

Machine Learning 2

D.1. Latent Dirichlet Allocation (LDA)

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)
Institute for Computer Science
University of Hildesheim, Germany

Syllabus

| | | | |
|------|-------|------|---|
| | | | A. Advanced Supervised Learning |
| Fri. | 24.4. | (1) | A.1 Generalized Linear Models |
| Fri. | 1.5. | — | — <i>Labour Day</i> — |
| Fri. | 8.5. | (2) | A.2 Gaussian Processes |
| Fri. | 15.5. | (3) | A.3 Advanced Support Vector Machines |
| | | | B. Ensembles |
| Fri. | 22.5. | (4) | B.1 Stacking & B.2 Boosting |
| Fri. | 29.5. | (5) | B.3 Mixtures of Experts |
| Fri. | 5.6. | — | — <i>Pentecoste Break</i> — |
| | | | C. Sparse Models |
| Fri. | 12.6. | (6) | C.1 Homotopy and Least Angle Regression |
| Fri. | 19.6. | (7) | C.2 Proximal Gradients |
| Fri. | 26.6. | (8) | C.3 Laplace Priors |
| Fri. | 3.7. | (9) | C.4 Automatic Relevance Determination |
| | | | D. Complex Predictors |
| Fri. | 10.7. | (10) | D.1 Latent Dirichlet Allocation (LDA) |
| Fri. | 17.7. | (11) | Q & A |

Outline

1. The LDA Model
2. Learning LDA via Gibbs Sampling
3. Learning LDA via Collapsed Gibbs Sampling
4. Learning LDA via Variational Inference
5. Supervised LDA

Outline

1. The LDA Model
2. Learning LDA via Gibbs Sampling
3. Learning LDA via Collapsed Gibbs Sampling
4. Learning LDA via Variational Inference
5. Supervised LDA

Documents / Finite Discrete Sequences

- ▶ instances $x_n \in \mathcal{A}^*$ are **discrete sequences**
 - ▶ $\mathcal{A} := \{1, \dots, A\}$ called **dictionary / alphabet** ($A \in \mathbb{N}$), where $a \in \mathcal{A}$ denotes the a -th **word / symbol / token**.
 - ▶ $\mathcal{A}^* := \bigcup_{\ell=1}^{\infty} \mathcal{A}^{\ell}$ called **documents / finite \mathcal{A} -sequences**.
 - ▶ $M_n := |x_n| := \ell$ called **length** (for $x_n \in \mathcal{A}^{\ell}$).
 - ▶ $x_{n,m}$ called **m -th word of x_n** .

- ▶ if there are no sequential effects (order does not matter), documents can be described by their **word frequencies (bag of words)**:

$$\tilde{x}_{n,a} := |\{m \in \{1, \dots, |x_n|\} \mid x_{n,m} = a\}|, \quad a \in \mathcal{A}$$

The LDA Model

$$p(x_{n,m} \mid z_{n,m} = k, \phi) := \text{Cat}(x_{n,m} \mid \phi_k), \quad n = 1, \dots, N, m = 1, \dots, M_n$$

$$p(z_{n,m} \mid \pi_n) := \text{Cat}(z_{n,m} \mid \pi_n), \quad n = 1, \dots, N, m = 1, \dots, M_n$$

$$p(\phi_k \mid \beta) := \text{Dir}(\phi_k \mid \beta \mathbf{1}_A), \quad k = 1, \dots, K$$

$$p(\pi_n \mid \gamma) := \text{Dir}(\pi_n \mid \gamma \mathbf{1}_K), \quad n = 1, \dots, N$$

- ▶ $z_{n,m} \in \{1, \dots, K\}$: topic the m -th word of document n belongs to.
- ▶ $\phi_k \in \Delta^A$: word probabilities of topic k .
- ▶ $\pi_n \in \Delta^K$: topic probabilities of document n .
- ▶ $\beta, \gamma \in \mathbb{R}^+$: priors of ϕ and π .

Note: $\Delta^K := \{z \in \mathbb{R}^K \mid z \geq 0, \sum_{k=1}^K z_k = 1\}$.

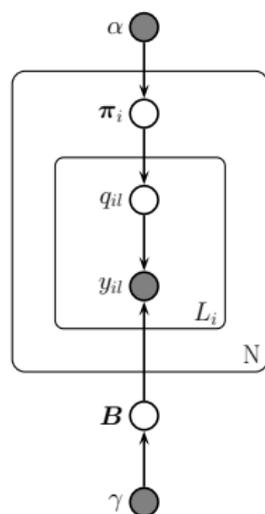
The LDA Model

$$p(x_{n,m} \mid z_{n,m} = k, \phi) := \text{Cat}(x_{n,m} \mid \phi_k),$$

$$p(z_{n,m} \mid \pi_n) := \text{Cat}(z_{n,m} \mid \pi_n),$$

$$p(\phi_k \mid \beta) := \text{Dir}(\phi_k \mid \beta \mathbf{1}_A),$$

$$p(\pi_n \mid \gamma) := \text{Dir}(\pi_n \mid \gamma \mathbf{1}_K),$$



- ▶ $z_{n,m} \in \{1, \dots, K\}$: topic the m -th word of document n belongs to.
- ▶ $\phi_k \in \Delta^A$: word probabilities of topic k .
- ▶ $\pi_n \in \Delta^K$: topic probabilities of document n .
- ▶ $\beta, \gamma \in \mathbb{R}^+$: priors of ϕ and π .

Note: $\Delta^K := \{z \in \mathbb{R}^K \mid z \geq 0, \sum_{k=1}^K z_k = 1\}$.

[Mur12, fig. 27.2]



Example $p(x_{n,m} | z_{n,m}, \phi)$

Topic 77

| word | prob. |
|-------------|-------|
| MUSIC | .090 |
| DANCE | .034 |
| SONG | .033 |
| PLAY | .030 |
| SING | .026 |
| SINGING | .026 |
| BAND | .026 |
| PLAYED | .023 |
| SANG | .022 |
| SONGS | .021 |
| DANCING | .020 |
| PIANO | .017 |
| PLAYING | .016 |
| RHYTHM | .015 |
| ALBERT | .013 |
| MUSICAL | .013 |

Topic 82

| word | prob. |
|-------------|-------|
| LITERATURE | .031 |
| POEM | .028 |
| POETRY | .027 |
| POET | .020 |
| PLAYS | .019 |
| POEMS | .019 |
| PLAY | .015 |
| LITERARY | .013 |
| WRITERS | .013 |
| DRAMA | .012 |
| WROTE | .012 |
| POETS | .011 |
| WRITER | .011 |
| SHAKESPEARE | .010 |
| WRITTEN | .009 |
| STAGE | .009 |

Topic 166

| word | prob. |
|-------------|-------|
| PLAY | .136 |
| BALL | .129 |
| GAME | .065 |
| PLAYING | .042 |
| HIT | .032 |
| PLAYED | .031 |
| BASEBALL | .027 |
| GAMES | .025 |
| BAT | .019 |
| RUN | .019 |
| THROW | .016 |
| BALLS | .015 |
| TENNIS | .011 |
| HOME | .010 |
| CATCH | .010 |
| FIELD | .010 |

[Mur12, fig. 27.4]

Example $x_{n,m}, z_{n,m}$

Document #29795

Bix beiderbecke, at age⁰⁶⁰ fifteen²⁰⁷, sat¹⁷⁴ on the slope⁰⁷¹ of a bluff⁰⁵⁵ overlooking⁰²⁷ the mississippi¹³⁷ river¹³⁷. He was listening⁰⁷⁷ to music⁰⁷⁷ coming⁰⁰⁹ from a passing⁰⁴³ riverboat. The music⁰⁷⁷ had already captured⁰⁰⁶ his heart¹⁵⁷ as well as his ear¹¹⁹. It was jazz⁰⁷⁷. Bix beiderbecke had already had music⁰⁷⁷ lessons⁰⁷⁷. He showed⁰⁰² promise¹³⁴ on the piano⁰⁷⁷, and his parents⁰³⁵ hoped²⁶⁸ he might consider¹¹⁸ becoming a concert⁰⁷⁷ pianist⁰⁷⁷. But bix was interested²⁶⁸ in another kind⁰⁵⁰ of music⁰⁷⁷. He wanted²⁶⁸ to play⁰⁷⁷ the cornet. And he wanted²⁶⁸ to play⁰⁷⁷ jazz⁰⁷⁷ ...

Document #1883

There is a simple⁰⁵⁰ reason¹⁰⁶ why there are so few periods⁰⁷⁸ of really great theater⁰⁸² in our whole western⁰⁴⁶ world. Too many things³⁰⁰ have to come right at the very same time. The dramatists must have the right actors⁰⁸², the actors⁰⁸² must have the right playhouses, the playhouses must have the right audiences⁰⁸². We must remember²⁸⁸ that plays⁰⁸² exist¹⁴³ to be performed⁰⁷⁷, not merely⁰⁵⁰ to be read²⁵⁴. (even when you read²⁵⁴ a play⁰⁸² to yourself, try²⁸⁸ to perform⁰⁶² it, to put¹⁷⁴ it on a stage⁰⁷⁸, as you go along.) as soon⁰²⁸ as a play⁰⁸² has to be performed⁰⁸², then some kind¹²⁶ of theatrical⁰⁸² ...

Document #21359

Jim²⁹⁶ has a game¹⁶⁶ book²⁵⁴. Jim²⁹⁶ reads²⁵⁴ the book²⁵⁴. Jim²⁹⁶ sees⁰⁸¹ a game¹⁶⁶ for one. Jim²⁹⁶ plays¹⁶⁶ the game¹⁶⁶. Jim²⁹⁶ likes⁰⁸¹ the game¹⁶⁶ for one. The game¹⁶⁶ book²⁵⁴ helps⁰⁸¹ jim²⁹⁶. Don¹⁸⁰ comes⁰⁴⁰ into the house⁰³⁸. Don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the game¹⁶⁶ book²⁵⁴. The boys⁰²⁰ see a game¹⁶⁶ for two. The two boys⁰²⁰ play¹⁶⁶ the game¹⁶⁶. The boys⁰²⁰ play¹⁶⁶ the game¹⁶⁶ for two. The boys⁰²⁰ like the game¹⁶⁶. Meg²⁸² comes⁰⁴⁰ into the house²⁸². Meg²⁸² and don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the book²⁵⁴. They see a game¹⁶⁶ for three. Meg²⁸² and don¹⁸⁰ and jim²⁹⁶ play¹⁶⁶ the game¹⁶⁶. They play¹⁶⁶...

[Mur12, fig. 27.5]



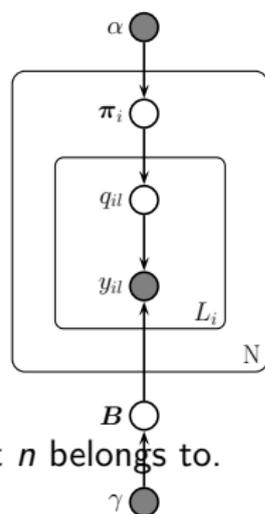
The LDA Model

$$p(x_{n,m} | z_{n,m} = k, \phi) := \text{Cat}(x_{n,m} | \phi_k)$$

$$p(z_{n,m} | \pi_n) := \text{Cat}(z_{n,m} | \pi_n)$$

$$p(\phi_k | \beta) := \text{Dir}(\phi_k | \beta \mathbf{1}_A)$$

$$p(\pi_n | \gamma) := \text{Dir}(\pi_n | \gamma \mathbf{1}_K)$$



- ▶ $z_{n,m} \in \{1, \dots, K\}$: topic the m -th word of document n belongs to.
- ▶ $\phi_k \in \Delta^A$: word probabilities of topic k .
- ▶ $\pi_n \in \Delta^K$: topic probabilities of document n .
- ▶ $\beta, \gamma \in \mathbb{R}^+$: priors of ϕ and π .

x_7 = "In Frankfurt, many banks are located at the banks of Main."

Q: How can LDA model that

- ▶ $x_{7,4}$ ="bank" denotes a credit institute and belongs to topic 1 "finance",
- ▶ $x_{7,9}$ ="bank" denotes a river bank and belongs to topic 2 "environment"?

Outline

1. The LDA Model
2. Learning LDA via Gibbs Sampling
3. Learning LDA via Collapsed Gibbs Sampling
4. Learning LDA via Variational Inference
5. Supervised LDA

Learning via Parameter Sampling

The posterior

$$p(\theta | \mathcal{D}) \propto p(\mathcal{D} | \theta)$$

describes the **distribution of the parameters given the data**.

If we can **sample parameters** from this distribution

$$\theta_1, \theta_2, \dots, \theta_S \sim p(\theta | \mathcal{D})$$

we can

- ▶ estimate **expected parameter values** and **their variances** from this parameter sample:

$$\hat{\theta} := E(\theta | \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S \theta_s, \quad V(\theta | \mathcal{D}) \approx \frac{1}{S-1} \sum_{s=1}^S (\theta_s - E(\theta | \mathcal{D}))^2$$

- ▶ **predict targets** for new instances x via model averaging:

$$p(y | x, \theta_{1:S}) = \frac{1}{S} \sum_{s=1}^S p(y | x, \theta_s)$$

Sampling

- ▶ for most closed-form distributions $p(x)$ there exist efficient sampling methods
 - ▶ categorical, normal, ...

- ▶ but most posteriors are not closed-form distributions.
 - ▶ but for example products thereof.

Gibbs Sampling

- ▶ task: sample from $p(x_1, \dots, x_N)$
- ▶ problem:
 - ▶ assume sampling from the joint distribution $p(x_1, \dots, x_N)$ is difficult.
 - ▶ assume sampling from marginals $p(x_n)$ or partial conditionals $p(x_n \mid \text{some } x_{n'})$ is also difficult.
 - ▶ assume sampling from all **full conditionals** $p(x_n \mid x_{-n})$ is easy.

Gibbs Sampling

- ▶ task: sample from $p(x_1, \dots, x_N)$
- ▶ problem:
 - ▶ assume sampling from the joint distribution $p(x_1, \dots, x_N)$ is difficult.
 - ▶ assume sampling from marginals $p(x_n)$ or partial conditionals $p(x_n \mid \text{some } x_{n'})$ is also difficult.
 - ▶ assume sampling from all **full conditionals** $p(x_n \mid x_{-n})$ is easy.

Gibbs sampling: given last sample x^s , sample x^{s+1} one variable at a time:

$$x_1^{s+1} \sim p(x_1 \mid x_{2:N} = x_{2:N}^s)$$

$$x_2^{s+1} \sim p(x_2 \mid x_{1:1} = x_{1:1}^{s+1}, x_{3:N} = x_{3:N}^s)$$

$$\vdots$$

$$x_n^{s+1} \sim p(x_n \mid x_{1:n-1} = x_{1:n-1}^{s+1}, x_{n+1:N} = x_{n+1:N}^s)$$

$$\vdots$$

$$x_N^{s+1} \sim p(x_N \mid x_{1:N-1} = x_{1:N-1}^{s+1})$$

Gibbs Sampling

- ▶ the distribution created by the Gibbs sampler eventually will converge to $p(x_1, \dots, x_N)$
- ▶ start Gibbs sampling with an arbitrary x^0
 - ▶ but ensure that $p(x^0) > 0$!
 - ▶ also consider restarts.
- ▶ throw away the first examples (**burn in**).
 - ▶ only after a while the chain has converged to the stationary distribution $p(x_1, \dots, x_N)$.
 - ▶ typical are 100-10,000 examples
- ▶ sometimes some variables can be marginalized out, improving the performance of the Gibbs sampler (**collapsed Gibbs sampling, Rao-Blackwellisation**)

Gibbs Sampling for LDA

$$\begin{aligned}
 p(x_{n,m} \mid z_{n,m} = k, \phi) &:= \text{Cat}(x_{n,m} \mid \phi_k) &&= \phi_{k,x_{n,m}} \\
 p(z_{n,m} \mid \pi_n) &:= \text{Cat}(z_{n,m} \mid \pi_n) &&= \pi_{n,z_{n,m}} \\
 p(\phi_k \mid \beta) &:= \text{Dir}(\phi_k \mid \beta \mathbf{1}_A) &&\propto \prod_{a=1}^A \phi_{k,a}^{\beta_a - 1} \\
 p(\pi_n \mid \gamma) &:= \text{Dir}(\pi_n \mid \gamma \mathbf{1}_K) &&\propto \prod_{k=1}^K \pi_{n,k}^{\gamma_k - 1}
 \end{aligned}$$

Full conditionals: 1. z

$$p(z_{n,m} = k \mid \phi, \pi_n) \propto p(x_{n,m} \mid z_{n,m} = k, \phi) p(z_{n,m} = k \mid \pi_n) = \phi_{k,x_{n,m}} \pi_{n,k}$$

Gibbs Sampling for LDA

$$p(x_{n,m} | z_{n,m} = k, \phi) := \text{Cat}(x_{n,m} | \phi_k) = \phi_{k,x_{n,m}}$$

$$p(z_{n,m} | \pi_n) := \text{Cat}(z_{n,m} | \pi_n) = \pi_{n,z_{n,m}}$$

$$p(\phi_k | \beta) := \text{Dir}(\phi_k | \beta \mathbf{1}_A) \propto \prod_{a=1}^A \phi_{k,a}^{\beta_a - 1}$$

$$p(\pi_n | \gamma) := \text{Dir}(\pi_n | \gamma \mathbf{1}_K) \propto \prod_{k=1}^K \pi_{n,k}^{\gamma_k - 1}$$

Full conditionals: 2. π

$$\begin{aligned}
 p(\pi_n | z_n, \phi) &\propto p(\pi_n | \gamma) \prod_{m=1}^{M_n} p(z_{n,m} = k | \pi_n) \\
 &\propto \prod_{k=1}^K \pi_{n,k}^{\gamma_k - 1} \prod_{m=1}^{M_n} \prod_{k=1}^K \pi_{n,k}^{\delta(z_{n,m}=k)}
 \end{aligned}$$

Q: Do you recognize this distribution?

Gibbs Sampling for LDA

$$p(x_{n,m} | z_{n,m} = k, \phi) := \text{Cat}(x_{n,m} | \phi_k) = \phi_{k,x_{n,m}}$$

$$p(z_{n,m} | \pi_n) := \text{Cat}(z_{n,m} | \pi_n) = \pi_{n,z_{n,m}}$$

$$p(\phi_k | \beta) := \text{Dir}(\phi_k | \beta \mathbf{1}_A) \propto \prod_{a=1}^A \phi_{k,a}^{\beta_a - 1}$$

$$p(\pi_n | \gamma) := \text{Dir}(\pi_n | \gamma \mathbf{1}_K) \propto \prod_{k=1}^K \pi_{n,k}^{\gamma_k - 1}$$

Full conditionals: 2. π

$$p(\pi_n | z_n, \phi) \propto p(\pi_n | \gamma) \prod_{m=1}^{M_n} p(z_{n,m} = k | \pi_n)$$

$$\propto \prod_{k=1}^K \pi_{n,k}^{\gamma_k - 1} \prod_{m=1}^{M_n} \prod_{k=1}^K \pi_{n,k}^{\delta(z_{n,m}=k)}$$

$$= \text{Dir}((\gamma_k + \sum_{m=1}^M \delta(z_{n,m} = k))_{k=1:K})$$

Gibbs Sampling for LDA

$$p(x_{n,m} | z_{n,m} = k, \phi) := \text{Cat}(x_{n,m} | \phi_k) = \phi_{k,x_{n,m}}$$

$$p(z_{n,m} | \pi_n) := \text{Cat}(z_{n,m} | \pi_n) = \pi_{n,z_{n,m}}$$

$$p(\phi_k | \beta) := \text{Dir}(\phi_k | \beta \mathbf{1}_A) \propto \prod_{a=1}^A \phi_{k,a}^{\beta_a - 1}$$

$$p(\pi_n | \gamma) := \text{Dir}(\pi_n | \gamma \mathbf{1}_K) \propto \prod_{k=1}^K \pi_{n,k}^{\gamma_k - 1}$$

Full conditionals: 3. ϕ

$$p(\phi_k | z, \pi) \propto \left(\prod_{n=1}^N \prod_{m=1}^{M_n} p(x_{n,m} = a | z_{n,m} = k, \phi_k) p(\phi_k | \beta) \right)_{a=1:A}$$

$$= \text{Dir} \left(\left(\beta_a + \sum_{n=1}^N \sum_{m=1}^{M_n} \delta(x_{n,m} = a, z_{n,m} = k) \right)_{a=1:A} \right)$$

Gibbs Sampling for LDA

- ▶ initialize randomly

$$\pi_n \sim \text{Dir}(\gamma \mathbf{1}_K), \quad \phi_k \sim \text{Dir}(\beta \mathbf{1}_A)$$

- ▶ sample iteratively:

$$z_{n,m} \sim \text{Cat}((\phi_{k,x_{n,m}} \pi_{n,k})_{k=1:K}), \quad \forall n \forall m$$

$$\pi_n \sim \text{Dir}((\gamma_k + \sum_{m=1}^M \delta(z_{n,m} = k))_{k=1:K}), \quad \forall n$$

$$\phi_k \sim \text{Dir}((\beta_a + \sum_{n=1}^N \sum_{m=1}^M \delta(x_{n,m} = a, z_{n,m} = k))_{a=1:A}), \quad \forall k$$

Outline

1. The LDA Model
2. Learning LDA via Gibbs Sampling
- 3. Learning LDA via Collapsed Gibbs Sampling**
4. Learning LDA via Variational Inference
5. Supervised LDA

Counts

$$c_{n,a,k} := \sum_{m=1}^{M_n} \delta(x_{n,m} = a, z_{n,m} = k)$$

$$c_{n,k} := \sum_{a=1}^A c_{n,a,k}$$

$$c_{a,k} := \sum_{n=1}^N c_{n,a,k}$$

$$c_k := \sum_{a=1}^A \sum_{n=1}^N c_{n,a,k}$$

$$c_n := \sum_{a=1}^A \sum_{k=1}^K c_{n,a,k} = M_n$$

Marginals over π

$$\begin{aligned}
 p(z_n | \gamma) &= \int \left(\prod_{m=1}^{M_n} \text{Cat}(z_{n,m} | \pi_n) \right) \text{Dir}(\pi_n | \gamma \mathbf{1}_K) d\pi_n \\
 &= \int \prod_{k=1}^K \pi_{n,k}^{c_{n,k}} \frac{\Gamma(K\gamma)}{\Gamma(\gamma)^K} \pi_{n,k}^\gamma d\pi_n \\
 &= \frac{\Gamma(K\gamma)}{\Gamma(\gamma)^K} \frac{\prod_{k=1}^K \Gamma(\gamma + c_{n,k})}{\Gamma(\sum_{k=1}^K \gamma + c_{n,k})} \int \underbrace{\frac{\Gamma(\sum_{k=1}^K \gamma + c_{n,k})}{\prod_{k=1}^K \Gamma(\gamma + c_{n,k})} \prod_{k=1}^K \pi_{n,k}^{\gamma + c_{n,k}}}_{=\text{Dir}(\pi_n | (\gamma + c_{n,k})_{k=1:K})} d\pi_n \\
 &= \frac{\Gamma(K\gamma)}{\Gamma(\gamma)^K} \frac{\prod_{k=1}^K \Gamma(c_{n,k} + \gamma)}{\Gamma(M_n + K\gamma)}
 \end{aligned}$$

Marginals over π and ϕ

$$\begin{aligned}
 p(z | \gamma) &= \prod_{n=1}^N \int \left(\prod_{m=1}^{M_n} \text{Cat}(z_{n,m} | \pi_n) \right) \text{Dir}(\pi_n | \gamma \mathbf{1}_K) d\pi_n \\
 &= \left(\frac{\Gamma(K\gamma)}{\Gamma(\gamma)^K} \right)^N \prod_{n=1}^N \frac{\prod_{k=1}^K \Gamma(c_{n,k} + \gamma)}{\Gamma(M_n + K\gamma)}
 \end{aligned}$$

Marginals over π and ϕ

$$\begin{aligned}
 p(z | \gamma) &= \prod_{n=1}^N \int \left(\prod_{m=1}^{M_n} \text{Cat}(z_{n,m} | \pi_n) \right) \text{Dir}(\pi_n | \gamma \mathbf{1}_K) d\pi_n \\
 &= \left(\frac{\Gamma(K\gamma)}{\Gamma(\gamma)^K} \right)^N \prod_{n=1}^N \frac{\prod_{k=1}^K \Gamma(c_{n,k} + \gamma)}{\Gamma(M_n + K\gamma)}
 \end{aligned}$$

$$\begin{aligned}
 p(x | z, \beta) &= \prod_{k=1}^K \int \left(\prod_{(n,m): z_{n,m}=k} \text{Cat}(x_{n,m} | \phi_k) \right) \text{Dir}(\phi_k | \beta \mathbf{1}_K) d\phi_k \\
 &= \left(\frac{\Gamma(A\beta)}{\Gamma(\beta)^A} \right)^K \prod_{k=1}^K \frac{\prod_{a=1}^A \Gamma(c_{a,k} + \beta)}{\Gamma(c_k + A\beta)}
 \end{aligned}$$

Conditional Probability for Single $z_{n,m}$

$$p(z \mid x, \beta, \gamma) \stackrel{\text{Bayes}}{=} \frac{p(x \mid z, \beta, \gamma) p(z \mid \beta, \gamma)}{p(x \mid \beta, \gamma)} \propto p(x \mid z, \beta) p(z \mid \gamma)$$

$$\begin{aligned} p(z \mid x, \beta, \gamma) &= p(z_{n,m} \mid z_{-(n,m)}, x, \beta, \gamma) p(z_{-(n,m)} \mid x, \beta, \gamma) \\ &= p(z_{n,m} \mid z_{-(n,m)}, x, \beta, \gamma) p(z_{-(n,m)} \mid x_{-(n,m)}, \beta, \gamma) \end{aligned}$$

\rightsquigarrow

$$\begin{aligned} p(z_{n,m} \mid z_{-(n,m)}, x, \beta, \gamma) &= \frac{p(z \mid x, \beta, \gamma)}{p(z_{-(n,m)} \mid x_{-(n,m)}, \beta, \gamma)} \\ &\propto \frac{p(x \mid z, \beta) p(z \mid \gamma)}{p(x_{-(n,m)} \mid z_{-(n,m)}, \beta) p(z_{-(n,m)} \mid \gamma)} \end{aligned}$$

Conditional Probability for Single $z_{n,m}$

$$p(z_{n,m} \mid z_{-(n,m)}, x, \beta, \gamma) \propto \frac{p(x \mid z, \beta) p(z \mid \gamma)}{p(x_{-(n,m)} \mid z_{-(n,m)}, \beta) p(z_{-(n,m)} \mid \gamma)}$$

Conditional Probability for Single $z_{n,m}$

$$p(z_{n,m} \mid z_{-(n,m)}, x, \beta, \gamma) \propto \frac{p(x \mid z, \beta) p(z \mid \gamma)}{p(x_{-(n,m)} \mid z_{-(n,m)}, \beta) p(z_{-(n,m)} \mid \gamma)}$$

Now let $c_{n,a,k}^-$ be the counts for the leave-one-out sample $x_{-(n,m)}, z_{-(n,m)}$ (all but m -th word of document n).

$$c_{n,a,k}^- = \begin{cases} c_{n,a,k} - 1, & \text{for } x_{n,m} = a, z_{n,m} = k \\ c_{n,a,k}, & \text{else} \end{cases}$$

- ▶ all terms other than for $x_{n,m} = a, z_{n,m} = k$ cancel out.
- ▶ terms for $x_{n,m} = a, z_{n,m} = k$ can be simplified via $\Gamma(x+1)/\Gamma(x) = x$

Conditional Probability for Single $z_{n,m}$

$$p(z_{n,m} \mid z_{-(n,m)}, x, \beta, \gamma) \propto \frac{p(x \mid z, \beta) p(z \mid \gamma)}{p(x_{-(n,m)} \mid z_{-(n,m)}, \beta) p(z_{-(n,m)} \mid \gamma)}$$

Now let $c_{n,a,k}^-$ be the counts for the leave-one-out sample $x_{-(n,m)}, z_{-(n,m)}$ (all but m -th word of document n).

$$c_{n,a,k}^- = \begin{cases} c_{n,a,k} - 1, & \text{for } x_{n,m} = a, z_{n,m} = k \\ c_{n,a,k}, & \text{else} \end{cases}$$

- ▶ all terms other than for $x_{n,m} = a, z_{n,m} = k$ cancel out.
- ▶ terms for $x_{n,m} = a, z_{n,m} = k$ can be simplified via $\Gamma(x+1)/\Gamma(x) = x$

$$p(z_{n,m} = k \mid z_{-(n,m)}, x, \beta, \gamma) \propto \frac{c_{x_{n,m},k}^- + \beta}{c_k^- + A\beta} \frac{c_{n,k}^- + \gamma}{M_n + K\gamma}$$

Collapsed LDA Implementation

- ▶ assign all $z_{n,m}$ randomly
- ▶ compute $c_{n,a,k}$
- ▶ for $s := 1, \dots, S$:
 - ▶ for $n := 1, \dots, N, \quad m := 1, \dots, M_n$:

$$c_{x_{n,m}, z_{n,m}} := c_{x_{n,m}, z_{n,m}} - 1$$

$$c_{n, z_{n,m}} := c_{n, z_{n,m}} - 1$$

$$c_{z_{n,m}} := c_{z_{n,m}} - 1$$

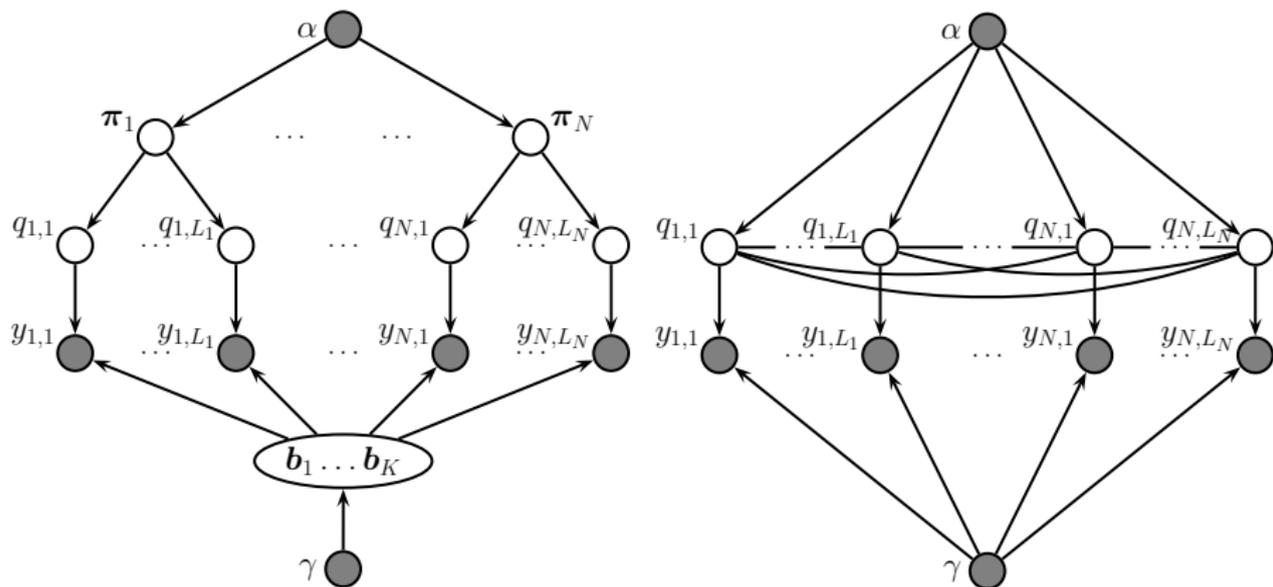
$$z_{n,m} \sim \text{Cat}\left(\left(\frac{c_{x_{n,m},k}^- + \beta}{c_k^- + A\beta} \frac{c_{n,k}^- + \gamma}{M_n + K\gamma}\right)_{k=1:K}\right)$$

$$c_{x_{n,m}, z_{n,m}} := c_{x_{n,m}, z_{n,m}} + 1$$

$$c_{n, z_{n,m}} := c_{n, z_{n,m}} + 1$$

$$c_{z_{n,m}} := c_{z_{n,m}} + 1$$

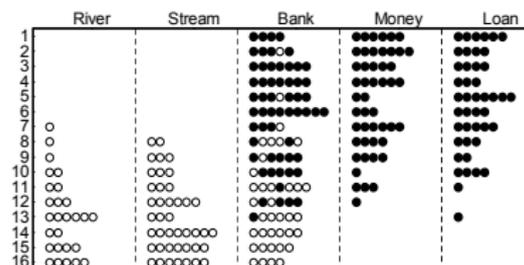
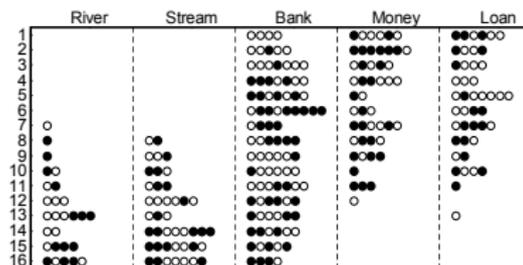
LDA vs Collapsed LDA



[Mur12, fig. 27.7]



Collapsed LDA / Example



$N = 16$ (rows), $A = 5$ (columns), $K = 2$ (colors)

[Mur12, fig. 27.8]

Outline

1. The LDA Model
2. Learning LDA via Gibbs Sampling
3. Learning LDA via Collapsed Gibbs Sampling
- 4. Learning LDA via Variational Inference**
5. Supervised LDA

Variational Inference via Mean Field Approximation

To solve the inference problem

$$\text{compute } p(x_1, \dots, x_N)$$

for intractable p , **approximate** p with a **fully factorized** density q

$$p(x_1, \dots, x_N) \approx q(x_1, \dots, x_N | \theta) := \prod_{n=1}^N q_n(x_n | \theta_n)$$

A good approximation should minimize the KL divergence of p and q :

$$(\theta_1, \dots, \theta_N) := \arg \min_{\theta_1, \dots, \theta_N} \text{KL}(q||p)$$

$$\text{KL}(q||p) := E_x \left(q(x) \log \frac{q(x)}{p(x)} \right)$$

which can be solved via coordinate descent:

$$\log q_n(x_n | \theta_n) = E_{x_{-n} \sim q_{-n}} (\tilde{p}(x_1, \dots, x_N)) + \text{const}$$

where \tilde{p} can be an unnormalized version of p .

Learning LDA via Mean Field Approximation

Mean field approximation

$$q(\pi_n | \tilde{\pi}_n) := \text{Dir}(\pi_n | \tilde{\pi}_n)$$

$$q(z_{n,m} | \tilde{z}_{n,m}) := \text{Cat}(z_{n,m} | \tilde{z}_{n,m})$$

in the E-step of EM leads to

E-step:

$$\tilde{z}_{n,m,k} = \phi_{x_{n,m},k} e^{\Psi(\tilde{\pi}_{n,k}) - \Psi(\sum_{k'} \tilde{\pi}_{n,k'})}$$

$$\tilde{\pi}_{n,k} = \gamma + \sum_m \tilde{z}_{n,m,k}$$

M-step:

$$\phi_{a,k} = \beta + \sum_n \sum_m \tilde{z}_{n,m,k} \delta(x_{n,m} = a)$$

Note: $E_{\pi_{n,k} \sim \text{Dir}(\tilde{\pi}_{n,k})}(\log \pi_{n,k}) = \Psi(\tilde{\pi}_{n,k}) - \Psi(\sum_{k'} \tilde{\pi}_{n,k'})$ with Ψ the digamma function.

Outline

1. The LDA Model
2. Learning LDA via Gibbs Sampling
3. Learning LDA via Collapsed Gibbs Sampling
4. Learning LDA via Variational Inference
5. Supervised LDA

Adding Further Information

- ▶ Add observed class information

$$y_n \in \mathcal{Y} := \{1, \dots, T\}, \quad n \in \{1, \dots, N\}$$

- ▶ goal now is either
 - ▶ to analyze x_n with an LDA model and predict targets y_n based on this analysis (supervised learning) or
 - ▶ to find topics that explain both, documents x_n and their classes y_n (unsupervised learning).
- ▶ Sometimes richer information is added, e.g., images.

Joint LDA and Logistic Regression

$$p(y_n | \pi_n, \theta) := \text{Cat}(y_n | \text{logistic}(\theta^T \pi_n))$$

$$p(\theta | \sigma^2) := \mathcal{N}(\theta | 0, \sigma^2)$$

$$p(x_{n,m} | z_{n,m} = k, \phi) := \text{Cat}(x_{n,m} | \phi_k), \quad n = 1, \dots, N, m = 1, \dots, M_n$$

$$p(z_{n,m} | \pi_n) := \text{Cat}(z_{n,m} | \pi_n), \quad n = 1, \dots, N, m = 1, \dots, M_n$$

$$p(\phi_k | \beta) := \text{Dir}(\phi_k | \beta \mathbf{1}_A), \quad k = 1, \dots, K$$

$$p(\pi_n | \gamma) := \text{Dir}(\pi_n | \gamma \mathbf{1}_K), \quad n = 1, \dots, N$$

Generative Supervised LDA

$$p(y_n | \bar{\pi}_n, \theta) := \text{Cat}(y_n | \text{logistic}(\theta^T \bar{\pi}_n)), \quad \bar{\pi}_{n,k} := \frac{1}{M_n} \sum_{m=1}^{M_n} \delta(z_{n,m} = k)$$

$$p(x_{n,m} | z_{n,m} = k, \phi) := \text{Cat}(x_{n,m} | \phi_k), \quad n = 1, \dots, N, m = 1, \dots, M_n$$

$$p(z_{n,m} | \pi_n) := \text{Cat}(z_{n,m} | \pi_n), \quad n = 1, \dots, N, m = 1, \dots, M_n$$

$$p(\phi_k | \beta) := \text{Dir}(\phi_k | \beta \mathbf{1}_A), \quad k = 1, \dots, K$$

$$p(\pi_n | \gamma) := \text{Dir}(\pi_n | \gamma \mathbf{1}_K), \quad n = 1, \dots, N$$

Discriminative Supervised LDA

$$\begin{aligned}
 p(x_{n,m} \mid z_{n,m} = k, \phi) &:= \text{Cat}(x_{n,m} \mid \phi_k), \quad n = 1, \dots, N, m = 1, \dots, M_n \\
 p(z_{n,m} \mid \pi_n, \mathbf{y}_n = \mathbf{t}) &:= \text{Cat}(z_{n,m} \mid \mathbf{A}_t \pi_n), \quad n = 1, \dots, N, m = 1, \dots, M_n \\
 p(\phi_k \mid \beta) &:= \text{Dir}(\phi_k \mid \beta \mathbf{1}_A), \quad k = 1, \dots, K \\
 p(\pi_n \mid \gamma) &:= \text{Dir}(\pi_n \mid \gamma \mathbf{1}_K), \quad n = 1, \dots, N
 \end{aligned}$$

- ▶ $\mathbf{A}_t \in \mathbb{R}^{K \times K}$ stochastic ($t = 1, \dots, T$)

Summary

- ▶ **Latent Dirichlet Allocation (LDA)** solves a clustering problem for sequence data (really: histograms)
 - ▶ clusters are called **topics**.
 - ▶ topics are described by **word/symbol probabilities**.
 - ▶ documents/sequences by **topic probabilities**.
 - ▶ a **latent variable “word topic”** for each word/element of each sequence.
 - ▶ semantically: disambiguation of the word (w.r.t. its topic)
- ▶ LDA can be learned via **Gibbs sampling**:
 - ▶ re-sample single variables from their full conditionals on all others in a round-robin fashion.
 - ▶ leads to sampling from categorical and Dirichlet distributions.
- ▶ LDA can be learned via **collapsed Gibbs sampling**:
 - ▶ integrate out word and document probabilities, leaving just the latent word topics.
 - ▶ leads to a way faster sampling from a categorical distribution only.

Summary (2/2)

- ▶ LDA can be learned via **variational inference** using **mean field approximation**.
 - ▶ approximate a distribution by a fully factorized distribution.
 - ▶ here: the distribution of the latent word topics and topic probabilities in the E-step of an EM algorithm for LDA.
 - ▶ leads to closed-form reestimation formulas.
- ▶ LDA can be extended different ways to take **document labels/classes** into account.
 - ▶ yielding a model for text classification.
 - ▶ joint LDA and logistic regression, generative supervised LDA, discriminative supervised LDA

Further Readings

- ▶ LDA:
 - ▶ [Mur12, chapter 27.3],
- ▶ Supervised LDA and other extensions:
 - ▶ [Mur12, chapter 27.4],

References



Kevin P. Murphy.

Machine learning: a probabilistic perspective.

The MIT Press, 2012.