

CORSO DI *WEB MINING E RETRIEVAL* *- INTRODUZIONE AL WM -*

Corso di Laurea in Informatica, Ing. Internet,
Ing. Informatica, Ing. Gestionale
(a.a. 2015-2016)

Roberto Basili

Overview

- Web Mining & Retrieval: Motivazioni e prospettive
 - Web, User-generated contents, Social Media
 - The role of learning
 - What is Machine Learning?
 - Data-driven algorithms: sources of complexity
- Main Applications
 - Intelligent Web Search
 - User Profiling for Marketing or Brand reputation management
 - Web Recommending
 - Spoken Dialogue Interaction in Robotics or in Web/mobile Interfaces

Do you know

More than
4,000 new books
are published every day



A Web of people and opinions

- **31.7%** of the more than 200 million bloggers worldwide blog about opinions on products and brands (Universal McCann, July 2009)
- **71%** of all active Internet users read blogs.
- 2009 Survey of **25,000** Internet users in **50** countries: **70%** of consumers trust opinions posted online by other consumers (Nielsen Global Online Consumer, 2010).

Social Media Analytics

- Complex process for Social Media Analytics are necessary whereas ...
 - Communities play a special role in circulating information
 - Search is the mostly collective function used (e.g. # in tweets)
 - Opinion Mining and Sentiment Analysis are important for individuals as well as organisations



WM&R: Motivazioni

- *Cos'è il Web Mining?*
- *Perché IR?*
- *Perché Apprendimento Automatico?*
- *Quale contributo l'IR fornisce alle tecnologie di sfruttamento delle informazioni del Web?*
- *Quali sono le prospettive per l'impiego di tali tecnologie?*

Cos'è il Web Mining?

- Web Mining attualmente si riferisce ad un insieme di tecnologie necessarie per lo sfruttamento delle informazioni pubblicamente disponibili nel Web
 - Contenuti: dati ma anche ... persone, luoghi, eventi, concetti, ...
 - Relazioni:
 - Link strutturali
 - Collegamenti tematici, concettuali e interpersonali
 - Ridondanze/analogie
 - Multilingualità
 - Trend e comportamenti collettivi
 - Opinioni

Perché IR?

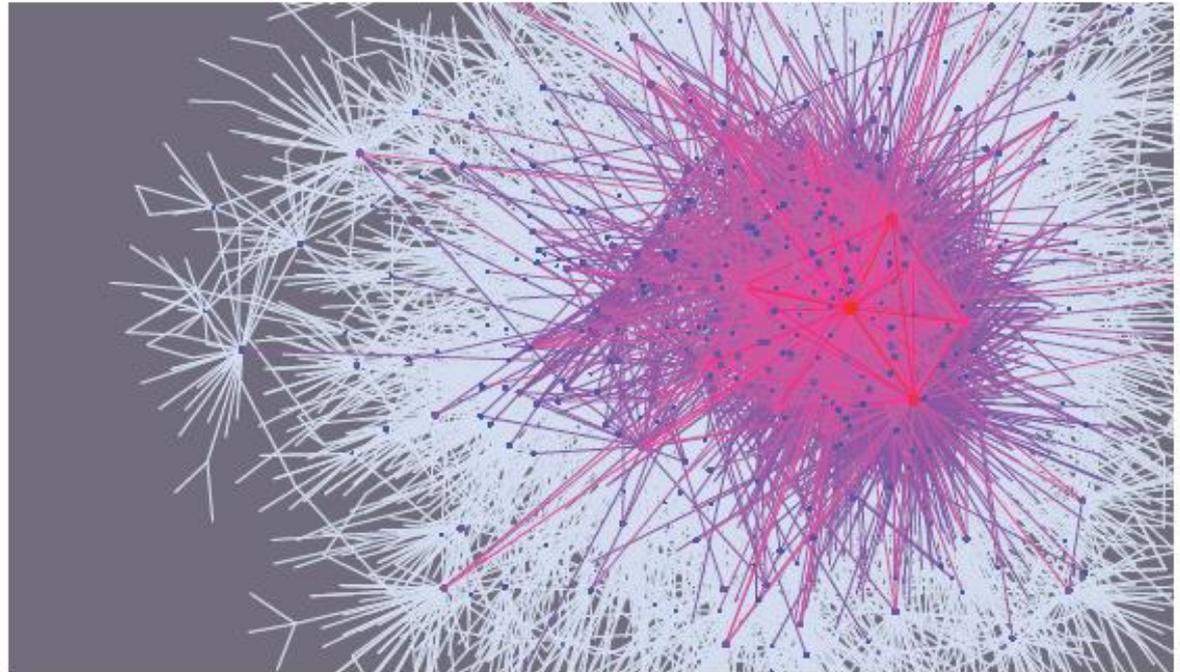
- La taglia delle informazioni in gioco pone il problema della *localizzazione*
- Accedere in modo automatico è possibile solo governando il problema di sapere **dove** si trova una informazione *rilevante*
- La ricerca corrisponde al calcolo di una funzione *aleatoria* di mapping tra requisiti e informazione utile

Machine Learning vs IR?

- La eterogeneità delle informazioni produce significativi effetti di incertezza nel processo di ricerca riguardo ad aspetti diversi del processo di IR
 - Incompletezza della informazione:
 - Query brevi come informazione (incomplete) sui fabbisogni informativi
 - Ricchezza di dati, formati e modalità di accesso
 - I contenuti sono sparsi in diverse forme nei dati
 - Requisiti vaghi
 - Spesso molte informazioni sono esplicite solo nel contesto
 - Aspetti soggettivi
 - La rilevanza dipende dallo user e non solo dal contenuto
 - Tempestività ed autorevolezza

ML vs. IR

- La pervasività degli elementi di incertezza rende impraticabile la ricerca di soluzioni esaustive (ottimi globali)
- “*Finding diamonds in the rough*”
(Fan Chung, UCSD)



ML vs. IR

- Le tecniche di ML propongono una ampia serie di algoritmi, strategie e tecniche per la produzione di soluzioni *sub-ottime* effettive
- Nel processo di *learning* i dati suggeriscono la ipotesi risolutiva per la funzione di *mapping*
- Tale ipotesi è attesa migliorare la prestazione complessiva del sistema di base
 - Accuratezza
 - Efficienza computazionale

Machine Learning

- (Langley, 2000): l'Apprendimento Automatico si occupa dei meccanismi attraverso i quali un agente intelligente migliora nel tempo le sue prestazioni P nell'effettuare un compito C .
- La prova del successo dell'apprendimento è quindi nella capacità di misurare l'incremento ΔP delle prestazioni sulla base delle esperienze E che l'agente è in grado di raccogliere durante il suo ciclo di vita.
- La natura dell'apprendimento è quindi tutta nella caratterizzazione delle nozioni qui primitive di *compito*, *prestazione* ed *esperienza*.

Esperienza ed Apprendimento

- L'esperienza, per esempio, nel gioco degli scacchi può essere interpretata in diversi modi:
 - i dati sulle vittorie (e sconfitte) pregresse per valutare la bontà (o la inadeguatezza) di strategie e mosse eseguite rispetto all'avversario.
 - valutazione fornita sulle mosse da un docente esterno (oracolo, guida).
 - Adeguatezza dei comportamenti derivata dalla auto-osservazione, cioè dalla capacità di analizzare partite dell'agente contro se stesso secondo un modello esplicito del processo (partita) e della sua evoluzione (comportamento, vantaggi, ...).

ML: una introduzione visuale

- See URL: http://www.r2d3.us/visual-intro-to-machine-learning-part-1/?imm_mid=0d76b4&cmp=em-data-na-na-newsltr_20150826

Algoritmi di Apprendimento

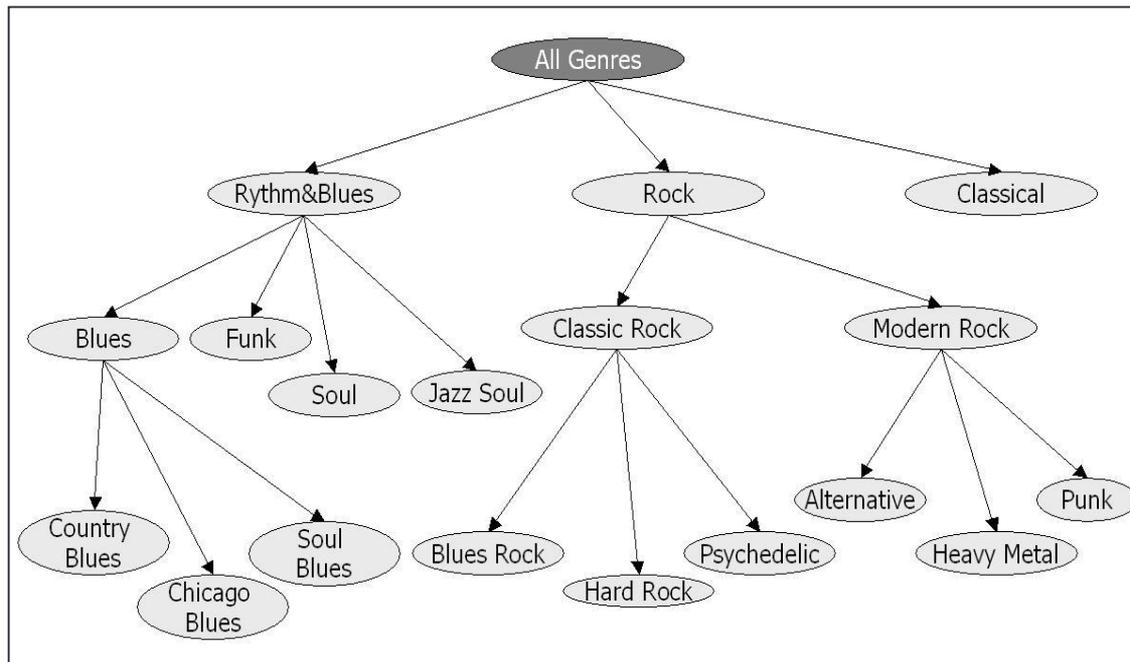
- Funzioni logiche booleane, (ad es., alberi di decisione).
- Funzione di Probabilità, (ad es., classificatore Bayesiano).
- Funzioni di separazione in spazi vettoriali
 - Non lineari: KNN, reti neurali multi-strato,...
 - Lineari, percettroni, Support Vector Machines,...
- Trasformazioni di spazi: embeddings, analisi spettrale

Apprendimento senza supervisione

- In assenza di un oracolo o di conoscenze sul task esistono ancora molti modi di migliorare le proprie prestazioni, ad es.
 - Migliorando il proprio modello del mondo (acquisizione/*discovery* della conoscenza)
 - Migliorando le proprie prestazioni computazionali (ottimizzazione)

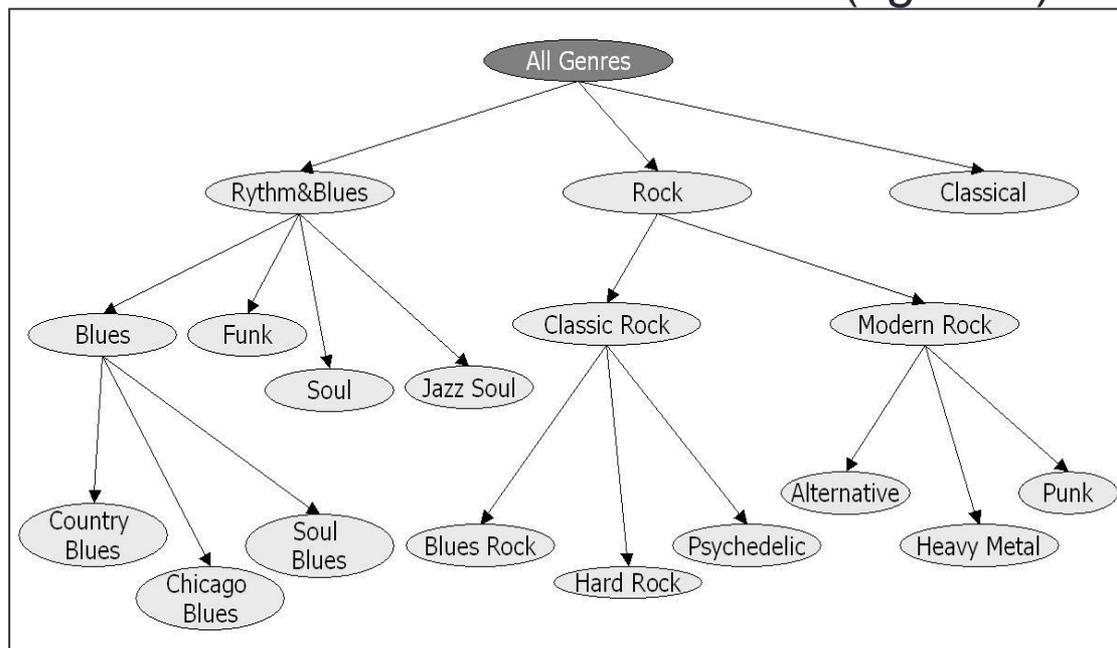
Apprendimento senza supervisione

- Esempio:
 - una collezione mp3 può essere organizzata in generi attraverso il raggruppamento di brani simili secondo proprietà audio (*clustering*): tale organizzazione è naturalmente gerarchica
 - Il miglioramento avviene quindi almeno rispetto agli algoritmi di ricerca: la organizzazione gerarchica consente di esaminare solo i membri dell'insieme in alcune classi (i generi).



Apprendimento senza supervisione

- Esempio: Al termine del processo di acquisizione il sistema dispone di un sistema di classi e relazioni indotti che migliora la sua interazione futura con l'ambiente operativo (ad es. l'utente)
- Il miglioramento avviene quindi almeno rispetto agli algoritmi di ricerca: la organizzazione gerarchica consente di esaminare solo i membri dell'insieme in alcune classi (i generi).



Web IR?

- I processi di IR studiati in domini antecedenti all'affermarsi del Web debbono essere estesi ed adattati rispetto alla maggiore ricchezza ed ai problemi maggiori che tali scenari presentano
 - Complessità strutturale: contenuti, topologia e uso
 - Affidabilità dell'informazione
 - Multimodalità, Multimedialità
 - Partecipazione (aspetti sociali)

Web IR

- Processing Web data: content detection, link detection, ...
- Web Crawling
- Web Search: indici, link analysis
- Ranking: weighting contents, links and formats, authority, timeliness
- Meta-search
- Link Analysis

Prospettive delle tecnologie WM&R

- Crescita esponenziale della taglia dei problemi
- Crescente interesse verso processi di IR agenti su dati complessi (multimediali, sociali)
- Web partecipativo: Web 2.0
- Ruolo crescente della mediazione degli strumenti informatici
 - Software as a Service
 - Personalizzazione
- Big Data challenges:
 - Scala
 - Opacità Semantica

Information, Web and the language

Web contents, characterized by rich multimedia information, are mostly **opaque from a semantic standpoint**

Today is 2011年11月13日 星期日 顯示器最佳分辨率 1024X768 今日天氣 加入最愛 設為首頁 大公網新版

www.takungpao.com.hk

2011 中国証券金紫荆獎 Golden Bauhinia Awards

首頁 國內 國際 港澳 兩岸 評論 財經 體育 教育 科技 醫學 娛樂 文化 副刊 軍事 生活 旅遊 圖片 博客

關鍵詞: 欄目: 全部 最近三個月 三個月之前 檢索

▶ 手機新聞 ▶ 手機博客 ▶ 漢語學習 ▶ 新聞點擊排行

滾動新聞:

胡總語特首:防範經濟金融風險
胡錦濤在夏威夷會見出席APEC峰會的曾蔭權。他祝賀香港區議會選舉成功,並充分肯定曾蔭權及港府工作,要求做好經濟金融風險防範

胡連會登場 共同宣示九二共識
胡錦濤第四次在APEC峰會期間會見連戰。他強調,認同「九二共識」是兩岸開展對話協商的必要前提,也是兩岸關係和平發展的重要基礎

西藏黨代會高調反「藏獨」 德國作家:外埠雜誌報道西藏
傳媒入日本福島核電站探險 英國大裁軍 傷兵難雜免
滇礦難已30死 13人生還 礦工講述內幕 事故並不意外
范徐麗泰認民望跌最不耐 選委再獲60提名表 累積逾千人
聖保羅中學本月底截止招新 選委再獲60提名表 累積逾千人
民調逆轉 藍高層:國親吵鬧地 秋門訴求多 向藍綠表不滿
世界新七奇觀 亞洲景佔四席 新奇觀選舉惹爭議

中國實體書店苦苦掙扎求存 加入TPP 台密集會談探險
香港人家/蔡仕榮 人生導師 活出自我 我香港人家/教導子女...
債務危機好 港ADR幾全線造藍 歐元反降 兌美元逼近1.38
入世十年/充分對接 華強北最輝煌 入世十年/挑戰「二次...
抽身除「雜軍」 工人險生國難 南亞漢命案 警日籍妻

即時新聞

- 組國/河南全國太極拳錦標賽賽
- 奧巴馬重申美不支持「台灣獨立
- 巴基斯坦西北部兩起襲擊 16人
- 圖文/胡錦濤會見美國總統奧巴
- 兩岸30對愛侶在廈門集體證婚
- 中日韓衛生部長會議在青島舉行
- 面向中國遊客中英雜誌紐約創刊
- 「CEO聖經」成內地官員考試內
- 斯特恩:經紀人是勞資談判的障
- 香港冀成爲人幣國際化關鍵角色
- 日學者提出地核物質形態新假設
- 中國影視機構向國際大師「取經

焦點關注

區議會選舉 香港

2011APEC 港黑金事件 2011

神八天宮對接 第七次陳江會 李

9.1衝擊事件 中國航母試航 辛

Information, Web and language

Hu meets KMT honorary chairman in Hawaii - People's Daily Online - Mozilla Firefox

File Modifica Visualizza Cronologia Segnalibri Yahoo! Strumenti Aiuto

Hu meets KMT honorary chairman

Indietro Avanti Download

Chinese President Hu Jintao (R) shakes hands with Honorary Chairman of the Chinese Kuomintang (KMT) Lien Chan, in Honolulu, Hawaii, the U.S., Nov. 11, 2011.
(Xinhua/Huang Jingwen)

HONOLULU, United States, Nov. 11 (Xinhua) -- Hu Jintao, general secretary of the Central

Latest News: • Indonesia to host European Higher Education Fair

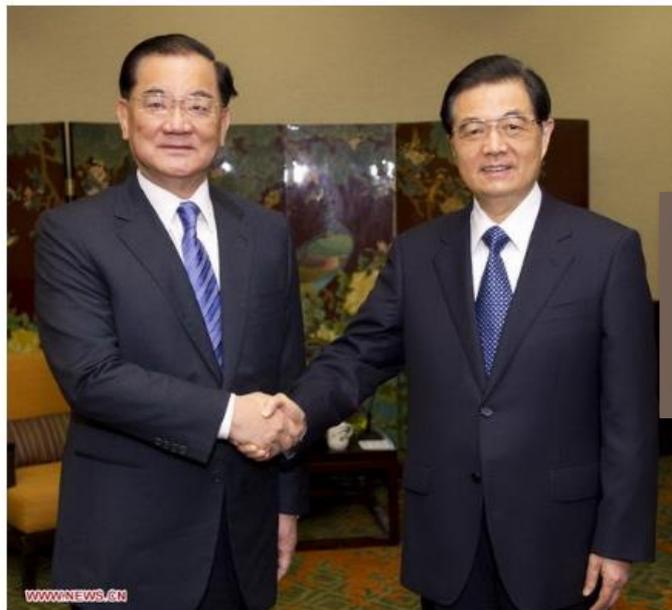
Beijing Sunny 15 / 1 City Forecast

Home >> China Politics

Hu meets KMT honorary chairman in Hawaii

(Xinhua)

11:10, November 12, 2011 🔍 +-



Chinese President Hu Jintao (R) shakes hands with Honorary Chairman of the Chinese Kuomintang (KMT) Lien Chan, in Honolulu, Hawaii, the U.S., Nov. 11, 2011.

Selections for you



Miao ethnic group celebrates Miao's New Year in SW China



World's first Angry Birds exclusive shop opens in Helsinki

Who is Hu Jintao?

Most Popular

- 1 Hu reaffirms support to Hong Kong's sta...
- 2 Hu meets KMT honorary chairman in Hawaii
- 3 China in APEC: a mutually beneficial en...
- 4 Night life in Shanghai
- 5 China's 2011 foreign trade to grow 20 p...
- 6 Beijing house prices stumble 5.1 pct as...
- 7 Lama students start school in Tibet Col...
- 8 Police in central China crack phoney ca...



Hu Jintao



Ricerca

Circa 725.000 risultati (0,09 secondi)

- Tutto
- Immagini
- Mappe
- Video
- Notizie
- Shopping
- PIÙ conte

Tutti i ri
Per argomento

Qualsiasi dimensione

- Grandi
- Medie
- Icone
- Maggiori di...
- Dimensioni esatte...

Qualsiasi colore

- A colori
- Bianco e nero



Qualsiasi tipo

- Volti
- Foto
- Clip art
- Disegni

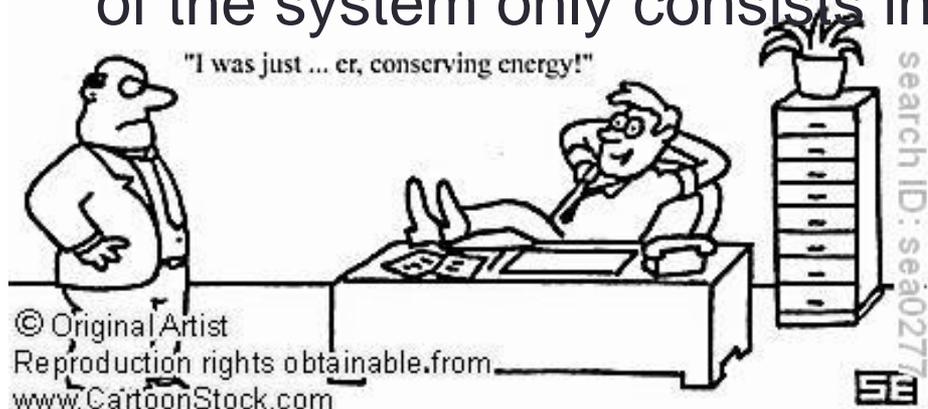
Visual standard

Mostra dimensioni



Benefits of a data-driven approach

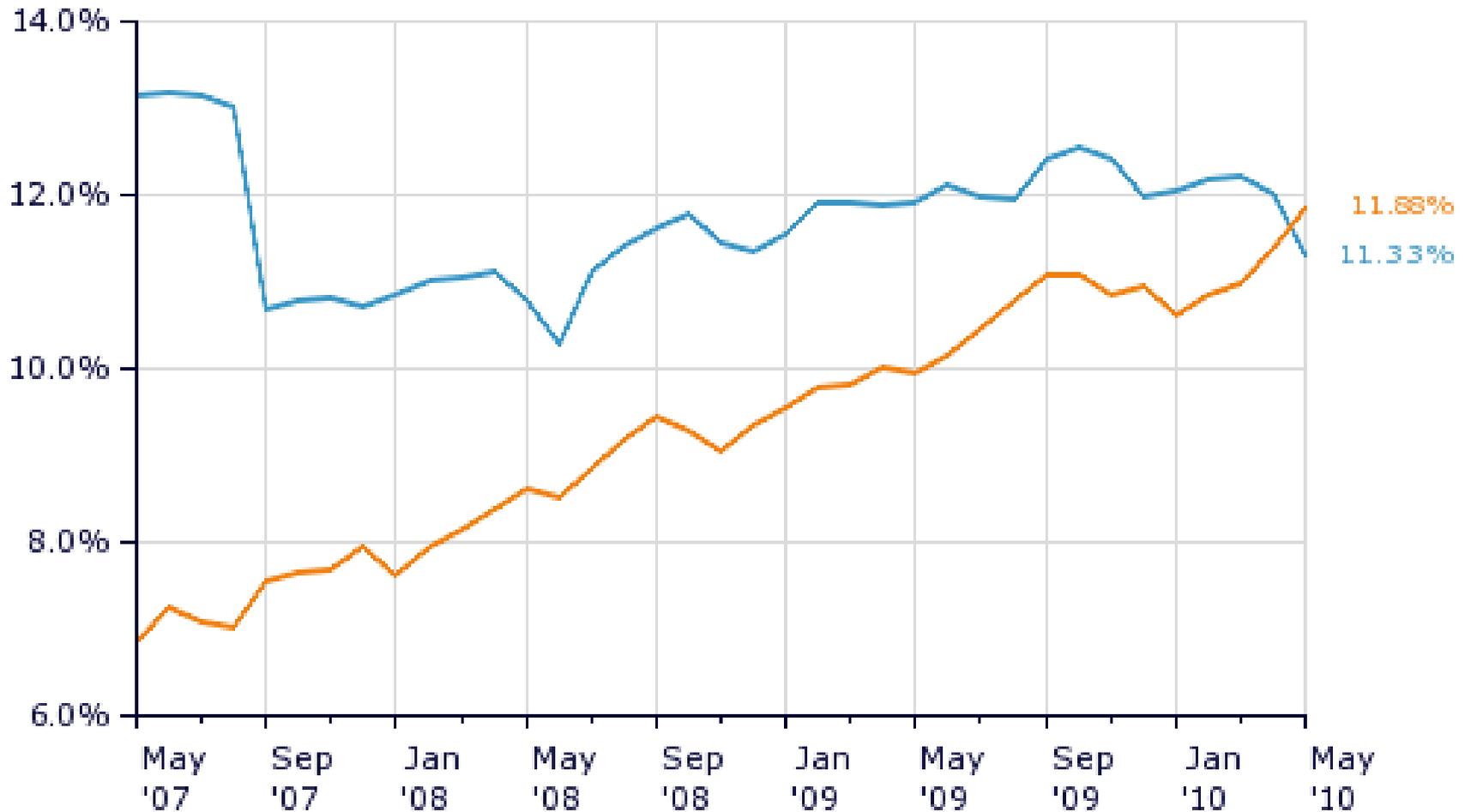
- Very effective learning algorithms available (e.g. Support Vector Machines)
- The ML technology is *portable* while imperative coding is *task (i.e. scenario) specific*
- Very accurate solutions can be obtained
- Gathering training data *much less expensive* than rule coding
- In dynamically evolving scenarios, incremental refinement of the system only consists in re-training



Data Mining: perspectives and benefits

- Technical advantages
 - Self adaptivity to changing operational conditions (i.e. domain)
 - Better SW management and incremental maintenance
 - More flexibility for special-purpose versioning:
 - No need for re-engineering or independent software developments
 - Just new domain-specific examples are needed
- Cost benefits
 - The data-driven approach has been shown to reduce the development costs up to 80-90% in several NLP tasks
- Market benefits
 - Reduced time-to-market
 - Competitive advantages: the lack of similar products makes the system targeted strongly competitive solutions

UK Internet visits to Social Networks and Search Engines



- Computers and Internet - Search Engines
- Computers and Internet - Social Networking and Forums

Monthly market share in 'All Categories', measured by visits, based on UK usage.

Created: 03/06/2010. © Copyright 1996-2010 Hitwise Pty. Ltd. Source: Experian Hitwise UK



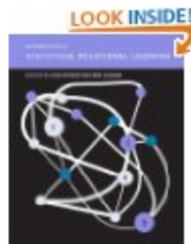
amazon.com

Recommended for You

Amazon.com has new recommendations for you based on [items](#) you purchased or told us you own.



[Collective Intelligence in Action](#)



[Introduction to Statistical Relational Learning \(Adaptive Computation and Machine Learning\)](#)



[Quantitative Methods In Linguistics](#)

[See More Recommendations](#)



[Collective Intelligence in Action](#)
by Satnam Alag
Average customer review: ★★★★★

[Add to cart](#)

[Add to Wish List](#)



Home

Preferences

Links

Documents

Contact

Roberto Basili

Search

All results (99)

- » [alessandro moschitti \(27\)](#)
- » [maria teresa pazienza \(25\)](#)
- » [tor vergata \(12\)](#)
- » [department of computer science systems \(9\)](#)
- » [artificial intelligence and human oriented \(6\)](#)
- » [database of individual seismogenic \(2\)](#)
- » [vigna murata 605 it 00143 \(2\)](#)
- » [fault mapper \(2\)](#)
- » [musical genre a machine learning \(2\)](#)
- » [roberto basili fabio \(2\)](#)

[More clusters](#)

YOU ARE IN "ALESSANDRO MOSCHITTI" CLUSTER WITH 27 DOCUMENTS

DBLP: ROBERTO BASILI

Paolo Annesi, **Roberto Basili**: Cross-Lingual Alignment of FrameNet Annotations through Hidden Markov Models. ... **Roberto Basili**, Cristina Giannone, Chiara Del Vescovo, Alessandro ...

<http://dblp.uni-trier.de/db/indices/a-tree/b/Basili:Roberto.html>

PUBZONE - ROBERTO BASILI

Roberto Basili, Danilo Croce, Cristina Giannone, Diego De Cao. ... **Roberto Basili**, Cristina Giannone, Chiara Del Vescovo, Alessandro Moschitti, Paolo Naggari. Kernel-Based ...

<http://www.pubzone.org/pages/publications/showAuthor.do?userId=79.2000>

LEARNING DOMAIN-SPECIFIC FRAMENETS FROM TEXTS

Marco Pennacchiotti, Diego De Cao, Paolo Marocco, **Roberto Basili**, ... Marco Pennacchiotti, Diego De Cao, **Roberto Basili**, Danilo Croce, Michael Roth, Automatic ...

<http://olp.dfki.de/olp3/Basili.pdf>

DIEGO DE CAO HOME PAGE

Roberto Basili, Danilo Croce, Diego De Cao, and Cristina Giannone. ... **Roberto Basili**, Diego De Cao, Danilo Croce, Bonaventura Coppola, and Alessandro Moschitti. ...

<http://art.uniroma2.it/decao/>



All results (98)

- » [computer](#) (33)
- » [learning algorithm for machine](#) (11)
- » [machine learning course](#) (6)
- » [machine learning group](#) (6)
- » [introduction to machine learning](#) (3)
- » [machine learning applications](#) (3)
- » [artificial intelligence](#) (7)
- » [applying machine learning](#) (3)
- » [challenges machine learning](#) (3)
- » [machine learning the study](#) (5)

[More clusters](#)

TOP 98 RESULTS OF RETRIEVED FOR THE QUERY MACHINE LEARNING

MACHINE LEARNING - WIKIPEDIA, THE FREE ENCYCLOPEDIA

Machine learning is a scientific discipline that is concerned with ... A major focus of **machine learning** research is to automatically learn to recognize complex ...

http://en.wikipedia.org/wiki/Machine_learning

MACHINE LEARNING: DEFINITION FROM ANSWERS.COM

machine learning (m??sh?n ?l?rni?) (computer science) The process or technique by which a device modifies its own behavior as the result of its

<http://www.answers.com/topic/machine-learning>

CMU 15-859 MACHINE LEARNING THEORY

Course description: This course will focus on theoretical aspects of **machine learning**. We will examine questions such as: What kinds of guarantees ...

<http://machinelearning.com/>

THE INTERNATIONAL MACHINE LEARNING SOCIETY - ABOUT

The International **Machine Learning** Society is a non-profit organisation whose main aim is to foster **machine learning** research and whose main activity ...

<http://www.machinelearning.org/>

MACHINE LEARNING: WEEKLY STUDY GUIDE

Weekly study guide for the course on **Machine Learning** taught by Vasant Honavar at Iowa

Semantics and News

Applicazioni Risorse Sistema mar 27 lug, 23.47 dan

Gmail ... x SRL_EN x Come ... x R Econo... x Googl... x Tanl It... x Frame... x SRL_EN x Econo... x

file:///home/danilo/Downloads/SRL_ITA/sorgente/Economia%20-%20Repubblica.it.html

Telefilm in stream... Telefilm in stream... Flash Forward pri... Telefilm in stream... Cronologia Altri Pr



L'ad punta a nuove regole sulla base del modello Pomigliano. L'annuncio, che prevede l'uscita da Federmeccanica, domani al vertice con il governo o giovedì con una lettera a Bombassei. Potrebbe avvenire assieme alla decisione di creare una new company per

Pomigliano di SALVATORE TROPEA

Cisl-Uil: "L'accordo di categoria non si tocca" di S. PAROLA

Sacconi: "Su Fiat partita aperta"

Nasce Fabbrica Italia Pomigliano

Si dimette il capo di Bp buonuscita un milione di sterline



Oggi l'annuncio: a Tony Hayward subentrerà il direttore esecutivo Robert Dudley. **I costi legati al disastro sono saliti a 32,2 miliardi di dollari**, ma la società li deterrà evitando di versare al fisco Usa 10 miliardi

Manager Usa, è Ellison di Oracle il più pagato del decennio



Ha guadagnato 1,84 miliardi di dollari. Nella classifica del *Wall Street Journal* sui leader delle società quotate, secondo con 1,14 miliardi il capo di Expedia, terzo Irani di Occidental Petroleum. Solo quarto Steve Jobs

Il nemico alle porte

La Consob e la mano invisibile

Altri articoli

PICCOLE GRANDI IMPRESE
DI LUCA PAGNI

La grande sfida del teleshopping

La crisi colpisce anche i porti turistici ma siamo sicuri che sia un male?

Altri articoli

PERCENTUALMENTE
DI ROSARIA AMATO

La prova del 9

L'export risolve il Pil, ma non le famiglie

Altri articoli

GLI ESPERTI RISPONDONO

CASA
A cura di Antonella Donati

Compenso extra, quando ne ha diritto l'amministratore

Mia moglie ed il fratello sono proprietari di un appartamento in condominio. Allo stato

Il tuo libro arriva dove
hai sempre sognato.

ilmiolibro.it

24ORE AGI

Roma 19:04
ACEA: NEL I SEMESTRE UTILE NETTO +52,1% A 2010, MLN

Parigi 18:42
AIR FRANCE-KLM: TORNA IN UTILE NEL PRIMO TRIMESTRE

← 3 → Le altre not

CREDITO ALLE IMPRESE

Microimprese: con la crisi aumenta il rischio di credito

IN COLLABORAZIONE CON

**WEB MINING & RETRIEVAL,
a.a. 2015-16**

Moduli	Argomenti	Lezioni
Basic ML	Introduzione all'algorithmica di ML. Probabilità e Metriche di similarità. Classificazione mediante algoritmi di base	Introduzione al Corso, ed al WM&R.
		Richiami ai metodi base di ML. Metodi Supervised vs. Unsupervised.
		Metodi probabilistici e generativi. Naive Bayes.
		Esercitazione. Decision Trees. Intro Weka (1)
Clustering and Probabilistic methods	Introduction to clustering algorithms. Generative Models. HMMs	Metodi algebrico-lineari. Metriche e similarità semantica.
		Metodi Unsupervised. Clustering. Metodi Agglomerativi: K-means
		Esercitazione. Metodi avanzati di Clustering. Uso di Weka (2)
		Modelli di Linguaggio. Processi Markoviani.
From PAC learnability to SVM	PAC learnability. VC-dimension. SVMs. Kernel methods	Modelli Generativi: HMM. (POS tagging)
		PAC Learnability. Perceptron
		SVM. Hard Margin.
		Soft margin SVM. La Nozione di Kernel.
Advanced ML topics: Neural Networks	Semi-supervised learning: ensemble methods, active learning, EM. On-line learning: Passive-Aggressive. Deep Neural network architectures.	Kernel polinomiali e RBF. Sequence & Tree Kernels.
		Use of Tree and Sequence Kernels in KELP
		Improving ML settings: Semisupervised Extensions
		Improving ML settings: <i>On-line Learning con esercitazione</i>
Basic IR topics	Ad hoc IR. Text Classification. IR models for texts and queries. Lexical Resources for IR.	From neural networks to deep learning: perceptrons and MLP
		Deep Learning over MLPs. Recurrent Neural Networks.
		Adopting Convolutional Neural Networks on images,
		IR tradizionale. Paradigmi e Architettura.
Advanced IR topics	Dimensionality Reduction for IR. Distributional Semantic models for IR. Word Embeddings with NNs.	IR tradizionale. Modelli algebrico-lineari
		Introduzione a Lucene: Perf Eval su Cranfield
		Query Operations/Expansion - Automatic Global Analysis - Thesaurus/SA
		Latent Semantic Analysis. LSA.
Web Information Retrieval	Web Applications & IR. Web search and Crawling. Random surfer models: Google PageRank	Latent Semantic Kernels & Semantic Kernels. Wordspaces
		Word Embeddings through Neural Networks.
		Word Space Exercise
		Introduzione al web and Social Medi Analysis
Social Media Analytics (*)	IR in Social Media. Community detection. User profiling and Recommending. Sentiment and Emotion Analysis.	Web Search: Rango e Relevance: PageRank. HITS
		Thematic Web Crawling project: Development of an indexed Web dump
		HITS. Web Crawling
		Social Media Analysis: Community Detection
Statistical Natural Language Processing (*)	Standard Natural Language Processing. Statistical Parsing. Evaluation of Statistical NLP tools.	Opinion Mining e Sentiment Analysis: the task
		Introduzione all'OM & SA: Twitter as a case study
		Development of a SA tool
		Social Media Analysis: Recommending
Statistical Natural Language Processing (*)	Standard Natural Language Processing. Statistical Parsing. Evaluation of Statistical NLP tools.	Richiami ai metodi di Elaborazione del Linguaggio Naturale: il TAL
		Risorse: Corpora, Lessici, Grammatiche e Basi di Conoscenza per il TAL.
		HMM-based POS tagging.
		Esercitazione: Stanford NLP chain. Evaluation.

WEB MINING & RETRIEVAL, a.a. 2015-16

Moduli	Argomenti	Lezioni
Basic ML	Introduzione all'algorithmica di ML. Probabilità e Metriche di similarità. Classificazione mediante algoritmi di base	Introduzione al Corso, ed al WM&R.
		Richiami ai metodi base di ML. Metodi Supervised vs. U
		Metodi probabilistici e generativi. Naive Bayes.
		Esercitazione. Decision Trees. Intro Weka (1)
Clusering and Probabilistic methods	Introduction to clustering algorithms. Generative Models. HMMs	Metodi Unsupervised. Clustering. Metodi Agglomerativi
		Esercitazione. Metodi avanzati di Clustering. Uso di W
		Modelli di Linguaggio. Processi Markoviani.
		Modelli Generativi: HMM. (POS tagging)
From PAC learnability to SVM	PAC learnability. VC-dimension. SVMs. Kernel methods	PAC Learnability. Perceptron
		SVM. Hard Margin.
		Soft margin SVM. La Nozione di Kernel.
		Kernel polinomiali e RBF. Sequence & Tree Kernels. Use of Tree and Sequence Kernels in KELP
Advanced ML topics: Neural Networks	Semi-supervised learning: ensemble methods, active learning, EM. On-line learning: Passive-Aggressive. Deep Neural network architectures.	Improving ML settings: Semisupervised Extensions
		Improving ML settings: <i>On-line Learning con esercitaz</i>
		From neural networks to deep learning: perceptrons and
		Deep Learning over MLPs. Recurrent Neural Networks. Adopting Convolutional Neural Networks on images,
Basic IR topics	Ad hoc IR. Text Classification. IR models for texts and queries. Lexical Resources for IR.	IR tradizionale. Paradigmi e Architettura.
		IR tradizionale. Modelli algebrico-lineari
		Introduzione a Lucene: Perf Eval su Cranfield
	Dimensionality Reduction for IR.	Query Operations/Expansion - Automatic Global Analy
		Latent Semantic Analysis. LSA.
		Latent Semantic Kernels & Semantic Kernels. Word

Clusering and Probabilistic methods	Introduction to clustering algorithms. Generative Models. HMMs	Metodi Unsupervised. Clustering. Metodi Agglomerativi: K-means Esercitazione. Metodi avanzati di Clustering. Uso di Weka (2 Modelli di Linguaggio. Processi Markoviani. Modelli Generativi: HMM. (POS tagging)
From PAC learnability to SVM	PAC learnability. VC-dimension. SVMs. Kernel methods	PAC Learnability. Perceptron SVM. Hard Margin. Soft margin SVM. La Nozione di Kernel. Kernel polinomiali e RBF. Sequence & Tree Kernels. Use of Tree and Sequence Kernels in KELP
Advanced ML topics: Neural Networks	Semi-supervised learning: ensemble methods, active learning, EM. On-line learning: Passive-Aggressive. Deep Neural network architectures.	Improving ML settings: Semisupervised Extensions Improving ML settings: <i>On-line Learning con esercitazione</i> From neural networks to deep learning: perceptrons and MLP Deep Learning over MLPs. Recurrent Neural Networks. Adopting Convolutional Neural Networks on images,
Basic IR topics	Ad hoc IR. Text Classification. IR models for texts and queries. Lexical Resources for IR.	IR tradizionale. Paradigmi e Architettura. IR tradizionale. Modelli algebrico-lineari Introduzione a Lucene: Perf Eval su Cranfield Query Operations/Expansion - Automatic Global Analysis - Th
Advanced IR topics	Dimensionality Reduction for IR. Distributional Semantic models for IR. Word Embeddings with NNs.	Latent Semantic Analysis. LSA. Latent Semantic Kernels & Semantic Kernels. Wordspaces Word Embeddings through Neural Networks. Word Space Exercise
Web Information Retrieval	Web Applications & IR. Web search and Crawling. Random surfer models: Google PageRank	Introduzione al web and Social Medi Analysis Web Search: Rango e Relevance: PageRank. HITS Thematic Web Crawling project: Development of an indexe HITS. Web Crawling
Social Media Analytics (*)	IR in Social Media. Community detection. User profiling and Recommending. Sentiment and Emotion Analysis.	Social Media Analysis: Community Detection Opinion Mining e Sentiment Analysis: the task Introduzione all'OM & SA: Twitter as a case study Development of a SA tool Social Media Analysis: Recommending
Statistical Natural Language Processing (*)	Standard Natural Language Processing. Statistical Parsing. Evaluation of Statistical NLP tools.	Richiami ai metodi di Elaborazione del Linguaggio Naturale: il Risorse: Corpora, Lessici, Grammatiche e Basi di Conoscenza p HMM-based POS tagging. Esercitazione: Stanford NLP chain. Evaluation.

Laboratori del Corso

- Nella'ambito dei Laboratori agli studenti saranno resi disponibili:
 - Piattaforme di Machine Learning: Weka, R, KELP
 - Motori di Ricerca: Lucene, Terrier
 - Strumenti di AI per l'elaborazione dei testi:
 - Recursive Neural Networks per l'apprendimento di lessici vettoriali
 - Parser grammaticali di linguaggi naturali (ita,eng)
 - Ambienti di Enterprise Semantic Search su Web

Kelp: Java-based kernel framework

The screenshot shows the website for the Semantic Analytics Group at the University of Rome, Tor Vergata. The page features a navigation menu with links for People, Research, Teaching, Publications, Projects, Demo & Software, and Contacts. Below the navigation is a large image of a glowing blue brain circuit. The main content area is titled "KeLP (Kernel-based Learning Platform)" and describes it as a machine learning platform developed within the SAG group, written in Java and focused on Kernel Machines. It includes a "Downloads" section with links to Maven and the installation page, and a "News" section with several recent updates.

Semantic Analytics Group @ Uniroma2
SAG is the Semantic Analytics Group at the University of Rome, Tor Vergata

People Research Teaching Publications Projects Demo & Software Contacts

KeLP (Kernel-based Learning Platform)

KeLP (Kernel-based Learning Platform) is a machine learning platform developed within the SAG group. It is entirely written in Java and it is strongly focused on *Kernel Machines*. It includes different Online and Batch Learning and Classification algorithms, Kernel functions, ranging from vector-based to structural kernels. **KeLP** allows to build complex kernel machine based systems, leveraging on the Java language and on a JSON interface to store and load classifiers configurations as well as to save the models to be reused.

For a deeper look, you can visit [What's inside KeLP page](#).

Downloads

KeLP is released under [Maven](#). To use it, please refer to the [Installation](#) page

To download **KeLP** source code you can go to the github [KeLP page](#).

Authentication

[Log In](#)

News

- [SAG's KeLP team ranked first at the SemEval 2016 Community Question Answering Task](#) February 16, 2016
- [KeLP 2.0.2 released!](#) February 16, 2016
- [KeLP 2.0.1 released](#) January 13, 2016
- [The ECIR 2016 paper has been accepted!](#) December 30, 2015
- [KeLP 2.0.0 released](#) December 4, 2015
- [SAG with Reveal @ Maker Faire 2015, Rome!!](#) October 16, 2015

`https://github.com/SAG-KeLP`

`http://sag.art.uniroma2.it/demo-software/kelp/`

KELP applications: cQA

General Description

Subtasks

Data and Tools

Important Dates

Results

Call for Papers

SemEval-2016 Task 3

Task 3: Community Question Answering

Building on the success of [SemEval 2015 Task 3](#) "Answer Selection in Community Question Answering" (see [the task description paper](#)), we propose an extension, which covers a full task on Community Question Answering (CQA) and which is, therefore, closer to a real application (see, e.g., [Qatar Living forum](#)).

CQA systems are gaining popularity online. Such systems are seldom moderated, quite open, and thus they have little restrictions, if any, on who can post and who can answer a question. On the positive side, this means that one can freely ask any question and expect some good, honest answers. On the negative side, it takes effort to go through all possible answers and to make sense of them. For example, it is not unusual for a question to have hundreds of answers, which makes it very time-consuming for the user to inspect and to winnow through them all. The present task could help to automate the process of finding good answers to new questions in a community-created discussion forum (e.g., by retrieving similar questions in the forum and by identifying the posts in the answer threads of those similar questions that answer the original question well).

In essence, the main CQA task can be defined as follows:

"given (i) a new question and (ii) a large collection of question-comment threads created by a user community, rank the comments that are most useful for answering the new question"

Contact Info

Organizers

- Preslav Nakov, Qatar Computing Research Institute, HBKU
- Lluís Màrquez, Qatar Computing Research Institute, HBKU
- Alessandro Moschitti, Qatar Computing Research Institute, HBKU
- Walid Magdy, Qatar Computing Research Institute, HBKU
- James Glass, CSAIL-MIT
- Bilal Randeree, Qatar Living

email : semeval-cqa@googlegroups.com

Other Info

Announcements

KELP applic

General Description

SemEval-20

Team ID	Team Affiliation
ConvKN	Qatar Computational
ECNU	East China Normal University
ICL00	Institute of Computing Technology, Chinese Academy of Sciences
ICRC-HIT	Intelligence and Information Technology Center, Harbin Institute of Technology
ITNLP-AiKF	Intelligence and Information Technology Center, Harbin Institute of Technology
Kelp	University of Cambridge
MTE-NN	Qatar Computational
overfitting	University of Cambridge
PMI-cool	Sofia University "St. Kliment Ohridski", Sofia
QAIIT	IIT Hyderabad
QU-IR	Qatar University
RDI.team	RDI Egypt, Cairo
SemanticZ	Sofia University "St. Kliment Ohridski", Sofia
SLS	MIT Computational Science and Technology Laboratory
SUper.team	Sofia University "St. Kliment Ohridski", Sofia
UH-PRHLT	Pattern Recognition and Computer Vision Group, Universitat Politècnica de Catalunya
UniMelb	The University of Melbourne
UPC_USMBA	Universitat Politècnica de Catalunya

Tal
shc
tea
the
rut

	Submission	MAP	AvgRec	MRR	P	R	F1	Acc
1	Kelp-primary	79.19₁	88.82₁	86.42₁	76.96₁	55.30₈	64.36₆	75.11₂
	ConvKN-contrastive1	78.71	88.98	86.15	77.78	53.72	63.55	74.95
	SUper.team-contrastive1	77.68	88.06	84.76	75.59	55.00	63.68	74.50
2	ConvKN-primary	77.66₂	88.05₃	84.93₄	75.56₂	58.84₆	66.16₂	75.54₁
	3 SemanticZ-primary	77.58₃	88.14₂	85.21₂	74.13₄	53.05₁₀	61.84₈	73.39₅
4	ConvKN-contrastive2	77.29	87.77	85.03	74.74	59.67	66.36	75.41
	4 ECNU-primary	77.28₄	87.52₅	84.09₆	70.46₆	63.36₄	66.72₁	74.31₄
5	SemanticZ-contrastive1	77.16	87.73	84.08	75.29	53.20	62.35	73.88
	5 SUper.team-primary	77.16₅	87.98₄	84.69₅	74.43₃	56.73₇	64.39₄	74.50₃
	MTE-NN-contrastive2	76.98	86.98	85.50	58.71	70.28	63.97	67.83
	SUper.team-contrastive2	76.97	87.89	84.58	74.31	56.36	64.10	74.34
	MTE-NN-contrastive1	76.86	87.03	84.36	55.84	77.35	64.86	65.93
	SLS-contrastive2	76.71	87.17	84.38	59.45	67.95	63.41	68.13
6	SLS-contrastive1	76.46	87.47	83.27	60.09	69.68	64.53	68.87
	6 MTE-NN-primary	76.44₆	86.74₇	84.97₃	56.28₉	76.22₁	64.75₃	66.27₈
7	7 SLS-primary	76.33₇	87.30₆	82.99₇	60.36₈	67.72₃	63.83₆	68.81₇
	ECNU-contrastive2	75.71	86.14	82.53	63.60	66.67	65.10	70.95
	SemanticZ-contrastive2	75.41	86.51	82.52	73.19	50.11	59.49	72.26
8	ICRC-HIT-contrastive1	73.34	84.81	79.65	63.43	69.30	66.24	71.28
	8 ITNLP-AiKF-primary	71.52₈	82.67₉	80.26₈	73.18₅	19.71₁₂	31.06₁₂	64.43₉
	ECNU-contrastive1	71.34	83.39	78.62	66.95	41.31	51.09	67.86
9	9 ICRC-HIT-primary	70.90₉	83.36₈	77.38₁₀	62.48₇	62.53₅	62.50₇	69.51₆
	10 PMI-cool-primary	68.79₁₀	79.94₁₀	80.00₉	47.81₁₂	70.58₂	57.00₉	56.73₁₂
11	UH-PRHLT-contrastive1	67.57	79.50	77.08	54.10	50.11	52.03	62.45
	11 UH-PRHLT-primary	67.42₁₁	79.38₁₁	76.97₁₁	55.64₁₀	46.80₁₁	50.84₁₁	63.21₁₀
	UH-PRHLT-contrastive2	67.33	79.34	76.73	54.97	49.13	51.89	62.97
12	12 QAIIT-primary	62.24₁₂	75.41₁₂	70.58₁₂	50.28₁₁	53.50₉	51.84₁₀	59.60₁₁
	QAIIT-contrastive2	61.93	75.22	69.95	49.48	49.96	49.72	58.93
	QAIIT-contrastive1	61.80	75.12	69.76	49.85	50.94	50.39	59.24
	Baseline 1 (IR)	59.53	72.60	67.83	—	—	—	—
	Baseline 2 (random)	52.80	66.52	58.71	40.56	74.57	52.55	45.26
	Baseline 3 (all 'true')	—	—	—	40.64	100.00	57.80	40.64
	Baseline 4 (all 'false')	—	—	—	—	—	—	59.36

Table 1: **Subtask A, English (Question-Comment Similarity):** results for all submissions. The first column shows the rank of the primary runs with respect to the official MAP score. The second column contains the team’s name and its submission type (primary vs. contrastive). The following columns show the results for the primary, and then for other, unofficial evaluation measures. The subindices show the rank of the

Natural Language Parsing tool: RevNLT

UK Economy News Headlines - FT.com - Mozilla Firefox

File Modifica Visualizza Cronologia Segnalibri Strumenti Aiuto

http://www.ft.com/world/uk/economy

Più visitati Corso: Basi di dati Gruppi Posta :: Benvenuto a H... ClustrMaps - map of vi... UniversitaCedol Tree Kernels in SVM-lig... Net RicercaAteneo Keysrc Calls EMEROTECA GEMS2010

The image displays a natural language parsing tree for the sentence: "Mortgage approvals fell sharply in June, lending yet more weight to the theory...". The tree is rooted at the top with a node labeled "1.00 V_PP". This node branches into three children: "1.00 V_Sog", "1.00 V_Adv", and "1.00 Adv_PP".

- The "1.00 V_Sog" node branches into "Mortgage_approvals" (type Nom) and "fell" (type VerFin, labeled as "Sentence").
- The "1.00 V_Adv" node branches into "sharply" (type Adv).
- The "1.00 Adv_PP" node branches into "in_June" (type Prep) and a "CongCo" node.
- The "in_June" node branches into "in" (type IN, morph invariante) and "'June'" (type NNP, morph mas.fem.sing.plur.).
- The "CongCo" node branches into "lending_yet_more_weight" (type Nom) and "to_the_theory" (type Prep).
- The "lending_yet_more_weight" node branches into "lending" (type NN, morph mas.fem.sing.), "yet" (type RB, morph invariante), and "more" (type JJR, morph mas.fem.plur.sing.).
- The "to_the_theory" node branches into "to" (type CongCo) and "the_theory" (type Nom).

Business Education Personal Finance Arts & Leisure Wealth In depth

Britain's place in the world, and how far it has travelled since 1947 - Jul-29

Gilts lose lustre for overseas investors
Flight from eurozone risk to UK government bonds is moderating - Jul-29

with Alex Barker and Jim Pickard

Mechanical & Electrical Engineering
Deputy Director of Finance
London Ambulance Service
RECRUITERS

http://www.ft.com/westminster

Italiano (Italia)

Today is: 2006-07-06 17:14:55

00:00:02:39

Timeline bar with play, stop, and volume controls.

Info	Transcription	Semantic Analysis	Content Analysis
00:06:36	chamonix america dove perde forza ma fa sempre paura l' uragano di mallarme italia andiamo nel centro		
00:06:41	che in florida riguardasse cento km di costa sull' atlantico si e' formata nel frattempo un' altra tempesta tropicale		
00:06:52	ha lasciato una riviera messicana dello jucker puntando verso la florida l' uragano delle corde wilma il dodicesimo ciclone di una stagione eco dell' atmosfera piu' di qualcuno lavatrici su strada a festeggiare lo scampato pericolo mentre dall' altra l' emergenza ha segnato l' inizio dei sa scarseggiano cibo e acqua si e' costretti a fare i conti con la sopravvivenza ad attraversare queste strade inondate sferzata dal		
00:07:22	vento la pioggia per raggiungere i centri della croce rossa vengono distribuiti ieri alla popolazione dino risi ma ha lasciato otto vittime soltanto migliaia di casi devastato la rete ospedaliera abbattuto centrali elettriche che ha causato danni a un milione di persone in florida e' attesa per		
00:07:44	e nelle isole di kiss e' gia' iniziata la grande fuga non bastasse sull' atlantico a sud di porto rico si e' formata falla venti di hemingway le		
00:07:52	nessuna tempesta tropicale della stagione la buona notizia che dovrebbe essere innocua la brutta notizia che la stagione degli uragani		
00:07:59	non e' ancora finita nulla fino al trentanove e c' e' stato una sciagura		
00:08:05	in nigeria		

- TG1 - 2005-10-23
- Other Classification
- Other Classification
- Other Classification
- Ambiente, Natura e Territorio
- Ambiente, Natura e Territorio**
- Other Classification
- Politica, Partiti, Istituzioni e Sindacati
- Politica, Partiti, Istituzioni e Sindacati
- Other Classification
- Other Classification
- Other Classification
- Other Classification
- Usi e costumi
- Other Classification
- Sanita' e Salute
- Giustizia, Criminalita' e Sicurezza
- Other Classification
- Other Classification
- Giustizia, Criminalita' e Sicurezza
- Other Classification
- Other Classification
- Musica e Spettacolo
- Sport
- Sport

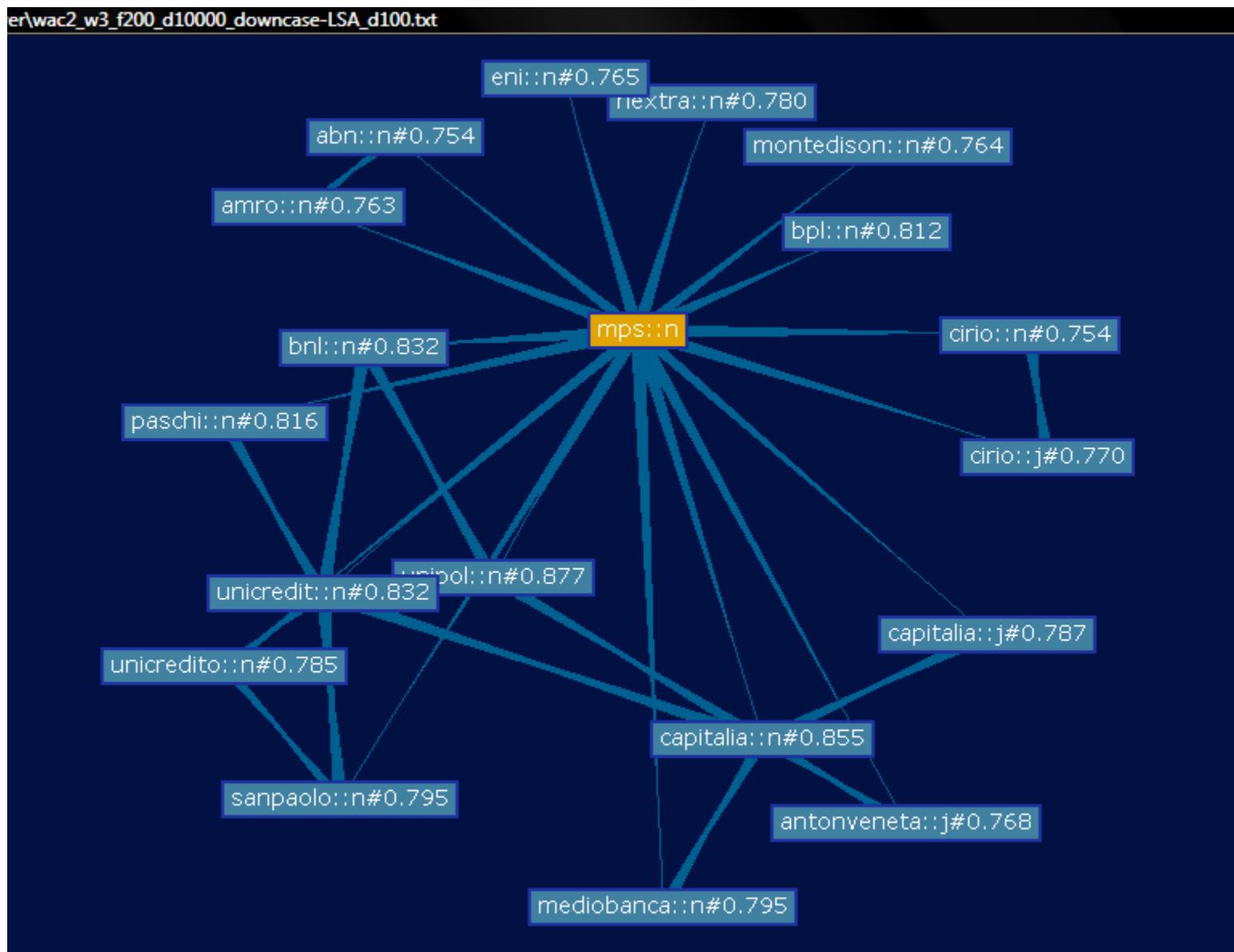
TIMELINE

00:06:36:10 00:06:51:07 00:06:54:20 00:06:56:01 00:06:57:10 00:07:00:08

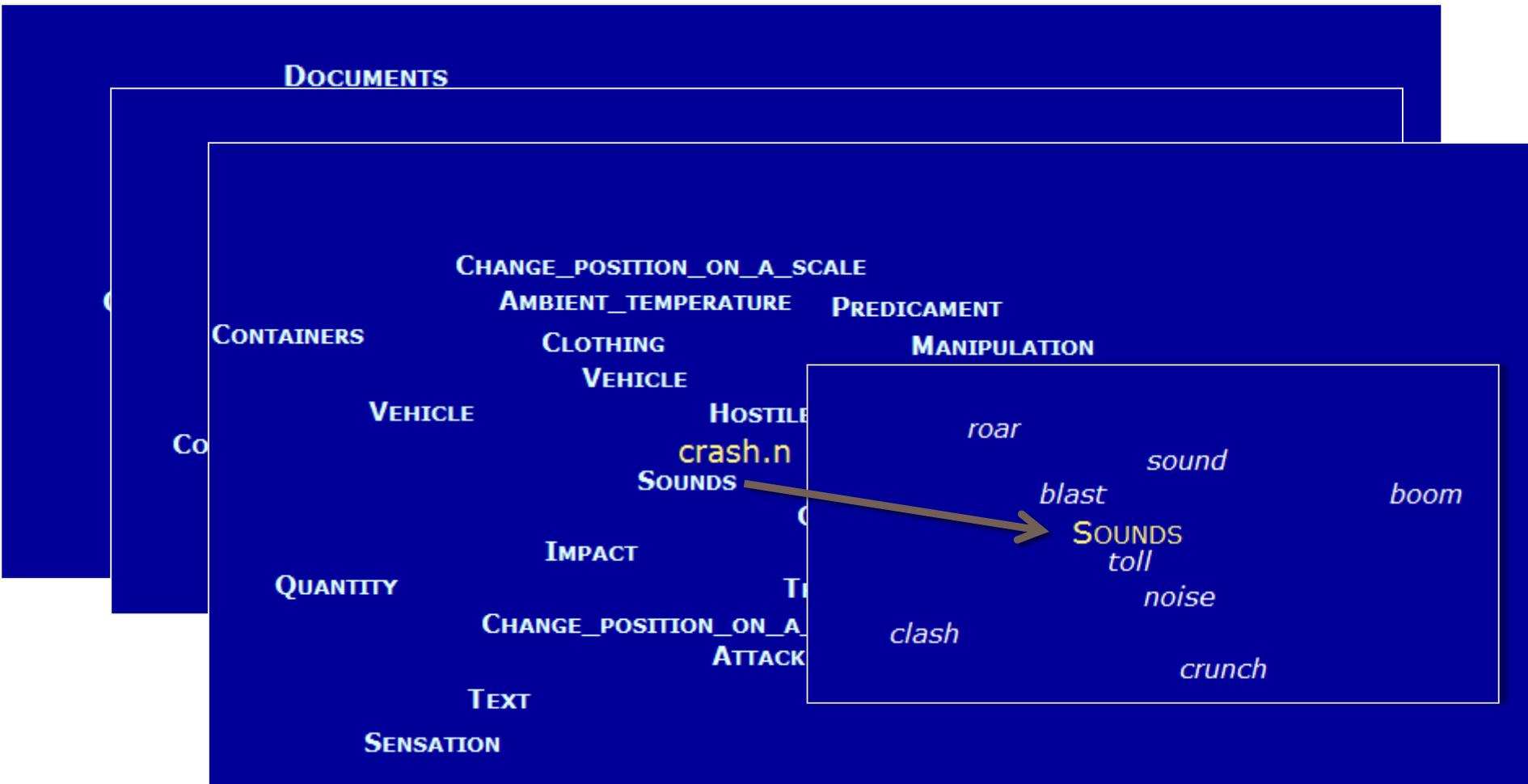
00:07:08:18 00:07:12:02 00:07:15:23 00:07:19:13 00:07:22:02 00:07:26:06

00:07:37:08 00:07:39:15 00:07:45:20 00:07:55:06 00:07:57:07 00:07:59:13

Vector Spaces for Lexical Semantics



Spaces for NL predicates



Sperimentazioni attive su

<http://mscoco.org/>



cocodataset@outlook.com

Home People Explore **Dataset** External

MS COCO image and annotations



a group of people point towards a green bus.

COCO-Text annotations

- 1) 'men',
legible
machine printed
English
- 2) 'transport',
legible
machine printed
English

OCR results:
CORRECT

COCO-Text

COCO-Text is for both text detection and recognition. The dataset annotates scene text with transcriptions along with attributes such as legibility, printed or handwritten text.



FM-IQA

The Freestyle Multilingual Image Question Answering (FM-IQA) dataset contains over 120,000 images and 250,000 freestyle Chinese question-answer pairs and their English translations.



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is laid under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?

VQA

VQA is a new dataset containing open-ended questions about images. These questions require an understanding of vision, language and commonsense knowledge to answer.

References

- Mitchell, Tom. M. 1997. *Machine Learning*. New York: McGraw-Hill.
- [Kernel machines, neural networks and graphical models](#), P. Frasconi, A. Sperduti, A. Starita, Rivista AI*IA Numero speciale per i “50 anni di IA”, 2007.
- Very good video lectures by Andrew Ng (Stanford) <http://academicearth.org/courses/machine-learning>