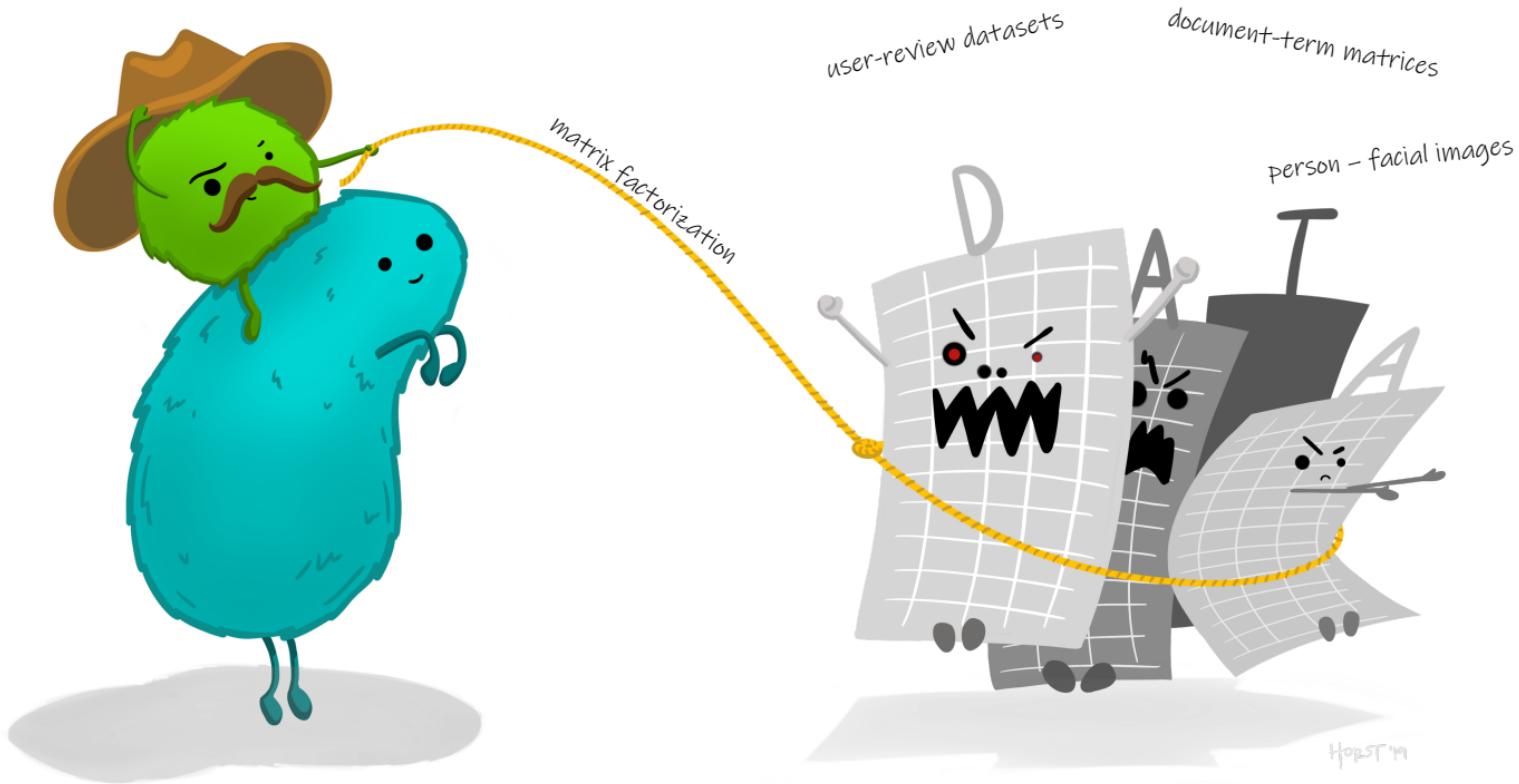




(Bayesian) non-negative matrix factorization

Greta Gašparac
Mentor: Assist. Prof. dr. Jana Faganeli Pucer



1.

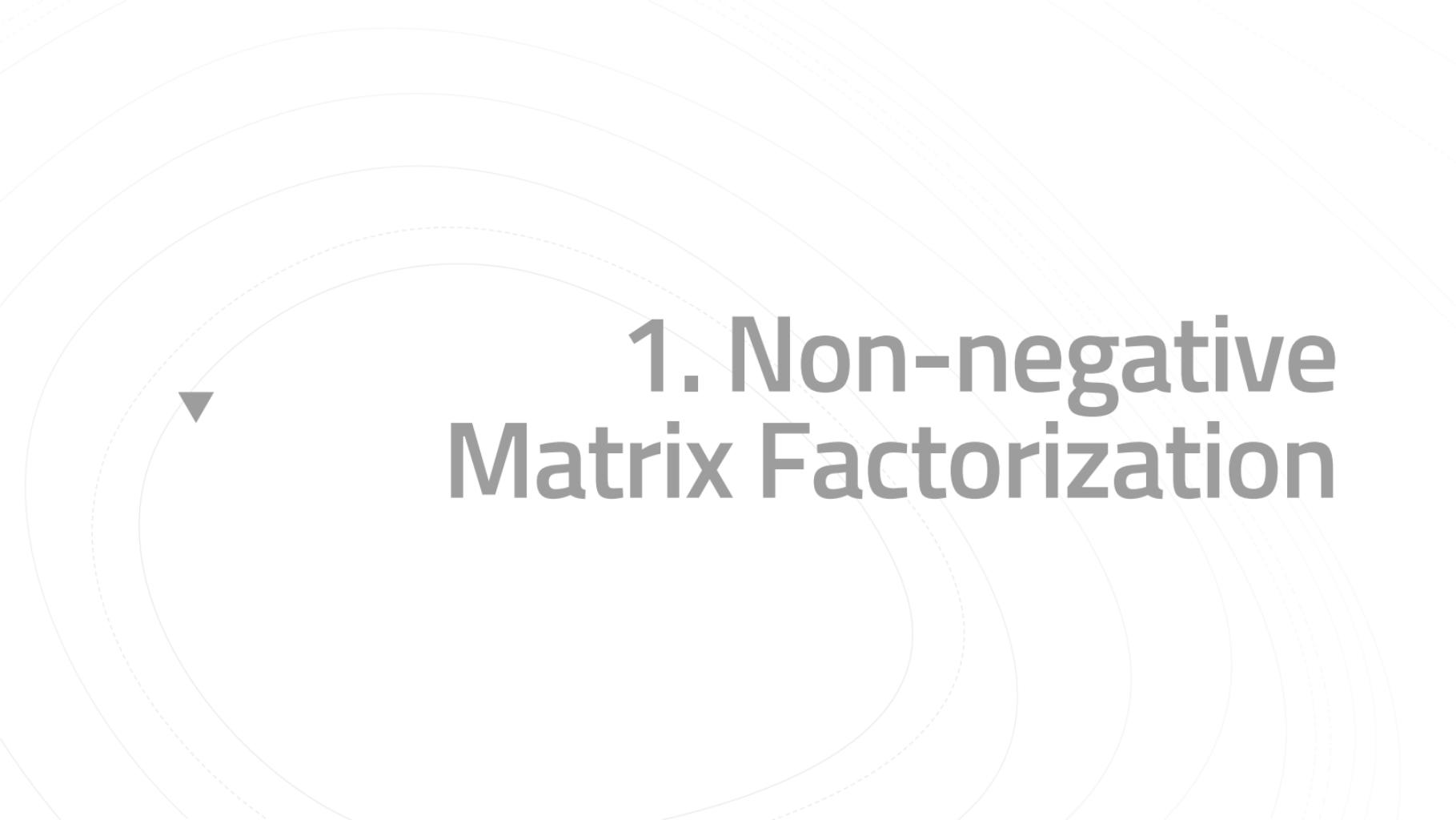
Non-negative Matrix Factorization

2.

An Overview of Bayesian Statistics & Gibbs Sampling

3.

Bayesian Non-negative Matrix Factorization



1. Non-negative Matrix Factorization

Non-negative matrix factorization (NMF)

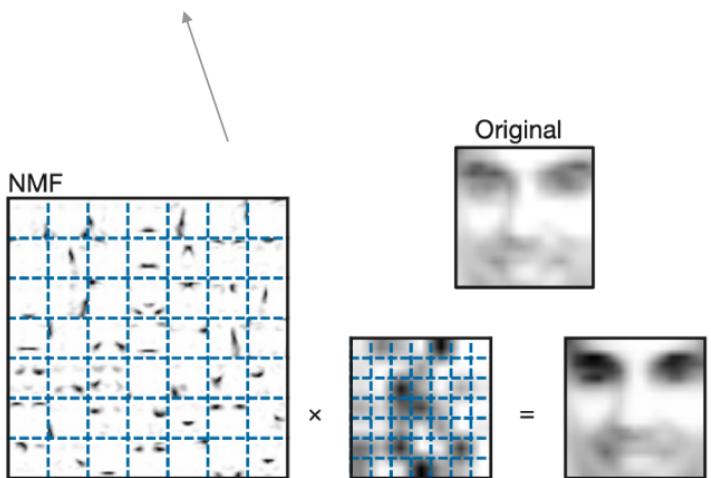
- Belongs to the family of linear dimensionality reduction techniques based on matrix factorization.
- Popularised by Lee and Seung in 1999.
- Approximates our original matrix \mathbf{X} with lower-rank matrices \mathbf{W} and \mathbf{H} , which are element-wise non-negative.
- The non-negativity constraint provides more intuitive factors.

$$\begin{matrix} u & & l \\ \left[\begin{array}{c|c} & \end{array} \right] & \approx & \left[\begin{array}{c|c} w_u & \end{array} \right] \times \left[\begin{array}{c|c} & h_l \end{array} \right] \\ X & & W & H \\ \left[\begin{array}{c|c} & \end{array} \right] & & \left[\begin{array}{c|c} & \end{array} \right] & \left[\begin{array}{c|c} & \end{array} \right] \\ n \times p & & n \times r & r \times p \end{matrix}$$

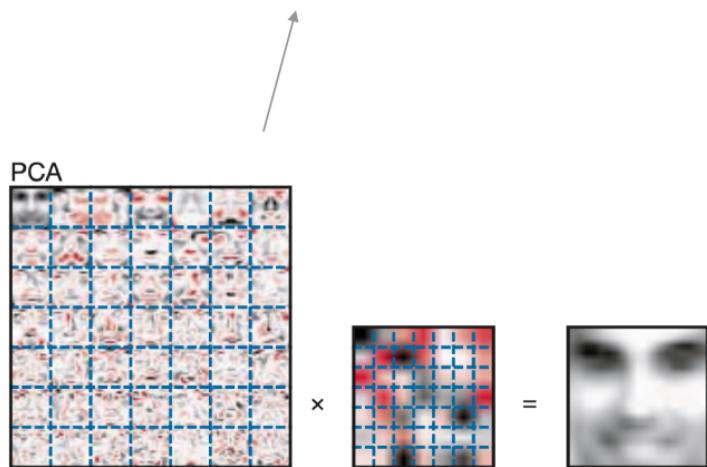
Example

$$\begin{bmatrix} & & \\ & \vdots & \\ X & & \end{bmatrix}_{n \times p} \approx \begin{bmatrix} w_1 & \dots & w_r \\ \vdots & & \vdots \\ W & & \end{bmatrix}_{n \times r} \times \begin{bmatrix} h_1 & \dots & h_p \\ \vdots & & \vdots \\ H & & \end{bmatrix}_{r \times p}$$

localized features, parts-based representation



"eigenfaces", resemble distorted faces from the database



What is NMF?

$$\begin{matrix} & & i \\ & u & \end{matrix} \quad \boxed{X} \quad \begin{matrix} & & i \\ & u & \end{matrix} \approx \begin{matrix} & & i \\ & w_u & \end{matrix} \quad \boxed{W} \quad \times \quad \begin{matrix} & & i \\ & h_i & \end{matrix} \quad \boxed{H}$$

$n \times p$ $n \times r$ $r \times p$

What makes a good latent representation?

Pick a measure to asses the quality of the approximation. Common choice: Frobenius norm.

$$\|X - WH\|_F^2 = \sum_{i,j} (X - WH)_{i,j}^2$$

Optimisation problem for a rank r factorization:

$$\min_{W \geq 0, H \geq 0} \|X - WH\|_F^2 + \text{---}$$

Problems



NMF is ill-posed.

- No guarantee there exists a single best decomposition.
- Results depend on the initialization.
- *Solutions:*
regularization terms in the objective function,
multiple reruns.

Choice of r .

- Trade-off.
- *Approaches:*
trial & error,
SVD decomposition,
domain expert knowledge.

NMF is NP-hard.

- Non-convex optimization problem.

NMF algorithms



Initialization

- Randomly
- Clustering
- SVD
- Use columns of X

Stopping criteria

- Value of objective function
- Time
- Maximum number of iterations
- Combination

Algorithm 1: Two-Block Coordinate Descent

1 Initialize matrices $\mathbf{W}^{(0)} \geq 0$ and $\mathbf{H}^{(0)} \geq 0$.

2 **for** $i = 1, 2, \dots$, *stopping criteria* **do**

3 $\mathbf{W}^{(i)} = \text{update}(\mathbf{X}, \mathbf{H}^{(i-1)}, \mathbf{W}^{(i-1)})$

4 $\mathbf{H}^{(i)} = \text{update}(\mathbf{X}, \mathbf{H}^{(i-1)}, \mathbf{W}^{(i)})$

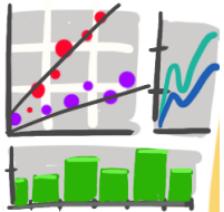
Multiplicative update rule proposed by Lee and Seung:

$$\mathbf{W} \leftarrow \mathbf{W} \circ \frac{\mathbf{X}\mathbf{H}^T}{\mathbf{W}\mathbf{H}\mathbf{H}^T},$$

$$\mathbf{H} \leftarrow \mathbf{H} \circ \frac{\mathbf{W}^T\mathbf{X}}{\mathbf{W}^T\mathbf{W}\mathbf{H}}$$

Applications

- Financial data mining



- Natural language processing



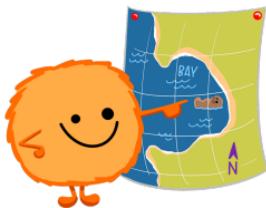
- Pattern detection in bioinformatics



- Recommender systems:
The Netflix prize



- Environmetrics





2. An Overview of Bayesian Statistics and Gibbs sampling

Bayesian Modeling

BAYES' THEOREM



$$p(\theta|y) = \frac{\underbrace{p(y|\theta) \cdot p(\theta)}_{\text{LIKELIHOOD PRIOR}}}{\underbrace{p(y)}_{\text{POSTERIOR}}} \propto p(y|\theta) \cdot p(\theta)$$

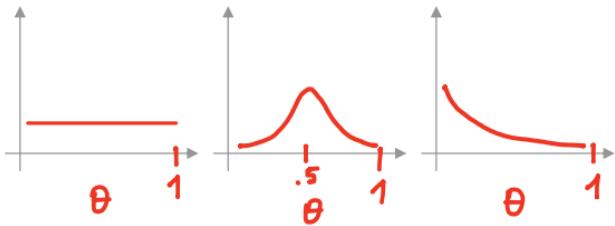
$\int \dots$

Coin toss example



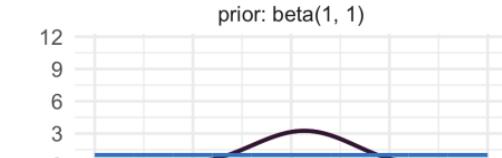
Bayesian inference steps:

1. Specify prior.
2. Identify the observed data. HTTHHTHHHHHTTHT
3. Construct a probabilistic model to represent the data.

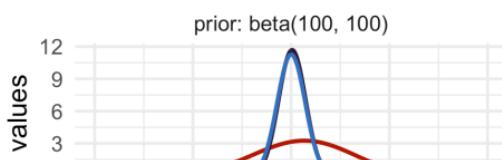


$$p(\theta|y) = \frac{p(y|\theta) \cdot p(\theta)}{p(y)} \propto p(y|\theta) \cdot p(\theta)$$

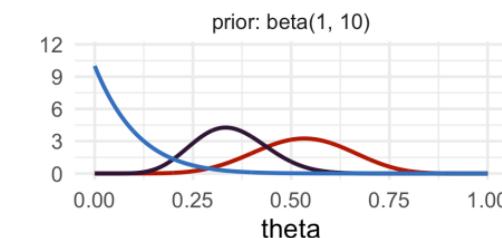
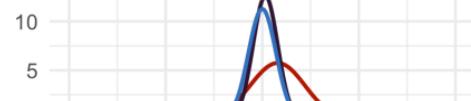
N = 15



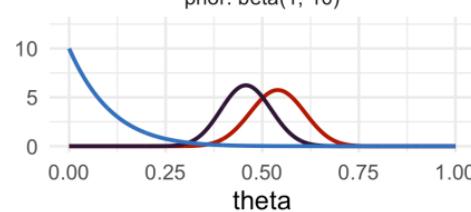
N = 50



prior: beta(100, 100)



prior: beta(1, 10)



type
— (scaled) likelihood — posterior — prior

type
— (scaled) likelihood — posterior — prior

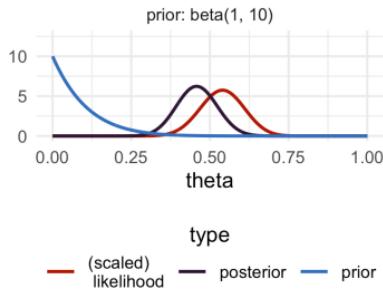
How do we get the posterior?

- Analytically

Problem: most of the models do not have nice conjugate posteriors!



$$p(\theta|y) = \frac{p(y|\theta) \cdot p(\theta)}{p(y)} \propto p(y|\theta) \cdot p(\theta)$$



- Using MCMC methods



Gibbs Sampling

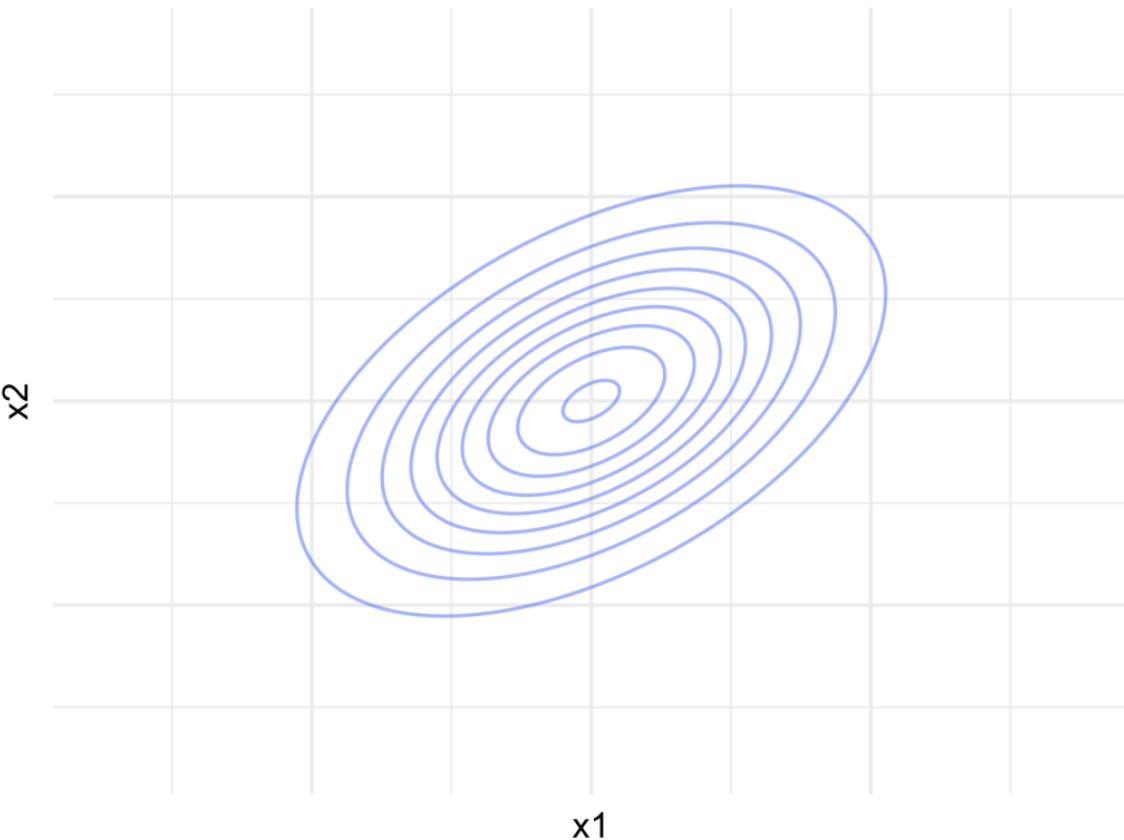
- Specific case of Metropolis-Hastings where the proposals are always accepted.
- Revolutionized the field of Bayesian statistics in the 80s and the 90s as it was the basis for the first generation of the Bayesian stats tools, such as BUGS, WinBUGS, JAGS.
- *Main idea behind:*
a joint distribution $p(x_1, x_2, \dots, x_n)$ can be characterized by its complete set of full conditional distributions:
 $p(x_1 | x_2, x_3, \dots, x_n), p(x_2 | x_1, x_3, \dots, x_n), \dots, p(x_n | x_1, x_2, \dots, x_{n-1})$.

Gibbs Sampling: bivariate normal example

$$\mathbf{X} = \begin{pmatrix} x_A \\ x_B \end{pmatrix} \sim N \left(\begin{pmatrix} \mu_A \\ \mu_B \end{pmatrix}, \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix} \right)$$

$$X_A | X_B \sim N(\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(x_B - \mu_B), \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA})$$

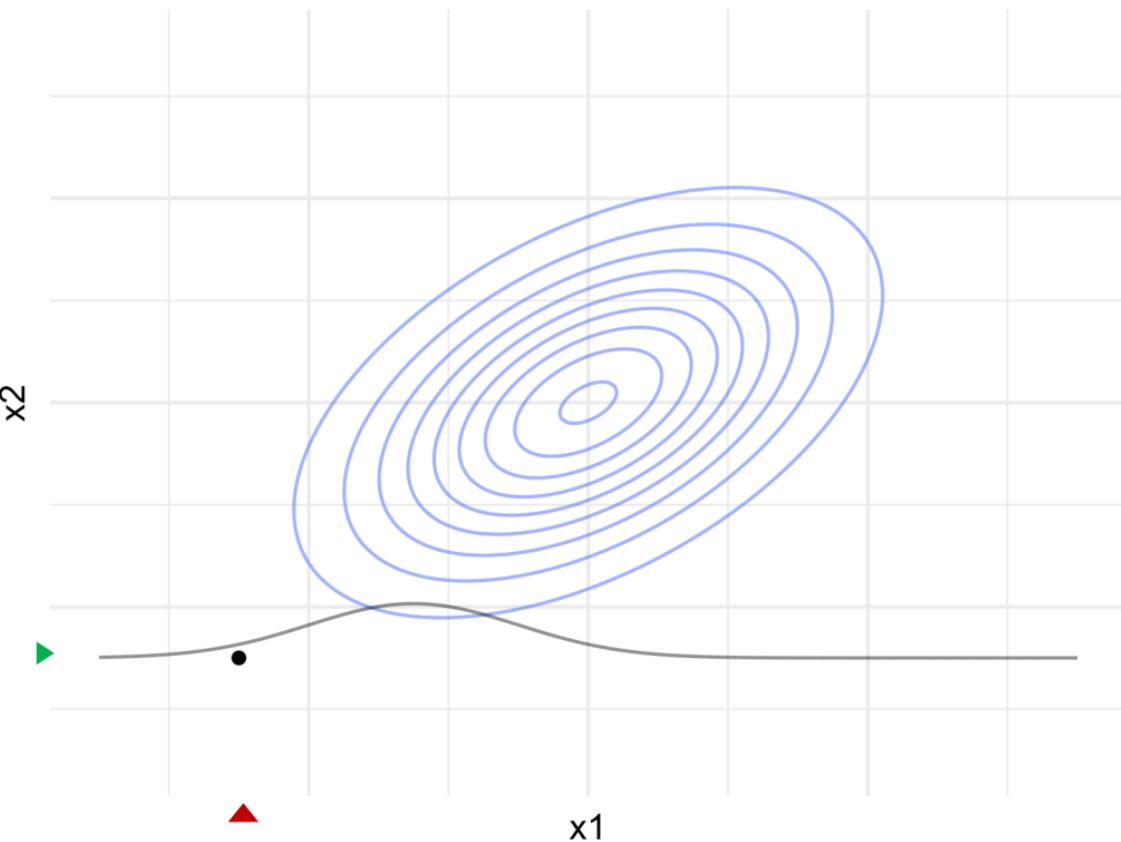
$$p(x_1, x_2)$$



Gibbs Sampling: bivariate normal example

$$x_0 = (\textcolor{red}{\blacktriangle}, \textcolor{green}{\triangleright})$$

$$i = 1, x_1 | x_2$$

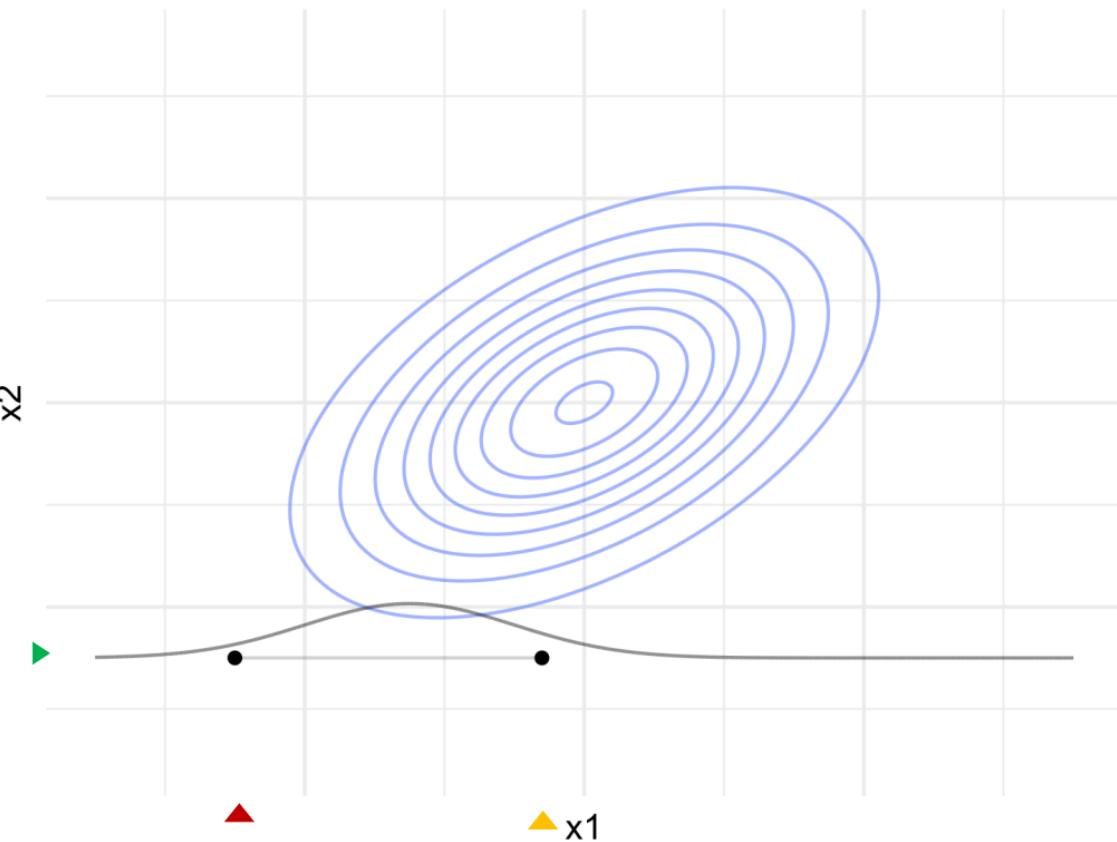


Gibbs Sampling: bivariate normal example

$$x_0 = (\textcolor{red}{\blacktriangle}, \textcolor{green}{\triangleright})$$

$$x_1 = (\textcolor{yellow}{\blacktriangle}, \textcolor{black}{\square})$$

$$i = 1, x_1 | x_2$$

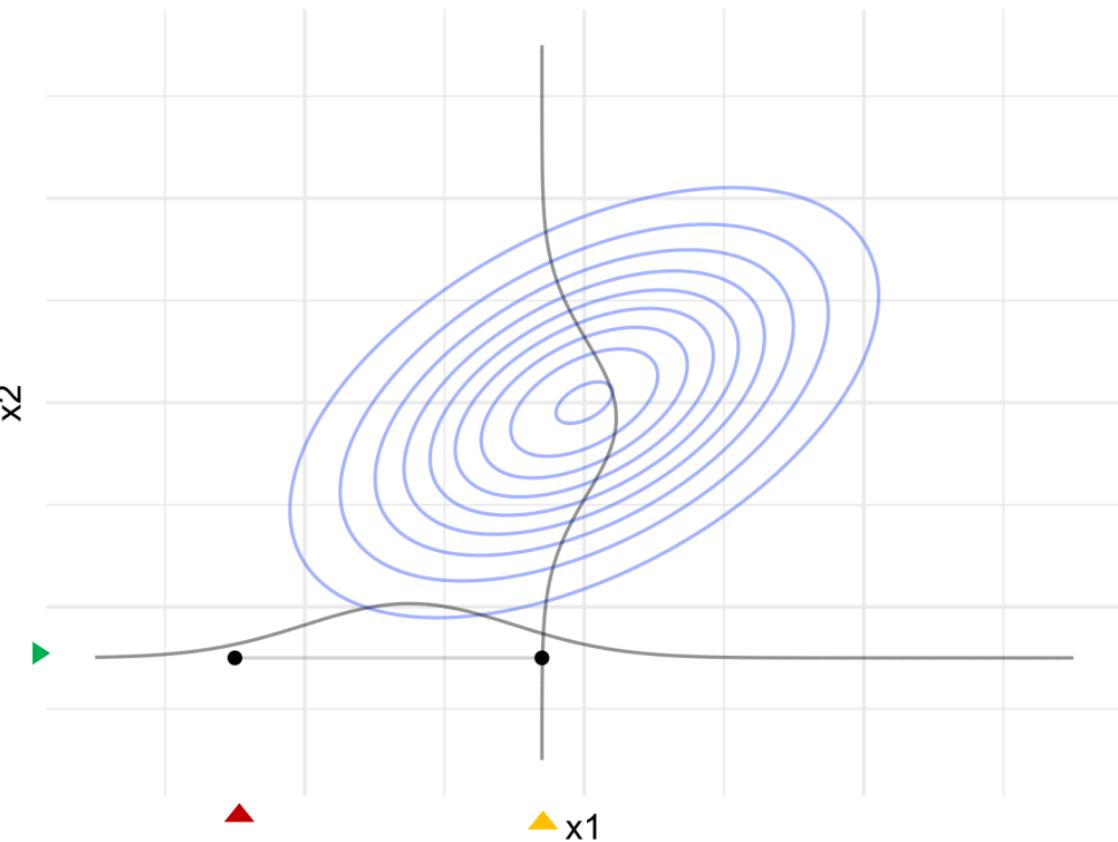


Gibbs Sampling: bivariate normal example

$$x_0 = (\textcolor{red}{\blacktriangle}, \textcolor{green}{\triangleright})$$

$$x_1 = (\textcolor{yellow}{\blacktriangle}, \textcolor{black}{\square})$$

$$i = 1, x_2 | x_1$$

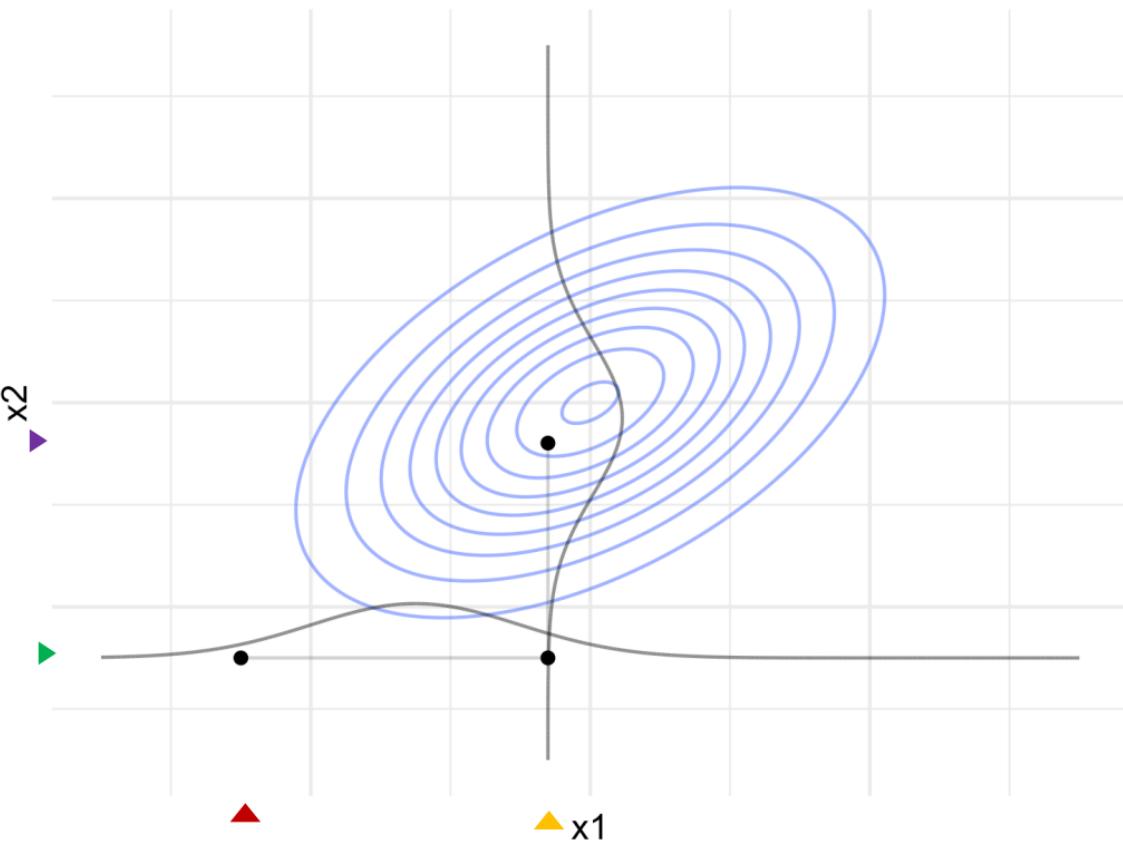


Gibbs Sampling: bivariate normal example

$$x_0 = (\textcolor{red}{\blacktriangle}, \textcolor{green}{\triangleright})$$

$$x_1 = (\textcolor{yellow}{\blacktriangle}, \textcolor{purple}{\triangleright})$$

$$i = 1, x_2 | x_1$$



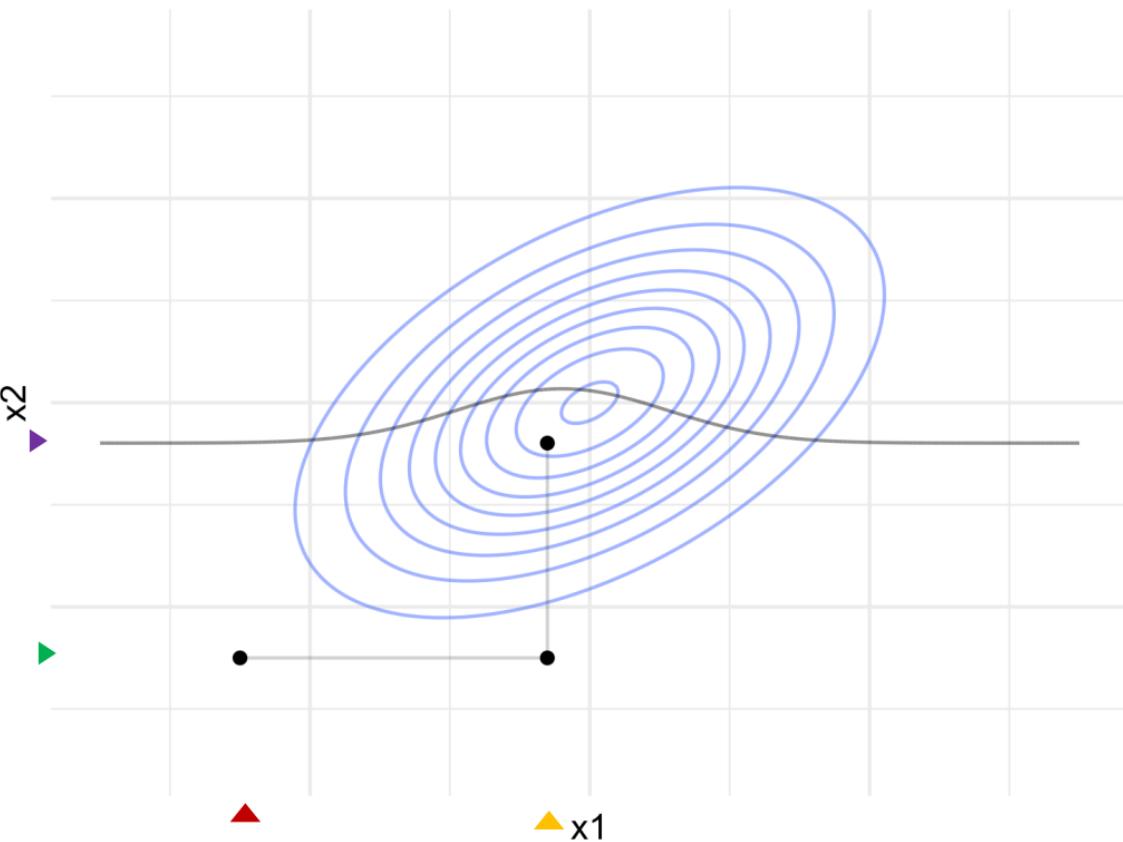
Gibbs Sampling: bivariate normal example

$x_0 = (\textcolor{red}{\blacktriangle}, \textcolor{green}{\triangleright})$

$x_1 = (\textcolor{yellow}{\blacktriangle}, \textcolor{purple}{\triangleright})$

$x_2 = (\textcolor{black}{\circ}, \textcolor{black}{\circ})$

$i = 2, x_1 | x_2$



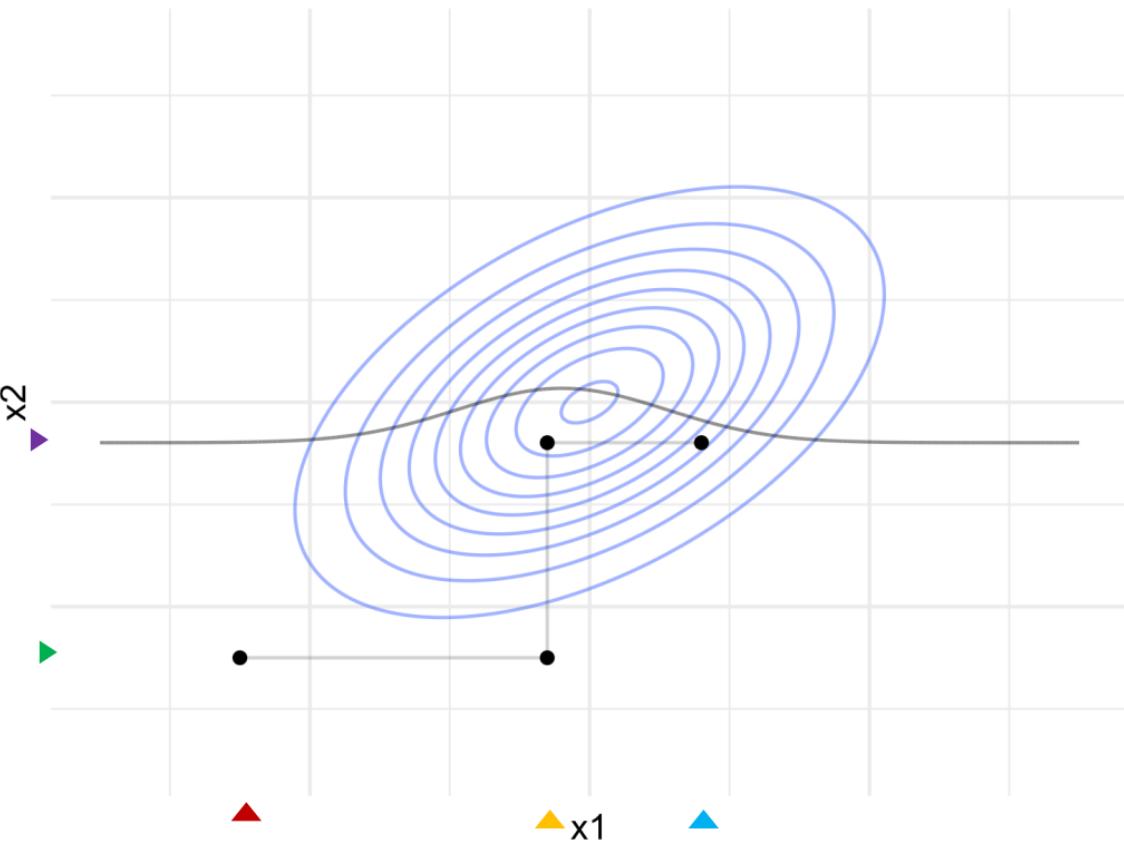
Gibbs Sampling: bivariate normal example

$x_0 = (\textcolor{red}{\blacktriangle}, \textcolor{green}{\triangleright})$

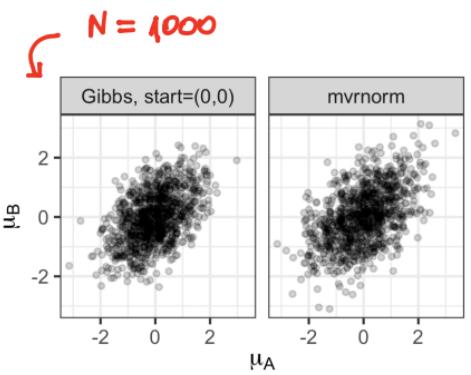
$x_1 = (\textcolor{yellow}{\blacktriangle}, \textcolor{purple}{\triangleright})$

$x_2 = (\textcolor{blue}{\blacktriangle}, \textcolor{black}{\triangleright})$

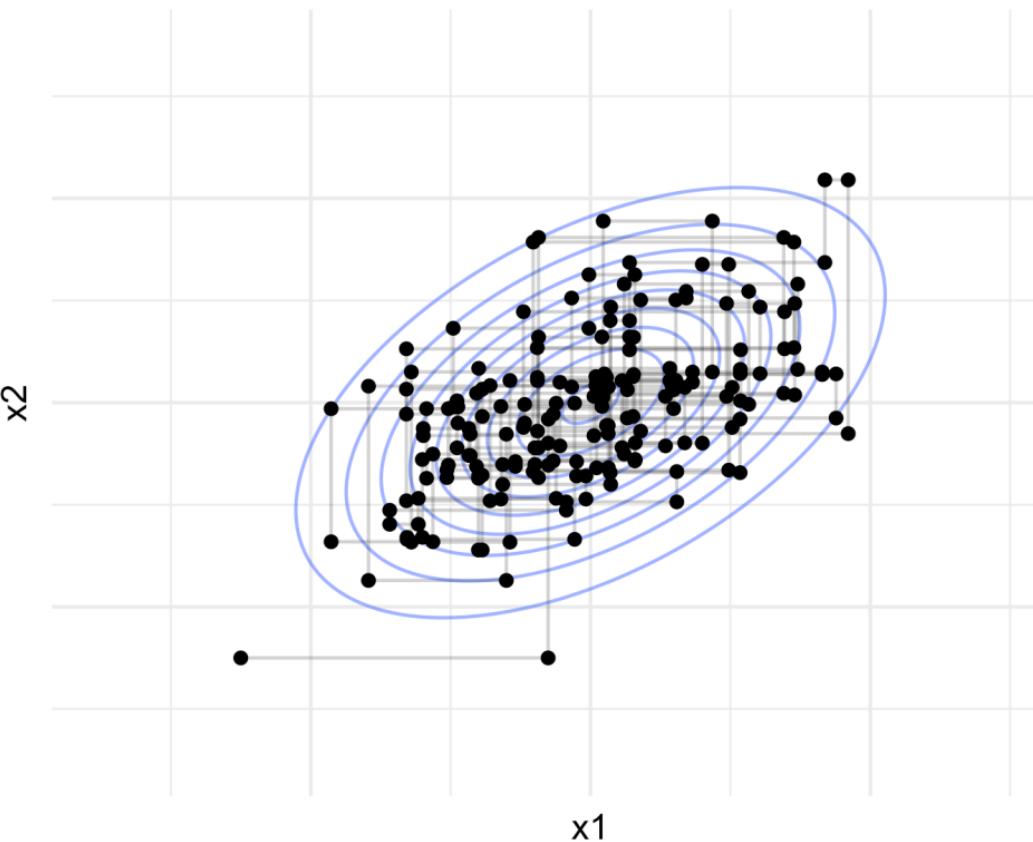
$$i = 2, x_1 | x_2$$



Gibbs Sampling: bivariate normal example



i = 100



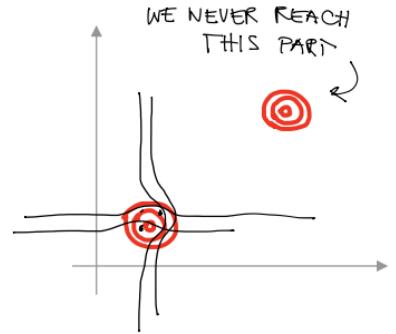
Gibbs Sampling

- We use Gibbs sampling when the joint posterior distribution is difficult to sample from directly and full conditional distributions of the parameters are known.

Algorithm 4: Gibbs sampler

- 1 Initialize $\mathbf{x}^{(0)}$.
- 2 **for** $i = 1:N$ **do**
- 3 $x_1^{(i)} \sim p(x_1 | x_2^{(i-1)}, x_3^{(i-1)}, \dots, x_n^{(i-1)})$
- 4 $x_2^{(i)} \sim p(x_2 | x_1^{(i)}, x_3^{(i-1)}, \dots, x_n^{(i-1)})$
- 5 \vdots
- 5 $x_n^{(i)} \sim p(x_n | x_1^{(i)}, x_2^{(i)}, \dots, x_{n-1}^{(i)})$
- 6 **Output:** $\{\mathbf{x}^{(i)}\}_{i=1}^N$

- Specific failure cases, diagnostics is important: check traceplots, autocorrelation, effective sample size.



3. Bayesian Non-negative Matrix Factorization



Bayesian NMF

- Bayesian matrix factorization approaches use a likelihood distribution to capture the noise in the data and place priors over factor matrices:

$$p(\mathbf{W}, \mathbf{H}, \boldsymbol{\theta} | \mathbf{X}) \propto p(\mathbf{X} | \mathbf{W}, \mathbf{H}, \boldsymbol{\theta}) p(\mathbf{W}, \mathbf{H}, \boldsymbol{\theta})$$

likelihood \approx cost function
priors \approx penalty terms

- Several approaches to Bayesian NMF:
 - Gibbs sampling with different prior-likelihood distribution combinations
 - Variational Bayesian inference
 - ...
- Our agenda: *NMF Gibbs sampler & Iterated conditional modes (ICM)* algorithm

$$p(\mathbf{W}, \mathbf{H}, \boldsymbol{\theta} | \mathbf{X}) \propto p(\mathbf{X} | \mathbf{W}, \mathbf{H}, \boldsymbol{\theta}) p(\mathbf{W}, \mathbf{H}, \boldsymbol{\theta})$$

NMF Gibbs sampler

The NMF problem can be rewritten as: $\mathbf{X} = \mathbf{W}\mathbf{H} + \mathbf{E}$, where \mathbf{E} is a residual matrix.

LIKELIHOOD

We assume that the residuals \mathbf{E} are **normally distributed** with mean 0 and variance σ^2 .

$$p(\mathbf{X} | \mathbf{W}, \mathbf{H}, \sigma^2) = \prod_{n,p} \mathcal{N}(\mathbf{X}_{n,p}; (\mathbf{W}\mathbf{H})_{n,p}, \sigma^2)$$

PRIORS

- \mathbf{W} and \mathbf{H} independently **exponentially distributed** with scales $\alpha_{n,r}, \beta_{r,p}$.
- The noise variance σ^2 has an **inverse gamma distribution** with shape k and scale θ .

$$p(\mathbf{W}) = \prod_{n,r} \mathcal{E}(\mathbf{W}_{n,r}; \alpha_{n,r})$$

$$p(\mathbf{H}) = \prod_{r,p} \mathcal{E}(\mathbf{H}_{r,p}; \beta_{r,p})$$

$$p(\sigma^2) = \mathcal{G}^{-1}(\sigma^2; k, \theta).$$

POSTERIOR $p(\mathbf{W}, \mathbf{H}, \sigma^2 | \mathbf{X})$

- Proportional to the product of likelihood and priors.
- However, we are interested in estimating the marginal density of the factors ... and we can't really integrate the posterior.



Gibbs sampling!

But we only need to know how to sample from the full conditionals!

NMF Gibbs sampler

Full conditionals:

$$p(W_{n,r} | \mathbf{X}, \mathbf{W}_{\setminus(n,r)}, \mathbf{H}, \sigma^2) \propto p(\mathbf{X} | \mathbf{W}, \mathbf{H}, \sigma^2) p(\mathbf{W}) \propto \dots \propto \mathcal{T}\mathcal{N}(W_{n,r}; \mu_{W_{nr}}, \sigma_{W_{nr}}^2)$$

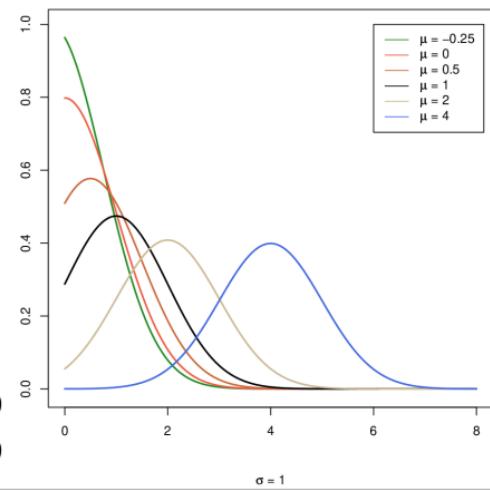
$$p(H_{r,p} | \mathbf{X}, \mathbf{W}, \mathbf{H}_{\setminus(r,p)}, \sigma^2) \propto p(\mathbf{X} | \mathbf{W}, \mathbf{H}, \sigma^2) p(\mathbf{H}) \propto \dots \propto \mathcal{T}\mathcal{N}(H_{r,p}; \mu_{H_{rp}}, \sigma_{H_{rp}}^2)$$

$$p(\sigma^2 | \mathbf{X}, \mathbf{W}, \mathbf{H}) \propto \dots \propto \mathcal{G}^{-1}(\sigma^2; k^*, \theta^*)$$

$$\mathcal{T}\mathcal{N}(x; \mu, \sigma^2) = \begin{cases} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$



Different Density Functions of Left Truncated Normal Distributions



NMF Gibbs sampler

Efficient implementation: the elements in each column of \mathbf{W} (row of \mathbf{H}) are conditionally independent, so we can sample an entire column (row) simultaneously.

Algorithm 3: NMF Gibbs sampler

```
1 Initialize  $\mathbf{W}^{(0)}$  and  $\mathbf{H}^{(0)}$ .
2 for  $i = 1:N$  do
3    $\mathbf{W} \leftarrow \mathbf{W}^{(i-1)}$ ,  $\mathbf{H} \leftarrow \mathbf{H}^{(i-1)}$ 
4   for  $j = 1$  to  $r$  do
5      $\mathbf{W}_{:,j} \leftarrow \mathcal{T}\mathcal{N}(\mu_{\mathbf{W}_{:,j}}, \sigma_{\mathbf{W}_{:,j}}^2)$ 
6      $\sigma^2 \leftarrow \mathcal{G}^{-1}(k^*, \theta_i^*)$ 
7     for  $j = 1$  to  $r$  do
8        $\mathbf{H}_{j,:} \leftarrow \mathcal{T}\mathcal{N}(\mu_{\mathbf{H}_{j,:}}, \sigma_{\mathbf{H}_{j,:}}^2)$ 
9    $\mathbf{W}^{(i)} \leftarrow \mathbf{W}$ ,  $\mathbf{H}^{(i)} \leftarrow \mathbf{H}$ 
10 Output:  $\{\mathbf{W}^{(i)}, \mathbf{H}^{(i)}\}_{i=1}^N$ 
```

Iterated conditional modes (ICM)

- Maximum-a-posteriori (MAP) estimate: $\theta_{MAP} = \max_{\theta} p(\theta | \mathbf{X})$
- For random variables $A \sim \mathcal{G}^{-1}(a, b)$ and $B \sim \mathcal{T}\mathcal{N}(\mu, \sigma^2)$ the modes are $\frac{b}{a+1}$ and $\max(0, \mu)$, respectively.

Algorithm 4: ICM

```
1 Initialize  $\mathbf{W}$  and  $\mathbf{H}$ .
2  $i = 1$ 
3 while stopping criteria not reached do
4   for  $j = 1$  to  $r$  do
5      $\mathbf{W}_{:,j} \leftarrow \max(0, \mu_{\mathbf{W}_{:,j}})$ 
6      $\sigma^2 = \frac{\theta_i^*}{k^* + 1}$ 
7     for  $j = 1$  to  $r$  do
8        $\mathbf{H}_{j,:} \leftarrow \max(0, \mu_{\mathbf{H}_{j,:}})$ 
9    $i = i + 1$ 
10 Output:  $\mathbf{W}, \mathbf{H}$ 
```

To Bayes or not to Bayes



Experiments

Convergence:

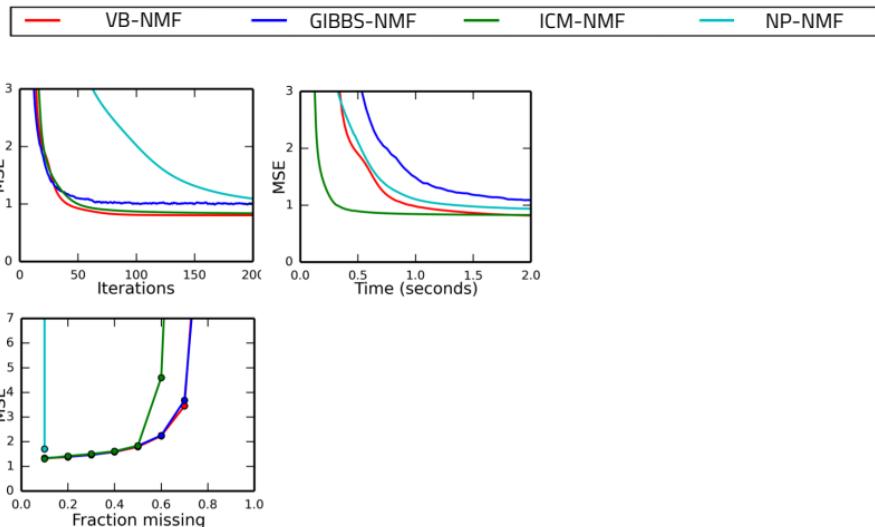
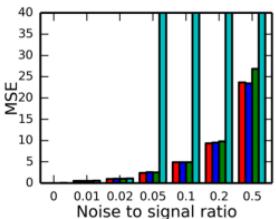
- Iterations: track the convergence rate of the error on the data against the number of iterations taken.
- Time: timing each method 10 times and taking the average.

Missing data:

- Testing the ability of the model to recover missing values as the fraction of unknown entries increases.

Noise:

- Adding different levels of Gaussian noise to the data.



Method	Estimate	Requires sampling	Speed of convergence	Robustness
Non-probabilistic	point	no	medium	low
ICM	point (MAP)	no	high	medium
Gibbs	full posterior	yes	low	high

To Bayes or not to Bayes



Choice of r

- In the Bayesian framework model selection can be performed by evaluating the marginal likelihood, $p(X)$.

Chib's method

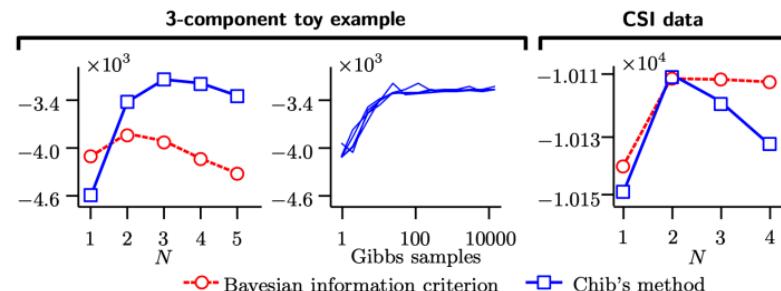
- Can be applied to Gibbs sampler output and is based on the relation: $p(X) = \frac{p(X|\theta)p(\theta)}{p(\theta|X)}$

$$\log p(X) = \log p(X|\theta^*) + \log p(\theta^*) - \log p(\theta^*|X)$$

- We can partition θ in arbitrary number of blocks, the only requirement: full-conditional sampling from each block needs to be possible!

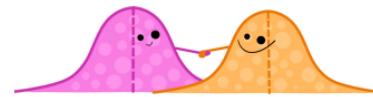
Example

- 3-component toy example:
 - Generated data: two random unit mean i.i.d. exponentially distributed matrices, where $n = 100$, $p = 20$, $r = 3$.
 - Added noise.
 - Computed marginal likelihood and BIC for r between 1 and 5.
- CSI data: chemical shift imaging data containing 2 spectral components.



Conclusion

- NMF is a dimensionality reduction method, providing more intuitive and meaningful interpretation, which is why it is so attractive for many different domains.
- Non-probabilistic vs. Bayesian:
 - Trade-offs in convergence, robustness to noise & sparsity.
 - Bayesian approach solves some of the issues of the non-probabilistic NMF.
- Nowadays we have more advanced methods for tasks, such as recommendation, however, when we are interested in interpretability, we still utilize methods based on the ideas presented today.



Materials & code: <https://github.com/ggapac/b-nmf>
Contact: greta.gasparac@fri.uni-lj.si