

# Identifying Genetic Material with Compressed Sensing Techniques

Jamie Dougherty

8th October 2018

# Acknowledgements

With thanks to

- The SRIM Scheme
- Dr Sergio Bacallado

# Outline

- 1 Introduction to Metagenomics
- 2 Problem Spec
- 3 Overview of Approximate Message Passing
- 4 Application to Matrix Factorisation

# Metagenomics



- Sequence DNA of mixture genomes in large groups of micro-organisms.

# Metagenomics



- Sequence DNA of mixture genomes in large groups of micro-organisms.
- Conventional cultivation methods have trouble in these environments.

# Metagenomics



- Sequence DNA of mixture genomes in large groups of micro-organisms.
- Conventional cultivation methods have trouble in these environments.
- Applications in medicine - treatment of infections.

# Metagenomics



- Sequence DNA of mixture genomes in large groups of micro-organisms.
- Conventional cultivation methods have trouble in these environments.
- Applications in medicine - treatment of infections.
- Still not well understood - room for new methods.

- Problem can be framed as the reconstruction of a noisy low-rank matrix with structured factors.
- Using prior information from factors can improve error. Still a challenge to derive algorithms that use this information.



# Approximate Message Passing (AMP)

- Message passing algorithms allow efficient marginalising of distributions with a certain structure of dependencies.

# Approximate Message Passing (AMP)

- Message passing algorithms allow efficient marginalising of distributions with a certain structure of dependencies.
- AMP does away with this structure. Just requires high number of dependencies.

# Approximate Message Passing (AMP)

- Message passing algorithms allow efficient marginalising of distributions with a certain structure of dependencies.
- AMP does away with this structure. Just requires high number of dependencies.
- Some limit theorems are then used to derive a different iterative procedure (AMP).

# Approximate Message Passing (AMP)

- Message passing algorithms allow efficient marginalising of distributions with a certain structure of dependencies.
- AMP does away with this structure. Just requires high number of dependencies.
- Some limit theorems are then used to derive a different iterative procedure (AMP).
- Some motivation comes from Statistical Physics (Replica and Cavity methods).

# Approximate Message Passing (AMP)

- Message passing algorithms allow efficient marginalising of distributions with a certain structure of dependencies.
- AMP does away with this structure. Just requires high number of dependencies.
- Some limit theorems are then used to derive a different iterative procedure (AMP).
- Some motivation comes from Statistical Physics (Replica and Cavity methods).

## Example

Problem: find sparse  $x$  given  $y = Ax$ ,  $A \in \mathbb{R}^{N \times n}$ ,  $n \ll N$

AMP iterate:

$$\begin{aligned}x^{t+1} &= \eta(x^t + A^T z^t) \\ z^t &= y - Ax^t + \eta'(x^{t-1} + A^T z^{t-1})\end{aligned}$$

# Interesting Points

- Structure of the problem is exploited.
- Updates are computationally efficient.

# Interesting Points

- Structure of the problem is exploited.
- Updates are computationally efficient.
- Engenders a separate set of equations known as **State Evolution**:
  - 1 Approximates error of iterative procedure.
  - 2 Can be used by algorithm to modify optimal mapping to target structure.

# Interesting Points

- Structure of the problem is exploited.
- Updates are computationally efficient.
- Engenders a separate set of equations known as **State Evolution**:
  - 1 Approximates error of iterative procedure.
  - 2 Can be used by algorithm to modify optimal mapping to target structure.

## Example with state evolution

Problem: find sparse  $x$  given  $y = Ax$ ,  $A \in \mathbb{R}^{N \times n}$ ,  $n \ll N$

AMP update:

$$\begin{aligned}x^{t+1} &= \eta(x^t + A^T z^t, \hat{\sigma}^t) \\z^t &= y - Ax^t + \eta'(x^{t-1} + A^T z^{t-1}, \hat{\sigma}^{t-1})\end{aligned}$$

State evolution:

$$\hat{\sigma}^t = \hat{\sigma}^t \eta'(A^T z^{t-1} + x^t, \hat{\sigma}^{t-1})$$



# AMP for matrix factorisation

## Problem

Find  $\mathbf{x}_0 \sim \mu_0$

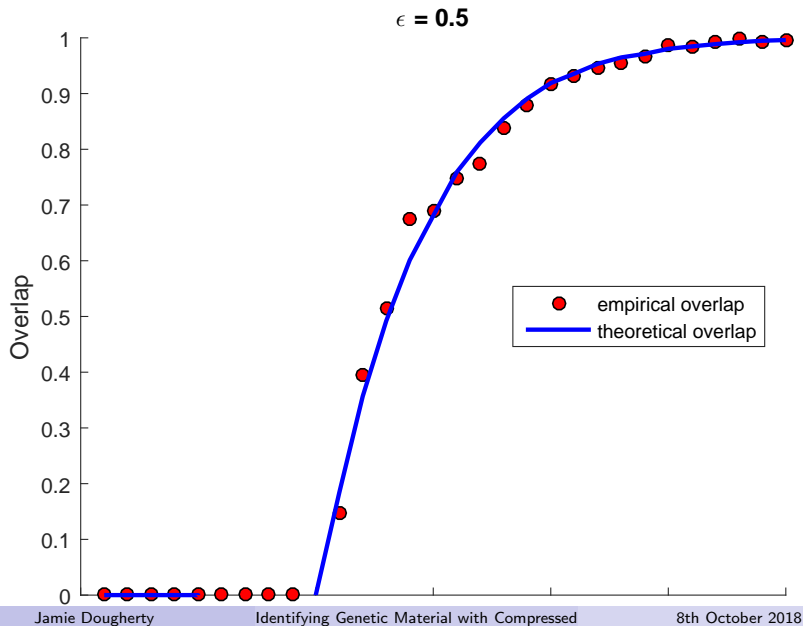
$$A = \frac{\lambda}{n} \mathbf{x}_0 \mathbf{x}_0^T + W$$

$W$  distributed as a Gaussian orthogonal ensemble.

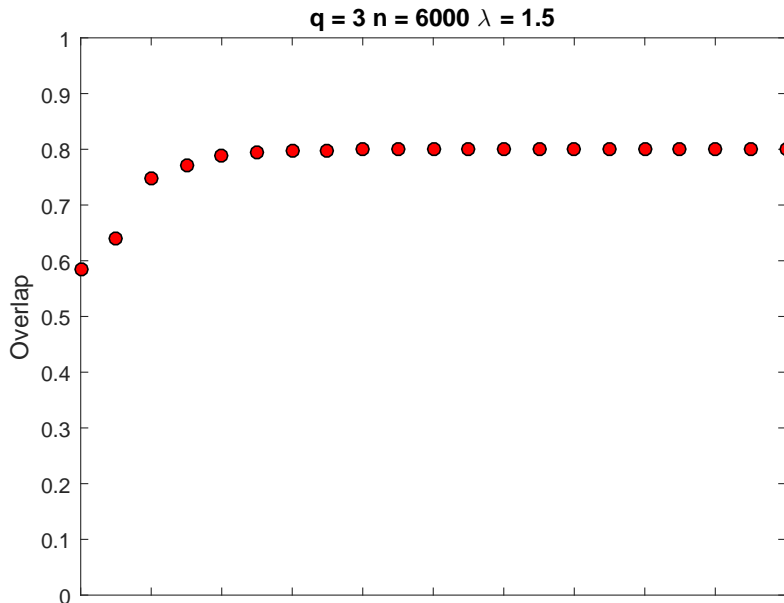
## AMP Solution

$$\mathbf{x}^{t+1} = A f_t(\mathbf{x}^t) - \mathbf{b}_t f_{t-1}(\mathbf{x}^{t-1})$$

# Results in rank one case



# How iterates improve over time



# Problems

What happens when we try and apply these techniques to our Metagenomic problem?

- Algorithm becomes more complex.

What happens when we try and apply these techniques to our Metagenomic problem?

- Algorithm becomes more complex.
- Much of the literature requires assumptions that don't apply.

# Problems

What happens when we try and apply these techniques to our Metagenomic problem?

- Algorithm becomes more complex.
- Much of the literature requires assumptions that don't apply.

Stochastically:

$$f_{\text{seq}=\text{seq}(t+1)} = f_{\text{seq}(t)} + 1 - \frac{1}{\sqrt{N}} \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \text{seq}(t+1) + O(1/\sqrt{N}), \quad (90)$$

$$c_{\text{seq}=\text{seq}(t+1)} = c_{\text{seq}(t)} + O\left(\frac{1}{\sqrt{N}}\right), \quad s_{\text{seq}=\text{seq}(t+1)} = s_{\text{seq}(t)} + 1 + O\left(\frac{1}{\sqrt{N}}\right), \quad (91)$$

$$\text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) = \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) + \frac{1}{\sqrt{N}} \text{seq}(t) \text{seq}(t) \text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) + O(1/\sqrt{N}). \quad (92)$$

From these equations we obtain the GAMP algorithm for matrix factorization

$$V_{\text{seq}}^{\text{seq}(t)} = \frac{1}{N} \sum_i \text{seq}(i) \text{seq}(i) + c_{\text{seq}(t)} \text{seq}(i) + s_{\text{seq}(t)}^2 \text{seq}(i), \quad (93)$$

$$c_{\text{seq}}^{\text{seq}(t)} = \frac{1}{N} \sum_i \text{seq}(i) \text{seq}(i) - \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \frac{1}{N} \sum_i \left\{ \text{seq}(i) \text{seq}(i) - 1 \right\} \text{seq}(i) + c_{\text{seq}(t)} \text{seq}(i) - 1 \text{seq}(i), \quad (94)$$

$$(V_{\text{seq}}^{\text{seq}(t)})^{-1} = \frac{1}{N} \sum_i \left\{ -\text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \left[ \text{seq}(i) + c_{\text{seq}(t)} \right] - \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \text{seq}(i) \right\}, \quad (95)$$

$$s_{\text{seq}}^{\text{seq}(t)} = \frac{1}{N} \sum_i \left\{ \frac{1}{N} \sum_j \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \text{seq}(i) - c_{\text{seq}(t)} \frac{1}{N} \sum_j \text{seq}(i) \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \right. \\ \left. - c_{\text{seq}(t)} - 1 \right\} \frac{1}{N} \sum_j \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \text{seq}(i) - c_{\text{seq}(t)} \frac{1}{N} \sum_j \text{seq}(i) \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)), \quad (96)$$

$$(V_{\text{seq}}^{\text{seq}(t)})^{-1} = \frac{1}{N} \sum_i \left\{ -\text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \left[ \text{seq}(i) + c_{\text{seq}(t)} \right] - \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \text{seq}(i) \right\}, \quad (97)$$

$$W_{\text{seq}}^{\text{seq}(t)} = s_{\text{seq}}^{\text{seq}(t)} \left( \frac{1}{N} \sum_i \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \text{seq}(i) - f_{\text{seq}(t)} \frac{1}{N} \sum_i \text{seq}(i) \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \right. \\ \left. - f_{\text{seq}(t)} - 1 \right) \frac{1}{N} \sum_i \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)) \text{seq}(i) - f_{\text{seq}(t)} \frac{1}{N} \sum_i \text{seq}(i) \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)), \quad (98)$$

$$s_{\text{seq}(t+1)} = f_{\text{seq}(t)} (V_{\text{seq}}^{\text{seq}(t)})^{-1}, \quad c_{\text{seq}(t+1)} = f_{\text{seq}(t)} (V_{\text{seq}}^{\text{seq}(t)}), \quad (99)$$

$$f_{\text{seq}(t+1)} = f_{\text{seq}(t)} (W_{\text{seq}}^{\text{seq}(t)}), \quad s_{\text{seq}(t+1)} = f_{\text{seq}(t)} (W_{\text{seq}}^{\text{seq}(t)}). \quad (100)$$

The initial condition for iterations are

$$s_{\text{seq}(t)} = 0, \quad c_{\text{seq}(t)} = 0, \quad f_{\text{seq}(t)} = 0, \quad s_{\text{seq}(t)} = 0, \quad (101)$$

$$f_{\text{seq}(t)} = \sqrt{N} \int \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)), \quad s_{\text{seq}(t)} = N \int \text{seq}(\text{seq}(t), \text{seq}(t), V_{\text{seq}}^{\text{seq}(t)}(t)). \quad (102)$$

In order to compute  $s_{\text{seq}(t)}^{\text{seq}(t)}$ ,  $c_{\text{seq}(t)}^{\text{seq}(t)}$  and  $W_{\text{seq}(t)}^{\text{seq}(t)}$  use the above equations as if  $f_{\text{seq}(t)} = 0$  and  $s_{\text{seq}(t)} = 0$ .  
The interpretation of the terms in the GAMP for matrix factorization is the following:  $s_{\text{seq}(t)}^{\text{seq}(t)}$  is the mean of the

# Solutions?

How do we resolve some of these problems?

# Solutions?

How do we resolve some of these problems?

- Don't try for Bayes optimal.



# Solutions?

How do we resolve some of these problems?

- Don't try for Bayes optimal.
- Choose non-linearities based on intuition.

# Solutions?

How do we resolve some of these problems?

- Don't try for Bayes optimal.
- Choose non-linearities based on intuition.
- Lot's of experimentation!

# Solutions?

How do we resolve some of these problems?

- Don't try for Bayes optimal.
- Choose non-linearities based on intuition.
- Lot's of experimentation!

Don't have working algorithm, more testing needed.

# Summary

# Summary

- Metagenomic problem seems tractable.

# Summary

- Metagenomic problem seems tractable.
- AMP is very efficient when used correctly.

# Summary

- Metagenomic problem seems tractable.
- AMP is very efficient when used correctly.
- However it is difficult to use correctly in this context.

# Summary

- Metagenomic problem seems tractable.
- AMP is very efficient when used correctly.
- However it is difficult to use correctly in this context.
- Better intuitions and more testing could generate a solution.