

本人が一番わかって いないLDA説明スライド

神田拓実

LDAとは？

- Latent Dirichlet Allocation (潜在的ディリクレ配分法)
- トピックモデルと呼ばれる手法の1つ
- 文章解析に用いられる
- 1つの文章が複数のトピックを持つものとして説明できる
- PLSIをベイズ化したもの

学がなさすぎてわからん

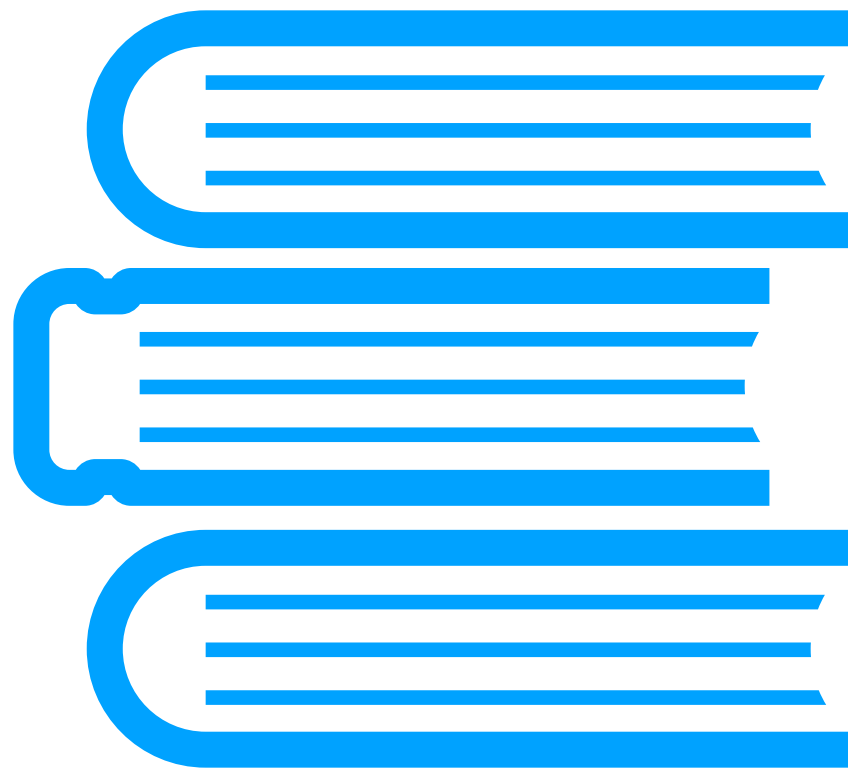


よくわかる解説

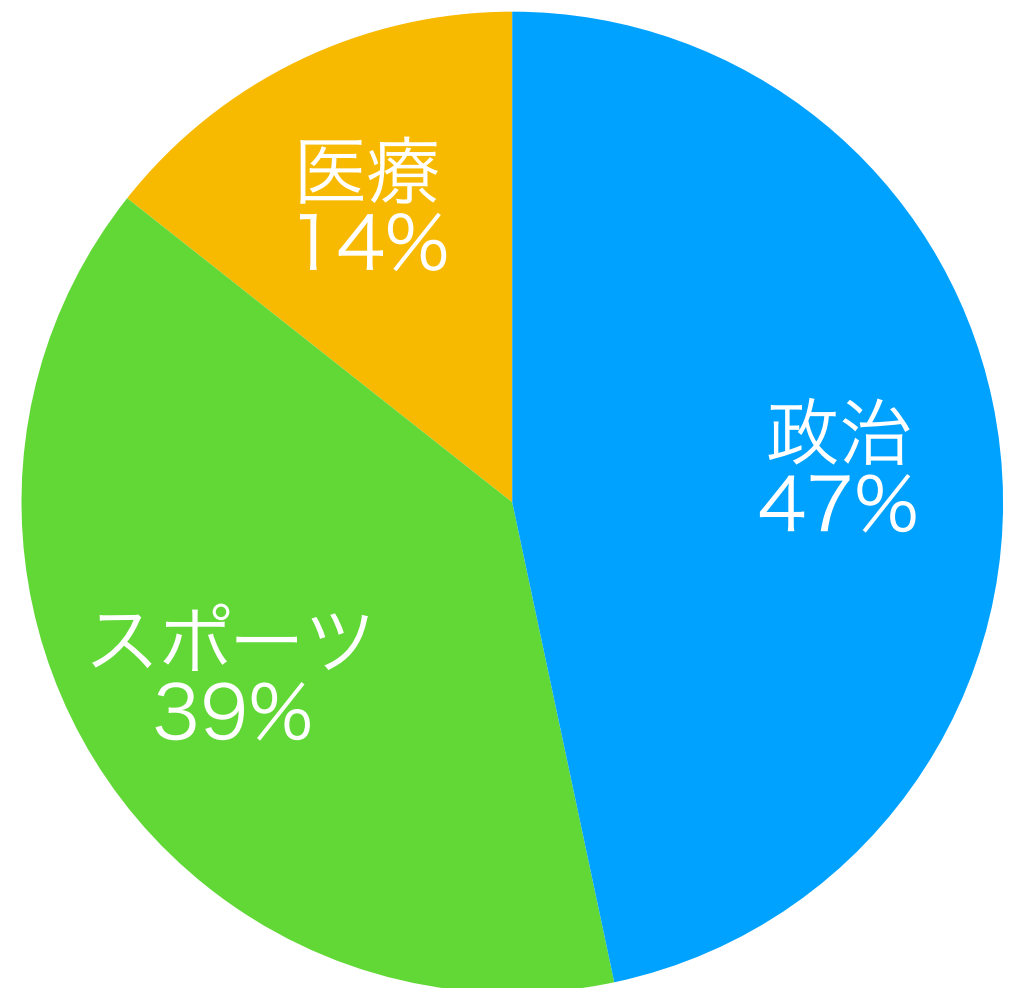
LDAは「トピック分布」と「単語分布」を用いて
文章をモデル化する

トピック分布とは

- 文章が持っているトピックを割合で表したもの

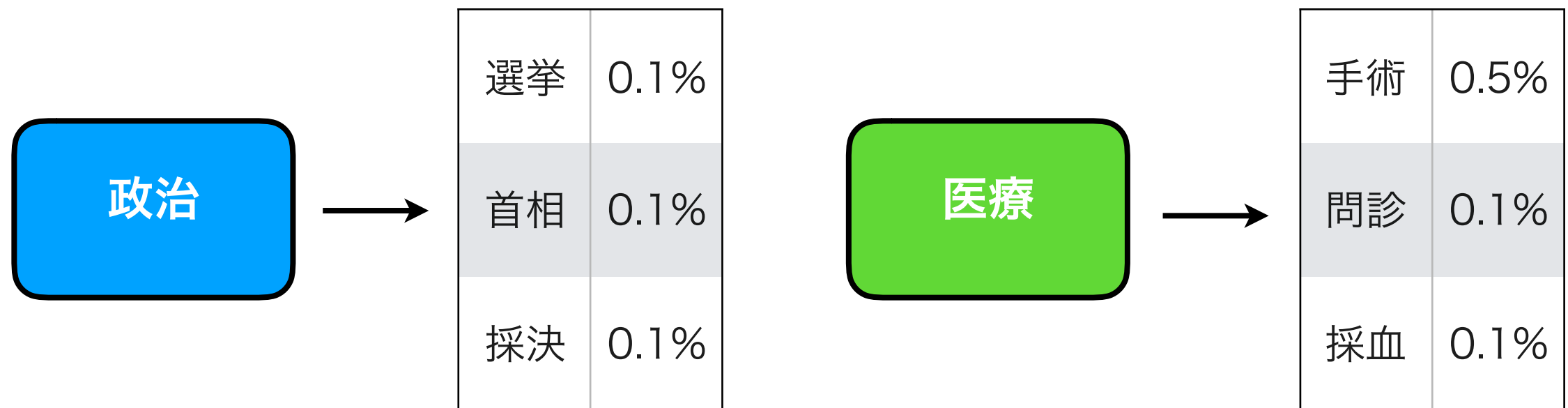


文章A



単語分布とは

- トピックに関連する単語を割合の形で表したもの



テキストコーパスを用いてLDAを「学習」することで、各文章のトピック分布と各トピックの単語分布を求めることができる

LDAの生成過程

1. 各文章の単語について、ランダムにトピックを割り当てる
2. 割り当てられたトピックから、文章ごとのトピック確率を計算する
3. 割り当てられたトピックから、トピックごとの単語確率を計算する
4. 2と3の積で計算される確率をもとに、各文章の単語にトピックを再び割り当てる
5. 2、3、4を収束条件まで繰り返す

LDAに関するあれこれ

- トピック分類は
「トピック 1 の割合40%、トピック 2 の割合60%」の
ような形で得られるため、
**トピック1が実際にどのような話題であるかは人間が
解釈する必要がある**
例：「トピック 1 はホームラン、盗塁などの単語が並ん
でいるから野球トピックだな」
- **LDAは教師なし学習である**ため、文章コーパスによっ
ては「標準語」「関西弁」のようなトピックが出現す
ることも考えられる