الإلكترون |
| ∞ | اينشتاين | ∞ | جسيمات |
| ∞ | الزاوي | ∞ | الموجية |
| ∞ | الكواركات | ∞ | ميكانيكا |
| ∞ | شرودينغر | ∞ | الفوتونات |
| ∞ | المدار | ∞ | الذرات |
| ∞ | الارتياح | ∞ | الأوتار |
| ∞ | الجزيئات | ∞ | الزمكان |
| ∞ | الحيود | ∞ | بلانك |

دقة استخلاص المصطلحات لقيم مختلفة من التكرار النسبي ومعامل الغرابة لمدونة الذكاء الصناعي.



دقة استخلاص المصطلحات لقيم مختلفة من التكرار النسبي ومعامل الغرابة لمدونة الفيزياء.





الخلاصة

- إن دقة الاستخلاص تتناسب طرديا مع معامل الغرابة ومع التكرار النسبي ولكن ازدياد قيم معامل الغرابة والتكرار النسبي يؤدي إلى قلة الكلمات التي ممن الممكن تحديدها كمصطلحات.
- إن حجم المدونة يؤثر على دقة الاستخلاص فكلما ازداد الحجم ازدادت الدقة.
- إن حجم مدونة المجالات العامة ومدى عموميته لهذه المجالات يؤثر على قيمة معامل الغرابة.