

Latent Dirichlet Allocation - Derivations & Proofs

Paul F. Roysdon, Ph.D.

Contents

1 Mathematical Derivations & Proofs	1
1.1 Introduction	1
1.2 Data and Notation	1
1.3 Model Formulation (Generative Process)	2
1.4 Conjugacy and Collapsed Marginals	2
1.5 Inference & Learning I: Mean-Field Variational Bayes	2
1.6 Inference & Learning II: Collapsed Gibbs Sampling	3
1.7 Hyperparameter Estimation (Optional)	4
1.8 Predictive Inference and Perplexity	4
1.9 Algorithmic Summary	4
1.10 Computational Aspects	5
1.11 Summary of Variables and Their Dimensions	5
1.12 Summary	5

1 Mathematical Derivations & Proofs

1.1 Introduction

Latent Dirichlet Allocation (LDA) is a generative probabilistic model of discrete data (e.g., text corpora), in which each document is modeled as a mixture over K latent topics, and each topic is a distribution over a vocabulary of size V . The key Bayesian ingredients are Dirichlet priors on (i) per-document topic proportions and (ii) per-topic word distributions, yielding conjugacy with categorical/multinomial likelihoods. We derive the full generative model, show Dirichlet–multinomial conjugacy, develop two standard inference/learning procedures—*mean-field variational Bayes* and *collapsed Gibbs sampling*—and provide complete algorithmic steps, proofs, and dimensions.

1.2 Data and Notation

Let the corpus consist of D documents. Document d has N_d tokens:

$$\mathbf{w}_d = (w_{d1}, \dots, w_{dN_d}), \quad w_{dn} \in \{1, \dots, V\}.$$

LDA posits K topics. For each token there is a latent *topic assignment*

$$\mathbf{z}_d = (z_{d1}, \dots, z_{dN_d}), \quad z_{dn} \in \{1, \dots, K\}.$$

Per-document topic proportions $\boldsymbol{\theta}_d \in \Delta^{K-1}$, per-topic word distributions $\boldsymbol{\beta}_k \in \Delta^{V-1}$. Stack topics as a matrix

$$\mathbf{B} = [\boldsymbol{\beta}_1 \cdots \boldsymbol{\beta}_K] \in [0, 1]^{V \times K}, \quad \sum_{v=1}^V \beta_{kv} = 1 \quad \forall k.$$

Priors:

$$\boldsymbol{\alpha} \in \mathbb{R}_{>0}^K \quad (\text{Dirichlet on } \boldsymbol{\theta}_d), \quad \boldsymbol{\eta} \in \mathbb{R}_{>0}^V \quad (\text{Dirichlet on each } \boldsymbol{\beta}_k).$$

Dimensions: D (docs), K (topics), V (vocab), N_d (# tokens in d), $N = \sum_d N_d$ (# tokens).

1.3 Model Formulation (Generative Process)

For each topic $k = 1, \dots, K$:

$$\boldsymbol{\beta}_k \sim \text{Dirichlet}(\boldsymbol{\eta}).$$

For each document $d = 1, \dots, D$:

$$\boldsymbol{\theta}_d \sim \text{Dirichlet}(\boldsymbol{\alpha}), \quad z_{dn} \mid \boldsymbol{\theta}_d \sim \text{Categorical}(\boldsymbol{\theta}_d), \quad w_{dn} \mid z_{dn}=k, \mathbf{B} \sim \text{Categorical}(\boldsymbol{\beta}_k), \quad n = 1, \dots, N_d.$$

Joint density (probability). With $\alpha_0 = \sum_k \alpha_k$, $\eta_0 = \sum_v \eta_v$, the complete joint density (probability) is

$$p(\mathbf{w}, \mathbf{z}, \{\boldsymbol{\theta}_d\}, \{\boldsymbol{\beta}_k\} \mid \boldsymbol{\alpha}, \boldsymbol{\eta}) = \left[\prod_{k=1}^K \text{Dir}(\boldsymbol{\beta}_k \mid \boldsymbol{\eta}) \right] \left[\prod_{d=1}^D \text{Dir}(\boldsymbol{\theta}_d \mid \boldsymbol{\alpha}) \prod_{n=1}^{N_d} \theta_{d,z_{dn}} \beta_{z_{dn},w_{dn}} \right], \quad (1)$$

where $\text{Dir}(\cdot)$ is the Dirichlet function.

1.4 Conjugacy and Collapsed Marginals

Dirichlet–multinomial integral (proof). For counts $\mathbf{n} = (n_1, \dots, n_K)$ with $\sum_k n_k = N$,

$$\int_{\Delta^{K-1}} \left(\prod_{k=1}^K \theta_k^{n_k} \right) \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \theta_k^{\alpha_k-1} d\boldsymbol{\theta} = \frac{B(\boldsymbol{\alpha} + \mathbf{n})}{B(\boldsymbol{\alpha})} = \frac{\Gamma(\alpha_0)}{\prod_k \Gamma(\alpha_k)} \cdot \frac{\prod_k \Gamma(\alpha_k + n_k)}{\Gamma(\alpha_0 + N)}.$$

This is the Dirichlet–multinomial (a.k.a. Polya) compound; $B(\cdot)$ is the multivariate Beta. ■

Collapsed joint $p(\mathbf{z}, \mathbf{w} \mid \boldsymbol{\alpha}, \boldsymbol{\eta})$. Let n_{dk} be the # of tokens in document d assigned to topic k ; n_{kv} the # of tokens using word v assigned to topic k ; $n_k = \sum_v n_{kv}$. Integrating Eqn. (1) over $\{\boldsymbol{\theta}_d\}$ and $\{\boldsymbol{\beta}_k\}$ yields

$$p(\mathbf{z}, \mathbf{w} \mid \boldsymbol{\alpha}, \boldsymbol{\eta}) = \prod_{d=1}^D \frac{\Gamma(\alpha_0)}{\prod_{k=1}^K \Gamma(\alpha_k)} \cdot \frac{\prod_{k=1}^K \Gamma(\alpha_k + n_{dk})}{\Gamma(\alpha_0 + N_d)} \times \prod_{k=1}^K \frac{\Gamma(\eta_0)}{\prod_{v=1}^V \Gamma(\eta_v)} \cdot \frac{\prod_{v=1}^V \Gamma(\eta_v + n_{kv})}{\Gamma(\eta_0 + n_k)}. \quad (2)$$

1.5 Inference & Learning I: Mean-Field Variational Bayes

We approximate the posterior with a factorized family

$$q(\{\boldsymbol{\beta}_k\}, \{\boldsymbol{\theta}_d\}, \mathbf{z}) = \left[\prod_{k=1}^K \text{Dir}(\boldsymbol{\beta}_k \mid \boldsymbol{\lambda}_k) \right] \left[\prod_{d=1}^D \text{Dir}(\boldsymbol{\theta}_d \mid \boldsymbol{\gamma}_d) \prod_{n=1}^{N_d} \text{Categorical}(z_{dn} \mid \boldsymbol{\phi}_{dn}) \right],$$

with variational parameters $\boldsymbol{\lambda}_k \in \mathbb{R}_{>0}^V$, $\boldsymbol{\gamma}_d \in \mathbb{R}_{>0}^K$, and $\boldsymbol{\phi}_{dn} \in \Delta^{K-1}$.

Evidence lower bound (ELBO).

$$\mathcal{L}(q) = \mathbb{E}_q[\log p(\mathbf{w}, \mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta} \mid \boldsymbol{\alpha}, \boldsymbol{\eta})] - \mathbb{E}_q[\log q(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta})]. \quad (3)$$

Coordinate ascent updates follow from setting derivatives of \mathcal{L} to zero.

Expected logs. For digamma $\psi(\cdot)$,

$$\mathbb{E}[\log \theta_{dk}] = \psi(\gamma_{dk}) - \psi\left(\sum_{t=1}^K \gamma_{dt}\right), \quad \mathbb{E}[\log \beta_{kv}] = \psi(\lambda_{kv}) - \psi\left(\sum_{u=1}^V \lambda_{ku}\right).$$

Coordinate ascent (CAVI) updates. For each document d and token n with word $v = w_{dn}$,

$$\phi_{dnk} \propto \exp\left(\mathbb{E}[\log \theta_{dk}] + \mathbb{E}[\log \beta_{kv}]\right) = \exp\left(\psi(\gamma_{dk}) - \psi(\sum_t \gamma_{dt}) + \psi(\lambda_{kv}) - \psi(\sum_u \lambda_{ku})\right), \quad (4)$$

$$\gamma_{dk} = \alpha_k + \sum_{n=1}^{N_d} \phi_{dnk}, \quad (5)$$

$$\lambda_{kv} = \eta_v + \sum_{d=1}^D \sum_{n=1}^{N_d} \phi_{dnk} \mathbf{1}\{w_{dn} = v\}. \quad (6)$$

Normalize ϕ_{dn} after Eqn. (4). Iterate Eqns. (4)–(6) until ELBO convergence.

Variational EM (ML topics). If we treat \mathbf{B} as *parameters* (not random) and maximize w.r.t. β_k under $\sum_v \beta_{kv} = 1$, then the M-step yields

$$\beta_{kv} \propto \sum_{d,n} \phi_{dnk} \mathbf{1}\{w_{dn} = v\} \quad (\text{optionally smoothed: } + \eta_v - 1), \quad (7)$$

followed by normalization across v .

Proof (CAVI). *Proof.* Taking $\partial \mathcal{L}/\partial \phi_{dnk}$ with normalization constraint gives $\log \phi_{dnk} \propto \mathbb{E}[\log \theta_{dk}] + \mathbb{E}[\log \beta_{kv}]$. Dirichlet–multinomial conjugacy in the complete conditionals yields Eqn. (5) and Eqn. (6). ■

1.6 Inference & Learning II: Collapsed Gibbs Sampling

Integrate out $\{\boldsymbol{\theta}_d\}$ and $\{\boldsymbol{\beta}_k\}$ to sample topic assignments \mathbf{z} from Eqn. (2). Let $v = w_{dn}$. Define “leave-one-out” counts that exclude token (d, n) :

$$n_{dk}^{\neg dn}, \quad n_{kv}^{\neg dn}, \quad n_k^{\neg dn} = \sum_v n_{kv}^{\neg dn}.$$

Collapsed conditional (derivation). From Eqn. (2), the ratio of the joint with $z_{dn}=k$ versus not depends only on factors involving n_{dk} and n_{kv} . Using $\Gamma(x+1) = x \Gamma(x)$,

$$\begin{aligned} p(z_{dn}=k \mid \mathbf{z}^{\neg dn}, \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\eta}) &\propto \frac{\Gamma(\alpha_k + n_{dk}^{\neg dn} + 1)}{\Gamma(\alpha_k + n_{dk}^{\neg dn})} \cdot \frac{\Gamma(\eta_v + n_{kv}^{\neg dn} + 1)}{\Gamma(\eta_v + n_{kv}^{\neg dn})} \cdot \frac{\Gamma(\eta_0 + n_k^{\neg dn})}{\Gamma(\eta_0 + n_k^{\neg dn} + 1)} \\ &= (\alpha_k + n_{dk}^{\neg dn}) \cdot \frac{\eta_v + n_{kv}^{\neg dn}}{\eta_0 + n_k^{\neg dn}}. \end{aligned} \quad (8)$$

After normalization over k ,

$$p(z_{dn}=k \mid \cdot) \propto (\alpha_k + n_{dk}^{\neg dn}) \cdot \frac{\eta_v + n_{kv}^{\neg dn}}{\eta_0 + n_k^{\neg dn}}. \quad (9)$$

Sampling scheme. Sweep all tokens, update counts in-place; after burn-in, thin/summarize $\{n_{kv}\}$ to estimate topics:

$$\hat{\beta}_{kv} = \frac{\eta_v + n_{kv}}{\eta_0 + n_k}, \quad \hat{\theta}_{dk} = \frac{\alpha_k + n_{dk}}{\alpha_0 + N_d}.$$

1.7 Hyperparameter Estimation (Optional)

Updating α (per-document Dirichlet). Maximize $\sum_d \log p(\theta_d \mid \alpha)$ under variational posterior by replacing $\log \theta_{dk}$ with $\mathbb{E}[\log \theta_{dk}]$:

$$\mathcal{J}(\alpha) = D \left(\log \Gamma(\alpha_0) - \sum_k \log \Gamma(\alpha_k) \right) + \sum_d \sum_k (\alpha_k - 1) \left(\psi(\gamma_{dk}) - \psi(\sum_t \gamma_{dt}) \right).$$

Gradient and Hessian:

$$g_k = D \left(\psi(\alpha_0) - \psi(\alpha_k) \right) + \sum_d \left(\psi(\gamma_{dk}) - \psi(\sum_t \gamma_{dt}) \right), \quad H_{kl} = D \left(\psi'(\alpha_0) - \delta_{kl} \psi'(\alpha_k) \right).$$

Perform Newton updates in the positive orthant (e.g. Minkas method), with line search to ensure ascent.

Updating η (per-topic Dirichlet). Analogously, with variational $q(\beta_k) = \text{Dir}(\lambda_k)$,

$$\mathcal{J}(\eta) = K \left(\log \Gamma(\eta_0) - \sum_v \log \Gamma(\eta_v) \right) + \sum_k \sum_v (\eta_v - 1) \left(\psi(\lambda_{kv}) - \psi(\sum_u \lambda_{ku}) \right),$$

and compute Newton steps with gradient g_v and Hessian $H_{uv} = K \left(\psi'(\eta_0) - \delta_{uv} \psi'(\eta_v) \right)$.

1.8 Predictive Inference and Perplexity

Held-out document likelihood (VB). Given topics $\{\lambda_k\}$, infer γ_d, ϕ_{dn} on held-out document d^* via Eqns. (4)–(5); approximate predictive with expected logs:

$$\log p(\mathbf{w}_{d^*} \mid \text{model}) \approx \sum_n \log \sum_k \exp \left(\psi(\gamma_{d^*k}) - \psi(\sum_t \gamma_{d^*t}) \right) \exp \left(\psi(\lambda_{k,w_{d^*n}}) - \psi(\sum_u \lambda_{ku}) \right).$$

Perplexity.

$$\text{Perplexity}(\mathcal{D}_{\text{test}}) = \exp \left(- \frac{\sum_{d \in \text{test}} \log p(\mathbf{w}_d)}{\sum_{d \in \text{test}} N_d} \right).$$

1.9 Algorithmic Summary

Mean-field Variational Bayes (VB) for LDA.

1. **Initialize** λ_k (e.g. $\eta + \text{random}$), for each d : $\gamma_d \leftarrow \alpha + N_d/K$, $\phi_{dn} \leftarrow 1/K$.
2. **Repeat** until convergence:
 - (a) For each document d : iterate Eqn. (4), Eqn. (5) over tokens n .
 - (b) Global update: Eqn. (6) for all (k, v) .
 - (c) (Optional) Update hyperparameters α, η by Newton steps.
3. **Outputs:** topics via $\hat{\beta}_{kv} \propto \lambda_{kv}$, proportions $\hat{\theta}_{dk} \propto \gamma_{dk}$.

Collapsed Gibbs Sampling for LDA.

1. **Initialize** z_{dn} uniformly at random; compute counts n_{dk}, n_{kv}, n_k .
2. **For** sweeps $t = 1, \dots, T$:
 - (a) For each token (d, n) , decrement counts, sample z_{dn} from Eqn. (9), increment counts.
3. **Estimate** topics and proportions from posterior means (averaged over samples): $\hat{\beta}_{kv} = (\eta_v + n_{kv}) / (\eta_0 + n_k)$, $\hat{\theta}_{dk} = (\alpha_k + n_{dk}) / (\alpha_0 + N_d)$.

1.10 Computational Aspects

VB per-iteration cost is $O(\sum_d N_d K + KV)$ due to updates of all ϕ_{dn} and λ_k . Collapsed Gibbs has cost $O(NK)$ per sweep (or less with sparse/alias-based sampling). For large V , top- S vocab pruning or sparse λ updates reduce memory/time.

1.11 Summary of Variables and Their Dimensions

- D (docs), K (topics), V (vocab size), N_d (# tokens in doc d), $N = \sum_d N_d$.
- $\mathbf{w}_d \in \{1, \dots, V\}^{N_d}$: tokens of doc d ; w_{dn} is the n th token index.
- $\mathbf{z}_d \in \{1, \dots, K\}^{N_d}$: topic assignments; z_{dn} is the topic for w_{dn} .
- $\boldsymbol{\theta}_d \in \Delta^{K-1}$: topic mixture for doc d ; prior $\text{Dir}(\boldsymbol{\alpha})$.
- $\boldsymbol{\beta}_k \in \Delta^{V-1}$: word distribution for topic k ; prior $\text{Dir}(\boldsymbol{\eta})$.
- $\mathbf{B} = [\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K] \in [0, 1]^{V \times K}$: topic matrix.
- Counts: n_{dk} (# tokens in d assigned to k), n_{kv} (# tokens of word v assigned to k), $n_k = \sum_v n_{kv}$.
- Variational params: $\phi_{dn} \in \Delta^{K-1}$, $\gamma_d \in \mathbb{R}_{>0}^K$, $\lambda_k \in \mathbb{R}_{>0}^V$.
- Hyperparameters: $\boldsymbol{\alpha} \in \mathbb{R}_{>0}^K$, $\boldsymbol{\eta} \in \mathbb{R}_{>0}^V$; sums α_0, η_0 .

1.12 Summary

From first principles, LDA defines a hierarchical Bayesian model with Dirichlet priors on per-document topic mixtures and per-topic word distributions, yielding conjugate Dirichlet–multinomial structure. Integrating out the Dirichlet variables gives the collapsed joint Eqn. (2), from which the Gibbs conditional Eqn. (9) follows. Alternatively, mean-field variational Bayes maximizes an ELBO with closed-form coordinate updates Eqns. (4)–(6) (and ML/smoothed updates Eqn. (7)). We have stated all distributions, objectives, updates, proofs of conjugacy, and algorithmic procedures, together with explicit variable dimensions, in a notation consistent with prior sections.