

# Dense Message Passing for Sparse Principal Component Analysis

Kevin Sharp    Magnus Rattray

University of Manchester

AIStats, Chia Laguna Resort, Sardinia, 14th May 2010

# Outline

- 1 Introduction
  - Motivating Application: Gene Regulation
- 2 Dense Message Passing
  - Model Description
  - Algorithm Description
  - Statistical Mechanics Theory
- 3 Results
  - Simulated Data
  - Gene Expression Data
  - Marginal Likelihood Estimation
- 4 Summary

# Motivation

- Gene regulation - inference of explanatory factors.
- Microarray data - 'Large  $p$  small  $N$ ' regime.
- Explanatory factors have truly sparse loadings.
- Zero-norm priors allocate probability mass to truly sparse solutions.
- Easy to encode prior knowledge of sparse structure.
- But, zero-norm priors are problematic for inference.

# Model - Probabilistic PCA

- For the  $n^{\text{th}}$  data point,  $\mathbf{y}_n$ , we assume:

$$\mathbf{y}_n = \mathbf{w}x_n + \epsilon_n ,$$

where  $x_n \sim \mathcal{N}(0, 1)$ .

- To simplify the description,  $\epsilon_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .
- Integrate out  $x$ :

$$P(\mathbf{y}_n | \mathbf{w}) = \mathcal{N}(\mathbf{y}_n | \mathbf{0}, \mathbf{I} + \mathbf{w}\mathbf{w}^T)$$

# Model - Probabilistic PCA

- For the  $n^{\text{th}}$  data point,  $\mathbf{y}_n$ , we assume:

$$\mathbf{y}_n = \mathbf{w}x_n + \epsilon_n ,$$

where  $x_n \sim \mathcal{N}(0, 1)$ .

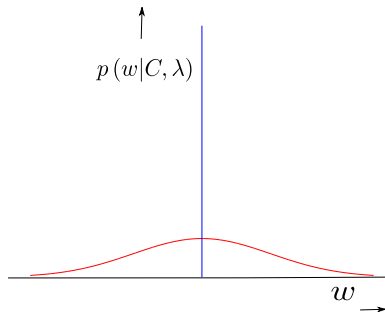
- To simplify the description,  $\epsilon_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .
- Integrate out  $x$ :

$$P(\mathbf{y}_n | \mathbf{w}) = \mathcal{N}(\mathbf{y}_n | \mathbf{0}, \mathbf{I} + \mathbf{w}\mathbf{w}^T)$$

# Sparsity Mixture Prior

We use a **spike and slab** mixture prior:

$$P(\mathbf{w}|C, \lambda) = \prod_{j=1}^p \left[ (1 - C)\delta(w_j) + C\mathcal{N}(w_j|0, (\lambda)^{-1}) \right]$$



$C$  - Fraction of non-zero  $w$ s  
 $\lambda$  - inverse width

# Sparsity Mixture Prior

Express in factorised form using binary variables  $z_j \in \{0, 1\}$ :

$$P(\mathbf{v}, \mathbf{z}) = \prod_{j=1}^p \left\{ (1 - C) \mathcal{N}(v_j | 0, 1) \right\}^{1-z_j} \left\{ C \mathcal{N}(v_j | 0, \lambda^{-1}) \right\}^{z_j}$$

where  $w_j = z_j v_j$  and  $v_j \in \mathcal{R}$

## Form of the Prior in High Dimensions

- $z_j \sim \text{Bernoulli}(C)$
- $\sum_j z_j \sim \text{Bin}(p, C)$
- For large dimension,  $p$ , the fraction of non-zero parameters is highly peaked at  $C$ .



# Form of the Prior in High Dimensions

- $P(w_j | z_j = 1) = \mathcal{N}(w_j | 0, \lambda^{-1})$
- For large dimension,  $p$ :

$$\|\mathbf{w}\|^2 \sim \mathcal{N}(0, Cp/\lambda)$$

- For large dimension,  $p$ , this distribution is **approximately spherical** with radius  $\sqrt{Cp/\lambda}$ .

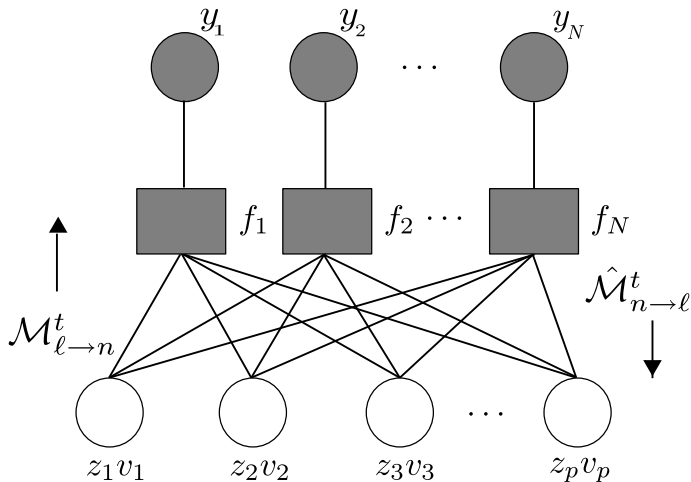
- Conclusion - A **constraint-based** prior:

$$p(\mathbf{w}, \mathbf{z} | C, \lambda) \propto \delta \left( \sum_{j=1}^p z_j - pC \right) \delta \left( \sum_{j=1}^p w_j^2 - \frac{pC}{\lambda} \right)$$

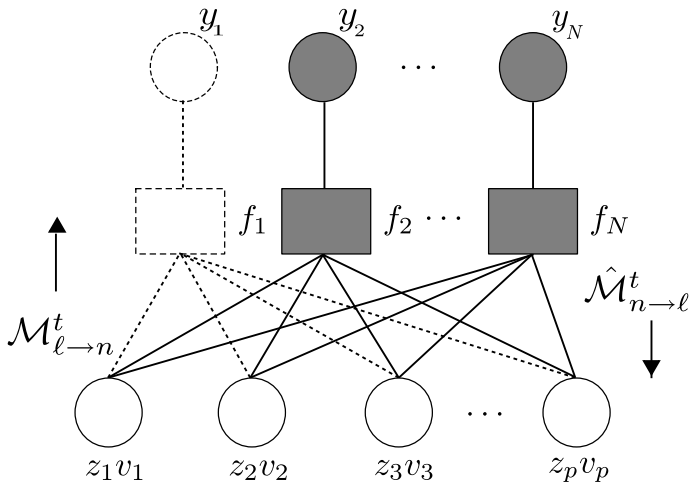
is almost equivalent to the mixture prior in high dimensions.

- This proves useful for developing the message passing algorithm.

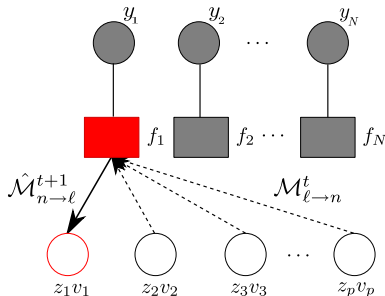
# Factor Graph Representation



# Belief Propagation

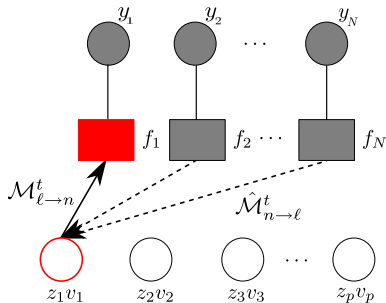


# Factor to Variable Messages



$$\hat{\mathcal{M}}_{n \rightarrow \ell}^{t+1}(v_\ell, z_\ell) \propto \int \prod_{j \neq \ell} dv_j \sum_{\mathbf{z} \setminus z_\ell} f_n(\mathbf{y}_n, \mathbf{z}, \mathbf{v}) \prod_{j \neq \ell} \mathcal{M}_{j \rightarrow n}^t(v_j, z_j)$$

# Variable to Factor Messages



$$\mathcal{M}_{\ell \rightarrow n}^t(v_\ell, z_\ell) \propto P(v_\ell, z_\ell) \prod_{m \neq n} \hat{\mathcal{M}}_{m \rightarrow \ell}^t(v_\ell, z_\ell)$$

# Marginal Beliefs

After  $t$  iterations, the approximate posterior marginal belief is:

$$p^t(z_\ell, v_\ell | \mathbf{Y}) = \frac{p(z_\ell, v_\ell) \prod_{m=1}^N \hat{\mathcal{M}}_{m \rightarrow \ell}^t(v_\ell, z_\ell)}{\int dv_\ell \sum_{z_\ell} p(z_\ell, v_\ell) \prod_{m=1}^N \hat{\mathcal{M}}_{m \rightarrow \ell}^t(v_\ell, z_\ell)}$$

where  $p(z_\ell, v_\ell)$  is the prior.

# Two Problems

Unfortunately,

1

$$\hat{\mathcal{M}}_{n \rightarrow \ell}^{t+1}(\mathbf{v}_\ell, \mathbf{z}_\ell) \propto \int \prod_{j \neq \ell} d\mathbf{v}_j \sum_{\mathbf{z} \setminus \mathbf{z}_\ell} f_n(\mathbf{y}_n, \mathbf{z}, \mathbf{v}) \prod_{j \neq \ell} \mathcal{M}_{j \rightarrow n}^t(\mathbf{v}_j, \mathbf{z}_j)$$

is hard to compute.

- 2 Belief propagation is not expected to converge for dense graphical models.



# Solutions

- 1 Exploit the high-dimensionality:

Use a **Gaussian approximation**.

- 2 Impose consistency requirements:

Use the constraint-based prior to **enforce sparsity and length constraints** self-consistently at each iteration.

Uda and Kabashima - Statistical Mechanical Development of a Sparse Bayesian Classifier, *J. Phys. Soc. Japan*, 2005

# Gaussian Approximation (1)

Notice that likelihood factors may be written as:

$$f_n(\mathbf{y}_n, \mathbf{z}, \mathbf{v}) = \frac{1}{\sqrt{(2\pi)^p (1 + \|\mathbf{w}\|^2)}} \exp\left(-\frac{\mathbf{y}_n^T \mathbf{y}_n}{2} + \Delta_n^2/2\right),$$

with  $\Delta_n$  defined by:  $\Delta_n = \frac{\sum_{j=1}^p y_j^n z_j v_j}{\sqrt{1 + \|\mathbf{w}\|^2}}$

For large dimension,  $p$ , Central Limit Theorem permits a Gaussian approximation.

## Gaussian Approximation (2)

For constant  $\|\mathbf{w}\|^2$ , we replace  $\Delta_n$  by:

$$\frac{y_\ell^n z_\ell v_\ell}{\sqrt{1 + Cp/\lambda}} + \underbrace{\frac{1}{\sqrt{1 + Cp/\lambda}} \sum_{j \neq \ell} y_j^n m_{j \rightarrow n}^t}_{\langle \Delta_{n \setminus \ell} \rangle_{n \setminus \ell}^t} + \sqrt{V_{n \setminus \ell}^t} u$$

where  $u \sim \mathcal{N}(0, 1)$ .

$m_{j \rightarrow n}^t$  is the mean of  $z_j v_j$  under the cavity distribution with the  $n^{\text{th}}$  data point removed.

## Gaussian Approximation (3)

- The variance,  $V_{n \setminus \ell}^t$  is given by:

$$\frac{1}{1 + Cp/\lambda} \sum_{j,k \neq \ell} y_j^n y_k^n \langle (z_j v_j - m_{j \rightarrow n}^t) (z_k v_k - m_{k \rightarrow n}^t) \rangle_{n \setminus \ell}^t$$

- For large dimension,  $p$ , fluctuations about the sample mean are  $\mathcal{O}\left(\frac{1}{\sqrt{p}}\right)$ :  $V_{n \setminus \ell}^t$  is *self-averaging*.

- $$V^t \approx \frac{1}{(1 + Cp/\lambda)} \left( Cp/\lambda - \sum_{j=1}^p (m_j^t)^2 \right)$$

# Consistency - Constraint-based Prior

- 1 The spike and slab prior can be written:

$$P(\mathbf{v}, \mathbf{z}) \propto \prod_{j=1}^p \exp\left(-\frac{1}{2}(1 - z_j + Gz_j)v_j^2 + \gamma z_j\right)$$

where  $\gamma = \ln\left(\frac{C\sqrt{\lambda}}{1-C}\right)$  and  $G = \lambda$ .

- 2 Adjust  $G$  and  $\gamma$  at each iteration to satisfy the constraint-based prior *on average*:

$$\sum_{j=1}^p \langle z_j \rangle^t = Cp \text{ and } \sum_{j=1}^p \langle z_j v_j^2 \rangle^t = Cp/\lambda$$

- 3 Note, after convergence,  $G \neq \lambda$  and  $\gamma \neq \ln\left(\frac{C\sqrt{\lambda}}{1-C}\right)$
- 4 Consistent with replica analysis.

# Consistency - Constraint-based Prior

- 1 The spike and slab prior can be written:

$$P(\mathbf{v}, \mathbf{z}) \propto \prod_{j=1}^p \exp\left(-\frac{1}{2}(1 - z_j + Gz_j)v_j^2 + \gamma z_j\right)$$

where  $\gamma = \ln\left(\frac{C\sqrt{\lambda}}{1-C}\right)$  and  $G = \lambda$ .

- 2 Adjust  $G$  and  $\gamma$  at each iteration to satisfy the constraint-based prior *on average*:

$$\sum_{j=1}^p \langle z_j \rangle^t = Cp \text{ and } \sum_{j=1}^p \langle z_j v_j^2 \rangle^t = Cp/\lambda$$

- 3 Note, after convergence,  $G \neq \lambda$  and  $\gamma \neq \ln\left(\frac{C\sqrt{\lambda}}{1-C}\right)$
- 4 Consistent with replica analysis.

# Consistency - Constraint-based Prior

- 1 The spike and slab prior can be written:

$$P(\mathbf{v}, \mathbf{z}) \propto \prod_{j=1}^p \exp\left(-\frac{1}{2}(1 - z_j + Gz_j)v_j^2 + \gamma z_j\right)$$

where  $\gamma = \ln\left(\frac{C\sqrt{\lambda}}{1-C}\right)$  and  $G = \lambda$ .

- 2 Adjust  $G$  and  $\gamma$  at each iteration to satisfy the constraint-based prior *on average*:

$$\sum_{j=1}^p \langle z_j \rangle^t = Cp \text{ and } \sum_{j=1}^p \langle z_j v_j^2 \rangle^t = Cp/\lambda$$

- 3 Note, after convergence,  $G \neq \lambda$  and  $\gamma \neq \ln\left(\frac{C\sqrt{\lambda}}{1-C}\right)$

- 4 Consistent with replica analysis.

# Consistency - Constraint-based Prior

- 1 The spike and slab prior can be written:

$$P(\mathbf{v}, \mathbf{z}) \propto \prod_{j=1}^p \exp\left(-\frac{1}{2}(1 - z_j + Gz_j)v_j^2 + \gamma z_j\right)$$

where  $\gamma = \ln\left(\frac{C\sqrt{\lambda}}{1-C}\right)$  and  $G = \lambda$ .

- 2 Adjust  $G$  and  $\gamma$  at each iteration to satisfy the constraint-based prior *on average*:

$$\sum_{j=1}^p \langle z_j \rangle^t = Cp \text{ and } \sum_{j=1}^p \langle z_j v_j^2 \rangle^t = Cp/\lambda$$

- 3 Note, after convergence,  $G \neq \lambda$  and  $\gamma \neq \ln\left(\frac{C\sqrt{\lambda}}{1-C}\right)$
- 4 Consistent with replica analysis.



## Replica Analysis (1)

- Compute average of the log marginal likelihood over all possible datasets for  $p \rightarrow \infty$
- $\alpha = N/p$  is held constant (where  $N$  is the sample size).
- Works well for  $\alpha \ll 1$  – ‘large  $p$  small  $N$ ’
- Not mathematically rigorous, but a useful tool.

## Replica Analysis (2)

- Derive expressions involving the posterior mean ( $^{\text{PM}}$ ) parameter vector,  $\mathbf{w}^{\text{PM}}$ :
  - squared length,  $\|\mathbf{w}^{\text{PM}}\|^2$
  - overlap with the true parameter vector,  $\mathbf{w}^{\text{PM}} \cdot \mathbf{w}^t$ .

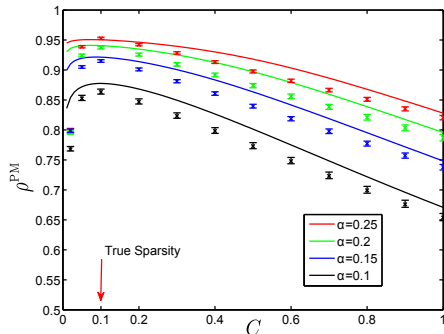
$$\mathbf{w}^t \sim \prod_{j=1}^p \left[ (1 - C_t) \delta(\mathbf{w}_j) + C_t \mathcal{N}(\mathbf{w}_j | 0, (\lambda_t)^{-1}) \right]$$

- Can show that the algorithm is consistent with this analysis.
- Can compare algorithm performance to theory using

$$\rho^{\text{PM}} = \frac{\mathbf{w}^{\text{PM}} \cdot \mathbf{w}^t}{\|\mathbf{w}^{\text{PM}}\| \|\mathbf{w}^t\|} .$$

# Simulated Data - DMP vs Theory

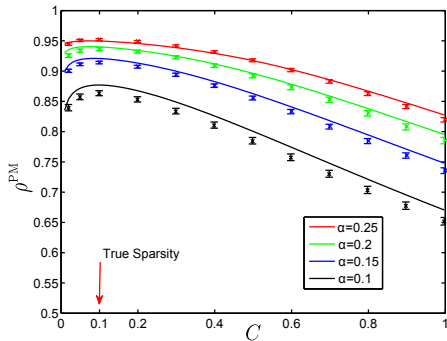
## DMP



$C$  - fraction of non-zero parameters;

$N = 200$  samples,  $\alpha = N/p$ ;

## Gibbs

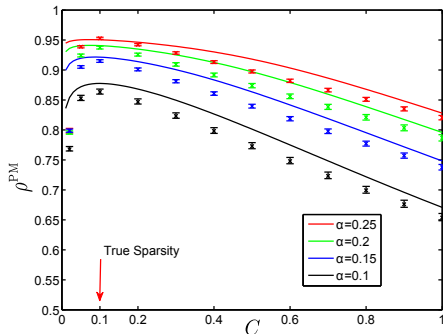


$\rho^{\text{PM}}$  cosine angle between  $\mathbf{w}^{\text{PM}}$  and  $\mathbf{w}^t$ .

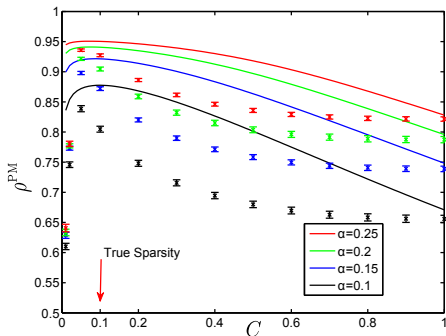
Results averaged over 50 sample datasets.

# Simulated Data - DMP vs emPCA

DMP



emPCA



$C$  - fraction of non-zero parameters;

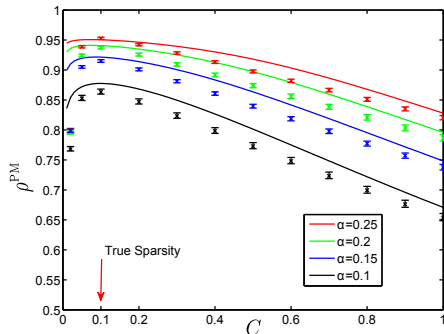
$N = 200$  samples,  $\alpha = N/p$ ;

$\rho^{\text{PM}}$  cosine angle between  $\mathbf{w}^{\text{PM}}$  and  $\mathbf{w}^t$ .

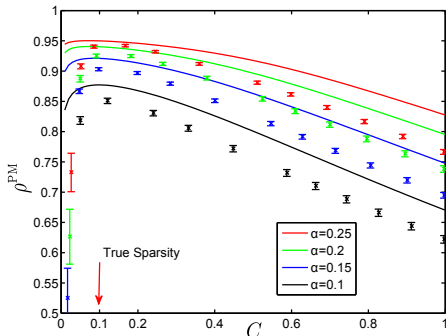
Results averaged over 50 sample datasets.

# Simulated data - DMP vs SPCA

DMP



SPCA



$C$  - fraction of non-zero parameters;

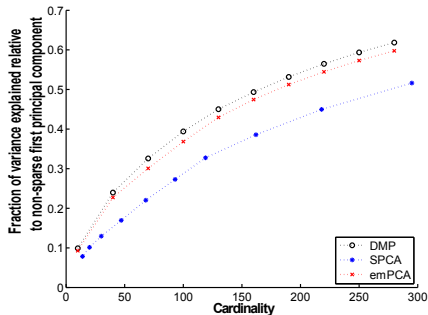
$N = 200$  samples,  $\alpha = N/p$ ;

$\rho^{\text{PM}}$  cosine angle between  $\mathbf{w}^{\text{PM}}$  and  $\mathbf{w}^t$ .

Results averaged over 50 sample datasets.

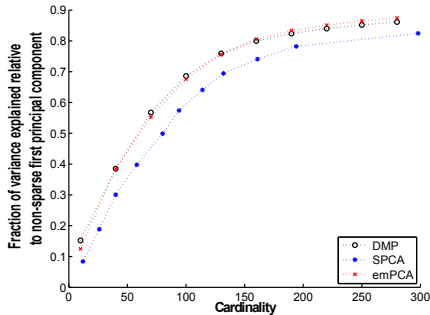
# Gene Expression Data - DMP vs emPCA and SPCA

Armstrong *et al.*



$p = 12582, N = 72$

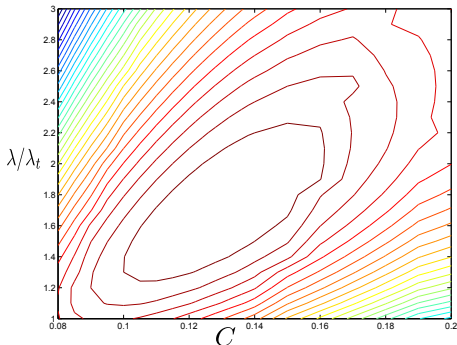
Ramaswamy *et al.*



$p = 16063, N = 144$

# Marginal Likelihood Estimation - Simulated Data

DMP

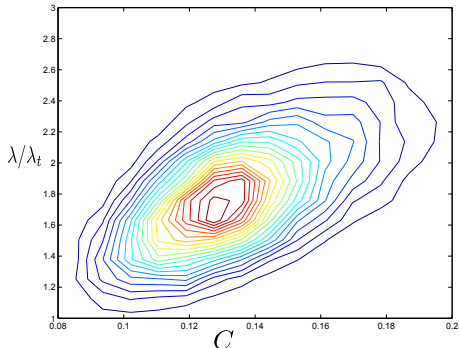


$C$  - fraction of non-zero parameters;

$\lambda$  - assumed signal precision;

True sparsity - 0.1.

Gibbs



$N = 200$  samples; dimension,  $p = 2000$ .

$\lambda^t$  - true signal precision.

# Summary

- Novel message passing algorithm for Sparse Bayesian PCA in high dimensions
- Message updates rendered tractable using a Gaussian approximation
- Convergence achieved by imposing consistency requirements derived from statistical mechanics analysis.
- Inference of posterior marginals exhibits near optimal performance compared to theory.
- Outperforms two other recently published algorithms.
- Approximation to Marginal Likelihood also available.



# Summary

- Novel message passing algorithm for Sparse Bayesian PCA in high dimensions
- Message updates rendered tractable using a Gaussian approximation
- Convergence achieved by imposing consistency requirements derived from statistical mechanics analysis.
- Inference of posterior marginals exhibits near optimal performance compared to theory.
- Outperforms two other recently published algorithms.
- Approximation to Marginal Likelihood also available.

# Summary

- Novel message passing algorithm for Sparse Bayesian PCA in high dimensions
- Message updates rendered tractable using a Gaussian approximation
- Convergence achieved by imposing consistency requirements derived from statistical mechanics analysis.
- Inference of posterior marginals exhibits near optimal performance compared to theory.
- Outperforms two other recently published algorithms.
- Approximation to Marginal Likelihood also available.

# Summary

- Novel message passing algorithm for Sparse Bayesian PCA in high dimensions
- Message updates rendered tractable using a Gaussian approximation
- Convergence achieved by imposing consistency requirements derived from statistical mechanics analysis.
- Inference of posterior marginals exhibits near optimal performance compared to theory.
- Outperforms two other recently published algorithms.
- Approximation to Marginal Likelihood also available.

# Summary

- Novel message passing algorithm for Sparse Bayesian PCA in high dimensions
- Message updates rendered tractable using a Gaussian approximation
- Convergence achieved by imposing consistency requirements derived from statistical mechanics analysis.
- Inference of posterior marginals exhibits near optimal performance compared to theory.
- Outperforms two other recently published algorithms.
- Approximation to Marginal Likelihood also available.

# Summary

- Novel message passing algorithm for Sparse Bayesian PCA in high dimensions
- Message updates rendered tractable using a Gaussian approximation
- Convergence achieved by imposing consistency requirements derived from statistical mechanics analysis.
- Inference of posterior marginals exhibits near optimal performance compared to theory.
- Outperforms two other recently published algorithms.
- Approximation to Marginal Likelihood also available.

# The Future

- Hyperparameter estimation using Marginal likelihood.
- Extension to multiple factors:
  - Relatively straightforward for orthogonal factors.  
(but will require efficient hyperparameter estimation).
  - For non-orthogonal factors the best approach is a subject of on-going research.

# The Future

- Hyperparameter estimation using Marginal likelihood.
- Extension to multiple factors:
  - Relatively straightforward for orthogonal factors.  
(but will require efficient hyperparameter estimation).
  - For non-orthogonal factors the best approach is a subject of on-going research.

## Explore further

Matlab code available from: <http://www.cs.man.ac.uk/~sharpk>