



Introduction

言語モデル

言語モデル

word2vec の概要

word2vec の概要

Theorem

Proof of Theorem

Algorithms

Conclusions

# word2vec Parameter Learning Explained

Shigeaki ONODA

Graduate School of Information System Engineering Dept  
t855005@st.pu-toyama.ac.jp

Tues., 12 January, 2016, Osaka Univ.



# Introduction

## Introduction

言語モデル

言語モデル

word2vec の概要

word2vec の概要

Theorem

Proof of Theorem

Algorithms

Conclusions

## Motivation

離散オブジェクトを数値として扱うことで和・差等の演算が可能となる。

自然言語をテンソルとして扱うためにはモデルとして以下が考案されている離散オブジェクト：  
単語，概念のような物理的に計測できる量を伴わず，通常記号を用いて離散的に表現するもの

▶ n-gram 言語モデル

## Agenda

- 言語モデルとは
- word2vec の概要
- 結論



# 言語モデルと n-gram

Introduction

言語モデル

言語モデル

word2vec の概要

word2vec の概要

Theorem

Proof of Theorem

Algorithms

Conclusions

## 言語モデル

人間が扱う自然言語で書かれた文や文書が生成される確率をモデル化したもの。  
例えば「おはようございます。 朝会を始めます。」には高い確率を与え、「おはようございます。懲戒を始めます。」といった出現確率の低い文には低い確率を与える

## n-gram モデル

隣り合って出現した  $n$  単語のことを  $n$ -gram と呼ぶ。  
nurture passes nature =>

- 1-gram: nurture, passes, nature
- 2-gram: nurture-passes, passes-nature
- 3-gram: nurture-passes-nature

このように 2 グラム以上だと前後の文脈を考慮したモデルとなる。

## Bag of Words

1-gram の場合の文書ベクトル表現のこと単純に文書中のその単語の頻度を返す。  
"nurture or nature? nurture passes nature."

だとその BoW は

$$x^d = (\text{nature}, \text{nurture}, \text{or}, \text{passes}) = (2, 2, 1, 1)$$

この手法でも言語処理ではある程度有効であるため、これを用いる研究も多い



# one hot ベクトル

Introduction

言語モデル

言語モデル

word2vec の概要

word2vec の概要

Theorem

Proof of Theorem

Algorithms

Conclusions

## one-hot ベクトル

要素のうち一つだけ 1 のベクトル

$$x = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} 0.01 & 0.02 & 0.03 & 0.04 \\ 0.06 & 0.01 & 0.08 & 0.10 \\ 0.02 & 0.12 & 0.14 & 0.15 \\ 0.01 & 0.02 & 0.03 & 0.04 \\ 0.21 & 0.22 & 0.23 & 0.24 \end{pmatrix} * \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0.02 \\ 0.01 \\ 0.12 \\ 0.02 \\ 0.22 \end{pmatrix} \quad (1)$$

## 言語処理における one hot

上の式の左行列を埋め込み行列とよび右辺のベクトルを埋め込みベクトルとよぶによって埋め込みベクトルが  $i$  番目の単語の情報を格納している

## 記号からテンソル, テンソルから記号への変換

単語数と同じ次元ベクトル  $\mathbf{o}$  を考えるそのベクトル中で最も大きな値になった要素番号 (インデックス) を返す近似式は以下となる

$$\text{softmax}_a(\mathbf{o}) = \frac{1}{\exp(a\mathbf{o}) \cdot \mathbf{1}} \exp(a\mathbf{o}) \quad (2)$$



# ニューラル言語モデル

Introduction

言語モデル

言語モデル

word2vec の概要

word2vec の概要

Theorem

Proof of Theorem

Algorithms

Conclusions

## ニューラル言語モデル

機械学習の権威ヒントンらが考案したニューラルネットを用いた言語モデルの総称 **n-gram** モデルと比べて言語モデルとしての性能が優れているただし計算量は **n-gram** に比べて膨大になる  
=>GPGPU の発展によりこの問題が解消されつつあるので近年フォーカスされている

## ニューラル言語モデルの定式化

単語の位置  $t$  より前に出現した  $t - a$  単語を  $\mathbf{Y}_{[a,t-1]} = (\mathbf{y}_a, \mathbf{y}_{a+1}, \dots, \mathbf{y}_{t-1})$  と表しこれが文脈に相当する。

このとき文  $\mathbf{Y}$  の生成確率  $P(\mathbf{Y})$  を各単語が生成される条件付き確率の積でモデル化する。ただし  $\mathbf{Y}$  は文中単語の **one hot** ベクトルが格納された集合である

$$P_{model}(\mathbf{Y}) = P(\mathbf{y}_0) \prod_{t=1}^{T+1} P_{model}(\mathbf{y}_t | \mathbf{Y}_{[a,t-1]}) \quad (3)$$

## ニューラル言語モデルの応用

これを応用（ただ学習した埋め込み行列を得る）することで分散表現を獲得できる。

分散表現とは分布仮設（単語の意味は周囲の単語によって決まる）に基づくものであり、離散オブジェクトを  $D$  次元ベクトルで表現したもの

これにより単語という離散オブジェクトに対して和や差、距離等のオペレーションが可能となる。



# word2vec

Introduction

言語モデル

言語モデル

word2vec の概要

word2vec の概要

Theorem

Proof of Theorem

Algorithms

Conclusions

## word2vec とは

前述したニューラル言語モデルにおける分散表現の獲得は精度が高いが計算量の問題もあったそこで計算量の問題を克服するため「最も簡単な言語モデル」として開発されたのが word2vec である。

word2vec とは単一の手法を指すわけではなく CBOW と skip-gram を合わせたモデル及びツールを指す。対数双線形モデルとも呼ばれる。

## CBOW

隠れ層の出力を計算するときに、CBOW は入力文脈ベクトルの平均を取る。導出手順は以下をとる。

$$\mathbf{h} = \frac{1}{C} \mathbf{W}^T (\mathbf{x}_1 + \mathbf{x}_2 + \dots + \mathbf{x}_c) = \frac{1}{C} (\mathbf{v}_{w_1} + \mathbf{v}_{w_2} + \dots + \mathbf{v}_{w_c})^T \quad (4)$$

ここで  $C$  は文脈中の単語の数、 $w_1 \dots w_c$  は文脈の単語、 $\mathbf{v}_w$  は単語  $w$  の入力ベクトルとする損失関数は以下となる

$$E = -\log p(w_O | w_{I,1}, \dots, w_{I,C}) \quad (5)$$

$$= -u_{j*} + \log \sum_{j=1}^V \exp(u_j) \quad (6)$$

$$= -\mathbf{v}_{w_O}^T \cdot \mathbf{h} + \log \sum_{j=1}^V \exp(\mathbf{v}_{w_j}^T \cdot \mathbf{h}) \quad (7)$$



## word2vec まとめ

基本的には更新式は通常のニューラルネットワークと同じだが文中の単語をまとめる部分だけ通常のネットワークと違う

Introduction

言語モデル

言語モデル

word2vec の概要

**word2vec の概要**

Theorem

Proof of Theorem

Algorithms

Conclusions



# Model (II)

結論 CBOW を用いると分散表現に特化したベクトル表現を得られる.

Introduction

言語モデル

言語モデル

word2vec の概要

**word2vec の概要**

Theorem

Proof of Theorem

Algorithms

Conclusions





# Theorem

[ Theorem ] for identifying parameters  $w_{hv}$ ,  $\mathbf{m}_{hv}$  and  $\Sigma_{hv}$

Vector  $\mathbf{z}_h^k$  obeys mixture normal distribution

$$p(\mathbf{z}_h^k | \theta_h) = \sum_{v=1}^V w_{hv} N_d(\mathbf{z}_{hv}^k, \phi_{hv})$$

thus parameters  $\theta_h = \{w_{hv}, \phi_{hv}\}$  can be derived by applying Expectation Maximization algorithms.

Because vector  $\mathbf{z}_{hv}^k$  obeys  $d$  dimensional normal distribution with mean vector  $\mathbf{m}_{hv}$ , and covariance matrix  $\Sigma_{hv}$  as

$$p(\mathbf{z}_{hv}^k | \phi_{hv}) = N_d(\mathbf{z}_{hv}^k, \phi_{hv}).$$

From the result, each estimation  $\hat{\theta}_h$  of parameter  $\theta_h$  can be derived by maximizing log-likelihood function for incomplete data,

$$L(\hat{\theta}_h) = \sum_{k=1}^K \log p(\mathbf{z}_h^k | \hat{\theta}_h).$$

Update rule of parameter  $\hat{\theta}$  is given by EM algorithms maximizing conditional expectation of log-likelihood function for complete data,

$$\begin{aligned} Q(\hat{\theta}_h | \hat{\theta}_h^{(t)}) &= \mathbb{E}[L(\hat{\theta}_h, v) | \mathbf{z}_h^k, \hat{\theta}_h^{(t)}] \\ &= \sum_{k=1}^K \sum_{v=1}^V h_{hv}^{(t)}(\mathbf{z}_h^k) \log p(\mathbf{z}_h^k, v | \hat{\theta}_l). \end{aligned}$$



# Proof of Theorem

They are iterative calculation form like

$$\begin{aligned}\hat{w}_{hv}^{(t+1)} &= \frac{1}{K} \sum_{k=1}^K h_{hv}^{(t)}(\mathbf{z}_h^k), \\ \hat{\mathbf{m}}_{hv}^{(t+1)} &= \sum_{k=1}^K \gamma_{hv}^{c(t)} \mathbf{z}_h^k, \\ \hat{\Sigma}_{hv}^{(t+1)} &= \sum_{k=1}^K \gamma_{hv}^{k(t)} (\mathbf{z}_h^k - \hat{\mathbf{m}}_{hv}^{(t+1)})(\mathbf{z}_h^k - \hat{\mathbf{m}}_{hv}^{(t+1)})^T\end{aligned}$$

where

$$\begin{aligned}h_{hv}^{(t)}(\mathbf{z}_h^k) &= \frac{\hat{w}_{hv}^{(t)} N_d(\mathbf{z}_{hv}^k, \hat{\phi}_{hv}^{(t)})}{\sum_{v=1}^V \hat{w}_{hv}^{(t)} N_d(\mathbf{z}_h^k, \hat{\phi}_{hv}^{(t)})}, \\ \gamma_{hv}^{k(t)} &= \frac{h_{hv}^{(t)}(\mathbf{z}_h^k)}{\sum_{k=1}^K h_{hv}^{(t)}(\mathbf{z}_h^k)}.\end{aligned}\quad \square$$

## Summary of Parameter Identification

- Observed data vector  $\mathbf{z}_h^k$  consists of faculty vector  $\mathbf{x}^k$  and choice result  $y_h^{k*}$ .
- Utility  $u_{hv}^k$  is given by conditional expectation  $\mathbb{E}[y_{hv}^k | \mathbf{x}^k]$  at value factor level.
- Probability of observed data obeys contaminated normal distribution with  $\theta_h$ .
- Model parameters  $\theta_h$  can be identified by EM algorithms.



# Algorithms

## Algorithms of Parameter Estimation for Modeling Decision Making

- ① Input data vector  $\mathbf{z}_h^k$  which consists of faculty vector  $\mathbf{x}^k$  and choice result  $y_h^{k*}$ .
- ② Give the initial values of parameters  $\theta_h^{(0)}$  that are set of  $\{w_{hv}^{(0)}, \phi_{hv}^{(0)}\}$  where  $\phi_{hv}^{(0)}$  is set of  $\{\mathbf{m}_{hv}^{(0)}, \Sigma_{hv}^{(0)}\}$ .
- ③ Calculate  $N_d(\mathbf{z}_h^k, \phi_{hv}^{(t)})$ ,  $h_{hv}^{(t)}(\mathbf{z}_h^k)$  and  $\gamma_{hv}^{k(t)}$ .
- ④ Update parameters  $\theta_h^{(t+1)}$ .
- ⑤ If the change,  $|\theta_h^{(t+1)} - \theta_h^{(t)}|$ , becomes smaller than given threshold  $\epsilon$ , then the procedure is terminated, otherwise return to Step 3.
- ⑥ From identified parameters  $\theta_h$ , decision making model parameters  $\beta_{hv}^k$  for deriving utility can be estimated.

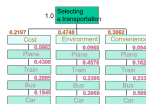


Fig. 9: Example of AHP  
(Case in (K,V,R)=(4, 1, 3))

By incorporating similar concept of AHP (Analytic Hierarchy Process), structure of discrete choice model will be more clearly.



# Conclusions

Introduction

言語モデル

言語モデル

word2vec の概要

word2vec の概要

Theorem

Proof of Theorem

Algorithms

Conclusions

## Model mapping from characteristic to selection

$$E[y_h^k | \mathbf{x}^k] = \sum_{v=1}^V \lambda_{hv}(\mathbf{x}^k) u_{hv}^k \quad \text{where} \quad \lambda_{hv}(\mathbf{x}^k) = \frac{w_{hv} p(\mathbf{x}^k | \phi_{hv}^x)}{\sum_{v'=1}^V w_{hv'} p(\mathbf{x}^k | \phi_{hv'}^x)}$$

## Consider ambiguity to decision making

$$u_{hv}^k = \beta_{hv}^{kT} \mathbf{x}'^k \quad \text{where} \quad \beta_{hv}^{kr} = f(\theta_h)$$

## Algorithms of parameter estimation for modeling discrete choice

$\theta_h = \{w_{hv}, \phi_{hv}\}$  can be estimated by Expectation Maximization algorithms

## Future plans

Apply the concept of discrete choice model.

- ▶ Development of attractive product and service on limited budget,
- ▶ Consensus building among stakeholder,
- ▶ Objective decision making and support mental training, etc.