

言語統計力学シンポジウム

単語分散表現の結合学習による 通時的な単語の意味変化の検出

相田太一

東京都立大学 自然言語処理研究室

背景：単語の意味・用法の変化

単語は以下の条件で異なる使われ方をする

- 異なる**時期**

- meat: 食べ物全般(**古英語**) → 動物の肉(**近代英語**)

- 了解: 理解(**戦前**) → 承知(**戦後**)

- 異なる**分野**

- interface: 境界(**一般**) → ソフトウェア(**情報**)

→この違いの検出は言語学・社会学・辞書学において有用

背景：単語の意味・用法の変化

単語は以下の条件で異なる使われ方をする

- 異なる**時期**

- meat: 食べ物全般(**古英語**) → 動物の肉(**近代英語**)
- 了解: 理解(**戦前**) → 承知(**戦後**)

- 異なる**分野**

- interface: 境界(**一般**) → ソフトウェア(**情報**)

→この**違いの検出**は言語学・社会学・辞書学において有用

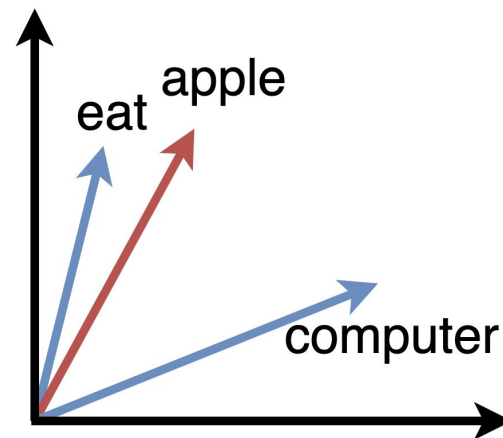
どのように？
例) 各単語の全用例を確認して判別...

背景：単語の意味・用法の変化

情報科学では**単語ベクトル**による調査が主流

周辺の単語情報から学習（例：apple 🍏 , 💻）

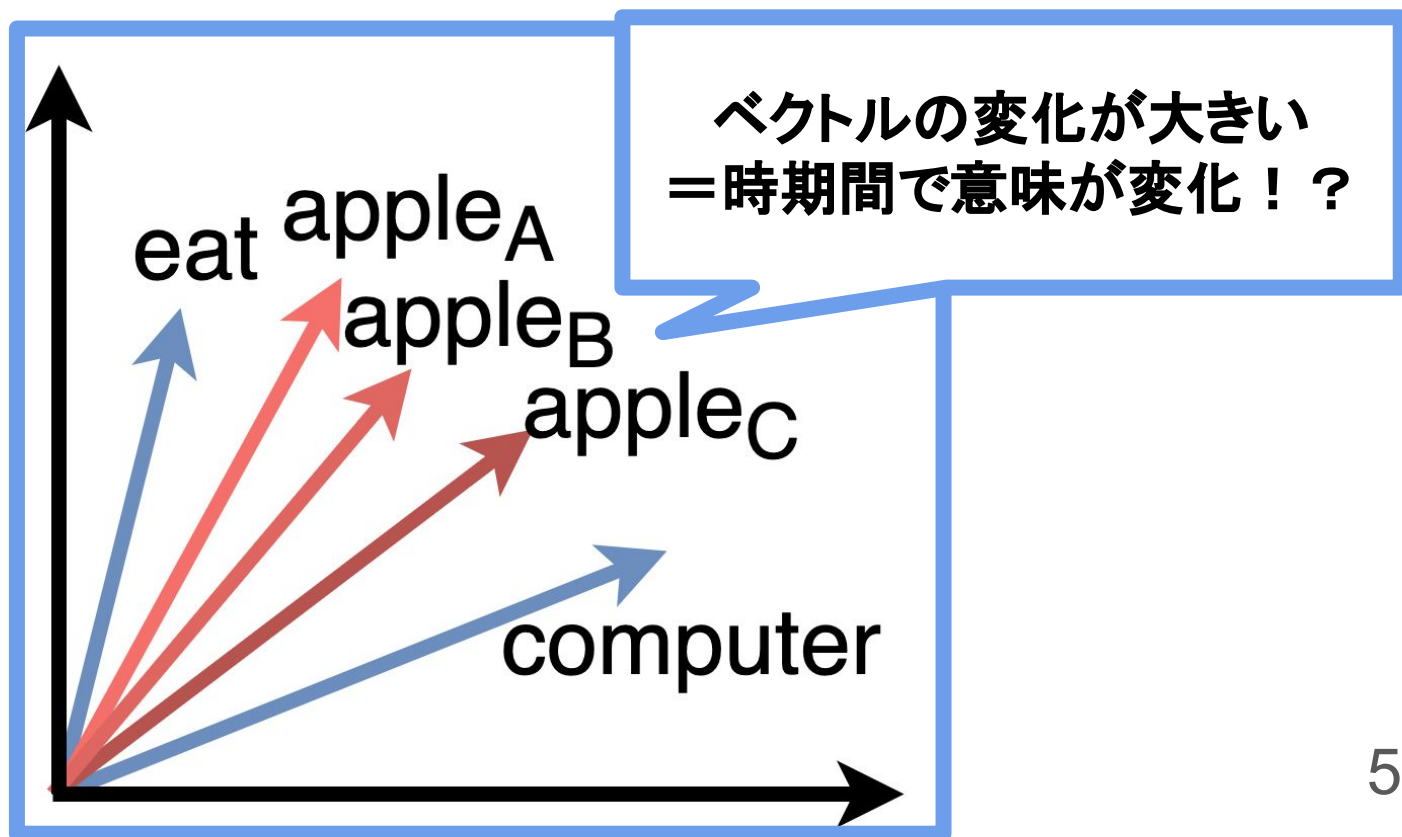
- I eat an **apple** every morning .
- He is eating an **apple** .
- I am an **apple** user.



背景：単語の意味・用法の変化

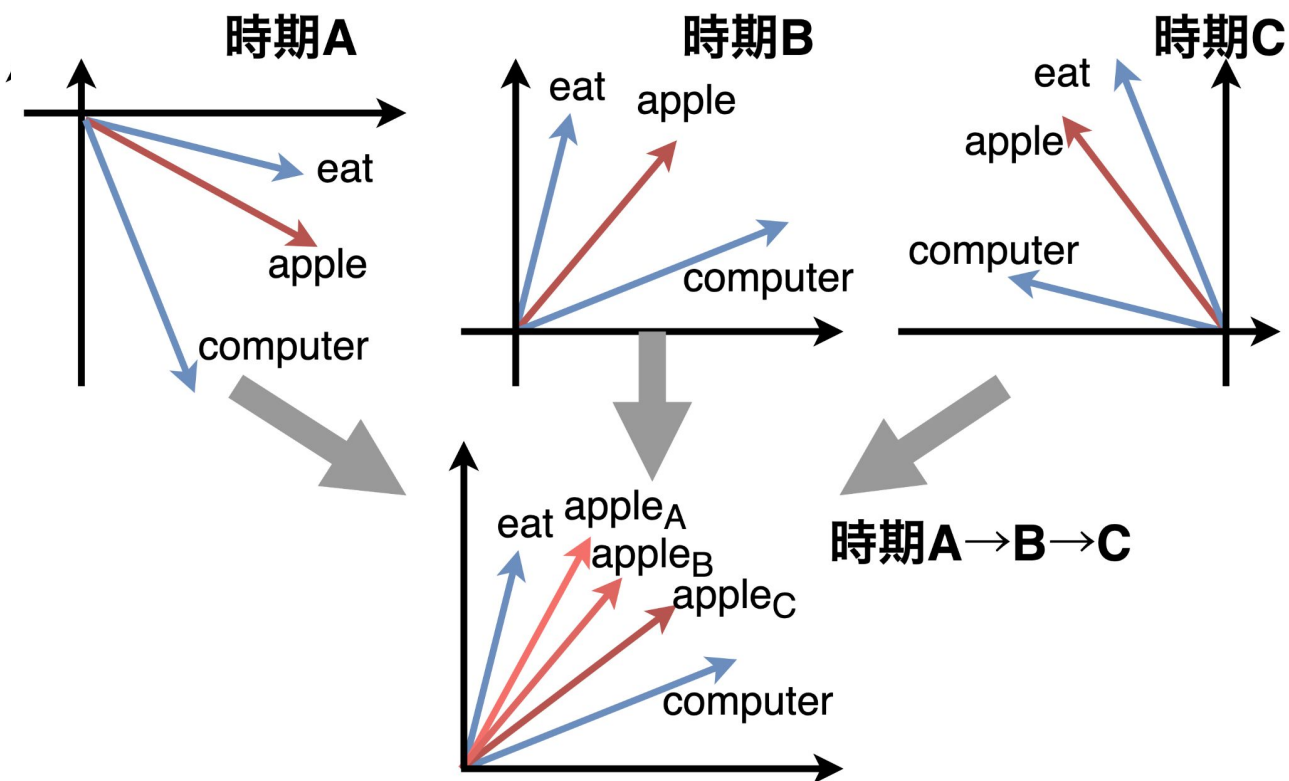
情報科学では**単語ベクトル**による調査が主流

例) 時期A→B→Cにおける単語 **“apple”** の変化



背景：単語の意味・用法の変化

単語ベクトルは文書ごと(時期A, B, C だと3つ)
→ 時期・分野間に対応した単語ベクトルを
どのように学習するか？

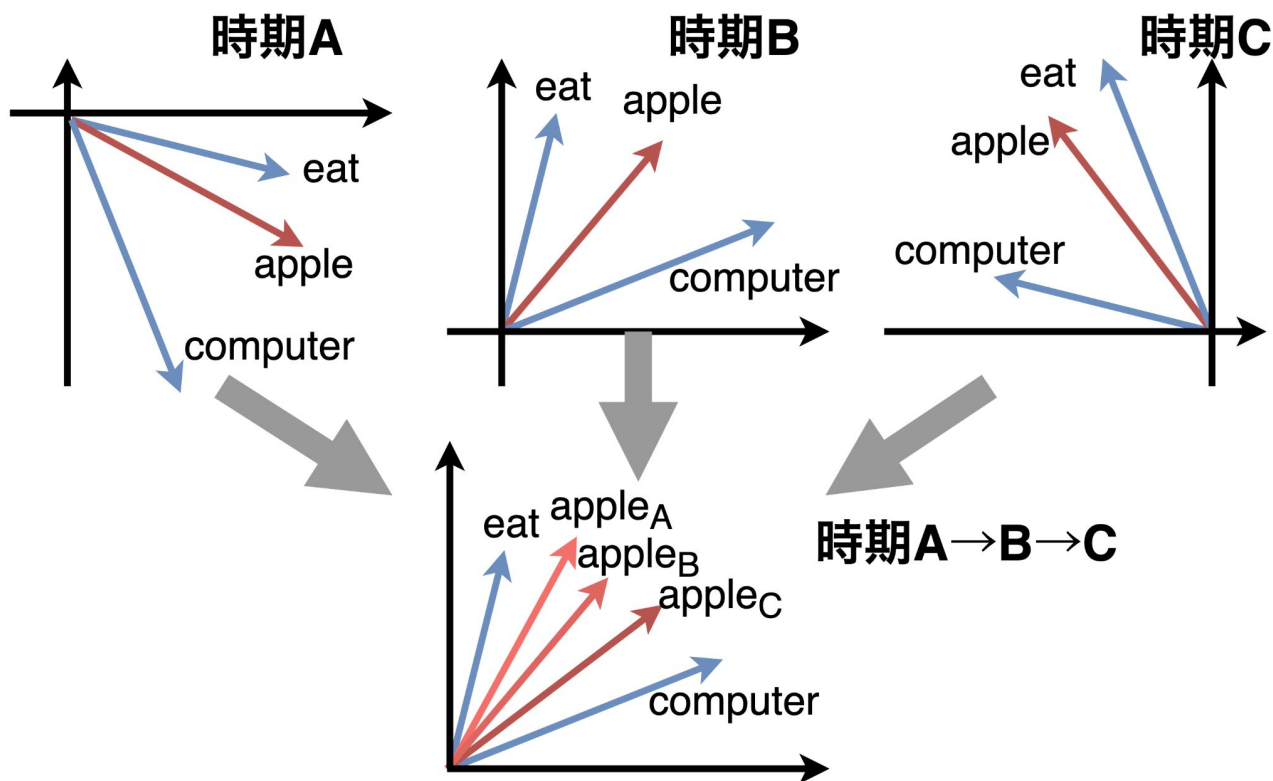


先行研究: alignment (Hamilton+16)

- 各時期で訓練し、**回転させて対応づけ**

✓ 簡単、計算コストが低い

✗ 強い仮定「各ベクトルは**線形変換で対応づけ**できる」



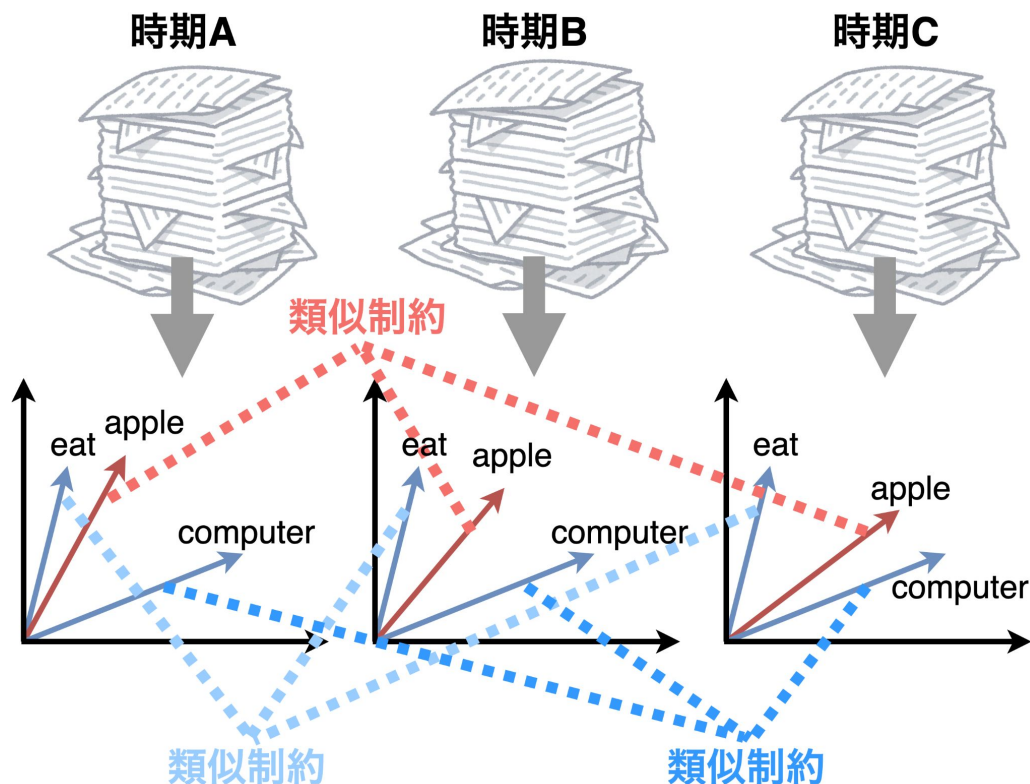
先行研究: alignment を避ける手法 (1)

Dynamic Word Embeddings: DWE (Yao+18)

- 時期・分野間で**ベクトルを同時に学習**する

✓ 同時に訓練することで、alignment が不要

✗ 3つの**ハイパーパラメータの調整**が必要



先行研究: alignment を避ける手法 (2)

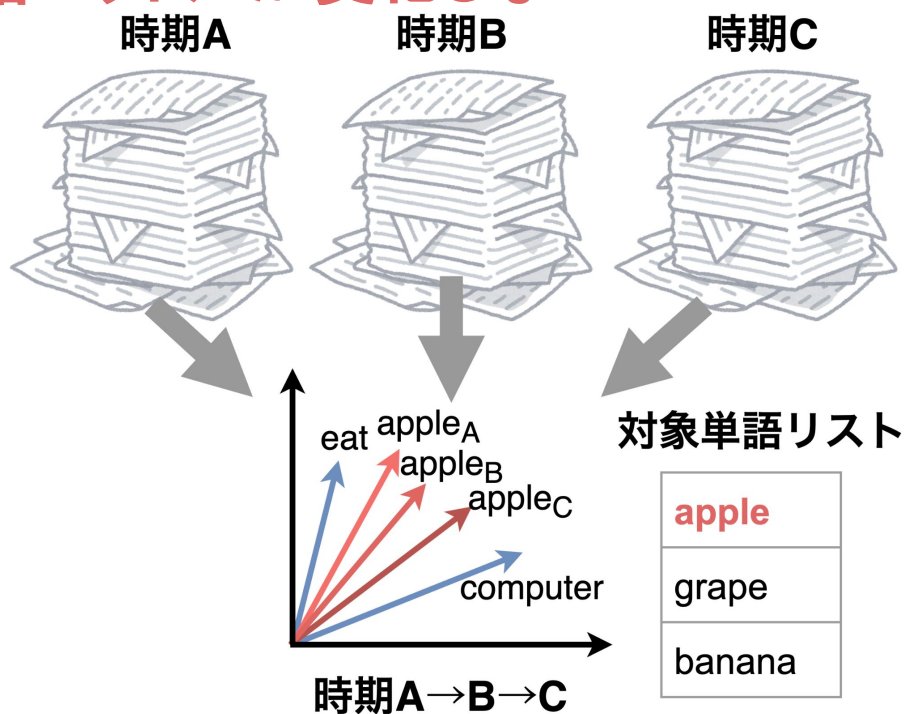
Temporal Referencing: TR (Dubossarsky+19)

- 文書を1つに結合し、**対象単語だけ区別**する

✓ 簡単、計算コストが低い

✗ **対象単語リスト**を用意する必要がある

✗ **周辺単語ベクトルが変化しない**



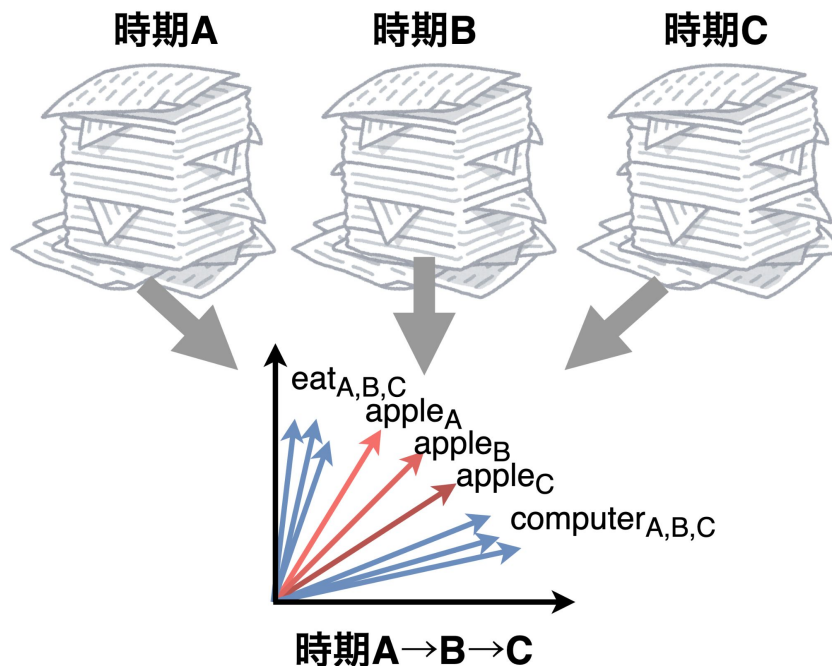
本研究: TR に対する2つの拡張

✗ **対象単語リスト**を用意する必要がある

→ 語彙に含まれる全ての単語を対象語にする

✗ **周辺単語ベクトル**が変化しない

→ 周辺単語ベクトルの変化を考慮する



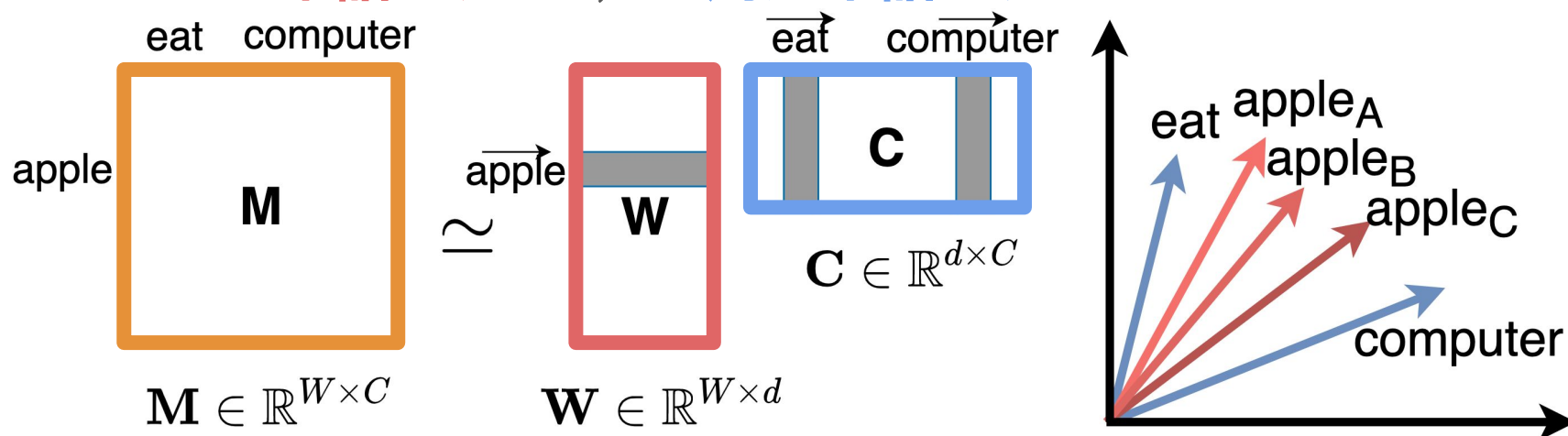
基盤となる手法: PMI-SVD (Levy+14)

1. 対象語-周辺単語 の共起行列から **PMI 行列** を計算

$$SPPMI(\text{apple}, \text{eat}) = \max\left(\log \frac{P(\text{apple}, \text{eat})}{P(\text{apple})P(\text{eat})} - \log k, 0\right)$$

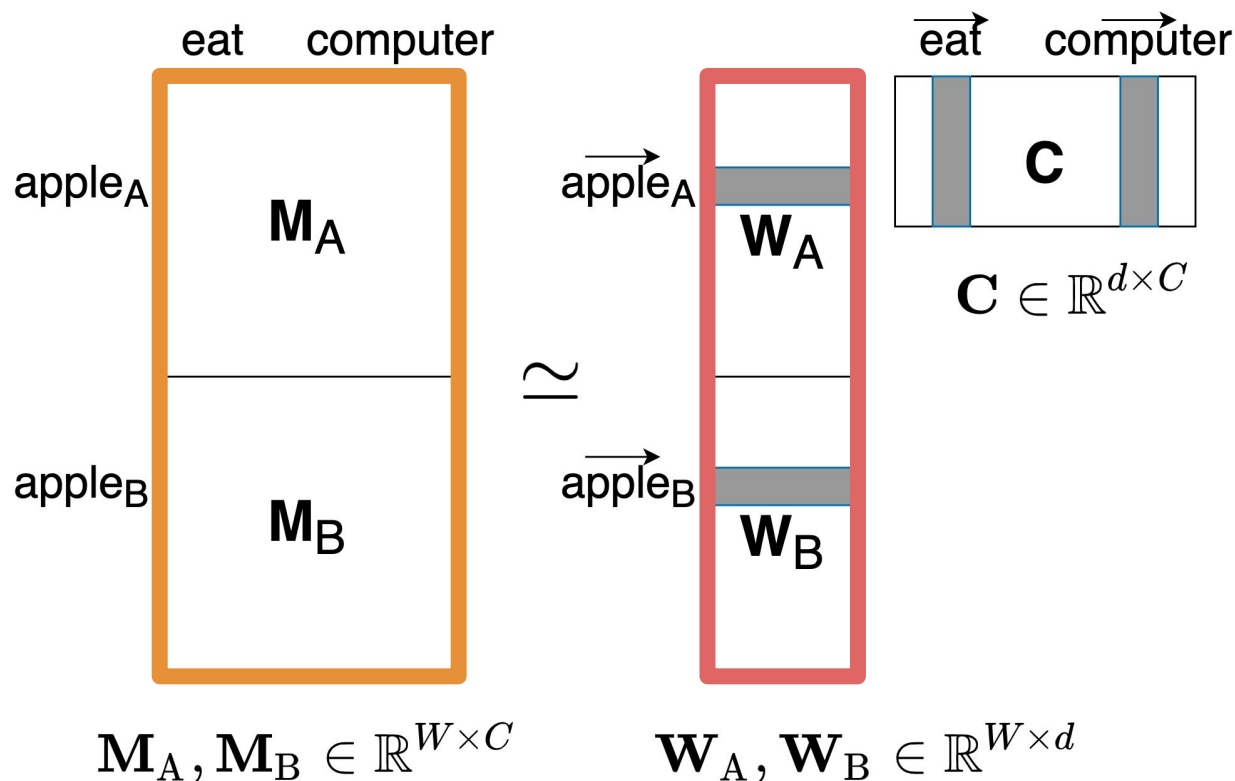
2. PMI 行列を SVD で次元削減

- **W: 単語ベクトル**, **C: 周辺単語ベクトル**



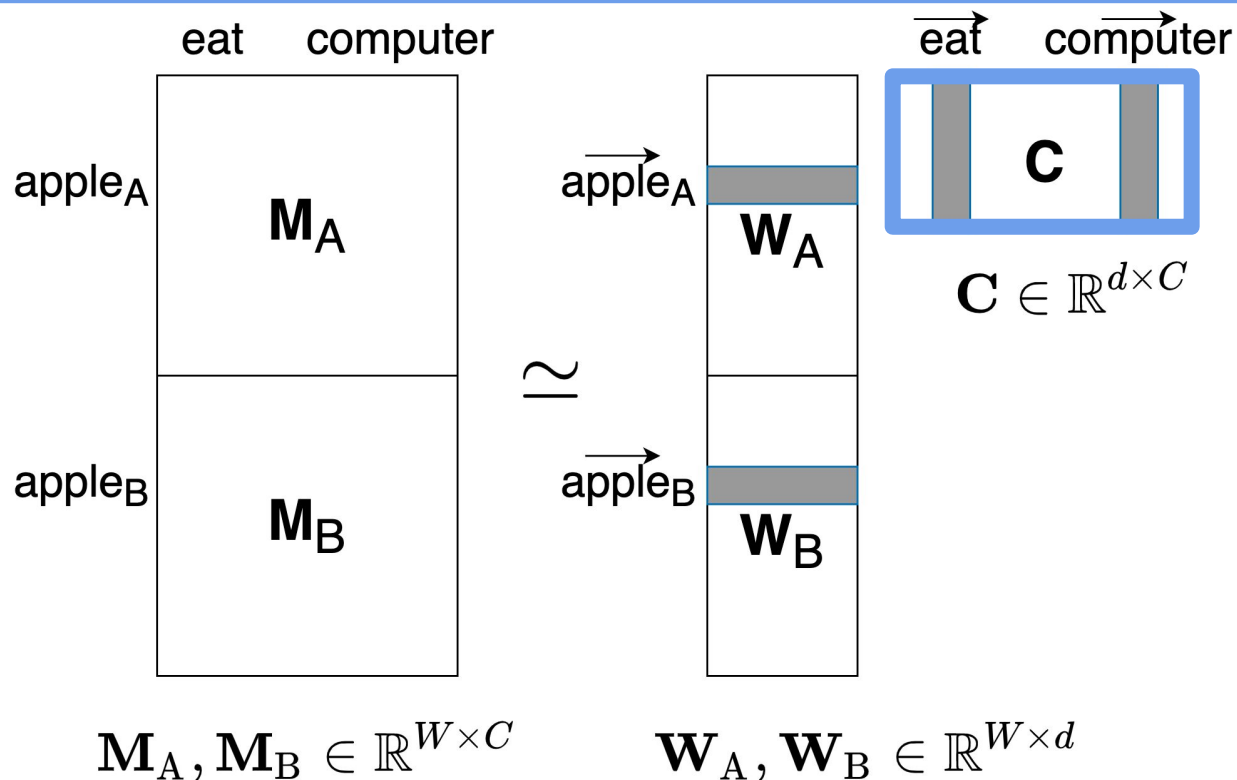
拡張1: 語彙に含まれる全ての単語を対象語にする (PMI-SVDjoint)

1. PMI 行列を各時期で計算
2. PMI 行列を結合し、同時に次元削減



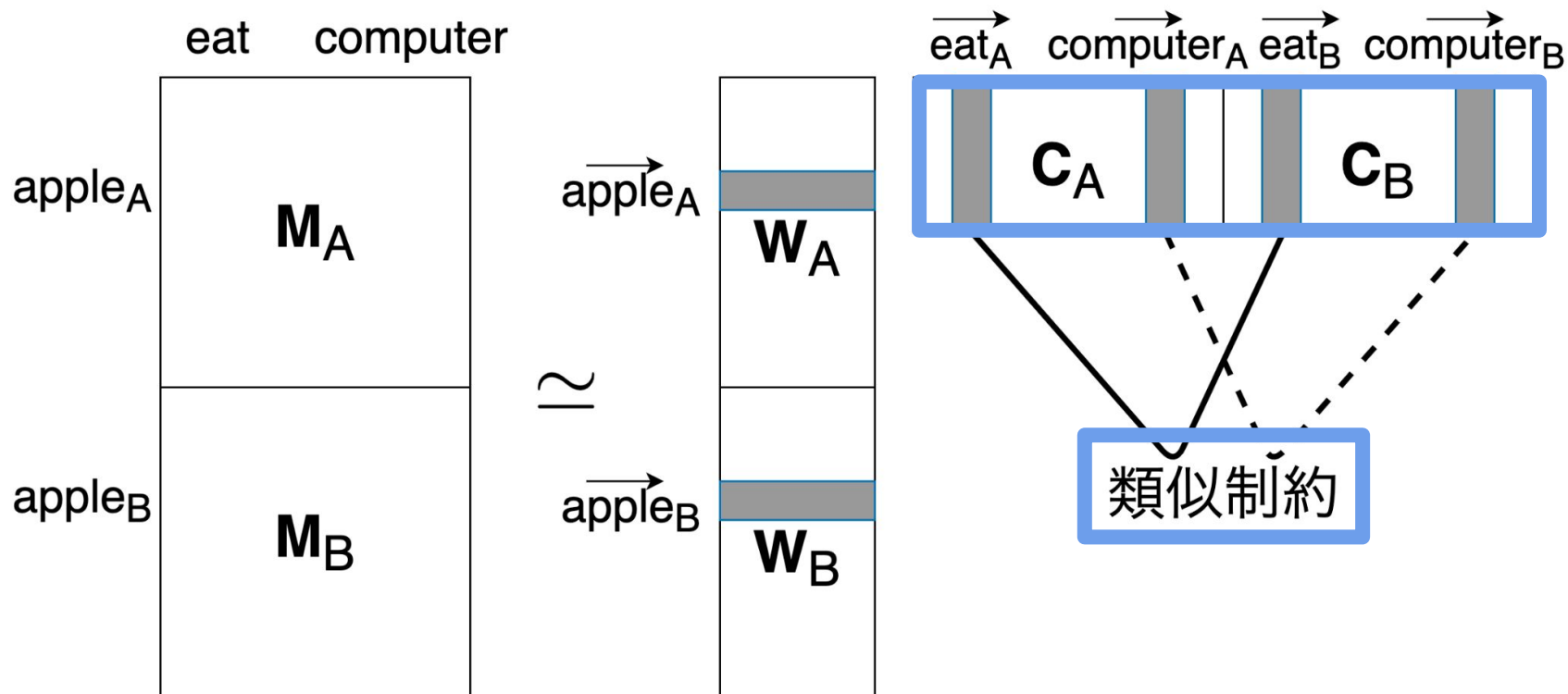
拡張1: 語彙に含まれる全ての単語を 対象語にする (PMI-SVDjoint)

1. この方法では 周辺単語ベクトルは1つ
2. → 文書間で周辺単語ベクトルが変化しない
(TR の2つ目の問題)



拡張2: 周辺単語ベクトルの変化を考慮 (PMI-SVDc)

- DWEと同様に**制約条件**を追加
- **ただ1つのハイパーパラメータを持つ**
(DWE は3つ)



実験：拡張の効果

- **先行研究 TR** に対する**2つ拡張**の効果を検証
 - **PMI-SVDtr**(先行研究)
 - **PMI-SVDjoint**
 - **PMI-SVDc**
- **データ・タスク: SemEval-2020 Task 1**
 - 2つの時期間で与えられた**対象単語の意味変化**を予測
 - **4つの言語**(英語、ドイツ語、ラテン語、スウェーデン語)
 - **2つのタスク**(意味変化の有無を分類、意味変化の度合いで並べ替え)

結果：意味変化の有無を分類

- 提案した拡張手法による性能の向上を確認

手法	分類正解率 (Acc.)				
	En	De	La	Sv	Avg.
PMI-SVD _{tr}	0.622	0.625	0.525	0.613	0.596
PMI-SVD _{joint}	0.649	0.708	0.525	0.677	0.640
PMI-SVD _c	0.649	0.667	0.650	0.613	0.645

結果:意味変化の度合いで順位付け

- 提案した拡張手法による性能向上
- **PMI-SVD_c はどの言語でも安定して機能する**

手法	スピアマンの順位相関 (ρ)				
	En	De	La	Sv	Avg.
PMI-SVD _{tr}	0.487	0.527	0.123	0.257	0.349
PMI-SVD _{joint}	0.438	0.540	0.141	0.478	0.399
PMI-SVD_c	0.424	0.597	0.328	0.328	0.433

結果：意味変化の度合いで順位付け

- ラテン語で大幅な性能の差
 - 時期は **-200~0年→0~2000年** と非常に広い
 - 意味が変わらなくとも、表現が変わるのでは
- **周辺単語の変化を考慮する必要性**

手法	スピアマンの順位相関 (ρ)				
	En	De	La	Sv	Avg.
PMI-SVD _{tr}	0.487	0.527	0.123	0.257	0.349
PMI-SVD _{joint}	0.438	0.540	0.141	0.478	0.399
PMI-SVD _c	0.424	0.597	0.328	0.328	0.433

実験: 語彙全体から意味変化した単語を検出 英語(1900s→1990s), 日本語(戦前→戦後)

- 提案手法

- PMI-SVD joint: 語彙に含まれる**全ての単語を対象語**にする
- PMI-SVD c: **周辺単語ベクトルの変化を考慮**する

- ベースライン

- PMI-SVD align (Hamilton+16)
 - Word2Vec align (Hamilton+16)
 - Dynamic word embeddings (Yao+18)
 - BERT (Martinc+20)
- alignment を用いた手法
- alignment を避ける手法 (同時学習)

alignment を避ける手法
(事前訓練済み言語モデル)

実験: 語彙全体から意味変化した単語を検出 英語(1900s→1990s), 日本語(戦前→戦後)

- 評価: 平均逆順位(MRR)

- 異なる時期、同じ単語の余弦類似度を計算し、類似度が低い順(=変化した順)に並べてリストを作成
- 評価セットに含まれる単語の MRR を算出

意味変化の度合いで
並べ替えたリスト

順位	単語
...	...
4	mouse
5	apple
...	...
14	web
...	...
100	computer

評価セット
(意味変化した単語)

apple
mouse
web

目的の単語をどれだけ上位で捉えられるか?
(高いほどよい)

$$MRR = \frac{1/5 + 1/4 + 1/14}{3}$$

結果：意味変化した単語の検出性能

- 2つの拡張手法 (PMI-SVD joint, c) は**両言語で優れた性能**

手法	英語	日本語	訓練時間
PMI-SVD _{joint}	0.00186	0.00131	2m58s
PMI-SVD _c	0.01045	0.00120	26m01s
Word2Vec _{align}	0.00040	0.00137	6m22s
PMI-SVD _{align}	0.00100	0.00091	3m26s
DWE	0.00047	0.00058	30h20m
<i>BERT*</i>	<i>0.00250</i>	<i>0.00163</i>	2h23m

結果：訓練時間の比較

- PMI-SVD joint が**最も高速**
- PMI-SVD c は**ハイパラ探索時間を大幅減**

手法	英語	日本語	訓練時間
PMI-SVD _{joint}	0.00186	0.00131	2m58s
PMI-SVD _c	0.01045	0.00120	26m01s
Word2Vec _{align}	0.00040	0.00137	6m22s
PMI-SVD _{align}	0.00100	0.00091	3m26s
DWE	0.00047	0.00058	30h20m
<i>BERT*</i>	<i>0.00250</i>	<i>0.00163</i>	2h23m

実験: 対象データで BERT の再訓練

- 今回使った BERT (BERT-base) は**大規模な外部データで事前訓練済み**
- 実験に使ったデータで**サイズの異なる BERT のモデルを訓練** (-tiny < -mini < -base)

手法	英語	日本語	訓練時間
PMI-SVD _{joint}	0.00186	0.00131	2m58s
PMI-SVD _c	0.01045	0.00120	26m01s
BERT-tiny	0.00100	0.00078	12days
BERT-mini	0.00135	0.00119	2weeks

結果: 対象データで BERT の再訓練

- 提案した拡張手法は**高い検出性能**を示す
- **訓練時間も BERT に比べ非常に高速**

手法	英語	日本語	訓練時間
PMI-SVD _{joint}	0.00186	0.00131	2m58s
PMI-SVD _c	0.01045	0.00120	26m01s
BERT-tiny	0.00100	0.00078	12days
BERT-mini	0.00135	0.00119	2weeks

分析

- 先行研究の分析：**意味変化が自明な単語のみ**（
“gay” 陽気な→同性愛）

× 対象データの選び方で変化する単語は変わるため、
幅広い分析が必要

- 本研究では、

- ✓ 各手法が意味変化を予測した単語
- ✓ 意味変化が自明な単語・自明でない単語

について、**網羅的に分析**

分析: 各手法が意味変化を予測した単語

- どちらの手法も**自明でない単語の変化を検出**

順位	BERT		PMI-SVD _c	
	単語	説明	単語	説明
1	若く	匹敵, 年齢 → 年齢	行い	振舞い → 振舞い, 実行
2	触れ	言及, 抵触 → 言及, 触る	かねて	以前 → 以前, 同時
3	行い	振舞い → 振舞い, 実行	おまけ	追加 → 追加, 減額
4	公明	公正 → 組織名, 公正	無論	(副詞的用法)
5	思い	思考, 感情 → 思考	年中	1年, 役職 → 1年
6	削除	文字の削除	キー	音楽, 人名 → 音楽, 鍵
7	在り	(物理的) → (概念的)	欠け	(物理的) → (概念的)
8	参議	参与 → 組織名	皆無	全然ないこと
9	欠け	(物理的) → (概念的)	馬場	人名, 芝
10	幼稚	幼い → 幼稚園, 幼い	反面	反対, 一方 → 一方

分析: 各手法が意味変化を予測した単語

- 捉える変化の大きさが違う → 文脈窓幅の違い

順位	BERT		PMI-SVD _c	
	単語	説明	単語	説明
1	若く	匹敵, 年齢 → 年齢	行い	振舞い → 振舞い, 実行
2	触れ	言及, 抵触 → 言及, 触る	かねて	以前 → 以前, 同時
3	行い	振舞い → 振舞い, 実行	おまけ	追加 → 追加, 減額
4	公明	公正 → 組織名, 公正	無論	(副詞的用法)
5	思い	思考, 感情 → 思考	年中	1年, 役職 → 1年
6	削除	文字の削除	キー	音楽, 人名 → 音楽, 鍵
7	在り	(物理的) → (概念的)	欠け	(物理的) → (概念的)
8	参議	参与 → 組織名	皆無	全然ないこと
9	欠け	(物理的) → (概念的)	馬場	人名, 芝
10	幼稚	幼い → 幼稚園, 幼い	反面	反対, 一方 → 一方

分析: 自明でない単語の近傍単語

「欠け(物理的→概念的)」

- どちらの手法も傾向を捉えている

BERT		PMI-SVD _c	
戦前	戦後	戦前	戦後
マイナス	欠如	切り	有し
決まり	乏しい	切ら	欠如
構え	不足	諦め	富ん
重み	崩れ	箸	づけ
当て	破れ	つける	把握

分析：自明な単語の近傍単語

「了解（理解→承諾）」

- どちらも変化の傾向を捉えている
 - **BERT は変化後の意味を事前に検出**
- **大規模な現代のデータによるバイアス**

BERT		PMI-SVD _c	
戦前	戦後	戦前	戦後
承諾	承諾	理解	承諾
承知	承知	納得	承知
納得	承認	推測	納得
理解	同意	判断	同意
断定	納得	断定	理解

結論・今後の展望

- 異なる時期・分野に対応した単語ベクトルを学習する既存の手法を拡張
 - 語彙中の全ての単語を調査対象単語にする
 - 周辺単語ベクトルの変化を考慮する
- 実験より、
 - 先行研究より高速に学習し、同等以上の性能
 - 意味変化が自明・自明でない単語を効果的に捉える
- 今後は以下2つを検討
 - 意味変化の種類を予測するには？
 - 意味変化とベクトル空間には関係があるか？