

Text Mining w R

WhyR? 2017 - warsztat

Norbert Ryciak

Politechnika Warszawska, Sages

4 października 2017



- Wydział Matematyki i Nauk Informacyjnych PW
 - ▶ doktorant (zastosowanie deep learningu w analizie tekstu)



- Instytut Podstaw Informatyki PAN
 - ▶ projekt CLARIN (zadanie analizy sentymentu)



- Sages
 - ▶ opiekun bootcampu Data Science, trener

Latent Dirichlet Allocation

Model LDA - wprowadzenie

Motywacja: przedstawienie tekstu jako mieszanki tematów.

Model LDA - wprowadzenie

Motywacja: przedstawienie tekstu jako mieszanki tematów.

Temat - rozkład prawdopodobieństwa na zbiorze słów.

Motywacja: przedstawienie tekstu jako mieszanki tematów.

Temat - rozkład prawdopodobieństwa na zbiorze słów.

Przykład:

- Mam gorączkę i katar.
- Graliśmy w siatkówkę.
- Sport to zdrowie.

Model LDA - wprowadzenie

Motywacja: przedstawienie tekstu jako mieszanki tematów.

Temat - rozkład prawdopodobieństwa na zbiorze słów.

Przykład:

- Mam gorączkę i katar.
- Graliśmy w siatkówkę.
- Sport to zdrowie.

Założmy, że mamy dwa tematy: „zdrowie” i „sport”. Wówczas:

Model LDA - wprowadzenie

Motywacja: przedstawienie tekstu jako mieszanki tematów.

Temat - rozkład prawdopodobieństwa na zbiorze słów.

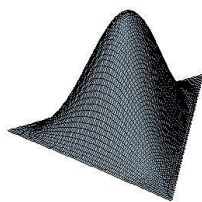
Przykład:

- Mam gorączkę i katar.
- Graliśmy w siatkówkę.
- Sport to zdrowie.

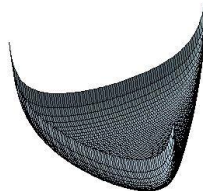
Założmy, że mamy dwa tematy: „zdrowie” i „sport”. Wówczas:

- Pierwsze zdanie = 100% zdrowie
- Drugie zdanie = 100% sport
- Trzecie zdanie = 50% sport + 50% zdrowie

Rozkład Dirichleta



(a) $\alpha = 3$



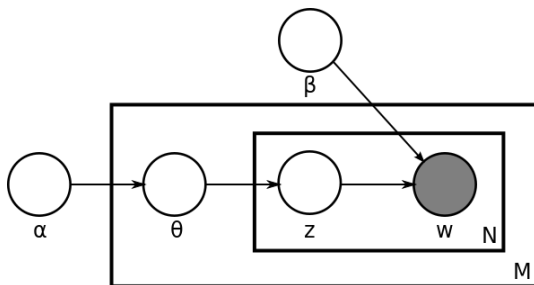
(b) $\alpha = 0.95$

Rysunek: Gęstość trójwymiarowego rozkładu Dirichleta $\text{Dir}(\alpha)$.

Wektor losowy (x_1, \dots, x_K) z K -wymiarowego rozkładu Dirichleta to punkt na $(K - 1)$ -wymiarowym sympleksie, czyli $x_1 + \dots + x_K = 1$, $x_i \geq 0$.

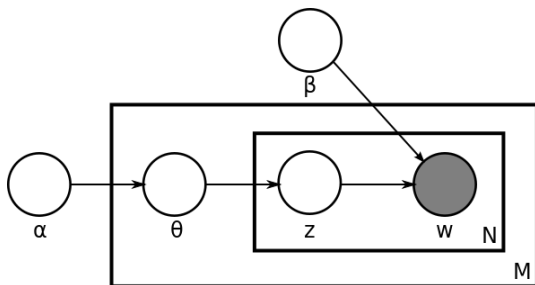
Model LDA - sformułowanie

LDA - Latent Dirichlet Allocation.



Model LDA - sformułowanie

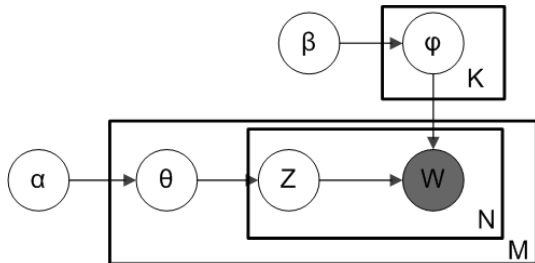
LDA - Latent Dirichlet Allocation.



Proces generowania dokumentu d :

- ❶ Ustal liczbę słów w dokumencie N_d .
- ❷ Losuj rozkład tematów w dokumencie $\theta_d \sim \text{Dir}(\alpha)$.
- ❸ Dla $i = \{1, \dots, N_d\}$:
 - ❶ Losuj temat z_{di} z rozkładu dyskretnego θ_d .
 - ❷ Losuj słowo w_{di} z rozkładu dyskretnego $\beta_{z_{di}}$

Wyglądający model LDA



$\varphi_k \sim \text{Dir}(\beta)$ dla $k = 1, \dots, K$,
 $\beta = (\beta_1, \dots, \beta_V)$, $\beta_i \in \mathbb{R}$.

Estymacja w modelu LDA - próbkowanie Gibbsa

Estymacja w modelu LDA - próbkowanie Gibbsa

- 1 Inicjalizacja: każde słowo przypisz do losowego tematu.

Estymacja w modelu LDA - próbkowanie Gibbsa

- ❶ Inicjalizacja: każde słowo przypisz do losowego tematu.
- ❷ Dla każdego słowa w_{dj} , $j = 1, \dots, N_d$, $d = 1, \dots, M$, powtarzaj:
 - ▶ Oblicz $P(z_{dj} = k | \mathbf{w}, \mathbf{z}_{-dj}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ dla $k = 1, \dots, K$.
 - ▶ Wylosuj z_{dj} z powyższego rozkładu.

Estymacja w modelu LDA - próbkowanie Gibbsa

- ❶ Inicjalizacja: każde słowo przypisz do losowego tematu.
- ❷ Dla każdego słowa w_{dj} , $j = 1, \dots, N_d$, $d = 1, \dots, M$, powtarzaj:
 - ▶ Oblicz $P(z_{dj} = k | \mathbf{w}, \mathbf{z}_{-dj}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ dla $k = 1, \dots, K$.
 - ▶ Wylosuj z_{dj} z powyższego rozkładu.

$$P(z_{dn} = k | \mathbf{z}_{-(dn)}, \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \propto \frac{n_{k, w_{dn}}^{(\cdot)} + \beta}{n_{k, (\cdot)}^{(\cdot)} + \beta V} \cdot \frac{n_{k, (\cdot)}^d + \alpha_k}{N_d + \sum_{k=1}^K \alpha_k},$$

$n_{k, w}^d$ - liczność przypisać słowa w do tematu k w dokumencie d .

Estymacja w modelu LDA - próbkowanie Gibbsa

- 1 Inicjalizacja: każde słowo przypisz do losowego tematu.
- 2 Dla każdego słowa w_{dj} , $j = 1, \dots, N_d$, $d = 1, \dots, M$, powtarzaj:
 - ▶ Oblicz $P(z_{dj} = k | \mathbf{w}, \mathbf{z}_{-dj}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ dla $k = 1, \dots, K$.
 - ▶ Wylosuj z_{dj} z powyższego rozkładu.

$$P(z_{dn} = k | \mathbf{z}_{-(dn)}, \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \propto \frac{n_{k, w_{dn}}^{(\cdot)} + \beta}{n_{k, (\cdot)}^{(\cdot)} + \beta V} \cdot \frac{n_{k, (\cdot)}^d + \alpha_k}{N_d + \sum_{k=1}^K \alpha_k},$$

$n_{k,w}^d$ - liczność przypisań słowa w do tematu k w dokumencie d .

$$\hat{\theta}_{di} = \frac{n_{i, (\cdot)}^d + \alpha_i}{N_d + \sum_{i=1}^K \alpha_i}, \quad i = 1, \dots, K, \quad d = 1, \dots, M,$$

$$\hat{\varphi}_{ij} = \frac{n_{ij}^{(\cdot)} + \beta}{n_{i, (\cdot)}^{(\cdot)} + \beta V}, \quad i = 1, \dots, K, \quad j = 1, \dots, V.$$

Próbkowanie Gibbsa:

- 1 Inicjalizacja tematami najbardziej prawdopodobnymi dla odpowiednich słów.
- 2 Analogiczne losowanie sekwencyjne: losujemy na podstawie

$$P(z_{dn} = k | \mathbf{z}_{-(dn)}, \mathbf{w}, \boldsymbol{\alpha}, \beta) \propto \frac{n_{k,(\cdot)}^d + \alpha_k}{N_d + \sum_{k=1}^K \alpha_k}.$$

- $\frac{n_{k,w_{dn}}^{(\cdot)} + \beta}{n_{k,(\cdot)}^{(\cdot)} + \beta V}$ jest ustalone.