$Topics$: $\beta_t \sim Dirichlet(\alpha_1, \dots, \alpha_{20000})$, $which\ is\ a\ prob\ over\ words\ for\ topics$

eg: $[1/4 \quad 1/4 \quad 1/2]$

$\quad$ W$_1$ $\qquad$ W$_2$ $\qquad$ W$_3$ ... W$_{50}$

Documents:

1. $\theta_d \sim Dirichlet(\gamma_1, \dots, \gamma_{50})$, which is distribution over topics

eg: $[1/4 \quad 1/2 \quad 1/4]$

$\quad$ t$_1$ $\qquad$ t$_2$ $\qquad$ t$_3$ ... t$_{50}$

2. $Each\ word\ is\ assigned\ to\ topic\ Z_{id}$ (topic for word i in doc d)

$Z_{id} \sim Multi(\theta_d)$, $which\ means\ prob(word\ i = topic\ 1) = \theta_{d,1}$

3. Each word is drawn from topic prob vertor $\omega_{id}$

$\omega_{id} \sim Multi(\beta_{Z_{id}})$

Then we can find:

$$\beta_1, \dots, \beta_{50}$$
$$\theta_1, \dots, \theta_N$$
$$Z_{i,1}, \dots, Z_{i,n_d} \ (i = 1, \dots N)$$

**3.11 ColumbiaX:** DS101X Statistical Thinking for Data Science and Analytics

**word distribution for each topic:**

$$\begin{array}{cccc} \text{city} & \text{town} & \text{Seattle} & \text{zoology} \end{array}$$

$$\text{topic } 1 \quad \beta_1 = [0.1 \;\; 0.08 \ldots 0.02 \ldots 10^{-12}]$$

$$\begin{array}{ccc} \text{cow} & \text{sheep} & \text{zoology} \end{array}$$

$$\text{topic } 50 \quad \beta_{50} = [0.01 \;\; 0.009 \ldots \ldots \ldots 10^{-13}]$$

**topic distribution for each document:**

$$\begin{array}{ccc} \text{topic } 27 & \text{topic } 43 & \text{topic } 2 \end{array}$$

$$doc\ 1: \;\; \theta_1 = [0.35 \;\; 0.21 \ldots \ldots 0.01]$$

$$\begin{array}{cc} \text{topic } 5 & \text{topic } 6 \end{array}$$

$$doc\ 3721: \;\; \theta_{3721} = [0.21 \;\; \ldots \ldots \ldots \ldots . 0.01]$$