

- 課題: n-gramのトレンド線は**複数の語彙的成分が混合し、解釈が難しい**
- 手法: n-gramの出現頻度を**カクテルパーティー問題に置き換え、ICAを用いた**
- 結果: **出現頻度が語彙別に分解されたことを定性的に確認した**

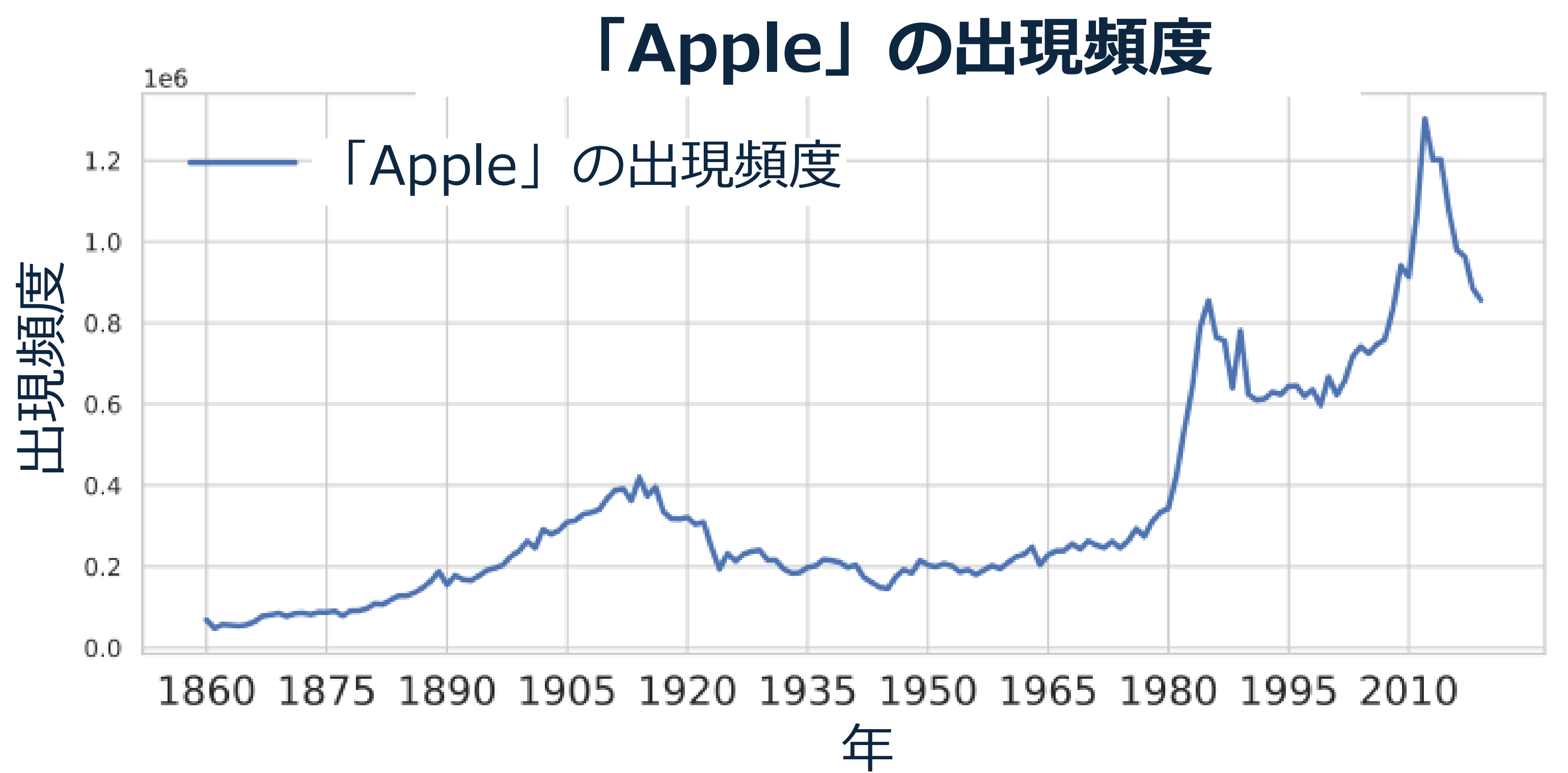
背景 : n-gramを活用したトレンド分析における課題

課題

n-gram の出現頻度はトレンド解析に広く用いられるが、表層形依存により**同型類語が混在してしまう**。そのため、語の起源の特定や高精度な流行語の把握が困難となる

目的

本研究はn-gramにおける**出現頻度を語彙別成分に分解**することを目的とする。



提案手法 : カクテルパーティー問題として捉えるn-gramにおける出現頻度

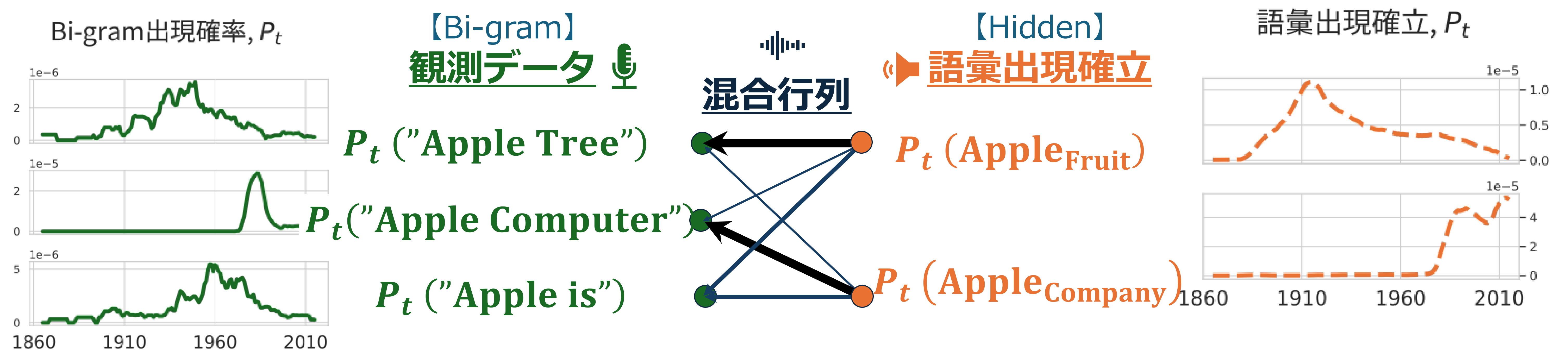
1. 観測データ = **独立成分**の線型結合

2. n-gramの出現確率を**語彙成分**の線結合と捉える

$$x_t = A S_t$$

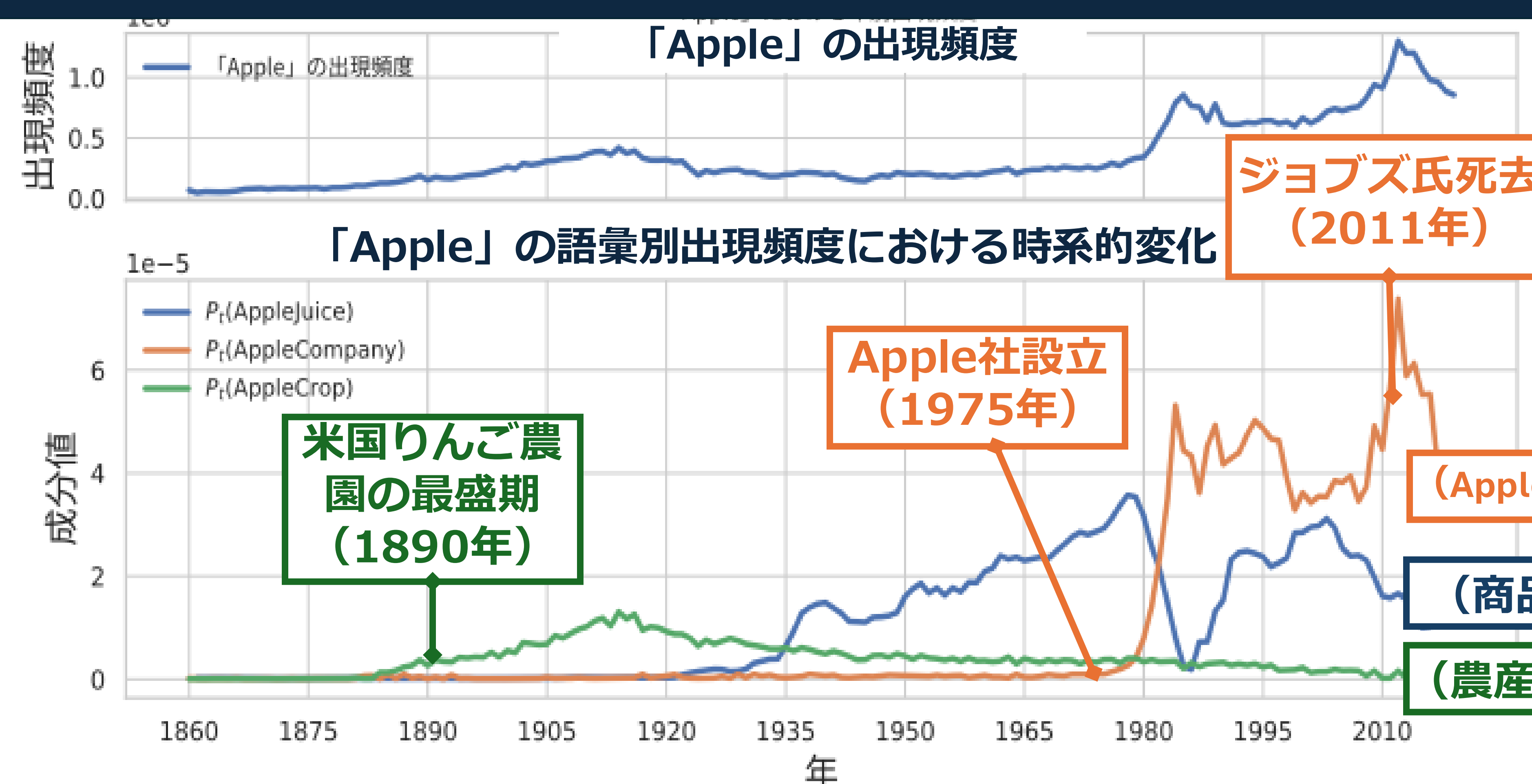
観測データ (ベクトル) 混合行列 独立成分 (ベクトル)

$$P_t(\text{"Apple Tree"}) = P(\text{"Apple Tree"} | \text{AppleFruit}) P_t(\text{AppleFruit}) + P(\text{"Apple Tree"} | \text{AppleCompany}) P_t(\text{AppleCompany})$$



独立成分分析 (ICA) : Zipff's Law考慮のためベータ事前分布付き二項 MAP ICAを活用。

結果 : トレンドデータを語彙別の成分に分解



結論・今後の展望

ICAの活用により、n-gramのトレンド線を語彙別に分解することが可能であることを定性的に確認した。

今後はこのようなデータを知識グラフに含めることにより、言葉の関係性を考慮した、トレンド線の回帰予測が可能なのかを検証して行く。