

ソーシャルメディア上のテキスト情報を 考慮したインフルエンサー検出モデル

東北大学大学院 経済学研究科 博士課程
五十嵐 未来

目次

1. はじめに
2. トピックモデル
3. 分析モデル
4. 分析結果
5. シミュレーション
6. Appendix

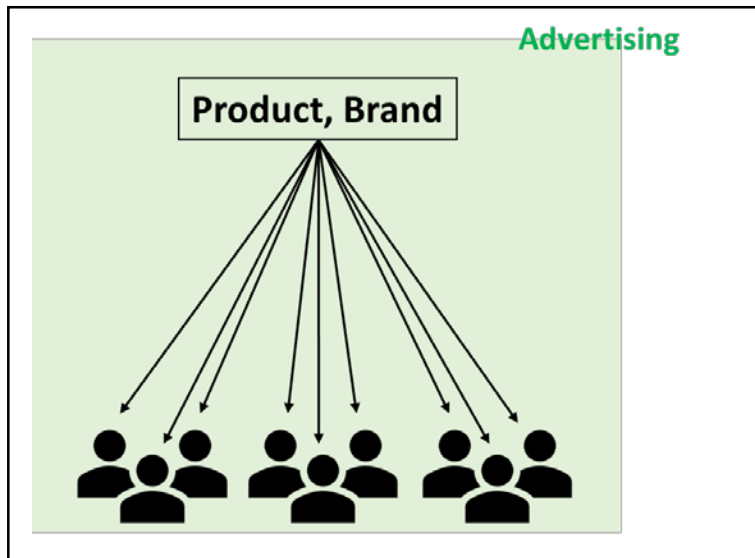
1.はじめに

バイラル・マーケティング

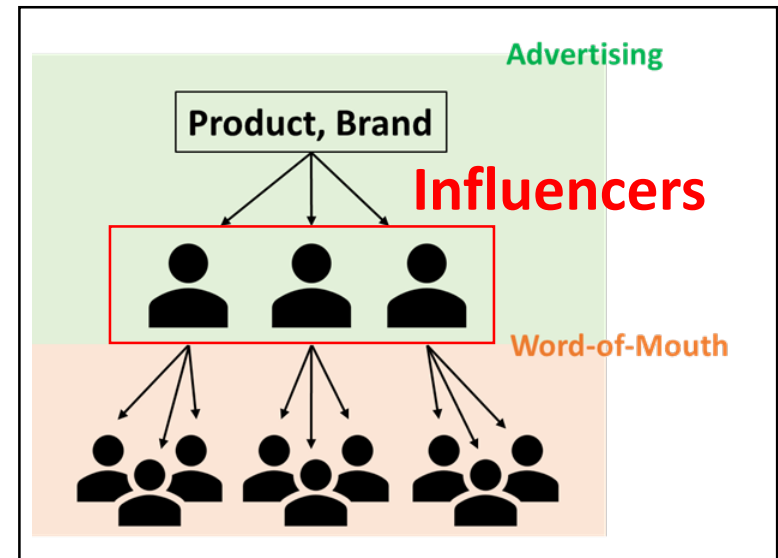
バイラル・マーケティング:

ユーザーの口コミを通じて商品やブランドに関する情報を拡散させるマーケティング手法

身近な知り合いによる口コミであるため、伝統的なマーケティング施策よりも強い影響力が期待される (Zhang et al. 2018)



一般的なプロモーション



Viral Marketing の概念図

関連研究

インフルエンサーの定義に関して、マーケティングの分野では大きく二つの考え方が存在する

- 他者への影響力を考慮

Chae et al. (2016):

インフルエンサーの口コミが持つ一般ユーザーの口コミへの影響

同商品に対しては促進、異ブランドや異カテゴリ商品に対しては減衰の効果を持つ

Gong et al. (2017):

インフルエンサーによるRT(情報の拡散)がテレビの閲覧数を増加させる効果を持つ

- ソーシャルネットワークのモデル化

Chen et al. (2017):

重み付きネットワークによるモデル化と重要度の推定

Bampo et al. (2008):

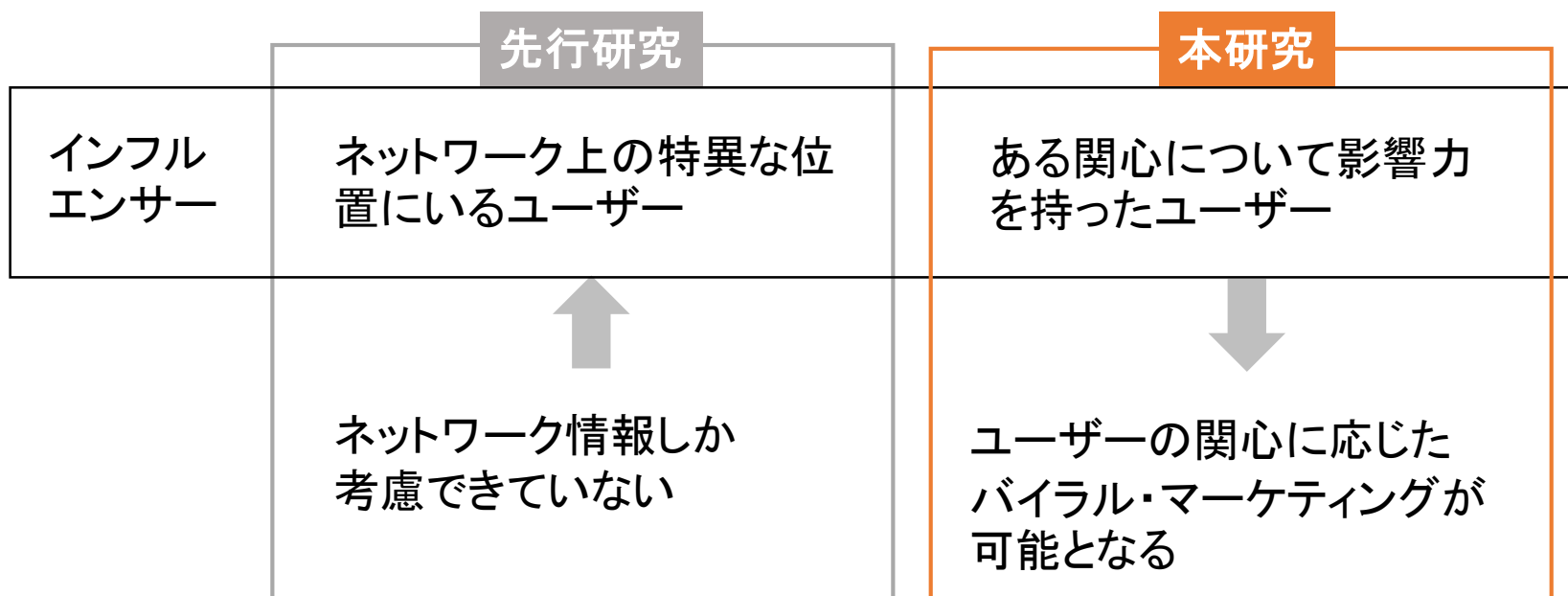
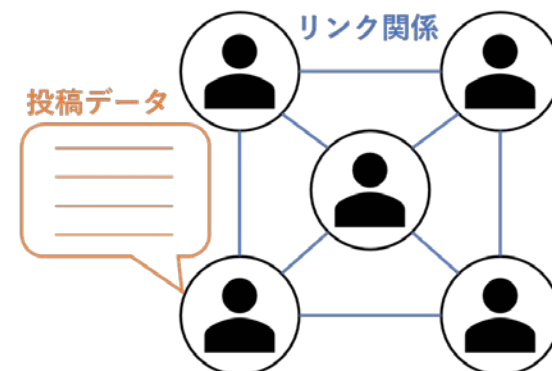
ランダムグラフとSIR型行動によるモデル化



先行研究ではユーザーがどのような関心を持っているかを考慮していない

本研究の特色

- ソーシャルネットワークのリンク関係
 - ソーシャルメディアに投稿されたテキスト情報
- トピックモデルによってモデル化
“どのような関心によってリンクしているのか”



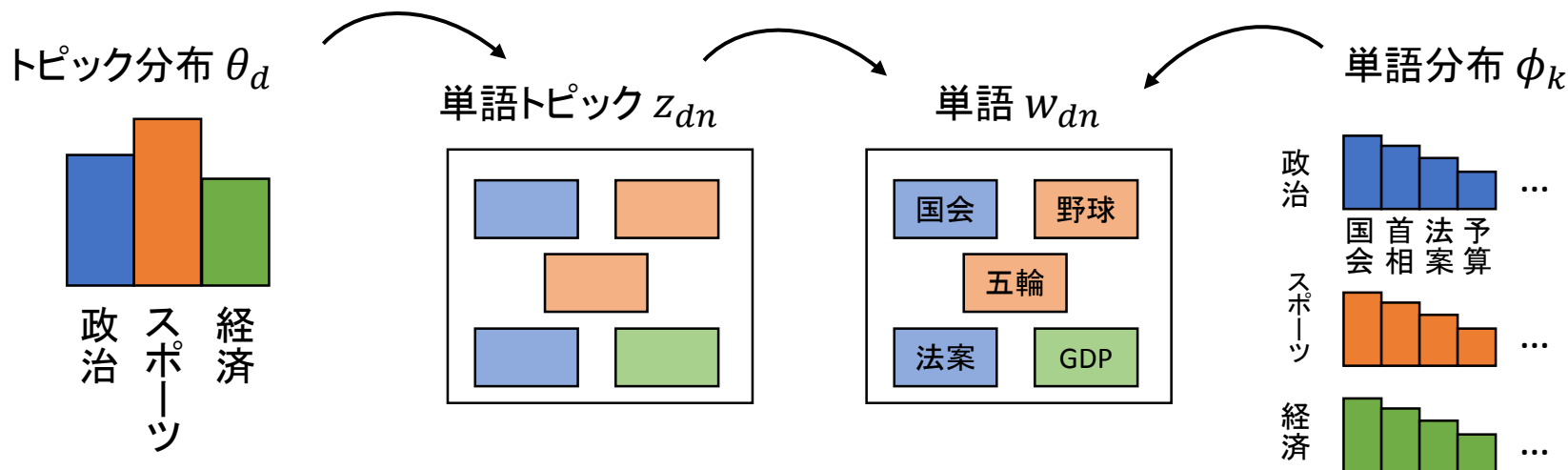
本研究の目的

1. トピックモデルによりリンク関係と投稿データを同時にモデル化し、ユーザー同士がどのような関心に応じて繋がっているのかを把握する
2. 拡散したい情報とユーザーの関心をマッチさせたバイラル・マーケティングを可能にする

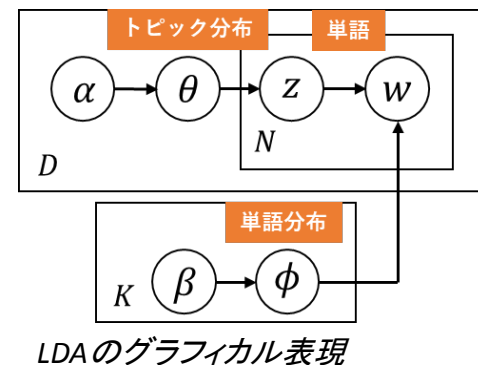
2. トピックスモデル

潜在ディリクレ配分法 (LDA)

Blei et al. 2003a



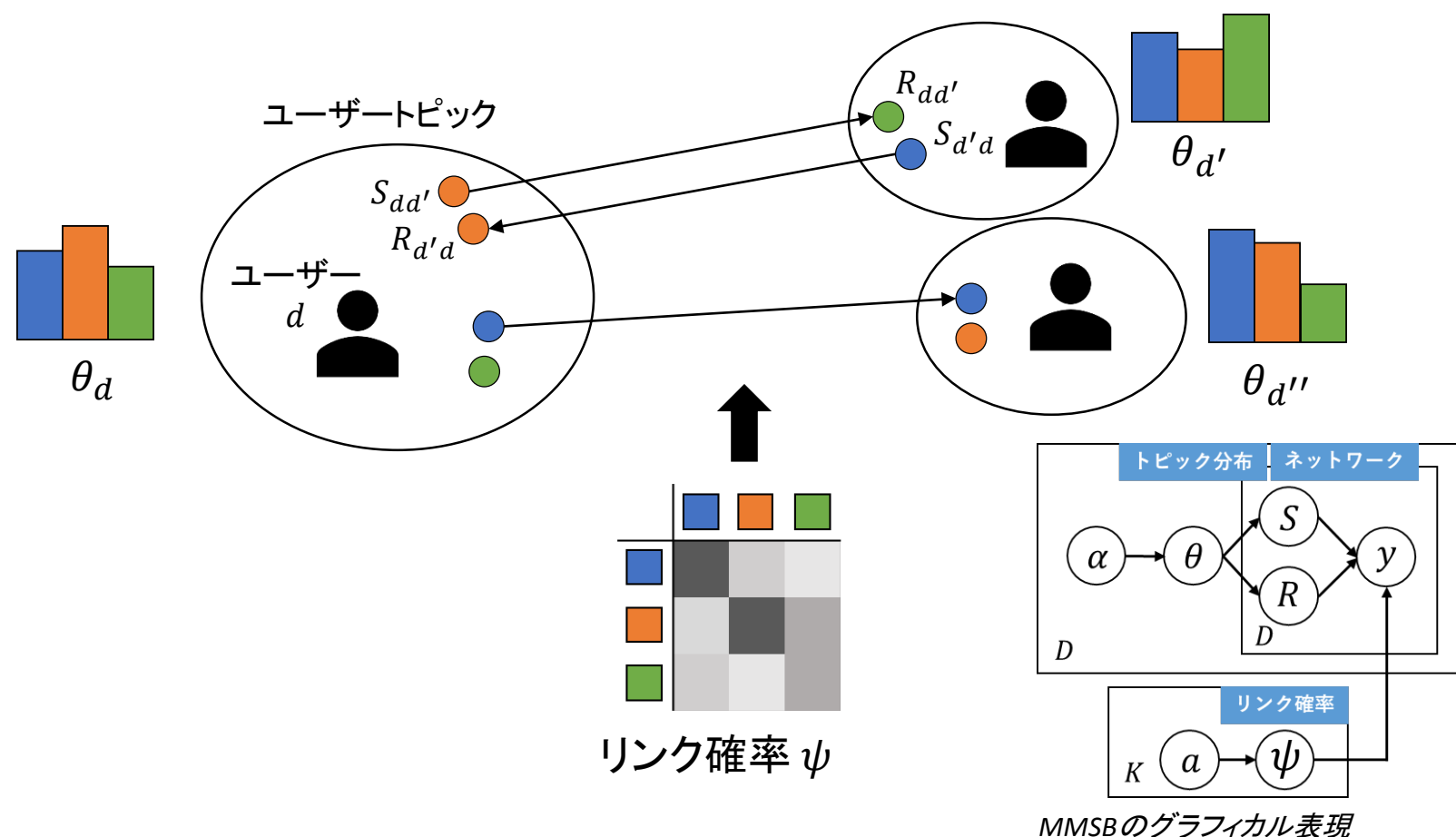
$$\begin{aligned}
 p(W|\Theta, \Phi) &= \prod_{d=1}^D \prod_{n=1}^{N_d} \sum_{k=1}^K p(z_{dn} = k | \theta_d) p(w_{dn} | \phi) \\
 &= \prod_{d=1}^D \prod_{n=1}^{N_d} \sum_{k=1}^K \theta_{dk} \phi_{kw_{dn}}
 \end{aligned}$$



→ 一つの文書が複数のトピックを持つと仮定する確率的潜在意味解析(PLSA)をベイズモデルに拡張したもの

混合メンバーシップ確率的ブロックモデル (MMSB) Airoldi et al. 2008

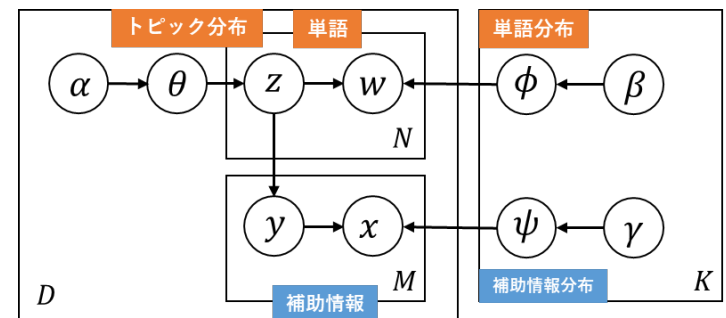
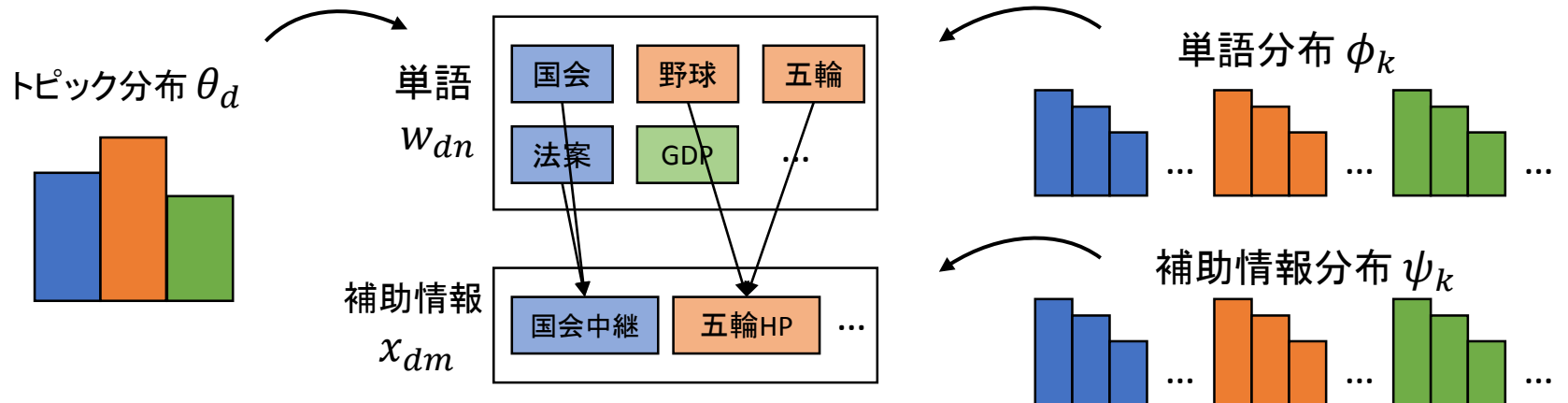
→ 一人のユーザーがリンクごとに複数のトピックを持つと仮定したもの
文書 \leftrightarrow ユーザー 単語 \leftrightarrow リンク



対応トピックモデル (CTM)

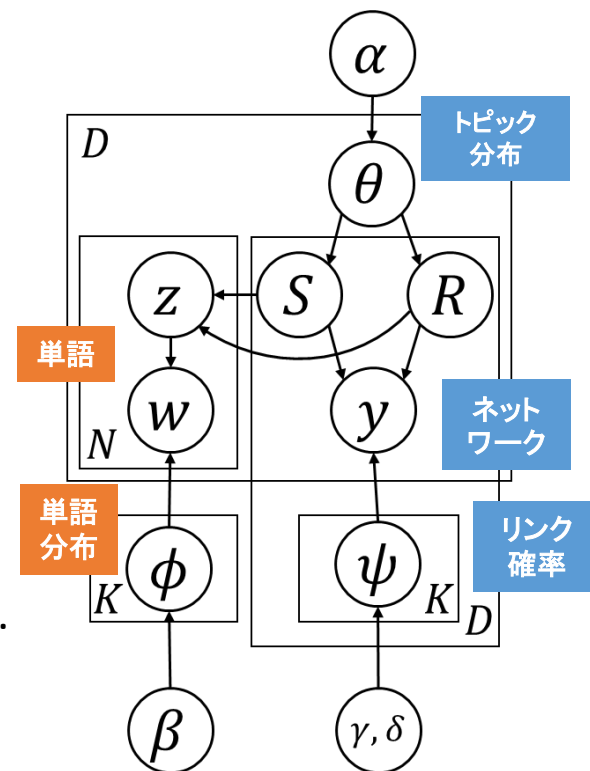
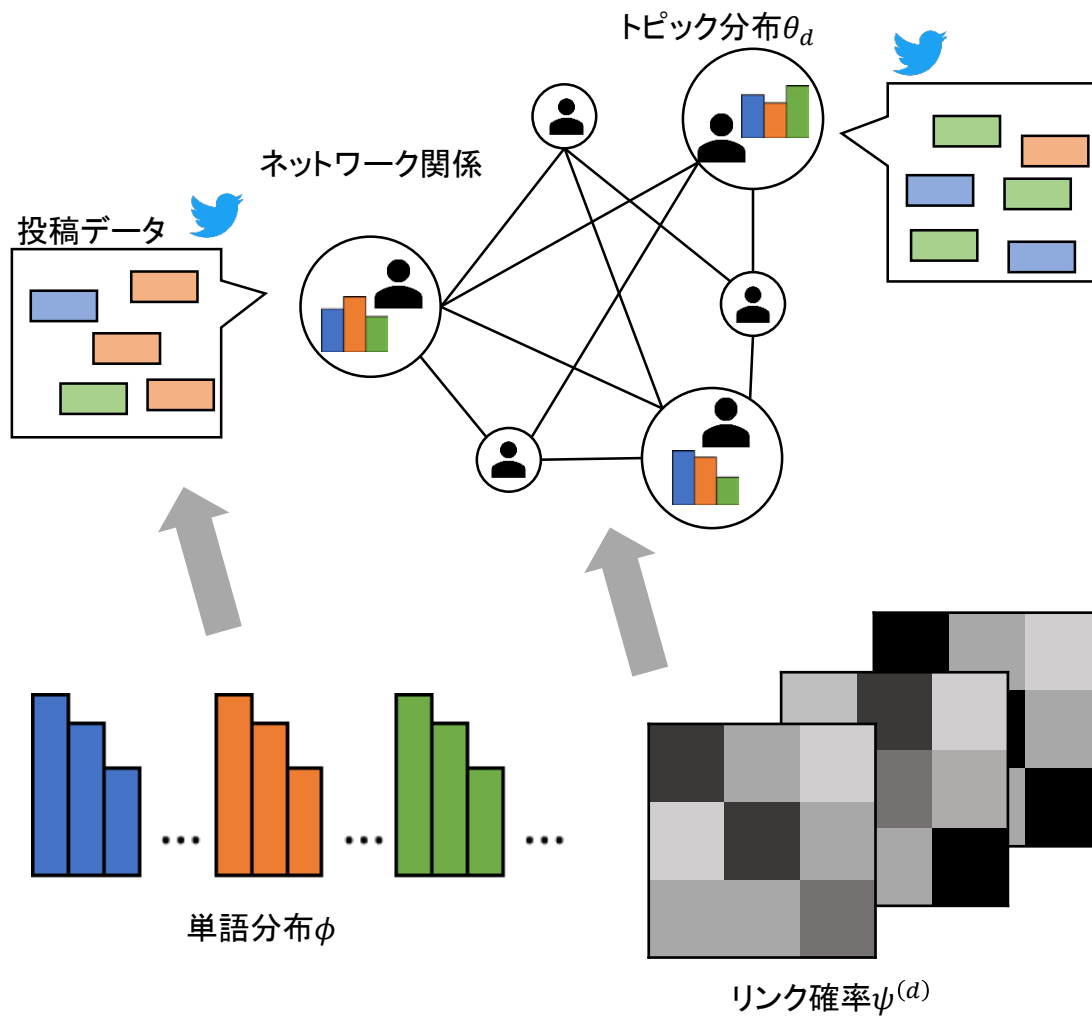
Blei et al. 2003b

→ 主情報と補助情報のトピックを対応させるモデル



3.分析モデル

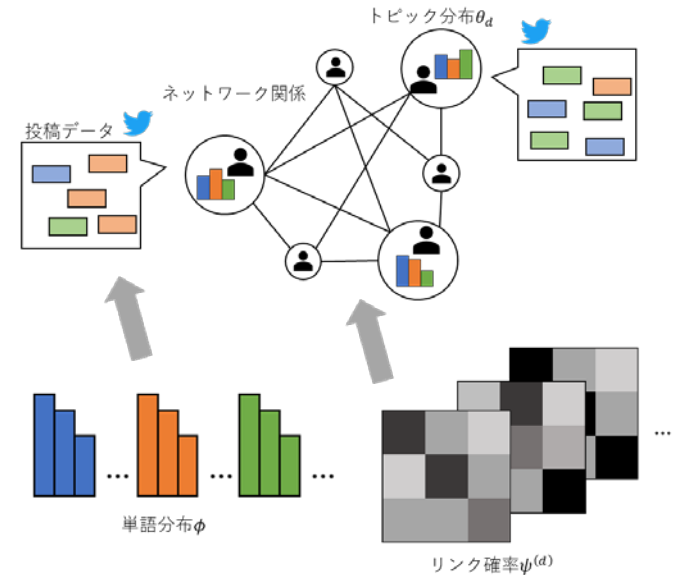
分析モデル



分析モデルのグラフィカル表現

分析モデル

本研究のモデルもCTMの一種
主情報: ネットワーク関係
補助情報: SNSの投稿データ
→ ネットワークデータからトピック
(=コミュニティ)を抽出しながら
テキスト情報によりコミュニティの
特色を表すことができる



モデルの特色

ネットワークを主情報としている

→ 機械学習の分野では文書のネットワーク構造に着目

(Chang & Blei 2010, Chen et al. 2015)

リンク確率に異質性を導入

→ ネットワーク上の特異な位置にいる個人を特定することが可能
(インフルエンサーの検出)

分析モデル

- モデルの推定には周辺化ギブスサンプリングを用いる
パラメータ(θ, ϕ, ψ)を積分消去してトピック(z, S, R)
のサンプリングのみを行う → サンプリング式などはAppendix参照
- トピック数は5, 10, 15, 20, 25の中からPerplexityが最小となるものを選択
Perplexityは負の対数尤度をテストデータの数で割って指数を取ったものであり、次の式で表される

$$Perplexity = \exp \left\{ - \frac{\log p(y^{test}, w^{test} | model)}{\sum_{d=1}^D N_d^{test}} \right\}$$

4.分析結果

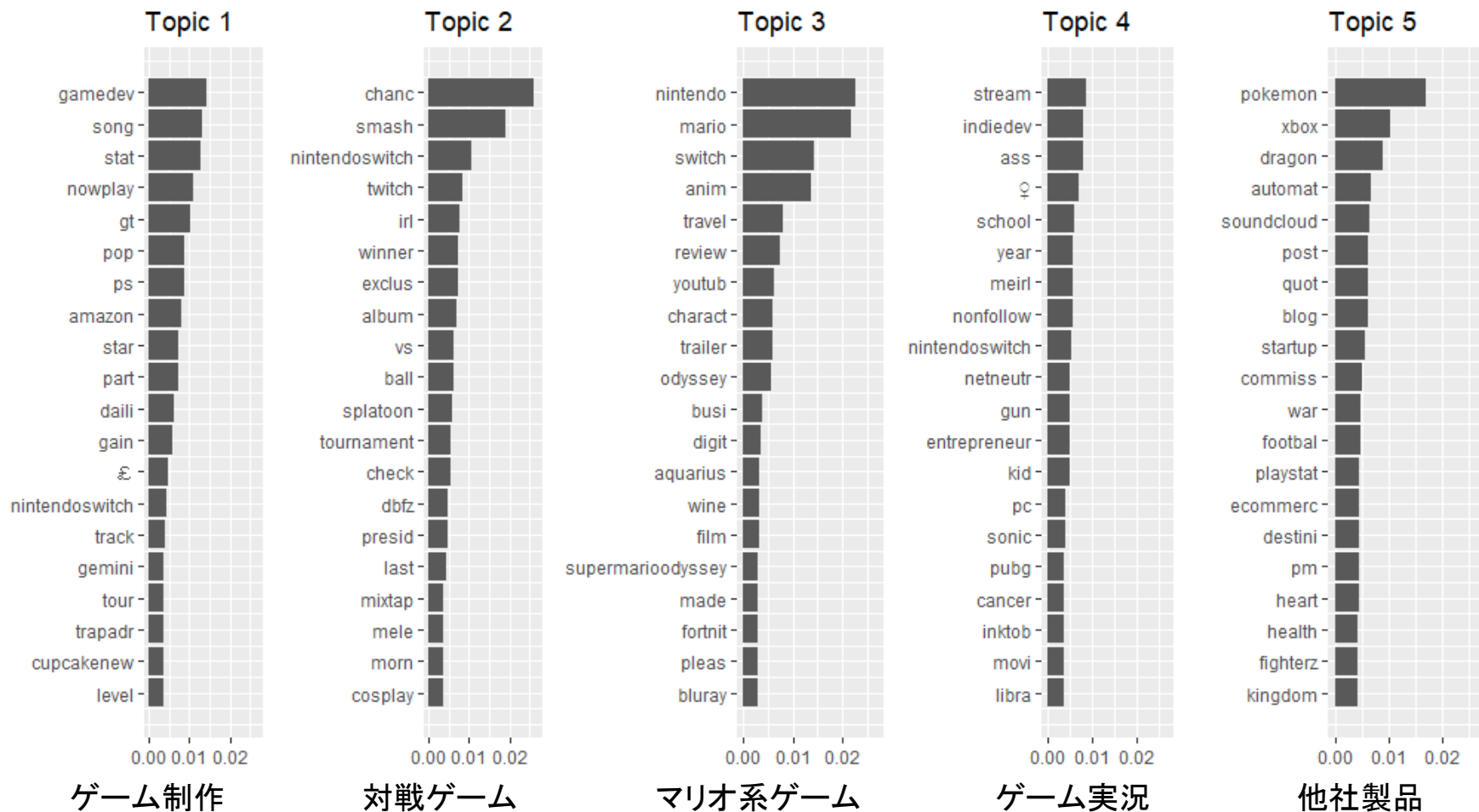
データ

- ソーシャルメディア : Twitter
- 対象 : Nintendo of Americaの公式アカウント (@NintendoAmerica) をフォローしているユーザー、及びそれらのユーザーをフォローしているユーザーからランダムにサンプリング
- 期間 : リンク関係は2018年5月1日時点
投稿データは2017年9月1日から2018年2月28日の6か月間
- 前処理 : 投稿データにはストップワードの除去やステミングなど前処理を行った(詳細はAppendix)
- 要約統計量 :

ユーザー数	語彙数	リンク数	平均単語数	平均リンク数	テスト比率
5,028	19,805	75,816	32.82	15.08	9:1

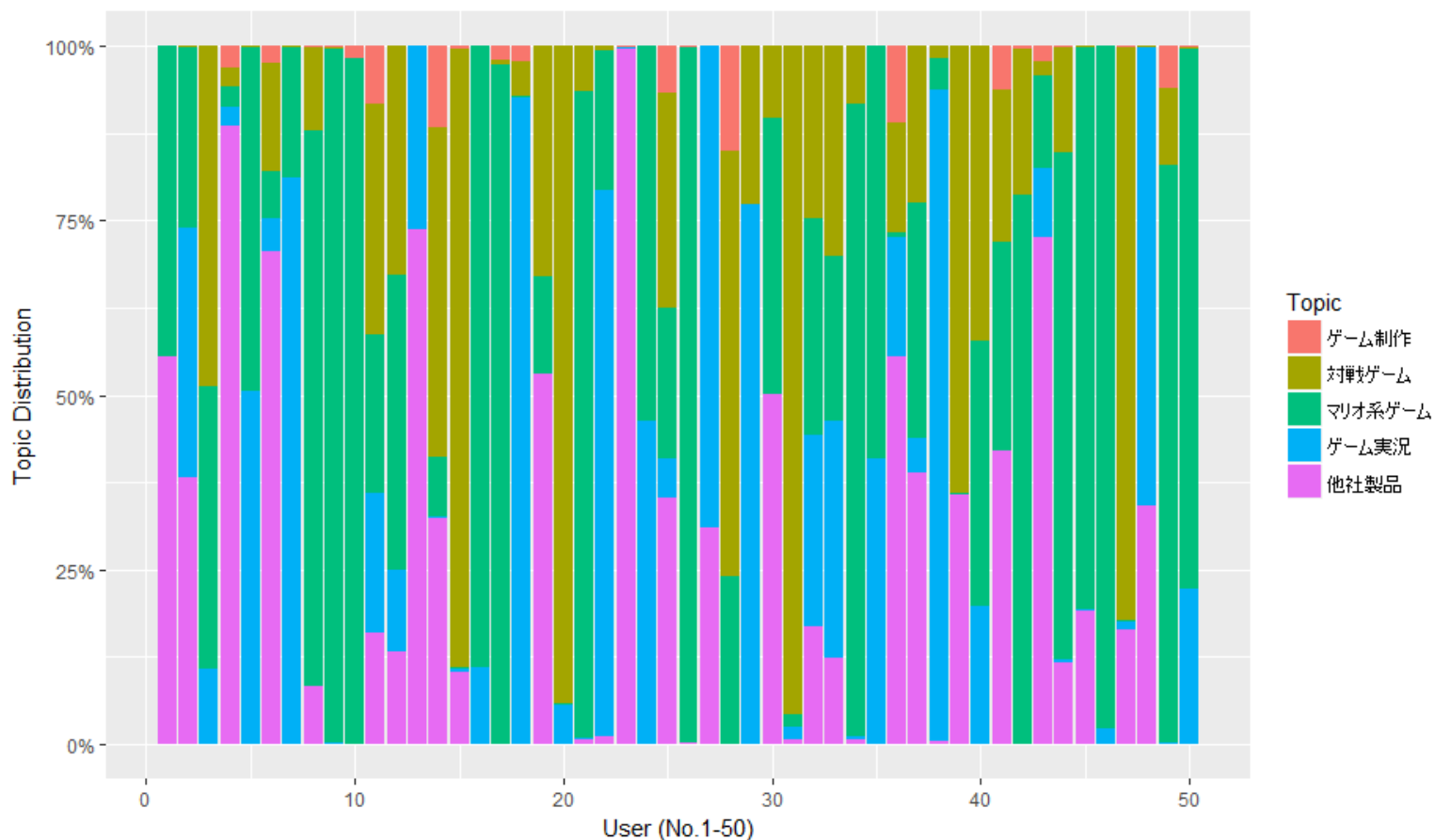
分析結果－パラメータの解釈－

- 単語分布 ϕ の解釈



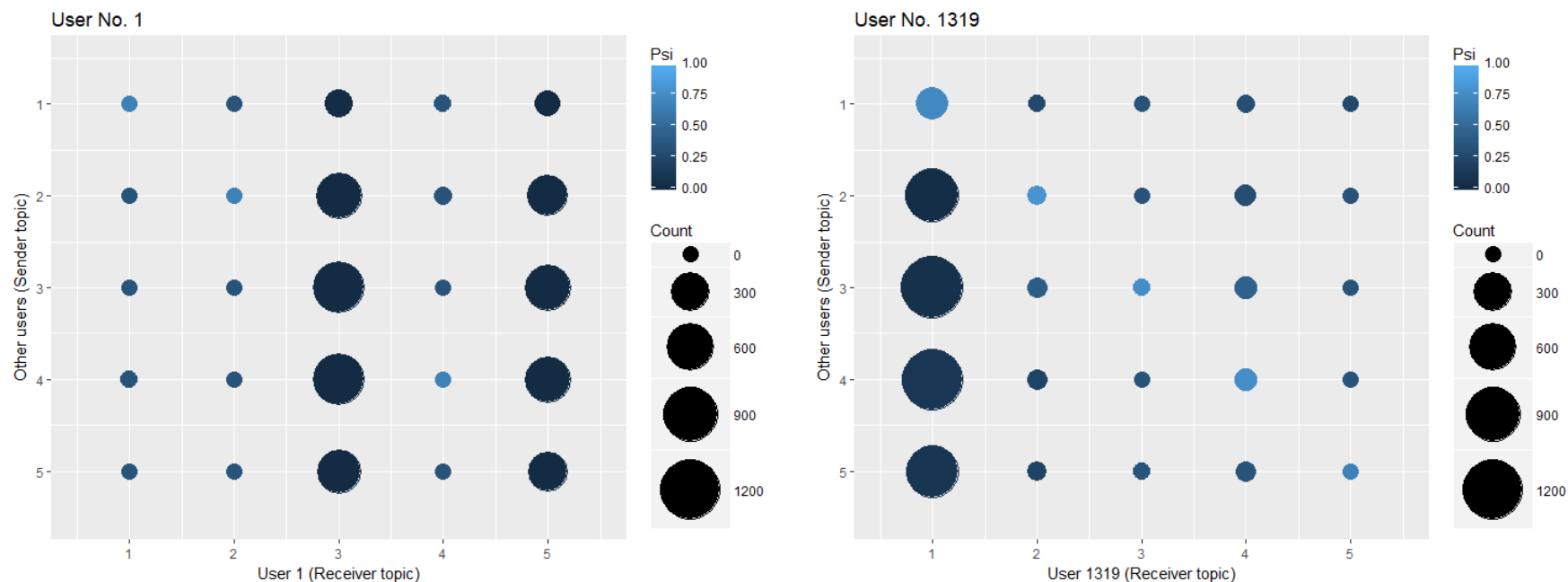
分析結果－パラメータの解釈－

- トピック分布 θ_d の解釈
→ ユーザーごとの潜在的なコミュニティ所属割合を把握できる



分析結果－パラメータの解釈－

- リンク確率 $\psi^{(d)}$ の解釈

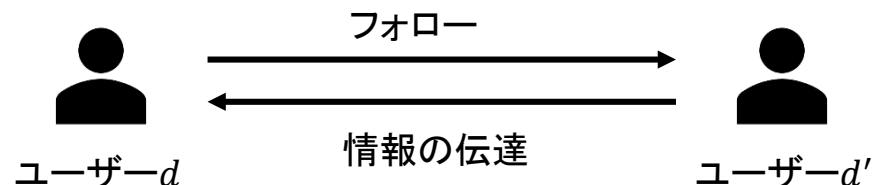


Psi(色の濃さ): 推定されたリンク確率 $\psi^{(d)}$

Count(円の大きさ): トピック $k \rightarrow k'$ へのリンク関係の数

5. シミュレーション

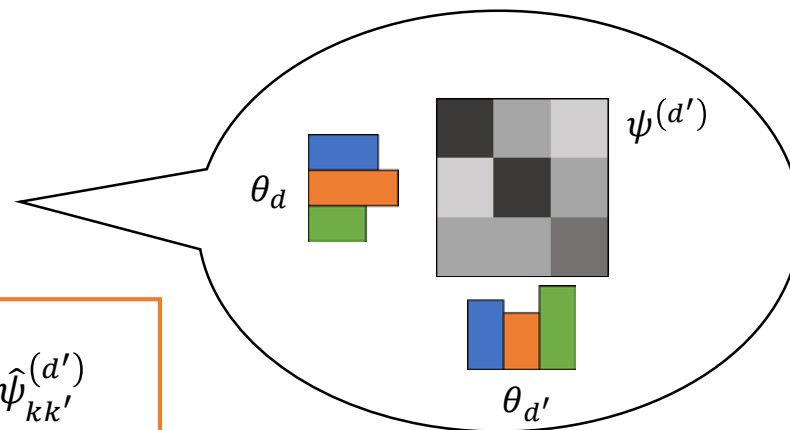
情報伝達の影響力



リンク確率 = d' が発信する情報を受け取りたい確率

$$\begin{aligned} &= \sum_{k=1}^K d' \text{が発信するトピック}k \text{に関する情報を受け取りたい確率} \\ &= p(y_{dd'} = 1 | \cdot) \\ &= \sum_{k=1}^K \sum_{k'=1}^K \hat{\theta}_{dk} \hat{\theta}_{d'k} \hat{\psi}_{kk'}^{(d')} \end{aligned}$$

$$d' \rightarrow d \text{の情報}k \text{に関する影響力} = \hat{\theta}_{dk} \sum_{k'=1}^K \hat{\theta}_{d'k} \hat{\psi}_{kk'}^{(d')}$$



リンク有でも影響力が小さければそのリンクによる情報伝達の効果は薄い

リンク無でも影響力が大きければ潜在的に高い情報伝達となる可能性がある

バイラル・マーケティングのシミュレーション

リンクごとの影響力を考慮することでユーザーの関心を加味したバイラル・マーケティングが可能となる

対象ツイート:

2018年3月に行われた新作ゲームの発表
→ トピック2のコミュニティが興味を持つと考えられる

実データ:

1061人(5028人中)がRT・コメント
→ 企業が介入し、推定された影響力の高いユーザー順にRT・コメントさせる場合と比較

仮定:

RTやゲームに関するコメントを見た際に、
影響力が閾値を超えていた場合、そのユーザーは関心を持った状態になると仮定する
→ RTやコメントを通じてユーザーの関心度がどのように高まっていくかをシミュレーション



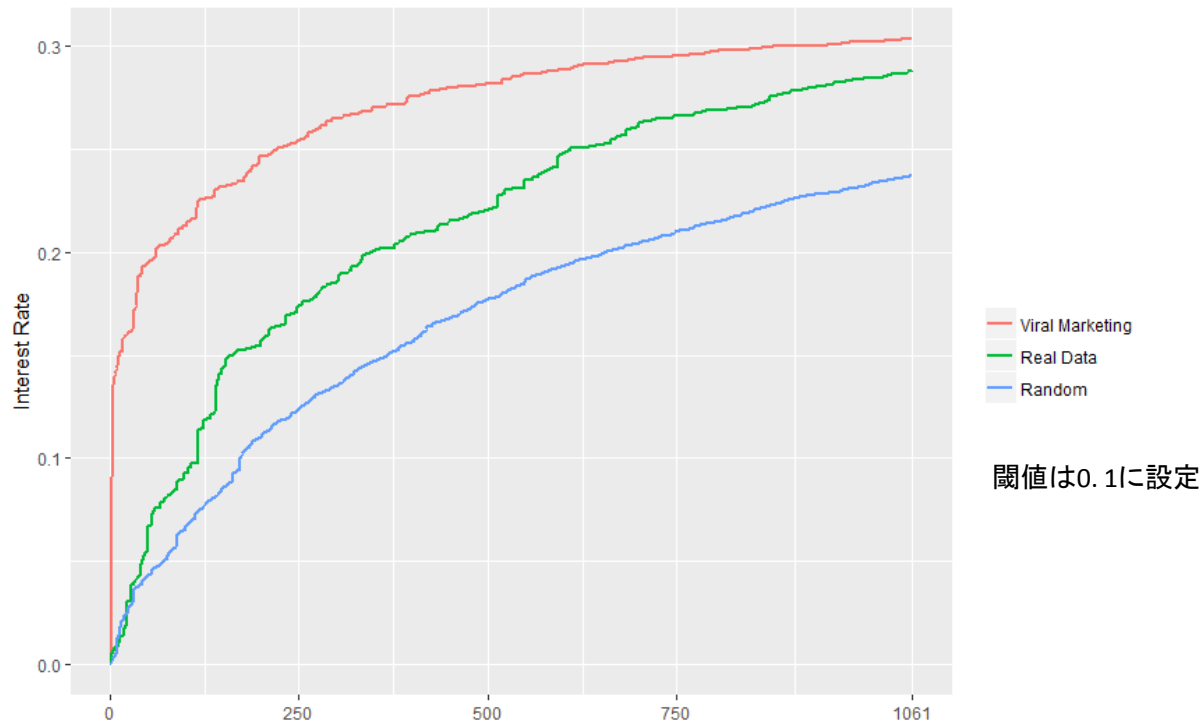
対象ツイート

バイラル・マーケティングのシミュレーション

Viral Marketing: 影響カスコアが高い順

Real Data: 実際に対象TweetをRTした順

Random: ランダムな試行を10回繰り返したときの平均



まとめ

研究目的

ユーザーがどのような関心でつながっているかを把握する

分析モデル

リンクトピックと単語トピックを対応させたトピックモデル
リンク確率に異質性を導入

分析結果

単語分布: 各トピックの内容
トピック分布: 各ユーザーのコミュニティ所属割合
リンク確率: 各ユーザーのリンク関係

シミュレーション

ユーザーごと・トピックごとの影響力を考慮した
バイラル・マーケティングが有効であるという示唆を得た

今後の課題

- モデルの有効性の検証

分析モデルを用いたバイラル・マーケティングの有効性はあくまでシミュレーションによって示唆されたただけである。

閾値の設定やRTによって実際に関心を持ったのかなど曖昧な点も多い。

実際に行われたバイラル・マーケティングのデータを用いるなどしてモデルの有効性を検証する必要がある。

- モデルの改良

トピック数の推定(階層ディリクレ過程)やハイパーパラメータの推定(不動点反復法)などモデルに改良の余地がある。また、ソーシャルネットワークのリンクデータは基本的に疎なネットワークであり、Airoldi et al. (2008)では Sparsity parameterを導入することを提唱している。これもまた改良の余地がある部分である。

参考文献

- Airoldi, E. M., Blei, D. M., Fienberg, S. E., & Xing, E. P. (2008). Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9(Sep), 1981-2014.
- Bampo, M., Ewing, M. T., Mather, D. R., Stewart, D., & Wallace, M. (2008). The effects of the social structure of digital networks on viral marketing performance. *Information systems research*, 19(3), 273-290.
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003a). Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan), 993-1022.
- Blei, D. M., & Jordan, M. I. (2003b, July). Modeling annotated data. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*(pp. 127-134). ACM.
- Chae, I., Stephen, A. T., Bart, Y., & Yao, D. (2016). Spillover effects in seeded word-of-mouth marketing campaigns. *Marketing Science*, 36(1), 89-104.
- Chen, N., Zhu, J., Xia, F., & Zhang, B. (2015). Discriminative relational topic models. *IEEE transactions on pattern analysis and machine intelligence*, 37(5), 973-986.
- Chang, J., & Blei, D. M. (2010). Hierarchical relational models for document networks. *The Annals of Applied Statistics*, 124-150.
- Chen, X., van der Lans, R., & Phan, T. Q. (2017). Uncovering the importance of relationship characteristics in social networks: Implications for seeding strategies. *Journal of Marketing Research*, 54(2), 187-201.
- Gong, S., Zhang, J., Zhao, P., & Jiang, X. (2017). Tweeting as a Marketing Tool: A Field Experiment in the TV Industry. *Journal of Marketing Research*, 54(6), 833-850.
- 岩田具治. (2015). トピックモデル (機械学習プロフェッショナルシリーズ).
- 佐藤一誠, 情報科学, & 奥村学. (2015). トピックモデルによる統計的潜在意味解析. コロナ社.

Appendix 分析モデルの生成過程

1. For トピック $k = 1, \dots, K$
 - a. 単語分布を生成 $\phi_k \sim \text{Dirichlet}(\beta)$
 - b. For トピック $k' = 1, \dots, K$
 - i. For ユーザー $d = 1, \dots, D$
 - リンク確率を生成 $\psi_{kk'}^{(d)} \sim \text{Beta}(\gamma, \delta)$
2. For ユーザー $d = 1, \dots, D$
 - a. トピック分布を生成 $\theta_d \sim \text{Dirichlet}(\alpha)$
 - b. For ユーザー $d' = 1, \dots, D$
 - i. リンクトピック(Sender)を生成 $S_{dd'} \sim \text{Categorical}(\theta_d)$
 - ii. リンクトピック(Receiver)を生成 $R_{dd'} \sim \text{Categorical}(\theta_{d'})$
 - iii. リンクを生成 $y_{dd'} \sim \text{Bernoulli}(\psi_{S_{dd'}, R_{dd'}}^{(d')})$
 - c. For 単語 $n = 1, \dots, N_d$
 - i. 単語トピックを生成 $z_{dn} \sim \text{Categorical}(\frac{N_{d\cdot}}{2(D-1)})$
 - ii. 単語を生成 $w_{dn} \sim \text{Categorical}(\phi_{z_{dn}})$

$N_{d\cdot} = (N_{d1}, \dots, N_{dK})$
 N_{dk} はユーザー d のリンクのうちトピック k が割り当てられたSender・Receiverの和

Appendix 周辺化ギブスサンプリングのサンプリング式

- $$p(S_{dd'} = k, R_{dd'} = k' | S_{\setminus dd'}, R_{\setminus dd'}, Y, Z, \alpha, \gamma, \delta)$$

$$\propto (N_{dk \setminus dd'} + \alpha_k) \times (N_{d'k' \setminus dd'} + \alpha_l)$$

$$\times \frac{\left(n_{kk' \setminus dd'}^{(+, d')} + \gamma_{kk'}\right)^{\mathbb{I}(y_{dd'}=1)} \left(n_{kk' \setminus dd'}^{(-, d')} + \delta_{kk'}\right)^{\mathbb{I}(y_{dd'}=0)}}{n_{kk' \setminus dd'}^{(+, d')} + n_{kk' \setminus dd'}^{(-, d')} + \gamma_{kk'} + \delta_{kk'}}$$

$$\times \left(\frac{N_{dk \setminus dd'} + 1}{N_{dk \setminus dd'}}\right)^{M_{dk}} \times \left(\frac{N_{d'k' \setminus dd'} + 1}{N_{d'k' \setminus dd'}}\right)^{M_{d'k'}}$$
- $$p(z_{dn} = k | Z_{\setminus dn}, W, S, R, \beta) \propto N_{dk} \times \left(\frac{M_{kw \setminus dn} + \beta_{w \setminus dn}}{\sum_v M_{kv \setminus dn} + \beta_v}\right)$$

記号

N_{dk}	ユーザー d のリンク関係のうちトピック k が割り当てられた数
$n_{kk'}^{(+, d)}$	ユーザー d のリンク関係のうちトピック k から k' に向かうリンク有の数
$n_{kk'}^{(-, d)}$	ユーザー d のリンク関係のうちトピック k から k' に向かうリンク無の数
M_{dk}	ユーザー d の文書のうちトピック k が割り当てられた単語の数
M_{kv}	語彙 v にトピック k が割り当てられた数
$\cdot \setminus dd' \cdot \cdot \setminus dn$	$y_{dd'}, w_{dn}$ に関するトピックの割り当て分をカウントから除く操作

ハイパーパラメータの設定

$\alpha = 1.0/K$ K はトピック数
$\beta = 1.0/V$ V は語彙数
$\gamma_{kk} = 2.0$ $\gamma_{kk'} = 1.0 (k \neq k')$
$\delta_{kk} = 1.0$ $\delta_{kk'} = 2.0 (k \neq k')$

Appendix Perplexityの導出

Perplexityを計算するためにはテストデータに対する予測確率 $p(y^*, w^* | model)$, $y^* \in y^{test}$, $w^* \in w^{test}$ を計算する必要がある。

しかし、これは解析的に求めることができないため、ギブスサンプリングによる各サンプル s によって以下のように近似する。

$$\begin{aligned} p(y^*, w^* | y^{train}, w^{train}, \alpha, \beta, \gamma, \delta) &= p(y^* | y^{train}, \alpha, \gamma, \delta) p(w^* | w^{train}, \beta) \\ p(y_d^* | y^{train}, \alpha, \gamma, \delta) &\sim \frac{1}{S} \sum_{s=1}^S \sum_{k=1}^K \sum_{k'=1}^K \frac{N_{dk}^{(s)} + \alpha_k}{\sum_l (N_{dl}^{(s)} + \alpha_l)} \frac{N_{d'k'}^{(s)} + \alpha_{k'}}{\sum_l (N_{d'l}^{(s)} + \alpha_l)} \frac{(n_{kk'}^{(+,d',s)} + \gamma_{kk'})^{\mathbb{I}(y_{dd'}=1)} (n_{kk'}^{(-,d',s)} + \delta_{kk'})^{\mathbb{I}(y_{dd'}=0)}}{n_{kk'}^{(+,d',s)} + n_{kk'}^{(-,d',s)} + \gamma_{kk'} + \delta_{kk'}} \\ p(w_d^* = v | w^{train}, \beta) &\sim \frac{1}{S} \sum_{s=1}^S \sum_{k=1}^K \frac{N_{dk}^{(s)}}{\sum_l N_{dl}^{(s)}} \frac{M_{kv}^{(s)} + \beta_v}{\sum_u (M_{ku}^{(s)} + \beta_u)} \end{aligned}$$

近似予測確率を用いてPerplexityは次のように計算される

$$Perplexity = \exp \left(\frac{\sum_{d=1}^D \log p(y_d^* | y^{train}, \alpha, \gamma, \delta)}{\sum_{d=1}^D N_d^{test}} + \frac{\sum_{d=1}^D \log p(w_d^* | w^{train}, \beta)}{\sum_{d=1}^D M_d^{test}} \right)$$

Appendix テキストデータの前処理

取得したテキストデータに対して以下の順に前処理を行った

1. ストップワードの除去

統計ソフトR内パッケージ“tm”に含まれている英語のストップワードに加えて、記号・URL文字列・screen name (@+英数字の文字列)などをストップワードとしてテキストデータから取り除いた。

2. ステミング(テキストの正規化)

全ての単語について小文字に統一し、活用形を直すステミングの処理を行った(COME, came, comingなどは全てcomeに統一される)。これには、“tm”の関数であるstemDocumentを使用した。

3. 形態素解析による名詞の抽出

単語の品詞を検出する形態素解析を行い、名詞以外の単語を取り除いた。これには、“TreeTagger”というソフトウェアを使用した。

4. TF-IDFによる単語の選定

トピックの解釈性向上と計算コストの削減を目的に、TF-IDFが上位1%の単語を選定した。なお、TF-IDFは以下の式で計算される単語の重要度である。

$$TF_{dv} = \frac{N_{dv}}{\sum_{v=1}^V N_{dv}}, \quad IDF_v = \log \left(\frac{D}{D_v} \right) + 1$$
$$TFIDF_{dv} = TF_{dv} \times IDF_v$$