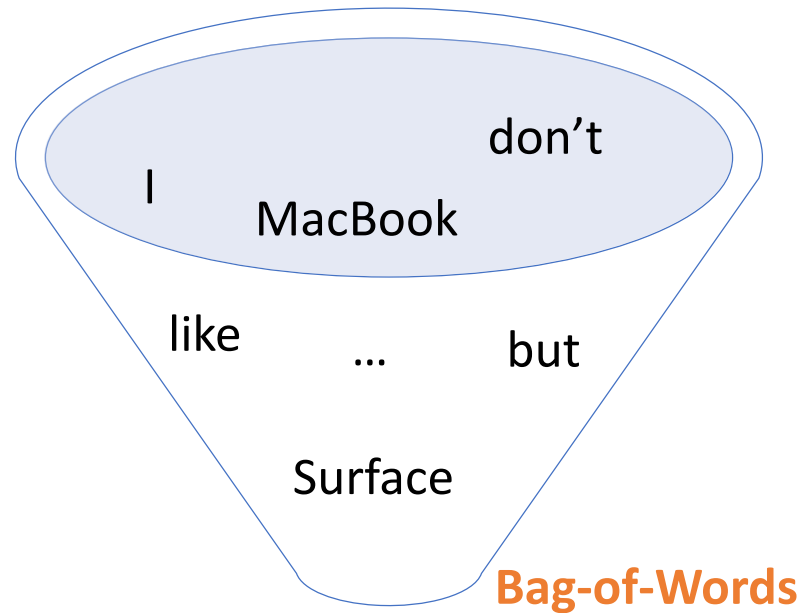# Beyond Bag-of-Words

Customer Review Analysis Using Word Embedding
Model Considering Text Topics and Sentiments

(work with N. Terui and P.K. Kannan)

# Introduction

**I like MacBook …
but don't like Surface …**

**I like Surface …
but don't like MacBook…**

don't

I

MacBook

like

…

but

Surface

**Bag-of-Words**

# Word2vec

→単語のベクトル表現を学習するニューラルネットワーク

$$\overrightarrow{w_i} = \{w_{i1}, \ldots, w_{iN}\}$$

**Skipgram**モデル

$$p(w_O|w_I) = \frac{\exp(\overrightarrow{w'_O} \cdot \overrightarrow{w_I})}{\sum_v \exp(\overrightarrow{w'_v} \cdot \overrightarrow{w_I})}$$

目的関数

$$\arg\max_{w,w'} p(w_{O1}, \ldots, w_{OC}|w_I) = \prod_{c=1}^{C} \frac{\exp(\overrightarrow{w'_{Oc}} \cdot \overrightarrow{w_I})}{\sum_v \exp(\overrightarrow{w'_v} \cdot \overrightarrow{w_I})}$$

# Word2vec

→ローカルな共起関係を考慮して単語をベクトル空間に射影する



Country and Capital Vectors Projected by PCA

Mikolov et al. (2013)

# LDA2vec



Moody (2016)

# LDA2vec

Context vector

$$\overrightarrow{c_{ih}} = \overrightarrow{w_i} + \overrightarrow{d_h}$$

Document vector

$$\overrightarrow{d_h} = \theta_{h1} \cdot \overrightarrow{t_1} + \cdots + \theta_{hK} \cdot \overrightarrow{t_K}$$

Loss functions

$$L = L^w + L^d$$

$$L^w = \sum_{i,j,h} \left\{ \log \sigma(\overrightarrow{c_{ih}} \cdot \overrightarrow{w_j}) + \sum_n \log \sigma(-\overrightarrow{c_{ih}} \cdot \overrightarrow{w_n}) \right\}$$

$$L^d = \sum_{h=1}^{H} \left\{ \lambda(\alpha - 1) \sum_{k=1}^{K} \log \theta_{hk} \right\}$$

# LDA2vec with rating regression

## Model

$$\overrightarrow{c_{ih}} = \overrightarrow{w_i} + \overrightarrow{d_h}$$

$$\overrightarrow{d_h} = \theta_{h1} \cdot \overrightarrow{t_1} + \cdots + \theta_{hK} \cdot \overrightarrow{t_K}$$

$$y_h = l \quad \text{if } \tau_{l-1} \leq \hat{y}_h < \tau_l, \quad \tau_0 = -\infty, \tau_L = +\infty$$

$$\hat{y}_h = \theta_h^T \beta + \text{other variables}$$

## Loss functions

$$L = L^w + L^d + L^y$$

$$L^w = \sum_{i,j,h} \left\{ \log \sigma(\overrightarrow{c_{ih}} \cdot \overrightarrow{w_j}) + \sum_n \log \sigma(-\overrightarrow{c_{ih}} \cdot \overrightarrow{w_n}) \right\}$$

$$L^d = \sum_{h=1}^{H} \left\{ \lambda(\alpha - 1) \sum_{k=1}^{K} \log \theta_{hk} \right\}$$

$$L^y = \sum_{h=1}^{H} \left\{ \sum_{l=1}^{y_h - 1} \psi(-\gamma_{hl}) + \sum_{l=y_h}^{L} \psi(\gamma_{hl}) \right\}, \qquad \psi(\cdot) = \{\text{hinge, logistic, expoential, etc}\}$$

# LDA2vec with rating regression

## Simulation experiments

| Topic 1 Business | Topic 2 Entertainment | Topic 3 Sports | Topic 4 Technology | Topic 5 Politics |
|---|---|---|---|---|
| bank | film | match | mobil | parti |
| growth | award | champion | technolog | labour |
| profit | actor | cup | user | elect |
| oil | album | coach | comput | blair |
| yuko | chart | rugbi | phone | howard |
| sharehold | nomin | chelsea | softwar | mp |
| airlin | song | ireland | onlin | lord |
| stock | oscar | victori | digit | brown |
| deficit | rock | injuri | blog | lib |

# LDA2vec with rating regression

## Simulation experiments

| thresholds ($\tau$) | 1-2 | 2-3 | 3-4 | 4-5 |
|---|---|---|---|---|
| true | -1.0 | 0.0 | 0.8 | 1.3 |
| estimates | -1.19 | -0.47 | 0.21 | 1.17 |

| coefficients ($\beta$) | business | entertainment | sport | technology | politics |
|---|---|---|---|---|---|
| true | 2.0 | 1.0 | -1.0 | -2.0 | 1.5 |
| estimates | 1.40 | 0.20 | -1.12 | -1.75 | 0.99 |

| | | True | | | | |
|---|---|---|---|---|---|---|
| | | 5 | 4 | 3 | 2 | 1 |
| | 5 | 51 | 0 | 0 | 0 | 0 |
| | 4 | 9 | 60 | 0 | 0 | 0 |
| Prediction | 3 | 0 | 0 | 43 | 0 | 0 |
| | 2 | 0 | 0 | 17 | 60 | 6 |
| | 1 | 0 | 0 | 0 | 0 | 51 |

# Sentiment LDA2vec

LDA2vec
Embedding space

Topic 2
(laptop)　MacBook

Topic 1
(laptop)

Surface

Topic vector

極性を考慮
(positive / negative)

# Sentiment LDA2vec

Embedding dimensionをpositiveとnegativeに分割

$$\overrightarrow{w_i} = \{0.23, \dots, 0.45, 0.78, \dots, 0.11\}$$

<span style="color:orange">**Positive**</span>　　<span style="color:blue">**Negative**</span>

半教師有学習（極性辞書を利用, Lin et al. 2016）

<span style="color:orange">**Positive単語**</span>　　$\{w_{i,1}, \dots, w_{i,N/2}, 0, \dots, 0\}$

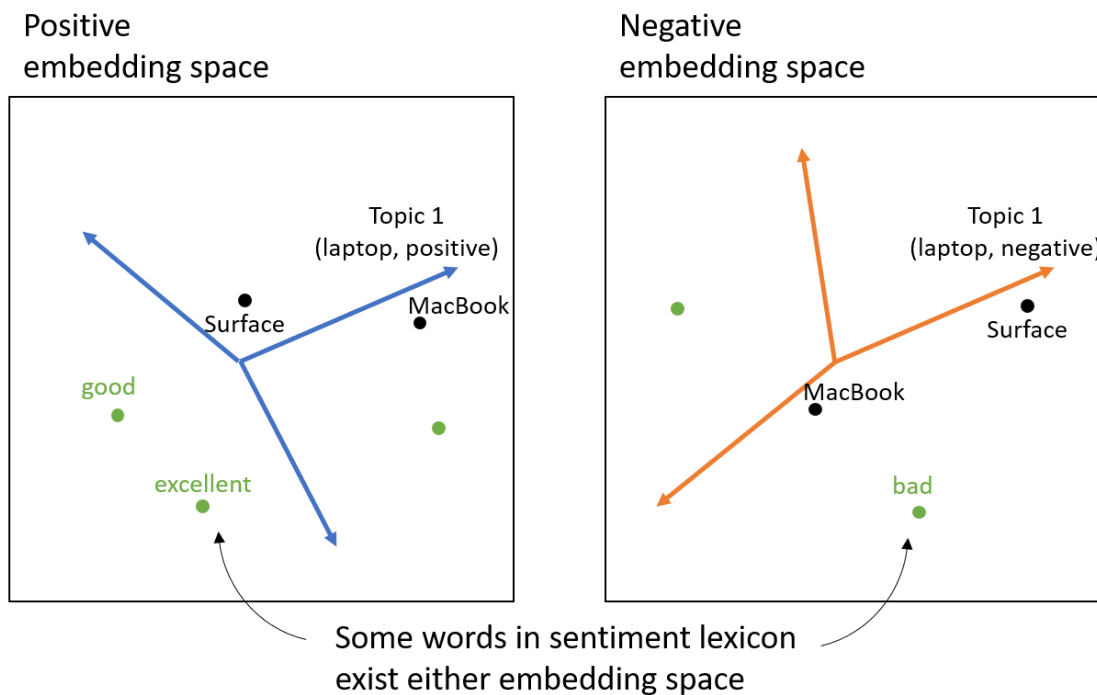<span style="color:blue">**Negative単語**</span>　　$\{0, \dots, 0, w_{i,N/2+1}, \dots, w_{i,N}\}$

<span style="color:green">**それ以外の単語**</span>　　$\{w_{i,1}, \dots, w_{i,N/2}, w_{i,N/2+1}, \dots, w_{i,N}\}$

# Sentiment LDA2vec

## Model

$$\overrightarrow{c_{ih}} = \overrightarrow{w_i} + f_h \cdot \overrightarrow{d_{h,pos}} + (1 - f_h) \cdot \overrightarrow{d_{h,neg}}$$

$$\overrightarrow{d_{hm}} = \theta_{hm1} \cdot \overrightarrow{t_{m1}} + \cdots + \theta_{hmK} \cdot \overrightarrow{t_{mK}}, \qquad m \in \{positive, negative\}$$

Positive
embedding space

Topic 1
(laptop, positive)

MacBook

Surface

good

excellent

Negative
embedding space

Topic 1
(laptop, negative)

Surface

MacBook

bad

Some words in sentiment lexicon
exist either embedding space

# Sentiment LDA2vec

## Rating regression

$$y_h = l, \qquad \text{if } \tau_{l-1} \leq \hat{y}_h < \tau_l, \qquad l = 1, \ldots, L$$

$$\hat{y}_h = f_h \cdot \theta_{h,pos}^T \beta_{pos} + (1 - f_h) \cdot \theta_{h,neg}^T \beta_{neg} + \text{other variables}$$

## Future work

- Simulation experiment

- Bayesian estimation