

نمو التكامل المعرفي من واقع توظيف الأنطولوجيات في أطار التنقيب عن البيانات: دراسة تحليلية.

Towards to Knowledge Integration by Ontology-based Web Mining:

An Analytical Study

إعداد

مؤمن النشرتي

مدرس مساعد بقسم المكتبات والوثائق والمعلومات

كلية الآداب – جامعة القاهرة

Mo'men Sayed Osman Elnasharty

Teaching assistant – Library, documents and Information Dept.

Faculty of Arts – Cairo University

Navigator001@gmail.com

المحتويات

1- الأطار النظري للدراسة.

- a. المقدمة.
- b. مشكلة الدراسة.
- c. أهداف الدراسة.
- d. منهج الدراسة.
- e. تساؤلات الدراسة.
- f. حدود الدراسة.

2- التكامل المعرفي على صعيد محتوى الويب:

- a. مفهوم التكامل المعرفي.
- b. واقع التكامل المعرفي على الويب.
- c. التحديات التي تواجه شبكة الويب في تحقيق التكامل المعرفي.

3- التنقيب عن البيانات The Data Mining:

- a. النشأة والمفهوم
- b. مراحل التنقيب على البيانات.
- c. معمارية التنقيب على البيانات.
- d. مهام التنقيب عن البيانات.
- e. آليات التنقيب عن البيانات.
- f. خوارزميات التنقيب عن البيانات.

4- التنقيب على الويب The Web Mining:

- a. الدوافع وراء التنقيب عن المحتوى على الويب.
- b. محدودية الويب في تحقيق التكامل المعرفي.
- c. مفهوم التنقيب على الويب.
- d. أقسام التنقيب على الويب.
- e. آليات التنقيب على الويب.
- f. خوارزميات التنقيب على الويب.

5- الانطولوجيات والتنقيب على الويب :

- a. مفهوم الانطولوجيات.
- b. البنية البنائية Syntax للانطولوجيا.
- c. لغات انطولوجيا الويب.
- d. مستويات انطولوجيا الويب.
- e. وظائف الانطولوجيات في أطار التنقيب على الويب.
- f. دور الانطولوجيا في تحقيق التكامل المعرفي على الويب.

6- نتائج الدراسة.
7- توصيات الدراسة.

المستخلص:

سعت هذه الدراسة إلى توضيح مفهوم التكامل المعرفي في سياق تقني، وأشهر المبادرات التقنية التي ساهمت في تحقيقه على صعيد المحتوى المتاح على الشبكة العنكبوتية، ثم تناولت الدراسة رصد القدرة على توظيف مبادرة الانطولوجيات في إطار التنقيب عن البيانات سعياً بذلك إلى الاستفادة منهما فيما يعرف بنظم الانطولوجيا للتنقيب على الويب Ontology-based Web mining لتحقيق التكامل المعرفي لمحتوى شبكة الويب.

اعتمدت الدراسة على المنهج المسحي لرصد التحديات الخوارزمية واللغوية التي تواجه شبكة الويب في تحقيق التكامل المعرفي، كذلك تعتمد الدراسة على المنهج التحليلي في رصد وتحليل واقع قدرات الانطولوجيا في أن توظف لتحقيق التكامل المعرفي للمحتوى المتاح على الويب.

وكان من أهم نتائج الدراسة هي قدرة التنقيب على الويب على تحقيق مستوى مرتفع من التشغيل المتبادل على صعيد محتوى الويب بأكمله لتوفير التكامل المعرفي.

This study aimed to clarify the concept of knowledge integration in the information science context, aimed to study the most famous technique initiatives that contributed to achieve knowledge integration at the level of the available content on the web, and monitoring ability to Ontology initiative within the framework of data mining in order to take advantage of them in Web Ontology - based Web mining system to achieve the knowledge integration.

This study based on survey method to monitor algorithm and linguistic challenges facing the web in the knowledge integration path.

As well as the study is based on the analytical method in analyzing the capabilities of the ontology to employ it to achieve the knowledge integration.

The most important result of the study was the ability of the web mining to achieve a high level of interoperability between web resources to provide knowledge integration.

1- الأطار النظري للدراسة:

المقدمة:

يأتي التكامل المعرفي كأحد أهم التطلعات التي سعت إليها الكثير من المبادرات التكنولوجية لتحقيقه على صعيد محتوى شبكة الويب، فقد أمست الويب (الشبكة العنكبوتية العالمية - The world wide web) منصة عمل فريدة من نوعها، لم تشهد نظم المعلومات مثيلاً لها من قبل، ومستودعا أساسيا لزخم متراكم من مصادر المعلومات، والتي وجدت طريقها للنشر والإتاحة في بيئة اتسمت بفجوة كبيرة في تحقيق التكاملية بين مصادرها المختلفة وبين احتياجات المستفيدين منها، موسومة في ذلك بنمو مطرد لمحتواها، وارتفاع ديناميكي في معدلات التغيير والتحديث، الأمر الذي جعل البعض يصف هذه الديناميكية بأن الحديث عن الويب بطرق موثقة أو شفوية لا يتسم بالإستقرار، فما يكاد أن يلبث إلا أن يتقدم محتوى هذا الحديث خلال فترة زمنية قصيرة، كما جاء تفرد الويب في طبيعتها المعمارية والتكوينية، فلم تصمم البنية المعمارية للويب على أن تعمل وفقاً لمنطق نظم إدارة قواعد البيانات Database Management Systems من حيث الهيكل والتنظيم المطرد للبيانات والمعلومات، ولم تصمم الويب أيضاً على أن تصبغ بهيئة الفهارس الببليوجرافية Online Public Access Catalogue من حيث ما تقدمه الثانية من عناصر تنظيمية واسترجاعية، وصيغ للبيانات وأشكال للاتصال Communication Forms كصيغة مارك MaRC تلك الصيغ التي كفلت القدرة على ضمان النسق والتوحيد في الإدخال للبيانات والمعالجة والتشغيل البيني بين الفهارس، فضلاً جذرية تغيير أحدثتها في سلوك المستفيدين منها تجاه تعاملهم معها ومع مصادرها.

بل صممت الويب لتتيح من خلالها كل شيء عن أي شيء، مما استتبع في أن تكون بمثابة الحياة البرية للمعلومات - (فقد قدر حجم المصادر المعلوماتية المتاحة على الويب في اغسطس عام 2000 بنحو 7 ملايين صفحة بعدد مستخدمين لها قدر بـ 500 مليون مستخدم، ليصل حجم الشبكة في اغسطس 2010 إلى نحو 7.74 مليار صفحة بعدد مستخدمين قدر بنحو 2 مليار مستخدم¹). فضلاً عن التباين والتنوع الموضوعي واللغوي والنوعي والشكلي والجغرافي لما تشمله من مصادر، ولعل المرجعية الأساس وراء ما تعانية الويب من تحديات يعود إلى الإزدواجية Duplication في تكوينها، فهي بيئة استرجاعية تعمل في نفس الوقت كبيئة للنشر والإتاحة الحرة، مما جعل أمر ضبط مصادرنا وتنظيمها أمراً يكاد أن يكون مستحيلًا في اكتماله، هذا الأمر كان بمثابة دافعا نحو التفكير في تطوير العديد من التقنيات التي تكفل القدرة على استثمار هذا الحجم الهائل من المحتوى لتحقيق التكامل المعرفي، وعلى هذا جاءت العديد من المبادرات التي تسعى إلى ضبط مصادر المعلومات المتاحة على شبكة الويب رامية بذلك نحو التكامل المعرفي، فجاءت مبادرة استخراج البيانات Information Extraction، ثم مبادرة التنقيب عن البيانات Web Mining، وتلي ذلك مبادرة الانطولوجيا Ontology والتي تعد من أبرز هذه المبادرات وأكثرهما توظيفاً.

¹Kunder, M. d. (n.d.). WorldWideWebSize.com | The size of the World Wide Web (The Internet). Retrieved September 21, 2011, from <http://www.worldwidewebsite.com>.

يشير مصطلح التنقيب على شبكة الويب Web Mining الى القدرة على اكتشاف المعرفة من واقع البيانات المتاحة على شبكة الويب، مرتكزة في ذلك على ثلاثة قطاعات رئيسية:

- ١- التنقيب عن محتوى الويب Web content mining.
- ٢- التنقيب اعتمادا على بنية الروابط المتاحة على الويب Web structure mining.
- ٣- التنقيب من واقع انماط الافادة والاستخدام على الويب Web usage mining.

وتلعب الانطولوجيا Ontology دورا رئيسيا في كلا من القطاعات الثلاث إذ تعمل على تمثيل المحتوى وفقا للمفاهيم والدلالات المراد التعبير عنها معتمدة في ذلك على انماط المنطق والإستدلال الرياضي ، كما تعمل على تحديد العلاقات والروابط بين البيانات المختلفة على الويب، وتعمل على تصنيف المحتوى لفئات وفقا لموضوعية استخدامه.

وفي هذا تعمل هذه الدراسة على رصد توظيف مبادرة الانطولوجيات في إطار التنقيب عن البيانات سعياً بذلك إلى الاستفادة منهما فيما يعرف بنظم الانطولوجيا للتنقيب على الويب Ontology-based Web mining كل من المبادرتين لتحقيق التكامل المعرفي لمحتوى شبكة الويب.

مشكلة الدراسة:

يمكن بلورة مشكلة الدراسة الأساسية في:

- 1- عجز تطبيقات وبرمجيات ومحركات شبكة الويب الحالية عن تحقيق التكامل المعرفي على صعيد محتواها، وذلك في ظل ما تتسم به شبكة الويب الحالية من قضايا شائكة ومشكلات متفاقمة.
- 2- دراسة امكانيات وقدرات الانطولوجيات على تحقيق التكامل المعرفي على صعيد كل من محتوى الويب وروابطه وأنماط الاستخدام.

أهداف الدراسة:

- 1- رصد التحديات الداخلية والخارجية التي تواجه تحقيق التكامل المعرفي على صعيد محتوى الويب من حيث القدرة على اكتشاف المصادر ذات العلاقة والصلة الموضوعية فيما بينها، ومن حيث معالجتها ومن حيث القدرة على التشغيل البيئي بين أدوات وانظمة البحث والاسترجاع المختلفة.
- 2- التأصيل النظري للانطولوجيات في إطار التنقيب عن البيانات.
- 3- التعرف على أبرز ملامح تقنيات وخوارزميات التنقيب عن البيانات والتي تمثل الأطار النظري لهذه الدراسة.
- 4- مدى قدرة الانطولوجيات على تحقيق التكامل المعرفي على صعيد محتوى الويب.

منهج الدراسة:

تعتمد الدراسة على المنهج المسحي لرصد التحديات الخوارزمية واللغوية التي تواجه شبكة الويب في تحقيق التكامل المعرفي، كذلك تعتمد الدراسة على المنهج التحليلي في رصد وتحليل واقع قدرات الانطولوجيا في أن توظف لتحقيق التكامل المعرفي للمحتوى المتاح على الويب.

تساؤلات الدراسة:

تحاول الدراسة الإجابة على التساؤلات التالية:

- ما هي أبرز التحديات الداخلية والخارجية التي تواجه برمجيات ونظم الويب في تحقيق التكامل المعرفي؟
- ماهي أبرز سمات وملامح التنقيب عن البيانات على الويب؟
- كيف يمكن للانطولوجيات أن تساهم في حل قضايا ومشكلات الويب في تحقيق التكامل المعرفي من خلال ما تم اقتراحه لها من معايير ومواصفات؟

حدود ومجال الدراسة:

تأتي الدراسة وفقاً للحدود الآتية:

1. الحدود الموضوعية:
تتناول الدراسة رصد قدرة الانطولوجيات في تحقيق التكامل المعرفي.
2. الحدود النوعية:
حيث تركز الدراسة على تناول الانطولوجيات في إطار التنقيب عن البيانات دون أن تتطرق إلى الأطر الأخرى - كالويب الدلالي - التي قامت هي الأخرى بتوظيف الانطولوجيات في سياقها.

التكامل المعرفي على صعيد محتوى الويب:

مفهوم التكامل المعرفي:

يجدر البيان، بأن قبيل التطرق إلى مفهوم التكامل المعرفي وواقع وجوده على صعيد محتوى الويب، ستقوم الدراسة بتوضيح مفهوم المعرفة ودور النظم في معالجتها وتيسير الإفادة منها.

كانت المهمة الأساس للحاسب الآلي وبرمجيات المعالجة الآلية للنصوص والأرقام ونظم البحث والاسترجاع، ومنصات التشغيل المختلفة على السواء، في القرن المنصرم، هو الحصول على المعلومات القيمة والقابلة للتطبيق من واقع تحويل أنماط البيانات إلى نسق من المعلومات، وعلى هذا أصطبغ هذا القرن بهوية "عصر المعلومات"، ولكن يأتي القرن الواحد والعشرون حاملاً معه تحولاً جوهرياً في تحويل تلك النظم والبرمجيات من مجرد معالج للبيانات ومقدم للمعلومات وعارض لها، إلى أنظمة لديها القدرة على استخراج المعرفة واستنباطها من واقع ما أتيح لها من معلومات لينشئ في هذا السياق ما يعرف بـ "عصر المعرفة".

وفي هذا، عمدت تقنيات ونظم عصر المعرفة إلى توجه جديد، وهو الانتقال بالمعلومات من ذاتية البرمجة وفردية النظم وعزلة المستودعات، إلى نطاقات التكامل والتقاسم والتشارك، هادفة بذلك إلى زيادة حجم

المعرفة وتوسيع عمق الاستدلال والاستنباط للمعلومات، هذا الأمر الذي يسفر عن أنماط جديدة من الذكاء الاصطناعي.

يعرف قاموس ويبستر المعرفة بأنها "مجموع تحصيل الإنسان للحقائق والمعلومات والمبادئ والقواعد التي اثمرت خلال فترة حياته"².

ويعرفها كل من Sunasee and Sewery الخبرة البشرية المكتسبة من واقع تجاربه وتفاعله وقرءاته وتعلمه في البيئة المحيطة به والمخزنة في العقل البشري"³.

تصنف المعرفة في الاسانيد الفكرية وادبيات الموضوع الى الأنماط الآتية:

- المعرفة القائمة على التصنيف.
- المعرفة الموجهة لاتخاذ القرارات.
- المعرفة الوصفية.
- المعرفة الإجرائية.
- المعرفة المنطقية.
- المعرفة الاستنباطية⁴.

وفي هذا الصدد، تعددت الرؤي في مفهوم التكامل المعرفي Knowledge Integration فيما بينها، وذلك لتعدد السياقات والمجالات المعرفية التي هدفت لتحقيق التكامل في مجالها، ففي سياق العلوم الإدارية ينظر للتكامل المعرفي بأنه قدرة المنظمة على الإستغلال الأمثل للمعلومات والقيم المتوافرة بغية استكشاف أنماط جديدة للمعلومات.

أما في سياق الذكاء الاصطناعي، فينظر اليه في كونه أحد أهم المهام والمتطلبات الأساسية لبناء نظام معرفي يكفل القدرة على تفاعل المستفيد مع الآلة بصورة ذكية.

يمكن النظر إلى مفهوم التكامل المعرفي بأنه عملية تنطوي على دمج وإدراج معلومات جديدة في كيان معرفي قائم يشتركا في بنية مفاهيمية واحدة أو درجة من الصلة أو العلاقة الموضوعية، حيث يعمل التكامل المعرفي على تحدد كيفية أن تتكامل المعلومات الجديدة مع الكيان المعرفي القائم، وكيف يمكن أن يتم تعديل هذا الكيان المعرفي لاستيعاب المعلومات الجديدة، وكيف يمكن تعديل المعلومات الجديدة في ضوء المعرفة الحالية⁵.

² Merriam-Webster (2001). Merriam-Webster, n.d. Web. 28 Sept. 2012. <<http://www.merriam-webster.com/>>.

³ Sunasee and Sewery (2002). Introduction to Knowledge Modeling Available at www.makhfi.com/KCM_intro.htm

⁴ Ibid.

⁵ Murray, K. S. (1996) KI: A tool for Knowledge Integration. Proceedings of the Thirteenth National Conference on Artificial Intelligence.

كذلك يمكن فهم مفهوم التكامل المعرفي في ضوء مقارنته بالتكامل المعلوماتي *information integration*، حيث ينطوي التكامل المعلوماتي على دمج المعلومات ذات مخططات البناء المختلفة في حين يهدف التكامل المعرفي إلى التركيز وبصورة أساسية على توليف وجهات النظر المختلفة القريبة والبعيدة لفهم موضوع محدد، أو بعبارة أخرى يركز التكامل المعرفي على العلاقات المفاهيمية والدلالية في حين يركز التكامل المعلوماتي على دراسة العلاقات الترابطية بين المعلومات بعضها البعض⁶.

ولكن يمكن أن ينظر إلى التكامل المعرفي في سياق المجرد بأنه " العملية التي تكفل القدرة على تجميع وتوليف وتلخيص نماذج من البيانات Data Models داخل نموذج مشترك أو نموذج واحد يوفر حالة معرفية فريدة"⁷.

واقع التكامل المعرفي على الويب:

يوضح التعريف السابق أن جوهر التكامل المعرفي يتجلى في اعتماده على نماذج البيانات Data Model والتي تعرف بأنها منهجية أو بنية Method or Structure تعمل على تنظيم البيانات وفقا لنسق وهيكل محدد هذا النسق يساعد على تخزين وتيسير البحث والاسترجاع عن المعلومات بصورة تتعالى فيها الكفاءة والفاعلية.

هذه النماذج تتباين في بنيتها وخوارزمياتها ومنهجية عملها وأغراض انشائها، ومن أمثلتها نماذج العلائقية للبيانات المستخدمة في نظم إدارة قواعد البيانات DBMS ونماذج الاتصال المعيارية في فهارس المكتبات كنموذج MARC21 ومخططات اللغات التكويدية XML Schemas، الأمر الذي يسفر عن أن يكون كل نموذج بمثابة جزيرة منفصلة بما يشمله من معلومات، ويلزم كل من يتعامل معه بالبحث والاسترجاع أن يعمل وفقا لقواعد النموذج وخوارزمياته.

وعلى الرغم من تعدد هذا النماذج وأهدافها إلا أن سمة قاسمين مشتركين بينهم يتمثلان في:

- كونهم قوالب تشتمل على المضمون والمحتوى والمعلومات.
- أن كل منهم – بالرغم من اختلاف بنيته - لديه القدرة على أن يعمل تحت منصة الويب.

هذا الأمر الذي كان بمثابة دافع أساسي في التفكير لإيجاد طريقة أو منهجية ما تعمل على توفير التشغيل البيئي بين نماذج البيانات المختلفة، وقد تمثلت هذه المنهجية في مبادرات متعاقبة، يأتي على رأسها مبادرة الشبكة العنكبوتية العالمية The World Wide Web.

في عام 1989 قام Tim Berners Lee بإنشاء العنكبوتية كمنصة عمل على شبكة الانترنت تخلص إلى نظام عالمي، حيث تعتمد على وضع معرفات فريدة للوثائق ونماذج البيانات وربطها بعضها البعض وتمثيلها

⁶ Ibid.

⁷ Linn, M. C. (2006) The Knowledge Integration Perspective on Learning and Instruction. R. Sawyer (Ed.). In The Cambridge Handbook of the Learning Sciences. Cambridge, MA. Cambridge University Press.

بصورة تسمح بالتشغيل البيئي، وفي عام 1990 قام باختبار رؤيته والذي اثبتت جدارة ونجاح، وفي عام 1991 قام بإنشاء أول خادم للعنكبوتية عرف بأسم world wide web ومنه اشتق اسم النظام، اما الدور الذي تميزت به الويب فتمثل في تكوين بنية معمارية اساس والتي تجلت فيها مفهوم التشغيل البيئي والمتبادل وهي سمة النجاح التي تميز بها هذا النظام ولكي تتحقق سمة التشغيل البيئي بين الأنظمة حدد منشئها ثلاثة مبادئ اساسية لها وهي:

- محددات الهوية Identification: ويرمز لها بـ URI والتي تهدف الى إكساب نماذج البيانات محددات ومعرفات فريدة على صعيد العالم بأسره بغية تحديد هويتها وتميزها عن بعضها البعض ويتم ذلك من خلال الاعتماد على أحد السمات المميزة لها وليكن موقعها على الشبكة (URL) او رقمها الدولي (DOI) او اسم مميز لها (ISBN) ك معرف لها.

- التفاعلية Interaction: ويقصد بها الأدوات التي تمكن اطراف النظام من التفاعل عبر عملية اتصالهم ببعضهم البعض، حيث يتم ذلك من خلال مجموعة من البوتكولات التي تحكم وتنظم قواعد الاتصال والتفاعل بين المرسل والمستقبل لمصدر المعلومات، ويعد أشهر نماذجها بروتوكول TCP/IP وبروتوكول HTTP.

- الصيغ أو التنسيقات Formats: ويقصد بها منهجيات تمثيل وتكويد المحتوى بصيغ تسمح بالتشغيل البيئي وتمكن بروتكولات النقل من نقلها، وتمكن نظم التشغيل المختلفة من حفظها، وبرامج التطبيقات من عرضها ويعد أشهر الصيغ لهذا الامر هي صيغ النصوص الفائقة Hypertext.

وبهذا تتراءى شبكة الويب على أنها نظام عالمي لنشر نماذج البيانات المختلفة، وشبكة للربط المصادر، ومستودع مفتوح للوثائق الموزعه، وارشيف متكامل للمحتوى، وعليه استلزم الامر لهذه المنظومة أن توفر أدوات البحث والاسترجاع لما تمده من محتوى ومعلومات ولتكون البوابة الأساسية لتحقيق التكامل المعرفي وعليه ظهر مفهوم البنية الاسترجاعية العنكبوتية للمعلومات أو ما يعرف بنظم استرجاع العنكبوتية للمعلومات Web Information Retrieval (Web IR).

وقد حُملت على هذه الأدوات الأمل في تحقيق التكاملية المعرفية للمحتوى المتاح على الويب، ولكن تضاعف لها الأمر فقد عملت كل من هذه الأدوات بمنأى عن بعضها البعض. فجاءت أدوات البحث والاسترجاع على المعلومات متخذة أشكال عدة وسمات مختلفة في بنائها وخوارزميات عملها فجاء منها:

- الأدلة الموضوعية. Web directories.
- محركات البحث على الشبكة العنكبوتية Web search engines.
- ما وراء محركات البحث Meta search engines
- بوابات الشبكة العنكبوتية Web portals
- فهارس العنكبوتية غير المرئية Invisible web catalogue.

○ فهرس المكتبات على العنكبوتية Web public access library catalogue .

وقد أوضح ضياء عبد الواحد⁸ ان المرجعية في التعددية بين أدوات البحث يعود إلى تقديم كل أداة للشكل من المعلومات فضلا عن التفاوت في خبرات المستخدمين في التعامل مع كل أداة وطبيعة احتياجاتهم البحثية التي تفرض التوجه للأداة محددة دون غيرها لكونها مناسبة لطبيعة بحثه. ويضيف الباحث على هذا، أن المسعى في تعددية هذه الأدوات يرجع إلى عجز كل أداة على حدى في تحقيق التكامل بين طبيعة ما تقدمه من بيانات ومعلومات ومحتوى وبين نمط آخر، فعلى سبيل المثال يعجز محرك البحث عن تقديم البيانات والمعلومات في صورة بيبليوجرافية كالتالي يقدمها الفهرس المتاح على الويب، كذلك يعجز فهرس العنكبوتية غير المرئية من تقديم المعلومات في صورة مبوبة ومتدرجة كالتالي تقدمها الأدلة الموضوعية.

وعلى الرغم من أنه قد تم أرساء شبكة الويب كمنصة عمل لاستخراج محتوى نماذج البيانات المختلفة (DBMS – MARC21... وغيرها) وتجميع مضامينه وإعادة هيكلته بغية التشغيل البيئي وتحقيقا للتكامل المعرفي من واقع ربط المضامين ذات الصلة بعضها البعض من خلال أدوات البحث المختلفة، إلا أن هذه أدوات البحث على الويب قد واجهت تحدي رئيسي يقف في طريق تحقيق التكامل المعرفي على صعيد الويب وهو أن شبكة الويب بجانب ما تحمله من نماذج للبيانات تيسر البحث والاسترجاع فيها ومن ثم تحقيق التكامل المعرفي، إلا أنها تحمل في بنيتها شكل آخر من أشكال قوالب ونماذج البيانات، وهو شكل لم يخضع بأي صورة للهيكلية أو التنظيم، هذا الشكل يعرف تعرف بإسم البيانات غير المهيكلة Unstructured Data والتي تمثل نحو 46% من إجمالي حجم محتوى الويب والمتمثلة في المواقع المحررة بلغة النص الفائق التشعبي Hypertext Markup Language وخطابات البريد الإلكتروني فضلا عن المحتوى ذات الطبيعة الانية Real-time content وهي معلومات تتسم بإنها ديناميكية في وجودها وفي تعطل محتواها، كجداول الطيران واسعار الاسهم وغيرها.

كان أحد أهم الأدوات التي اعتمد عليها مخترع الويب Tim Berners Lee في تحقيق رؤيته هي لغات النص الفائق HTML والتي سرعان ما أتضح أنها أعظم التحديات التي تقف في طريق تحقيق التكامل المعرفي، نظرا لعدم صلاحية هذه اللغة في تحقيق هيكلية للبيانات في ظل انتمائها إلى عائلة لغات التمثيل غير المهيكل للبيانات والتي يصعب من خلالها استثمار المحتوى.

وفي هذا يجدر الإشارة إلى طبيعة بنية البيانات في الويب تتخذ نمطين أساسيين هما:

- البيانات المهيكلة: وهي تلك البيانات التي تسكن نماذج البيانات وهي بذلك خضعت لتنظيم ما وفقا لمعرفات أو محددات منطقية بحيث تسمح بتحديد أجزاء معينة من المعلومات مما يسهل عملية التعرف عليها والبحث عنها واسترجاعها كما هو الحال في قواعد البيانات العلائقية، وتعتمد هذه النوعية على سمات وخصائص تعكس من

⁸ضياء عبد الواحد (2005). محركات البحث المتخصصة دراسة تجريبية . القاهرة: جامعة حلوان قسم المكتبات والوثائق والمعلومات. رسالة دكتوراه. ص

خلالها المحتوى والمعنى والاستخدام وتشكل هذه النسبة 15% من محتوى الويب وزعد لغة structure query language أكثر اللغات المستخدمة لهذا الغرض.

- البيانات غير المهيكلة: تشير إلى البيانات التي لم تخضع لأي صيغة من الهيكلة أو التنظيم أو التحديد لها بشكل مسبق عند انشائها كما لا يظهر فيها المعرفات أو المؤشرات مما يجعل طبيعتها بمثابة كتلة من البيانات المجمعة لها بداية ونهاية دون أية إشارة إلى أقسامها أو القطاعات التي تشتمل عليها وتمثل هذه الفئة نسبة 85% من إجمالي محتوى الويب ، ويشكل البريد الإلكتروني والمدونات والمجموعات الإخبارية والمنتديات والمحادثات والدراسات المسحية والأبحاث والعروض التقديمية والمواقع نسبة 65% من هذه البنية، أما عن النسبة الباقية 35% فهي اسيرة التنظيم في الجداول الإلكترونية وبرامج تحليل البيانات النصية وغيرها، وتعد لغة html أحد أكثر اللغات التكويدية انتشارا - فقد بلغت عدد المواقع المكوّدة بلغة html نحو 85% من إجمالي مواقع الويب - والتي تستخدم لغرض التمثيل دون الهيكلة أو التنظيم. فمصممي ومنشئي المحتوى على العنكبوتية يلجأون إلى تنظيم محتوهم بشكل انطباعي ذاتي أو بطرقهم الخاصة دون اللجوء إلى هيكلة منطقية لها والذي يعد أمرا غير مستحب في معالجة المحتوى، ورغم ما توفره لغة HTML من منهجيات أولية لهيكلة البيانات كالتيجان TAGS ونقاط الارساء ANCHORS إلى أنها تتسم بالمحدودية وعدم الجدوى نسبيا في هيكلة المحتوى.

أما عن أوجه القصور التي تعتلي لغة Hypertext Markup Language فيتحقيق التكامل المعرفي فتتمثل في:

- لا تحمل هذه اللغة أية نماذج مفاهيمية من شأنها أن تكفل التنظيم أو الدلالة أو المعنى بل يقتصر دورها على التمثيل البصري للمعلومات.
- أن هذه اللغة بمثابة كتلة من البيانات غير المهيكلة، حيث لا توفر هذه اللغة توصيف لبنية المحتوى، فهي تمثيل أصم للبيانات.
- تعد هذه اللغة بمثابة عالم فوضوي من الروابط الصماء والتي لا تعني للبرامج والتطبيقات البحثية إلا جملة "هذا الموقع مرتبط مع هذا الموقع" فقط من دون أن تحمل هذه الجملة أي دلالة كأن يشير الرابط إلى أن الموقع مرتبط بآخر بعلاقة "هو نوع من" أو "جزء من" وغيرها من العلاقات الوصفية (رغم ذلك تعتمد أقوى محركات البحث على هذا العنصر في بناء خوارزمياتها وترتيب نتائجها للمستفيد).
- قصور التكامل بين البيانات في ظل التطوير البرمجي المستقل للبرامج المستخدمة لتحريير النصوص بلغة .html
- ان غالبية صفحات html تبنى بشكل غير معياري بمعنى انها منشؤها لايلتزمون بمواصفات الانشاء والتحرير معتمدون في ذلك على قدرة المتصفحات في دعم المواصفات الناقصة ولكن الامر يختلف بشكل كبير عند المعالجة في أدوات البحث لتحقيق التكامل المعرفي.

وعليه، كان على أدوات البحث المختلفة على الويب أن تفاضل فيما بين الفئتين السابقتين من نماذج البيانات (البيانات المهيكلة – البيانات غير المهيكلة)، وذلك من حيث أيهما يمثل حقل عمل هذه الأدوات؟ وأيها ستركز عليه في التجميع والتنظيم والإفادة؟ ومن ثم التكامل المعرفي؟ وقد كان الترويج لصالح البيانات غير المهيكلة ومن ثم فإن أدوات البحث الحالية مهما كانت تعمل على تطوير قدراتها في البحث الذكي عن المعلومات وتحليلها والتكامل المعرفي فإنها محدودة في النهاية بقيود البحث الشكلي وفقا لطبيعة تكويد محتوى الوثائق بلغة النص الفائق التشعبي، ومن هنا كانت محدودية هذه اللغة منبع للتحديات التي تواجه التكامل المعرفي في ظل افتقار برمجيات وتطبيقات الويب النماذج المفاهيمية التي من شأنها ضمان الاتساق بين المحتوى.

التحديات التي تواجه شبكة الويب في تحقيق التكامل المعرفي:

1- منهجيات خداع محتوى Content Spandexing:

يمكن القول بأن بعض التحديات والمشكلات التي تواجه الويب في التكامل المعرفي ارتبطت بظهور ووجود الويب نفسها، والبعض الآخر كان بمثابة تغيرات مستحدثة لمشكلات حظيت من قبل بالدراسة والاهتمام العلمي في أدبيات معالجة البيانات ومجال استرجاع المعلومات والشبكات والدراسات البيبليومترية أي قبل ظهور الويب.

ولقد ارتبطت عمليات البحث على المعلومات على الويب من خلال أدوات البحث المختلفة من فهارس وأدلة وكشافات بتقنيات التنقيب عن المحتوى ومدى أسهامها في دقة النتائج المسترجعة من قبل أدوات البحث، وقد اتسم نمط تعامل وسلوك المستخدمين في إيجاد والبحث عن المعلومات إلى اللجوء والاعتماد بصورة قصوى على النتائج الأولى داخل الصفحات الأولى من نتائج نظام استرجاع المعلومات وقصر التصفح عليها دون النظر إلى باقي النتائج في الصفحات الأخرى، فقد أوضحت Silverstein أن 85% من عمليات البحث يقتصر فيها التصفح على النتائج الأولى فقط، ولعل المرجعية في ذلك إلى قناعة ذاتية من جانب المستخدمين في أن النتائج ذات الصلة باستفساره لا بد وأن تظهر أولاً، وأنه كلما توغل في النتائج الأخرى ابتعدت به عن مجال استفساره⁹.

وعليه ادرك مديرون المواقع والقائمين على إدارة المحتوى تبعية منطقية مفادها في أن تضمين المواقع داخل النتائج العشر الأولى لأدوات البحث يؤدي إلى ما يعرف بارتفاع معدل المرور إلى محتوى الموقع *traffic to web site*¹⁰، وعلى النقيض فإن استثناء أو استبعاد المواقع من الشاشة الأولى أو النتائج الأولى للبحث يسمح لعدد محدود من المستخدمين من رؤية محتوى الموقع أو تصفحه.

⁹C. Silverstein, M. R. Henzinger, J. Marais, and M. Moricz.(1999)."Analysis of a very large AltaVista query log." ACM SIGIR Forum, 33:P6-12.

* ظاهرة المرور إلى المواقع Web traffic تعرف بانها حجم البيانات التي ترسل وتستقبل من قبل زائري المواقع والتي تقدر وفقا للمعادلة بحساب عدد الزائرين وعدد الصفحات التي تم تصفحها من خلالهم ويعد هذا المقياس مؤشر مهما للقائمين على إدارة المواقع في تحديد جدوى صفحاتهم.

وبالتالي يعمل أغلب موفري المحتوى على الويب – خاصة المواقع التجارية – على جعل محتوى صفحات مواقعهم تحظى بترتيب طبقي عالي في نتائج أدوات داخل محركات البحث العامة والأكثر استخداما . ويتم تحقيق ذلك عن طريق منهجية واضحة ومباشرة وهي تحسين جودة صفحات الموقع سواء كان في المحتوى أو الشكل من خلال ما يعرف بمنهجيات تحسين الأداء (SEO- Search Engines Optimization) ولكن تكتنف هذه الطريقة الكثير من المال والوقت والجهد والموارد الأخرى، و عوضاً عن كل ذلك لجأ بعض موفري المحتوى إلى طرق أخرى مختصرة تتمثل في التلاعب في نتائج محركات البحث من خلال استخدام أساليب غير اخلاقية عند بناء محتوى المواقع أو تصميم الصفحات، فيما تعرف هذه المحاولة بتضليل خوارزمية الترتيب الطبقي لنتائج محركات البحث بمنهجيات خداع محركات البحث search engine spam.

تعرف منهجيات خداع محركات البحث ب "المنهجية أو اسلوب الذي يعتمد على استخدام بعض الاليات المصممة عمدا لرفع ترتيب المواقع أو الصفحات في نتائج محركات البحث"¹¹.

أما عن أوجه الضرر التي قد تلحقه منهجيات خداع المحتوى بتقنيات التنقيب عن الويب فتتمثل في:

- تصعب مهمة المستفيدين في ارضاء حاجتهم المعلوماتية information needs.
- إرساء خلفية سيئة تجاه محرك البحث ونتائجه وعدم مصداقيته.
- تلويث كشف محركات البحث بمواقع زائفة.
- حرق النطاق الترددي للزواحف أو ما يعرف بBandwidth.
- تشوية نتائج محركات البحث.
- قدرة الخادعات على توليد العديد من شبيهاتها.
- قد يمتد الضرر إلى مزودي محتوى الويب ومديري المواقع، فقد يدفع مدراء المواقع إلى التخلي عن جهودهم في توفير صفحات ومحتوى ذو جودة مرتفعة، اخذين على عاتقهم مهمة انشاء محتوى يستطيع ان يتعامل مع أدوات البحث الملوثة مما يستتبع أن يجعل محتوى الويب موجه إلى أدوات البحث بدلا من أن يوجه إلى مستفيدي شبكة الويب الامر الذي يشكل كارثة في نهاية المطاف على الشبكة العنكبوتية¹².

فئات برمجيات خداع محركات البحث:

يمكن أن تصنف خادعات محركات البحث من خلال تقنياتها وأهدافها إلى:

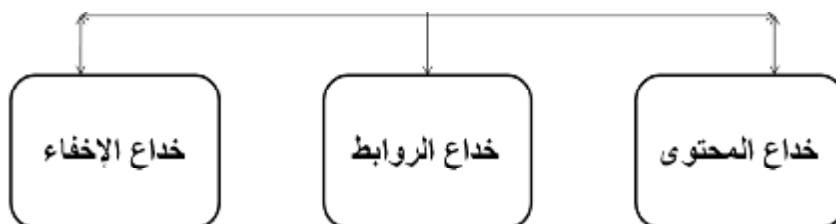
§ خداع محتوى المواقع.

§ خداع روابط الصفحات.

¹¹ Z. Gyongyi and H. Garcia-Molina.(2005) Web spam taxonomy.In First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb). Retrieved July 21, 2011, from <http://www.airweb.cse.lehigh.edu/2005/gyongyi.pdf>

¹² A. Ntoulas, M. Najork, M. Manasse, and D. Fetterly.(2006). Detecting spam web pages through content analysis. In Proceedings of the World Wide Web conference, pages 83-92. available at <http://www.research.microsoft.com/apps/pubs/default.aspx?id=65140>

§ خداع معتمد على إخفاء أو تضمين الصفحات.



شكل رقم (1) يوضح أنواع منهجيات خداع المحتوى.

1- خداع المحتوى:

تعتمد هذه الفئة على تغيير محتوى صفحات الويب للحصول على رتبة أعلى في نتائج محركات البحث، إن أغلب خوارزميات محركات البحث المستخدمة في ترتيب النتائج تعتمد وبشكل أساسي على استخدام أنماط مختلفة من المعادلات والحسابات المختلفة ولكن تعد خوارزمية tf-idf القاسم المشترك بين أغلب هذه الخوارزميات وتحسب العلاقة بين الوثيقة والاستفسار من خلال المعادلة الآتية والتي سبق الإشارة إليها:

$$tf-idf(p, q) = \sum_{t \in p, t \in q} tf(t) \cdot idf(t)$$

وفي هذا يعتمد الخادعون spammers على هذه المعادلة لتحقيق هدفين أساسيين:

جعل صفحاتهم تتصل بعدد كبير من الاستفسارات وذلك من خلال جعل صفحاتهم لا تحصل على الرقم صفر كنتاج لهذه المعادلة ويتم تحقيق ذلك من خلال تضمين مجموعة كبيرة وضخمة من المصطلحات بنية الوثيقة.

جعل الوثيقة أكثر صلة باستفسار محدد وبالتالي تحظى برتبة عالية في نتائج محركات البحث من خلال ان تحصل على درجة عالية في المعادلة السابقة تكاد تصل إلى رقم 1 ويتم تحقيق ذلك من خلال تكرار مصطلحات محددة في بنية الوثيقة.

2- الويب غير المرئي The Invisible Web:

يمكن النظر إلى الويب وفقا لمنظور التنقيب على الويب Web Mining على أنها تنقسم إلى قسمين:

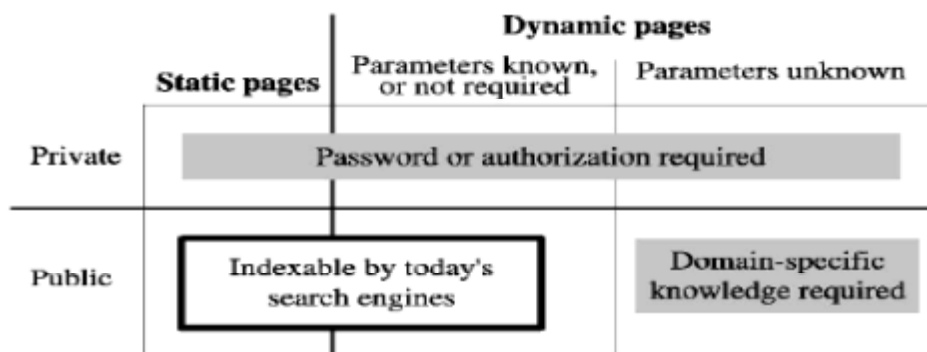
- القسم الأول: الويب السطحية Surface web:

أو الويب القابلة للتكشيف **Indexable web Publically** أو الويب المفتوحة **Open Web**: وهي ذلك المحتوى والصفحات والمواقع التي تخضع للتجميع والتكشيف والاسترجاع والتحليل والمعالجة من قبل تقنيات التنقيب على الويب وذلك من خلال تتبع الروابط.

- القسم الثاني: الويب غير المرئية **Invisible web**:

وتعرف أيضا بالويب العميقة **the deep web**، أو الويب المظلمة **Dark net** أو ما تحت الويب **Under net**: وهي شبكة الويب التي تشتمل على مصادر المعلومات لإدوات البحث أو تقنيات التنقيب على الويب من الوصول إليها ومن ثم أخضاعها للتجميع أو التكشيف أو التسكين في مستودعاتها¹³.

فجز تقنيات التنقيب على الويب وأدوات البحث لا يتمثل في عدم قدرتها على التعامل مع الويب الخفية وحسب، بل يمتد ليتضح أن تغطيتها لما هو متاح على الويب المرئية لا يتجاوز من 20% إلى 50%¹⁴.



شكل رقم (2) رسم توضيحي يوضح فئات الوثائق من خلال أنواع الويب¹⁵.

ويمكن تحديد خصائص الويب غير المرئية على النحو الآتي:

- يقدر حجم الويب غير المرئية من 400 - 500 مرة من الويب المرئية.
- تشتمل الويب على ما يقرب من 70 ألف تيرابايت من المعلومات مقارنة بـ 9 آلاف تيرابايت متاحة على العنكبوتية القابلة للتكشيف.
- تشتمل الويب غير المرئية على ما يقرب من 550 بليون وثيقة منفردة مقارنة ببليون وثيقة متاحة على نظيرتها.

¹³Sherman, C., & Price, G. (2001). *The Invisible Web: uncovering information sources search engines can't see*. Medford, N.J.: CyberAge Books.

¹⁴ MICHAEL K. BERGMAN .The Deep Web: Surfacing Hidden Value. available: <http://brightplanet.com/images/uploads/12550176481-deepwebwhitepaper.pdf>

¹⁵ Ibid.

- يشتمل الويب غير المرئية على أكثر المصادر تتمتع بمصداقية وموثوقية عما هو متاح على العنكبوتية السطحية.
- تشتمل غير المرئية على أغلب المصادر التي يمكن أن تحقق نسبة مرتفعة من تلبية الحاجات المعلوماتية.
- 95% من معلومات الويب غير المرئية متاح مجاناً¹⁶.

3- سلوك المستخدمين البحثي:

ذهب Broder إلى أن نشأة شبكة الويب حمل معها العديد من الأنماط الجديدة من الحاجات المعلوماتية للمستخدمين سواء كان على صعيد البحث أو الاسترجاع ، وقد امتدت هذه الانماط لتشكيل فئات جديدة من الاستفسارات المعلوماتية لدى المستخدمين، ولذا كان لزاماً على تقنيات التنقيب على الويب ومبادرات التكامل المعرفي لمحتوى الويب التوجه إلى حصر هذه الحاجات المعلوماتية تمهيدا لتوظيف التقنيات التي تعمل على أشباعها، ففي سياق الويب لا تقتصر الحاجة من وراء الاستفسار على استرجاع الوثائق فحسب - كما هو الحال في نظم استرجاع المعلومات وقواعد البيانات- بل امتدت هذه الحاجة لتشمل 3 اقسام تتمثل في:

- الحاجة الملاحية Navigational need.
- الحاجة المعلوماتية Informational need.
- الحاجة الاجرائية Transactional need .
- الحاجة الاستنتاجية Deductive need¹⁷.

1- الحاجة الملاحية:

يعد نمطا حديثا من الاستفسارات، فالغرض منه هو الوصول إلى موقع أو مصدر معلومات معين لدى المستخدم، فالحاجة من وراء هذا الاستفسار تتمثل في أن المستخدم ربما قد زار أو تعرض لهذا الموقع من قبل، أو على فرضية أن مثل هذا الموقع يمكن ان يكون موجودا على الويب. هذا النمط من الحاجات أو البحث كان يشار اليه سالفاً بالبحث عن "المفردة المعروفة" بمعنى أن يكون المستخدم على دراية وعلم بأن هذه المفردة - سواء كانت هذه المفردة وثيقة أو ملف صوتي أو موقع .. الخ - موجودة على الويب وجزير بالإشارة إلى أن هذا النمط من الحاجات قد ظهر في نظم الاسترجاع الكلاسيكية ولكن كان الغرض منه يتمثل في تقييم أداة الاسترجاع.

2- الحاجة المعلوماتية:

¹⁶ MICHAEL K. BERGMAN .The Deep Web: Surfacing Hidden Value. available:

<http://brightplanet.com/images/uploads/12550176481-deepwebwhitepaper.pdf>

¹⁷ A. Broder, "A taxonomy of web search," presented at SIGIR Forum, 2002. Available at:

<http://www.sigir.org/forum/F2002/broder.pdf>

يعد هذا النمط تقليديا إلى حد ما فالغرض منه لا يختلف كثيرا عن الغرض من الحاجات المعلوماتية من نظم الاسترجاع التقليدية وهو ايجاد المعلومات التي يفترض أن تكون متاحة على العنكبوتية في صورة ثابتة ولكن سمة امتياز امتازت به العنكبوتية تمثل في قدرتها على المزج بين البيانات المتاحة من خلال المصادر المختلفة في صورة اقرب لتحقيق التكاملية بين مصادر العنكبوتية فمثلا من خلال البحث عن مدينة القاهرة في بعض محركات البحث تكفل خوارزميات هذه المحركات القدرة على توفير مصادر المعلومات التي تتناول القاهرة ثم تستطرد لتسترجع أسماء فنادق القاهرة والمعالم الاثرية بها... الخ .

3- الحاجة الاجرائية:

والغرض من طبيعة هذه الحاجة هو الوصول إلى مواقع تتعالى فيها سمة التفاعل سواء كان بين المستخدمين بعضهم البعض (كالشبكات الاجتماعية) أو تفاعل المستخدمين مع الالة (كتحميل الملفات والتنزيل الهابط وكالمعاملات التجارية مع البنوك).

4- الحاجة الاستنتاجية:

نمط آخر من الحاجات المعلوماتية يعتمد عليه متخذي القرار، وهو القدرة على الاستدلال والاستنتاج من خلال ما هو متاح من معطيات وبيانات ومعلومات للحصول على نتائج جديدة ومستحدثة.

Arabic Speaking Internet Users					
COUNTRIES	Population (2009 Est.)	Internet Users, Latest Data	Penetration (% Population)	User Growth 2000-2009	% User in Table
Algeria	34,178,188	4,100,000	12.0 %	8,100.0 %	6.
Bahrain	728,709	402,900	55.3 %	907.3 %	0.
Comoros	752,438	23,000	3.1 %	1,433.3 %	0.
Djibouti	724,622	19,200	2.6 %	1,271.4 %	0.
Egypt	78,866,635	16,636,000	21.1 %	3,596.9 %	27.
Iraq	28,945,569	300,000	1.0 %	2,300.0 %	0.
Jordan	6,269,285	1,595,200	25.4 %	1,153.4 %	2.
Kuwait	2,692,526	1,000,000	37.1 %	566.7 %	1.
Lebanon	4,017,095	945,000	23.5 %	215.0 %	1.
Libya	6,324,357	323,000	5.1 %	3,130.0 %	0.
Mauritania	73,129,486	60,000	1.9 %	1,100.0 %	0.
Morocco	31,285,174	10,442,500	33.4 %	10,342.5 %	17.
Oman	3,418,085	557,000	16.3 %	518.9 %	0.
Qatar	833,285	436,000	52.3 %	1,353.3 %	0.
Saudi Arabia	28,686,633	7,761,800	27.1 %	3,780.3 %	12.
Somalia	9,832,017	102,000	1.0 %	50,900.0 %	0.
Sudan	41,087,825	4,200,000	10.2 %	13,900.0 %	7.
Syria	21,762,978	3,565,000	16.4 %	11,783.3 %	5.
Tunisia	10,486,339	3,500,000	33.4 %	3,400.0 %	5.
United Arab Emirates	4,798,491	3,558,000	74.1 %	384.1 %	5.
Palestine	2,461,267	355,500	14.4 %	915.7 %	0.
Yemen	22,858,238	370,000	1.6 %	2,366.7 %	0.
TOTAL	344,139,242	60,252,100	17.5 %	2,297.7 %	100

جدول رقم (1) يوضح الهوية المعلوماتية العربية على الويب من حيث عدد المستخدمين ومعدلات النمو في الاستخدام¹⁸

يمكن أجمال التحديات التي تتعلق بتفاعل المستفيدين مع محتوى الويب والتي تعوق تحقيق التكامل المعرفي في العناصر الآتية:

1- حجم الاستفسارات:

فحجم استفسارات المقدمة إلى الويب في نمو مستمر فقد بلغت نحو 10 مليار استفسار عام 2008 موجه إلى 5 محركات بحث (Google, Yahoo!, MSN, AOL, Ask Jeeves). كما أن عدد الاستفسارات الموجهة لمحرك البحث جوجل أكثر من 200 مليون استفسار يوميا عام 2003، لتصل هذه النسبة في عام 2010 إلى

¹⁸ Internet Statstics (2012). Available at : www.internetworldstats.com

3,5 مليار استفسار يوميا بمعدل يصل إلى 40,000 استفسار في الثانية، وبالتالي تعد هذه المؤشرات بمثابة تحد كبير يواجه محركات البحث¹⁹.

2- صياغة الاستفسار:

عادة ما يميل المستفيدون إلى تقديم استفسارات أعم بكثير من الحاجة الفعلية إلى المحتوى، والمرجعية وراء ذلك تكمن في ظن المستفيد بأن أداة البحث يمكن أن تعمل على مستوى أعرض وليس على مستوى التخصص، هذا من ناحية ومن ناحية أخرى كشفت الدراسات المقارنة – القائمة على تحليل سجلات نظم أدوات البحث – أن سلوك المستفيدين في البحث على المعلومات ينطوي على أن معدلات الاستفسارات تتسم بالقصر. فمعظم المستفيدين يدخل ما بين مصطلحين إلى ثلاثة مصطلحات في الاستفسار الواحد، ويدخل ما بين استفسارين إلى ثلاثة استفسارات في عملية البحث الواحدة²⁰.

3- معدلات توزيع مفردات ومصطلحات البحث:

أوضحت إحدى الدراسات البحثية أن معدل توزيع تردد مصطلحات الاستفسار يتسم بالانحراف من إجمالي عمليات البحث، فقد كشفت هذه الدراسة أن بعض المصطلحات تستخدم بشكل متكرر في مقابل الكثير من المصطلحات التي تستخدم مرة واحدة فقط فقد تم الكشف عن أن 63 مصطلح حظي بتردد ظهور بلغ أكثر من 100 مرة، في حين كون هذه المصطلحات تمثل أقل من 1% من إجمالي المصطلحات القابلة للبحث. مما يظهر طبيعة أن البحث على الويب يمكن وصفه بأنه بحث بمصطلحات تتسم بنسبة منخفضة في تردد الظهور، مقابل مصطلحات تتسم بتردد عالي في الظهور.

4- إعادة صياغة الاستفسار:

فمعظم الباحثين على أدوات البحث يلجأون إلى صياغة حاجتهم البحثية في صورة استفسار واحد فقط دون اللجوء إلى إعادة صياغته مرة أخرى بمعدل بلغ اثنين من 3 باحثين وبصورة عامة بلغت نسبة المستفيدين الذين يعتمدون على تعديل استفسارتهم في البحث على العنكبوتية نحو 44% بينما بلغت نسبة من يقوم بتوجيه أكثر من 3 استفسارات في عملية البحث الواحدة نحو 25%.

5- نمط الاستفسار:

إن 50% من إجمالي الاستفسارات المقدمه لأحد محركات البحث كانت تبدأ بعبارات استفهامية مثل "Where do I find . . . ?" كيف أجد" وأن 25% من المستفيدين يبدأون بحثهم بعبارات طلبية مثل "احضر لي

¹⁹Jansen, B. J., & Spink, A. (2003). An analysis of web information seeking and use: Documents retrieved versus documents viewed. In Proceedings of the 4th International Conference on Internet Computing, pp. 65-69. Las Vegas, Nevada. 23-26 June.

²⁰Ibid

معلومات عن **Get me information**، في حين أن معظم محركات البحث تعتمد في المقام الأول على الكلمات المفتاحية مما يسفر في الحالات السابقة عن الكثير من النتائج غير المرضي عنها من قبل المستفيد.

6- الاعتماد على الروابط البولينية والبحث المتقدم:

بلغت درجة الاعتماد على الروابط البولينية في عمليات البحث على درجة تكاد تصل إلى الندرة في عمليات البحث داخل أدوات البحث، فشخص واحد من إجمالي 18 شخص يلجأ إلى الاعتماد على الروابط البولينية فضلاً عن شخصان من أصل 3 اشخاص يستخدم هذه الروابط بصورة خاطئة.²¹

²¹ Jansen, B. J., & Spink, A. (2005). How are we searching the World Wide Web? A comparison of nine large search engine transaction logs. *Information Processing and Management*, 42(1), 248-263.

2- التنقيب عن البيانات Data Mining:

النشأة والمفهوم:

يأتي مفهوم التنقيب عن البيانات على الويب Web Mining كأحد تطبيقات مفهوم التنقيب في قواعد البيانات Data Mining، وكلاهما ينتميان إلى مجال أكتشاف المعرفة من داخل قواعد البيانات Knowledge Discovery from Database وبالتالي يعد من الضروري قبيل التطرق إلى التنقيب عن البيانات على الويب Web Mining، التعرف على مفهوم التنقيب في قواعد البيانات Data Mining والذي يعد جوهر استخراج المعرفة من على الويب.

طور مفهوم التنقيب عن البيانات Data Mining في سياق مجال أكتشاف المعرفة في قواعد البيانات "Knowledge Discovery in Databases" وهو مجال معرفي عمل على توظيف الإحصاء والذكاء الاصطناعي ونظم قواعد البيانات والتعلم الآلي في معالجة محتوى قواعد البيانات، حيث هدف إلى الخروج من مجموعة البيانات المخزنة في قواعد البيانات بمعلومات ذات دلالات ومؤشرات، وتحويلها إلى بنية مفهومة تمهيدا لاستثمارها مرة أخرى والتنبؤ من خلالها بعوامل وفرضيات جديدة في ظل المعطيات المتوافرة، إذ يعمل هذا المجال على تحديد وتعيين البيانات في أدنى مستوى لها وتحويلها إلى معلومات أكثر إحكاما وتجريدا وفائدة حاملة معها القدرة على توليد نموذج تنبؤي لتقدير الأوضاع في المستقبل لمجال ما أو منظمة ما²².

تاريخياً، حظيت فكرة استثمار محتوى قواعد البيانات بالاهتمام في مجال صناعة المعلومات فضلا عن اجتذابه لقدر كبير من الإهتمام من جانب مجالات معرفية مختلفة، نظرا لما تمتاز به هذه البيانات من طبيعة هيكلية وتنظيمية تسمح بتشغيلها في سياقات ونسق مختلفة، ونتيجة لتزايد عدد قواعد البيانات وضخامة حجم بيانات المخزنة بها، ووجود حاجة ملحة لتحويل هذه البيانات إلى معلومات يمكن أن تستثمر للخروج بها إلى معرفة جديدة مجردة، ظهر مفهوم التنقيب عن البيانات Data Mining.

حظي مفهوم استخراج المعرفة من قواعد البيانات بالعديد من المسميات والتي تنوعت بين (استخراج المعرفة – أكتشاف المعلومات – جمع المعلومات – معالجة انماط البيانات – التنقيب عن البيانات)، إلى أن تم الإستقرار على مسمى واحد لها وهو Data Mining ليدل على مفهوم واحد وهو اكتشاف المعرفة في قواعد البيانات KDD في ولكن يجب الإشارة إلى أن الأستخدام الأولي لمفردة التنقيب عن البيانات Data Mining قد ورد في سياق فعاليات ورشة العمل الأولى لاكتشاف المعرفة في قواعد البيانات على هامش اجتماعات اللجنة الدولية المشتركة للذكاء الاصطناعي International Joint Conferences on Artificial Intelligence والمنعقدة في ولاية ميتشجان في الولايات المتحدة.

²² U. Fayyad, G. P.-Shapiro, and P. Smyth. From data mining to knowledge discovery in databases. *AI Magazine*, Vol. 17 No. 3, pp. 37-54, Fall 1996.

يوضح كل من U. Fayyad, G. P.-Shapiro, and P. Smyth أن العلاقة بين مجال اكتشاف البيانات KDD وبين Data Mining هي علاقة أشتمال من الأولى للثانية، حيث أن كينونة KDD تتضح في أنها عملية تهدف إلى اكتشاف المعرفة من واقع البيانات، في حين يشير التنقيب عن البيانات إلى كونه خطوة أو مرحلة معينة في إطار العملية السابقة.

وفي ضوء ذلك، يعرف التنقيب عن البيانات Data Mining بأنه:

" آلية تعمل على تحليل مجموعات من البيانات المجردة من خلال توظيف خوارزميات معينة تعمل على استخراج انماط محددة لإيجاد علاقات غير متوقعة أو لتلخيص البيانات بطرق جديدة وبصورة مفهومة وقابلة للاستثمار في اتخاذ القرار"²³.

لا تعد مرحلة التنقيب عن البيانات هي المرحلة الوحيدة أو الفردية في عملية اكتشاف المعرفة في قواعد البيانات KDD إذ ينطوي مجال اكتشاف المعرفة على عدد من الخطوات والمراحل تحدد من خلالها كيفية اكتشاف المعرفة من البيانات المتاحة، وماهي منهجية حفظ وتخزين البيانات، وما هي أنسب الطرق لأتاحة هذه البيانات، وكيف يمكن أن تقوم خوارزمياتها بتعزيز المعرفة، بحيث تكفل هذه المراحل اكتشاف وتوفير وأشتقاق المعرفة بصورة سليمة من واقع البيانات المتوفرة.

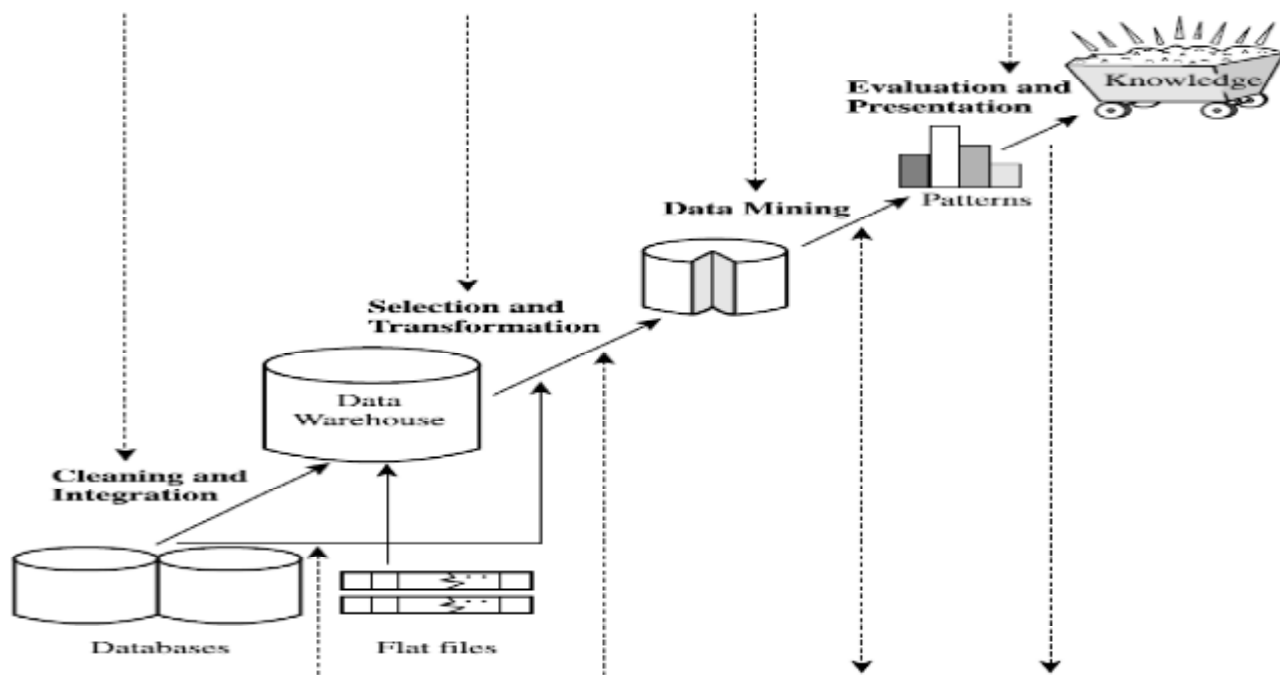
مراحل التنقيب عن البيانات:

وردت العديد من الرؤي في تحديد مراحل اكتشاف المعرفة في قواعد البيانات (موقع التنقيب عن البيانات في منظومة اكتشاف المعرفة) ولكن تعد الرؤية المعيارية في هذا الصدد هي الرؤية التي قام بوضعها كل من Jiawei Han, Micheline Kamber, Jian Pei فعملية اكتشاف المعرفة لديهم هي عملية تفاعلية وتكرارية، تنطوي على خطوات محددة تتمثل في:

- 1- Data cleaning (تنقية البيانات): تختص هذه العملية بحذف البيانات غير المهمة والمكررة.
- 2- Data Integration (تكامل البيانات): وتعمل هذه المرحلة على جمع البيانات من المصادر المختلفة.
- 3- Data Selection (إختيار البيانات): وفي هذه المرحلة يتم إختيار البيانات التي سيتم عليها التحليل.
- 4- Data Transformation (تحويل البيانات): أي عملية توحيد/إدماج البيانات في أشكال محددة مناسبة حتى تكون مهيئة لعملية التنقيب.
- 5- Data mining (التنقيب عن البيانات): وتعد الخطوة الجوهرية حيث يتم استخدام التقنيات الذكية لاستنباط أنماط مفيدة من المعلومات.
- 6- Pattern Evaluation (تقييم الأنماط): أي تحديد الأنماط التي تمثل المعرفة وفق المقاييس المعطاة.

²³ David Hand, Heikki Mannila, and Padhraic Smith, *Principles of Data Mining*, MIT Press, Cambridge, MA, 2001

7- Knowledge Representation (تمثيل المعرفة): وهي المرحلة الأخيرة حيث يتم اكتشاف المعرفة الجديدة، وفي هذه المرحلة يتم استخدام التقنيات المصورة Visualization لمساعدة المستخدمين على فهم وتفسير النتائج المستخرجة²⁴.



شكل رقم (3) يمثل خطوات اكتشاف المعرفة من واقع تحليل البيانات، ويوضح موقع دور التنقيب عن البيانات كمرحلة جوهرية داخل هذه العملية²⁵.

معمارية التنقيب عن البيانات:

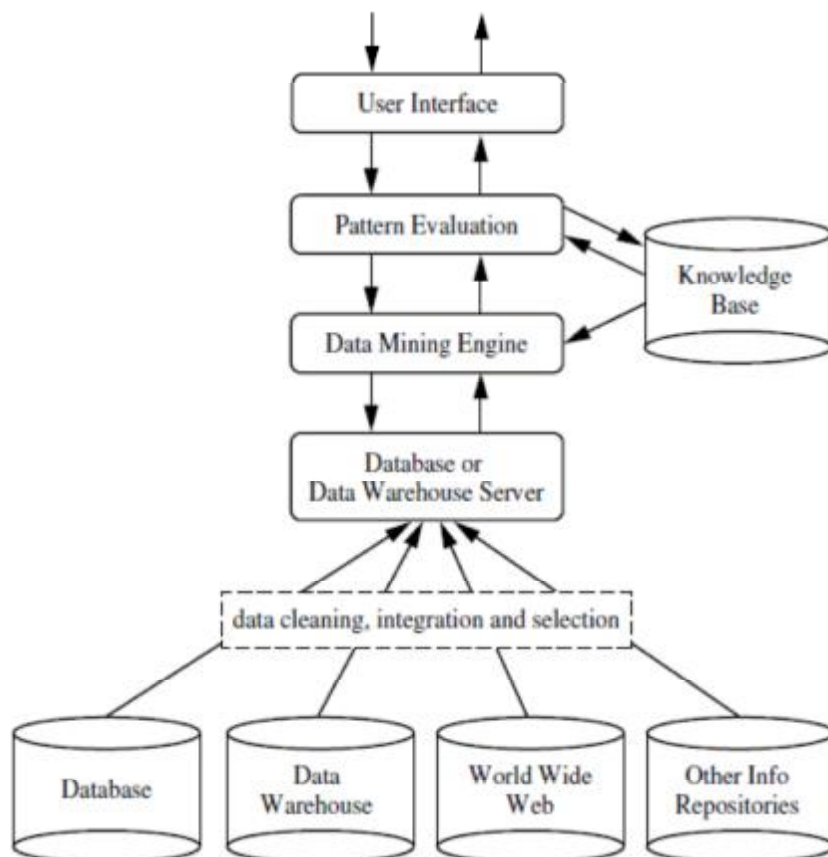
وفي ضوء هذه المراحل والخطوات اقترح كل من Jiawei Han, Micheline Kamber, Jian Pei تصورا لمعمارية مقترحة لما يجب أن تكون عليه نظم التنقيب عن البيانات وجاءت وحداتها على النحو التالي:

- 1- مستودعات البيانات Data Repository: وتتمثل في قواعد البيانات Data Base و شبكة الويب العالمية World Wide Web ومخازن البيانات Data Warehouse وتلعب كلا من مرحلة تنقية وتكامل البيانات الدور الرئيسي في هذه الوحدة.
- 2- خوادم مستودعات البيانات: وهي وحدة مسؤولة عن جلب وإحضار البيانات ذات الصلة بناء على طلب المستخدم في إطار التنقيب عن البيانات.

²⁴ Han, J., & Kamber, M. (2001). *Data mining: concepts and techniques*. San Francisco: Morgan Kaufmann Publishers.

²⁵ Ibid.

- 3- قواعد المعرفة Knowledge Base: ويقصد بها المجال المعرفي المستخدم في توجيه عملية البحث وتقييم درجة الصلة الناتجة عن استرجاع النماذج والأنماط، وتعمل هذه الوحدة على إنتاج الأنماط التي سيتم في ضوءها تحليل البيانات واستنباط المعرفة كما تعمل هذه الوحدة على تنظيم المعرفة في مستويات مجردة الصورة الهرمية التي تتدرج فيها من العام إلى الخاص، ويعد أشهر نماذج المعرفة التي تستخدم في هذه الوحدة هي الميتاداتا.
- 4- محرك التنقيب عن البيانات (Data mining engine): والذي يعد قلب وجوهر نظم التنقيب عن البيانات، وويتألف من وحدات وظيفية فرعية تعمل بالقيام بالمهم الآتية: التوصيف – الربط – تحليل الروابط – التصنيف – التنبؤ- التحليل العنقودي للبيانات – تحليل النماذج – تحليل الاستدلال.
- 5- وحدة تقييم الأنماط: تهدف هذه الوحدة على توظيف مقاييس تهدف إلى تحديد أي الموضوعات تحظى بالاهتمام داخل مجموعات البيانات المختلفة.
- 6- وحدة واجهة المستخدم: وهي الوحدة المسؤولة عن تحقيق التفاعل بين المستخدمين ونظم التنقيب عن البيانات، وذلك من خلال توفير القدرة للمستخدم أن يطرح الاستفسارات والتي سيتم ضوءها التنقيب عن البيانات، كما تعمل هذه الوحدة على توفير إمكانية السماح للمستخدمين بتصفح مخططات قواعد البيانات المحللة، ومستودعات البيانات وبنية البيانات²⁶.



شكل رقم (4) يوضح البنية المعمارية الأساس لنظام التنقيب عن البيانات²⁷.

مهام التنقيب عن البيانات:

حدد كل من Fayyad, Usama; Piatetsky-Shapiro, Gregory; Smyth, Padhraic المهام الرئيسية التي يمكن أن تضطلع بها نظم وتطبيقات التنقيب عن البيانات في ستة قطاعات رئيسية:

- 1- الكشف الخاص Anomaly detection: وتشير هذه المهمة الى الكشف عن أنماط في مجموعة من البيانات والتي تتطلب مزيدا من التحقيق والتحليل نظرا لعدم أنساقها مع بعضها البعض.
- 2- الربط وفقا لقواعد محددة Association rule learning: أو تعرف أحيانا بالتمذجة ذات التبعية Dependency modeling وتشير هذه المهمة إلى البحث عن القات بين المتغيرات المختلفة فعلي سبيل المثال تعمل هذه المهمة على جمع البيانات عن عادات وأنماط سلوك المستخدمين في البحث عن البيانات

²⁷ Ibid.

- وتحليلها وفي ظل هذه المهمة يقوم النظام بدفع بعض المواقع ومصادر المعلومات للمستفيدين من واقع تحليل سجلات تاريخ البحث للمستفيدين وربطها معا.
- 3- العنقدة Clustering: تعمل هذه المهمة على ان تكتشف داخل مجموعات البيانات عن الصلات أو درجات التشابه بين البيانات تمهيدا لجمعها معا في نسق محدد يعرف بالعنقدة.
- 4- التصنيف Classification: تعمل هذه المهمة على تعميم بنية محددة لتطبق على البيانات الجديدة لتسكن وفقا لهويتها أو موضوعها فعلي سبيل المثال يستخدم البريد الالكتروني نموذج محدد وهيكل معين لتحديد وتسكين البريد الوارد في قطاعين المشروع والمزرج.
- 5- الإنحدار Regression : حيث تشير هذه المهمة في فهم كيفية تأثير المتغيرات الثابتة على قيم المتغيرات المتغيرة أو بمعنى آخر يساعد على فهم تبعية التفسير التي تحدث في المتغير التابع من جراء التغيير في المتغير المستقل.
- 6- التلخيص Summarization: أذ تعمل هذه المهمة على تحديد وتمثيل البيانات في صورة أكثر ايجازا تمهيدا لتوليد التقارير لمساعدة متخذي القرار.²⁸

آليات التنقيب عن البيانات:

- أوضح كل من Fayyad, Usama; Piatetsky-Shapiro أن مختلف خوارزميات التنقيب عن البيانات تدور في فلك 3 عوامل أساسية:
- 1- تمثيل النماذج Model representation: ويقصد بها اللغة التي تستخدم لوصف أنماط التنقيب التي تعمل على اكتشاف المعرفة وتحليلها، او بمعنى آخر هي لغة صياغة انماط Patterns.
- 2- معايير تقييم النماذج Model-evaluation criteria: وهي منهجية كمية أو مجموعة من الدوال التي تهدف الى تقييم مدى نجاح الأنماط المحددة في نظام التنقيب عن البيانات، أو بمعنى آخر تعمل هذه الدالات على اختبار ما إذا كانت الانماط Parameters تلبي أهداف قواعد بيانات اكتشاف المعرفة.
- 3- منهجيات البحث Search Method: وتعمل هذه المنهجية على تكوين عنصرين مهمين في نظم التنقيب عن البيانات وهما نماذج البحث ومحددات البحث والذان يتم توظيفهما بغية تحقيق التكامل المعرفي على صعيد البيانات المحللة.

خوارزميات التنقيب عن البيانات.

في عام 2006 قام معهد مهندسي الكهرباء والالكترونيات Institute of Electrical and Electronics Engineers بتكليف فريق بحثي لإعداد دراسة تهدف الى تحديد أشهر الخوارزميات الأساسية

²⁸ Fayyad, Usama; Piatetsky-Shapiro, Gregory; Smyth, Padhraic (1996). "From Data Mining to Knowledge Discovery in Databases". Retrieved 17 December 2008. Available at: <http://www.kdnuggets.com/gpspubs/aimag-kdd-overview-1996-Fayyad.pdf>

المستخدمة في مجال التنقيب عن البيانات ذات التأثير والاستخدام داخل الأوساط البحثية، لتوثق هذه الدراسة في فعاليات المؤتمر الدولي للتنقيب عن البيانات International Conference on Data Mining (ICDM) وقد جاءت هذه الخوارزميات على النحو الآتي:

1- خوارزمية الموجهات في الفراغ Vector space model:

طورت هذه الخوارزمية على يد Gerard Salton عام 1975 في إطار تطويره لنظام أحصائي لتحليل البيانات عرف بإسم SMART (System for the Mechanical Analysis and Retrieval of Text) ، قدم هذا النموذج إطار عمل جديدًا بحيث ينظر إلى محتوى الوثيقة على أنه حقيبة كلمات Bag of words بمعنى أن محتوى الوثيقة يشتمل على مصطلحات غير مرتبة وذات ترددية غير منتظمة داخل محتوى الوثيقة.

كان الاستخدام الأول لهذه الخوارزمية في سياق نظم استرجاع المعلومات حيث تمثلت رؤيته في حساب درجة صلة محتوى الوثيقة بالاستفسار من خلال تحديد درجة التشابه بينهما، حيث يمثل كلا من محتوى الوثيقة والاستفسار في صورة موجهات في فراغ متعدد الأبعاد:

$$d_j = (w_{1,j}, w_{2,j}, \dots, w_{t,j})$$

$$q = (w_{1,q}, w_{2,q}, \dots, w_{t,q})$$

حيث ينطوي كل موجه على اوزان غير ثنائية للمصطلحات الكشفية في كلا من محتوى الوثيقة والاستفسارات والتي يشار إليها بالرمز w_1

وتحسب درجة الصلة للوثائق من خلال مقارنة انحراف الزوايا بين كل من موجه الوثيقة وموجه الاستفسار كما هو موضح من خلال المعادلة الآتية:

$$\cos \theta = \frac{d_2 \cdot q}{\|d_2\| \|q\|}$$

- الأساس الرياضي الذي يعتمد عليه هذا النموذج:

يمكن أن توصف العلاقة بين محتوى الوثيقة **D** والمصطلح **T** من خلال المصفوفة **tf-idf** كمعيار كمي يشتمل على محورين أساسيين:

○ المحور الأول: هو تردد المصطلح **TF** ويشير إلى عدد مرات ظهور المصطلح **t** في محتوى الوثيقة **d** وتأتي المعادلة لحساب تردد المصطلح على هذا النحو:

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}}$$

§ حيث تشير $tf_{i,j}$ إلى حساب تردد المصطلح.

§ تشير $n_{i,j}$ إلى عدد مرات ظهور المصطلح t_i في محتوى الوثيقة d_j .

§ وتشير $\sum_k n_{k,j}$ إلى مجموع عدد المصطلحات في إجمالي الوثيقة.

مثال: إذا افترضنا أن وثيقة ما تتكون من 100 مصطلح، ويظهر مصطلح المكتبات 4 مرات في الوثيقة فإن المعادلة ستكون $0.04 = (4/100)$

○ المحور الثاني: هو تردد الوثيقة العكسي، والذي يعمل على حساب نسبة إجمالي عدد الوثائق المخزنة في النظام إلى عدد الوثائق التي تشتمل على المصطلح T وتظهر معادلته على هذا النحو:

$$idf_i = \log \frac{|D|}{|\{d : t_i \in d\}|}$$

§ تشير idf_i إلى حساب تردد الوثيقة العكسي.

§ بينما تشير \log إلى حساب لوغاريتم ناتج القسمة.

§ وتشير $|D|$ إلى إجمالي عدد الوثائق في النظام.

§ وتشير $|\{d : t_i \in d\}|$ إلى عدد الوثائق التي يظهر فيها المصطلح t_i .

وتبعاً للمثال السابق، فإذا افترضنا أن عدد الوثائق المخزنة في النظام تبلغ **1000000** وثيقة ويظهر مصطلح المكتبات في 1000 وثيقة من إجمالي عدد الوثائق وبالتالي يحسب $3 = \log(1000000/1000)$.
ويحسب معدل التردد العام للوثيقة من خلال حاصل ضرب تردد المصطلح X تردد الوثيقة المعكوس المعادلة الآتية:

$$d_t = TF(d, t) * IDF(t)$$

ومن خلال المثال السابق تكون المعادلة $0.12 = 0.04 * 3$ أي أن رتبة الوثيقة يساوي 0.12، ولعل من الملاحظ أن إجمالي القيم ستأتي منحصرة بين رقمي 1 و 0.

$$TF(d, t) = \begin{cases} 0 & \text{if } n(d, t) = 0 \\ 1 + \log(1 + \log(n(d, t))) & \text{otherwise} \end{cases}$$

وعليه يحسب جيب الزاوية الخاصة بالتشابه بين الوثيقة والاستفسار من خلال المعادلة الآتية:

$$sim(d_j, q) = \frac{d_j \cdot q}{\|d_j\| \|q\|} = \frac{\sum_{i=1}^N w_{i,j} * w_{i,q}}{\sqrt{\sum_{i=1}^N w_{i,j}^2} * \sqrt{\sum_{i=1}^N w_{i,q}^2}}$$

-2 خوارزمية C4.5 and beyond:

تأتي هذه الخوارزمية كأحد أشهر وأهم الخوارزميات المستخدمة في بناء نظم التصنيفات هذه النظم تعتمد مجموعة من الحالات Cases كمدخل لها، حيث أن كل حالة تنتمي الى عدد قليل الفئات وتوصف من خلال قيمها لتعمل هذه الخوارزمية على بناء هيكل أو شجرة قرارات تساهم على التنبؤ بتحديد وتسكين أي من الحالات قد ينتمي لفئة من الفئات، طورت هذه الخوارزمية على يد روس كوينلان Ross Quinlan.

تقوم منهجية هذه الخوارزمية على افتراض أن البيانات تعيين في مجموعة حيث $S = s_1, s_2, \dots$ والتي تصنف بالفعل في عينات Sample بحيث تشتمل كل عينة على مجموعة من المحددات Vector لتأخذ الشكل الآتي $s_i = x_1, x_2, \dots$ حيث أن كل محدد x_1, x_2, \dots يمثل خواص عينة من العينات، تتمثل المرحلة التالية في تزويد البيانات بمحددات تعمل على تسكينها داخل الفئات التي تنتمي إليها.

3- خوارزمية The k-means algorithm:

تنتمي هذه الخوارزمية الى منهجيات التحليل العنقودي لمجموعات البيانات أذ تعمل على تقسيم مجموعات البيانات المحددة الى مجموعات عنقودية محددة تعرف ب K وقد مر تطوير هذه الخوارزمية بالعديد من المراحل لما يقرب من عشرة سنوات.

تعتمد هذه الخوارزمية توافر مجموعة من البيانات (x_1, x_2, \dots, x_n) حيث يتم إعطاء لكل بيان محدد ذات أبعاد d-dimensional real vector ليتم تسكين هذه المجموعات في عنقاويد K وفقا لانتمائها الموضوعي.

4- خوارزمية Naive Bayes:

وهي أحد الخوارزميات التصنيف الاحتمالي المعتمدة على نظرية تعرف بنظرية Bayes' theorem تعد هذه النظرية احد نظريات مجال الاحتمالات في علم الرياضيات والتي تعني بقياس العلاقة بين احتمالين شرطييين والذي عادة ما يعلوهما التناقض فيما بينهم، أن المنهجية الأساسية التي تستند عليها هذه الخوارزمية يتم في أن يتم توصيف الكيانات المختلفة حيث أن كل كيان ينتمي إلى فئة معروفة، وكل كيان لديه مجموعة من المحددات المعروفة بحيث تهدف هذه الخوارزمية الى بناء قاعدة تسمح بتحديد وتعيين الكيانات المستقبلية داخل المجموعات في ظل ما توافر من محدّدات ومتغيرات.

3- التنقيب على الويب Web Mining:

الدوافع وراء التوجه للتنقيب عن المحتوى على الويب:

اتضح من واقع دراسة البنية المعمارية لنظم التنقيب عن البيانات ضرورة توافر منصات الحفظ ومستودعات التخزين للبيانات كركيزة أساس لتلك النظم، فهي بمثابة الحاويات التي تستورد منها نظم التنقيب عن البيانات مدخلاتها وتقوم علي تقديم أوجه المعالجة المختلفة لاكتشاف واستخراج المعرفة منها ومن ثم تحقيق التكامل المعرفي، ولما كان الويب يمثل أكبر مستودعات للبيانات في العالم²⁹ سواء المهيكّل منها أو غير المهيكّل، كان

²⁹ Liu, B. (2007). *Web data mining exploring hyperlinks, contents, and usage data*. Berlin: Springer.

لزما على الباحثين والعاملين في مجال التنقيب عن البيانات التطرق بالبحث والدراسة لاستثمار هذا الكم بغية تحقيق التكامل المعرفي على صعيد محتواه وأنماط استخدامه وبنيته.

يرى تيم بيرنرز Tim Berners lee - مبتكر شبكة الويب - أن المشكلات والتحديات التي تواجه الويب في تحقيق التكامل المعرفي يكمن في طبيعة أدوات البحث (من أدلة موضوعية ومحركات بحث... وغيرها) ، فلقد صممت هذه الأدوات للإجابة على تساؤل واحد " ماهي الوثائق التي تشتمل على الكلمات أو الجمل الآتية" دون النظر إلى اعتبارات أخرى كالعلاقات بين المحتوى ذات الصلة ومصداقية وموثوقية المحتوى أو تكاملية المعرفة، وقد أوضح ذلك في قوله...

"... إذا استطاع محرك البحث على الويب أن يجمع في تقنياته بين محرك الاستدلال reasoning engine ومحرك البحث Search engine فإننا يمكن أن نحظى بالتكامل المعرفي على صعيد شبكة الويب...." ³⁰

"اننا لانزال في عصر ما قبل الويب".

"أن الشبكة الويب الآن بصورتها الحالية مفهومة لنا نحن البشر ولكنها بالنسبة إلى الحاسب الآلي عبارة عن صفحات ممثلة بالصفحة الواحد ولا تعني لها شيئا، إن ما نريده من مبادرات التكامل المعرفي أن يجعل التطبيقات والبرمجيات قادرة على فهم ماذا تعني محتويات الصفحة التي يعرضها وما هو موجود في الويب من معلومات يمكن من خلالها تحقيق التكامل المعرفي."

"لا يمكن أن ينكر أحد أن الويب قد أضفت سمة من التغيير الشامل في طرق اكتساب المعرفة والتواصل والنجاح على مختلف الأصعدة وأنه كان المحرك الأساس للاتجاه نحو اقتصاديات المعرفة ومجتمعات المعرفة بما قد ضمنه من قدرات وإمكانات كفلت نقل المحتوى بيسر واتحته بسهولة، ولكن لم تصل الويب من خلال هذا المحتوى أو به إلى مرحلة النضج فهي مازالت قائمة على تطبيقات ومستودعات للمحتوى منعزلة ومنفصلة تتمثل في كيانات مستقلة تعمل دون تكاملية في المحتوى، فضلا عما يعلوها من عجز في التعامل مع هذا المحتوى وفهمه وتحليله بناء على دلالاته ومضمونه".

ويرى ريكاردو بيزا Ricardo Baeza-Yates أن المشكلة الرئيسية في الوصول للتكامل المعرفي على الويب يكمن في قضية التنبؤ بتحديد أي من الوثائق قد يتسم بالصلة بعضها البعض وأي منها لا يتسم بالصلة³¹.

كذلك أوضح كلا من كارلو تاسو & بيتر بيرسيل كوف PETER BRUSILOVSKY و CARLO TASSO أن كافة التحديات التي تواجه الويب في الوصول الى التكامل المعرفي تدور في فلك عاملين أساسيين هما: 1- المعالجة اللغوية: حيث أن غالبية برمجيات الويب تفقر في معالجتها لمحتوى الويب على وجود أو

³⁰ Alesso, H. P., & Smith, C. F. (2006). Thinking on the Web: Berners-Lee, Go del, and Turing. Hoboken, N.J.: Wiley-Interscience.(p.67).

³¹ Yates, R., & Neto, B. (2011). Modern information retrieval: the concepts and technology behind search (Second ed.). New York: Addison Wesley.(p.11).

غياب الكلمات المفتاحية في النص دون أية محاولة لتحليل المحتوى أو تحديد المفاهيم المشار إليها في النص وهو السبب الرئيسي وراء انخفاض الدقة في عمليات البحث والاسترجاع فضلا عن الظواهر اللغوية الأخرى كالترادف والتجانس. و2- محدودية الآليات والخوارزميات: وتتجلى هذه المحدودية في عجز التطبيقات والبرمجيات من فهم المحتوى والمعلومات المقدمة نظراً لعدم وجود نسق عالمي للتشغيل المتبادل فضلا عن اتساع الفجوة بين الجانبين الأساس لشبكة الويب المتمثلان في المحتوى (وهو المضمون الذي تشمله صفحات المعلومات على شبكة الويب) والبرمجيات (المسئولة عن معالجة هذا المحتوى واسترجاعه)، وبالأحرى افتقار البرامج لمعالجة المحتوى. فهما يفتقران إلى التكامل لينسجا معا نسيج الويب³².

بينما أوضح ماركوف Markov أن تحديات التي تواجه الويب في الوصول الى التكامل المعرفي منبعا يعود إلى المحدودية الدلالية لشبكة الويب ذاتها، فصفحات الويب لا تحمل دلالة لمحتواها ولكن تحمل تنسيقا جيد وتمثيل عظيم للبيانات، أما الروابط فتكاد تنعدم دلالاتها على الويب والدلالة الوحيدة التي تحملها في أطرافها هي أن الموقع هذا يرتبط بالموقع ذاك دون أية تحديد لدلالة الربط أو نوع الارتباط³³.

ويرى كلا من فان هيرميلين و سنتكشميدت Van Harmelen & Stuckenschmidt أن التحديات التي تواجهها تكمن في افتقار الويب إلى النماذج المفاهيمية لمصادر المعلومات وعدم وضوح حدود وملامح الويب في ظل ديناميكيته المفرطة³⁴.

محدودية الويب في تحقيق التكامل المعرفي:

بالرغم من أنصاح بنية الويب المعمارية والتكوينية، إلا أن تحقيق التكامل المعرفي على صعيد محتواه يتسم الصعوبة للأسباب الآتية³⁵:

1- ضخامة حجم المحتوى المتاح على الويب:

فقد قدر حجم المصادر المعلوماتية المتاحة والقابلة للتكشيف على الويب في اغسطس عام 2000 بنحو 7 ملايين صفحة بعدد مستخدمين لها قدر بـ 500 مليون مستخدم، ليصل حجم الشبكة في اغسطس 2010 إلى نحو 7.74 مليار صفحة وبعدد روابط قدر بنحو 4 مليار رابط وبعدد مستخدمين قدر بنحو 2 مليار مستخدم³⁶ - فقد قامت الويب بالسماح

³² Peter Brusilovsky, Carlo Tasso, Preface to Special Issue on User Modeling for Web Information Retrieval, User Modeling and User-Adapted Interaction, v.14 n.2-3, p.147-157, June 2004.

³³Markov, Z., & Larose, D. T. (2007). Data mining the Web: uncovering patterns in Web content, structure, and usage. New York: Wiley-Interscience.

³⁴Stuckenschmidt, H., & Harmelen, F. v. (2005). Information sharing on the semantic web: New York: Springer.

³⁵ Berners-Lee, T. (n.d.). The Semantic Web: Scientific American. *Science News, Articles and Information | Scientific American*. Retrieved August 2, 2011, from <http://www.scientificamerican.com/article.cfm?id=the-semantic-web>

³⁶Kunder, M. d. (n.d.). WorldWideWebSize.com | The size of the World Wide Web (The Internet). Retrieved September 21, 2011, from <http://www.worldwidewebsite.com>.

لمختلف الأدوار والأفراد بإضافة المحتوى والتحليل والنقد، مما أثمر عن حجم هائل من المحتوى ذات التنوع اللغوي والموضوعي والشكلي، الأمر الذي أسفر عن تحديات غير مسبوقة على صعيد ضبطه وتنظيمه.

2- افتقار بنية شبكة الويب الداخلية إلى البنية السليمة لتمثيل المعلومات والمحتوى:

تلك البنية التي تقف وراء عجزها الكامل عن تحقيق التكامل المعرفي على صعيد محتواها سواء النصي أو ذو الوسائط المتعددة، ولعل المرجعية الأساس وراء هذا العجز والافتقار يعود إلى اعتماد شبكة الويب على لغات تمثيل ومواصفات يقتصر دورها فقط على التمثيل الأصم للمحتوى plain text، مع افتقار كامل إلى المعاني meaningful tags، والمؤشرات mark-up indicating، التي تيسر من معالجة المحتوى واستثماره وتحقيق التكامل على صعيده³⁷.

3- تعقد احتياجات المستفيدين والباحثين في التعامل مع محتوى الويب:

سواء كان ذلك على صعيد التأليف والتحرير أو على صعيد البحث والاسترجاع، فكثير من المستفيدين يقومون بطرح استفساراتهم بصورة قلما يعبروا من خلالها عن حاجاتهم البحثية الفعلية، نظرا لافتقارهم إلى الفهم الصحيح للموضوع قيد البحث، أو عدم الإلمام الشامل بمفردات حاجاتهم المعلوماتية، أو تقديم استفسارات أعم بكثير من الحاجة الفعلية إلى المعلومات، فضلا عن حجم الاستفسارات المقدمة إلى الويب والذي يتسم بالنمو المتزايد والمستمر فقد أوضحت إحدى الدراسات عن أن حجم الاستفسارات على الويب قد بلغت نحو 10 مليار استفسار عام 2008 موجه إلى 5 محركات بحث (Google, Yahoo!, Bing, AOL). وأن معدل توزيع تردد مصطلحات الاستفسار يتسم بالانحراف من إجمالي عمليات البحث فقد كشفت هذه الدراسة أن بعض المصطلحات تستخدم بشكل متكرر في مقابل الكثير من المصطلحات التي تستخدم مرة واحدة فقط، فقد تم الكشف عن أن 63 مصطلح حظي بتردد ظهور بلغ أكثر من 100 مرة، في حين كون هذه المصطلحات تمثل أقل من 1% من إجمالي المصطلحات. مما يظهر طبيعة أن البحث على الويب يمكن وصفه بأنه بحث بمصطلحات تتسم بنسبة منخفضة في تردد الظهور مقابل مصطلحات تتسم بتردد عالي في الظهور³⁸.

4- الغموض المعلوماتي لمحتوى الويب Ambiguity of information:

وهي تلك الظاهرة التي أرسنها بنية روابط شبكة الويب الحالية والناجئة عن ضعف توصيف وتحديد طبيعة الروابط والعلاقات بين المصادر المتاحة على الويب، فما تشير إليه مواقع الويب في شبكة الويب الحالية هو ان الموقع (أ) يرتبط بالموقع (ب) دون وجود لدلالة أو توصيف لطبيعة هذه العلاقة فهل هي علاقة ابوة وبنوة أم علاقة اشتغال أم اكتمال أو غيرها من انماط العلاقات بين الكيانات، وبمعنى آخر يمكن القول بأن الروابط في الويب هي روابط صماء.

5- صبغ عملية إنتاج المحتوى على الويب بصبغة الهدف Target:

³⁷Lim, E.P., Sun, A. (2005) : Web Mining - The Ontology Approach. In: Proceedings of The International Advanced Digital Library Conference (IADLC 2005), Nagoya, Japan (August 2005) available at: [http://iadlc.nul.nagoya-u.ac.jp/archives/IADLCpresen/Lim.pdf](http://iadlc.nul.nagoya-
http://iadlc.nul.nagoya-u.ac.jp/archives/IADLCpresen/Lim.pdf) date: 7/3/2012.

³⁸Ibid .

- فإننتاج المحتوى على الويب يتسم بطابع الهدف Target بمعنى أن المحتوى ينقسم إلى فئتين وفقا للهدف المرجو منه:
3. الفئة الأولى: يتمثل في المعلومات التي يتم انتاجها اساسا للاستخدام من قبل البشر والمستخدمين، المتمثلة في الرسالة الفكرية لدى المؤلفين والتي يقوم بنقلها من خلال الويب.
 4. الفئة الثانية: فيتمثل في المحتوى الذي ينتج بهدف أن يستخدم من قبل البرمجيات واجهزة الحاسب ولا يمكن للبشر أن يستفيدوا منه، ويتمثل في بروتوكولات الاتصال ولغات البرمجة البيئية ونظم ادارة قواعد البيانات.
- هذا الأمر الذي أسفر عن عجز التطبيقات والبرمجيات من فهم المحتوى والمعلومات المقدمة نظراً لعدم وجود نسق عالمي للتشغيل المتبادل³⁹
- 6- تحمل اللغات الطبيعية المعبر بها عن محتوى الويب الكثير من القضايا الشائكة:
التي تتعلق بالترادف والأضداد والتروية والجناس اللفظي، تلك القضايا التي تنسحب بطبيعة الحال على مجال المعالجة الآلية للمحتوى من جانب، وعلى تفاعل المستخدمين مع أنظمة البحث والاسترجاع من جانب آخر في ظل الاعتماد على التعبير اللفظي والنصي عن الحاجة المعلوماتية.
 - 7- محدودية المعالجة الذكية لمحتوى الويب:
فالمعالجة الآلية لمحتوى الويب تعتمد بالأساس على المعالجة المفردة للكلمات الأمر الذي يعرف بالمعالجة المعجمية واللغوية للمحتوى Lexicon Handling دون ان تعمل على المعالجة وفقاً للمفاهيم والدلالات والسياقات الواردة بها كلمات النص الأمر الذي يصعب القدرة على تحقيق التكامل المعرفي المعتمد في الأساس على معالجة المعاني وفقاً لفهم المضمون والسياق.
 - 8- صعوبة تحديد المصداقية والموثوقية لمحتوى الويب:
فكثير من محتوى الويب يعطوه سمة خلط بين الحقيقي وما هو زيف، ولا يمكن تحديد درجة أو معيار للموثوقية للمحتوى. على الرغم من كون الفكرة الأساسية للتكامل المعرفي هو القدرة تبادل المعلومات والتكامل البيئي مطمئنة في ذلك الى مصدقيتها.
 - 9- التعامل الذكي.
تقف الويب بمعماريتها وتقنياتها ثابتة أما ابواب التعامل الذكي مع المعلومات والمحتوى فليس لديها ملكة الذكاء في تحليله بناء على مضمونه وكذلك تتلاشى قدرات الربط بين المعلومات واستنباط نتائج جديدة من واقع مما هو متاح.

³⁹Sanjib kumar (march 2009), "TOWARDS SEMANTIC WEB BASED SEARCH ENGINES" National Conference on "Advances in Computer Networks & Information Technology (NCACNIT-09) March 24-25, available at http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=5974163 at date: 7/2/2012.

10- المحتوى المغلق.

على الرغم من أن الويب قد اتسمت بانها قاتل للتطبيقات Application killer بما قامت بتوفيره من برامج الويكي والمدونات وغيرها من التطبيقات التي كفلت التحرير الحي للمحتوى على الويب، إلا أن هذا المحتوى عادة ما ينشر هذا المحتوى على الويب في صورة اقرب للصومعات مما يعيق عملية اعادة استخدامها أو تدويرها من قبل برمجيات اخرى على الويب.

11- افتقار الدلالة في مخططات ونماذج الوصف:

فالويب تشتمل على نماذج توصيفية (كالميتاداتا) تعتمد على وصف المصادر من خلال مخططات (كالدبلن كور) وتفتقر هذه المخططات والنماذج إلى الدلالة والمعاني، فيقتصر دورها فقد على توصيف المحتوى دون عبارات حرة، بصورة قد لاتشير بالشكل الكامل إلى محتوى الموقع أو الصفحة فضلا عن عجزها في امكانية تحقيق التشغيل المتبادل.

12- ديناميكية البيانات كتحد لمركات البحث Data Dynamic:

تعد احد السمات التي يتميز بها محتوى العنكبوتية انه محتوى ديناميكي الوجود، بسبب ديناميكية العنكبوتية والانترنت حيث ترتفع معدلات التعديل في محتوى العنكبوتية، ان المرجعية الاساس وراء ديناميكية محتوى العنكبوتية يعود إلى 3 اسباب:

- المحتوى ذات الطبيعة الحساسة للوقت: كالمواقع التي تشتمل على اسعر الاسهم وعاوين الاخبار.
- المحتوى المولد وفقا بصورة مخصصة لمستفيد ما: وتتمثل فيما يعرف بالمواقع ذات سمات الخاصة لمستفيد ما personalization من (حيث الشكل والسمات والبنية والمحتوى) لتتناسب مستخدم ما.
- المحتوى المولد بناء على المدخلات: ويتمثل في المواقع التي تعتمد على استقبال مدخلات من قبل المستفيد كشاشات قواعد البيانات.

الجانب الآخر هو معدلات التعديل سواء كان بالحذف او الاضافة وقد قدر بنحو 80% يوميا، وفي دراسة قام بها Ntous & Olson عن تحديث صفحات العنكبوتية وجدا أن 320 مليون صفحة جديدة تضاف اسبوعيا كما ان 20% من صفحات العنكبوتية اليوم سوف تختفي خلال عام واحد . كما ان 50% من محتوى تلك الصفحات سوف يتغير خلال نفس الفترة. ولا يقتصر الامر على المحتوى فحسب بل تمتد هذه الديناميكية لتشمل هيكله البيانات وقوابها والتي اتسمت بعدم الاستقرار والديناميكية والمرجعية في ذلك تعود إلى وفرة برامج ادارة المحتوى الرقمي على العنكبوتية، يمكن اجمال ديناميكية العنكبوتية من خلال المؤشرات الاتية:

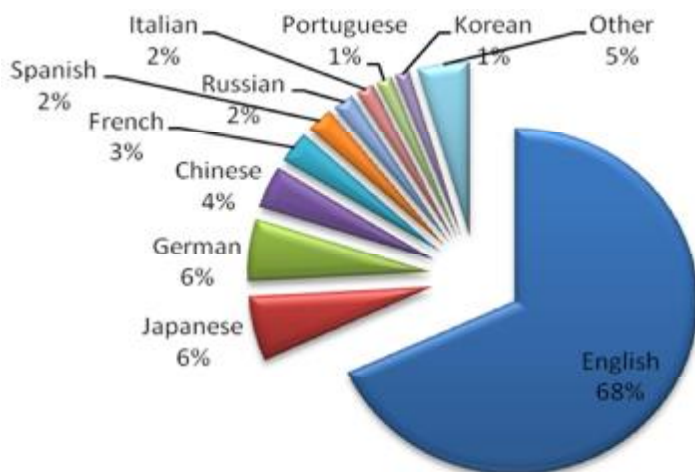
- لايزيد عمر 50% من اجمالي صفحات العنكبوتية عن 100 يوم بينما نسبة 25% من صفحات العنكبوتية تزيد عمرها عن سنة واحدة.

- 40% من الصفحات ذات النطاق .com تتغير كل يوم، بينما المتوسط العمري للصفحات ذات النطاق .gov و .edu لا يزيد عن أربعة أشهر.
- متوسط عمر محتوى العنكبوتية عامان.
- متوسط عمر المحددات الفريدة للمواقع لا يتجاوز 4 سنوات.
- متوسط عمر bookmarks لا يتجاوز الشهران⁴⁰.

13- التباين اللغوي:

فضلا عن التباين في التكويد وصيغ التكويد داخل الويب تتفاقم ظاهرة أخرى تحول بين التكامل المعرفي لمحتوى الويب وهي تعدد لغات المحتوى بين العربية والعبرية والفرنسية والصينية واليابانية وغيرها من اللغات الحية المنطوقة. الأمر الذي استتبع أن تقوم برمجيات معالجة المحتوى بالتركيز على لغة محتوى بعينه دون الآخر لتظهر بعد ذلك طوائف البرمجيات وفقا للغات ووفقا للنطاق الجغرافي للتركيز على المحتوى في لغات محددة وفي مناطق جغرافية محددة،⁴¹

أما عن واقع اللغة العربية على الويب فيتمثل في حجم المحتوى العربي المتاح والذي يشغل نسبة 0.2% من حجم المحتوى المتاح على الانترنت بواقع 100 مليون صفحة⁴².



شكل رقم (5) يوضح واقع حجم لغات المحتوى على الويب.

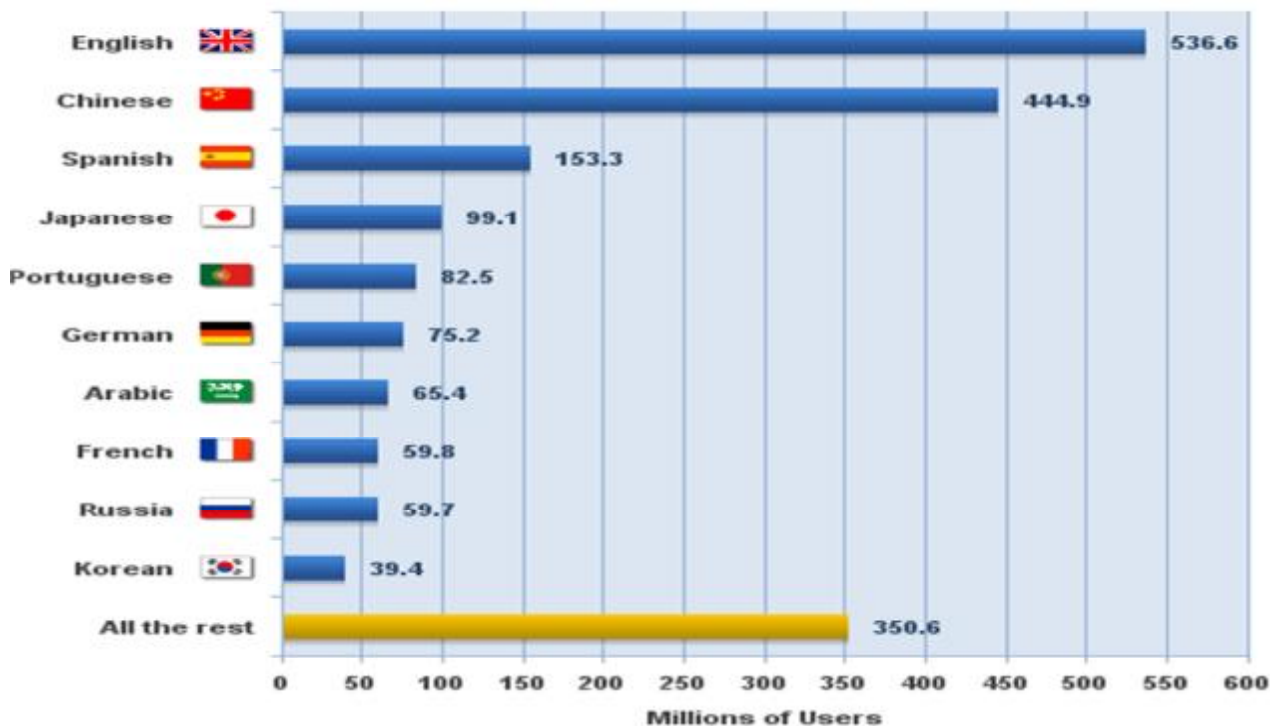
⁴⁰ Terrence A. Brooks. Web search: how the Web has changed information retrieval. Information Research, 8(3):(paper no. 154), April 2003.

⁴¹ Andrew Hammond (2004). Arabic search engine may boost content available at . <http://www.abc.net.au>.

⁴² <http://www.translate-to-success.com/online-language-web-site-content.html>

اللغة	جمالي نسبة المحتوى
English	34.40%
Japanese	9.90%
German	8.80%
Chinese	9.90%
French	10.00%
Spanish	14.40%
Russian	9.90%
Italian	6.60%
Portuguese	4.40%
Arabic	3.30%
Abic	6.60%

جدول رقم (1) يوضح حجم المحتوى المصاغ باللغة العربية على صعيد العالم.⁴³



شكل رقم (6) يوضح ترتيب اللغة العربية من حيث عدد مستخدميها على الانترنت⁴⁴

⁴³ Top Ten Internet Languages - World Internet Statistics. (n.d.). *Internet World Stats - Usage and Population Statistics*. Retrieved July 20, 2011, from <http://www.internetworldstats.com/stats7.htm>

⁴⁴ Ibid

14- الروابط الصماء:

على الرغم من ان بنية الويب الحالية تمتاز بالروابط بين نصوصها ومحتواها إلا أنها تعد من أكثر المشكلات في الويب، فالروابط داخل بنية الويب لاتحدد طبيعة أو نمط الارتباط بين المصادر والكيانات فهي تشير فقط أن الموقع (أ) يرتبط بالموقع (ب) دون تحديد لنوع الربط ودرجته أو الموثوقية فيه، فوفقا لتعريف Tim Lee للنص الفائق فان أي شئ يمكن أن يرتبط بأي شئ وبالتالي ووفقا لهذا المفهوم يمكن أن ترتبط معلومات أكاديمية بمعلومات تجارية⁴⁵.

وعلى هذا يتضح أن شبكة الويب الحالية تواجه العديد من القضايا الشائكة والمشكلات المتفاقمة التي كانت دافعا للتوجه الى التنقيب الى البيانات Web Data Mining بغية التكامل المعرفي.

⁴⁵ Berners-Lee, T. (n.d.). The Semantic Web: Scientific American. *Science News, Articles and Information | Scientific American*. Retrieved August 2, 2011, from <http://www.scientificamerican.com/article.cfm?id=the-semantic-web>