

Topic Models

Why Topic Model

- Enormous amount of data is being generated every minute.
- Dimension reduction is necessary for effective use of the data.



User as Document

Consider that a company named M\$ has collected user data from many sources. They are stored as documents (bags of words) under user ids.



Find Relevant User

Toyota wants to promote their new compact car by showing ads to users who are interested in cars.

M\$ needs to extract users' interested topics from:

- Search queries
- Tweets
- Posts
- And any sort of text records.

Gordon Roqué @gordonroque · May 27

Today, I am #thankful for my car. I drive a red 2008 #Toyota #Scion xD. Its compact size allows it... [instagram.com/p/3NcvlgQUsk/](https://www.instagram.com/p/3NcvlgQUsk/)

8:14 PM - 27 May 2015 - Details



Topics:

- Vehicle
- Compact Car

Strong interest in car.

imgur



upload images



I watched a police **car** crash into a public transit **bus**. Officers then engaged everyone onboard in a **firefight**.



Topics:

- Gaming
- Vehicle
- Battle

Doesn't show interest in car.

Topic Model

- Users' activities are diverse and numerous.
Dimension reduction is necessary for effective use of the data.
- Topic model is a type of statistical models that find abstract "topics" from a collection of documents, so that each document can be represented by few topics.
- With the topic representation, we can do things like
 - Information Retrieval (find potential car buyer)
 - Classification (or prediction)

- Introduction
- **Dimension Reduction on Texts**
 - TF-IDF (1975)
 - LSI (1990)
 - Aspect Model (1999)
 - LDA (2003)
 - Deep Learning (2006)
- Probabilistic Modeling
 - Graphical Models
 - Multinomial and Dirichlet Distributions
- Latent Dirichlet Allocation
 - Generative Model
 - Inference
 - Applications

TF-IDF (1975)

Salton, Wong and Yong proposed the famous Term Frequency - Inverse Document Frequency formula to represent each document by a fixed-length vector.

$$\text{tfidf}(t, d, D) = \frac{\text{frequency of term } t \text{ in document } d}{\text{frequency of term } t \text{ in corpus } D}$$

- Dimension reduction is small if vocabulary V is large.
- TF-IDF reveals little about statistical structure of the document.

Latent Semantic Indexing (1990)

SVD (Singular Value Decomposition) is the general dimension reduction technique on real-valued vectors.

(Deerwester et. al, 1990) applies SVD to identify the principle components of the tf-idf vectors and use them to represent the documents.

- Deerwester et al. argue that the derived features can capture some aspects of basic linguistic notions such as synonymy.
- LSI is not a generative model.

Aspect Model (1999)

(Hofmann 1999) proposed pLSI (probabilistic LSI) model, aka. aspect model, as an alternative to LSI. pLSI models each word w as a sample from a mixture model.

$$P(w | d) = \sum_{z \in Z} P(w | z)P(z | d)$$

The mixture components are multinomial random variables z that can be viewed as "topics".

Aspect Model (1999)

The generative model is defined as follows

1. Select a document with probability $P(d)$, $d \in D$
2. Pick a latent topic z acc. to topic proportion $P(z | d)$
3. Generate a word w acc. to word proportion $P(w | z)$

$$P(w | d) = \sum_{z \in Z} P(w | z)P(z | d)$$

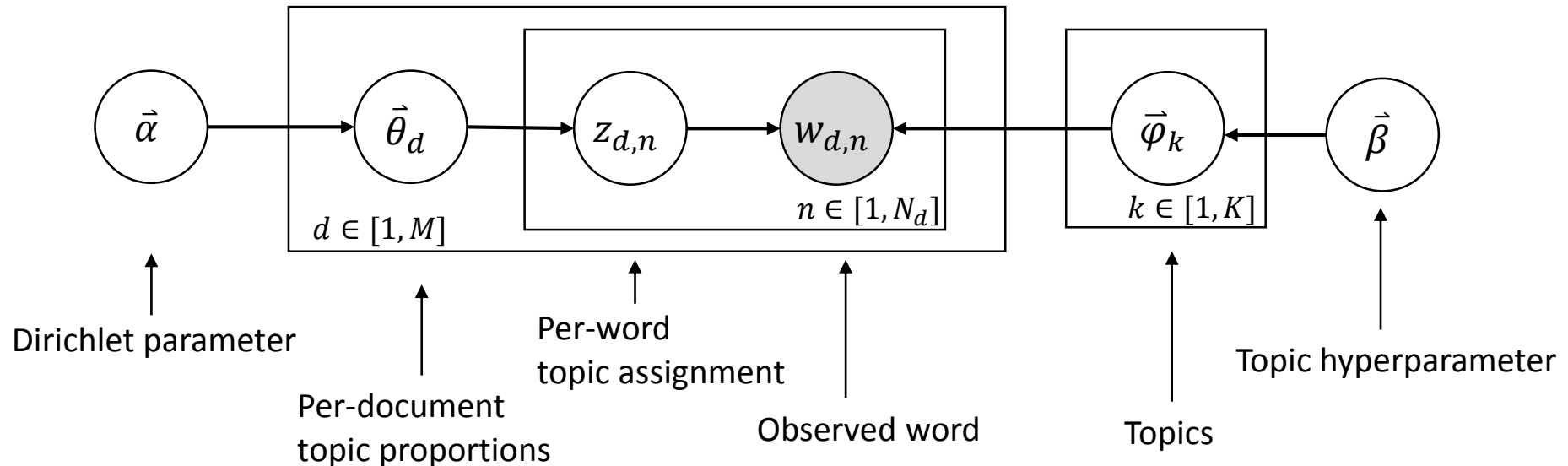
$$P(w, d) = P(w | d)P(d)$$

No probabilistic model at the level of documents.

Latent Dirichlet Allocation (2003)

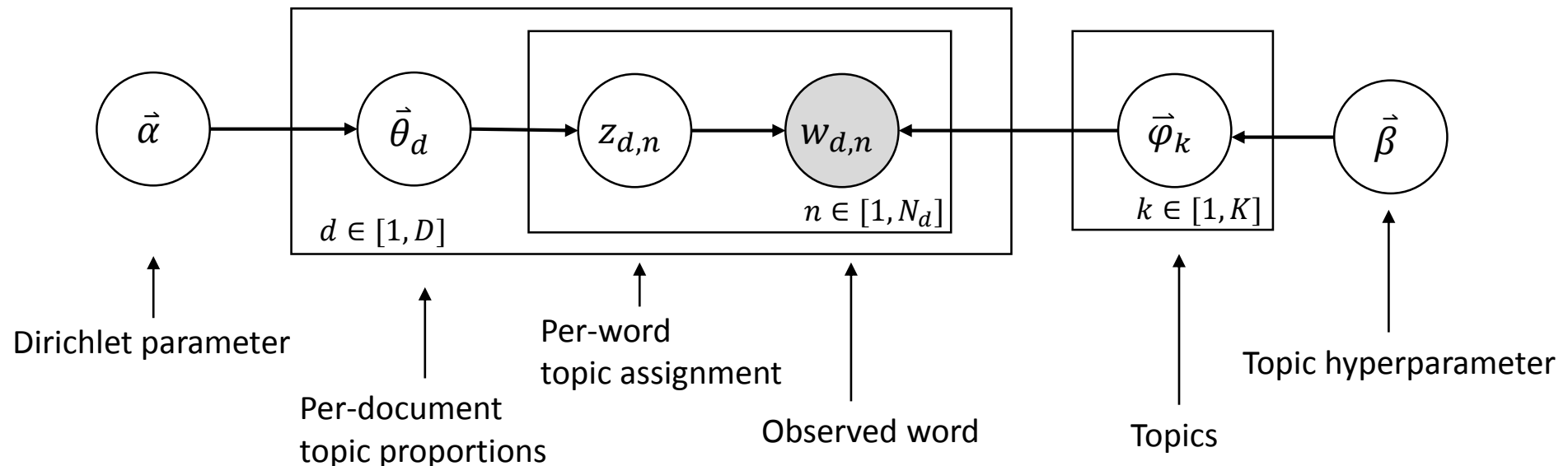
Add Dirichlet priors on aspect model, we get LDA model

$$P(\vec{w}_d, \vec{z}_d, \vec{\theta}_d, \vec{\varphi} | \vec{\alpha}, \vec{\beta}) = \underbrace{\prod_{n=1}^{N_d} \text{Mult}(w_{d,n} | \vec{\varphi}_{z_{d,n}}) \text{Mult}(z_{d,n} | \vec{\theta}_d)}_{\text{word likelihood}} \cdot \underbrace{\text{Dir}(\vec{\theta}_d | \vec{\alpha})}_{\text{document prior}} \underbrace{\text{Dir}(\vec{\varphi} | \vec{\beta})}_{\text{topic prior}}$$



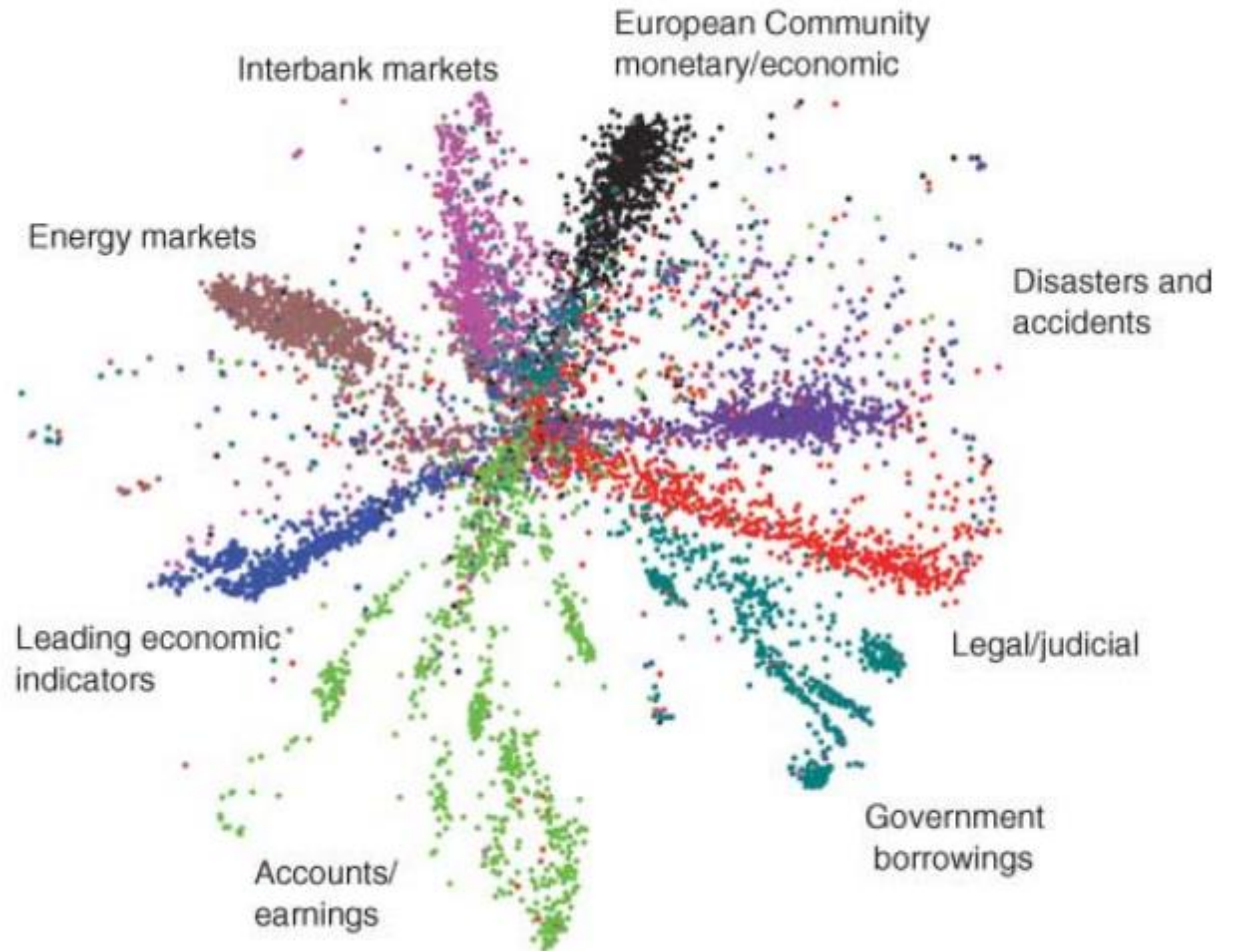
Latent Dirichlet Allocation (Blei et. al. 2003)

1. For each topic $k \in [1, K]$, sample word proportion $\vec{\phi}_k \sim \text{Dir}(\vec{\beta})$
2. For each document $d \in [1, D]$, sample topic proportion $\vec{\theta}_d \sim \text{Dir}(\vec{\alpha})$
 1. For each word $n \in [1, N_d]$, sample topic index $z_{d,n} \sim \text{Mult}(\vec{\theta}_d)$
 2. For each word $n \in [1, N_d]$, sample term $w_{d,n} \sim \text{Mult}(\vec{\phi}_{z_{d,n}})$



Deep Learning (2006)

G. E. Hinton and R. R. Salakhutdinov converts high-dimensional data to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors.



- Introduction
- Dimension Reduction on Texts
 - TF-IDF (1975)
 - LSI (1990)
 - Aspect Model (1999)
 - LDA (2003)
 - Deep Learning (2006)
- **Probabilistic Modeling**
 - Graphical Models
 - Multinomial and Dirichlet Distributions
- Latent Dirichlet Allocation
 - Generative Model
 - Inference
 - Applications

Probabilistic Modeling

1. Treat data as observations that arise from a generative probabilistic process that includes hidden variables.
 - For documents, the hidden variables reflect the thematic structure of a document.
2. Infer the hidden structure (topics) using posterior inference.
 - MCMC (Markov Chain Monte Carlo)

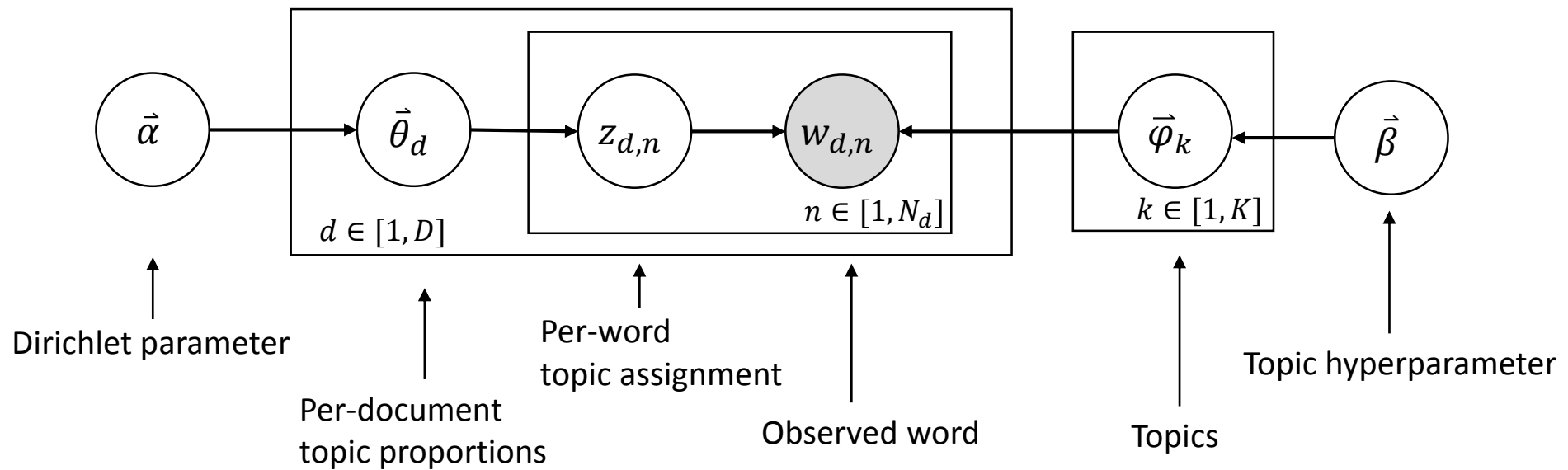
$$P(\theta | X) = \frac{P(X | \theta)P(\theta)}{P(X)}$$

$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{evidence}}$$

3. Situate new data into the estimated model.

Graphical Models

- Nodes are random variables
- Edges denote possible dependence.
- Observed variables are shaded.
- Replicate structures are boxed.



Multinomial and Dirichlet Distributions

- Multinomial distribution conditioned on word proportion \vec{p}

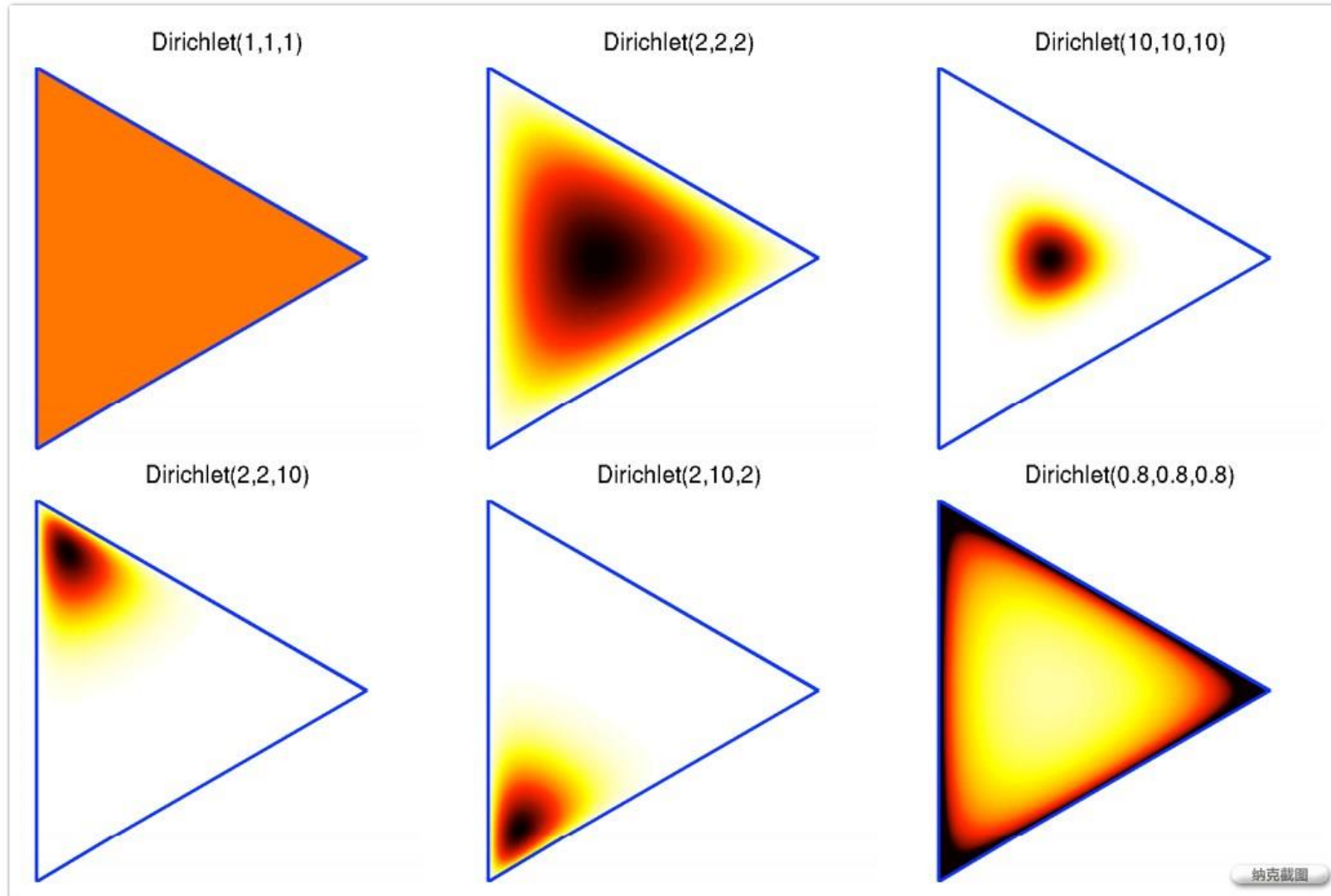
$$P(\vec{w} | \vec{p}) = \text{Mult}(\vec{n} | \vec{p}) = \frac{\Gamma(\sum_t n^{(t)} + 1)}{\prod_t \Gamma(n^{(t)} + 1)} \prod_{t=1}^V p_t^{n^{(t)}}$$

- Dirichlet distribution conditioned on pseudo-count $\vec{\alpha}$

$$P(\vec{p} | \vec{\alpha}) = \text{Dir}(\vec{p} | \vec{\alpha}) = \frac{1}{\Delta(\vec{\alpha})} \prod_{t=1}^V p_t^{\alpha_t - 1}$$

$$\Delta(\vec{\alpha}) = \frac{\prod_t \Gamma(\alpha_t)}{\Gamma(\sum_t \alpha_t)}$$

Dirichlet Pseudo-Count when $K=3$

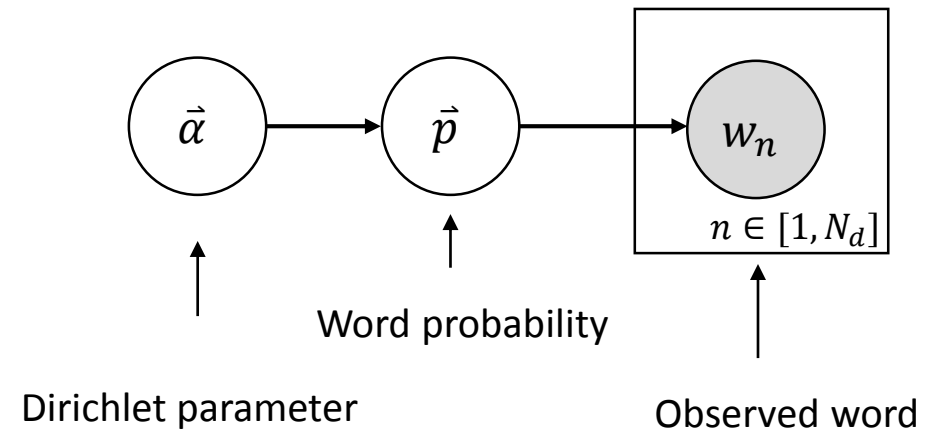


Multinomial and Dirichlet Distributions

- Thanks to conjugacy, the Dirichlet posterior given multinomial observation \vec{w} and pseudo-counts $\vec{\alpha}$ is

$$P(\vec{p} \mid \vec{w}, \vec{\alpha}) = \frac{\text{Mult}(\vec{w} \mid \vec{p}) \text{Dir}(\vec{p} \mid \alpha)}{P(\vec{w} \mid \vec{\alpha})} = \text{Dir}(\vec{p} \mid \vec{\alpha} + \vec{n})$$

- Likelihood(\vec{p}) \sim $\text{Mult}(\vec{w} \mid \vec{p})$
- Prior(\vec{p}) \sim $\text{Dir}(\vec{p} \mid \vec{\alpha})$
- Having new observation is equivalent with adjusting pseudo-count.

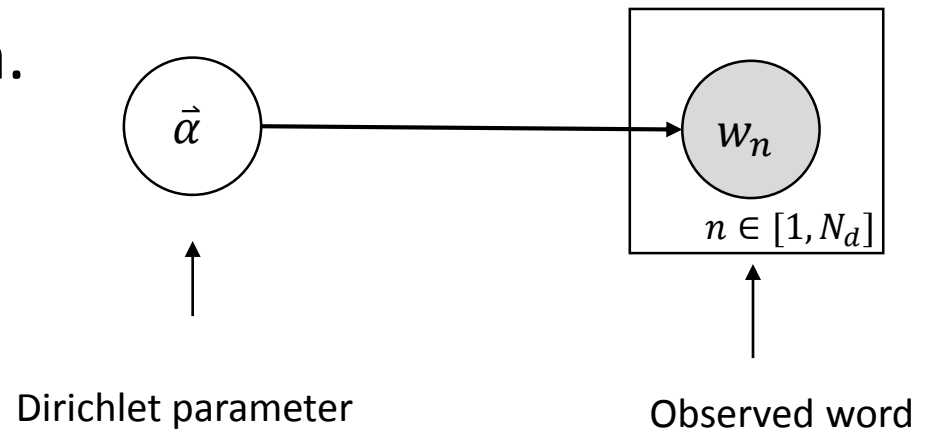


Multinomial and Dirichlet Distributions

- Integrate out p to get

$$\begin{aligned} P(\vec{w} | \vec{\alpha}) &= \int_{\vec{p}} \text{Mult}(\vec{w} | \vec{p}) \text{Dir}(\vec{p} | \alpha) d\vec{p} \\ &= \frac{\Delta(\vec{n} + \vec{\alpha})}{\Delta(\vec{\alpha})} \end{aligned}$$

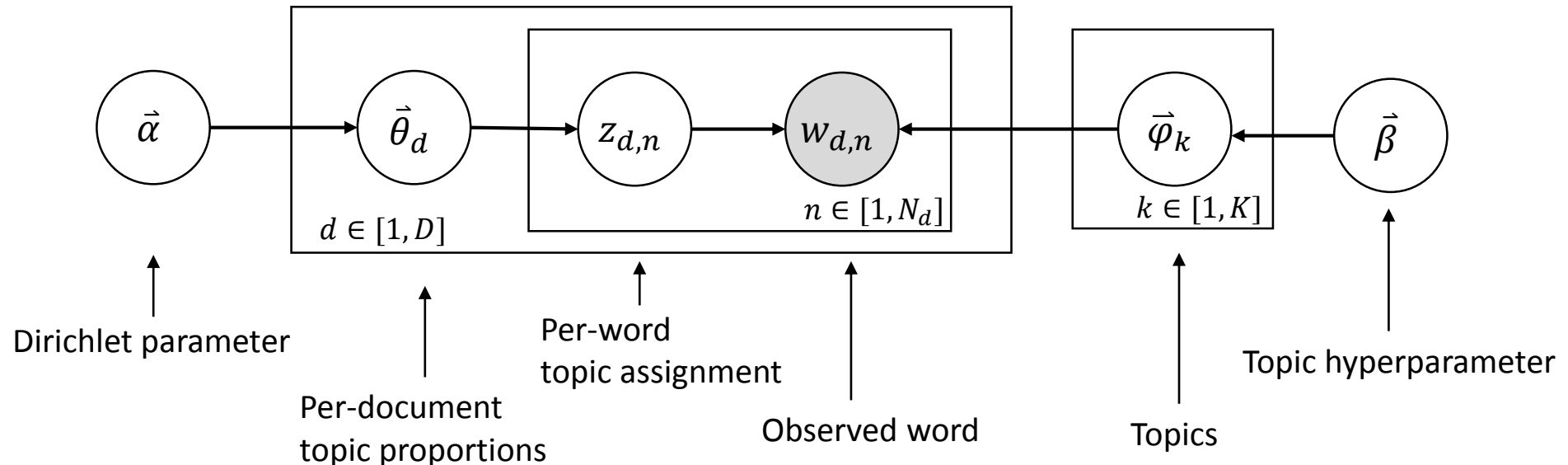
- We'll use this repeatedly in later derivation.



- Introduction
- Dimension Reduction on Texts
 - TF-IDF (1975)
 - LSI (1990)
 - Aspect Model (1999)
 - LDA (2003)
 - Deep Learning (2006)
- Probabilistic Modeling
 - Graphical Models
 - Multinomial and Dirichlet Distributions
- **Latent Dirichlet Allocation**
 - Generative Model
 - Inference
 - Applications

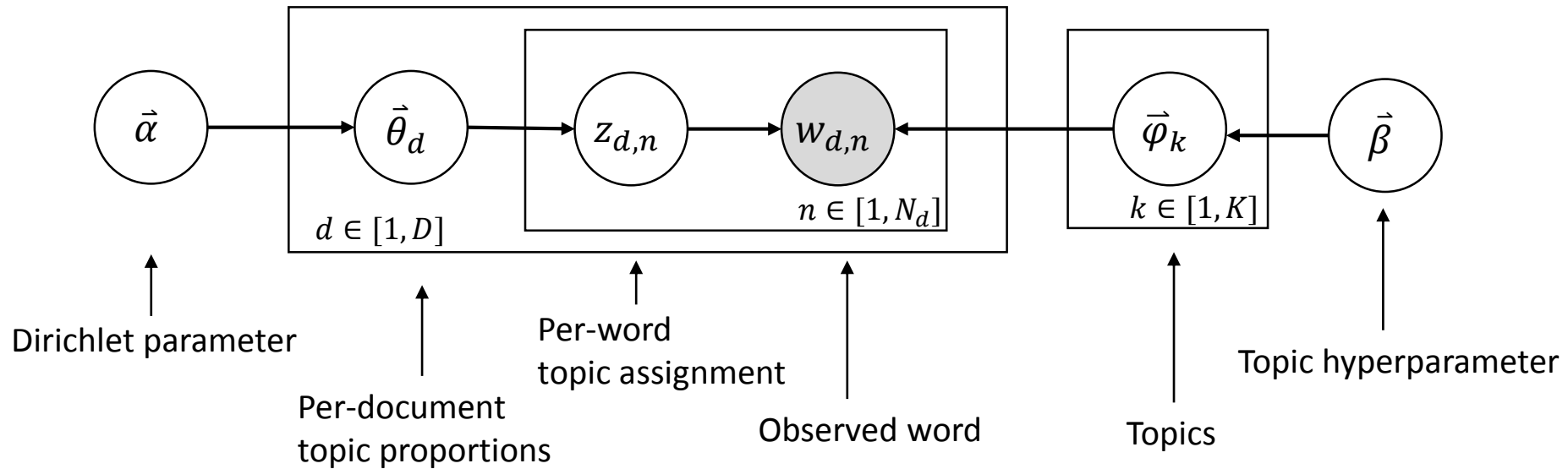
Generative Model

1. For each topic $k \in [1, K]$, sample word proportion $\vec{\phi}_k \sim \text{Dir}(\vec{\beta})$
2. For each document $d \in [1, D]$, sample topic proportion $\vec{\theta}_d \sim \text{Dir}(\vec{\alpha})$
 1. For each word $n \in [1, N_d]$, sample topic index $z_{d,n} \sim \text{Mult}(\vec{\theta}_d)$
 2. For each word $n \in [1, N_d]$, sample term $w_{d,n} \sim \text{Mult}(\vec{\phi}_{z_{d,n}})$



Generative Model

$$P(\vec{w}_d, \vec{z}_d, \vec{\theta}_d, \vec{\varphi} | \vec{\alpha}, \vec{\beta}) = \underbrace{\prod_{n=1}^{N_d} \text{Mult}(w_{d,n} | \vec{\varphi}_{z_{d,n}}) \text{Mult}(z_{d,n} | \vec{\theta}_d)}_{\text{word likelihood}} \cdot \underbrace{\text{Dir}(\vec{\theta}_d | \vec{\alpha})}_{\text{document prior}} \underbrace{\text{Dir}(\vec{\varphi} | \vec{\beta})}_{\text{topic prior}}$$



Inference

- Given observations W , we want to infer the hidden structure (topics) from posterior probability.
- Posterior probability is computationally intractable. Approximate inference is used.
 - MCMC (Markov Chain Monte Carlo): Gibbs sampler, collapse Gibbs sampler.
 - Variational methods: replace sampling with optimization.
 - “Distributed Algorithms for Topic Models” by D. Newman, A. Asuncion et. al., 2009.

Inference - Strategy

- Given document \vec{w} , sample topic assignment \vec{z} from posterior

$$P(\vec{z} | \vec{w}, \vec{\alpha}, \vec{\beta}) = \frac{P(\vec{w} | \vec{z}, \vec{\beta})P(\vec{z} | \vec{\alpha})}{P(\vec{w} | \vec{\alpha}, \vec{\beta})}$$

- Estimate topic proportion and word proportion by Dirichlet posterior

$$P(\vec{\theta}_d | \vec{z}_d, \vec{\alpha}) = \text{Dir}(\vec{\theta}_d | \vec{\alpha} + \vec{n}_d)$$
$$P(\vec{\phi}_k | \vec{z}, \vec{w}, \vec{\beta}) = \text{Dir}(\vec{\phi}_k | \vec{\beta} + \vec{n}_k)$$
$$\mathbf{E}(\theta_{d,k}) = \frac{n_d^{(k)} + \alpha_k}{\sum_{k'=1}^K n_d^{(k')} + \alpha_{k'}}$$
$$\mathbf{E}(\phi_{k,t}) = \frac{n_k^{(t)} + \beta_t}{\sum_{t'=1}^V n_k^{(t')} + \beta_{t'}}$$

Inference – Likelihood

- Given document \vec{w} , sample topic assignment \vec{z} from posterior

$$P(\vec{z} | \vec{w}, \vec{\alpha}, \vec{\beta}) = \frac{P(\vec{w} | \vec{z}, \vec{\beta})P(\vec{z} | \vec{\alpha})}{P(\vec{w} | \vec{\alpha}, \vec{\beta})}$$

- Likelihood

$$\begin{aligned} P(\vec{w} | \vec{z}, \vec{\beta}) &= \int P(\vec{w} | \vec{z}, \vec{\phi})P(\vec{\phi} | \vec{\beta}) d\vec{\phi} \\ &= \int \prod_{k=1}^K \text{Mult}(\vec{n}_k | \vec{\phi}_k) \text{Dir}(\vec{\phi}_k | \vec{\beta}) d\vec{\phi} \\ &= \prod_{k=1}^K \frac{\Delta(\vec{n}_k + \vec{\beta})}{\Delta(\vec{\beta})} \end{aligned}$$

Inference – Prior

- Given document \vec{w} , sample topic assignment \vec{z} from posterior

$$P(\vec{z} | \vec{w}, \vec{\alpha}, \vec{\beta}) = \frac{P(\vec{w} | \vec{z}, \vec{\beta}) P(\vec{z} | \vec{\alpha})}{P(\vec{w} | \vec{\alpha}, \vec{\beta})}$$

- Prior

$$\begin{aligned} P(\vec{z} | \vec{\alpha}) &= \int P(\vec{z} | \vec{\theta}) P(\vec{\theta} | \vec{\alpha}) d\vec{\theta} \\ &= \int \prod_{d=1}^D \text{Mult}(\vec{n}_d | \vec{\theta}_d) \text{Dir}(\vec{\theta}_d | \vec{\alpha}) d\vec{\theta} \\ &= \prod_{d=1}^D \frac{\Delta(\vec{n}_d + \vec{\alpha})}{\Delta(\vec{\alpha})} \end{aligned}$$

Inference – Collapsed Conditional

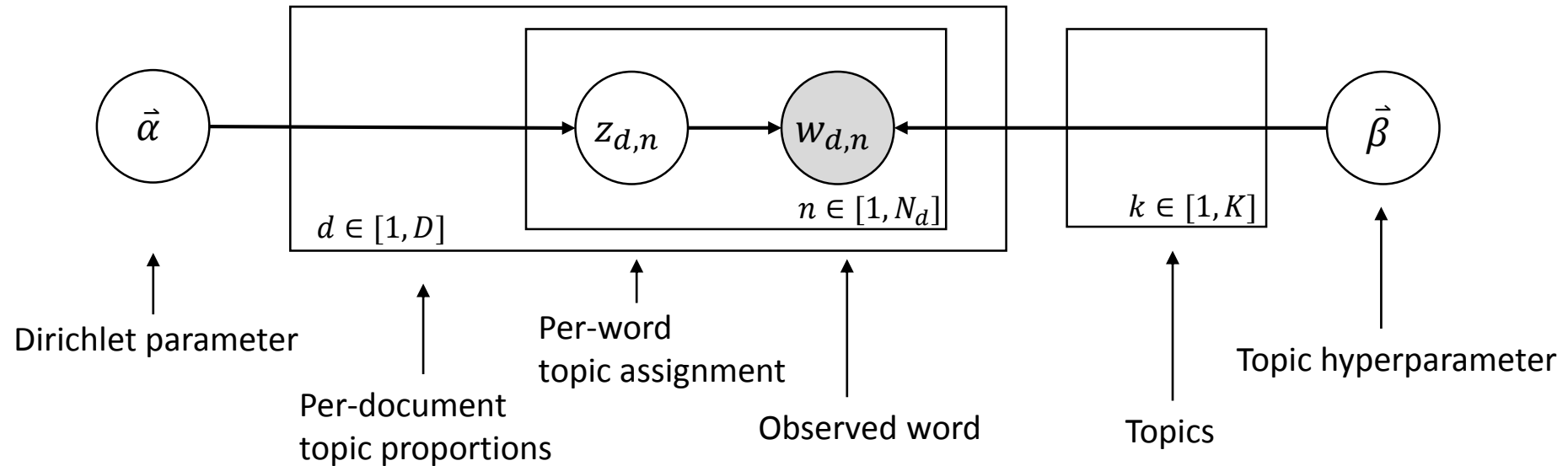
- Given document \vec{w} , sample topic assignment \vec{z} from posterior

$$P(\vec{z} | \vec{w}, \vec{\alpha}, \vec{\beta}) = \frac{P(\vec{w} | \vec{z}, \vec{\beta})P(\vec{z} | \vec{\alpha})}{P(\vec{w} | \vec{\alpha}, \vec{\beta})}$$

- For $\vec{w} = \{w_i = t, \vec{w}_{-i}\}$ and $\vec{z} = \{z_i = k, \vec{z}_{-i}\}$

$$\begin{aligned} P(z_i = k | \vec{z}_{-i}, \vec{w}, \vec{\alpha}, \vec{\beta}) &= \frac{P(\vec{w} | \vec{z}, \vec{\beta})P(\vec{z} | \vec{\alpha})}{P(\vec{w}, \vec{z}_{-i} | \vec{\alpha}, \vec{\beta})} = \frac{P(\vec{w} | \vec{z}, \vec{\beta})}{P(\vec{w}_{-i} | \vec{z}_{-i}, \vec{\beta})P(w_i)} \frac{P(\vec{z} | \vec{\alpha})}{P(\vec{z}_{-i} | \vec{\alpha})} \\ &\propto \frac{\Delta(\vec{n}_z + \vec{\beta})}{\Delta(\vec{n}_{z,-i} + \vec{\beta})} \frac{\Delta(\vec{n}_d + \vec{\alpha})}{\Delta(\vec{n}_{m,-i} + \vec{\alpha})} \\ &\propto \frac{n_{k,-i}^{(t)} + \beta_t}{\sum_{t'=1}^V n_{k,-i}^{(t')} + \beta'_t} \left(n_{d,-i}^{(k)} + \alpha_k \right) \end{aligned}$$

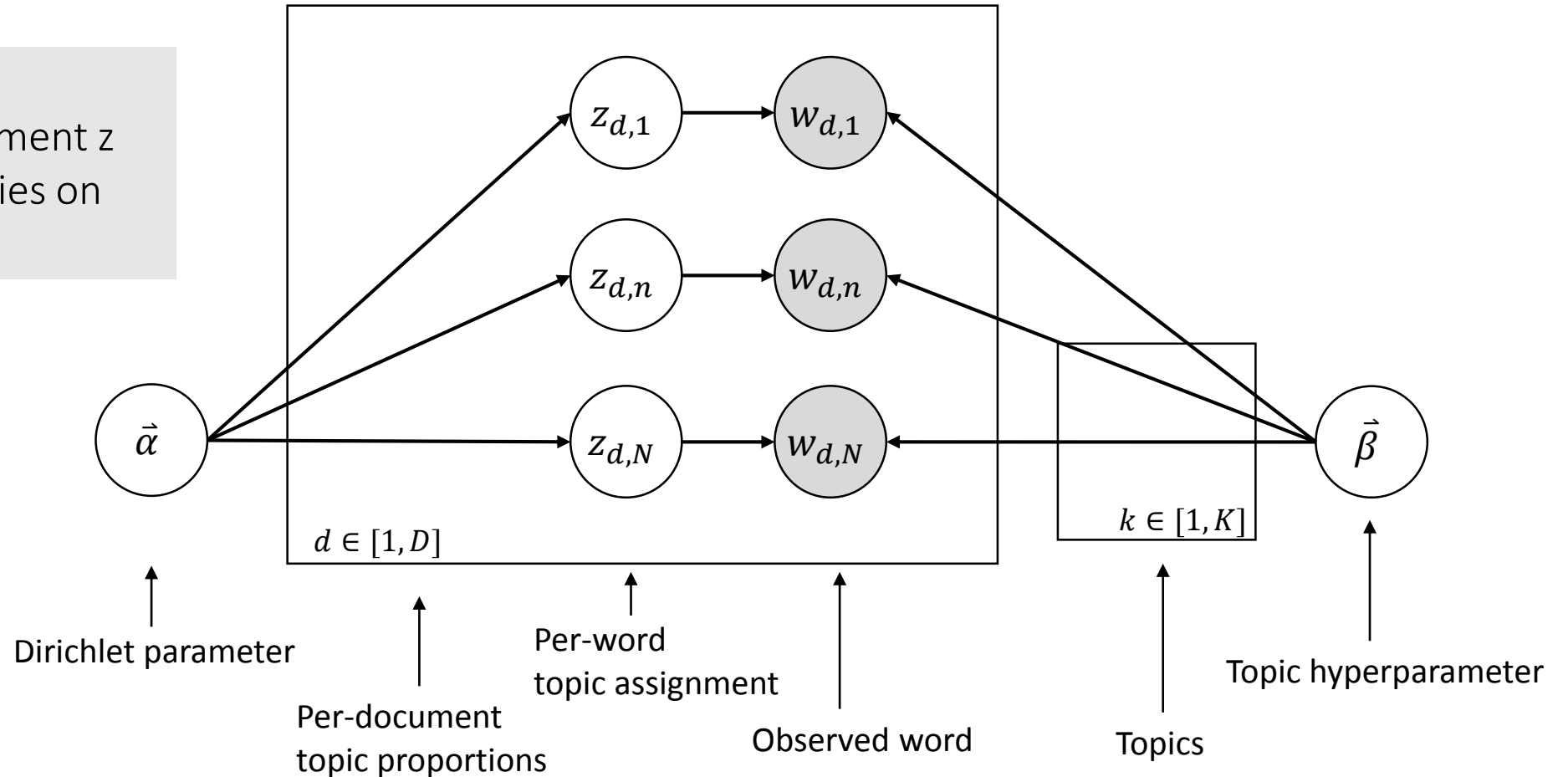
Inference – Collapsed Conditional



Inference – Gibbs Sampling

Stage 1 Initialization:

Initialize topic assignment z with equal probabilities on topics.

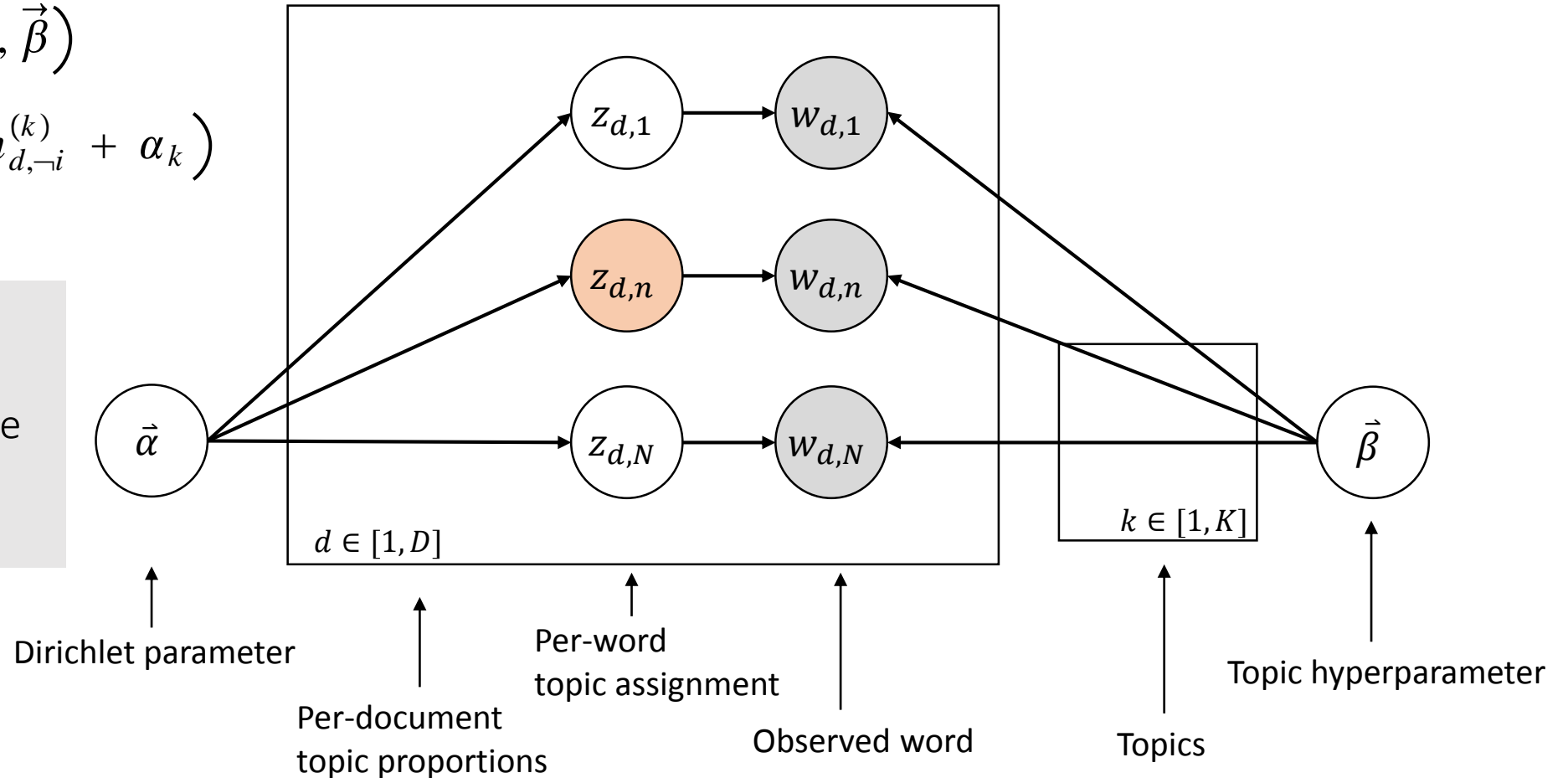


Inference – Gibbs Sampling

$$P(z_i = k | \vec{z}_{-i}, \vec{w}, \vec{\alpha}, \vec{\beta})$$

$$\propto \frac{n_{k,-i}^{(i)} + \beta_t}{\sum_{t'=1}^V n_{k,-i}^{(t')} + \beta'_t} (n_{d,-i}^{(k)} + \alpha_k)$$

Stage 2 Burn-In Period:
Randomly select topic assignment z . Update the assignment acc. to the collapsed conditional.

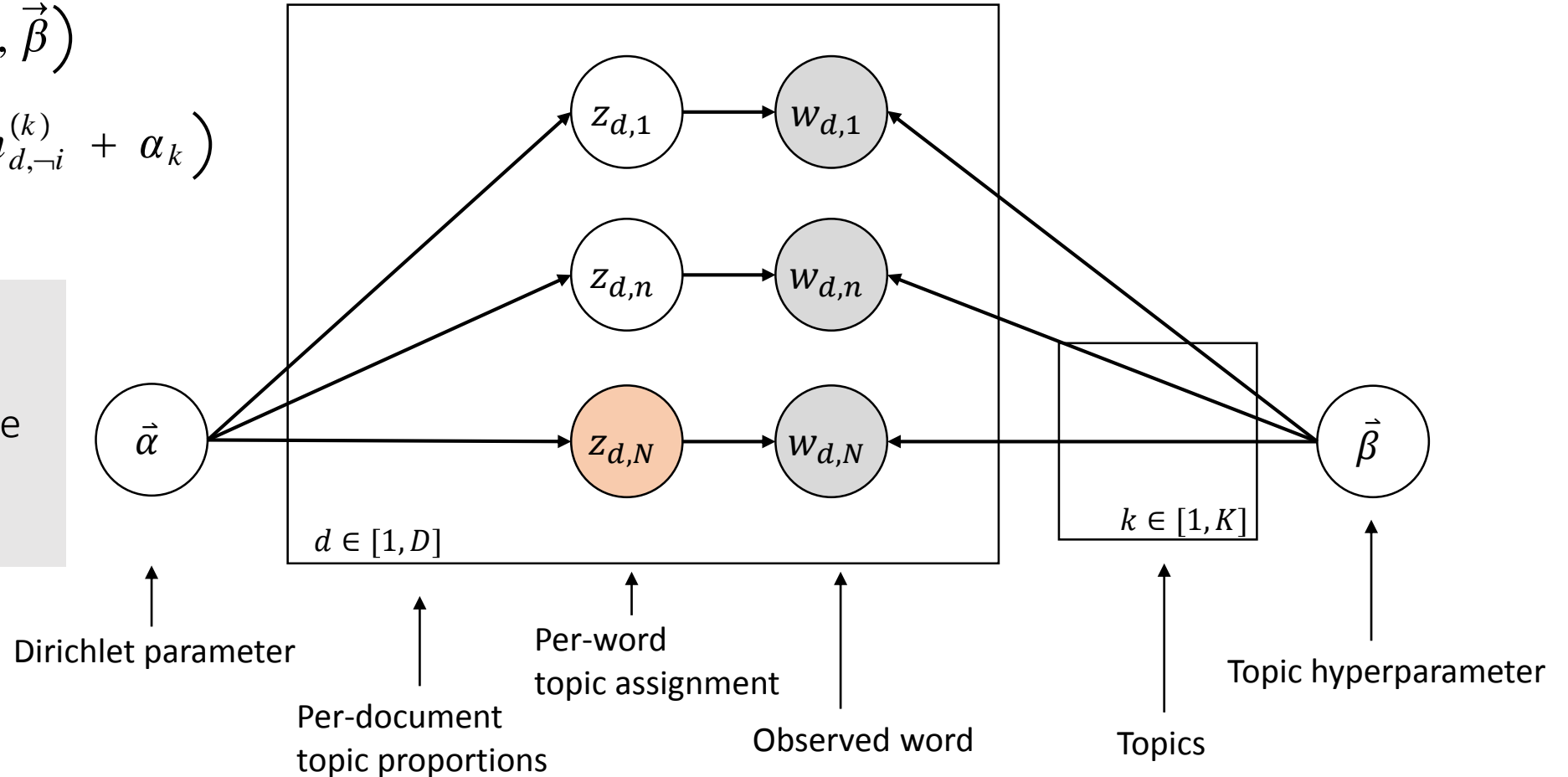


Inference – Gibbs Sampling

$$P(z_i = k | \vec{z}_{-i}, \vec{w}, \vec{\alpha}, \vec{\beta})$$

$$\propto \frac{n_{k,-i}^{(i)} + \beta_t}{\sum_{t'=1}^V n_{k,-i}^{(t')} + \beta'_t} (n_{d,-i}^{(k)} + \alpha_k)$$

Stage 2 Burn-In Period:
Randomly select topic assignment z . Update the assignment acc. to the collapsed conditional.

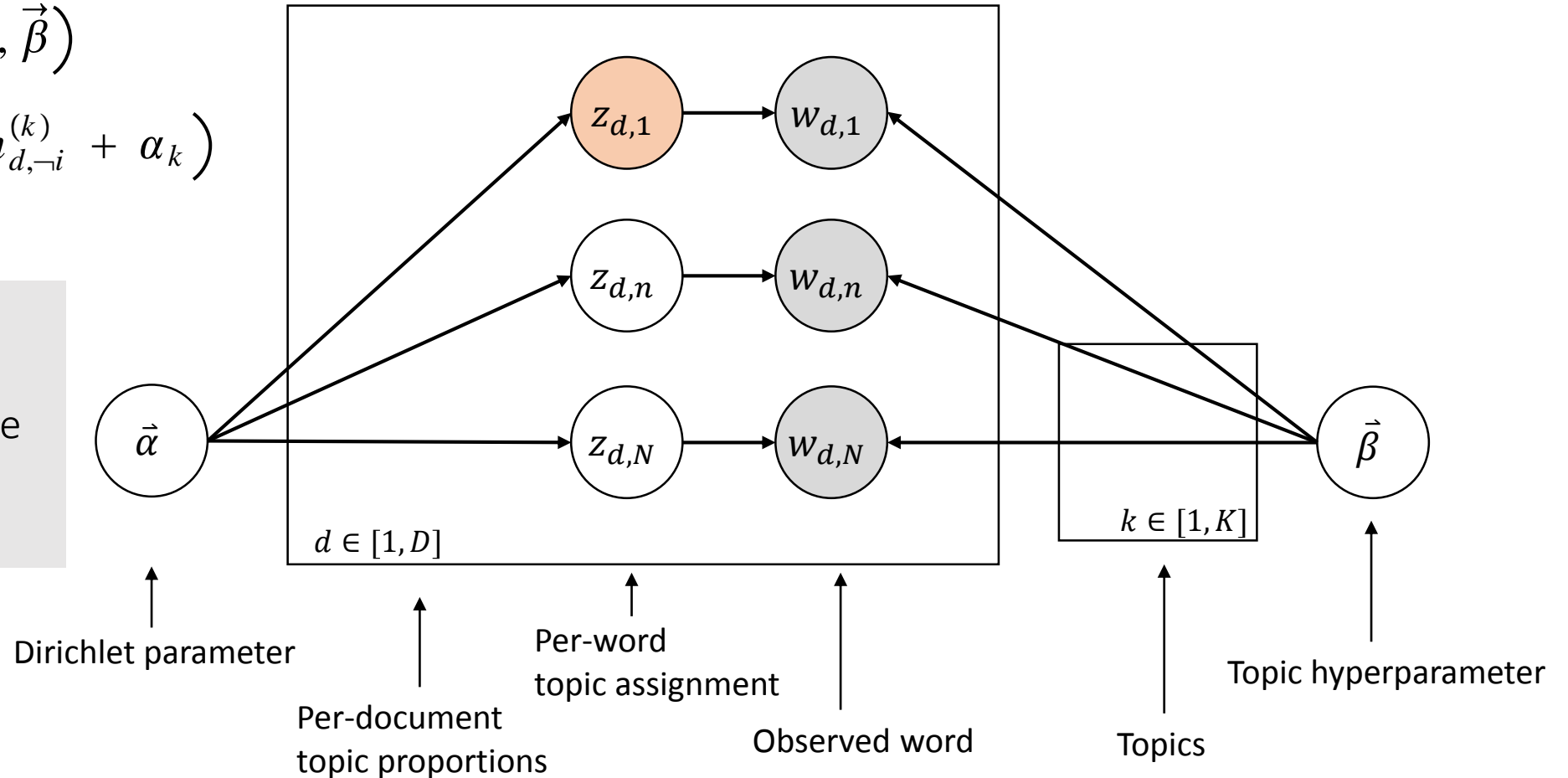


Inference – Gibbs Sampling

$$P(z_i = k | \vec{z}_{-i}, \vec{w}, \vec{\alpha}, \vec{\beta})$$

$$\propto \frac{n_{k,-i}^{(i)} + \beta_t}{\sum_{t'=1}^V n_{k,-i}^{(t')} + \beta'_t} (n_{d,-i}^{(k)} + \alpha_k)$$

Stage 2 Burn-In Period:
Randomly select topic assignment z . Update the assignment acc. to the collapsed conditional.



Inference – Gibbs Sampling

Stage 3 Terminate:

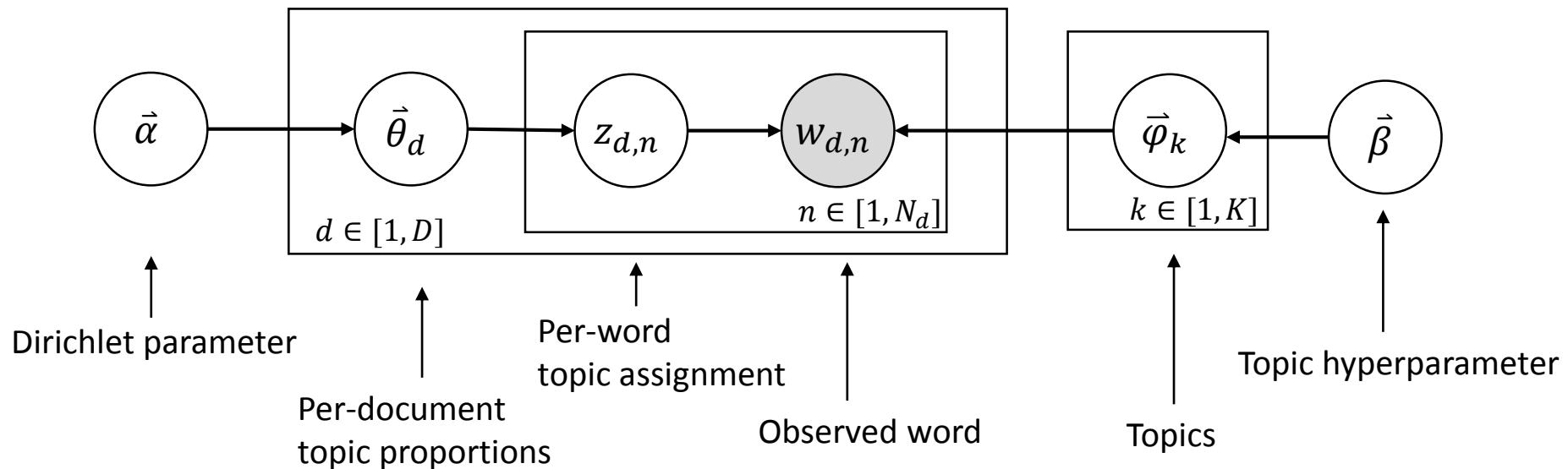
When stop criteria is met, read out word proportion and topic proportion from topic assignment counts.

$$P(\vec{\theta}_d | \vec{z}_d, \vec{\alpha}) = \text{Dir}(\vec{\theta}_d | \vec{\alpha} + \vec{n}_d)$$

$$P(\vec{\varphi}_k | \vec{z}, \vec{w}, \vec{\beta}) = \text{Dir}(\vec{\varphi}_k | \vec{\beta} + \vec{n}_k)$$

$$\mathbf{E}(\theta_{d,k}) = \frac{n_d^{(k)} + \alpha_k}{\sum_{k'=1}^K n_d^{(k')} + \alpha_{k'}}$$

$$\mathbf{E}(\varphi_{k,t}) = \frac{n_k^{(t)} + \beta_t}{\sum_{t'=1}^V n_k^{(t')} + \beta_{t'}}$$



Inference – Gibbs Sampler Initialization

Algorithm $\{\vec{\phi}, \vec{\theta}, \vec{z}\} = \text{LDA_Gibbs}(\{\vec{w}\}, \vec{\alpha}, \vec{\beta}, K)$

//initialize counts and topic assignment \vec{z}

sample topic indices $\vec{z} \sim \text{Mult}(1/K)$ for every word in $\{\vec{w}\}$

set document-topic count $n_d^{(k)} = \|\{z_{d,n} = k \mid n \in [1, N_d]\}\|$

set topic-term count $n_k^{(t)} = \|\{w_{d,n} = t \mid z_{d,n} = k\}\|$

Inference – Gibbs Sampler Burn-In Period

```
while not finished do  
  for all documents  $d \in [1, D]$  do  
    //for the current assignment of topic  $k$  to term  $t$  for word  $w_{d,n}$   
    decrease counts:  $n_d^{(k)} -= 1$ ,  $n_k^{(t)} -= 1$   
    //multinomial sampling for each component  $i$   
    sample topic index  $k' \sim P(z_i | \vec{z}_{-i}, \vec{w}, \vec{\alpha}, \vec{\beta})$   
    //for the new assignment of  $z_{d,n}$  to the term  $t$  for word  $w_{d,n}$   
    increment counts:  $n_d^{(k')} += 1$ ,  $n_{k'}^{(t)} += 1$   
  end for  
  
  if converged or  $L$  sampling iterations since last read out then  
    read out parameter estimates  $\mathbf{E}(\theta_{d,k})$  and  $\mathbf{E}(\varphi_{k,t})$ .  
  end if  
end while
```

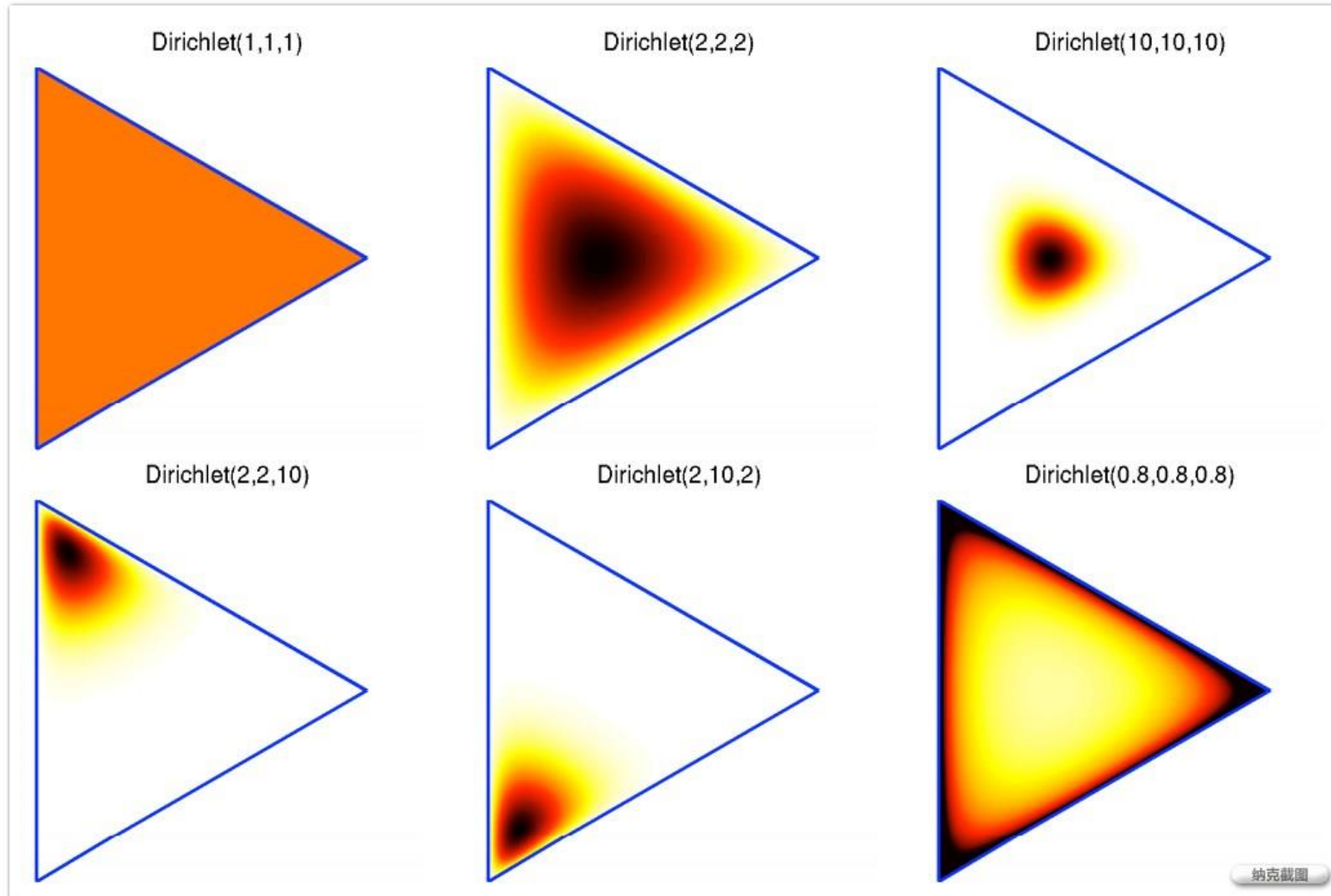
Inference - Dirichlet Parameters

- Dirichlet parameters describe our belief in the outcome.
- Dirichlet parameters less than one penalize dense topic/word proportions.
- Suggested value [Griffiths&Steyvers 2004]

$$\alpha = 50/K$$

$$\beta = 0.01$$

Dirichlet Pseudo-Count when $K=3$



Applications

- Given new document \tilde{w} , Gibbs sampling update:

$$P(\tilde{z}_i = k \mid \tilde{w}_i = t, \tilde{z}_{-i}, \tilde{w}_{-i}, \vec{\phi}, \vec{\theta}) \propto \varphi_{k,t} * \left(n_{d',-i}^{(k)} + \alpha_k \right)$$

$$\mathbf{E}(\theta_{d',k}) = \frac{n_{d'}^{(k)} + \alpha_k}{\sum_{k'=1}^K n_{d'}^{(k')} + \alpha_{k'}}$$

- Thus we obtain the topic proportion for documents and queries. Information retrieval and classification can be carried out on top of this representation.

References

- Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." *the Journal of machine Learning research* 3 (2003)
- Heinrich, Gregor. "Parameter estimation for text analysis." Technical report, 2005.
- Hofmann, Thomas. "Probabilistic latent semantic indexing." *Proceedings of the 22nd annual international ACM SIGIR*, 1999.
- http://videolectures.net/mlss09uk_blei_tm by David Blei