# A generic approach to topic models

Gregor Heinrich

CTO, vsonix GmbH, Darmstadt
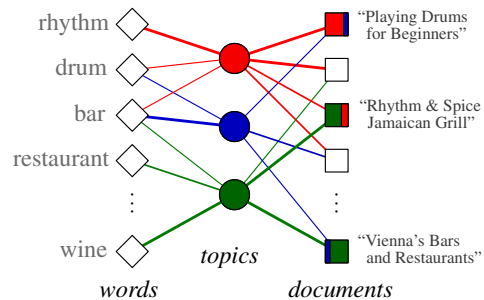
`http://vsonix.com`

`http://arbylon.net`

Research Seminar, Institute for Statistics and Mathematics
Vienna University of Economics and Business
`http://www.wu.ac.at/statmath/resseminar`

Vienna, 1 June 2012

---

## Overview

- Topic models – motivation and review
- Networks of mixed membership (NoMMs)
- Inference – a Gibbs "meta-sampler"
- NoMM typology and design
- Application to tag-enhanced expertise finding
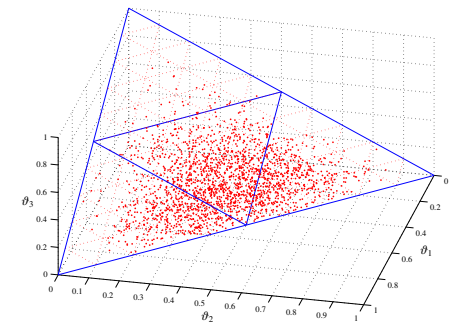- Conclusions and outlook

---

## Topic models



- Probabilistic representations of grouped discrete data
- Illustrative for text: words grouped in documents
  - Latent topics (a.k.a. concepts, components) = cluster semantically related words (Landauer and Dumais 1997; Griffiths et al. 2007)
  - Language = semantic meaning (topics) + noise
- → Reduce vocabulary problem by discovery of semantic relations
- → Reduce sparsity problem by dimensionality reduction ↔ discrete principal components analysis (Buntine and Jakulin 2005)
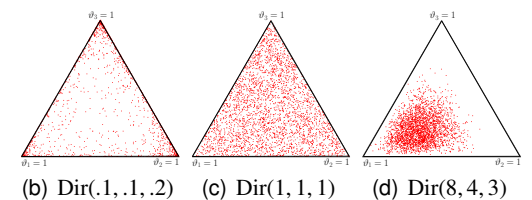
---

## Towards Bayesian topic models: the Dirichlet distribution

Bayesian methodology:

- Parameters generated from *prior* distributions
- Language data: popular prior for the multinomial / discrete distribution: Dirichlet distribution
  - Conjugacy: straight-forward mathematical form
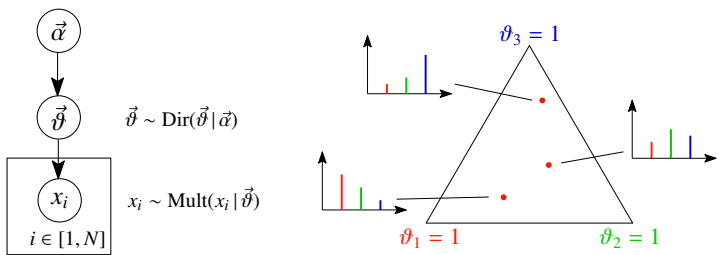- Bayesian topic model: Latent Dirichlet allocation (Blei et al. 2003)



(a) $\mathrm{Dir}(4, 4, 2)$

(b) $\mathrm{Dir}(.1, .1, .2)$   (c) $\mathrm{Dir}(1, 1, 1)$   (d) $\mathrm{Dir}(8, 4, 3)$

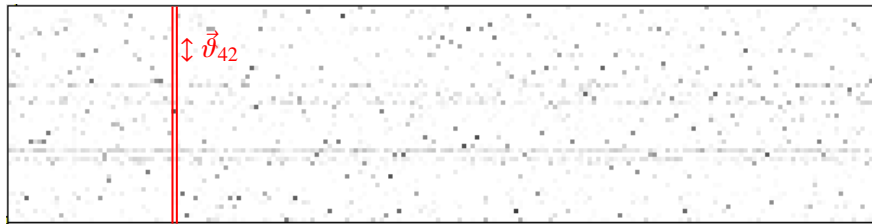## Bayesian networks: Dirichlet-generated multinomials



Bayesian networks:

- Graphical modelling of joint probability distributions
- Node: random variable
- Edge: conditional probability distribution
- Plate: repeated i.i.d. samples

## Latent Dirichlet Allocation



Draw word from term distribution of topic 2, "learning"

## Example document–topic distributions

Document $m = 42$ (column): Traditional machine learning relies on the availability of a large amount of data to train a model, which is then applied to test data in the same feature space. However, labeled data are often scarce and expensive to obtain...
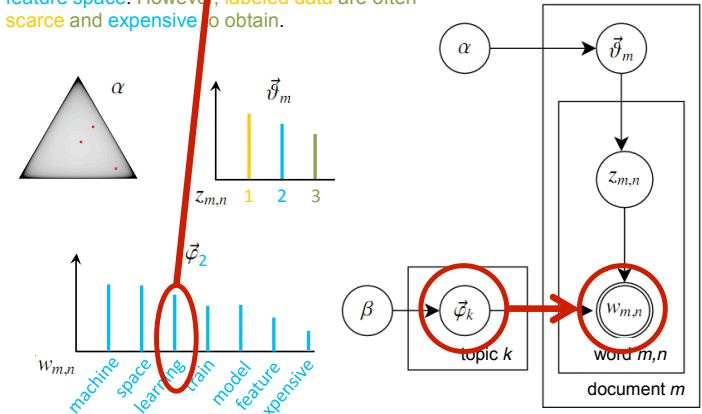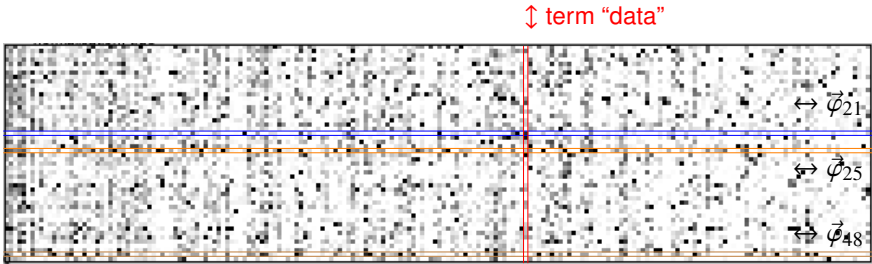
Strongest topics: $k = \{25, 21, 48, \ldots\}$



transposed view: rows = topics, columns = documents

Figure: Excerpt from document–topic matrix $\vartheta$ ($M = 200, K = 50$).

## Example topic–term distributions

Topic $k = 21$ (row): data word feature label data scarce obtain...
Topic $k = 25$ (row): machine learning train model test feature space...
Topic $k = 48$ (row): computing support grant project system method...

$\updownarrow$ term "data"
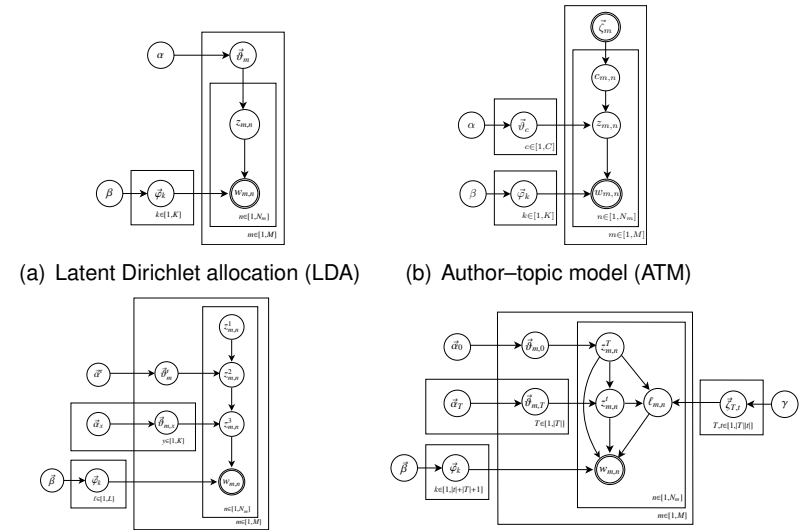


rows = topics, columns = terms

Figure: Excerpt from topic–term matrix $\varphi$ ($V = 200, K = 50$).

## Example: Text mining for semantic clusters

| Topic label | Most likely terms according to $\varphi_{k,t} = p(\text{word}|\text{topic})$ |
|---|---|
| Politische Parteien | CDU Partei Kohl Aufklärung Schäuble Zeitung Union Krise Wahrheit Affäre Christdemokraten Glaubwürdigkeit Konsequenzen |
| Bundesliga | FC SC München Borussia SV VfL Kickers SpVgg Uhr Köln Bochum Freiburg VfB Eintracht Bayern Hamburger Bayern+München |
| Polizei / Unfall | Polizei verletzt schwer Auto Unfall Fahrer Angaben schwer+verletzt Menschen Wagen Verletzungen Lawine Mann vier Meter Straße |
| Tschetschenien | Rebellen russischen Grosny russische Tschetschenien Truppen Kaukasus Moskau Angaben Interfax tschetschenischen Agentur |
| Politik / Hessen | FDP Koch Hessen CDU Koalition Gerhardt Wagner Liberalen hessischen Westerwelle Wolfgang Roland+Koch Wolfgang+Gerhardt |
| Wetter | Grad Temperaturen Regen Schnee Süden Sonne Wetter Wolken Deutschland zwischen Nacht Wetterdienst Wind |
| Politik / Kroatien | Parlament Partei Stimmen Mehrheit Wahlen Wahl Opposition Kroatien Präsident Parlamentswahlen Mesic Abstimmung HDZ |
| Die Grünen | Grünen Parteitag Atomausstieg Trittin Grüne Partei Trennung Mandat Ausstieg Amt Roestel Jahren Müller Radcke Koalition |
| Russische Politik | Russland Putin Moskau russischen russische Jelzin Wladimir Tschetschenien Russlands Wladimir+Putin Kreml Boris Präsidenten |
| Polizei / Schulen | Polizei Schulen Schüler Täter Polizisten Schule Tat Lehrer erschossen Beamten Mann Polizist Beamte verletzt Waffe |

Bigram LDA topics, 18400 German news messages, Jan. 2000 (Heinrich et al. 2005)

## Topic models: Example structures



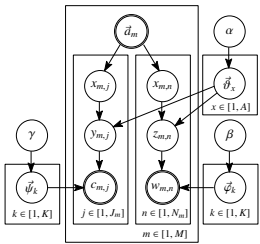(a) Latent Dirichlet allocation (LDA)     (b) Author–topic model (ATM)

(c) Pachinko allocation model (PAM4)     (d) Hierarchical PAM (hPAM)

(Blei et al. 2003; Rosen-Zvi et al. 2004; Li and McCallum 2006; Li et al. 2007)

## Typical derivation method (Is it really that complex?)



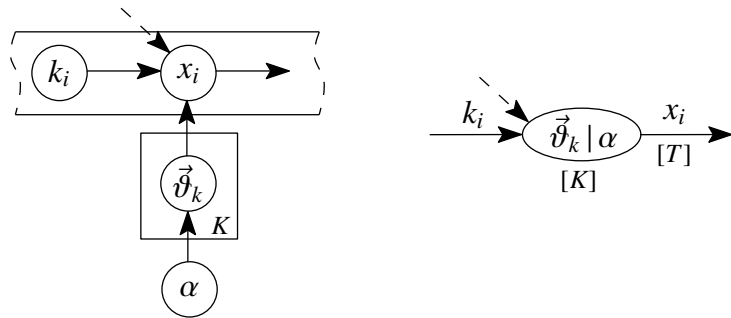(e) Expert–tag–topic model (ETT)

(Heinrich 2011)

## Topic models – bottom line

- Expanding research field with practical relevance
- No existing analysis as generic model class

$\rightarrow$ Conjecture:
  - Important properties generic across models
  - Simplifications in the derivation of model properties, inference algorithms and design methods
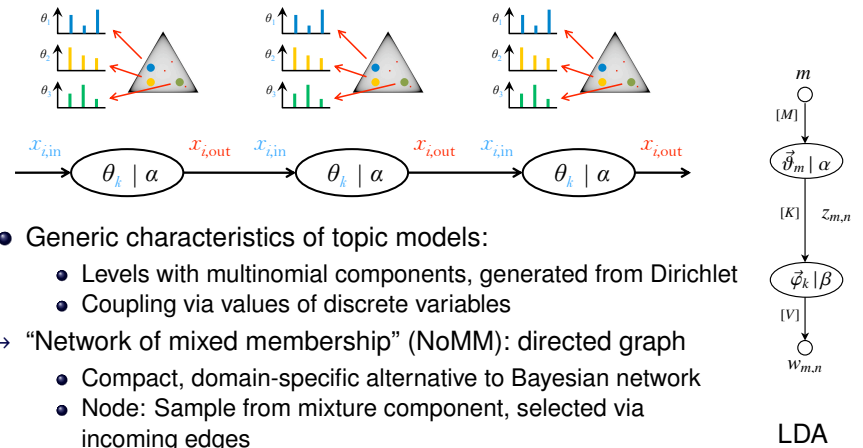
## Overview

- Topic models – motivation and review
- Networks of mixed membership (NoMMs)
- Inference – a Gibbs "meta-sampler"
- NoMM typology and design
- Application to tag-enhanced expertise finding
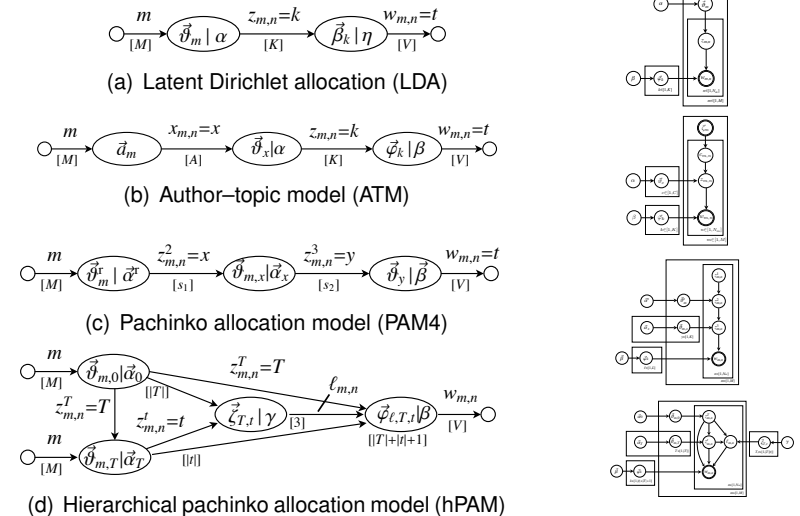- Conclusions and outlook

## Generic topic models – "NoMMs"



- Generic characteristics of topic models:
  - Levels with multinomial components, generated from Dirichlet
  - Coupling via values of discrete variables
- → "Network of mixed membership" (NoMM): directed graph
  - Compact, domain-specific alternative to Bayesian network
  - Node: Sample from mixture component, selected via incoming edges
  - Terminal node: observation
  - Edge: Propagation of discrete values to children

LDA

## NoMM level notation



parameters + hyperparameters ⇔ nodes

variables ⇔ edges

plates ⇔ indices + dimensions

## Topic models in NoMM representation



(a) Latent Dirichlet allocation (LDA)

(b) Author–topic model (ATM)

(c) Pachinko allocation model (PAM4)

(d) Hierarchical pachinko allocation model (hPAM)
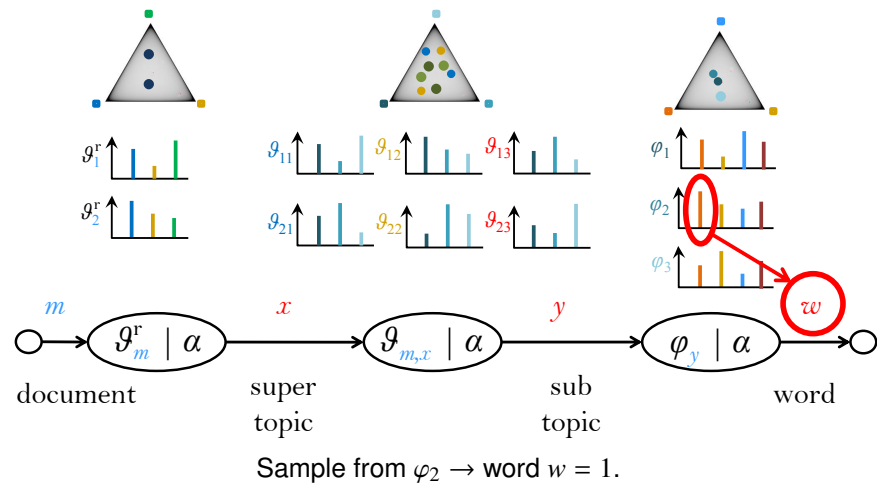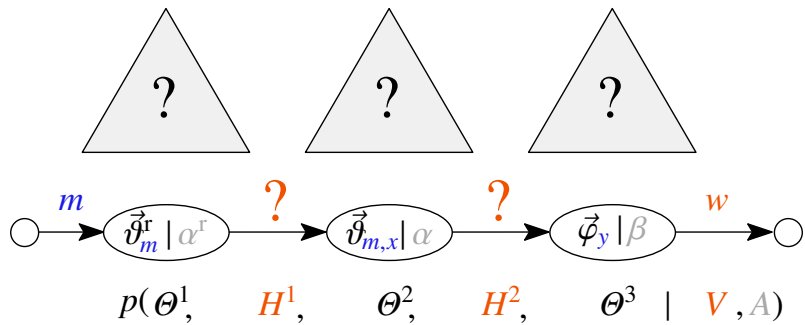
(Blei et al. 2003; Rosen-Zvi et al. 2004; Li and McCallum 2006; Li et al. 2007)

## Example NoMM generative process: PAM4



Sample from $\varphi_2 \to$ word $w = 1$.

## Overview

- Topic models – motivation and review
- Networks of mixed membership (NoMMs)
- Inference – a Gibbs "meta-sampler"
- NoMM typology and design
- Application to tag-enhanced expertise finding
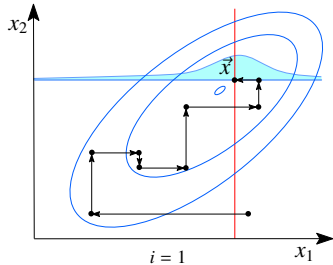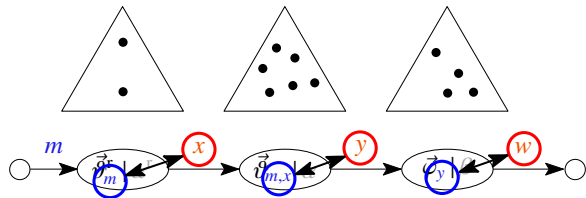- Conclusions and outlook

## Bayesian inference problem

- Bayesian inference: "Reverse generative process"
- Estimate (distributions over) parameters $\Theta$ and latent variables ("topics") $H$ given observations $V$ and hyperparameters $A$.
- $\to$ Find posterior distribution $p(H, \Theta \,|\, V, A) \to$ exponential complexity!



$$p(\Theta^1, \quad H^1, \quad \Theta^2, \quad H^2, \quad \Theta^3 \quad | \quad V, A)$$

## Collapsed Gibbs sampling



- Collapsed Gibbs sampling: stochastic EM / MCMC:
  - NoMMs: parameters $\Theta$ correlated with $H \to$ integrated out
  - For each data token $i$: Sample latent variables $H_i = (y_i, z_i, \dots)$, given all other data, latent $H_{\neg i}$ and visible $V$:

$$H_i \sim p(H_i \,|\, H_{\neg i}, V, A). \qquad (1)$$

  - Stationary state: full conditional simulates posterior
- Faster absolute convergence for NoMMs than, e.g., variational Bayes (Heinrich and Goesele 2009)
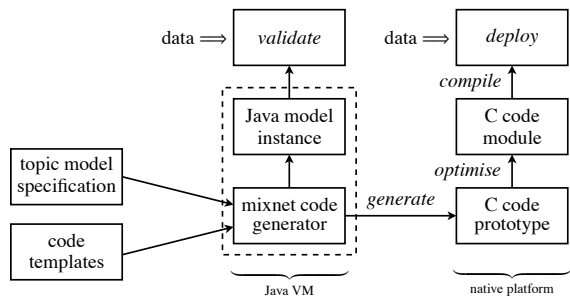
# Collapsed Gibbs full conditionals



- NoMM full conditionals can be generically derived (Heinrich 2009)
- Typical case leads to weights with straight-forward factor structure:

$$p(H_i \mid H_{\neg i}, V, A) \propto \prod_\ell \left[ \frac{n_{k,t}^{\neg i} + \alpha}{n_k^{\neg i} + T\alpha} \right]^{[\ell]} . \tag{2}$$

- $n_{k,t}$ = count of co-occurrences between input and output values of a NoMM level $\ell$
- More generally: $p(H_i \mid \cdot) \propto \prod_\ell [q(k,t)]^{[\ell]}$ with $t$ = *set* of values/edges
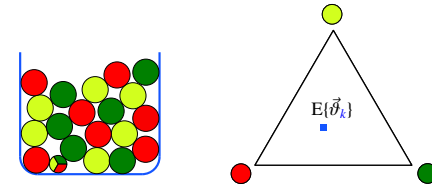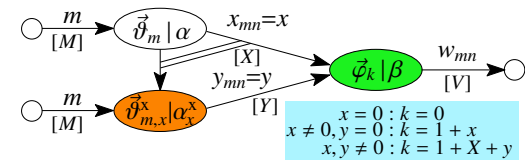
# $q$-functions and Pólya urn



Figure: Pólya urn and multinomial parameters.

$$q(k,t) \triangleq \frac{B(\vec{n}_k + \alpha)}{B(\vec{n}_k^{\neg i} + \alpha)} \stackrel{|t|=1}{=} \frac{n_{k,t}^{\neg t_i} + \alpha}{n_k^{\neg t_i} + T\alpha} = \text{smoothed ratio of co-occurrence counts}$$

$$\stackrel{t=\{u,v\}}{=} \frac{n_{k,u}^{\neg u_i} + \alpha}{n_k^{\neg u_i} + T\alpha} \cdot \frac{n_{k,v}^{\neg v_i} + \alpha + \delta(u-v)}{n_k^{\neg v_i} + T\alpha + 1} \triangleq q(k, u \oplus v)$$

$$\dots$$

# Implementation: Gibbs "meta-sampler"



- Code generator for topic models in Java and C
- Separation of knowledge domains: topic model applications vs. machine learning vs. computing architecture

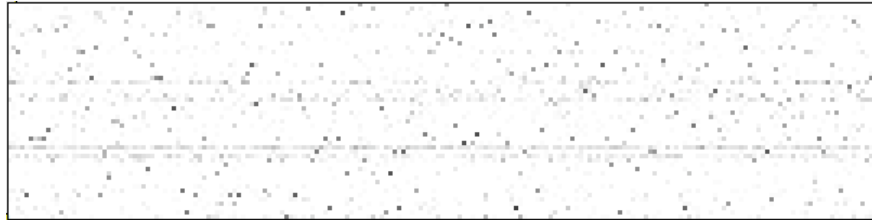# Example NoMM script and generated kernel: hPAM2



```
model = HPAM2

description:
    Hierarchical PAM model 2 (HPAM2)

sequences:
    # variables sampled for each (m,n)
    w, x, y : m, n

network:
    # each line one NoMM node
    m    >>  theta   | alpha        >>  x
    m,x  >>  thetax  | alphax[x]    >>  y
    x,y  >>  phi[k]                 >>  w
    # java code to assign k
    k : {
        if (x==0) { k = 0; }
        else if (y==0) k = 1 + x;
        else k = 1 + X + y;
    }.
```

```
// hidden edge
for (hx = 0; hx < X; hx++) {
    // hidden edge
    for (hy = 0; hy < Y; hy++) {
        mxsel = X * m + hx;
        mxjsel = hx;
        if (hx == 0)
            ksel = 0;
        else if (hy == 0)
            ksel = 1 + hx;
        else
            ksel = 1 + X + hy;
        pp[hx][hy] = (nmx[m][hx] + alpha[hx])
            * (nmxy[mxsel][hy] + alphax[mxjsel][hy])
            / (nmxysum[mxsel] + alphaxsum[mxjsel])
            * (nkw[ksel][w[m][n]] + beta)
            / (nkwsum[ksel] + betasum);
        psum += pp[hx][hy];
    } // for h
} // for h
```

## Example document–topic distributions



$t = 500$, converged

Figure: Excerpt from document–topic matrix $\vartheta$ ($M = 200, K = 50$).

---

## Fast sampling: hybrid acceleration methods

Serial:

- Exploit saliency of few weights, e.g., generalising (Porteous et al. 2008): compute only few weights on average + estimate normalisation term
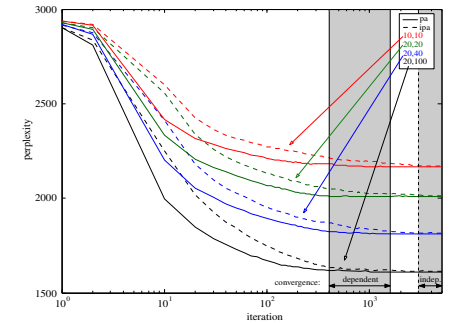- Complex data structures, especially for larger models

Parallel:

- Distribute local parameters (document-specific etc.)
- Need to sync global parameters: different methods, e.g., generalising (Newman et al. 2009)
- Occupancy: balance communication and computation (architecture-spec.)

Independence assumption:

- Reduce complexity: $\prod_\ell T^\ell \gg \sum_\ell T^\ell$

| method | model | parameters | speedup (iter., converge) | |
|--------|-------|------------|---------|---------|
| S×P4 | LDA | $K = 100$ | 6.3 | |
| S×P4 | LDA | $K = 500$ | 30.2 | |
| I | PAM4 | $K, L = 40, 40$ | 21.8 | 7.4 |
| P4×I | PAM4 | $K, L = 40, 40$ | 78.7 | 24.1 |
| S×P4×I | PAM4 | $K, L = 40, 40$ | 163.2 | 49.8 |
| S×P4×I | PAM4 | $K, L = 20, 100$ | 143.6 | 43.5 |



$\to$ Extend code generation to more complex implementations

---

## Overview

- Topic models – motivation and review
- Networks of mixed membership (NoMMs)
- Inference – a Gibbs "meta-sampler"
- NoMM typology and design
- Application to tag-enhanced expertise finding
- Conclusions and outlook

---

## $q$-functions and Pólya urn revisited



Figure: Pólya urn and multinomial parameters.

$$q(k, t) \triangleq \frac{\mathrm{B}(\vec{n}_k + \alpha)}{\mathrm{B}(\vec{n}_k^{\neg i} + \alpha)} \overset{|t|=1}{=} \frac{n_{k,t}^{\neg t_i} + \alpha}{n_k^{\neg t_i} + T\alpha} = \text{smoothed ratio of co-occurrence counts}$$

$$\overset{t=\{u,v\}}{=} \frac{n_{k,u}^{\neg u_i} + \alpha}{n_k^{\neg u_i} + T\alpha} \cdot \frac{n_{k,v}^{\neg v_i} + \alpha + \delta(u - v)}{n_k^{\neg v_i} + T\alpha + 1} \triangleq q(k, u \oplus v)$$

$\cdots$

## NoMM sub-structure typology

**N1. Dirichlet–multinomial parameters**

$$q(a, z)\, q(z, b)$$

**E2. Autonomous edges**

$$q(a, x \oplus y)\, q(x, b)\, q(y, c)$$

**C2. Combined indices**

$$q(a, x)\, q(b, y)\, q(k, c),\ k = f(x, y)$$

**N2. Observed parameters**

$$\vartheta^c_{a,z}\, q(z, b)$$

**E3. Coupled edges**

$$q(a, z)\, q(z, b)\, q(z, c)$$

**C3. Interleaved indices**

$$\approx q(a, z^1)\, q(b, z^2)\, q(z^1, c \oplus \tilde{c})\, q(z^2, \tilde{c} \oplus c)$$

Gibbs full conditional assembled via:

$$p(H_i \mid \cdot) \propto \prod_\ell \left[ q(k, t) \right]^\ell \tag{3}$$

---

---

## Towards a design process



Figure: NoMM design process.

---

## Overview

- Topic models – motivation and review
- Networks of mixed membership (NoMMs)
- Inference – a Gibbs "meta-sampler"
- NoMM typology and design
- Application to tag-enhanced expertise finding
- Conclusions and outlook

## Define evidence

- Expertise finding in digital libraries
  - Find authors from document content
  - Semantic tags to disambiguate word meaning and provide additional retrieval method



- Example: scientific community of *Neural Information Processing Systems* (NIPS) conference

Tags: *probabilistic methods, variational inference, learning algorithms*

## Define tasks + metrics; set up terminals

- Retrieval of experts $a$ for term queries $\vec{w}$ and tag queries $\vec{c}$: query likelihood model: $p(\vec{w}|a)$ and $p(\vec{c}|a)$ → measure retrieval precision
- Topic quality → measure coherence score
- Baseline: Author–topic model ATM (Rosen-Zvi et al. 2004), LDA (Blei et al. 2003)



Figure: Model design: Terminals.

## Modelling assumptions



(a) Expertise of authors weighted by the portion of authorship $a_{m,a}$.

(b) Expertise semantics expressed by topics $z$. Each author has a single field of expertise (topic distribution).

(c) Tag semantics expressed by topics $y$. Tag topics $y$ could be $\equiv z$.

## Compose model



$$p(\ldots | \vec{a}, \vec{w}, \vec{c}) \propto \ldots$$

Starting from terminals

## Compose model



$$p(x, \dots \mid \cdot) \propto a_{m,x}\, q(x, \dots) \dots$$

Up-stream evidence $\vec{a}_m$
$\rightarrow$ observed parameter node samples word author $x$

## Compose model



$$p(x, z, \dots \mid \cdot) \propto a_{m,x}\, q(x, z) \dots$$

Each author only one field of expertise (topic distribution)
$\rightarrow$ $q$-term $q(x, z)$ assigns topics to sampled author $x$ (cf. ATM)

## Compose model



$$p(x, z, \dots \mid \cdot) \propto a_{m,x}\, q(x, z)\, q(z, w) \dots$$

Topic distribution over words $\rightarrow$ can connect directly via $q(z, w)$

## Compose model



$$p(x, z \mid \cdot) \propto a_{m,x}\, q(x, z)\, q(z, w)\, q(z, c)$$

Incorporate tags via $q(z, c)$ conditioned on the same topic
$\rightarrow$ Problem: How to determine tag $c_{m,n}$ for word?

## Compose model



$$p(x, z, y \,|\, \cdot) \propto a_{m,x} \, q(x, z \oplus y) \, q(z, w) \, q(y, c)$$

→ Incorporate tag topics $y_{m,j}$ on separate sequence $(m, j)$
→ Tag boosting: adjust tag influence via tag sequence length $J_m$

## ETT1 model



Assembled $q$-terms:

$$p(x, z, y \,|\, \cdot) \propto a_{m,x} \, q(x, z \oplus y) \, q(z, w) \, q(y, c) \tag{4}$$

Easy expansion to standard Gibbs full conditionals:

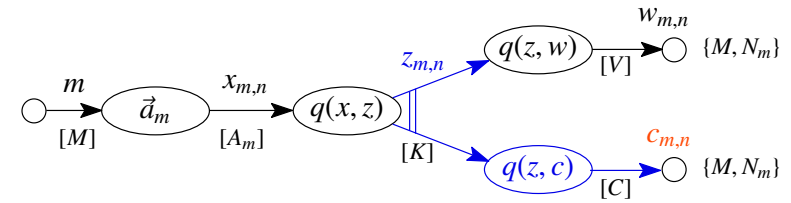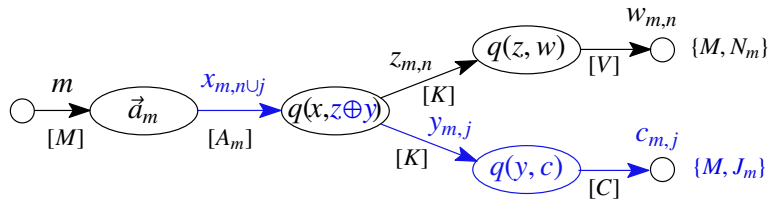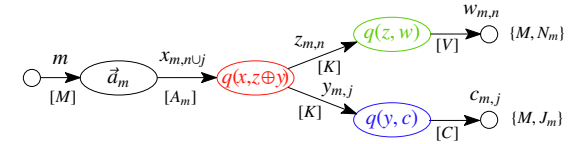$$p(x_{m,n}{=}x, z_{m,n}{=}z \,|\, \cdot) \propto a_{m,x} \cdot \frac{n_{x,z}^{\neg\{x,z\}_{m,n}} + \alpha}{n_x^{\neg\{x,z\}_{m,n}} + K\alpha} \cdot \frac{n_{z,w}^{\neg z_{m,n}} + \beta}{n_z^{\neg z_{m,n}} + V\beta} \tag{5}$$

$$p(x_{m,j}{=}x, y_{m,j}{=}y \,|\, \cdot) \propto a_{m,x} \cdot \frac{n_{x,y}^{\neg\{x,y\}_{m,j}} + \alpha}{n_y^{\neg\{x,y\}_{m,j}} + K\alpha} \cdot \frac{n_{y,c}^{\neg y_{m,j}} + \gamma}{n_y^{\neg y_{m,j}} + C\gamma} \tag{6}$$

Retrieval via query likelihood model:

$$p(\vec{w} \,|\, a) = \prod_{w \in \vec{w}} \sum_z \vartheta_{a,z} \varphi_{z,w} \qquad p(\vec{c} \,|\, a) = \prod_{c \in \vec{c}} \sum_y \vartheta_{a,y} \psi_{y,c} \,. \tag{7}$$

## ETT1 model



Assembled $q$-terms:

$$p(x, z, y \,|\, \cdot) \propto a_{m,x} \, q(x, z \oplus y) \, q(z, w) \, q(y, c) \tag{4}$$

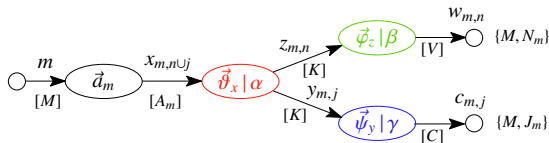Easy expansion to standard Gibbs full conditionals:

$$p(x_{m,n}{=}x, z_{m,n}{=}z \,|\, \cdot) \propto a_{m,x} \cdot \frac{n_{x,z}^{\neg\{x,z\}_{m,n}} + \alpha}{n_x^{\neg\{x,z\}_{m,n}} + K\alpha} \cdot \frac{n_{z,w}^{\neg z_{m,n}} + \beta}{n_z^{\neg z_{m,n}} + V\beta} \tag{5}$$

$$p(x_{m,j}{=}x, y_{m,j}{=}y \,|\, \cdot) \propto a_{m,x} \cdot \frac{n_{x,y}^{\neg\{x,y\}_{m,j}} + \alpha}{n_y^{\neg\{x,y\}_{m,j}} + K\alpha} \cdot \frac{n_{y,c}^{\neg y_{m,j}} + \gamma}{n_y^{\neg y_{m,j}} + C\gamma} \tag{6}$$

Retrieval via query likelihood model:

$$p(\vec{w} \,|\, a) = \prod_{w \in \vec{w}} \sum_z \vartheta_{a,z} \varphi_{z,w} \qquad p(\vec{c} \,|\, a) = \prod_{c \in \vec{c}} \sum_y \vartheta_{a,y} \psi_{y,c} \,. \tag{7}$$

## Typical derivation method (Is it really that complex?)



(a) Expert–tag–topic model 1 (ETT1)
(Heinrich 2011)

## Model evaluation



Figure: ETT1 example query in community browser.

---

## Truncated average precision



$$AP@5 = \frac{\frac{1/2}{} \quad +2/4 \; +3/5}{3} = 0.533$$

$$AP@5 = \frac{1/1 \quad +2/2 \qquad +3/5}{3} = 0.867$$

Figure: Average precision at 5 (3 relevant documents in corpus).

---

## Retrieval and clustering results



(a) Term queries

(b) Tag queries

- Term retrieval improved by tag influence during *training* time
- Mutual information between a-priori tag clusterings $p(c\,|\,a)$ and topic clusterings $p(z\,|\,a)$: ETT1 $\geq$ 1.002 vs. ATM = 0.865.
- Semi-supervised features: find relevant items with missing tags
- Tag strength: bias towards strong tags in combinations

---

## Topic coherence results

Topic coherence (Mimno et al. 2011):
- $\approx$ How often do top-ranked topic terms co-occur in documents?
- Re-enacts human judgement in topic intrusion experiments (Chang et al. 2009; Heinrich 2011)



Words in topic (choose worst match (A-F) in every group):

| 1. | 2. | 3. |
|---|---|---|
| A. orientation | A. likelihood | A. risk |
| B. cortex | B. mixture | B. return |
| C. visual | C. theorem | C. stock |
| D. ocular | D. density | D. trading |
| E. acoustic | E. em | E. processor |
| F. eye | F. prior | F. prediction |

| 4. | 5. | 6. |
|---|---|---|
| A. language | A. circuit | A. validation |
| B. word | B. bayesian | B. set |
| C. stress | C. analog | C. variance |
| D. grammar | D. voltage | D. regression |
| E. neural | E. vlsi | E. selection |
| F. syllable | F. chip | F. bias |

(a) Topic intrusion experiment

(b) Coherence scores

## Overview

- Topic models – motivation and review
- Networks of mixed membership (NoMMs)
- Inference – a Gibbs "meta-sampler"
- NoMM typology and design
- Application to tag-enhanced expertise finding
- Conclusions and outlook

## Conclusions

- Networks of mixed membership:
  Domain-specific compact representation
- Inference:
  - Generic Gibbs sampling: $q$-functions as central quantity in model behaviour
  - Gibbs meta-sampler: simplify implementation
  - Hybrid acceleration methods
  - Alternatives: variational Bayes (Heinrich and Goesele 2009), collapsed VB
- Typology and design method:
  - Model structure types: literature + novel
  - Building blocks for design with predictable properties
- Application:
  - Expert–tag–topic model demonstrates design
  - Tags improve retrieval and topic coherence

## Towards an R-based Gibbs meta-sampler

- R environment becoming popular for topic models, e.g.:
  - `topicmodels` package implementing general and various special cases (Grün and Hornik 2011), based on text mining package `tm`
  - `lda` package with LDA, supervised, relational topic models (Blei et al. 2003; Blei and McAuliffe 2007; Chang and Blei 2009)
- Vision: Use Gibbs meta-sampler as front-end to create R-based high-performance code ↔ use R as experimental front-end



- Extend to non-parametric distributions, e.g., based on `DPpackage` (Jara et al. 2012):
  - NoMMs as polymorphism of parametric and non-parametric models (with different Bayesian networks)

Q+A

`http://arbylon.net/resources.html`

# References I

References

Blei, D. and J. McAuliffe (2007).
Supervised topic models.
In *Advances in Neural Information Processing Systems.*

Blei, D., A. Ng, and M. Jordan (2003, January).
Latent Dirichlet allocation.
*Journal of Machine Learning Research 3*, 993–1022.

Buntine, W. and A. Jakulin (2005).
Discrete principal components analysis.
In *Proc. ECML.*

Chang, J. and D. M. Blei (2009).
Relational topic models for document networks.
In *AISTATS.*

Chang, J., J. Boyd-Graber, S. Gerrish, C. Wang, and D. Blei (2009).
Reading tea leaves: How humans interpret topic models.
In *Proc. Neural Information Processing Systems (NIPS).*

Griffiths, T. L., J. B. Tenenbaum, and M. Steyvers (2007).
Topics in semantic representation.
*Psychological Review 114*(2), 211–244.

# References II

Grün, B. and K. Hornik (2011).
topicmodels: An R package for fitting topic models.
*Journal of Statistical Software 43*(13).

Heinrich, G. (2009).
A generic approach to topic models.
In *Proc. European Conf. on Mach. Learn. / Principles and Pract. of Know. Discov. in Databases (ECML/PKDD), Part 1*, pp. 517–532.

Heinrich, G. (2011).
Typology of mixed-membership models: Towards a design method.
In *Proc. European Conf. on Mach. Learn. / Principles and Pract. of Know. Discov. in Databases (ECML/PKDD).*

Heinrich, G. and M. Goesele (2009).
Variational Bayes for generic topic models.
In *Proc. 32nd Annual German Conference on Artificial Intelligence (KI2009).*

Heinrich, G., J. Kindermann, C. Lauth, G. Paaß, and J. Sanchez-Monzon (2005).
Investigating word correlation at different scopes – a latent concept approach.
In *Workshop Lexical Ontology Learning at Int. Conf. Mach. Learning.*

# References III

Jara, A., T. Hanson, F. A. Quintana, P. Mueller, and G. L. Rosner (2012, Feb.).
DPpackage: Bayesian nonparametric modeling in R.
software documentation.

Landauer, T. K. and S. T. Dumais (1997).
Solution to Plato's problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge.
*Psych. Rev. 104*(2), 211–240.
Cognitive view on LSA.

Li, W., D. Blei, and A. McCallum (2007).
Mixtures of hierarchical topics with pachinko allocation.
In *International Conference on Machine Learning.*

Li, W. and A. McCallum (2006).
Pachinko allocation: DAG-structured mixture models of topic correlations.
In *ICML '06: Proceedings of the 23rd international conference on Machine learning*, New York, NY, USA, pp. 577–584. ACM.

Mimno, D., H. M. Wallach, E. Talley, M. Leenders, and A. McCallum (2011, July).
Optimizing semantic coherence in topic models.
In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, Edinburgh, UK*, pp. 262272.

# References IV

Newman, D., A. Asuncion, P. Smyth, and M. Welling (2009, August).
Distributed algorithms for topic models.
*JMLR 10*, 1801–1828.

Porteous, I., D. Newman, A. Ihler, A. Asuncion, P. Smyth, and M. Welling (2008).
Fast collapsed Gibbs sampling for latent Dirichlet allocation.
In *KDD '08: Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, New York, NY, USA, pp. 569–577. ACM.

Rosen-Zvi, M., T. Griffiths, M. Steyvers, and P. Smyth (2004).
The author-topic model for authors and documents.
In *Proc. 20th Conference on Uncertainty in Artificial Intelligence (UAI).*