# Advanced Topics in Information Retrieval

# 9. Social Media

**Vinay Setty
(vsetty@mpi-inf.mpg.de)**

Jannik Strötgen
(jtroetge@mpi-inf.mpg.de)

# Outline

max planck institut
informatik

MAX-PLANCK-GESELLSCHAFT

# 9.1. What is Social Media?

‣ Content creation is **supported by software** (no need to know HTML, CSS, JavaScript)

‣ Content is **user-generated** (as opposed to by big publishers) or **collaboratively-edited** (as opposed to by a single author)

‣ **Web 2.0** (if you like –outdated– buzzwords)

‣ Examples:

   ‣ **Blogs** (e.g., Wordpress, Blogger, Tumblr)

   ‣ **Social Networks** (e.g., facebook, Google+)

   ‣ **Wikis** (e.g., Wikipedia but there are many more)
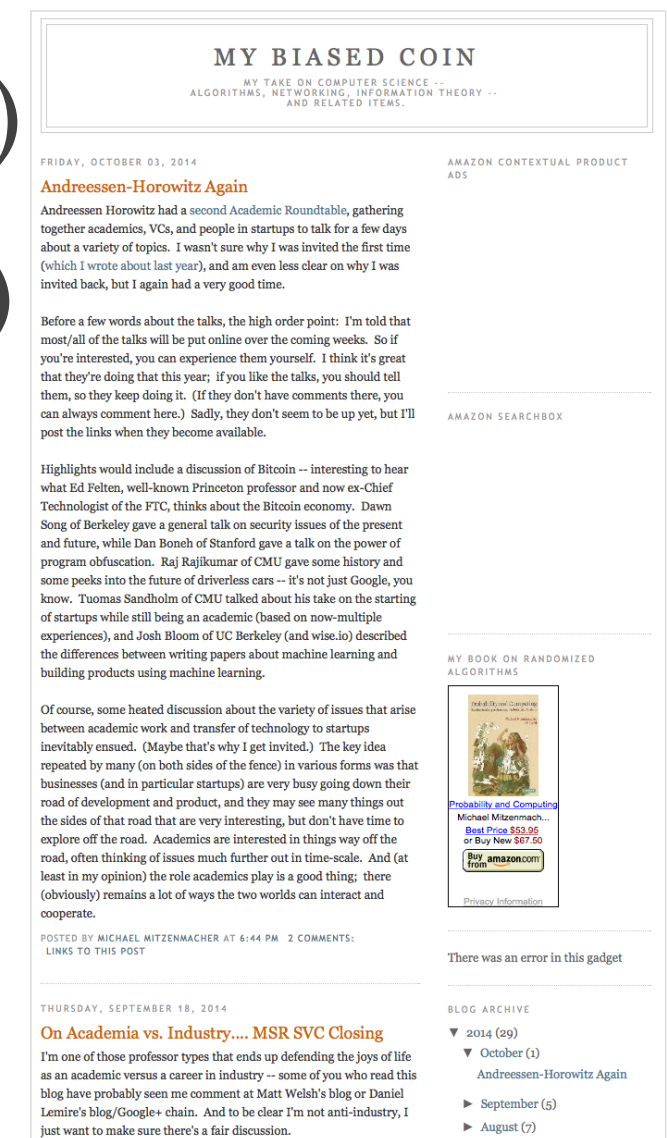
   ‣ …

= ?!? =

max planck institut
informatik

# Weblogs, Blogs, the Blogosphere

▸ **Journal-like website**, editing supported by software, self-hosted or as a service

▸ Initially often run by **enthusiasts**, now also common in the **business world**, and some bloggers make their living from it

▸ **Reverse chronological order** (newest first)

▸ **Blogroll** (whose blogs does the blogger read)

▸ **Posts** of varying length and topics

▸ **Comments**

▸ Backed by **XML feed** (e.g., RSS or Atom) for **content syndication**

# Weblogs, Blogs, the Blogosphere



http://mybiasedcoin.blogspot.de

- WordPress.com

  - ~ 60M blogs

  - ~ 50M posts/month

  - ~ 50M comments/month

- Tumblr.com (by Yahoo!)

  - ~ 208M blogs

  - ~ 95B posts

  - ~ 100M posts/day

# Twitter



- Micro-blogging service created in March '06

- Posts (tweets) limited to **140 characters**

- **271M** monthly active **users**

- 500M tweets/day = ~6K tweets/second

- **2B queries** per day

- 77% of accounts are outside of the U.S.

- Hashtags (#atir2016)

- Messages (@vinaysetty)

- Retweets

# Facebook, Twitter, LinkedIn, Pinterest, …

# Challenges & Opportunities

▸ Content

    ▸ **plenty of context** (e.g., publication timestamp, relationships between users, user profiles, comments, external urls)

    ▸ **short posts** (e.g., on Twitter), **colloquial/cryptic language**

    ▸ **spam** (e.g., splogs, fake accounts)

▸ Dynamics

    ▸ **up-to-date content** – real-world events covered as they happen

    ▸ **high update rates** pose severe engineering challenges (e.g., how to maintain indexes and collection statistics)

max planck institut informatik

# How do People Search Blogs?

‣ Mishne and de Rijke [8] analyzed a **month-long query log** from a blog search engine ([blogdigger.com](blogdigger.com)) and found that

> ‣ queries are **mostly informational** (vs. transactional or navigational)
>
> > ‣ **contextual**: in which context is a specific **named entity** (i.e., person, location, organization) mentioned, for instance, to find out opinions about it
> >
> > ‣ **conceptual**: which blogs cover a specific **high-level concept or topic** (e.g., stock trading, gay rights, linguists, islam)
> >
> > ‣ contextual more common than conceptual both for **ad-hoc and filtering queries**
>
> ‣ most popular topics: **technology, entertainment, and politics**
>
> ‣ many queries (15–20%) **related to current events**

# How do People Search Twitter?

- Teevan et al. [10] **conducted a survey** (54 MS employees), compared **query logs** from web search and Twitter, finding that queries on Twitter

    - are often related to **celebrities, memes, or other users**

    - are **often repeated** to monitor a specific topic

    - are on average **shorter** than web queries (1.64 vs. 3.08 words)

    - tend to return **results** that are **shorter** (19.55 vs. 33.95 words), **less diverse**, and more often relate to **social gossip** and **recent events**

- People also directly **express information needs** using Twitter: **17% of tweets** in the analyzed data correspond to **questions**

# What Data?

- **Feeds** (e.g., blog, twitter user, facebook page)

- **Posts** (e.g., blog posts, tweets, facebook posts)

- We'll consider

  - **textual content** of posts

  - **publication timestamps** of posts

  - **hyperlinks** contained in posts

- We'll ignore

  - other links (e.g., friendship, follower/followee)

  - hashtags, images, comments

# Tasks

‣ **Meme tracking** grouping of memes to track them over period of time

‣ **Post retrieval** identifies posts relevant to a specific information need (e.g., how is life in Iceland?)

‣ **Opinion retrieval** finds posts **relevant** to a specific **named entity** (e.g., a company or celebrity) which **express an opinion** about it

‣ **Feed distillation** identifies feeds relevant to a topic, so that the user can subscribe to their posts (e.g., who tweets about C++?)

‣ **Top-story identification** leverages social media to determine the most important news stories (e.g., to display on front page)

# Outline

9.1. What is Social Media?

9.2. Tracking Memes

9.3. Opinion Retrieval

9.4. Feed Distillation

9.5. Top-Story Identification

# 9.2. Tracking Memes

▸ Leskovec et al. [5] track **memes** (e.g., "lipstick on a pig") and visualize their volume in traditional news and blogs



▸ Demo: http://www.memetracker.org

# Phrase Graph Construction

▸ <u>Problem</u>: Memes are **often modified** as they spread, so that first **all mentions of the same meme** need to be identified

▸ Construction of a **phrase graph** G(V, E):

  ‣ **vertices** V correspond to **mentions of a meme**
    that are reasonably long and occur often enough

  ‣ **edge** (u,v) exists if meme mentions u and v

    ‣ u is **strictly shorter** than v

    ‣ <u>either</u>: have **small directed token-level edit distance**
      (i.e., u can be transformed into v by adding at most ε tokens)

    ‣ <u>or</u>: have a **common word sequence** of length at least k

  ‣ **edge weights** based on **edit distance** between u and v
    and how often v occurs in the document collection

# Meme Phrase Graph

# Phrase Graph Partitioning

▸ Phrase graph is an **directed acyclic graph** (DAG) by construction

▸ Partition G(V, E) by **deleting a set of edges having minimum total weight,** so that each resulting **component is single-rooted**

▸ Phrase graph partitioning is *NP*-hard, hence addressed by **greedy heuristic algorithm**

# Applications

‣ **Clustering of meme mentions** allows for insightful analyses, e.g.:

  ‣ **volume of meme** per time interval

  ‣ **peak time** of meme in traditional news and social media

  ‣ **time lag** between peek times in traditional news and social media

# Outline

9.1. What is Social Media?

9.2. Tracking Memes

**9.3. Opinion Retrieval**

9.4. Feed Distillation

9.5. Top-Story Identification

max planck institut
informatik

# 9.3. Opinion Retrieval

‣ **Opinion retrieval** finds posts **relevant** to a specific **named entity** (e.g., a company or celebrity) which **express an opinion** about it

‣ <u>Examples</u>: (from TREC Blog track 2006)

> ‣ macbook pro
>
> ‣ jon stewart
>
> ‣ whole foods
>
> ‣ mardi gras
>
> ‣ cheney hunting

> **Title:**
> whole foods
>
> **Description:**
> Find opinions on the quality, expense, and value of purchases at Whole Foods stores.
>
> **Narrative:**
> All opinions on the quality, expense and value of Whole Foods purchases are relevant. Comments on business and labor practices or Whole Foods as a stock investment are not relevant. Statements of produce and other merchandise carried by Whole Foods without comment are not relevant.

‣ **Standard retrieval models** can help with finding relevant posts; but how to determine **whether a post expresses an opinion?**

# Opinion Retrieval Task Example

```
<top>
<num> Number: 863

<title> netflix

<desc> Description:
Identify documents that show customer opinions
of Netflix.

<narr> Narrative:
A relevant document will indicate subscriber
satisfaction with Netflix.  Opinions about
the Netflix DVD allocation system, promptness
or delay in mailings are relevant.
Indications of having been or
intent to become a Netflix subscriber that do
not state an opinion are not relevant.
</top>
```

max planck institut
informatik

# Opinion Dictionary

‣ What if we had a **dictionary of opinion words?**
(e.g., like, good, bad, awesome, terrible, disappointing)

‣ Lexical resources with **word sentiment information**

  ‣ **SentiWordNet** (http://sentiwordnet.isti.cnr.it/)

  | | |
  |---|---|
  | unspeakable#2  terrible#2  painful#3  dreadful#2  **awful**#1  atrocious#2  abominable#2 | 01126291 |

  exceptionally bad or displeasing; "atrocious taste"; "abominable workmanship"; "an awful voice"; "dreadful manners"; "a painful performance"; "terrible handwriting"; "an unspeakable odor came sweeping into the room"

  P: 0 O: 0.125 N: 0.875

  Feedback on SentiWordNet values: They are OK. | Suggest your values.

  | | |
  |---|---|
  | awing#1  **awful**#6  awesome#1  awe-inspiring#1  amazing#2 | 01282510 |

  inspiring awe or admiration or wonder; "New York is an amazing city"; "the Grand Canyon is an awe-inspiring sight"; "the awesome complexity of the universe"; "this sea, whose gently awful stirrings seem to speak of some hidden soul beneath"- Melville; "Westminster Hall's awing majesty, so vast, so high, so silent"

  P: 0.875 O: 0 N: 0.125

  Feedback on SentiWordNet values: They are OK. | Suggest your values.

  ‣ **General Inquirer** (http://www.wjh.harvard.edu/~inquirer/)

  ‣ **OpinionFinder** (http://mpqa.cs.pitt.edu)

# Opinion Dictionary

- He et al. [4] construct an **opinion dictionary** from training data

  - consider only words that are neither too frequent (e.g., and, or) nor too rare (e.g., aardvark) in the post collection D

  - let $D_{rel}$ be a set of **relevant posts** (to any query in a workload) and $D_{relopt} \subset D_{rel}$ be the subset of **relevant opinionated posts**

  - two options to **measure opinionatedness of a word** v

    - **Kullback-Leibler Divergence**

$$op_{KLD}(v) = P[\,v \mid D_{relopt}\,] \log_2 \frac{P[\,v \mid D_{relopt}\,]}{P[\,v \mid D_{rel}\,]}$$

    - **Bose Einstein Statistics**

$$op_{BO}(v) = tf(v, D_{relopt}) \log_2 \frac{1 + \lambda}{\lambda} + \log_2(1 + \lambda) \quad \text{with} \quad \lambda = \frac{tf(v, D_{rel})}{|D_{rel}|}$$

# Re-Ranking

‣ He et al. [4] **measure opinionatedness of a post** d as follows

  ‣ consider the set $Q_{opt}$ of k **most opinionated words** from the dictionary

  ‣ issue $Q_{opt}$ as a query (e.g., using Okapi BM25 as a retrieval model)

  ‣ the retrieval status value $score(d, Q_{opt})$ measures how opinionated d is

‣ Posts are ranked in response to query Q (e.g., whole foods) according to a (linear) **combination of retrieval scores**

$$score(d) = \alpha \cdot score(d, Q) + (1 - \alpha) \cdot score(d, Q_{opt})$$

with $0 \leq \alpha \leq 1$ as a **tunable mixing parameter**

# Sentiment Expansion

▸ Huang and Croft [5] **expand the query** with **query-independent** (Q<sub>I</sub>) and **query-dependent** (Q<sub>D</sub>) opinion words; posts are then ranked according to

$$score(d) = \alpha \cdot score(d, Q) + \beta \cdot score(d, Q_I)$$
$$+ (1 - \alpha - \beta) \cdot score(d, Q_D)$$

with $0 \leq \alpha, \beta \leq 1$ as a **tunable mixing parameters**
and retrieval scores based on **language model divergences**

▸ **Query-independent opinion words** are obtained as

  ‣ **seed words** (e.g, good, nice, excellent, poor, negative, unfortunate, …)

  ‣ **most frequent words** in opinionated corpora (e.g., movie reviews)

# Sentiment Expansion (Query Independent)

▸ Examples: (of most frequent words in different corpora)

   ‣ Cornell movie reviews: like, even, good, too, plot

   ‣ MPQA opinion corpus: against, minister, terrorism, even, like

   ‣ Blog06(op): like, know, even, good, too

▸ Observation: Query-independent opinion words are very general (e.g., like, good) or specific to the corpus (e.g., minister, terrorism)

# Sentiment Expansion (Query Dependent)

▸ **Query-dependent opinion words** are obtained as words that frequently co-occur with query terms in **pseudo-relevant documents** (following the approach by Lavrenko and Croft [6])

▸ Given a query q, identify the set of R of top-k pseudo-relevant documents, and top-n words having highest probability

$$P[\,w \mid R\,] \propto \sum_{d \in R} P[\,w \mid d\,] \prod_{v \in q} P[\,v \mid d, w\,]$$

$$P[\,v \mid d, w\,] = \begin{cases} \frac{tf(v,d)}{\sum_u tf(u,d)} & : & w \in d \\ 0 & : & \text{otherwise} \end{cases}$$

with parameter set as k = 5 and n = 20 in practice

# Sentiment Expansion

▶ <u>Examples</u>: (of query-dependent opinion words)

‣ **mozart** → (like, good, too, even, death, best, great, genius)

‣ **allianz** → (best, premium, great, value, traditional, fidelity)

‣ **wikipedia** → (like, open, good, know, free, great, knowledge)

# Outline

9.1. What is Social Media?

9.2. Tracking Memes

9.3. Opinion Retrieval

**9.4. Feed Distillation**

9.5. Top-Story Identification

# 9.4. Feed Distillation

- **Feed distillation** identifies feeds (e.g., blogs, Twitter users) that are **relevant** to a **specific (typically rather broad) topic**

- Examples: (from TREC Blog track 2007)

  - movie review

  - firearm control

  - baseball

  - garden

  - mobile phone

  ---

  **Title:**
  baseball

  **Description:**
  Blogs with recurring interests in Major League Baseball, or lesser leagues, for example, giving news or analysis of games or player moves.

  **Narrative:**
  Relevant blogs will have news or analysis from the major league baseball and other leagues. Blogs listing only product reviews, or with other nonsensical information are not relevant.

- Challenges: How to capture whether a blog **consistently covers** the given topic? How to bridge **vocabulary gap** to posts?

# Language Models

▸ Weerkamp et al. [11] develop two approaches to feed distillation estimating **language models** for **entire blog(ger)s** and **individual posts**, respectively

▸ <u>Notation</u>:

  ▸ a blog b is a set of posts; |b| is the number of posts by b

  ▸ a post p is a bag of terms

  ▸ tf(v, p) denotes the term frequency of term v in post p

  ▸ B denotes a virtual post concatenating all posts from all blogs

# Blogger Model (BM)

▸ Estimates a language model **for each blog(ger)** b

$$P[\,q \mid \theta_b\,] = \prod_{v \in q} P[\,v \mid \theta_b\,]^{\,tf(v,q)}$$

▸ Smooths probability estimates using the collection of blogs B

$$P[\,v \mid \theta_b\,] = (1 - \lambda_b) \cdot P[\,v \mid b\,] + \lambda_b \cdot P[\,v \mid B\,]$$

with **blog-specific smoothing parameter**

$$\lambda_b = \frac{\beta}{(1/|b| \cdot \sum_{p \,\in\, b} \sum_v tf(v,p)) + \beta}$$

thus smoothing blogs with **shorter posts more aggressively**

# Blogger Model

▸ **Two-step generation** of term v from blog b

$$P[\,v \mid b\,] = \sum_{p \,\in\, b} P[\,v \mid p, b\,]\, P[\,p \mid b\,]$$

assuming **conditional independence** of terms given blog

$$P[\,v \mid b\,] = \sum_{p \,\in\, b} P[\underbrace{\,v \mid p\,}_{\substack{\text{2. Draw term} \\ \text{from post}}}] \, P[\underbrace{\,p \mid b\,}_{\substack{\text{1. Draw post} \\ \text{from blog}}}]$$

▸ **Uniform probability** of posts given blog (i.e., equal importance)

$$P[\,p \mid b\,] = 1/|b|$$

▸ **Maximum-likelihood estimate** $\quad P[\,v \mid p\,] = \dfrac{tf(v, p)}{\sum_w tf(w, p)}$

# Posting Model (PM)

▸ Estimates a language model **for each individual post** p

$$P[\,v \mid \theta_p\,] = (1 - \lambda_p) \cdot P[\,v \mid p\,] + \lambda_p \cdot P[\,v \mid B\,]$$

with **post-specific smoothing parameter**

$$\lambda_p = \frac{\beta}{\left(\sum_w tf(w, p)\right) + \beta}$$

thus smoothing **short posts more aggressively**

▸ **Maximum-likelihood estimate** $P[\,v \mid p\,] = \dfrac{tf(v, p)}{\sum_w tf(w, p)}$

# Posting Model

▸ Likelihood of generating query q from language model of post p

$$P[\,q \mid \theta_p\,] = \prod_{v \in q} P[\,v \mid \theta_p\,]^{\,tf(v,q)}$$

▸ **Two-step generation** of query q from blog b

$$P[\,q \mid b\,] = \sum_{p \in b} P[\,\underbrace{q \mid \theta_p}\,]\,P[\,\underbrace{p \mid b}\,]$$

2. Generate query    1. Draw post
from post         from blog

▸ **Uniform probability** of posts given blog (i.e., equal importance)

$$P[\,p \mid b\,] = 1/|b|$$

# Query Expansion for Vocabulary Gap

▸ Elsass et al. [3] proposed the highly similar **Large Document Model** (~BM) and **Small Document Model** (~PM) approaches

▸ Focus on bridging the **vocabulary gap** between high-level topic descriptions (e.g., garden) and posts (e.g., seed, flower, crop)

▸ **Query expansion** with terms from **pseudo-relevant documents** retrieved from different corpora

  ‣ **Blogs** (MAP 0.266 compared to small document model 0.315)

  ‣ **Posts** (MAP 0.282)

  ‣ **Wikipedia articles** (MAP 0.314)

  ‣ **Wikipedia passages** (MAP 0.313)

No Improvement!

# Query Expansion for Vocabulary Gap

▸ **Query expansion** based on **anchor phrases** in Wikipedia

- ▸ **issue original query** q against Wikipedia articles as corpus

- ▸ **consider** top-k and top-n (k < n) **results** returned by query

- ▸ **score every anchor phrase** a occurring in any top-n result and pointing to a document d from the top-k result as

$$score(a) = \sum_{(a,d)} (k - rank(d))$$

anchor phrase **a** from top-**n** article
pointing to top-**k** article **d**



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

united states

united states of america

america

land of the free

the states

favoring **frequent anchor phrases** pointing to **highly ranked articles**

- ▸ **expand query** with top-m anchor phrases (MAP 0.361)

IMPROVEMENT!

# Outline

9.1. What is Social Media?

9.2. Tracking Memes

9.3. Opinion Retrieval

9.4. Feed Distillation

**9.5. Top-Story Identification**

# Online News Media



Thousands of news articles generated each day!

# Google News

# News Aggregators

Portal:Current events

digg

YAHOO! News

reddit

Flipboard

news360

NEWSTIFY
NEWS IS WHAT WE AIM FOR

max planck institut
informatik

# Wikipedia Current Events Portal

## July 7, 2016 (Thursday)

edit   history   watch

### Disasters and accidents

- Super Typhoon Nepartak
    - The first major typhoon of 2016 threatens Taiwan, China and northern Luzon, Philippines. Thousands of people have been evacuated in Taiwan. (The Weather Channel), (ABC News)
    - Typhoon Nepartak is expected to make landfall on mainland China on Friday and will make flooding worse. Nearly 200 people have died in flood waters in China in the past week with 41 people missing, 1.6 million relocated and almost 50000 houses collapsed. (*The Telegraph*)

### Law and crime

- A group of suspected radical Islamists hurl homemade bombs at police officers in the Kishoreganj District in central Bangladesh killing at least one officer and injuring several others. (AP via ABC News)

### Politics and elections

- Australian federal election, 2016
    - Australian Prime Minister Malcolm Turnbull's Liberal/National coalition, behind Bill Shorten's Labor Party in the first 48 hours following Saturday's election, is now ahead of Labor in the Lower House, 74-71 seats, just two seats shy of the minimum needed to form a government. Minor parties and independents have won five seats; mail-in and absentee votes are still being counted. Turnbull is on the road today seeking support from a small handful of independent and small party lawmakers. (Reuters) (*The Australian*) (*Daily Mail*)

---

**Time:** 11:33 UTC | **Day:** 7 July | Purge

| << | **July 2016** | | | | | >> |
|---|---|---|---|---|---|---|
| S | M | T | W | T | F | S |
| | | | | | 1 | 2 |
| 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| 17 | 18 | 19 | 20 | 21 | 22 | 23 |
| 24 | 25 | 26 | 27 | 28 | 29 | 30 |
| 31 | | | | | | |

More July 2016 events...

About this page · Suggest a headline

News about Wikipedia

## Ongoing events

### Business

- United Kingdom withdrawal from the European Union
- 2016 Bangladesh Bank heist
- Panama Papers

### Disasters

- 2016 California wildfires
- 2016 Fort McMurray wildfire

### Health

- Flint water crisis
- Zika virus outbreak

42

# Top-Story Identification

▶ **Top-story identification** (another task within the TREC Blog track) aims to identify the **most important news stories for a specific day d** based on their **coverage in the blogosphere**

  ‣ **real-time** (online, limited statistics, time critical: small lag)

  ‣ **retrospective**: (offline, full statistics)

▶ <u>Notation:</u>

  ‣ d denotes the day of interest

  ‣ $B_d$ is the set of posts published at day d; p denotes a post

  ‣ n denotes a news article (consisting of headline and content)

  ‣ $tf(v,p)$ is the term frequency of term v in post p

# Top-Story Identification

▸ Lee and Lee [7] address retrospective top-story identification using **language models** estimated from news and blogs

▸ <u>Intuition</u>: *"News article important if discussed by many posts"*

$$Importance(n,d) \propto KL(\theta_n \parallel \theta_{B_d})$$

<span style="color:magenta">LM representing news article **n**</span>     <span style="color:magenta">LM representing posts published at day **d**</span>

(Note: This is a simplified version of the approach described in [7])

▸ Only articles **published -1/+1 around the day of interest** d are considered as candidates and ranked by the approach

# Top-Story Identification Workflow

# Blog Post Language Model

▸ Language model for **blog posts published at** d is estimated as

$$P[\, v \mid \theta_{B_d} \,] = \frac{tf(v, B_d) + \mu \cdot \frac{tf(v,B)}{\sum_w tf(w,B)}}{(\sum_w tf(w, B_d)) + \mu}$$

using Dirichlet smoothing with the collection of all posts B

# News-Story Language Model

▸ <u>Option 1</u>: Estimate **directly from content** of news article

$$P[\,v \mid \theta_n\,] = \frac{tf(v,n) + \mu \cdot \frac{tf(v,N)}{\sum_w tf(w,N)}}{\left(\sum_w tf(w,n)\right) + \mu}$$

VOCABULARY GAP?!?

using Dirichlet smoothing with the entire news collection N

▸ <u>Option 2</u>: Estimate from top-k **pseudo-relevant blog posts** $B_n$ retrieved using **headline** as query and **published within -1/+1 month** of the news article; again using Dirichlet smoothing with the collection of all posts B

▸ <u>Option 3</u>: **Interpolate language models** estimated from news article content and top-k pseudo-relevant blog posts

# Summary

- ▸ **Meme tracking**
  grouping variants of memes to track them over time

- ▸ **Opinion retrieval**
  finds posts expressing an opinion about a specific named entity

- ▸ **Feed distillation**
  identifies feeds worth following for a given high-level topic

- ▸ **Top-story identification**
  spots most important news articles based on coverage in blogs

- ▸ **Vocabulary gaps**
  are a common obstacle in IR but can often be bridged

- ▸ **Language models**
  are versatile and can be used to address many (if not most) tasks

# References

[1]  A. Dong, R. Zhang, P. Kolari, J. Bai, F. Diaz, Y. Chang, Z. Zheng:
*Time is of the Essence: Improving Recency Ranking Using Twitter Data*,
WWW 2010

[2]  M. Efron:
*Information Search and Retrieval in Microblogs*,
JASIST, 62(6):996–1008, 2011

[3]  J. Elsass, J. Arguello, J. Callan, J. G. Carbonell:
*Retrieval and Feedback Models for Blog Feed Search*,
SIGIR 2008

[4]  B. He, C. Macdonald, J. He, Iadh Ounis:
*An Effective Statistical Approach for Blog Post Opinion Retrieval*,
CIKM 2008

[5]  X. Huang and W. B. Croft:
*A Unified Relevance Model for Opinion Retrieval*,
CIKM 2009

[6]  V. Lavrenko and W. B. Croft:
*Relevance-Based Language Models*,
SIGIR 2001

# References

[7]     Y. Lee and J.-H. Lee:
        *Identifying top news stories based on their popularity in the blogosphere,*
        Information Retrieval 17:326–350, 2014

[8]     G. Mishne and M. de Rijke:
        *A Study of Blog Search,*
        ECIR 2006

[9]     R. L. T. Santos, C. Macdonald, R. McCreadie, I. Ounis:
        *Information Retrieval on the Blogosphere,*
        FTIR 6(1):1–125, 2012

[10]    J. Teevan, D. Ramage, M. R. Morris:
        *#TwitterSearch: A Comparison of Microblog Search and Web Search,*
        WSDM 2011

[11]    W. Weerkamp, K. Balog, M. de Rijke:
        *Blog feed search with a post index,*
        Information Retrieval 14:515–545, 2011