Contributions to probabilistic non-negative matrix factorization Maximum marginal likelihood estimation and temporal Markovian models

Louis FILSTROFF IRIT, Univ. Toulouse, CNRS

Joint work with Cédric Févotte, Alberto Lumbreras, Olivier Gouvert, Olivier Cappé

November 29th, 2019









Louis Filstroff

INRA-MIAT Seminar

1/39





2 MMLE in the Gamma-Poisson model





Outline



- 2 MMLE in the Gamma-Poisson model
- 3 Temporal NMF

4 Conclusions and perspectives

In many situations, data are available as matrices

- In many situations, data are available as matrices
- Consider N samples \mathbf{v}_n in \mathbb{R}^F (i.e., described by F features)

- In many situations, data are available as matrices
- Consider N samples \mathbf{v}_n in \mathbb{R}^F (i.e., described by F features)
- Samples stored column-wise, yielding an $F \times N$ data matrix **V**

- In many situations, data are available as matrices
- Consider N samples \mathbf{v}_n in \mathbb{R}^F (i.e., described by F features)
- Samples stored column-wise, yielding an F × N data matrix V

V represents	f	п	Typical <i>F</i>
A corpus of documents	Words	Documents	$10^{4}-10^{5}$
A collection of grayscale images	Pixels	Images	$10^4 - 10^6$
The spectrogram of an audio signal	Frequencies	Time frames	$10^{3} - 10^{4}$
Ratings	Items	Users	$10^{6} - 10^{8}$

Table 1: Examples of data available as matrices

Matrix factorization (1/3)

Matrix factorization (MF)

MF aims at finding a decomposition of the data matrix ${\bf V}$ as the product of two matrices

$$\mathbf{V} \simeq \mathbf{W} \mathbf{H},$$
 (1)

where **W** is of size $F \times K$, and **H** is of size $K \times N$

Matrix factorization (1/3)

Matrix factorization (MF)

MF aims at finding a decomposition of the data matrix ${\bf V}$ as the product of two matrices

$$\mathbf{V} \simeq \mathbf{W} \mathbf{H},$$
 (1)

where **W** is of size $F \times K$, and **H** is of size $K \times N$

• $K \ll \min(F, N)$: low-rank approximation

Matrix factorization (1/3)

Matrix factorization (MF)

MF aims at finding a decomposition of the data matrix ${\bf V}$ as the product of two matrices

$$\mathbf{V} \simeq \mathbf{W} \mathbf{H},$$
 (1)

where **W** is of size $F \times K$, and **H** is of size $K \times N$

- $K \ll \min(F, N)$: low-rank approximation
- Linear dimensionality reduction technique

$$\mathbf{v}_n \simeq \sum_{k=1}^K h_{kn} \mathbf{w}_k \tag{2}$$

Matrix factorization (2/3)

- ► W is called the dictionary. Columns represent characteristic or recurring patterns of the data
- ► H is called the activation coefficients. The *n*-th column represents how much of each pattern is needed to represent v_n



Matrix factorization (3/3)

MF can be written as an optimization problem

$$\min_{\mathbf{W},\mathbf{H}} D(\mathbf{V}|\mathbf{W}\mathbf{H}) = \sum_{f,n} d(v_{fn}|[\mathbf{W}\mathbf{H}]_{fn})$$
(3)

- D is a separable measure of fit ("divergence")
- Additional constraints over W and H for interpretability

Matrix factorization (3/3)

MF can be written as an optimization problem

$$\min_{\mathbf{W},\mathbf{H}} D(\mathbf{V}|\mathbf{W}\mathbf{H}) = \sum_{f,n} d(v_{fn}|[\mathbf{W}\mathbf{H}]_{fn})$$
(3)

- D is a separable measure of fit ("divergence")
- ► Additional constraints over **W** and **H** for interpretability
- Ubiquitous example : principal component analysis (PCA) [Pearson, 1901, Hotelling, 1933]...

Matrix factorization (3/3)

MF can be written as an optimization problem

$$\min_{\mathbf{W},\mathbf{H}} D(\mathbf{V}|\mathbf{W}\mathbf{H}) = \sum_{f,n} d(v_{fn}|[\mathbf{W}\mathbf{H}]_{fn})$$
(3)

- D is a separable measure of fit ("divergence")
- Additional constraints over W and H for interpretability
- Ubiquitous example : principal component analysis (PCA) [Pearson, 1901, Hotelling, 1933]...
- ... but does not take into account the support of the data

Non-negative matrix factorization

Non-negative matrix factorization (NMF)

NMF aims at finding a decomposition of a non-negative data matrix ${\bf V}$ as the product of two non-negative matrices

$$\min_{\mathbf{W} \ge 0, \mathbf{H} \ge 0} D(\mathbf{V} | \mathbf{W} \mathbf{H})$$
(4)

[Paatero and Tapper, 1994, Lee and Seung, 1999]

Non-negative matrix factorization

Non-negative matrix factorization (NMF)

NMF aims at finding a decomposition of a non-negative data matrix ${\bf V}$ as the product of two non-negative matrices

$$\min_{\mathbf{N} \ge 0, \mathbf{H} \ge 0} D(\mathbf{V} | \mathbf{W} \mathbf{H})$$
(4)

[Paatero and Tapper, 1994, Lee and Seung, 1999]

► Non-negativity constraints improve interpretability :

- $\blacktriangleright~ \textbf{W} \geq 0$: direct interpretation of the columns of W
- $\mathbf{H} \ge 0$: induces part-based representations

Non-negative matrix factorization

Non-negative matrix factorization (NMF)

NMF aims at finding a decomposition of a non-negative data matrix ${\bf V}$ as the product of two non-negative matrices

$$\min_{\mathbf{W} \ge 0, \mathbf{H} \ge 0} D(\mathbf{V} | \mathbf{W} \mathbf{H})$$
(4)

[Paatero and Tapper, 1994, Lee and Seung, 1999]

- Non-negativity constraints improve interpretability :
 - $\blacktriangleright~ \textbf{W} \geq 0$: direct interpretation of the columns of W
 - $\blacktriangleright~\textbf{H} \geq 0$: induces part-based representations
- Many application fields :
 - Audio signal processing (source separation [Virtanen, 2007], music transcription [Smaragdis and Brown, 2003])
 - Text information retrieval (topic modeling [Xu et al., 2003])
 - Hyperspectral imaging (unmixing [Bioucas-Dias et al., 2012])

Divergences and statistical models

Name	d(x y)
Squared Euclidean distance (SED)	$\frac{1}{2}(x-y)^2$
Kullback-Leibler (KL) divergence	$x \log\left(\frac{x}{y}\right) - x + y$
Itakura-Saito (IS) divergence	$\frac{x}{y} - \log\left(\frac{x}{y}\right) + 1$

Table 2: Typical divergences used in NMF

For many usual cost functions, the minimization problem is equivalent to the joint maximum likelihood estimation of ${\bf W}$ and ${\bf H}$

$$\min_{\mathbf{W} \ge 0, \mathbf{H} \ge 0} D(\mathbf{V} | \mathbf{W} \mathbf{H}) \Leftrightarrow \max_{\mathbf{W}, \mathbf{H}} p(\mathbf{V}; \mathbf{W}, \mathbf{H})$$
(5)

The probabilistic framework

NMF problem	Equivalent likelihood
SED-NMF KL-NMF	$egin{aligned} & v_{fn} \sim \mathcal{N}([\mathbf{WH}]_{fn}, \sigma^2) \ & v_{fn} \sim \mathrm{Poisson}([\mathbf{WH}]_{fn}) \end{aligned}$
IS-NMF	$v_{\mathit{fn}} \sim Exp\left(rac{1}{[WH]_{\mathit{fn}}} ight)$

Table 3: Equivalences with statistical models

The probabilistic framework

NMF problem	Equivalent likelihood
SED-NMF KL-NMF	$v_{fn} \sim \mathcal{N}([\mathbf{WH}]_{fn}, \sigma^2)$ $v_{fn} \sim \text{Poisson}([\mathbf{WH}]_{fn})$ $v_{fn} \in \mathbb{E}$
13-INIVIE	$V_{fn} \sim \text{Exp}\left(\overline{[\mathbf{WH}]_{fn}}\right)$

Table 3: Equivalences with statistical models

Probabilistic NMF

Learning (estimation and/or inference) tasks in statistical models of the form

$$\mathbf{v}_n \sim p(.; \mathbf{Wh}_n, \mathbf{\Psi}),$$
 (6)

i.e., parametrized by the dot product \mathbf{Wh}_n (and Ψ)

NB : Most of the time we have $\mathbb{E}(\mathbf{v}_n) = \mathbf{W}\mathbf{h}_n$

8/39



(1) Frequentist NMF – Maximum likelihood estimation



► (1) Frequentist NMF – Maximum likelihood estimation

(3) Bayesian NMF – Infer posterior distribution $p(\mathbf{W}, \mathbf{H} | \mathbf{V})$



- (1) Frequentist NMF Maximum likelihood estimation
- (2) Semi-Bayesian NMF Our setting
- ► (3) Bayesian NMF Infer posterior distribution $p(\mathbf{W}, \mathbf{H} | \mathbf{V})$

Two estimation paradigms

Two estimation paradigms

Maximizing joint likelihood estimation (MJLE) :

 $\max_{\mathbf{W},\mathbf{H}} \log p(\mathbf{V},\mathbf{H};\mathbf{W}) = \log p(\mathbf{V}|\mathbf{H};\mathbf{W}) + \log p(\mathbf{H})$ (7)

Estimation of FK + KN parameters

Two estimation paradigms

Maximizing joint likelihood estimation (MJLE) :

 $\max_{\mathbf{W},\mathbf{H}} \log p(\mathbf{V},\mathbf{H};\mathbf{W}) = \log p(\mathbf{V}|\mathbf{H};\mathbf{W}) + \log p(\mathbf{H})$ (7)

Estimation of FK + KN parameters

Maximizing marginal likelihood estimation (MMLE) :

$$\max_{\mathbf{W}} \log p(\mathbf{V}; \mathbf{W}) = \log \int_{\mathbf{H}} p(\mathbf{V} | \mathbf{H}; \mathbf{W}) p(\mathbf{H}) d\mathbf{H}$$
(8)

Estimation of FK parameters

Two estimation paradigms

Maximizing joint likelihood estimation (MJLE) :

 $\max_{\mathbf{W},\mathbf{H}} \log p(\mathbf{V},\mathbf{H};\mathbf{W}) = \log p(\mathbf{V}|\mathbf{H};\mathbf{W}) + \log p(\mathbf{H})$ (7)

Estimation of FK + KN parameters

Maximizing marginal likelihood estimation (MMLE) :

$$\max_{\mathbf{W}} \log p(\mathbf{V}; \mathbf{W}) = \log \int_{\mathbf{H}} p(\mathbf{V} | \mathbf{H}; \mathbf{W}) p(\mathbf{H}) d\mathbf{H}$$
(8)

Estimation of FK parameters

MMLE is a better-posed approach because the number of parameters to be estimated is fixed w.r.t. the number of samples N

 Empirical comparison of the two paradigms [Dikmen and Févotte, 2011, Dikmen and Févotte, 2012] in the Poisson and Exponential models

- Empirical comparison of the two paradigms [Dikmen and Févotte, 2011, Dikmen and Févotte, 2012] in the Poisson and Exponential models
 - MMLE tends to automatically prune the columns of W, while MJLE makes use of all K columns

- Empirical comparison of the two paradigms [Dikmen and Févotte, 2011, Dikmen and Févotte, 2012] in the Poisson and Exponential models
 - MMLE tends to automatically prune the columns of W, while MJLE makes use of all K columns
 - Favorable behavior that was left unexplained

- Empirical comparison of the two paradigms [Dikmen and Févotte, 2011, Dikmen and Févotte, 2012] in the Poisson and Exponential models
 - MMLE tends to automatically prune the columns of W, while MJLE makes use of all K columns
 - Favorable behavior that was left unexplained
- Related works of the literature

- Empirical comparison of the two paradigms [Dikmen and Févotte, 2011, Dikmen and Févotte, 2012] in the Poisson and Exponential models
 - MMLE tends to automatically prune the columns of W, while MJLE makes use of all K columns
 - Favorable behavior that was left unexplained
- Related works of the literature
 - Integration of nuisance parameters [Berger et al., 1999]

- Empirical comparison of the two paradigms [Dikmen and Févotte, 2011, Dikmen and Févotte, 2012] in the Poisson and Exponential models
 - MMLE tends to automatically prune the columns of W, while MJLE makes use of all K columns
 - Favorable behavior that was left unexplained
- Related works of the literature
 - Integration of nuisance parameters [Berger et al., 1999]
 - Noisy ICA [Moulines et al., 1997]

- Empirical comparison of the two paradigms [Dikmen and Févotte, 2011, Dikmen and Févotte, 2012] in the Poisson and Exponential models
 - MMLE tends to automatically prune the columns of W, while MJLE makes use of all K columns
 - Favorable behavior that was left unexplained
- Related works of the literature
 - Integration of nuisance parameters [Berger et al., 1999]
 - Noisy ICA [Moulines et al., 1997]
 - Latent Dirichlet allocation (LDA) [Blei et al., 2003]

Overview

We will focus on two sub-cases :
Overview

We will focus on two sub-cases :

Independent priors on H

$$p(\mathbf{H}) = \prod_{n=1}^{N} p(\mathbf{h}_n)$$
(9)

- Poisson likelihood + Gamma prior
- (Exponential likelihood + Inverse Gamma prior)

Overview

We will focus on two sub-cases :

Independent priors on H

$$p(\mathbf{H}) = \prod_{n=1}^{N} p(\mathbf{h}_n)$$
(9)

- Poisson likelihood + Gamma prior
- (Exponential likelihood + Inverse Gamma prior)
- ► Temporal priors on **H**

$$p(\mathbf{H}) = p(\mathbf{h}_1) \prod_{n \ge 2} p(\mathbf{h}_n | \mathbf{h}_{n-1})$$
(10)

Design of a meaningful prior





2 MMLE in the Gamma-Poisson model

3 Temporal NMF

4 Conclusions and perspectives

Model and objective

The Gamma-Poisson (GaP) model [Canny, 2004] $h_{kn} \sim \text{Gamma}(\alpha_k, \beta_k)$ (11) $v_{fn} | \mathbf{h}_n \sim \text{Poisson}([\mathbf{WH}]_{fn})$ (12)

Model and objective

The Gamma-Poisson (GaP) model [Canny, 2004]

 $h_{kn} \sim \text{Gamma}(\alpha_k, \beta_k)$ (11) $v_{fn} | \mathbf{h}_n \sim \text{Poisson}([\mathbf{WH}]_{fn})$ (12)

- Different application fields for this observation model
 - Text information retrieval (observation model very close to LDA) [Canny, 2004, Buntine and Jakulin, 2006]
 - Recommender systems ("Poisson factorization") [Gopalan et al., 2015]
 - Image processing [Cemgil, 2009]

Model and objective

The Gamma-Poisson (GaP) model [Canny, 2004]

 $h_{kn} \sim \text{Gamma}(\alpha_k, \beta_k)$ (11) $v_{fn} | \mathbf{h}_n \sim \text{Poisson}([\mathbf{WH}]_{fn})$ (12)

- Different application fields for this observation model
 - Text information retrieval (observation model very close to LDA) [Canny, 2004, Buntine and Jakulin, 2006]
 - Recommender systems ("Poisson factorization") [Gopalan et al., 2015]
 - Image processing [Cemgil, 2009]
- MMLE amounts to

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}) = -\log \int p(\mathbf{V}|\mathbf{H}; \mathbf{W}) p(\mathbf{H}) d\mathbf{H}$$
(13)

Hyperparameters lpha and eta may also be optimized

Data augmentation

Superposition property of the Poisson distribution

The augmented GaP model	
$h_{kn}\sim Gamma(lpha_k,eta_k)$	(14)
$c_{fkn} \sim Poisson(w_{fk}h_{kn})$	(15)
$v_{fn} = \sum c_{fkn}$	(16)
k	

Data augmentation

Superposition property of the Poisson distribution

The augmented GaP model	
$h_{kn} \sim Gamma(lpha_k,eta_k)$	(14)
$c_{fkn} \sim Poisson(w_{fk}h_{kn})$	(15)
$v_{fn} = \sum_{k} c_{fkn}$	(16)

• **C** denotes the $F \times K \times N$ tensor with entries c_{fkn}

Data augmentation

Superposition property of the Poisson distribution

(14)
(15)
(16)

- **C** denotes the $F \times K \times N$ tensor with entries c_{fkn}
- ► Thanks to the conjugacy between the Poisson and the Gamma distribution, *h_{kn}* can be marginalized out from Eqs. (14)-(15)

New formulation of GaP

$$\mathbf{c}_{kn} \sim \mathsf{NM}\left(\alpha_{k}, \left[\frac{w_{1k}}{\sum_{f} w_{fk} + \beta_{k}}, \dots, \frac{w_{Fk}}{\sum_{f} w_{fk} + \beta_{k}}\right]^{\mathsf{T}}\right)$$
(17)
$$\mathbf{v}_{n} = \sum_{k} \mathbf{c}_{kn}$$
(18)

where $\mathbf{c}_{kn} = [c_{1kn}, \dots, c_{Fkn}]^{\mathsf{T}}$ is a vector of size F

New formulation of GaP

$$\mathbf{c}_{kn} \sim \mathsf{NM}\left(\alpha_{k}, \left[\frac{w_{1k}}{\sum_{f} w_{fk} + \beta_{k}}, \dots, \frac{w_{Fk}}{\sum_{f} w_{fk} + \beta_{k}}\right]^{\mathsf{T}}\right)$$
(17)
$$\mathbf{v}_{n} = \sum_{k} \mathbf{c}_{kn}$$
(18)

where $\mathbf{c}_{kn} = [c_{1kn}, \dots, c_{Fkn}]^{\mathsf{T}}$ is a vector of size F

The vector c_{kn} has a so-called negative multinomial (NM) distribution, known in closed form

New formulation of GaP

$$\mathbf{c}_{kn} \sim \mathsf{NM}\left(\alpha_{k}, \left[\frac{w_{1k}}{\sum_{f} w_{fk} + \beta_{k}}, \dots, \frac{w_{Fk}}{\sum_{f} w_{fk} + \beta_{k}}\right]^{\mathsf{T}}\right)$$
(17)
$$\mathbf{v}_{n} = \sum_{k} \mathbf{c}_{kn}$$
(18)

where $\mathbf{c}_{kn} = [c_{1kn}, \dots, c_{Fkn}]^{\mathsf{T}}$ is a vector of size F

- The vector c_{kn} has a so-called negative multinomial (NM) distribution, known in closed form
- GaP can therefore be seen as a composite NM model

New formulation of GaP

$$\mathbf{c}_{kn} \sim \mathsf{NM}\left(\alpha_{k}, \left[\frac{w_{1k}}{\sum_{f} w_{fk} + \beta_{k}}, \dots, \frac{w_{Fk}}{\sum_{f} w_{fk} + \beta_{k}}\right]^{\mathsf{T}}\right)$$
(17)
$$\mathbf{v}_{n} = \sum_{k} \mathbf{c}_{kn}$$
(18)

where
$$\mathbf{c}_{kn} = [c_{1kn}, \dots, c_{Fkn}]^{\mathsf{T}}$$
 is a vector of size F

- The vector c_{kn} has a so-called negative multinomial (NM) distribution, known in closed form
- ► GaP can therefore be seen as a composite NM model
- Alternative characterization with the multinomial distribution

Closed-form marginal likelihood

Closed-form marginal likelihood

Denote by C_V the set of "admissible components"

$$\mathcal{C}_{\mathbf{V}} = \{ \mathbf{C} \in \mathbb{N}^{F \times K \times N} \, | \, \forall (f, n), \sum_{k} c_{fkn} = v_{fn} \}.$$
(19)

Closed-form marginal likelihood

Denote by C_V the set of "admissible components"

$$\mathcal{C}_{\mathbf{V}} = \{ \mathbf{C} \in \mathbb{N}^{F \times K \times N} \, | \, \forall (f, n), \sum_{k} c_{fkn} = v_{fn} \}.$$
(19)

The marginalization of C yields

$$p(\mathbf{V};\mathbf{W}) = \sum_{\mathbf{C}\in\mathcal{C}_{\mathbf{V}}} p(\mathbf{C};\mathbf{W}) = \sum_{\mathbf{C}\in\mathcal{C}_{\mathbf{V}}} \prod_{k,n} \underbrace{p(\mathbf{c}_{kn};\mathbf{w}_k)}_{NM}$$
(20)

After computation, we obtain

$$-\frac{1}{N}\mathcal{L}(\mathbf{W}) = -\frac{1}{N}\log\left(\sum_{\mathbf{C}\in\mathcal{C}_{\mathbf{V}}}f(\mathbf{C};\mathbf{W})\right)$$
(21)
$$+\sum_{k}\alpha_{k}\log(||\mathbf{w}_{k}||_{1}+\beta_{k}) + \operatorname{cst}$$
(22)

After computation, we obtain

$$-\frac{1}{N}\mathcal{L}(\mathbf{W}) = -\frac{1}{N}\log\left(\sum_{\mathbf{C}\in\mathcal{C}_{\mathbf{V}}}f(\mathbf{C};\mathbf{W})\right)$$
(21)
$$+\sum_{k}\alpha_{k}\log(||\mathbf{w}_{k}||_{1}+\beta_{k}) + \operatorname{cst}$$
(22)

"Data-fitting term" + "regularization term"

After computation, we obtain

$$-\frac{1}{N}\mathcal{L}(\mathbf{W}) = -\frac{1}{N}\log\left(\sum_{\mathbf{C}\in\mathcal{C}_{\mathbf{V}}}f(\mathbf{C};\mathbf{W})\right)$$
(21)
+ $\sum_{k}\alpha_{k}\log(||\mathbf{w}_{k}||_{1}+\beta_{k}) + \operatorname{cst}$ (22)

- "Data-fitting term" + "regularization term"
- ► Term of the form R(x) = ∑_k log(|x_k| + ϵ) is known to be sparsity-inducing [Candès et al., 2008]

After computation, we obtain

$$-\frac{1}{N}\mathcal{L}(\mathbf{W}) = -\frac{1}{N}\log\left(\sum_{\mathbf{C}\in\mathcal{C}_{\mathbf{V}}}f(\mathbf{C};\mathbf{W})\right)$$
(21)
$$+\sum_{k}\alpha_{k}\log(||\mathbf{w}_{k}||_{1}+\beta_{k}) + \operatorname{cst}$$
(22)

- "Data-fitting term" + "regularization term"
- ► Term of the form R(x) = ∑_k log(|x_k| + ϵ) is known to be sparsity-inducing [Candès et al., 2008]
- Provides a deeper understanding of the self-regularization phenomenon

17/39

 EM algorithm [Dempster et al., 1977] iteratively optimizes the following functional

$$Q(\mathbf{W}; \tilde{\mathbf{W}}) = \int_{\mathbf{Z}} \log p(\mathbf{V}, \mathbf{Z}; \mathbf{W}) p(\mathbf{Z} | \mathbf{V}; \tilde{\mathbf{W}})$$
(23)

 ${\boldsymbol{\mathsf{Z}}}$ is the set of latent variables, $\tilde{{\boldsymbol{\mathsf{W}}}}$ the current value

 EM algorithm [Dempster et al., 1977] iteratively optimizes the following functional

$$Q(\mathbf{W}; \tilde{\mathbf{W}}) = \int_{\mathbf{Z}} \log p(\mathbf{V}, \mathbf{Z}; \mathbf{W}) p(\mathbf{Z} | \mathbf{V}; \tilde{\mathbf{W}})$$
(23)

 ${\boldsymbol{\mathsf{Z}}}$ is the set of latent variables, $\tilde{{\boldsymbol{\mathsf{W}}}}$ the current value

- Three possible choices :
 - ► **Z** = {**C**, **H**} : known from [Dikmen and Févotte, 2012]
 - ► **Z** = {**H**} : known from [Dikmen and Févotte, 2012]

► **Z** = {**C**} : novel

 EM algorithm [Dempster et al., 1977] iteratively optimizes the following functional

$$Q(\mathbf{W}; \tilde{\mathbf{W}}) = \int_{\mathbf{Z}} \log p(\mathbf{V}, \mathbf{Z}; \mathbf{W}) p(\mathbf{Z} | \mathbf{V}; \tilde{\mathbf{W}})$$
(23)

 ${\boldsymbol{\mathsf{Z}}}$ is the set of latent variables, $\tilde{{\boldsymbol{\mathsf{W}}}}$ the current value

- Three possible choices :
 - ► **Z** = {**C**, **H**} : known from [Dikmen and Févotte, 2012]
 - $Z = \{H\}$: known from [Dikmen and Févotte, 2012]

In all cases the posterior of the latent variables p(Z|V; W̃) is not tractable

¹Majorization-Minimization

Louis Filstroff

▶ We resort to Monte Carlo EM [Wei and Tanner, 1990]

¹Majorization-Minimization

Louis Filstroff

- ▶ We resort to Monte Carlo EM [Wei and Tanner, 1990]
- Sampling from $p(\mathbf{C}, \mathbf{H} | \mathbf{V}; \tilde{\mathbf{W}})$ with a Gibbs sampling procedure

¹Majorization-Minimization

Louis Filstroff

- ▶ We resort to Monte Carlo EM [Wei and Tanner, 1990]
- Sampling from $p(\mathbf{C}, \mathbf{H} | \mathbf{V}; \tilde{\mathbf{W}})$ with a Gibbs sampling procedure
- Optimizes instead

$$\hat{Q}(\mathbf{W}) = \frac{1}{J} \sum_{j} \log p(\mathbf{V}, \mathbf{Z}^{(j)}; \mathbf{W})$$
(24)

Can be carried out in closed form for EM-CH and EM-C, EM-H requires a $\rm MM^1\text{-}based$ procedure

¹Majorization-Minimization

Experimental work (1/3)

Synthetic dataset \bm{V}_1 (4 \times 100) generated from the GaP model with \bm{W}_1^\star



Figure 1: Speed of convergence comparison of the three algorithms on dataset \mathbf{V}_1 ($\mathcal{K}=3$)

Experimental work (2/3)

- Synthetic dataset V₂ (4 × 100) generated from the GaP model with W^{*}₂ = 100 × W^{*}₁
- Over-dispersed, non-sparse



Figure 2: Speed of convergence comparison of the three algorithms on dataset ${\bf V}_2$ (${\cal K}=3)$

Experimental work (3/3)

- Taste Profile dataset (1509 × 805)
- Over-dispersed, sparse



Figure 3: Speed of convergence comparison of the three algorithms the Taste Profile dataset (K = 10)

Louis Filstroff

Outline



2 MMLE in the Gamma-Poisson model



4 Conclusions and perspectives

Temporal prior

Matrices V where the samples are correlated

Temporal prior

- Matrices V where the samples are correlated
- ► Add correlation to the models by lifting the independence assumption on the columns of **H**

Temporal prior

- Matrices V where the samples are correlated
- ► Add correlation to the models by lifting the independence assumption on the columns of **H**
- Markov structure on the columns + independence of the rows

$$p(\mathbf{H}) = \prod_{k} p(h_{k1}) \prod_{n \ge 2} p(h_{kn} | h_{k(n-1)})$$
(25)

So-called temporal models
Temporal prior

- Matrices V where the samples are correlated
- ► Add correlation to the models by lifting the independence assumption on the columns of **H**
- Markov structure on the columns + independence of the rows

$$p(\mathbf{H}) = \prod_{k} p(h_{k1}) \prod_{n \ge 2} p(h_{kn} | h_{k(n-1)})$$
(25)

So-called temporal models

 Non-negative Markov chains, in relation with the Gamma distribution

Gamma Markov Chains of the literature (1/3)

Chaining on the rate parameter

$$h_{kn}|h_{k(n-1)} \sim \text{Gamma}\left(\alpha, \frac{\beta}{h_{k(n-1)}}\right)$$
 (26)

[Févotte et al., 2009, Févotte, 2011]

Gamma Markov Chains of the literature (1/3)

Chaining on the rate parameter

$$h_{kn}|h_{k(n-1)} \sim \operatorname{Gamma}\left(\alpha, \frac{\beta}{h_{k(n-1)}}\right)$$
 (26)

[Févotte et al., 2009, Févotte, 2011]

•
$$\mathbb{E}(h_{kn}|h_{k(n-1)}) = \frac{\alpha}{\beta}h_{k(n-1)}$$

No well-defined stationary distribution

Gamma Markov Chains of the literature (1/3)

Chaining on the rate parameter

$$h_{kn}|h_{k(n-1)} \sim \operatorname{Gamma}\left(\alpha, \frac{\beta}{h_{k(n-1)}}\right)$$
 (26)

[Févotte et al., 2009, Févotte, 2011]

•
$$\mathbb{E}(h_{kn}|h_{k(n-1)}) = \frac{\alpha}{\beta}h_{k(n-1)}$$

No well-defined stationary distribution



Gamma Markov Chains of the literature (2/3)

Chaining on the rate parameter with an auxiliary variable

$$z_{kn}|h_{k(n-1)} \sim \mathsf{Gamma}\left(\alpha_z, \beta_z h_{k(n-1)}\right)$$
 (27)

$$h_{kn}|z_{kn} \sim \text{Gamma}\left(\alpha_h, \beta_h z_{kn}\right)$$
 (28)

[Cemgil and Dikmen, 2007]

Gamma Markov Chains of the literature (2/3)

Chaining on the rate parameter with an auxiliary variable

$$z_{kn}|h_{k(n-1)} \sim \mathsf{Gamma}\left(\alpha_z, \beta_z h_{k(n-1)}\right)$$
 (27)

$$h_{kn}|z_{kn} \sim \text{Gamma}\left(\alpha_h, \beta_h z_{kn}\right)$$
 (28)

[Cemgil and Dikmen, 2007]

$$\blacktriangleright \mathbb{E}(h_{kn}|h_{k(n-1)}) = \frac{\beta_z \alpha_h}{\beta_h(\alpha_z - 1)} h_{k(n-1)}$$

No well-defined stationary distribution

Gamma Markov Chains of the literature (2/3)

Chaining on the rate parameter with an auxiliary variable

$$z_{kn}|h_{k(n-1)} \sim \mathsf{Gamma}\left(lpha_z, eta_z h_{k(n-1)}
ight)$$
 (27)

$$h_{kn}|z_{kn} \sim \text{Gamma}\left(\alpha_h, \beta_h z_{kn}\right)$$
 (28)

[Cemgil and Dikmen, 2007]

$$\blacktriangleright \mathbb{E}(h_{kn}|h_{k(n-1)}) = \frac{\beta_z \alpha_h}{\beta_h(\alpha_z - 1)} h_{k(n-1)}$$

No well-defined stationary distribution



Gamma Markov Chains of the literature (3/3)

Chaining on the shape parameter

$$h_{kn}| h_{k(n-1)} \sim \text{Gamma}\left(\alpha h_{k(n-1)}, \beta\right)$$
 (29)

[Acharya et al., 2015, Schein et al., 2016]

Gamma Markov Chains of the literature (3/3)

Chaining on the shape parameter

$$h_{kn}|h_{k(n-1)} \sim \text{Gamma}\left(\alpha h_{k(n-1)}, \beta\right)$$
 (29)

[Acharya et al., 2015, Schein et al., 2016]

$$\blacktriangleright \mathbb{E}(h_{kn}|h_{k(n-1)}) = \frac{\alpha}{\beta}h_{k(n-1)}$$

No well-defined stationary distribution

Gamma Markov Chains of the literature (3/3)

Chaining on the shape parameter

$$h_{kn}| h_{k(n-1)} \sim \text{Gamma}\left(\alpha h_{k(n-1)}, \beta\right)$$
 (29)

[Acharya et al., 2015, Schein et al., 2016]

$$\blacktriangleright \mathbb{E}(h_{kn}|h_{k(n-1)}) = \frac{\alpha}{\beta}h_{k(n-1)}$$

No well-defined stationary distribution



▶ All the chains proposed in the NMF literature are based on $\mathbb{E}(h_{kn}|h_{k(n-1)}) \propto h_{k(n-1)}$

- ► All the chains proposed in the NMF literature are based on $\mathbb{E}(h_{kn}|h_{k(n-1)}) \propto h_{k(n-1)}$
- All share the same drawback : the absence of a well-defined stationary distribution

- ► All the chains proposed in the NMF literature are based on $\mathbb{E}(h_{kn}|h_{k(n-1)}) \propto h_{k(n-1)}$
- All share the same drawback : the absence of a well-defined stationary distribution
- Leads to degenerate realizations of the chain
- Difficult to interpret from a generative perspective

- ► All the chains proposed in the NMF literature are based on $\mathbb{E}(h_{kn}|h_{k(n-1)}) \propto h_{k(n-1)}$
- All share the same drawback : the absence of a well-defined stationary distribution
- Leads to degenerate realizations of the chain
- Difficult to interpret from a generative perspective
- We propose to use Markov chains with a well-defined stationary distribution

First-order autoregressive Beta-Gamma process – BGAR(1)

$$h_{k1} \sim \text{Gamma}(\alpha, \beta)$$
 (30)

$$h_{kn} = b_{kn}h_{k(n-1)} + \epsilon_{kn} \tag{31}$$

where $b_{kn} \in [0,1]$ and $\epsilon_{kn} \geq 0$ are i.i.d. r.v. such that

$$b_{kn} \sim \text{Beta}(\alpha \rho, \alpha(1-\rho))$$
 (32)

$$\epsilon_{kn} \sim \text{Gamma}(\alpha(1-\rho),\beta)$$
 (33)

[Lewis et al., 1989]

First-order autoregressive Beta-Gamma process – BGAR(1)

$$h_{k1} \sim \text{Gamma}(\alpha, \beta)$$
 (30)

$$h_{kn} = b_{kn}h_{k(n-1)} + \epsilon_{kn} \tag{31}$$

where $b_{kn} \in [0,1]$ and $\epsilon_{kn} \ge 0$ are i.i.d. r.v. such that

$$b_{kn} \sim \text{Beta}(\alpha \rho, \alpha(1-\rho))$$
 (32)

$$\epsilon_{kn} \sim \text{Gamma}(\alpha(1-\rho),\beta)$$
 (33)

[Lewis et al., 1989]

We have

$$\mathbb{E}(h_{kn}|h_{k(n-1)}) = \rho h_{k(n-1)} + \frac{\alpha(1-\rho)}{\beta}$$
(34)

First-order autoregressive Beta-Gamma process – BGAR(1)

$$h_{k1} \sim \text{Gamma}(\alpha, \beta)$$
 (30)

$$h_{kn} = b_{kn}h_{k(n-1)} + \epsilon_{kn} \tag{31}$$

where $b_{kn} \in [0,1]$ and $\epsilon_{kn} \ge 0$ are i.i.d. r.v. such that

$$b_{kn} \sim \text{Beta}(\alpha \rho, \alpha(1-\rho))$$
 (32)

$$\epsilon_{kn} \sim \text{Gamma}(\alpha(1-\rho),\beta)$$
 (33)

[Lewis et al., 1989]

We have

$$\mathbb{E}(h_{kn}|h_{k(n-1)}) = \rho h_{k(n-1)} + \frac{\alpha(1-\rho)}{\beta}$$
(34)

h_{kn} is marginally Gamma distributed

INRA-MIAT Seminar

More precisely...

 $\blacktriangleright \ \alpha$ and β control the marginal distribution

- $\blacktriangleright \ \alpha$ and β control the marginal distribution
- ρ controls the correlation between two successive values

- $\blacktriangleright \ \alpha$ and β control the marginal distribution
- ρ controls the correlation between two successive values
- ▶ $\rho \rightarrow 0$: i.i.d. random variables, $\rho \rightarrow 1$: deterministic process

- α and β control the marginal distribution
- \blacktriangleright ρ controls the correlation between two successive values
- ▶ $\rho \rightarrow 0$: i.i.d. random variables, $\rho \rightarrow 1$: deterministic process



The BGAR-NMF model

$$\underline{\mathbf{h}}_{k} \sim \mathsf{BGAR}(\rho_{k}, \alpha_{k}, \beta_{k})$$
(35)

$$v_{fn} | \mathbf{h}_n \sim \text{Poisson}([\mathbf{WH}]_{fn})$$
 (36)

The BGAR-NMF model

$$\underline{\mathbf{h}}_{k} \sim \mathsf{BGAR}(\rho_{k}, \alpha_{k}, \beta_{k})$$
(35)

$$|v_{fn}| |\mathbf{h}_n \sim \mathsf{Poisson}([\mathbf{WH}]_{fn})$$
 (36)

▶ W is left to be a deterministic variable to be estimated

The BGAR-NMF model

$$\mathbf{\underline{h}}_{k} \sim \mathsf{BGAR}(\rho_{k}, \alpha_{k}, \beta_{k}) \tag{35}$$

$$|v_{fn}| |\mathbf{h}_n \sim \mathsf{Poisson}([\mathbf{WH}]_{fn})$$
 (36)

- ▶ W is left to be a deterministic variable to be estimated
- α , β , ho are treated as fixed hyperparameters

The BGAR-NMF model

$$\underline{\mathbf{h}}_{k} \sim \mathsf{BGAR}(\rho_{k}, \alpha_{k}, \beta_{k})$$
(35)

$$|v_{fn}||\mathbf{h}_n \sim \mathsf{Poisson}([\mathbf{WH}]_{fn})$$
 (36)

- ▶ W is left to be a deterministic variable to be estimated
- α , β , ho are treated as fixed hyperparameters
- ▶ V and H define a hidden Markov model [Cappé et al., 2005]



Maximize the marginal likelihood

$$\max_{\mathbf{W}} p(\mathbf{V}; \mathbf{W}) = \int_{\mathbf{H}, \mathbf{B}} p(\mathbf{V}, \mathbf{H}, \mathbf{B}; \mathbf{W}) \mathrm{d}\mathbf{H} \mathrm{d}\mathbf{B}$$
(37)

Using ${\boldsymbol{\mathsf{B}}}$ as auxiliary variables

Maximize the marginal likelihood

$$\max_{\mathbf{W}} p(\mathbf{V}; \mathbf{W}) = \int_{\mathbf{H}, \mathbf{B}} p(\mathbf{V}, \mathbf{H}, \mathbf{B}; \mathbf{W}) \mathrm{d}\mathbf{H} \mathrm{d}\mathbf{B}$$
(37)

Using ${\boldsymbol{\mathsf{B}}}$ as auxiliary variables

 Amounts to estimating the static parameters of the HMM [Kantas et al., 2015]

Maximize the marginal likelihood

$$\max_{\mathbf{W}} p(\mathbf{V}; \mathbf{W}) = \int_{\mathbf{H}, \mathbf{B}} p(\mathbf{V}, \mathbf{H}, \mathbf{B}; \mathbf{W}) \mathrm{d}\mathbf{H} \mathrm{d}\mathbf{B}$$
(37)

Using **B** as auxiliary variables

- Amounts to estimating the static parameters of the HMM [Kantas et al., 2015]
- MCEM algorithm whose sampling step is carried out with sequential Monte Carlo (SMC)



Figure 4: Evolution of the norm of the columns of **W** w.r.t. the number of EM iterations on the NIPS dataset (11463×29).

- Method seemingly works on small dimensioned datasets
- Fails to produce exploitable results on real datasets
 - Samples of poor quality ?
 - Label switching



Figure 4: Evolution of the norm of the columns of W w.r.t. the number of EM iterations on the NIPS dataset (11463 \times 29).

- Method seemingly works on small dimensioned datasets
- Fails to produce exploitable results on real datasets
 - Samples of poor quality ?
 - Label switching
- Need for an alternative estimation paradigm in the BGAR-NMF model

Louis Filstroff

INRA-MIAT Seminar

MAP estimation in the BGAR-NMF model (1/3)

 MAP estimation amounts to the minimization of the following function

$$C(\mathbf{W}, \mathbf{H}, \mathbf{B}) = -\log p(\mathbf{H}, \mathbf{B} | \mathbf{V}; \mathbf{W})$$
(38)

$$= -\log p(\mathbf{V}|\mathbf{H};\mathbf{W}) - \log p(\mathbf{H},\mathbf{B}) + \mathrm{cst} \quad (39)$$

MAP estimation in the BGAR-NMF model (1/3)

 MAP estimation amounts to the minimization of the following function

$$C(\mathbf{W}, \mathbf{H}, \mathbf{B}) = -\log p(\mathbf{H}, \mathbf{B} | \mathbf{V}; \mathbf{W})$$
(38)

$$= -\log p(\mathbf{V}|\mathbf{H};\mathbf{W}) - \log p(\mathbf{H},\mathbf{B}) + \mathrm{cst} \quad (39)$$

We resort to a MM-based scheme. Only – log p(V|H; W) needs to be majorized. Standard scheme in the NMF literature [Lee and Seung, 2000, Févotte and Idier, 2011]

MAP estimation in the BGAR-NMF model (1/3)

 MAP estimation amounts to the minimization of the following function

$$C(\mathbf{W}, \mathbf{H}, \mathbf{B}) = -\log p(\mathbf{H}, \mathbf{B} | \mathbf{V}; \mathbf{W})$$
(38)

$$= -\log p(\mathbf{V}|\mathbf{H};\mathbf{W}) - \log p(\mathbf{H},\mathbf{B}) + \mathrm{cst} \quad (39)$$

- ► We resort to a MM-based scheme. Only log p(V|H; W) needs to be majorized. Standard scheme in the NMF literature [Lee and Seung, 2000, Févotte and Idier, 2011]
- ► For **H** and **B**, leads to order-3 polynomial equations to solve
- We have to restrict ourselves to certain values of hyperparameters


Figure 5: Evolution of $\underline{\mathbf{h}}$ w.r.t. ρ on a synthetic dataset

Comparison of all the presented models in a MAP framework

- Comparison of all the presented models in a MAP framework
- ▶ Prediction problem on the NIPS dataset. 80/10/10 split

- Comparison of all the presented models in a MAP framework
- Prediction problem on the NIPS dataset. 80/10/10 split

Method	ℓ_1 error	$\ell_2 \text{ error}$	KL error
GaP	17.14 ± 0.47	4011 ± 446	267031 ± 9859
Rate	13.07 ± 0.58	4652 ± 1907	206574 ± 12839
Rate + Aux	9.34 ± 0.29	912 ± 225	$\textbf{136412} \pm \textbf{6212}$
Shape	12.63 ± 0.3	1946 ± 240	192849 ± 5753
BGAR	$\textbf{9.3} \pm \textbf{0.17}$	$\textbf{839} \pm \textbf{119}$	138351 ± 6262

Table 4: Prediction results on the NIPS dataset

Outline



- 2 MMLE in the Gamma-Poisson model
- 3 Temporal NMF



We tackled maximum marginal likelihood estimation in semi-Bayesian NMF models

Study of two particular instances : the GaP model and the IGCN model

- Study of two particular instances : the GaP model and the IGCN model
 - Rewriting of the models free of H, which led to an expression of the marginal likelihood

- Study of two particular instances : the GaP model and the IGCN model
 - Rewriting of the models free of H, which led to an expression of the marginal likelihood
 - \blacktriangleright The expression revealed a penalty term on ${\bf W}$

- Study of two particular instances : the GaP model and the IGCN model
 - Rewriting of the models free of H, which led to an expression of the marginal likelihood
 - The expression revealed a penalty term on ${\bf W}$
 - We tackled the optimization of the likelihood with (MC)-EM algorithms

- Study of two particular instances : the GaP model and the IGCN model
 - Rewriting of the models free of H, which led to an expression of the marginal likelihood
 - \blacktriangleright The expression revealed a penalty term on ${\bf W}$
 - We tackled the optimization of the likelihood with (MC)-EM algorithms
- Study of temporal Markovian NMF models

- Study of two particular instances : the GaP model and the IGCN model
 - Rewriting of the models free of H, which led to an expression of the marginal likelihood
 - ► The expression revealed a penalty term on **W**
 - We tackled the optimization of the likelihood with (MC)-EM algorithms
- Study of temporal Markovian NMF models
 - Thorough review of the literature, which revealed that all considered Markov chains shared the same drawback

- Study of two particular instances : the GaP model and the IGCN model
 - Rewriting of the models free of H, which led to an expression of the marginal likelihood
 - \blacktriangleright The expression revealed a penalty term on ${\bf W}$
 - We tackled the optimization of the likelihood with (MC)-EM algorithms
- Study of temporal Markovian NMF models
 - Thorough review of the literature, which revealed that all considered Markov chains shared the same drawback
 - We proposed the use of an overlooked model from the time series literature, BGAR(1)

- Study of two particular instances : the GaP model and the IGCN model
 - Rewriting of the models free of H, which led to an expression of the marginal likelihood
 - ► The expression revealed a penalty term on **W**
 - We tackled the optimization of the likelihood with (MC)-EM algorithms
- Study of temporal Markovian NMF models
 - Thorough review of the literature, which revealed that all considered Markov chains shared the same drawback
 - We proposed the use of an overlooked model from the time series literature, BGAR(1)
 - MMLE tackled with SMC : not satisfying

- Study of two particular instances : the GaP model and the IGCN model
 - ► Rewriting of the models free of **H**, which led to an expression of the marginal likelihood
 - The expression revealed a penalty term on ${f W}$
 - We tackled the optimization of the likelihood with (MC)-EM algorithms
- Study of temporal Markovian NMF models
 - Thorough review of the literature, which revealed that all considered Markov chains shared the same drawback
 - We proposed the use of an overlooked model from the time series literature, BGAR(1)
 - MMLE tackled with SMC : not satisfying
 - MAP estimation tackled with an MM-based algorithm showed better performance on prediction tasks

Model aspects

Break out of conjugate prior distributions

- Break out of conjugate prior distributions
- Break out of composite models

- Break out of conjugate prior distributions
- Break out of composite models
- Carry out the analysis in a family of distributions

- Break out of conjugate prior distributions
- Break out of composite models
- Carry out the analysis in a family of distributions
- Other Markov chains with stationary Gamma distribution

Model aspects

- Break out of conjugate prior distributions
- Break out of composite models
- Carry out the analysis in a family of distributions
- Other Markov chains with stationary Gamma distribution

Optimization aspects

- Break out of conjugate prior distributions
- Break out of composite models
- Carry out the analysis in a family of distributions
- Other Markov chains with stationary Gamma distribution
- Optimization aspects
 - Break out of MC-EM : alternative schemes ?

Model aspects

- Break out of conjugate prior distributions
- Break out of composite models
- Carry out the analysis in a family of distributions
- Other Markov chains with stationary Gamma distribution

Optimization aspects

- Break out of MC-EM : alternative schemes ?
- Direct optimization of the likelihood

Associated publications

- Filstroff, L., Lumbreras, A., and Févotte, C. (2018).
 Closed-form Marginal Likelihood in Gamma-Poisson Matrix Factorization. In Proceedings of the International Conference of Machine Learning (ICML).
- Filstroff, L., and others (2019) Temporal Non-negative Matrix Factorization with Gamma Markov Chains. In preparation. IEEE Transactions on Signal Processing.
- Xia, R., Tan, V.Y.F., <u>Filstroff, L.</u>, and Févotte, C. (2019).
 A Ranking Model Motivated by Nonnegative Matrix Factorization with Applications to Tennis Tournaments.

In Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML-PKDD).

Lumbreras, A., <u>Filstroff, L.</u>, and Févotte, C. (2018).
 Bayesian mean-parameterized nonnegative binary matrix factorization.
 In revision. Data Mining and Knowledge Discovery.

Thank you for your attention.

Additional slides



- Variants of the EM algorithm
- Links between PF and LDA
- NB and NM distributions
- Examples of self-reg.
- The IGCN model
- MM algorithm for BGAR-NMF
- MAP estimation

The Majorization-Minimization (MM) framework

- Majorize the function f by an auxiliary function g
- Minimize g instead
- g is such that $g(x; \tilde{x}) \ge f(x)$ and $g(\tilde{x}; \tilde{x}) = f(\tilde{x})$



Variants of the EM algorithm

- ► MCEM : converges as N_{iter} → +∞ AND J → +∞ At iteration t, uses J samples for maximization
- SAEM [Delyon et al., 1999, Kuhn and Lavielle, 2004] Converges as N_{iter} → +∞

$$\hat{Q}_t(\boldsymbol{\theta}) = (1 - \gamma_t)\hat{Q}_{t-1}(\boldsymbol{\theta}) + \frac{\gamma_t}{J_t}\sum_j \log p(\mathbf{V}, \mathbf{Z}^{(j)}; \boldsymbol{\theta})$$
(40)

 $(\gamma_l)_{l\geq 1}$ is a sequence of positive step sizes decreasing to 0 At iteration *t*, uses all past samples for maximization

On-line EM [Cappé et al., 2005]
 Stochastic approximation "one sample at a time"

◀ Go back

Equivalence between Poisson Factorization and LDA

$$v_{fn} \sim \text{Poisson}([\mathbf{WH}]_{fn})$$
 (41)

is equivalent to

v

$$L_n \sim \text{Poisson}\left(\sum_f w_{fk} \sum_k h_{kn}\right)$$
(42)
$$u_n | L_n \sim \text{Mult}(L_n, \lambda_n)$$
(43)

with

$$\lambda_{fn} = \frac{\sum_{k} w_{fk} h_{kn}}{\sum_{f} w_{fk} \sum_{k} h_{kn}}$$
(44)

Imposing $\sum_{f} w_{fk} = 1$ and $\sum_{k} h_{kn} = 1$ leads to the observation model of LDA (difference in budget)

◀ Go back

NB and NM distributions

Negative Binomial (NB) distribution

With $\alpha > 0$ and $p \in [0, 1]$. For all $c \in \mathbb{N}$

$$\mathbb{P}(X=c) = \frac{\Gamma(\alpha+c)}{\Gamma(\alpha)c!} (1-p)^{\alpha} p^{c}$$
(45)

$$\begin{cases} \lambda \sim \mathcal{G}(\alpha, \beta) \\ X | \lambda \sim \mathsf{Poisson}(\lambda) \end{cases} \Leftrightarrow X \sim \mathsf{NB}\left(\alpha, \frac{1}{\beta + 1}\right) \tag{46}$$

Negative Multinomial (NM) distribution

With $\alpha > 0$ and $p_f \in [0, 1]$ and $\sum_f p_f \le 1$. For all $c_1, \dots, c_F \in \mathbb{N}^F$ $\mathbb{P}(X_1 = c_1, \dots, X_n = c_F) = \frac{\Gamma(\alpha + \sum_f c_f)}{\Gamma(\alpha) \prod_f c_f!} p_0^{\alpha} \prod_f p_f^{c_f}$ (47)

Examples of self-regularization



The IGCN model and objective

The IGCN model

$$h_{kn} \sim \mathcal{IG}(\alpha_k, \beta_k)$$
(48)

$$\kappa_{fn} | \mathbf{h}_n \sim \mathcal{CN}(0, [\mathbf{WH}]_{fn})$$
(49)

- CN denotes the complex normal distribution
- Standard STFT model in audio signal processing ("Gaussian composite model")
 [Févotte et al., 2009, Hoffman et al., 2010]
- Inverse Gamma prior for practical reasons : conjugacy with the normal distribution of known mean
- MMLE amounts to

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}) = -\log \int p(\mathbf{X}|\mathbf{H}; \mathbf{W}) p(\mathbf{H}) d\mathbf{H}$$
(50)

Marginalization of the activation coefficients

- Same augmentation with variables C as in the GaP model (superposition property of the normal distribution)
- IGCN can be written as a composite complex Student's t model
- Does not lead to a closed-form expression of the likelihood, because we end up with an intractable integral

$$p(\mathbf{X}; \mathbf{W}) = \int_{\mathbf{C} \in \mathcal{C}} p(\mathbf{C}; \mathbf{W}) \mathrm{d}\mathbf{C} = \int_{\mathbf{C} \in \mathcal{C}} \prod_{k,n} p(\mathbf{c}_{kn}; \mathbf{w}_k) \mathrm{d}\mathbf{C} \quad (51)$$

Still exhibits a term of the form ∑_{f,k} log(w_{fk}), which promotes "local" sparsity

Optimization of the marginal likelihood

- In [Dikmen and Févotte, 2011], the optimization was tackled with a variational algorithm
- We have proposed three novel EM algorithms based on three choices of the latent variables
- E-step based on a shared Gibbs sampling procedure
- M-step in closed form (EM-CH), or tackled with MM-based schemes (EM-H, EM-C)
Experimental results

- Real audio decomposition task on a short piano sequence
- Performance compared with the standard IS-NMF
- No obvious advantage in this case (similar audio accuracy, dictionary not especially sparse, computationally prohibitive)
- Conceptually interesting, but shows the limitations of the method

▲ Go back

MM : Constraints

In BGAR, we have

$$h_{kn} = b_{kn}h_{k(n-1)} + \epsilon_{kn} \tag{52}$$

$$h_{k(n+1)} = b_{k(n+1)}h_{kn} + \epsilon_{k(n+1)}$$
(53)

Leads to

$$b_{kn}h_{k(n-1)} \le h_{kn} \le \frac{h_{k(n+1)}}{b_{k(n+1)}}$$
 (54)

 and

$$0 \le b_{kn} \le \min\left(1, \frac{h_{kn}}{h_{k(n-1)}}\right) \tag{55}$$

◀ Go back

MM : Hyperparameter constraints



Figure 6: Hyperparameter values of the parameters α_k and ρ_k ensuring a well-posed MAP estimation in the BGAR-NMF model

MAP Estimation

Objective function

$$C(\mathbf{W}, \mathbf{H}) = \underbrace{-\log p(\mathbf{V}|\mathbf{H}; \mathbf{W})}_{\text{Majorize}} - \log p(\mathbf{H})$$
(56)

Standard majorization in the Poisson case

$$G(\mathbf{H}; \tilde{\mathbf{H}}) = -\sum_{k,n} p_{kn} \log h_{kn} + \sum_{k,n} q_k h_{kn}$$
(57)

- Rate : order-2 polynomials
- Rate + Aux : order-1 polynomials
- Shape : Newton's method
- BGAR : order-3 polynomials

References I

 [Acharya et al., 2015] Acharya, A., Ghosh, J., and Zhou, M. (2015).
 Nonparametric Bayesian Factor Analysis for Dynamic Count Matrices.
 In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS), pages 1–9.

[Berger et al., 1999] Berger, J. O., Liseo, B., and Wolpert, R. L. (1999). Integrated Likelihood Methods for Eliminating Nuisance Parameters. *Statistical Science*, 14(1):1–28.

[Bioucas-Dias et al., 2012] Bioucas-Dias, J. M., Plaza, A., Dobigeon, N., Parente, M., Du, Q., Gader, P., and Chanussot, J. (2012).
Hyperspectral Unmixing Overview: Geometrical, Statistical, and Sparse Regression-Based Approaches.
IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, 5(2):354–379.

[Blei et al., 2003] Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent Dirichlet Allocation. Journal of Machine Learning Research, 3(Jan):993–1022.

[Buntine and Jakulin, 2006] Buntine, W. and Jakulin, A. (2006). Discrete component analysis.

In Subspace, Latent Structure and Feature Selection, pages 1–33. Springer.

References II

[Candès et al., 2008] Candès, E. J., Wakin, M. B., and Boyd, S. P. (2008). Enhancing sparsity by reweighted ℓ1 minimization. Journal of Fourier Analysis and Applications, 14(5-6):877–905.

[Canny, 2004] Canny, J. (2004).

GaP: A Factor Model for Discrete Data. In Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval, pages 122–129.

[Cappé et al., 2005] Cappé, O., Moulines, E., and Rydén, T. (2005). Inference in Hidden Markov Models. Springer.

[Cemgil, 2009] Cemgil, A. T. (2009). Bayesian Inference for Nonnegative Matrix Factorisation Models. Computational Intelligence and Neuroscience, (Article ID 785152).

[Cemgil and Dikmen, 2007] Cemgil, A. T. and Dikmen, O. (2007). Conjugate Gamma Markov random fields for modelling nonstationary sources. In Proceedings of the International Conference on Independent Component Analysis and Signal Separation (ICA), pages 697–705.

[Delyon et al., 1999] Delyon, B., Lavielle, M., and Moulines, E. (1999). Convergence of a stochastic approximation version of the em algorithm. *The Annals of Statistics*, 27(1):94–128.

References III

[Dempster et al., 1977] Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. Journal of the Royal Statistical Society, Series B (Methodological), 39(1):1–38.

[Dikmen and Févotte, 2011] Dikmen, O. and Févotte, C. (2011). Nonnegative dictionary learning in the exponential noise model for adaptive music signal representation.

In Advances in Neural Information Processing Systems (NIPS), pages 2267–2275.

[Dikmen and Févotte, 2012] Dikmen, O. and Févotte, C. (2012).

 ${\sf Maximum}$ marginal likelihood estimation for nonnegative dictionary learning in the Gamma-Poisson model.

IEEE Transactions on Signal Processing, 60(10):5163–5175.

[Févotte, 2011] Févotte, C. (2011).

Majorization-Minimization Algorithm for Smooth Itakura-Saito Nonnegative Matrix Factorization.

In Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1980–1983.

[Févotte et al., 2009] Févotte, C., Bertin, N., and Durrieu, J.-L. (2009). Nonnegative matrix factorization with the itakura-saito divergence: With application to music analysis. *Neural Computation*, 21(3):793–830.

References IV

[Févotte and Idier, 2011] Févotte, C. and Idier, J. (2011). Algorithms for nonnegative matrix factorization with the β -divergence. *Neural Computation*, 23(9):2421–2456.

[Gopalan et al., 2015] Gopalan, P., Hofman, J. M., and Blei, D. M. (2015).
 Scalable recommendation with hierarchical poisson factorization.
 In Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI), pages 326–335.

[Hoffman et al., 2010] Hoffman, M. D., Blei, D. M., and Cook, P. R. (2010). Bayesian Nonparametric Matrix Factorization for Recorded Music. In Proceedings of the International Conference on Machine Learning (ICML), pages 439–446.

[Hotelling, 1933] Hotelling, H. (1933).

Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology*, 24(6):417–441.

[Kantas et al., 2015] Kantas, N., Doucet, A., Singh, S. S., Maciejowski, J., and Chopin, N. (2015).
On Particle Methods for Parameter Estimation in State-space Models. *Statistical Science*, 30(3):328–351.

References V

[Kuhn and Lavielle, 2004] Kuhn, E. and Lavielle, M. (2004). Coupling a stochastic approximation version of em with an mcmc procedure. ESAIM: Probability and Statistics, 8:115–131.

[Lee and Seung, 1999] Lee, D. D. and Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791.

[Lee and Seung, 2000] Lee, D. D. and Seung, H. S. (2000). Algorithms for Non-negative Matrix Factorization. In Advances in Neural Information Processing Systems (NIPS), pages 556–562.

[Lewis et al., 1989] Lewis, P., McKenzie, E., and Hugus, D. K. (1989). Gamma processes. Communications in Statistics. Stochastic Models, 5(1):1–30.

[Moulines et al., 1997] Moulines, E., Cardoso, J.-F., and Gassiat, E. (1997). Maximum likelihood for blind separation and deconvolution of noisy signals using mixture models.

In Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 3617–3620.

References VI

[Paatero and Tapper, 1994] Paatero, P. and Tapper, U. (1994).

Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(2):111–126.

[Pearson, 1901] Pearson, K. (1901).

On lines and planes of closest fit to systems of points in space.

The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science, 2(11):559–572.

[Schein et al., 2016] Schein, A., Wallach, H. M., and Zhou, M. (2016). Poisson-Gamma Dynamical Systems.

In Advances in Neural Information Processing Systems (NIPS), pages 5005–5013.

[Smaragdis and Brown, 2003] Smaragdis, P. and Brown, J. C. (2003). Non-negative matrix factorization for polyphonic music transcription. In Proceedings of the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA), pages 177–180.

[Virtanen, 2007] Virtanen, T. (2007).

Monaural Sound Source Separation by NonnegativeMatrix Factorization With Temporal Continuityand Sparseness Criteria.

IEEE Transactions on Audio, Speech, and Language Processing, 15(3):1066–1074.

[Wei and Tanner, 1990] Wei, G. C. G. and Tanner, M. A. (1990). A Monte Carlo Implementation of the EM Algorithm and the Poor Man's Data Augmentation Algorithms.

Journal of the American Statistical Association, 85(411):699–704.

 [Xu et al., 2003] Xu, W., Liu, X., and Gong, Y. (2003).
 Document Clustering Based On Non-negative MatrixFactorization.
 In Proceedings of the International ACM SIGIR Conference on Research and Development in Informaion Retrieval, pages 267–273.