

Topic Detection & Tracking

Advanced Social Computing

Department of Computer Science
University of Massachusetts, Lowell
Spring 2020

Hadi Amiri
hadi@cs.uml.edu



Announcement

- **Project Presentations**
 - Date: 3/29, 3:30 – 6:00 PM

- **Project Final Report**
 - Due Date: 5/15, 11:59 PM

Lecture Topics

- Topic Detection
- Topic Tracking
- Early Prediction

Matrix Factorization for Topic Detection

Topic Detection

Tweets

Computer technology: 2-Tone L.E.D. to Simplify Screens

Stock Market: A Better Deal For Investors Isn't Simple. Large Sale 03/02

The Shape of Cinema, Transformed At the Click of a Mouse. Movie production.

The three big Internet portals begin to distinguish among themselves as shopping malls

Topic Detection

Topics

Computer	0.02
Technology	0.03
System	0.04
Internet	0.01
...	

Sale	0.02
Product	0.03
Market	0.02
Consumer	0.04
...	

Film	0.05
Movie	0.04
Theater	0.02
Production	0.04
...	

Tweets

Computer technology: 2-Tone L.E.D. to Simplify Screens

Stock Market: A Better Deal For Investors Isn't Simple. Large Sale 03/02

The Shape of Cinema, Transformed At the Click of a Mouse. Movie production.

The three big Internet portals begin to distinguish among themselves as shopping malls

A topic is a distribution over words

Topic Detection

Topics

Computer	0.02
Technology	0.03
System	0.04
Internet	0.01
...	

Sale	0.02
Product	0.03
Market	0.02
Consumer	0.04
...	

Film	0.05
Movie	0.04
Theater	0.02
Production	0.04
...	

Tweets

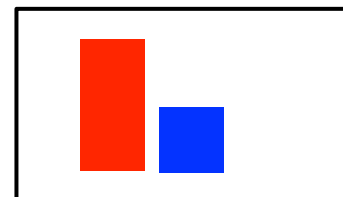
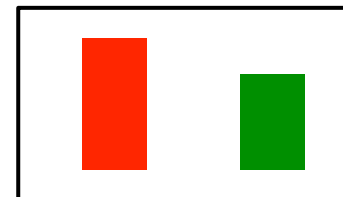
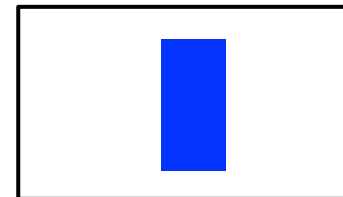
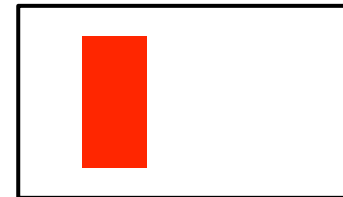
Computer technology: 2-Tone L.E.D. to Simplify Screens

Stock Market: A Better Deal For Investors Isn't Simple. Large Sale 03/02

The Shape of Cinema, Transformed At the Click of a Mouse. Movie production.

The three big Internet portals begin to distinguish among themselves as shopping malls

Assignments



A **topic** is a distribution over words

A **tweet** is a mixture of topics / distribution over topics

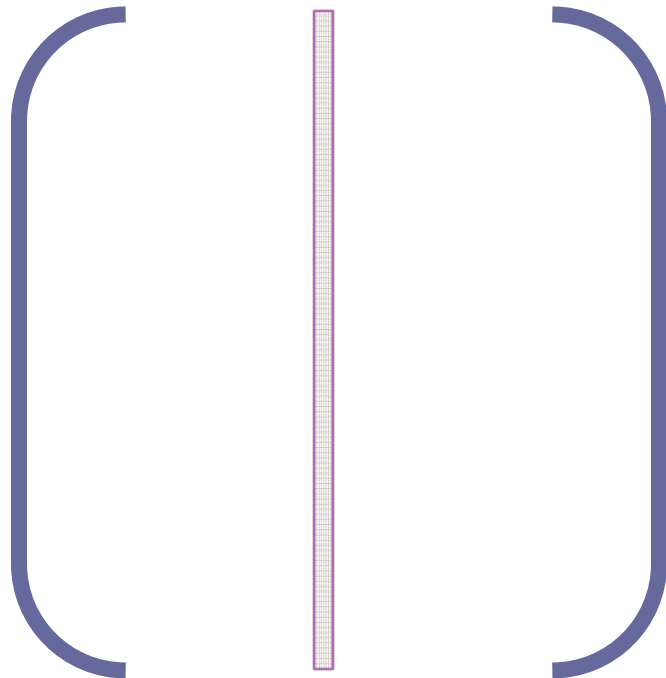
Topic Detection

- Different learning techniques
- Matrix factorization methods
 - LU decomposition
 - Singular Value Decomposition(SVD)
 - Probabilistic Matrix Factorization(PMF)
 - (Online) Non-negative Matrix Factorization(NMF)
 - Etc.

Topic Detection - NMF

m : # **terms** in the dataset
 n : # **docs** in the dataset
 k : # **topics** in the dataset

dataset

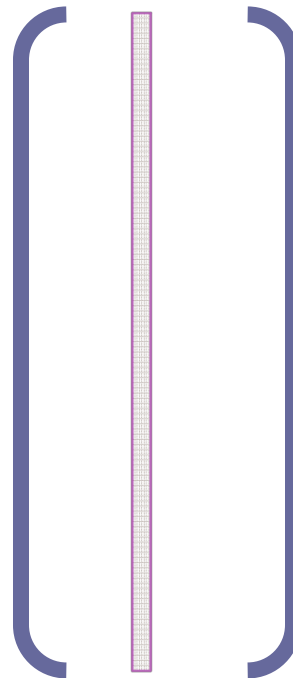


$m \times n$

a sample doc vector

\approx

topics

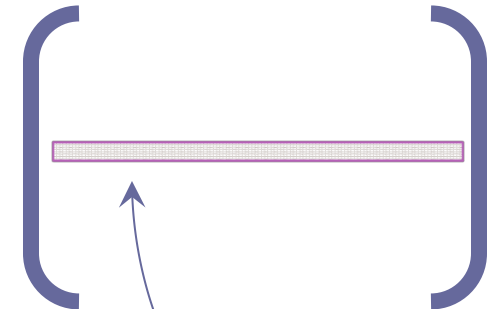


$m \times k$

a sample topic vector

\times

topic assignment



$k \times n$

the topic assignment to docs

Topic Detection - NMF

m : # **terms** in the dataset
 n : # **docs** in the dataset
 k : # **topics** in the dataset

dataset

S

\approx

topics

D

\times

topic assignment

X

$k \times n$

$m \times n$

$m \times k$

Topic Detection - NMF

m : # **terms** in the dataset
 n : # **docs** in the dataset
 k : # **topics** in the dataset

dataset

S

$m \times n$

\approx

topics

D

$m \times k$

\times

topic assignment

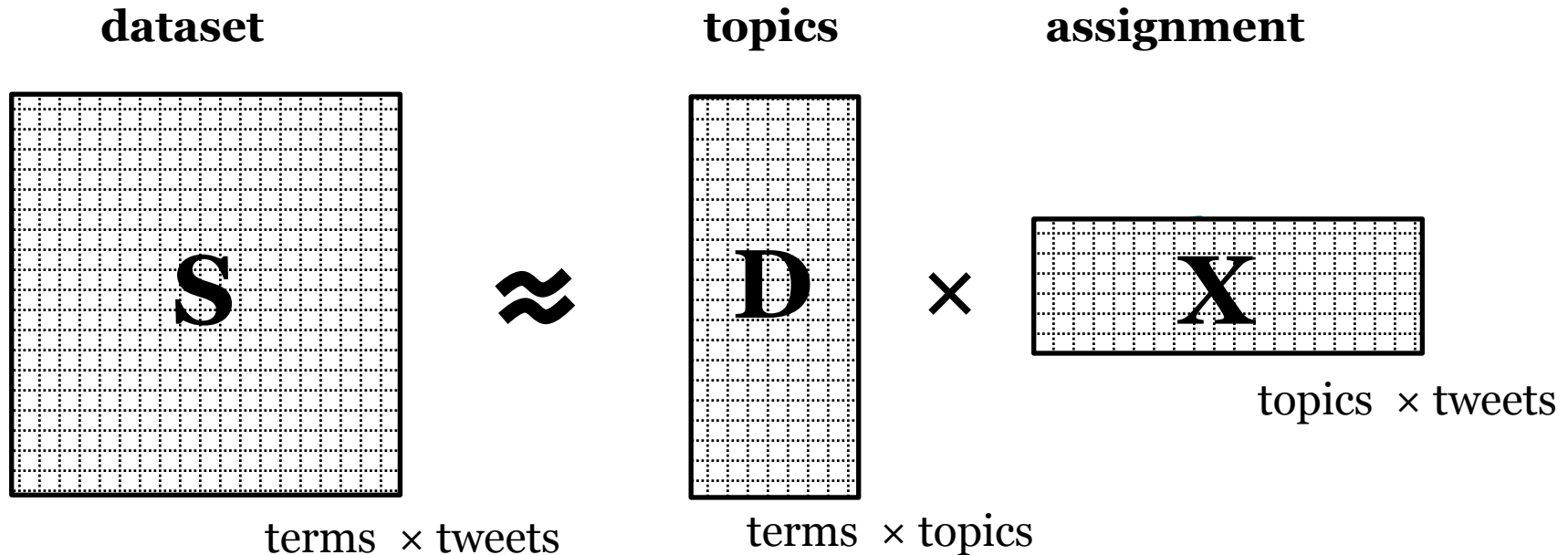
X

$k \times n$

$$(\mathbf{D}, \mathbf{X}) = \arg \min_{\mathbf{D}, \mathbf{X}} \left\| \mathbf{S} - \mathbf{D} \mathbf{X} \right\|_F^2 + \lambda \left\| \mathbf{X} \right\|_1$$

s.t. $\mathbf{X} \geq \mathbf{0}$, $\mathbf{D} \geq \mathbf{0}$, $\|d_i\| = 1$ for $i = \{1, \dots, k\}$

Topic Detection - NMF



$$(\mathbf{D}, \mathbf{X}) = \arg \min_{\mathbf{D}, \mathbf{X}} \|\mathbf{S} - \mathbf{DX}\|_F^2 + \lambda \|\mathbf{X}\|_1$$

s.t. $\mathbf{X} \geq \mathbf{0}$, $\mathbf{D} \geq \mathbf{0}$, $\|d_i\| = 1$ for $i = \{1, \dots, k\}$

Topic Detection - NMF

$$(\mathbf{D}, \mathbf{X}) = \arg \min_{\mathbf{D}, \mathbf{X}} \left\| \mathbf{S} - \mathbf{D} \mathbf{X} \right\|_F^2 + \lambda \left\| \mathbf{X} \right\|_1$$

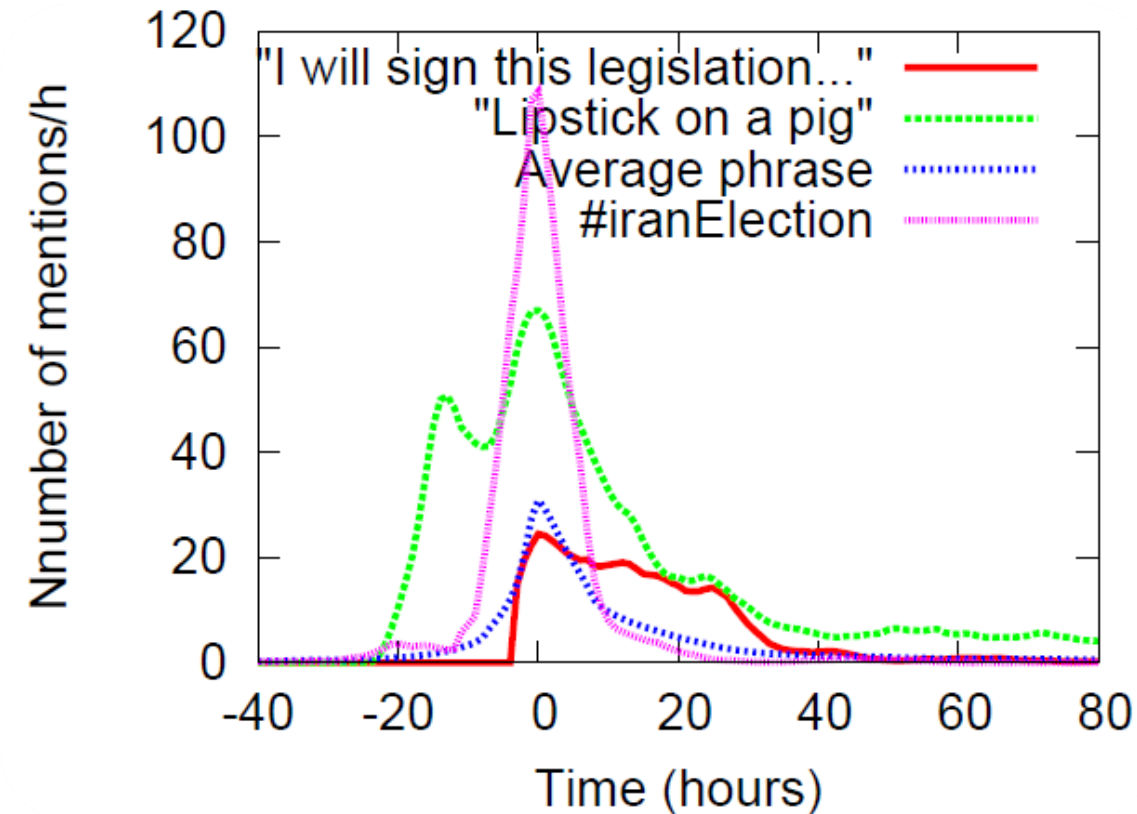
s.t. $\mathbf{X} \geq \mathbf{0}$, $\mathbf{D} \geq \mathbf{0}$, $\|\mathbf{d}_i\| = 1$ for $i = \{1, \dots, k\}$

- Non-convex optimization problem.
 - many local optimum.
- But, if one of the variables, either \mathbf{D} or \mathbf{X} , is known, optimization wrt the other will be convex.
 - Solution:
 - Iteratively optimize the objective function
 - Alternatively optimize wrt \mathbf{D} and \mathbf{X} while holding the other fixed!

Topic Tracking

Topic Tracking

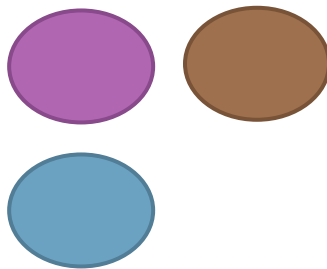
- Smooth evolution of topics through time



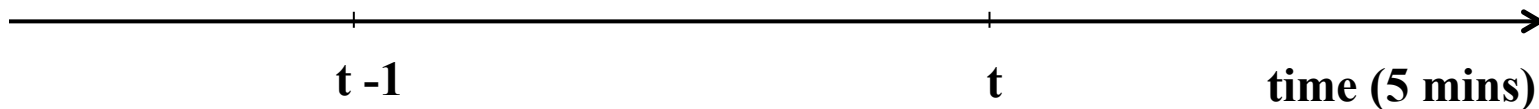
Topic Tracking

- Incremental Clustering

Incoming data at t



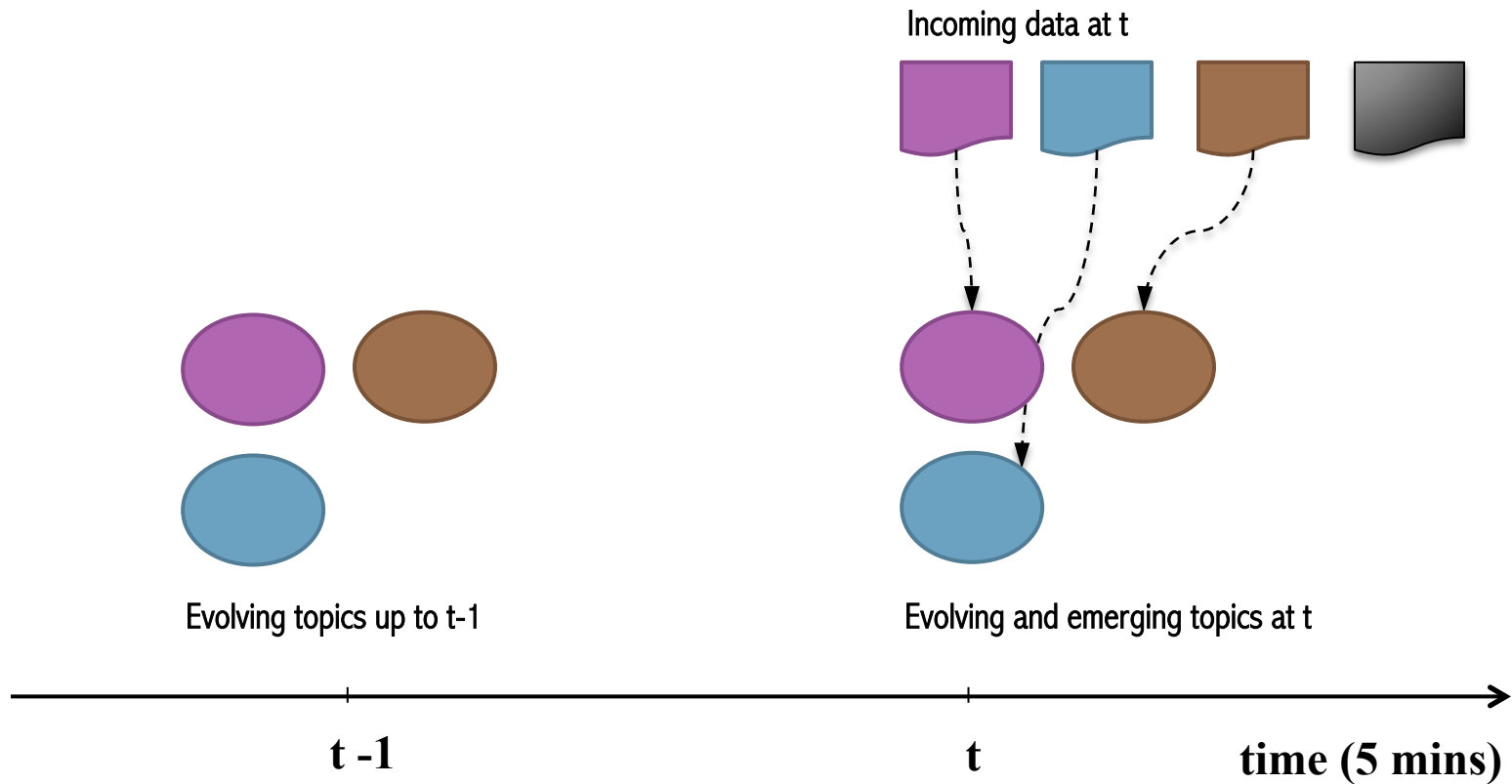
Evolving topics up to $t-1$



Evolving topic: a previously identified topic.
Emerging topic: new topics

Topic Tracking

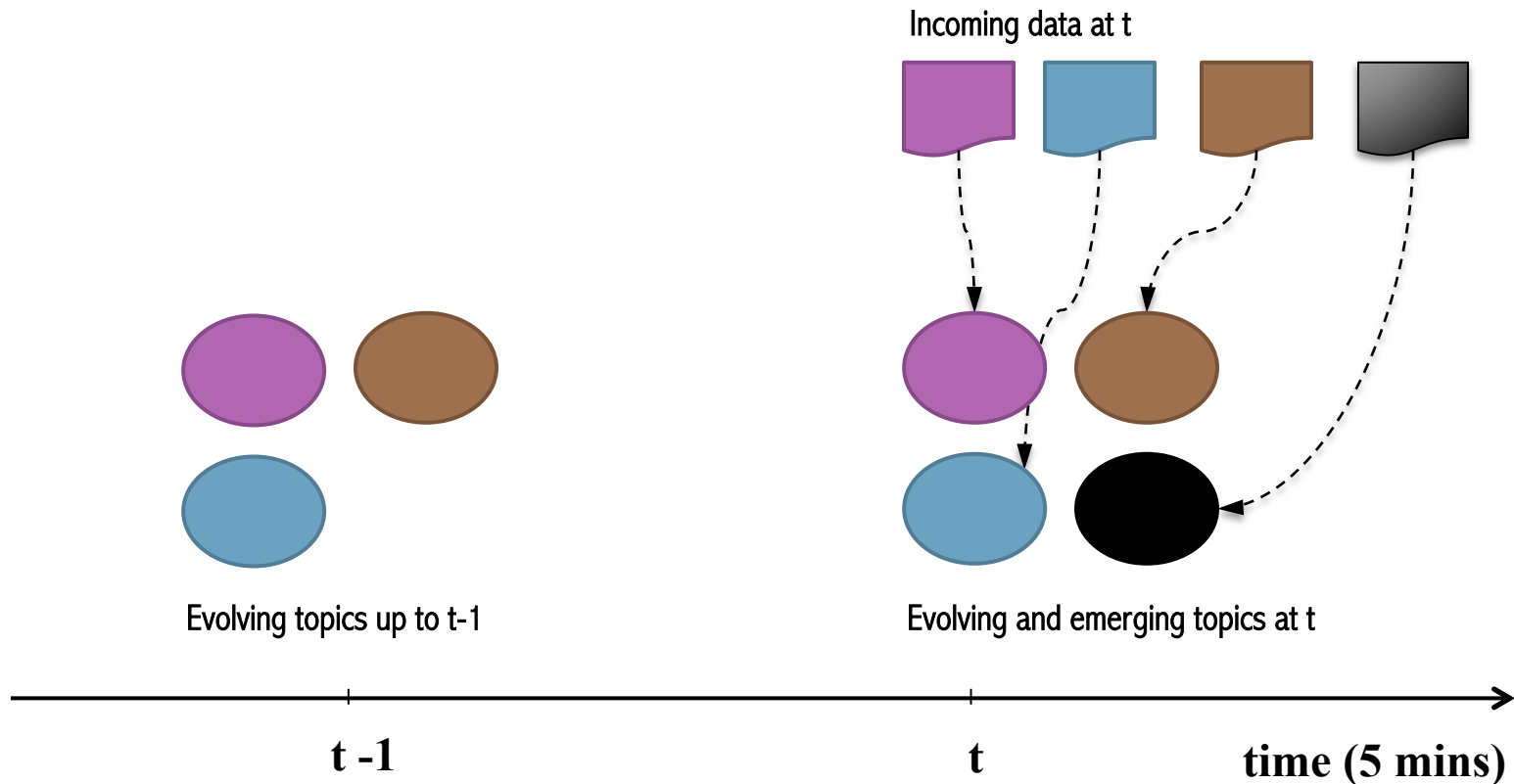
- Incremental Clustering



Evolving topic: a previously identified topic.
Emerging topic: new topics

Topic Tracking

- Incremental Clustering



Evolving topic: a previously identified topic.
 Emerging topic: new topics

Topic Tracking

• Incremental Clustering for Topic Discovery

- Compute similarity btw each incoming tweet and each cluster center.
- If the maximum similarity value is greater than τ , assign the tweet to the cluster and update cluster center.
- Otherwise, generate a new cluster and cluster center.
- *faster* approach: Minhash or LSH

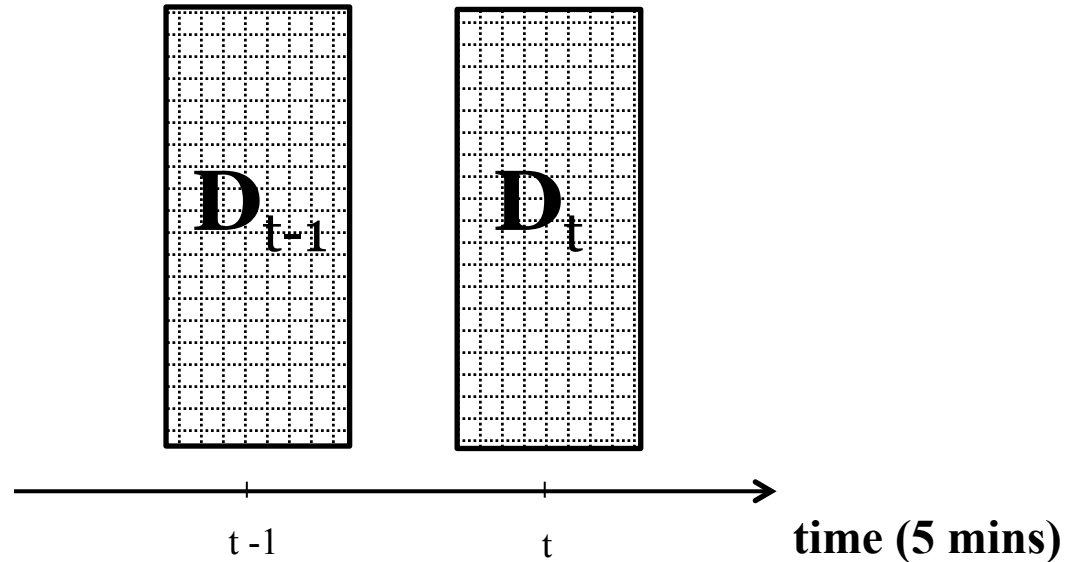
```

1: Input: tweet sets  $D$ , topic cluster set  $C$ , cluster center set  $Center$ , and threshold  $\tau$ .
2: Output: update topic clusters  $C$ , and update cluster centers  $Center$ .
3: Process:
4: if  $C = \emptyset$  then
5:   random select  $N$  tweets from  $D$  and add into  $C$  and  $Center$ .
6: end if
7: initialize  $max$ ,  $tmp_C$ ,  $tmp_{center}$ .
8: for  $d_i \in D$  do
9:   for  $center_j \in Center$  do
10:    compute Cosine Similarity  $sim$  between  $center_j$  and  $d_i$ .
11:    if  $sim > max$  then
12:       $max = sim$ ,  $tmp_C = C_j$ ,  $tmp_{center} = center_j$ .
13:    end if
14:  end for
15:  if  $max > \tau$  then
16:    distribute  $d_i$  to cluster  $tmp_C$ , and update  $tmp_{center}$ .
17:  else
18:    new cluster and centroid and add to  $C$  and  $Center$ .
19:  end if
20: end for
21: return  $C$  and  $Center$ .

```

Topic Tracking

- **Key Idea:** Temporal Coherence, smooth evolution



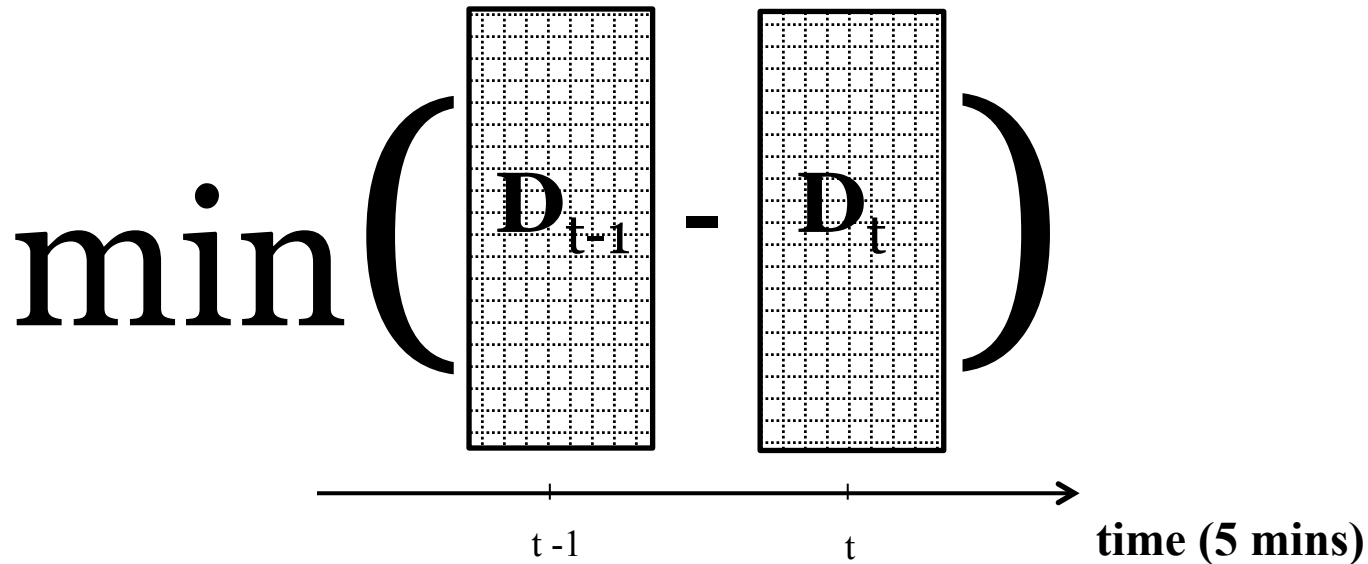
\mathbf{D} at t to be a smooth evolution of \mathbf{D} at $t-1$

No dramatic change in distribution over words for the same **evolving** topic in consecutive time stamps.

The nature of the topic remains the same.

Topic Tracking

- **Key Idea:** Temporal Coherence, smooth evolution



\mathbf{D} at t to be a smooth evolution of \mathbf{D} at $t-1$

No dramatic change in distribution over words for the same **evolving** topic in consecutive time stamps.

The nature of the topic remains the same.

Topic Tracking

$$\mathcal{L}(\mathbf{D}) = \|\mathbf{S} - \mathbf{D}\mathbf{X}\|_F^2 + \lambda \|\mathbf{X}\|_1 + \mu \|\mathbf{D} - \mathbf{D}^{t-1}\|_F^2$$

$$\mathcal{H}[\mathcal{L}(\mathbf{D})] = \mathbf{X}\mathbf{X}^T + 2\mu\mathbf{I}_k \quad \mathbf{D}_{i+1} = P \left[\mathbf{D}_i - \alpha_i \nabla_{\mathbf{D}} \mathcal{L}(\mathbf{D})_{[\mathbf{D}_i, \mathbf{X}]} \right]$$

Algorithm 5.2. Computing \mathbf{D}^t and \mathbf{X}^t at time t , see TL in Figure 4

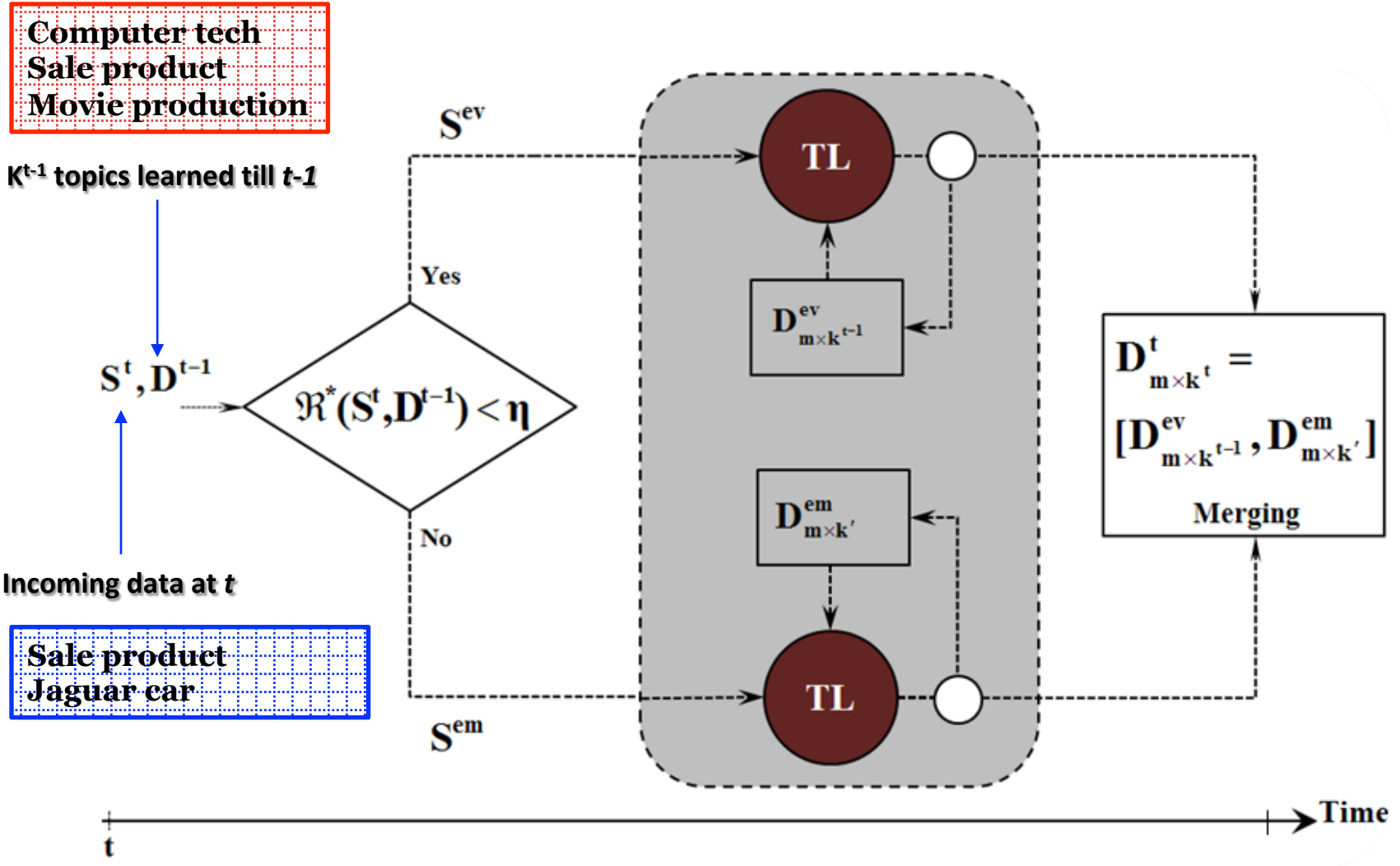
Input: \mathbf{S}^t , \mathbf{D}^{t-1} , itr: number of iterations

Output: \mathbf{D}^t , \mathbf{X}^t

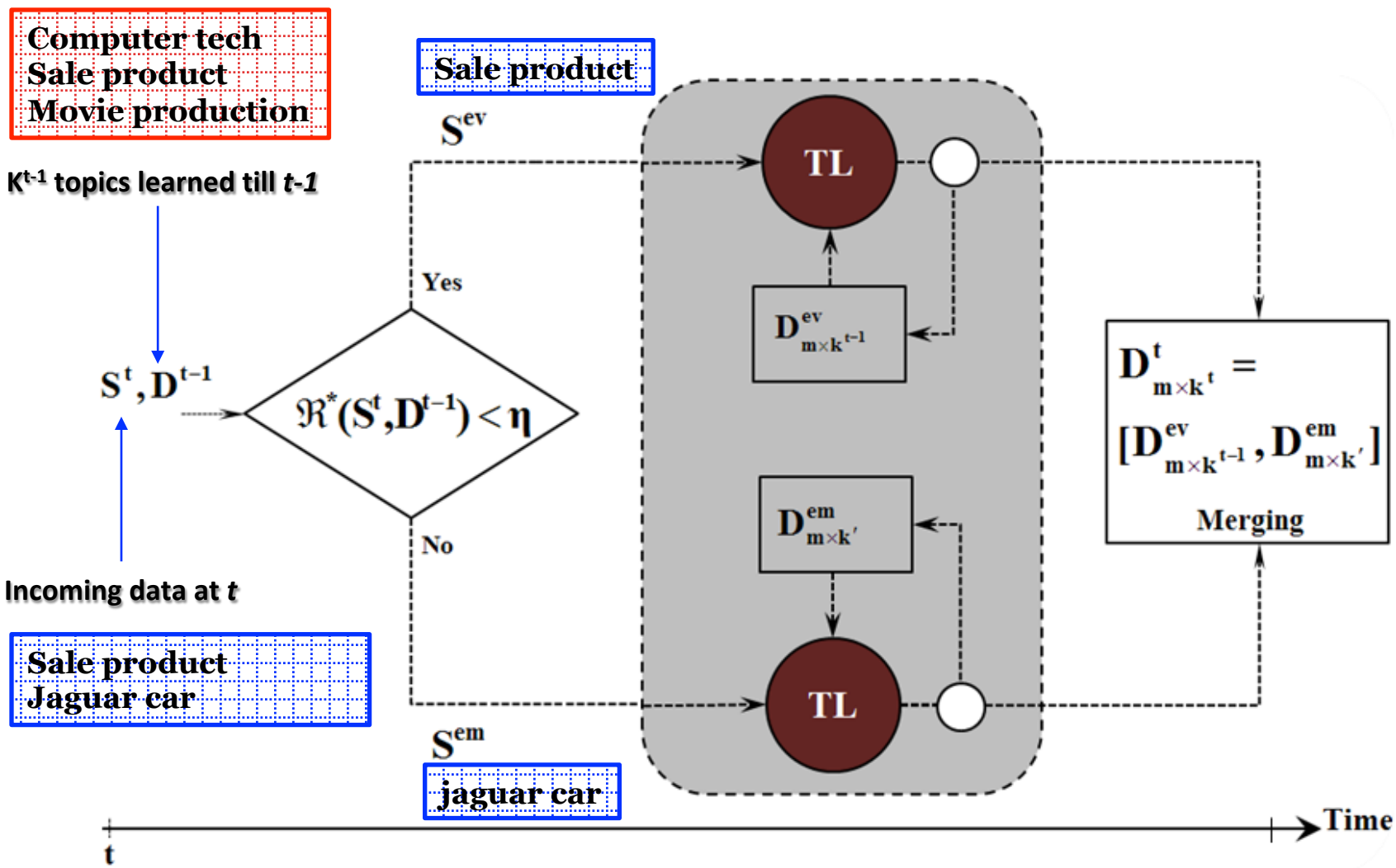
1. Compute \mathbf{X}^t using \mathbf{S}^t and \mathbf{D}^{t-1}
2. $\mathbf{D}_0^t = \mathbf{D}^{t-1}$
3. for $i=1 : \text{itr}$ do
4. compute $\nabla_{\mathbf{D}} \mathcal{L}(\mathbf{D}_{i-1}^t)$
5. $\mathbf{U} = \nabla_{\mathbf{D}} \mathcal{L}(\mathbf{D}_{i-1}^t) \text{diag}^{-1}(\mathcal{H}[\mathcal{L}(\mathbf{D})]_{[\mathbf{X}^t]}) + \mathbf{D}_{i-1}^t$
6. $\mathbf{D}_i^t = \max(\mathbf{0}, \mathbf{U})$
7. end for

[1] Julien Mairal, Francis Bach, Jean Ponce, Guillermo Sapiro: *Online Learning for Matrix Factorization and Sparse Coding*. Journal of Machine Learning Research 11: 19-60 (2010)

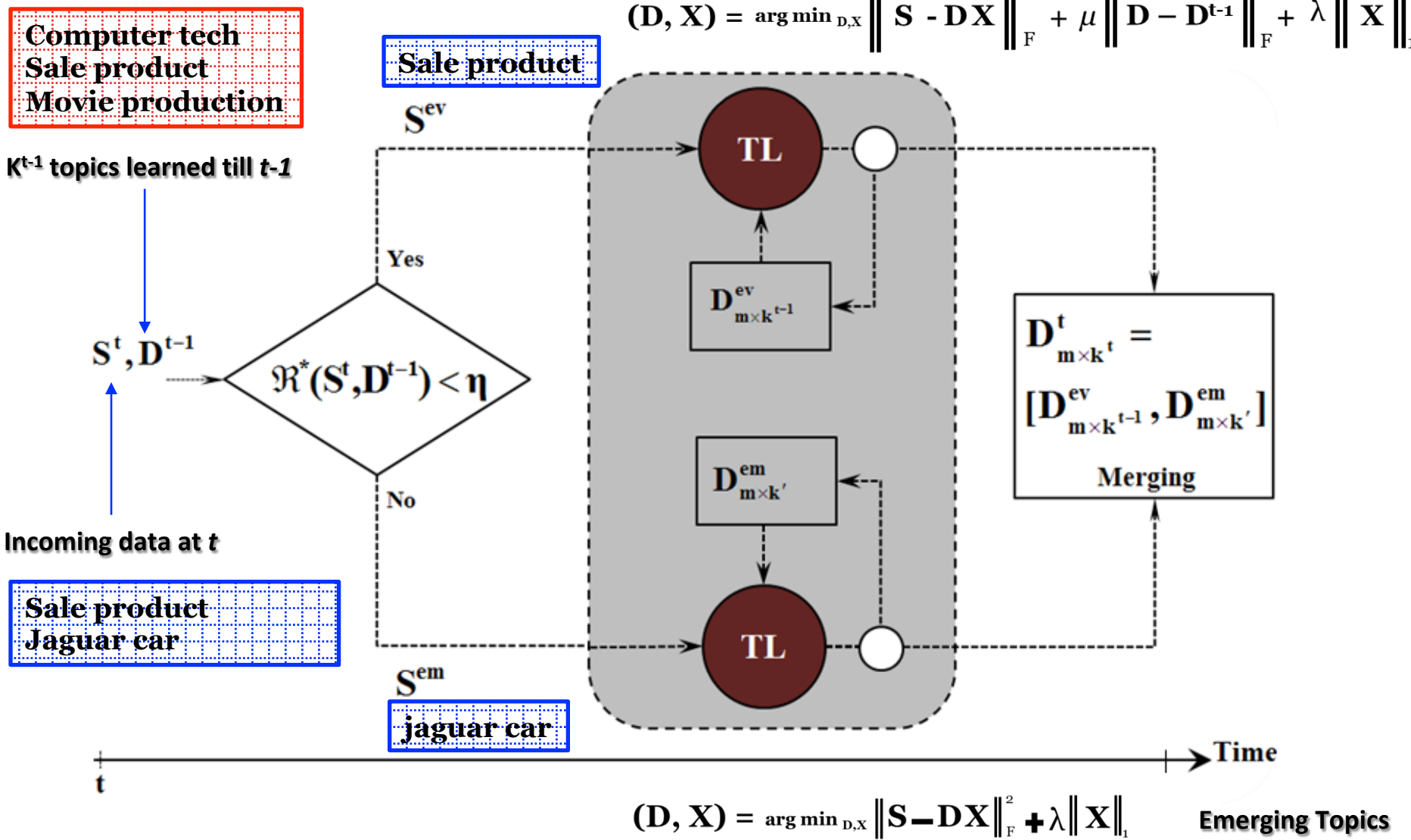
Temporal Coherence



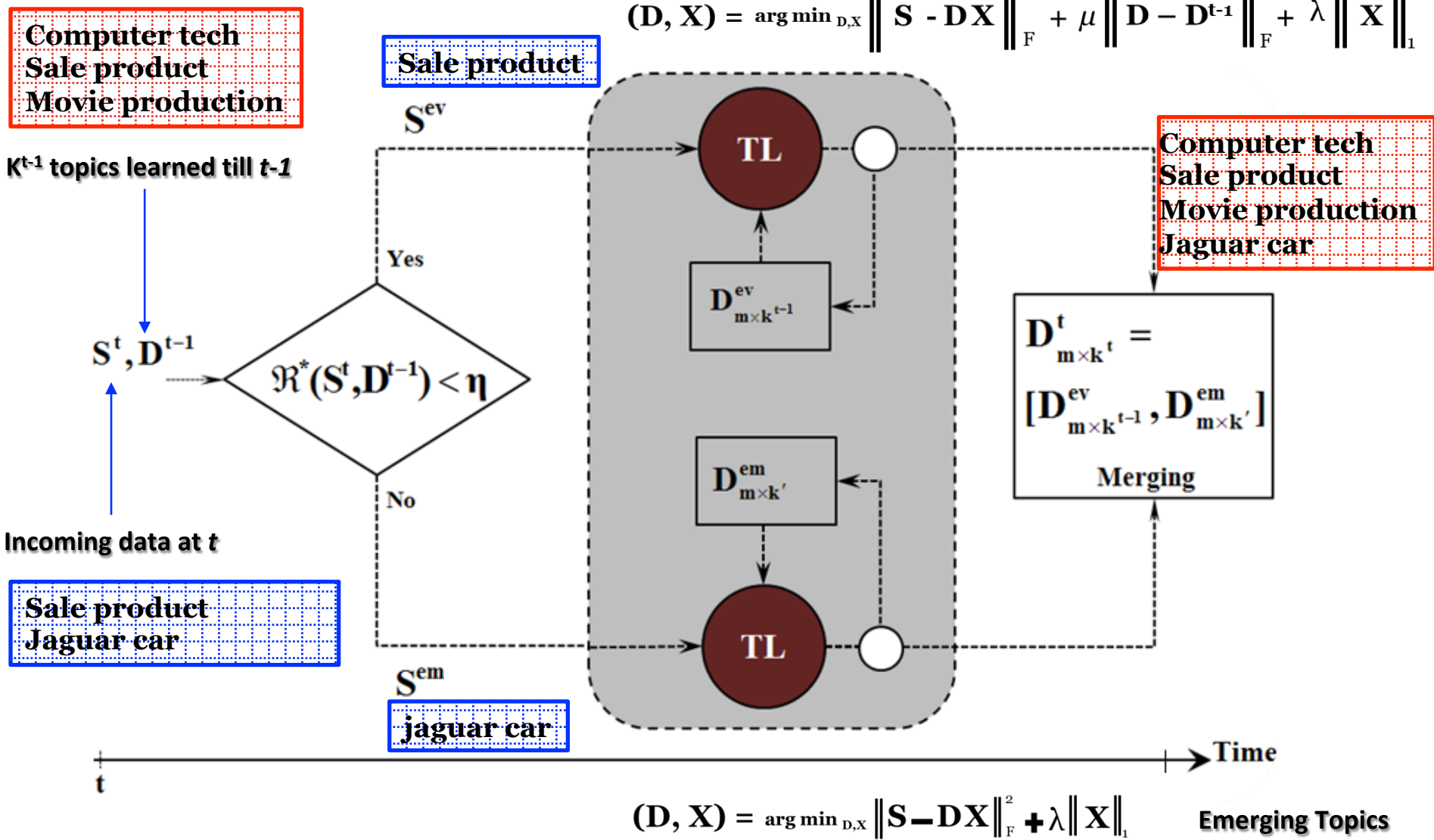
Temporal Coherence



Temporal Coherence



Temporal Coherence



Topic Tracking- Cnt.

- Temporal Coherence constraint for topic learning:
 - \mathbf{D}^{ev} to be a smooth evolution of \mathbf{D}^{t-1}

$$(\mathbf{D}, \mathbf{X}) = \arg \min_{\mathbf{D}, \mathbf{X}} \left\| \mathbf{S} - \mathbf{D}\mathbf{X} \right\|_F^2 + \lambda \left\| \mathbf{D} - \mathbf{D}^{t-1} \right\|_F^2 + \lambda \left\| \mathbf{X} \right\|_1$$

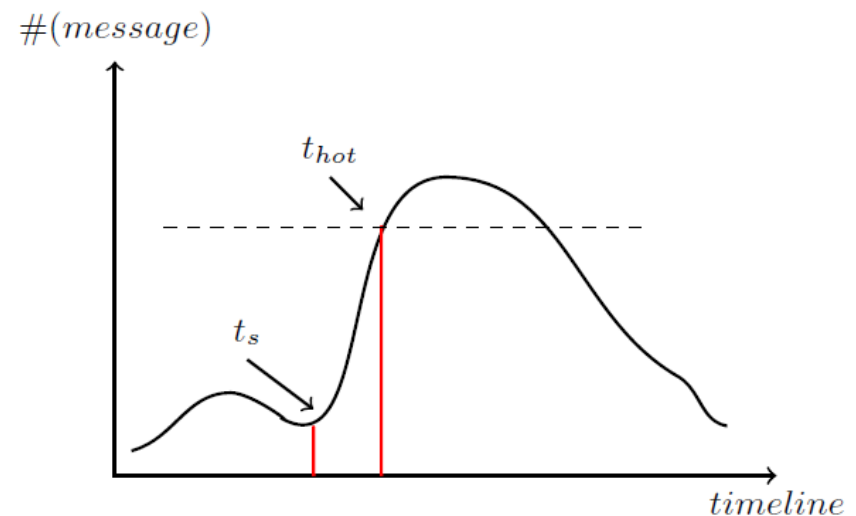
s.t. $\mathbf{X} \geq \mathbf{0}$, $\mathbf{D} \geq \mathbf{0}$, $\|\mathbf{d}_i\| = \mathbf{1}$ for $i = \{1, \dots, k\}$

- Can be solved efficiently
 - Space: $O(n*m)$, given that $m \gg k$
 - Running time: $O(n)$

Early Detection of Emerging Topics

Early Detection of Topics

- Evolution of a hot topic
 - t_s topic detection time
 - t_{hot} the time by which topic becomes major.
 - tweets number exceeds a threshold.
- We aim to predict if an already-detected topic will be major in the near future!



Early Detection

View 1: rate indicators

- Rate of increase in #users
- Rate of increase in #tweets
- Rate of increase in #re-tweets

View 2: overlap indicators

- Overlap btw users posted about topic and influential users
- Overlap btw topic keywords and top influential keywords

Co-training (Co): Two SVM classifiers trained on the above two orthogonal views of features

Ensemble Learner (En): Ensemble of three classifiers (Decision Tree, SVM, and Naive Bayesian) vote for each unlabeled topic.

Early Detection

- User authority / user influence against the topic
- Tweet authority / derived from topical user auth.

- f_1 is the rate of increase of user number,

$$f_1 = \frac{|U^t|}{\sum_{x=0}^t \frac{1}{t-x+1} |U^x|}. \quad (6)$$

- f_2 is the rate of increase of tweets number,

$$f_2 = \frac{|Tw^t|}{\sum_{x=0}^t \frac{1}{t-x+1} |Tw^x|}. \quad (7)$$

- f_3 is the rate of increase of re-tweets number,

$$f_3 = \frac{|Rtw^t|}{\sum_{x=0}^t \frac{1}{t-x+1} |Rtw^x|}. \quad (8)$$

- f_4 is the overlap between org keyusers and top N influential topic users,

$$f_4 = \frac{\#(ku_{tp} \cap ku)}{\#ku_{tp}}. \quad (9)$$

- f_5 is the overlap between org keywords and top N influential topic keywords, and

$$f_5 = \frac{\#(kw_{tp} \cap kw)}{\#kw_{tp}}. \quad (10)$$

- f_6 represents the rate of increase of influence of the accumulated weight of tweets,

$$f_6 = \frac{|A^t|}{\sum_{x=0}^t \frac{1}{t-x+1} |A^x|}, \quad (11)$$

where $A = \frac{\sum_{tw \in Tw_{tp}} auth_{tp}(tw)}{|Tw_{tp}|}$.

Evaluation

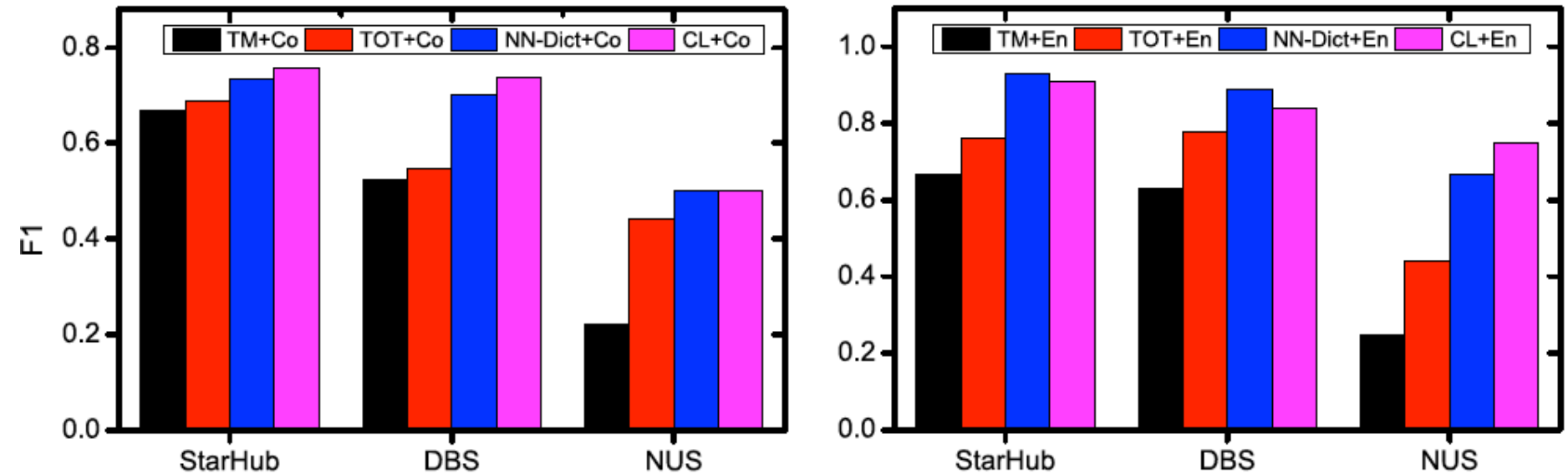


Figure 5: Performance of emerging topic Detection when $T_L = t_{hot}$

CL: Incremental clustering

Co: Co-training

En: Ensemble Learner

Evaluation

Table 2: Performance of emerging topic detection when $T_L = t_{hot}$

Methods	Organization	recall	precision	F_1
CL+En	<i>StarHub</i>	0.93	0.87	0.90
CL+TSVM		0.86	0.75	0.80
CL+Semi-NB		0.86	0.71	0.77
CL+En	<i>DBS</i>	0.89	0.80	0.84
CL+TSVM		0.89	0.73	0.80
CL+Semi-NB		0.89	0.67	0.70
CL+En	<i>NUS</i>	1.00	0.60	0.75
CL+TSVM		1.00	0.50	0.67
CL+Semi-NB		1.00	0.42	0.73

CL: Incremental clustering

Co: Co-training

En: Ensemble Learner

Evaluation

Table 3: Performance of emerging topic detection when $T_L = t_{mid}$

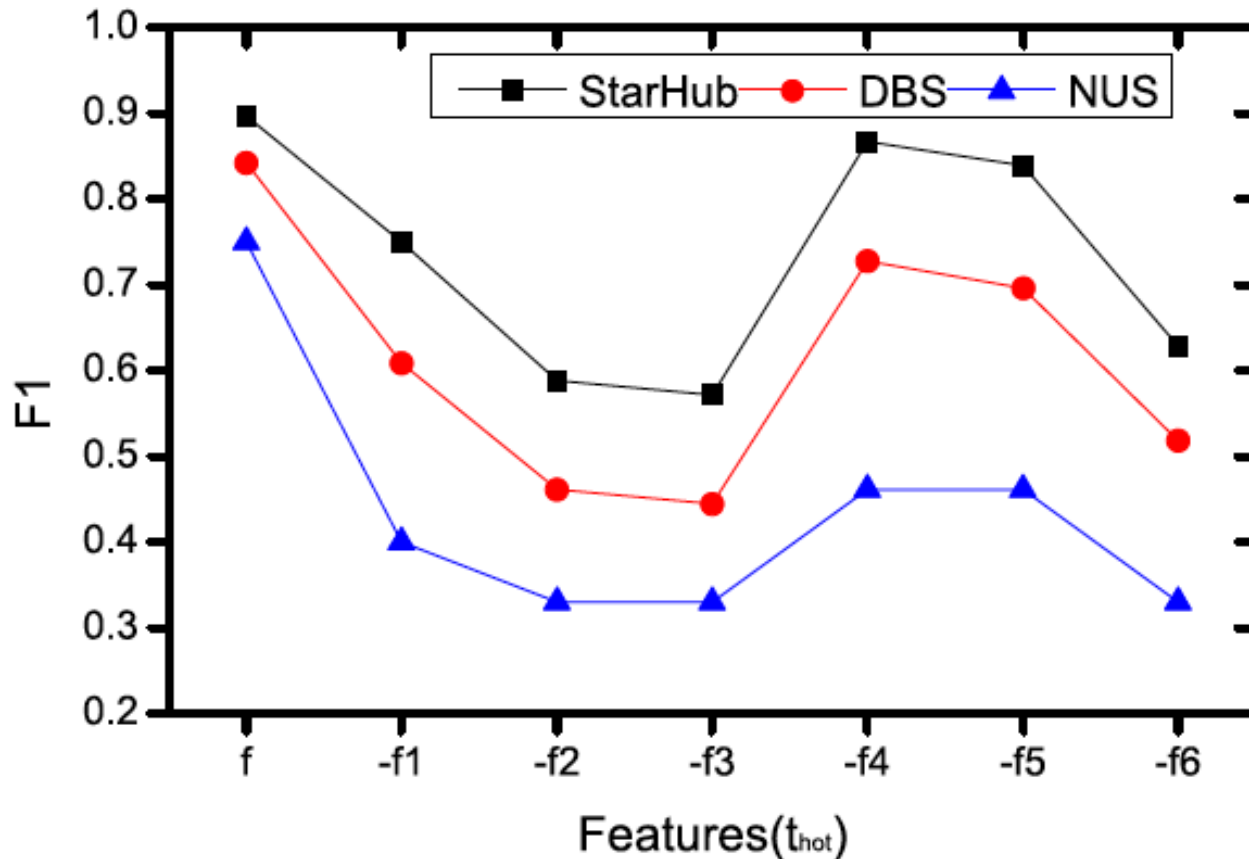
Methods	Organization	recall	precision	F_1
CL+En	<i>StarHub</i>	0.71	0.83	0.77
CL+TSVM		0.71	0.71	0.71
CL+Semi-NB		0.71	0.67	0.69
CL+En	<i>DBS</i>	0.78	0.78	0.78
CL+TSVM		0.78	0.70	0.74
CL+Semi-NB		0.78	0.64	0.70
CL+En	<i>NUS</i>	0.67	0.50	0.57
CL+TSVM		0.67	0.40	0.50
CL+Semi-NB		0.67	0.40	0.50

CL: Incremental clustering

Co: Co-training

En: Ensemble Learner

Evaluation



f2: rate of increase in #tweets

f3: rate of increase in #re-tweets,

f6: overall accumulated influence of tweets

Summary

- Effective NMF model with temporal coherence constraint
 - Improves topic tracking in streaming data.
- Effective framework for early prediction of emerging topics.
 - Rate and overlap features

Reading

- Emerging topic detection for organizations from microblogs. Chen, Y., et al. SIGIR'13.
- Learning evolving and emerging topics in social media. Saha, A. et al. WSDM'12
- Community detection in social networks considering topic correlations. Wang, Y., et al. AAAI'19.