

Machine Learning 2 8. Latent Dirichlet Allocation (LDA)

Lars Schmidt-Thieme

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

《日》《聞》《臣》《臣》 된言 '오오오

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□

Syllabus



		A. Advanced Supervised Learning		
Tue. 9.12.	(1)	A.1 Generalized Linear Models		
Wed. 10.12.	(2)	A.2 Gaussian Processes		
Tue. 16.12.	(3)	A.3 Advanced Support Vector Machines		
Wed. 17.12.	(4)	A.4 Neural Networks		
Tue. 6.1.	(5)	A.5 Ensembles		
Wed. 7.1.	Wed. 7.1. (6) A.5b Ensembles (ctd.)			
Tue. 13.1.	ie. 13.1. (7) A.6 Sparse Linear Models — L1 regularization			
Wed. 14.1.	(8)	A.6b Sparse Linear Models — L1 regularization (ctd.)		
Tue. 20.1.	(9)	A.7. Sparse Linear Models — Further Methods		
		B. Complex Predictors		
Wed. 21.1.	(10)	B.1 Latent Dirichlet Allocation (LDA)		
Tue. 27.1.	(11)	B.1b Latent Dirichlet Allocation (LDA; ctd.)		
Wed. 28.1.	(12)	_		
Tue. 3.2.	(13)			
Wed. 4.2.	(14)			

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

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

・ロト・四ト・王ト・王ト 別に ろくの

Documents / Finite Discrete Sequences

Juniversiter Hildeshein

- ► instances $x_n \in A^*$ are **discrete sequences**
 - A := {1,..., A} called dictionary / alphabet (A ∈ N), where a ∈ A denotes the a-th word / symbol / token.
 - $\mathcal{A}^* := \bigcup_{\ell=1}^{L} \mathcal{A}^{\ell}$ 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):

$$\Phi_{n,a} := |\{\ell \in \{1, \dots, |x_n|\} \mid x_{n,\ell} = a\}$$

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

◆□▶ ◆□▶ ★∃▶ ★∃▶ ★目★ 少々で

The LDA Model



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

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

$$p(\phi_k | \beta) := Dir(\phi_k | \beta 1_A), \quad k = 1, ..., K$$

$$p(\pi_n | \gamma) := Dir(\pi_n | \gamma 1_K), \quad n = 1, ..., N$$

- ► $z_{n,m} \in \{1, ..., K\}$: 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 .

《日》《國》《王》《王》 法正 ろくの

The LDA Model

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

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

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

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

- ► $z_{n,m} \in \{1, ..., K\}$: 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 .

[Mur12, fig. 27.2] 《마》《문》《문》 문》 분들 것으(~

Topic 77

Example



word

MUSIC

MUSICAL

prob.

.090

.013

Topic 166

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

POEM	.034	DANCE
POETRY	.033	SONG
POET	.030	PLAY
PLAYS	.026	SING
POEMS	.026	SINGING
PLAY	.026	BAND
LITERARY	.023	PLAYED
WRITERS	.022	SANG
DRAMA	.021	SONGS
WROTE	.020	DANCING
POETS	.017	PIANO
WRITER	.016	PLAYING
SHAKESPEARE	.015	RHYTHM
WRITTEN	.013	ALBERT

Topic 82

word pro

LITERATURE

[Mur12, fig. 27.4] Sac

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

STAGE

Example



Document #29795

Bix beiderbecke, at age⁶⁶⁰ 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 eart¹⁹. It was jazz⁰⁷⁷. Bix beiderbecke had already had music⁰⁷⁷ lessons⁰⁷⁷. He showed⁰⁰² promise¹³⁴ on the piano⁰⁷⁷, and his parents³³⁵ hoped²⁶⁸ he might consider¹¹⁸ becoming a concert⁰⁷⁷ pianist⁰⁷⁷. But bix was interested²⁶⁸ in another kind⁰⁵⁰ of music⁰⁷⁷. He wanted²⁶⁸ to play⁰⁷⁷ he cornet. And he wanted²⁶⁸ to play⁰⁷⁷.

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 for two. The 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 jlay¹⁶⁶ the game¹⁶⁶. They play¹⁶⁶ the game¹⁶⁶.

[Mur12, fig. 27.5]

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

・ロト ・部ト ・王ト ・王ト シスペ



Learning via Parameter Sampling The loglikelihood

$$\mathit{p}(heta \mid \mathcal{D}) \propto \mathit{p}(\mathcal{D} \mid heta)$$

describes the **distribution of the parameters given the data**. If we can **sample parameters** from this distribution

$$\theta_1, \theta_2, \dots, \theta_S \sim p(\theta \mid D)$$

we can

estimate expected parameter values and their variances from this parameter sample:

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

predict targets for new instances x via model averaging:

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

Sampling



- ▶ for most closed-form distributions p(x) there exist efficient sampling methods
 - categorical, normal, ...
- but most loglikelihoods are not closed-form distributions.
 - but for example products thereof.

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□

Gibbs Sampling

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

◆□▶ ◆□▶ ★∃▶ ★∃▶ ★目★ 少々で



Gibbs Sampling

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

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

$$\begin{aligned} x_1^{s+1} &\sim p(x_1 \mid x_{2:N} = x_{2:N}^s) \\ x_2^{s+1} &\sim p(x_2 \mid x_{1:1} = x_{1:1}^{s+1}, x_{3:N} = x_{3:N}^s) \\ &\vdots \\ x_n^{s+1} &\sim p(x_n \mid x_{1:n-1} = x_{1:n-1}^{s+1}, x_{n+1:N} = x_{n+1:N}^s) \\ &\vdots \\ x_N^{s+1} &\sim p(x_N \mid x_{1:N-1} = x_{1:N-1}^{s+1}) \end{aligned}$$



Gibbs Sampling



- ► the distribution created by the Gibbs sampler eventually will converge to p(x₁,...,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**; 100-10,000).
 - only after a while the chain has converged to the stationary distribution $p(x_1, \ldots, x_N)$.
 - ► 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

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

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

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

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

Full conditionals:

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





Gibbs Sampling for LDA

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

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

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



Full conditionals:

$$p(\phi_k \mid z, \pi) \propto \prod_{n=1}^N \prod_{m=1}^{M_n} p(z_{n,m} = k \mid \pi_n) p(\phi_k \mid \beta)$$
$$= \operatorname{Dir}((\beta_a + \sum_{n=1}^N \sum_{m=1}^M \delta(x_{n,m} = a, z_{n,m} = k))_{a=1:A})$$



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

《日》《聞》《臣》《臣》 王曰 '오�?'

Counts



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

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

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

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

シック 비門 《파》《파》《西》《日》



Marginals over π and ϕ

$$p(z \mid \beta) = \prod_{n=1}^{M_n} \int (\prod_{m=1}^{M_n} \operatorname{Cat}(z_{n,m} \mid \pi_n)) \operatorname{Dir}(\pi_n \mid \gamma \mathbf{1}_K) d\pi_n$$
$$= \left(\frac{\Gamma(K\gamma)}{\Gamma(\gamma)^K}\right)^N \prod_{n=1}^N \frac{\prod_{k=1}^K \Gamma(c_{n,k} + \gamma)}{\Gamma(M_n + K\gamma)}$$

シック 正則 スポッスポッス モッ



Marginals over π and ϕ

$$p(z \mid \beta) = \prod_{n=1}^{M_n} \int (\prod_{m=1}^{M_n} \operatorname{Cat}(z_{n,m} \mid \pi_n)) \operatorname{Dir}(\pi_n \mid \gamma \mathbf{1}_K) d\pi_n$$
$$= \left(\frac{\Gamma(K\gamma)}{\Gamma(\gamma)^K}\right)^N \prod_{n=1}^N \frac{\prod_{k=1}^K \Gamma(c_{n,k} + \gamma)}{\Gamma(M_n + K\gamma)}$$

$$p(x \mid z, \beta) = \prod_{k=1}^{K} \int (\prod_{(n,m):z_{n,m}=k} \operatorname{Cat}(x_{n,m} \mid \phi_{k})) \operatorname{Dir}(\phi_{k} \mid \beta 1_{K}) d\phi_{k}$$
$$= \left(\frac{\Gamma(A\beta)}{\Gamma(\beta)^{A}}\right)^{K} \prod_{k=1}^{K} \frac{\prod_{a=1}^{A} \beta(c_{a,k} + \beta)}{\Gamma(c_{k} + A\beta)}$$

・ロト < 団ト < 三ト < 三ト < 三ト < ロト

Shiversiter Hildeshein

Marginals over π and ϕ

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

◇▷▷ 비로 《토》《토》 《팀》 ◇□ >

Marginals over π and ϕ



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

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 $\Gamma(x+1)/\Gamma(x) = x$

(今夕) 비로 《王》《王》《臣》《曰》

Marginals over π and ϕ



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

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 $\Gamma(x+1)/\Gamma(x) = x$

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

Collapsed LDA Implementation

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

 $\begin{aligned} c_{x_{n,m},z_{n,m}} &:= c_{x_{n,m},z_{n,m}} - 1 \\ c_{n,z_{n,m}} &:= c_{n,z_{n,m}} - 1 \\ c_{z_{n,m}} &:= c_{z_{n,m}} - 1 \\ z_{n,m} &:= k \sim \frac{c_{x_{n,m},k}^{-} + \gamma}{c_{k}^{-} + A\gamma} \frac{c_{n,k}^{-} + \beta}{M_{n} + K\beta} \\ c_{x_{n,m},z_{n,m}} &:= c_{x_{n,m},z_{n,m}} + 1 \\ c_{n,z_{n,m}} &:= c_{n,z_{n,m}} + 1 \\ c_{z_{n,m}} &:= c_{z_{n,m}} + 1 \end{aligned}$





LDA vs Collapsed LDA





Collapsed LDA / Example



	River	Stream	Bank	Money	Loan
1		1	0000	000000	000000
2		1	00000	0000000	0000
3			00000000	00000	0000
4		1	0000000	000000	000
5			0000000	e 0	0000000
6		1	000000000	080	0000
7	0		0000	660060	00000
8	•	0	000000	0000	660
9		000	000000	0000	0
10			000000	•	0000
11		000	00000000		•
12	000	0000000	000000	0	1
13	0000000	000	000000		0
14	00	00000000	000000	!	1
15	0000	6000000	00000		i
16	00000	000000	0000		1

River	Stream	Bank	Money	Loan
1 2 3 4 5 6 7 8 0 8 0 0 10 0 0 11 0 0 12 0 0 12 0 0 12 0 0 12 0 0 12 0 0 12 0 0 12 0 0 12 0 0 12 0 0 12 0 0 12 0 0 0 0 12 0 0 0 12 0 0 0 0 0 0 0 12 0 0 0 0 0 0 0 0 0 0 0 0 0	00 000 000 000 000000 0000000 0000000 0000	Contraction Contracti	000000 000000 000000 000 000 000 0000 0000	

[Mur12, fig. 27.8] イロトイクトイミトイミト ミニ クヘベ

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

Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany

◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□ ● ◆○



Variational Inference via Mean Field Approximation To solve the inference problem

compute
$$p(x_1, \ldots, x_N)$$

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

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

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

which can be solved via coordinate descent:

$$\log q_n(x_n \mid \theta_n) = E_{x_{-1} \sim q_{-n}}(\tilde{p}(x_1, \ldots, x_N)) + \text{const}$$

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



Learning LDA via Mean Field Approximation Mean field approximation

$$q(\pi_n \mid \tilde{\pi}_n) := \mathsf{Dir}(\pi_n \mid \tilde{\pi}_n)$$

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

in the E-step of EM leads to

E-step:

$$\tilde{z}_{n,m,k} = \Psi(\tilde{\pi}_{n,k}) - \Psi(\sum_{k'} \tilde{\pi}_{n,k'})$$
$$\tilde{\pi}_{n,k} = \gamma + \sum_{m} \tilde{z}_{n,m,k}$$

M-step:

$$\hat{\phi}_{a,k} = \beta + \sum_{n} \sum_{m} \tilde{z}_{n,m,k} \delta(x_{n,m} = a)$$

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 $\Psi_{\alpha} = \bigoplus_{k \in \mathcal{D}} \bigoplus_$

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

ィロト イラト イラト イラト ショー シークへで Lars Schmidt-Thieme, Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, Germany Machine Learning 2 5. Supervised LDA

. . .



Further Readings



- ► L1 regularization: [Mur12, chapter 13.3–5], [HTFF05, chapter 3.4, 3.8, 4.4.4], [Bis06, chapter 3.1.4].
 - ► LAR, LARS: [HTFF05, chapter 3.4.4], [Mur12, chapter 13.4.2],
- ► Non-convex regularizers: [Mur12, chapter 13.6].
- ► Automatic Relevance Determination (ARD): [Mur12, chapter 13.7], [HTFF05, chapter 11.9.1], [Bis06, chapter 7.2.2].
- ► Sparse Coding: [Mur12, chapter 13.8].

- 《日》 《日》 《日》 《日》 《日》

References





Christopher M. Bishop.

Pattern recognition and machine learning, volume 1. springer New York, 2006.



Trevor Hastie, Robert Tibshirani, Jerome Friedman, and James Franklin.

The elements of statistical learning: data mining, inference and prediction, volume 27. Springer, 2005.



Kevin P. Murphy.

Machine learning: a probabilistic perspective. The MIT Press, 2012.