

C. STATO DELL'ARTE

A partire dal problema che cerchiamo di risolvere, è importante iniziare a calarsi nel concreto, rendendosi conto a quali campi ispirarsi per costruire un'architettura funzionale alla soluzione della questione.

Web semantico: con questo termine, coniato dal suo ideatore Tim Berners-Lee, si intende la trasformazione del World Wide Web in un ambiente dove i documenti pubblicati (pagine HTML, file, immagini, e così via) sono associati ad informazioni e dati che ne specificano il contesto semantico in un formato adatto all'interrogazione e l'interpretazione (es. tramite motori di ricerca) e, più in generale, all'elaborazione automatica. [Wikipedia] Questo non significa limitare le numerosissime forme in cui è possibile presentare i contenuti sul web, bensì un modo per fornire una struttura capace di rispondere a delle query, permettendo la lettura automatica di pagine e documenti. Risorse con questa caratteristica sono enormemente utili per la facilità con la quale è possibile estrarre informazioni e adattarle alle proprie necessità; sono quindi queste la scelta più ovvia come semplice punto di partenza per la ricerca di contenuti atti a stimolare la reminiscenza. Una delle realtà più grosse che opera nel web semantico è DBpedia, il cui approccio è quello di parsare periodicamente i dump rilasciati da Wikipedia, estrarre le uniche informazioni strutturate delle pagine, le cosiddette infobox, unendo alcuni attributi delle versioni localizzate dello stesso articolo e rappresentando le risorse tramite Resource Description Framework (RDF). Ogni risorsa è definita tramite un URI, un identificatore unico dell'entità presentata. Il lavoro di DBpedia è eccellente, ma è involontariamente ostacolato dagli editori di Wikipedia, che per la compilazione delle infobox non hanno dei template ben definiti, finendo per rappresentare lo stesso attributo in articoli diversi con nomi diversi; DBpedia è interrogabile tramite SPARQL, un linguaggio SQL-like costruito per leggere RDF. Dopo il rilascio del primo dataset al pubblico nel 2007, la situazione attuale è quella di un grosso grafo di altri set di dati collegati tra di loro tramite l'accoppiamento di risorse appartenenti a insiemi diversi, ma rappresentanti la stessa entità. Pubblicare questi Linked Data rende molto più facile la ricerca di informazioni sul web, rendendole molto più precise e complete. Alcuni altri progetti di web semantico sono: Freebase, che a differenza di DBpedia è un progetto proprietario e orientato al profitto; Uniprot, una knowledge base contenente dati liberamente accessibili sulle sequenze di proteine; GeoNames, un database geografico contenente oltre 10 milioni di luoghi, accessibile e scaricabile sotto licenza Creative Commons. A settembre 2011, lo stato dei collegamenti tra datasets semantici è rappresentato dal grafo visualizzabile a (<http://lod-cloud.net/versions/2011-09-19/lod-cloud.html>). Un approccio diverso nell'attribuzione di un significato ai contenuti presenti sul web è quello dei microformati, un'estensione di markup che, tramite l'utilizzo degli attributi HTML class, rel e rev, consente l'attribuzione di regole semantiche a normali pagine web: date le informazioni sul contatto di Nicola Parrello

```
<div>
  <div>Nicola Parrello</div>
  <div>University of Trento</div>
  <div>123-123123</div>
```

```
<a href="mailto:hi@everybody.com">Email</a>
</div>
```

queste, se riscritte con il microformato hCard (specifico per i contatti), diventano

```
<div class="vcard">
  <div class="fn">Nicola Parrello</div>
  <div class="org">University of Trento</div>
  <div class="tel">123-123123</div>
  <a class="email" href="mailto:hi@everybody.com">Email</a>
</div>
```

Attraverso questi microformati, ad esempio, un software come un browser può estrarre facilmente informazioni e navigare le relazioni tra oggetti diversi, mantenendo comunque la normale leggibilità di una pagina web. Oltre ad hCard, solo hCalendar è stato formalizzato: altri microformati, come hAtom, hMedia e hNews (rispettivamente per feed Atom, contenuti multimediali e notizie) sono solo più o meno abbozzati, e non rappresentano quindi uno standard.

Web search e contestualizzazione dell'informazione: Il passo successivo nella ricerca di una soluzione al problema che è il tema di questo documento, dopo aver scelto da dove reperire i dati, è quello di decidere come permettere a un eventuale utente di accedere ai contenuti che vogliamo proporre. L'approccio che sembra più naturale è quello di costruire un piccolo motore di ricerca, in modo tale da creare una maschera che permetta una richiesta di risorse in maniera uniforme, anche se da fonti diverse; inoltre, per avere una ricerca rapida e reattiva, il buon senso suggerisce che è una buona idea quella di indicizzare il materiale ricercabile. In questo campo è difficile non considerare Google come lo stato dell'arte, essendo loro il search engine più utilizzato al mondo. Per gli utenti da ogni parte del globo possano utilizzare un servizio veloce ed efficiente, Google utilizza sette componenti dinamiche, che salvano e leggono dati in altre strutture; il procedimento può essere riassunto così: Il Server URL parte da un URL e, leggendo il Document Index, invia gli URL ai Crawler. I Crawler scaricano le pagine web e le mandano nello Store Server. Lo Store Server comprime le pagine nel formato zlib (RFC1950), riducendo la loro dimensione a un terzo dell'originale, che vengono poi di un docId univoco e immagazzinate nel repository. Il repository viene letto dall'Indexer, che decompime i documenti e li parse, assegnando a ogni parola (a cui viene assegnato un wordId univoco) informazioni su posizione, grandezza del testo, numero di occorrenze e altro; ogni voce viene poi aggiunta ad un indice parzialmente ordinato. Dalle pagine lette, l'Indexer estrae anche i dati dei link in esse contenuti, oltre a servirsi del testo analizzato per costruire un lessico, utile per la vera e propria funzione di ricerca. Lo URL Resolver combina i documenti con i dati dei link, per costruire tabelle di coppie di docId, utilizzati per calcolare il PageRank. Il Sorter riordina l'indice (ordinato per docId) in un indice inverso (ordinato per wordId), aggiungendo altri dati per rendere la ricerca più precisa. Infine viene calcolato il PageRank, una misura dell'importanza di una pagina web calcolata in base a quali e quante sono le pagine che hanno dei link che puntano a essa.

Negli anni Google Ã diventato colmo di pubblicitÃ, mostrata e scelta in base ai dati raccolti dalle ricerche degli utenti, ma certamente non Ã lâunica possibilitÃ in quanto a motori di ricerca: un esempio su tutti e quello di Duck Duck Go che, a differenza di Google, Ã un Semantic Search Engine. Il sistema quindi si occupa di valutare lâeffettivo significato dei termini di ricerca, eliminando in maniera piÃ precisa i risultati irrilevanti. Ma quello che lo contraddistingue in maniera maggiore dagli altri Ã la sua attenzione alla privacy: Duck Duck Go infatti non conserva e non vende a terzi nessuna informazione sulle ricerche, permettendo una ricerca anonima e al sicuro da data leak, richieste legali da parte delle istituzioni e disonesti. (<https://duckduckgo.com/privacy>) Quello che puÃ infastidire gli utenti, cioÃ gli annunci di cui sopra, sono perÃ molto interessanti per la nostra ricerca, in quanto introducono un altro aspetto molto importante: perchÃ la ricerca sia semplice da utilizzare, e soprattutto piacevole, puÃ essere necessario renderla per cosÃ dire automatica, nascondendo lâazione manuale dellâinserimento di parametri e mostrando i risultati direttamente come contesto di qualche altra azione. Premiata nel 2012 dal magazine Popular Science come innovazione dellâanno (<http://www.popsci.com/bown/2012/product/google-now>), Google Now Ã unâestensione dellâapplicazione mobile Google Search che, oltre ad essere un assistente vocale, analizza abitudini, ricerche e posizioni ricorrenti per fornire dati e risultati contestuali, mostrando allâutente le informazioni prima che egli ne faccia richiesta esplicita allâapp.

Algoritmi spazio-temporali: Abbiamo definito la raccolta e la ricerca, ma un ultimo punto rimane fumoso: in base a quale criterio scegliere un set di risorse rispetto a un altro, con la condizione che queste siano effettivamente rilevanti ed efficaci nella stimolazione della reminiscenza? Il vissuto di un individuo puÃ essere riassunto in una lista di eventi, di storie di vita; allo stesso tempo, un evento puÃ essere identificato da una coppia $\langle t, s \rangle$, dove t Ã la coordinata temporale e s descrive la posizione nello spazio. Utilizzando questo modello, il criterio di cui sopra diventa quello della distanza spazio-temporale tra due entitÃ, e il calcolo di questa distanza diventa il mezzo per effettuare unâindicizzazione preliminare delle risorse.

Un interessante esempio di sfruttamento estensivo delle informazione su tempo e spazio al fine di restituire dei risultati a una query Ã TimeTrails. TimeTrails Ã un sistema per lâestrazione e lâesplorazione di coordinate spazio-temporali che si possono trovare nei documenti di testo, composto di tre componenti principali: (1) una pipeline che, dopo aver letto i dati forniti da dei moduli che estraggono documenti di testo da varie sorgenti (e.g. la vetrina di Wikipedia (<http://it.wikipedia.org/wiki/Wikipedia:Vetrina>)), si occupa di estrarre le date e i luoghi contenuti nel testo, normalizzarli (e.g. da â25 luglio 1991â a â25/07/1991â per le prime e da âIPergine Valsuganaâ a â46.06853620, 11.23528960â per i luoghi), e calcolare il numero e la posizione delle occorrenze trovate, in modo da verificare se la coordinata dello spazio e quella del tempo identificano un evento ben preciso oppure se sono due riferimenti senza alcuna correlazione. Il risultato dellâelaborazione, cioÃ il documento originale unito a una sequenza ordinata di tuple $\langle t, s \rangle$, viene poi salvato in (2) un database ottimizzato per

contenere dati su luoghi (<http://postgis.net/>), dal quale (3) unâinterfaccia permette la ricerca testuale di documenti, con un risultato che verrÃ mostrato su una mappa sotto forma di traiettoria, rappresentata da tutte le tuple lette dal db. Nel caso la query dovesse restituire piÃ di un documento, puÃ risultare interessante vedere dove e quando le traiettorie si incrociano: se i documenti sono biografie, lâintersezione delle traiettorie in uno dei punti segnati dalle tuple $\langle t, s \rangle$ puÃ indicare che i personaggi si sono incontrati.