

Statistical Methods for Integration and Analysis of Online Opinionated Text Data

ChengXiang (“Cheng”) Zhai

**Department of Computer Science
University of Illinois at Urbana-Champaign**

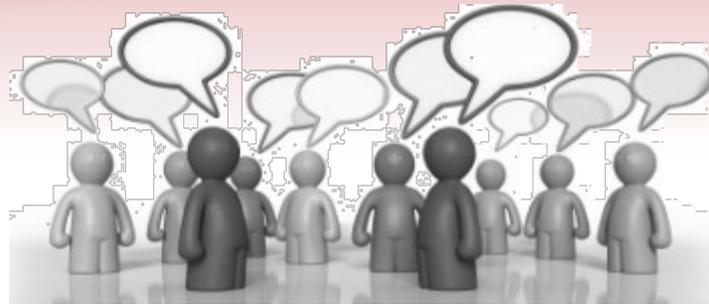
<http://www.cs.uiuc.edu/homes/czhai>

Joint work with Yue Lu, Qiaozhu Mei, Kavita Ganesan, Hongning Wang, and others

Online opinions cover all kinds of topics

Topics:

- People
- Events
- Products
- Services, ...



Sources:

- Blogs
- Microblogs
- Forums
- Reviews, ...

45M reviews ↑ 53M blogs ↑ 65M msgs/day ↑ 115M users ↑
 1307M posts ↑ 10M groups ↑



Great opportunities for many applications

Opinionated Text Data



Topics:
People
Events
Products
Services, ...



Sources:
Blogs
Microblogs
Forums
Reviews, ...



Decision Making & Analytics

“Which cell phone should I **buy**?”

“What are the winning features of **iPhone** over **blackberry**?”

“How do people like this new **drug**?”

“How is Obama’s health care **policy** received?”

“Which presidential candidate should I **vote** for?”

...

However, it's not easy to for users to make use of the online opinions

How can I collect all opinions?

How can I digest them all?

How can I ...?

How can I ...?



Research Questions

- How can we integrate scattered opinions?
- How can we summarize opinionated text articles?
- How can we analyze online opinions to discover patterns and understand consumer preferences?
- How can we do all these in a general way with no or minimum human effort?
 - **Must work for all topics**
 - **Must work for different natural languages**

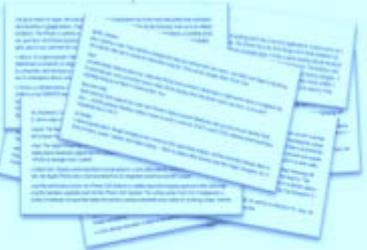
Solutions: Knowledge-Learn Statistical Methods (Statistical Language Models)

Lots of related work (usually not as general):

Bing Liu, *Sentiment Analysis and Opinion Mining*, Morgan & Claypool Publishers, 2012

Rest of the talk: general methods for

Opinionated Text Data



Topics: People, Events, Products, Services, ...



Sources: Blogs, Microblogs, Forums, Reviews, ...



1. Opinion Integration



2. Opinion Summarization

Query: *Dell Laptop*

	positive	negative	neutral
Topic 1 (Price)	- it is the best one and they show Dell coupon code as early as possible	- Even though Dell's price is cheaper, we still don't want it.	- that got us one previous price compare. - (DELL is trading at \$24.88)
Topic 2 (Battery)	- One thing I really like about this Dell battery is the express charge feature.	- The Dell battery sucks. - stupid old laptop battery	- I still want a free battery from dell.

3. Opinion Analysis



Decision Making & Analytics

"Which cell phone should I buy?"

"What are the winning features of iPhone over blackberry?"

"How do people like this new drug?"

"How is Obama's health care policy received?"

"Which presidential candidate should I vote for?"

...

Outline

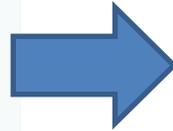
Opinionated Text Data



Topics: People, Events, Products, Services, ...



Sources: Blogs, Microblogs, Forums, Reviews, ...



1. Opinion Integration



2. Opinion Summarization

Query: *Dell Laptop*

	positive	negative	neutral
Topic 1 (Price)	- it is the best one and they show Dell coupon code as early as possible	- Even though Dell's price is cheaper, we still don't want it.	- most people are picky about a price compare. - DELL is trading at \$24.88
Topic 2 (Battery)	- One thing I really like about this Dell battery is the Express Charge feature.	- The Dell battery sucks. - stupid old laptop battery	- I still want a free battery from dell.



3. Opinion Analysis

Decision Making & Analytics

"Which cell phone should I buy?"

"What are the winning features of iPhone over blackberry?"

"How do people like this new drug?"

"How is Obama's health care policy received?"

"Which presidential candidate should I vote for?"

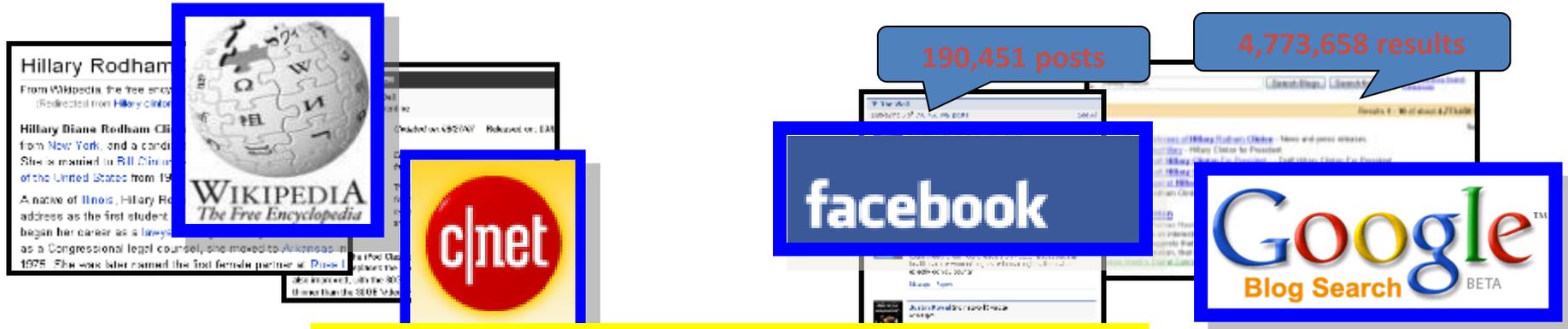
...

How to digest all scattered opinions?

Need tools to automatically integrate all scattered opinions

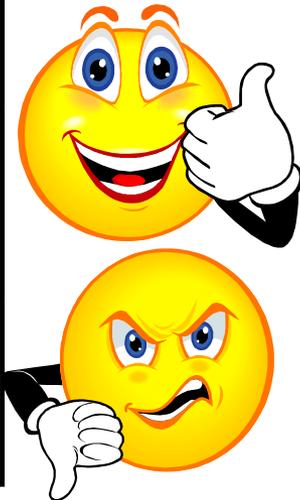


Observation: two kinds of opinions



Can we combine them?

Expert opinions
<ul style="list-style-type: none"> • CNET editor's review • Wikipedia article
<ul style="list-style-type: none"> • Well-structured • Easy to access
<ul style="list-style-type: none"> • Maybe biased • Outdated soon



Ordinary opinions
<ul style="list-style-type: none"> • Forum discussions • Blog articles
<ul style="list-style-type: none"> • Represent the majority • Up to date
<ul style="list-style-type: none"> • Hard to access • fragmented

Opinion Integration Strategy 1

[Lu & Zhai WWW 08]

Align scattered opinions with well-structured
expert reviews

Yue Lu, ChengXiang Zhai. Opinion Integration Through Semi-supervised Topic Modeling, *Proceedings of the World Wide Conference 2008 (WWW'08)*, pages 121-130.

Review-Based Opinion Integration

Input

Topic: iPod



Expert review with aspects

Design
Battery
Price..

Text collection of ordinary opinions, e.g. Weblogs



Output

Review Aspects
Extra Aspects

Design
Battery
Price

iTunes	... easy to use...
warranty	...better to extend..

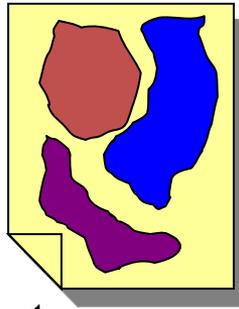
Similar opinions	Supplementary opinions
cute... tiny...	..thicker..
last many hrs	die out soon
could afford it	still expensive

Integrated Summary

Solution is based on probabilistic latent semantic analysis (PLSA) [Hofmann 99]

Topic model
 = unigram language model
 = multinomial distribution

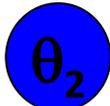
Document



Topics

battery 0.3
life 0.2..

design 0.1
screen 0.05



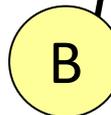
...



price 0.2
purchase 0.15

$1 - \lambda_B$

λ_B



is 0.05
the 0.04
a 0.03 ..

Collection background

Generate a word in a document

$$p_d(w) = \lambda_B p(w|\theta_B) + (1 - \lambda_B) \sum_{j=1}^k [\pi_{d,j} p(w|\theta_j)]$$

Basic PLSA: Estimation

Generate a word
in a document

$$p_d(w) = \lambda_B p(w|\theta_B) + (1 - \lambda_B) \sum_{j=1}^k [\pi_{d,j} p(w|\theta_j)]$$

Log-likelihood of
the collection

Count of word in
the document

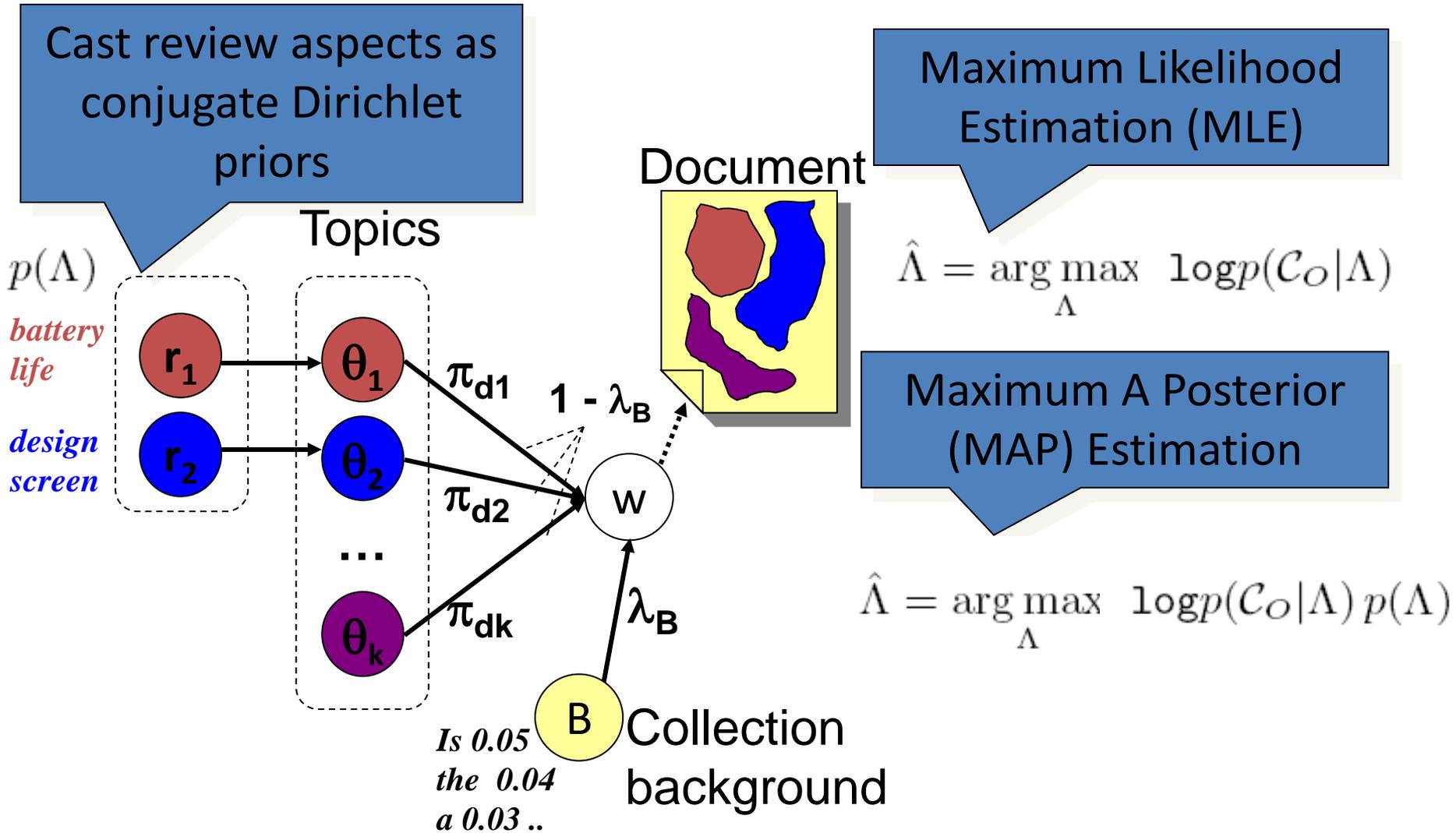
$$\log p(\mathcal{C}_O|\Lambda) = \sum_{d \in \mathcal{C}_O} \sum_{w \in V} \{c(w, d) \times \log p_d(w)\}$$

$$\boxed{p(w|\theta_j)} \quad \boxed{\pi_{d,j}}$$

- Parameters estimated with Maximum Likelihood Estimator (MLE) through an EM algorithm

$$\hat{\Lambda} = \arg \max_{\Lambda} \log p(\mathcal{C}_O|\Lambda)$$

Semi-supervised Probabilistic Latent Semantic Analysis



Results: Product (iPhone)

- Opinion Integration with review aspects

Review article	Similar opinions	Supplementary opinions
<p>You can make emergency calls, but you can't use any other functions...</p> <p>Activation</p>	<p>N/A</p> <p>Confirm the opinions from the review</p>	<p>... methods for unlocking the iPhone... on the iPhone for a few weeks, although they involve tinkering with the iPhone hardware...</p> <p>Unlock/hack iPhone</p>
<p>rated battery life of 8 hours talk time, 24 hours of music playback, 7 hours of video playback, and 6 hours of Internet use.</p> <p>Battery</p>	<p>Up to 8 Hours of Talk Time, 6 Hours of Internet Use, 7 Hours of Video Playback or 24 Hours of Audio Playback</p>	<p>Playing relatively high bitrate VGA H.264 videos, our iPhone lasted almost exactly 9 freaking hours of continuous playback with cell and WiFi on (but Bluetooth off).</p> <p>Additional info under real usage</p>

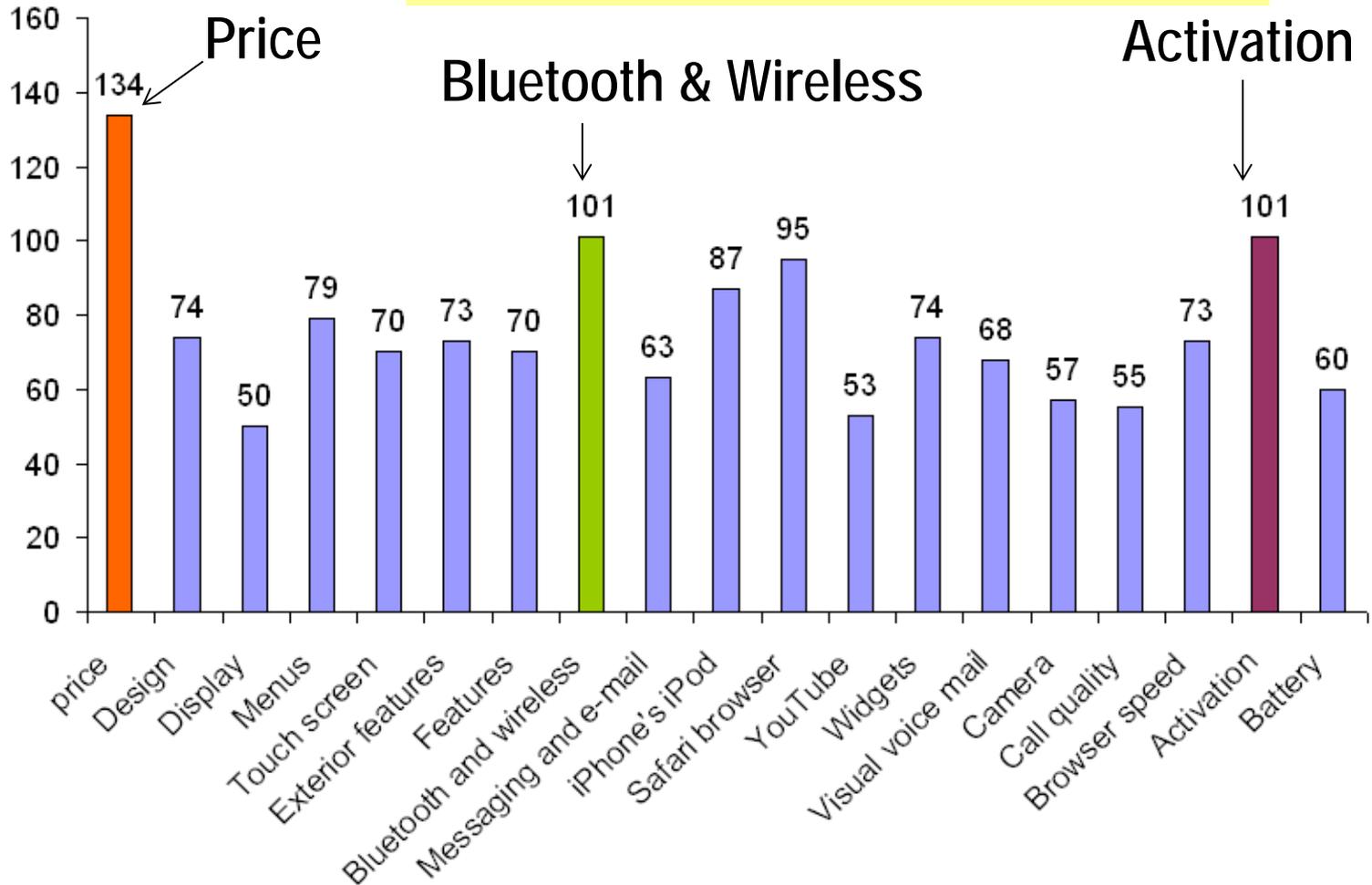
Results: Product (iPhone)

- Opinions on extra aspects

support	Supplementary opinions on extra aspects
15	<p>You may have heard of iASign , an iPhone app that allows you to activate your phone with iTunes rigamarole.</p> <p>Another way to activate iPhone</p>
13	<p>Cisco has owned the trademark on the name "iPhone" since 2000, when it acquired InfoGear, which originally registered the name.</p> <p>iPhone trademark originally owned by Cisco</p>
13	<p>With the imminent availability of the iPhone, a look at 10 things current smartphones like the Nokia N95 have been able to do that the iPhone can't currently match...</p> <p>A better choice for smart phones?</p>

As a result of integration...

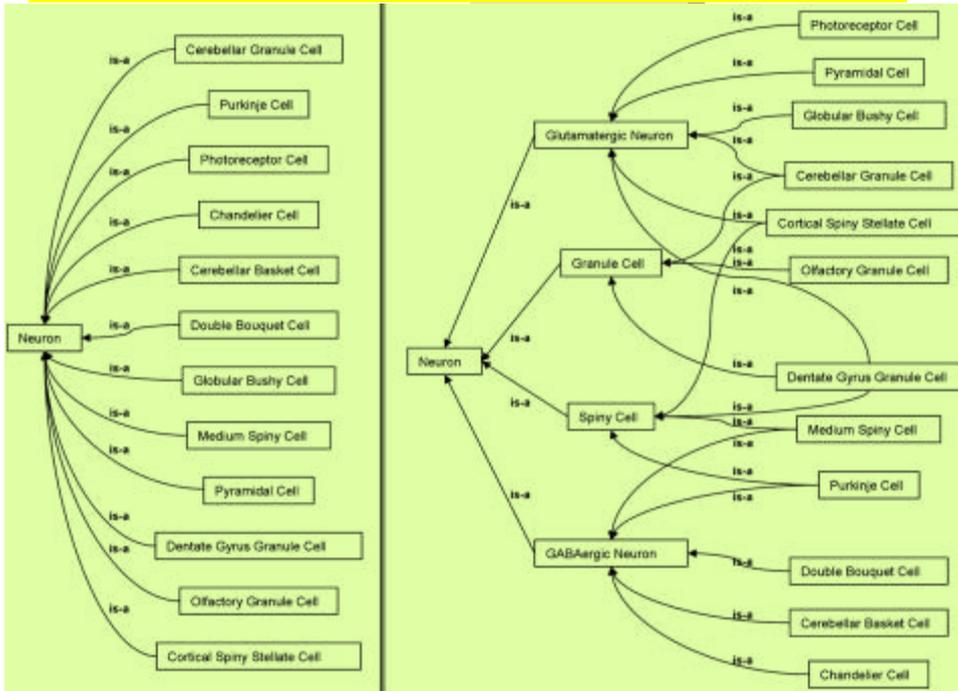
What matters most to people?



What if we don't have expert reviews?

How can we organize scattered opinions?

Exploit online ontology!



Ordinary opinions

- Forum discussions
- Blog articles
- Represent the majority
- Up to date
- Hard to access
- fragmented

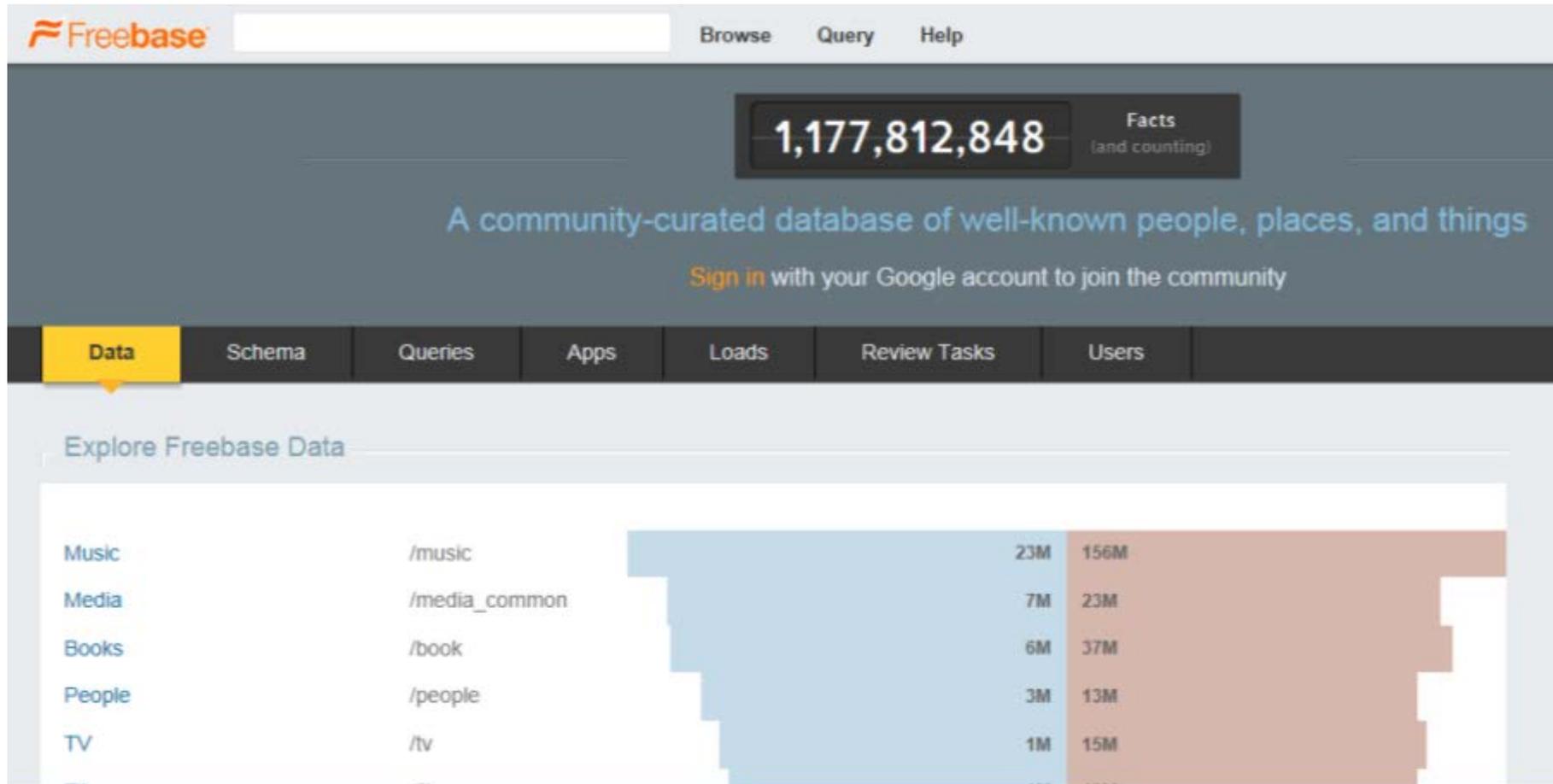
Opinion Integration Strategy 2

[Lu et al. COLING 10]

Organize scattered opinions using an ontology

Yue Lu, Huizhong Duan, Hongning Wang and ChengXiang Zhai. Exploiting Structured Ontology to Organize Scattered Online Opinions, *Proceedings of COLING 2010 (COLING 10)*, pages 734-742.

Sample Ontology: Freebase™



The screenshot shows the Freebase website interface. At the top left is the Freebase logo. To its right is a search bar and navigation links for "Browse", "Query", and "Help". Below this is a large dark grey banner with a white box containing the number "1,177,812,848" and the text "Facts (and counting)". Below the banner is the text "A community-curated database of well-known people, places, and things" and a "Sign in" link. A navigation bar below the banner contains tabs for "Data", "Schema", "Queries", "Apps", "Loads", "Review Tasks", and "Users". The "Data" tab is selected. Below the navigation bar is a section titled "Explore Freebase Data" containing a table with two columns: "Category" and "Count". The table lists categories like Music, Media, Books, People, and TV with their respective counts.

Category	URI	Count	Count
Music	/music	23M	156M
Media	/media_common	7M	23M
Books	/book	6M	37M
People	/people	3M	13M
TV	/tv	1M	15M

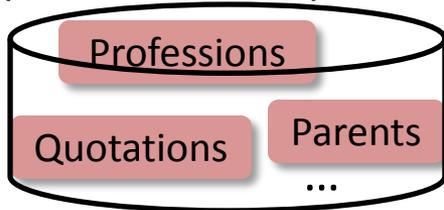
Ontology-Based Opinion Integration

Two key tasks: 1. Aspect Selection. 2. Aspect Ordering

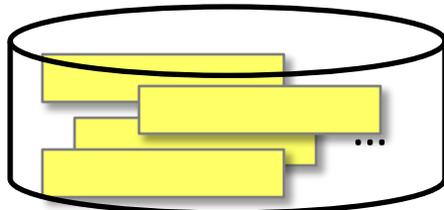
Topic = "Abraham Lincoln"
(Exists in ontology)



Aspects from Ontology
(more than 50)

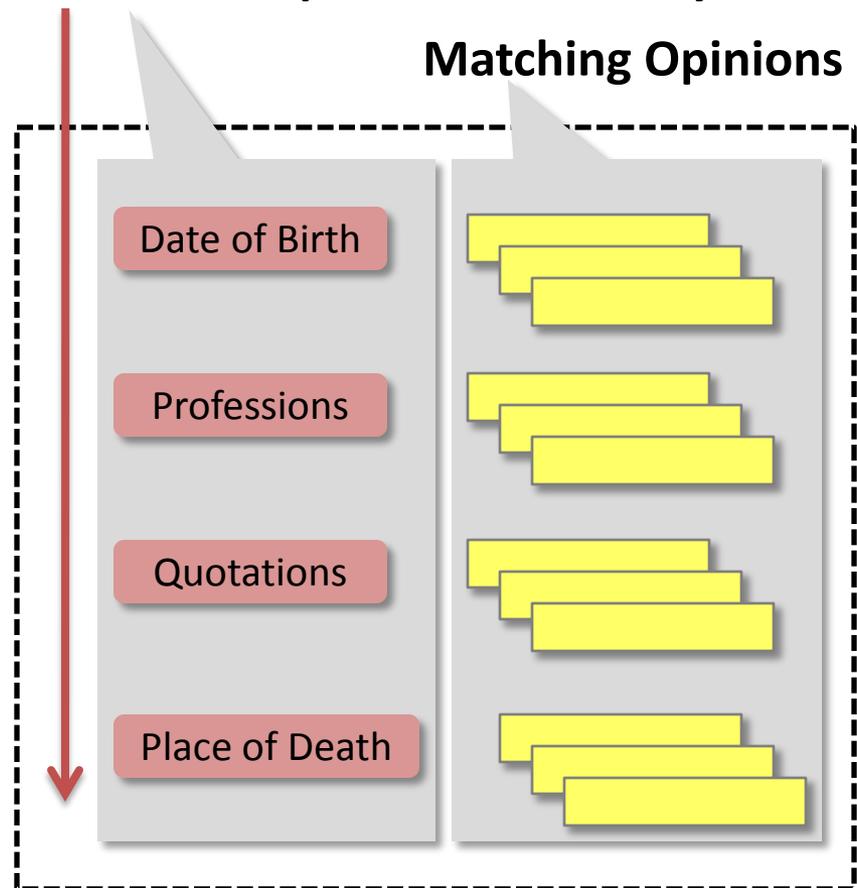


Online Opinion Sentences



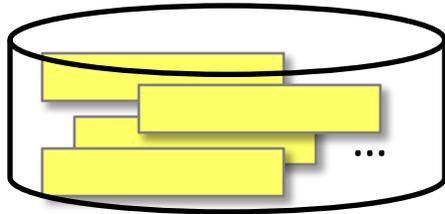
Subset of Aspects
Ordered to optimize readability

Matching Opinions

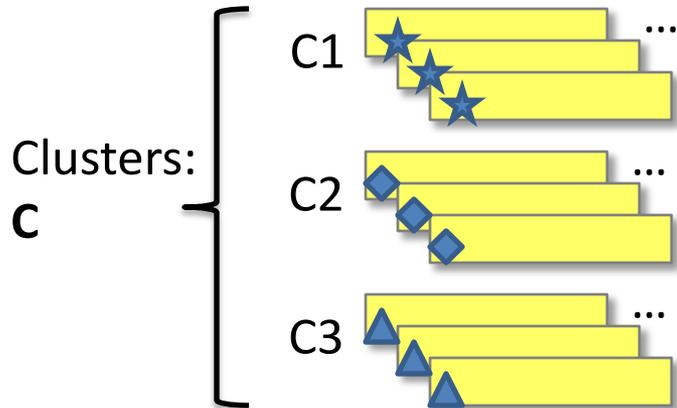


1. Aspect Selection: Conditional Entropy-based Method

Collection:

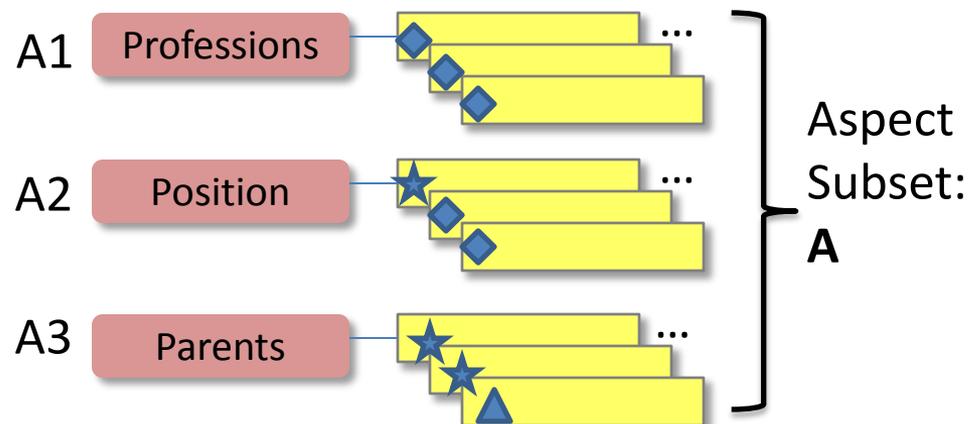


K-means Clustering

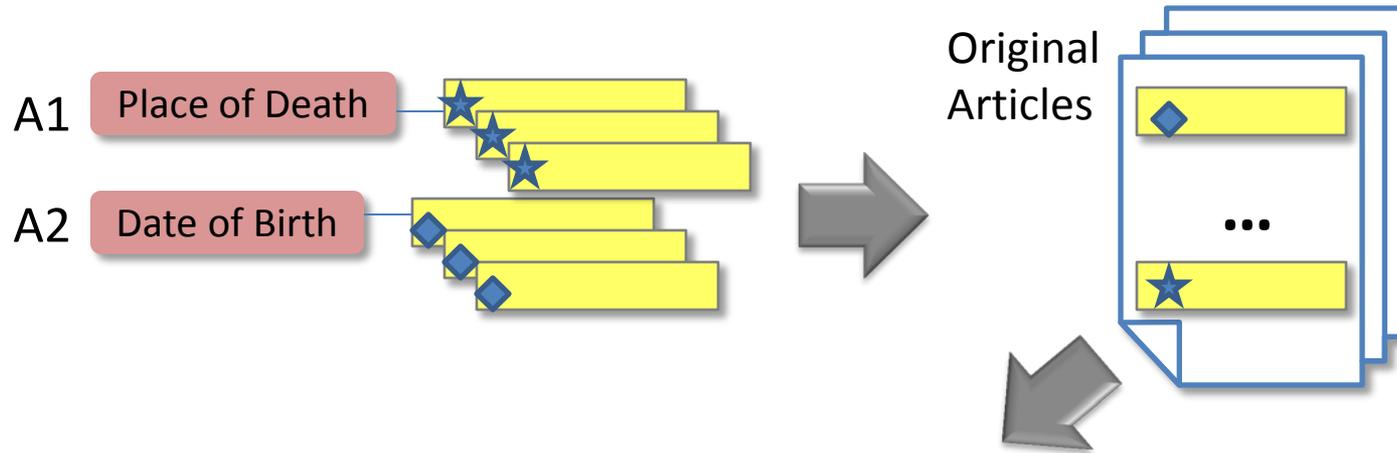


$$A = \operatorname{argmin} H(C|A)$$

$$= \operatorname{argmin} - \sum_i p(A_i, C_i) \log \frac{p(A_i, C_i)}{p(A_i)}$$



2. Aspect Ordering: Coherence Order



$\text{Coherence}(A1, A2) \leftarrow \#(\star \text{ is before } \diamond)$

$\text{Coherence}(A2, A1) \leftarrow \#(\diamond \text{ is before } \star)$

So, $\text{Coherence}(A2, A1) > \text{Coherence}(A1, A2)$

$$\Pi(A) = \operatorname{argmax} \sum_{A_i \text{ before } A_j} \text{Coherence}(A_i, A_j)$$

Sample Results: Sony Cybershot DSC-W200

Freebase Aspects	sup	Representative Opinion Sentences
Format: Compact	13	Quality pictures in a compact package. ...amazing is that this is such a small and compact unit but packs so much power
Supported Storage Types: Memory Stick Duo	11	This camera can use Memory Stick Pro Duo up to 8 GB Using a universal storage card and cable (c'mon Sony)
Sensor type: CCD	10	I think the larger ccd makes a difference. but remember this is a small CCD in a compact point-and-shoot.
Digital zoom: 2X	47	once the digital :smart" zoom kicks in you get another 3x of zoom. I would like a higher optical zoom, the W200 does a great digital zoom translation...

More opinion integration results are
available at:

<http://sifaka.cs.uiuc.edu/~yuelu2/opinionintegration/>

Outline

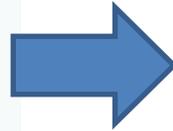
Opinionated Text Data



Topics: People, Events, Products, Services, ...



Sources: Blogs, Microblogs, Forums, Reviews, ...



1. Opinion Integration



2. Opinion Summarization

Query: *Dell Laptop*

	positive	negative	neutral
Topic 1 (Price)	- it is the best one and they show Dell coupon code as best as possible	- Even though Dell's price is cheaper, we still don't want it.	- most people like price point, a price compare. - DELL is trading at \$24.88
Topic 2 (Battery)	- One thing I really like about this Dell battery is the Express Charge feature.	- The Dell battery sucks - stupid old laptop battery	- I still want a free battery from dell.



3. Opinion Analysis

Decision Making & Analytics

"Which cell phone should I buy?"

"What are the winning features of iPhone over blackberry?"

"How do people like this new drug?"

"How is Obama's health care policy received?"

"Which presidential candidate should I vote for?"

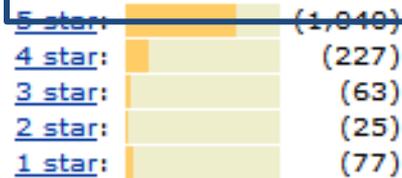
...

Need for opinion summarization

Customer Reviews

Average Customer Rating

★★★★☆ (1,432 customer reviews)



1,432 customer reviews

How can we help users digest these opinions?

Most Helpful Customer Review

3,677 of 3,770 people

★★★★★ **WARN**

By [Hassan B. Bn Hadhram](#) - [See all my reviews](#)

REAL NAME

Amazon Verified Purchase (What's this?)

This review is from: **Apple iPod touch 8 GB (2nd Generation--with iPhone OS 3.1 Software Installed) [NEWEST MODEL]** (Electronics)

Before i start let me just tell you "what's New" with the iPod touch Third generation" :

- Faster Cpu/Double the ram/Better graphic (faster Boot time/faster loading is all what i did notice)
- Double the storage for the same old price
- Voice control (I'll explain it in a second)
- Latest firmware for free

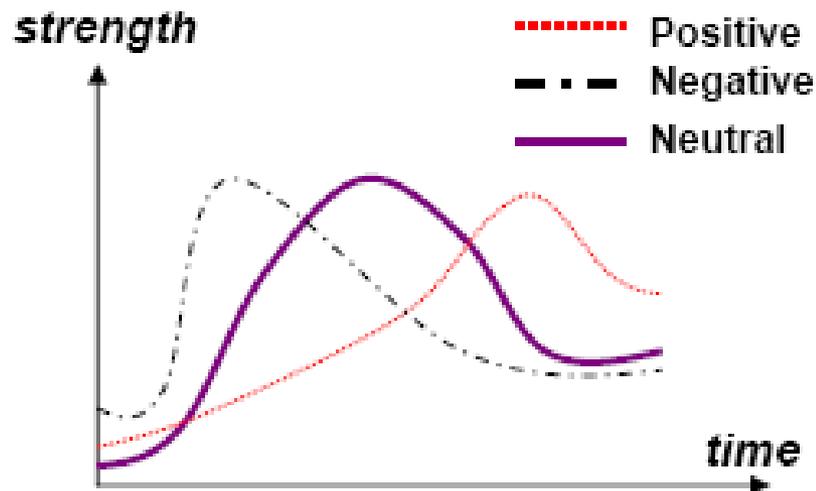
Nice to have....

Topic-sentiment summary

Query: *Dell Laptop*

	positive	negative	neutral
Topic 1 <i>(Price)</i> 	<ul style="list-style-type: none"> • It is the best site and they show Dell coupon code as early as possible 	<ul style="list-style-type: none"> • Even though Dell's price is cheaper, we still don't want it. • 	<ul style="list-style-type: none"> • mac pro vs. dell precision: a price compar... • DELL is trading at \$24.66
Topic 2 <i>(Battery)</i> 	<ul style="list-style-type: none"> • One thing I really like about this Dell battery is the Express Charge feature. 	<ul style="list-style-type: none"> • my Dell battery sucks • Stupid Dell laptop battery • 	<ul style="list-style-type: none"> • I still want a free battery from dell.. •

Topic-sentiment dynamics (Topic = *Price*)



Can we do this in a general way?

Opinion Summarization 1:

[Mei et al. WWW 07]

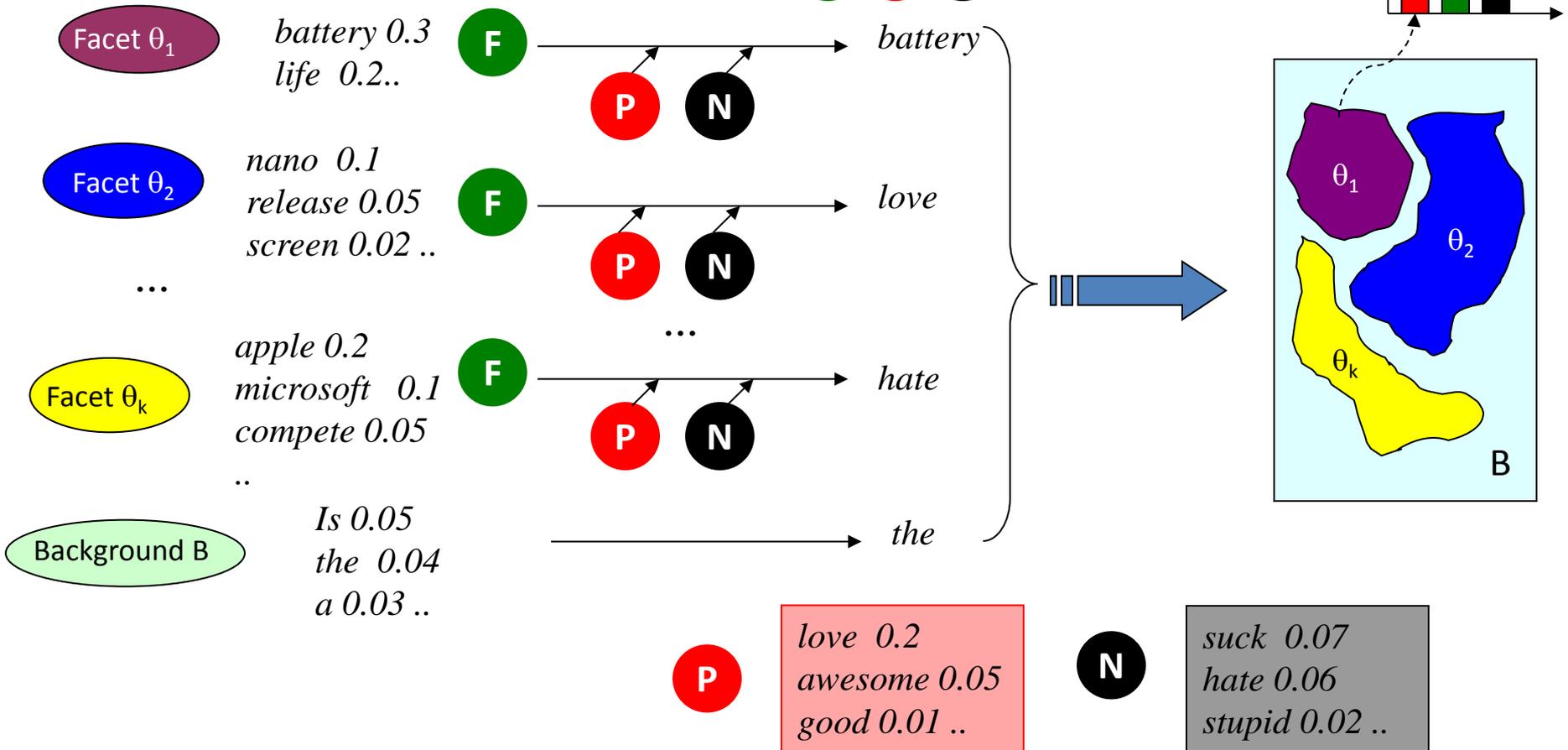
Multi-Aspect Topic Sentiment Summarization

Qiaozhu Mei, Xu Ling, Matthew Wondra, Hang Su, ChengXiang Zhai, Topic Sentiment Mixture: Modeling Facets and Opinions in Weblogs, *Proceedings of the World Wide Conference 2007 (WWW'07)*, pages 171-180

A Topic-Sentiment Mixture Model

Choose a facet (subtopic) θ_i

Draw a word from the mixture of topics and sentiments (**F** **P** **N**)



The Likelihood Function

$$\log p(C) = \sum_{d \in C} \sum_{w \in V} c(w, d) \log [\lambda_B p(w | B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{dj} (\delta_{j,d,F} p(w | \theta_j) + \delta_{j,d,P} p(w | \theta_P) + \delta_{j,d,N} p(w | \theta_N))]$$

Count of word w
in document d

Generating w
using the neutral topic model

Generating w
using the background model

Generating w
using the positive sentiment model

Generating w
using the negative sentiment model

Choosing
a faceted opinion

Two Modes for Parameter Estimation

- **Training Mode: Learn the sentiment model**

$$\log(C) = \sum_{d \in C} \sum_{w \in V} c(w, d) \log[\lambda_B p(w | B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{dj} \times (\delta_{j,d,F} p(w | \theta_j) - \delta_{j,d,P} p(w | \theta_P) + \delta_{j,d,N} p(w | \theta_N))]$$

Fixed for each d

One of them is zero for d

- **Testing Mode: Extract the Topic models**

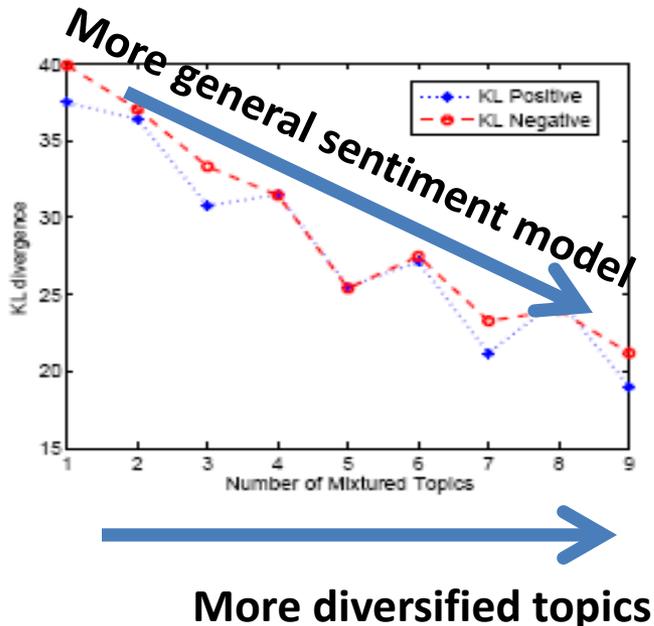
$$\log(C) = \sum_{d \in C} \sum_{w \in V} c(w, d) \log[\lambda_B p(w | B) + (1 - \lambda_B) \sum_{j=1}^k \pi_{dj} \times (\delta_{j,d,F} p(w | \theta_j) + \delta_{j,d,P} p(w | \theta_P) + \delta_{j,d,N} p(w | \theta_N))]$$

Feed strong prior on sentiment models

EM algorithm can be used for estimation

Results: General Sentiment Models

- Sentiment models trained from diversified topic mixture v.s. single topics



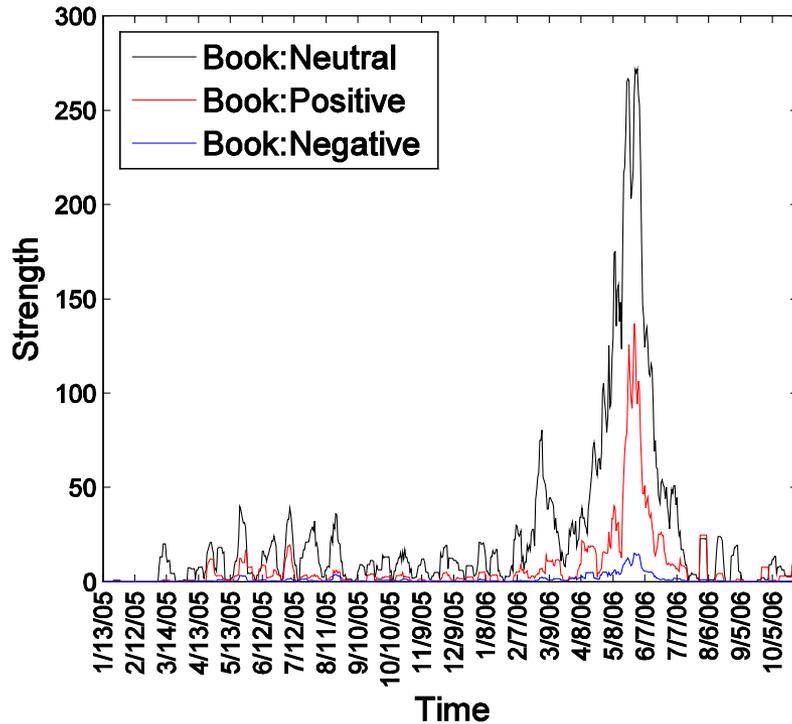
Pos-Mix	Neg-Mix	Pos-Cities	Neg-Cities
love	suck	beautiful	hate
awesome	hate	love	suck
good	stupid	awesome	people
miss	ass	amaze	traffic
amaze	fuck	live	drive
pretty	horrible	good	fuck
job	shitty	night	stink
god	crappy	nice	move
yeah	terrible	time	weather
bless	people	air	city
excellent	evil	greatest	transport

Multi-Faceted Sentiment Summary (query="Da Vinci Code")

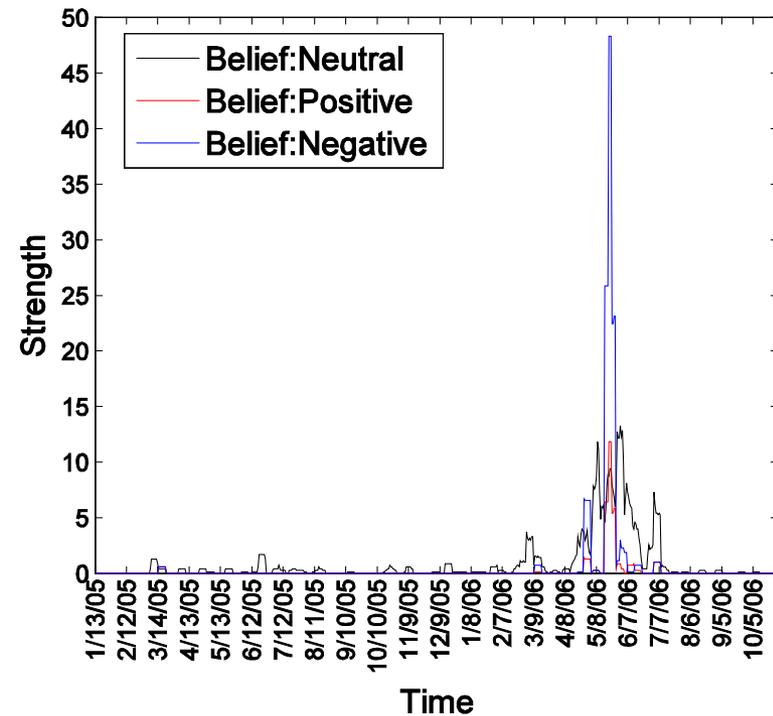
	Neutral	Positive	Negative
Facet 1: Movie	... Ron Howards selection of Tom Hanks to play Robert Langdon.	Tom Hanks stars in the movie,who can be mad at that?	But the movie might get delayed, and even killed off if he loses.
	Directed by: Ron Howard Writing credits: Akiva Goldsman ...	Tom Hanks, who is my favorite movie star act the leading role.	protesting ... will lose your faith by ... watching the movie.
	After watching the movie I went online and some research on ...	Anybody is interested in it?	... so sick of people making such a big deal about a FICTION book and movie.
Facet 2: Book	I remembered when i first read the book, I finished the book in two days.	Awesome book.	... so sick of people making such a big deal about a FICTION book and movie.
	I'm reading "Da Vinci Code" now. ...	So still a good book to past time.	This controversy book cause lots conflict in west society.

Separate Theme Sentiment Dynamics

“book”



“religious beliefs”



Can we make the summary more concise?

	Neutral	Positive	Negative
Facet 1: Movie	... Ron Howards selection of Tom Hanks to play Robert Langdon.	Tom Hanks stars in the movie,who can be mad at that?	But the movie might get delayed, and even killed off if he loses.
	Directed by: Ron Howard Writing credits: Akiva Goldsman ...	Tom Hanks, who is my favorite movie star act the leading role.	protesting ... will lose your faith by ... watching the movie.
What if the user is using a smart phone?			
	research on ...		FICTION book and movie .
Facet 2: Book	I remembered when i first read the book, I finished the book in two days.	Awesome book.	... so sick of people making such a big deal about a FICTION book and movie .
	I'm reading "Da Vinci Code" now. ...	So still a good book to past time.	This controversy book cause lots conflict in west society.

Opinion Summarization 2:

[Ganesan et al. WWW 12]

“Micro” Opinion Summarization

Kavita Ganesan, Chengxiang Zhai and Evelyne Viegas, Micropinion Generation: An Unsupervised Approach to Generating Ultra-Concise Summaries of Opinions, *Proceedings of the World Wide Conference 2012 (WWW'12)*, pages 869-878, 2012.

Micro Opinion Summarization

- **Generate a set of non-redundant phrases:**
 - Summarizing **key opinions** in text
 - Short (2-5 words)
 - Readable

 **Micropinions**

Micropinion summary for a restaurant:

“Good service”
“Delicious soup dishes”

- **Emphasize (1) ultra-concise nature of phrases; (2) abstractive summarization**

“Room is large”
“Room is clean”



“large clean room”

A general unsupervised approach

- **Main idea:**
 - use **existing words** in original text to compose meaningful summaries
 - leverage Web-scale n-gram language model to assess meaningfulness
- **Emphasis on 3 desirable properties of a summary:**
 - **Compactness**
 - summaries should use as **few words** as possible
 - **Representativeness**
 - summaries should reflect **major opinions** in text
 - **Readability**
 - summaries should be fairly **well formed**

Optimization Framework to capture compactness, representativeness & readability

$$M = \arg \max_{\{m_1 \dots m_k\}} \sum_{i=1}^k \underline{S_{rep}(m_i)} + \underline{S_{read}(m_i)}$$

Micropinion Summary, M

- 2.3 very clean rooms
- 2.1 friendly service
- 1.8 dirty lobby and pool
- 1.3 nice and polite staff

subject to

$$\sum_{i=1}^k |m_i| \leq \sigma_{ss} \quad \text{Size of summary}$$

$$S_{rep}(m_i) \geq \sigma_{rep} \quad \text{Minimum rep.}$$

$$S_{read}(m_i) \geq \sigma_{read} \quad \text{\& readability}$$

$$sim(m_i, m_j) \leq \sigma_{sim} \forall i, j \in [1, k] \quad \text{Redundancy}$$

Representativeness scoring: $S_{rep}(m_i)$

- 2 properties of a **highly representative phrase**:
 - Words should be **strongly associated** in text
 - Words should be **sufficiently frequent** in text
- Captured by **modified** pointwise mutual information

$$pmi'(w_i, w_j) = \log_2 \frac{p(w_i, w_j) \times c(w_i, w_j)}{p(w_i) \times p(w_j)} \leftarrow \text{Add frequency of occurrence within a window}$$

$$pmi_{local}(w_i) = \left[\frac{1}{2C} \sum_{j=i-C}^{i+C} pmi'(w_i, w_j) \right]$$

$$S_{rep}(w_1..w_n) = \frac{1}{n} \sum_{i=1}^n pmi_{local}(w_i)$$

Readability scoring, Sread(mi)

- Phrases are constructed from seed words, thus we can have **new phrases** not in original text
- Readability scoring based on N-gram language model (normalized **probabilities** of phrases)
 - **Intuition:** A phrase is more readable if it occurs more frequently on the web

$$S_{read}(w_k \dots w_n) = \frac{1}{K} \log_2 \prod_{k=q}^n p(w_k | w_{k-q+1} \dots w_{k-1})$$

Ungrammatical

“sucks life battery” -4.51

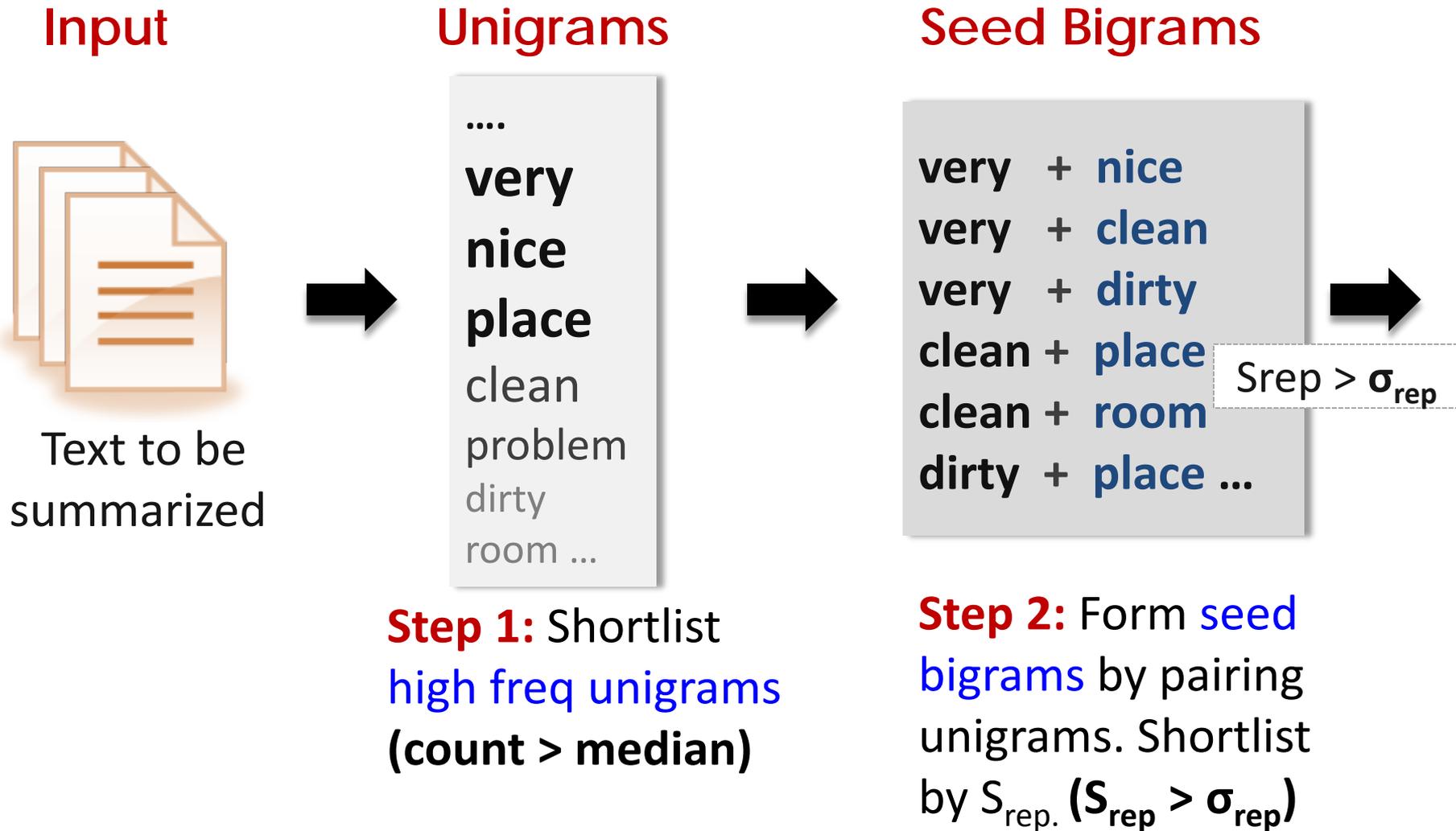
“life battery is poor” -3.66

Grammatical

“battery life sucks” -2.93

“battery life is poor” -2.37

Overview of summarization algorithm



Overview of summarization algorithm

Higher order n-grams

Summary

Candidates + Seed Bi-grams = New Candidates

very clean + clean rooms
clean bed = very clean rooms
very clean bed

very dirty + dirty room
dirty pool = very dirty room
very dirty pool

very nice + nice place
nice room = very nice place
very nice room

$S_{rep} < \sigma_{rep}; S_{read} < \sigma_{read}$



0.9 very clean rooms
0.8 friendly service
0.7 dirty lobby and pool
0.5 nice and polite staff
.....
.....

Sorted Candidates

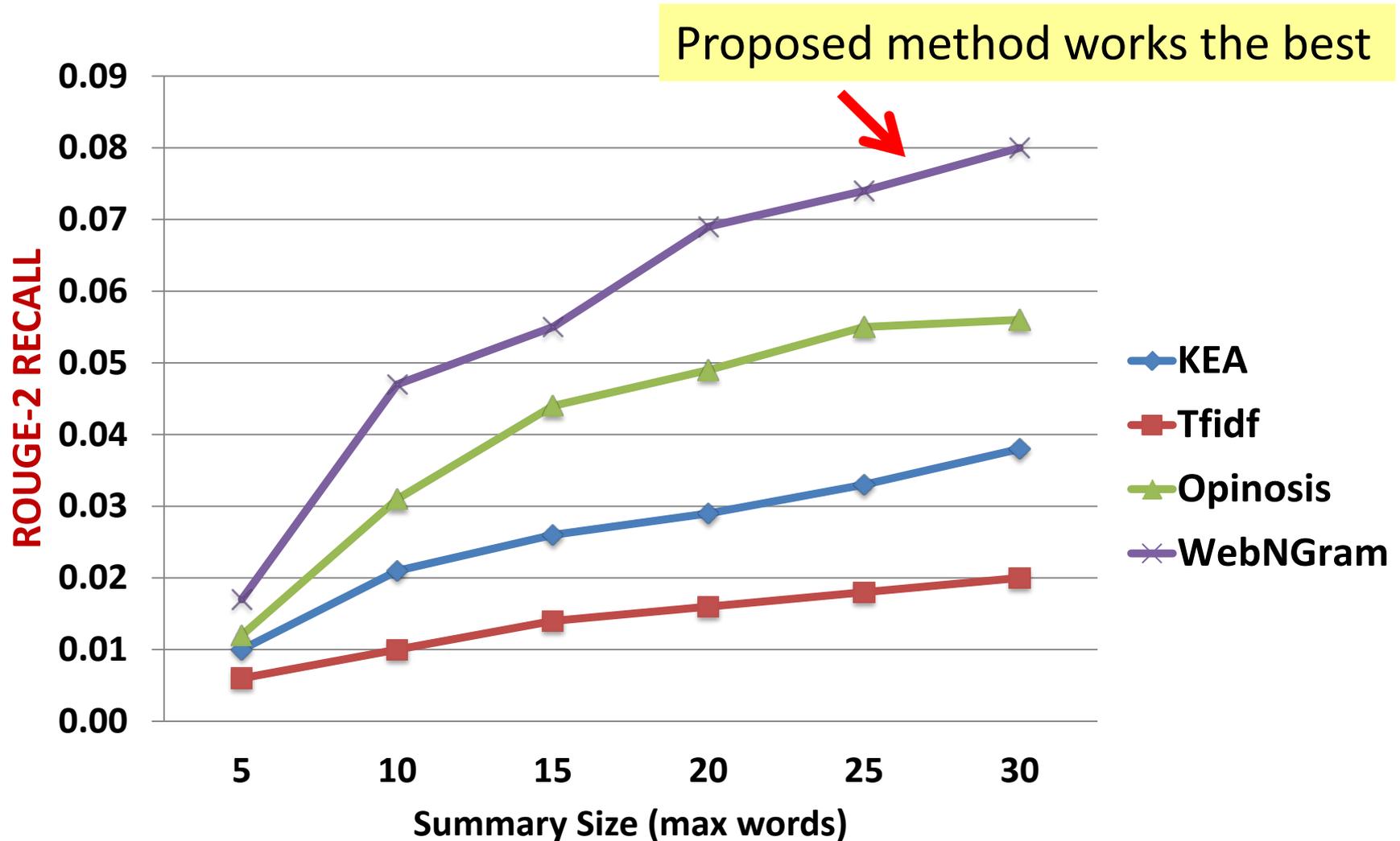
Step 3: Generate higher order n-grams.

- Concatenate existing candidates + seed bigrams
- Prune non-promising candidates (S_{rep} & S_{read})
- Eliminate redundancies ($\text{sim}(m_i, m_j)$)
- Repeat process on shortlisted candidates (until no possibility of expansion)

Step 4: Final summary.

Sort by objective function value. Add phrases until $|M| < \sigma_{ss}$

Performance comparisons (reviews of 330 products)



The program can generate meaningful novel phrases

Example:

Unseen N-Gram (Acer AL2216 Monitor)

“wide screen lcd monitor is bright”

readability : -1.88

representativeness: 4.25



“...plus the **monitor** is very **bright**...”

“...it is a **wide screen**, great color, great quality...”

“...this **lcd monitor** is quite **bright** and clear...”

**Related
snippets in
original text**

A Sample Summary

Canon Powershot SX120 IS

Easy to use
Good picture quality
Crisp and clear
Good video quality



Useful for pushing opinions
to devices where the **screen**
is small



**E-reader/
Tablet**



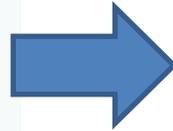
**Smart
Phones**



Cell Phones

Outline

Opinionated Text Data



1. Opinion Integration



2. Opinion Summarization

Query: *Dell Laptop*

	positive	negative	neutral
Topic 1 (Price)	- it is the best one and they show Dell laptop are as good as possible	- Even though Dell's price is cheaper, we still don't want it.	- most people prefer dell, a price compare. - DELL is trading at \$24.88
Topic 2 (Battery)	- One thing I really like about this Dell battery is the express charge feature.	- The Dell battery sucks. - stupid dell laptop battery	- I still want a free battery from dell.



3. Opinion Analysis

Decision Making & Analytics

"Which cell phone should I buy?"

"What are the winning features of iPhone over blackberry?"

"How do people like this new drug?"

"How is Obama's health care policy received?"

"Which presidential candidate should I vote for?"

...

Motivation

Hotel Palomar Chicago: Traveler Reviews

“ Great location+spacious room =happy traveler ”



★★★★★

leos_10 3 contributions
Boston

Jul 11, 2010 | Trip type: Couples **NEW**

Save Review



My ratings for this hotel

★★★★☆ Value
★★★★★ Rooms
★★★★★ Location
★★★★★ Cleanliness

★★★★★ Service
★★★★★ Sleep Quality

Stayed for a weekend in July. Walked everywhere, enjoyed the comfy bed and quiet hallways. [more](#)

“ terrific service and gorgeous facility ”



★★★★★

ahickling 1 contribution
Greensboro, North Carolina

Jul 7, 2010 | Trip type: Family **NEW**

Save Review



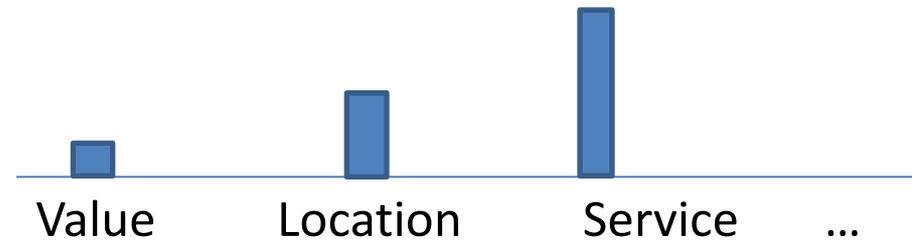
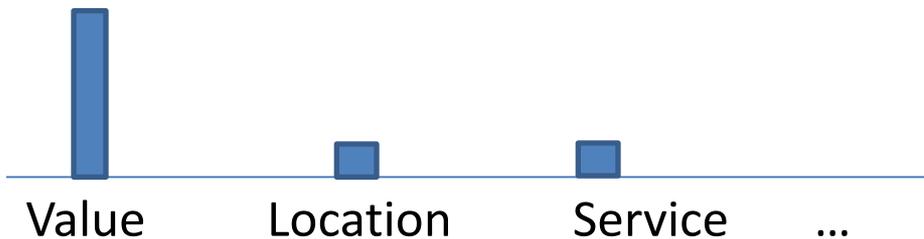
My ratings for this hotel

★★★★☆ Value
★★★★★ Rooms
★★★★★ Location
★★★★★ Cleanliness

★★★★★ Service
★★★★★ Sleep Quality

I stayed at the Palomar with my young daughter for three nights June 17-20, 2010 and absolutely loved the hotel. The room was one of the nicest I've ever stayed in (My daughter loved the Fuji jetted tub so much that she wanted to take 2 baths a day!) in terms of decor, design, and size. (It compared favorably to... [more](#)

How to infer aspect weights?



Opinion Analysis:

[Wang et al. KDD 2010] & [Wang et al. KDD 2011]

Latent Aspect Rating Analysis

Hongning Wang, Yue Lu, ChengXiang Zhai. Latent Aspect Rating Analysis on Review Text Data: A Rating Regression Approach, *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'10)*, pages 115-124, 2010.

Hongning Wang, Yue Lu, ChengXiang Zhai, Latent Aspect Rating Analysis without Aspect Keyword Supervision, *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'11)*, 2011, pages 618-626.

Latent Aspect Rating Analysis

- **Given a set of review articles about a topic with overall ratings**
- **Output**
 - Major aspects commented on in the reviews
 - Ratings on each aspect
 - Relative weights placed on different aspects by reviewers
- **Many applications**
 - Opinion-based entity ranking
 - Aspect-level opinion summarization
 - Reviewer preference analysis
 - Personalized recommendation of products
 - ...

Solving LARA in two stages: Aspect Segmentation + Rating Regression

Aspect Segmentation

+

Latent Rating Regression

Reviews + overall ratings

Aspect segments

Term Weights

Aspect Rating

Aspect Weight

W_{di}

location:1
amazing:1
walk:1
anywhere:1

β_i

0.0
2.9
0.1
0.9
0.1
1.7
0.1
3.9
2.1
1.2
1.7
2.2
0.6

S_i

3.9
4.8
5.8

α_d

0.2
0.2
0.6

r_d

Save Review

“Loved, Loved, Loved it”

Hotel Palomar Chicago



Tifplace 3 contributions
Queens, New York

Jul 7, 2010 | Trip type: Friends getaway

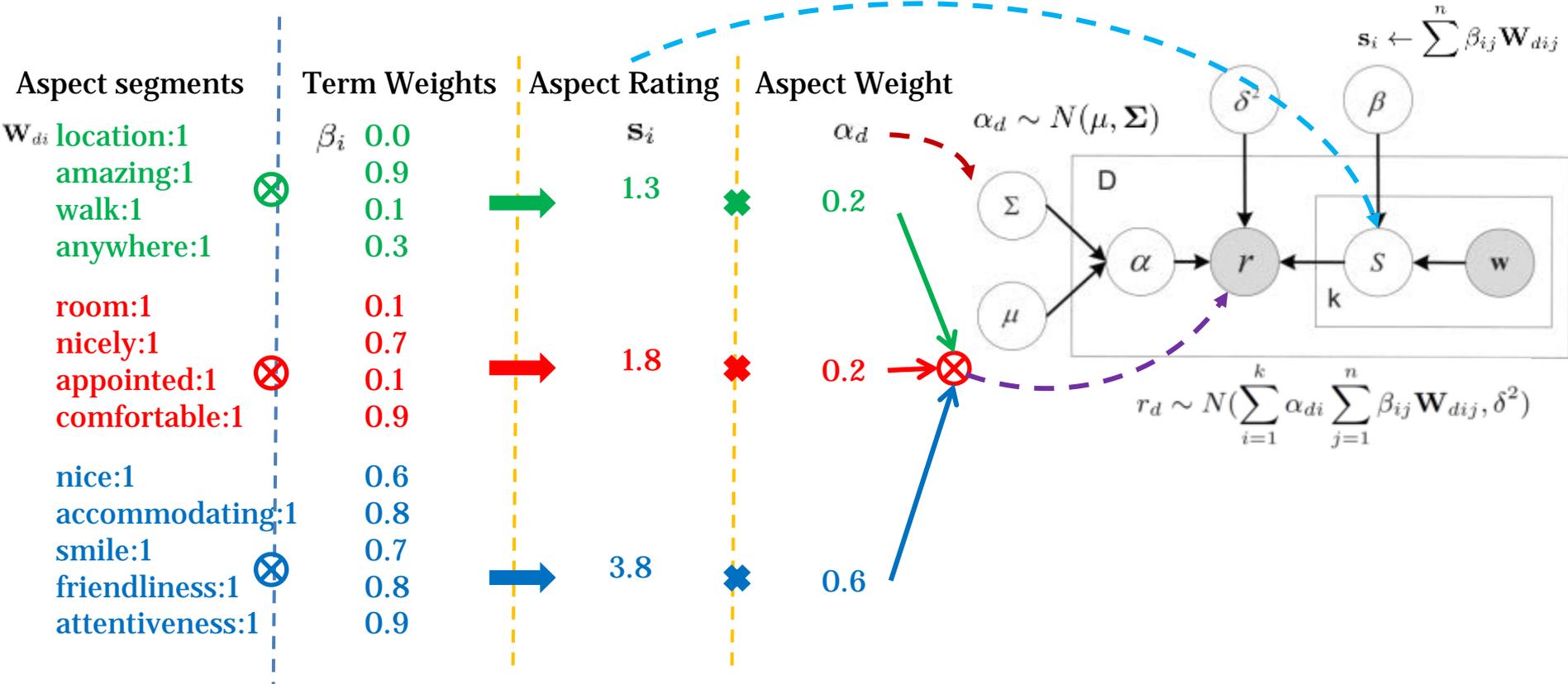
1 person found this review helpful

A friend and I stayed at the Hotel Palomar for the fourth of July Weekend. The hotel was very nice. The location was amazing. We could walk almost anywhere. The room was very nicely appointed and the bed was sooo comfortable. Eventhough the bathroom door did not close all the way, it was still pretty private. My friend and I were not out on by the door. I really liked the Sangria during the cocktail hour in the living room. But what I liked best about the Palomar was the staff. They were soooo nice and accomodating from my boy D money at the door, to the ladies at reception, to my new friend Greg in housekeeping and my other friend Ricky. Any questions or request we had were answered and fulfilled. They had us smiling and laughing the whole time. We really appreciated all the information they provided us with about where to go and what to do. When I come back to Chicago I will definitely stay at the Palomar again. I am sure there are other nice hotels in Chicago but I am not sure if you would get the same level of friendliness and attentiveness from their staff. If you stay at the hotel and the doorman D is there tell him that G Money sent you. Lol.

Observed

Latent!

Latent Rating Regression



Conditional likelihood

$$P(r|d) = P(r_d | \mu, \Sigma, \delta^2, \beta, W_d)$$

$$= \int p(\alpha_d | \mu, \Sigma) p(r_d | \sum_{i=1}^k \alpha_{di} \sum_{j=1}^n \beta_{dij} W_{dij}, \delta^2) d\alpha_d$$

A Unified Generative Model for LARA

Entity



Aspects



Review

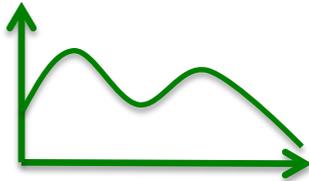


Aspect Rating

Aspect Weight

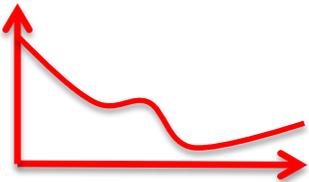
Location

location
amazing
walk
anywhere



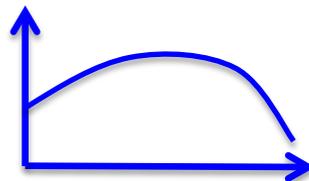
Room

room
dirty
appointed
smelly



Service

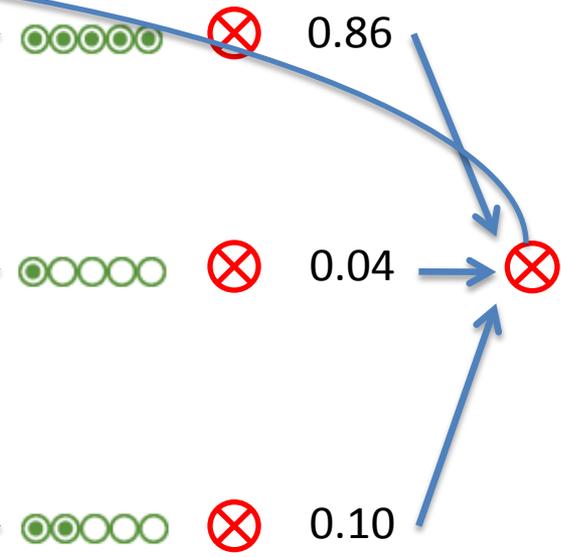
terrible
front-desk
smile
unhelpful



“Spend your money elsewhere”

○○○○○ Reviewed September 19, 2010

Excellent location in walking distance to Tiananmen Square and shopping streets. That’s the best part of this hotel! The rooms are getting really old. Bathroom was nasty. The fixtures were falling off, lots of cracks and everything looked dirty. I don’t think it worth the price. Service was the most disappointing part, especially the door men. this is not how you treat guests, this is not hospitality.



Latent Aspect Rating Analysis Model

• Unified framework

“Spend your money elsewhere”

●●○○○ Reviewed September 19, 2010

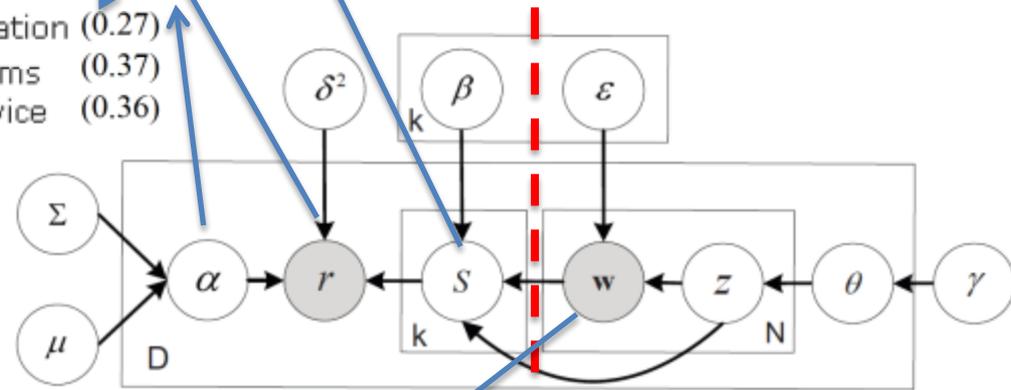
Excellent location in walking distance to Tiananmen Square and shopping streets. That’s the best part of this hotel! The rooms are getting really old. Bathroom was nasty. The fixtures were falling off, lots of cracks and everything looked dirty. I don’t think it worth the price. Service was the most disappointing part, especially the door men. this is not how you treat guests, this is not hospitality.

●●●○○ Location (0.27)
 ●●○○○ Rooms (0.37)
 ●○○○○ Service (0.36)

$$r \sim N\left(\sum_{i=1}^k \alpha_i \sum_{n=1}^{|d|} \beta_{ij} \Delta[w_n = v_j, z_n = i], \delta^2\right)$$

$$s_i = \sum_{n=1}^{|d|} \beta_{ij} \Delta[w_n = v_j, z_n = i]$$

$$\alpha \sim N(\mu, \Sigma)$$



$$p(\mathbf{W}, \mathbf{z}, \theta | \gamma, \epsilon) = p(\theta | \gamma) \prod_{n=1}^{|d|} p(w_n | z_n, \epsilon) p(z_n | \theta)$$

Rating prediction module Aspect modeling module

Sample Result 1: Rating Decomposition

- Hotels with the same overall rating but different aspect ratings

(All 5 Stars hotels, ground-truth in parenthesis.)

Hotel	Value	Room	Location	Cleanliness
Grand Mirage Resort	<u>4.2</u> (4.7)	3.8(3.1)	4.0(4.2)	4.1(4.2)
Gold Coast Hotel	4.3(4.0)	3.9(3.3)	3.7(3.1)	<u>4.2</u> (4.7)
Eurostars Grand Marina Hotel	3.7(3.8)	4.4(3.8)	4.1(4.9)	<u>4.5</u> (4.8)

- Reveal detailed opinions at the aspect level

Sample Result 2: Comparison of reviewers

- **Reviewer-level Hotel A**
 - Different reviewers' ratings

Reviewer	Value
Mr.Saturday	3.7(4.0)
Salsrug	<u>5.0(5.0)</u>

- Reveal differences in opinions

“ Good Price for what we got ”

Riu Palace Punta Cana



salsrug  13 contributions
Marylander

Save Review

Oct 27, 2008 | Trip type: Family

1 person found this review helpful

We stayed for six days, five nights. Overall, we had a very good time. The was pretty good and the staff was very friendly. They definitely do not skimp on the free alcohol. The room was a little smelly, which we had read on trip advisor so we bought a candle with us - no problem. They only thing I had an issue with was the little bugs. They were like gnats or fleas but they weren't either. I had some candy and popcorn which we brought from the States to munch on. I left it out on the table and within 40 minutes, the bag was infested. DO NOT keep any open food in your room. Also we ended up having to wash all of our clothes (clean and dirty) and airing out our luggage when we got home because we could still smell the room on them. For the price we paid, we really did have an excellent time besides those small things. The pool was awesome and the beach was spectacular. Out of the nearby resorts that we saw, Riu Palace Punta Cana was the best (it was also the nicest out of the other Riu's on Punta Cana). We went on the 1/2 day Outback Safari and had a great time. We got coffee and souvenirs cheaper than other places and the hotel. General - not good or bad just things that we noticed - There were a lot of topless sunbathers. The crowd is middle aged (35 - 55) so we were on the younger side and the majority of the people were European or Brazilian. It helps to know some spanish but it's not a necessity.

Liked — The beach was excellent.

Disliked — Room smell and little bugs.

My ratings for this hotel

 Value

 Rooms

 Location

 Cleanliness

 Check in / front desk

 Service

 Business service (e.g., internet access)

Sample Result 3: Aspect-Specific Sentiment Lexicon

<i>Value</i>	<i>Rooms</i>	<i>Location</i>	<i>Cleanliness</i>
resort 22.80	view 28.05	restaurant 24.47	clean 55.35
value 19.64	comfortable 23.15	walk 18.89	smell 14.38
excellent 19.54	modern 15.82	bus 14.32	linen 14.25
worth 19.20	quiet 15.37	beach 14.11	maintain 13.51
<i>bad -24.09</i>	<i>carpet -9.88</i>	<i>wall -11.70</i>	<i>smelly -0.53</i>
<i>money -11.02</i>	<i>smell -8.83</i>	<i>bad -5.40</i>	<i>urine -0.43</i>
<i>terrible -10.01</i>	<i>dirty -7.85</i>	<i>road -2.90</i>	<i>filthy -0.42</i>
<i>overprice -9.06</i>	<i>stain -5.85</i>	<i>website -1.67</i>	<i>dingy -0.38</i>

Uncover sentimental information directly from the data

Sample Result 4:

Validating preference weights

- **Analysis of hotels preferred by different types of reviewers**

<i>City</i>	<i>AvgPrice</i>	<i>Group</i>	<i>Val/Loc</i>	<i>Val/Rm</i>	<i>Val/Ser</i>
Amsterdam	241.6	top-10	190.7	214.9	221.1
		bot-10	270.8	333.9	236.2
Barcelona	280.8	top-10	270.2	196.9	263.4
		bot-10	330.7	266.0	203.0
San Francisco	261.3	top-10	214.5	249.0	225.3
		bot-10	321.1	311.1	311.4
Florence	272.1	top-10	269.4	248.9	220.3
		bot-10	298.9	293.4	292.6

- Reviewers emphasizing the ‘value’ aspect more would prefer cheaper hotels

Application 1: Rated Aspect Summarization

<i>Aspect</i>	<i>Summary</i>	<i>Rating</i>
Value	Truly unique character and a great location at a reasonable price Hotel Max was an excellent choice for our recent three night stay in Seattle.	3.1
	Overall not a negative experience, however considering that the hotel industry is very much in the impressing business there was a lot of room for improvement.	1.7
Location	The location, a short walk to downtown and Pike Place market, made the hotel a good choice.	3.7
	When you visit a big metropolitan city, be prepared to hear a little traffic outside!	1.2
Business Service	You can pay for wireless by the day or use the complimentary Internet in the business center behind the lobby though.	2.7
	My only complaint is the daily charge for internet access when you can pretty much connect to wireless on the streets anymore.	0.9

(Hotel Max in Seattle)

Application 2: Discover consumer preferences

- **Amazon reviews: no guidance**

Table 2: Topical Aspects Learned on MP3 Reviews

Low Overall Ratings			High Overall Ratings		
unit	jack	service	files	player	vision
usb	headphone	charge	format	music	video
battery	warranty	problem	included	download	player
charger	replacement	support	easy	headphones	quality
reset	problem	hours	convert	button	great
time	player	months	mp3	set	product
hours	back	weeks	videos	hours	sound
work	months	back	file	buds	radio
thing	buy	customer	wall	volume	accessory
wall	amazon	time	hours	ear	fm

battery life accessory service file format volume video

Application 3: User Rating Behavior Analysis

	<i>Expensive Hotel</i>		<i>Cheap Hotel</i>	
	<i>5 Stars</i>	<i>3 Stars</i>	<i>5 Stars</i>	<i>1 Star</i>
Value	0.134	0.148	0.171	0.093
Room	0.098	0.162	0.126	0.121
Location	0.171	0.074	0.161	0.082
Cleanliness	0.081	0.163	0.116	0.294
Service	0.251	0.101	0.101	0.049

People like expensive hotels because of good service

People like cheap hotels because of good value

Application 4: Personalized Ranking of Entities

Table 10: Personalized Hotel Ranking

Query: 0.9 value
0.1 others

Non-Personalized



Personalized



Hotel	Overall Rating	Price	Location
Majestic Colonial	5.0	339	Punta Cana
Agua Resort	5.0	753	Punta Cana
Majestic Elegance	5.0	537	Punta Cana
Grand Palladium	5.0	277	Punta Cana
Iberostar	5.0	157	Punta Cana
Elan Hotel Modern	5.0	216	Los Angeles
Marriott San Juan Resort	4.0	354	San Juan
Punta Cana Club	5.0	409	Punta Cana
Comfort Inn	5.0	155	Boston
Hotel Commonwealth	4.5	313	Boston

Open Questions

- **How can we combine all these methods in a general unified decision-support system?**
 - What are the basic common functions required by all applications?
 - How do we support users to interact with the system?
- **How far can we go with such pure statistical approaches?**
 - How can we maximize the benefit of unsupervised learning?
Continuous learning from the Web?
 - How can we combine unsupervised learning naturally with supervised learning through user interactions?
- **How can we incorporate linguistic resources & knowledge?**
 - How can we build a sentiment analyzer to take advantage all the resources available today?
 - Can we automatically construct sophisticated features for sentiment analysis? Deep learning?

Demo: FindILike System

Opinionated Text Data



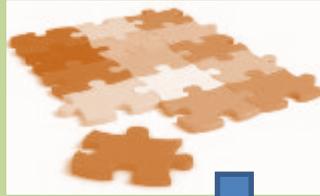
Topics: People, Events, Products, Services, ...



Sources: Blogs, Microblogs, Forums, Reviews, ...



1. Opinion Integration



2. Opinion Summarization

www.findilike.com

Query: Dell Laptop

	positive	negative	neutral
Topic 1 (Price)	it is the most best price that could use the coupon code as well as possible	price is high could be an option as the best price is	price per se can be a good price because it will be saving at \$2.48
Topic 2 (Battery)	standing longer than most other laptops in the region large feature	could be an option as the best price is	also a free battery and



3. Opinion Analysis

Decision Making & Analytics

"Which cell phone should I buy?"



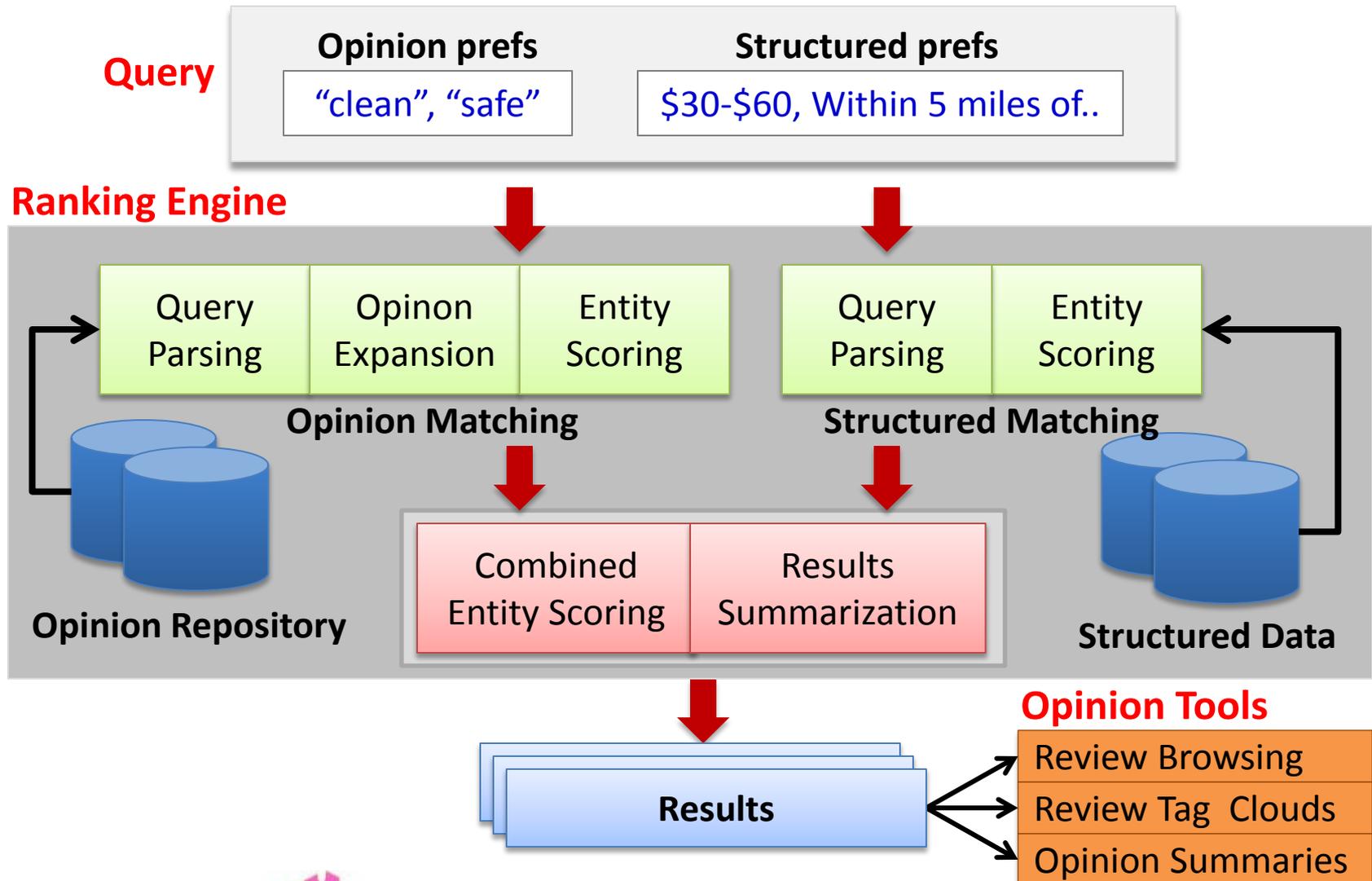
res of
y drug?"
e policy

"Which presidential candidate should I vote for?"

...

Findilike: Opinion-Based Decision-Support

www.findilike.com



Opinion-Based Entity Ranking



Query = "near ohare airport, free internet"

YourMatch 5/5

CHICAGO O'HARE HOTEL
★ ★ ★ ★ ☆ (680 reviews)
8201 W Higgins
gym . internet ac

Opinions ★ ★ ★ ★ ★
[near ohare airport](#) ★ ★ ★ ★ ★
Conceinently located by **Ohare** ... hotels.com
[free internet](#) ★ ★ ★ ★ ★
Good location if you're traveling from Chicago O'Hare airport
(Free shuttle every 30 minutes!) excellent restaurant (Bistro
90) on the premisses, **free wireless internet** access

- Room was clean (7)
- Room was nice (3)
- Great room for the price (30)
- Staff was helpful (6)
- Bed was comfortable (5)
- Internet is only in the lobby (3)
- Great food (6)

What's Buzzing...
People Think...
See On Map

hotel to airport (4) chicao o'hare hotel (6) room the price (8)
good location (11) service was friendly (2)
place to stay (8) available in lobby (5) internet in lobby (7)
room was good (4) **close the airport (12)**
good room (11) staff was helpful (4)
bed was comfortable (6) shuttle to airport (6) hotel was clean (2)
service to airport (4) o'hare garden hotel (6) price was good (4)
location was good (4) **front desk (15)**
room was clean (7) chicao o'hare (10)
service was good (6) shuttle service good (3)
good price (16) garden hotel (10)
room was nice (4) breakfast was donuts (4)

Map Review

Firefox Search, analyze and decide bas... findilike.cs.illinois.edu

Chicago Hotels - Showing 20 results

view: List, Simple, Map (checked); Sort: 20, YourMatch

- 1. **Chicago O'hare Garden** Ho... \$74
★☆☆☆☆ (680 guest ratings)
5.0 Opinions ★★★★★
- 2. **Aloft Chicago O'hare** \$89
★★★★☆ (365 guest ratings)
4.5 Opinions ★★★★★
- 3. **Hilton Rosemont Chicago ...** \$99
★★★★☆ (348 guest ratings)
4.5 Opinions ★★★★★
- 4. **Radisson Hotel Chicago O...** \$89
★★★★☆ (292 guest ratings)
4.5 Opinions ★★★★★
- 5. **Travelodge Hotel Downtown** \$69
★★★★☆ (417 guest ratings)
4.5 Opinions ★★★★★

O'Hare Airport

7:39 AM 1/28/2013

Acknowledgments

- Collaborators: Yue Lu, Qiaozhu Mei, Kavita Ganesan, Hongning Wang, and many others
- Funding



Thank You!

Questions/Comments?

More information can be found at <http://timan.cs.uiuc.edu/>