

# Assegnamento Automatico Di Macrocategorie Agli Articoli Di Wikipedia

## Tesi di Laurea Triennale

Tesi di: Jacopo Farina, matricola 713091

Relatore: Prof. Marco Colombetti

Correlatori: Ing. Riccardo Tasso, Ing. David Laniado

Sessione di laurea del 21 settembre 2010

- Enciclopedia consultabile via Web il cui contenuto è **liberamente modificabile** da chiunque
- Lanciata il 15 gennaio 2001, ha avuto una crescita esplosiva
- Disponibile in oltre 270 lingue, di cui un centinaio attive
- La versione in inglese conta più di **3.4 milioni** di voci, ed è la più grande in assoluto seguita da quella in tedesco con 1.2 milioni
- La versione in italiano ne ha “solo” 720.000

# Le categorie di Wikipedia

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegna-  
mento

Risultati

Nel maggio 2004 sono state introdotte le categorie per organizzare meglio la quantità sempre maggiore di articoli esistenti.

- Attualmente ne esistono circa 500mila
- Ogni pagina può appartenere a una o più categorie, scelte arbitrariamente dagli utenti
  - In media una pagina appartiene a **2,68** categorie
  - Il 64% degli articoli appartiene a una o due categorie, il 93% a meno di 7
  - La pagina *Winston Churchill* ha il numero record di 70 categorie di appartenenza, *Albert Einstein* ne ha 56.
- Ogni categoria appartiene ad altre categorie e ne può contenere altre, senza vincoli particolari
- **non rappresentano una tassonomia**, anche se di solito sono organizzate in maniera gerarchica, semmai dei legami semantici

# Le categorie di Wikipedia

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegna-  
mento

Risultati

Nel maggio 2004 sono state introdotte le categorie per organizzare meglio la quantità sempre maggiore di articoli esistenti.

- Attualmente ne esistono circa 500mila
- Ogni pagina può appartenere a una o più categorie, scelte arbitrariamente dagli utenti
  - In media una pagina appartiene a **2,68** categorie
  - Il 64% degli articoli appartiene a una o due categorie, il 93% a meno di 7
  - La pagina *Winston Churchill* ha il numero record di 70 categorie di appartenenza, *Albert Einstein* ne ha 56.
- Ogni categoria appartiene ad altre categorie e ne può contenere altre, senza vincoli particolari
- **non rappresentano una tassonomia**, anche se di solito sono organizzate in maniera gerarchica, semmai dei legami semantici

# Le categorie di Wikipedia

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegna-  
mento

Risultati

Nel maggio 2004 sono state introdotte le categorie per organizzare meglio la quantità sempre maggiore di articoli esistenti.

- Attualmente ne esistono circa 500mila
- Ogni pagina può appartenere a una o più categorie, scelte arbitrariamente dagli utenti
  - In media una pagina appartiene a **2,68** categorie
  - Il 64% degli articoli appartiene a una o due categorie, il 93% a meno di 7
  - La pagina *Winston Churchill* ha il numero record di 70 categorie di appartenenza, *Albert Einstein* ne ha 56.
- Ogni categoria appartiene ad altre categorie e ne può contenere altre, senza vincoli particolari
- **non rappresentano una tassonomia**, anche se di solito sono organizzate in maniera gerarchica, semmai dei legami semantici

# Una pagina di Wikipedia

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegnamento

Risultati

## Milan

From Wikipedia, the free encyclopedia

Coordinates: 45°27′51″N 09°﻿ / ﻿

*"Milano" redirects here. For other uses, see Milano (disambiguation).*

*For other uses, see Milan (disambiguation).*

**Milan** (Italian: Milano, About this sound listen<sup>(help·info)</sup> Italian pronunciation: [miˈlaːno]; Western Lombard: Milan, About this sound listen<sup>(help·info)</sup>) is a city in Italy<sup>(help·info)</sup> and the capital of the region of Lombardy<sup>(help·info)</sup> and of the province of Milan<sup>(help·info)</sup>. The city proper has a population of about 1,310,000, while the urban area is the first in Italy and the fifth largest in the European Union<sup>(help·info)</sup> with a population of 4,345,000 over an area of 2,370 km2 (915 sq mi).[2] The Milan metropolitan area, by far the largest in Italy, is estimated by the OECD<sup>(help·info)</sup> to have a population of 7,400,000.[3]

- The Milan Garden of the Righteous
- Milan In Hebrew מילאנו

<span><span><span></span></span></span> Regional capitals of Italy	
<span><span></span></span> Lombardy <span> </span> • <i>Comuni</i> of the Province of Milan	

Categories:Cities and town in Lombardy|Communes of the province of Milan|Populated places established in 1st millennium



# Una categoria di Wikipedia

Assegnamento automatico delle macro-categorie

Tesi di:  
Jacopo Farina,  
matricola 713091

Introduzione


Il problema dell'assegnamento

Risultati

Category: [Discussion](#) [read](#) [edit](#) [view history](#) [Get help](#)

## Category:Populated places established in the 1st millennium BC

From Wikipedia, the free encyclopedia

 Wikimedia Commons has media related to: *Settlements established in the 1st millennium BC*

### Subcategories

This category has the following 10 subcategories, out of 10 total.

<ul style="list-style-type: none"><li><span>[x]</span> <a href="#">Populated places established in the 1st century BC (13 P)</a></li><li><span>[+]</span> <a href="#">Populated places established in the 2nd century BC (2 C, 13 P)</a></li><li><span>[x]</span> <a href="#">Populated places established in the 3rd century BC (18 P)</a></li><li><span>[x]</span> <a href="#">Populated places established in the 4th century BC (33 P)</a></li></ul>	<b>cont.</b> <ul style="list-style-type: none"><li><span>[x]</span> <a href="#">Populated places established in the 5th century BC (14 P)</a></li><li><span>[x]</span> <a href="#">Populated places established in the 6th century BC (16 P)</a></li><li><span>[x]</span> <a href="#">Populated places established in the 7th century BC (22 P)</a></li><li><span>[x]</span> <a href="#">Populated places established in the 8th century BC (16 P)</a></li></ul>	<b>cont.</b> <ul style="list-style-type: none"><li><span>[x]</span> <a href="#">Populated places established in the 9th century BC (6 P)</a></li><li><span>[x]</span> <a href="#">Populated places established in the 10th century BC (1 P)</a></li></ul>
--	--	---

### Pages in category "Populated places established in the 1st millennium BC"

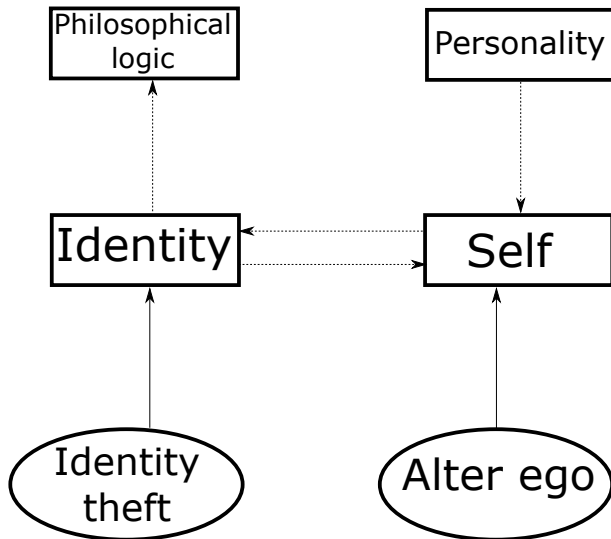
The following 28 pages are in this category, out of 28 total. This list may not reflect recent changes ([learn more](#)).

<b>G</b> <ul style="list-style-type: none"><li><a href="#">Geneva</a></li></ul>	<b>S</b> <ul style="list-style-type: none"><li><a href="#">Seleucia Pieria</a></li><li><a href="#">Stena</a></li><li><a href="#">Soli, Cilicia</a></li></ul>	<b>Y</b> <ul style="list-style-type: none"><li><a href="#">Yangzhou</a></li><li><a href="#">Yevpatoria</a></li></ul>
<b>H</b> <ul style="list-style-type: none"><li><a href="#">Herculaneum</a></li></ul>		
<b>K</b> <ul style="list-style-type: none"><li><a href="#">Ksar-el-Kebir</a></li></ul>		
<b>M</b> <ul style="list-style-type: none"><li><a href="#">Marseille</a></li></ul>		

[Categories: Populated places by year of establishment | 1st-millennium BC establishments](#)

Navigation icons: back, forward, search, etc.

# Un esempio reale



Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

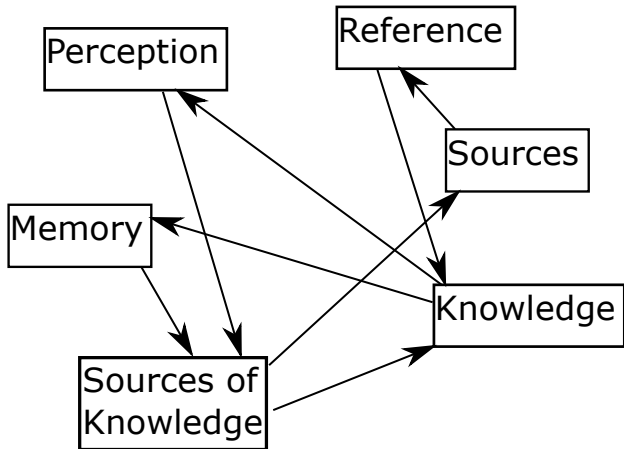
Introduzione

Il problema  
dell'assegna-  
mento

Risultati



# Un ciclo interessante



- Gli archi non rappresentano sempre una relazione di tipo *isa*
- I cicli si trovano con l'algoritmo di Tarjan per le strutture fortemente connesse

Wikipedia contiene quantità **enormi** di testo **liberamente accessibile** su qualsiasi argomento, ed è quindi di grande interesse nel campo dell'Intelligenza Artificiale:

- Mappatura tra termini di ontologie (come **Cyc**) e articoli di Wikipedia per estrarre conoscenza
- Analisi automatica del testo per estrarre legami tra parole (sinonimi, iperonimi, iponimi...)
- Ricerca di legami semantici tra i termini
- Analisi statistica degli argomenti trattati dal progetto e dai singoli utenti

Wikipedia contiene quantità **enormi** di testo **liberamente accessibile** su qualsiasi argomento, ed è quindi di grande interesse nel campo dell'Intelligenza Artificiale:

- Mappatura tra termini di ontologie (come **Cyc**) e articoli di Wikipedia per estrarre conoscenza
- Analisi automatica del testo per estrarre legami tra parole (sinonimi, iperonimi, iponimi...)
- Ricerca di legami semantici tra i termini
- Analisi statistica degli argomenti trattati dal progetto e dai singoli utenti

# Utilizzi nell'IA

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione e

Il problema  
dell'assegna-  
mento

Risultati

Wikipedia contiene quantità **enormi** di testo **liberamente accessibile** su qualsiasi argomento, ed è quindi di grande interesse nel campo dell'Intelligenza Artificiale:

- Mappatura tra termini di ontologie (come **Cyc**) e articoli di Wikipedia per estrarre conoscenza
- Analisi automatica del testo per estrarre legami tra parole (sinonimi, iperonimi, iponimi...)
- Ricerca di legami semantici tra i termini
- Analisi statistica degli argomenti trattati dal progetto e dai singoli utenti

Wikipedia contiene quantità **enormi** di testo **liberamente accessibile** su qualsiasi argomento, ed è quindi di grande interesse nel campo dell'Intelligenza Artificiale:

- Mappatura tra termini di ontologie (come **Cyc**) e articoli di Wikipedia per estrarre conoscenza
- Analisi automatica del testo per estrarre legami tra parole (sinonimi, iperonimi, iponimi...)
- Ricerca di legami semantici tra i termini
- Analisi statistica degli argomenti trattati dal progetto e dai singoli utenti

Wikipedia contiene quantità **enormi** di testo **liberamente accessibile** su qualsiasi argomento, ed è quindi di grande interesse nel campo dell'Intelligenza Artificiale:

- Mappatura tra termini di ontologie (come **Cyc**) e articoli di Wikipedia per estrarre conoscenza
- Analisi automatica del testo per estrarre legami tra parole (sinonimi, iperonimi, iponimi...)
- Ricerca di legami semantici tra i termini
- Analisi statistica degli argomenti trattati dal progetto e dai singoli utenti

# Il problema dell'assegnamento

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione e

Il problema  
dell'assegnamento

Risultati

- Dato un set di argomenti, chiamati **macrocategorie**, che partiziona idealmente la conoscenza umana, come stabilire automaticamente quale (o quali) di essi è più adatto a contenere un certo articolo ?
- Kittur (2008) ha analizzato il grafo delle categorie di Wikipedia come se fosse una **rete semantica**, abbinando ogni articolo alla macrocategoria con la minore distanza topologica dall'articolo tra quelle di un set di 11 elementi.
- Il criterio di Kittur è ancora valido dopo l'aumento delle dimensioni del grafo? E se si scelgono più macrocategorie?
- Esistono criteri più efficaci, ossia che assegnino automaticamente le macrocategorie agli articoli in maniera più simile a quanto farebbe un essere umano?

# Il problema dell'assegnamento

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione e

Il problema  
dell'assegnamento

Risultati

- Dato un set di argomenti, chiamati **macrocategorie**, che partiziona idealmente la conoscenza umana, come stabilire automaticamente quale (o quali) di essi è più adatto a contenere un certo articolo ?
- Kittur (2008) ha analizzato il grafo delle categorie di Wikipedia come se fosse una **rete semantica**, abbinando ogni articolo alla macrocategoria con la minore distanza topologica dall'articolo tra quelle di un set di 11 elementi.
- Il criterio di Kittur è ancora valido dopo l'aumento delle dimensioni del grafo? E se si scelgono più macrocategorie?
- Esistono criteri più efficaci, ossia che assegnino automaticamente le macrocategorie agli articoli in maniera più simile a quanto farebbe un essere umano?



# Il problema dell'assegnamento

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione e

Il problema  
dell'assegnamento

Risultati

- Dato un set di argomenti, chiamati **macrocategorie**, che partiziona idealmente la conoscenza umana, come stabilire automaticamente quale (o quali) di essi è più adatto a contenere un certo articolo ?
- Kittur (2008) ha analizzato il grafo delle categorie di Wikipedia come se fosse una **rete semantica**, abbinando ogni articolo alla macrocategoria con la minore distanza topologica dall'articolo tra quelle di un set di 11 elementi.
- Il criterio di Kittur è ancora valido dopo l'aumento delle dimensioni del grafo? E se si scelgono più macrocategorie?
- Esistono criteri più efficaci, ossia che assegnino automaticamente le macrocategorie agli articoli in maniera più simile a quanto farebbe un essere umano?

# Il problema dell'assegnamento

Assegnamento  
automatico  
delle macro-  
categorie

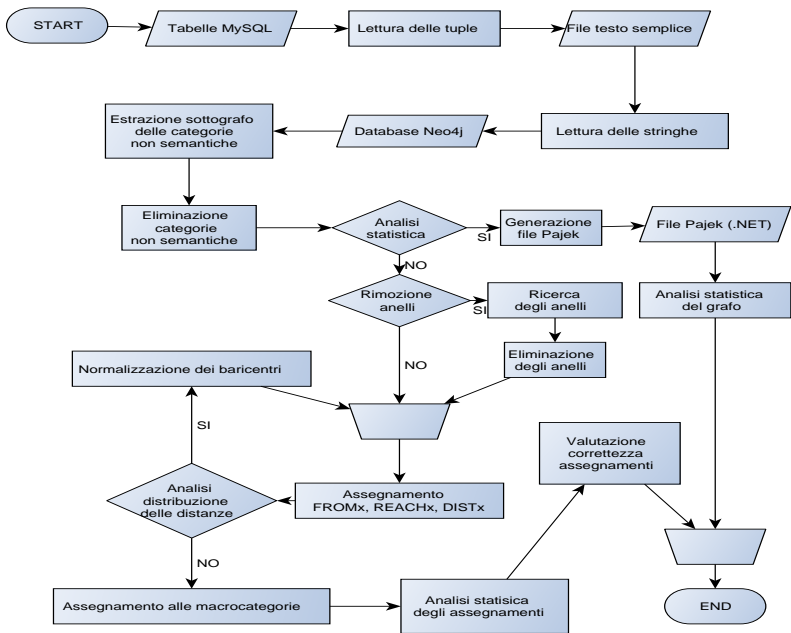
Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegnamento

Risultati

- Dato un set di argomenti, chiamati **macrocategorie**, che partiziona idealmente la conoscenza umana, come stabilire automaticamente quale (o quali) di essi è più adatto a contenere un certo articolo ?
- Kittur (2008) ha analizzato il grafo delle categorie di Wikipedia come se fosse una **rete semantica**, abbinando ogni articolo alla macrocategoria con la minore distanza topologica dall'articolo tra quelle di un set di 11 elementi.
- Il criterio di Kittur è ancora valido dopo l'aumento delle dimensioni del grafo? E se si scelgono più macrocategorie?
- Esistono criteri più efficaci, ossia che assegnino automaticamente le macrocategorie agli articoli in maniera più simile a quanto farebbe un essere umano?



# Formato dei risultati e valutazione

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegna-  
mento

Risultati

- Per ogni pagina ottengo delle percentuali di assegnamento alle macrocategorie

Politecnico\_di\_Milano: Education:87,5; History and events:12,5;

- Valuto la precisione dei risultati confrontando le quote assegnate con quelle stabilite manualmente da un umano
- Valutazione di 50 pagine scelte a caso tra quelle assegnate
- Viene usato come valore di similarità il coseno dei vettori degli assegnamenti (automatico ed umano)

$$\cos(A, B) = \frac{A \bullet B}{\|A\| * \|B\|}$$

# Dimensione delle macrocategorie

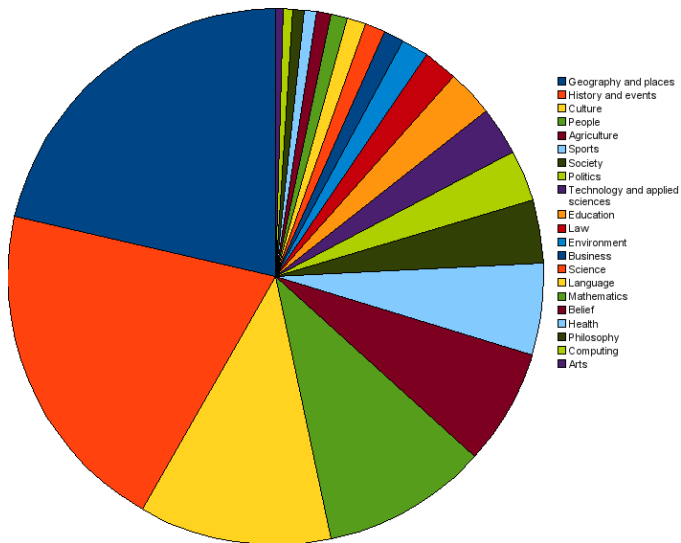
Assegnamento automatico delle macrocategorie

Tesi di:  
Jacopo Farina,  
matricola  
713091

Introduzione

Il problema dell'assegnamento

Risultati



# Alcuni assegnamenti

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

introduzione

Il problema  
dell'assegnamento

Risultati

Milan: History and events:100;  
Italy: History and events:70; Culture:30;  
Politecnico\_di\_Milano: Education:87,5; History and events:12,5;  
Java(software\_platform): Computing:50; Technology and applied sciences:50;  
Java: Geography and places:100; (L'isola principale dell'Indonesia)  
Earth: Geography and places:100;  
Albert\_Einstein: Technology and applied sciences:18,53; Science:18,53; People:24,78; Mathematics:18,53; Politics:11,39; History and events:5,31; Society:2,93;  
Nintendo: History and events:71,67; Agriculture:6,67; Culture:6,67; Geography and places:6,67; Philosophy:6,67; Technology and applied sciences:1,67; (famosa azienda produttrice di videogiochi)

- Perché Nintendo viene assegnata ad Agriculture più che a Technology and applied sciences ?

# I percorsi da Nintendo

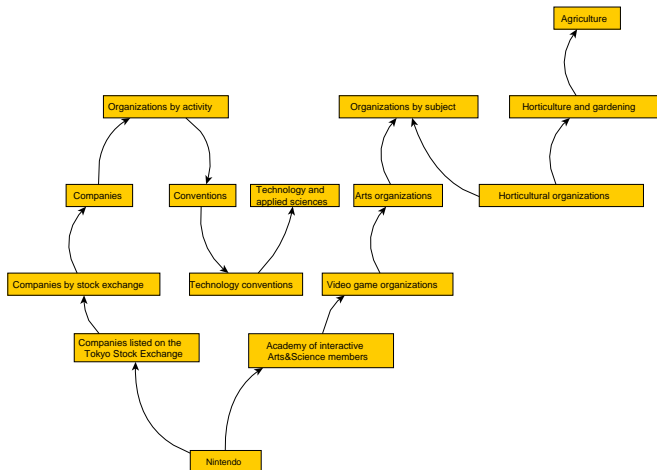
Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

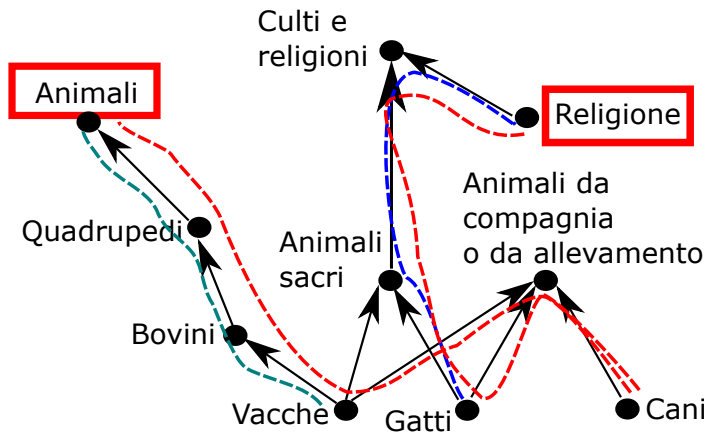
Introduzione

Il problema  
dell'assegnamento

Risultati



# Percorsi diretti e indiretti



- Percorsi orientati: non sempre arrivo a una macrocategoria
- Percorsi non orientati: posso giungere a degli errori



# Distribuzione con i costi differenziati

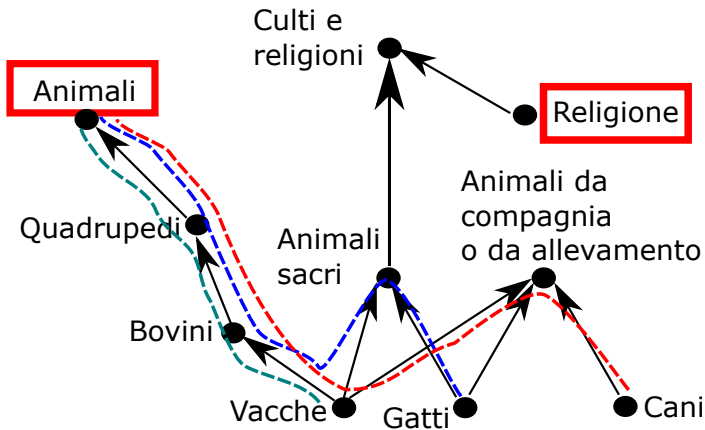
Assegnamento automatico delle macro-categorie

Tesi di:  
Jacopo Farina,  
matricola  
713091

Introduzione

Il problema dell'assegnamento

Risultati



- Ibrido: mi muovo in entrambi i modi preferendo i percorsi orientati

# Considerare la connettività

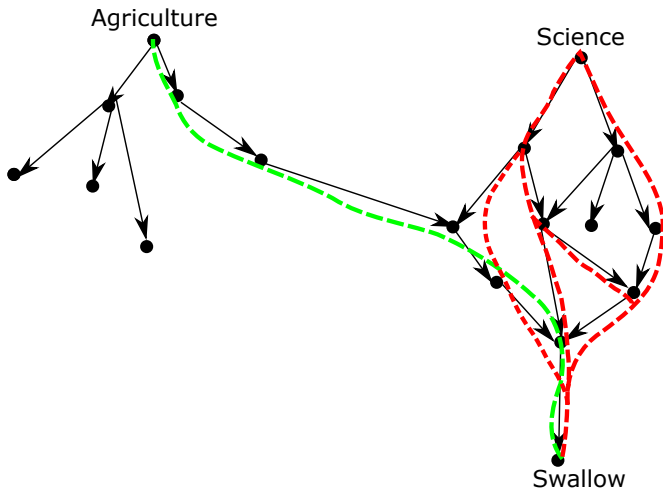
Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

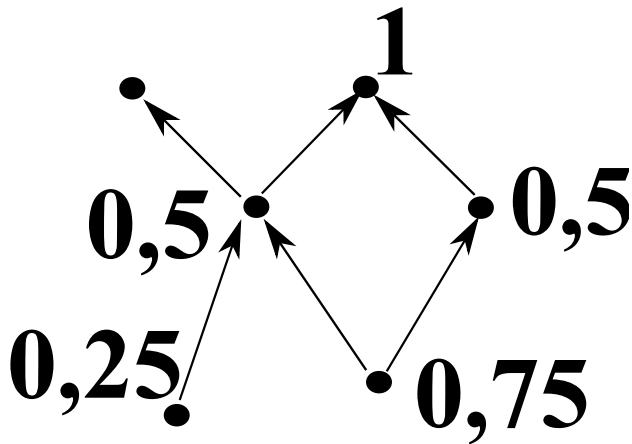
Il problema  
dell'assegnamento

Risultati



Potrebbe essere utile considerare più percorsi possibili alla volta.

# Ripartizione di punteggi



- Ripartisco il punteggio di una categoria tra quelle contenute, eventualmente sommandolo a quelli già assegnati.

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegnamento

Risultati

# Probabilità di raggiungimento

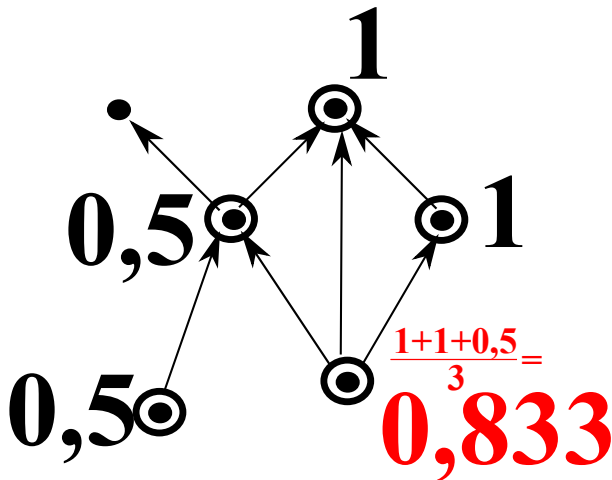
Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegnamento

Risultati



Marco i nodi raggiungibili, ogni categoria ha una probabilità uguale alla somma delle probabilità di quelle che la contengono divisa per gli archi uscenti.

# Dimensione delle macrocategorie con assegnamento probabilistico

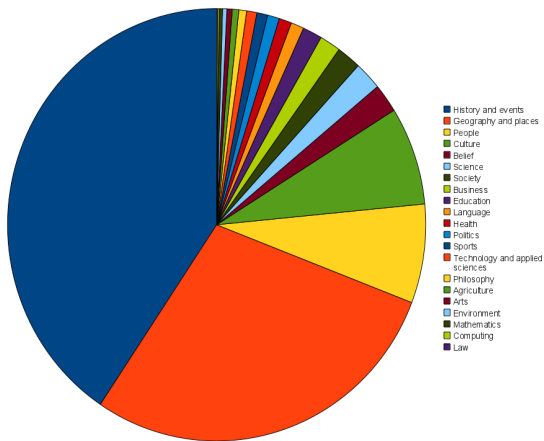
Assegnamento automatico delle macrocategorie

Tesi di:  
Jacopo Farina,  
matricola  
713091

Introduzione

Il problema dell'assegnamento

Risultati



- Si nota che Agriculture si è molto ridimensionata, assumendo una quota più plausibile.

# I metodi esaminati

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegnamento

Risultati

Metodo	valore di correttezza	% pagine valutate
Caso base (Kittur, Holloway)	0.34	100%
Solo percorsi diretti	0.35	65%
Costi differenziati	<b>0.37</b>	100%
Ripartizione dei punteggi	0.35	62%
Probabilità di raggiungimento	0.36	62%

# Conclusioni

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione e

Il problema  
dell'assegnamento

Risultati

- Il metodo più efficace tra quelli esaminati è quello a costi di attraversamento differenziati.
- Anche gli algoritmi basati sui percorsi multipli funzionano meglio del criterio basato sulla semplice distanza topologica
- Il metodo di Kittur fornisce assegnamenti correlati a quelli forniti manualmente da un valutatore umano anche se aumenta la complessità del grafo e si scelgono più macrocategorie.

# Sviluppi futuri

Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegnamento

Risultati

- Utilizzare non solo le appartenenze alle categorie ma anche i collegamenti tra le pagine per generare il grafo
- Elaborare tecniche del calcolo dei percorsi multipli che classifichino il 100% delle pagine seguendo gli archi in tutte le direzioni
  - Molte delle tecniche del calcolo del flusso esistenti non sono applicabili su grafi così grossi
  - Si potrebbero utilizzare approcci stocastici per calcolare approssimativamente la raggiungibilità di un nodo
- Utilizzare gli assegnamenti ottenuti per applicazioni pratiche
  - Mappare gli articoli con documenti testuali per calcolarne la vicinanza semantica
  - Stabilire la macrocategoria di un documento analizzandone le singole parole



# Ripartizione di punteggi

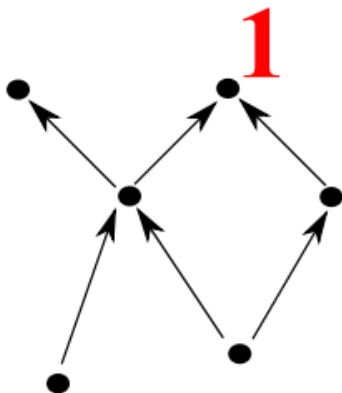
Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

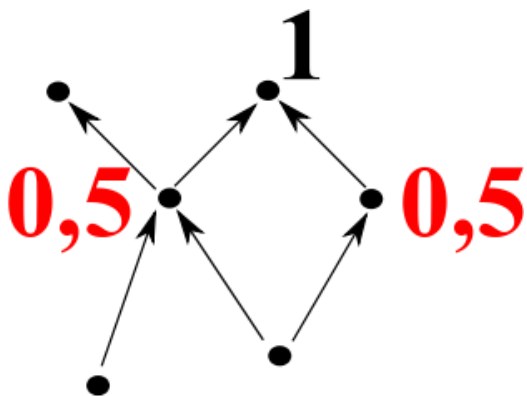
Introduzione

Il problema  
dell'assegna-  
mento

Risultati



# Ripartizione di punteggi



Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione e

Il problema  
dell'assegnamento

Risultati

# Ripartizione di punteggi

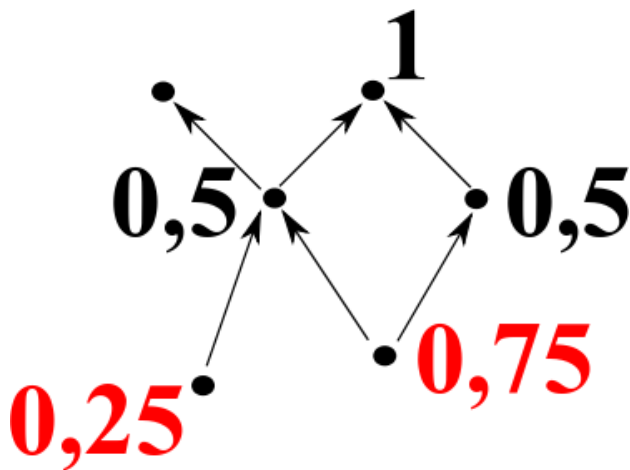
Assegnamento automatico delle macro-categorie

Tesi di:  
Jacopo Farina,  
matricola  
713091

Introduzione

Il problema dell'assegnamento

Risultati



# Probabilità di raggiungimento

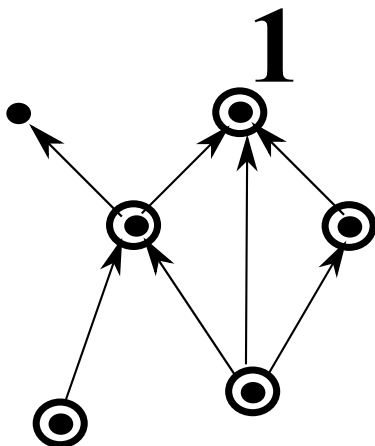
Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

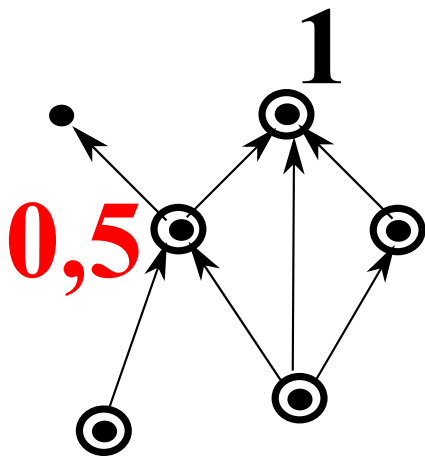
Il problema  
dell'assegna-  
mento

Risultati



I nodi raggiungibili da una macrocategoria sono stati marcati, eventualmente da una tecnica precedente.

# Probabilità di raggiungimento



Si attraversa il grafo secondo la logica depth-first

# Probabilità di raggiungimento

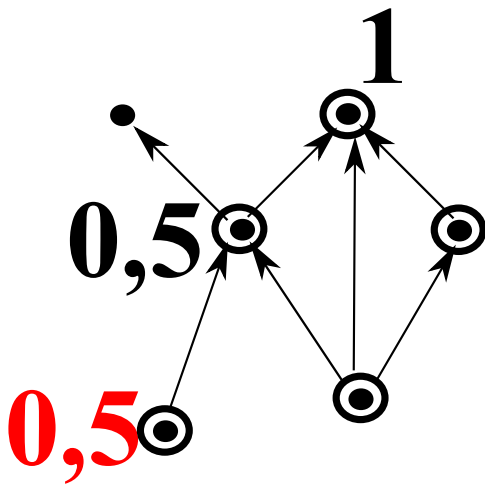
Assegnamento automatico delle macro-categorie

Tesi di:  
Jacopo Farina,  
matricola  
713091

Introduzione

Il problema dell'assegnamento

Risultati



# Probabilità di raggiungimento

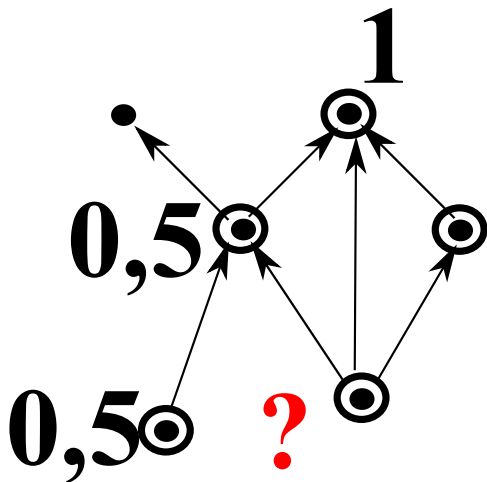
Assegnamento  
automatico  
delle macro-  
categorie

Tesi di:  
Jacopo  
Farina,  
matricola  
713091

Introduzione

Il problema  
dell'assegnamento

Risultati



Se un nodo ha archi uscenti verso nodi marcati ma senza una probabilità assegnata, lo salto: ci passerò dopo

# Probabilità di raggiungimento

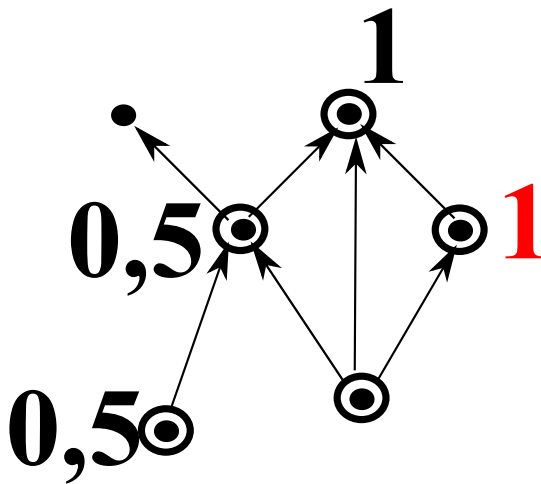
Assegnamento automatico delle macro-categorie

Tesi di:  
Jacopo Farina,  
matricola  
713091

Introduzione

Il problema dell'assegnamento

Risultati





# Probabilità di raggiungimento

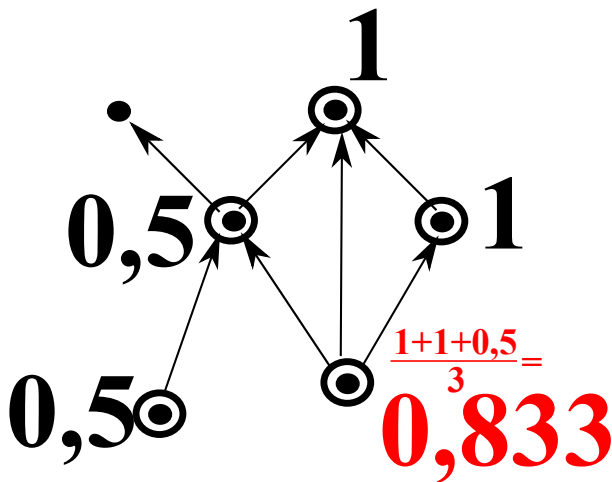
Assegnamento automatico delle macro-categorie

Tesi di:  
Jacopo Farina,  
matricola  
713091

Introduzione

Il problema dell'assegnamento

Risultati



Finalmente tutti i nodi marcati sono stati elaborati, posso assestare la probabilità di raggiungimento.

- Per maggiori dettagli si veda la pagina dell'AirWiki
  - <http://airwiki.ws.dei.polimi.it/index.php/WikipediaCategoryGr>