



Non-negative matrix factorization with Gaussian process priors

Mikkel N. Schmidt

Technical University of Denmark



Outline

■ Non-negative matrix factorization (NMF)

Examples of applications

- DNA microarray analysis
- Monaural audio separation

■ Gaussian processes (GP)

■ NMF with GP-priors

Example of application

- Chemical shift brain imaging



Non-negative matrix factorization

■ Non-negative bi-linear decomposition

$$x_{i,j} \approx \sum_{k=1}^K d_{i,k} \cdot h_{k,j} \quad \text{s.t. } d_{i,k}, h_{k,j} \geq 0$$

■ In matrix notation

$$\mathbf{X} \approx \mathbf{D}\mathbf{H} \quad \text{s.t. } \mathbf{D}, \mathbf{H} \geq 0$$

$$\begin{matrix} & \text{N} \\ \text{M} & \mathbf{X} \end{matrix} \approx \begin{matrix} & \text{K} \\ \text{M} & \mathbf{D} \end{matrix} \times \begin{matrix} & \text{N} \\ \text{K} & \mathbf{H} \end{matrix}$$



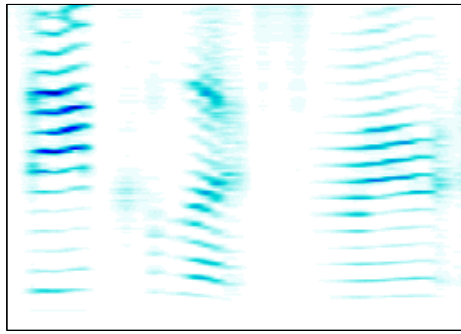
Why non-negativity?

- **Many signals are non-negative by nature**
 - Pixel intensities
 - Amplitude spectra
 - Occurrence counts
 - Discrete probabilities
 - etc.
- **Non-subtractive model**
 - No terms cancel out
 - **Parts-based: The whole is modeled as a sum of parts**

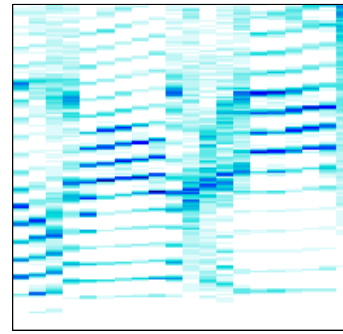


NMF, PCA, and VQ

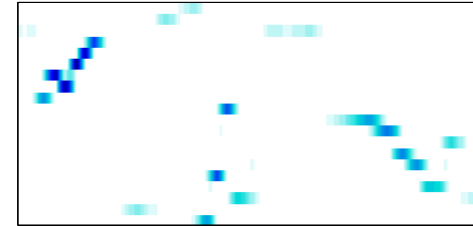
NMF



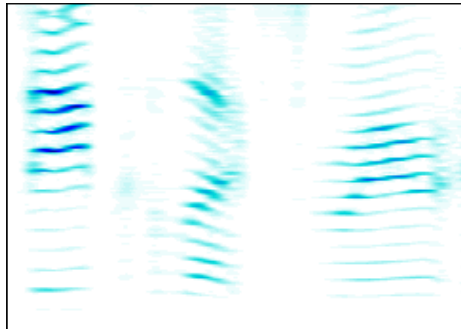
=



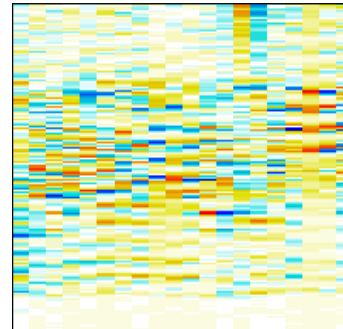
×



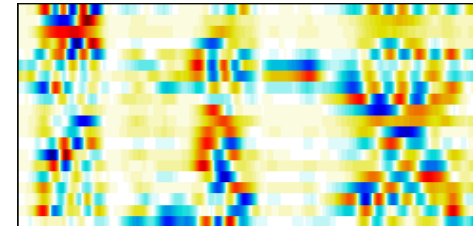
PCA



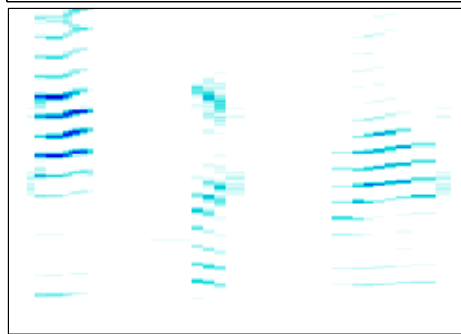
=



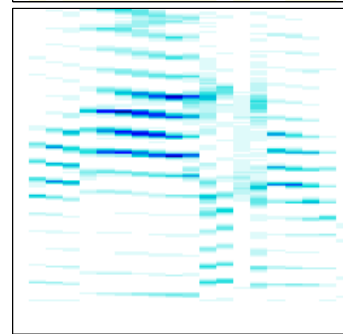
×



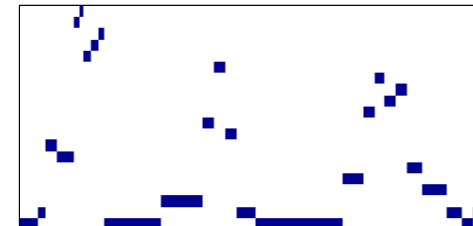
VQ



=



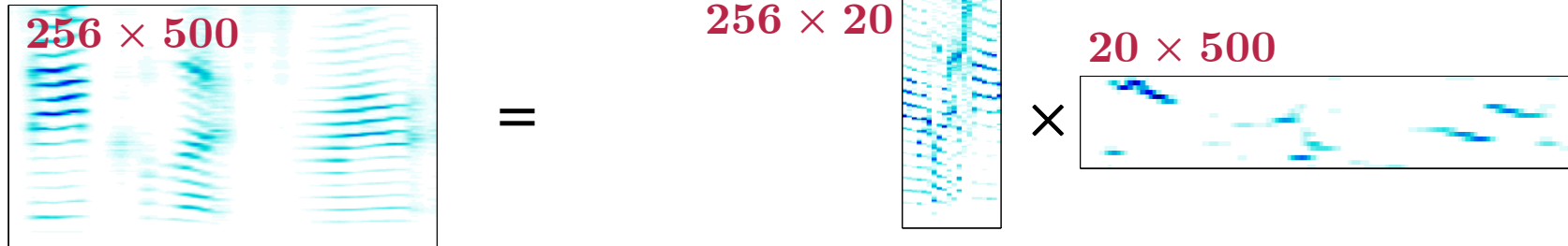
×



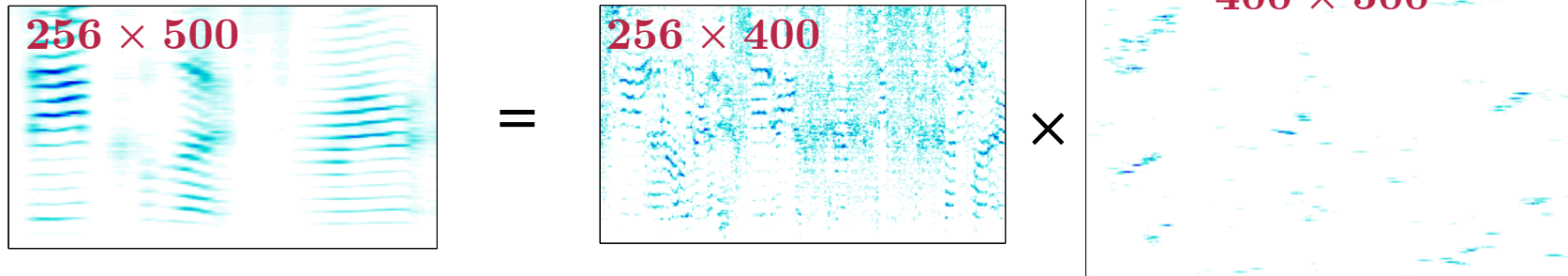


Over-complete decomposition

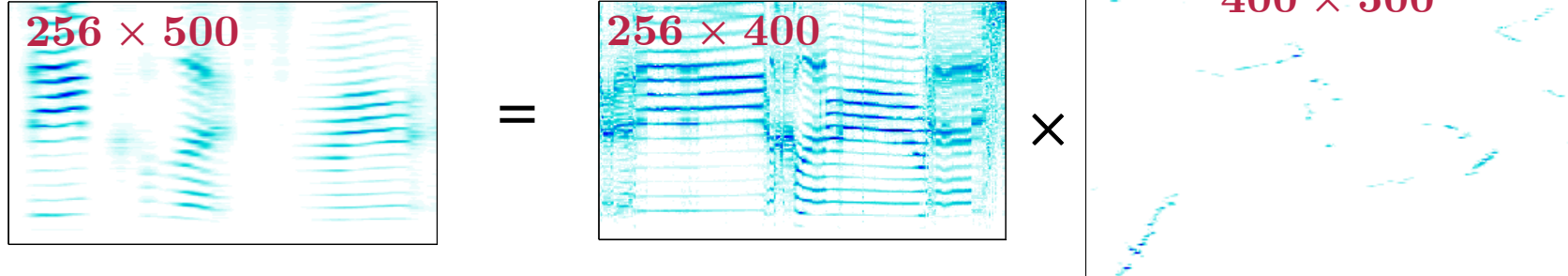
Under-complete



Over-complete



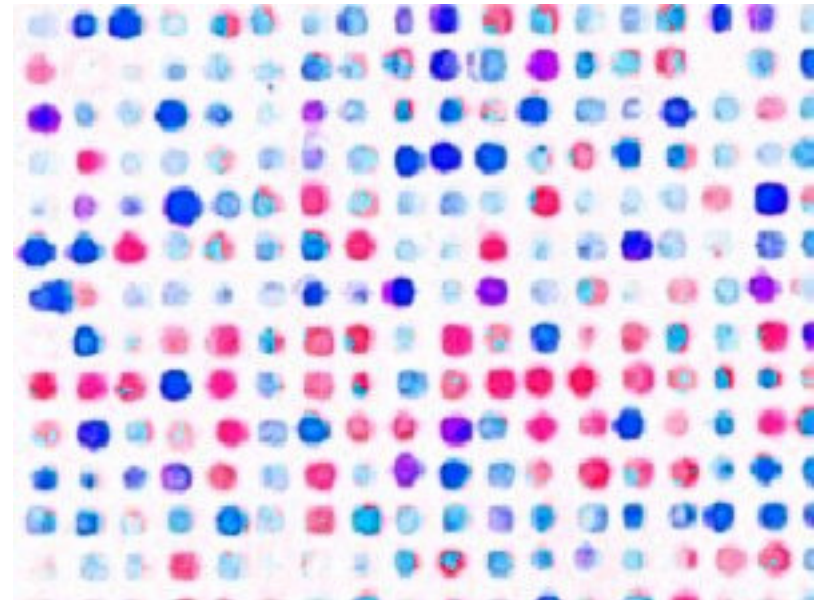
Sparse and over-complete





DNA microarray analysis

- DNA microarray technology enables parallel analysis of thousands of genes
- Data can be represented in non-negative matrix e.g. **gene** × **experiment**





Two types of leukemia

$$X \approx D \times H$$

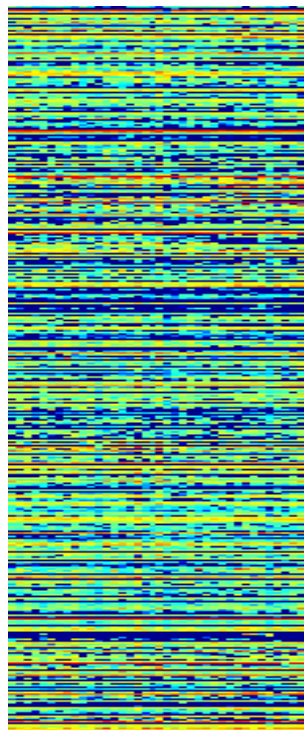
38 samples

2 classes

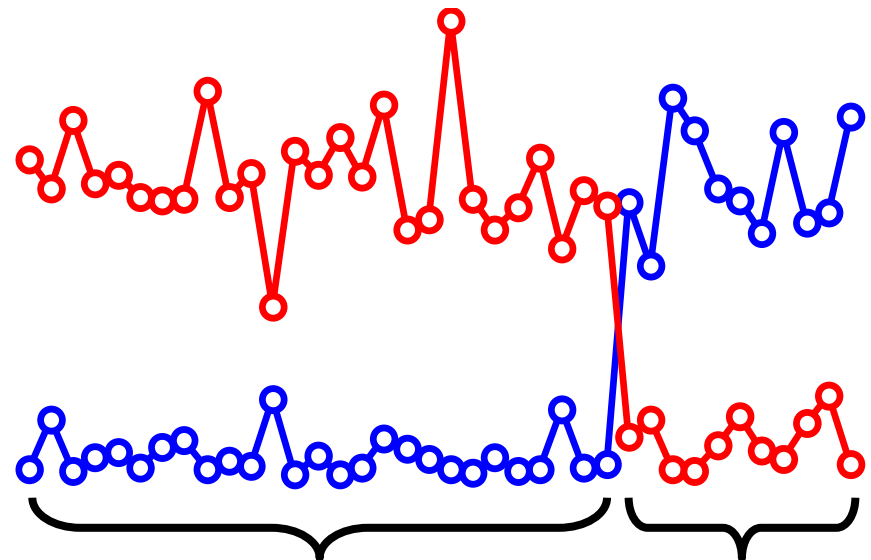
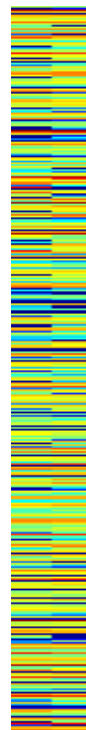
38 samples

2 classes

7129 genes



7129 genes



Class 1
27 samples

Class 2
11 samples



Monaural audio separation

- **Problem:** Separate audio sources using only one-microphone recording of mixture

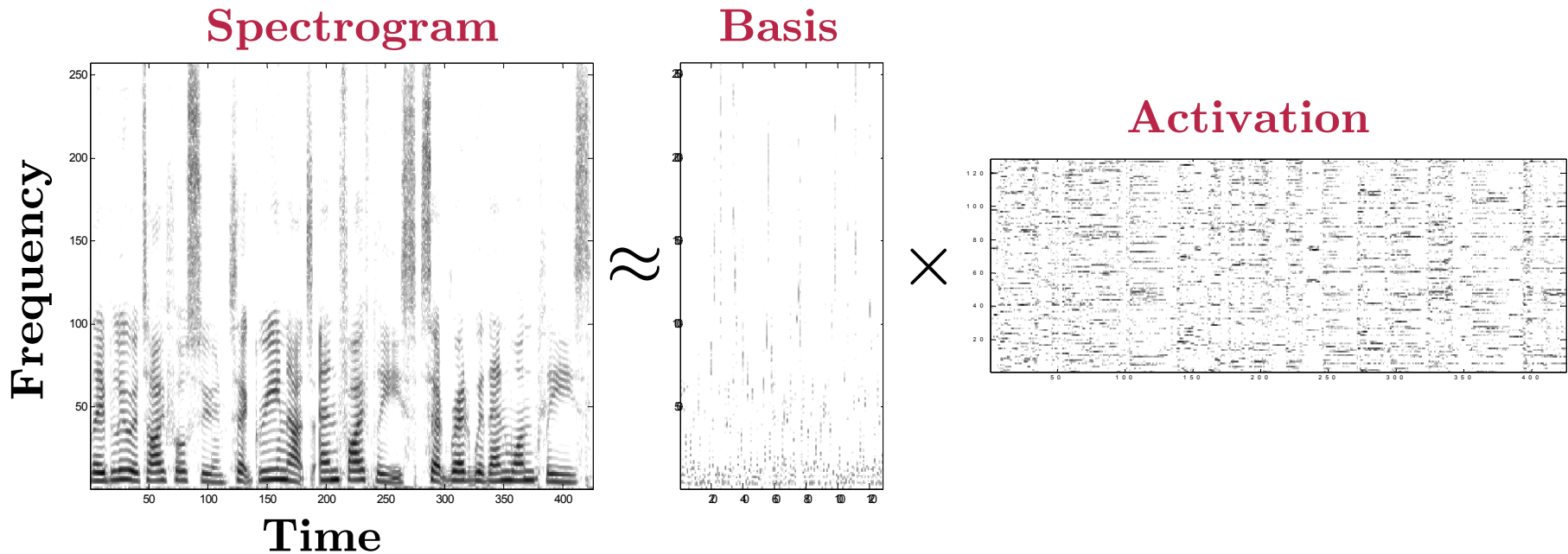




Amplitude spectrogram

- Audio represented as a non-negative matrix

$$X \approx DH$$





Audio separation with NMF

1. Learn a basis for each source

$$D_1, D_2, \dots, D_N$$

2. Compute activation for mixture

$$X \approx DH = [D_1, D_2, \dots, D_N]$$

$$\begin{bmatrix} H_1 \\ H_2 \\ \vdots \\ H_N \end{bmatrix}$$

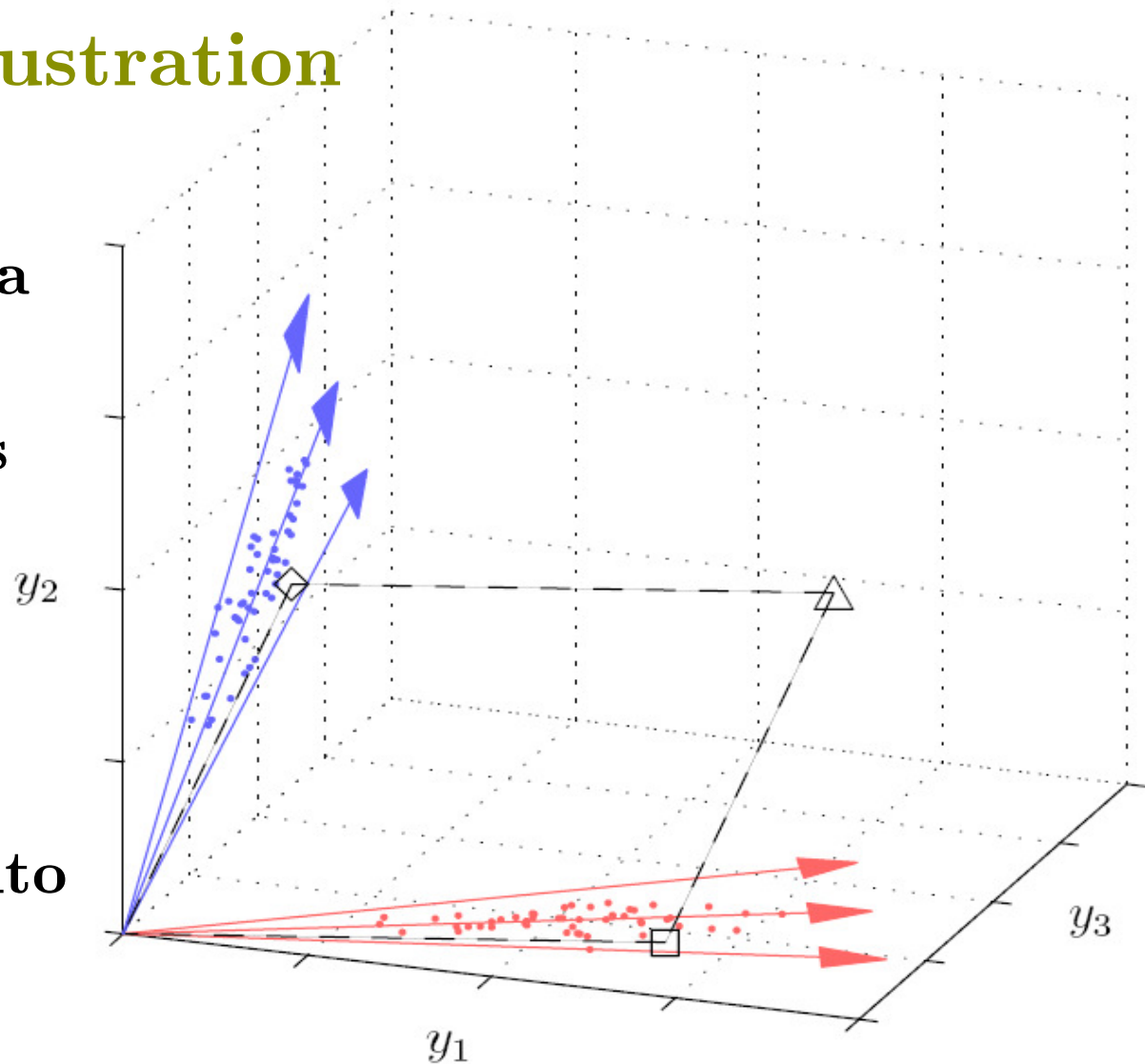
3. Reconstruct each source separately

$$\widehat{X}_n = D_n H_n$$



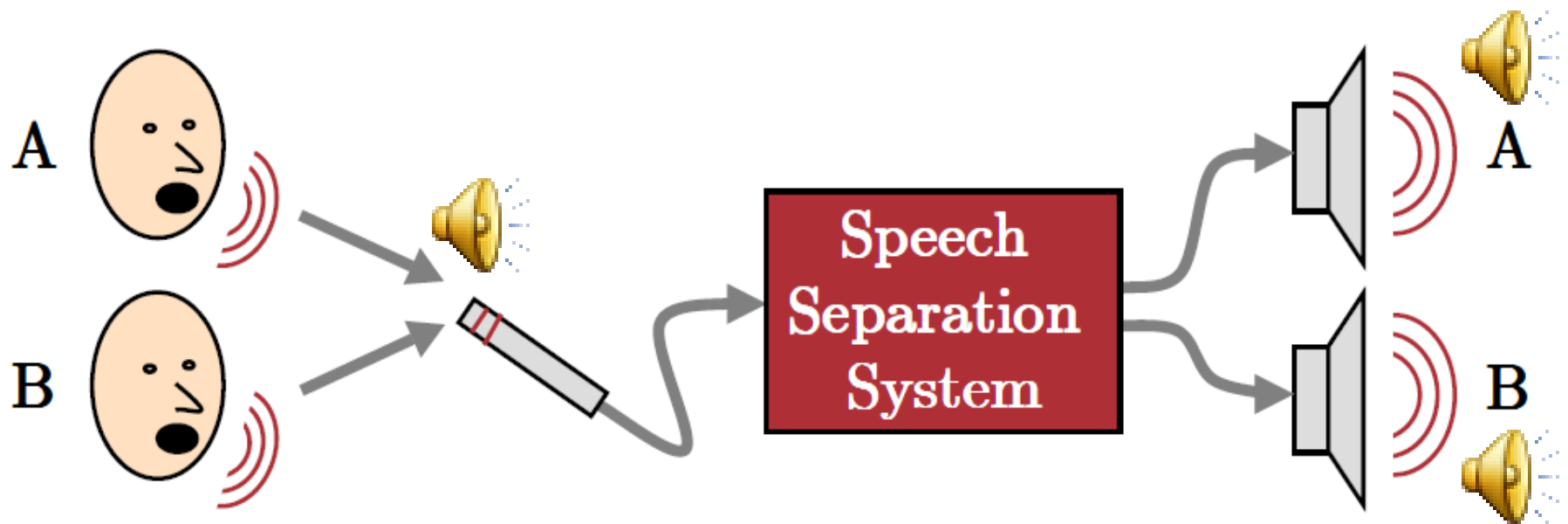
Subspace illustration

- Training data
- ↗↘↙ Basis vectors
- △ Observed mixture
- ◇ Mappings onto
- subspaces





Audio demonstration





Computing NMF

- Minimize divergence between X and DH , possibly plus some regularization

$$\min_{D, H \geq 0} \mathcal{D}(X, DH) + \mathcal{R}(D, H)$$

- **Constrained minimization problem:** A variety of algorithms for different divergence measures and regularizations
- Choice of divergence measure corresponds to assumptions about the data/noise/factor distribution



Maximum likelihood NMF

Example: Gaussian i.i.d. noise

■ Likelihood function

$$p(\mathbf{X} | \mathbf{D}, \mathbf{H}) = \prod_{i,j} \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(\mathbf{X}_{i,j} - [\mathbf{DH}]_{i,j})^2}{2\sigma^2}\right)$$

■ Negative log-likelihood (serves as divergence)

$$-\log p(\mathbf{X} | \mathbf{D}, \mathbf{H}) \propto \sum_{i,j} (\mathbf{X}_{i,j} - [\mathbf{DH}]_{i,j})^2$$



Maximum a posteriori NMF

■ Posterior distribution

$$p(\mathbf{D}, \mathbf{H} | \mathbf{X}) = \frac{p(\mathbf{X} | \mathbf{D}, \mathbf{H}) p(\mathbf{D}, \mathbf{H})}{p(\mathbf{X})}$$

■ Negative log posterior

$$-\log p(\mathbf{D}, \mathbf{H} | \mathbf{X}) \propto \underbrace{-\log p(\mathbf{X} | \mathbf{D}, \mathbf{H})}_{\text{Likelihood}} \underbrace{-\log p(\mathbf{D}, \mathbf{H})}_{\text{Prior}}$$

Likelihood

(divergence)

Prior

(regularization)



Bayesian NMF

- **Posterior distribution**

$$p(\mathbf{D}, \mathbf{H} | \mathbf{X}) = \frac{p(\mathbf{X} | \mathbf{D}, \mathbf{H})p(\mathbf{D}, \mathbf{H})}{p(\mathbf{X})}$$

- **In general, integrals over the posterior are intractable → approximate, e.g., using MCMC.**



Probabilistic model of factors

■ Standard NMF

$$X \approx DH \quad \text{s.t. } D, H \geq 0$$

- Factors constrained to be non-negative
- No other assumptions about the factors

■ Prior distribution over factors

$$p(D, H)$$

- Prior distribution captures **non-negativity** as well as other properties, such as **sparseness, smoothness, symmetries**, etc.



Which prior distributions to use?

- **Distribution over non-negative reals**
 - Rectified Gaussian: L_2 norm regularization
 - One-sided exponential: L_1 norm regularization
- **Gaussian process mapped to the non-negative reals**
 - Flexible, principled, and practical approach
 - Sparseness, smoothness, symmetries, etc.



Gaussian Processes

- A stochastic process which generates samples, x_i , such that any linear combination of x_i is Gaussian
- Characterized by its mean and covariance function
- Defines a distribution over functions



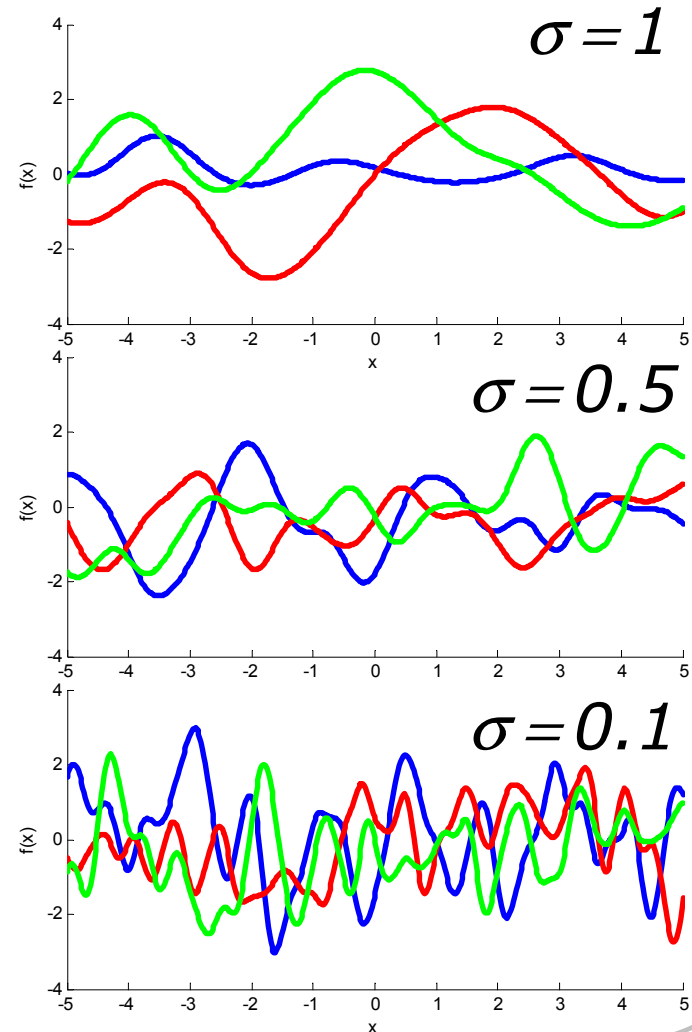
Example of Gaussian process

■ Mean function

$$m(x) = 0$$

■ Covariance function

$$k(x_1, x_2) = \exp\left(-\frac{(x_1 - x_2)^2}{2\sigma^2}\right)$$



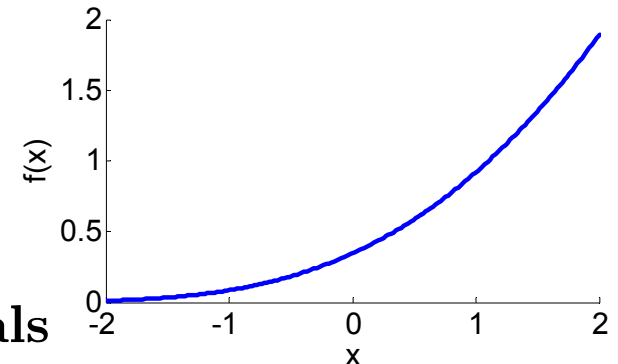


NMF with GP priors

$$X \approx DH$$

$$D = f_D(G_D) \quad H = f_H(G_H)$$

- G_D and G_H are Gaussian processes
- Link functions
 - Strictly increasing
 - Maps the reals to the non-negative reals





Link Function

- Map marginal distribution of Gaussian process to desired marginal distribution

$$H = f_H(\mathbf{G}_H) = P_H^{-1} \left(P_{G_H}(\mathbf{G}_H) \right)$$

- **Example: Gaussian-to-Exponential**

$$f_H(\mathbf{G}_H) = -\frac{1}{\lambda} \log \left(\frac{1}{2} - \frac{1}{2} \Phi \left(\frac{\mathbf{G}_H}{\sqrt{2}\sigma_i} \right) \right)$$



Change of Variable

- New variables δ and η

$$D = f_D(\mathbf{C}_D^\top \delta) \quad H = f_H(\mathbf{C}_H^\top \eta)$$

- \mathbf{C} is Cholesky decomposition of covariance matrix
- Variables δ and η are i.i.d. Gaussian

- **MAP estimate: $p(\delta, \eta | X)$**

- More robust and less local minima
- Unconstrained optimization (use e.g. conj. grad.)



Illustration of NMF with GP priors

■ Full matrix data

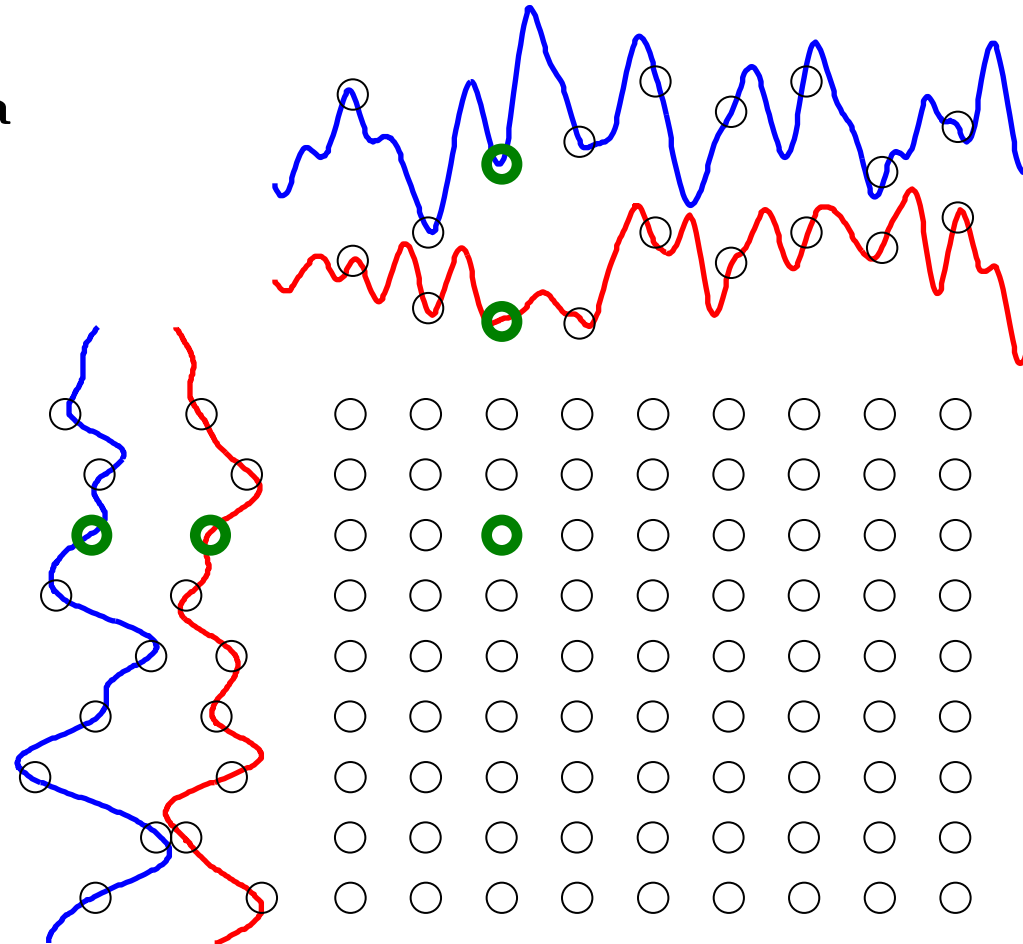




Illustration of NMF with GP priors

- Full matrix data
- Missing values

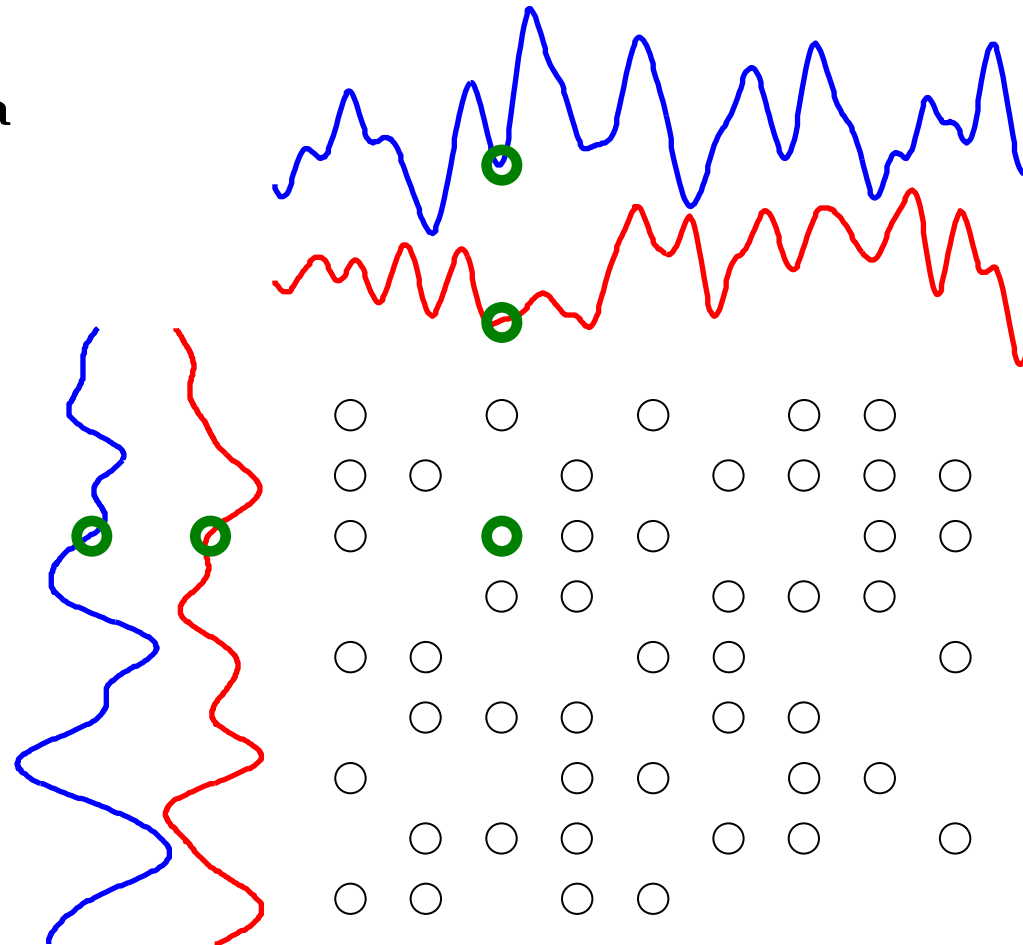
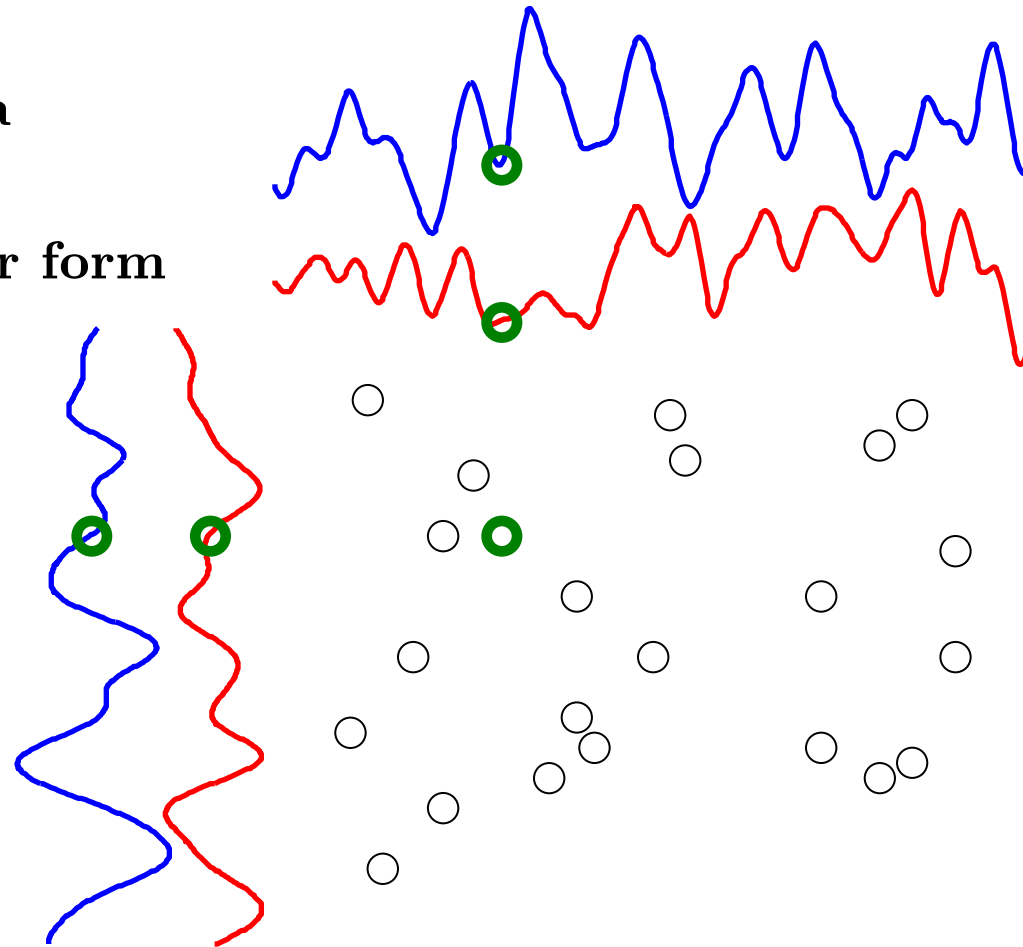




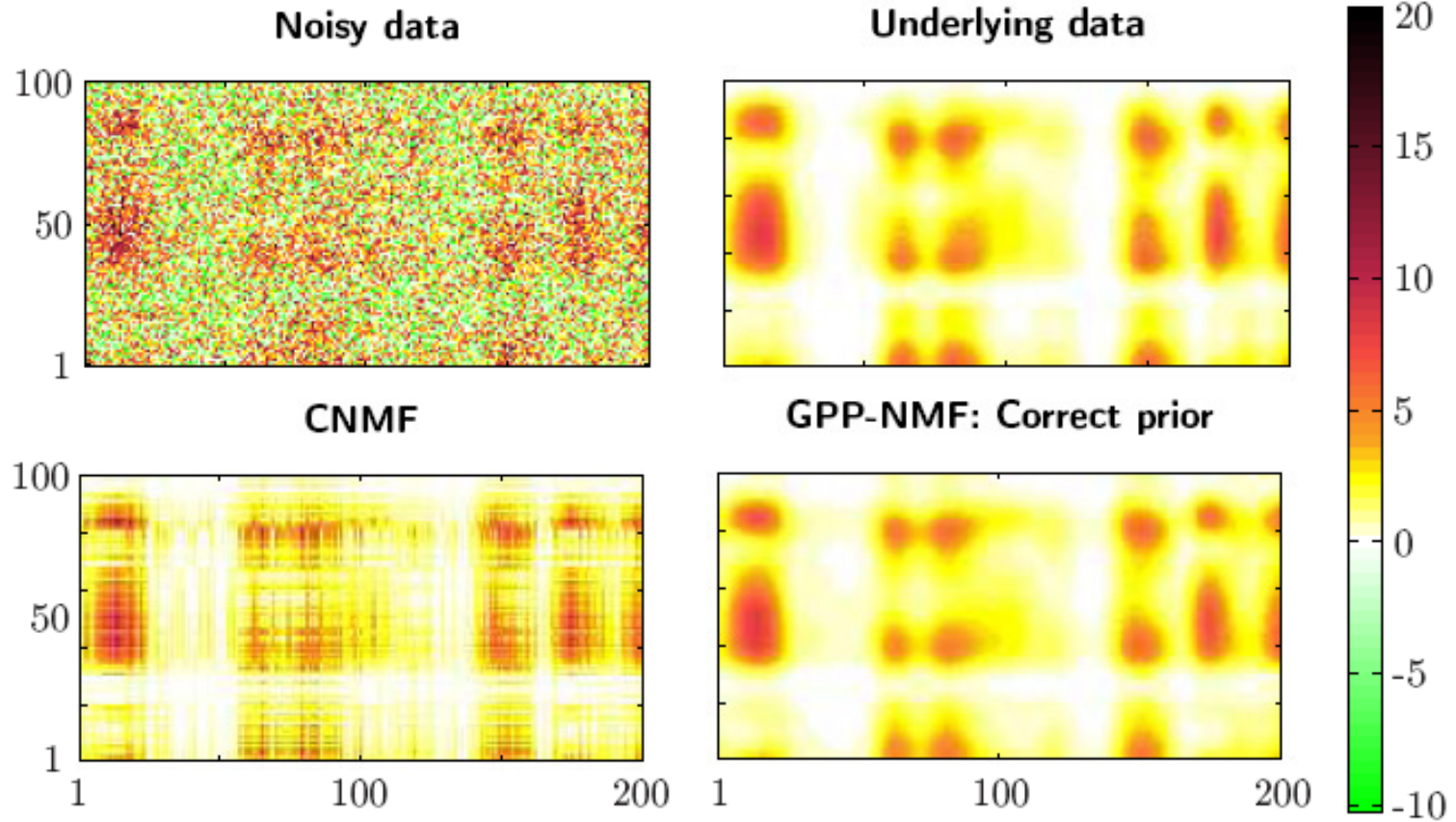
Illustration of NMF with GP priors

- Full matrix data
- Missing values
- General bi-linear form



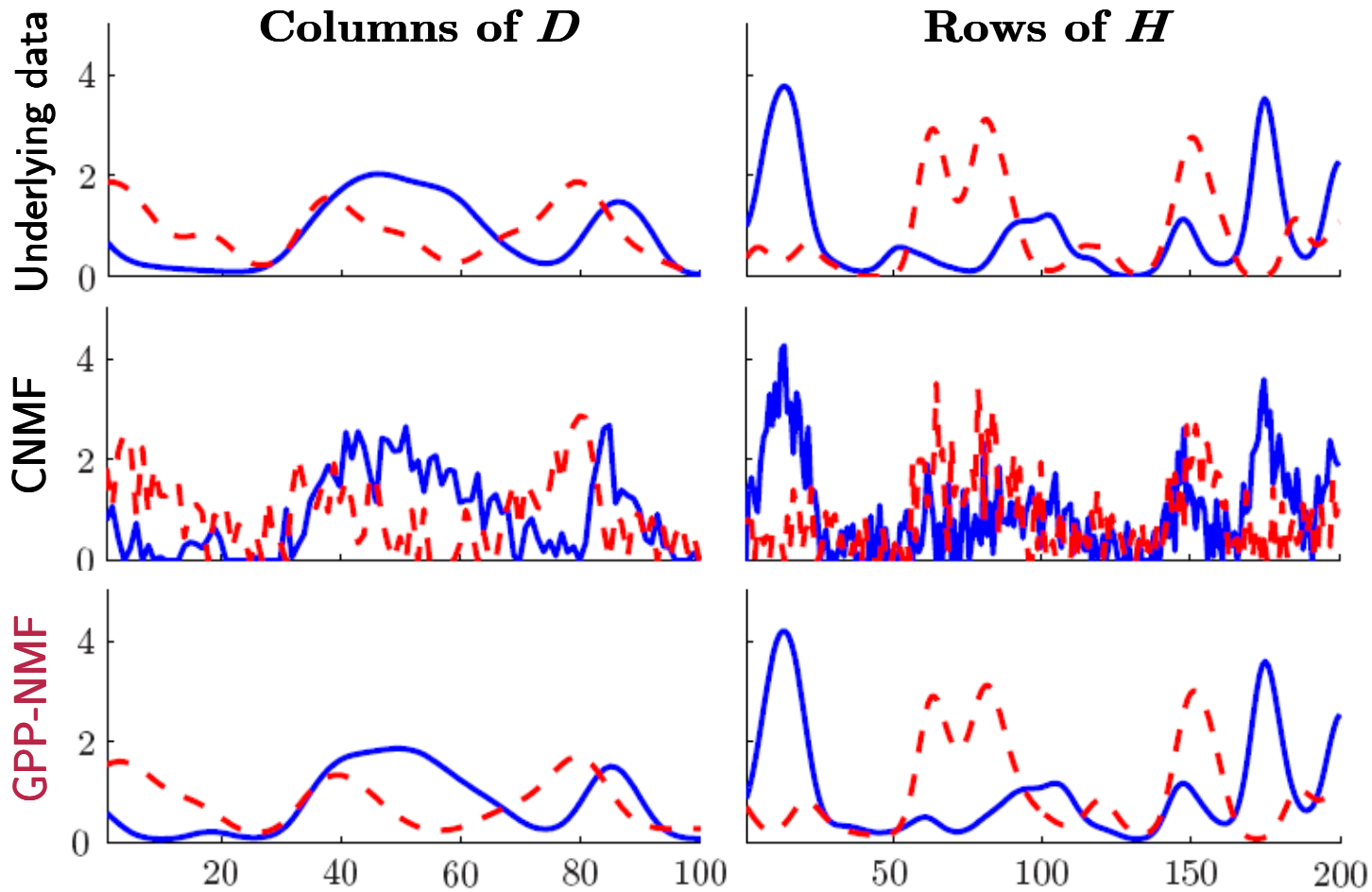


Toy Example





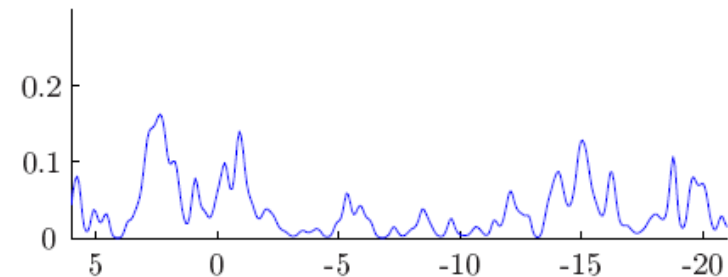
Toy Example





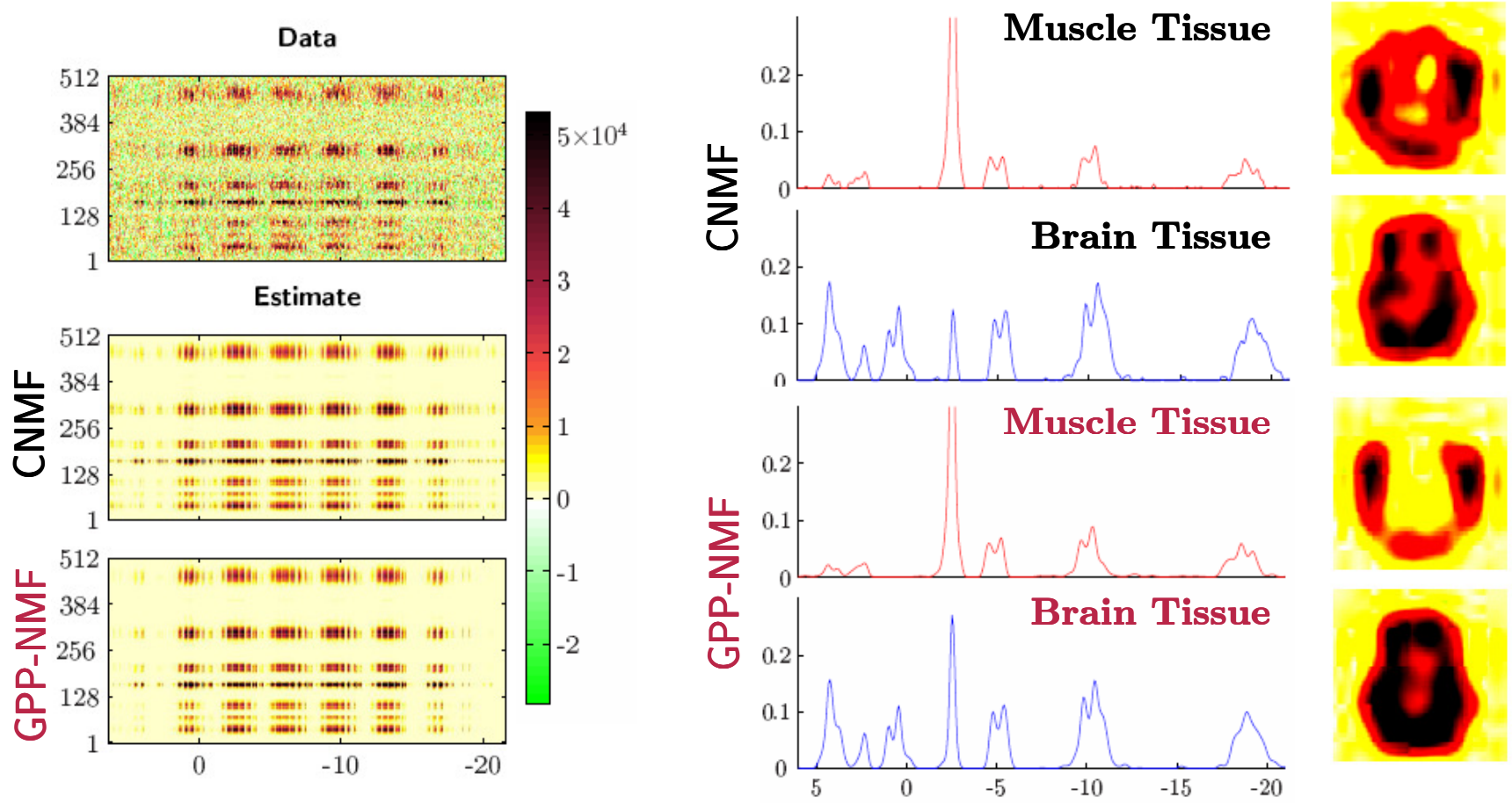
Chemical Shift Brain Imaging

- **Data:** 369 point chemical shift spectra measured at 512 positions ($8 \times 8 \times 8$ grid) in human skull.
- **Task:** Distinguish between brain and muscle tissue
- **Prior for spectra:**
Smooth, exponentially distributed.
- **Prior for activations in skull:**
Smooth in 3D, exponentially distributed, left-to-right symmetric.





Chemical Shift Brain Imaging





Conclusions

■ NMF

- General and versatile method
- Can be used to analyze a variety of problems including
 - DNA microarray analysis
 - Audio signal separation

■ NMF with GP-priors

- Extends the NMF framework by adding prior information
- Can improve the quality of non-negative factorizations



Future work

- Full Bayesian treatment of the model (MCMC)
- Learn parameters of kernel function
- Learn link functions from data
- Learn number of components (nonparametric Bayes)



References

Non-negative Matrix Factorization

- D. D. Lee and H. S. Seung. **Learning the parts of objects by non-negative matrix factorization.** *Nature*, 401(6755):788–791, 1999.
- P. Paatero and U. Tapper. **Positive matrix factorization: A nonnegative factor model with optimal utilization of error-estimates of data values.** *Environmetrics*, 5:111–126, 1994.
- Michael W. Berry, Murray Browne, Amy N. Langville, V. Paul Pauca, and Robert J. Plemmons. **Algorithms and applications for approximate nonnegative matrix factorization.** *Computational Statistics and Data Analysis*, 2006.

DNA Microarray Analysis

- Jean-Philippe Brunet, Pablo Tamayo, Todd R. Golub, and Jill P. Mesirov. **Metagenes and molecular pattern discovery using matrix factorization.** *Proceedings of the National Academy of Sciences (PNAS)*, 101(12):4164–4169, Mar 2004.



References

Speech Separation with NMF

- Mikkel N. Schmidt and Rasmus K. Olsson, **Linear Regression on Sparse Features for Single-Channel Speech Separation**. Applications of Signal Processing to Audio and Acoustics, IEEE Workshop on, 2007.
- Mikkel N. Schmidt and Rasmus K. Olsson, **Single-Channel Speech Separation using Sparse Non-Negative Matrix Factorization**. Spoken Language Processing, ISCA International Conference on, 2006.

Gaussian Processes

- Carl Edward Rasmussen and Christopher K. I. Williams, **Gaussian Processes for Machine Learning**. MIT Press, 2006.

NMF with GP-priors

- Mikkel N. Schmidt and Hans Laurberg, **Non-negative matrix factorization with Gaussian process priors**. Computational Intelligence and Neuroscience, 2008.