

WWW勉強会

佃 光撰（産総研）

2016/6/25

Tweet Properly: Analyzing Deleted Tweets to Understand and Identify Regrettable Ones

Lu Zhou, Wright State University

Wenbo Wang, Wright State University

Keke Chen, Wright State University

- **投稿したことを後悔して削除されたツイートを分析**
 - 削除された4,000ツイートを人手で分析
 - 削除されたツイートの18%は投稿を後悔したことが理由
 - 後悔理由を人手で10カテゴリに分類 (e.g. alcohol, drug, violence)

- **ツイートが後悔して削除されるか否かの**分類器を構築****
 - Content-based features + personalized history-based features
 - 精度0.78を達成
 - ツイート投稿前に警鐘を鳴らす

- **先行研究**

- SNS上で後悔を理由に削除された投稿の**種類や理由**を分析

- **本研究**

- 後悔を理由に削除された投稿を**自動的に特定**

ユーザ数	30,000
総ツイート数	17,587,816
総非削除ツイート数	14,325,871
総削除ツイート数	3,261,945
1回以上ツイートを削除したユーザ数	26,543
1回以上ツイートを投稿したユーザ数	28,778

- リツイートは含まない
- タイポや言い回しの修正のために削除したツイートは含めない
 - 削除直後のツイートとの編集距離からタイポや言い回しの修正かどうかを判定

削除された4,000ツイートをランダムに選択し以下の2種類に分類

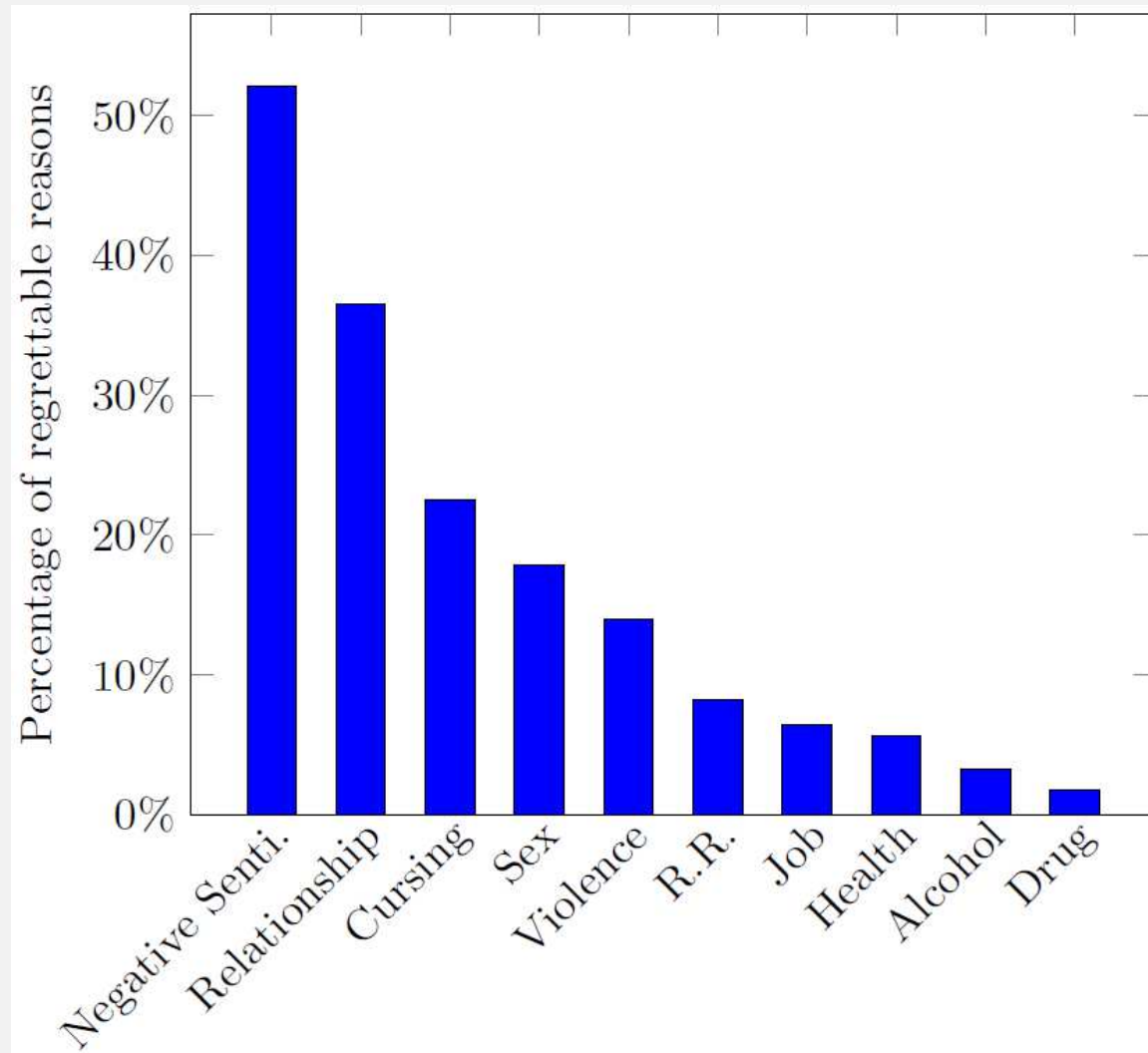
● **Regrettable tweet**

- 何かしらに後悔したことが理由で削除されたツイート
- 更に10種類のカテゴリに分類（ひとつのツイートの複数カテゴリを許す）
negative sentiment, cursing, sex, alcohol, drug, violence, health, racial and religion, job, relationship

● **Unsure tweet**

- 削除の理由が後悔ではないツイート
- 例 1 : my birthdays in 5 more days
- 例 2 : Yay! My cousins are visiting me next weekend for the @Dodgers game!

削除されたツイートの**18%**は何かしらに後悔したことが理由
(Facebookでは削除された投稿の6%が削除に値する理由有り)



ツイートごとに以下の2種類の素性を求める

● Content-based (10次元)

- 10カテゴリそれぞれの特徴語を含む辞書を作成 (WordNet等使用)
- ツイートに含まれるカテゴリのみ 1 を持つベクトル作成
- 例: **職場**に行く**と頭痛**がする → (0,0,0,0,0,0,1,0,1,0)
job health

● User-based (12次元)

- ユーザ毎の10カテゴリそれぞれのツイートが削除された割合
 - 例: u_1 のjobに関する削除ツイート数 / u_1 のjobに関する総ツイート数
- 上記の10カテゴリをまとめて考慮したときの削除割合
- カテゴリを考慮しない場合の削除割合

● データセット

- 5,000ユーザから4ツイートずつ収集
- 正例：削除されたツイート2件
- 負例：削除されなかったツイート2件
- 正例、負例ともに10カテゴリの少なくともひとつに属するツイート

● 分類手法

- Naive Bayes
- SVM
- J48
- AdaBoost

分類精度

		Content-only	Content+User-history
NB	Precision	0.552 ± 0.031	0.775 ± 0.046
	Recall	0.349 ± 0.078	0.486 ± 0.055
	F1-Score	0.427 ± 0.043	0.598 ± 0.035
SVM	Precision	0.536 ± 0.041	0.753 ± 0.042
	Recall	0.478 ± 0.049	0.626 ± 0.054
	F1-Score	0.505 ± 0.041	0.683 ± 0.034
J48	Precision	0.537 ± 0.027	0.711 ± 0.072
	Recall	0.593 ± 0.030	0.716 ± 0.081
	F1-Score	0.563 ± 0.019	0.714 ± 0.081
AdaBoost	Precision	0.541 ± 0.055	0.731 ± 0.048
	Recall	0.434 ± 0.077	0.696 ± 0.068
	F1-Score	0.482 ± 0.048	0.713 ± 0.055

影響の大きい素性ランキング

	Content-only	Content+User-history
1	job	total deletion ratio
2	R.R.	total regrettable deletion ratio
3	sex	negative senti. deletion ratio
4	drug	cursing deletion ratio
5	alcohol	sex deletion ratio
6	negative senti.	R.R. deletion ratio
7	health	job deletion ratio
8	violence	alcohol deletion ratio
9	relationship	violence deletion ratio
10	cursing	drug deletion ratio

- **Content-basedな素性はほぼ効果なし**

- 後悔が理由で削除されたツイートから辞書を作成したわけではないので、この素性で分類しようとしたこと自体無理があったのでは

- **全素性の中で最も有効な素性はtotal deletion ratio**

- よく削除するユーザのツイートを全部削除に分類することで精度が上がっているだけ？

Mining Aspect-Specific Opinion using a Holistic Lifelong Topic Model

Shuai Wang, University of Illinois

Zhiyuan Chen, University of Illinois

Bing Liu, University of Illinois

- レビューから以下の**4要素を抽出するモデル**を提案

- 観点 (e.g. screen)
- 意見 (e.g. clear, great)
- 意見の極性 (positive or negative)
- 意見が観点到に依存か非依存か (e.g. clearはscreen依存、greatは非依存)

- **2種類のモデル**を提案

- JASTモデル：対象レビューのドメイン (e.g. テレビ) の知識のみ使用
- LASTモデル：他ドメインの知識も利用してJASTモデルを強化

- **先行研究**

- 4要素のうちの一つのみ考慮したモデルを提案

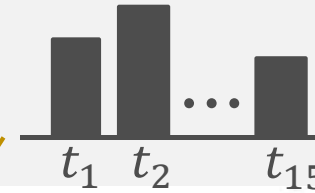
- **本研究**

- 4要素の全てを同時に考慮した初めてのモデルを提案
- 同時に考慮することで要素間の相関も考慮したより良いモデルになる

Posの観点非依存意見語の分布



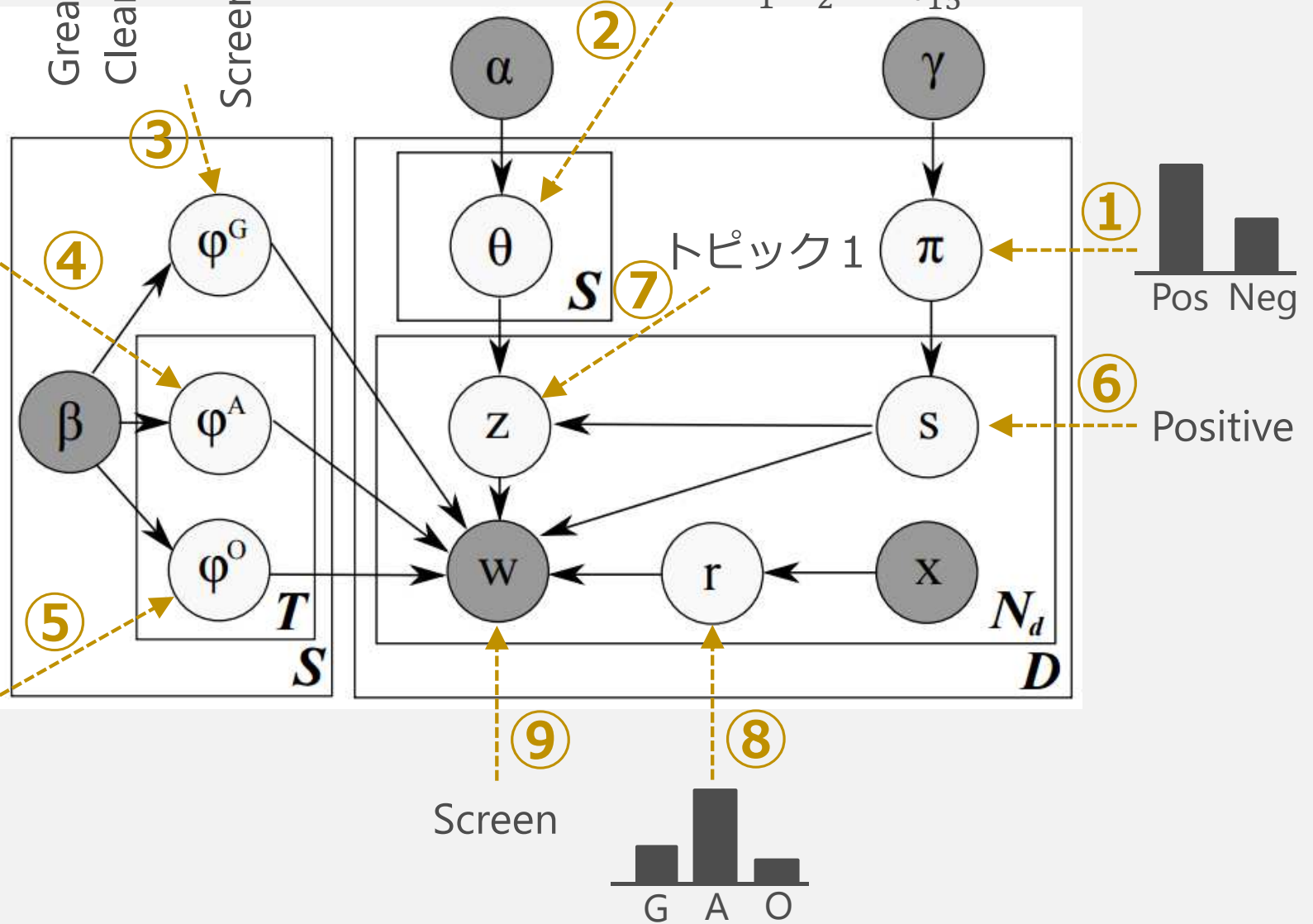
Posのトピック分布

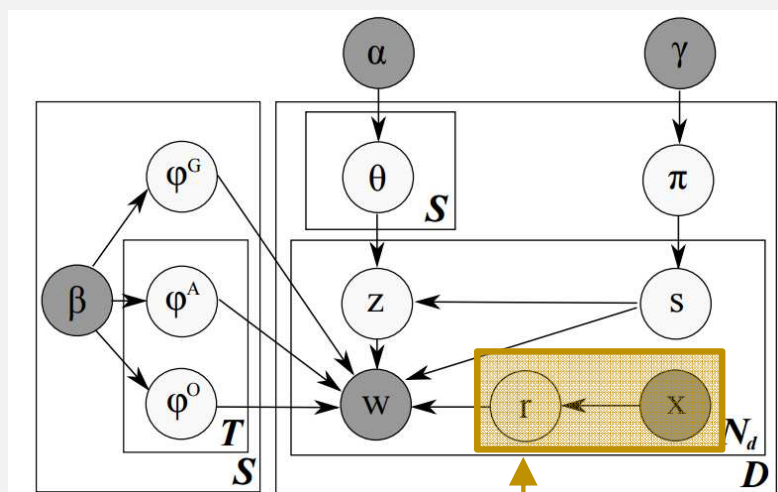


Posのトピック1の観点語の分布



Posのトピック1の観点依存意見語の分布





- G : 観点非依存の意見語 (e.g. great)
- A : 観点 (e.g. screen)
- O : 観点依存の意見語 (e.g. clear)

- Webにある**意見語辞書**を利用

- wが辞書に含まれる : $x=1$
- wが辞書に含まれない : $x=0$

- **$x=1$** なら**GとOの確率が高くなり**
 $x=0$ なら**Aの確率が高くなるように**
Gibbs Samplingに組み込む

● JASTモデルの問題点

- | |
|----------------|
| 観点非依存の意見語が観点依存 |
| 観点依存の意見語が観点非依存 |

と推定される（特に低頻度の語）
- 低頻度の語の観点依存性を数少ない共起関係から無理やり推定するため

● LASTモデルの解決策

- 他のドメインのJASTモデルの推定結果を取り入れることで精度改善
- 観点語の確率分布が十分に似たドメインから取り入れる
- PCドメインではタッチスクリーン対応PCが少ないので意見「smooth」が観点「screen」依存ではなく観点非依存と推定される
- スマホドメインではsmoothはscreen依存の意見と推定できる
- スマホドメインの結果をPCドメインに取り入れてsmoothを観点依存に

青字 : 観点非依存の意見語
 赤字 : 不正解

Battery (Negative)			Shipping&Order (Positive)		
LAST	JAST	ASUM	LAST	JAST	ASUM
die	old	<i>problem</i>	new	free	<i>great</i>
dead	die	hot	free	<i>happy</i>	<i>good</i>
short	fail	<i>bad</i>	fast	fast	quickly
drain	<i>suck</i>	die	quick	<i>pleased</i>	<i>well</i>
fail	useless	<i>original</i>	refund	refund	<i>love</i>
old	hassle	old	promptly	<i>recommend</i>	<i>perfect</i>
hassle	<i>bad</i>	<i>new</i>	original	new	<i>nice</i>
<i>wrong</i>	<i>concern</i>	<i>long</i>	correct	works	<i>perfectly</i>
useless	bother	break	works	quick	new
<i>complain</i>	<i>nervous</i>	<i>hate</i>	accurate	promptly	fast

Battery			Shipping&Order		
LAST	JAST	ASUM	LAST	JAST	ASUM
battery	battery	charge	order	arrive	<i>screen</i>
charge	charge	battery	receive	receive	receive
hour	life	recharge	arrive	order	arrive
life	hour	<i>iphone</i>	shipping	purchase	order
power	<i>device</i>	<i>sd</i>	ship	<i>expect</i>	<i>privacy</i>
charger	cable	<i>card</i>	today	send	cost
recharge	<i>phone</i>	receive	delivery	ship	money
<i>night</i>	<i>ipad</i>	replacement	usual	shipping	<i>monitor</i>
outlet	power	<i>purchase</i>	<i>expect</i>	back	purchase
aaa	plug	<i>star</i>	<i>manner</i>	<i>seller</i>	<i>seller</i>

Battery、Shipping&Order : 人が名付けたトピック名

- LASTモデルの観点依存の意見語のPrecision@10は0.6程度
- LASTモデルの観点語のPrecision@10は0.8程度