

RICERCA SEMANTICA MULTIDOMINIO SUL WEB

By

Gianluca Di Fiore

SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE OF
MASTER OF SCIENCE

AT

DALHOUSIE UNIVERSITY
HALIFAX, NOVA SCOTIA

JANUARY 2003

© Copyright by Gianluca Di Fiore, 2003

DALHOUSIE UNIVERSITY
DEPARTMENT OF
MATHEMATICS STATISTICS AND COMPUTING SCIENCE

The undersigned hereby certify that they have read and recommend to the Faculty of Graduate Studies for acceptance a thesis entitled “**Ricerca semantica multidominio sul WEB**” by **Gianluca Di Fiore** in partial fulfillment of the requirements for the degree of **Master of Science**.

Dated: January 2003

Supervisor:

Readers:

DALHOUSIE UNIVERSITY

Date: **January 2003**

Author: **Gianluca Di Fiore**

Title: **Ricerca semantica multidominio sul WEB**

Department: **Mathematics Statistics and Computing Science**

Degree: **M.Sc.** Convocation: **May** Year: **2003**

Permission is herewith granted to Dalhousie University to circulate and to have copied for non-commercial purposes, at its discretion, the above title upon the request of individuals or institutions.

Signature of Author

THE AUTHOR RESERVES OTHER PUBLICATION RIGHTS, AND NEITHER THE THESIS NOR EXTENSIVE EXTRACTS FROM IT MAY BE PRINTED OR OTHERWISE REPRODUCED WITHOUT THE AUTHOR'S WRITTEN PERMISSION.

THE AUTHOR ATTESTS THAT PERMISSION HAS BEEN OBTAINED FOR THE USE OF ANY COPYRIGHTED MATERIAL APPEARING IN THIS THESIS (OTHER THAN BRIEF EXCERPTS REQUIRING ONLY PROPER ACKNOWLEDGEMENT IN SCHOLARLY WRITING) AND THAT ALL SUCH USE IS CLEARLY ACKNOWLEDGED.

A mia madre e Daria

Table of Contents

Table of Contents	v
List of Tables	viii
List of Figures	ix
Introduzione	1
0.1 Organizzazione della tesi	3
1 Stato dell'arte	5
1.1 Una definizione di Web Mining	5
1.2 Una tassonomia per il Web Mining	6
1.3 Web content mining	8
1.3.1 Approccio Agent-Based	8
1.3.2 Database Approach	9
1.4 Web Usage Mining	10
1.4.1 Pattern discovery tools	11
1.4.2 Pattern analysis tools	11
1.5 L'architettura del sistema WebTMS	12
1.5.1 Gathering agent	13
1.5.2 Preprocessing Agent	13
1.5.3 Categorization Agent	14
1.5.4 Clustering Agent	15
1.5.5 MDDAE (Multi-Dimension Document Analysis Engine)	15
1.5.6 User Interface Agent (UIA)	16
1.6 Il meta-motore GSA	16
1.6.1 L'architettura di GSA	17
1.6.2 Il ranking delle pagine	18
1.6.3 Come GSA impara dall'utente	21

1.7	Il motore di ricerca Mondou, basato sul Text Mining	22
1.7.1	L'architettura di Mondou	23
1.7.2	Il data mining nel motore di ricerca Mondou	24
1.8	Il progetto Harvest	28
2	Approccio teorico	29
2.1	Gli Agenti Software	29
2.1.1	Tipologie di Agenti	30
2.2	Ontologia	58
2.2.1	Semantic Network	59
2.3	Modello formale	65
2.3.1	Domini, concetti e parole	66
2.3.2	Dominio unico	71
2.3.3	Domini indipendenti distinti	77
2.3.4	Domini distinti non indipendenti	79
2.4	La metrica	81
2.4.1	GPrS e posizione mediata	81
2.4.2	Contributo sintattico	84
2.4.3	Contributo semantico	86
2.4.4	Peso dei vocaboli di una pagina	86
2.4.5	Voto complessivo	92
3	Il sistema	93
3.1	Intelligent Search Agent	93
3.1.1	La struttura ad agenti	101
3.2	Architettura	109
3.2.1	User Interface	111
3.2.2	Search Engine Wrapper	113
3.2.3	Web Spider	115
3.2.4	Document Preprocessor	115
3.2.5	Semantic Knowledge Base	117
3.2.6	Miner	120
3.3	Struttura del Web Repository	121
3.4	Struttura della Semantic Knowledge Base	122
4	Sperimentazione	124
4.1	Parametri del sistema	124
4.2	Ipotesi preliminari	126
4.3	La prove effettuate	128

4.3.1	Query: Zuccherò	128
4.3.2	Query: Morandi	134
4.3.3	Query: Bongiorno	139
5	Interfaccia utente del sistema SSA	140
5.1	Gestione Motori	141
5.2	Nuova Query	145
5.3	Carica Query	149
5.4	Gestione Rete Semantica	150
5.4.1	Gestione Concetti	152
5.4.2	Gestione Domini	153
5.4.3	Gestione Rete	154
5.5	Gestione Dizionario	157
5.6	Impostazioni	157
	Bibliography	160

List of Tables

1.1	Esempi di <i>associative keywords</i>	25
1.2	Le dimensioni degli spazi delle keywords	26
4.1	Link restituiti e link giudicati rilevanti nella query Zuccherò nel dominio Musica	129
4.2	Link restituiti e link giudicati rilevanti nella query Morandi nel dominio Musica	134
4.3	Link restituiti e link giudicati rilevanti nella query Bongiorno nel dominio Televisione	139

List of Figures

1.1	Tassonomia per il Web Mining	7
1.2	L'architettura del sistema WebTMS	12
1.3	Il Gathering Agent	13
1.4	Il Preprocessing Agent	14
1.5	Il Categorization Agent	14
1.6	Il Clustering Agent	15
1.7	Un esempio di Document Cube	16
1.8	L'architettura del sistema GSA	17
1.9	L'architettura del sistema Mondou	23
1.10	Relazioni tra attributi differenti	27
2.1	Una tipologia di agenti	31
2.2	Classificazione di Agenti Software	33
2.3	Architettura Pleiades	36
2.4	Interface Agents	38
2.5	Architettura Telescript	44
2.6	Information Agent	46
2.7	Brooks Subsumption Architecture	49
2.8	The InteRRap Hybrid Architecture	53
2.9	Sistema ad agenti eterogenei	56
2.10	Legame tra due concetti	60
2.11	Distanza probabilistica tra due concetti	61

2.12	Una rete semantica per i concetti del dominio musicale	62
2.13	Distanza probabilistica lungo un percorso	63
2.14	Legami tra alcuni concetti del dominio automobilistico	64
2.15	Esempio di insieme monoconcetto	71
2.16	Esempio di insieme multiconcetto	72
2.17	Esempio di mappa dei domini	80
2.18	Esempio di risultato di ricerca con Google	82
3.1	Architettura del sistema	110
3.2	Menu Iniziale	111
3.3	Search Engine Wrapper	112
3.4	Esempio di Query	113
3.5	Web Spider	114
3.6	Document Preprocessor	116
3.7	Semantic Knowledge Base	118
3.8	Esempio di Query	119
3.9	Modello E-R del Web Repository	121
3.10	Modello E-R del Semantic Knowledge Base	123
4.1	Query sottoposta al sistema: zucchero	129
4.2	Risultato del sistema alla query: zucchero	130
4.3	Risultato di Google alla query: zucchero	131
4.4	Risultato di Altavista alla query: zucchero	132
4.5	Risultato di Lycos alla query: zucchero	133
4.6	Risultato del sistema alla query: morandi	135
4.7	Risultato di Google alla query: morandi	136
4.8	Risultato di Altavista alla query: morandi	137
4.9	Risultato di Lycos alla query: morandi	138
5.1	Interfaccia principale del sistema SSA	140
5.2	Finestra di dialogo per la gestione dei motori di ricerca	142

5.3	Interfaccia Nuova Query	144
5.4	Esempio di link di livello 0	146
5.5	Esempio di link di livello 1	146
5.6	Interfaccia Carica Query	148
5.7	Interfaccia della Gestione della Rete Semantica	151
5.8	Finestra di dialogo per la gestione dei concetti	152
5.9	Finestra di dialogo per la gestione dei domini	153
5.10	Finestra di dialogo per la gestione della Semantic Network	155
5.11	Finestra di dialogo per la gestione del Dizionario contenente le parole conosciute dal sistema	156
5.12	Finestra di dialogo la modifica delle impostazioni	158

Introduzione

Il *World Wide Web* è diventato progressivamente una vasta risorsa di informazioni; a causa della complessità nell'organizzazione dei dati e della quantità di informazioni presenti, la ricerca sul WEB di informazioni davvero utili per un utente è diventata decisamente complessa: anche se un motore di ricerca Web assiste nella scoperta di risorse, esso è lontano dal soddisfare, per la sua scarsa precisione, le richieste di una ricerca mirata. Infatti trovare le informazioni richieste nei risultati di un motore di ricerca si rivela fruttuoso solo in presenza di argomenti di una certa notorietà e importanza, e di query molto precise; negli altri casi questo lavoro può implicare una considerevole perdita di tempo che un utente può pagare in termini di costi di collegamento alla rete.

Questo avviene perché i motori di ricerca tradizionali effettuano ricerche di tipo **sintattico**: essi restituiscono le pagine che contengono le keywords presenti nelle query degli utenti, indipendentemente dal contesto in cui esse vengono utilizzate. Se ciò da un lato è conveniente in termini di velocità di reperimento delle pagine e restituzione dei risultati, dall'altro lato porta spesso a risultati errati o imprecisi, dato

che vengono restituite molte pagine non attinenti al contesto della query dell'utente.

Ad esempio, un utente che voglia cercare pagine relative alla pesca, inserendo in un motore come keyword la parola "pesca", oltre a trovare pagine attinenti al dominio ittico troverà sicuramente pagine inerenti al dominio ortofrutticolo. Sarebbe conveniente avere a disposizione un sistema in grado di capire di cosa parla una pagina valutando la sua attinenza con i domini di interesse per l'utente. Una ricerca di tal tipo è detta **semantica** in quanto non restituisce semplicemente pagine che contengono le keywords, ma pagine che hanno anche un contenuto semantico aderente al dominio desiderato dall'utente.

Un problema analogo, relativo però alle informazioni contenute nei database, si ha con il *Data Mining*; infatti esso è usato per identificare modelli validi, nuovi, potenzialmente utili e comprensibili, da una raccolta di dati in una collettività di database.

Il lavoro che invece si interessa delle informazioni non strutturate ed eterogenee presenti sul WEB viene spesso identificato con il nome di **Web Mining**. Anche se non c'è un accordo generale sulla definizione di Web Mining possiamo dire, in analogia con la definizione di Data Mining, che esso *rappresenta l'attività di identificazione di modelli impliciti in una grande raccolta di documenti presenti sul WEB*.

Il progetto sviluppato si colloca nella teoria del Web Mining e si basa sugli **Intelligent Search Agents (ISA)**. Questi sono programmi usati per filtrare le informazioni

ritrovate sul WEB in base agli interessi dell'utente, memorizzandone le scelte e costruendone il profilo. Ad essi si stanno interessando molte tra le più importanti aziende del settore informatico, da Netscape a Microsoft, da IBM a AT&T.

Tale lavoro si prefigge lo scopo di estendere un prodotto già realizzato precedentemente, che valuta l'attinenza delle pagine reperite ad un dominio applicativo prefissato, al fine di renderlo multidominio, cioè capace di apprendere nuovi domini e valutare l'attinenza delle pagine a tutti i domini conosciuti, considerando anche eventuali gradi di correlazione tra i domini. Il sistema, come il prodotto precedente, consente di reperire pagine Web interfacciandosi con i motori di ricerca più utilizzati dalla comunità Web ma consente di aggiungere nuovi motori di ricerca nel caso in cui si presenti la necessità.

0.1 Organizzazione della tesi

Nel primo capitolo è descritto lo stato dell'arte del Web Mining ed è mostrata una sua possibile tassonomia; vengono altresì illustrati alcuni sistemi già realizzati.

Nel secondo capitolo viene spiegato l'approccio teorico sul quale si basa il seguente lavoro di tesi, in particolare

- si chiarisce il concetto di *Agente Software* e viene presentata una panoramica sulle diverse tipologie di agenti esistenti in letteratura
- si definisce cosa si intende per *Semantic Network*

- si introduce il *modello formale* alla base del sistema realizzato
- si definisce la *metrica* mediante la quale vengono giudicate le pagine Web

Nel terzo capitolo viene illustrato il sistema realizzato in ogni sua componente. In tale capitolo vengono anche presentati gli agenti sviluppati per mezzo di uno specifico linguaggio di descrizione degli agenti: il *Service Description Language (SDL)*. Vengono infine mostrate le strutture del *Web Repository* e della *Semantic Knowledge Base*.

Nel quarto capitolo vengono illustrati i risultati sperimentali ed i confronti effettuati con i risultati dei motori di ricerca.

Chapter 1

Stato dell'arte

Nel corso del capitolo, dopo aver spiegato cosa si intende per Web Mining, verranno illustrati alcuni sistemi che affrontano il problema dell'identificazione di informazioni nascoste nel WEB.

1.1 Una definizione di Web Mining

Il **Web Mining** è un nuovo problema di ricerca sotto discussione, che ha molto interesse per molte comunità. Attualmente non c'è accordo su cosa debba intendersi con il termine Web Mining, ma una sua possibile definizione è la seguente [4]

Definition 1.1.1 (Web Mining). Il Web Mining è l'attività di identificazione di modelli impliciti in una grande raccolta di documenti C , che può essere denotata con un mapping $\xi : C \rightarrow p$.

Poiché il Web Mining deriva dal Data Mining, la sua definizione è molto simile alla nota definizione di Data Mining¹.

¹Si definisce Data Mining l'attività usata per identificare modelli validi, nuovi, potenzialmente utili e comprensibili da una collettività di database [4].

Ciononostante, il Web Mining ha molte caratteristiche uniche se confrontato con il Data Mining:

- Le sorgenti del Web Mining sono documenti Web
- Il WEB è assimilabile ad un grafo orientato i cui nodi sono i documenti ed i cui archi sono i link. Per questo motivo, il modello identificato può riguardare il contenuto dei documenti oppure la struttura WEB
- I documenti Web sono semi-strutturati o non strutturati, mentre i dati su cui opera il Data Mining sono strutturati in database²

1.2 Una tassonomia per il Web Mining

Spesso c'è una forte ambiguità quando si parla di Web Mining, poiché con la stessa locuzione si può far riferimento a due campi di applicazione molto diversi tra loro:

- **Web Content Mining:** descrive il processo di ricerca di informazioni o risorse da milioni di fonti attraverso il Web
- **Web Usage Mining:** è il processo di mining dei file log di accesso al Web su uno o più siti

²Si noti che il Web Mining è differente dal Web Information Retrieval in quanto hanno differenti obiettivi. Sebbene il Web Mining è più del Web Information Retrieval esso non è inteso per rimpiazzarlo, invece sono due tecnologie che si integrano. D'altro canto il Web Mining può essere utilizzato per aumentare la precisione del Web Information Retrieval e migliorare l'organizzazione del risultato del recupero delle informazioni [4].

In questi due campi, poi, occorre fare un ulteriore sforzo di classificazione per poter capire bene tutte le possibili applicazioni cui spesso ci si riferisce semplicemente con la locuzione "Web Mining".

Una possibile tassonomia è rappresentata nella figura 1.1.

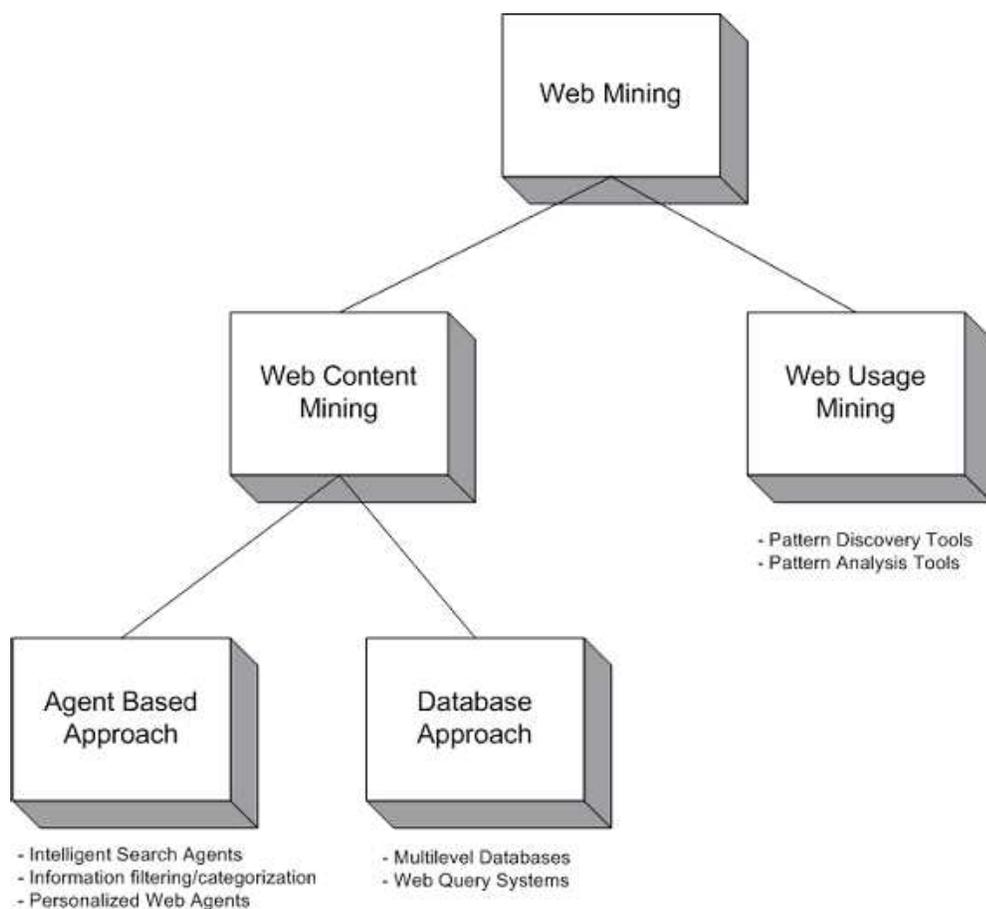


Figure 1.1: Tassonomia per il Web Mining

1.3 Web content mining

L'eterogeneità e la mancanza di struttura che caratterizza le fonti di informazione sul WEB ne rende difficile la scoperta automatica, l'organizzazione e la gestione delle informazioni. Gli strumenti tradizionali di ricerca e indicizzazione di Internet e del World Wide Web, come Lycos, Altavista, WebCrawler e altri, costituiscono un buon supporto per l'utente, ma generalmente non forniscono informazioni strutturate e non categorizzano, filtrano o interpretano i documenti.

Negli anni recenti c'è stato un grosso impulso da parte della ricerca per sviluppare strumenti di information retrieval più "intelligenti", come, ad esempio, agenti Web intelligenti o l'estensione di tecniche di database e di Data Mining per fornire un più alto livello di organizzazione per i dati semi-strutturati disponibili sul WEB.

Gli approcci generalmente proposti in letteratura sono quelli "Agent Based" e quelli "Data Base Approach". Nel seguito verranno discussi i due approcci.

1.3.1 Approccio Agent-Based

Di solito, i sistemi di **Web Mining Agent-Based** si possono dividere in tre categorie:

- **Agenti di ricerca intelligente:** tali agenti ricercano le informazioni rilevanti sulla base di caratteristiche del "*dominio*" e dei "*profili utente*" per organizzare e interpretare le informazioni trovate. Esistono agenti per informazioni di domini pre-specificati per cercare particolari tipi di documenti (ad esempio **Harvest** [2] [41], **FAQ-Finder** [42], **Information Manifold** [43]) ed agenti che imparano la struttura di informazioni sconosciute (ad esempio **ShopBot** [44])
- **Agenti per il filtraggio/categorizzazione delle informazioni:** usano varie

tecniche di *Information Retrieval* e le caratteristiche ipertestuali del WEB per reperire, filtrare e categorizzare i documenti (ad esempio **HyPursuit** [45] che clusterizza i documenti sulla base di informazioni insite nelle strutture dei link e del contenuto dei documenti)

- **Agenti Web personalizzati:** imparano le preferenze dell'utente e sulla base di esse (e di preferenze di utenti con interessi analoghi) ricercano le informazioni sul WEB (ad esempio **WebWatcher** [46])

1.3.2 Database Approach

Tale approccio consente di organizzare i dati semi-strutturati del WEB in collezioni di risorse più strutturate, usando i meccanismi standard di interrogazione dei database e tecniche di Data Mining per l'analisi. Sono state proposte due soluzioni:

- **Database multilivello:** in tali database il livello più basso contiene informazioni semi-strutturate memorizzate in diversi *Web Repositories*, ad esempio documenti ipertestuali. Ai livelli più alti sono estratti i metadati che vengono organizzati in collezioni strutturate di dati, come i database relazionali. Ad esempio, il sistema **ARANEUS** [60] estrae informazioni rilevanti dai documenti ipertestuali e li integra in *Iperesti Web Derivati (Derived Web Hypertexts)* di più alto livello che rappresentano generalizzazioni del concetto di vista di un database

- **Sistemi Web Query:** utilizzano i linguaggi standard di interrogazione dei database, ad esempio SQL, informazioni strutturali e perfino il linguaggio naturale per sottomettere query sul WEB. Ad esempio, **W3QL** [47] combina query riguardanti la struttura, basate sull'organizzazione degli ipertesti, e query riguardanti il contenuto dei documenti, basate su tecniche di Information Retrieval.

1.4 Web Usage Mining

Il **Web Usage Mining** consiste nella scoperta automatica dei pattern di accesso degli utenti ai siti Web a partire dai *log* memorizzati nei server Web. Le organizzazioni hanno a disposizione ingenti quantità di dati generate automaticamente dai server Web e raccolte nei *log*. Il Web Usage Mining consiste nell'applicazione di tecniche di Data Mining agli *usage logs* di grandi repository di dati Web in modo da analizzare come un sito viene utilizzato. L'analisi di utilizzo include statistiche dirette, come la frequenza di accesso alle pagine, e forme più sofisticate di analisi, come trovare *paths* di attraversamento comuni attraverso un sito Web (*common traversal paths*).

Le informazioni sull'utilizzo possono essere usate per ristrutturare un sito Web in modo da servire meglio l'utente:

- Percorsi di attraversamento lunghi e complessi o il basso utilizzo di una pagina con informazioni importanti possono suggerire che i link del sito e le informazioni non sono posizionate in modo intuitivo

- Il progetto del layout fisico dei dati o lo schema di caching per un server Web possono essere migliorati dalla conoscenza di come gli utenti tipicamente navigano attraverso il sito
- Le informazioni sull'uso possono anche essere usate per aiutare direttamente la navigazione del sito fornendo una lista di destinazioni preferite da una particolare pagina Web

Distinguiamo due tipologie di tools: **pattern discovery tools** e **pattern analysis tools**.

1.4.1 Pattern discovery tools

Sono tools che usano tecniche sofisticate di intelligenza artificiale, Data Mining, psicologia e teoria dell'informazione per "scovare" conoscenza da collezioni di dati. Ad esempio, **WEBMINER** [48] è in grado di scovare automaticamente regole d'associazione e pattern sequenziali dai log.

1.4.2 Pattern analysis tools

Sono tools per interpretare e visualizzare i pattern trovati. Alcuni tools fanno uso di *tecniche OLAP* come i *data cube* per semplificare l'analisi. **WEBMINER** [48] propone un meccanismo di query simile a SQL per interrogare i pattern scoperti.

1.5 L'architettura del sistema WebTMS

Per meglio comprendere gli approcci descritti, nel seguito illustriamo un certo numero di sistemi presentati in letteratura.

WebTMS [38] (*Web Text Mining System*), sviluppato dal Dipartimento di Computer Science and Technology dell'Università di Nanjing, è un sistema multi-agente che combina il Text Mining con l'analisi multidimensionale dei documenti per aiutare l'utente nel mining dei documenti HTML reperibili sul WEB. Lo scopo è quello di identificare pattern validi e potenzialmente utili in una collezione di dati presenti in una Web Warehouse.

L'architettura del sistema viene mostrata sinteticamente in figura 1.2.

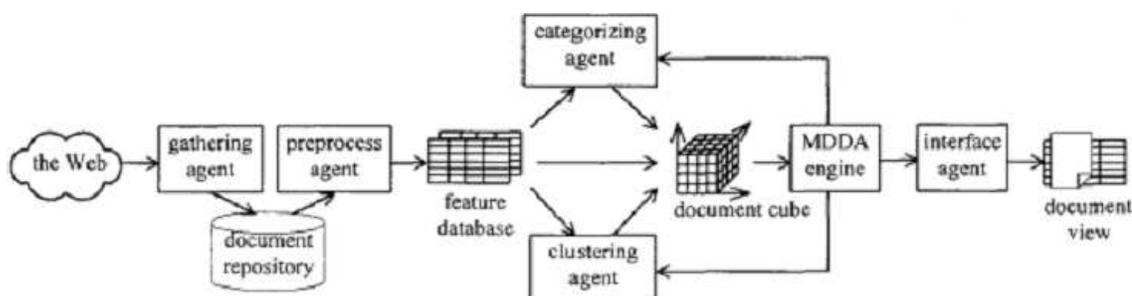


Figure 1.2: L'architettura del sistema WebTMS

Analizziamo di seguito i singoli componenti dell'architettura. Si noti che ognuno degli agenti può completare il lavoro in modo relativamente autonomo. Questi agenti possono essere messi in un computer o in diversi computer distribuiti in una rete. Inoltre, nuovi componenti possono essere facilmente aggiunti in WebTMS vista l'alta modularità.

1.5.1 Gathering agent

Tale agente, illustrato in figura 1.3, recupera i documenti, che possono essere distribuiti su parecchi Server Web, su cui effettuare il mining.

Tali documenti sono immagazzinati in opportuni repository del sistema. I documenti vengono raccolti a partire da una lista iniziale di URL immessi dall'utente.

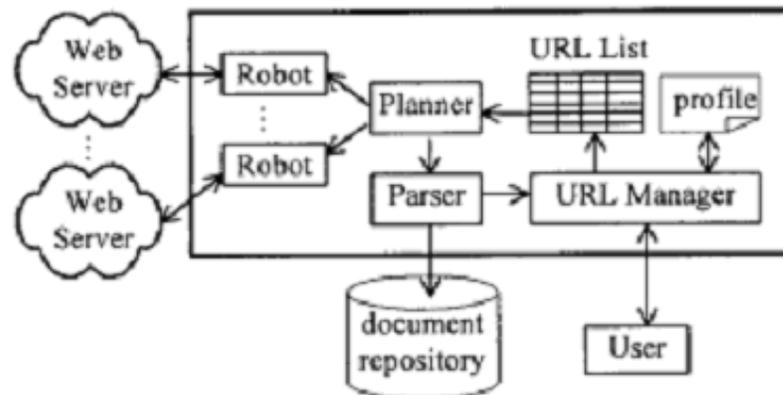


Figure 1.3: Il Gathering Agent

1.5.2 Preprocessing Agent

Nella figura 1.4 è rappresentato lo schema del Preprocessing Agent.

A partire da regole euristiche, vengono estratti metadati che poi vengono utilizzati per dedurre una rappresentazione del documento.

I metadati e la rappresentazione dei documenti sono tutti memorizzati, come dati strutturati nel database delle caratteristiche, per un uso futuro nel Text Mining.

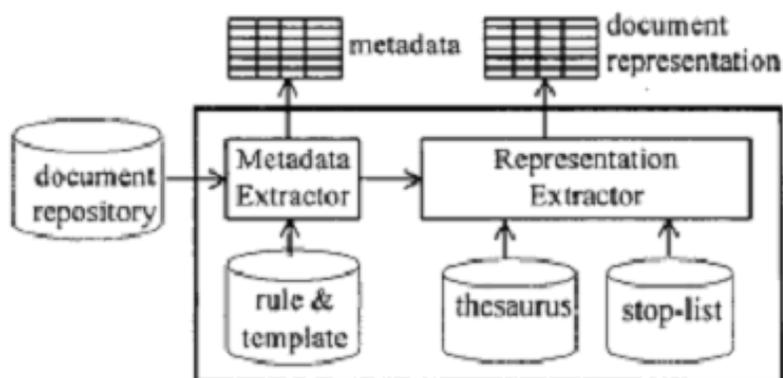


Figure 1.4: Il Preprocessing Agent

1.5.3 Categorization Agent

L'agente è mostrato in figura 1.5.

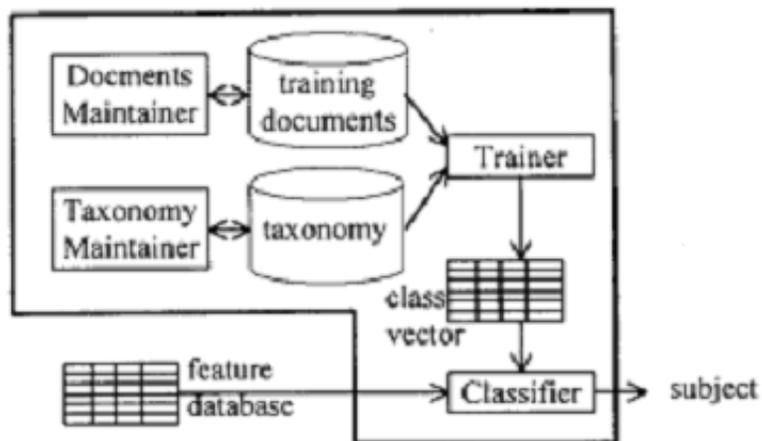


Figure 1.5: Il Categorization Agent

A partire da una tassonomia predefinita vengono divisi i documenti in classi, viene in altri termini realizzata una categorizzazione.

1.5.4 Clustering Agent

L'agente è raffigurato in figura 1.6. A differenza del *categorization agent*, tale agente non fa uso di una tassonomia predefinita per clusterizzare i documenti.

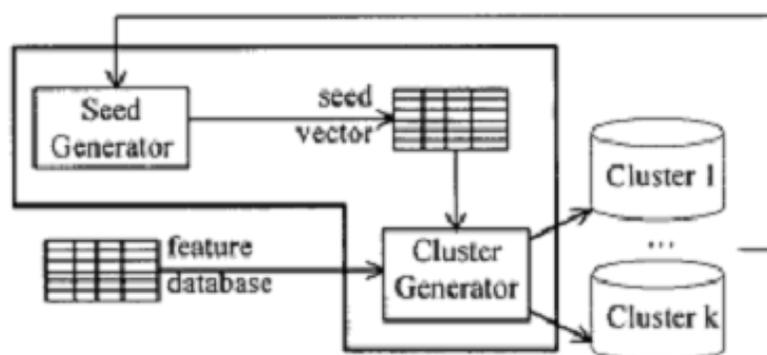


Figure 1.6: Il Clustering Agent

L'obiettivo è di dividere una raccolta di documenti in un insieme di cluster tali che è massimizzata la somiglianza intra-cluster e minimizzata quella inter-cluster.

1.5.5 MDDAE (Multi-Dimension Document Analysis Engine)

WebTMS usa una tecnologia di analisi multidimensionale basata su un *document cube*, per fornire viste multi-dimensione dei documenti per gli utenti.

Definition 1.5.1 (Document cube). Dato un documento D da analizzare, si definisce "document cube" la tupla:

$$CD = (d_1, d_2, \dots, d_m, D)$$

dove D è il centro attorno al quale i metadati (di dimensioni d_i) girano.

Un esempio di document cube è mostrato in figura 1.7.

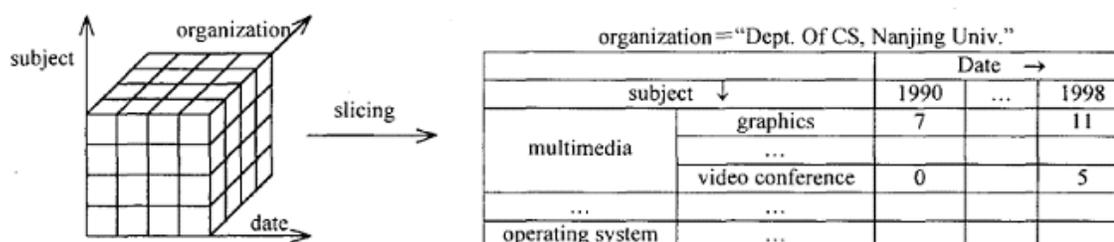


Figure 1.7: Un esempio di Document Cube

Si possono effettuare operazioni di slicing, dicing, rotation, drill down, roll up per dare all'utente diversi punti di vista del documento.

MDDAE offre anche funzionalità di analisi statistica per rivelare le distribuzioni o le tendenze nei documenti.

1.5.6 User Interface Agent (UIA)

È un'interfaccia visuale tra gli utenti e il MDDAE. Esso traduce le richieste degli utenti in un linguaggio specializzato e le passa ai moduli MDDAE, Document Categorization Agent o Document Clustering Agent, e poi mostra le viste dei documenti multi-dimensione e i documenti ritornati per gli utenti.

1.6 Il meta-motore GSA

GSA (Global Search Agent) [1] è un meta-motore di ricerca semi-adattativo sviluppato dal *DEIS* dell'Università della Calabria congiuntamente con le società

chiave inserite dall'utente.

ME estrae i singoli risultati dai motori di ricerca appena essi sono disponibili e chiede a SE di iniziare lo spidering delle pagine relative agli URL restituiti dai motori di ricerca. SE si occupa dello spidering e del parsing degli URL e di assegnare ad ogni pagina un punteggio (rank). Inoltre, SE decide se l'URL è abbastanza interessante così da mostrare soltanto gli URL che sembrano rilevanti per l'utente e decide se è il caso di effettuare ulteriori operazioni di spidering di pagine a profondità maggiore.

La ricerca termina quando ME ed SE non hanno più documenti da analizzare, cioè quando non occorre più effettuare alcuna operazione di spidering.

Il *Feedback Engine (FE)* lavora off-line. L'utente può specificare la sua opinione su quali sono i documenti effettivamente di suo interesse e quali non lo sono. FE a questo punto restituisce un set di parole rilevanti così che l'utente possa rifinire la ricerca.

1.6.2 Il ranking delle pagine

GSA utilizza una funzione di ranking molto interessante, che ha dato buoni risultati in fase di sperimentazione, pertanto è utile illustrarla, anche perché costituisce un buon confronto con la metrica che sarà definita più avanti per il sistema sviluppato.

Siano:

- D un documento
- W l'insieme delle parole chiave
- $w = \text{card}(W)$
- K_0 e K_1 due costanti (il cui ruolo sarà chiarito in seguito)

- K_2 , K_3 e K_4 opportuni pesi dei tre contributi specificati in seguito
- K_5 la massima distanza (in parole) considerata significativa per due occorrenze di termini in W
- N_p la somma dei valori di *presenza* di ogni termine (per *presenza* si intende la massima *somiglianza* trovata per un certo termine nel testo considerato, così realizzando una particolare tecnica di stemming)
- N_t la somma delle occorrenze delle parole di W , ognuna pesata con il corrispondente valore di *somiglianza*
- N_s il numero di parole di W con un significativo valore di *presenza*
- $d(i, j)$ la minima distanza trovata tra due occorrenze significative delle parole i e j

Si calcolino le seguenti quantità:

$$1. \mathbf{b}_0 = K_2 \cdot \frac{N_p}{w}$$

$$2. \mathbf{f}_0 = K_3 \cdot \frac{N_t}{w}$$

$$3. \mathbf{d}_0 = K_4 \cdot \frac{K_5 - \frac{\sum_{i=0}^{N_s-1} \sum_{j=i+1}^{N_s} \min(d(i,j), k_5)}{\sum_{k=0}^{N_s-1} (N_p - k)}}{K_5}$$

$$4. \mathbf{f} = b_0 + f_0 + d_0$$

Ed infine la funzione di ranking relativa al documento D ed al set di keywords W viene calcolata come segue:

$$\mathbf{r}(\mathbf{D}, \mathbf{W}) = k_0 \cdot \frac{f}{f + k_1}$$

Le costanti k_0 e k_1 hanno l'unico scopo di controllare la forma della funzione di ranking. Osserviamo che la funzione di ranking di GSA:

- Incorpora una tecnica di "stemming"
- Esprime le distanze in parole e non in caratteri
- Assume valori compresi tra 0 e K_0

Lo **stemming** è il processo di trasformazione di una parola nella sua radice ("stem"). Per esempio, considerando la lingua inglese, i diversi termini *learns*, *learned*, *learning* si possono ricondurre alla stessa radice learn.

1.6.3 Come GSA impara dall'utente

L'utente di GSA può esprimere una preferenza di tipo booleano (interessante o non interessante) su ogni documento restituito dall'agente, oppure ignorarne alcuni; dopodiché può chiedere al sistema di tener conto delle sue preferenze.

GSA effettua un parsing dei documenti interessanti (*hot*) e non (*cold*) ed estrae due insiemi di vocaboli: un insieme *good* ed un insieme *bad*.

Per capire se un termine è *good* o *bad* viene applicata un'euristica, dopo aver eliminato le *stop-words* dai documenti, ossia parole molto comuni che non inficiano il contenuto semantico del documento.

L'idea alla base dell'euristica è una funzione di ranking delle singole parole e di considerare *good* soltanto i termini il cui *rank* supera una certa soglia.

Siano:

- **G** l'insieme dei documenti *hot*
- **B** l'insieme dei documenti *cold*
- **t** il termine di cui si vuole calcolare il rank
- **\mathbf{o}_i** il numero di occorrenze del termine *t* nel documento *i*
- **W** l'insieme delle parole della query di partenza
- **$\mathbf{w}_k \in W$** il k-esimo termine della query di partenza
- **$\mathbf{D}(\mathbf{w}_k, \mathbf{t})$** la minima distanza (in parole) tra una occorrenza significativa del termine w_k ed il termine *t*
- **$\mathbf{w} = \text{card}(\mathbf{W})$**

- $\mathbf{g} = \text{card}(\mathbf{G})$
- $\mathbf{b} = \text{card}(\mathbf{B})$

Si calcolino le seguenti quantità:

1. $\mathbf{O}_i = \min_{k=1}^w (d(w_k, t) \cdot \sum_{j=0}^{O_i-1} (\frac{1}{2^j}))$
2. $\mathbf{r} = \sum_{i=1}^g \frac{O_i}{g} - \sum_{i=1}^b \frac{O_i}{b}$

in cui r è il rank del termine t . Ogni termine incrementa il suo punteggio se appare in un documento *hot* e decrementa il suo punteggio se appare in un documento *cold*. Altre occorrenze di un termine in un documento oltre la prima non alterano più di tanto il valore di r .

1.7 Il motore di ricerca Mondou, basato sul Text Mining

Il motore di ricerca giapponese **Mondou** [39] è stato sviluppato dal *Department of Systems Science* presso la Graduate School of Informatics dell'Università di Kyoto.

Nonostante si tratti di un motore di ricerca di prima generazione, le tecniche di Data Mining sono adottate sin dal 1995. Una delle tecniche più importanti utilizzate dal motore è quella basata sulle regole d'associazione.

Il motore dà agli utenti una lista di *associative keywords*, che costituiscono un feedback per l'utente mediante il quale può raffinare la sua ricerca.

1.7.1 L'architettura di Mondou

Mondou è costituito da tre moduli fondamentali: il *Database*, il *Query Server* ed i *Web Robots*.

In figura 1.9 è mostrata l'architettura del sistema.

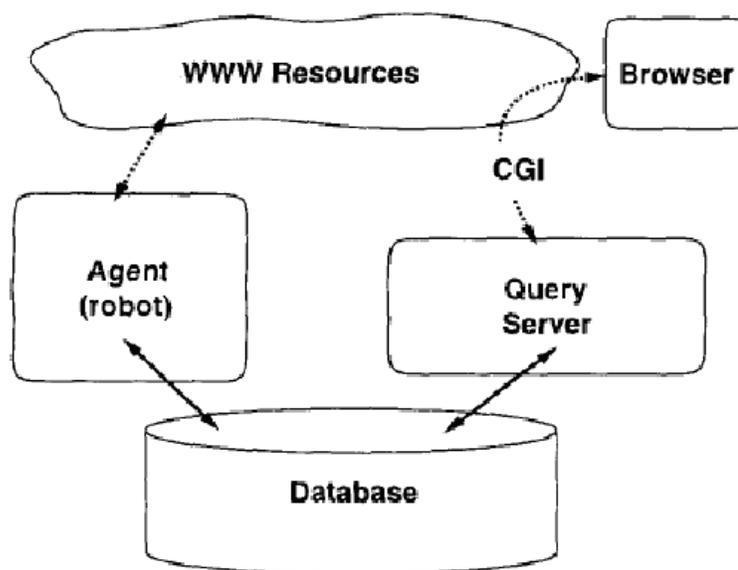


Figure 1.9: L'architettura del sistema Mondou

Il *modulo Database* non memorizza semplicemente il testo delle pagine Web, ma anche le relazioni tra i link.

Il *Query Server* è il modulo CGI più importante e si occupa di restituire i risultati delle ricerche, le associative keywords e i valori degli attributi per l'interfaccia grafica.

Il *modulo Web Robot* colleziona i dati dai server Web su Internet o su Intranet e li memorizza nel database. Oltre alle solite funzioni dei normali Web Robots, questo robot effettua il parsing dei documenti e sceglie le parole chiave più importanti mediante operazioni euristiche basate su analisi testuali del linguaggio naturale e dà

pesi appropriati ai dati.

1.7.2 Il data mining nel motore di ricerca Mondou

Mondou applica molti algoritmi di data mining per rendere più efficace la ricerca di informazioni sul WEB. Il più importante è quello volto a trovare le *associative keywords*.

Definition 1.7.1. *Dato un set di documenti, ognuno costituito da un certo numero di keywords, una **regola d'associazione** è un'espressione*

$$X \rightarrow Y$$

con X e Y insiemi di keywords.

Il significato intuitivo di un tal tipo di regola è che i documenti che contengono le keywords X tendono a contenere anche le Y . X si dice *testa* della regola o *premessa*; Y si dice *corpo* della regola o *conseguenza* [40].

La validità di una regola d'associazione è misurabile mediante due quantità:

- **Supporto:** è la probabilità che in una osservazione siano presenti sia la testa sia il corpo di una regola
- **Confidenza:** è la probabilità che in una osservazione sia presente la testa di una regola, essendo già presente il suo corpo.

Intuitivamente, il supporto misura l'importanza di una regola (quanto spesso premessa e conseguenza sono presenti contemporaneamente) mentre la confidenza ne

misura l'affidabilità (quante volte, essendo presente la premessa, è presente anche la conseguenza).

Keywords	Number of URLs	Related Keywords
Applications	975	windows, speech, computer, internet, helper, commercial, user, audio, network
Informatics	462	genome, medical, departement, center, school, faculty
Japan	20,366	japan, country, webec, japan, japanese, back
Engine	2,836	lycos, yahoo, dragon, search, creative, japan
Disney	183	tokyo, disney land, walt, world, parks, records
Mining	134	akita, geology, college, net, data
Mondou	28	rcaau, search

Table 1.1: Esempi di *associative keywords*

Mondou estende l'algoritmo di ricerca delle regole di associazione per considerare la frequenza e l'importanza delle keywords che appaiono nei documenti Web. Il Web Robot di Mondou dà l'appropriato valore di peso alle keywords in base alla frequenza ed al tipo di tag html. Esempi di associative keywords sono mostrati in Tabella 1.1.

In particolare, Mondou suddivide i tipi di attributi in tre categorie:

- **Characteristic** include keywords che mostrano le caratteristiche fondamentali di un documento e che sono specificate dall'autore o dall'editore del documento stesso, come ad esempio titolo e keywords

- **Description** include frasi che descrivono il contenuto, come ad esempio sommario, titolo, etc.
- **Supplement** include informazioni aggiuntive, come ad esempio autore, editore, dimensioni, data, etc.

Category	Keyword space	INSPEC	Average
Characteristics	Small	Title Keywords	11 30
Description	Small ~ Large	Abstract	183
Supplement	Small	Author	3

Table 1.2: Le dimensioni degli spazi delle keywords

Come mostrato dalla tabella 1.2, lo spazio delle keywords delle categorie *Characteristic* e *Supplement* include un numero di keywords troppo piccolo per poter effettuare significative operazioni di mining basate sulle regole d'associazione.

La soluzione proposta dai ricercatori giapponesi è quella di estendere l'algoritmo base aumentando il numero delle keywords introducendo relazioni tra attributi differenti, così come mostrato in figura 1.10.

Ad esempio, si introduce una relazione tra il titolo A_u e l'autore A_v per espandere le keywords del titolo, considerando titoli di altri documenti dello stesso autore.

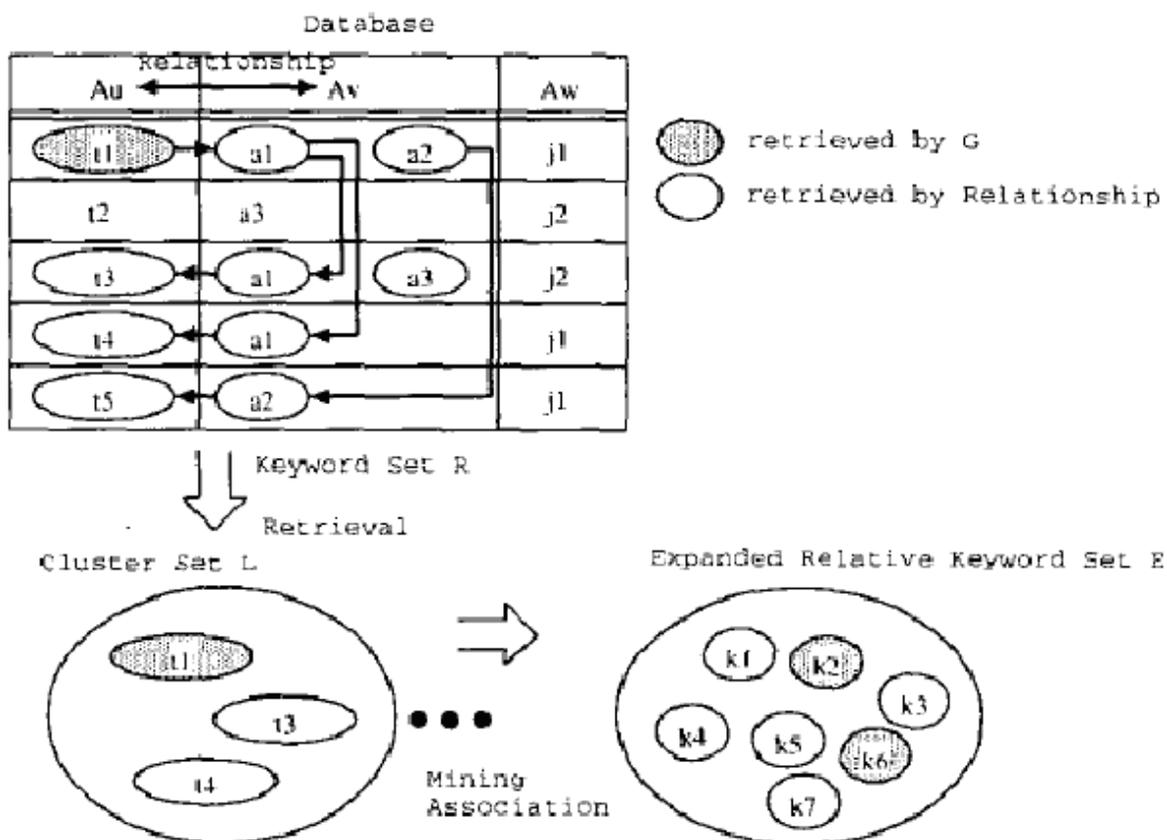


Figure 1.10: Relazioni tra attributi differenti

1.8 Il progetto Harvest

Scopo del progetto **Harvest Information Discovery and Access System** [2][41] non è soltanto il reperimento dei dati, ma anche il filtraggio degli stessi: filtrare enormi masse di dati per trarne informazioni concise riguardo il loro contenuto. Harvest è stato annunciato nel novembre del 1994 da un team di ricercatori universitari e del mondo dell'industria e fu rilasciato nell'agosto del 1996.

Lo scopo primario del progetto era di facilitare l'accesso alle diverse risorse di informazioni disponibili su Internet. L'idea di base era un'indicizzazione delle risorse disponibili sul WEB, dando la possibilità di cercare negli indici in modo flessibile e replicare tali indici su Internet così da ridurre il traffico ridondante di informazioni inutili sulla rete e facilitare il reperimento delle sole informazioni di interesse dell'utente.

Alcuni prodotti commerciali hanno fatto uso dei tool di Harvest, soprattutto in virtù della sua buona abilità nel caching delle pagine. Alcuni esempi sono il *Netscape Catalogue Server*, oppure *Glimpse*, un'applicazione per Unix per la creazione di motori di ricerca.

Il progetto, però, nasceva con ambizioni ben più grandi ed uno degli auspici dei suoi promotori era quello di essere utilizzato in larga scala dai maggiori motori di ricerca. Ciò non avvenne principalmente a causa del fatto che Harvest fu rilasciato troppo tardi, quando i motori di ricerca più importanti già esistevano ed avevano la propria tecnologia e non erano ben disposti a cambiarla in modo così radicale.

Chapter 2

Approccio teorico

2.1 Gli Agenti Software

Nei decenni passati gli sviluppatori di software hanno realizzato un'ingente quantità di programmi software monolitici caratterizzati dal fatto che ciascuno di essi realizzava un ampio e variegato insieme di task. Negli ultimi anni, al contrario, c'è stato un cambiamento nella modalità di sviluppo: invece di programmi con milioni di linee di codice, ci si è orientati verso pezzi di codice più piccoli e modulari, dove ogni modulo realizza un task ben definito (o, al limite, un insieme limitato di task) piuttosto che migliaia di task differenti, come si era soliti fare, per esempio, nei legacy systems.

Gli agenti software rappresentano l'ultimo ritrovato tecnologico in questa tendenza a suddividere sistemi software complessi in componenti. Un agente software è un "pezzo" di software avente le seguenti caratteristiche [5]:

- È in grado di agire in modo tale da portare a termine task per conto degli utenti
- Mette a disposizione uno o più servizi che altri agenti possono usare sotto certe condizioni

- Include una descrizione dei servizi offerti
- È in grado di agire autonomamente senza eventuali indicazioni da parte di un essere umano
- Descrive in che modo un agente decide quale azione intraprendere (anche se questa descrizione può essere tenuta nascosta agli altri agenti)
- È in grado di interagire con altri agenti (incluso esseri umani)

È da notare che non tutti gli agenti software debbono possedere le proprietà menzionate anche se ogni paradigma di programmazione basata su agenti dovrebbe dare la possibilità di sviluppare agenti con tutte o, eventualmente, alcune di queste proprietà.

2.1.1 Tipologie di Agenti

In questo paragrafo si tenterà di collocare i vari agenti esistenti in diverse classi in modo tale da poter definire una precisa tipologia degli agenti software [2] .

In primo luogo, gli agenti possono essere classificati in base alla loro mobilità (intesa come capacità di "muoversi" su una qualche rete): in base a ciò gli agenti possono suddividersi in **Statici** o **Mobili**.

In secondo luogo, possono suddividersi in **Deliberativi** o **Reattivi**. I primi derivano dal paradigma di pensiero deliberativo: posseggono al loro interno un modello di ragionamento simbolico e sono in grado di agire, pianificare e negoziare in modo tale da coordinarsi con altri agenti. Gli agenti di tipo reattivo, al contrario, non posseggono alcun modello simbolico dell' ambiente in cui si trovano e agiscono sulla base di un modello di comportamento stimolo/risposta, ossia rispondendo allo stato presente dell'ambiente in cui sono collocati [6].

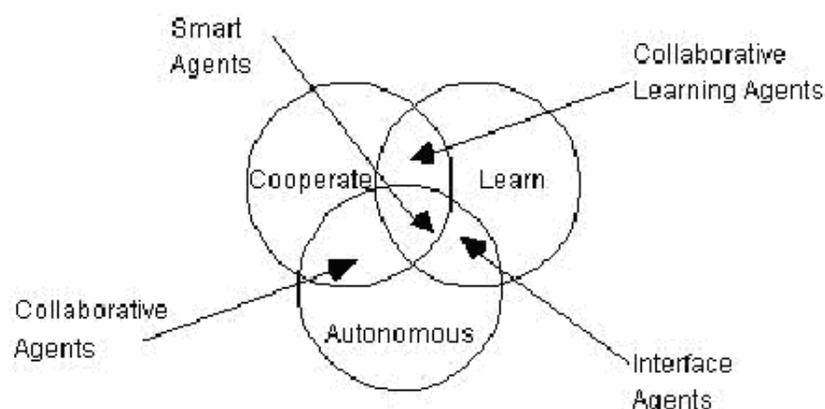


Figure 2.1: Una tipologia di agenti

In terzo luogo, gli agenti possono essere classificati in base ad alcuni fondamentali attributi che essi dovrebbero possedere. Tre sono gli attributi principali: *Autonomia*, *Apprendimento* e *Cooperazione*.

Con il termine **autonomia** si intende la capacità degli agenti di agire senza la guida dell'uomo. Tali agenti posseggono degli stati interni e degli obiettivi e agiscono in modo tale da realizzare i loro scopi nell'interesse dei propri utenti. Elemento chiave dell'autonomia è la cosiddetta "*proactiveness*", ossia la capacità di prendere iniziative piuttosto che agire semplicemente in risposta agli stimoli dell'ambiente circostante [7].

Con il termine **cooperazione**, si intende la capacità di interagire con altri agenti. È un attributo fondamentale per poter avere agenti multipli piuttosto che uno solo che lavora su un determinato task. Per poter cooperare, gli agenti devono possedere quella che viene definita capacità sociale ossia la capacità di interagire con altri agenti o esseri umani tramite qualche linguaggio di comunicazione.

Ultimo attributo è l'**apprendimento**; con esso si intende la capacità di un agente di imparare come esso reagisce e/o interagisce con l'ambiente esterno.

Sulla base di tali caratteristiche si può derivare una prima tipologia di agenti. Essi sono rappresentati nella figura 2.1 e sono¹:

1. **Collaborative agents**
2. **Collaborative learning agents**
3. **Interface agents**
4. **Smart agents**

In quarto luogo, gli agenti sono talvolta classificati sulla base dei ruoli svolti; per esempio, si può parlare di **World Wide Web Information Agents** intendendo così agenti che permettono di gestire le informazioni presenti su Internet, oppure di **Testing Agents** intendendo agenti che facilitano il testing del software. Altri tipi di agenti sono i **Presentation Agents** e gli **Analysis and Design Agents**, tutti agenti classificati in base ai compiti svolti.

Un' ultima tipologia di agenti sono gli **agenti ibridi** i quali uniscono due o più diverse caratteristiche in un singolo agente.

Oltre agli attributi menzionati, ce ne sono altri che vengono però considerati secondari ai fini della determinazione di una tipologia di agenti, quali:

- **Versatilità**, ossia la capacità di avere più obiettivi
- **Mendacità**, ossia la capacità di essere o no sempre veritieri

¹Tali agenti verranno spiegati nel dettaglio successivamente

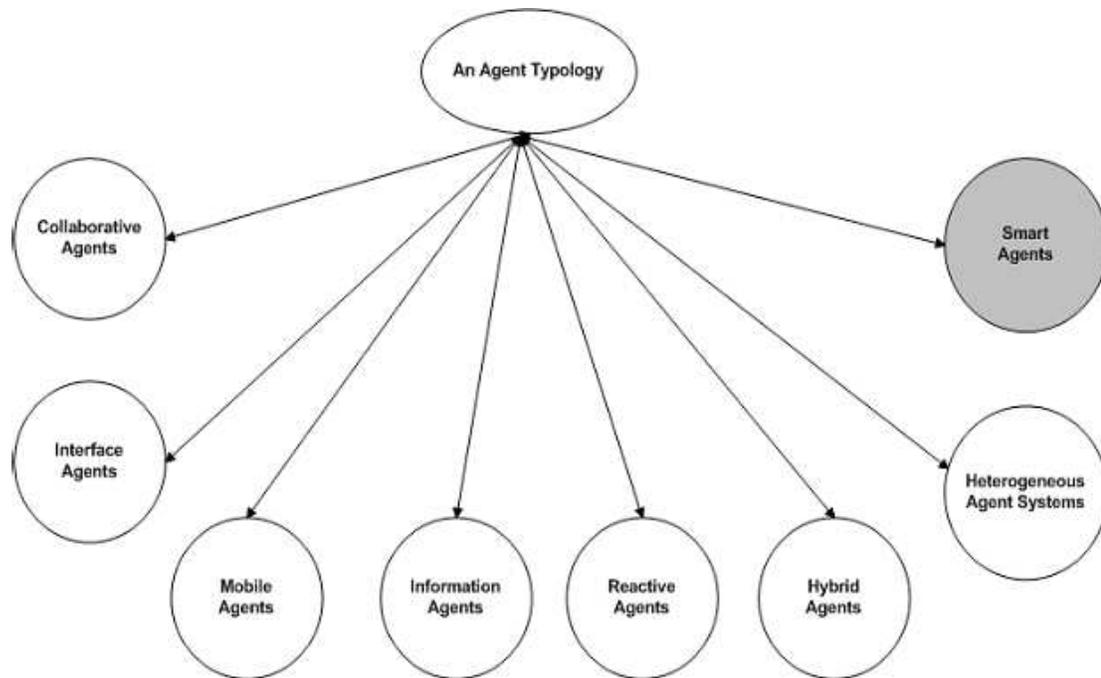


Figure 2.2: Classificazione di Agenti Software

Sulla base di quanto esposto, è possibile identificare sette tipi di agenti, i quali verranno spiegati nel dettaglio nei paragrafi successivi:

1. **Collaborative agents**
2. **Interface agents**
3. **Mobile agents**
4. **Information/Internet agents**
5. **Reactive agents**
6. **Hybrid agents**

7. Heterogeneous agents system

8. Smart agents

Si noti che rispetto alla figura 2.2 non sono presenti i *Collaborative Learning Agents*. Ciò è motivato dal fatto che non si è ancora a conoscenza di agenti che collaborano e apprendono, ma non sono autonomi.

C'è da dire che esistono diverse applicazioni che combinano agenti appartenenti a due o più categorie; a tali applicazioni ci si riferisce con il termine *Sistemi ad Agenti Eterogenei* (*heterogenous agent systems*). Tali applicazioni, però, non sono così diffuse come quelle elencate precedentemente.

Un'altra nota è che gli agenti non devono essere necessariamente benevoli gli uni con gli altri; infatti, è possibile avere agenti in competizione tra loro o, addirittura, antagonisti. Comunque, si può vedere agli agenti competitivi come ad una potenziale sottoclasse di quelli visti. In tal modo è possibile avere *Collaborative Agents* di tipo competitivo, *Interface Agents* di tipo competitivo e via dicendo.

Collaborative Agents

Come mostrato in figura 2.1, tali agenti pongono l'accento sulle caratteristiche di autonomia e cooperazione (con altri agenti) al fine di realizzare task per gli utenti [8].

Essi possono eventualmente anche apprendere ma ciò non è un aspetto fondamentale. Al fine di collaborare, possono avere necessità di effettuare negoziazioni per accordarsi su diversi aspetti dei task in cui sono coinvolti.

Alcuni ricercatori nel campo dell'*Intelligenza Artificiale* (AI) stanno tentando di

fornire delle definizioni più adeguate per tali agenti ponendo l'accento su alcuni attributi di tipo mentale come convinzioni, desideri e intenzioni (beliefs, desires e intentions, da cui l'acronimo di collaborative agents di tipo BDI).

Brevemente, le caratteristiche fondamentali di tali agenti sono *autonomia*, *capacità sociale*, *comprensione* e *proactiveness*. In tal modo sono in grado di agire in maniera razionale ed autonoma in un ambiente multi-agente. In genere sono statici e possono essere benevoli e veritieri o una qualche combinazione di essi.

Ci sono diversi motivi che possono portare allo sviluppo di sistemi basati su tale tipo di agenti:

- Risolvere problemi troppo complessi da realizzare tramite un sistema basato su un singolo agente centralizzato
- Permettere l'interconnessione di legacy system esistenti come sistemi esperti o sistemi di supporto alle decisioni
- Risolvere problemi di natura distribuita come sistemi di sensori distribuiti o di controllo del traffico aereo [9]
- Affrontare problemi con sorgenti di informazione distribuite
- Ottenere modularità (che riduce la complessità), velocità (dovuta al parallelismo), affidabilità (dovuta alla ridondanza), flessibilità (nuovi task possono essere affrontati tramite nuovi agenti).

Un esempio di sistema basato su agenti di tipo collaborativo è quello sviluppato presso il CMU: il **Pleiades Project** [10][11]. Tale sistema applica tale tipo di agenti

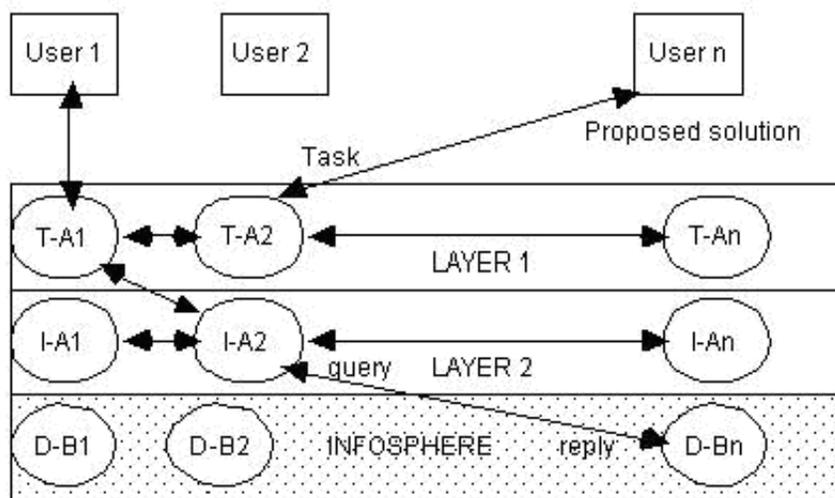


Figure 2.3: Architettura Pleiades

al problema di effettuare decisioni di tipo organizzativo. Esso è raffigurato in figura 2.3 e si basa su un'architettura ad agenti a due livelli: il primo livello contiene i *Tasks-Specific Collaborative Agents*, che realizzano specifici task come organizzare appuntamenti e incontri.

I *Tasks-Specific Collaborative Agents* collaborano tra di loro per risolvere conflitti e integrare informazioni. Per ottenere le informazioni necessarie, essi le richiedono agli agenti di secondo livello: gli *Information-Specific Agents*. Anch'essi possono collaborare tra loro ed attingono informazioni da specifici database di sistema. Alla fine, i *Tasks-Specific Collaborative Agents* propongono soluzioni ai diversi utilizzatori del sistema.

Sebbene ci siano altri esempi di sistemi basati su agenti di tipo collaborativo, questi ultimi sono stati poco sviluppati anche se la situazione sta mutando, come dimostrano i sistemi di controllo dei malfunzionamenti dello Space Shuttle o i simulatori di volo

usati presso la Royal Australian Airforce [12]. Ci sono, infatti, alcuni problemi che si sta cercando di risolvere. Se ne possono menzionare alcuni anche se essi non sono tipici solamente di tale categoria di agenti:

- Cercare di ingegnerizzare lo sviluppo dei sistemi [13]
- Migliorare la coordinazione Inter-Agenti; la coordinazione è essenziale affinché gruppi di agenti risolvano in maniera corretta i problemi. Senza un'adeguata teoria sulla coordinazione, c'è il rischio che anarchia o arresti non desiderati prendano il sopravvento [14]
- Stabilire eventuali tecniche di apprendimento che non portino all'instabilità dei sistemi
- Stabilire tecniche e metodologie per integrare agenti e legacy system
- Stabilire dei metodi per la verifica e la validazione dei sistemi che permettano di capire se le specifiche funzionali sono rispettate

In conclusione si può affermare che, nonostante i problemi ancora da risolvere, sono molti i campi di interesse in cui possono trovare applicazione tali tipi di agenti.

Interface Agents

Gli **Interface Agents**, come si nota nella figura 2.1, pongono l'accento sulle caratteristiche di autonomia e apprendimento al fine di portare a termine i task per i loro utilizzatori.

Pattie Maes, che rappresenta uno dei maggiori ricercatori in questo campo, ha evidenziato quella che è la metafora alla base di tali agenti: *un assistente personale*

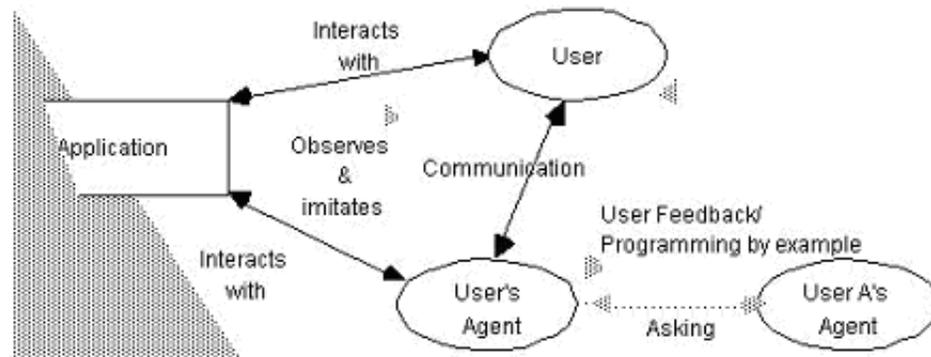


Figure 2.4: Interface Agents

che collabora con l'utente nello stesso ambiente di lavoro [15]. Da ciò si evince la sottile differenza tra questo tipo di agenti e quelli descritti nel paragrafo precedente: questi ultimi collaborano tra di loro mentre gli Interface Agents collaborano esclusivamente con gli utenti. In figura 2.4 viene mostrata la modalità di funzionamento di tali agenti.

In sintesi, tali agenti forniscono generalmente assistenza ad un utente che sta utilizzando una particolare applicazione (per esempio, fogli di calcolo o sistemi operativi). Tali agenti monitorano le azioni intraprese dagli utenti, apprendono nuovi shortcuts e suggeriscono in tal modo le migliori modalità per portare a termine un lavoro. Perciò tali agenti possono essere visti come assistenti personali che collaborano con l'utente nella realizzazione dei task in un'applicazione.

Tali tipi di agenti posseggono la capacità di apprendere al fine di assistere meglio gli utenti; ciò viene fatto in 4 modi [15]:

1. Osservando e imitando gli utenti
2. Ricevendo feedback positivi e negativi dagli utenti

3. Ricevendo esplicite istruzioni dagli utenti
4. Chiedendo ad altri agenti

Riassumendo, si può affermare che, sotto certe ipotesi, un tale tipo di agente è in grado di programinarsi ossia di acquisire la conoscenza necessaria per assistere gli utenti. L'agente è dotato di una conoscenza minima che può aumentare osservando gli utenti oppure chiedendo ad altri agenti [15].

La cooperazione con altri agenti è limitata esclusivamente alla richiesta di consigli e non ad eventuali negoziazioni al fine di collaborare al compimento di un task (come si verificava nei Collaborative Agents).

L'osservazione alla base dello sviluppo di tali agenti, è che le interfacce della maggior parte dei sistemi sono di tipo passivo, ossia eseguono compiti solamente quando sono sollecitate. Con lo sviluppo di tali agenti si cerca di avere, invece di un utente che impartisce comandi ad un'interfaccia, una modalità di funzionamento basata sulla collaborazione uomo-agente, in cui entrambi possono iniziare una comunicazione con l'altro, realizzare task e monitorare eventi. La chiave, quindi, è questa cooperazione tra l'uomo e l'agente [16]. Pertanto, l'obiettivo è la migrazione da una metafora basata sulla diretta manipolazione ad una che deleghi alcuni task a degli Interface Agents magari per aiutare utenti alle prime armi.

Sono tre i principali benefici che derivano dall'utilizzo di tali agenti [15]. In primo luogo rendono meno faticoso il lavoro degli utenti in quanto task noiosi o faticosi possono essere delegati ad essi (per esempio la gestione delle mail vecchie, la schedulazione degli appuntamenti, la ricerca di informazioni da un grosso insieme). In secondo luogo, l'agente può adattarsi nel tempo alle abitudini e le preferenze dei propri utenti.

In ultimo, la conoscenza all'interno di una comunità di utenti può essere condivisa (ciò accade, per esempio, quando un agente acquisisce conoscenza da un suo pari). Ci sono diversi esempi di sistemi basati su tali agenti:

- *Calendar Agent* [51], sviluppato da Kozierok & Maes, che assiste gli utenti nella schedulazione degli appuntamenti apprendendo le loro preferenze
- *Letizia Agent* [?], sviluppato da Liebermann, che assiste gli utenti nella navigazione sul WEB effettuando ricerche di pagine Web basandosi sui comportamenti degli utenti
- *The Remembrance Agent* [19], sviluppato da Rhodes & Starner, che restituisce, durante la scrittura di una e-mail, le cinque e-mail più attinenti a quella che si sta scrivendo contenute nel proprio archivio
- *The Kasbah Agent* [20], sviluppato da Chavez & Maes, che filtra gli annunci pubblicitari sul WEB e li seleziona in base agli interessi degli utenti

Concludiamo il paragrafo indicando quelli che sono i principali oggetti della ricerca nel campo degli Interface Agents:

- Dimostrare che la conoscenza acquisita tramite tali agenti sia davvero utile a ridurre il lavoro degli utenti e dimostrare, invece, che tali utenti, davvero li vogliono come collaboratori
- Portare a compimento centinaia di esperimenti usando varie tecniche di apprendimento per capire quale sia la migliore su un certo dominio e perchè

- Analizzare gli effetti delle diverse tecniche di apprendimento sulle risposte fornite dagli agenti agli utenti
- Dare la possibilità a tali agenti di negoziare con altri alla pari
- Estendere il campo di applicazione di tali agenti in nuove aree innovative come l'intrattenimento, come gli agenti ALIVE e HOMR stanno facendo

Al di là di tali osservazioni, non si può negare il fatto che gli Interface Agents sono destinati ad un continuo sviluppo e ad una continua applicazione poichè sono semplici, lavorano in domini limitati e, in generale, non richiedono cooperazione con altri agenti.

Mobile Agents

Gli **Agenti Mobili** [2] sono processi software in grado di muoversi all'interno di **WAN (Wide Area Networks)** oppure nel **WWW (World Wide Web)**; sono capaci di interagire con altri hosts, di ottenere eventuali informazioni e di ritornare "alla base" dopo aver portato a termine i loro compiti.

Tali compiti possono variare, per esempio, dalla prenotazione di un volo alla gestione di una rete di telecomunicazioni. Comunque, la mobilità non è né una condizione necessaria né sufficiente per essere considerati agenti.

Gli Agenti Mobili sono agenti in quanto godono delle proprietà di autonomia e cooperazione, sebbene in maniera differente dai Collaborative Agents. Per esempio possono cooperare e comunicare con altri agenti rendendo loro noti la locazione di qualche oggetto interno oppure i loro metodi; nel fare ciò un agente scambia dei dati o qualche informazione con altri senza necessariamente fornire tutte le sue informazioni.

L'ipotesi base di questo tipo di agenti è che essi non devono essere stazionari: ciò, in certe applicazioni, può portare dei benefici di tipo non funzionale rispetto all'utilizzo di agenti statici.

Per esempio [21], si consideri il caso in cui si voglia scrivere un'applicazione che permetta da casa di poter prenotare un volo accedendo ai diversi database di sistema delle compagnie aeree nei quali si gestisce tale operazione. L'utente dovrebbe immettere dati quali: classe, fascia oraria di partenza, città di partenza e arrivo, numero di scali; un agente di tipo statico, per fare ciò, potrebbe richiedere ai vari database un elenco di tutti i voli compresi nella fascia oraria impostata e questo potrebbe costare in termini di kilobytes.

Tale agente potrebbe richiedere ancora altre informazioni quali una lista di tutti i collegamenti tra le due città impostate e procedere sino a restringere la ricerca. Tutte queste azioni coinvolgono una serie di operazioni spesso inutili che non fanno altro che affollare la rete, rendere più lunga l'operazione con eventuali costi di collegamento alla rete che aumentano.

Lo scenario opposto è quello che prevede l'uso di Agenti Mobili. L'utente potrebbe incapsulare, in uno stile orientato agli oggetti, l'intera applicazione in un agente che occuperebbe probabilmente meno di 2 kilobytes. Tale agente provvede a muoversi attraverso i sistemi di prenotazione dei voli delle varie compagnie, interroga localmente i loro database e ritorna all'utente una semplice scheda con il volo scelto che l'utente può accettare o rifiutare.

Si capisce come ciò possa ridurre i costi di collegamento e i kilobytes di informazioni inutili scaricati localmente. In sostanza, tali agenti offrono una serie di vantaggi di tipo pratico, sebbene non funzionale che li distinguono dalle loro controparti

di tipo statico. Tali vantaggi possono essere:

- Riduzione dei costi di comunicazione dovuta all'eliminazione di informazioni inutili da scaricare localmente. L'agente restituisce all'host locale solo i dati necessari selezionandoli direttamente sull'host remoto. Ciò è un grosso vantaggio nel caso la connessione alla rete venga pagata in base al tempo di utilizzo
- Possibilità di realizzare operazioni in modo asincrono: mentre l'agente lavora, l'utente può fare altre cose o, addirittura, spegnere il terminale ricevendo i risultati in un secondo momento nella propria mailbox
- Possibilità di superare i limiti dovuti a risorse locali limitate; infatti la capacità di calcolo e di immagazzinamento di una macchina può essere molto limitata ma ciò può essere superato tramite l'uso degli Agenti Mobili
- Possibilità di realizzare architetture distribuite molto potenti

Un'applicazione di tali agenti mobili è l'architettura **Telescript** [21][22] mostrata in figura 2.5.

Telescript è un linguaggio di programmazione interpretato, object-oriented e remoto per la scrittura di applicazioni distribuite. L'ambiente di sviluppo è rappresentato dal Telescript engine integrato tramite API al sistema operativo.

La primitiva più importante è rappresentata dalla "go"; essa permette la comunicazione tra i processi: due o più processi-agente possono comunicare per usufruire dei servizi dell'altro instaurando un canale di comunicazione.

La primitiva richiede uno spazio di destinazione ed il Telescript engine incapsula l'agente con tutti i suoi dati e lo invia; la destinazione può anche trovarsi su una

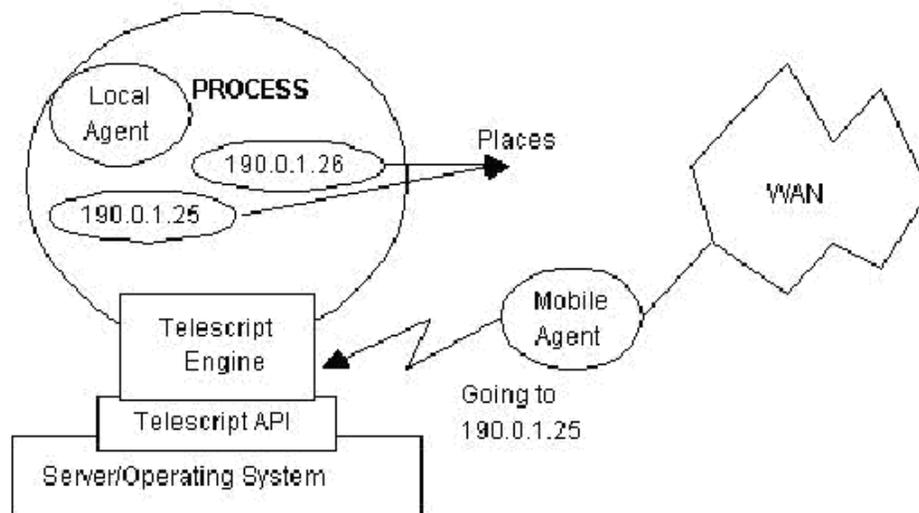


Figure 2.5: Architettura Telescript

WAN (come mostrato in figura). A destinazione, l'altro Telescript engine autentica l'agente permettendogli di eseguire i suoi task e al termine l'agente ritorna sull'host di partenza.

Concludiamo il paragrafo indicando quelli che sono i principali oggetti della ricerca nel campo dei Mobile Agents [23]:

- Stabilire dei metodi affidabili di autenticazione per stabilire se un agente che intende entrare in un sistema è davvero chi dice di essere
- Stabilire dei metodi sicuri per gestire la segretezza dei dati che viaggiano tramite tali agenti
- Stabilire dei metodi per la sicurezza di un sistema in cui accede un agente per evitare che questo vada in loop consumando cicli di CPU
- Stabilire degli standard in modo tale che un agente scritto in un certo linguaggio

di programmazione possa funzionare su un motore scritto in un altro

Si evince che uno degli argomenti di ricerca più importante è rappresentato dalla sicurezza. In molte aziende sono installati opportuni firewall che impediscono l'accesso a tali tipi di agenti, indipendentemente dal fatto che essi siano o meno pericolosi.

Altra tecnica per aumentare la sicurezza di chi deve accettare un agente mobile nel proprio sistema, è quella di impedire la realizzazione di agenti che possano scrivere nella memoria o sui dischi del sistema stesso.

Information/Internet Agents

Gli **Information Agents** [2] sono nati a causa della forte domanda di applicazioni per la gestione della sempre crescente mole di dati presente oggi in rete. Essi hanno il compito di gestire, manipolare e collezionare informazioni da sorgenti distribuite.

L'ipotesi alla base degli Information Agents è che la rete è sovraccaricata da dati e informazioni e, quindi, è richiesta una migliore e più efficiente gestione di essi.

Tom Henry, vice presidente dell'azienda SandPoint [24], ha asserito che *una delle sfide più affascinanti dell'informatica è quella di realizzare un sistema dove la richiesta e l'ottenimento di informazioni utili usando tali agenti sia naturale come usare il telefono o leggere un giornale.*

In teoria, usando tali agenti ed addestrandoli opportunamente, un utente potrebbe ricevere sul proprio terminale il suo giornale preferito all'ora stabilita oppure raccogliere in maniera personalizzata tutta una serie di informazioni utili all'interno di un dominio.

Sono due i principali motivi che spingono allo sviluppo di Information/Internet Agents; il primo, come già sottolineato, è che a causa della crescente mole di dati

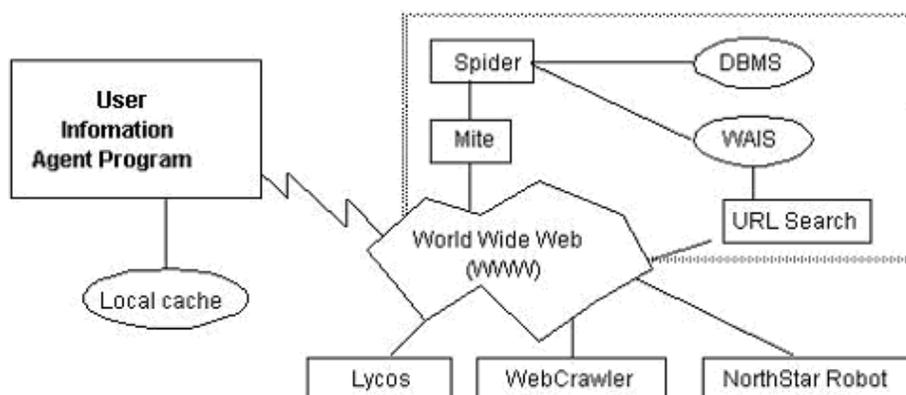


Figure 2.6: Information Agent

presenti in rete, si sente l'esigenza di avere applicazioni per gestire efficacemente questa crescente quantità.

Il secondo motivo è di natura economica: si ritiene che chi riuscirà a realizzare un sistema ad agenti di tipo proattivo, dinamico, adattativo e cooperativo per la gestione delle informazioni sul WWW, facilmente potrà arricchirsi poichè sempre maggiori sono le spese che vengono effettuate sul WEB. Soffermandoci sulle caratteristiche che tali agenti debbono possedere, si può affermare che essi in genere sono di tipo mobile, possono essere non cooperativi e possono o non apprendere.

Esistono anche casi di Information Agent di tipo non mobile, come evidenziato nella figura 2.6

L'Information Agent è integrato nel browser [24] e si appoggia a sistemi di spidering o a motori di ricerca per ottenere le informazioni desiderate.

Uno Spider è un sistema in grado di muoversi nel WEB e di memorizzare nel database di sistema la topologia dei siti visitati. L'agente, al fine di ottenere informazioni su un certo argomento, si interfaccia con i motori di ricerca o lo Spider

(oppure si interfaccia con la cache locale se una certa richiesta era già stata effettuata ed i risultati memorizzati).

Si possono citare diversi esempi di Information Agents sviluppati di recente. Etzioni & Weld [25] nel 1994 hanno sviluppato un sistema detto **Internet Softbot**: si tratta di un agente che permette di effettuare query complesse (con congiunzioni, disgiunzioni, quantificazioni) e che è capace di muoversi nel WWW al fine di inferire conoscenza e di soddisfare le richieste. È in grado di tollerare ambiguità, omissioni ed errori da parte dell'utente. È capace, inoltre, non solo di inferire conoscenza, ma anche di filtrare e-mail, schedulare appuntamenti e realizzare task di gestione del sistema su cui opera.

Un tale tipo di agente è classificato come Information Agent e non come Interface Agent poiché l'apprendimento non è un fattore cruciale, sebbene possa implementare tecniche di apprendimento elementari. Per esempio, è in grado di capire, impostata una query, se una simile è già stata sottoposta. Altro esempio è l'agente sviluppato da Davies e Weeks nel 1995 presso i BT Labs: il **Jasper Agent (Joint Access to Stored Pages with Easy Reatrieval)** [26]. Esso è in grado di immagazzinare, ritrovare e indicizzare pagine Web sulla base delle query dell'utente. In più, è capace di interagire con altri agenti presenti sul WEB al fine di migliorare le ricerche. Offre la possibilità di estrarre delle keywords da ogni pagina Web in modo tale da restituire queste pagine nel caso in cui un utente effettui delle ricerche usando le keywords estratte.

Per quanto concerne la ricerca sugli Information Agents, si può affermare che gli oggetti principali sono simili a quelli già trattati quando si è parlato degli Interface

Agents o dei Mobile Agents; infatti se l'agente è di tipo statico, ad esso si può applicare tutto quello già visto per gli Interface Agents (tecniche di apprendimento, negoziazione), mentre se è di tipo mobile si può applicare quello visto per i Mobile Agents (autenticazione, sicurezza, segretezza). Si tornerà in seguito su questi agenti, quando si parlerà più approfonditamente degli **Intelligent Search Agents (ISA)**.

Reactive Agents

Gli **Agenti Reattivi** [2] rappresentano una particolare categoria di agenti caratterizzati dal fatto che non posseggono al loro interno modelli simbolici del loro ambiente; al contrario, essi agiscono in uno stile sollecitazione/risposta allo stato presente dell'ambiente in cui si trovano. Si tratta di agenti in genere relativamente semplici ed in grado di interagire con altri in maniera elementare.

Tre sono le ipotesi principali alla base di tali agenti [27]; la prima è la *semplicità*, in quanto non c'è nessuna specificazione a priori del comportamento di un insieme di agenti reattivi.

La seconda è la possibilità di decomporre in task un'applicazione; infatti un agente reattivo è visto come una collezione di moduli ognuno operante in maniera autonoma e responsabile di specifici task (rilevazioni, computazioni). La comunicazione tra i moduli è ridotta al minimo e di basso livello.

La terza è che tali agenti lavorano su semplici modelli dell'ambiente circostante che sono quelli forniti, per esempio, da sensori in contrasto con modelli simbolici di alto livello adottati negli altri agenti discussi in precedenza.

L'elemento chiave del successo di tali agenti, è proprio questo nuovo modo di vedere l'ambiente circostante; infatti, i principali teorici nel campo dell'intelligenza

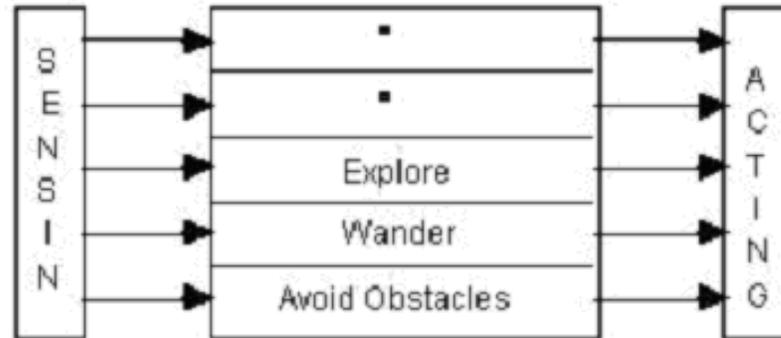


Figure 2.7: Brooks Subsumption Architecture

artificiale sostenevano che per costruire un sistema intelligente si dovesse fornire una rappresentazione simbolica del mondo, mentre ora, questa viene bandita insieme ai modelli di ragionamento simbolico [28]. L'ambiente, ora, viene visto "per quello che è" tramite una serie di sensori ed attuatori. In tal modo gli agenti reattivi si dimostrano semplici, facili da comprendere e con costi di memorizzazione ridotta dato che devono "ricordare" poco. Essi non programmano nulla in anticipo, ma le loro azioni dipendono da quello che accade in quel momento [6].

I benefici che spingono ad uno sviluppo di sistemi basati su agenti reattivi, oltre a quelli già menzionati, sono:

- Maggiore robustezza e tolleranza ai guasti: un agente può danneggiarsi senza effetti catastrofici
- Maggiore flessibilità e adattabilità, in contrasto alla rigidità, agli alti tempi di risposta e all'instabilità dei classici sistemi di intelligenza artificiale

Un architettura basata su agenti di tipo reattivo è la **Brooks Subsumption Architecture** [29] raffigurata in figura 2.7.

Si tratta di un'architettura usata per costruire robot meccanici: essa è formata da una serie di moduli descritti come una macchina a stati finiti. Un modulo entra in azione se il suo segnale di input supera una soglia fissata. Tali moduli sono le uniche unità di processazione nell'architettura, in quanto non ci sono simboli come quelli usati nei classici sistemi di Intelligenza Artificiale.

I moduli sono organizzati a livelli e ciascun modulo ha uno scopo ben preciso (per esempio, evitare ostacoli o abilitare e controllare il movimento).

Un'applicazione interessante è quella della simulazione artificiale. Ferber [6] nel 1994 realizzò un'applicazione basata su agenti reattivi in grado di simulare colonie di formiche; ogni formica era modellata tramite un agente ed egli riuscì a simulare un piccolo ecosistema. Altro simulatore è quello realizzato da Nwana [30] nel 1993 dove si simulava un gruppo di bambini in un campo da gioco.

Perciò gli agenti reattivi possono far diventare il computer un "laboratorio virtuale" dove i ricercatori possono far variare molteplici parametri e validare i più svariati modelli.

Un altro campo di applicazione è rappresentato dall'industria dei giochi e dell'intrattenimento, come dimostrano le ricerche effettuate negli ultimi anni nei laboratori della Philips [31].

Concludiamo il paragrafo indicando quelli che sono gli oggetti principali di ricerca nel campo degli agenti reattivi:

- Cercare di ampliare la tipologia ed il numero di applicazioni basate su tali agenti in quanto attualmente esse sono concentrate principalmente ai simulatori e ai

giochi

- Sviluppare un'adeguata metodologia di progettazione delle applicazioni e, quindi, teorie, architetture e linguaggi
- Affrontare in maniera ingegneristica i problemi di scalabilità, estendibilità e valutazione delle performance

Comunque, a discapito di queste osservazioni finali, c'è da dire che ci si aspetta uno sviluppo notevole nei prossimi anni di applicazioni basate su questo tipo di agenti.

Hybrid Agents

Sino ad ora si sono analizzati cinque tipi di agenti: collaborativi, di interfaccia, mobili, di Internet e i reattivi; ampio è il dibattito su quale tra essi sia il migliore. Poiché ciascuno di essi ha i suoi pregi ed i suoi difetti, si è cercato di massimizzare i primi e ridurre i secondi adattando al meglio ciascuno di essi ai propri scopi.

Un modo per massimizzare i propri obiettivi è quello che comunemente viene denominato approccio ibrido [2] che cerca di unire i pregi sia degli agenti basati sul paradigma deliberativo (e quindi i collaborativi, quelli di interfaccia, i mobili e di internet) sia di quelli basati sul paradigma reattivo. Perciò, con il termine **Agente Ibrido** ci si riferisce ad agenti che fondono caratteristiche di due o più tipi di agenti in uno solo.

L'ipotesi alla base di tali agenti è che spesso, in diverse applicazioni, si possono ottenere performance migliori combinando diverse filosofie in un unico agente piuttosto che affidarsi ad un unico agente che si basi su una singola filosofia tra quelle già viste. Ovviamente, i benefici derivanti dall'uso di tale tipo di agenti nascono dall'unione dei benefici dei diversi tipi di agenti adottati. Per esempio, si supponga di realizzare un agente che unisca le caratteristiche di collaborazione (e, quindi, un agente di tipo deliberativo) e di reazione; la componente reattiva, che potrebbe avere la precedenza sull'altra, potrebbe portare i seguenti benefici: robustezza, adattabilità e tempi di risposta minori. La componente deliberativa potrebbe occuparsi degli scopi a lungo termine, aumentando la flessibilità del sistema. In genere i sistemi che si basano su un

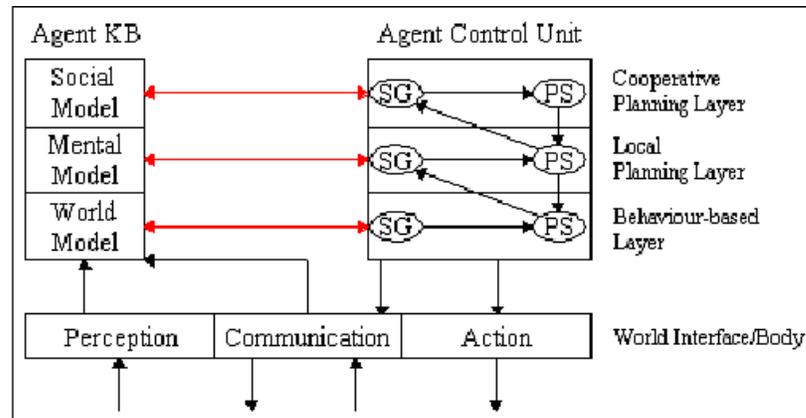


Figure 2.8: The InteRRaP Hybrid Architecture

approccio ibrido presentano un'architettura a livelli; un esempio è il sistema sviluppato presso il Centro di Ricerca Tedesco sull'Intelligenza: **The InteRRaP Hybrid Architecture** [32][33]. Esso è rappresentato in figura 2.8. Si tratta di una struttura a livelli per la realizzazione di agenti e robot autonomi e che sposa la filosofia reattiva con quella deliberativa.

La parte reattiva (*Behaviour-based Layer - BBL*) del sistema è basata su agenti con regole di tipo stimolo/reazione e che contribuiscono alle proprietà di efficienza e robustezza. Il livello intermedio (*Local Planning Layer - LPL*) è costituito da agenti con comportamenti di tipo goal driven mentre l'ultimo livello (*Cooperative Planning Layer - CPL*) è costituito da agenti che permettono la cooperazione con altri in modo da realizzare una struttura multi-agente.

Un'altra architettura è la **Touring Machine** [34] sviluppata da Ferguson nel 1992; si tratta ancora di un sistema a livelli: un livello *reattivo*, uno di *pianificazione* ed un altro di *modellazione*.

Si tratta di un'architettura di tipo orizzontale, al contrario di quella precedentemente esposta che è di tipo verticale; infatti nella prima tutti i livelli hanno accesso ai sensori e possono contribuire all'azione, mentre nell'applicazione vista in precedenza solo il livello più basso riceve gli stimoli dai sensori e agisce.

Nonostante i molti vantaggi che può portare un'architettura basata su agenti ibridi, sono ancora relativamente poche quelle realmente progettate; infatti due delle critiche mosse è che un approccio di tal genere può portare a sistemi troppo specifici e poco generali e che l'approccio teorico e ingegneristico non è ben specificato in letteratura. Perciò, gli obiettivi dei ricercatori in questo campo sono molto simili a quelli dei ricercatori nel campo degli agenti di tipo reattivo.

Inoltre ci si aspetta di vedere sistemi ibridi che non uniscano solo le due filosofie deliberativa/reattiva ma anche altre; per esempio combinare in un unico agente, le caratteristiche degli Interface Agents e dei Mobile Agents, entrambi appartenenti alla categoria degli agenti di tipo deliberativo.

Heterogeneous Agents System

I sistemi basati su Agenti Eterogenei [2], a differenza dei sistemi ibridi descritti in precedenza, sono dei sistemi integrati di due o più tipi di agenti tra quelli discussi nel capitolo.

Mentre gli agenti ibridi combinavano filosofie diverse in un'unica struttura, ora vengono combinati direttamente diverse tipologie di agenti.

Uno dei principali motivi che hanno portato allo sviluppo di tali sistemi è l'ampia raccolta di prodotti software ciascuno dei quali offrono una vasta quantità di servizi [35]. Sebbene essi lavorino in modalita stand alone, c'è una crescente richiesta di

interoperabilità tra di essi nella speranza di ottenere un valore aggiunto dalla loro collaborazione. Al fine di ingegnerizzare e standardizzare l'interoperabilità tra diversi agenti software è nata una nuova branca dell'informatica: l'*Agent-Based Software Engineering*.

Uno dei suoi principali obiettivi è stato quello di avere un linguaggio di comunicazione tra agenti (*Agent Communication Language - ACL*) comune in modo tale che agenti software differenti potessero comunicare in una piattaforma comune.

Diversi sono i benefici che si possono ottenere dalla realizzazione di tali sistemi [35]:

- Le applicazioni di tipo stand alone possono aumentare il numero dei servizi offerti inserendole in un contesto eterogeneo
- I legacy system possono vedere ridotti i costi dovuti alla reingegnerizzazione del software se operanti in collaborazione con altri sistemi
- Nuovi orizzonti nella ricerca della progettazione software e nella sua implementazione e manutenzione

Da quanto detto sino ad ora, si evince che uno dei problemi da affrontare per avere sistemi ad agenti eterogenei è rappresentato dal protocollo di comunicazione. Sono stati studiati diversi meccanismi al fine di facilitare ciò:

- Affiancare ad ogni agente un trasduttore il quale riceve i messaggi dagli altri agenti e li traduce in un linguaggio comprensibile al suo agente
- Utilizzare opportuni wrappers che, come affermano Genesereth e Ketchpel [35], "iniettino codice in un programma per consentirgli di comunicare". In tal modo

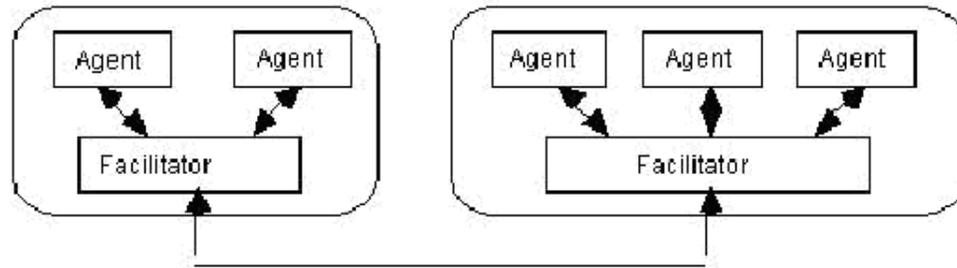


Figure 2.9: Sistema ad agenti eterogenei

ogni agente ha associato un pezzo di codice che è espresso in un comune linguaggio usato dagli altri agenti

- Riscrivere da zero il codice degli agenti; tale alternativa si dimostra essere molto dispendiosa
- Usare, come proposto da Wiederhold [36], un apposito mediatore con il quale si interfacciano tutti gli agenti del sistema e provvede allo smistamento dei messaggi

Nella figura 2.9 è mostrato un esempio di sistema basato su agenti eterogenei [35].

Nel sistema sono presenti cinque agenti distribuiti su due macchine, una con due agenti e l'altra con tre.

Gli agenti non comunicano tra di loro in maniera diretta, bensì tramite dei facilitator che fungono da opportuni mediatori. Gli agenti cedono una parte della loro autonomia ai mediatori i quali sono in grado di mettersi in contatto con gli altri agenti nella rete per richiederne i servizi. I mediatori si occupano anche di stabilire le connessioni nella rete e di assicurare una corretta comunicazione tra gli agenti.

Un prototipo che si basa sull'architettura presentata nella precedente figura è il **PACT**: *Palo Alto Collaborative Testbed* [37]. Si tratta di un sistema per integrare quattro legacy system in un'architettura comune. Sono coinvolti trentuno agenti che girano su quindici workstations e microcomputers. Gli agenti sono organizzati gerarchicamente e comunicano tramite opportuni mediatori.

Il lavoro di sviluppo di sistemi basati su agenti eterogenei è in continua evoluzione ed è molto sentita dai ricercatori l'esigenza di metodologie, tools, tecniche e standards al fine di raggiungere gli obiettivi di interoperabilità tra fonti di informazioni eterogenee. Gli oggetti principali di ricerca sono:

- Fissare un comune linguaggio di comunicazione tra gli agenti
- Stabilire in che modo gli agenti possono comunicare sfruttando un eventuale linguaggio comune
- Fissare degli standards con i quali progettare architetture che permettano l'interoperabilità tra gli agenti

Nel seguito si riprenderanno alcuni dei concetti ora esposti al fine di spiegare nel dettaglio cosa si intende per **Agente di Ricerca Intelligente (Intelligent Search Agent - ISA)** e verranno presentati gli agenti progettati nella nostra architettura².

²Non si è descritto il tipo di agente ombreggiato in figura 2.2 (*Smart Agents*) poiché esso rappresenta un'aspirazione dei ricercatori piuttosto che una realtà. Ricordiamo che un tale tipo di agente è uno che gode contemporaneamente degli attributi di Autonomia, Apprendimento e Cooperazione.

2.2 Ontologia

Un'ontologia è la specificazione esplicita di una concettualizzazione [55]. La parola "ontologia" sembra generare controversie nelle discussioni riguardanti l'intelligenza artificiale. È possibile definire un'ontologia come una descrizione dei concetti e delle relazioni che intercorrono tra di essi. Essa permette di lavorare con un insieme strutturato di concetti, dove le relazioni risultino chiare e, soprattutto, significative a livello semantico. Una ontologia può essere specificata informalmente mediante proposizioni in linguaggio naturale. Una ontologia formale è invece specificata da una raccolta di nomi per i concetti che si vogliono specificare e dalle relazioni che intercorrono tra di loro.

Le ontologie sono spesso considerate alla stregua delle gerarchie tassonomiche di classi, ma le prime, nella visione di T. Gruber [55] (studioso dell'università di Stanford), non si limitano a dare definizioni conservative, ovvero definizioni nel senso tradizionale di introdurre terminologia senza aggiungere conoscenza a proposito del mondo. Infatti, per mantenere un sistema "aperto", l'ontologia non è basata su una gerarchia fissa di categorie, ma su un insieme di distinzioni, a partire dalle quali la gerarchia è generata automaticamente. È importante considerare gli obiettivi di un'ontologia. Se ne progetta una per permettere la condivisione della conoscenza e il suo riutilizzo. Pragmaticamente, si sceglie di scrivere un'ontologia come un insieme strutturato di definizioni di un vocabolario formale. Sebbene questo non sia l'unico modo di specificare una concettualizzazione, esso ha l'interessante proprietà di permettere la condivisione della conoscenza tra software di intelligenza artificiale. È sostanzialmente un accordo sull'uso di un certo vocabolario (ovvero effettuare interrogazioni e fare affermazioni) in modo che esso sia consistente (ma non completo)

rispetto alla teoria specificata dall'ontologia. A partire dalle ontologie, si possono costruire agenti che fanno affidamento su di esse. Le ontologie sono quindi progettate allo scopo di condividere la conoscenza con e tra gli agenti. Non è richiesto comunque che un tale agente risponda a tutte le domande che possono essere formulate nel vocabolario condiviso. Infatti, un accordo sull'uso di una ontologia comune è una garanzia di consistenza ma non di completezza rispetto alle interrogazioni e alle affermazioni che possono essere fatte usando il vocabolario definito dall'ontologia stessa.

2.2.1 Semantic Network

Dati un insieme di concetti³ in qualche modo correlati tra loro, nasce la necessità di definire una **Semantic Network - Rete Semantica**, per poter effettuare una analisi semantica soddisfacente.

Definition 2.2.1 (Semantic Network). *Dati un insieme di n concetti C_1, \dots, C_n si definisce Semantic Network il grafo pesato i cui nodi rappresentano i concetti C_i e gli archi non orientati le relazioni tra essi. I pesi associati agli archi rappresentano il **grado di correlazione** tra i concetti interessati.*

Le relazioni rappresentate dagli archi possono essere di similitudine, di generalizzazione/specializzazione o di semplice attinenza.

I pesi associati agli archi costituiscono una misura del grado di parentela tra i concetti collegati dall'arco stesso.

³Una definizione di concetto sarà data successivamente quando verrà illustrato il modello formale al quale il sistema fa riferimento

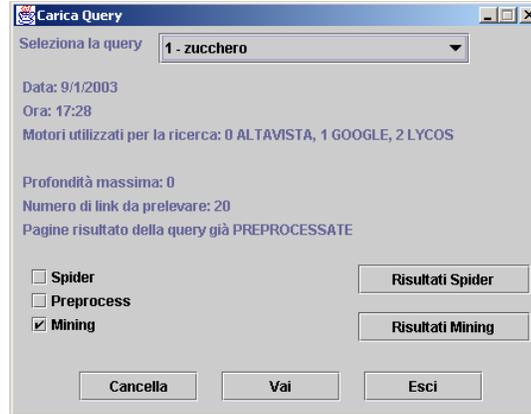


Figure 2.10: Legame tra due concetti

Più precisamente, con riferimento alla figura 2.10, considerati due concetti C_1 e C_2 tali che esistano due archi diretti rispettivamente da C_1 a C_2 e da C_2 a C_1 :

- P_{12} è la probabilità che il concetto C_2 possa soddisfare una richiesta del concetto C_1
- P_{21} è la probabilità che il concetto C_1 possa soddisfare una richiesta del concetto C_2

Definiti i pesi sugli archi, possiamo dare, con riferimento alla figura 2.11, la definizione di distanza probabilistica tra due concetti.

Definition 2.2.1 (Distanza probabilistica). Dato un insieme di n concetti C_1, \dots, C_n e considerati due di tali *concetti* C_i e C_j , collegati da uno o più percorsi, si definisce **distanza probabilistica tra i concetti C_i e C_j** la quantità

$$DP(C_i, C_j) = \max_{k \in 1..n} \left(\prod_{t=1}^{na} P_{k,t} \right) \quad (2.2.1)$$

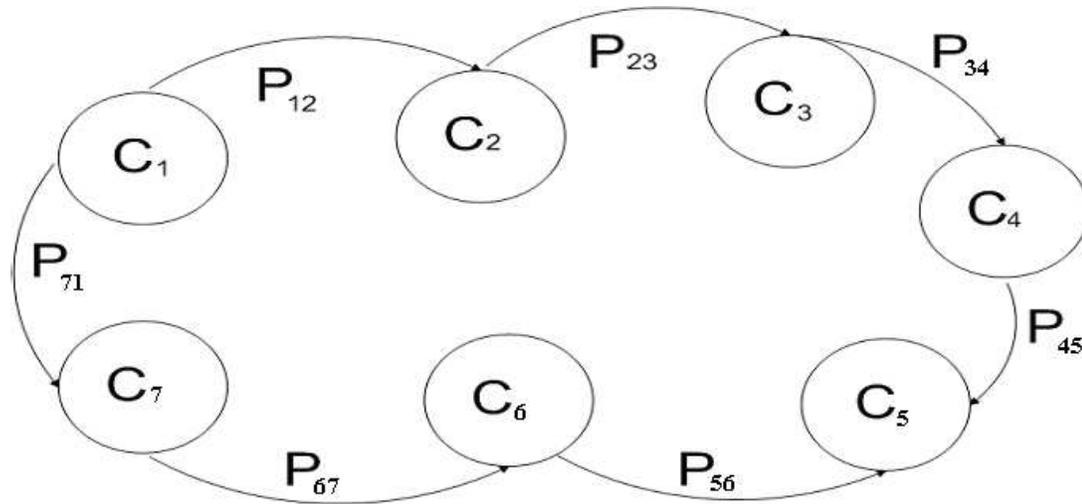


Figure 2.11: Distanza probabilistica tra due concetti

dove n è il numero di percorsi esistenti tra i concetti C_i e C_j , n_k è il numero di archi del k -simo percorso e $P_{k,t}$ è il peso associato al t -simo arco del k -simo percorso.

Si noti che $DP(C_i, C_j)$ può ancora essere interpretata come la probabilità che il concetto C_j possa soddisfare una richiesta del concetto C_i e coincide con P_{ij} se esiste un arco diretto da C_i a C_j e non esistono tra gli stessi concetti percorsi migliori di quello diretto. Due concetti sono tanto più vicini dal punto di vista semantico quanto più la loro distanza probabilistica è prossima ad 1.

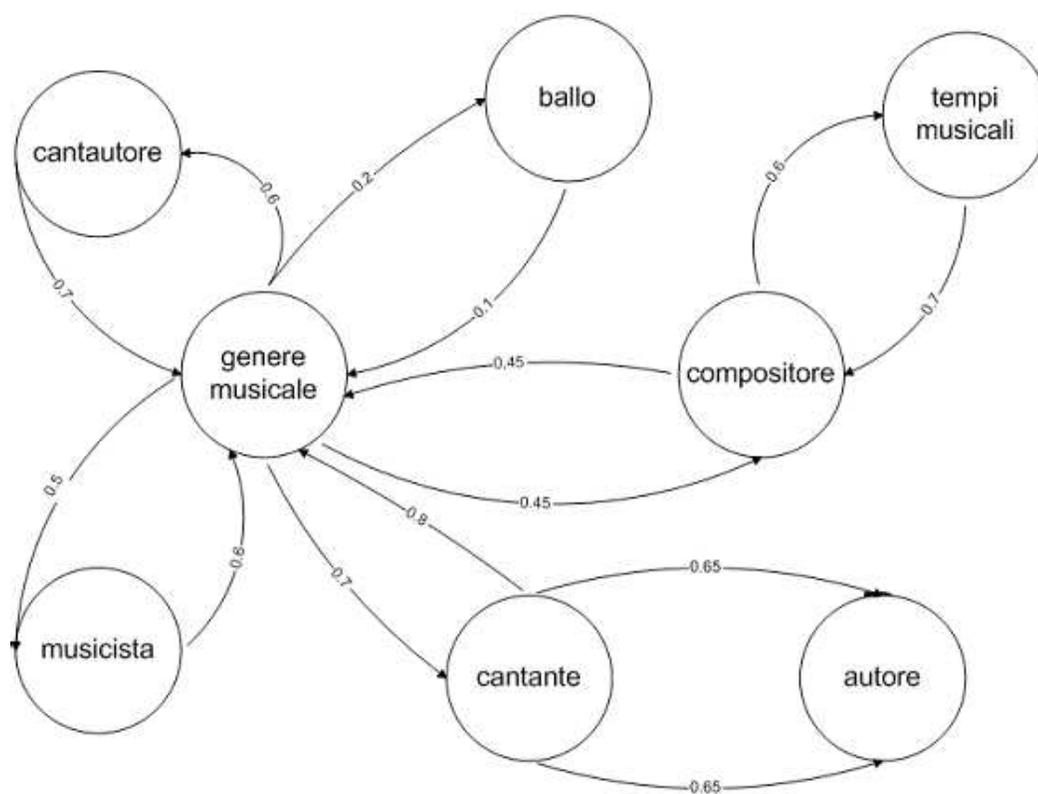


Figure 2.12: Una rete semantica per i concetti del dominio musicale

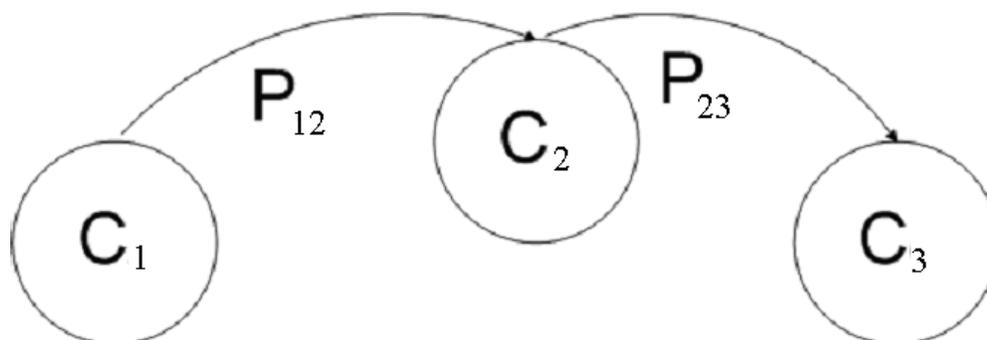


Figure 2.13: Distanza probabilistica lungo un percorso

Example 2.2.1. *Con riferimento alla figura 2.12, che rappresenta una possibile rete semantica specializzata sul dominio musicale, alcune delle distanze probabilistiche individuabili tra i concetti sono:*

- $DP(\text{genere musicale}, \text{cantante}) = 0.8$
- $DP(\text{cantante}, \text{genere musicale}) = 0.7$
- $DP(\text{cantautore}, \text{compositore}) = 0.7 \cdot 0.45 = 0.315$
- $DP(\text{tempi musicali}, \text{musicista}) = 0.7 \cdot 0.45 \cdot 0.5 = 0.1575$
- $DP(\text{musicista}, \text{tempi musicali}) = 0.6 \cdot 0.45 \cdot 0.6 = 0.162$

Facciamo ora alcune considerazioni per giustificare la distanza introdotta. Motiviamo prima di tutto la presenza del simbolo di produttoria. Con riferimento al semplice esempio di figura 2.13, la probabilità che il concetto C_3 possa soddisfare una richiesta del concetto C_1 è pari alla probabilità dell'evento " C_3 soddisfa una richiesta del concetto C_2 e C_2 soddisfa una richiesta del concetto C_1 ", ovvero alla probabilità

dell'intersezione di due eventi statisticamente indipendenti. Com'è noto dalla teoria della probabilità, dati due eventi statisticamente indipendenti A e B risulta:

$$P(A \cdot B) = P(A) \cdot P(B)$$

Inoltre bisogna chiarire perché si assume come distanza tra due concetti la probabilità associata al percorso probabilisticamente più corto, ossia il percorso cui è associata la distanza probabilistica più prossima ad 1. Se tra due concetti esiste più di un percorso, il percorso probabilisticamente più corto è chiaramente quello che esprime le relazioni più esplicite tra i concetti in questione e risulta dunque naturale assumere la probabilità ad esso associata come misura di similitudine tra gli stessi.

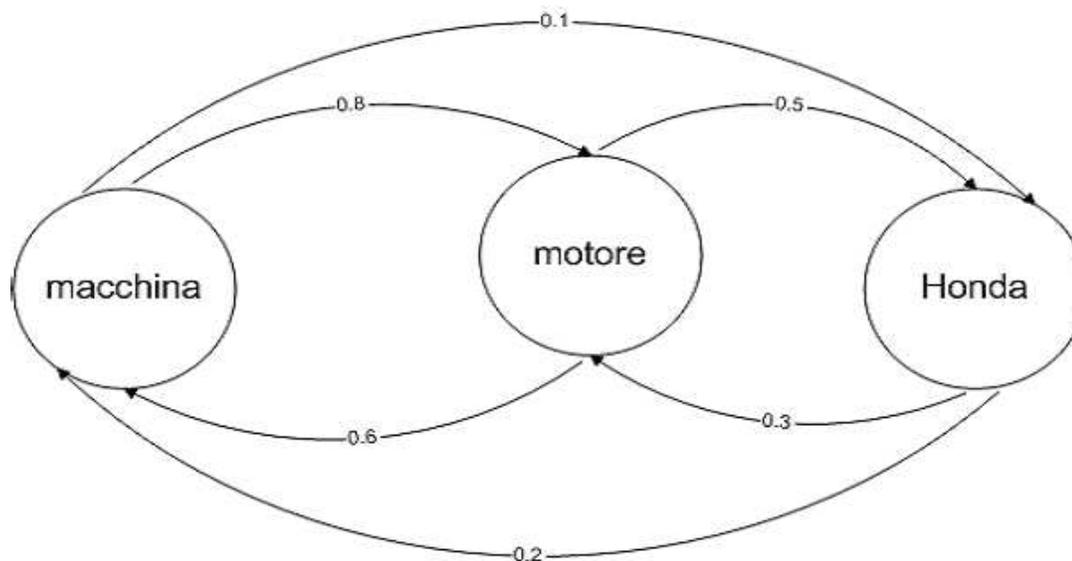


Figure 2.14: Legami tra alcuni concetti del dominio automobilistico

Example 2.2.2. *Con riferimento alla figura 2.14, che rappresenta i legami tra alcuni concetti del dominio automobilistico, si ha che*

- $DP(\text{macchina}, \text{Honda}) = 0.8 \cdot 0.5 = 0.4$
- $DP(\text{Honda}, \text{macchina}) = 0.2$

Si noti come nel primo caso il percorso diretto macchina-Honda non sia quello probabilisticamente più corto. Nel secondo caso accade, invece, l'esatto contrario.

Per concludere, osserviamo che la *distanza probabilistica* soddisfa le seguenti proprietà:

- $DP(C_i, C_i) = 1$ $\forall C_i \in \text{insieme dei concetti}$
- $DP(C_i, C_j) \neq DP(C_j, C_i)$
- $DP(C_i, C_j) \neq DP(C_i, C_k) \cdot DP(C_k, C_j)$ $\forall C_i, C_j, C_k \in \text{insieme dei concetti}$
- $DP(C_i, C_j) = 0$ $\iff \text{non esiste percorso tra } C_i \text{ e } C_j$

2.3 Modello formale

Consideriamo un insieme di elementi A: nella *teoria classica degli insiemi* l'appartenenza di ogni elemento x ad A è indicata in maniera certa. Ci sono però dei casi in cui la teoria classica non riesce a classificare in maniera adeguata gli elementi di un sottoinsieme dell'universo del discorso. In questi casi può essere opportuno non

indicare la semplice appartenenza di un elemento all'insieme, ma la probabilità percentuale che l'elemento appartenga o meno all'insieme: questa probabilità percentuale di appartenenza sarà un valore compreso nell'intervallo $[0,1]$ dei numeri reali ed ovviamente l'appartenenza certa si avrà con la probabilità percentuale pari ad 1, mentre in corrispondenza della probabilità percentuale pari a 0 si avrà la non appartenenza all'insieme.

2.3.1 Domini, concetti e parole

Iniziamo con il definire cosa sono un vocabolo ed una lingua.

Definition 2.3.1 (Vocabolo). *Si definisce vocabolo un insieme di caratteri isolato, colto fuori da un contesto e da ogni legame grammaticale o logico, nella sua individualità lessicale [62].*

Example 2.3.1. *Alcuni vocaboli sono, ad esempio:*

Riserva, House, Maman, Pfad, Elicóptero,...

I vocaboli sono la parte fondamentale di una lingua.

Definition 2.3.2 (Lingua). *Si definisce una lingua come l'insieme dei vocaboli utilizzati da un popolo, insieme alle regole che li governano [62].*

Example 2.3.2. *Alcune possibili lingue possono essere:*

ITALIANO, TEDESCO, INGLESE, SPAGNOLO,...

Definition 2.3.3 (Regole sintattiche). *Si definiscono regole sintattiche (o sintassi) l'insieme delle norme che studiano le relazioni che i vocaboli hanno in una frase (gruppo di vocaboli) [62].*

Definition 2.3.4 (Regole semantiche). *Si definiscono regole semantiche (o semantica) l'insieme delle norme che si occupano del significato dei vocaboli e dei cambiamenti di essi [62].*

Ci occuperemo di seguito della lingua Italiana e delle sue regole sintattiche e semantiche. Definiamo ora cosa è un dominio ed un concetto.

Definition 2.3.5 (Dominio). *Detto S l'universo del discorso si definisce dominio D un sottoinsieme di S , ovvero*

$$D \subseteq S.$$

Example 2.3.3. *Possibili esempi di domini, nella lingua Italiana, potrebbero essere il dominio della MUSICA, il dominio della PITTURA, ecc.*

Definition 2.3.6 (Concetto). *Per concetto linguistico si intende una parola, frase, acronimo o nome ricco di significato, che è stato estratto da componenti non strutturati del testo, incluso blocchi isolati di testo, sommari, intestazioni, paragrafi etc. Esso è definito anche come termine linguistico, mentre ogni documento si riferisce a un gruppo di concetti o termini [61].*

Dato un dominio D e detti C_i i concetti in esso definiti si ha che l'unione di tutti i concetti è una copertura per il dominio stesso, ovvero

$$D = \bigcup_{i=1}^n C_i$$

Example 2.3.4. *Se consideriamo il dominio della MUSICA, possibili esempi di concetti di tale dominio potrebbero essere il concetto STRUMENTI, il concetto CANTANTI, ecc.*

Ogni concetto è costituito da **vocaboli** della lingua alla quale si ci riferisce. Possiamo, però considerare che i vocaboli non appartengano ad un concetto necessariamente con una probabilità pari ad 1. Soprattutto nei casi in cui si vuole rappresentare un dominio con un insieme sintetico di concetti, può accadere che alcuni vocaboli facenti parte del dominio, siano inseriti in concetti ai quali non appartengono pienamente. Inoltre un vocabolo può appartenere a più concetti e pare opportuno, in questi casi, immaginare che non appartenga a pieno ad entrambi i concetti.

Example 2.3.5. *Consideriamo ad esempio i due vocaboli della lingua italiana: basso*

e fisarmonica. Possiamo immaginare che fisarmonica appartenga al concetto strumento del dominio musicale con una probabilità maggiore rispetto a basso, in quanto quest'ultimo crea maggiori ambiguità.

Per tale motivo introduciamo il concetto di centralità di un vocabolo rispetto ad un concetto.

Definition 2.3.7 (Centralità). *Definiamo centralità la probabilità percentuale di appartenenza del vocabolo lessicale ad un concetto.*

Appare evidente che più il valore della centralità è prossimo ad 1 e più è probabile che il vocabolo appartenga al concetto associato.

Definition 2.3.8 (Peso). *Definiamo peso il grado con cui un vocabolo lessicale rappresenta il dominio al quale è associato.*

Anche il peso è un valore appartenente all'intervallo $[0, 1]$ e più è grande, maggiore è il livello di rappresentazione, da parte del vocabolo, del dominio associato.

Definition 2.3.9 (Parola). *Definiamo parola un vocabolo di una lingua, a cui è associato un concetto C, il dominio D a cui il concetto appartiene e due indici centralità e peso, quindi una parola è la quadrupla*

<vocabolo lessicale, concetto, dominio, centralità, peso>

Example 2.3.6 (Parole). *Possibili esempi di parole possono essere i seguenti.*

<i>Vocabolo</i>	<i>Concetto</i>	<i>Dominio</i>	<i>Centralità</i>	<i>Peso</i>
<i>Chitarra</i>	<i>Strumenti</i>	<i>Musica</i>	<i>0.8</i>	<i>0.6</i>
<i>Olio</i>	<i>Ricambi</i>	<i>Motori</i>	<i>0.6</i>	<i>0.55</i>
<i>Aereo</i>	<i>Mezzi di trasporto</i>	<i>Turismo</i>	<i>0.9</i>	<i>0.7</i>
<i>Presentatore</i>	<i>Televisione</i>	<i>Spettacolo</i>	<i>0.85</i>	<i>0.6</i>
<i>Bacio</i>	<i>Amore</i>	<i>Sentimenti</i>	<i>0.85</i>	<i>0.9</i>

Indicato con S sempre l'universo del discorso, è possibile anche definire una **funzione F** tale che data una *parola* x ed un concetto C appartenente ad un dominio D , il valore

$$F_C(x) \in [0, 1]$$

rappresenta la centralità di x rispetto al concetto C , con la convenzione che se il vocabolo inserito nella parola x non appartenga al concetto C la funzione F restituisca 0 ⁴.

In maniera analoga alla centralità è possibile definire una **funzione P** tale che dato una *parola* x ed un dominio D , il valore

$$P_D(x) \in [0, 1]$$

rappresenta il peso di x rispetto al Dominio D .

⁴In questo caso è come dire che il vocabolo associato alla parola x appartiene al concetto con una probabilità nulla, cioè non vi appartiene.

2.3.2 Dominio unico

Consideriamo innanzitutto l'esistenza di un solo **dominio** e diamo in tal caso la definizione di insieme monoconcetto e multiconcetto.

Definition 2.3.10 (Insieme monoconcetto). Dato il **dominio** D ed un suo **concetto** C , l'insieme di parole A si definisce **monoconcetto** se

$$\forall x \in \mathcal{A} \rightarrow F_C(x) \neq 0. \quad (2.3.1)$$

Example 2.3.7. Un possibile esempio di insieme monoconcetto A , riferito al concetto *Repulsione* appartenente al dominio *Sentimenti* è illustrato nella figura 2.15

	<i>Detestare</i>	<i>Repulsione</i>	<i>Sentimenti</i>	<i>0.9</i>	<i>0.9</i>	
	<i>Disprezzo</i>	<i>Repulsione</i>	<i>Sentimenti</i>	<i>0.6</i>	<i>0.7</i>	
$A = <$	<i>Avversione</i>	<i>Repulsione</i>	<i>Sentimenti</i>	<i>0.9</i>	<i>0.5</i>	$>$
	<i>Scontentezza</i>	<i>Repulsione</i>	<i>Sentimenti</i>	<i>0.3</i>	<i>0.3</i>	
	<i>Grasso</i>	<i>Repulsione</i>	<i>Sentimenti</i>	<i>0.2</i>	<i>0.1</i>	

Figure 2.15: Esempio di insieme monoconcetto

Definition 2.3.11 (Insieme multiconcetto). Dato il dominio D ed i concetti C_1, \dots, C_n di D , l'insieme di parole A è un insieme **multiconcetto** se

$$\forall x \in \mathcal{A} \rightarrow \exists C_i, \text{ con } i = 1..n : F_{C_i}(x) \neq 0. \quad (2.3.2)$$

dove la funzione $F_{C_i}(x)$, come detto in precedenza, restituisce la centralità della parola x rispetto al concetto C_i .

Example 2.3.8. *Un possibile esempio di insieme multiconcetto è illustrato nella figura 2.16*

	<i>Disprezzo</i>	<i>Repulsione</i>	<i>Sentimenti</i>	<i>0,6</i>	<i>0,7</i>	
	<i>Indebolire</i>	<i>Violenza</i>	<i>Sentimenti</i>	<i>0,3</i>	<i>0,2</i>	
$A = <$	<i>Vendicare</i>	<i>Violenza</i>	<i>Sentimenti</i>	<i>0,8</i>	<i>0,5</i>	$>$
	<i>Attaccare</i>	<i>Violenza</i>	<i>Sentimenti</i>	<i>0,9</i>	<i>0,8</i>	
	<i>Attaccare</i>	<i>Conflitto</i>	<i>Sentimenti</i>	<i>0,8</i>	<i>0,7</i>	

Figure 2.16: Esempio di insieme multiconcetto

Si noti che uno stesso vocabolo lessicale può appartenere all'insieme in due *parole* differenti, ognuna delle quali si riferisce ad un *concetto*, come nel caso del vocabolo **attaccare**.

Sia nel caso di insieme monoconcetto, che in quello di insieme multiconcetto, è possibile dare una definizione di **centralità dell'insieme**; ovviamente nel caso di insieme monoconcetto la centralità la si potrà calcolare rispetto al concetto al quale si riferisce (negli altri casi appare evidente che la centralità sarà nulla), mentre nel caso di insieme multiconcetto potremo calcolare la centralità rispetto ad uno qualsiasi dei concetti al quale si riferisce.

Definition 2.3.12 (Centralità di un insieme monoconcetto). *Dato un insieme monoconcetto A riferito al concetto C , si definisce **centralità di A rispetto a C** la media delle centralità delle parole appartenenti all'insieme A .*

$$CEN_C(A) = \frac{\sum_{x \in A} (F_C(x))}{N}. \quad (2.3.3)$$

dove N è il numero di parole appartenenti ad A

Example 2.3.9. Consideriamo ad esempio l'insieme monoconcetto A nella figura 2.15, in tal caso avremo che la sua centralità di rispetto alla **Repulsione**, sarà 0,58; ovvero

$$CEN_{REPULSIONE}(A) = 0,58$$

Per un insieme multiconcetto è, come detto, è possibile definire la centralità che esso ha nei vari concetti al quale fa riferimento.

Definition 2.3.13 (Centralità di un insieme multiconcetto). Dato un insieme multiconcetto A riferito agli n concetti C_1, \dots, C_n si definisce **centralità di A rispetto a C_i** la media delle centralità delle parole appartenenti all'insieme A ed associate al concetto C_i , ovvero

$$CEN_{C_i}(A) = \frac{\sum_{x \in A} (F_{C_i}(x))}{N_i}. \quad (2.3.4)$$

dove N_i è il numero di parole appartenenti ad A ed associate al concetto C_i .

Example 2.3.10. Consideriamo ad esempio l'insieme multiconcetto A della figura 2.16, in tal caso la centralità di quest'insieme rispetto al concetto **Repulsione** è 0.6, rispetto al concetto **Violenza** è 0.66, mentre rispetto al concetto **Conflitto** è 0.8, ovvero

$$CEN_{REPULSIONE}(A) = 0.6$$

$$CEN_{VIOLENZA}(A) = 0.66$$

$$CEN_{CONFLITTO}(A) = 0.8$$

È possibile anche definire delle operazioni tra gli insiemi monoconcetto o multiconcetto.

Definition 2.3.14 (Unione di insiemi). *Detti A e B due insiemi (monoconcetto o multiconcetto), l'insieme C si definisce **unione tra A e B** , e si indica nel seguente modo*

$$C = \text{Unione}(A, B)$$

se esso è formato dalle parole appartenenti ad A e B . Nel caso in cui una parola fosse contenuta sia in A che in B questa viene inserita una sola volta in C .

$$C = \{x \in \text{insieme delle parole} \mid x \in A \text{ oppure } x \in B\} \quad (2.3.5)$$

In relazione all'operazione di unione valgono le seguenti **proprietà**:

- se A e B sono due insiemi monoconcetto che si riferiscono allo stesso concetto, allora la loro unione sarà ancora un insieme monoconcetto.
- se A e B sono due insiemi monoconcetto che si riferiscono a due concetti differenti, allora la loro unione sarà un insieme multiconcetto.
- se A e B sono due insiemi multiconcetto, allora la loro unione sarà comunque un insieme multiconcetto.

Dal punto di vista analitico è ovvio che

$$\text{Unione}(A, A) = A$$

$$\text{Unione}(A, B) = \text{Unione}(B, A)$$

$$\text{Unione}(\text{Unione}(A, B), C) = \text{Unione}(A, \text{Unione}(B, C))$$

Definition 2.3.15 (Dimensione di un insieme). *Dato un insieme A (monoconcetto o multiconcetto) definiamo la **dimensione di A** e la indichiamo con **$DIM(A)$** il numero di **PAROLE** contenute all'interno dell'insieme.*

Example 2.3.11. *Consideriamo sempre l'insieme multiconcetto A di figura 2.16: si avrà che la dimensione di A è pari a cinque, cioè*

$$DIM(A) = 5$$

Definition 2.3.16 (Peso di un insieme). *Dato un dominio D composto dai concetti C_1, \dots, C_n ed un insieme A (monoconcetto o multiconcetto), si definisce **peso di A** il valore ottenuto facendo la media dei valori dei pesi delle parole dell'insieme A .*

$$Peso(A) = \sum_{j=1}^{DIM(A)} \left(\frac{P_D(x_j)}{DIM(A)} \right) \quad (2.3.6)$$

Example 2.3.12. *Considerato sempre l'insieme multiconcetto A precedentemente illustrato in figura 2.16 si avrà che il peso di A sarà $2.9/5$.*

Si noti che poiché il peso è una caratteristica che riguarda il dominio, non ci sono differenza tra insiemi monoconcetto o multiconcetto per quanto riguarda il calcolo del loro peso. In relazione al peso vale la seguente disuguaglianza: dati due insiemi (monoconcetto o multiconcetto) A e B si ha che

$$\text{Peso}(\text{Unione}(\mathbf{A}, \mathbf{B})) \neq \frac{\text{Peso}(A) + \text{Peso}(B)}{2} \quad (2.3.7)$$

Dato un dominio D costituito da un insieme di concetti C_1, \dots, C_n è possibile definire una grado di associazione tra i concetti C_i e C_j , con $i, j = 1..n$

$$\langle \mathbf{C}_i, \mathbf{C}_j, \text{gradodiassociazione} \rangle$$

Si può costruire in questo modo una **Semantic Network** tra i concetti, in cui i nodi rappresentino i concetti e l'arco dal nodo i al nodo j il grado di associazione tra C_i e C_j .

Definition 2.3.17 (Voto di un insieme monodominio). *Dato un dominio D composto dai concetti C_1, \dots, C_n ed un insieme A si definisce **Voto di A** la somma di due contributi uno pari al peso di A, e l'altro funzione della **distanza probabilistica DP** tra tutti i concetti interessati da A.*

$$\mathbf{Voto}(\mathbf{A}) = \text{Peso}(\mathbf{A}) + f(\mathbf{DP}) \quad (2.3.8)$$

In particolare la $f(\mathbf{DP})$ è la somma delle distanze probabilistiche tra i concetti di tutte le coppie di parole relative ad A moltiplicate per la media delle centralità delle parole nei relativi concetti, che però superano un valore limite.

2.3.3 Domini indipendenti distinti

Consideriamo adesso l'esistenza di m domini D_1, \dots, D_m indipendenti.

Definition 2.3.18 (Insieme multidominio). *Dati m domini D_1, \dots, D_m ognuno dei quali composto dai concetti C_{d1}, \dots, C_{dn_d} , dove d è l'indice dei domini mentre n_d è il numero di concetti del dominio d , (si noti che può capitare che uno stesso concetto C appartenga a due domini diversi) si definisce **insieme multidominio** quell'insieme di Parole i cui elementi appartengono a domini differenti.*

Un insieme multidominio A può essere visto come l'unione degli insiemi (monoconcetto o multiconcetto) ottenuti facendo la proiezione di A sui domini D_1, \dots, D_m . Indichiamo la proiezione di A sul dominio D_i con

$$A|D_i = \text{proiezione di } A \text{ sul dominio } D_i$$

È possibile in questo modo estendere le definizioni di centralità e peso ad un insieme multidominio.

Definition 2.3.19 (Centralità di un insieme multidominio). *La centralità di un insieme multidominio A rispetto ad uno dei concetti C_{ij} appartenenti al dominio D_i , si definisce come la centralità dell'insieme (monoconcetto o multiconcetto) $A|D_i$ rispetto al concetto C_{ij} .*

$$\text{CEN}_{C_{ij}}(A) = \text{CEN}_{C_{ij}}(A|D_i) \quad \text{con } C_{ij} \in D_i$$

Definition 2.3.20 (Peso di un insieme multidominio rispetto ad un dominio D_i). Dati m domini D_1, \dots, D_m ognuno dei quali composto dai concetti C_{d1}, \dots, C_{dn_d} , dove d è l'indice dei domini mentre n_d è il numero di concetti del dominio d , e considerato un insieme multidominio A si definisce **peso di un insieme multidominio A rispetto ad un dei domini D_i** , la sommatoria dei pesi delle parole appartenenti all'insieme dell'insieme (monoconcetto o multiconcetto) $A|D_i$ diviso però per la dimensione di A .

$$\text{Peso}_{D_i}(A) = \sum_{j=1}^{DIM(A|D_i)} \left(\frac{P_{D_i}(x_j)}{DIM(A)} \right) \quad (2.3.9)$$

Come detto in precedenza, avremo una *mappa dei concetti* per ogni dominio, ma poiché è anche possibile che stessi concetti appartengano a domini differenti, potrebbe capitare che esista una *distanza probabilistica* tra i concetti C_k e C_t in due o più domini; indicata con $DP_{D_i}(C_k, C_t)$ la distanza probabilistica tra C_k e C_t in relazione al dominio D_i , consideriamo la seguente situazione:

$$\begin{aligned} & DP_{D_1}(C_k, C_t) \\ & DP_{D_2}(C_k, C_t) \\ & \vdots \\ & DP_{D_m}(C_k, C_t) \end{aligned}$$

in questo caso la distanza probabilistica tra C_k e C_j può essere calcolata con la seguente formula

$$\text{DP}(C_k, C_j) = \frac{1}{m} \sum_{i=1}^m (DP_{D_i}(C_k, C_j)) \quad (2.3.10)$$

La motivazione della precedente formula è la seguente: dalla teoria della probabilità è noto che dato un evento E ed una partizione D_1, \dots, D_m dell'evento certo S per la **legge della probabilità totale**

$$\mathbf{P}(\mathbf{E}) = \sum_{i=1}^m (P(D_i) \cdot P(E|D_i))$$

Nel caso in esame i domini D_1, \dots, D_m costituiscono una partizione dell'universo lessicale (evento certo) e sono equiprobabili con probabilità $\frac{1}{m}$, inoltre la distanza probabilistica $DP_{D_i}(C_k, C_t)$ può essere considerata come la probabilità condizionata dato il dominio D_i .

Nel caso di m domini indipendenti il voto di un insieme multidominio rispetto ad uno degli m domini può essere calcolato con una formula analoga a quella definita nel caso di insiemi monodominio.

$$\mathbf{Voto}(\mathbf{A}) = \mathbf{Peso}(\mathbf{A}) + \mathbf{f}(\mathbf{DP}) \quad (2.3.11)$$

2.3.4 Domini distinti non indipendenti

Consideriamo infine l'esistenza di m *domini* D_1, \dots, D_m non indipendenti tra loro. In questo caso è possibile definire una **mappa dei domini**.

Definition 2.3.21 (mappa dei domini). *Dati m domini D_1, \dots, D_m ognuno dei quali composto dai concetti C_{d1}, \dots, C_{dn_d} , dove d è l'indice dei domini mentre n_d è il numero di concetti del dominio d , si definisce mappa dei domini il grafo pesato i cui nodi rappresentano i domini D_i e gli archi non orientati le relazioni tra essi. I pesi associati agli archi rappresentano il **grado di correlazione** tra i domini interessati: si indica con $GA(D_i, D_j)$ il grado di associazione tra D_i e D_j . Si assume $GA(D_i, D_j) = 0$ se non esiste un arco diretto tra D_i e D_j .*

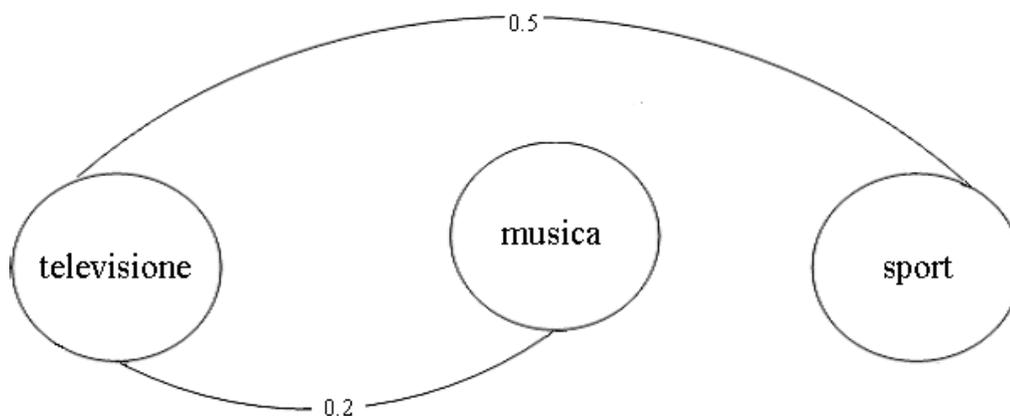


Figure 2.17: Esempio di mappa dei domini

Anche nel caso di m domini non indipendenti è possibile calcolare il voto di un insieme multidominio rispetto ad uno degli m domini utilizzando la definizione data. In questo caso ovviamente la funzione $f(\mathbf{DP})$ sarà più complessa dovendo tener conto delle relazioni tra i domini.

2.4 La metrica

Per poter discriminare le pagine di interesse per uno degli m domini esistenti, da quelle che non lo sono, occorre definire una metrica che consenta al sistema di dare un voto alle pagine a seconda del loro contenuto sintattico e semantico.

La metrica qui proposta prende in considerazione cinque informazioni che fanno riferimento alla pagina a cui vogliamo dare un giudizio, due delle quali legate alla sintattica e tre alla semantica delle pagine:

- *Profondità della pagina rispetto al motore di ricerca*
- *Posizione mediata della pagina*
- *Peso dei vocaboli nella pagina*
- *Premio dell'accoppiata di vocaboli che si riferiscono a concetti correlati tra loro, pesato con la distanza dei vocaboli stessi*
- *Premio dovuto all'associazione tra domini correlati*

2.4.1 GPrS e posizione mediata

Consideriamo Π l'insieme delle pagine WEB. Dati n motori di ricerca, viene assegnato ad ogni motore un peso w_i , compreso tra 0 e 1, tale che:

$$\sum_{i=1}^n (w_i) = 1 \quad (2.4.1)$$

Definition 2.4.1 (Funzione posizione). *Definiamo funzione di posizione*

$$p_i : \Pi \rightarrow N$$

come quella funzione intera che ad ogni pagina WEB x appartenenti all'insieme Π associa la sua posizione all'interno dell'elenco dei risultati del motore di ricerca i , se la pagina viene restituita dal motore stesso, 0 nel caso in cui la pagina non fosse trovata dal motore.

Google [Ricerca avanzata](#) [Preferenze](#) [Strumenti per le lingue](#) [Suggerimenti per la ricerca](#)

spagna

Cerca nel Web Cerca solo le pagine in Italiano

[Web](#) [Immagini](#) [Gruppi](#) [Directory](#)

Google ha cercato **spagna** nelle pagine in **Italiano**. Risultati **1 - 10** di circa **418,000**. Durata della ricerca: **0.05** secondi.

Categoria: [World](#) > [Italiano](#) > [Regionale](#) > [Europa](#) > [Spagna](#)

[Spagna.com :: informazioni, itinerari, città e località di mare ...](#)
 informazioni, itinerari, città, località e suggerimenti per viaggiare
 in **Spagna**. per ricevere la newsletter basta la tua email ...
[www.spagna.com/](#) - 10k - 17 Gen 2003 - [Copia cache](#) - [Pagine simili](#)

[...: Ivana Spagna - Official Site - www.ivanaspagna.it ::...](#)
[www.ivanaspagna.it/](#) - 2k - [Copia cache](#) - [Pagine simili](#)

[Spagna contemporanea](#)
Spagna contemporanea: rivista semestrale di storia, cultura e bibliografia; rivista
 de historia y cultura de la España contemporánea; Journal of History and ...
 Descrizione: Rivista di storia, cultura e bibliografia nata nel 1992 per iniziativa congiunta di un gruppo di studiosi...
 Categoria: [World](#) > [Italiano](#) > [Scienza](#) > [Scienze Sociali](#) > [Studi Geoculturali](#)
[www.spagnacontemporanea.it/](#) - 5k - [Copia cache](#) - [Pagine simili](#)

Collegamenti sponsorizzati

[Free Internet Access](#)
 Connect to the Internet in Spain
 totally unlimited and for free.
[www.gonuts4free.com](#)
 Livello di interesse:

[Per visualizzare il vostro
 messaggio...](#)

Figure 2.18: Esempio di risultato di ricerca con Google

Example 2.4.1. *Considerato l'esempio di figura 2.18, si ha che la posizione della pagina [www.spagnacontemporanea.it/](#) ha una posizione pari a 3, rispetto al motore di ricerca Google.*

Definition 2.4.2 (GPrS). *Data una pagina WEB x appartenente all'insieme Π e siano A l'insieme di numeri naturali i tali che il motore di ricerca i -esimo restituisca la pagina x , e B l'insieme di numeri naturali j tali che il motore di ricerca j -esimo non restituisca la pagina x . Sia inoltre $a = \text{card}(A)$, si definisce **GPrS** della pagina x la quantità:*

$$\mathbf{GPrS}(x) = \frac{\sum_{i \in A} (w_i)}{\sum_{i \in A} p_i(x) \cdot (w_i + \frac{\sum_{j \in B} w_j}{a})} \quad (2.4.2)$$

ovviamente risulta $\mathbf{GPrS}(x) \in]0, 1]$.

La situazione migliore si ha quando $\mathbf{GPrS}(x) = 1$, infatti in tal caso tutti i motori di ricerca hanno trovato la pagina x nella posizione numero 1, man mano che la pagina viene trovata in posizioni maggiori o da meno motori di ricerca, $\mathbf{GPrS}(x)$ si avvicina sempre più a 0.

Definition 2.4.3 (Posizione mediata). *Data una pagina WEB x appartenente all'insieme Π , si definisce posizione mediata di x , e la si indica con $\mu(x)$, la quantità inversa di $\mathbf{GPrS}(x)$, cioè*

$$\mu(x) = \frac{1}{\mathbf{GPrS}(x)} \quad (2.4.3)$$

ovviamente risulta $\mu(x) \in [1, \infty]$.

Si noti che se tutti i motori di ricerca restituissero la pagina x , la posizione mediata sarebbe la media pesata delle posizioni nei singoli motori, con i pesi attribuiti ai motori di ricerca.

2.4.2 Contributo sintattico

Definition 2.4.4 (Quantità $q(x)$). *Data una pagina WEB x appartenente all'insieme Π , e detto L il numero di link prelevati da ogni motore di ricerca, si definisce $q(x)$ la quantità*

$$q(x) = \frac{L}{(L + GPrS(x) - 1)} \quad (2.4.4)$$

ovviamente si ha che $q(x) \in]0, 1]$.

Si noti che la situazione migliore si ha con $q(x) = 1$: tale valore si ottiene nel caso in cui $GPrS(x)=1$ ovvero quando il valore di L tende ad ∞ .

Definition 2.4.5 (Contributo sintattico). *Detta b la profondità della pagina x , si definisce contributo sintattico, e lo si indica con $Sin(x)$, la quantità*

$$Sin(x) = \frac{q(x)}{1 + b \cdot K_b} \quad (2.4.5)$$

dove K_b è una costante da trovare sperimentalmente per definire quanto la profondità di una pagina debba incidere sul suo contributo sintattico. Ovviamente si ha che $Sin(x) \in]0, 1]$ e la situazione migliore si ha con $Sin(x) = 1$.

Se consideriamo una pagina WEB x appartenente all'insieme Π , è possibile che al suo interno siano presenti dei link ad altre pagine WEB sempre appartenenti all'insieme Π .

Definition 2.4.6 (Profondità di una pagina). *Data una pagina WEB x appartenente all'insieme Π e detta y una pagina WEB appartenente a Π il cui link si trova all'interno di x , definiamo la profondità della pagina y come la profondità della pagina x aumentata di 1. Si definisce inoltre che la profondità delle pagine restituite da qualche motore di ricerca sia pari a θ^5 .*

Definition 2.4.7 (Contributo sintattico). *Data una pagina WEB x appartenente all'insieme Π , detta b la sua profondità e K_b un valore reale, si definisce contributo sintattico, e lo si indica con **Sin(x)**, la quantità*

$$\text{Sin}(x) = \frac{q(x)}{1 + b \cdot K_b} \quad (2.4.6)$$

ovviamente si ha che $\text{Sin}(x) \in]0, 1]$.

Si noti che anche per il contributo sintattico la situazione migliore si ha quando il suo valore è pari ad 1, ciò avviene per le pagine che hanno $q(x)$ uguale ad 1 ed appartengono all'insieme delle pagine restituite da qualche motore di ricerca.

Si noti ancora che il valore della costante della profondità K_b deve essere trovato

⁵*Il motivo della definizione della profondità verrà meglio chiarita in seguito nella spiegazione del funzionamento del sistema*

sperimentalmente; tale valore indica quanto la profondità di una pagina debba incidere sul suo contributo sintattico; nel caso in cui K_b fosse pari a 0 si avrebbe che la profondità non influenza il contributo sintattico.

2.4.3 Contributo semantico

Il contributo semantico, che possiamo indicare con V , è somma di tre quantità:

- *Peso dei vocaboli trovati in uno dei dizionari presenti nel sistema*
- *Contributo funzione delle distanze probabilistiche dei concetti trovati*
- *Peso dovuto all'associazione tra domini correlati*

Si noti che d'ora in poi faremo riferimento in generale ad una pagina anche se tutto può essere riferito separatamente al contenuto, al titolo, alla descrizione o alle keywords della pagina stessa.

2.4.4 Peso dei vocaboli di una pagina

Definition 2.4.8 (Peso di una pagina in un dominio D_i). *Data una pagina x appartenente all'insieme Π e m domini D_1, \dots, D_m , detto inoltre N_i il numero di vocaboli presenti in x ed associati al dominio D_i e NT il numero di vocaboli presenti in x ed appartenenti ad uno degli m domini, cioè*

$$NT = \sum_{i=1}^m N_i$$

si definisce peso di una pagina nel dominio D_i , e lo si indica con $Peso_{D_i}(x)$, la somma di tutti i pesi $t_{i,j}$, relativi al dominio D_i , associati agli N vocaboli presenti in

una pagina (ovviamente se un vocabolo non è presente nel dominio D_i il suo peso vale 0).

$$\mathbf{Peso}_{D_i}(\mathbf{x}) = \frac{\sum_{j=1}^N t_{i,j}}{NT} \quad (2.4.7)$$

Si noti che il *Peso di una pagina* porta in se anche un significato sintattico in quanto il peso viene dato in base alla forma del vocabolo (e quindi in base alla sua sintassi).

Contributo funzione delle distanze probabilistiche dei concetti presenti nella pagina

Abbiamo visto che il peso di una pagina dipende dai vocaboli trovati al suo interno. Poiché ad ogni vocabolo è associato un concetto e la relativa centralità possiamo parlare anche di concetti, e centralità associate, trovati in una pagina.

Definition 2.4.9 (Distanza tra due concetti). *Data una pagina x appartenente all'insieme Π ed una coppia di concetti C_i e C_j presente all'interno della pagina si definisce distanza tra i concetti, e la si indica con $DIST(C_i, C_j)$ la differenza tra le posizioni dei concetti nella pagina.*

$$\mathbf{DIST}(C_i, C_j) = |\mathit{Posizione}(C_i) - \mathit{Posizione}(C_j)| \quad (2.4.8)$$

Definition 2.4.10 (Centralità media tra due concetti). *Data una pagina x appartenente all'insieme Π ed una coppia di concetti C_i e C_j presente all'interno della pagina si definisce Centralità media tra due concetti, e la si indica con $CENM(C_i, C_j)$ la media delle centralità associate ai concetti C_i e C_j .*

$$CENM(C_i, C_j) = \frac{CEN(C_i) + CEN(C_j)}{2} \quad (2.4.9)$$

Definition 2.4.11 (Distanza probabilistica normalizzata). *Data una pagina x appartenente all'insieme Π ed una coppia di concetti C_i e C_j trovata all'interno della pagina si definisce distanza probabilistica normalizzata rispetto al dominio D_i , e la si indica con $DPN_{D_i}(C_i, C_j)$ la distanza probabilistica che i due concetti hanno nel dominio D_i diviso la distanza tra i due concetti nella pagina.*

$$DPN_{D_i}(C_i, C_j) = \frac{DP_{D_i}(C_i, C_j)}{DIST(C_i, C_j)} \quad (2.4.10)$$

Definition 2.4.12 (Contributo diretto delle DP). *Dati m domini D_1, \dots, D_m ognuno dei quali composto dai concetti C_{d1}, \dots, C_{dn_d} , dove d è l'indice dei domini mentre n_d è il numero di concetti del dominio d , ed una pagina x appartenente all'insieme Π si definisce contributo diretto delle DP relativo al dominio D_i , e lo si indica con $CdDP_{D_i}(x)$, la somma delle $DPN_{D_i}(C_{ik}, C_{ih})$ che superano una certa soglia, tra tutte le coppie di concetti trovati in x ed associati al dominio D_i , moltiplicate per la centralità media $CENM(C_{ik}, C_{ih})$.*

$$\mathbf{CdDP}_{D_i}(\mathbf{x}) = \sum_{h=1}^{NT} \sum_{k=h+1}^{NT} (DPN_{D_i}(C_{ik}, C_{ih}) \cdot CENM(C_{ik}, C_{ih})) \cdot X_{k,h} \quad (2.4.11)$$

dove $X_{k,h} = 1$ se $DPN_{D_i}(C_{ik}, C_{ih}) \geq \mathbf{soglia}$ mentre vale 0 negli altri casi.

Definition 2.4.13 (Contributo indiretto delle DP). Dati m domini D_1, \dots, D_m ognuno dei quali composto dai concetti C_{d1}, \dots, C_{dn_d} , dove d è l'indice dei domini mentre n_d è il numero di concetti del dominio d , ed una pagina x appartenente all'insieme Π si definisce contributo indiretto delle DP reativo al dominio D_i , e lo si indica con $CiDP_{D_i}(x)$, la somma delle $DPN_{D_i}(C_{jh}, C_{jk})$ che superano una certa soglia, tra tutte le coppie di concetti trovati in x e associati a qualche dominio D_j (con $j \neq i$) moltiplicate per la centralità media $CENM(C_{jh}, C_{jk})$ ed il grado di correlazione, $GC(D_i, D_j)$, tra il dominio D_i e il dominio D_j .

$$\mathbf{CiDP}_{D_i}(\mathbf{x}) = \sum_{k=1}^{NT} \sum_{h=k+1}^{NT} DPN_{D_i}(C_{jk}, C_{jh}) \cdot CENM(C_{jk}, C_{jh}) \cdot GC(D_i, D_j) \cdot X_{k,h} \quad (2.4.12)$$

dove $X_{k,h} = 1$ se $DPN_{D_i}(C_{jk}, C_{jh}) \geq \mathbf{soglia}$ mentre vale 0 negli altri casi.

Peso dovuto all'associazione tra domini correlati

Definition 2.4.14 (Peso indiretto). *Dati m domini D_1, \dots, D_m ognuno dei quali composto dai concetti C_{d1}, \dots, C_{dn_d} , dove d è l'indice dei domini mentre n_d è il numero di concetti del dominio d , ed una pagina x appartenente all'insieme Π si definisce Peso indiretto di x relativo al dominio D_i , e lo si indica con $PesoI_{D_i}(x)$, la somma dei pesi $Peso_{D_j}(x)$ della pagina x nei domini D_j con $j \neq i$ per il grado di correlazione, $GC(D_i, D_j)$, tra D_i e D_j .*

$$\mathbf{PesoI}_{D_i}(\mathbf{x}) = \sum_{k=1}^m (Peso_{D_k}(x) \cdot GC(D_i, D_j)) \quad \text{con } j \neq i \quad (2.4.13)$$

Contributo semantico relativo al dominio D_i

Il contributo semantico di una pagina è dipende dal **titolo**, dal **contenuto** della pagina, dalle **keywords** e della **descrizione**. Indichiamo le precedenti parti di interesse rispettivamente con i pedici **t**, **c**, **k** e **d**.

Definition 2.4.15 (Contributo semantico al dominio D_i). *Dati m domini D_1, \dots, D_m ognuno dei quali composto dai concetti C_{d1}, \dots, C_{dn_d} , dove d è l'indice dei domini mentre n_d è il numero di concetti del dominio d , ed una pagina x appartenente all'insieme Π si definisce Contributo semantico relativo al dominio D_i della pagina x , e lo si indica con V_{sem} , la somma pesata dei diversi contributi diretti e indiretti delle DP relativa a titolo, contenuto, descrizione e keywords.*

$$\mathbf{V}_{sem} = \sum_{p \in (t, c, d, k)} w_p \cdot (CdDP_{D_i}(x) + CiDP_{D_i}(x))_p \quad (2.4.14)$$

dove w_c, w_d, w_t, w_k sono rispettivamente i pesi del contenuto, della descrizione, del titolo e delle keyword della pagina x , tali che

$$w_c + w_d + w_t + w_k = 1 \quad (2.4.15)$$

Contributo semantico-sintattico relativo al dominio D_i

Anche il contributo semantico di una pagina dipende dal **titolo**, dal **contenuto** della pagina, dalle **keywords** e della **descrizione**.

Definition 2.4.16 (Contributo semantico-sintattico al dominio D_i). *Dati m domini D_1, \dots, D_m ognuno dei quali composto dai concetti C_{d1}, \dots, C_{dn_d} , dove d è l'indice dei domini mentre n_d è il numero di concetti del dominio d , ed una pagina x appartenente all'insieme Π si definisce Contributo semantico-sintattico relativo al dominio D_i della pagina x , e lo si indica con V_{ss} , la somma pesata dei diversi pesi diretti e indiretti del titolo, contenuto, descrizione e keywords.*

$$\mathbf{V}_{ss} = \sum_{p \in (t,c,d,k)} w_p \cdot (\text{Peso}_{D_i}(x) + \text{Peso}I_{D_i}(x))_p \quad (2.4.16)$$

dove w_c, w_d, w_t, w_k sono sempre i pesi del contenuto, della descrizione, del titolo e delle keyword della pagina x .

2.4.5 Voto complessivo

Definition 2.4.17 (Voto complessivo di una pagina). *Dati m domini D_1, \dots, D_m ognuno dei quali composto dai concetti C_{d1}, \dots, C_{dn_d} , dove d è l'indice dei domini mentre n_d è il numero di concetti del dominio d , ed una pagina x appartenente all'insieme Π si definisce Voto complessivo di x relativo al dominio D_i , e lo si indica con $Voto_{D_i}(x)$, la quantità*

$$\mathbf{Voto}_{D_i}(\mathbf{x}) = \frac{K_{sint} \cdot Sin(x) + K_{sem} \cdot V_{sem} + K_{ss} \cdot V_{ss}}{K_{sint} + K_{sem} + K_{ss}} \quad (2.4.17)$$

dove $K_{sint}, K_{sem}, K_{ss}$ sono tre costanti chiamate rispettivamente **costante sintattica**, **costante semantica** e **costante semantico-sintattica**.

Chapter 3

Il sistema

Nel corso del capitolo verrà presentata l'architettura del sistema realizzato

Il sistema SSA è stato progettato secondo il paradigma degli "agenti". Verrà quindi definito cos'è un agente software e verranno descritti gli agenti progettati mediante il Service Description Language.

3.1 Intelligent Search Agent

Gli **Intelligent Search Agents** sono programmi che filtrano le informazioni ritrovate sul WEB in base agli interessi dell'utente, memorizzandone le scelte e costruendone il profilo. Ad essi si stanno interessando molte tra le più importanti aziende del settore informatico, da Netscape a Microsoft, da IBM a AT&T.

Di norma, un buon agente dovrebbe permettere di utilizzare meno tempo per il reperimento delle informazioni e più tempo per l'analisi di quelle già reperite. Il nome di questi strumenti software richiede un piccolo chiarimento [3]:

- **Intelligent**: sono chiamati così in quanto non si limitano solamente a restituire pagine che contengono le keywords immesse dall'utente (come fanno i motori di

ricerca), ma anche pagine con parole di significato simile alle keywords assicurando che siano utilizzate nel giusto contesto evitando di restituire link inutili;

- **Search:** chiamati così in quanto realizzano, anche se in modo diverso, lo stesso task dei comuni motori di ricerca ossia il reperimento di pagine sul WEB;
- **Agent:** viene chiamato così qualcosa (nel nostro caso, un software) che agisce in base agli interessi di qualcuno. L'idea alla base di questo concetto è che un utente potrebbe fornire all'agente un set di istruzioni ed esso potrebbe eseguirle e riportare in seguito i risultati.

La maggior parte degli strumenti di Information Retrieval su Web non prescindono da un utente estremamente "impegnato" nell'attività di ricerca dell'informazione utile e presuppongono una ricerca in qualche modo "occasionale", non fissata a priori nei tempi. Nel momento specifico in cui una determinata informazione serve, ci si mette a cercarla e la ricerca richiede che vengano compiuti, ogni volta ed in prima persona, un certo numero di azioni alcune delle quali senz'altro ripetitive: per esempio, collegarsi alla pagina di Google (www.google.com) oppure a quella di Altavista (www.altavista.com), impostare i parametri della ricerca, restare collegato mentre si aspettano i risultati ed infine filtrare le risposte del motore. Sarebbe molto comodo, in questi casi, disporre di una sorta di "segreteria intelligente" che conosca in linea di massima i propri interessi, sappia prevedere ed anticipare le richieste, ci sostituisca, prendendo autonomamente le decisioni più opportune, sia per le azioni più ripetitive che per quelle meno ripetitive.

L'idea di *Agente di Ricerca Intelligente* cerca di avvicinarsi a questo ideale: in

sostanza, si tratta di incaricare un programma di svolgere, anche a intervalli prefissati, determinate ricerche, chiedendogli di reagire autonomamente ai risultati della ricerca stessa (magari "filtrandoli" attraverso l'uso di criteri che potrebbero essere difficili se non impossibili da impostare direttamente su un motore di ricerca e raggruppandoli in base al dominio di appartenenza).

Addirittura, se l'agente non è fisicamente localizzato sul proprio computer poiché situato su un server remoto o distribuito fra più server remoti, gli si potrebbero affidare ricerche e compiti da svolgere anche mentre si è scollegati dalla rete. Per esempio, un utente potrebbe chiedere all'agente di cercare informazioni su un certo argomento ed appartenenti ad un certo dominio, e l'agente potrebbe ricercare le pagine Web più aggiornate restituendo un report con le risorse disponibili sul WWW. Tali azioni potrebbero richiedere un giorno intero e l'utente potrebbe, una volta sottomessa la query, continuare i suoi lavori e ottenere i risultati in seguito. A molti, ciò potrebbe risultare sconveniente rispetto ad un motore di ricerca, dove i risultati vengono forniti in brevissimo tempo: in effetti, per veloci ricerche niente è meglio dei motori di ricerca, ma gli agenti di ricerca intelligenti offrono altri vantaggi. Per esempio, un utente potrebbe volere tutte le informazioni più recenti su diversi argomenti e potrebbe volerle ogni mattina sul proprio terminale, alla stregua dei servizi di news personalizzati. Esse verrebbero presentate esattamente come l'utente vuole e l'agente addirittura potrebbe apprendere nel tempo quali sono le sue preferenze in modo da migliorare il servizio.

Per certi aspetti, alcuni strumenti oggi disponibili sul WEB possiedono già, almeno in parte, qualcuna delle caratteristiche esposte: è il caso, ad esempio, dei servizi

di net filtering che permettono di impostare una ricerca da ripetere ad intervalli regolari comunicando i risultati via posta elettronica. Oppure programmi client per ricevere informazioni attraverso meccanismi di *information pushing*: una volta impostati, saranno loro ad occuparsi di collegarsi al server o ai server remoti e a scaricare dati. Il problema, in questi casi, è la carenza, o meglio l'assenza, della capacità di prendere decisioni realmente autonome, reagendo dinamicamente alle caratteristiche dell'universo informativo nel quale si muovono.

La realizzazione [49] di un agente di questo tipo richiede di includere meccanismi per poter effettuare ricerche sulla base di keywords sia multiple sia in combinazione tra di esse, per poter gestire l'esclusione di informazioni non rilevanti rispetto ai domini conosciuti dal sistema.

Data una query, l'agente lavora per soddisfarla considerando più percorsi possibili sul WEB. In genere si parte da una locazione di partenza (location seed) fornita dall'utente e l'agente esplora tutte le altre locazioni collegate alla radice che possono soddisfare la richiesta. Diversi sono i parametri che si possono impostare: la query (per esempio, delle keywords o una frase), da dove partire, quante iterazioni effettuare, sino a che profondità spingersi, quanto tempo impiegare e che algoritmi di ricerca impiegare nella navigazione sul WEB. Sono due le principali classi di algoritmi di ricerca [50]:

- di tipo informato
- di tipo non informato

I primi, o *algoritmi di tipo euristico*, si basano su alcune informazioni quali il

costo dei cammini effettuati o il numero di passi da effettuare a partire dal punto di partenza. Tali algoritmi sono fondati su tali informazioni per potare l'albero o il grafo di ricerca eliminando opportunamente diversi rami. I secondi, invece, non sfruttano nessuna di queste informazioni e si limitano a esplorare esaustivamente l'albero di ricerca o il grafo di ricerca sino a raggiungere l'obiettivo. In genere gli algoritmi di tipo euristico permettono di ottenere performance migliori ed è per questo che vengono sfruttati dalla maggior parte degli agenti di ricerca.

La progettazione si basa su quattro fasi principali:

- *Inizializzazione*
- *Percezione*
- *Azione*
- *Effetto*.

Nella fase di *inizializzazione* vengono fissate le variabili, le strutture e tutti i dati necessari per effettuare la ricerca. Vengono impostate la query, gli obiettivi della ricerca, il numero massimo di pagine Web da visitare, una locazione di partenza ed un metodo di ricerca se ne sono disponibili più di uno.

La fase di *percezione* sfrutta i dati forniti dall'utente per contattare un sito e prelevare da esso le informazioni. In tale fase l'agente è in grado di capire se la pagina soddisfa la query e di identificare percorsi verso altri siti Web (per esempio link ad altre pagine).

Nella fase di *azione* l'agente raccoglie tutte le informazioni disponibili e determina

se l'obiettivo è stato raggiunto. Se non è stato raggiunto, l'agente decide dove proseguire la ricerca: proprio in questo sta l'intelligenza dell'agente e il metodo di ricerca usato indica quanto esso sia abile.

Nella fase di *effetto* viene presentata all'utente una lista delle migliori locazioni trovate.

L'unione di tali fasi ed una connessione ad Internet permettono di realizzare un Agente di Ricerca Intelligente, o, come viene spesso chiamato in letteratura, un **Web Hunter** [51].

Durante la fase di inizializzazione vengono create apposite strutture dati per gestire le informazioni sui siti che verranno visitati: in genere è una lista a puntatori o un array. Ogni elemento o nodo contiene informazioni sull'URL della pagina e la sua locazione e un valore che rappresenta il costo per raggiungerla a partire dalla locazione di partenza (se si usano algoritmi di ricerca euristici). Il nodo radice contiene un valore detto seme che rappresenta la locazione da cui si parte nella ricerca: esso può essere un indirizzo IP messo a caso ma in genere è consigliabile partire dall'URL di un motore di ricerca o dall'indirizzo IP di una pagina attinente all'argomento della ricerca.

Effettuata l'inizializzazione, il Web Hunter entra nella fase di percezione aprendo una socket sulla porta HTTP (la porta 80) verso l'URL specificato. Stabilita la connessione, si sfrutta il comando HTTP "GET" per richiedere al server la pagina [?]; essa viene sottoposta a parsing al fine di trovare i link in essa contenuti e per determinare se essa è attinente alla query dell'utente. Il parser solitamente si basa sulla ricerca dei tag HTML "HREF=" e "SRC=" che indicano la presenza di un link nella pagina. È da notare che il parser può essere reso flessibile impostando diversi

metodi di ricerca all'interno della pagina: si possono, per esempio, ottenere file di grafica semplicemente cercando file con estensione "JPG" oppure "GIF". Le pagine Web possono contenere un numero molto elevato di link ad altre: per questo può essere utile, se l'agente lo consente, impostare un upper bound sul numero di link da prelevare da ogni pagina.

Una volta catturata la pagina ed estratte le informazioni più importanti, si entra nella fase in cui il Web Hunter manifesta la sua intelligenza; infatti, se l'obiettivo è stato raggiunto (confrontando, ad esempio, il numero di siti Web visitati con quello impostato dall'utente), il suo lavoro è terminato con successo. Se ciò non si è verificato, l'agente decide qual è la successiva azione da intraprendere: il cuore della decisione è rappresentato da quale algoritmo di ricerca utilizzare; infatti, i link estratti da una pagina e che, quindi, devono ancora essere visitati, vengono aggiunti alla lista spiegata in precedenza assieme al loro costo. La modalità di inserimento e, soprattutto, di prelievo è funzione dell'algoritmo di ricerca. Più il metodo è informato, maggiori sono le probabilità di raggiungere gli obiettivi prefissati. Ci sono diversi elementi per valutare il costo di una pagina, oltre a quelli già esposti in precedenza:

- Valutazione del suo round-trip-time tramite il comando "ping". Quanto più è basso, tanto più sarà basso il costo della pagina rispetto alla radice
- Valutazione della nazionalità della pagina. In una ricerca che si basi esclusivamente su siti, per esempio, americani, si potrebbe dare a pagine che non abbiano estensione compresa tra "COM", "NET", "EDU", "ORG", "GOV" o "MIL" (che si suppone essere pagine americane) un valore di costo elevato per penalizzarle

Se lo spazio di ricerca si esaurisce (non ci sono altri nodi da esplorare), l'agente

può richiedere un nuovo seme e ripartire oppure semplicemente terminare, entrando nella fase di effetto dove viene presentata all'utente una lista dei risultati ottenuti.

Il Web Hunter lavora muovendosi dalla fase di inizializzazione a quella di effetto, passando per quelle di percezione e azione. Tutto ciò sino a quando l'obiettivo non viene raggiunto. Diverse sono le applicazioni in cui può trovare uso un agente di questo tipo; infatti, dotandolo di un opportuno dizionario, può diventare un potente motore di ricerca domain-dependant oppure può essere modificato per trasformarsi da Web Hunter in Document Hunter in modo tale da cercare file di documenti in locale piuttosto che pagine Web.

Quando si testa un software di questo tipo è preferibile farlo localmente sino a quando non sia stabilita con certezza la sua stabilità [51]. Inoltre è preferibile che tali agenti si identifichino sul WEB utilizzando il campo HTTP "User-agent" e si registrino presso il "*Web Robots Database*" [53] in modo tale che possa essere identificato su un server protetto nel caso in cui debba effettuare su di esso diverse operazioni. Un'ultima nota che deve essere tenuta presente durante lo sviluppo e l'uso di agenti di questo tipo è la presenza su diversi siti Web di un file denominato *robot.txt*. Gli agenti di ricerca intelligente dovrebbero testare sempre la sua presenza, in quanto contiene informazioni rilevanti quali l'accettazione o meno da parte del sito di questo tipo di agenti oppure la politica di utilizzo delle informazioni in esso contenute. Altre pagine possono contenere tag HTML quali `<Meta Name="Robots" Content="Nofollow">` che indicano che è vietato il parsing della pagina stessa.

Per concludere il paragrafo, c'è da dire che di programmi che rispecchiano fedelmente quanto scritto ne esistono ben pochi, almeno in versioni "avanzate". È prevedibile, però, che il settore degli *Agenti di Ricerca Intelligente* conoscerà nei prossimi

anni un'evoluzione tale da far sembrare i primi strumenti di oggi solo rozze e primitive approssimazioni di applicazioni assai più sofisticate, potenti ed autonome.

3.1.1 La struttura ad agenti

Il sistema SSA è stato progettato secondo il paradigma degli "agenti". Abbiamo già illustrato cosa è un agente software e nel seguito di questo paragrafo verranno descritti gli agenti progettati mediante il *Service Description Language*.

Il Service Description Language (SDL)

Gli agenti progettati nel nostro sistema verranno di seguito descritti tramite uno specifico linguaggio di descrizione: il **Service Description Language (SDL)**.

Si tratta di un linguaggio che permette di descrivere gli agenti dal punto di vista dei servizi offerti. La descrizione di un servizio si basa sulla specifica dei seguenti elementi [11][15]:

- **Nome del servizio:** è un'espressione costituita da un verbo e da un nome (nella forma verbo: nome) che specifica il servizio offerto
- **Input:** è un elenco degli ingressi richiesti dal servizio
- **Output:** è un elenco delle uscite fornite dal servizio
- **Attributi:** alcuni servizi possono avere degli attributi associati. Per esempio, costo del servizio e tempo medio di esecuzione

Quando si devono elencare gli input e gli output, è necessario specificare le variabili ed i tipi ad esse associati. Ciò viene fatto tramite un apposito formalismo:

Definition 3.1.1 (item $s: \tau$). *Se s è una variabile avente valori nell'insieme τ , allora l'espressione $s: \tau$ è detta **item**.*

Ogni servizio può richiedere zero, uno o più input. Alcuni di essi sono obbligatori (il servizio non può essere portato a termine senza di essi), mentre altri sono facoltativi (non sono necessari, ma possono, se forniti, aumentare l'efficienza e la qualità del servizio offerto).

Prima di passare alla definizione dei servizi offerti dagli agenti progettati, è necessario fornire qualche definizione.

Definition 3.1.2 (item atomici). *Se $s: \tau$ è un item, allora l'espressione*

$$\langle I \rangle s : \tau \langle I \rangle \text{ (rispettivamente } \langle MI \rangle s : \tau \langle MI \rangle \text{)}$$

indica un item di ingresso atomico (rispettivamente item di ingresso obbligatorio atomico), mentre l'espressione

$$\langle O \rangle s : \tau \langle \emptyset \rangle$$

indica un item di uscita atomica.

Definition 3.1.3 (lista di item). *Se $s_1: \tau_1, \dots, s_n: \tau_n$ sono degli item, allora l'espressione*

$$\langle I \rangle s_1 : \tau_1, \dots, s_n : \tau_n \langle I \rangle^1$$

¹abbreviazione di $\langle I \rangle s_1 : \tau_1 \langle I \rangle \dots \langle I \rangle s_n : \tau_n \langle I \rangle$

indica una lista di item di ingresso, mentre le espressioni

$$\langle MI \rangle s_1 : t_1, \dots, s_n : t_n \quad \langle O \rangle s_1 : t_1, \dots, s_n : t_n$$

indicano una lista item di ingressi obbligatori ed una lista di uscite.

Definition 3.1.4 (Service Description). Siano s_n il nome di un servizio, i_1, \dots, i_n item di ingresso atomici, mi_1, \dots, mi_n item di ingresso obbligatori e o_1, \dots, o_n item di uscita atomici. Viene chiamata **Descrizione del Servizio (Service Description)**, la seguente espressione:

$$\begin{aligned} \langle S \rangle & \quad s_n \\ & \quad mi_1 \dots mi_k \\ & \quad i_1 \dots i_n \\ & \quad o_1 \dots o_r \\ \langle \backslash S \rangle & \end{aligned}$$

Tramite questo formalismo, è possibile presentare gli agenti dal punto di vista dei servizi da essi offerti.

²abbreviazioni, rispettivamente, di $\langle MI \rangle s_1 : \tau_1 \langle \backslash MI \rangle \dots \langle MI \rangle s_n : \tau_n \langle \backslash MI \rangle$ e $\langle O \rangle s_1 : \tau_1 \langle \backslash O \rangle \dots \langle O \rangle s_n : \tau_n \langle \backslash O \rangle$

Struttura ad agenti del sistema progettato

Il sistema SSA (Semantic Search Agent) progettato può essere visto come costituito da due agenti principali:

- *Spider Agent*
- *Mining Agent*

Spider Agent

È l'agente che si occupa di reperire le pagine Web sulla base delle query immesse dall'utente e di memorizzarle nel sistema. I servizi offerti, descritti tramite il *Service Description Language*, sono:

```
< S >  Adapt: query
        < MI > Query: String < \MI >
        < O > AdaptedQuery: String < \O >
< \S >
```

Questo servizio accetta come input la *query* sottoposta dall'utente (costituita da una stringa contenente singole keywords o intere frasi racchiuse tra virgolette). Compito del servizio è adattare la query immessa dall'utente alla sintassi adottata dai diversi motori di ricerca interrogati.

```
< S >  Submit: query
        < MI > AdaptedQuery: String < \MI >
        < MI > MaxRes: Integer < \MI >
        < O > Results: WebPage < \O >
< \S >
```

Questo servizio accetta come input la *query adattata*, insieme al massimo numero di link che ciascun motore deve prelevare dal WEB, e la sottopone al motore di ricerca corrispondente.

L'output sarà costituito dalla pagina restituita dal motore di ricerca contenente i link trovati.

```
< S >   ResultsParse: Results
          < MI > Results: WebPage < \MI >
          < O > Links: Strings < \O >
< \S >
```

Questo servizio accetta come input la pagina Web restituita dal motore di ricerca e fornisce come output i link in essa contenuti.

```
< S >   WebCatch: Links
          < MI > Link: String < \MI >
          < MI > MaxProf: Integer < \MI >
          < O > Page: WebPage < \O >
          < O > Average Position: Float < \O >
          < O > Depth: Integer < \O >
< \S >
```

Questo servizio accetta come input uno dei link contenuti all'interno delle pagine Web, insieme alla profondità massima che ogni pagina Web deve avere rispetto al motore di ricerca affinché essa venga prelevata dal sistema, e fornisce come output la corrispondente WebPage. Vengono inoltre restituiti, di ogni pagina Web, la posizione mediata e la profondità rispetto al motore di ricerca (parametri necessari ad alcuni servizi offerti dal Mining Agent). I concetti di profondità e posizione mediata verranno

ampiamente spiegati nel corso del capitolo successivo.

```
< S > DocumentParse: WebPage
      < MI > Document: WebPage < \MI >
      < O > Links: Strings < \O >
< \S >
```

Questo servizio accetta come input una delle pagine prelevate dal WEB e restituisce in output i link in esse contenuti.

```
< S > Construct: Repository
      < MI > Document: WebPage < \MI >
      < MI > Link: String < \MI >
      < O > Web Repository: DataBase < \O >
< \S >
```

Questo servizio accetta come input i link trovati ed i vari documenti prelevati dal WEB e costruisce il Web Repository di sistema.

Mining Agent

È l'agente che si occupa di analizzare le pagine Web prelevate dallo Spider, al fine di trovare gli elementi necessari all'attribuzione del voto semantico, un valore che misura l'attinenza ad un certo dominio delle pagine reperite. Si occupa, inoltre, dell'attribuzione del voto stesso. I servizi offerti, sono:

```

< S >   Get: Content
          < MI > Document: WebPage < \MI >
          < O > Content: Text < \O >
< \S >

```

Questo servizio accetta come input le varie pagine Web prelevate dal sistema e fornisce come uscita il loro contenuto semantico, ossia le pagine senza i tag HTML.

```

< S >   Preprocess: Content
          < MI > Content: Text < \MI >
          < MI > Knowledge Base: Database < \MI >
          < O > Adapted SemanticContent: Text < \O >
< \S >

```

Questo servizio permette di preprocessare il contenuto semantico estrapolato da ogni pagina Web al fine di sostituire eventuali caratteri di escape con i corrispondenti caratteri ASCII. I caratteri di escape sono preceduti dal simbolo & e servono per codificare in HTML speciali caratteri (per esempio, in HTML la "e commerciale" o &, viene codificata con la sequenza &).

Tale servizio, inoltre, consente di eliminare, all'interno del contenuto semantico estrapolato, eventuali stop-words, ossia quelle parole che non danno contributo alla semantica del documento. Per effettuare queste operazioni, è necessario accedere ad una base di conoscenza (Knowledge Base) ove è contenuta una lista dei caratteri di escape noti (con le rispettive corrispondenze) ed una lista delle stop-words.

```

< S >   Get: Title
          < MI > Document: WebPage < \MI >
          < O > WebTitle: Text < \O >
< \S >

```

Questo servizio accetta come input le varie pagine Web prelevate dal sistema e fornisce come uscita il contenuto del tag HTML Title.

```

< S >   Get: Meta Description
          < MI > Document: WebPage < \MI >
          < O > Description: Text < \O >
< \S >

```

Questo servizio accetta come input le varie pagine Web prelevate dal sistema e fornisce come uscita il contenuto del Meta Tag Description.

```

< S >   Get: Meta Keywords
          < MI > Document: WebPage < \MI >
          < O > Keywords: Text < \O >
< \S >

```

Questo servizio accetta come input le varie pagine Web prelevate dal sistema e fornisce come uscita il contenuto del Meta Tag Keywords.

```

< S >   Get: Semantic Vote
          < MI > MaxRes: Integer < \MI >
          < MI > Average Position: Float < \MI >
          < MI > Depth: Integer < \MI >
          < MI > Adapted Content: Text < \MI >
          < MI > WebTitle: Text < \MI >
          < MI > Description: Text < \MI >
          < MI > Keywords: Text < \MI >
          < MI > Knowledge Base: Database < \MI >
          < O > Semantic Vote: Floats < \O >
< \S >

```

Questo servizio accetta come input il massimo numero di link che l'utente ha chiesto allo Spider Agent di prelevare all'atto della sottomissione della query, la posizione mediata di una pagina Web, la sua profondità rispetto al motore di ricerca, il suo contenuto semantico, il suo titolo ed il contenuto dei Meta Tag Description e Keywords e restituisce in uscita il voto semantico. Per effettuare questa operazione, è necessario accedere ad una base di conoscenza (**Knowledge Base**) ove è contenuta una lista delle parole e dei concetti conosciuti dal sistema.

3.2 Architettura

L'architettura del sistema **SSA (Semantic Search Agent)** è rappresentata in figura 3.1; illustriamo dettagliatamente i componenti che ne fanno parte.

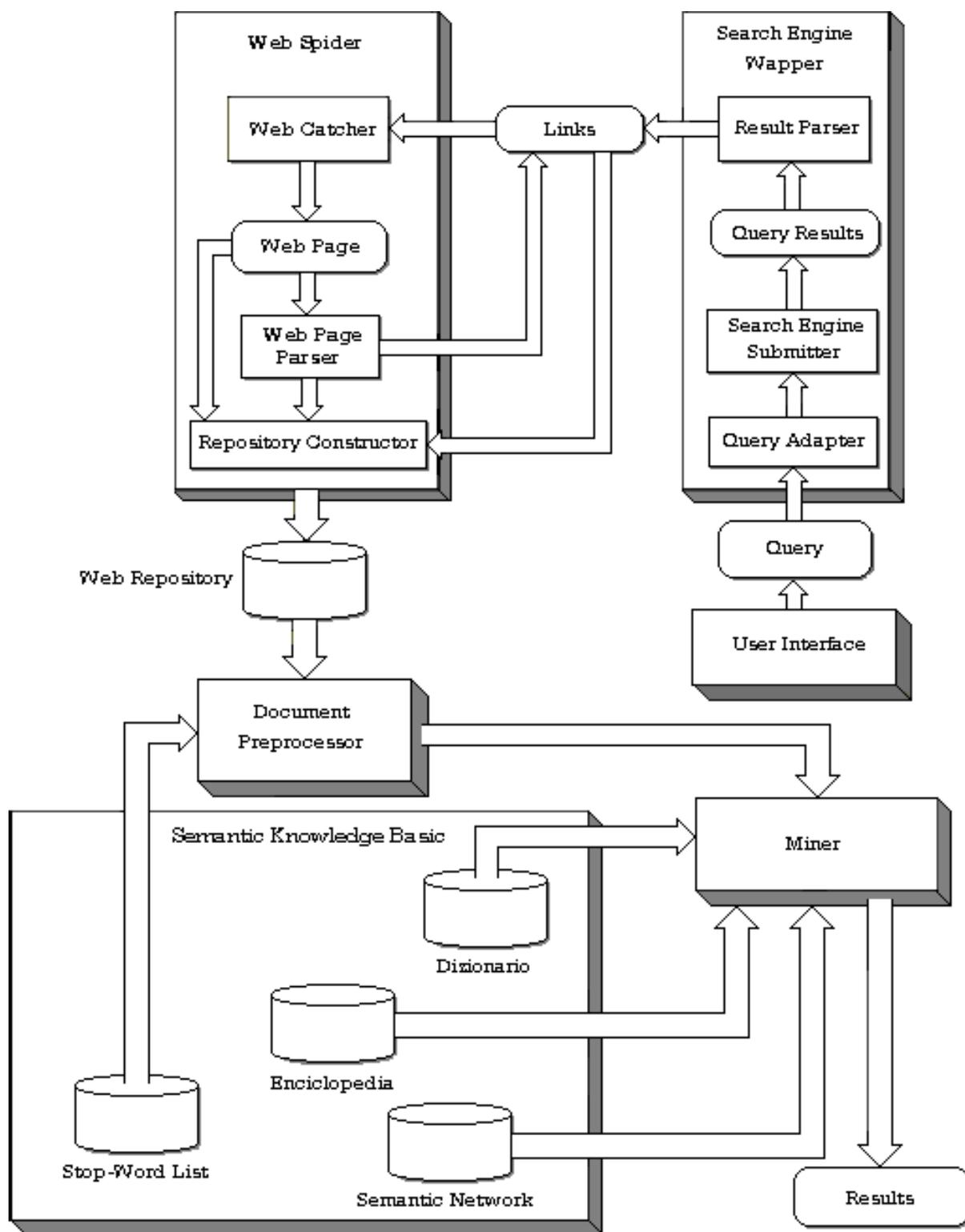


Figure 3.1: Architettura del sistema

3.2.1 User Interface

La **User Interface - Interfaccia Utente** è delegata a chiedere all'utente, tra le altre cose, la query con le parole chiave che egli vuole cercare sul WEB.



Figure 3.2: Menu Iniziale

Essa si presenta come in figura 3.2 ed i suoi pulsanti sono:

- *Gestione Motori* per inserire o eliminare uno dei motori di ricerca utilizzati dal sistema
- *Nuova Query* per la gestione di una nuova query da sottoporre al sistema
- *Carica Query* per la gestione delle query già sottoposte al sistema

- *Gestione Rete Semantica* per la gestione della Rete Semantica usata dal sistema per la valutazione dell'attinenza delle pagine ai domini
- *Gestione Dizionario* per la gestione della lista delle parole di interesse di tutti i domini
- *Impostazioni* per la gestione dei parametri dell'applicazione
- *Esci* per terminare l'applicazione

Tali funzionalità verranno spiegate in maniera dettagliata successivamente.

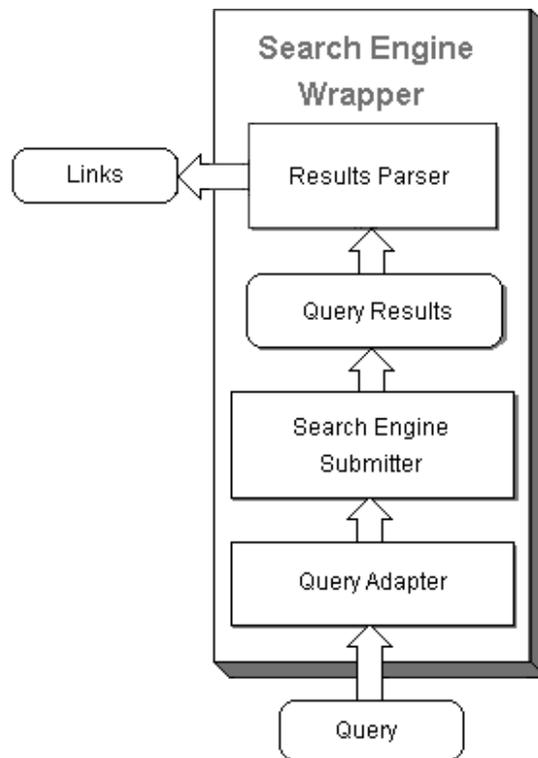


Figure 3.3: Search Engine Wrapper

3.2.2 Search Engine Wrapper

Il modulo **Search Engine Wrapper** è rappresentato in figura 3.3

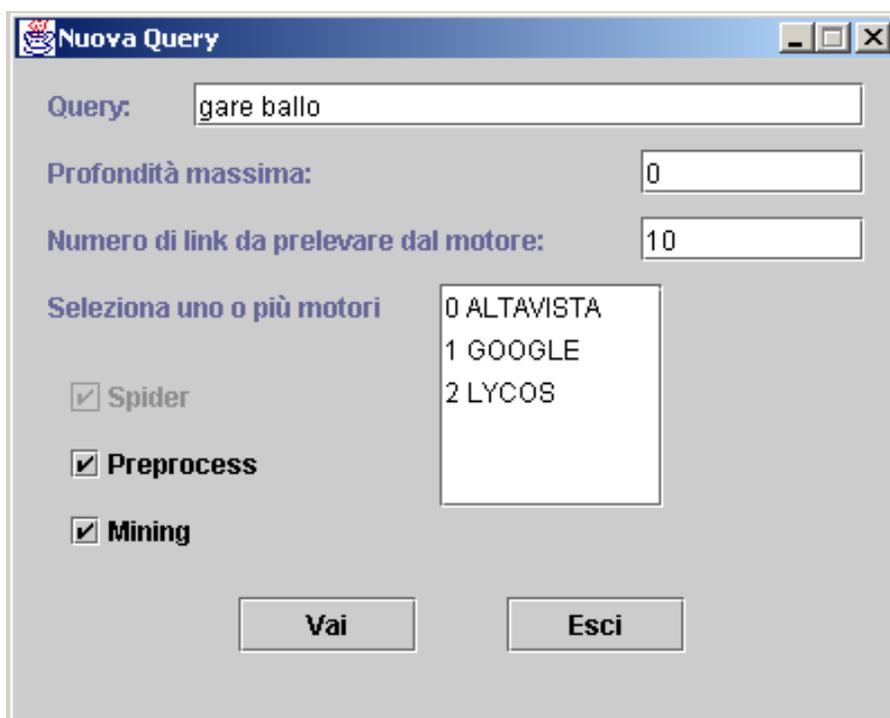


Figure 3.4: Esempio di Query

Questo modulo prende la query inserita mediante l'interfaccia utente, mostrata in figura 3.4, e la adatta alle specifiche sintassi dei motori di ricerca, mediante il *Query Adapter*, che, partendo dalla query effettuata nella sintassi del sistema, crea la stringa di interrogazione per i singoli motori di ricerca da interrogare. Siccome i motori di ricerca utilizzano regole differenti per l'inserimento della query, viene definita una sintassi unica a cui dovrà attenersi l'utente del sistema. Tale sintassi è simile a quella utilizzata dal motore di ricerca Google [56][57] ed è definita come segue

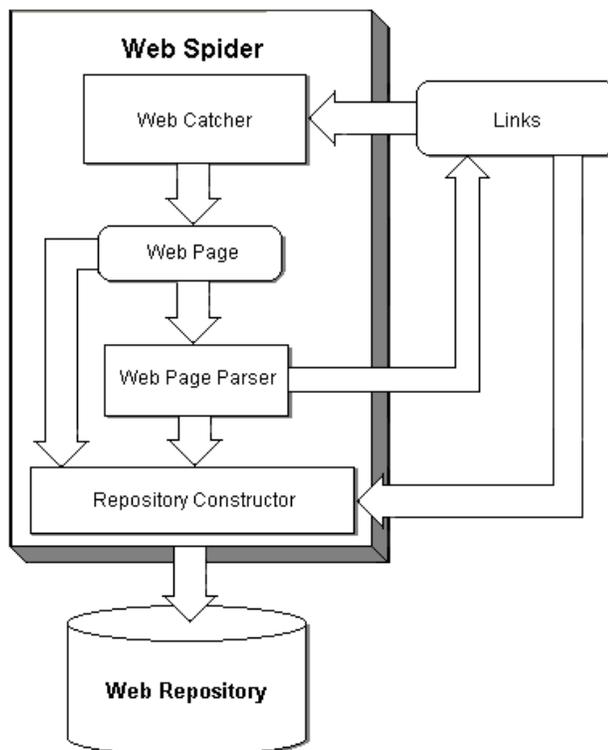


Figure 3.5: Web Spider

Definition 3.2.1 (Sintassi della query). *Le parole immesse sono tutte in AND tra loro, ossia devono essere tutte presenti nel testo della pagina affinché essa venga restituita dal motore di ricerca, ed inoltre è possibile immettere query contenenti intere frasi racchiudendo tali frasi tra virgolette*

Allo scopo di rendere del tutto trasparente per l'utente l'interrogazione dei motori di ricerca, il *Search Engine Wrapper* si occupa della sottomissione delle query adattate dal *Query Adapter* ai motori di ricerca e di prelevarne i risultati. Il *Search Engine Submitter* sottopone la query ad un motore di ricerca ottenendo la pagina con i

risultati da esso restituiti ed il *Results Parser* analizza tale pagina allo scopo di prelevarne i link in essa contenuti.

3.2.3 Web Spider

Il modulo **Web Spider** è rappresentato in figura 3.5

Le pagine corrispondenti ai link restituiti dai motori di ricerca vengono prelevate e memorizzate nel repository locale dal *Web Spider*. Di ogni pagina viene effettuato il prelievo dal *Web Catcher* ed un parsing finalizzato alla ricerca dei link contenuti nella pagina da parte del *Web Page Parser*. In tal modo è possibile seguire la struttura di un sito Web fino alla profondità scelta dall'utente. Le pagine prelevate ed i link trovati vengono inseriti nel *Web Repository* che contiene le informazioni necessarie per ricostruire localmente le pagine prelevate e la struttura del sito. La costruzione del *Web Repository* è delegata al *Repository Constructor*, che provvede all'inserimento delle pagine Web e dei link trovati nel database di sistema.

3.2.4 Document Preprocessor

Il modulo **Document Preprocessor** è rappresentato in figura 3.6.

Le pagine Web sono costituite da tre tipi di informazioni:

- **Tag HTML**, il cui significato riguarda essenzialmente la formattazione del testo nella pagina
- **Meta Tag HTML**, in cui l'autore della pagina può inserire una descrizione

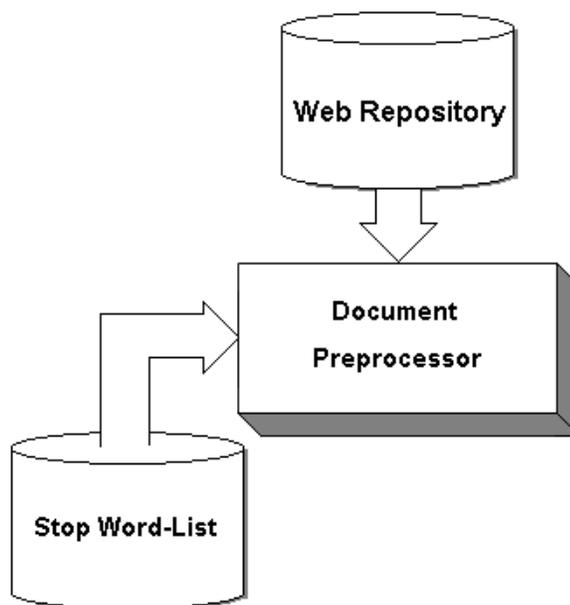


Figure 3.6: Document Preprocessor

sintetica del sito, le parole chiave che lo caratterizzano, l'autore del sito ed, eventualmente, altre informazioni. La presenza dei Meta Tag non è obbligatoria ed è cura dell'autore inserirli. Peraltro, l'utente è spinto ad immettere un contenuto coerente in tali campi in quanto alcuni motori di ricerca ne fanno uso per semplificare l'indicizzazione delle pagine e, talvolta, premiano le pagine in cui i Meta Tag sono particolarmente significativi

- **Testo**

Il contenuto semantico di una pagina è chiaramente concentrato nel testo, ma anche i Meta Tag hanno la loro importanza, in quanto sintetizzano il contenuto semantico della pagina.

Nella metrica utilizzata dal sistema realizzato, che verrà ripresa successivamente,

vengono considerati i seguenti campi della pagina HTML:

- *Testo*
- *Titolo*
- *Meta Tag "Description"*
- *Meta Tag "Keywords"*

Per estrarre tali informazioni dalla pagina è necessario analizzarne i tag HTML. Tale operazione viene effettuata dal *Document Preprocessor*, che estrapola il contenuto semantico dalla pagina HTML, lo separa nelle quattro componenti succitate e le memorizza nel database di sistema. Inoltre il *Document Preprocessor* provvederà anche ad una pulizia del documento dalle Stop Words e alla sostituzione delle stringhe che servono per la formattazione del documento in HTML con le relative combinazioni di caratteri.

3.2.5 Semantic Knowledge Base

Il cuore dell'agente di *Mining* è costituito dalla **Semantic Knowledge Base**, rappresentato in figura 3.7

In essa sono contenute le informazioni necessarie per analizzare il contenuto semantico di una pagina e misurarne l'attinenza con un certo dominio. A tal fine occorre valutare una metrica che consente di attribuire ad ogni pagina un voto che misura il grado di appartenenza della pagina ai diversi domini in esame.

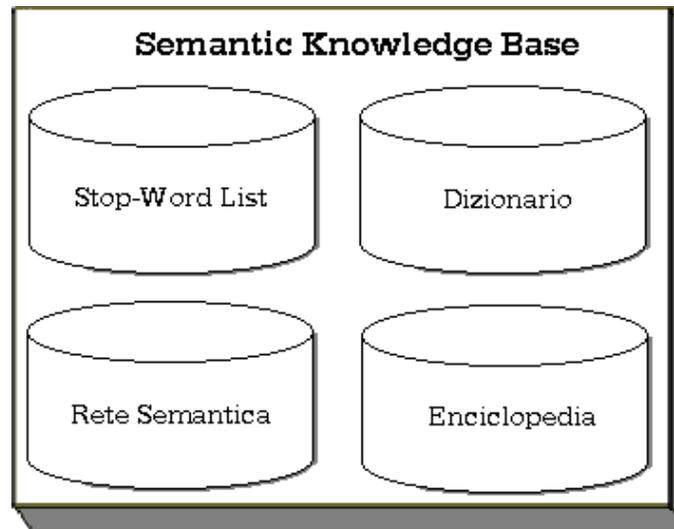


Figure 3.7: Semantic Knowledge Base

La metrica usata dal sistema realizzato, come detto in precedenza, considera cinque tipi di informazioni:

- Profondità della pagina rispetto al motore di ricerca
- Posizione mediata della pagina
- Peso dei vocaboli presenti
- Premio dell'accoppiata di vocaboli che si riferiscono a concetti tra loro correlati, pesato con la distanza dei vocaboli
- Premio dovuto a domini con un grado di correlazione diverso da zero

I primi due contributi vengono valutati dal *Repository Constructor* all'atto della costruzione del *Web Repository*.

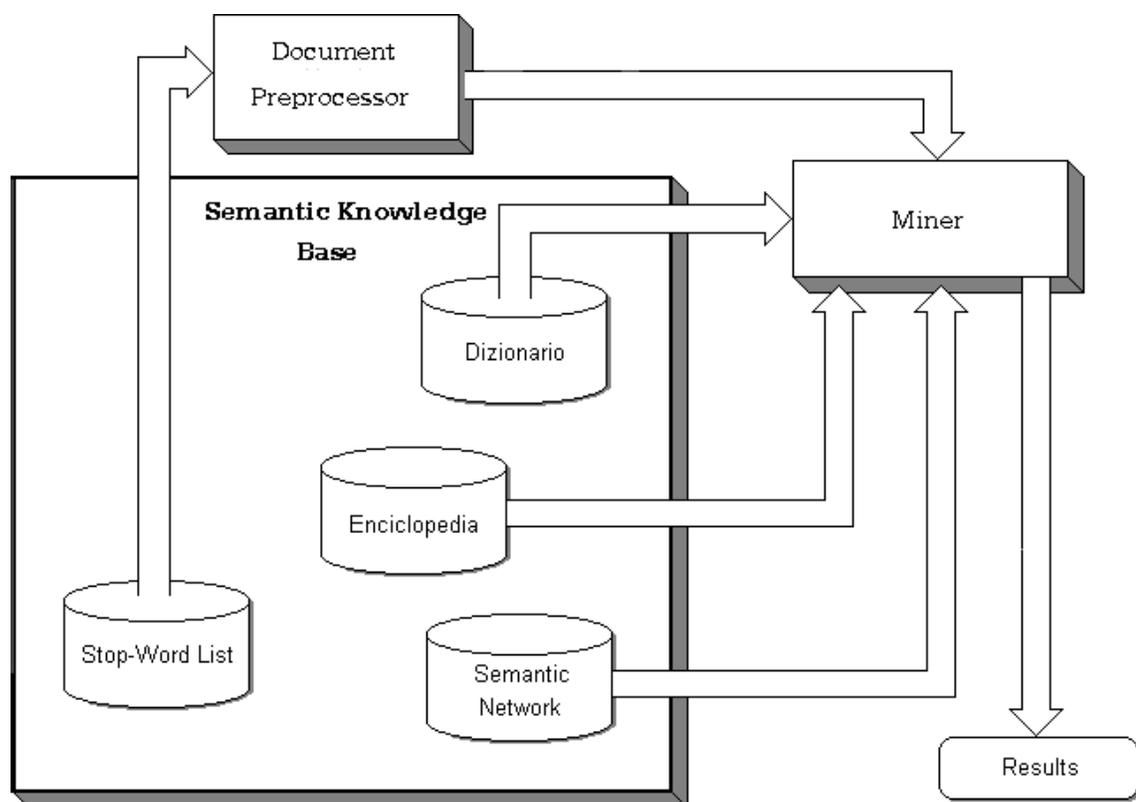


Figure 3.8: Esempio di Query

Gli ultimi tre contributi, invece, devono essere valutati dal *Mining Agent*. A tal fine, la prima operazione che esso deve effettuare è l'eliminazione delle stop-words, ossia di quelle parole che non danno contributo alla semantica del documento. Tali parole sono tipiche di ogni lingua e, poiché nel nostro sistema viene considerata esclusivamente la lingua italiana, la lista delle stop-words della lingua italiana sono memorizzate nel Database di sistema nella *Stop-Words List*.

Ottenuto il documento depurato dalle stop-words, occorre valutare i tre contributi relativi alla presenza dei vocaboli, all'accoppiata di vocaboli relativi a concetti legati in qualche modo tra loro e al grado di associazione tra i domini.

Il primo contributo è valutato confrontando i vocaboli contenuti nella pagina con le parole presenti nel dizionario, il quale contiene le parole appartenenti a tutti i domini conosciuti dal sistema sistema.

Il secondo contributo, invece, viene valutato a partire da una tabella nella quale sono presenti il concetto associato ad ogni parola trovata e la posizione del vocabolo all'interno del documento. A partire dalla tabella, si vanno a considerare le accoppiate dei vocaboli corrispondenti a concetti correlati tra loro e se ne calcola la distanza probabilistica. Tale distanza è valutata a partire dalla *Semantic Network*, che è una struttura a grafo che descrive tutte le possibili relazioni (di similitudine, di generalizzazione, specializzazione, etc.) esistenti tra concetti noti.

Infine, l'ultimo contributo è valutato tenendo presente il peso di una pagina rispetto ad un dominio, ed assegnando, in base a questo, un premio ai domini associati. Al fine di assegnare questi ultimi tre contributi la Semantic Knowledge Base contiene quattro strutture in cui sono memorizzate in maniera permanente le stop-words, il Dizionario, la rete semantica tra i concetti e la rete delle associazioni tra i domini.

3.2.6 Miner

Il **Miner**, consultando la *Semantic Knowledge Base*, valuta la metrica di attinenza delle pagine ai diversi domini conosciuti dal sistema, fornendo i risultati all'utente. Esso è rappresentato in figura 3.8

3.3 Struttura del Web Repository

Il **Web Repository** è la struttura dati che contiene le informazioni relative alle pagine Web prelevate dallo *Spider* e preprocessate dal *Document Preprocessor*. Il diagramma Entità-Relazioni è mostrato in figura 3.9.

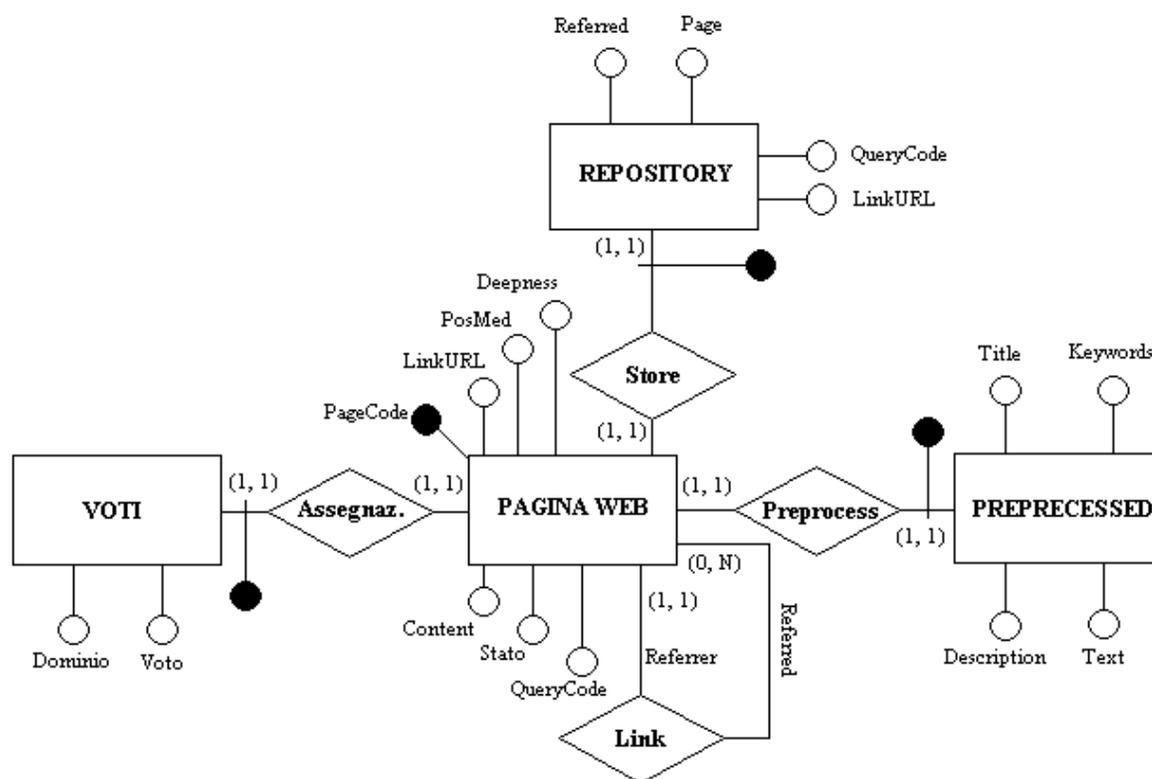


Figure 3.9: Modello E-R del Web Repository

A tale modello concettuale corrisponde il seguente schema di Database:

WEBPAGE(PageCode, LinkURL, Referrer, Deepness, PosMed, Content, Stato)

PREPROCESSED(PageCode, Title, Keywords, Description, Text)

REPOSITORY(PageCode, LinkURL, QueryCode, Page, Referrer)

3.4 Struttura della Semantic Knowledge Base

La **Semantic Knowledge Base** è composta dalla **Stop-Words List**, dal **Dizionario**, dalla **Semantic Network** e dalla **mappa dei domini**. Per quanto riguarda la lista delle stop-words, essa verrà immagazzinata in una tabella apposita con un unico campo atto a contenere le parole considerate stop-words. Per costruire lo schema delle tabelle atte a contenere il *Dizionario*, la *Semantic Network* e la *mappa dei domini* è di aiuto il diagramma Entità-Relazioni mostrato in figura 3.10

PAROLA(Name, ConceptCode, DomCode, Peso, centralità)

CONCETTO(ConceptCode, Name)

DOMINIO(DomCode, Nome)

LINK(Source, Target, costo, dominio)

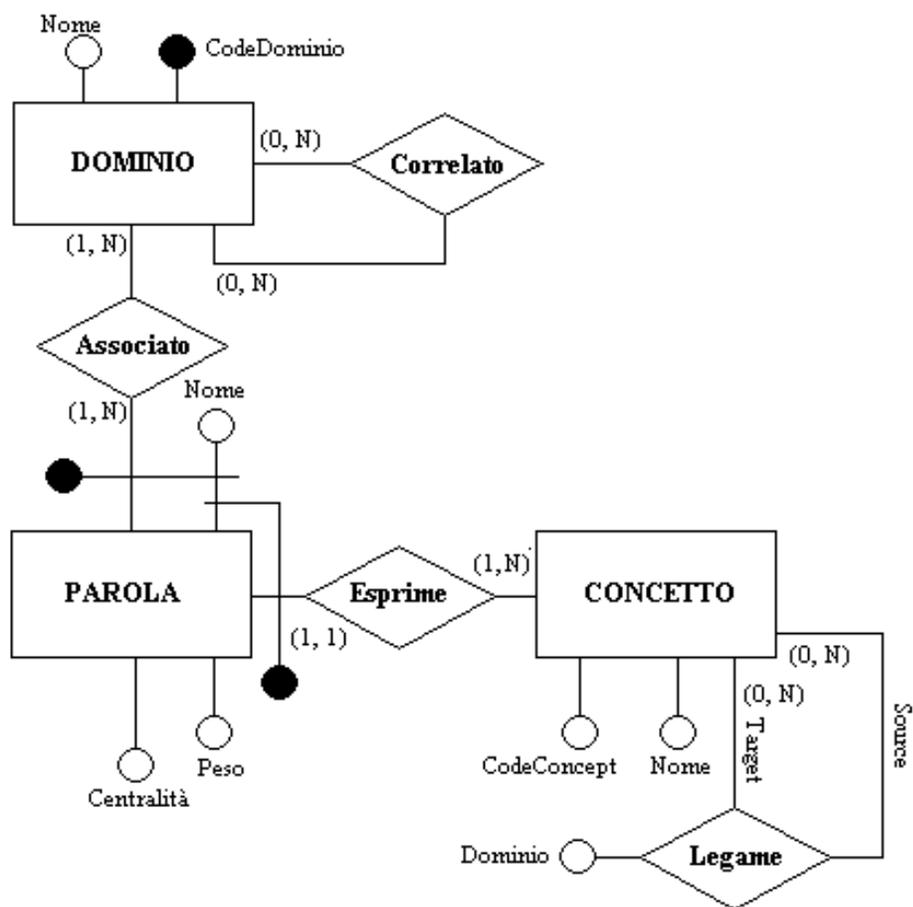


Figure 3.10: Modello E-R del Semantic Knowledge Base

Chapter 4

Sperimentazione

In questo capitolo viene mostrata una prova sperimentale effettuata sul sistema progettato. La prova è stata effettuata utilizzando la seguente configurazione della macchina:

- Processore Intel Pentium II 400 MHz
- 512 MB di memoria
- S.O. Microsoft Windows 2000 Professional
- Oracle Enterprise Edition 8.1

4.1 Parametri del sistema

Come illustrato nel capitolo relativo agli aspetti teorici, le formule che permettono di attribuire il voto semantico ad una pagina web sono le seguenti:

1. $\mathbf{Sin}(\mathbf{x}) = \frac{q(x)}{1+b \cdot K_b}$

2. $\mathbf{V}_{ss} = \sum_{p \in (t,c,d,k)} w_p \cdot (Peso_{D_i}(x) + PesoI_{D_i}(x))_p$
3. $\mathbf{V}_{sem} = \sum_{p \in (t,c,d,k)} w_p \cdot (CdDP_{D_i}(x) + CiDP_{D_i}(x))_p$
4. $\mathbf{Voto}_{D_i}(\mathbf{x}) = \frac{K_{sint} \cdot Sin(x) + K_{sem} \cdot V_{sem} + K_{ss} \cdot V_{ss}}{K_{sint} + K_{sem} + K_{ss}}$

I parametri che l'utente può inserire per gestire la modalità di votazione delle pagine sono i seguenti:

- **Costante sintattica \mathbf{k}_{sint}** , che misura l'importanza data al contributo sintattico nella valutazione della pagina
- **Costante semantica \mathbf{k}_{sem}** , che misura l'importanza data al contributo semantico nella valutazione della pagina
- **Costante sintattico-sintattica \mathbf{k}_{ss}** , che misura l'importanza data al contributo sintattico-sintattico nella valutazione della pagina
- **Soglia semantica \mathbf{T}** , che rappresenta il minimo valore di attinenza di una coppia di token ritrovata nel testo affinché la loro accoppiata possa essere considerata significativa
- **Peso del titolo \mathbf{w}_t** , che misura l'importanza data al titolo di una pagina web nel calcolo del voto semantico
- **Peso della descrizione \mathbf{w}_d** , che misura l'importanza data al contenuto del Meta Tag description di una pagina web nel calcolo del voto semantico
- **Peso delle keywords \mathbf{w}_k** , che misura l'importanza data al contenuto del Meta Tag keywords di una pagina web nel calcolo del voto semantico

- **Peso del contenuto w_c** , fattore che misura l'importanza data al contenuto semantico di una pagina web nel calcolo del voto semantico
- **Costante di penalizzazione della profondità k_b** , fattore che controlla la penalizzazione delle pagine dovuta alla loro profondità
- **Peso dei motori di ricerca**, che misurano l'importanza dei risultati restituiti da un motore di ricerca

È da notare che la somma dei pesi dati ai motori di ricerca deve essere necessariamente pari ad uno e che tutti i pesi devono essere maggiori di zero.

4.2 Ipotesi preliminari

I motori di ricerca inseriti nel sistema sono stati Google [56], Altavista [58] e Lycos [59]¹. I loro pesi sono stati scelti in base alla qualità dei motori; nella comunità web si ritiene che la qualità dei risultati restituiti da Google sia superiore di quella dei risultati restituiti da Altavista e Lycos. Pertanto, nella prova successiva, sono stati utilizzati i seguenti valori:

- Peso di Google = 0.4
- Peso di Altavista = 0.3
- Peso di Lycos = 0.3

¹Si ricordi che è comunque possibile aggiungere o eliminare motori di ricerca

La costante di penalizzazione della profondità k_b è stata fissata pari a 0.25. Tale parametro è di difficile calibrazione poiché le query a profondità maggiore di 0 sono particolarmente onerose in termini di tempi di elaborazione. Pertanto si è supposto di penalizzare il voto di ogni pagina a profondità b decrementandolo del $b \cdot 25\%$.

Inoltre si è considerato che volendo fare una analisi semantica, bisognava dare maggior peso al contenuto di una pagina e eliminare, quanto più possibile, il contributo semantico delle keywords². In aggiunta è parso opportuno assegnare un peso maggiore al titolo nei confronti della descrizione, in quanto spesso il titolo è una sintetica frase che però descrive il contenuto della pagina. Per tale motivo sono state assegnati i seguenti pesi³

- $w_c = 0.5$
- $w_t = 0.2$
- $w_d = 0.15$
- $w_k = 0.15$

Si è infine ipotizzato di ritenere non attinenti al dominio considerato quei link la cui votazione non superasse il valore 0.025.

²Essendo le keywords un insieme di parole e non una frase, dovrebbero avere un contributo semantico pressoché nullo

³Anche tenendo presente i precedenti sviuzzi del sistema

Dopo una attività di calibrazione volta all'ottimizzazione di un fattore di merito costituito dal rapporto tra numero di link restituiti dal sistema e numero di link giudicati rilevanti, si sono ottenuti i seguenti parametri ottimali:

- $K_{sint} = 1$
- $K_{sem} = 70$
- $K_{ss} = 5$
- $T = 0.1$

4.3 La prove effettuate

La prove effettuate sono state scelte utilizzando query quanto più ambigue. Per tale motivo sono state scelte effettuare le seguenti interrogazioni:

1. **Zucchero**
2. **Morandi**
3. **Bongiorno**

4.3.1 Query: Zucchero

La query zucchero è stata fatta a profondità 0, prendendo 20 link per ogni motore e utilizzando tutti i motori del sistema, come mostrato nella figura 4.1

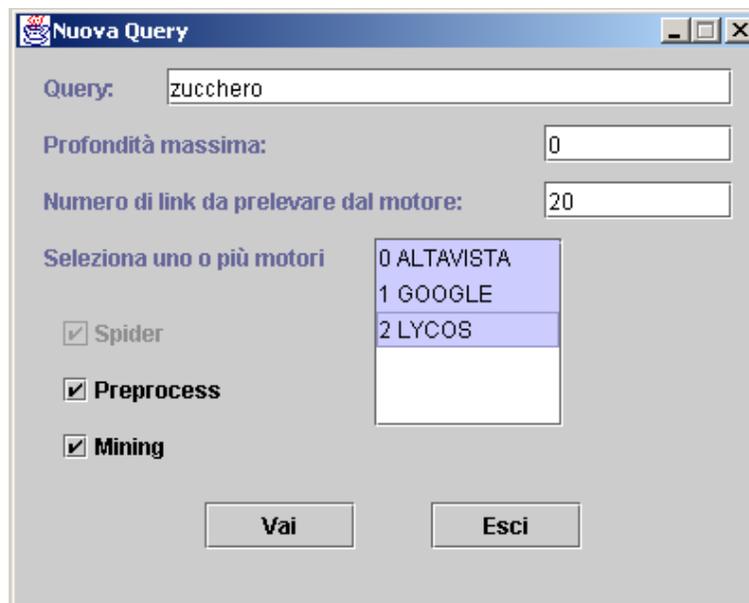


Figure 4.1: Query sottoposta al sistema: zucchero

Il numero dei link restituiti dal sistema (ossia quelli che hanno ottenuto un voto maggiore di 0.25), quello dei link giudicati rilevanti tra quelli restituiti e quello dei link rilevanti ma non restituiti sono riportati in tabella 4.1.

Nella figura 4.2 viene riportato il risultato del sistema alla query effettuata.

Per apprezzare maggiormente il risultato ottenuto si riportano, nelle figure 4.3, 4.4 e 4.5, le prime pagine di risultati della query, ottenuti dai tre motori di ricerca.

Query	Link restituiti	Link rilevanti	Link non restituiti
Zucchero	17	16	0

Table 4.1: Link restituiti e link giudicati rilevanti nella query Zucchero nel dominio Musica

Google [Ricerca avanzata](#) [Preferenze](#) [Strumenti per le lingue](#) [Suggerimenti per la ricerca](#)

zucchero

Cerca nel Web Cerca solo le pagine in Italiano

[Web](#) [Immagini](#) [Gruppi](#) [Directory](#)

Google ha cercato **zucchero** nelle pagine in **Italiano**. Risultati **1 - 10** di circa **125.000**. Durata della ricerca: **0.28** secondi.

Categorie: [World > Italiano > Arte > Musica > Artisti > Z > Zucchero](#)

ZUCCHERO
 Descrizione: Sito ufficiale: biografia, discografia, news, una fanzine, il fun club e la presentazione dell'ultimo...
 Categoria: [World > Italiano > Arte > Musica > Artisti > Z > Zucchero](#)
[www.zucchero.it/](#) - 5k - [Copia cache](#) - [Pagine simili](#)

ZUCCHERO - intro
 skip intro.
[www.zucchero.it/english/](#) - 2k - [Copia cache](#) - [Pagine simili](#)
 [Altri risultati in [www.zucchero.it](#)]

Zucchero Web Site - Shop
 Descrizione: Vendita di gadget dedicati a Zucchero Fornaciari: magliette, adesivi, berretti, occhiali, bandane,...
 Categoria: [World > Italiano > Acquisti > Musica](#)
[store.zucchero.it/](#) - 2k - [Copia cache](#) - [Pagine simili](#)

Dolcificanti Bustine di zucchero Stuzzicadenti da Personal ...
 Articoli personalizzati e prodotti per la ristorazione: bustine di zucchero, tovaglioli, tovaglette, blocchi comode. Distribuzione ...
 Descrizione: [Poggio Berni, RN] Offre prodotti per la ristorazione, come bustine di zucchero e dolcificante personaliz...
 Categoria: [World > Italiano > Affari > Alimentazione e Prodotti Correlati](#)
[www.personalzucchero.com/](#) - 4k - [Copia cache](#) - [Pagine simili](#)

ASSOCIAZIONE
 Associazione collezionisti - bustine e zollette di zucchero -Amici di Zuccheromania. ... Associazione Italiana collezionisti di bustine e zollette di zucchero...
 Descrizione: [Macerata] Dedicato agli 'Amici di zuccheromania', raccoglie informazioni e consigli per gli appassionati...
 Categoria: [World > Italiano > ... > Bustine da Zucchero](#)
[www.amicidizuccheromania.it/](#) - 12k - [Copia cache](#) - [Pagine simili](#)

Zucchero Fornaciari by Composizioni.it - artisti italiani ...
 Zucchero by Composizioni.it - Zucchero Fornaciari testi accordi midi e karaoke canzoni testi musicali. ... Zucchero Fornaciari Adelmo ...
[www.compore.it/zuccheroit.htm](#) - 42k - [Copia cache](#) - [Pagine simili](#)

Un Po' di Zucchero - Unofficial Zucchero fan club (non ufficiale)
 Un po' Di Zucchero. The Unofficial Zucchero Fornaciari fan club site.
 Il fan club non ufficiale di Zucchero. "Un po' di Zucchero ...
 Descrizione: Fan club site includes artist biography, discography, news, pictures, questions and answers, and...
 Categoria: [Arts > Music > Bands and Artists > Z > Zucchero](#)
[www.unpodizucchero.net/](#) - 3k - [Copia cache](#) - [Pagine simili](#)

Zucchero: Shake
 A 46 anni appena compiuti, Zucchero "Sugar" Fornaciari può tracciare un bilancio assai lusinghiero della propria carriera: sedici milioni di dischi venduti ...
[www.italica.rai.it/principali/argomenti/musica/zucchero2.htm](#) - 18k - [Copia cache](#) - [Pagine simili](#)

Zucchero
 Zucchero. ... L'Italia è il terzo produttore europeo di barbietola da zucchero, tuttavia siamo costretti ad importarne circa il 20 per cento del fabbisogno. ...
[www.milleunoalimenti.it/zucche.htm](#) - 9k - [Copia cache](#) - [Pagine simili](#)

Info Special
 ... mucca pazza virus bse soia biotech zucchero mercato caucci' farine animali carne fondi decreti pomodoro imprese agricoltura biodinamica transgenico vino ...
[www.dietachef.com/infospecial.asp?c=zucchero](#) - 16k - [Copia cache](#) - [Pagine simili](#)

Link non rilevanti

Figure 4.3: Risultato di Google alla query: zucchero

altavista Web Immagini MP3/Audio Video Directory Avanzata Fu

zucchero **TROVA** Ma

RICERCA: Tutto il mondo Italia **RISULTATO:** Tutte le lingue Italiano, Inglese

Perfeziona la tua ricerca con AltaVista Prisma Fai clic su un termine per perfezionare la tua ricerca e ripristinare la ricerca.

[Bustine Di Zucchero >>](#) [Articoli >>](#) [Discografia >>](#)
[Zucchero In Bustine >>](#) [Blues >>](#) [Dolciaria >>](#)
[Zucchero Sugar Fornaciari >>](#) [Commercializzazione >>](#) [Italia >>](#)

Sfondi per il tuo Desktop!!!
Questo programma ha rilevato gli sfondi più belli del momento! Fai [Entra](#) per visualizzarli subito!

Sponsor [Informazioni su](#)

Zucchero - MP3 e MIDI da scaricare !!!
Dal Pop al Rock, Hip Hop, R & B, Dance... Scegli le canzoni fra centinaia di gruppi musicali e cantanti. Classifiche italiane e straniere dei brani più scaricati.

Ingaggia per un concerto Zucchero
Café Carosello Entertainment, leader in Italia e in Europa nell'organizzazione di eventi, convention e feste aziendali in locations, discoteche, teatri e hotel esclusivi, con vip, comici e attrazioni.

AltaVista ha trovato 75.489 risultati [Informazioni su](#)

ZUCCHERO
italiano con Intro senza Intro english with Intro without Intro Attenzione! Ogni giorno ti aspettanno novità nel sito ... ricordati di svuotare la "cache" (memoria) del tuo browser per non vedere ...
[www.zucchero.it/](#) • [Aggiornato nelle ultime 48 ore](#) • [Traduci](#)
[Ulteriori pagine in www.zucchero.it](#)

ZUCCHERO - OFFICIAL SITE
The website is optimized for Internet Explorer 5.0 or higher
[www.zucchero.it/english/english.htm](#) • [Traduci](#)
[Ulteriori pagine in www.zucchero.it](#)

ASSOCIAZIONE
Associazione collezionisti bustine e zollette di **zucchero** -Amici di Zucchero mania ... Italiana collezionisti di bustine e zollette di **zucchero** **ULTIMO AGGIORNAMENTO** 13.01.2003 Versione Italiana ...
[www.amicidizucchermania.it/](#) • [Aggiornato nelle ultime 48 ore](#) • [Traduci](#)
[Ulteriori pagine in www.amicidizucchermania.it](#)

Dolcificanti Bustine di zucchero Stuzzicadenti da Personal Zucchero
Prodotti monodose per la ristorazione: **zucchero** in bustine, salviette detergenti, ... Personal **zucchero** distribuisce prodotti monodose valorizza la tua immagine con articoli personalizzati per la ...
[www.personalzucchero.com/](#) • [Aggiornato nelle ultime 48 ore](#) • [Traduci](#)
[Ulteriori pagine in www.personalzucchero.com](#)

Zucchero e Zuccheri F.lli Pinin Pero C.
Sito aziendale della Figli di Pinin Pero C. azienda produttrice di **Zucchero** Zuccheri di ... Ottimizzato 800 X 600
[www.pininpero.com/](#) • [Aggiornato nelle ultime 48 ore](#) • [Traduci](#)
[Ulteriori pagine in www.pininpero.com](#)

Zucchero - Artisti - Musica Italiana .com
Notizie, recensioni, interviste, date concerti, audioclips e molto altro su **Zucchero** ... 1a 73b6 Sei qui: Home > Artisti > **Zucchero** Mappa del sito Home News Recensioni Concerti Interviste Artisti ...
[www.musicaitaliana.com/artisti/zucchero/](#) • [Traduci](#)
[Ulteriori pagine in www.musicaitaliana.com](#)

FunkyGallo: The Zucchero Sugar Fornaciari Unofficial Site
funkygallo.supereva.it/ • [Aggiornato nelle ultime 48 ore](#) • [Traduci](#)
[Ulteriori pagine in funkygallo.supereva.it](#)

Zucchero: Shake
... Scrivici Informazioni Cerca Argomenti Shake di **Zucchero** Indice Biografia Discografia Il sito ufficiale ... Blues A 46 anni appena compiuti, **Zucchero** "Sugar" Fornaciari può tracciare un bilancio assai ...
[www.italica.rai.it/principali/argomenti/...a/zucchero2.htm](#) • [Aggiornato nelle ultime 48 ore](#) • [Traduci](#)
[Ulteriori pagine in www.italica.rai.it](#)

Microsoft Word - Accordo nazionale trasporto barbabietole da zucchero campag...
Tipo di file: PDF - [Scarica PDF Reader](#)
1. ACCORDO NAZIONALE TRASPORTO BARBABIETOLE DA **ZUCCHERO** .CAMPAGNA 2002/03. Nel quadro dell'autonomia imprenditoriale ed organizzativa dei vettori e degli utenti e della conseguente libera...
[www.confartigianatrasp.com/VARI%20PDF/...a%202002-03.pdf](#)

Zucchero Fornaciari by Composizioni.it - artisti italiani Zucchero Fornaciari spartiti testi canzoni tabs partiture accordi ...
Zucchero by Composizioni.it - **Zucchero** Fornaciari testi accordi midi e karaoke canzoni ... 3 9/Nov/2001 (Biagio Antonacci) 4 Shake (**Zucchero**) 5 Swing when you're winning (Robbie Williams) 6 ...
[www.composizioni.it/zuccheroit.htm](#) • [Aggiornato nelle ultime 48 ore](#) • [Traduci](#)
[Ulteriori pagine in www.composizioni.it](#)

Link non rilevanti

Figure 4.4: Risultato di Altavista alla query: zucchero

LYCOS
La tua guida personale su internet

Guarda.
Cerca un hotel dentro: città

Home Cerca E-mail Chat Mobile Love@Lycos Il tuo sito

La tua ricerca per pagine in italiano **Trova!** Aiuto | Ricerca avanzata | Categorie

Risultati della ricerca

Risultati dal web italiano 10 di 235348
Ulteriori risultati della ricerca: [1 2 3 4 5] >>

<< Indietro **Avanti** >>

- [Zucchero Web Site - Shop](#)
Descrizione: Vendita di gadget dedicati a Zucchero Fornaciari: magliette, adesivi, berretti, occhiali, bandane, spille e poster. Spedizioni in tutta Italia, pagamento contrassegno.
<http://store.zucchero.it>
- [Dolcificanti Bustine di zucchero Stuzzicadenti da Personal Zucchero](#)
... dolcificanti, stuzzicadenti Personal **zucchero** distribuisce prodotti monodose ... ristorazione come bustine di **zucchero**, tovaglioli, tovaglette, blocchi ... detergenti, stuzzicadenti, ...
Descrizione: Articoli personalizzati e prodotti per la ristorazione: bustine di zucchero, tovaglioli, tovaglette, blocchi comande. Distribuzione prodotti monodose: salviette detergenti, ...
<http://www.personalzucchero.com>
- [zucchero latin site \(ZLS\)](#)
Descrizione: la unica pagina de Zucchero en español
<http://www.flv.to/zucchero> | ...altri risultati da questo dominio
- [ASSOCIAZIONE](#)
... Associazione Italiana collezionisti di bustine e zollette di **zucchero** ULTIMO AGGIORNAMENTO 19.12.2008 Versione Italiana English Version Deutsche Ausführung ...
Descrizione: Associazione collezionisti - bustine e zollette di **zucchero** -Amici di Zuccheromania
<http://www.amici dizuccheromania.it>
- [Zucchero e Zuccheri F.lli Pinin Pero & C.](#)
Descrizione: Sito aziendale della Figli di Pinin Pero & C. azienda produttrice di Zucchero & Zuccheri di altissima qualità
<http://www.pininpero.com>
- [ZUCCHERO - OFFICIAL SITE](#)
Descrizione: Il sito è ottimizzato per Internet Explorer 5.0 e versioni successive
<http://www.zucchero.it/italiano.htm> | ...altri risultati da questo dominio
- [Camillo Shop - dischi rari e da collezione - zucchero "sugar" fornaciari](#)
... rare records, vinyl, cd s for collectors **Zucchero** ARTISTA TITOLO SUPP ANNO ETICHETTA STAMPA NOTE COND ?
ZUCCHERO D'oro incenso e birra CDS3" 90 POLYDOR ... Diavolo in me -12" ...
Descrizione: Camillo Shop, dischi rari e da collezione
<http://www.camilloshop.com/zucchero.htm>
- [Humanitasonline](#)
Humanitasonline Consulta il dizionario Inserisci una parola per trovare informazioni e proposte sulla salute **Zucchero**
zucchero termine con cui si indica commercialmente il saccarosio, disaccaride ...
Descrizione: Un nuovo punto di riferimento per la salute!Tanti check up su misura, consulenze mediche qualificate per un altro parere clinico, dizionario medico, prenotazione on line . tutte ...
<http://www.humanitasonline.com/HOL/index.cfm?ID=26181>
- [ALIMENTAZIONE](#)
... come il saccarosio (il normale **zucchero** che usiamo in cucina o al bar ... scissi nelle singole molecole di **zucchero**.A livello intestinale queste ... costituiti da una sola molecola di ...
Descrizione: Wayfitness è il portale del benessere in movimento: dal fitness al body building, dallo step alla palestra. Per mantenersi in forma allenando il corpo e rilassando la mente.
http://www.wayfitness.net/root/175_384.html
- [Confetture, marmellate, prodotti senza zucchero - Naso&Gola](#)
... Tisane, infusi e the Olio Aceto balsamico Funghi Prodotti senza **zucchero** Pasticceria Antiche ricette di salute IN EVIDENZA... BALSAMO ... acquisterà su Naso&Gola (vedi regolamento). ...
<http://www.nasoegola.com/ita/shop/sugarless.html>

Ulteriori risultati della ricerca: [1 2 3 4 5] >>

<< Indietro **Avanti** >>

Link non rilevanti

Figure 4.5: Risultato di Lycos alla query: zucchero

4.3.2 Query: Morandi

La query morandi è stata fatta a profondità 0, prendendo 20 link per ogni motore e utilizzando tutti i motori del sistema.

Il numero dei link restituiti dal sistema, quello dei link giudicati rilevanti e quello dei link rilevanti ma non restituiti dal sistema sono riportati in tabella 4.2

Nella figura 4.6 viene riportato il risultato del sistema alla query effettuata ⁴.

Query	Link restituiti	Link rilevanti	Link non restituiti
Morandi	5	5	2

Table 4.2: Link restituiti e link giudicati rilevanti nella query Morandi nel dominio Musica

Per apprezzare maggiormente il risultato ottenuto si riportano, anche in questo caso,, nelle figure 4.7, 4.8 e 4.9, le prime pagine di risultati della query, ottenuti dai tre motori di ricerca.

⁴Si noti che una delle due pagine non restituite ha una pagina di redirect.



Query = Morandi

Tempo di Spidering= 0ms
 Tempo di Preprocessing= 0ms
 Tempo di Mining= 485067ms
 Tempo Totale= 485067ms

Risultati nel dominio = MUSICA

- 1 - Come fa bene il Giannino - **Voto: 0.2354**
 &COME FA BENE IL GIANNINO& | NEWS | LA STORIA | DISCOGRAFIA | GIANNI AL CINEMA | TESTI E ACCORDI | TOUR | C'ERA UN RAGAZZO | UNO DI NOI | | RACCONTI (CON FOTO) | VOTA IL GIANNINO | GIANNI-ARCHIVIO | MESSAGGI | LINKS | GUESTBOOK | CONTATTAMI | SITO CREATO INTERAMENTE DA: GABRIELE...
- 2 - A.e.P. Amati-arrangiatori per Morandi Ramazzotti ed altri - **Voto: 0.03949**
 AP BEAT MUSIC ON WEB HOME > FRIENDS BIO MP3 SERVIZI STUDIO FRIENDS CHI E' LA SUA PAGINA CONTIENE NOTIZIE SU: MARCO ADAMI AUTOINTERVISTA FRANCESCO DESMAELE MORANDI, PERCENTONETTO, BARBARA COLA, GIULIA, GROSSOMODO, A. SAMMARINI GROSSOMODO EMERGENZA ROCK, ROBERTO GATTO DANIELE MARCELLI ESACORDO...
- 3 - Morandi - **Voto: 0.03026**
 1971 (RIST. 2000) - UN MONDO DI DONNE - GIANNI MORANDI - RCA PLS 10517 CANTA: LA CANZONE DI MARINELLA CLICK TO ENLARGE COVER BY ORESTE CLICK TO ENLARGE COVER BY ORESTE
- 4 - Gianni Morandi - Artisti - Musica Italiana .com - **Voto: 0.0278**
 SEI QUI: HOME > ARTISTI > GIANNI MORANDI MAPPA DEL SITO HOME NEWS RECENSIONI CONCERTI INTERVISTE ARTISTI SHOPPING COMMUNITY ...:CLUBS ...:CHAT ...:FORUM ...:ISCRIVITI ...:MODIFICA DATI ...:E.M.M.A. TOP20 CLICCANTE FOTO LINK MISS & MISTER SANREMO NETWORK ----- DICCI TUTTO! CERCA: IN ...
- 5 - Marco Morandi - Artisti - Musica Italiana .com - **Voto: 0.0255**
 SEI QUI: HOME > ARTISTI > MARCO MORANDI MAPPA DEL SITO HOME NEWS RECENSIONI CONCERTI INTERVISTE ARTISTI SHOPPING COMMUNITY ...:CLUBS ...:CHAT ...:FORUM ...:ISCRIVITI ...:MODIFICA DATI ...:E.M.M.A. TOP20 CLICCANTE FOTO LINK MISS & MISTER SANREMO NETWORK ----- DICCI TUTTO! CERCA: IN ...
- ***** VALORI AL DI SOTTO DELLA SOGLIA *****
- 6 - Recensioni e critica musicale - Guida di superEva - **Voto: 0.02256**
 SITE MAP FAI DI SUPEREVA LA TUA HOMEPAGE INTRATTENIMENTO E SPETTACOLO RECENSIONI E CRITICA MUSICALE DI ANGY PLATANIA , UNA DELLE 500 GUIDE NELLA GUIDA NEL CANALE INTRATTENIMENTO E SPETTACOLO IN SUPEREVA IN ITALIA NEL MONDO HOME NEWSLETTER FORUM CHAT LA GUIDA RISPONDE ...
- 7 - Camillo Shop - dischi rari e da collezione - gianni morandi - **Voto: 0.02057**
 DISCHI RARI, VINILE, CD DA COLLEZIONE RARE RECORDS, VINYL, CD S FOR COLLECTORS GIANNI MORANDI ARTISTA TITOLO SUPP ANNO ETICHETTA STAMPA NOTE COND --MORANDI GIANNI ANDAVO A CENTO ALL'ORA CD55" 99 BMG 74321-64977-2 ITA 2 TRACKS P/S INCL. "LOREDANA" M 11,00 MORANDI GIANNI BANANE E LAMPONE CD55" 93 RCA...
- 8 - Amati-Appbeat. Arrangiatori per Morandi, Ramazzotti, Mietta. - **Voto: 0.01986** ← **Pagina interessante**
- 9 - Hotel Morandi - Sanremo - **Voto: 0.01658**
 HOTEL MORANDI CORSO MATUZIA 51 - 18038 SANREMO ...
- 10 - L'Opera di Giorgio Morandi (Bologna 1890 - 1964) arriva a... - **Voto: 0.01404**
 GIORGIO MORANDI (BOLOGNA 1890 - 1964) CHERASCO - PALAZZO SALMATORIS 12 OTTOBRE 05 DICEMBRE 2002 L'OPERA DI GIORGIO MORANDI (BOLOGNA 1890 - 1964) ARRIVA A PALAZZO SALMATORIS DI CHERASCO DOPO LE RETROSPETTIVE DEDICATE A PABLO PICASSO, ANTONIO LIGABUE, MASSIMO CAMPILGI, FILIPPO DE PISIS, GIORGIO DE...
- 11 - Collezione Peggy Guggenheim - Artisti - Giorgio Morandi... - **Voto: 0.01287**
 HOME NEWS! INFORMAZIONE MOSTRE PADIGLIONE USA COLLEZIONE ARTISTI LA STORIA MEMBERSHIP INTRAPRES/E INTERNSHIP SHOP CAFE UFFICIO STAMPA LINKS I SPEAK ENGLISH INTRODUZIONE ARTISTI A-C ALBERS , ANNI ALBERS , JOSEF ALECHINSKY , PIERRE APOLLONIO , MARINA APPEL , KAREL ARCHIPENKO , ALEXANDER ARMAN...
- 12 - Museo Morandi - **Voto: 0.01226**
 IL NUOVO SITO DEL MUSEO MORANDI È VISIBILE ALL'INDIRIZZO: HTTP://WWW.MUSEOMORANDI.IT/ THE MUSEO MORANDI SITE IS NOW AT: HTTP://WWW.MUSEOMORANDI.IT/
- 13 - Museo Morandi - **Voto: 0.01204**
 LA PAGINA CORRENTE UTILIZZA I FRAME. QUESTA CARATTERISTICA NON È SUPPORTATA DAL BROWSER IN USO.
- 14 - Gianni Morandi - **Voto: 0.01172**
 SCANSERVICE MUSICA PRESENTA GIANNI MORANDI CELESTE AZZURRO E BLU I SUCCESSI DI MORANDI MORANDI IN TEATRO MORANDI MORANDI QUESTA È LA STORIA: SCENDE LA PIOGGIA QUESTA È LA STORIA QUESTA È LA STORIA: ANDAVO A 100 ALL'ORA TUTTI I SUCCESSI VOL.1 TUTTI I SUCCESSI VOL.2 VARIETA 30VOLTEMORANDI ...
- 15 - Autori - **Voto: 0.01121**
 UNIVERSITÀ DI BOLOGNA IST. DI ANATOMIA UMANA NORMALE GLI AUTORI LELLI MANZOLINI MORANDI SUSINI ASTORRI ALTRI ERCOLE LELLI (BOLOGNA 1702-1766) F IGLIO DI UN ARCHIBUGIERE, DOPO AVER TRASCORSO UN PERIODO DI APPRENDISTATO PRESSO LA BOTTEGA PATERNA, ERCOLE LELLI INTRAPRESE LA CARRIERA ARTISTICA...
- 16 - Museo Morandi - **Voto: 0.0111**
 LA PAGINA CORRENTE UTILIZZA I FRAME. QUESTA CARATTERISTICA NON È SUPPORTATA DAL BROWSER IN USO.
- 17 - Emilio Morandi - **Voto: 0.01056**
 VIVE ED OPERA A PONTE NOGSA. SI OCCUPA DI PITTURA, È OPERATORE DI INSTALLAZIONI E PERFORMER.FA PARTE DEL GRUPPO "PULS/PLUS" - GRUPPO DI RICERCA ARTE VISIVA.COLLABORA ALL'ARCHIVIO SONORO "V.E.C." DI MAASTRICHT - OLANDA.DIRIGE LA GALLERIA "ARTESTUDIO" DI PONTE NOGSA (BG), SPAZIO APERTO AD ARTISTI CHE...
- 18 - Marco Morandi - Sito non ufficiale - **Voto: 0.009891**
 CLICK HERE TO CONTINUE. MARCO MORANDI - SITO NON UFFICIALE IL PRIMO SITO NON UFFICIALE DI MARCO MORANDI MARCO MORANDI.
- 19 - Comunicato Stampa -Giorgio Morandi- - **Voto: 0.009654**
 COMUNE DI FORL&IGRAVE GIORGIO MORANDI PITTORE E INCISORE MOSTRA : GIORGIO MORANDI, PITTORE E INCISORE LUOGO : FORL&IGRAVE, ORATORIO DI SAN SEBASTIANO, PIAZZA GUIDO DA MONTEFELTRO, TEL. 0543/21424 INAUGURAZIONE : 24 APRILE, ORE 17,30 DURATA : 25 APRILE - 20 LUGLIO 1997 ORARI ...
- 20 - Italy Posters and Prints, Movie Posters, Cuisine Posters,... - **Voto: 0.009568**
 SEARCH: MOST POPULAR ITALIAN POSTERS ITALIAN VINTAGE POSTERS ITALIAN LANDSCAPES PANORAMIC VIEWS OF ITALY ITALIAN OPERA POSTERS ITALIAN MOVIE POSTERS ITALIAN ARTISTS' GALLERIES MOVIEGOODS ORIGINALS FRAMED FINE ART OF ITALY POSTERIDEA ARTWORKS ON CANVAS 2003 ITALY'S CALENDARS ITALIAN LANDSCAPES...
- 21 - Morandi Tappeti - tappeti orientali, Tappeti... - **Voto: 0.009264**
 IL GUSTO ITALIANO PER IL TAPPETO PERSIANO.

Figure 4.6: Risultato del sistema alla query: morandi

Google [Ricerca avanzata](#) [Preferenze](#) [Strumenti per le lingue](#) [Suggerimenti per la ricerca](#)

morandi

Cerca nel Web Cerca solo le pagine in Italiano

Web Immagini Gruppi Directory

Google ha cercato **morandi** nelle pagine in Italiano. Risultati 1 - 10 di circa 52,500. Durata della ricerca: 0.07 secondi.

Categorie: [World > Deutsch > ... > M > Morandi, Giorgio](#) [World > Italiano > Consultazione > Musei > Arte e Intrattenimento](#)

Museo Morandi
 Il nuovo sito del Museo **Morandi** è visibile all'indirizzo: <http://www.museomorandi.it/>.
 The Museo **Morandi** Site is now at: <http://www.museomorandi.it/>.
 Descrizione: Ospitato nel Palazzo Comunale di Bologna, offre informazioni sull'artista, sugli orari del museo,...
 Categoria: [World > Italiano > ... > Arte e Intrattenimento > Musei](#)
www.comune.bologna.it/bologna1/Cultura/Museicomun/Morandi/ - 1k - [Copia cache](#) - [Pagine simili](#)

Museo Morandi
www.museomorandi.it/ - 5k - [Copia cache](#) - [Pagine simili](#)

Museo Morandi
www.museomorandi.it/index_exp.htm - 2k - [Copia cache](#) - [Pagine simili](#)
 [Altri risultati in www.museomorandi.it/]

MorandiMania
www.morandimania.it/ - 2k - [Copia cache](#) - [Pagine simili](#)

Nuova pagina 1 - [[Traduci questa pagina](#)]
 Hotel. **MORANDI ALLA GROCCETTA**. Firenze. Via Laura 50, 50121 Firenze ITALY. Tel. 055 234 4 747 Fax 055 248 0 954. e-mail: welcomed@hotelmorandi.it. www.hotelmorandi.it/ ...
 Descrizione: Description, some images and reservation request. Housed in a former convent. Close to the Duomo.
 Categoria: [Regional > Europe > ... > Travel and Tourism > Lodging > Hotels](#)
www.hotelmorandi.it/ - 10k - [Copia cache](#) - [Pagine simili](#)

Edizione Morandi-Libri sul i paesaggi el'arte
 Una biografia nuovo di Giorgio **Morandi** sul i paesaggi e più di 250 pagine di testo, illustrazione, e le interviste recente con contemporani di **Morandi**...
 Descrizione: I paesaggi di **Morandi** e più di 250 pagine di testo, illustrazione, e le interviste recente con contempor...
 Categoria: [World > Italiano > Acquisti > Pubblicazioni > Libri](#)
www.morandieditions.com/index_italia.html - 5k - [Copia cache](#) - [Pagine simili](#)

Morandi Tappeti - tappeti orientali, Tappeti persiani, antichi, ...
Morandi Tappeti - Vendita tappeti antichi,prove di ambientazione e consulenza, lavaggi manuali, controlli e manutenzione, laboratorio di restauro interno ...
 Descrizione: [Castelvetro Piacentino, PC] Vendita, consulenza, manutenzione e restauro tappeti antichi. E' presente...
 Categoria: [World > Italiano > ... > Antiquariato e Collezionismo](#)
www.moranditappeti.it/ - 2k - [Copia cache](#) - [Pagine simili](#)

Titolo pagina
 contatti. servizi. home. a chi è rivolto. entra.
 Descrizione: [Litta Parodi, AL] Offre servizi di project management per l'industria elettromeccanica ed impiantistica...
 Categoria: [World > Italiano > Affari > Beni e Servizi per l'Industria](#)
spaziowind.libero.it/projectmanager/ - 6k - [Copia cache](#) - [Pagine simili](#)

PML Morandi: rubinetteria, miscelatori monocomando e articoli per ...
 La PML di **Morandi** produce e distribuisce articoli idrosanitari, comprendente rubinetteria tradizionale, miscelatori monocomando e articoli per doccia. ...
 Descrizione: Produce e distribuisce articoli idrosanitari, comprendente rubinetteria tradizionale, miscelatori...
 Categoria: [World > Italiano > Affari > Edilizia > Impianti](#)
www.pml-morandi.it/ - 5k - [Copia cache](#) - [Pagine simili](#)

Caffè Morandi Basket
 CaffèMorandi.it. Sito in costruzione. Tutto sulla squadra di basket del Caffè **Morandi**. Le partite. I giocatori. Foto. Scrivete al webmaster. Hosted by.
 Descrizione: Sito ufficiale di una società di basket amatoriale di Bologna. Le schede dei giocatori ei tabellini...
 Categoria: [World > Italiano > Sport > Pallacanestro > Società](#)
www.caffemorandi.it/ - 2k - [Copia cache](#) - [Pagine simili](#)

Link non rilevanti

Figure 4.7: Risultato di Google alla query: morandi

altavista Web Immagini MP3/Audio Video Directory

morandi **TROVA** Maggiore precisione

RICERCA: Tutto il mondo Italia RISULTATO: Tutte le lingue Italiano, Inglese

Perfeziona la tua ricerca con AltaVista Prisma Fai clic su un termine per perfezionare la tua ricerca. Fai clic su >> per ripi

Gianni Morandi >> Museo Morandi >> Musica Italiana >>
 Giorgio Morandi >> Il Sito >> Anni >>
 Marco Morandi >> La Mostra >> Attività >>

Bellissimi sfondi per il tuo desktop!!! @ 100% (RGB)
www.SfondiFantastici.com stop it!

AltaVista ha trovato 31.296 risultati [Informazioni su](#)

Museo Morandi
 Questo sito è ottimizzato per Microsoft Internet Explorer 4.0 (o versione successiva) e Netscape Communicator 4.0 (o versione successiva). Inoltre il browser deve essere abilitato ad eseguire ...
www.museomorandi.it/ • Aggiornato nelle ultime 48 ore • Traduci
 Ulteriori pagine in www.museomorandi.it

Museo Morandi
 Progetto culturale realizzato nell'ambito delle iniziative per Bologna 2000 Città Europea della Cultura
 Coordinamento e Gestione del Progetto: Museo **Morandi** - Bologna. Titolarità del progetto...
www.museomorandi.it/index_exp.htm • Traduci
 Ulteriori pagine in www.museomorandi.it

Museo Morandi
 Il nuovo sito del Museo **Morandi** è visibile all'indirizzo: <http://www.museomorandi.it/> The Museo **Morandi** Site is now at: <http://www.museomorandi.it/>
www.comune.bologna.it/Cultura/Museicomun...ditaliano.html • Aggiornato nelle ultime 48 ore • Traduci
 Ulteriori pagine in www.comune.bologna.it

Gianni Morandi - Artisti - Musica Italiana .com
 Notizie, recensioni, interviste, date concerti, audioclips e molto altro su Gianni **Morandi** ... 1a.7112 Sei qui: Home > Artisti > Gianni **Morandi** Mappa del sito Home News Recensioni Concerti Interviste ...
www.musicaitaliana.com/artisti/giannimorandi/ • Traduci
 Ulteriori pagine in www.musicaitaliana.com

Morandi Tappeti - tappeti orientali, Tappeti persiani, antichi, caucasici, restauro
Morandi Tappeti - Vendita tappeti antichi, prove di ambientazione, lavaggi manuali, ... la fiducia accordata, proseguire nella ricerca del meglio da offrirvi. Fabio **Morandi** Host & Development by A&A
www.moranditappeti.it/ • Traduci
 Ulteriori pagine in www.moranditappeti.it

Giorgio Morandi
www.romagna.net/gmorandi/uno.htm • Aggiornato nelle ultime 48 ore • Traduci
 Ulteriori pagine in www.romagna.net

Marco Morandi - Sito non ufficiale
 IL PRIMO SITO NON UFFICIALE DI MARCO **MORANDI** ... Click here to continue. Marco **Morandi** - Sito non ufficiale
 IL PRIMO SITO NON UFFICIALE DI MARCO **MORANDI** Marco **Morandi**
www.marcomorandi.cjb.net/ • Traduci

Autori
 Università di Bologna Ist. di Anatomia Umana Normale Gli Autori Lelli Manzolini **Morandi** Susini Astorri Altri Ercole Lelli (Bologna 1702-1766) Figlio di un archibugiere, dopo aver trascorso un ...
biocfarm.unibo.it/museocere/autori.htm • Aggiornato nelle ultime 48 ore • Traduci
 Ulteriori pagine in biocfarm.unibo.it

PML Morandi: rubinetteria, miscelatori monocomando e articoli per doccia
 La PML di **Morandi** produce e distribuisce articoli idrosanitari, comprendente rubinetteria tradizionale, miscelatori monocomando e articoli per doccia.
www.pml-morandi.it/ • Traduci
 Ulteriori pagine in www.pml-morandi.it

La Repubblica/spettacoli e cultura: Morandi batte De Filippi Con Fiorello non c'è partita
 ... tiene comunque testa al concorrente **Morandi** batte De Filippi Con ... è toccato a Gianni **Morandi** spuntare la percentuale più alta ... (24%). Il programma di **Morandi**, nel periodo di sovrapposizione ...
www.repubblica.it/online/spettacoli_e_cu...morandidue.html • Aggiornato nelle ultime 48 ore • Traduci
 Ulteriori pagine in www.repubblica.it

Link non rilevanti

Figure 4.8: Risultato di Altavista alla query: morandi

LYCOS
La tua guida personale su internet

SMS
349 | 1234567 | Invio sms Animati e 160 caratteri!!!

Home Cerca E-mail Chat Mobile Love@Lycos Il tuo sito

La tua ricerca per pagine in italiano **Trova!** [Aiuto](#) | [Ricerca avanzata](#) | [Categorie](#)

Risultati della ricerca

Ultimissime novità >> [Ulteriori Notizie](#)
[Sanremo, Gerini e Autieri sul palco accanto a Baudo](#)
La Repubblica - Giovedì, 1.16.2003 - 20:35 Ore

Risultati dal web italiano 10 di 112421
Ulteriori risultati della ricerca: [1 2 3 4 5] >>

- [Edizione Morandi-Libri sul i paesaggi e l'arte](#)
... quassu I PAESAGGI DI GIORGIO MORANDI SETTANTA TONALITE DI VERDE ... Finalmente una biografia intimo di **Morandi** con efasi su i paesaggi. Tradotto ... primo tempo. Commentario su ...
Descrizione: I paesaggi di Morandi e più di 250 pagine di testo, illustrazione, e le interviste recente con contemporanei di Morandi. Informazioni e ordini per questa edizione limitata.
http://www.morandieditions.com/index_italia.html
- [Museo Morandi](#)
Il nuovo sito del Museo **Morandi** è visibile all'indirizzo: <http://www.museomorandi.it/> The Museo **Morandi** Site is now at: <http://www.museomorandi.it>
Descrizione: [Bologna] Ospitato nel Palazzo Comunale, offre informazioni sull'artista, sugli orari del museo, ed alcune immagini della collezione.
<http://www.comune.bologna.it/bologna1/Cultura/Museicomun/Mor...> | ...altri risultati da questo dominio
- [Morandi Tappeti - tappeti orientali, Tappeti persiani, antichi, caucasici, restauro](#)
Morandi Tappeti - Tappeti antichi, manutenzione e restauro Consulenza di arredamento gratuita ... Ringrazio per la fiducia accordata, proseguirò nella ricerca del meglio da offrirvi. Fabio ...
Descrizione: Morandi Tappeti - Vendita tappeti antichi, prove di ambientazione e consulenza, lavaggi manuali, controlli e manutenzione, laboratorio di restauro interno, grandi restauri eseguiti ...
<http://www.moranditappeti.it>
- [PML Morandi: rubinetteria, miscelatori monocomando e articoli per doccia](#)
PML **Morandi**: rubinetteria, miscelatori monocomando e articoli per doccia PML **Morandi**: rubinetteria, miscelatori monocomando e articoli per doccia ... miscelatori monocomando e articoli per doccia. ...
Descrizione: La PML di Morandi produce e distribuisce articoli idrosanitari, comprendente rubinetteria tradizionale, miscelatori monocomando e articoli per doccia.
<http://www.pml-morandi.it>
- [Benvenuto dalla Cartaria Morandi S.p.A.](#)
Descrizione: [Milano] Recupero e commercio di carta da macero, presentazione delle infrastrutture e dei servizi offerti.
<http://www.cartaria-morandi.com>
- [Caffè Morandi Basket](#)
CaffèMorandi.it Sito in costruzione Tutto sulla squadra di basket del Caffè **Morandi** Le partite I giocatori-Foto Scrivete al webmaster Hosted by
Descrizione: Sito ufficiale della squadra di basket amatoriale Caffè **Morandi** Bologna
<http://www.caffemorandi.it>
- [Nuova pagina 1](#)
Hotel **MORANDI** ALLA CROCETTA Firenze Via Laura 50, 50121 Firenze ITALY Tel. 055 234 4 ...
Descrizione: Description, some images and reservation request. Housed in a former convent. Close to the Duomo.
<http://www.hotelmorandi.it>
- [Morandi Maurizio & Luigi](#)
Morandi snc: idea, segno, forma nel ferro, lavorazione artistica metallo
Descrizione: Morandi snc: idea, segno, forma nel ferro, lavorazione artistica metallo
<http://www.morandisnc.it>
- [Sito non ufficiale di Gianni Morandi](#)
Un invito a donare sangue, soprattutto nei mesi estivi Gianni **Morandi** e Laura Fogli Webdesign: Manuela Cianfarini G I A N N M O R A N D I U N O F F I C I A L S I T E M O R A N D I U N O F F I C I A L S I T E
Descrizione: Tutto quello che volete sulla vita di Gianni Morandi
<http://www.giannimorandi.info>
- [Discoteca Morandi unofficial](#)
<http://www.morandi.too.it>

Ulteriori risultati della ricerca: [1 2 3 4 5] >>

<< Indietro **Avanti** >>

Link non rilevanti

Figure 4.9: Risultato di Lycos alla query: morandi

4.3.3 Query: Bongiorno

La query bongiorno è stata fatta a profondità 0, prendendo 10 link per ogni motore e utilizzando tutti i motori del sistema.

Il numero dei link restituiti dal sistema, quello dei link giudicati rilevanti e quello dei link non restituiti anche se giudicati rilevanti sono riportati in tabella 4.3

Query	Link restituiti	Link rilevanti	Link non restituiti
Bongiorno 2	2	0	

Table 4.3: Link restituiti e link giudicati rilevanti nella query Bongiorno nel dominio Televisione

Chapter 5

Interfaccia utente del sistema SSA

In questo capitolo viene presentata l'interfaccia utente del sistema e vengono spiegati singolarmente la funzione dei vari pulsanti in essa contenuta.

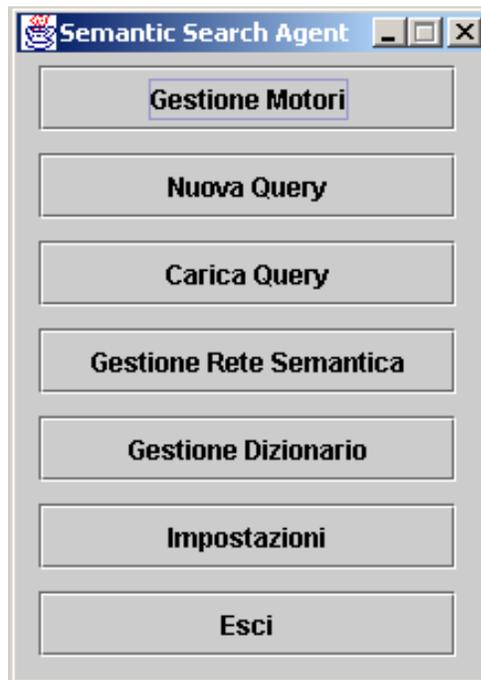


Figure 5.1: Interfaccia principale del sistema SSA

La suddetta interfaccia si presenta come nella figura 5.1.

In essa sono presenti i seguenti pulsanti:

- **Gestione Motori:** per la gestione dei motori di ricerca utilizzati dal sistema
- **Nuova Query:** per la gestione di una nuova query da sottoporre al sistema
- **Carica Query:** per la gestione delle query già effettuate e memorizzate nel repository
- **Gestione Rete Semantica:** per la gestione delle reti semantiche utilizzata dal sistema
- **Gestione Dizionario:** per la gestione del dizionario multi-dominio inserito nel sistema
- **Impostazioni:** per la gestione dei parametri dell'applicazione
- **Esci:** per terminare l'applicazione

5.1 Gestione Motori

La pressione del pulsante Gestione Motori provoca l'apertura della finestra di dialogo illustrata nella figura 5.2.

Nella parte in alto della finestra di dialogo è presente l'elenco dei motori di ricerca conosciuti dal sistema. Selezionandone uno è cliccando sul pulsante *Cancella motore* è possibile eliminare un motore di ricerca dal sistema.

Nella parte inferiore è possibile inserire le informazioni per l'aggiunta di un motore di ricerca tra quelli conosciuti dal sistema; tali informazioni sono le seguenti:

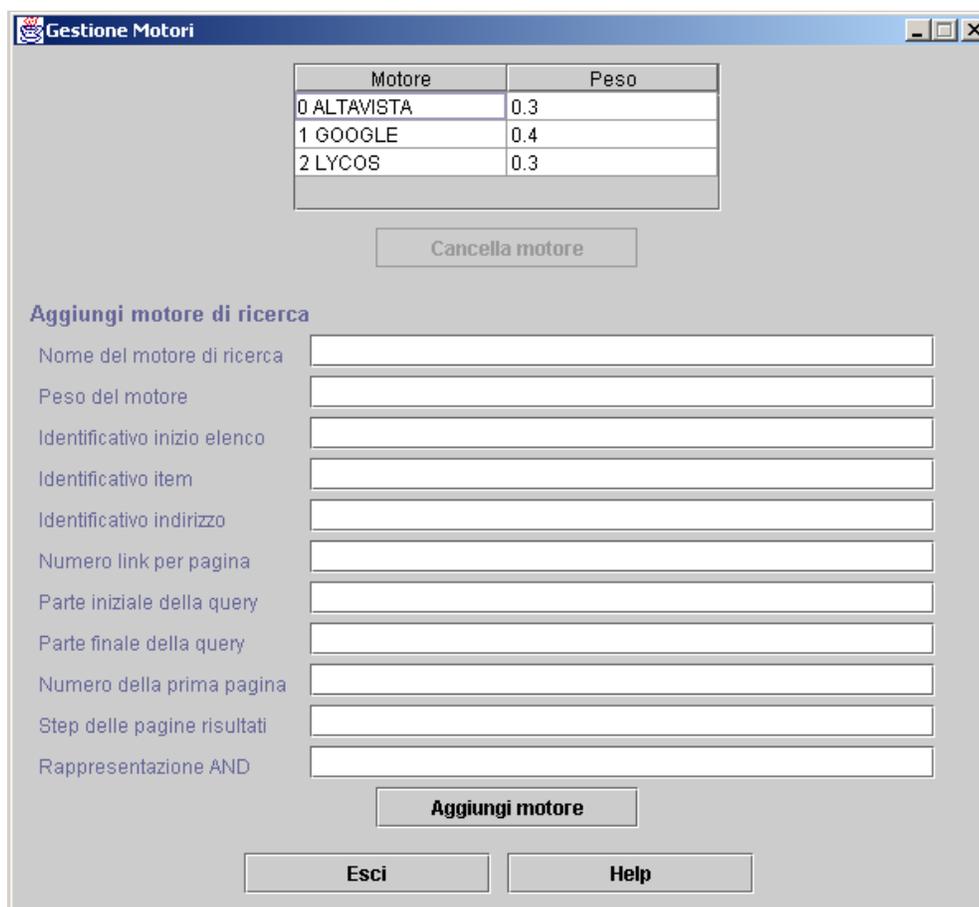


Figure 5.2: Finestra di dialogo per la gestione dei motori di ricerca

- **Nome del motore di ricerca:** tale informazione potrebbe essere anche un identificativo simbolico e non il nome reale del motore di ricerca.
- **Peso del motore:** tale valore rappresenta il peso da dare ai risultati della ricerca del motore. Una volta inserito il peso bisogna cambiare quello degli altri motori affinché la somma sia 1¹.
- **Identificativo inizio elenco:** stringa presente nella pagina HTML, risultato della ricerca del motore, che ci permetta di individuare in modo univoco la parte in cui inizia l'elenco delle pagine trovate
- **Identificativo item:** stringa presente nella pagina HTML, risultato della ricerca del motore, che ci permetta di individuare in modo univoco un punto elenco delle pagine trovate
- **Identificativo indirizzo:** stringa presente nella pagina HTML, risultato della ricerca del motore, che ci permetta di individuare in modo univoco un indirizzo dell'elenco delle pagine trovate
- **Numero di link per pagina:** numero di link trovati per ogni pagina risultato del motore. Ogni motore restituisce una pagina con un certo numero di risultati e poi un link alla successiva pagina con altri risultati
- **Parte iniziale della query:** stringa che rappresenta la query di ricerca fino al punto in cui ci sono le parole cercate adattate.

¹Si veda anche la sezione relativa alle Impostazioni.

- **Parte finale della query:** stringa che rappresenta la query di ricerca dal punto in cui ci sono le parole cercate adattate, fino alla posizione in cui inserire il numero della pagina contenente i risultati
- **Numero della prima pagina:** ogni motore identifica la prima pagina contenente i risultati della ricerca con un intero. Tale intero è solitamente 0
- **Step delle pagine risultato:** passo che intercorre tra gli identificativi (numeri interi) di due pagine contenenti i risultati, restituite dal motore
- **Rappresentazione AND:** stringa che indica come il motore rappresenta l'AND nella query adattata (solitamente +)

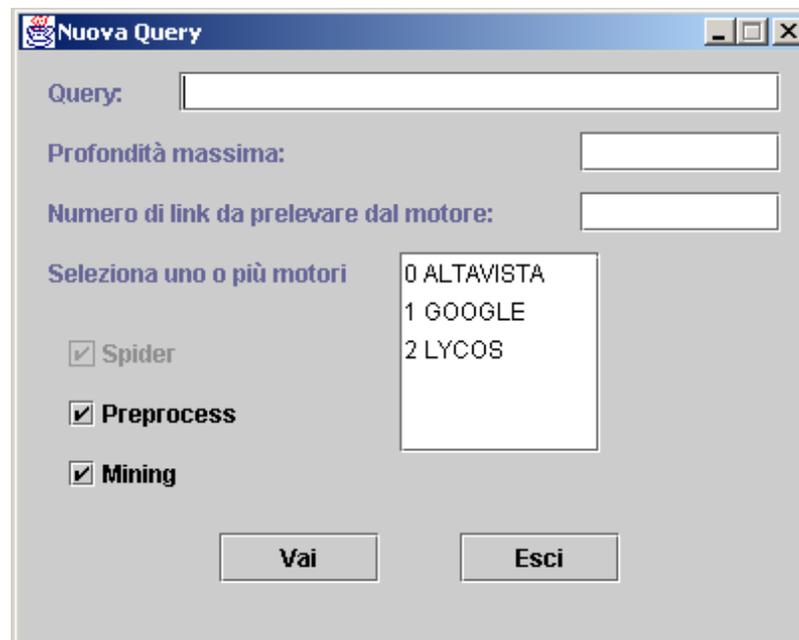


Figure 5.3: Interfaccia Nuova Query

5.2 Nuova Query

La pressione del pulsante Nuova Query provoca l'apertura della finestra di dialogo illustrata nella figura 5.3.

- **Query:** vengono inserite le keywords che verranno utilizzate dai motori di ricerca per il reperimento delle pagina Web. La sintassi adottata è quella di Google[57,60] che prevede che tutte le keywords immesse siano in AND tra loro, ossia devono essere presenti tutte in una pagina affinché la stessa venga restituita dal motore. È anche possibile immettere query contenenti frasi (o gruppi di parole) racchiuse tra virgolette; la stringa di caratteri racchiuse tra virgolette deve essere presente nella pagina affinché venga restituita dal motore.
- **Profondità massima:** deve essere inserita la profondità massima che una pagina deve avere, rispetto al motore di ricerca, affinché questa venga presa in considerazione. Si usa la convenzione che le pagine restituite dal motore abbiano profondità 0, mentre le pagine ottenute perchè linkate da qualche altra pagina abbiano profondità pari alla profondità della pagina linkante aumentata di 1. Il concetto di profondità è illustrato nella figura 5.4 e nella figura 5.5².
- **Numero di link da prelevare da ogni motore di ricerca:** deve essere inserito il numero massimo di link che si desidera considerare tra quelli restituiti dai singoli motori di ricerca utilizzati.
- **Motori di ricerca:** è necessario infine selezionare almeno un motore tra quelli conosciuti dal sistema. È possibile anche selezionare più motori di ricerca.

²Si noti che la pagina di figura 5.5 è presente tra i link di livello 0 della figura 5.4

Google Ricerca avanzata Preferenze Strumenti per le lingue Suggestimenti per la ricerca

IEEE Cerca con Google

Cerca nel Web Cerca solo le pagine in Italiano

Web Immagini Gruppi Directory

Google ha cercato IEEE nelle pagine in Italiano. Risultati 1 - 10 di circa 57,000. Durata della ricerca: 0.03 secondi.

[IEEE Central & South Italy Section](#)
Home page of the IEEE C&S Section. Useful links to other IEEE sites and lot of news for IEEE members. ...
[iee.ing.uniroma1.it/ - 22k - Copia cache - Pagine simili](#)

[IEEE Student Branch Trieste](#)
IEEE, IEEE Student Branch Trieste The Institute of Electrical and Electronics Engineers - Trieste Student Branch, Università degli Studi di Trieste, ...
[www.univ.trieste.it/~ieeesb/ - 9k - Copia cache - Pagine simili](#)

[IV2000 Home Page](#)
[www.ce.unipr.it/iv2000/ - 1k - Copia cache - Pagine simili](#)

[IEEE Student Branch Catania - COS'E' L'IEEE](#)
... In tutti questi anni l'IEEE è divenuto catalizzatore e punto di riferimento per tutto ciò che riguarda l'innovazione tecnologica nell'ambito delle ...
[iee.dees.unict.it/iee/ - 9k - Copia cache - Pagine simili](#)

Link di livello 0

Figure 5.4: Esempio di link di livello 0

 **IEEE Student Branch Trieste**
The Institute of Electrical and Electronics Engineers - Trieste Student Branch 

NEWS: IEEE Student Branch Trieste_ Ricerca: GO Language: Italiano

IEEE	Attività	Servizi	Links
L'istituto	Seminari	Proposte per tesi di laurea	Dipartimento - DEE
Lo Student Branch di Trieste	Visite guidate	Proposte di lavoro	Facoltà di Ingegneria
Perche' iscriversi?		Piani di studio ed ordinamenti	Università di Trieste
HomePage IEEE		Materiale didattico (tutor)	Trieste (HomePage Comune)
IEEE Xplore		NewsGroup (solo lettura)	Altri links

Link di livello 1

Qualsiasi domanda riguardante lo Student Branch può essere rivolta al comitato: ieeesb@ing.univ.trieste.it
IEEE Student Branch Trieste HomePage - <http://www.univ.trieste.it/~ieeesb>
Pagine a cura del comitato dello Student Branch di Trieste

Visite dal 01/05/1998: 7 1 2 1 - Ultimo aggiornamento: 16/1/2003

Risoluzione ottimale per la visione: 800x600 - Si consiglia l'uso di un Browser che supporti JavaScript

[\[Switch to english\]](#)

Figure 5.5: Esempio di link di livello 1

Una volta immessi i dati necessari alla ricerca e selezionati i motori da utilizzare, è possibile scegliere quali operazioni far compiere al sistema. Poiché si sta sottoponendo una nuova query, la casella di Spidering risulta già spuntata e non modificabile. Tale situazione è giustificata dal fatto che non è possibile effettuare le altre operazioni previste senza che le pagine siano state prelevate e memorizzate nel sistema (che è compito dello Spider).

Le operazioni che si possono compiere (oltre allo Spidering) sono le seguenti:

- **Preprocess:** per l'analisi delle pagine web al fine di ottenerne il contenuto semantico (ossia la pagina senza i tag HTML), il contenuto del campo Title e il contenuto dei Meta Tag Description e Keywords
- **Mining:** per l'analisi semantica delle pagine web al fine di ricavarne il voto semantico che rappresenta il grado di attinenza della pagina al dominio considerato

È da notare che l'operazione di *Mining* può avvenire solo a valle dell'operazione di *Preprocess* in quanto la prima utilizza, per la valutazione delle pagine web, i risultati della seconda. Pertanto, non è consentito effettuare il Mining senza che si sia anche scelto di effettuare il *Preprocess* delle pagine.

Alla pressione del **tasto Vai**, se non si è scelto di effettuare né il *Preprocess* né il *Mining*, viene restituita una pagina HTML dove viene presentato l'elenco dei link trovati ordinati in base alla profondità e alla posizione. Quest'ultimo valore dipende dal tipo di link: se è uno di livello zero, rappresenta la posizione mediata all'interno dei motori selezionati (come già spiegato); se è di livello maggiore, rappresenta la posizione mediata del suo migliore link referrer.

Nel caso in cui venga selezionata anche l'operazione di *Mining*, viene restituita una diversa pagina HTML contenente, stavolta, i link ordinati in base al voto attribuito dal sistema alle pagine web corrispondenti. Di ogni pagina web si riporta il titolo (se presente) e parte del testo. All'interno della pagina HTML costruita, si riportano anche i tempi (sempre in millisecondi) impiegati per effettuare le operazioni di Spidering, Preprocess e Mining.

Se si sceglie solo l'operazione di *Preprocess*, viene restituita la stessa pagina HTML costruita nel caso in cui si opti solo per lo Spidering e, in aggiunta, viene effettuato il preprocess delle pagine web (come già spiegato) andando a riempire le specifiche tabelle del database di sistema.

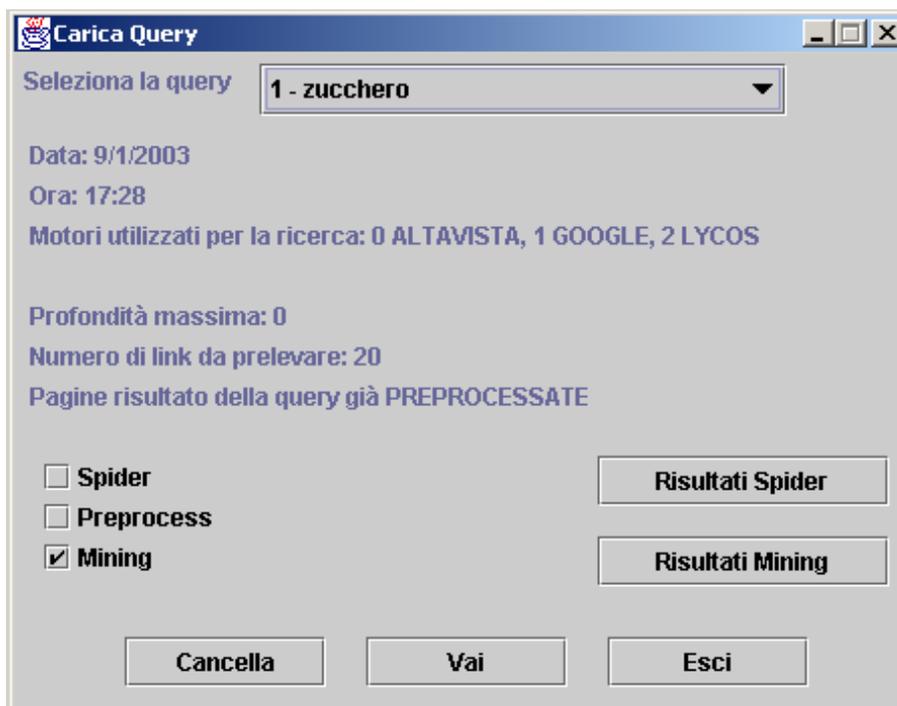


Figure 5.6: Interfaccia Carica Query

5.3 Carica Query

La funzione Carica Query è preposta alla gestione delle query già sottoposte al sistema. Essa si presenta come illustrato nella figura 5.6

Viene presentato, tramite una casella a scorrimento, l'elenco delle query già effettuate.

Selezionata una di esse, vengono restituite le seguenti informazioni:

- Data e ora dell'ultimo Spidering effettuato
- I motori di ricerca utilizzati
- La massima profondità di ricerca
- Il numero massimo di link da prelevare
- Se le pagine web associate alla query sono state o meno preprocessate

È possibile decidere se effettuare nuovamente lo *Spidering* oppure il *Preprocess* o il *Mining*. La prima operazione risulta utile nel caso in cui si ritenga che la query selezionata sia datata e, quindi, necessiti di un nuovo Spidering; la seconda e la terza permettono rispettivamente di preprocessare e di effettuare l'analisi semantica delle pagine web nel caso in cui ciò non era stato fatto quando la query veniva sottomessa per la prima volta al sistema. L'opzione di Mining risulta utile anche nel caso in cui si siano apportate delle modifiche ai parametri utilizzati dal sistema per l'attribuzione del voto alle pagine oppure modifiche alla Rete Semantica o alla Token List e, quindi, è richiesta una nuova analisi semantica delle pagine. È da notare che, nel caso in cui si voglia effettuare lo Spidering modificando i motori utilizzati oppure modificando il

numero di link da prelevare o la profondità di ricerca, ciò deve essere fatto andando ad impostare una nuova query tramite la funzionalità già descritta.

Allorquando si sceglie di effettuare un nuovo Spidering della query selezionata, tutti i vecchi dati ad essa associati vengono cancellati dal sistema e aggiornati. È da notare ancora che la casella Preprocessing risulta già spuntata nel caso in cui le pagine associate alla query siano ancora da preprocessare e che non è consentito effettuare il Mining se a monte non è stato prima effettuato il Preprocess delle pagine.

Come esposto in Nuova Query, alla pressione del tasto Vai, a seconda se si è scelto di effettuare solo lo Spidering oppure anche il Mining, viene costruita una specifica pagina HTML con i risultati dell'operazione. All'interno della finestra sono presenti i tre pulsanti: Risultati Spider, Risultati Miner e Cancella.

Il primo consente di vedere la pagina HTML che riporta i risultati dell'ultimo Spidering effettuato. Il secondo è attivo soltanto se le pagine restituite dalla query risultano già analizzate semanticamente attraverso la rete semantica e la token list e, quindi, presentano un voto associato. Alla pressione di tale tasto viene presentata la pagina HTML che riporta i risultati dell'ultima operazione di mining effettuata. Il terzo pulsante consente di cancellare tutte le informazioni memorizzate nel sistema associate alla query selezionata.

5.4 Gestione Rete Semantica

Alla pressione del pulsante Gestione Rete Semantica viene presentata l'interfaccia mostrata nella figura 5.7 in cui sono presenti i seguenti pulsanti:

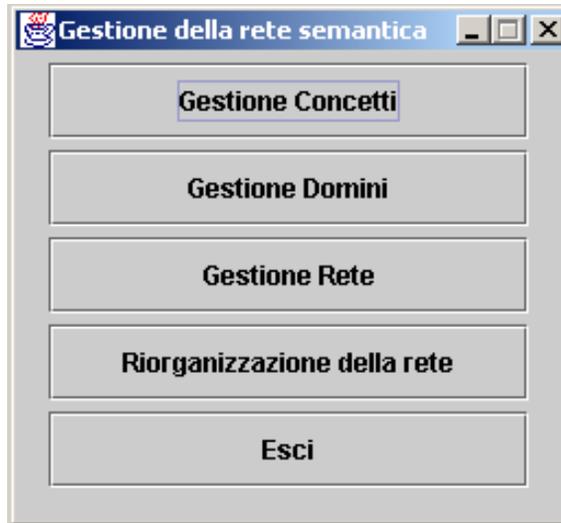


Figure 5.7: Interfaccia della Gestione della Rete Semantica

- **Gestione Concetti:** per la gestione dei concetti (eliminazione o aggiunta)
- **Gestione Domini:** per la gestione dei domini (eliminazione, aggiunta o definizione della loro correlazione)
- **Gestione Rete:** per la gestione della Semantic Network dei concetti.
- **Riorganizzazione della rete:** per riorganizzare la Semantic Network e la mappa dei domini dopo una modifica
- **Esci:** per uscire dal menù.³

³Si noti che se si tenta di uscire senza aver fatto la riorganizzazione viene comunque chiesto se si desidera farla.

5.4.1 Gestione Concetti

La funzionalità per la gestione dei concetti si presenta con la finestra di dialogo mostrata nella figura 5.8.

Nella parte alta della finestra è presente un elenco dei concetti presenti. Da tale elenco è possibile selezionarne uno e, premendo il pulsante *Cancella concetto*, eliminarlo.



Figure 5.8: Finestra di dialogo per la gestione dei concetti

Nella parte bassa è possibile scrivere il nome di un nuovo concetto e, premendo il pulsante *Aggiungi*, inserirlo nel sistema.

5.4.2 Gestione Domini

La funzionalità per la gestione dei domini si presenta con la finestra di dialogo mostrata nella figura 5.9.

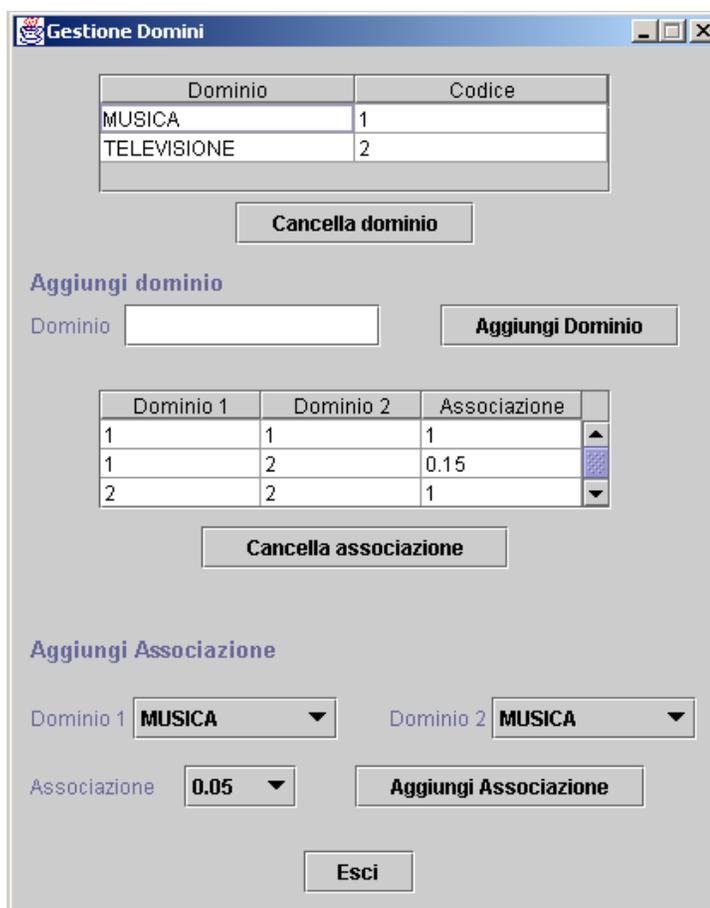


Figure 5.9: Finestra di dialogo per la gestione dei domini

Nella parte alta della finestra è presente un elenco dei domini presenti. Da tale elenco è possibile selezionarne uno e, premendo il pulsante *Cancella dominio*, eliminarlo.

Immediatamente sotto è possibile scrivere il nome di un nuovo dominio e, premendo il pulsante *Aggiungi dominio*, inserirlo nel sistema.

Ancora più in basso vi è un elenco delle associazioni tra i concetti. Da tale elenco è possibile selezionare una associazione e, premendo il pulsante *Cancella associazione*, eliminarla.

Infine, nella parte bassa della finestra c'è la possibilità di aggiungere una nuova associazione tra i domini scegliendo i due domini interessati, il grado di associazione e premendo il pulsante *Aggiungi associazione*⁴.

5.4.3 Gestione Rete

La finestra per la gestione della Semantic Network tra i concetti si presenta come in figura 5.10.

Innanzitutto è necessario scegliere il dominio della Semantic Network che si vuole gestire⁵.

Nella parte centrale c'è l'elenco degli archi della rete. Da tale elenco è possibile selezionare un arco ed eliminarlo mediante la pressione del pulsante *Cancella arco*.

Nella parte inferiore della finestra è possibile scegliere due concetti, da altrettanti elenchi a discesa, e un peso per inserire un nuovo arco nella rete semantica del dominio selezionato.

⁴Si ricordi che la mappa dei domini è un grafo i cui archi non sono orientati quindi non può essere inserito un arco dal dominio D_i al dominio D_j se esiste già l'arco da D_j a D_i

⁵Si noti che anche se nel Database le reti semantiche vengono gestite come una unica rete i cui archi pesati portano anche l'informazione del dominio in cui hanno senso, dal punto di vista dell'utente è comodo pensare che ci siano diverse reti semantiche una per ogni dominio.

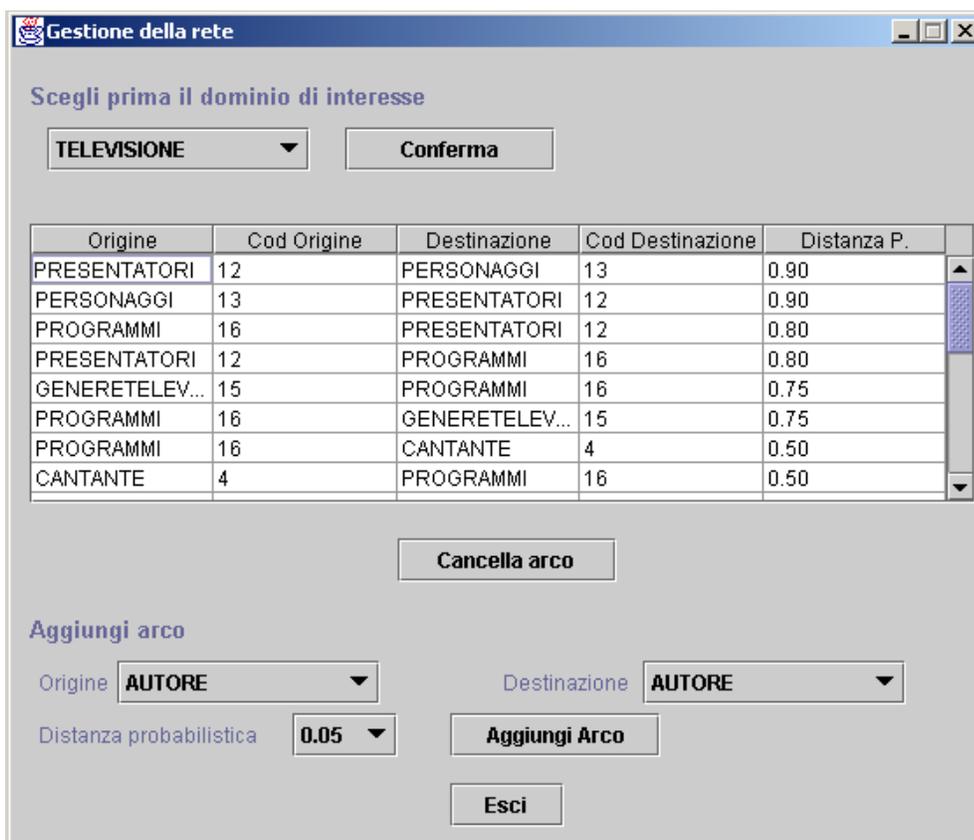


Figure 5.10: Finestra di dialogo per la gestione della Semantic Network

Gestione Dizionario

Nuovo termine da aggiungere

Concetto associato **AUTORE** Dominio di appartenenza **MUSICA**

Peso Centralità

Aggiungi parola

Termine	Concetto	Dominio	Centralità	Peso
ACCORDI	TERMINITECNICI	MUSICA	0.85	0.55
ACCORDO	TERMINITECNICI	MUSICA	0.85	0.55
ADAGIO	TEMPIMUSICALI	MUSICA	0.90	0.35
ADAMS	CANTANTE	MUSICA	0.90	0.40
AFRE	CASADISCOGRA...	MUSICA	0.85	0.45
ALBUM	TERMINITECNICI	MUSICA	0.80	0.70
ALLEGRETTO	TEMPIMUSICALI	MUSICA	0.85	0.65
ALLEGRO	TEMPIMUSICALI	MUSICA	0.70	0.40
ALLEMANDA	BALLO	MUSICA	0.90	0.80
AMADEUS	PRESENTATORI	TELEVISIONE	0.85	0.80
AMBRA	PERSONAGGI	TELEVISIONE	0.80	0.40

Esci **Cancella la parola**

Figure 5.11: Finestra di dialogo per la gestione del Dizionario contenente le parole conosciute dal sistema

5.5 Gestione Dizionario

La finestra per la gestione del Dizionario si presenta come in figura 5.11.

Nella parte alta della finestra c'è la sezione mediante la quale è possibile aggiungere una nuova parola nel dizionario. Per fare tale operazione è necessario inserire il vocabolo, il concetto associato, il dominio di appartenenza, la centralità del vocabolo rispetto al concetto e il peso del vocabolo nel dominio. Una volta inseriti le precedenti informazioni la nuova parola può essere aggiunta al Dizionario premendo il pulsante *Aggiungi parola*.

Nella parte inferiore della finestra vi è l'elenco di tutte le parole presenti nel dizionario. È possibile selezionarne una e cancellarla mediante la pressione del pulsante *Cancella Parola*.

5.6 Impostazioni

La funzionalità *Impostazioni* è preposta alla gestione dei parametri dell'applicazione. La funzione si presenta come nella figura 5.12:

La finestra mostra le ultime impostazioni salvate dall'utente. Si possono impostare i seguenti parametri:

- **Pesi dei motori di ricerca:** la somma di tali pesi deve necessariamente essere 1. Tali pesi misurano l'importanza dei risultati restituiti dai singoli motori.
- **Peso del contenuto:** fattore che misura l'importanza data al testo di una pagina web nel calcolo del voto semantico

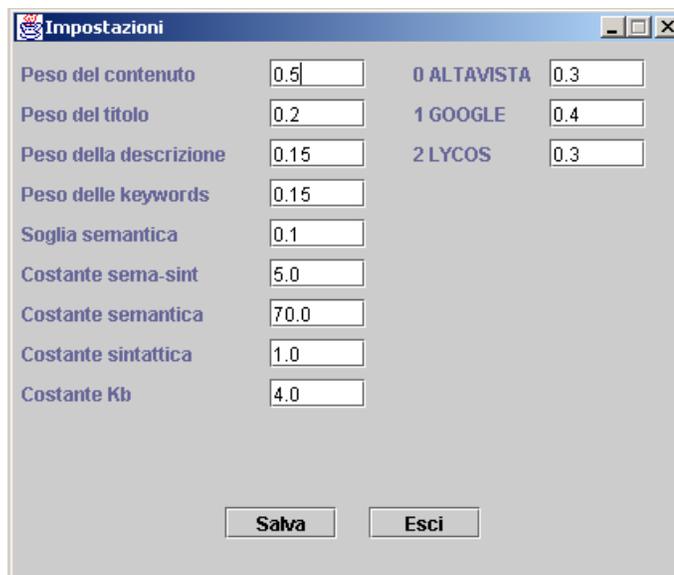


Figure 5.12: Finestra di dialogo la modifica delle impostazioni

- **Peso del titolo:** fattore che misura l'importanza data al titolo di una pagina web nel calcolo del voto semantico
- **Peso della descrizione:** fattore che misura l'importanza data al contenuto del Meta Tag description di una pagina web nel calcolo del voto semantico
- **Peso delle keywords:** fattore che misura l'importanza data al contenuto del Meta Tag keywords di una pagina web nel calcolo del voto semantico
- **Soglia semantica:** minimo valore di attinenza di una coppia di vocaboli ritrovata nel testo affinché la loro accoppiata possa essere considerata significativa.

Si noti che la somma dei pesi relativi al *contenuto*, *titolo*, *descrizione* e *keywords* deve essere necessariamente pari ad uno; inoltre tutti i pesi devono essere maggiori di zero.

Alla pressione del tasto Salva tutti i valori inseriti verranno memorizzati all'interno del sistema.

Bibliography

- [1] A. Castellucci, G. Ianni, D. Vasile, S. Costa - "*Searching and surfing the web using a semi-adaptative meta-engine*" - IEEE, 2001
- [2] H. S. Nwana - "*Software Agents: An Overview*" - Cambridge University, 1996. - <http://www.sce.carleton.ca/netmanage/docs/AgentsOverview/ao.html>
- [3] Darren Greaves - "*Intelligent Search Agents*" - <http://krapplets.org/report/report.html>
- [4] R. Cooley, B. Mobasher, J. Srivastava - "*Web Mining: Information and Pattern Discovery on the World Wide Web*" - IEEE, 1997 - <http://maya.cs.depaul.edu/mobasher/webminer/survey/survey.html>
- [5] V. S. Subrahmanian, Piero Bonatti, Jurgen Dix, Thomas Eiter, Sarit Kraus, Fatma Ozcan, Robert Ross - "*Heterogeneous Agent System*" - The MIT Press Cambridge, Massachusetts, London, England.
- [6] Ferber, J. - "*Simulating with Reactive Agents*" - In Hillebrand, E. & Stender, J. (Eds.), *Many Agent Simulation and Artificial Life*, Amsterdam: IOS Press, 8-28, 1994
- [7] M. Wooldridge, N. Jennings - "*Intelligent Agents: Theory and Practice*" - The Knowledge Engineering Review 10 (2), 115-152, 1995

- [8] M. N. Huhns - M. P. Singh - "*Distributed Artificial Intelligence for Information Systems*" - CKBS-94 Tutorial, June 15, University of Keele, UK, 1994
- [9] E. H. Durfee - V. R. Lesser - D. Corkill - "*Coherent Cooperation among Communicating Problem Solvers*" - IEEE Transactions on Computers C-36(11), 1275-1291, 1987
- [10] <http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-5/www/pleiades.html>
- [11] K. Sycara - "*Intelligent Agents and the Information Revolution*" - UNICOM Seminar on Intelligent Agents and their Business Applications, 8-9 November 1995, London, 143-159,
- [12] M. Georgeff - "*Agents with Motivation: Essential Technology for Real World Applications*" - The First International Conference on the Practical Applications of Intelligent Agents and Multi-Agent Technology, London, UK, 24th April 1996
- [13] R. Smith - Personal Communication - 1996
- [14] H. S. Nwana, L. Lee, N. R. Jennings - "*Coordination in Software Agent Systems*" - British Telecommunications Technology Journal 14 (4), October 1996
- [15] P. Maes - "*Agents that Reduce Work and Information Overload*" - Communications of the ACM 37 (7), 31-40, 1994
- [16] P. Maes - "*Intelligent Software*" - Scientific American 273 (3), September 1995
- [17] R. Kozierok, P. Maes - "*A Learning Interface Agent for Scheduling Meetings*" - Proceedings of the ACM-SIGCHI International Workshop on Intelligent User Interfaces, Florida, 81-93, 1993
- [18] H. Lieberman - "*Letizia: An Agent that Assists Web Browsing*" - In Proceedings of IJCAI 95, AAAI Press, 1995

- [19] B. J. Rhodes, T. Starner - "*Remembrance Agent: A Continuously Automated Information Retrieval System*" - In Proceedings the First International Conference on the Practical Application of Intelligent Agents and Multi-Agent Technology (PAAM 96), London, April 1996
- [20] A. Chavez, P. Maes - "*Kasbah: An Agent Marketplace for Buying and Selling Goods*" - In Proceedings the First International Conference on the Practical Application of Intelligent Agents and Multi-Agent Technology (PAAM 96), London, April 1996
- [21] P. Wayner - "*Agents Unleashed: A Public Domain Look at Agent Technology*" - Boston, MA: AP Professional, 1995
- [22] <http://www.genmagic.com>
- [23] P. Wayner - "*Free Agents*" - Byte, March, 1995
- [24] K. Indermaur - "*Baby Steps*" - Byte, March, 1995
- [25] O. Etzioni, D. Weld - "*A Softbot-Based Interface to the Internet*" - Communications of the ACM 37 (7), 72-76, 1994
- [26] N. J. Davies, R. Weeks - "*Jasper: Communicating Information Agents*" - in Proceedings of the 4th International Conference on the World Wide Web, Boston, USA, December 1995
- [27] P. Maes - "*Designing Autonomous Agents: Theory and Practice from Biology to Engineering and Back*" - London: The MIT press, 1991
- [28] R. A. Brooks - "*Intelligence without Representation*" - Artificial Intelligence 47, 139-159, 1991
- [29] R. A. Brooks - "*A Robust Layered Control System for a Mobile Robot*" - IEEE Journal of Robotics and Automation 2 (1), 14-23, 1986

- [30] H. S Nwana - "*Simulating a Childrens Playground in ABLE*" - Working Report, Department of Computer Science, Keele University, UK, 1993
- [31] P. Wavish, M. Graham - "*A Situated Action Approach to Implementing Characters in Computer Games*" - Applied AI Journal, 1995
- [32] K. Fisher, J. P. Muller, M. Pischel - "*Unifying Control in a Layered Agent Architecture*" - Technical Report TM-94-05, German Research Center for AI - (DFKI GmbH), 1996
- [33] J. P. Muller, M. Pishel, M. Thiel - "*Modelling Reactive Behaviour in Vertically Layered Agent Architectures*" - In Wooldridge, M. & Jennings, N. (eds.) (1995b), Intelligent Agents, Lecture Notes in Artificial Intelligence 890, Heidelberg: Springer Verlag, 261-276, 1995
- [34] I. A. Ferguson - "*TouringMachines: An Architecture for Dynamic, Rational, Mobile Agents*" - PhD Thesis, Computer Laboratory, University of Cambridge, UK, 1992
- [35] M. R. Genesereth, S. P. Ketchpel - "*Software Agents*" - Communications of the ACM 37 (7), 48-53, 1994
- [36] G. Wiederhold - "*Mediators in the Architecture of Future Information Systems*" - IEEE Computer 25 (3), 38-49, 1992
- [37] M. R. Cutkosky, R. S. Engelmores, R. E. Fikes, M. R. Genesereth, T. R. Gruber, J. M. Tenenbaum, J. C. Weber - "*PACT: An Experiment in Integrating Concurrent Engineering Systems*" - IEEE Computer 1, January, 28-37, 1993
- [38] Wang Jicheng, Huang Yuan, Wu Gangshan, Zhang Fuyan - "*Web Mining: Knowledge Discovery on the Web*" - Department of Computer Science and Technology, Nanjing University

- [39] H. Kawano - "*Overview of Mondou web search engine using text mining and information visualizing technologies*" - IEEE, 2001
- [40] P. Atzeni, S. Ceri, S. Paraboschi, R. Torlone - "*Basi di dati - Seconda Edizione*" - McGraw-Hill, 1999
- [41] C. M. Brown, B. B. Danzig, D. Hardy, U. Manber, M. F. Schwartz - "*The Harvest information discovery and access system*" - Proc. 2nd International World Wide Web Conference, 1994
- [42] K. Hammond, R. Burke, C. Martin, S. Lytinen - "*Faq-finder: A case-based approach to knowledge navigation*" - In Working Notes of the AAAI Spring Symposium: Information Gathering from Heterogeneous, Distributed Environments. AAAI Press, 1995.
- [43] T. Kirk, A. Y. Levy, Y. Sagiv, D. Sristava - "*The information manifold*" - In Working Notes of the AAAI Spring Symposium: Information Gathering from Heterogeneous, Distributed Environments. AAAI Press, 1995.
- [44] R.B.Doorenbos, O.Etzioni, D.S.Weld - "*A scalable comparison shopping agent for world wide web*" - Technical Report 96-01-03, University of Washington, Dept. Of Computer Science and Engineering, 1996
- [45] R. Weiss, B. Velez, M. A. Sheldon, C. Namprempre, P. Szlagyi, A. Duda, D. K. Gifford - "*Hypursuit: a hierarchical network search engine that exploits content-link hypertext clustering*" - In Hypertext'96: The Seventh ACM Conference on Hypertext, 1996
- [46] R. Armstrong, d. Freitag, T. Joachims, T. Mitchell - "*Webwatcher: A learning apprentice for the world wide web*" - In Proc. AAAI Spring Symposium on Information Gathering from Heterogeneous, Distributed Environments, 1995.

- [47] Konopnicki, O. Shmueli - "*W3qs: A query system for the world wide web*" - In Proc. Of the 21th VLDB Conference, pagine 54-65, Zurigo, 1995
- [48] B. Mobobasher, N. Jain, E. Han, J. Srivastava - "*Web mining: Pattern Discovery from world wide web transactions*" - Technical Report TR 96-050, University of Minnesota, Dept. of Computer Science, Minneapolis, 1996
- [49] G. M. Youngblood "*Web Hunting: Designe of a Simple Intelligent Web Search Agent*" - 1999 - <http://www.acm.org/crossroads/xrds5-4/webhunting.html>
- [50] Russell, Stuart, P. Norvig - "*Artificial Intelligence A Modern Approach*" - Prentice Hall, Upper Saddle River, N.J., 1995
<http://www.amazon.com/exec/obidos/ISBN=0131038052/acmcrossroadsstu>
- [51] Koster, Martijn - "*Guidelines for Robot Writers*"
<http://info.webcrawler.com/mak/projects/robots/guidelines.html>,
30 Nov. 1998.
- [52] Berners-Lee, R. Fielding, H. Frystyk - "*Hypertext Transfer Protocol - HTTP/1.0*" - RFC:1945, May 1996
<http://info.internet.isi.edu:80/in-notes/rfc/files/rfc1945.txt>
- [53] Koster, Martijn - "*The Web Robots Database*" - 30 Nov. 1998 -
<http://info.webcrawler.com/mak/projects/robots/active.html>
- [54] V. S. Subrahmanian, T. Eiter, T. Lukasiewicz, J. J. Lu - "*Probabilistic Object Bases*" - Infsys Research Report 1843-99-11, November 1999 & April 2001
- [55] T. Gruber - "*What is an ontology?*" -
<http://ksl-web.stanford.edu/kst/what-is-an-ontology.html>
- [56] *Google web site* - www.google.com

- [57] *Sintassi del motore di ricerca Google* -
<http://www.google.it/intl/it/help/basics.html>
- [58] *Altavista web site* - www.altavista.com
- [59] *Lycos web site* - www.lycos.it
- [60] P. Merialdo, P. Atzeni, G. Mecca - "*Semistructured and structured data in the web: Going back and forth*" - in Proceedings of the Workshop on the Management of Semistructured Data, 1997
- [61] C. Latiri, S. Ben Yahia - *Generating implicit association rules from Textual Data*
- IEEE
- [62] *Dizionario della lingua italiana* - Centri Garzanti, 1980