

Categorical Distributions in Natural Language Processing

Version 0.1

MURAWAKI Yugo

12 May 2016

Categorical distribution

Suppose random variable x takes one of K values. x is generated according to categorical distribution $\text{Cat}(\theta)$, where

$$\theta = (0.1, 0.6, 0.3).$$

R Y G

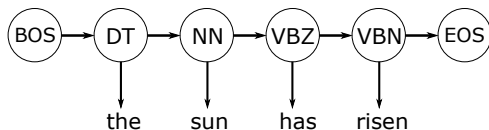
In many task settings, we do not know the true θ and need to infer it from observed data $\mathbf{x} = (x_1, \dots, x_N)$. Once we infer θ , we often want to predict new variable x' .

NOTE: In Bayesian settings, θ is usually integrated out and x' is predicted directly from \mathbf{x} .

Categorical distributions are a building block of natural language models

- N-gram language model (predicting the next word)
- POS tagging based on a Hidden Markov Model (HMM)
- Probabilistic context-free grammar (PCFG)
- Topic model (Latent Dirichlet Allocation (LDA))

Example: HMM-based POS tagging



Let K be the number of POS tags and V be the vocabulary size (ignore BOS and EOS for simplicity).

The transition probabilities can be computed using K categorical distributions ($\theta_{DT}^{TRANS}, \theta_{NN}^{TRANS}, \dots$), with the dimension K .

$$\theta_{DT}^{TRANS} = (0.21, 0.27, 0.09, \dots)$$

NN NNS ADJ

Similarly, the emission probabilities can be computed using K categorical distributions ($\theta_{DT}^{EMIT}, \theta_{NN}^{EMIT}, \dots$), with the dimension V .

$$\theta_{NN}^{EMIT} = (0.012, 0.002, 0.005, \dots)$$

sun rose cat

- Categorical and multinomial distributions
- Conjugacy and posterior predictive distribution
- LDA (Latent Dirichlet Allocation) as an application
- Gibbs sampling for inference

Suppose θ is known.

The probability of generating random variable $x \in \{1, \dots, K\}$ is

$$p(x = k|\theta) = \theta_k.$$

$0 \leq \theta_k \leq 1$ and $\sum_{k=1}^K \theta_k = 1$.

For example, if $\theta = (0.1, 0.6, 0.3)$, then $p(x = 2|\theta) = \theta_2 = 0.6$.

The probability of generating a sequence of random variables with length N , $\mathbf{x} = (x_1, \dots, x_N)$, is

$$p(\mathbf{x}|\boldsymbol{\theta}) = \prod_{i=1}^N \theta_{x_i} = \prod_{k=1}^K \theta_k^{n_k},$$

where n_k is the number of times value k is observed in \mathbf{x} ($\sum_{k=1}^K n_k = N$).

NOTE: $p(\mathbf{x}|\boldsymbol{\theta})$ does not depend on the ordering of \mathbf{x} but only on the total number of observed values in the sequence (sufficient statistics).

Multinomial distribution: the probability of counts

Replacing $\mathbf{x} = (x_1, \dots, x_N)$ with the counts $\mathbf{n} = (n_1, \dots, n_K)$, we obtain the multinomial distribution:

$$\begin{aligned}\text{Multi}(\mathbf{n}|\boldsymbol{\theta}) &= \binom{N}{n_1 \dots n_K} \prod_{k=1}^K \theta_k^{n_k} \\ &= \frac{N!}{n_1! \dots n_K!} \prod_{k=1}^K \theta_k^{n_k} \\ &= \frac{\Gamma(N+1)}{\prod_{k=1}^K \Gamma(n_k+1)} \prod_{k=1}^K \theta_k^{n_k}.\end{aligned}$$

Note $\Gamma(n) = (n-1)!$. The second term is the same as $p(\mathbf{x}|\boldsymbol{\theta})$. The first term is the combinatorial number of \mathbf{x} mapped to \mathbf{n} .

NOTE: In NLP, the categorical distribution is often called a multinomial distribution.

Estimating θ from N observations

If θ is unknown, we may want to infer it from x .

First, the likelihood function is defined as follows:

$$L(\theta; \mathbf{x}) = \text{Cat}(\mathbf{x}|\theta).$$

In maximum likelihood estimation (MLE), we estimate θ^{ML} that satisfies:

$$\theta^{\text{ML}} = \underset{\theta}{\operatorname{argmax}} L(\theta; \mathbf{x}) = \underset{\theta}{\operatorname{argmax}} \text{Cat}(\mathbf{x}|\theta).$$

It has the following analytical solution (use the method of Lagrange multipliers):

$$\theta_k^{\text{ML}} = \frac{n_k}{N}.$$

If $n_k = 0$, then $\theta_k^{\text{ML}} = 0$ (zero frequency problem).

Would larger observed data fix the problem? Not necessarily. Natural language symbols usually follow a power law (cf. Zipf's law). There are always low-frequency symbols.

Bayes' theorem

Now we introduce $p(\theta)$, a prior distribution over θ . A prior distribution is an (arbitrary) distribution that expresses our beliefs about θ .

Bayes' theorem:

$$\begin{aligned} p(\theta|\mathbf{x}) &= \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} \\ &\propto p(\mathbf{x}|\theta)p(\theta). \end{aligned}$$

$p(\theta|\mathbf{x})$ is a posterior distribution, $p(\mathbf{x}|\theta)$ is a likelihood, and $p(\theta)$ is a prior distribution.

$p(\theta|\mathbf{x})$ can be interpreted as the distribution of θ *after* observing data.

A prior for categorical distributions

For a categorical distribution, we usually set the Dirichlet distribution $\text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha})$ as its prior because it has nice analytical properties.

$$p(\boldsymbol{\theta}|\mathbf{x}, \boldsymbol{\alpha}) \propto \text{Cat}(\mathbf{x}|\boldsymbol{\theta})\text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha}),$$

where

$$\text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha}) = \frac{\Gamma(A)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k-1},$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$, $\alpha_k > 0$, and $A = \sum_{k=1}^K \alpha_k$. $\boldsymbol{\alpha}$ is a parameter given a priori (hyperparameter). We usually set $\alpha_i = \alpha_j$. α_k can be a real number (e.g. 1.5).

Note that the Dirichlet distribution resembles a categorical distribution:

$$\text{Cat}(\mathbf{x}|\boldsymbol{\theta}) = \prod_{k=1}^K \theta_k^{n_k}.$$

Posterior distribution $p(\boldsymbol{\theta}|\mathbf{x}, \boldsymbol{\alpha}) \propto \text{Cat}(\mathbf{x}|\boldsymbol{\theta})\text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha})$ is (proof omitted):

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{x}, \boldsymbol{\alpha}) &= \frac{\Gamma(N + A)}{\Gamma(n_1 + \alpha_1) \cdots \Gamma(n_K + \alpha_K)} \prod_{k=1}^K \theta_k^{n_k + \alpha_k - 1} \\ &= \text{Dir}(\boldsymbol{\theta}|\mathbf{n} + \boldsymbol{\alpha}). \end{aligned}$$

The posterior has the same form as the prior. The parameter $\boldsymbol{\alpha}$ is replaced with $\mathbf{n} + \boldsymbol{\alpha}$. This property is called conjugacy.

Maximum a posteriori (MAP) estimation

Maximum a posteriori (MAP) estimation employs θ^{MAP} that maximizes the posterior probability:

$$\begin{aligned}\theta^{\text{MAP}} &= \operatorname{argmax}_{\theta} p(\theta|\mathbf{x}) \\ &= \operatorname{argmax}_{\theta} p(\mathbf{x}|\theta)p(\theta).\end{aligned}$$

For Dirichlet-categorical parameter $\theta^{\text{MAP}} = \operatorname{argmax}_{\theta} \operatorname{Dir}(\theta|\mathbf{n} + \alpha)$, the solution is

$$\theta_k^{\text{MAP}} = \frac{n_k + \alpha_k - 1}{\sum_{k=1}^K (n_k + \alpha_k - 1)} = \frac{n_k + \alpha_k - 1}{N + A - K}.$$

Compare it with ML estimate $\theta_k^{\text{ML}} = \frac{n_k}{N}$. The additional term $\alpha_k - 1$ is used to smooth the estimate (additive smoothing).

Instead of using one estimate of θ , we now consider all possible values of θ . To do so, we integrate out θ to obtain marginal likelihood:

$$p(\mathbf{x}) = \int_{\theta} p(\mathbf{x}|\theta)p(\theta) d\theta.$$

The probability of generating new variable x' given \mathbf{x} (predictive probability) is:

$$\begin{aligned} p(x'|\mathbf{x}) &= \int_{\theta} p(x'|\theta)p(\theta|\mathbf{x}) d\theta \\ &= \int_{\theta} p(x'|\theta) \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} d\theta \end{aligned}$$

Thanks to conjugacy, we have analytical solutions for Dirichlet-categorical (next slide).

$$\begin{aligned}
 p(\mathbf{x}|\boldsymbol{\alpha}) &= \int_{\boldsymbol{\theta}} \text{Cat}(\mathbf{x}|\boldsymbol{\theta})\text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta} \\
 &= \int_{\boldsymbol{\theta}} \left(\prod_{k=1}^K \theta_k^{n_k} \right) \left(\frac{\Gamma(A)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \right) d\boldsymbol{\theta} \\
 &= \frac{\Gamma(A)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \int_{\boldsymbol{\theta}} \prod_{k=1}^K \theta_k^{n_k + \alpha_k - 1} d\boldsymbol{\theta} \\
 &= \frac{\Gamma(A)}{\prod_{k=1}^K \Gamma(\alpha_k)} \frac{\prod_{k=1}^K \Gamma(n_k + \alpha_k)}{\Gamma(N + A)} \\
 &= \frac{\Gamma(A)}{\Gamma(N + A)} \prod_{k=1}^K \frac{\Gamma(n_k + \alpha_k)}{\Gamma(\alpha_k)}
 \end{aligned}$$

Posterior predictive for Dirichlet-categorical

The probability of generating x' after observing \mathbf{x} is

$$\begin{aligned} p(x' = k' | \mathbf{x}, \boldsymbol{\alpha}) &= \frac{p(x' = k', \mathbf{x} | \boldsymbol{\alpha})}{p(\mathbf{x} | \boldsymbol{\alpha})} \\ &= \frac{\frac{\Gamma(A)}{\Gamma(N+1+A)} \prod_{k=1}^K \frac{\Gamma(n_k + I(k=k') + \alpha_k)}{\Gamma(\alpha_k)}}{\frac{\Gamma(A)}{\Gamma(N+A)} \prod_{k=1}^K \frac{\Gamma(n_k + \alpha_k)}{\Gamma(\alpha_k)}} \\ &= \frac{\Gamma(N+A)}{\Gamma(N+1+A)} \frac{\Gamma(n_{k'} + 1 + \alpha_{k'})}{\Gamma(n_{k'} + \alpha_{k'})} \\ &= \frac{n_{k'} + \alpha_{k'}}{N + A}, \end{aligned}$$

where $I(\text{statement})$ is an indicator function. It gives 1 if statement is true; 0 otherwise.

In analogy with n_k (the count of observations), α_k is called a pseudo-count.

Sequential updates of x

Let $\mathbf{x}_{1:N} = (x_1, \dots, x_N)$. $p(\mathbf{x}_{1:N}|\boldsymbol{\alpha})$ is the product of the probabilities of generating x one by one.

$$\begin{aligned} p(\mathbf{x}_{1:N}|\boldsymbol{\alpha}) &= p(x_N|\mathbf{x}_{1:N-1}, \boldsymbol{\alpha})p(\mathbf{x}_{1:N-1}|\boldsymbol{\alpha}) \\ &= \frac{n_{x_N}(\mathbf{x}_{1:N-1}) + \alpha_{x_N}}{N - 1 + A} p(\mathbf{x}_{1:N-1}|\boldsymbol{\alpha}) \\ &= \frac{n_{x_N}(\mathbf{x}_{1:N-1}) + \alpha_{x_N}}{N - 1 + A} p(x_{N-1}|\mathbf{x}_{1:N-2}, \boldsymbol{\alpha})p(\mathbf{x}_{1:N-2}|\boldsymbol{\alpha}) \\ &= \frac{n_{x_N}(\mathbf{x}_{1:N-1}) + \alpha_{x_N}}{N - 1 + A} \frac{n_{x_{N-1}}(\mathbf{x}_{1:N-2}) + \alpha_{x_{N-1}}}{N - 2 + A} p(\mathbf{x}_{1:N-2}|\boldsymbol{\alpha}) \\ &= \frac{\prod_{k=1}^K \prod_{i=1}^{n_k(\mathbf{x}_{1:N})} i - 1 + \alpha_k}{\prod_{i=1}^N i - 1 + A} \\ &= \frac{\Gamma(A)}{\Gamma(N + A)} \prod_{k=1}^K \frac{\Gamma(n_k(\mathbf{x}_{1:N}) + \alpha_k)}{\Gamma(\alpha_k)}, \end{aligned}$$

where $n_k(\mathbf{x})$ is the number of times k appears in \mathbf{x} .

As stated earlier, $p(\mathbf{x}|\theta)$ does not depend on the ordering of \mathbf{x} but on the total number of times each value is observed in the sequence (sufficient statistics). Let us create \mathbf{x}' by swapping an arbitrary pair of variables x_i, x_j in \mathbf{x} . Since the ordering does not matter, the following is hold:

$$p(\mathbf{x}'|\alpha) = p(\mathbf{x}|\alpha)$$

This property is called exchangeability.

As we will see later, exchangeability makes inference easy.

Example: Sequential updates of x

Suppose random variable x takes one of three values: R, Y, G. The probability of generating the sequence $\mathbf{x} = (\text{R}, \text{R}, \text{Y})$ is

$$p(\mathbf{x}|\boldsymbol{\alpha}) = \frac{0 + \alpha_{\text{R}}}{0 + A} \frac{1 + \alpha_{\text{R}}}{1 + A} \frac{0 + \alpha_{\text{Y}}}{2 + A}$$

Verify that reordering does not affect the probability.

For Dirichlet-categorical, θ is first generated from a Dirichlet distribution with parameter α and then each random variable x_n is generated from a categorical distribution with parameter θ . This process can be summarized as follows:

$$\theta|\alpha \sim \text{Dir}(\alpha)$$

$$x_n|\theta \sim \text{Cat}(\theta)$$

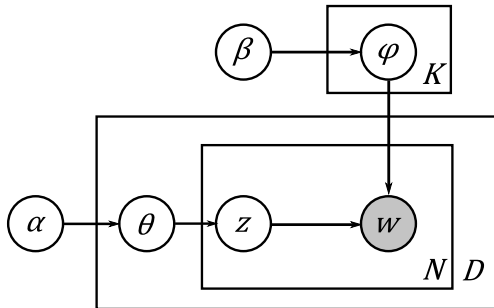
Application: LDA (Latent Dirichlet allocation)

NOTE: We use the modified model by Griffiths et al. (2004), not the original one by Blei et al. (2002).

We are given D documents. Document i contains N_i words, and $w_{i,j}$ is the j -th word in document i . LDA assumes that each document is a mixture of K topics. Document i is associated with a mixing proportion θ_i . Each topic k has a corresponding vocabulary distribution $\text{Cat}(\varphi_k)$. Each word $w_{i,j}$ is tied with a latent variable $z_{i,j} = k'$. $z_{i,j}$ is generated from $\text{Cat}(\theta_i)$, and then $w_{i,j}$ is generated from $\text{Cat}(\varphi_{k'})$.

Application: LDA (Latent Dirichlet allocation)

The generative story can be represented by a directed graph called a graphical model.



Each arrow indicates a dependency between variables. Each plate represents repetition with the number of repetition is shown in the corner. The shaded node represents an observed variable.

The generative story is as follows:

- 1 For each topic $k \in \{1, \dots, K\}$, draw vocabulary distribution $\varphi_k \sim \text{Dir}(\beta)$
- 2 For each document $i \in \{1, \dots, D\}$, draw topic distribution $\theta_i \sim \text{Dir}(\alpha)$
- 3 For each word $j \in \{1, \dots, N_i\}$ in document i ,
 - 1 Draw topic assignment $z_{i,j} \sim \text{Cat}(\theta_i)$, and
 - 2 Draw word $w_{i,j} \sim \text{Cat}(\varphi_{z_{i,j}})$

Let \mathbf{W} be the set of all observed words in the collection of documents and \mathbf{Z} be the corresponding latent topic assignments. Given \mathbf{W} , we want to infer \mathbf{Z} .

Application: LDA (Latent Dirichlet allocation)

We begin with the joint probability:

$$p(\mathbf{W}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\varphi} | \boldsymbol{\alpha}, \boldsymbol{\beta}) = \prod_{k=1}^K p(\boldsymbol{\varphi}_k | \boldsymbol{\beta}) \times \prod_{i=1}^D \left(p(\boldsymbol{\theta}_i | \boldsymbol{\alpha}) \prod_{j=1}^{N_i} p(z_{i,j} | \boldsymbol{\theta}_i) p(w_{i,j} | \boldsymbol{\varphi}, z_{i,j}) \right).$$

$\boldsymbol{\theta}_i$ and $\boldsymbol{\varphi}_k$ can be integrated out to obtain marginal likelihood.

$$\begin{aligned} p(\mathbf{W}, \mathbf{Z} | \boldsymbol{\alpha}, \boldsymbol{\beta}) &= \int_{\boldsymbol{\theta}} \int_{\boldsymbol{\varphi}} p(\mathbf{W}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\varphi} | \boldsymbol{\alpha}, \boldsymbol{\beta}) d\boldsymbol{\theta} d\boldsymbol{\varphi} \\ &= \prod_{i=1}^D \left(\frac{\Gamma(\alpha_{(\cdot)})}{\Gamma(n_{i,(\cdot)}^{(\cdot)} + \alpha_{(\cdot)})} \prod_{k=1}^K \frac{\Gamma(n_{i,(\cdot)}^k + \alpha_k)}{\Gamma(\alpha_k)} \right) \\ &\quad \times \prod_{k=1}^K \left(\frac{\Gamma(\beta_{(\cdot)})}{\Gamma(n_{(\cdot),(\cdot)}^k + \beta_{(\cdot)})} \prod_v \frac{\Gamma(n_{(\cdot),v}^k + \beta_v)}{\Gamma(\beta_v)} \right), \end{aligned}$$

where $n_{i,v}^k$ is the number of times word type v in document i is tied with topic k . We can see that LDA is simply a composite of $D + K$ Dirichlet-categorical distributions.

Application: LDA (Latent Dirichlet allocation)

Now consider sequential updates of topic assignments and words. It is easy to verify that exchangeability is hold.

Doc ID	Topic	Word	Prob.
1	Politics	Obama	$\frac{0+\alpha_{\text{Politics}}}{0+\alpha_{(\cdot)}} \times \frac{0+\beta_{\text{Obama}}}{0+\beta_{(\cdot)}}$
	Economy	FRB	$\frac{0+\alpha_{\text{Economy}}}{1+\alpha_{(\cdot)}} \times \frac{0+\beta_{\text{FRB}}}{0+\beta_{(\cdot)}}$
	Politics	election	$\frac{1+\alpha_{\text{Politics}}}{2+\alpha_{(\cdot)}} \times \frac{0+\beta_{\text{election}}}{1+\beta_{(\cdot)}}$
	Politics	Obama	$\frac{2+\alpha_{\text{Politics}}}{3+\alpha_{(\cdot)}} \times \frac{1+\beta_{\text{Obama}}}{2+\beta_{(\cdot)}}$
2	Economy	Yellen	$\frac{0+\alpha_{\text{Economy}}}{0+\alpha_{(\cdot)}} \times \frac{0+\beta_{\text{Yellen}}}{1+\beta_{(\cdot)}}$
	Economy	FRB	$\frac{1+\alpha_{\text{Economy}}}{1+\alpha_{(\cdot)}} \times \frac{1+\beta_{\text{FRB}}}{2+\beta_{(\cdot)}}$

We want to infer Z that satisfies

$$\operatorname{argmax}_Z p(Z|W, \alpha, \beta) = \operatorname{argmax}_Z p(W, Z|\alpha, \beta)$$

The number of possible assignments is $K^{|Z|}$. The search space is too large for exhaustive search.

We resort to approximate inference using random walk. Specifically, we often use Gibbs sampling as a theoretically-grounded algorithm.

Sampling is the procedure of obtaining a sample from a probability distribution. It is typically implemented by modifying the values of the rand function r ($0 \leq r < 1$).

Suppose x is a sample from $\theta = (0.1, 0.6, 0.3)$. Using r , we can create a sample x as follows:

$$x = \begin{cases} 1 & 0 \leq r < 0.1, \\ 2 & 0.1 \leq r < 0.7, \text{ and} \\ 3 & 0.7 \leq r < 1. \end{cases}$$

Recall the predictive distribution of Dirichlet-categorical is

$$p(x' = k' | \mathbf{x}, \boldsymbol{\alpha}) = \frac{n_{k'}(\mathbf{x}) + \alpha_{k'}}{n(\mathbf{x}) + \alpha_{(\cdot)}}.$$

Let $\mathbf{x} = (\text{R}, \text{R}, \text{Y}, \text{R}, \text{R})$, and $\boldsymbol{\alpha} = (\alpha_{\text{R}}, \alpha_{\text{Y}}, \alpha_{\text{G}}) = (1, 1, 1)$. Then

$$\begin{aligned} p(x' = \text{R} | \mathbf{x}, \boldsymbol{\alpha}) &= \frac{4 + 1}{5 + 3} = 0.625 \\ p(x' = \text{Y} | \mathbf{x}, \boldsymbol{\alpha}) &= \frac{1 + 1}{5 + 3} = 0.25 \\ p(x' = \text{G} | \mathbf{x}, \boldsymbol{\alpha}) &= \frac{0 + 1}{5 + 3} = 0.125. \end{aligned}$$

Draw a sample from $\boldsymbol{\theta} = (0.625, 0.25, 0.125)$, and that is the sample we want to obtain.

LDA is a composite of Dirichlet-categorical distributions $p(\mathbf{Z}|\mathbf{W}, \alpha, \beta)$. Directly sampling from this complex distribution is difficult.

However, it is easy to draw a sample from $p(z_i|\mathbf{W}, \mathbf{Z}^{-i}, \alpha, \beta)$, where $\mathbf{Z}^{-i} = \mathbf{Z} \setminus \{z_i\}$. i is now a corpus-wide index.

Gibbs sampling is a method of obtaining samples of \mathbf{Z} by repeatedly sampling from $p(z_i|\mathbf{W}, \mathbf{Z}^{-i}, \alpha, \beta)$.

The algorithmic overview of Gibbs sampling is as follows:

- 1 initialize \mathbf{Z}
- 2 For each iteration $\tau = 1, \dots, T$,
 - randomly select i
 - sample z'_i from $p(z_i|\mathbf{W}, \mathbf{Z}^{-i}, \alpha, \beta)$
 - update \mathbf{Z} by replacing z_i with z'_i

After a sufficient number of iterations, \mathbf{Z} can be seen as samples from $p(\mathbf{Z}|\mathbf{W}, \alpha, \beta)$.

Gibbs sampling is applicable if we can easily draw a sample from $p(z_i | \mathbf{Z}^{-i})$.

Due to exchangeability, the probability of observing the sequence $(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_N)$ is the same as the probability of observing the sequence $(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_N, z_i)$.

This means that $p(z_i | \mathbf{Z}^{-i})$ is the posterior predictive distribution of z_i after observing $(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_N)$, the sequence excluding z_i .

$$\begin{aligned}
 & p(z_i = k' | \mathbf{Z}^{-i}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \\
 \propto & p(z_i = k', \mathbf{Z}^{-i}, \mathbf{W} | \boldsymbol{\alpha}, \boldsymbol{\beta}) \\
 = & \prod_{j=1}^D \left(\frac{\Gamma(\alpha_{(\cdot)})}{\Gamma(n_{j,(\cdot)}^{\cdot}(\mathbf{Z}^{-i}) + I(w_i \in \mathbf{W}_j) + \alpha_{(\cdot)})} \right. \\
 & \left. \prod_{k=1}^K \frac{\Gamma(n_{j,(\cdot)}^k(\mathbf{Z}^{-i}) + I(k' = k \ \& \ w_i \in \mathbf{W}_j) + \alpha_k)}{\Gamma(\alpha_k)} \right) \\
 \times & \prod_{k=1}^K \left(\frac{\Gamma(\beta_{(\cdot)})}{\Gamma(n_{(\cdot),(\cdot)}^k(\mathbf{Z}^{-i}) + I(k' = k) + \beta_{(\cdot)})} \right. \\
 & \left. \prod_v \frac{\Gamma(n_{(\cdot),v}^k(\mathbf{Z}^{-i}) + I(k' = k \ \& \ w_i = v) + \beta_v)}{\Gamma(\beta_v)} \right).
 \end{aligned}$$

\mathbf{W}_j is the sequence of words in the j -th document. Most terms are irrelevant with z_i and can be dropped.

$$\begin{aligned}
 & p(z_i = k' | \mathbf{Z}^{-i}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \\
 \propto & \frac{\Gamma(n_{j',(\cdot)}^{(\cdot)}(\mathbf{Z}^{-i}) + \alpha_{(\cdot)})}{\Gamma(n_{j',(\cdot)}^{(\cdot)}(\mathbf{Z}^{-i}) + 1 + \alpha_{(\cdot)})} \times \frac{\Gamma(n_{j',(\cdot)}^{k'}(\mathbf{Z}^{-i}) + 1 + \alpha_{k'})}{\Gamma(n_{j',(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \alpha_{k'})} \\
 & \times \frac{\Gamma(n_{(\cdot),(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \beta_{(\cdot)})}{\Gamma(n_{(\cdot),(\cdot)}^{k'}(\mathbf{Z}^{-i}) + 1 + \beta_{(\cdot)})} \times \frac{\Gamma(n_{(\cdot),v'}^{k'}(\mathbf{Z}^{-i}) + 1 + \beta_{v'})}{\Gamma(n_{(\cdot),v'}^{k'}(\mathbf{Z}^{-i}) + \beta_{v'})} \\
 = & \frac{n_{j',(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \alpha_{k'}}{n_{j',(\cdot)}^{(\cdot)}(\mathbf{Z}^{-i}) + \alpha_{(\cdot)}} \times \frac{n_{(\cdot),v'}^{k'}(\mathbf{Z}^{-i}) + \beta_{v'}}{n_{(\cdot),(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \beta_{(\cdot)}},
 \end{aligned}$$

where $v' = w_i$, and j' is the index of the document w_i belongs to. The first term is the posterior predictive probability of adding topic assignment with value k' to document j' (the denominator can be dropped because it is constant wrt z_i). The second term is the posterior predictive probability of drawing word type v' from the vocabulary distribution of topic k' .

Intuitive explanation of Gibbs sampling for LDA

$$p(z_i = k' | \mathbf{Z}^{-i}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\beta}) \propto \frac{n_{j',(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \alpha_{k'}}{n_{j',(\cdot)}^{(\cdot)}(\mathbf{Z}^{-i}) + \alpha_{(\cdot)}} \times \frac{n_{(\cdot),v'}^{k'}(\mathbf{Z}^{-i}) + \beta_{v'}}{n_{(\cdot),(\cdot)}^{k'}(\mathbf{Z}^{-i}) + \beta_{(\cdot)}}$$

- To make the first term larger, we need to select k' that appears often in document j' .
- To make the second term larger, we need to select k' from which $w_i = v'$ is often drawn.
- z_i reflects a balance between the first and second terms.
- After repeating the procedure, \mathbf{Z} will generally satisfy these conditions.