

# Real-Time Top-R Topic Detection on Twitter with Topic Hijack Filtering

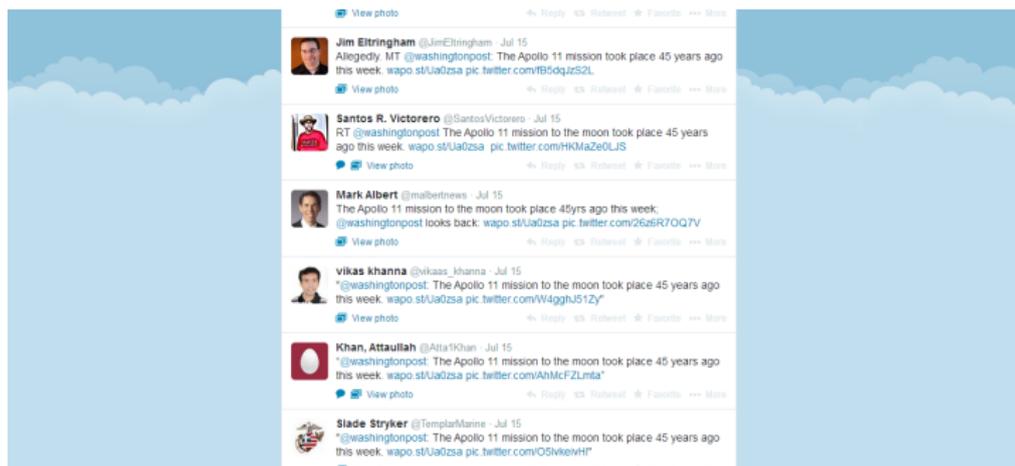
Kohei Hayashi (NII & ERATO)

August 3, 2015

Joint work with

- Takanori Maehara (Shizuoka Univ)
- Masashi Toyoda (Univ Tokyo)
- Ken-ichi Kawarabayashi (NII & ERATO)

# Twitter



SNS with short messages (**tweets**)

**Big data** 41M users, 1.4B interactions

**Diversity** Covering any topics: news, politics, TV, ...

**Rapid** 1 tweet  $\leq$  140 chars

$\Rightarrow$  Low latency

# Automatic Trend Detection on Twitter



# Automatic Trend Detection on Twitter



A promising data resource for **topic detection**

- Find word clusters by *word co-occurrence*
- May discover breaking news and events **even faster than news media**

# Two Challenges

- Topic Detection in Real-time
- Noise Filtering

# Topic Detection in Real-time



An ultimate goal of topic detection on Twitter

- Have to deal with **0.27M tweets/min**
- Words rarely co-occur  
⇒ **Severely degrade the quality of topics**

# Noise Filtering

Many spam tweets generated by **not human**

- e.g. “tweet buttons”

The image shows a news article from The Washington Post titled "The Apollo 11 Mission to the Moon". A red circle highlights a "Tweet" button in the share options. A red arrow points from this button to a vertical list of tweets on the right. These tweets are all variations of a spam message: "Allegedly. MT @washingtonpost. The Apollo 11 mission took place 45 years ago this week. wapo.st/Ua0zsa pic.twitter.com/IB5dqJzS2L". The tweets are generated by different accounts, some with profile pictures of the Washington Post logo, and all include a "View photo" link. This illustrates how a single click on a "tweet button" can generate a large volume of spam tweets.

Exaggerate co-occurrence and **“hijack”** important topics

# Contributions

**A streaming topic detection algorithm** based on *non-negative matrix factorization (NMF)*

- ① **Highly scalable:**  
Able to deal with a  $20\text{M} \times 1\text{M}$  sparse matrix/sec
- ② **Automatic topic hijacking detection & elimination**

# Contributions

**A streaming topic detection algorithm** based on *non-negative matrix factorization (NMF)*

- ① **Highly scalable:**  
Able to deal with a  $20M \times 1M$  sparse matrix/sec
- ② **Automatic topic hijacking detection & elimination**

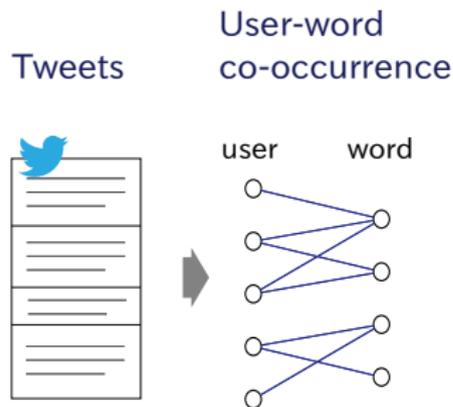
## Technical Points

- ① Reformulate NMF in a **stochastic manner**
  - Stochastic gradient descent (SGD) updates with  $O(\text{NNZ})$  time
- ② Use of **statistical testing**
  - Assume normal topics follow *power law*

# Streaming NMF

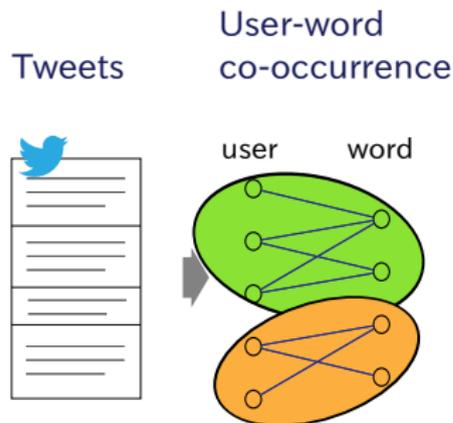
# Topic Detection by NMF

Consider to obtain  $R$  topics from all past tweets



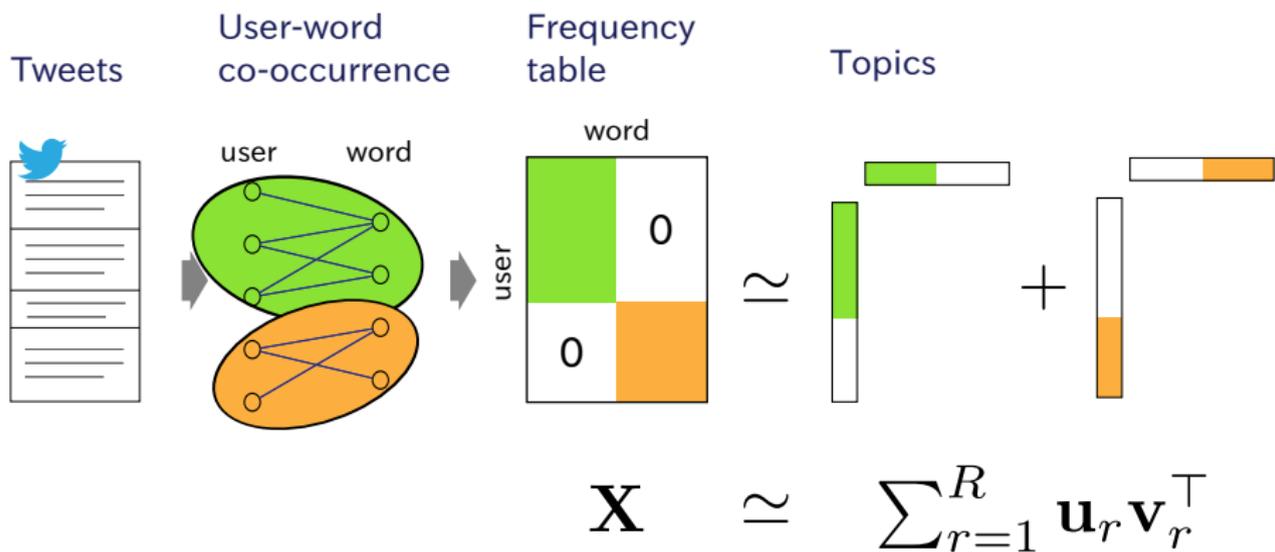
# Topic Detection by NMF

Consider to obtain  $R$  topics from all past tweets



# Topic Detection by NMF

Consider to obtain  $R$  topics from all past tweets



## Problem

$$\min_{\mathbf{U} \in \mathbb{R}_+^{I \times R}, \mathbf{V} \in \mathbb{R}_+^{J \times R}} f_\lambda(\mathbf{X}; \mathbf{U}, \mathbf{V}),$$

$$f_\lambda(\mathbf{X}; \mathbf{U}, \mathbf{V}) = \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{V}^\top\|_F^2 + \frac{\lambda}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)$$

## Batch Algorithm

Repeat until convergence:

- 1  $\mathbf{U} \leftarrow [\mathbf{U} - \eta \nabla_{\mathbf{U}} f_\lambda(\mathbf{X}; \mathbf{U}, \mathbf{V})]_+$
- 2  $\mathbf{V} \leftarrow [\mathbf{V} - \eta \nabla_{\mathbf{V}} f_\lambda(\mathbf{X}; \mathbf{U}, \mathbf{V})]_+$

Guaranteed converging to stationary points

# Stochastic Formulation

- Now we observe  $\mathbf{X}^{(t)}$  for each time  $t$
- Keep track to  $\bar{\mathbf{X}}^{(t)} = \frac{1}{t} \sum_{s=1}^t \mathbf{X}^{(s)}$ 
  - Efficient updates from  $\bar{\mathbf{X}}^{(t)}$  to  $\bar{\mathbf{X}}^{(t+1)}$ ?

# Stochastic Formulation

- Now we observe  $\mathbf{X}^{(t)}$  for each time  $t$
- Keep track to  $\bar{\mathbf{X}}^{(t)} = \frac{1}{t} \sum_{s=1}^t \mathbf{X}^{(s)}$ 
  - Efficient updates from  $\bar{\mathbf{X}}^{(t)}$  to  $\bar{\mathbf{X}}^{(t+1)}$ ?

**Key idea:** decompose  $f_\lambda$  for each  $t$

$$\|\bar{\mathbf{X}}^{(t)} - \mathbf{UV}^\top\|_{\text{F}}^2 = \frac{1}{t} \sum_s \|\mathbf{X}^{(s)} - \mathbf{UV}^\top\|_{\text{F}}^2 + \text{const.}$$

# Stochastic Formulation

- Now we observe  $\mathbf{X}^{(t)}$  for each time  $t$
- Keep track to  $\bar{\mathbf{X}}^{(t)} = \frac{1}{t} \sum_{s=1}^t \mathbf{X}^{(s)}$ 
  - Efficient updates from  $\bar{\mathbf{X}}^{(t)}$  to  $\bar{\mathbf{X}}^{(t+1)}$ ?

**Key idea:** decompose  $f_\lambda$  for each  $t$

$$\|\bar{\mathbf{X}}^{(t)} - \mathbf{UV}^\top\|_{\text{F}}^2 = \frac{1}{t} \sum_s \|\mathbf{X}^{(s)} - \mathbf{UV}^\top\|_{\text{F}}^2 + \text{const.}$$

If  $\mathbf{X}^{(t)}$  is i.i.d. random variable,

$$\|\mathbb{E}[\mathbf{X}^{(t)}] - \mathbf{UV}^\top\|_{\text{F}}^2 = \mathbb{E}\|\mathbf{X}^{(t)} - \mathbf{UV}^\top\|_{\text{F}}^2 + \text{var}[\mathbf{X}^{(t)}]$$

Now we can use **stochastic optimization**

## Streaming Algorithm

For  $t = 1, \dots, T$ :

- 1  $\mathbf{A}_U \leftarrow \nabla_{\mathbf{U}} \nabla_{\mathbf{U}} f_\lambda, \quad \mathbf{A}_V \leftarrow \nabla_{\mathbf{V}} \nabla_{\mathbf{V}} f_\lambda$  (metrics)
- 2  $\mathbf{U}^{(t)} \leftarrow [\mathbf{U}^{(t-1)} - \eta_t \nabla_{\mathbf{U}} f_\lambda(\mathbf{X}^{(t)}; \mathbf{U}, \mathbf{V}^{(t-1)}) \mathbf{A}_U^{-1}]_+$
- 3  $\mathbf{V}^{(t)} \leftarrow [\mathbf{V}^{(t-1)} - \eta_t \nabla_{\mathbf{V}} f_\lambda(\mathbf{X}^{(t)}; \mathbf{U}^{(t)}, \mathbf{V}) \mathbf{A}_V^{-1}]_+$

Guaranteed converging to the stationary points of  $f_\lambda(\bar{\mathbf{X}}^{(t)}; \mathbf{U}, \mathbf{V})$  for some  $\eta_n$  and i.i.d.  $\{\mathbf{X}^{(t)}\}$

# Comparing with the Batch NMF ...

## Much faster

- Update cost: depends on  $\text{NNZ}(\mathbf{X}^{(t)}) \ll \text{NNZ}(\bar{\mathbf{X}}^{(t)})$

## Memory efficient

- Able to discard  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(t-1)}$

# Comparing with the Batch NMF ...

## Much faster

- Update cost: depends on  $\text{NNZ}(\mathbf{X}^{(t)}) \ll \text{NNZ}(\bar{\mathbf{X}}^{(t)})$

## Memory efficient

- Able to discard  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(t-1)}$

## Smoothing effect

$$\begin{aligned}\mathbf{U}^{(t)} &= [(1 - \eta_t)\mathbf{U}^{(t-1)} + \eta_t \mathbf{X}^{(t)} \mathbf{V}^{(t-1)} \mathbf{A}_{\mathbf{U}}^{-1}]_+ \\ &= [(1 - \eta_t)\mathbf{U}^{(t-1)} + \eta_t \underset{\mathbf{U}}{\text{argmin}} f_\lambda(\mathbf{X}^{(t)}; \mathbf{U}, \mathbf{V}^{(t-1)})]_+\end{aligned}$$

- A weighted average of the prev solution and the NMF solution of  $\mathbf{X}^{(t)}$
- Mitigates the sparsity of  $\mathbf{X}^{(t)}$

# Topic Hijacking Detection

# Problem Setting

**Goal:** Find hijacked topics

**Idea:** Check word distributions

- The word dist of *a hijacked topic* should be different from the word dist of *a normal topic*

# Problem Setting

**Goal:** Find hijacked topics

**Idea:** Check word distributions

- The word dist of a *hijacked topic* should be different from the word dist of a *normal topic*
- NMF estimates topic-specific word dists as  $\mathbf{V}$

$$\begin{aligned}X_{ij} &\propto p(\text{user}_i, \text{word}_j) \\ &\propto \sum_r p(\text{topic}_r) p(\text{user}_i | \text{topic}_r) p(\text{word}_j | \text{topic}_r) \\ &\propto \sum_r u_{ir} v_{jr}\end{aligned}$$

# Defining Normal/Hijacked Topics

## Normal topics:

- Many users involve
  - ⇒ Mixing many different vocabularies
  - ⇒ Heavy-tailed (Zipf's law)

# Defining Normal/Hijacked Topics

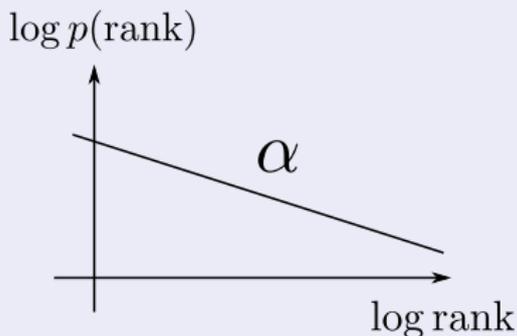
## Normal topics:

- Many users involve
  - ⇒ Mixing many different vocabularies
  - ⇒ Heavy-tailed (Zipf's law)

## Definition: A Normal Topic

Topic  $r$  is normal if

$$p(\text{rank}(\text{word}) \mid \text{topic}_r) \\ = \text{power}(\alpha)$$



# Defining Normal/Hijacked Topics (Cont'd)

## Hijacked topics:

- Few users involve
  - ⇒ The same vocabulary is repeatedly used
  - ⇒ Uniform probs & (almost) no tail

# Defining Normal/Hijacked Topics (Cont'd)

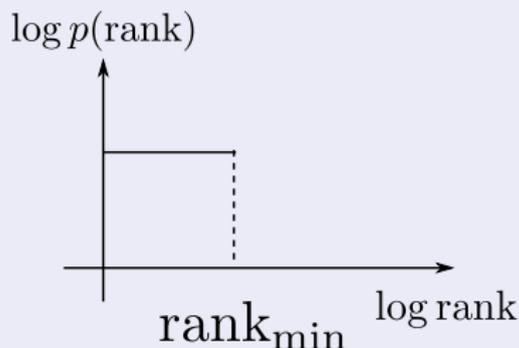
## Hijacked topics:

- Few users involve
  - $\Rightarrow$  The same vocabulary is repeatedly used
  - $\Rightarrow$  **Uniform probs & (almost) no tail**

## Definition: A Hijacked Topic

Topic  $r$  is **hijacked** if

$$p(\text{rank}(\text{word}) \mid \text{topic}_r) \\ = \text{step}(\text{rank}_{\min})$$



# Log-likelihood Ratio Test

$$\mathcal{L}(\text{rank}_{\min}) = \sum_j \log \frac{\text{step}(\text{rank}_j \mid \text{rank}_{\min})}{\text{power}(\text{rank}_j \mid \hat{\alpha})}$$

## Theorem (Asymptotic normality [Vuong'89])

Let  $N$  be # of observed words. Then,  $\mathcal{L}(\text{rank}_{\min})/\sqrt{N}$  converges in distribution to  $N(0, \sigma^2)$  where

$$\sigma^2 = \frac{1}{N} \sum_{j=1}^N \left( \log \frac{\text{step}(\text{rank}_j \mid \text{rank}_{\min})}{\text{power}(\text{rank}_j \mid \hat{\alpha})} \right)^2 - \left( \frac{1}{N} \mathcal{L}(\text{rank}_{\min}) \right)^2.$$

# Log-likelihood Ratio Test

$$\mathcal{L}(\text{rank}_{\min}) = \sum_j \log \frac{\text{step}(\text{rank}_j \mid \text{rank}_{\min})}{\text{power}(\text{rank}_j \mid \hat{\alpha})}$$

## Theorem (Asymptotic normality [Vuong'89])

Let  $N$  be # of observed words. Then,  $\mathcal{L}(\text{rank}_{\min})/\sqrt{N}$  converges in distribution to  $N(0, \sigma^2)$  where

$$\sigma^2 = \frac{1}{N} \sum_{j=1}^N \left( \log \frac{\text{step}(\text{rank}_j \mid \text{rank}_{\min})}{\text{power}(\text{rank}_j \mid \hat{\alpha})} \right)^2 - \left( \frac{1}{N} \mathcal{L}(\text{rank}_{\min}) \right)^2.$$

For  $r = 1, \dots, R$ :

- Estimate  $\hat{\alpha} = \text{argmax}_{\alpha} \log \text{power}(\text{rank}(\mathbf{u}_k) \mid \alpha)$
- For  $\text{rank}_{\min} = 1, \dots, 140$ :
  - Compute  $\mathcal{L}(\text{rank}_{\min})$
  - Topic  $r$  is hijacked if  $p\text{-val} < 0.05$

# Experiments

# Data

## Japanese Twitter stream

- April 15–16, 2013
  - 417K users
  - 1.98M words
  - 15.3M tweets
  - 69.4M co-occurrences
- Generated  $\mathbf{X}^{(t)}$  for each 10K co-occurrences

# Runtime

	NMF	1%	10%	100%
Batch		27.0m	1.9h	16.4h
Online	[Cao+ 07]	5.6h	9.2h	17.8h
Dynamic	[Saha+ 12]	16.7h	>24h	>24h
Streaming	[proposed]	4.0m	21.7m	3.6h
w/ Filter	[proposed]	5.3m	24.1m	3.8h

# Runtime

	NMF	1%	10%	100%
Batch		27.0m	1.9h	16.4h
Online	[Cao+ 07]	5.6h	9.2h	17.8h
Dynamic	[Saha+ 12]	16.7h	>24h	>24h
Streaming	[proposed]	4.0m	21.7m	3.6h
w/ Filter	[proposed]	5.3m	24.1m	3.8h

- ✓ Streaming NMF:  $\times 5$ –250 faster!
  - 67K tweets/m  
⇒ The real-time speed of all jp tweets!
- ✓ Filtering cost is ignorable

# Perplexity

- Similarity between a topic dist and a target dist
  - Used Y! headlines<sup>1</sup> for the target term dist
- Lower is better

	NMF	1%	10%	100%
Batch		9.01E+09	6.32E+06	4.51E+04
Online	[Cao+ 07]	1.06E+04	3.25E+04	1.83E+04
Dynamic	[Saha+ 12]	6.53E+04	N/A	N/A
Streaming	[proposed]	5.65E+07	3.25E+04	8.71E+03
w/ Filter	[proposed]	5.25E+09	2.40E+04	7.90E+03

---

<sup>1</sup><http://news.yahoo.co.jp/list>

# Perplexity

- Similarity between a topic dist and a target dist
  - Used Y! headlines<sup>1</sup> for the target term dist
- Lower is better

	NMF	1%	10%	100%
Batch		9.01E+09	6.32E+06	4.51E+04
Online	[Cao+ 07]	1.06E+04	3.25E+04	1.83E+04
Dynamic	[Saha+ 12]	6.53E+04	N/A	N/A
Streaming	[proposed]	5.65E+07	3.25E+04	8.71E+03
w/ Filter	[proposed]	5.25E+09	2.40E+04	7.90E+03

- ✓ Streaming NMF: best at 100% data
- ✓ Topic hijacking filter improves perplexity

<sup>1</sup><http://news.yahoo.co.jp/list>



# Detected Hijacking Phrases

- シネマトウデイ 妖精トム 映画 ハリウッド 戦闘 クルーズ 主演
- 発売 2013 情報 開催 参加 商品 入荷
- 限定 情報 開催 イベント 好評 商品 @\*\*\*\*\*
- リプライ 富士山 樋口 フェスティバル 興味 早稲田 開催 現在 不問 募集 スタッフ 大

広告系

- リフォロー アカウント 19836 フォロワー Only 希望 支援 696382 交流
- 無料 日本 拡散 希望 2013/04 リプ

拡散系

- 特徴 燃料 格闘 操縦 射撃 装甲 整備 評価 機動
- 所有 City ポイント 前回 Tweet アクセス 獲得 Intel
- 自動 だれ 宣伝 AutoTweet オートツイート Twitter 定期 設定 サイト

ボット系

- とノω・ □ | □□ | □ □ □ | □ □ □ | □
- □ | う— □ □ Å・ | お ● □ □ □ | よ。) □。・
- \*・ \*・ :。...: \*・! \*・!° :。...:

顔文字系

# Summary

Proposed **the streaming algorithm for Twitter topic detection**

- Works in real time  
(would handle all jp tweets in theory)
- Automatically filters spam topics

# Summary

Proposed **the streaming algorithm for Twitter topic detection**

- Works in real time  
(would handle all jp tweets in theory)
- Automatically filters spam topics

Thank you!

# Integrated Twitter Topic Detection System

For  $t = 1, \dots, T$ :

- ① Generate  $\mathbf{X}^{(t)}$  from tweets  $\notin$  Blacklist  $O(N_t)$
- ② Update  $\mathbf{U}, \mathbf{V}$  by SGD  $O(N_t R^2 + R^3)$
- ③ With some intervals,  
Detect Topic Hijacking and update Blacklist  $O(J)$