

# Automatic Topic Detection in Persian Blogs

**Dr. Karine Megerdooomian and Dr. Ali Hadjarian**  
**The MITRE Corporation**

**8 February 2011**



© 2011 The MITRE Corporation. All rights Reserved.



# Why Blogs?

- **Bloggers provide a view into societies where polls may not be available**
  - Audience analysis for strategic communication
  - Understanding culture and beliefs
  - Public opinion
- **Cyber-communities emerge in the blogosphere**
  - Analyze communities of interest
  - Identify influential bloggers in the network
  - See what issues and ideas bring people together
- **Trends analysis**
  - Study idea propagation by tracking issues through time
- **But blogs are not always representative of the society**
  - Depends on country's infrastructure, literacy rate, government filtering, popularity of blogs

# Persian Blogosphere

- Estimates place Persian among the top ten languages of the global blog community
- Current blogs estimated at 700,000 with ~70,000 active blogs (updated at least once a week)
- The Persian blogosphere represents a wide spectrum of groups and opinions



Афсурдаги вакте ба сабк табдил мешавад  
Ёдам меояд вакте Самарқандро тақ мекардам, дидишам доштам ва ба  
хамаи сохтаҳои худ барои мекардам, ба дустон ва хам  
гирифтаам.



**زن شناسی**

چرا با مهریه مخالفیم؟!  
نویسنده: آیتام - ساعت ۱۰:۰۷ و ۱۰:۰۸ روز ۱۳۸۸/۱۲/۲۰

آیا شما به مساوات و عدالت یاور دارید؟ منظوری حقوق مساوی بین زن و مرد است. اصلن بیایید را به دست بگیرید و تمامی باورهای درست و نادرست خود از جنس علمی، مذهبی، فرهنگی را ساده مطرح می کنیم: تساوی حقوق زن و مرد در یک رابطه مشترک به چه معنی است؟ اجرای از

**کانون وبلاگ نویسندگان افغانستان**

من | نمایندگی های کانون | چگونه به کانون بپیوندیم | مطالب مربوط به کانون | مطالب آموزشی وبلاگ

بنگاه ویژه به پاسان یام دنیا  
Friday, September 24, 2010

هر بیست جنگجوی افغانی که بعد از شهر بیست و یکم فرمانده شان دلیل موجهی داشت تا با جیب، یا قاطر از اکتبر ۲۰۰۱ از مرز پاکستان به سوی سرزمین شان گذشتند ن راه ها را ببیماید. بیش از ده سال پیش از آن، در مامور

# Challenges of Persian Blogs

## ■ Persian *diglossia*

- Two variants of the language, *literary* and *conversational*, are used side by side
- The conversational variant affects the lexicon, word forms, and even syntax
- Differences between major dialects (Farsi vs. Dari)

## ■ Non-standard orthography

- Spelling more directly represents the pronunciation of the words (e.g., *downloadz*, *dunno* in English)
- Variation in spelling the same word across blogs or within a single blog

## ■ Neologisms

- Widespread use of borrowings, newly created words by bloggers, and new official terms by the Persian Language Academy

## ■ Blog formats and encoding issues

- The format and encoding used are not consistent across blogs or within a single blog

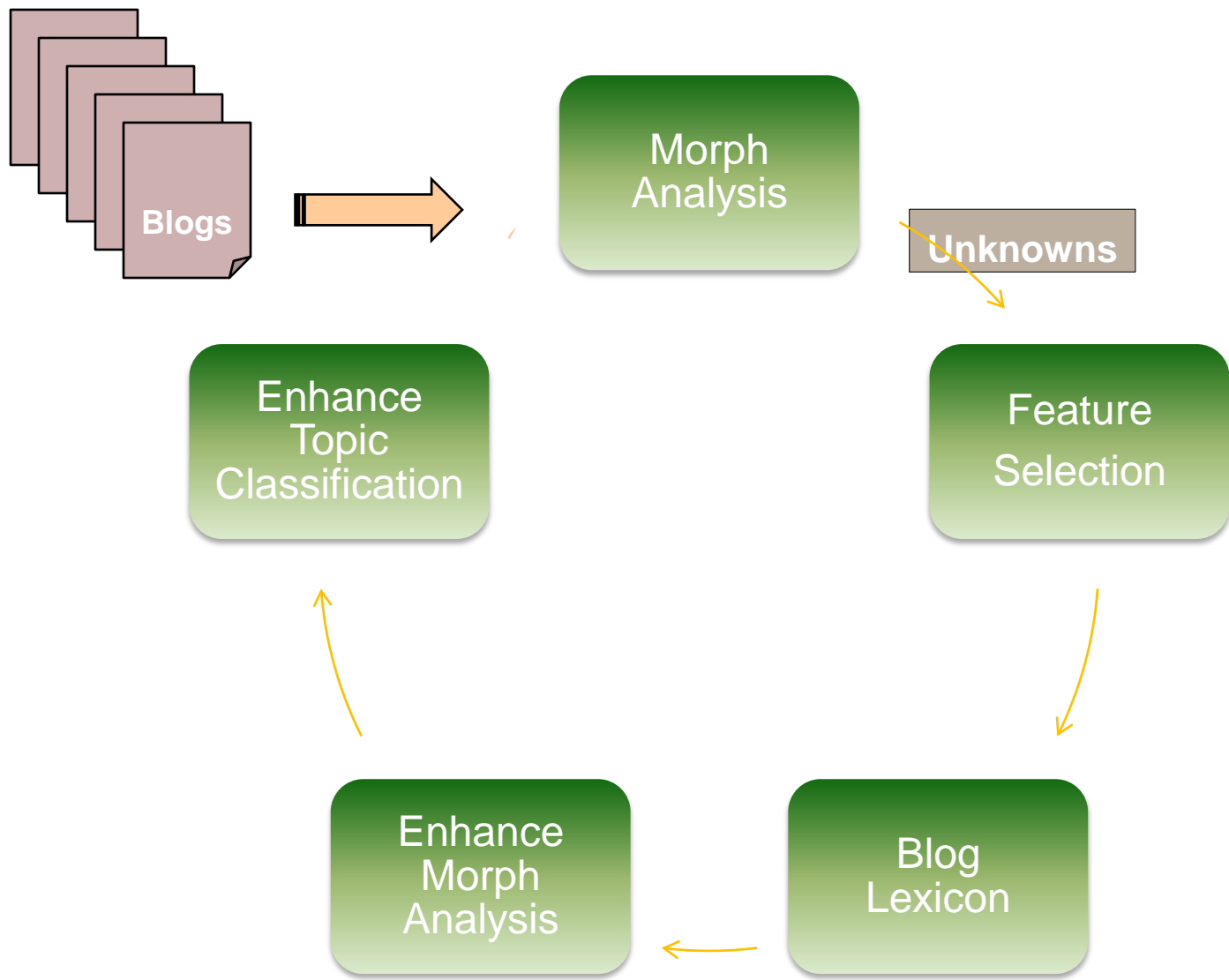
# Automatic Topic Detection

- **Develop automatic techniques for detecting topically related material in Persian language blogs.**
- **Combine language and text analytics with social network analysis to build a model of topic clusters**
- **Enhance system components:**
  - Lexicon of blog terms
  - Morphological analyzer to process blog language
  - Topic classification system



# Technical Approach

- **Automatically extract neologisms**
  - We employ morphological analysis in conjunction with a profile-based classification technique to extract a pertinent candidate list for identifying new word-level constructions in blogs.
- **Build a Blog Lexicon**
  - These neologisms are used to build a *blog lexicon* that is then integrated within the morphological analyzer.
- **Enhance the topic classification system**
  - We use the enhanced morphological analyzer and lexicon to extract a new set of features per topic category.



# Blog Data

<b>Topic</b>	<b>Number of blogs</b>	<b>Median size</b>	<b>Average number of words</b>
<b>Computers and Internet</b>	<b>497</b>	<b>14 kb</b>	<b>986</b>
<b>Cinema and theatre (sinama va ta'atr)</b>	<b>255</b>	<b>18 kb</b>	<b>1380</b>
<b>Political (siyasat-e-rooz)</b>	<b>500</b>	<b>22 kb</b>	<b>2171</b>
<b>Medical (pezeshki)</b>	<b>499</b>	<b>27 kb</b>	<b>2285</b>
<b>Sports (varzesh)</b>	<b>498</b>	<b>19 kb</b>	<b>1528</b>

# Linguistic Parsing

- **Morphological Analyzer (Amtrup 2003)**
  - Typed feature structures with unification
  - Full inflectional morphology, developed for MT application
  - Unanalyzed words are tagged as unknowns
  - 97% coverage and 93% accuracy on a 7 MB corpus collected from online news sources
- **Lexicon used for lookup**
  - About 40,000 entries in citation form and 5,000 proper names
  - Developed in 1999-2002 for coverage of online news articles
- **Unanalyzed tokens**
  - Conversational forms
  - Misspellings
  - Proper nouns
  - Specialized domain vocabulary
  - New words not in the lexicon

# Linguistic Parsing

The screenshot displays the Shooka UI interface, which is used for linguistic parsing. The main window is titled "Shooka UI" and contains several tabs: "Configuration", "Morphology", "Lexicon", and "TextAnalysis". The "TextAnalysis" tab is active, showing a text input field with the file path "C:\Home\Shooka\testfiles\Test1\BBCTestfile.txt" and an "Analyze" button. Below the input field, the text "C:\Home\Shooka\outFiles\BBCTestfile.txt.xml" is displayed, along with a "Load" button.

The main text area contains a Persian paragraph about the seizure of weapons in Pakistan. The text is as follows:

سپاه پاسداران را تکذیب کرد. یک افسر پلیس پاکستان به نام دادالرحمن به خبرنگاری اسوشیندپرس گفته است بازداشت شدگان تحت بازجویی قرار دارند و دو وسیله نقلیه آنها هم ضبط شده است. خبرنگاری رویترز به نقل از یک مقام امنیتی پاکستان که نخواست نامش فاش شود نوشت: "این مسئله بسیار جدی است. ما در حال بررسی این موضوع هستیم که چرا این افراد وارد خاک ما شده اند." چرا سربازان وارد خاک پاکستان شدند؟ در حال حاضر جزئیات ورود سربازان سپاه به خاک پاکستان مشخص نیست. شبکه پرس تیوی پس از اعلام خبر بازداشت سربازان ایرانی، آنها را "در تلاش برای به دام انداختن قاچاقچیان سوخت" توصیف کرد. اما خبرنگاری رویترز به نقل از یک مقام امنیتی دیگر در پاکستان اعلام کرد نیروی مرزی ایران به مقام های پاکستانی گفته اند که ورود به خاک پاکستان تصادفی بوده و "پس از آن رخ داد که سپاه پاسداران در نزدیکی مرز با پاکستان عملیاتی را علیه اعضا جندالله آغاز کردند." سپاه پاسداران در ماه های اخیر مسئولیت کنترل امنیت در سیستان و بلوچستان را به عهده داشته است. محمد پاکپور فرمانده نیروی زمینی سپاه پاسداران ساعاتی پیش از اعلام خبر بازداشت سربازان ایرانی در خاک پاکستان، در یک کنفرانس خبری به "قول مساعد پاکستان" برای مبارزه با جندالله اشاره کرد و گفت: "با توجه به این مسئله نیازی به وارد شدن به خاک پاکستان نیست." هر چند پیمان **فروزش**، نماینده زاهدان در مجلس پس از کشته شدن چند فرمانده سپاه در جریان یک بمبگذاری در جنوب استان سیستان و بلوچستان (در روز 26 مهر) گفته بود: "ما خواستار آن هستیم تا عملیات متقابلی از سوی سپاه در خاک پاکستان انجام شود." ایران، **جندالله**، پاکستان خبر دستگیری یازده سرباز ایران در خاک پاکستان تنها سه روز پس از سفر سر تیپ مصطفی محمد نجار، وزیر کشور ایران به پاکستان صورت می گیرد که در این سفر دو کشور درباره همکاری های امنیتی و مرزی گفتگو کردند. در این دیدار آقای نجار پس از ورود به خاک پاکستان از وجود "**شواهدی**" مبنی بر وجود عبدالمالک ریگی، رهبر گروه جندالله در خاک پاکستان صحبت به میان آورد و خواستار تحویل این فرد به ایران شد. این خواسته پس از آن مطرح شد که در روز 18 اکتبر (26 مهر) در اثر یک بمبگذاری در جنوب سیستان و بلوچستان که بعداً جندالله مسئولیت آن را به عهده گرفت، چند تن از فرماندهان سپاه پاسداران (از جمله جانشین فرمانده نیروی زمینی سپاه) کشته شدند. این حمله توسط شورای امنیت سازمان ملل متحد محکوم شد و مقام های پاکستانی نیز درباره آن ابراز تاسف کردند. اما در هنگام سفر وزیر کشور ایران به پاکستان، رحمان مالک همتای پاکستانی او حضور عبدالمالک ریگی در خاک پاکستان را رد کرد. جنبش مقاومت ملی ایران - جندالله - که رهبری آن را عبدالمالک ریگی برعهده دارد ظاهراً در سال 1382 در مخالفت با سیاست جمهوری اسلامی در قبال آنچه که تبعیض علیه پیروان تسنن در سیستان و بلوچستان می خواند تشکیل شد. این گروه مسئول یک رشته حملات مسلحانه، عملیات انتحاری، بمب گذاری و گروگانگیری معرفی شده است.

Below the text, the morphological analysis of the word "شودن" is shown:

Morphology: کتابهايمان  
Analyze  Show only longest analyses  
[0-11] +Noun+pl+clit.1pl+NPB [book:]

Lexicon: حزب  
حزب  
حزب اقلیت  
حزب اکثریت  
حزب جمهوری خواه  
حزب جمهوریخواه  
حزب حاکم  
حزب دیمکرات  
حزب محافظه کار

Hzb Hzb  
Irregular Plural  
POS: Noun  
Event:  Act  Become  Cause  
Headword: Hzb  
Translation: party;political party;  
Number:  Plural

LightVerb+Ind+perf.past+3sg [be introduced:]

# Information Gain

- Information Gain based feature selection
- Each attribute has a binary value which signifies the presence or absence of a given unknown term in the document

$$IG(D, t) = Entropy(D) - \frac{|D_t|}{|D|} Entropy(D_t) - \frac{|\overline{D}_t|}{|D|} Entropy(\overline{D}_t)$$

- For a term to be selected, it not only needs to have a high IG, but it needs to be present in a higher proportion of positive examples than the negative ones
  - This prevents the selection of terms that while are good descriptors of the negative class and thus carry a high IG, are not necessarily pertinent to the positive class (i.e., the topic category under consideration). So IG of a term not meeting the above constraint is effectively set to zero.
- Extract neologisms that would have the most discriminatory power in distinguishing between topics

# Weighted Unknowns for Computers

Transliteration		Weight	Translation
ویندوز	vyndvz	0.100033	Windows
دانلود	danlvd	0.080559	download
فایل	fayl	0.058319	file
کاربران	karbran	0.051595	users
جاوا	Java	0.048287	Java
کلیک	klyk	0.048180	click
ياهو	yahv	0.044999	Yahoo
نوکیا	nvkya	0.044807	Nokia
فلش	flG	0.042718	Flash
مرورگر	mrvrgr	0.041374	browser
هک	hk	0.041074	hack
مسنجر	msnJr	0.040853	Messenger
چت	Ct	0.039987	chat
پسورد	psvrd	0.039213	password
کد	kd	0.035936	code

# Blog Lexicon

## 1. Automatically Extract Neologisms

- Collected a corpus of Persian language blogs in 5 major topics (politics, cinema, computers, medicine, sports)
- Used a morphological analyzer to tag unknowns
- Generated a weight-ordered list of unknown words for each topic category using information gain
- Words with English translations extracted from Wikipedia (Burger 2009)

# Blog Lexicon

## 2. Manually collected terms

- **Collected words from existing glossaries**
  - Nuclear terms
  - Slang terms
  - Newly created technical and scientific terms from the Persian Language Academy
- **Performed corpus analysis**
  - Studied political and computer blogs to detect new usages that were not found in traditional lexicons using concordance tools

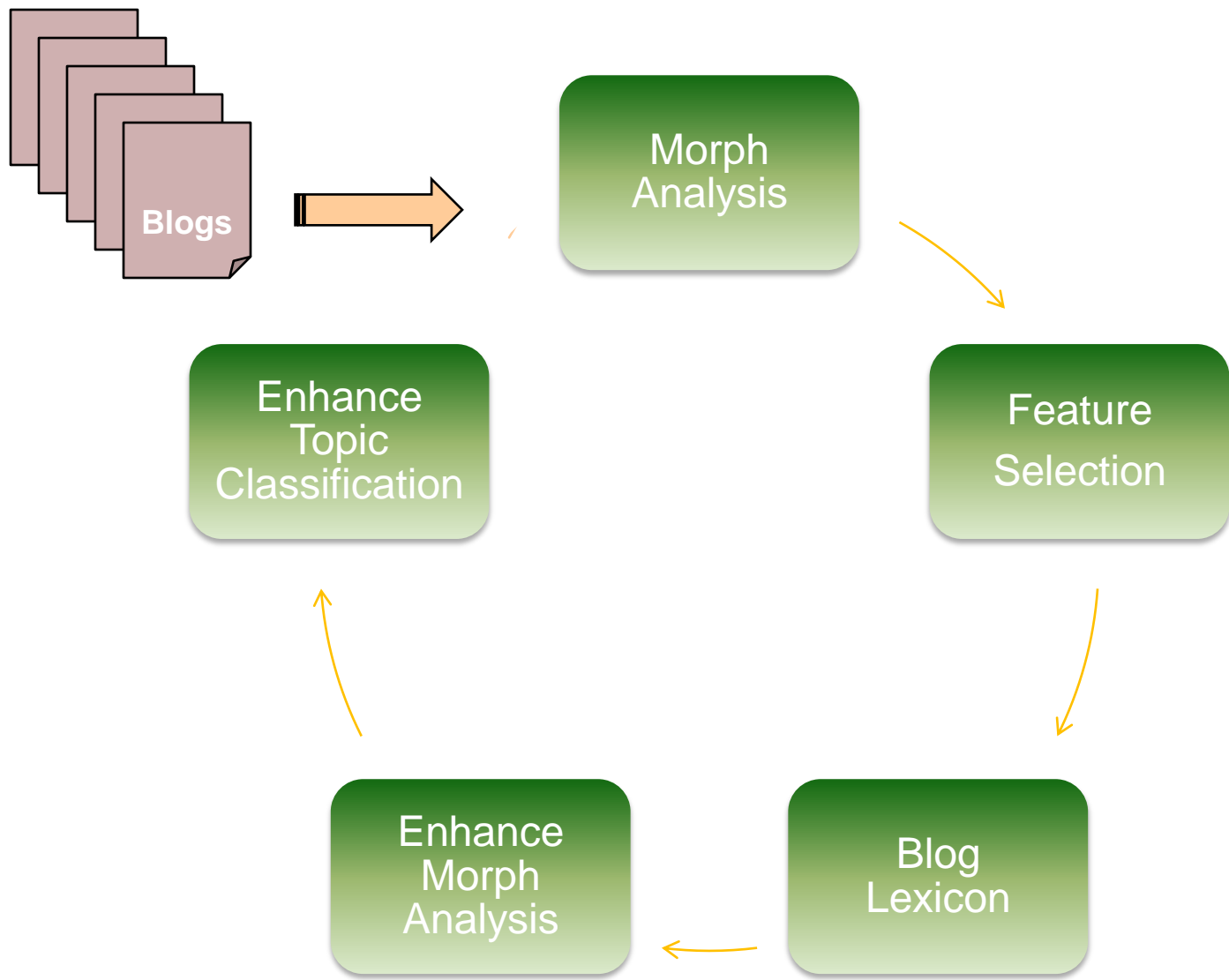
# Blog Lexicon

- Contains about 3,000 terms in 10 distinct categories
- Includes the Persian script, transcription, English translation and linguistic category for each term (e.g., English loan, abbreviation)

Topic	Number	Topic	Number
Technical	1840	Science	89
General	332	Medical	58
Politics	257	Film	30
Nuclear	150	Slang	31
Military	116	Religion	26

# Blog Lexicon Examples

Persian	Transcription	Formation	Translation	Domain
مرورگر	<i>morurgar</i>	browse + -er	Browser	Computer
دگر باش	<i>degarbâsh</i>	other + being	Queer	Politics/Society
بالگرد	<i>bâlgard</i>	wing + turn	Helicopter	Military
نمایشگر	<i>namâyeshgar</i>	show + -er	Monitor	Computer
مونیتور	<i>monitor</i>	(English loan)	Monitor	Computer
آوردنمایی	<i>âvardnamâyi</i>	input + showing	Battlefield Visualization	Military
فتنه گر	<i>fetnegar</i>	sedition + -er	Seditious	Politics
وبلاگستان	<i>veblâgestân</i>	weblog + -stan	Blogosphere	Computer
پیشدستانه	<i>spishdastâne</i>	pre + handed + ly	Preemptive	Military
چتیدن	<i>chatidan</i>	chat + (verb)	To chat	Computer
اس ام اس زدن	<i>es-em-es zadan</i>	SMS + hit	To send a text message	Computer
ذو الحقوق	<i>zolhoquq</i>	(Arab loan)	Righteous	Politics



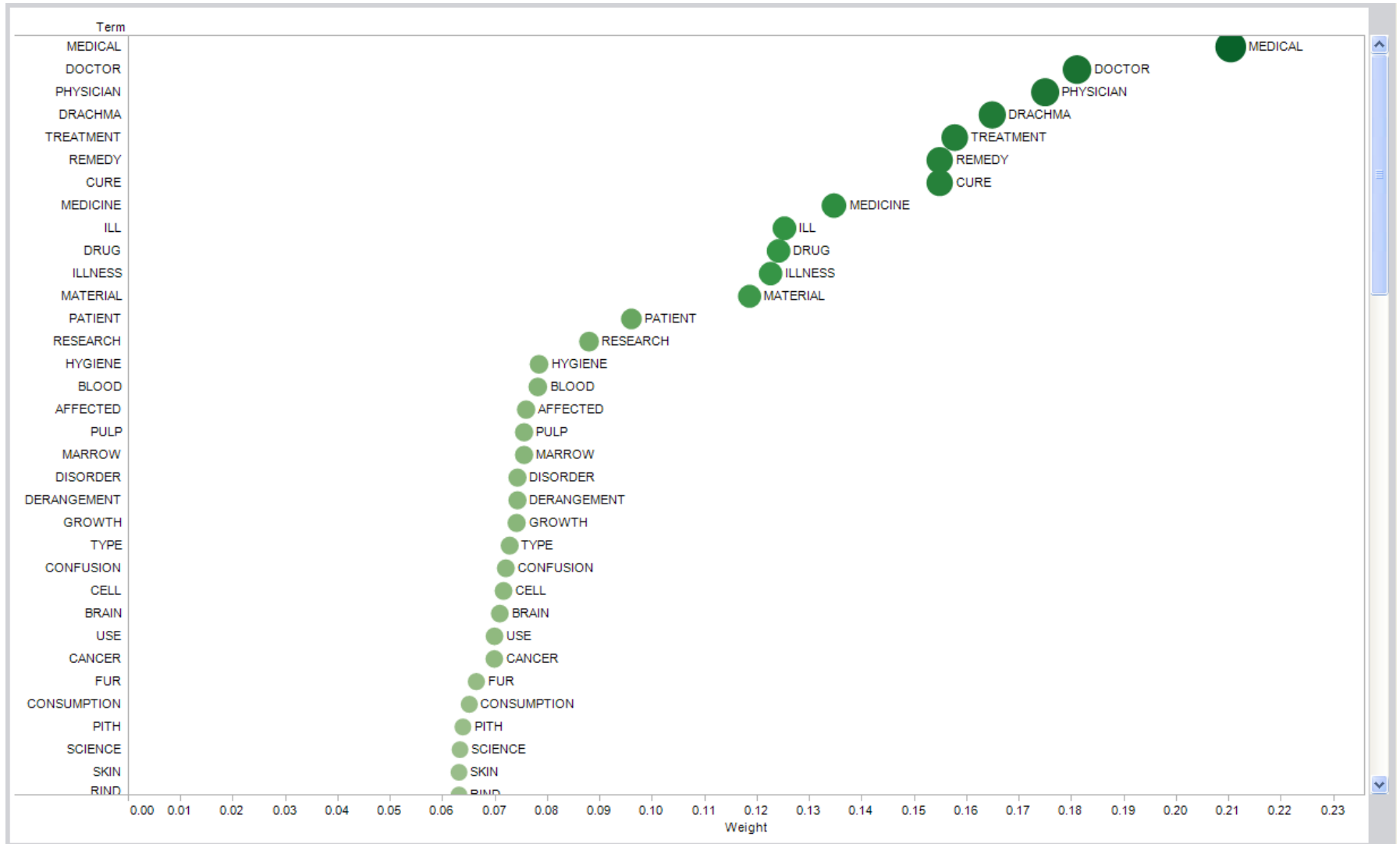
# Weighted Terms per Category

1	Medical	11	Material
2	Doctor	12	Patient
3	Physician	13	Research
4	Treatment	14	Hygiene
5	Remedy	15	Blood
6	Cure	16	Affected
7	Medicine	17	Marrow
8	Ill	18	Pulp
9	Drug	19	Disorder
10	Illness	20	Derangement

**Top 20 terms for Medicine**

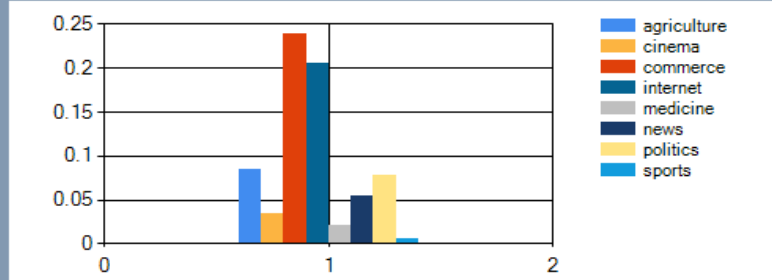
- Run blogs through enhanced morphological analyzer to obtain lemma or citation forms for each topic category
- Profile-based classifier trained on the terms (words) in the documents for each category
- Weighted terms used to identify the most significant key terms per topic
- Category of a blog post can then be determined by the similarity of its vector representation to the prototype vector for each topic

# Features for the Medical Domain



Choose File to Classify: C:\Projects\UIDelivery\outFiles\testBlogFA\test.txt.xml

Topic Match Results: commerce (0.238280972954)



commerce(0.238280972954)

دکتر الوانی در کتاب **مدیریت عمومی** خود نیز در ضمیمهٔ شش کتاب به تئوری SHE و تئوری مقابل آن یعنی HE اشاراتی نموده است. برخی از نظریات پروفسور پیتر دراگر (دروگر) نیز نشأت گرفته از همین کتاب است و چشمه هائی از این نظریه در کتاب '**مدیریت در جامعهٔ آینده**' دراگر فقیه به چشم می خورد. دیگر کتاب های رابرتسون از طریق این لینک قابل دسترس است. برخی دیگر از کتاب های او را می توانید به رایگان دانلود نمایید.

+ نوشته شده در 23/9/2008 ساعت 17:4 توسط امیر | 5 نظر

سخنی با خوانندگان این وبلاگ

در این پست از وبلاگ قصد دارم به برخی از سؤالاتی که در بخش نظرات این وبلاگ و یا توسط ایمیل با من در میان گذاشته می شود پاسخگو باشم. اغلب سؤالات و نظرات راجع به موضوع کارشناسی **ارشد** است و از آنجائی که فرصت کافی برای پاسخگویی به تک تک شما عزیزان را ندارم نظراتن را به خواندن این پست از بلاگ جلب می کنم. سوال ۲. من از چه راهی می توانم جزوات موسساتی را که کلاسهای کنکور برگزار می کنند را داشته باشم؟ فرستندگان: حامد ، سارا ، سمیه جواب: از طریق موسسات منتشر کننده جزوات اقدام نمایید. این جزوات تحت قانون کپی رایت هستند و نسخه برداری از آنها بدون اجازه موسسات منتشرکننده غیرقانونی است. **فروش آنها** نیز به همین شکل سوال ۳. سلام ، از مطالبتون بسیار متشکرم بسیار مطالب خوب و جامعی رو اشاره کردید اگر ممکنه سایتی رو معرفی کنید که بتونیم از نمونه سوالات مربوطه به رشته **مدیریت** استفاده کنیم با تشکر فراوان فرستنده: دریا

من سایت به خصوصی در این رابطه سراخ ندارم و در جستجو روی گوگل نیز چنین سایتی را نیافتم. می توانید از کتاب های موجود در **بازار** که همراه با جواب ارائه شده اند استفاده کنید.

سوال ۴. سلام از راهنمایی تون ممنون میخوامستم بدونم استفاده صرف از جزوات پ.ا.رسه تا چه حد تاثیر داره؟ فرستنده: نازنین به اندازهٔ اینکه درس را متوجه شوید و لا **هدف** 'برخی' اساتید موسسات( اسم نمی برم ) و جزواتشان بیشتر در راستای آموختن به دانشجو است و نه تئوری در **ارشد** ، البته با این نگرششون کاملا موافقم چرا که دانشجویی که تنها درس را بالاخص دروس مدیریتی را برای قبولی **ارشد** بخواند همان بهتر که **ارشد** قبول نشود. البته این نوع نگرش در بین اساتیدی که دروس **مدیریت** و گاه مهندسی صنایع درس می دهند بیشتر به چشم می خورد. سوال ۵. سلام اگر برائون مقدور **منابع** کارشناسی **ارشد** مدیریت صنعتی رو برام میل کنید ودر ضمن بگید که مینونم این کتابارو از طریق اینترنت خرید کنم. فرستنده: مهدیه ، فرناز **منابع** رشته **صنعتی** در همین وبلاگ و پست های مرتبط با رشته **مدیریت** صریحا بیان شده. وبلاگ را کامل بخوانید تا اطلاعات خود را بدست آورید. می توانید با کلیک روی هر کتاب که از این وبلاگ لینکی برای آنها ایجاد شده نیاز خود را با خرید از اینترنت بر طرف نمایید. سوال ۶. همیشه **منابع** فلان رشته رو معرفی کنید. فرستنده: هر کی میاد تو این وبلاگ این درخواست رو داره

جواب: من قادر به معرفی تمامی **منابع** نیستم ، اگر بخواهم چنین کنم بابستی ساعت ها و روزها به دنبال آنها باشم. بهترین **روش** برای یافتن **منابع** این است که از سه یا چهار استاد به نام در رشته ای که قصد **شرکت** در کنکور **ارشد** آن را دارید سوال **اساتید** مشتراکاً تکید نمودند رو حتما بخوانید. چه بهتر که از اساتیدی کمک بگیرید که با طراحان سوال ارتباط دارند. سوال ۷. سلام

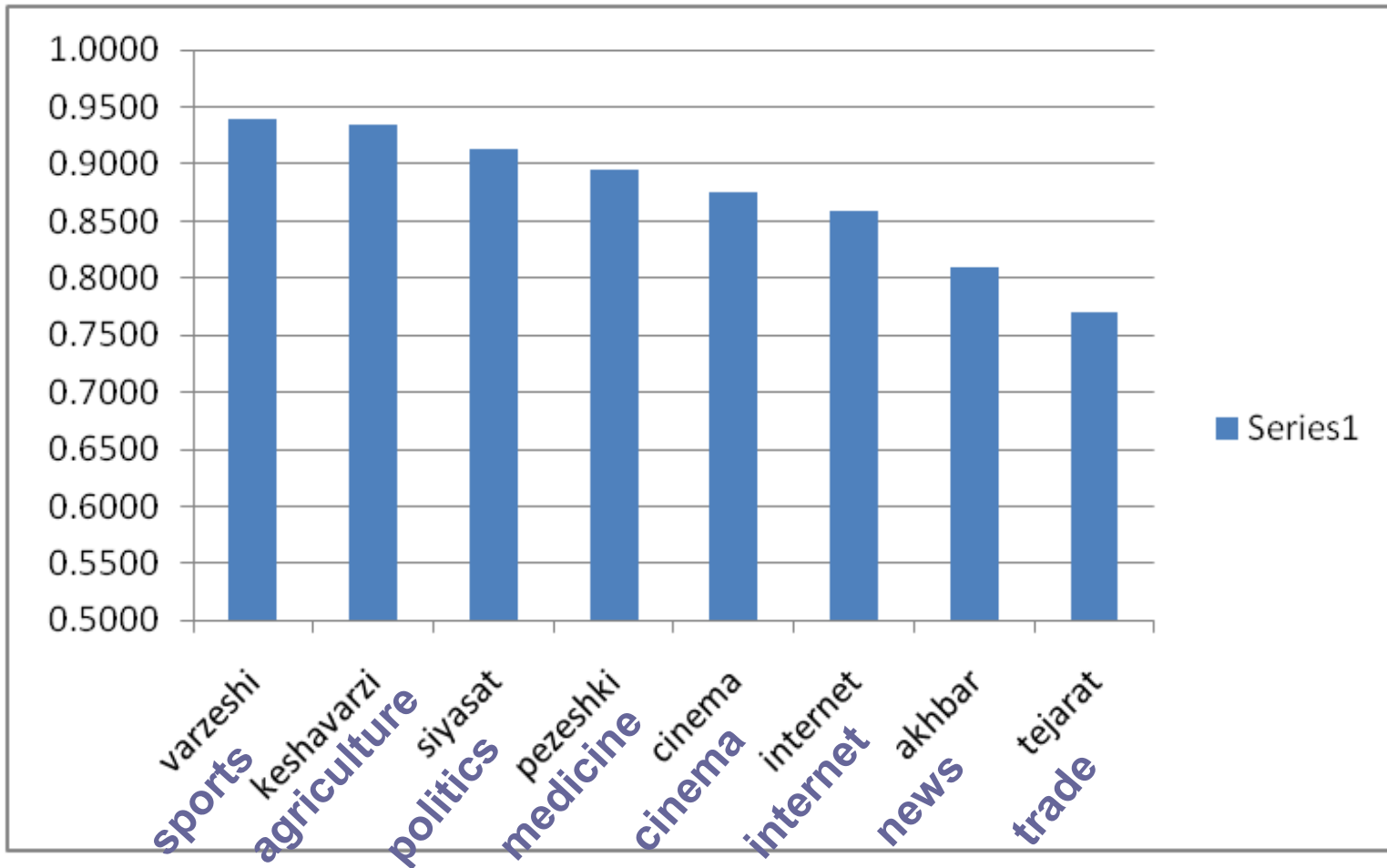
از وبلاگ خوبیت ممنون خیلی کمک خوبیه

به خواهش اژت داشتم و اون اینکه من می خوام بین **مدیریت** بازرگانی و مدیریت صنعتی یکی رو انتخاب کنم (برای کنکور **ارشد**) می تونی راهنمایی کنی؟ فرستنده: صادق جواب: در انتخاب گرایش به علاقه های خودتان رجوع کنید. در هر رشته ای که هستید شما هر گزایشی هم که بخوانید و حتی اگر دکتری هم داشته باشید تا اطلاعات نداشته باشید **شغلی** مناسبی نخواهید یافت مگر آنکه مثل بعضی ها از طریق لینک هائی به جانی برسید. چه بهتر که به آنچه علاقه دارید بپردازید ، انسان با پرداختن به علاقه هایش به رضایت از زندگی دست می یابد. 'سوالات مشابه در این مورد زیاد بوده' سوال ۸. تصمیم دارم امسال در کنکور کارشناسی **ارشد** مدیریت اجرائی **شرکت** کنم. خواهش میکنم اگر ممکنه برای برنامه ریزی درسی و **مطالعه** نوی این چند ماهی که فرصت دارم من رو راهنمایی کنید. باز هم سیاستگذارم فرستنده: گیتنا

جواب: حداقل امکان از کتابخانه برای درس خواندن استفاده کنید البته نه کتابخانه ای که پر از بچه های کنکوری باشد. رمز یادگیری در دوره هست این موضوع را هرگز فراموش نکنید. دوره های خود را زمان بندی کنید. اگر از کلاس استفاده می کنید بعد از کلاس حتما درس را مرور کنید(اینو دیگه شگهی می دونید). کنکور **ارشد** در رشته **مدیریت** نیاز به اطلاعات وسیعی دارد لذا حتی خواندن **مقالات** به روز هم کمک کننده خواهد بود مخصوصا رشته اجرائی که فقط با اطلاعات به روز سروکار دارد. سخن آخر: دکتر استن کای به عنوان یکی از مشاورین **ارشد شرکت های** بزرگ آمریکائی و به عنوان یکی از اعجازگران گران عرصهٔ **مدیریت**(Management guru) روی این موضوع بسیار تاکید دارد که افراد خوششان باید به خود بیاموزند لذا دوست عزیز این فقط خود تو هستی که یاد می گیری چگونه یاد بگیری ، استاد ، کتاب ، جزوه ، **موسسه** و .... فقط ببیند درصدم تاثیرگذار است. روی سخن بیشتر با آنهاست که می گویند شما نترانید به خوش به حالتان ، جزوه دارید ، کلاس دارید.

با بهترین آرزوها - امیر

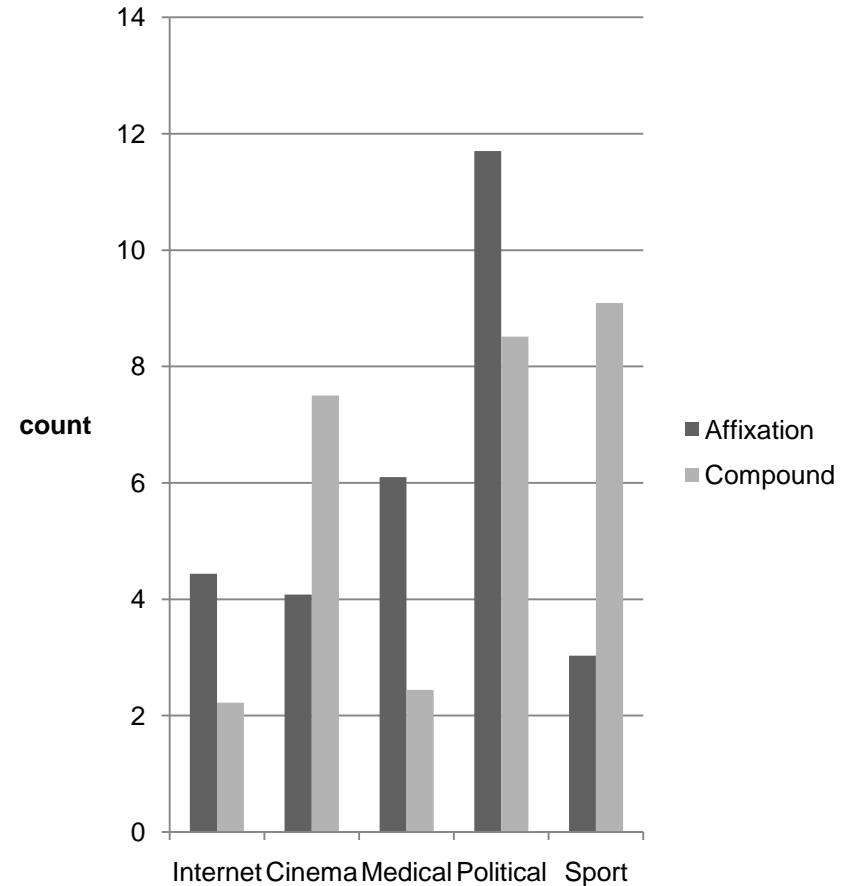
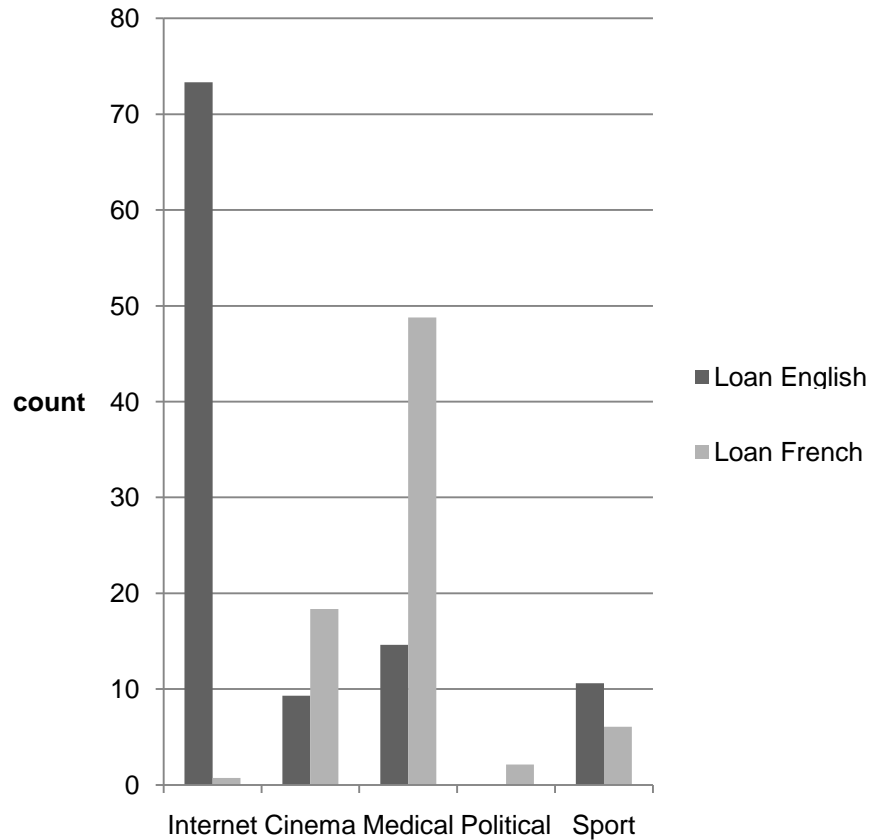
# Topic Classification



# Language Correlations

- **Weighted features selected by classifier correlates with topic**
- **Type of neologism used correlates with topic**
  - Internet and computer blogs have lots of English loans
  - Political blogs create words based on Persian rules
- **Type of neologism used correlates with author**
  - Younger bloggers use English loans rather than use the official technical terms
  - It seems that younger bloggers use conversational language more freely

# Language Correlations



# Link Structure Analysis

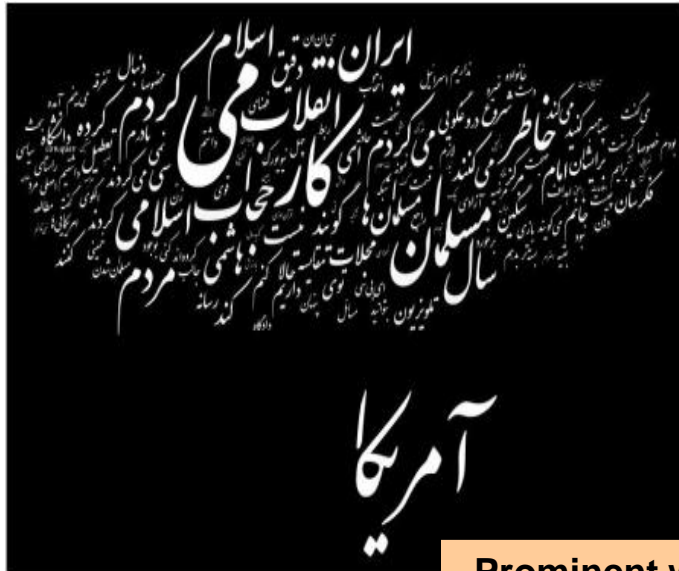
- **Background**

- Investigating how blogs are connected through hyperlinks
- Blogs that link to each other share issues and ideologies

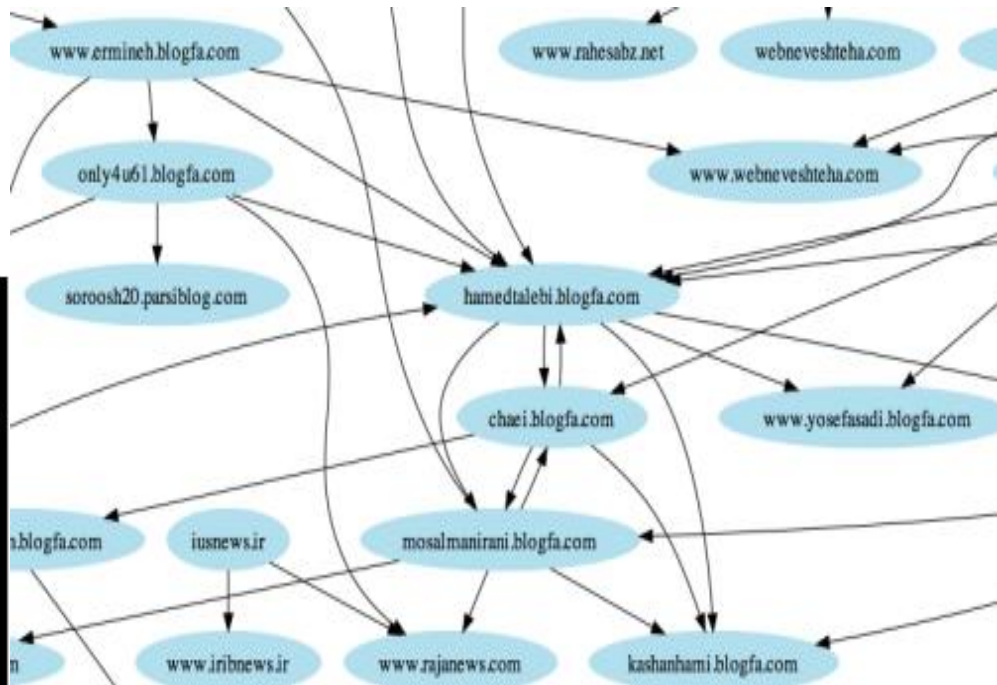
- **Automatically extracted links from blogs**

- **Combined with the results of classification to identify virtual clusters based on topic**

# Politics Topic Cluster



Prominent words for political cluster



Link Analysis of political blogs

*America, work, muslim, revolution, year, justice, sedition, sacrifice, Iran, hijab, people, soft war, trial, Islamic, magazines, Hashemi, comparison, Imam...*

# Conclusion

- **Built a Blog Lexicon**
- **Extended a morphological analyzer to process word forms encountered in blog language**
- **Developed a profile-based topic classifier for Persian blogs**
- **We combine the topic classifier with the results of link structure analysis to detect virtual topic clusters**
- **Future Work**
  - **Integrate results of sentiment analysis to identify bloggers' opinion(s) within topic clusters**
  - **Temporal analysis of topics in blog clusters**