

A High-bias Low-variance Introduction to Approximate Bayesian Inference

Xiangming Meng

Postdoctoral Researcher

Approximate Bayesian Inference (ABI) Team
RIKEN Center for Advanced Intelligence Project

Oct. 11, 2019

Tokyo Institute of Technology

Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

Physics Reports

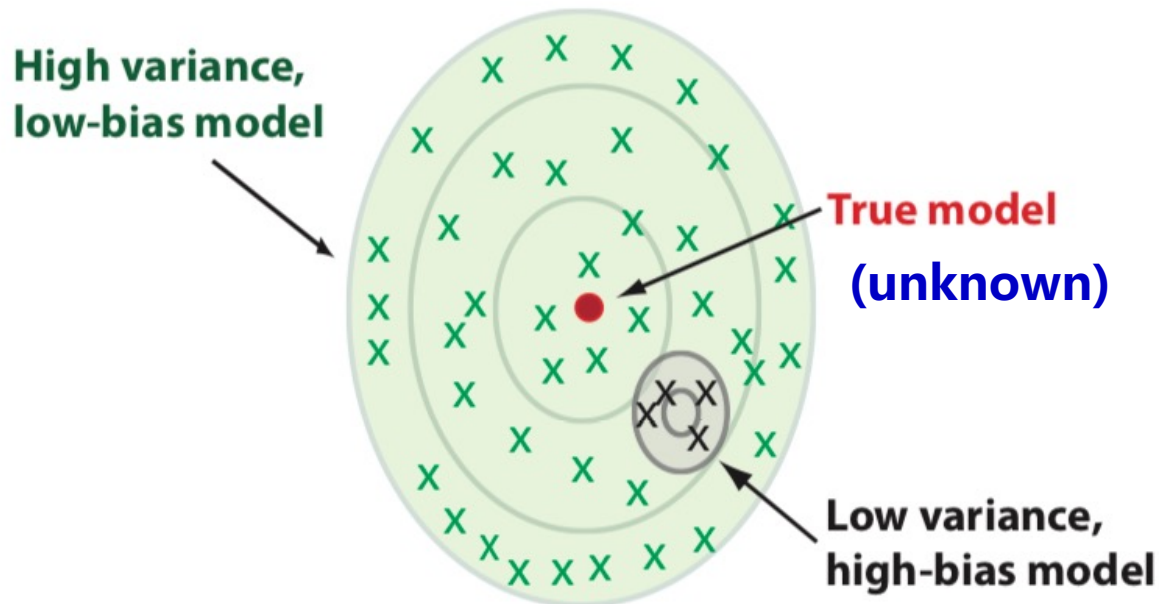
journal homepage: www.elsevier.com/locate/physrep

A high-bias, low-variance introduction to Machine Learning for physicists



Pankaj Mehta^a, Marin Bukov^{b,*}, Ching-Hao Wang^a, Alexandre G.R. Day^a,
Clint Richardson^a, Charles K. Fisher^c, David J. Schwab^d

100 pages !



Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

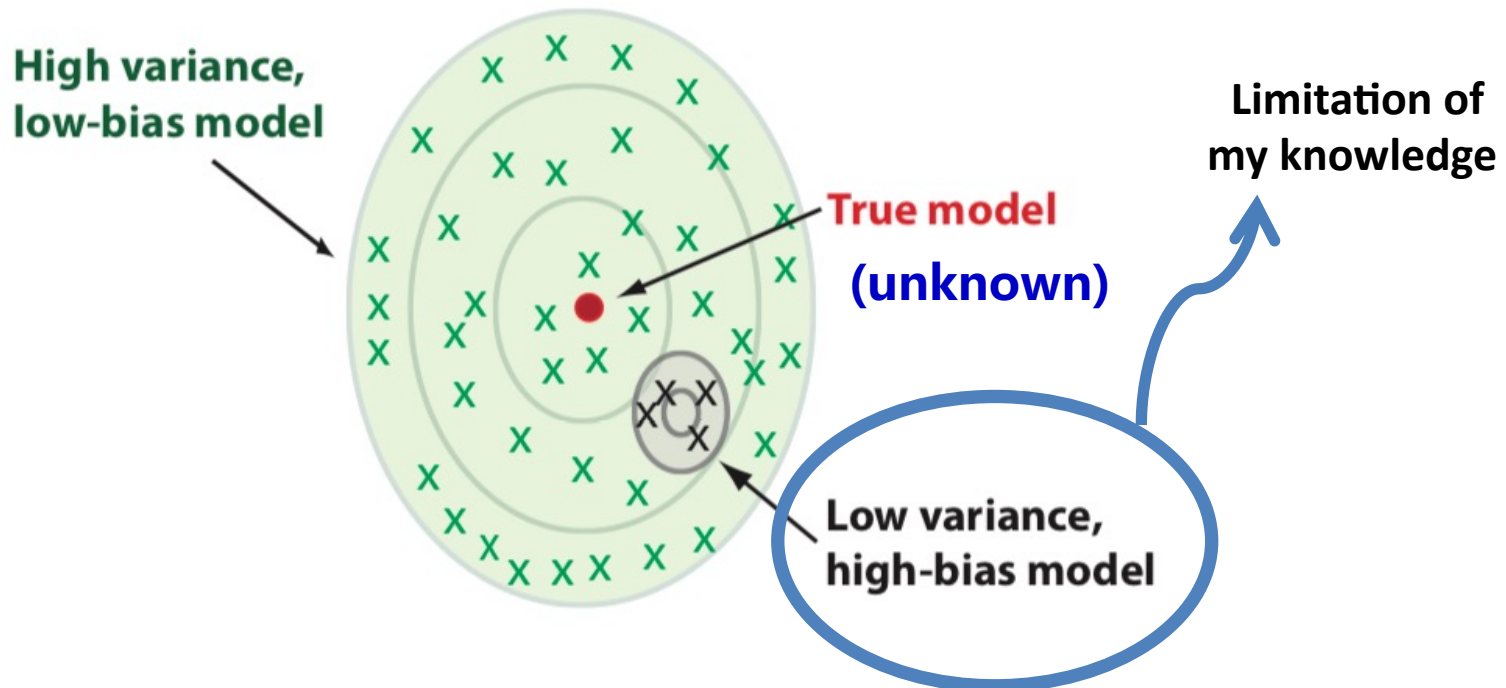
Physics Reports

journal homepage: www.elsevier.com/locate/physrep

A high-bias, low-variance introduction to Machine Learning for physicists



Pankaj Mehta^a, Marin Bukov^{b,*}, Ching-Hao Wang^a, Alexandre G.R. Day^a,
Clint Richardson^a, Charles K. Fisher^c, David J. Schwab^d



Outline

- **Background**
- **Variational Inference**
- **Expectation Propagation**
- **A Unified EP Perspective on AMP and its extensions**
- **Conclusion**

Background

- 3 little princes from 3 planets

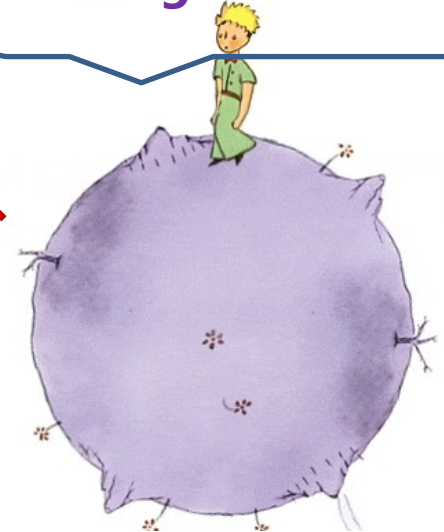
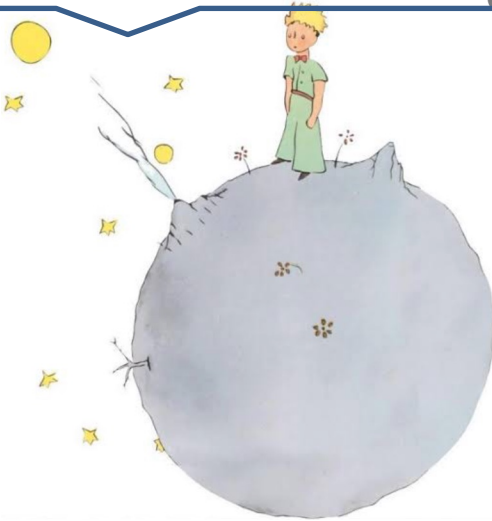


Hi, I study disordered systems in statistical physics.

Hi, I study coding and compressed sensing.

Hi, I study machine learning.

Statistical Physics Planet A



Information theory Planet B

Computer Science Planet C

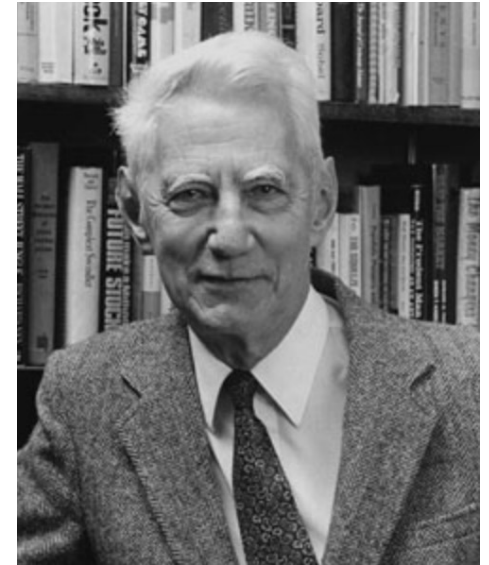
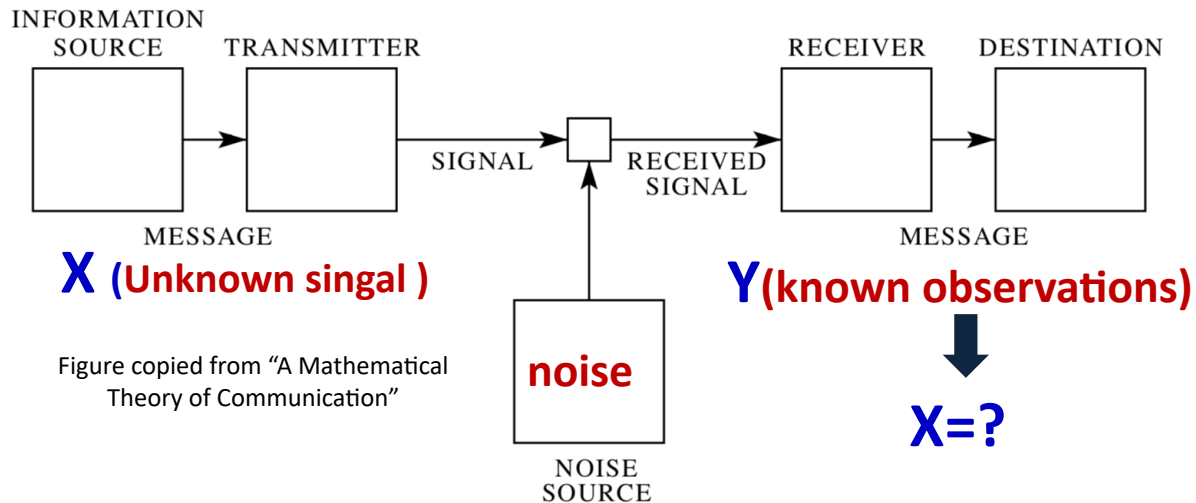
After a long time discussion, it turns out that they are studying similar problems using different languages

Background

□ Communication

*“The fundamental problem of communication is that of reproducing at one point **either exactly or approximately** a message selected at another point”*

—Shannon (1948)



Claude Elwood Shannon
(1916-2001)

Fig 1. Schematic diagram of a general communication system

Background

□ Communication

*“The fundamental problem of communication is that of reproducing at one point **either exactly or approximately** a message selected at another point”*

—Shannon (1948)

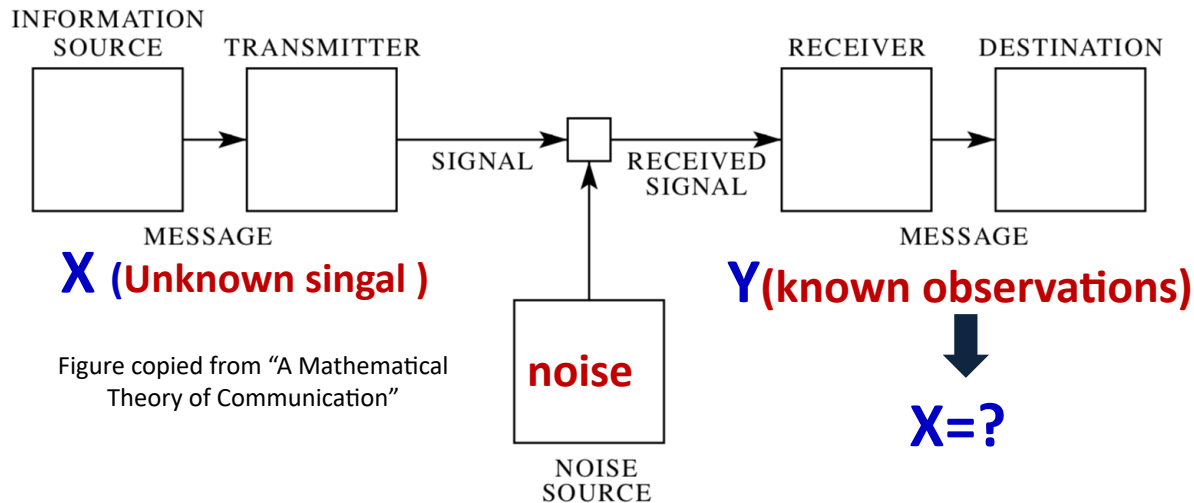
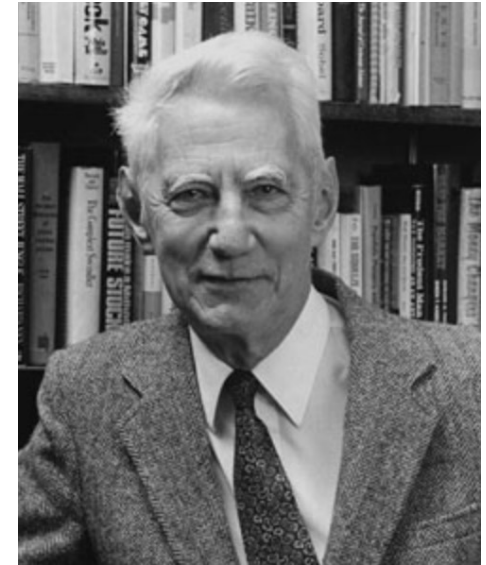


Figure copied from “A Mathematical Theory of Communication”



Claude Elwood Shannon
(1916-2001)

Fig 1. Schematic diagram of a general communication system

- **Q1: How to quantize information?**
Entropy $H = -\sum_k p_k \log p_k$
- **Q2: What is the capacity of a communication system?**
Shannon Formula: $C = W \cdot \log(1+S/N)$ maximum rate
- **Q3: How to approach the capacity?**
Channel coding (Turbo code, LDPC code, Polar code in 5G)

'You should call it entropy, for two reasons. In the first place your uncertainty function **has been used in statistical mechanics** under that name, so it already has a name. In the second place, and more important, **no one really knows what entropy really is**, so in a debate you will always have the advantage.'

—John von Neumann

Background

□ Communication

Received
Message y

Tokye Institote of Technalogy

?

Background

□ Communication

Received
Message y

Tokye Institote of Technalogy

Corrected
Message x

Tokyo Institute of Technology

Background

□ Communication

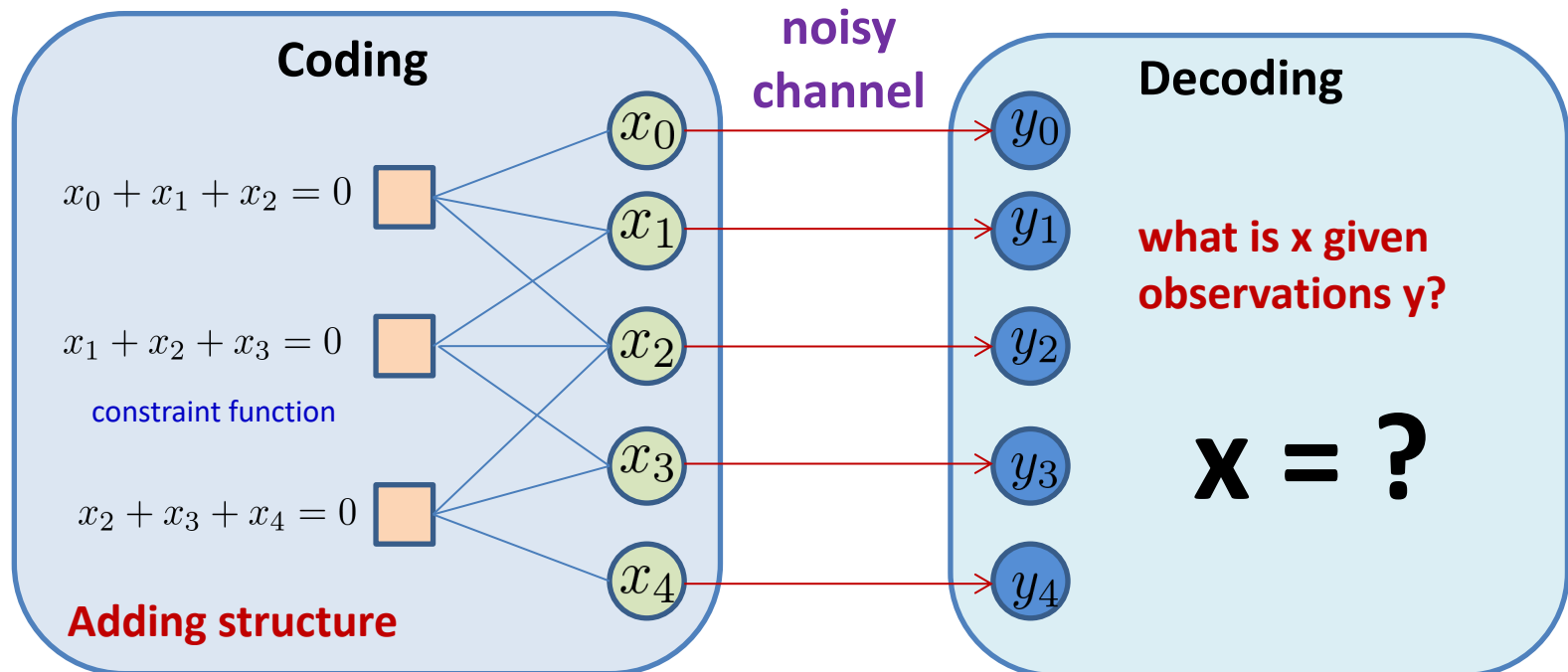
Received
Message y

Tokye Institote of Technalogy

Corrected
Message x

Tokyo Institute of Technology

There is structure within the transmitted codes.



Background

□ Compressed Sensing



Raw: 15MB



JPEG: 150KB

- **Massive data acquisition**
- **Most of the data is redundant**
- **Wasteful measurements**
- **Could we acquire images using less/efficient measurements?**

Background

Compressed Sensing

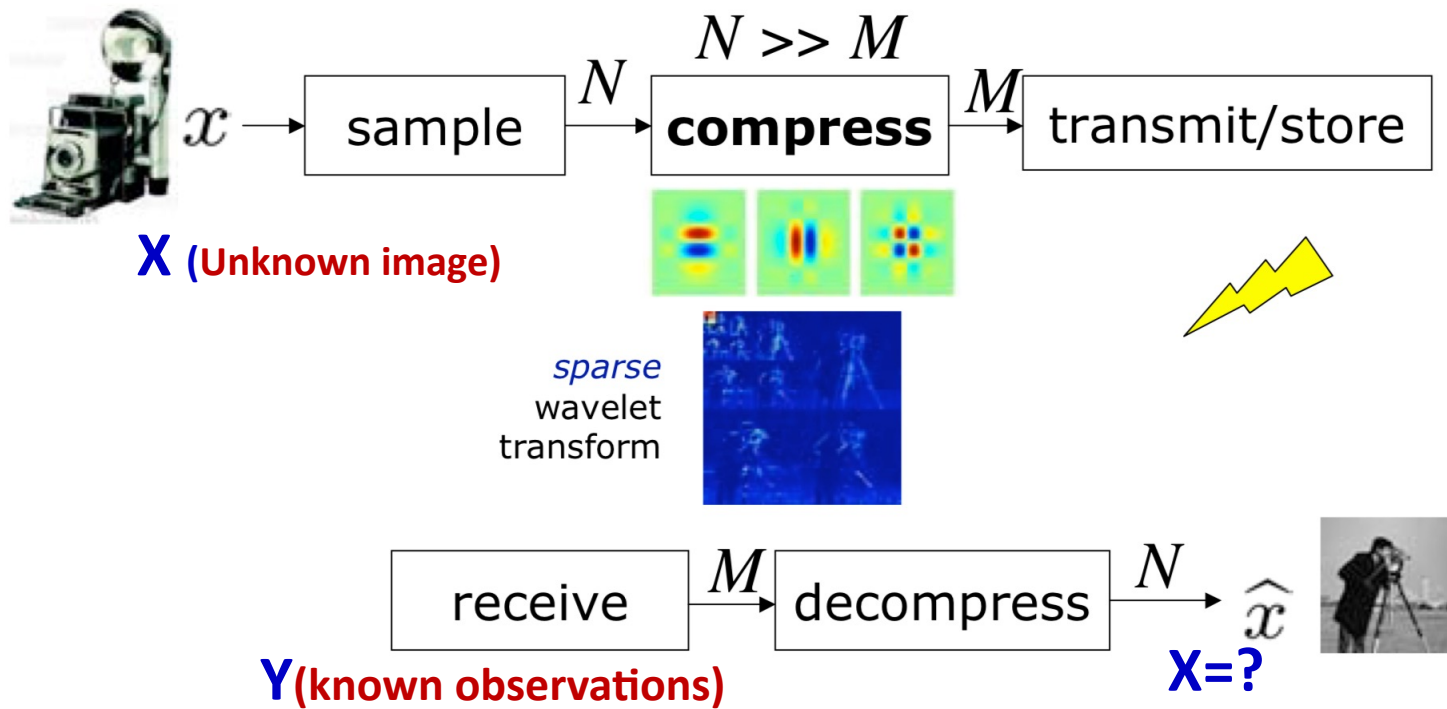


Raw: 15MB



JPEG: 150KB

- Massive data acquisition
- Most of the data is redundant
- Wasteful measurements
- **Could we acquire images using less/efficient measurements?**



Background

□ Bayesian Perspective

$\mathbf{x} \sim p(\mathbf{x})$
Unknown Signal



\mathbf{y}

Known Observations

what is \mathbf{x} ?



Background

□ Bayesian Perspective

Prior distribution

$$\mathbf{x} \sim p(\mathbf{x})$$

Unknown Signal

likelihood distribution

$$p(\mathbf{y}|\mathbf{x})$$

\mathbf{y}

Known Observations

Posterior distribution

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}$$

Bayes' rule

evidence
(partition
function)

$$\mathbf{Z} \triangleq p(\mathbf{y}) = \sum_{\mathbf{x}} p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$$



Thomas Bayes (1702-1761)

Background

□ Bayesian Perspective

Prior distribution

$$\mathbf{x} \sim p(\mathbf{x})$$

Unknown Signal

likelihood distribution

$$p(\mathbf{y}|\mathbf{x})$$

\mathbf{y}

Known Observations

Posterior distribution

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}$$

Bayes' rule

evidence
(partition
function)

$$\mathbf{Z} \triangleq p(\mathbf{y}) = \sum_{\mathbf{x}} p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$$



Thomas Bayes (1702-1761)

□ Goal

marginal distribution $p(x_i|\mathbf{y}) = \sum_{x_j \neq x_i} p(\mathbf{x}|\mathbf{y})$

$i = 1, \dots, N$

posterior mean $\hat{x}_i = E(x_i|\mathbf{y}) = \sum_{x_i} x_i p(x_i|\mathbf{y})$

Background

Bayesian Perspective

Prior distribution

$$\mathbf{x} \sim p(\mathbf{x})$$

Unknown Signal

likelihood distribution

$$p(\mathbf{y}|\mathbf{x})$$

\mathbf{y}

Known Observations

Posterior distribution

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}$$

Bayes' rule

evidence
(partition
function)

$$p(\mathbf{y}) = \sum_{\mathbf{x}} p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$$



Thomas Bayes (1702-1761)

Goal

marginal distribution $p(x_i|\mathbf{y}) = \sum_{x_j \neq x_i} p(\mathbf{x}|\mathbf{y})$

posterior mean $\hat{x}_i = E(x_i|\mathbf{y}) = \sum_{x_i} x_i p(x_i|\mathbf{y})$

Curse of
Dimensionality!

e.g., N spins, $O(2^N)$

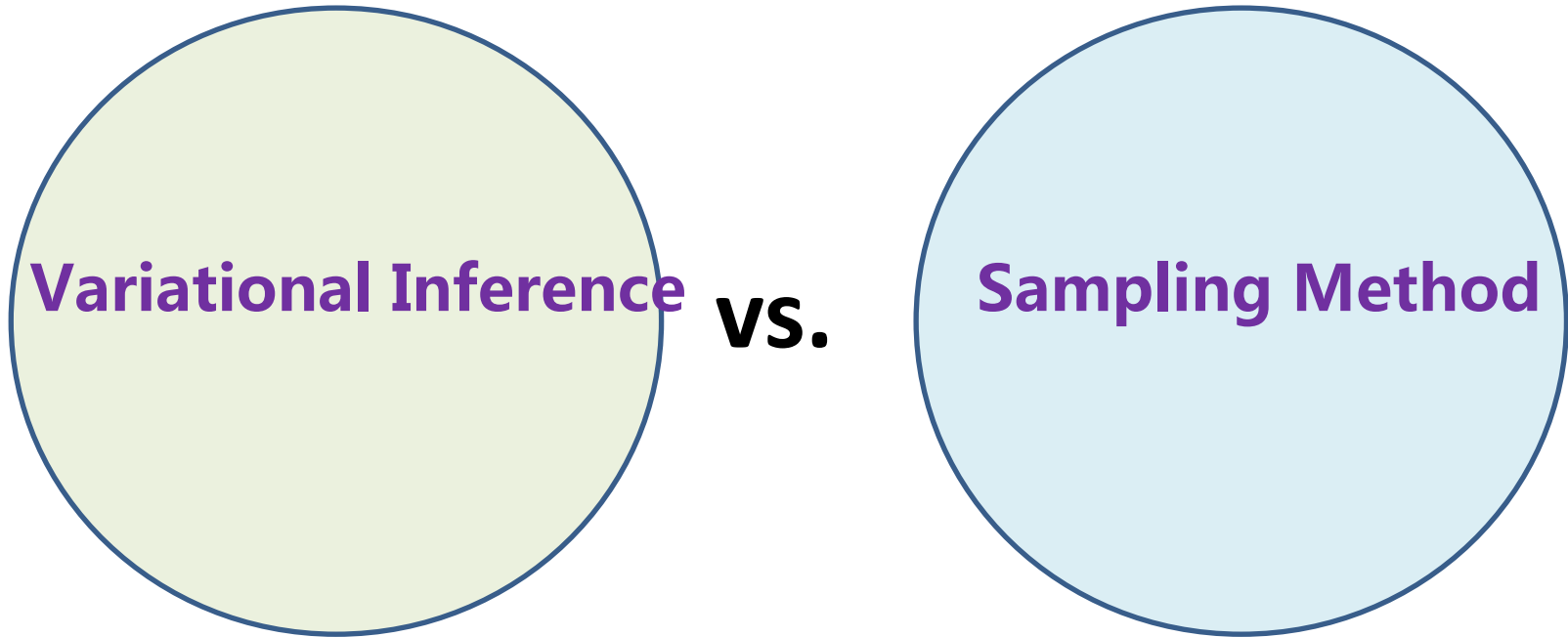
We have to resort to approximate Bayesian Inference methods

Outline

- Background
- **Variational Inference**
- Expectation Propagation
- A Unified EP Perspective on AMP and its extensions
- Conclusion

Variational Inference

□ Two Common Approaches of Approximate Inference

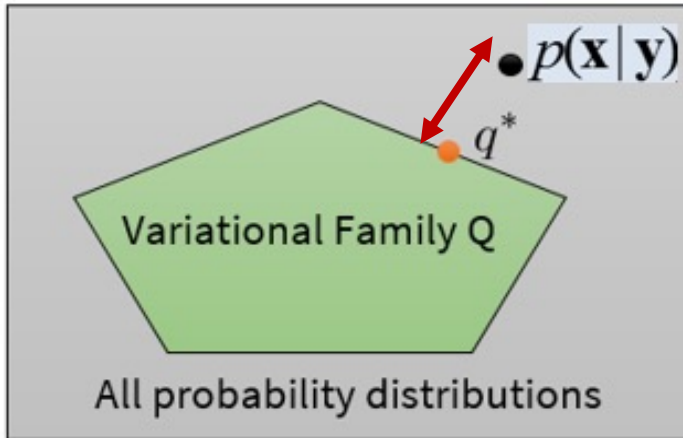


- **Deterministic**
- **Biased**
- **Scalable**

- **Stochastic**
- **Unbiased**
- **Non-scalable**

Variational Inference

□ Basic Principle

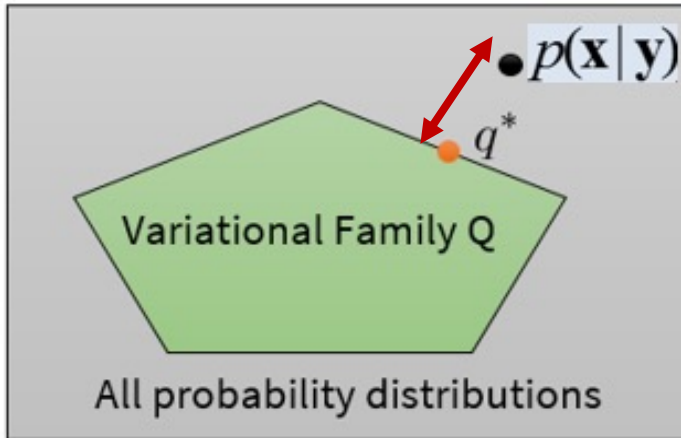


To approximate complicated target distribution p with a *simple distribution* q *as close to p as possible*

$$q \approx p$$

Variational Inference

□ Basic Principle



To approximate complicated target distribution p with a **simple distribution** q **as close to p as possible**

$$q \approx p$$

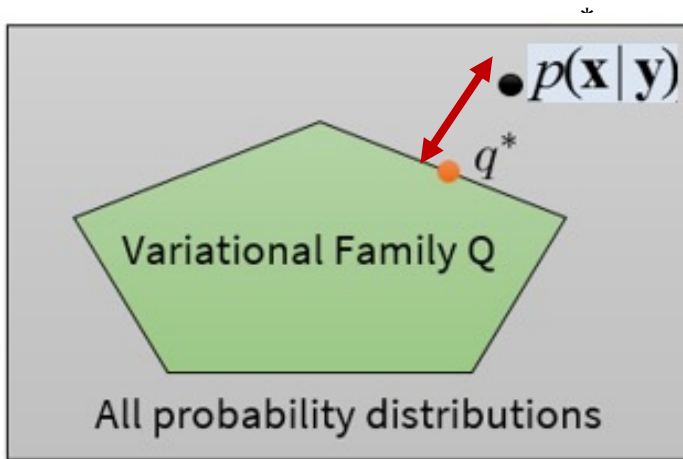
Optimization problem $q^* = \arg \min_{q \in Q} KL(q(\mathbf{x}) || p(\mathbf{x}|\mathbf{y}))$

KL divergence
"distance"

$$KL(q||p) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

Variational Inference

□ Basic Principle



$$\min_{q \in Q} KL(q(\mathbf{x}) || p(\mathbf{x} | \mathbf{y}))$$

To approximate complicated target distribution p with a simple distribution q as close to p as possible

$$q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

$$q \approx p$$

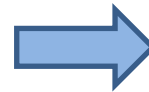
Optimization problem

KL divergence

“distance”

- Non-negativity of KL

$KL(p || q) \geq 0$ and $KL(p || q) = 0$ if and only if $p = q$



“Gibbs inequality”

- Non-symmetry of KL

$KL(p || q)$ is *not* equal to $KL(q || p)$



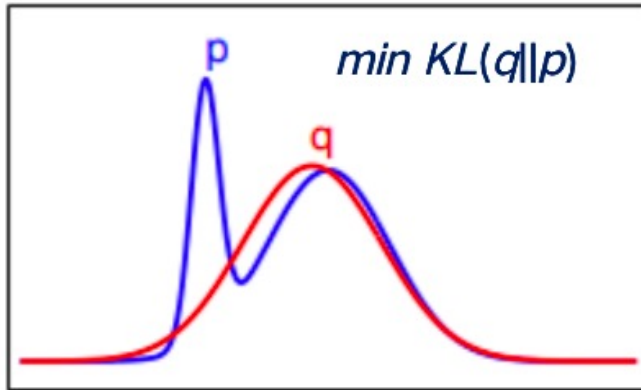
qseudo distance

Variational Inference

□ Basic Principle

- KL divergence

$$KL(q||p) = \int q(x) \log \frac{q(x)}{p(x)}$$



≠

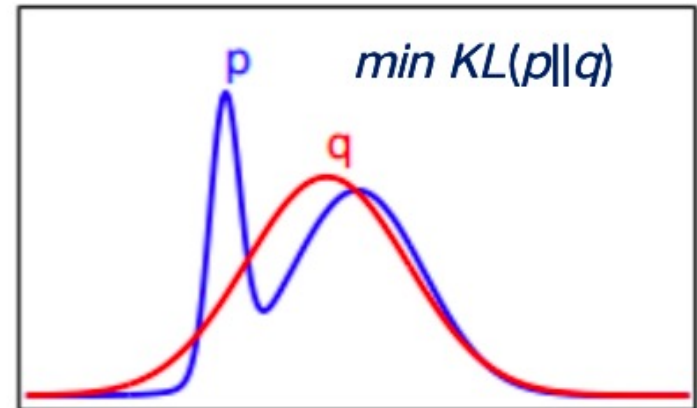


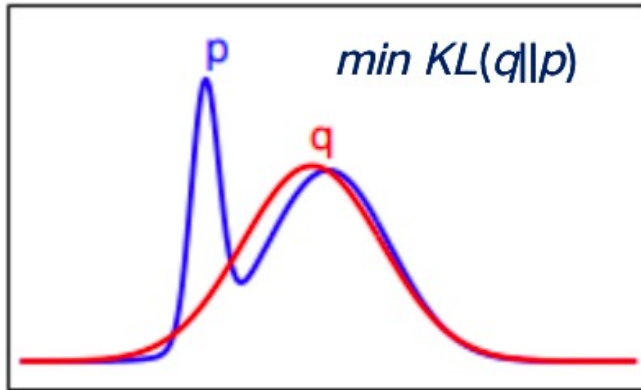
figure copied from [Bishop06]

Variational Inference

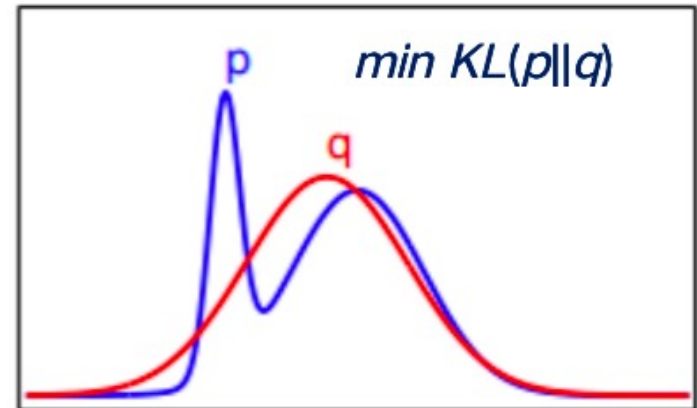
□ Basic Principle

- KL divergence

$$KL(q||p) = \sum q(x) \log \frac{q(x)}{p(x)}$$



≠



goal

Remember that VI uses $KL(q||p)$

figure caption

dilemma!

To calculate the KL divergence, **we must know the target distribution in advance**, which is our primary goal!

Variational Inference

□ ELBO bound

$$\begin{aligned} KL(q(\mathbf{x})||p(\mathbf{x}|\mathbf{y})) &= \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x}|\mathbf{y})} = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x},\mathbf{y})} && \text{Bayes' Rule} \\ &= \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x}) - \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x},\mathbf{y}) + \log p(\mathbf{y}) && \text{Expansion} \\ &\geq 0 && \text{"Gibbs inequality"} \end{aligned}$$

Variational Inference

□ ELBO bound

$$\begin{aligned} KL(q(\mathbf{x})||p(\mathbf{x}|\mathbf{y})) &= \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x}|\mathbf{y})} = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x},\mathbf{y})} && \text{Bayes' Rule} \\ &= \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x}) - \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x},\mathbf{y}) + \log p(\mathbf{y}) && \text{Expansion} \\ &\geq 0 && \text{"Gibbs inequality"} \end{aligned}$$

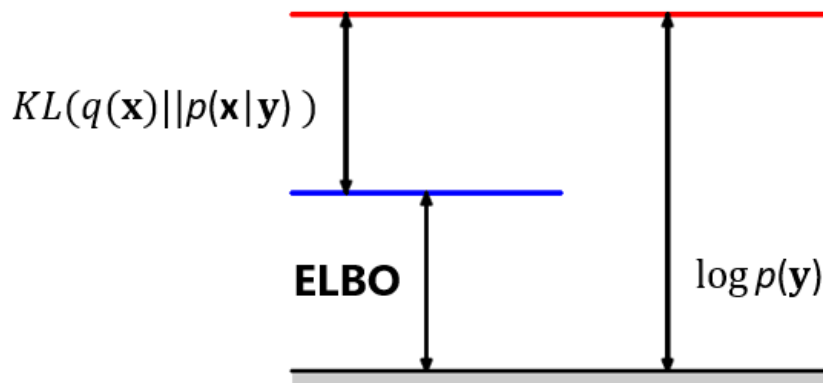
$$\underbrace{\log p(\mathbf{y})}_{\text{Log Partition function}} \geq \underbrace{\sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x},\mathbf{y}) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})}_{\text{Evidence Lower Bound (ELBO)}}$$

Variational Inference

□ ELBO bound

$$\begin{aligned} KL(q(\mathbf{x})||p(\mathbf{x}|\mathbf{y})) &= \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x}|\mathbf{y})} = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x},\mathbf{y})} && \text{Bayes' Rule} \\ &= \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x}) - \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x},\mathbf{y}) + \log p(\mathbf{y}) && \text{Expansion} \\ &\geq 0 && \text{"Gibbs inequality"} \end{aligned}$$

$$\underbrace{\log p(\mathbf{y})}_{\substack{\text{Log} \\ \text{Partition function}}} \geq \underbrace{\sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x},\mathbf{y}) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})}_{\text{Evidence Lower Bound (ELBO)}}$$

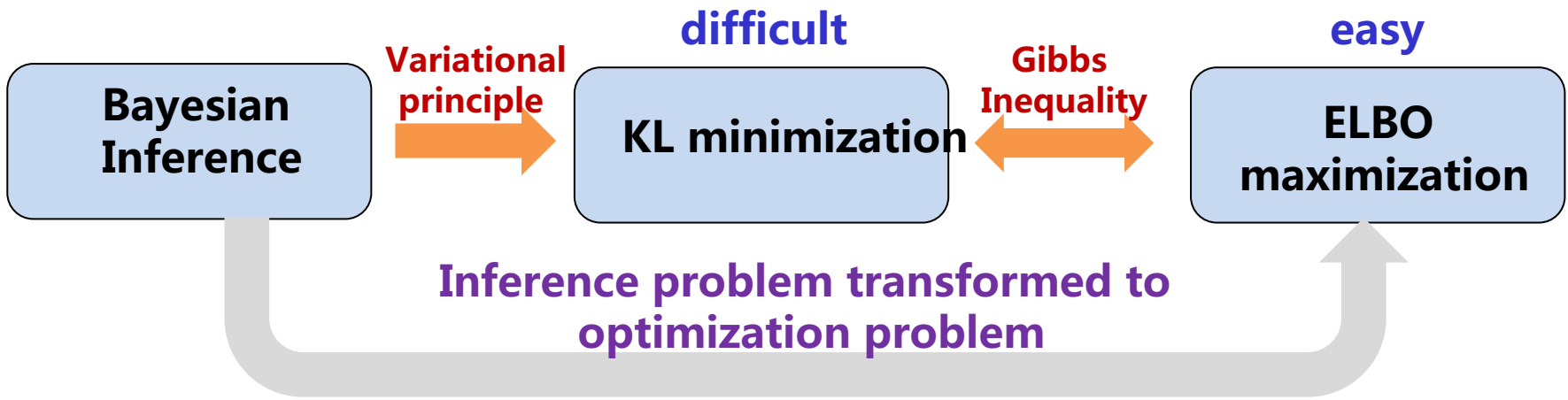


minimize KL = maximize ELBO

Variational Inference

□ ELBO bound

Big Picture of VI



Variational Inference

□ Analogy between different planets

Computer Science Planet

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}$$

evidence lower bound ELBO $\triangleq \langle \log p(\mathbf{x}, \mathbf{y}) \rangle_q + H(q(\mathbf{x})) \leq \log p(\mathbf{y})$

Statistical physics Planet

$$p(\mathbf{s}|\beta, J) = \frac{e^{-\beta E(\mathbf{s}, J)}}{Z(\beta, J)}$$

free energy

variational free energy $\beta F_q(\mathbf{J}) \triangleq \beta \langle E(\mathbf{s}, J) \rangle_q - H(q(\mathbf{s})) \geq -\log Z(\beta, J)$

Statistical Physics

Computer Science/Information Theory

Spins/degrees of freedom \mathbf{s}	Hidden variables/signal of interest \mathbf{x}
Couplings/quenched disorder \mathbf{J}	Data observations \mathbf{y}
Boltzmann factor $e^{-\beta E(\mathbf{s}, J)}$	Joint distribution $p(\mathbf{x}, \mathbf{y})$
Partition function $Z(\beta, J)$	Evidence $p(\mathbf{y})$
Energy $\beta E(\mathbf{s}, J)$	Negative log-joint distribution $-\log p(\mathbf{x}, \mathbf{y})$
Free Energy $-\log Z(\beta, J)$	Negative log evidence $-\log p(\mathbf{y})$
Variational distribution $q(\mathbf{s})$	Variational distribution $q(\mathbf{x} \mathbf{y})$
Variational free energy $\beta F_q(\mathbf{J})$	Negative ELBO -ELBO

Table modified from Table I in [Mehta et al 19]

Variational Inference

□ Why transforming inference to optimization?

$$\max \text{ELBO} = \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x}, y) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})$$

There are a bunch of optimization methods we could leverage!

- **different choice of q**

- ✓ **structure:** mean-field, Bethe, etc.
- ✓ **parametric:** Gaussian, neural network, etc.

- **different optimization methods**

- ✓ **coordinate descent**
- ✓ **gradient descent**
- ✓ **stochastic gradient descent**
- ✓ **natural gradient descent**
- ✓

Different combinations lead to different inference algorithms

Mean-filed Approximation

□ Mean Field Approximation

$$\max \text{ELBO} = \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x}, \mathbf{y}) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})$$

Mean Field structure

$$q(\mathbf{x}) = \prod_i q(x_i)$$

different variables
are independent

ELBO

$$= \sum_{\mathbf{x}} \prod_i q(x_i) \log p(\mathbf{x}, \mathbf{y}) - \sum_i q(x_i) \log q(x_i)$$

Mean-field Approximation

□ Mean Field Approximation

$$\max \text{ELBO} = \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x}, \mathbf{y}) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})$$

Mean Field structure

$$q(\mathbf{x}) = \prod_i q(x_i)$$

different variables
are independent

Using **coordinate descent optimization**, we obtain the **variational message passing (VMP) algorithm**:

ELBO

$$= \sum_{\mathbf{x}} \prod_i q(x_i) \log p(\mathbf{x}, \mathbf{y}) - \sum_i q(x_i) \log q(x_i)$$

Input: A model $p(\mathbf{x}, \mathbf{y})$, a dataset

Output: $q(\mathbf{x}) = \prod_i q(x_i)$

- 1: Initialize variational factors $q(\mathbf{x})$
- 2: **while** the ELBO has not converged **do**
- 3: **for** $\tilde{l} \in 1, 2, \dots, d$ **do**
- 4: $q(x_{\tilde{l}}) \propto \exp \left\{ \mathbb{E}_{\prod_{j \neq \tilde{l}} q(x_j)} [\log p(\mathbf{x}, \mathbf{y})] \right\}$
- 5: **end for**
- 6: Compute ELBO
- 7: **end while**

Bethe Approximation

□ Bethe approximation/Kikuchi Approximation

$$\max \text{ELBO} = \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x}, y) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})$$

Bethe Approximation

$$q(\mathbf{x}) = \frac{\prod_a q(\mathbf{x}_a)}{\prod_i [q(x_i)]^{d_i-1}}$$

pair-wise correlations

$$\text{ELBO} = \sum_a \sum_{\mathbf{x}_a} q(\mathbf{x}_a) \log f_a(q(x_i)) - \sum_i (d_i - 1) \sum_{x_i} q(x_i) \log q(x_i)$$

ELBO with Bethe Approximation

$$- \sum_a \sum_{\mathbf{x}_a} q(\mathbf{x}_a) \log q(\mathbf{x}_a) + \sum_i (d_i - 1) \sum_{x_i} q(x_i) \log q(x_i)$$

s.t.

Bethe Approximation

□ Bethe approximation/Kikuchi Approximation

$$\max \text{ELBO} = \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x}, y) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})$$

Bethe Approximation

$$q(\mathbf{x}) = \frac{\prod_a q(\mathbf{x}_a)}{\prod_i [q(x_i)]^{d_i-1}}$$

pair-wise correlations

ELBO with Bethe Approximation

$$-\sum_a \sum_{\mathbf{x}_a} q(\mathbf{x}_a) \log q(\mathbf{x}_a) + \sum_i (d_i - 1) \sum_{x_i} q(x_i) \log q(x_i)$$

s.t.

$$m_{a \rightarrow i}(x_i) = \sum_{\mathbf{x}_a \setminus x_i} f_a(\mathbf{x}_a) \prod_j m_{j \rightarrow a}(x_i)$$

Belief Propagation(BP)

Lagrange Multiplier



$m_{i \rightarrow a}(x_i) = \prod_{b \neq a} m_{b \rightarrow i}(x_i)$
Message passing on the factor graph

Belief Propagation

□ A Toy Example

There are 4 random discrete variables, each taking 10 possible values randomly.

The joint distribution is

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

Question: How to compute the marginal distribution $p(x_2)$

Belief Propagation

□ A Toy Example

There are 4 random discrete variables, each taking 10 possible values randomly.
The joint distribution is

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

Question: How to compute the marginal distribution $p(x_2)$

Direct Answer: marginalize out all other variables of the joint distribution

$$p(x_2) = \sum_{x_1, x_3, x_4} p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

Number of values of x_2	10
For each combination x_1, x_3, x_4	3 multiplications
Number of combinations x_1, x_3, x_4	$10 \cdot 10 \cdot 10 = 10^3$

Total Number of Multiplications :	$10 \cdot 3 \cdot 10^3 = 3 \cdot 10^4$
Total Number of Additions :	$10 \cdot (10^3 - 1) \approx 10^4$

$\mathcal{O}(10^4)$

Belief Propagation

□ A Toy Example

There are 4 random discrete variables, each taking 10 possible values randomly.
The joint distribution is

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

Question: How to compute the marginal distribution $p(x_2)$

Direct Answer: marginalize out all other variables of the joint distribution

$$p(x_2) = \sum_{x_1, x_3, x_4} p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

Number of values of x_2	10
For each combination x_1, x_3, x_4	3 multiplications
Number of combinations x_1, x_3, x_4	$10 * 10 * 10 = 10^3$

Total Number of Multiplications :	$10 * 3 * 10^3 = 3 * 10^4$
Total Number of Additions :	$10 * (10^3 - 1) \approx 10^4$

$\mathcal{O}(10^4)$

The structure of the joint distribution is totally ignored!

Belief Propagation

□ A Toy Example

The distributive law

$$ab + ac = a(b + c)$$

Belief Propagation

□ A Toy Example

The distributive law

$$ab + ac = a(b + c)$$

Alternative Answer:

$$\begin{aligned} p(x_2) &= \sum_{x_1, x_3, x_4} p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2) \\ &= \underbrace{\sum_{x_1} p(x_1) p(x_2|x_1)}_{m_A(x_2)} \underbrace{\sum_{x_3} p(x_3|x_2)}_{m_B(x_2)} \underbrace{\sum_{x_4} p(x_4|x_2)}_{m_C(x_2)} \end{aligned}$$

The distributive law

Belief Propagation

□ A Toy Example

The distributive law

$$ab + ac = a(b + c)$$

Alternative Answer:

$$p(x_2) = \sum_{x_1, x_3, x_4} p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

The original problem is divided into small sub-problems

$$= \underbrace{\sum_{x_1} p(x_1) p(x_2|x_1)}_{m_A(x_2)} \underbrace{\sum_{x_3} p(x_3|x_2)}_{m_B(x_2)} \underbrace{\sum_{x_4} p(x_4|x_2)}_{m_C(x_2)}$$

The distributive law

Total Number of Multiplications : $10 \cdot (10+2) = 120$

Total Number of Additions : $10 \cdot (9+9+9) = 270$

$$\mathcal{O}(10^2)$$

Belief Propagation

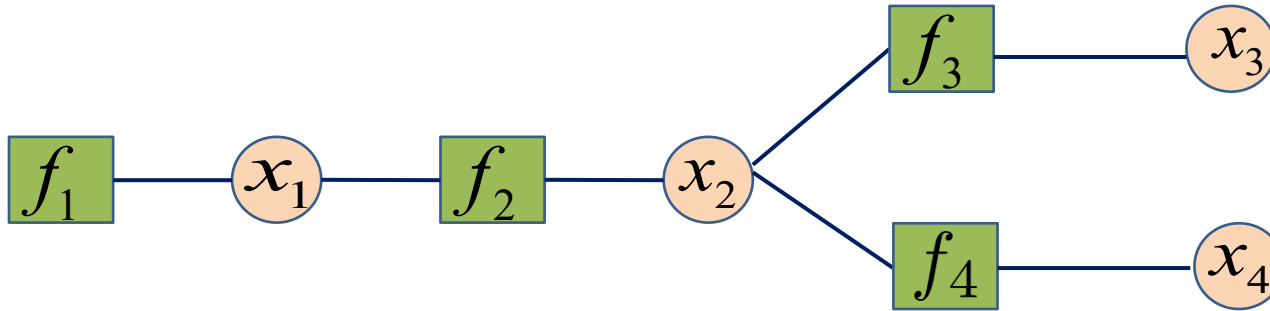
□ A Toy Example

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

f_1 f_2 f_3 f_4

Factor graph

- circle nodes represent random variables
- square nodes represent factorizing functions
- function node f connects variable node x if and only if x is one of argument of f



Belief Propagation

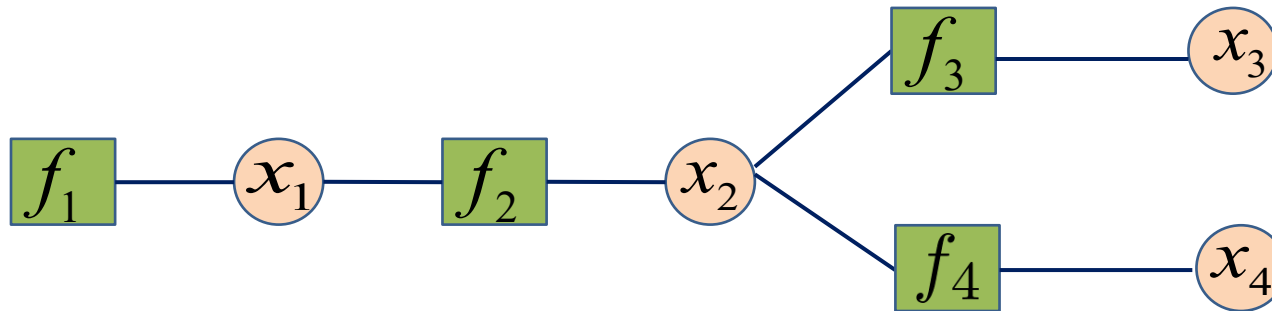
□ A Toy Example

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

f_1 f_2 f_3 f_4

Factor graph

- circle nodes represent random variables
- square nodes represent factorizing functions
- function node f connects variable node x if and only if x is one of argument of f



$$p(x_2) = \sum_{x_1} p(x_1) p(x_2|x_1) \sum_{x_3} p(x_3|x_2) \sum_{x_4} p(x_4|x_2)$$

Inference process



Message Passing on graph

Belief Propagation

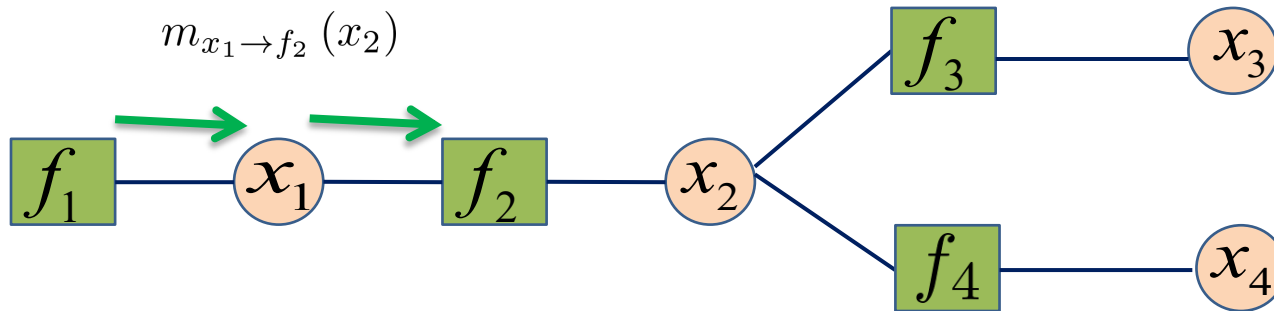
□ A Toy Example

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

f_1 f_2 f_3 f_4

Factor graph

- circle nodes represent random variables
- square nodes represent factorizing functions
- function node f connects variable node x if and only if x is one of argument of f



$$p(x_2) = \sum_{x_1} \overbrace{p(x_1) p(x_2|x_1)}^{m_{x_1 \rightarrow f_2}(x_2)} \sum_{x_3} p(x_3|x_2) \sum_{x_4} p(x_4|x_2)$$

Belief Propagation

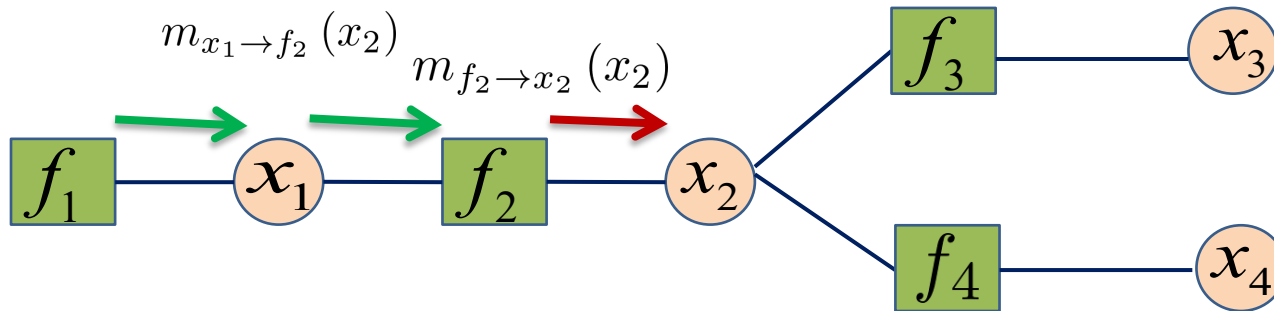
□ A Toy Example

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

f_1 f_2 f_3 f_4

Factor graph

- circle nodes represent random variables
- square nodes represent factorizing functions
- function node f connects variable node x if and only if x is one of argument of f



$$p(x_2) = \sum_{x_1} \overbrace{p(x_1) p(x_2|x_1)}^{m_{x_1 \rightarrow f_2}(x_2)} \sum_{x_3} p(x_3|x_2) \sum_{x_4} p(x_4|x_2)$$

$\underbrace{\hspace{10em}}_{m_{f_2 \rightarrow x_2}(x_2)}$

Belief Propagation

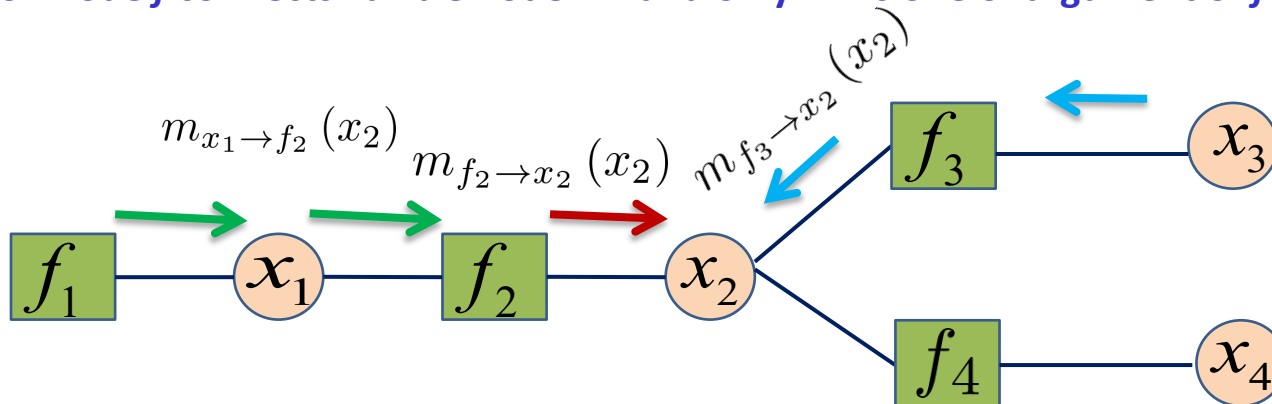
□ A Toy Example

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

f_1 f_2 f_3 f_4

Factor graph

- circle nodes represent random variables
- square nodes represent factorizing functions
- function node f connects variable node x if and only if x is one of argument of f



$$p(x_2) = \sum_{x_1} \overbrace{p(x_1) p(x_2|x_1)}^{m_{x_1 \rightarrow f_2}(x_2)} \sum_{x_3} \underbrace{p(x_3|x_2)}_{m_{x_3 \rightarrow f_3}(x_2)} \sum_{x_4} p(x_4|x_2)$$

$m_{f_2 \rightarrow x_2}(x_2)$ $m_{f_3 \rightarrow x_2}(x_2)$

Belief Propagation

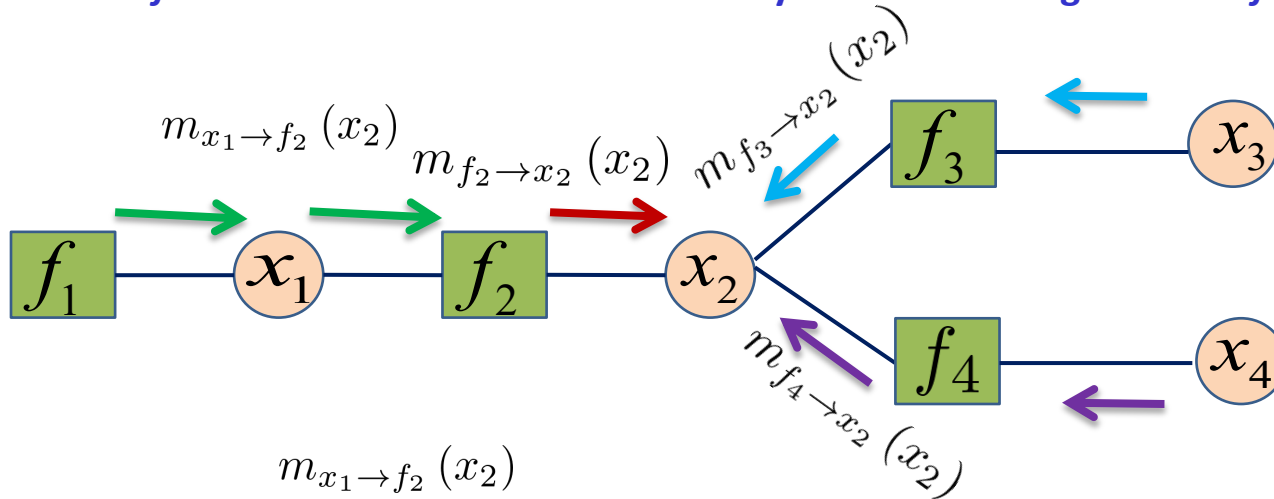
□ A Toy Example

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

f_1 f_2 f_3 f_4

Factor graph

- circle nodes represent random variables
- square nodes represent factorizing functions
- function node f connects variable node x if and only if x is one of argument of f



$$p(x_2) = \sum_{x_1} \overbrace{p(x_1) p(x_2|x_1)}^{m_{x_1 \rightarrow f_2}(x_2)} \sum_{x_3} \underbrace{p(x_3|x_2)}_{m_{f_3 \rightarrow x_2}(x_2)} \sum_{x_4} \underbrace{p(x_4|x_2)}_{m_{f_4 \rightarrow x_2}(x_2)}$$

Belief Propagation

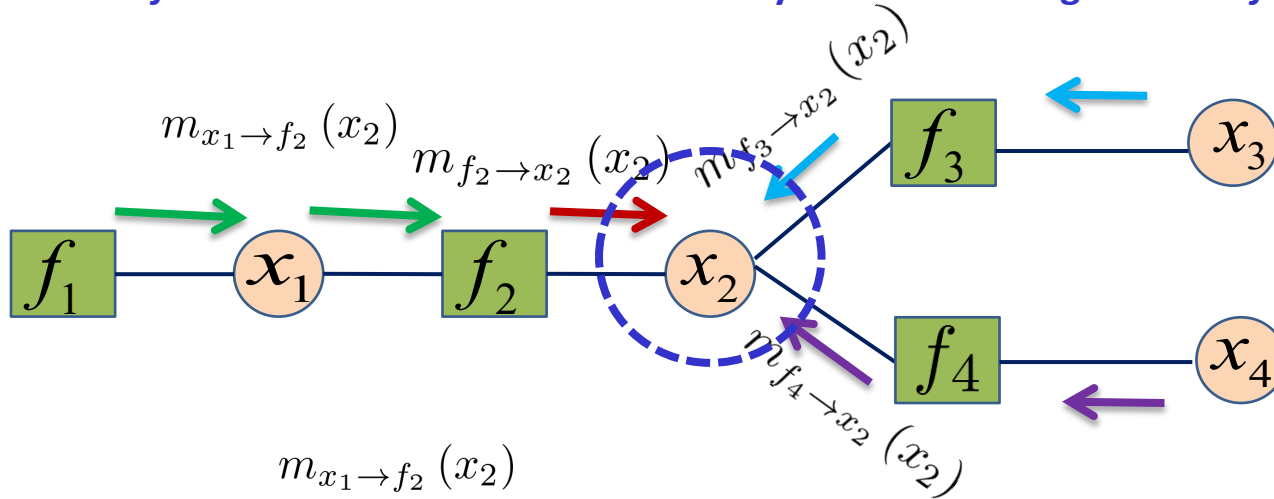
□ A Toy Example

$$p(x_1, x_2, x_3, x_4) = p(x_1) p(x_2|x_1) p(x_3|x_2) p(x_4|x_2)$$

f_1 f_2 f_3 f_4

Factor graph

- circle nodes represent random variables
- square nodes represent factorizing functions
- function node f connects variable node x if and only if x is one of argument of f



$$p(x_2) = \sum_{x_1} \overbrace{p(x_1) p(x_2|x_1)}^{m_{x_1 \rightarrow f_2}(x_2)} \underbrace{\sum_{x_3} p(x_3|x_2)}_{m_{f_2 \rightarrow x_2}(x_2)} \underbrace{\sum_{x_4} p(x_4|x_2)}_{m_{f_4 \rightarrow x_2}(x_2)}$$

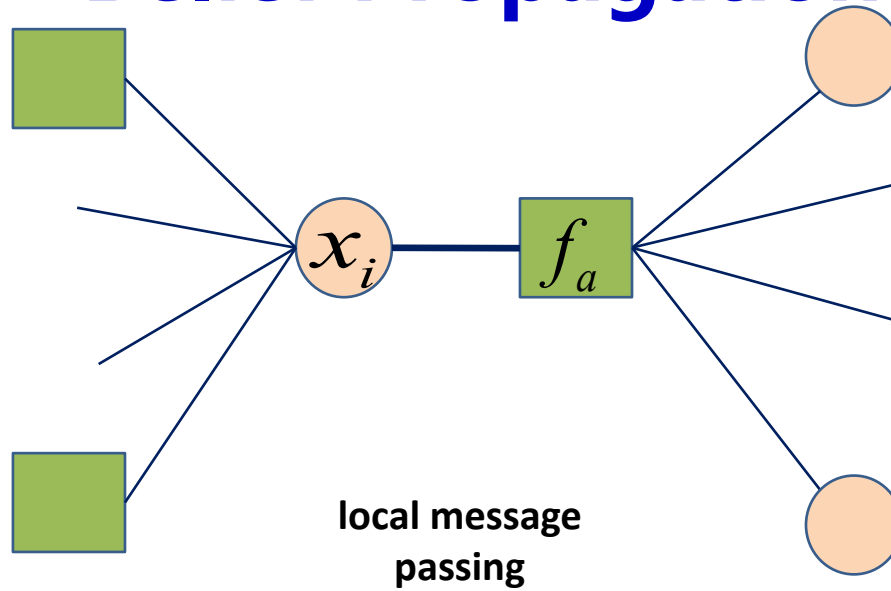
$m_{f_3 \rightarrow x_2}(x_2)$ $m_{f_4 \rightarrow x_2}(x_2)$

product of all incoming messages

Belief Propagation

□ Factor Graph

BP on
general graph



Loopy Belief Propagation (LBP)

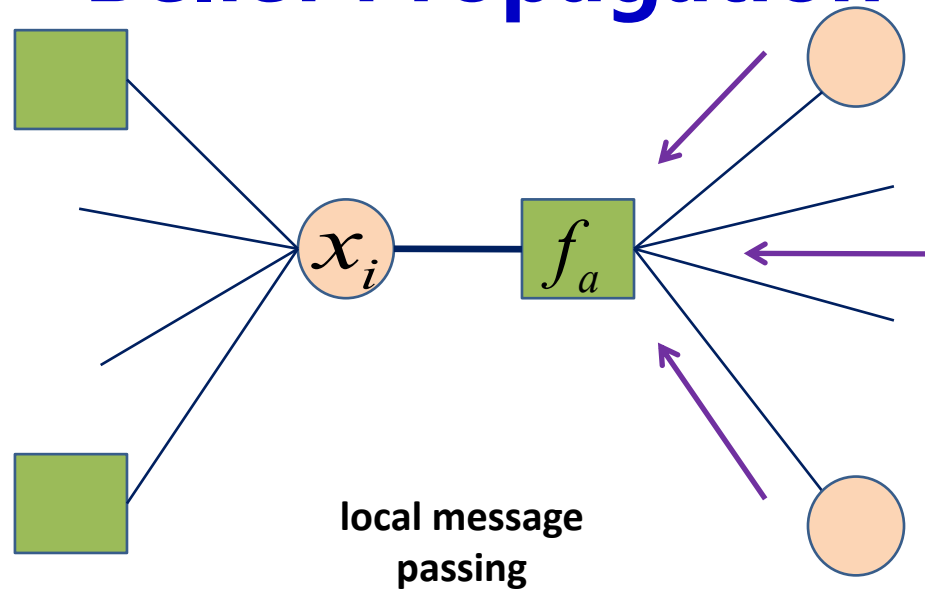
Factor to variable

Variable to factor

Belief Propagation

□ Factor Graph

BP on
general graph



Loopy Blief Propagation (LBP)

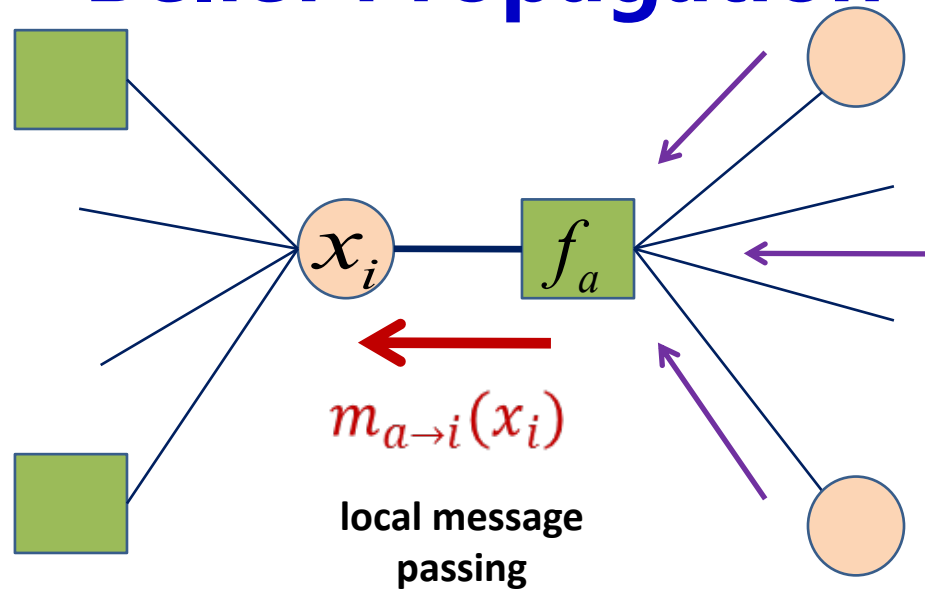
Factor to variable

Variable to factor

Belief Propagation

□ Factor Graph

BP on
general graph



Loopy Belief Propagation (LBP)

Factor to variable

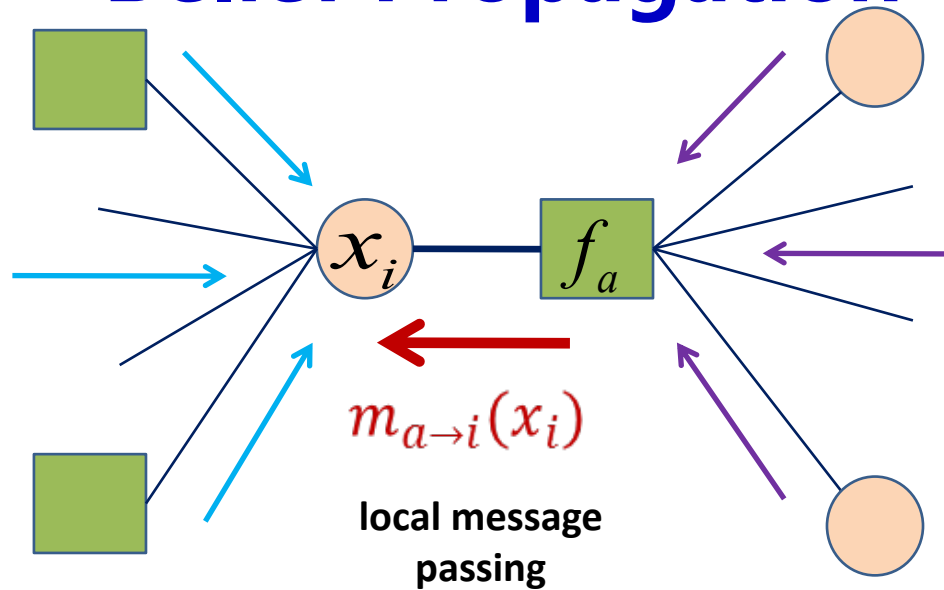
$$m_{a \rightarrow i}(x_i) = \sum_{x_j, j \neq i} f_a(\mathbf{x}_a) \prod_{j \neq i} m_{j \rightarrow a}(x_j)$$

Variable to factor

Belief Propagation

□ Factor Graph

BP on
general graph



Loopy Belief Propagation (LBP)

Factor to variable

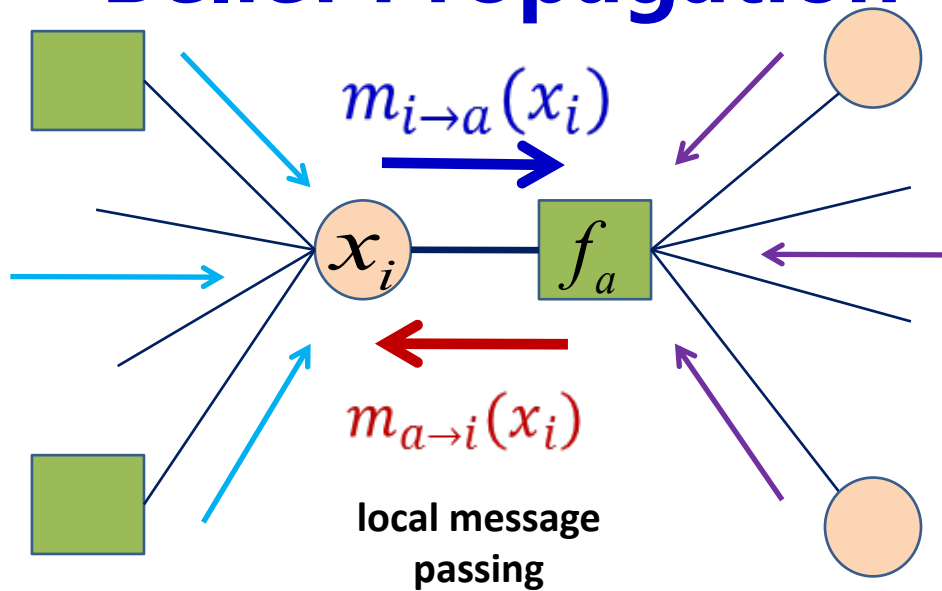
$$m_{a \rightarrow i}(x_i) = \sum_{x_j, j \neq i} f_a(\mathbf{x}_a) \prod_{j \neq i} m_{j \rightarrow a}(x_j)$$

Variable to factor

Belief Propagation

□ Factor Graph

BP on
general graph



Loopy Blief Propagation (LBP)

Factor to variable

$$m_{a \rightarrow i}(x_i) = \sum_{x_j, j \neq i} f_a(\mathbf{x}_a) \prod_{j \neq i} m_{j \rightarrow a}(x_j)$$

Variable to factor

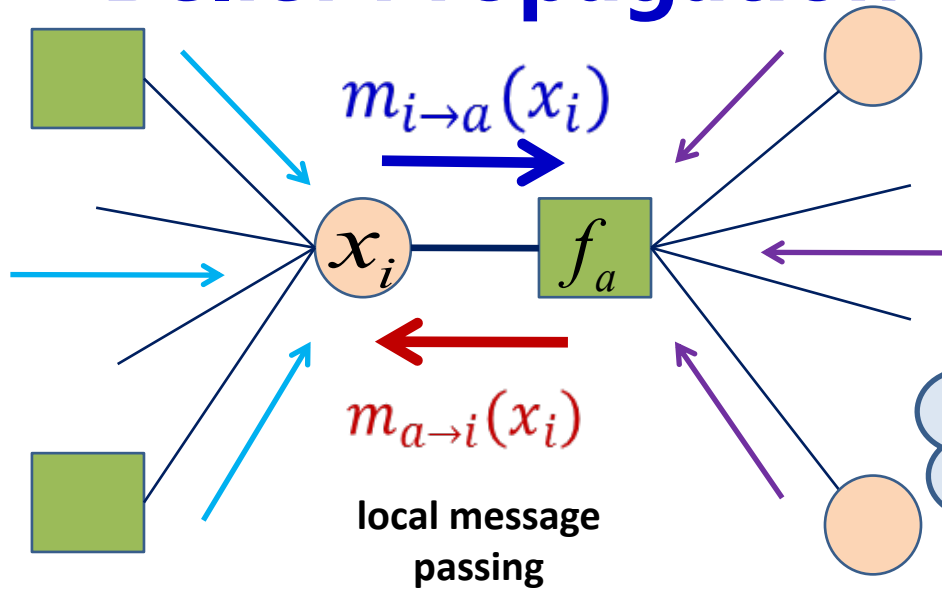
$$m_{i \rightarrow a}(x_i) = \prod_{b \neq a} m_{b \rightarrow i}(x_i)$$

Excluding incoming message itself

Belief Propagation

□ Factor Graph

BP on
general graph



BP is exact for
graph without loops!

Loopy Belief Propagation (LBP)

Factor to variable

$$m_{a \rightarrow i}(x_i) = \sum_{x_j, j \neq i} f_a(\mathbf{x}_a) \prod_{j \neq i} m_{j \rightarrow a}(x_j)$$

Iterations

(graph with loops)

Variable to factor

$$m_{i \rightarrow a}(x_i) = \prod_{b \neq a} m_{b \rightarrow i}(x_i)$$

Excluding incoming
message itself

Belief Propagation



figure copied from http://computerrobotvision.org/2009/tutorial_day/crv09_belief_propagation_v2.pdf

Message passing is a beautiful algorithmic framework to tackle difficult problems using **divide and conquer by **local computation** and **information sharing****

Parametric Approximation

□ Parameterization

$$\max \text{ELBO} = \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x}, y) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})$$

Parameterization $q(\mathbf{x}) \equiv q(\mathbf{x}; \boldsymbol{\phi})$

- **Exponential family** $q(\mathbf{x}; \boldsymbol{\phi}) = \exp\{\langle \boldsymbol{\phi}, \boldsymbol{\eta}(\mathbf{x}) \rangle - A(\boldsymbol{\phi})\}$ E.g., Gaussian, Bernoulli, exponential...
- **Deep neural network**, such as that used in variational auto-encoder (VAE)

Parametric Approximation

□ Parameterization

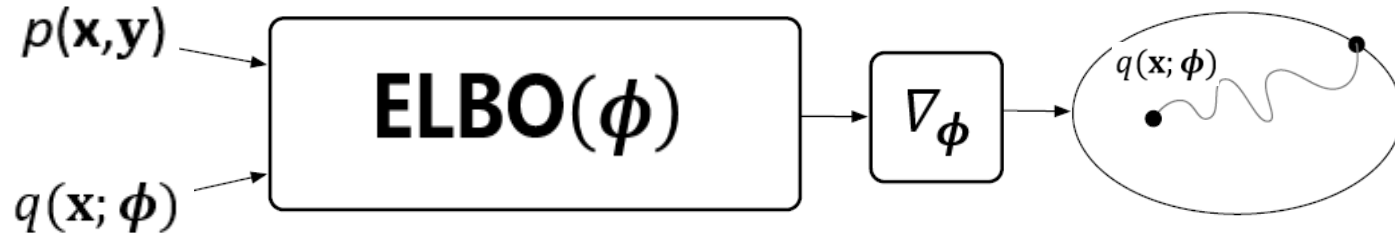
$$\max \text{ELBO} = \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x}, \mathbf{y}) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})$$

Parameterization $q(\mathbf{x}) \equiv q(\mathbf{x}; \boldsymbol{\phi})$

- **Exponential family** $q(\mathbf{x}; \boldsymbol{\phi}) = \exp\{\langle \boldsymbol{\phi}, \boldsymbol{\eta}(\mathbf{x}) \rangle - A(\boldsymbol{\phi})\}$ E.g., Gaussian, Bernoulli, exponential...
- **Deep neural network**, such as that used in variational auto-encoder (VAE)

New Optimization Objective

$$\max_{\boldsymbol{\phi}} \text{ELBO}(\boldsymbol{\phi}) = \sum_{\mathbf{x}} q(\mathbf{x}; \boldsymbol{\phi}) \log p(\mathbf{x}, \mathbf{y}) - \sum_{\mathbf{x}} q(\mathbf{x}; \boldsymbol{\phi}) \log q(\mathbf{x}; \boldsymbol{\phi})$$



Stochastic variational inference framework

The variational parameters are optimized using SGD

[Hoffman et al 2013]

Parametric Approximation

□ Parameterization

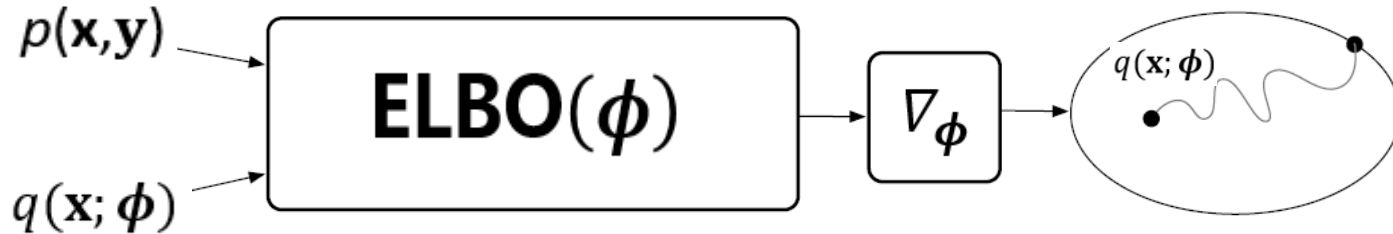
$$\max \text{ELBO} = \sum_{\mathbf{x}} q(\mathbf{x}) \log p(\mathbf{x}, \mathbf{y}) - \sum_{\mathbf{x}} q(\mathbf{x}) \log q(\mathbf{x})$$

Parameterization $q(\mathbf{x}) \equiv q(\mathbf{x}; \boldsymbol{\phi})$

- **Exponential family** $q(\mathbf{x}; \boldsymbol{\phi}) = \exp\{\langle \boldsymbol{\phi}, \boldsymbol{\eta}(\mathbf{x}) \rangle - A(\boldsymbol{\phi})\}$ E.g., Gaussian, Bernoulli, exponential...
- **Deep neural network**, such as that used in variational auto-encoder (VAE)

New Optimization Objective

$$\max_{\boldsymbol{\phi}} \text{ELBO}(\boldsymbol{\phi}) = \sum_{\mathbf{x}} q(\mathbf{x}; \boldsymbol{\phi}) \log p(\mathbf{x}, \mathbf{y}) - \sum_{\mathbf{x}} q(\mathbf{x}; \boldsymbol{\phi}) \log q(\mathbf{x}; \boldsymbol{\phi})$$



Stochastic variational inference framework

The variational parameters are optimized using SGD

The original integration problem builds down to derivative problem

Outline

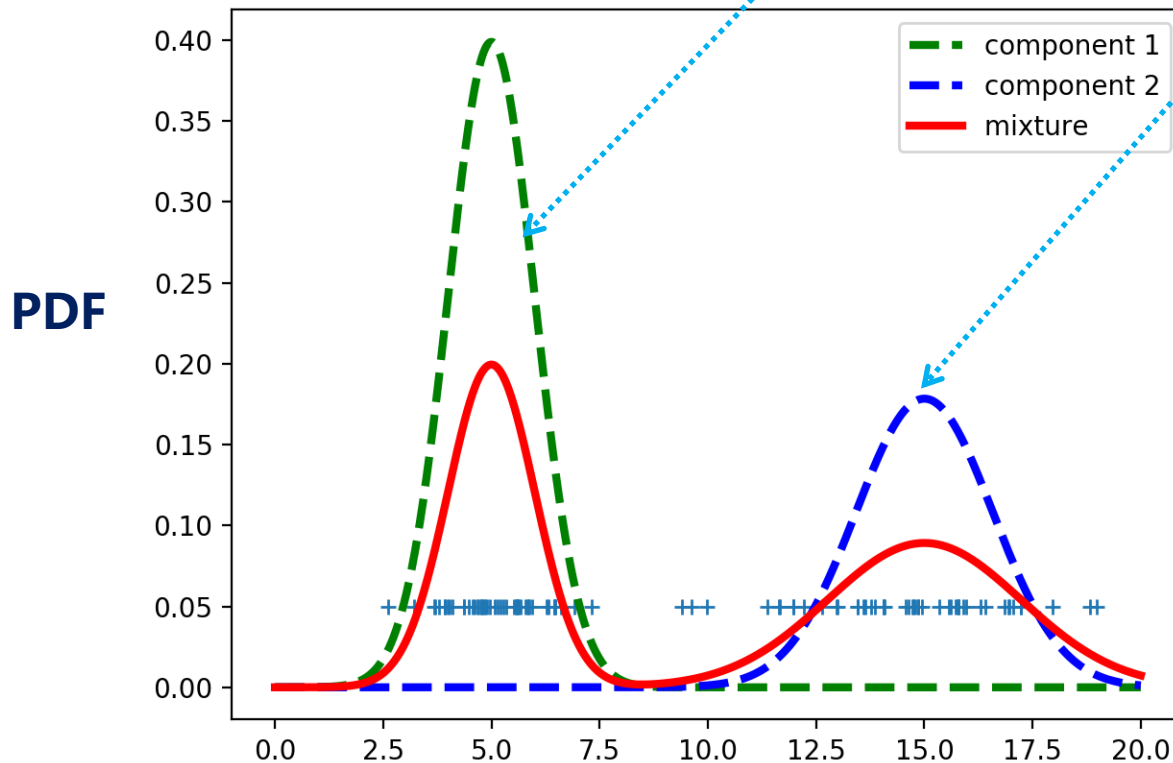
- Background
- Variational Inference
- **Expectation Propagation**
- A Unified EP Perspective on AMP and its extensions
- Conclusion

A Toy Problem

□ Problem Statement

we obtain a sequence of data points $y_i, i = 1 \dots N$

$$p(y_i|x) = \underbrace{0.5\mathcal{N}(y_i; x, 1)}_{\text{component one}} + \underbrace{0.5\mathcal{N}(y_i; x + 10, 5)}_{\text{component two}}$$



What is the value of x ?

This example is modified from example in [Minka01b]

A Toy Problem

□ Probabilistic Modeling

- **prior distribution** $p(x) = \mathcal{N}(x; 0, 100)$ **Guassian prior**
- **likelihood distribution** $p(y_i|x) = 0.5\mathcal{N}(y_i; x, 1) + 0.5\mathcal{N}(y_i; x + 10, 5)$

After obtaining N observations, the joint distribution could be written as

$$p(x, \mathbf{y}) = p(x) \prod_{i=1}^N p(y_i|x)$$

- **posterior distribution**

$$p(x|\mathbf{y}) = \frac{p(x) \prod_{i=1}^N p(y_i|x)}{p(\mathbf{y})}$$

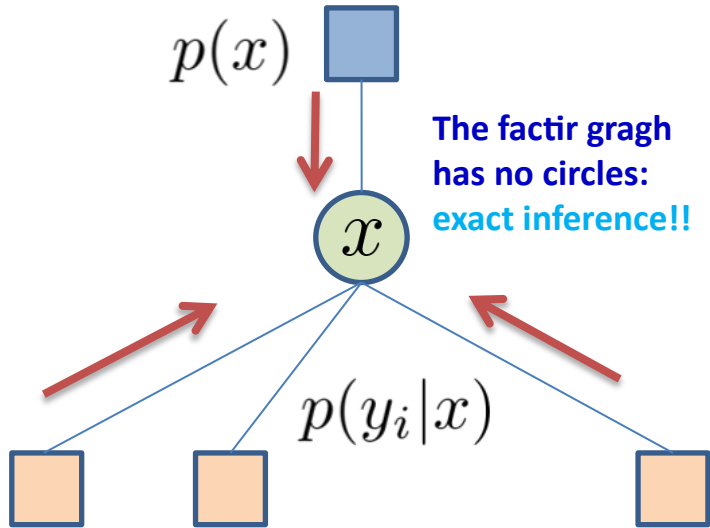
We could perform Bayesian inference to compute the posterior distribution

All the codes for this toy example are available:
<https://github.com/mengxiangming/ep-demo>

A Toy Problem

Factor Graph and Belief Propagation

$$p(x, \mathbf{y}) = p(x) \prod_{i=1}^N p(y_i|x)$$



Belief Propagation

factor to variable: $m_{i \rightarrow x}(x) = p(y_i|x)$

variable update: $q(x) = p(x) \prod_{i=1}^N m_{i \rightarrow x}(x)$

Already Done?

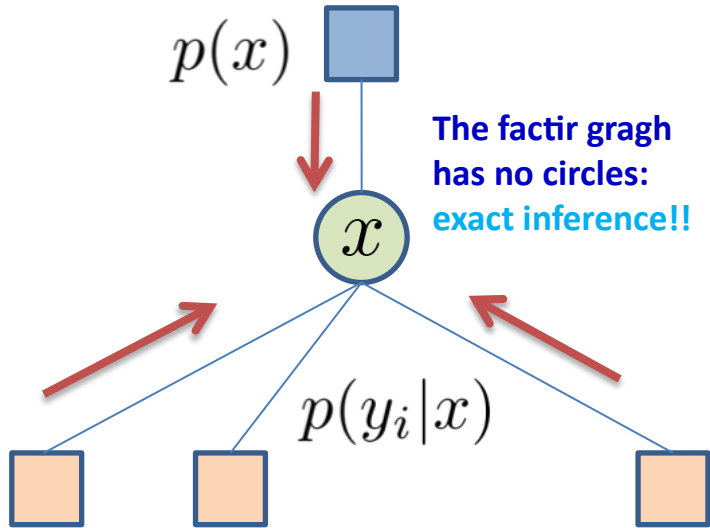
$$p(y_i|x) = 0.5\mathcal{N}(y_i; x, 1) + 0.5\mathcal{N}(y_i; x + 10, 5)$$

$$p(x) = \mathcal{N}(x; 0, 100)$$

A Toy Problem

Factor Graph and Belief Propagation

$$p(x, \mathbf{y}) = p(x) \prod_{i=1}^N p(y_i|x)$$



Belief Propagation

factor to variable: $m_{i \rightarrow x}(x) = p(y_i|x)$

variable update: $q(x) = p(x) \prod_{i=1}^N m_{i \rightarrow x}(x)$

Already Done?

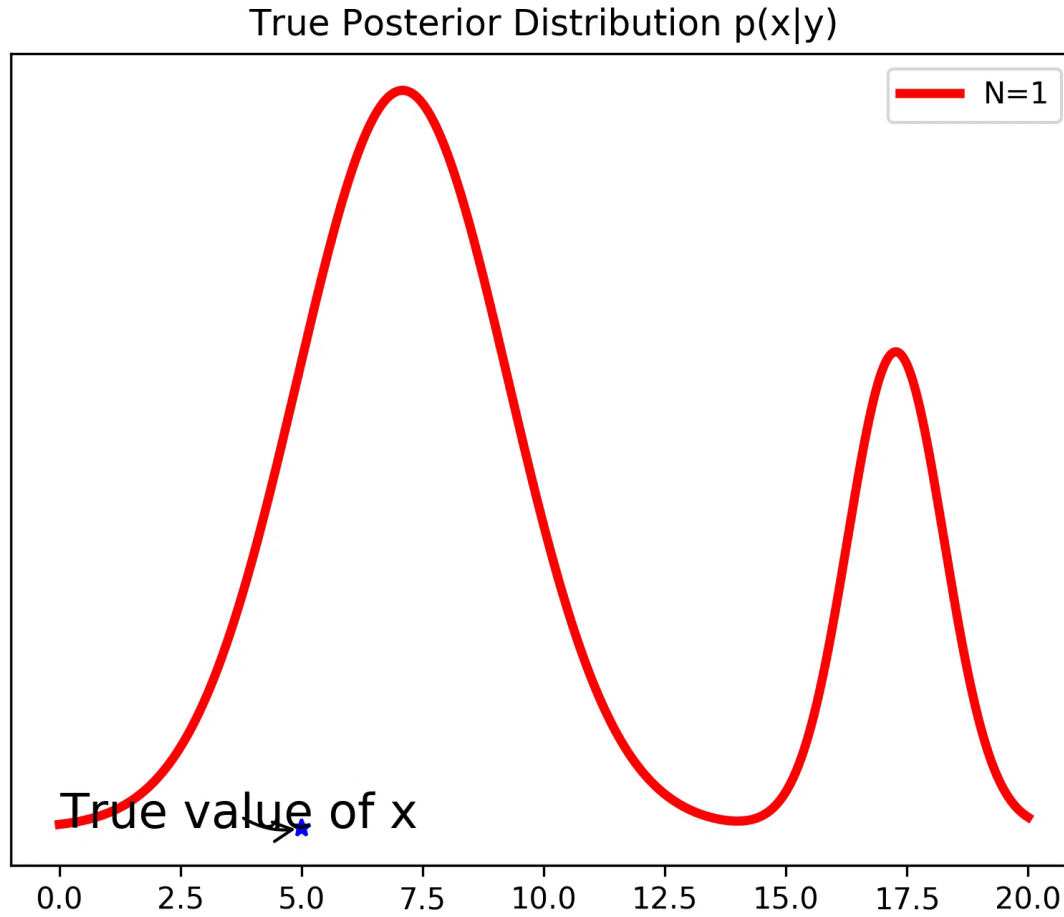
$$p(y_i|x) = 0.5\mathcal{N}(y_i; x, 1) + 0.5\mathcal{N}(y_i; x + 10, 5)$$

$$p(x) = \mathcal{N}(x; 0, 100)$$

- The posterior distribution is a mixture of Ω^N Gaussians.
- The computational complexity is exponential with N

A Toy Problem

□ The True Posterior



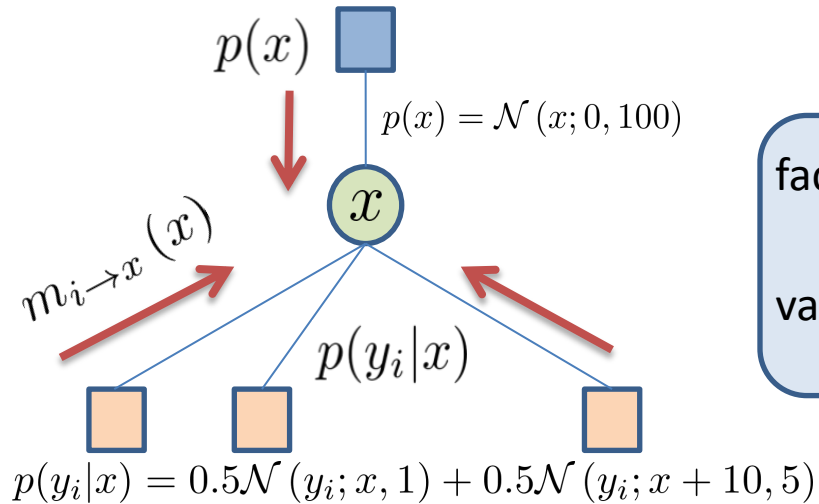
A Toy Problem

□ The True Posterior

**Approximating the posterior as
one Gaussian distribution**

A Toy Problem

□ A Naive Approximation



Approximating each BP message itself
as Gaussian distribution independently

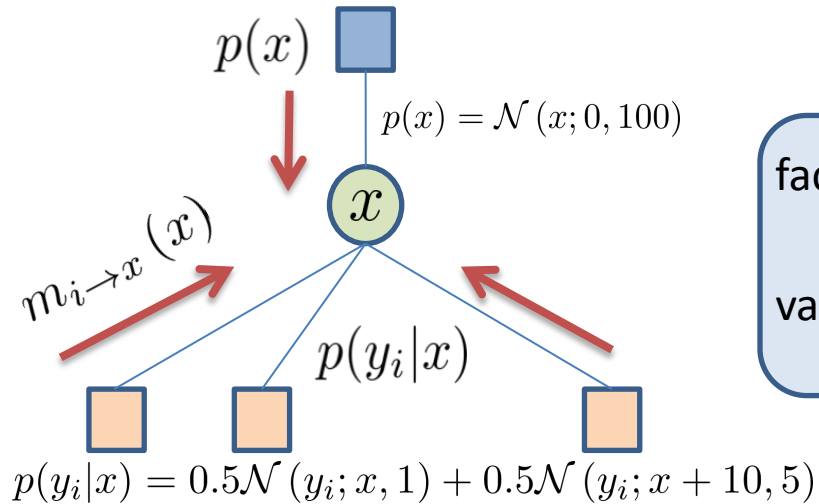
factor to variable: $m_{i \rightarrow x}(x) \approx \mathcal{N}(x; m_i, v_i)$

variable update: $q(x) \approx p(x) \prod_{i=1}^N \mathcal{N}(x; m_i, v_i)$

Naive Gaussian Message Approximation

A Toy Problem

□ A Naive Approximation



Approximating each BP message itself as Gaussian distribution independently

factor to variable: $m_{i \rightarrow x}(x) \approx \mathcal{N}(x; m_i, v_i)$

variable update: $q(x) \approx p(x) \prod_{i=1}^N \mathcal{N}(x; m_i, v_i)$

Naive Gaussian Message Approximation

True Posterior

$$p(x|\mathbf{y}) \propto p(x) \prod_{i=1}^N p(y_i|x)$$

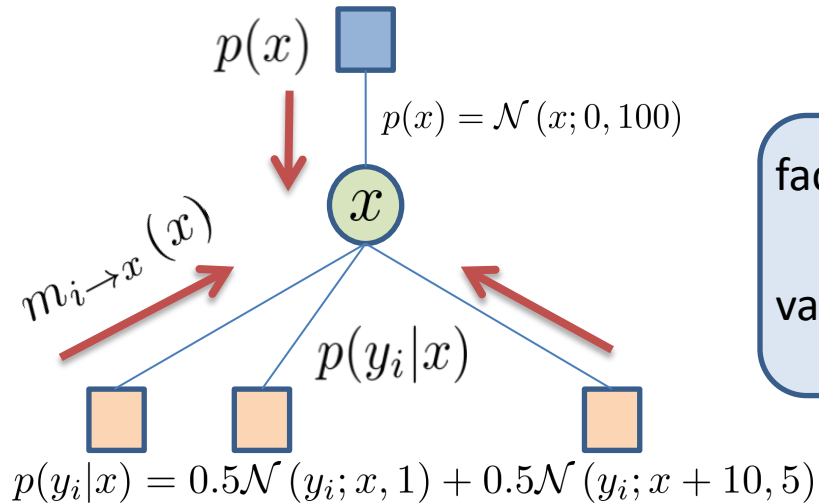
Approximate

$$p(x|\mathbf{y}) \propto p(x) \prod_{i=1}^N \mathcal{N}(x; m_i, v_i)$$

Each non-Gaussian likelihood is approximated as a Gaussian factor

A Toy Problem

□ A Naive Approximation



Approximating each BP message itself as Gaussian distribution independently

factor to variable: $m_{i \rightarrow x}(x) \approx \mathcal{N}(x; m_i, v_i)$

variable update: $q(x) \approx p(x) \prod_{i=1}^N \mathcal{N}(x; m_i, v_i)$

Naive Gaussian Message Approximation

True Posterior

$$p(x|\mathbf{y}) \propto p(x) \prod_{i=1}^N p(y_i|x)$$

Approximate

$$p(x|\mathbf{y}) \propto p(x) \prod_{i=1}^N \mathcal{N}(x; m_i, v_i)$$

Each non-Gaussian likelihood is approximated as a Gaussian factor

The posterior will be also Gaussian due to the product rule of Gaussian

$$\mathcal{N}(x; m, v) \propto \mathcal{N}(x; m_1, v_1) \mathcal{N}(x; m_2, v_2)$$

$$\frac{1}{v} = \frac{1}{v_1} + \frac{1}{v_2}$$

$$\frac{m}{v} = \frac{m_1}{v_1} + \frac{m_2}{v_2}$$

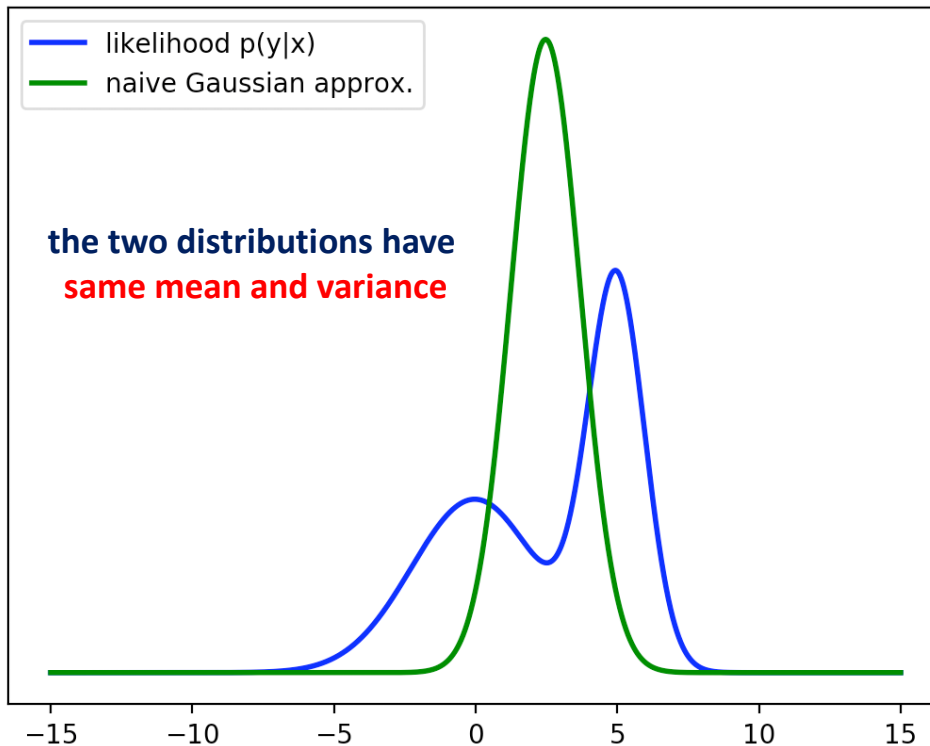
A Toy Problem

□ A Naive Approximation

For each non-Gaussian message

$$p(y_i|x) = 0.5\mathcal{N}(y_i; x, 1) + 0.5\mathcal{N}(y_i; x + 10, 5)$$

Gaussian Approximation $\tilde{t}_i(x) \triangleq \text{Proj}[p(y_i|x)]$ **Gaussian Projection Operator**



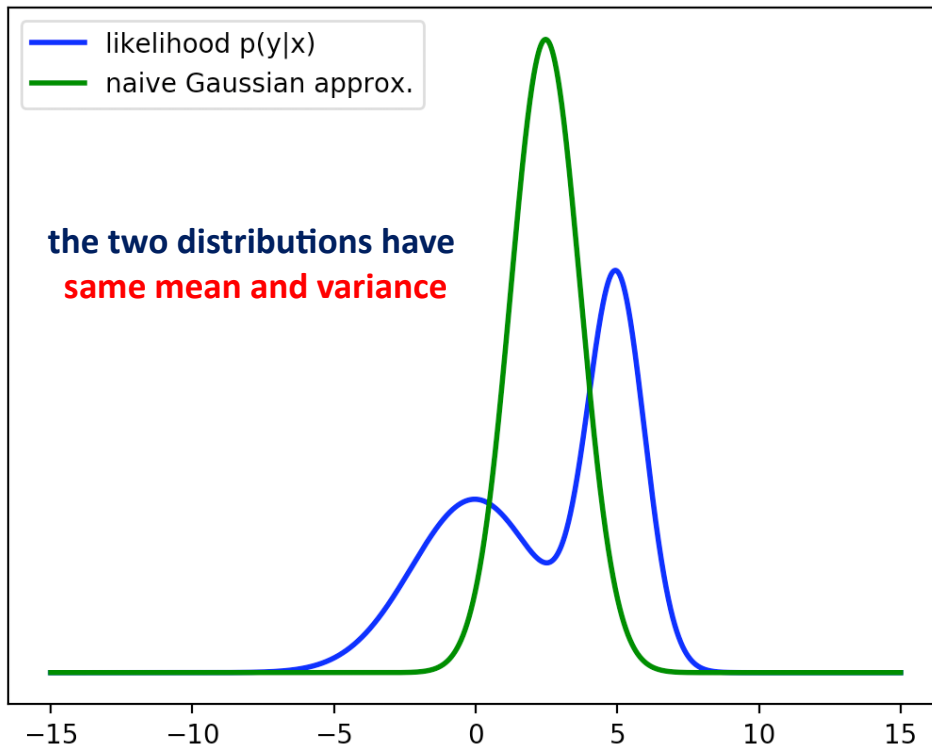
A Toy Problem

□ A Naive Approximation

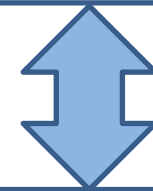
For each non-Gaussian message

$$p(y_i|x) = 0.5\mathcal{N}(y_i; x, 1) + 0.5\mathcal{N}(y_i; x + 10, 5)$$

Gaussian Approximation $\tilde{t}_i(x) \triangleq \text{Proj}[p(y_i|x)]$ Gaussian Projection Operator



$$q(x) = \arg \max_{q \in \text{Gaussian}} KL(p(x) || q(x))$$



equivalent

$$q(x) \triangleq \text{Proj}[p(x)] \quad m = \mathbf{E}_{p(x)}(x) \\ = \mathcal{N}(x; m, v) \quad v = \text{Var}_{p(x)}(x)$$

Moment Matching

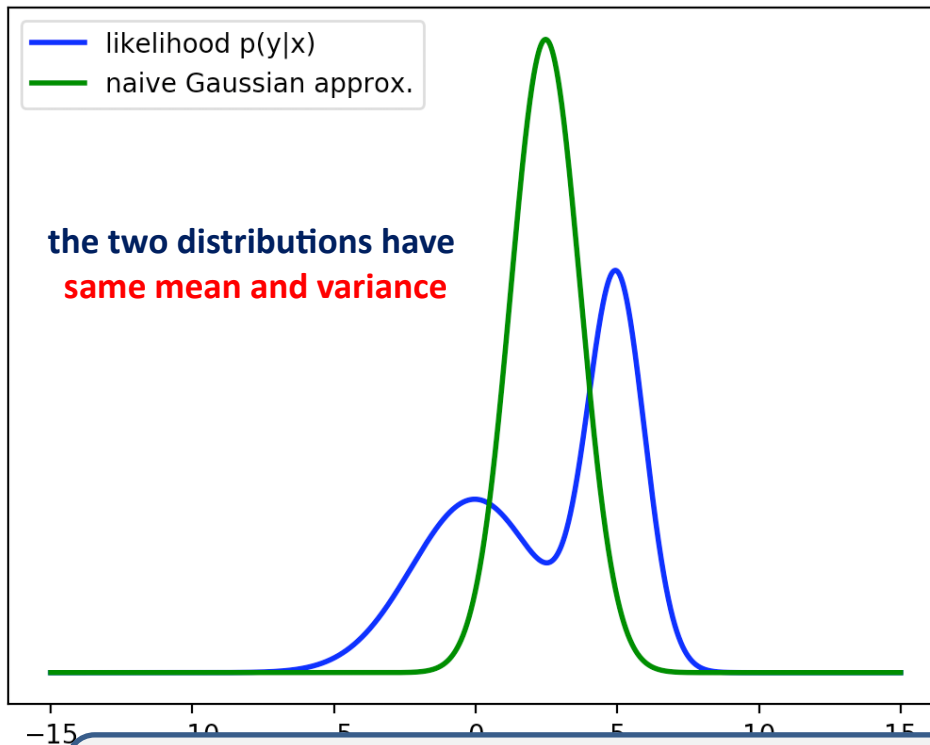
A Toy Problem

□ A Naive Approximation

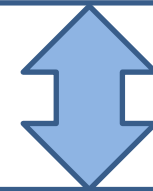
For each non-Gaussian message

$$p(y_i|x) = 0.5\mathcal{N}(y_i; x, 1) + 0.5\mathcal{N}(y_i; x + 10, 5)$$

Gaussian Approximation $\tilde{t}_i(x) \triangleq \text{Proj}[p(y_i|x)]$ Gaussian Projection Operator



$$q(x) = \arg \max_{q \in \text{Gaussian}} KL(p(x) || q(x))$$



equivalent

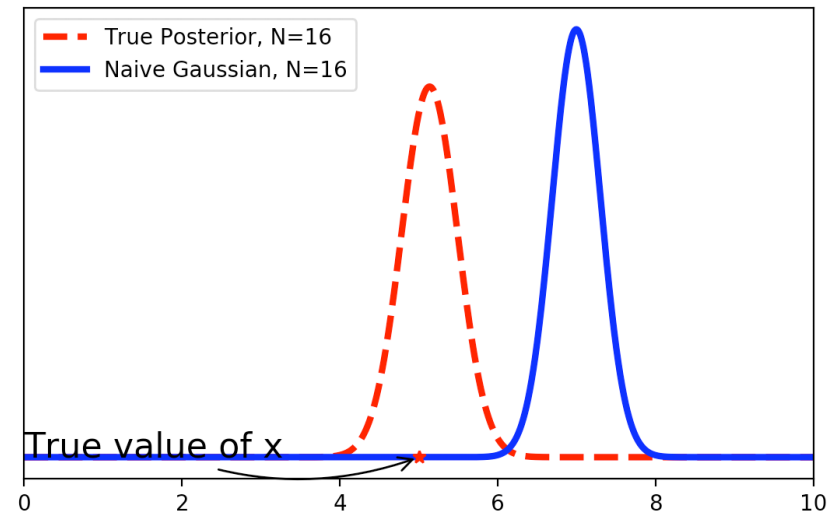
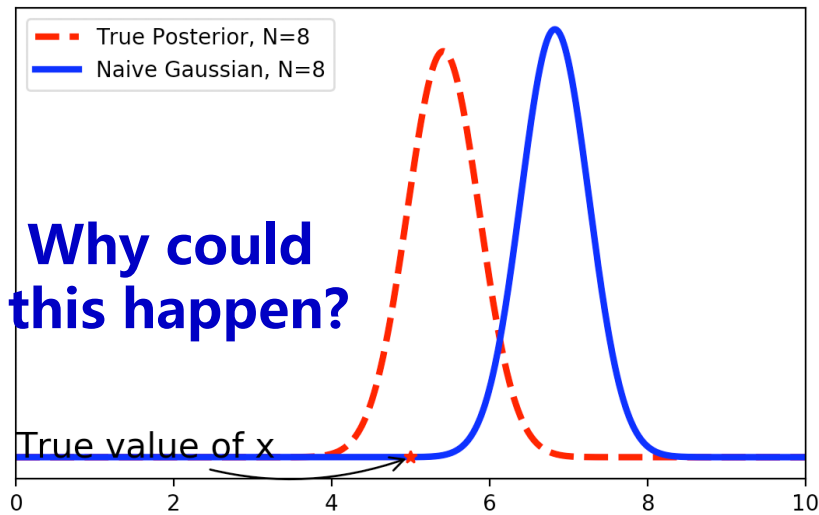
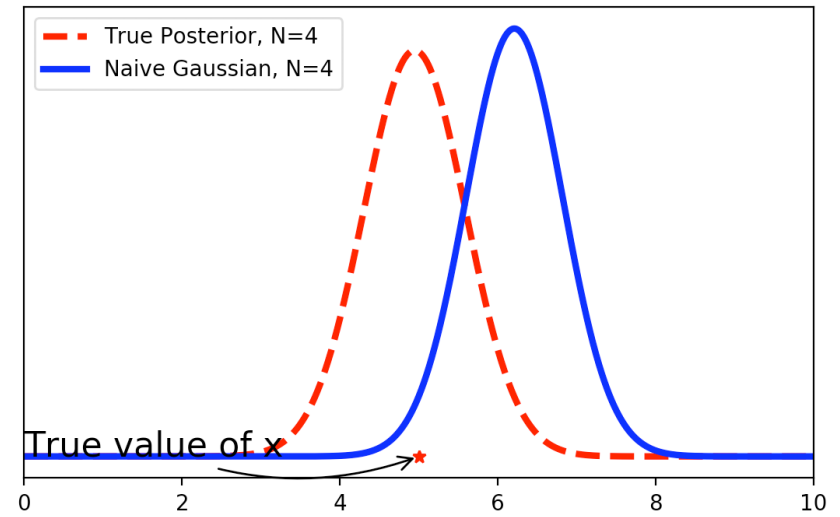
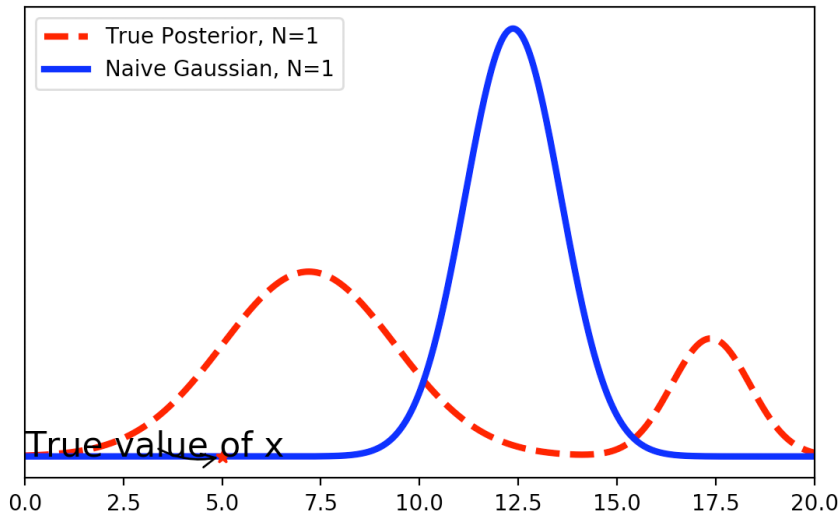
$$q(x) \triangleq \text{Proj}[p(x)] \quad m = \mathbf{E}_{p(x)}(x) \\ = \mathcal{N}(x; m, v) \quad v = \text{Var}_{p(x)}(x)$$

Moment Matching

Then, how will the posterior be like?

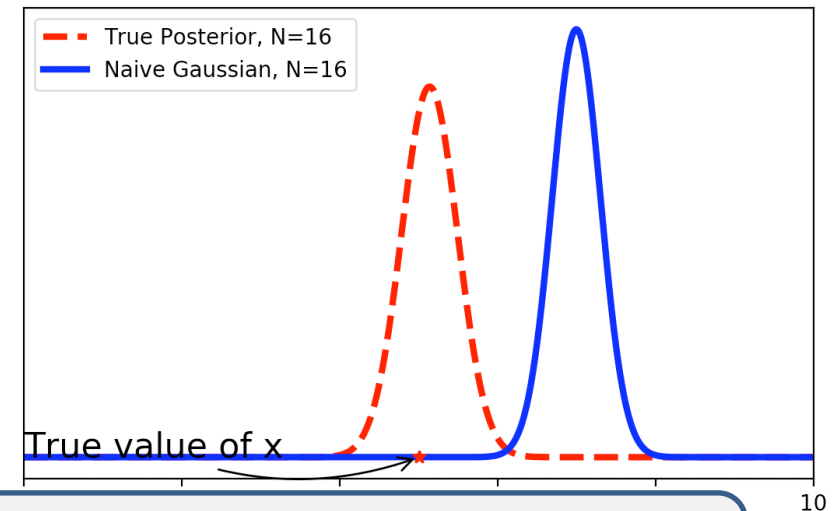
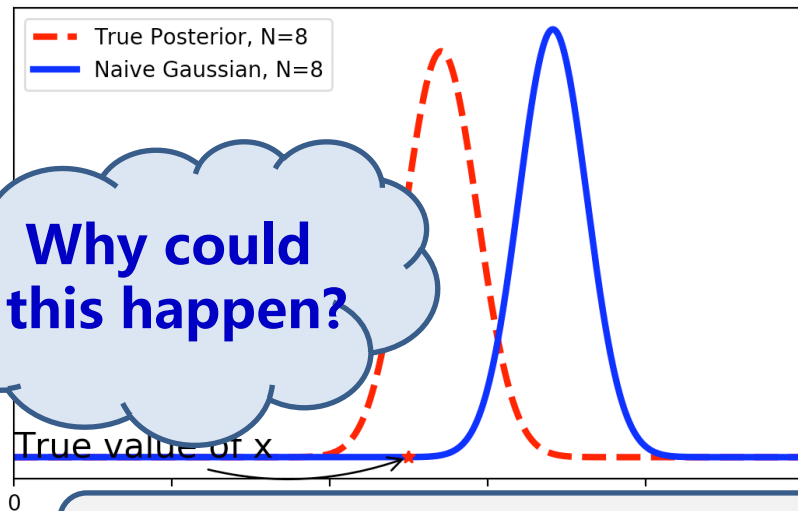
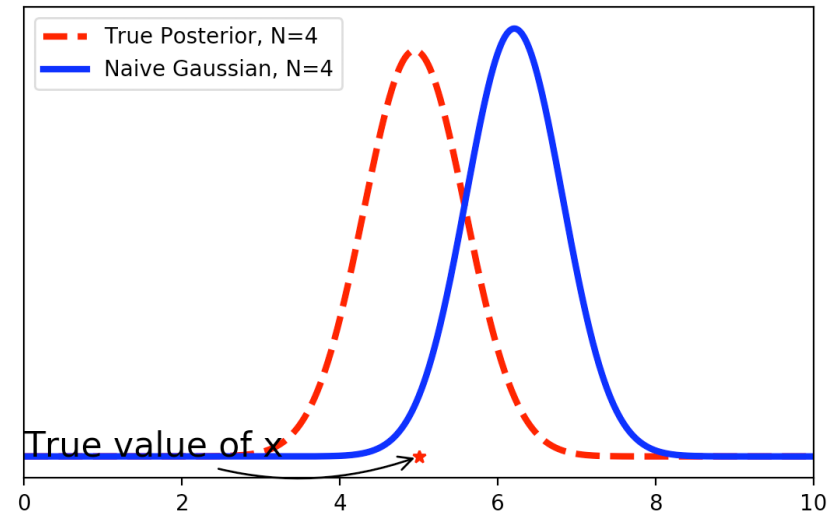
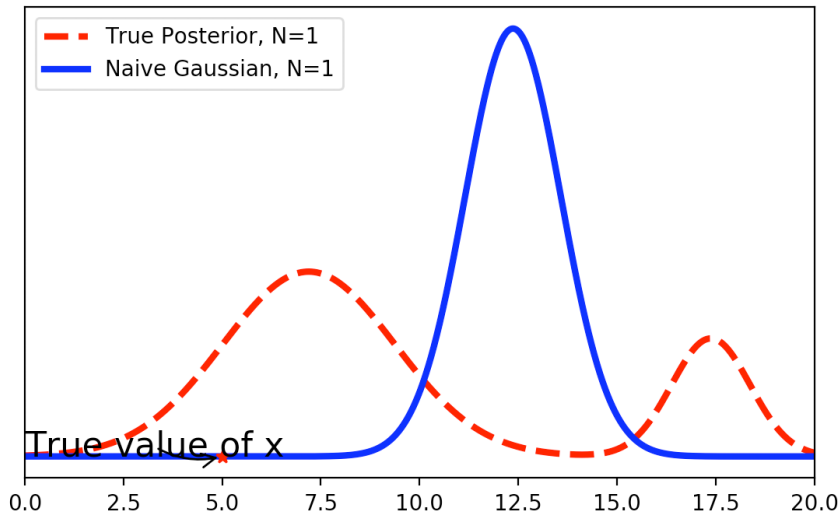
A Toy Problem

□ A Naive Approximation



A Toy Problem

□ A Naive Approximation



Why could this happen?

There is still a big discrepancy between the true posterior and naive Gaussian approximation, **even when the true posterior approaches Gaussian!**

A Toy Problem

- A Naive Approximation

Because it is **naive** selfish

A Toy Problem

□ A Naive Approximation

Because it is **naive** selfish

Each factor (message) only cares about itself
when making approximations
while forgetting the ultimate goal is to make a good
approximation to the global posterior

A Toy Problem

□ An Alternative Gaussian Approximation

Consider the simple case of $N = 1$ (only one observation)

True posterior $p(x|\mathbf{y}) \propto \underbrace{p(x)}_{\text{Gauss}} \underbrace{p(y_1|x)}_{\text{Non-Gauss}}$

- **Step 1: Approximating the product** $p(x)p(y_1|x)$ as Gaussian Proj $[p(x)p(y_1|x)]$

A Toy Problem

□ An Alternative Gaussian Approximation

Consider the simple case of $N = 1$ (only one observation)

True posterior $p(x|\mathbf{y}) \propto \underbrace{p(x)}_{\text{Gauss}} \underbrace{p(y_1|x)}_{\text{Non-Gauss}}$

- **Step 1: Approximating the product** $p(x)p(y_1|x)$ as Gaussian $\text{Proj} [p(x)p(y_1|x)]$
- **Step 2: Divide the Gaussian** $\text{Proj} [p(x)p(y_1|x)]$ by $p(x)$ to obtain a Gaussian

$$\tilde{t}_1(x) = \frac{\text{Proj} [p(x)p(y_1|x)]}{p(x)}$$

taking care of $p(x)$
when approximating
 $p(y_1|x)$

A Toy Problem

□ An Alternative Gaussian Approximation

Consider the simple case of $N = 1$ (only one observation)

True posterior $p(x|y) \propto \underbrace{p(x)}_{\text{Gauss}} \underbrace{p(y_1|x)}_{\text{Non-Gauss}}$

- **Step 1: Approximating the product** $p(x)p(y_1|x)$ as Gaussian $\text{Proj} [p(x)p(y_1|x)]$
- **Step 2: Divide the Gaussian** $\text{Proj} [p(x)p(y_1|x)]$ by $p(x)$ to obtain a Gaussian

$$\tilde{t}_1(x) = \frac{\text{Proj} [p(x)p(y_1|x)]}{p(x)}$$

taking care of $p(x)$ when approximating $p(y_1|x)$

$$q(x|y_1) \propto p(x) \tilde{t}_1(x)$$

Posterior Gauss approximation

A Toy Problem

□ Assumed Density Filtering (ADF)

Consider general case of N observations

True posterior $p(x|\mathbf{y}) \propto p(x)p(y_1|x)p(y_2|x)p(y_3|x)\cdots p(y_N|x)$

Approximate posterior $q(x|\mathbf{y}) \propto p(x)\tilde{t}_1(x)\tilde{t}_2(x)\tilde{t}_3(x)\cdots\tilde{t}_N(x)$

A Toy Problem

□ Assumed Density Filtering (ADF)

Consider general case of N observations

True posterior $p(x|\mathbf{y}) \propto p(x)p(y_1|x)p(y_2|x)p(y_3|x)\cdots p(y_N|x)$

Approximate posterior $q(x|\mathbf{y}) \propto p(x)\tilde{t}_1(x)\tilde{t}_2(x)\tilde{t}_3(x)\cdots\tilde{t}_N(x)$

ADF Algorithm

- Initialize $q^0(x) = p(x)$
- For each new observation y_i

Inclusion $\hat{p}(x) = \frac{q^{i-1}(x)p(y_i|x)}{\int q^{i-1}(x)p(y_i|x)dx}$

Projection $q^i(x) = \text{Proj}[\hat{p}(x)]$

$$\tilde{t}_i(x) \propto \frac{q^i(x)}{q^{i-1}(x)} \quad \text{only implicitly made}$$

A Toy Problem

□ Assumed Density Filtering (ADF)

Consider general case of N observations

True posterior
$$p(x|\mathbf{y}) \propto p(x)p(y_1|x)p(y_2|x)p(y_3|x)\cdots p(y_N|x)$$

Approximate posterior
$$q(x|\mathbf{y}) \propto p(x)\tilde{t}_1(x)\tilde{t}_2(x)\tilde{t}_3(x)\cdots\tilde{t}_N(x)$$

ADF Algorithm

- Initialize $q^0(x) = p(x)$
- For each new observation y_i

Inclusion
$$\hat{p}(x) = \frac{q^{i-1}(x)p(y_i|x)}{\int q^{i-1}(x)p(y_i|x)dx}$$

Projection
$$q^i(x) = \text{Proj}[\hat{p}(x)]$$

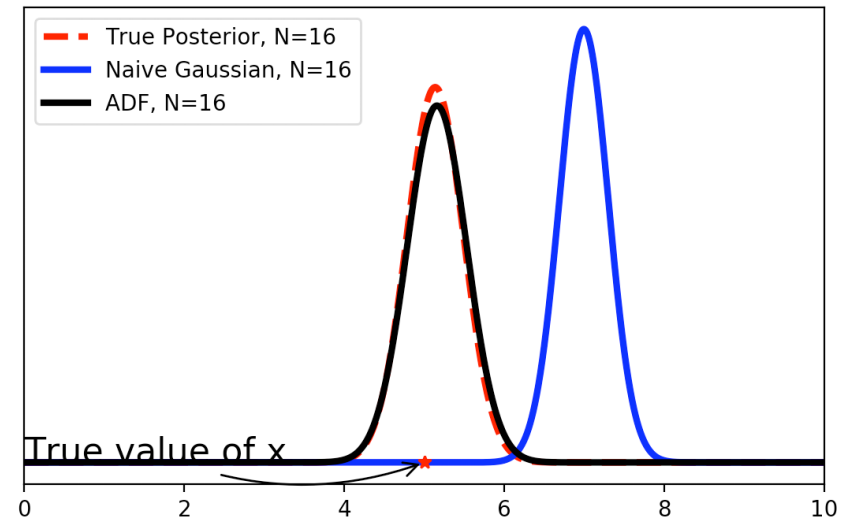
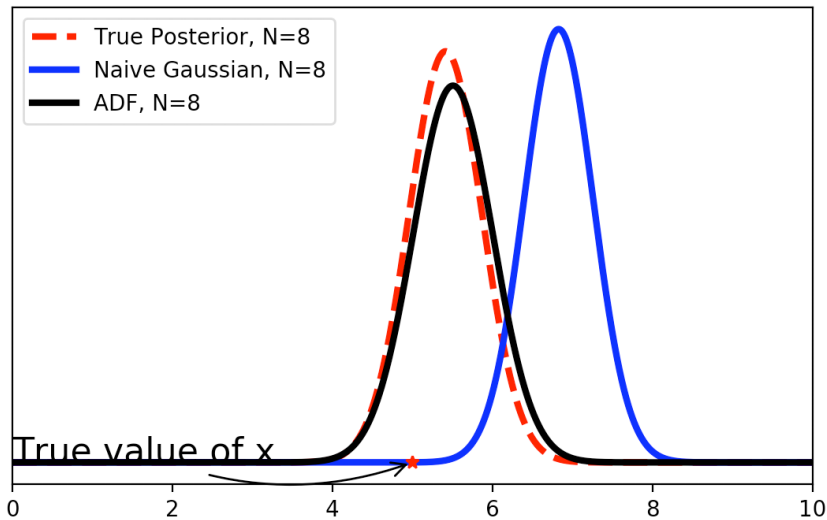
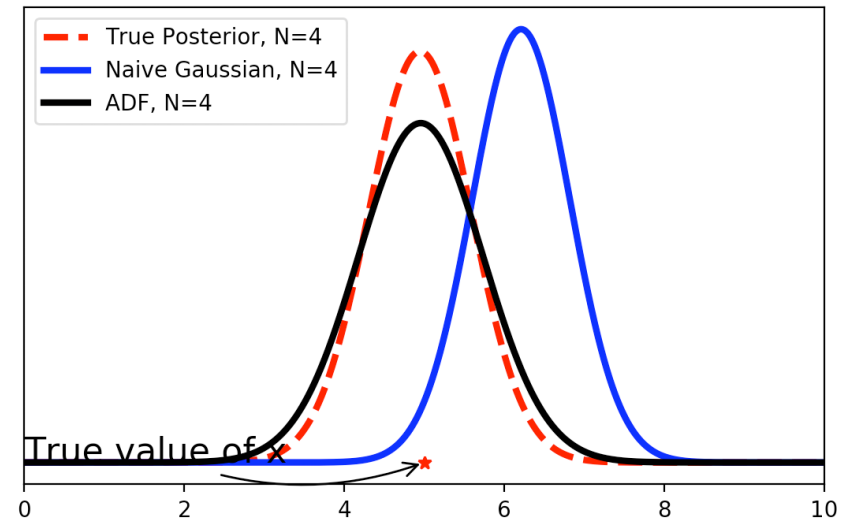
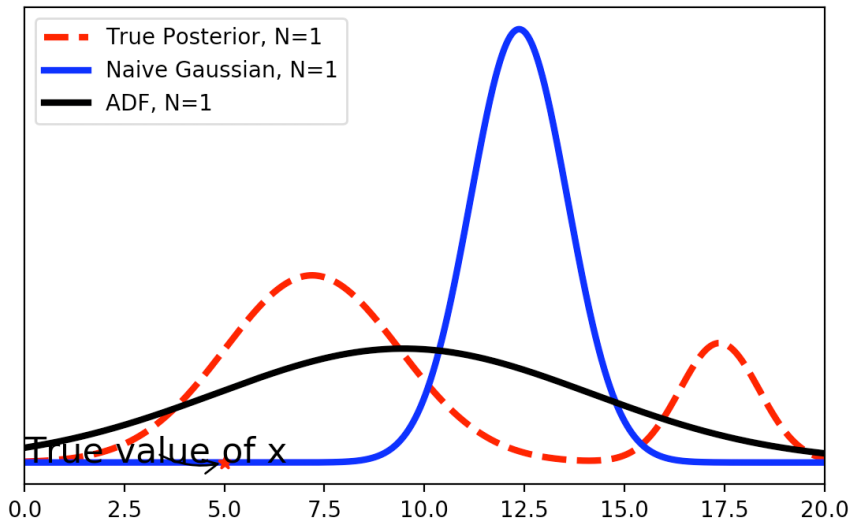
$$\tilde{t}_i(x) \propto \frac{q^i(x)}{q^{i-1}(x)} \quad \text{only implicitly made}$$

ADF is one kind of sequential Gaussian Projector [Minka01b]

$$\begin{array}{c} \xrightarrow{q^{i-1}(x)} \\ \tilde{t}_1(x) \quad \tilde{t}_2(x) \quad p(y_3|x) \quad p(y_4|x) \quad p(y_5|x) \\ \downarrow \\ \tilde{t}_1(x) \quad \tilde{t}_2(x) \quad \tilde{t}_3(x) \quad p(y_4|x) \quad p(y_5|x) \\ \xrightarrow{\quad\quad\quad} q^i(x) \end{array}$$

A Toy Problem

Assumed Density Filtering (ADF)

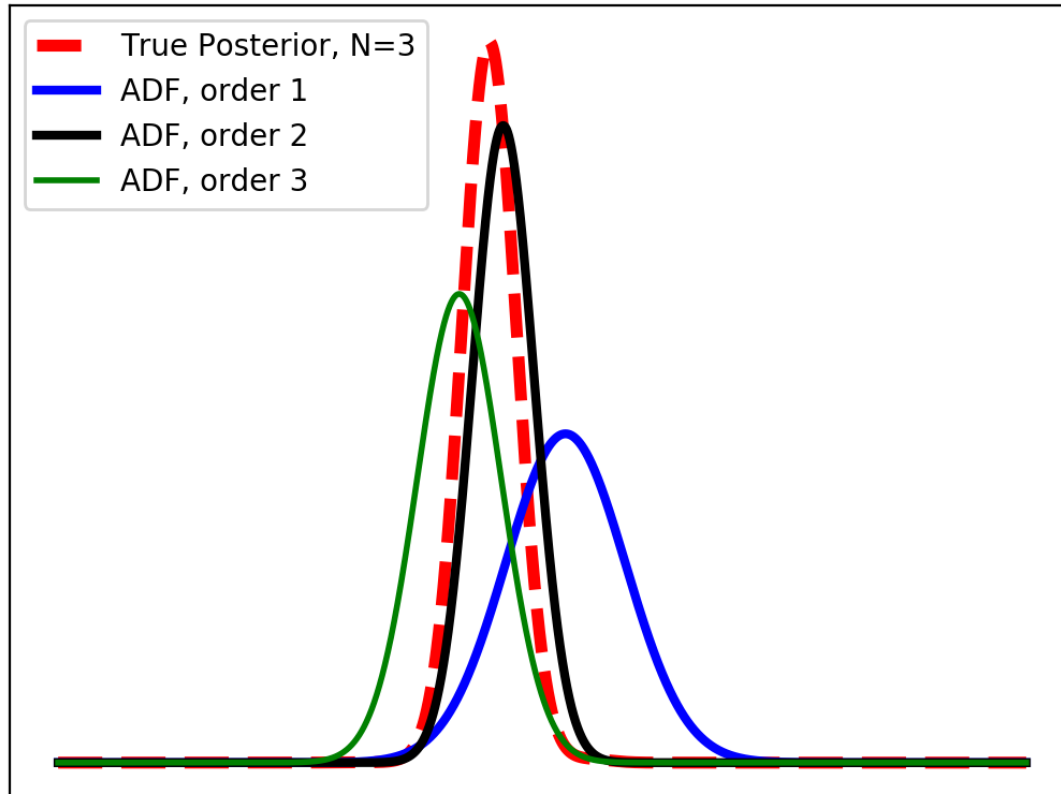


ADF is much better than naive Gauss approximation

A Toy Problem

□ Assumed Density Filtering (ADF)

However, ADF is **sensitive to the order of approximations** !



How to avoid the effect of different ordering

A Toy Problem

□ Expectation Propagation

Expectation Propagation = ADF + Iteratively Refine

A Toy Problem

□ Expectation Propagation

Expectation Propagation = ADF + Iteratively Refine

True posterior $p(x|\mathbf{y}) \propto p(x)p(y_1|x)p(y_2|x)p(y_3|x)\cdots p(y_N|x)$

Approximate posterior $q(x|\mathbf{y}) \propto p(x)\tilde{t}_1(x)\tilde{t}_2(x)\tilde{t}_3(x)\cdots\tilde{t}_N(x)$

EP Algorithm

- Initialize $\tilde{t}_i(x), i = 1\dots N, q(x) = p(x) \prod_i \tilde{t}_i(x)$
- For iter = 1... Num_iter
For i in 1... N

division $q^{\setminus i}(x) \propto \frac{q(x)}{\tilde{t}_i(x)} = p(x) \prod_{j \neq i} \tilde{t}_j(x)$

inclusion $\hat{p}(x) = \frac{q^{\setminus i}(x) p(y_i|x)}{\int q^{\setminus i}(x) p(y_i|x) dx}$

projection $q(x) = \text{Proj}[\hat{p}(x)]$

refinement $\tilde{t}_i(x) \propto \frac{q(x)}{q^{\setminus i}(x)}$

A Toy Problem

□ Expectation Propagation

Expectation Propagation = ADF + Iteratively Refine

True posterior $p(x|\mathbf{y}) \propto p(x)p(y_1|x)p(y_2|x)p(y_3|x)\cdots p(y_N|x)$

Approximate posterior $q(x|\mathbf{y}) \propto p(x)\tilde{t}_1(x)\tilde{t}_2(x)\tilde{t}_3(x)\cdots\tilde{t}_N(x)$

EP Algorithm

- Initialize $\tilde{t}_i(x), i = 1\dots N, q(x) = p(x) \prod_i \tilde{t}_i(x)$
- For iter = 1... Num_iter

For i in 1... N

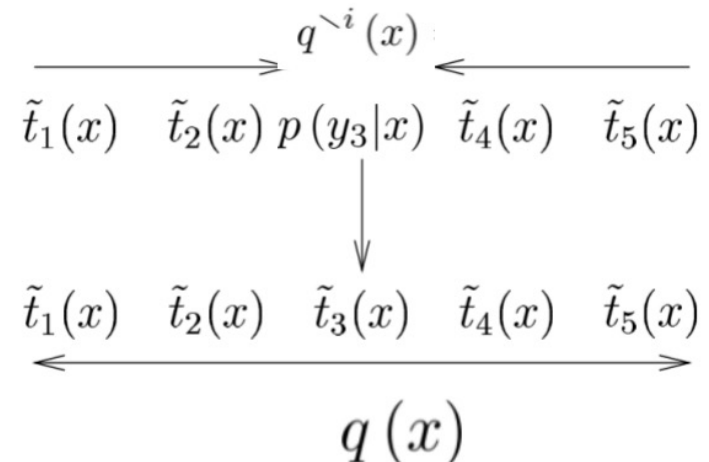
division $q^{\setminus i}(x) \propto \frac{q(x)}{\tilde{t}_i(x)} = p(x) \prod_{j \neq i} \tilde{t}_j(x)$

inclusion $\hat{p}(x) = \frac{q^{\setminus i}(x) p(y_i|x)}{\int q^{\setminus i}(x) p(y_i|x) dx}$

projection $q(x) = \text{Proj}[\hat{p}(x)]$

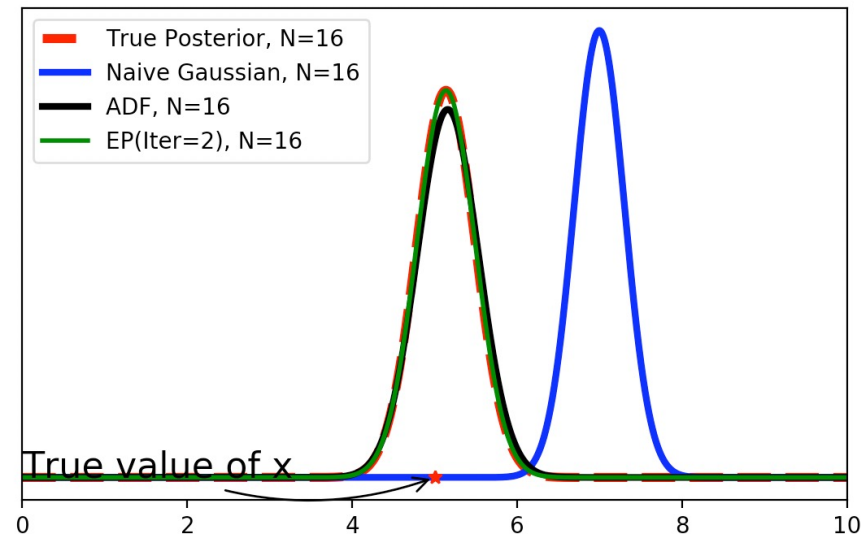
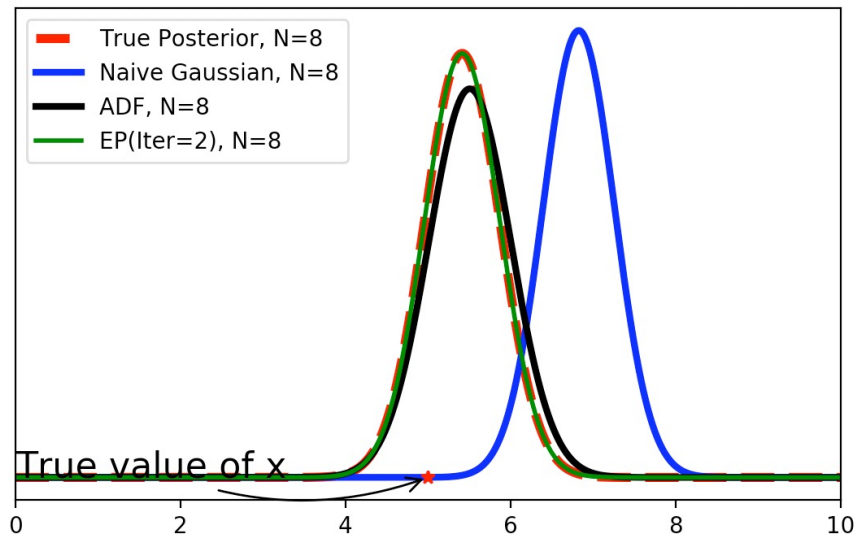
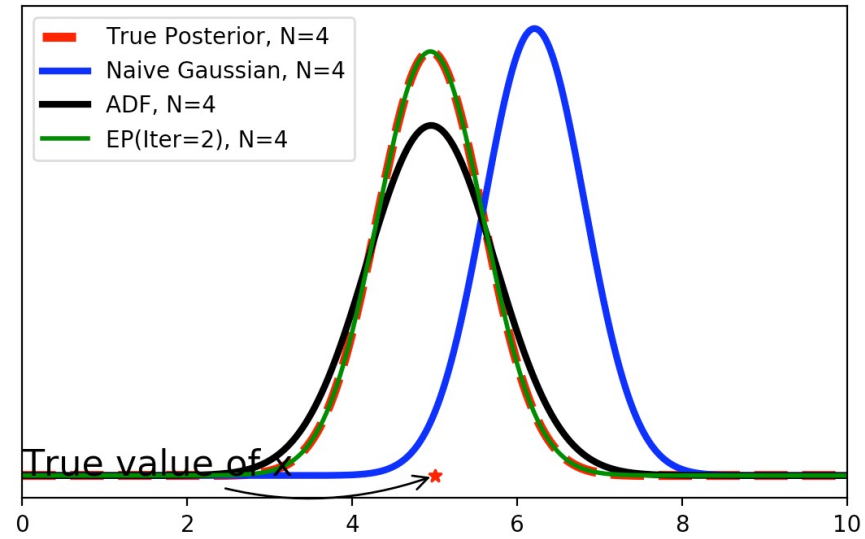
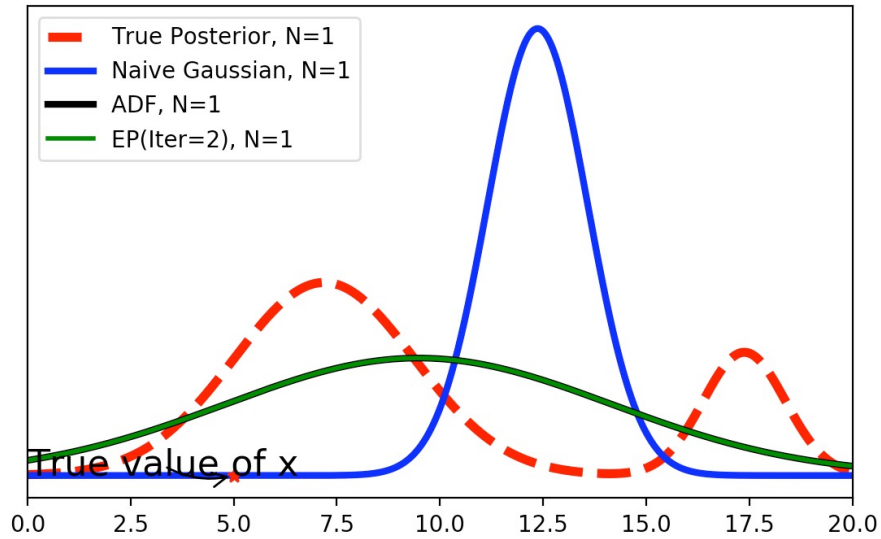
refinement $\tilde{t}_i(x) \propto \frac{q(x)}{q^{\setminus i}(x)}$

EP is an iterative refinement of ADF and is not affected by order [Minka01b]



A Toy Problem

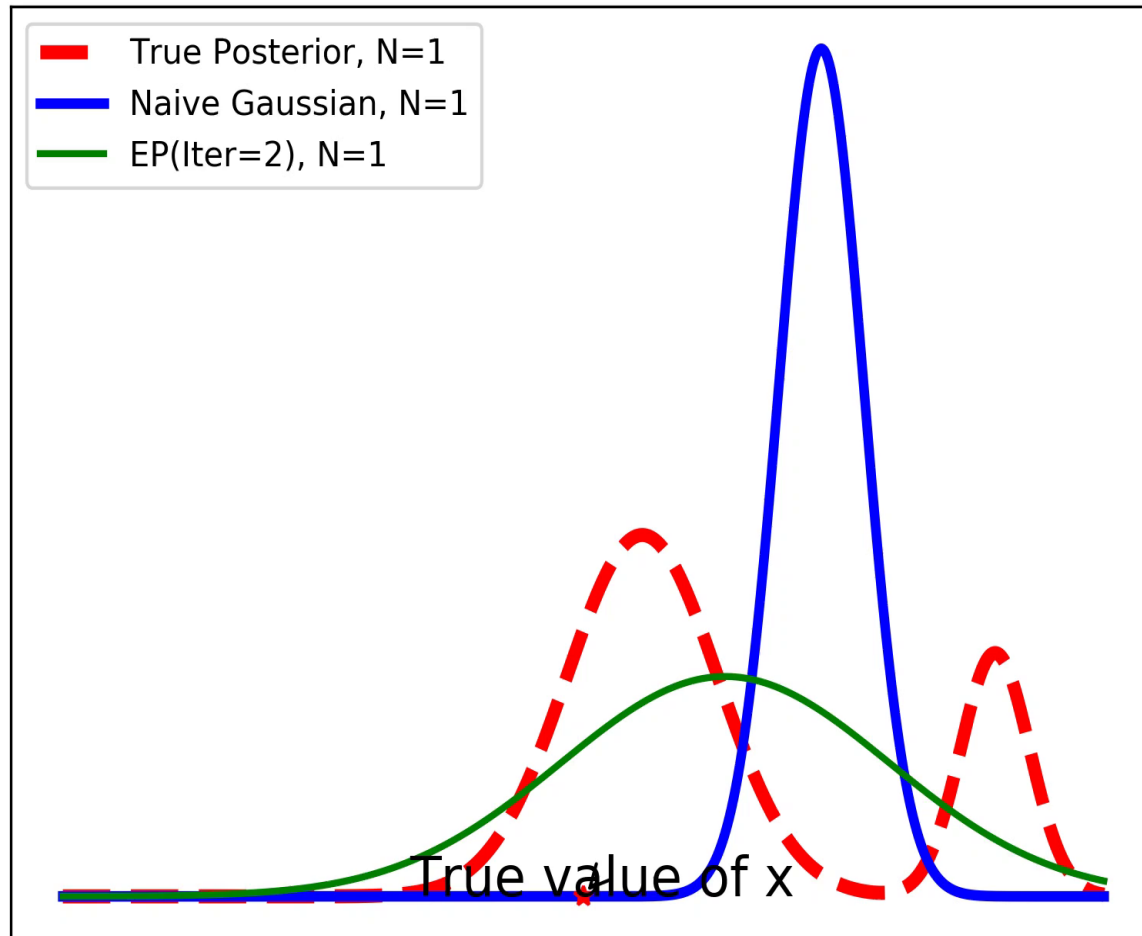
□ Expectation Propagation



EP approximation is close to the true posterior !

A Toy Problem

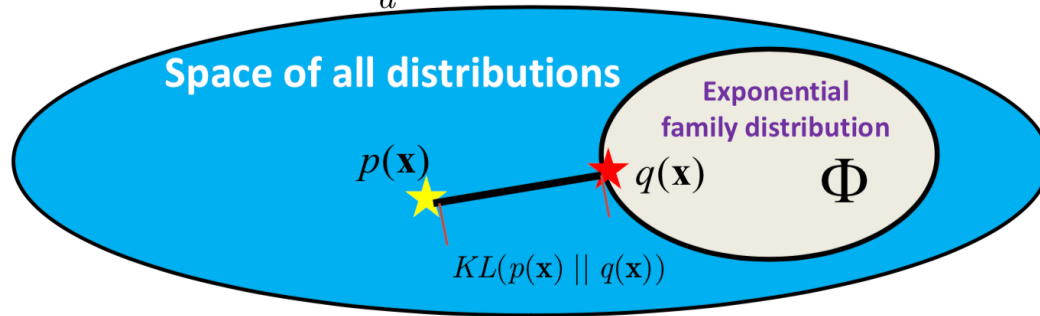
□ Expectation Propagation



EP as Optimization

□ Expectation Propagation (EP) [Minka01] [Opper05]

$$p(\mathbf{x}) = \prod_a f_a(\mathbf{x}) \xrightarrow{\text{approximated as}} q(\mathbf{x}) = \prod_a \tilde{f}_a(\mathbf{x})$$

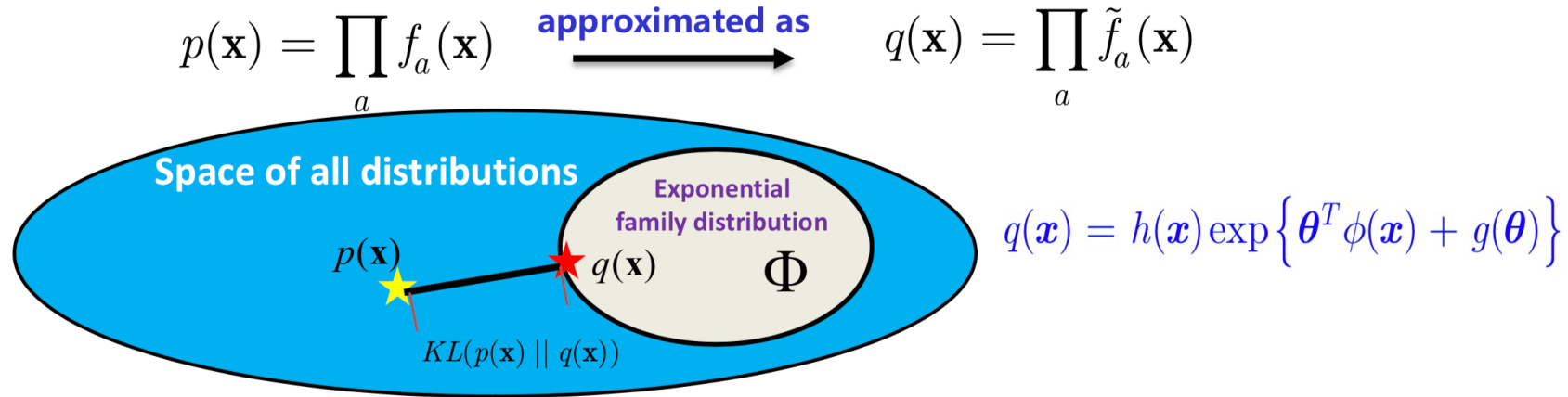


$$q(\mathbf{x}) = h(\mathbf{x}) \exp\{\boldsymbol{\theta}^T \phi(\mathbf{x}) + g(\boldsymbol{\theta})\}$$

Optimization objective: $\min KL(p(\mathbf{x}) || q(\mathbf{x}))$

EP as Optimization

□ Expectation Propagation (EP) [Minka01] [Opper05]



Optimization objective: $\min KL(p(\mathbf{x}) || q(\mathbf{x}))$

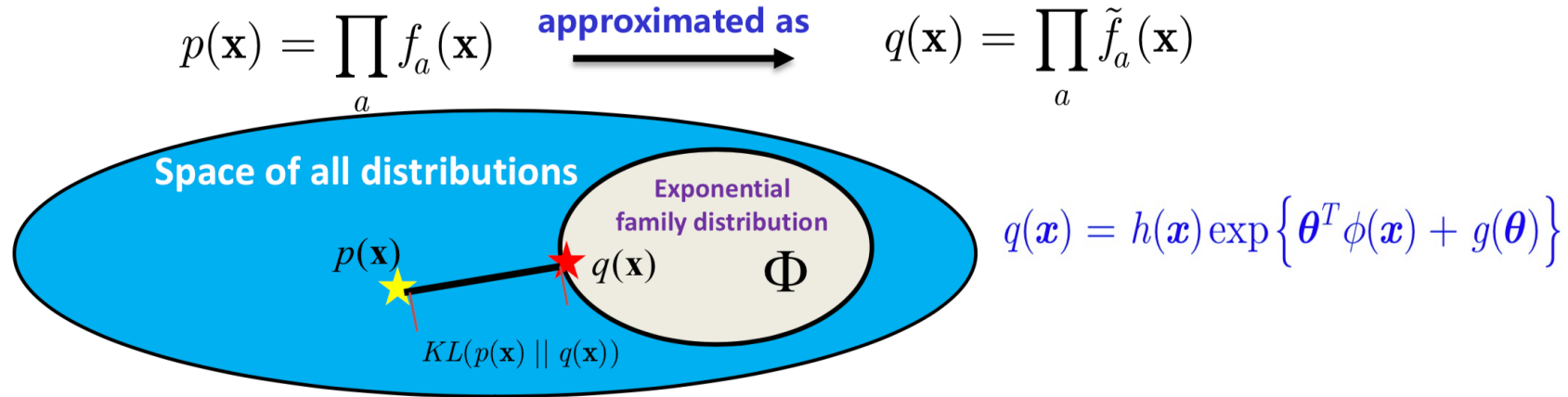
Iterative local optimization

Iteratively refine each factor

$$\tilde{f}_a(\mathbf{x}) = \arg \min_{t(\mathbf{x}) \in \Phi} KL(f_a(\mathbf{x}) \prod_{b \neq a} \tilde{f}_b(\mathbf{x}) || t(\mathbf{x}) \prod_{b \neq a} \tilde{f}_b(\mathbf{x}))$$

EP as Optimization

□ Expectation Propagation (EP) [Minka01] [Opper05]



Optimization objective: $\min KL(p(\mathbf{x}) || q(\mathbf{x}))$

Iterative local optimization

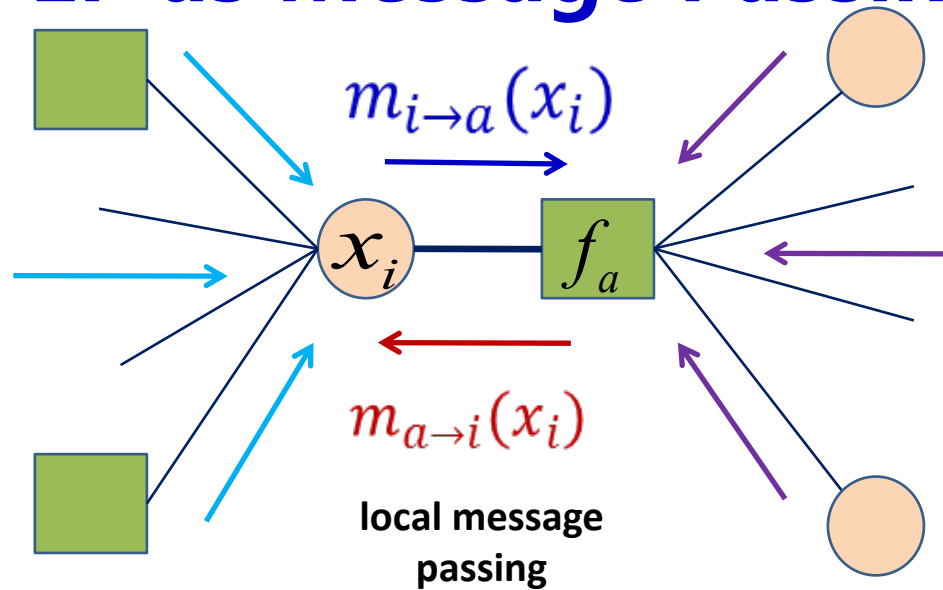
Iteratively refine each factor

$$\tilde{f}_a(\mathbf{x}) = \arg \min_{t(\mathbf{x}) \in \Phi} KL(f_a(\mathbf{x}) \prod_{b \neq a} \tilde{f}_b(\mathbf{x}) || t(\mathbf{x}) \prod_{b \neq a} \tilde{f}_b(\mathbf{x}))$$

- BP minimizes $KL(q || p)$ while EP minimizes $KL(p || q)$
- EP can be also implemented as message passing on factor graph

EP as Message Passing

□ Factor Graph



Expectation Propagation

Factor to variable

$$m_{a \rightarrow i}(x_i) = \frac{\text{Proj} \left[m_{i \rightarrow a}(x_i) \sum_{x_j, j \neq i} f_a(\mathbf{x}_a) \prod_{j \neq i} m_{j \rightarrow a}(x_j) \right]}{m_{i \rightarrow a}(x_i)}$$

Iterations

Variable to factor

$$m_{i \rightarrow a}(x_i) = \prod_{b \neq a} m_{b \rightarrow i}(x_i)$$

Excluding incoming message itself

EP vs. BP

Expectation Propagation

VS.

Belief Propagation

- minimize KL ($p||q$)
- A generalization of BP
- Discrete & continuous variable
- Might iterative without loop

- minimize KL ($q||p$)
- EP with fully factorization
- Discrete variable
- Non-iterative without loop

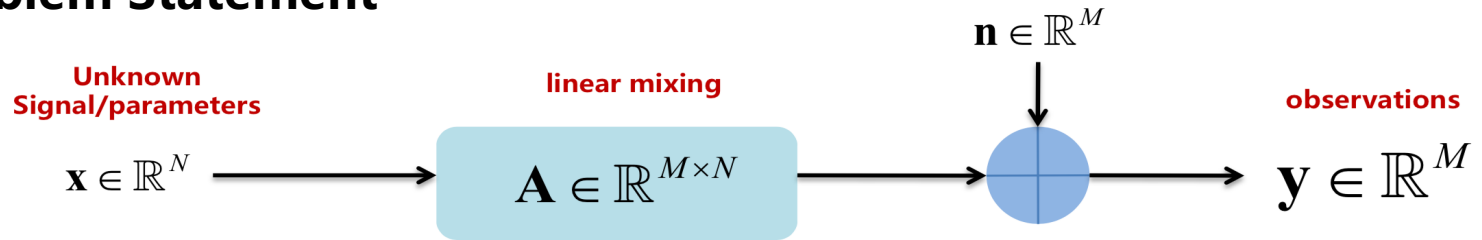
- EP is related to the cavity method in physics [Mezard et al 87] [Opper&Saad 01]

Outline

- Background
- Variational Inference
- Expectation Propagation
- **A Unified EP Perspective on AMP and its extensions**
- Conclusion

Linear Observations

□ Problem Statement



$$\mathbf{x} \sim p_0(\mathbf{x})$$

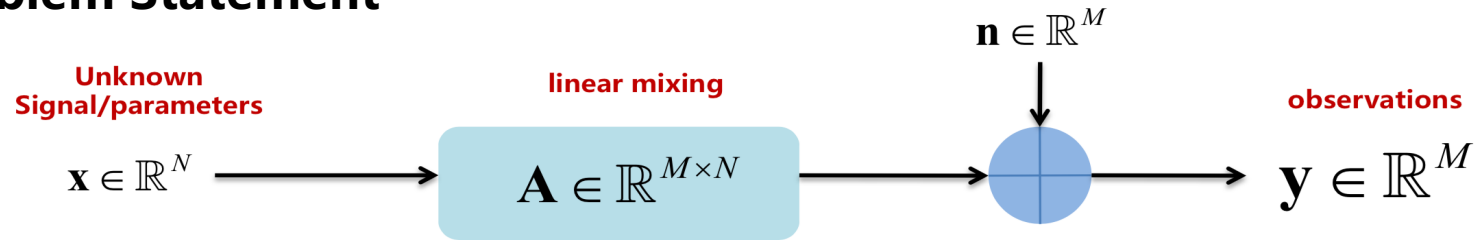
$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}$$

$$\mathbf{n} \sim \mathcal{N}(0, \sigma^2 I)$$

- The goal is to recover signal \mathbf{x} given the observations \mathbf{y} .
- A fundamental problem in communication, compressed sensing, statistics

Linear Observations

□ Problem Statement



$$\mathbf{x} \sim p_0(\mathbf{x})$$

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}$$

$$\mathbf{n} \sim \mathcal{N}(0, \sigma^2 I)$$

- The goal is to recover signal \mathbf{x} given the observations \mathbf{y} .
- A fundamental problem in communication, compressed sensing, statistics

First, we write the the joint distribution can be written as follows

$$p(\mathbf{x}, \mathbf{y}) = p_0(\mathbf{x})p(\mathbf{y}|\mathbf{x})$$

$$= p_0(\mathbf{x}) \frac{1}{(2\pi\sigma^2)^{\frac{M}{2}}} e^{-\frac{(\mathbf{y} - \mathbf{A}\mathbf{x})^T (\mathbf{y} - \mathbf{A}\mathbf{x})}{2\sigma^2}}$$

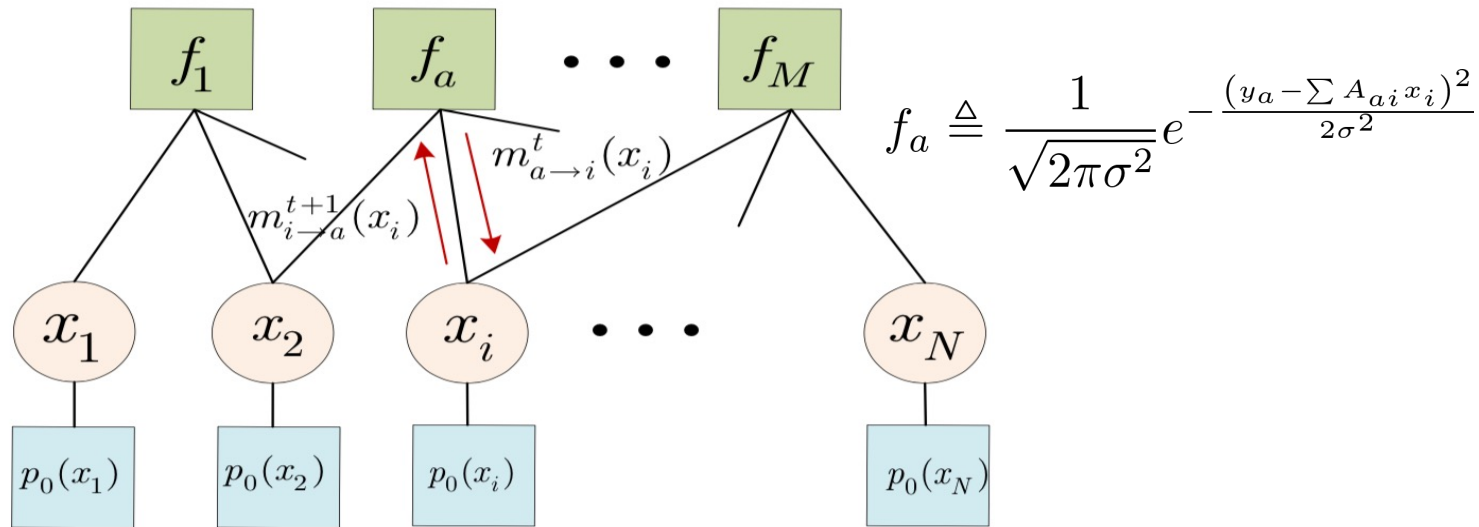
Vector-Form

$$= \prod_{i=1}^N p_0(x_i) \prod_{a=1}^M \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_a - \sum A_{ai}x_i)^2}{2\sigma^2}}$$

Fully-Factorized

Linear Observations

□ Fully-Factorized Factor Graph



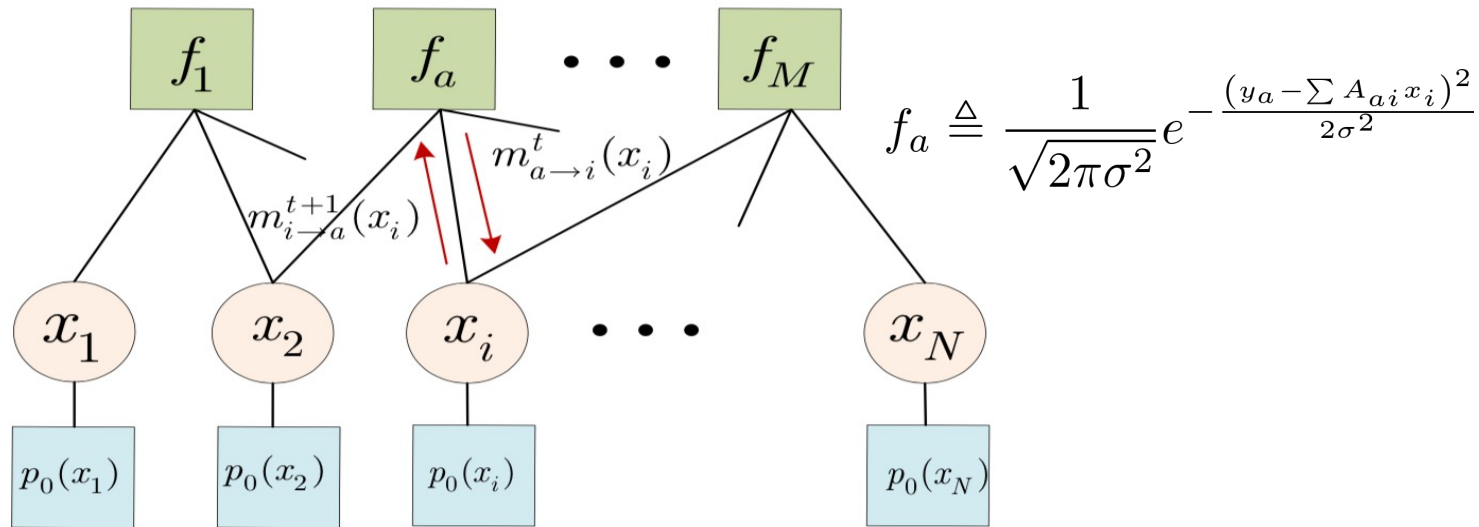
Expectation Propagation (EP)

$$m_{a \rightarrow i}^t(x_i) \propto \frac{\text{Proj}_{\Phi} \left[m_{i \rightarrow a}^t(x_i) \int \prod_{j \neq i} m_{j \rightarrow a}^t(x_j) p(y_a | \mathbf{x}) \right]}{m_{i \rightarrow a}^t(x_i)}$$

$$m_{i \rightarrow a}^{t+1}(x_i) \propto \frac{\text{Pro}_{\Phi} \left[p_0(x_i) \prod_b m_{b \rightarrow i}^t(x_i) \right]}{m_{a \rightarrow i}^t(x_i)}$$

Linear Observations

□ Fully-Factorized Factor Graph



Expectation Propagation (EP)

$$m_{a \rightarrow i}^t(x_i) \propto \frac{\text{Proj}_{\Phi} \left[m_{i \rightarrow a}^t(x_i) \int \prod_{j \neq i} m_{j \rightarrow a}^t(x_j) p(y_a | \mathbf{x}) \right]}{m_{i \rightarrow a}^t(x_i)}$$

$$m_{i \rightarrow a}^{t+1}(x_i) \propto \frac{\text{Pro}_{\Phi} \left[p_0(x_i) \prod_b m_{b \rightarrow i}^t(x_i) \right]}{m_{a \rightarrow i}^t(x_i)}$$

It seems quite easy ?

Linear Observations

□ An EP Perspective on AMP

$$m_{a \rightarrow i}^t(x_i) \propto \mathcal{N}(x_i; \hat{x}_{a \rightarrow i}^t, v_{a \rightarrow i}^t)$$

$$m_{i \rightarrow a}^{t+1}(x_i) \propto \mathcal{N}(x_i; \hat{x}_{i \rightarrow a}^{t+1}, v_{i \rightarrow a}^{t+1})$$

where

$$V_{a \rightarrow i}^t = \sum_{j \neq i} |A_{aj}|^2 \nu_{j \rightarrow a}^t \quad Z_{a \rightarrow i}^t = \sum_{j \neq i} A_{aj} \hat{x}_{j \rightarrow a}^t$$

$$\hat{x}_{a \rightarrow i}^t = \frac{y_a - Z_{a \rightarrow i}^t}{A_{ai}}, \quad v_{a \rightarrow i}^t = \frac{\sigma^2 + V_{a \rightarrow i}^t}{|A_{ai}|^2}$$

$$\Sigma_i^t = \left[\sum_a \frac{|A_{ai}|^2}{\sigma^2 + V_{a \rightarrow i}^t} \right]^{-1}, \quad R_i^t = \Sigma_i^t \sum_a \frac{A_{ai}^* (y_a - Z_{a \rightarrow i}^t)}{\sigma^2 + V_{a \rightarrow i}^t}$$

$$\hat{x}_i^{t+1} = f_a(R_i^t, \Sigma_i^t) \quad \hat{\nu}_i^{t+1} = f_c(R_i^t, \Sigma_i^t)$$

$$\frac{1}{\nu_{i \rightarrow a}^{t+1}} = \frac{1}{\nu_i^{t+1}} - \frac{|A_{ai}|^2}{\sigma^2 + V_{a \rightarrow i}^t},$$

$$\hat{x}_{i \rightarrow a}^{t+1} = \nu_{i \rightarrow a}^{t+1} \left(\frac{\hat{x}_i^{t+1}}{\nu_i^{t+1}} - \frac{A_{ai}^* (y_a - Z_{a \rightarrow i}^t)}{\sigma^2 + V_{a \rightarrow i}^t} \right).$$

Linear Observations

□ An EP Perspective on AMP

$$m_{a \rightarrow i}^t(x_i) \propto \mathcal{N}(x_i; \hat{x}_{a \rightarrow i}^t, v_{a \rightarrow i}^t)$$

$$m_{i \rightarrow a}^{t+1}(x_i) \propto \mathcal{N}(x_i; \hat{x}_{i \rightarrow a}^{t+1}, v_{i \rightarrow a}^{t+1})$$

where

$$V_{a \rightarrow i}^t = \sum_{j \neq i} |A_{aj}|^2 v_{j \rightarrow a}^t \quad Z_{a \rightarrow i}^t = \sum_{j \neq i} A_{aj} \hat{x}_{j \rightarrow a}^t$$

$$\hat{x}_{a \rightarrow i}^t = \frac{y_a - Z_{a \rightarrow i}^t}{A_{ai}}, \quad v_{a \rightarrow i}^t = \frac{\sigma^2 + V_{a \rightarrow i}^t}{|A_{ai}|^2}$$

$$\Sigma_i^t = \left[\sum_a \frac{|A_{ai}|^2}{\sigma^2 + V_{a \rightarrow i}^t} \right]^{-1}, \quad R_i^t = \Sigma_i^t \sum_a \frac{A_{ai}^* (y_a - Z_{a \rightarrow i}^t)}{\sigma^2 + V_{a \rightarrow i}^t}$$

$$\hat{x}_i^{t+1} = f_a(R_i^t, \Sigma_i^t) \quad \hat{v}_i^{t+1} = f_c(R_i^t, \Sigma_i^t)$$

$$\frac{1}{v_{i \rightarrow a}^{t+1}} = \frac{1}{v_i^{t+1}} - \frac{|A_{ai}|^2}{\sigma^2 + V_{a \rightarrow i}^t},$$

$$\hat{x}_{i \rightarrow a}^{t+1} = v_{i \rightarrow a}^{t+1} \left(\frac{\hat{x}_i^{t+1}}{v_i^{t+1}} - \frac{A_{ai}^* (y_a - Z_{a \rightarrow i}^t)}{\sigma^2 + V_{a \rightarrow i}^t} \right)$$

Still Too Complicated!

- However, the number of messages are $O(MN)$, which is still intractable for high-dimensional problems

Linear Observations

□ An EP Perspective on AMP

$$m_{a \rightarrow i}^t(x_i) \propto \mathcal{N}(x_i; \hat{x}_{a \rightarrow i}^t, v_{a \rightarrow i}^t)$$

where

$$m_{i \rightarrow a}^{t+1}(x_i) \propto \mathcal{N}(x_i; \hat{x}_{i \rightarrow a}^{t+1}, v_{i \rightarrow a}^{t+1})$$

• However, the **number of messages are $O(MN)$** , which is still intractable for high-dimensional problems

• To reduce the number of messages, **neglect the high-order terms** in large system limit.

$$Z_a^t = \sum_i A_{ai} \hat{x}_{i \rightarrow a}^t \quad V_a^t = \sum_i |A_{ai}|^2 v_{i \rightarrow a}^t$$

$$Z_{a \rightarrow i}^t = Z_a^t - A_{ai} \hat{x}_{i \rightarrow a}^t \quad \rightarrow \text{Be careful!}$$

$$V_{a \rightarrow i}^t = V_a^t - |A_{ai}|^2 v_{i \rightarrow a}^t \quad \rightarrow V_{a \rightarrow i}^t \approx V_a^t$$

$$v_{i \rightarrow a}^{t+1} \approx v_i^{t+1} \quad \rightarrow V_a^t \approx \sum_i |A_{ai}|^2 v_i^t$$

$$V_{a \rightarrow i}^t = \sum_{j \neq i} |A_{aj}|^2 v_{j \rightarrow a}^t \quad Z_{a \rightarrow i}^t = \sum_{j \neq i} A_{aj} \hat{x}_{j \rightarrow a}^t$$

$$\hat{x}_{a \rightarrow i}^t = \frac{y_a - Z_{a \rightarrow i}^t}{A_{ai}}, \quad v_{a \rightarrow i}^t = \frac{\sigma^2 + V_{a \rightarrow i}^t}{|A_{ai}|^2}$$

$$\Sigma_i^t = \left[\sum_a \frac{|A_{ai}|^2}{\sigma^2 + V_{a \rightarrow i}^t} \right]^{-1} \quad R_i^t = \Sigma_i^t \sum_a \frac{A_{ai}^* (y_a - Z_{a \rightarrow i}^t)}{\sigma^2 + V_{a \rightarrow i}^t}$$

$$\hat{x}_i^{t+1} = f_a(R_i^t, \Sigma_i^t) \quad \hat{v}_i^{t+1} = f_c(R_i^t, \Sigma_i^t)$$

$$\frac{1}{v_{i \rightarrow a}^{t+1}} = \frac{1}{v_i^{t+1}} - \frac{|A_{ai}|^2}{\sigma^2 + V_{a \rightarrow i}^t},$$

$$\hat{x}_{i \rightarrow a}^{t+1} = v_{i \rightarrow a}^{t+1} \left(\frac{\hat{x}_i^{t+1}}{v_i^{t+1}} - \frac{A_{ai}^* (y_a - Z_{a \rightarrow i}^t)}{\sigma^2 + V_{a \rightarrow i}^t} \right)$$

Still Too Complicated!

Linear Observations

□ An EP Perspective on AMP

$$m_{a \rightarrow i}^t(x_i) \propto \mathcal{N}(x_i; \hat{x}_{a \rightarrow i}^t, v_{a \rightarrow i}^t)$$

where

$$m_{i \rightarrow a}^{t+1}(x_i) \propto \mathcal{N}(x_i; \hat{x}_{i \rightarrow a}^{t+1}, v_{i \rightarrow a}^{t+1})$$

• However, the **number of messages are $O(MN)$** , which is still intractable for high-dimensional problems

• To reduce the number of messages, **neglect the high-order terms** in large system limit.

$$Z_a^t = \sum_i A_{ai} \hat{x}_{i \rightarrow a}^t \quad V_a^t = \sum_i |A_{ai}|^2 v_{i \rightarrow a}^t$$

$$Z_{a \rightarrow i}^t = Z_a^t - A_{ai} \hat{x}_{i \rightarrow a}^t \quad \rightarrow \text{Be careful!}$$

$$V_{a \rightarrow i}^t = V_a^t - |A_{ai}|^2 v_{i \rightarrow a}^t \quad \rightarrow V_{a \rightarrow i}^t \approx V_a^t$$

$$v_{i \rightarrow a}^{t+1} \approx v_i^{t+1} \quad \rightarrow V_a^t \approx \sum_i |A_{ai}|^2 v_i^t$$

After some algebra, the number of messages is reduced to **$O(M+N)$** and we obtain AMP

$$V_{a \rightarrow i}^t = \sum_{j \neq i} |A_{aj}|^2 v_{j \rightarrow a}^t \quad Z_{a \rightarrow i}^t = \sum_{j \neq i} A_{aj} \hat{x}_{j \rightarrow a}^t$$

$$\hat{x}_{a \rightarrow i}^t = \frac{y_a - Z_{a \rightarrow i}^t}{A_{ai}}, \quad v_{a \rightarrow i}^t = \frac{\sigma^2 + V_{a \rightarrow i}^t}{|A_{ai}|^2}$$

$$\Sigma_i^t = \left[\sum_a \frac{|A_{ai}|^2}{\sigma^2 + V_{a \rightarrow i}^t} \right]^{-1} \quad R_i^t = \Sigma_i^t \sum_a \frac{A_{ai}^* (y_a - Z_{a \rightarrow i}^t)}{\sigma^2 + V_{a \rightarrow i}^t}$$

$$\hat{x}_i^{t+1} = f_a(R_i^t, \Sigma_i^t) \quad \hat{v}_i^{t+1} = f_c(R_i^t, \Sigma_i^t)$$

$$\frac{1}{v_{i \rightarrow a}^{t+1}} = \frac{1}{v_i^{t+1}} - \frac{|A_{ai}|^2}{\sigma^2 + V_{a \rightarrow i}^t}$$

$$\hat{x}_{i \rightarrow a}^{t+1} = v_{i \rightarrow a}^{t+1} \left(\frac{\hat{x}_i^{t+1}}{v_i^{t+1}} - \frac{A_{ai}^* (y_a - Z_{a \rightarrow i}^t)}{\sigma^2 + V_{a \rightarrow i}^t} \right)$$

Still Too Complicated!

Initialization AMP Algorithm

Loop: For $t = 1, \dots, T$

Factor node update $\left\{ \begin{array}{l} V_a^t = \sum_i A_{ai}^2 v_i^t \\ Z_a^t = \sum_i A_{ai} \hat{x}_i^t - \frac{(y_a - Z_a^{t-1})}{\sigma^2 + V_a^{t-1}} V_a^t \end{array} \right.$ **Onsager term**

Variable node update $\left\{ \begin{array}{l} \Sigma_i^t = 1 / \sum_a \frac{A_{ai}^2}{\sigma^2 + V_a^t} \\ R_i^t = \hat{x}_i^t + \Sigma_i^t \sum_a \frac{A_{ai} (y_a - Z_a^t)}{\sigma^2 + V_a^t} \end{array} \right.$ **Linear Complexity $O(MN)$**

End $\hat{x}_i^{t+1} = E(x_i | R_i^t, \Sigma_i^t), \hat{v}_i^{t+1} = \text{Var}(x_i | R_i^t, \Sigma_i^t)$

Relation to AMP

□ An EP Perspective on AMP

AMP iteratively decouples the original **vector inference** problem to **scalar inference** problems

$$y = Ax + n \quad \xrightarrow{\text{decoupled}} \quad \begin{cases} R_1 = x_1 + \tilde{n}_1 \\ \vdots \\ R_N = x_N + \tilde{n}_N \end{cases} \quad \text{decoupling principle}$$

• Comments

- ✓ The first AMP-like method was derived by Kabashima for CDMA detection [Kabashima 03] and later derived by Donoho et. al for compressed sensing [DMM09].
- ✓ For i.i.d. Gaussian \mathbf{A} , AMP is proved to be asymptotically Bayesian optimal and rigorously analyzed via state evolution (SE) [BM11]
- ✓ For general matrices \mathbf{A} , AMP may diverge [BM11]

Relation to AMP

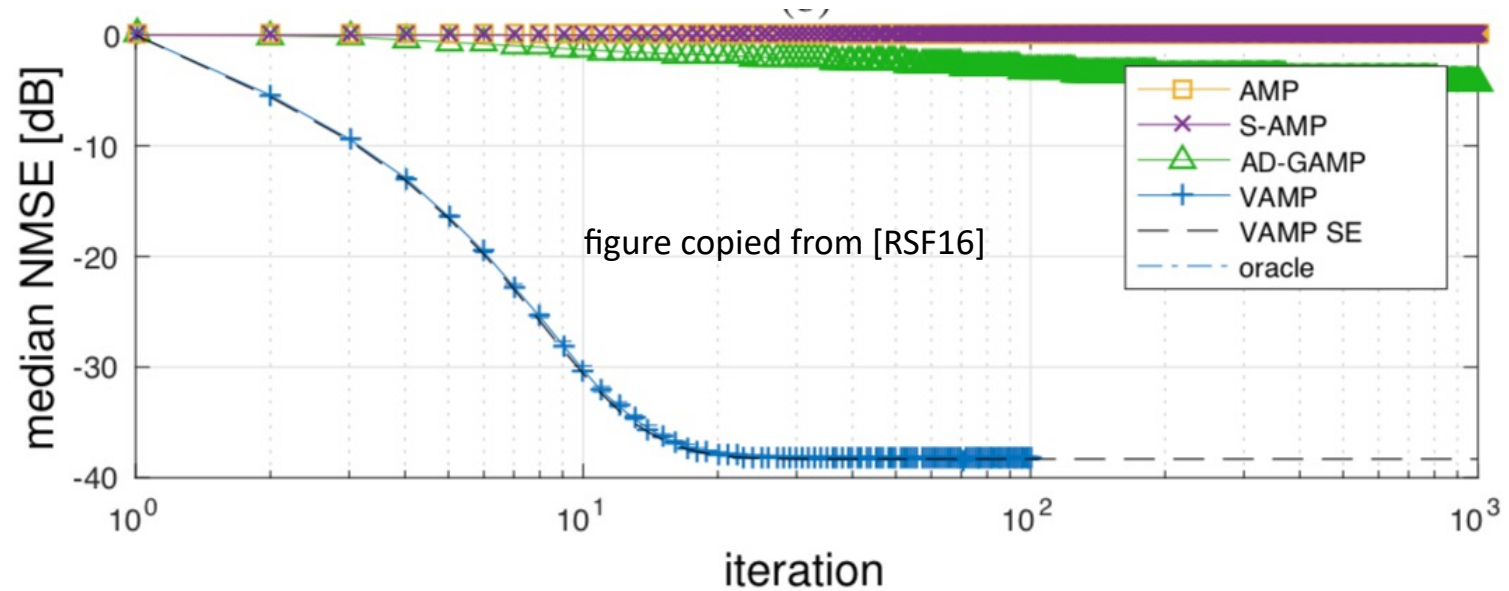
□ An EP Perspective on AMP

AMP iteratively decouples the original **vector inference** problem to **scalar inference** problems

$$y = Ax + n \quad \xrightarrow{\text{decoupled}} \quad \begin{cases} R_1 = x_1 + \tilde{n}_1 \\ \vdots \\ R_N = x_N + \tilde{n}_N \end{cases} \quad \text{decoupling principle}$$

• Comments

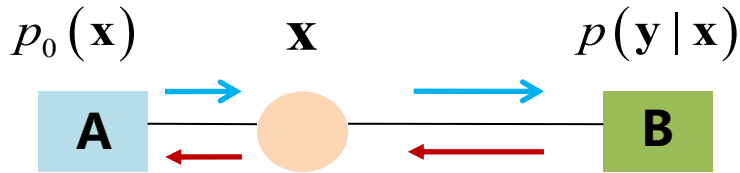
- ✓ The first AMP-like method was derived by Kabashima for CDMA detection [Kabashima 03] and later derived by Donoho et. al for compressed sensing [DMM09].
- ✓ For i.i.d. Gaussian **A**, AMP is proved to be asymptotically Bayesian optimal and rigorously analyzed via state evolution (SE) [BM11]
- ✓ For general matrices **A**, AMP may diverge [BM11]
- ✓ **Vector AMP (VAMP) converges for right-rotationally invariant matrices** [RSF16]



EP Perspective on VAMP

□ Vector-form Factor Graph

$$p(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^N p_0(x_i) \prod_{a=1}^M \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_a - \sum A_{ai}x_i)^2}{2\sigma^2}}$$



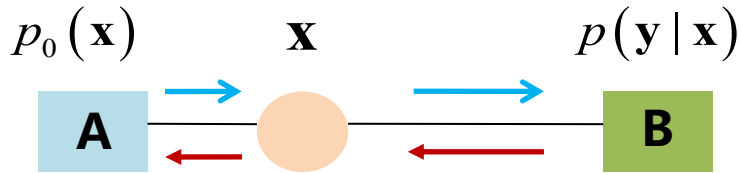
vector-form factor graph

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{\frac{M}{2}}} e^{-\frac{(\mathbf{y}-\mathbf{Ax})^T(\mathbf{y}-\mathbf{Ax})}{2\sigma^2}}$$

EP Perspective on VAMP

□ Vector-form Factor Graph

$$p(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^N p_0(x_i) \prod_{a=1}^M \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_a - \sum A_{ai}x_i)^2}{2\sigma^2}}$$



vector-form factor graph

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{\frac{M}{2}}} e^{-\frac{(\mathbf{y}-\mathbf{Ax})^T(\mathbf{y}-\mathbf{Ax})}{2\sigma^2}}$$

$$m_{A \rightarrow x}(x_i) = \frac{\text{Proj}[p_0(x_i) m_{x \rightarrow A}(x_i)]}{m_{x \rightarrow A}(x_i)} = \mathcal{N}(x_i; m_{i \rightarrow A}, v_{i \rightarrow A})$$

$$m_{x \rightarrow B}(x_i) = m_{A \rightarrow x}(x_i)$$

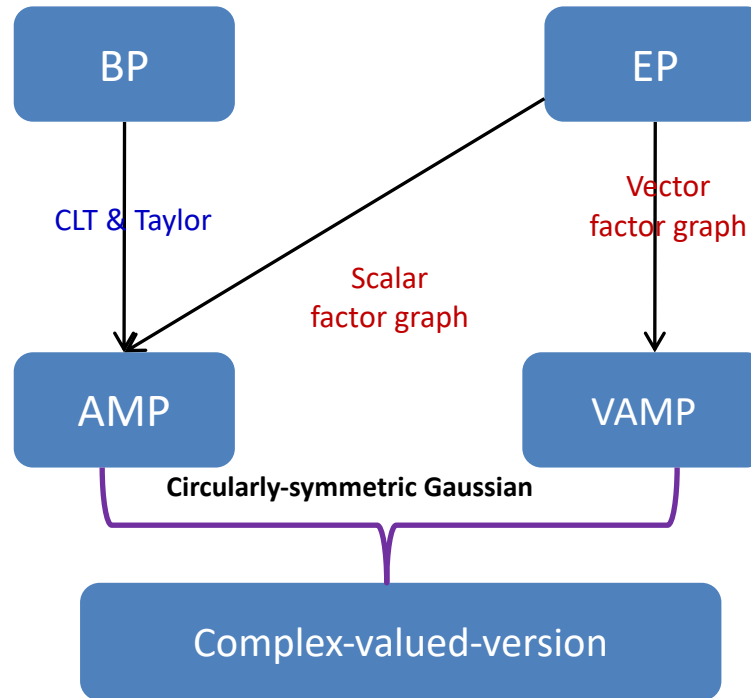
$$m_{B \rightarrow x}(x_i) = \int \mathcal{N}(\mathbf{y}; \mathbf{Ax}, \sigma^2 \mathbf{I}) \prod m_{x \rightarrow B}(x_i) dx_{j \neq i} = \mathcal{N}(x_i; m_{B \rightarrow i}, v_{B \rightarrow i})$$

$$m_{x \rightarrow A}(x_i) = m_{B \rightarrow x}(x_i)$$

**This is exactly the
MMSE form of VAMP**
[RSF16]

A Unified Perspective

□ An EP Perspective on AMP

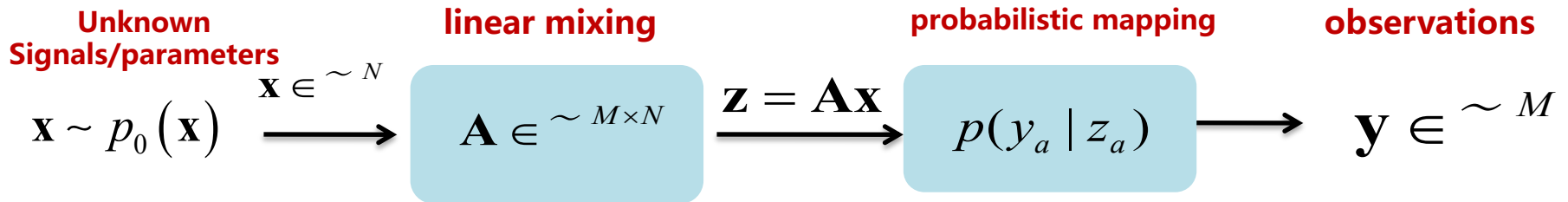


• The EP perspective of AMP and VAMP:

- ✓ Explicitly **establishing the relationship between AMP**
- ✓ Simplifying the extension of AMP to the **complex-valued AMP** (simply using circularly-symmetric Gaussian) [MWKL15b]
- ✓ **Providing a unified view of AMP and VAMP** (derived from scalar EP [MWKL15a] and vector EP [RSF16], respectively)

NonLinear Observations

□ Background



- The measurements are often obtained **in a nonlinear way**
 - one-bit (quantized) compressed sensing
 - phase retrieval
 - logistic regression
 -

Inference on Generalized linear model (GLM)

NonLinear Observations

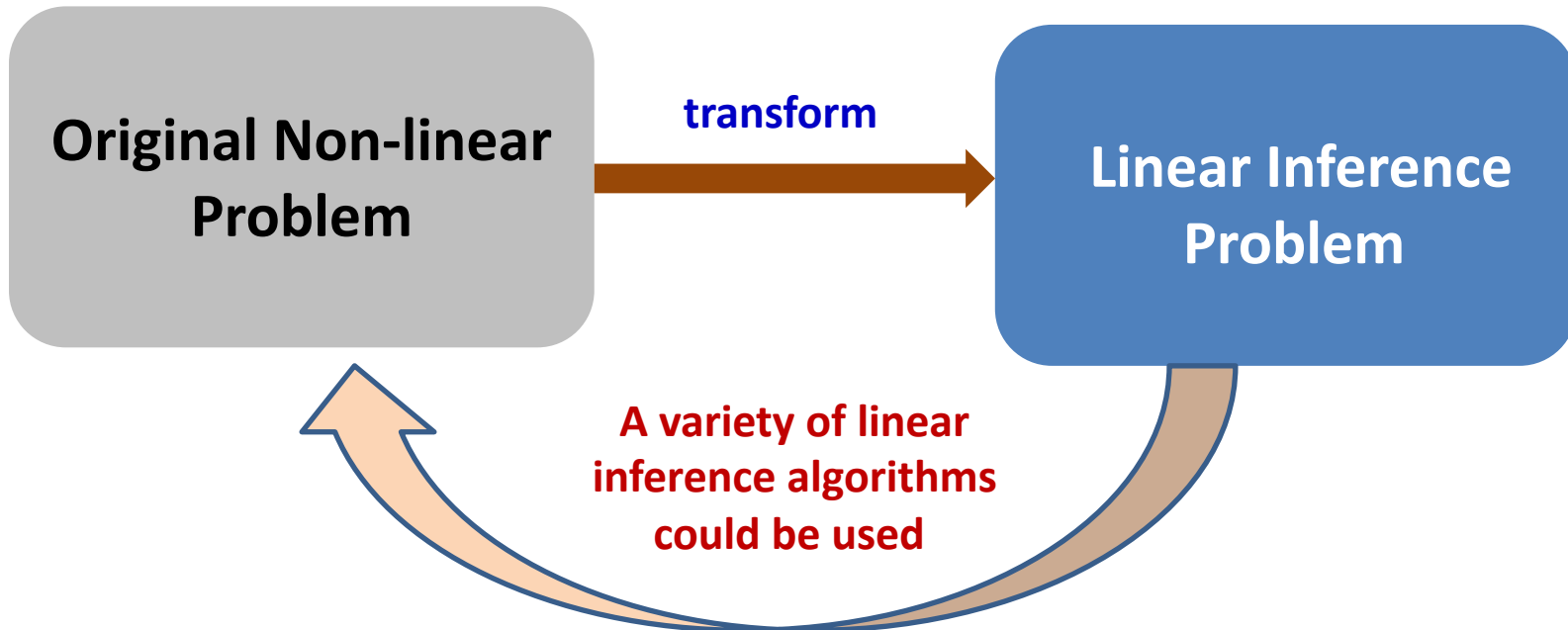
Basic Idea:

Is it possible to transform the nonlinear inference problem to linear inference problems?

NonLinear Observations

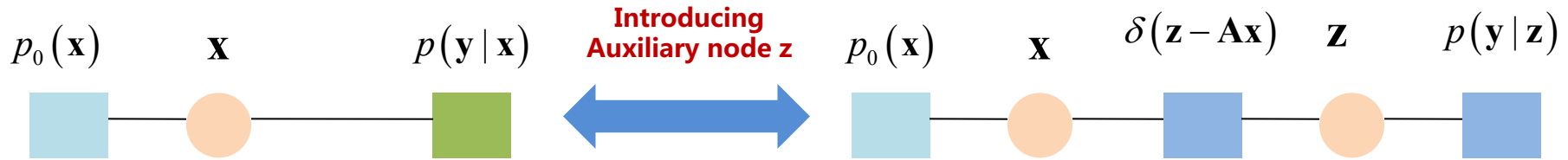
Basic Idea:

Is it possible to transform the nonlinear inference problem to linear inference problems?



A Unified Inference Framework for GLM

□ Two Equivalent Factor Graphs of GLM

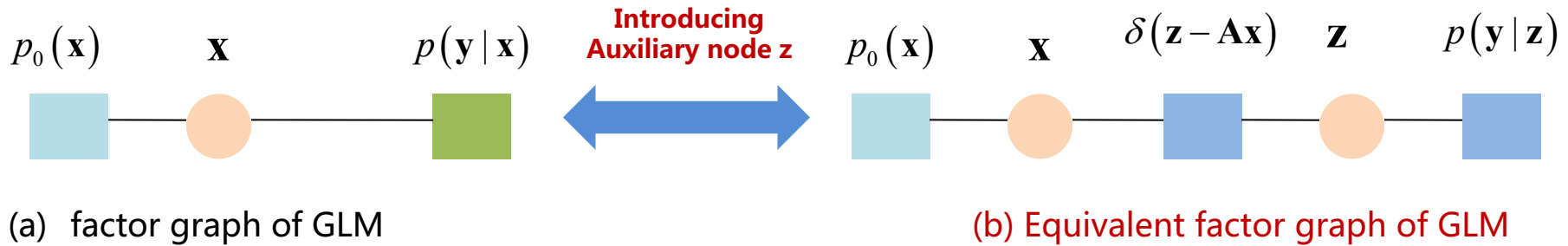


(a) factor graph of GLM

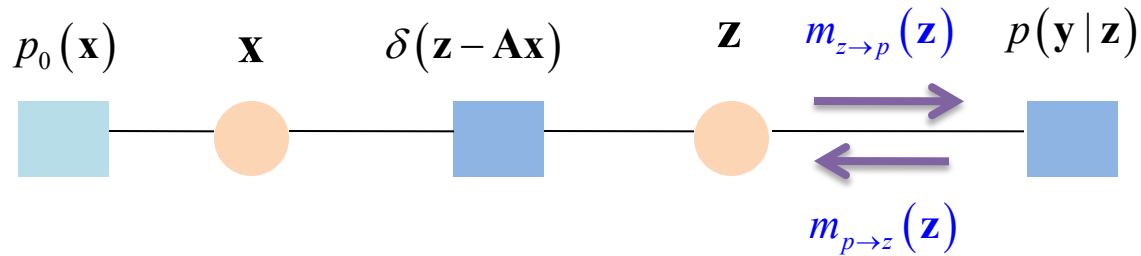
(b) Equivalent factor graph of GLM

A Unified Inference Framework for GLM

□ Two Equivalent Factor Graphs of GLM



□ Decoupling GLM into SLM via EP



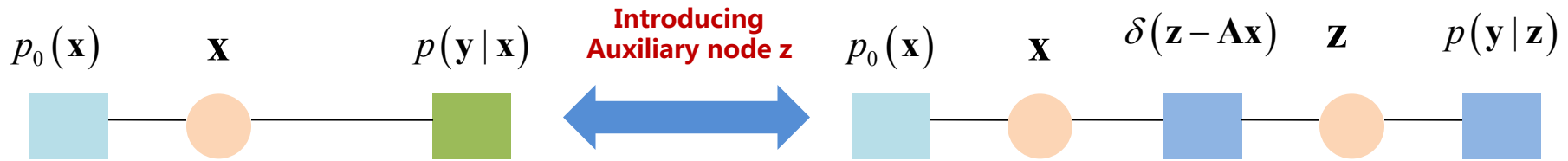
$$m_{z \rightarrow p}^{t-1}(\mathbf{z}) \propto N\left(\mathbf{z}; z_A^{ext}(t-1), v_A^{ext}(t-1)I\right)$$

EP message passing
(t -th iteration)

$$m_{p \rightarrow z}^t(\mathbf{z}) \propto \frac{\text{Proj}_{\Phi}\left(p(\mathbf{y}|\mathbf{z})m_{z \rightarrow p}^{t-1}(\mathbf{z})\right)}{m_{z \rightarrow p}^{t-1}(\mathbf{z})} \propto N\left(\mathbf{z}; z_B^{ext}(t), v_B^{ext}(t)I\right)$$

A Unified Inference Framework for GLM

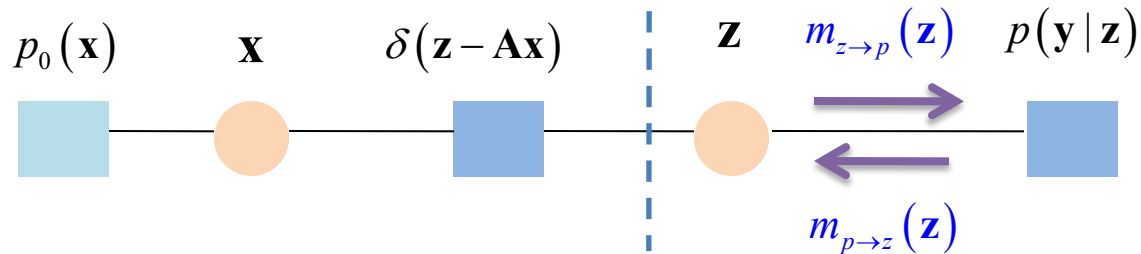
□ Two Equivalent Factor Graphs of GLM



(a) factor graph of GLM

(b) Equivalent factor graph of GLM

□ Decoupling GLM into SLM via EP



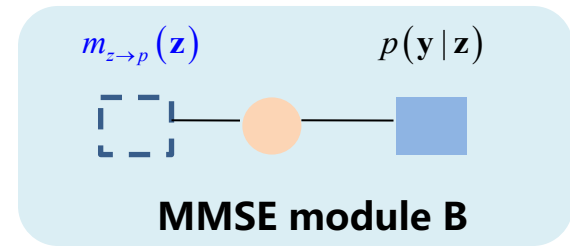
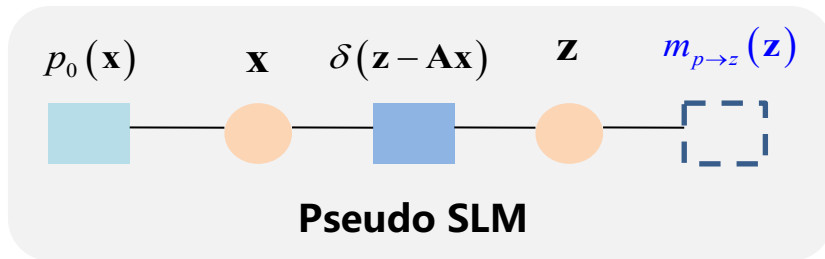
$$m_{z \rightarrow p}^{t-1}(\mathbf{z}) \propto N\left(\mathbf{z}; \mathbf{z}_A^{ext}(t-1), \mathbf{v}_A^{ext}(t-1)I\right)$$

EP message passing
(t -th iteration)

$$m_{p \rightarrow z}^t(\mathbf{z}) \propto \frac{\text{Proj}_{\Phi}\left(p(\mathbf{y}|\mathbf{z})m_{z \rightarrow p}^{t-1}(\mathbf{z})\right)}{m_{z \rightarrow p}^{t-1}(\mathbf{z})} \propto N\left(\mathbf{z}; \mathbf{z}_B^{ext}(t), \mathbf{v}_B^{ext}(t)I\right)$$

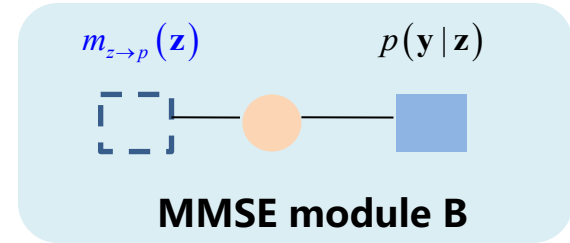
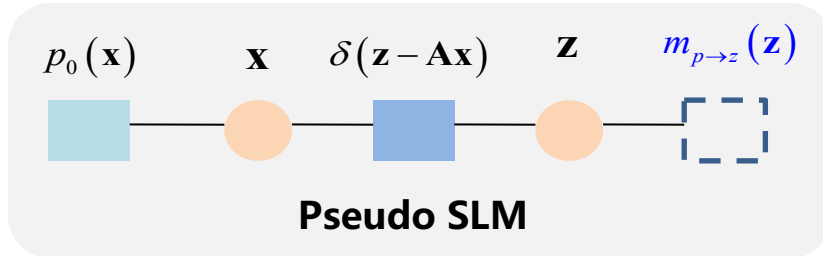
A Unified Inference Framework for GLM

□ Decoupling GLM into SLM via EP

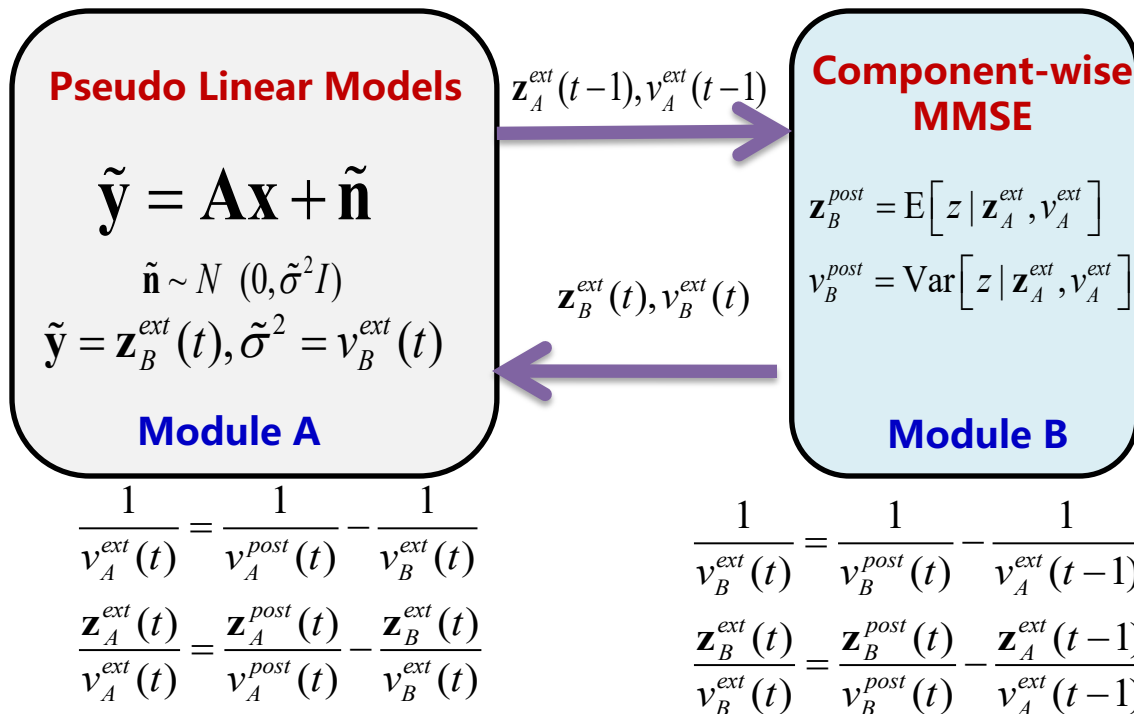


A Unified Inference Framework for GLM

□ Decoupling GLM into SLM via EP



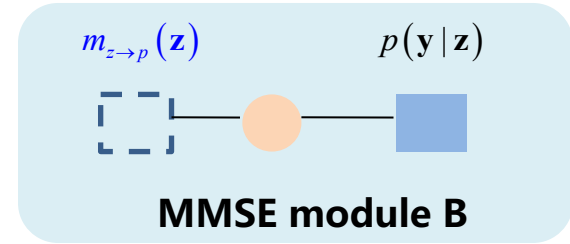
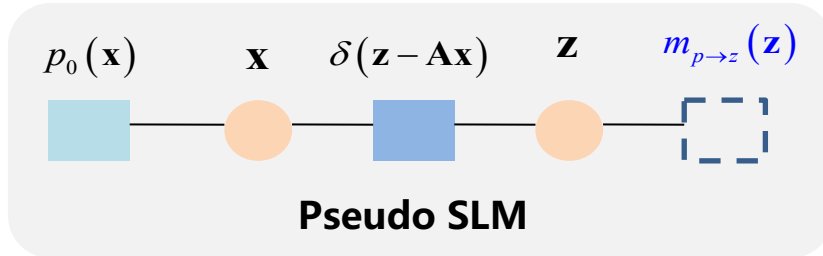
- The original GLM is **iteratively decoupled** into a sequence of simple SLM problems



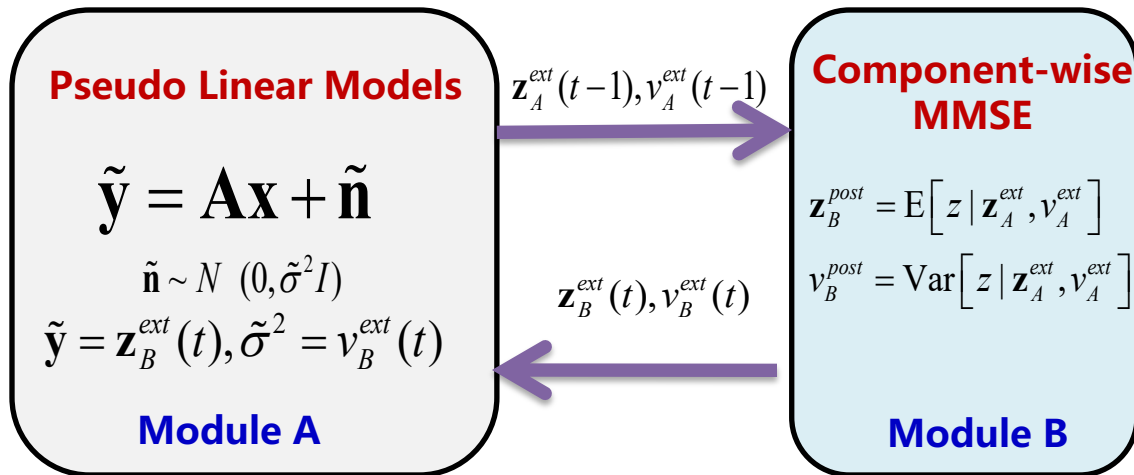
Note: The computation of posterior mean and variance of \mathbf{z} in module A may differ for different SLM inference methods.

A Unified Inference Framework for GLM

□ Decoupling GLM into SLM via EP



- The original GLM is **iteratively decoupled into a sequence of simple SLM problems**



Pseudo Linear Models

$$\tilde{\mathbf{y}} = \mathbf{A}\mathbf{x} + \tilde{\mathbf{n}}$$

$$\tilde{\mathbf{n}} \sim N(0, \tilde{\sigma}^2 \mathbf{I})$$

$$\tilde{\mathbf{y}} = \mathbf{z}_B^{ext}(t), \tilde{\sigma}^2 = v_B^{ext}(t)$$

Module A

Component-wise MMSE

$$\mathbf{z}_B^{post} = \mathbb{E}[\mathbf{z} | \mathbf{z}_A^{ext}, v_A^{ext}]$$

$$v_B^{post} = \text{Var}[\mathbf{z} | \mathbf{z}_A^{ext}, v_A^{ext}]$$

Module B

Universal Algorithm Design [MWZ18]

Unified Inference Framework for GLM

- Initialization $\mathbf{z}_A^{ext}(0), v_A^{ext}(0)$
- For $t = 1: T$, Do
 1. Perform component-wise MMSE
 2. Update $\mathbf{z}_B^{ext}(t), v_B^{ext}(t)$
 3. Perform **SLM inference** one or more iterations
 4. Compute $\mathbf{z}_A^{post}(t), v_A^{post}(t)$ and then update $\mathbf{z}_A^{ext}(t), v_A^{ext}(t)$

$$\frac{1}{v_A^{ext}(t)} = \frac{1}{v_A^{post}(t)} - \frac{1}{v_B^{ext}(t)}$$

$$\frac{\mathbf{z}_A^{ext}(t)}{v_A^{ext}(t)} = \frac{\mathbf{z}_A^{post}(t)}{v_A^{post}(t)} - \frac{\mathbf{z}_B^{ext}(t)}{v_B^{ext}(t)}$$

$$\frac{1}{v_B^{ext}(t)} = \frac{1}{v_B^{post}(t)} - \frac{1}{v_A^{ext}(t-1)}$$

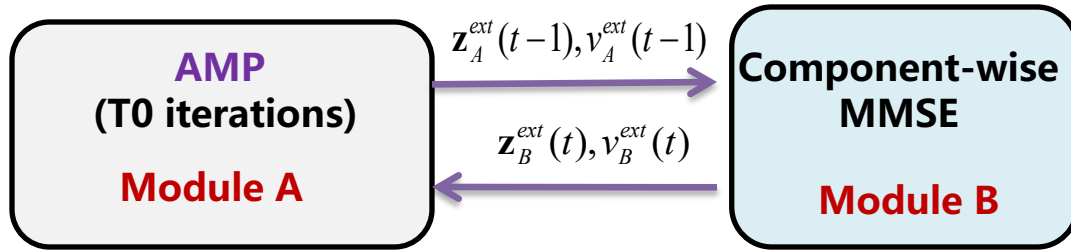
$$\frac{\mathbf{z}_B^{ext}(t)}{v_B^{ext}(t)} = \frac{\mathbf{z}_B^{post}(t)}{v_B^{post}(t)} - \frac{\mathbf{z}_A^{ext}(t-1)}{v_A^{ext}(t-1)}$$

Note: The computation of posterior mean and variance of \mathbf{z} in module A may differ for different SLM inference methods.

[MWZ18] X. Meng, S. Wu and J. Zhu, "A unified Bayesian inference framework for generalized linear model," IEEE Signal Processing Letters, vol. 25, no. 3, Mar. 2018.

A Unified Inference Framework for GLM

□ From AMP to Gr-AMP

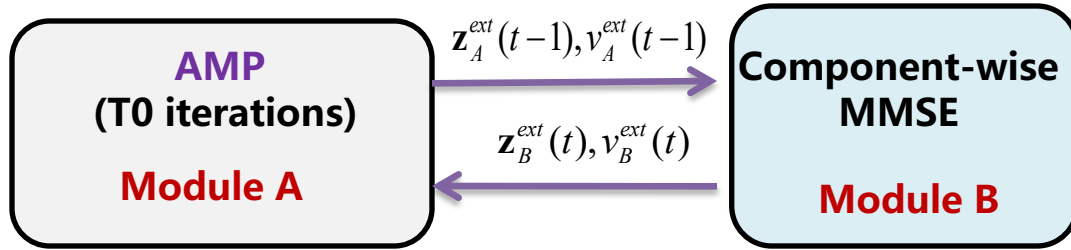


The Gr-AMP Algorithm

- Initialization $\mathbf{z}_A^{ext}(0), v_A^{ext}(0)$
- For $t = 1: T$, Do
 1. Perform component-wise MMSE
 2. Update $\mathbf{z}_B^{ext}(t), v_B^{ext}(t)$
 3. Perform AMP for T0 iterations
 4. Compute $\mathbf{z}_A^{post}(t), v_A^{post}(t)$ and then update $\mathbf{z}_A^{ext}(t), v_A^{ext}(t)$

A Unified Inference Framework for GLM

□ From AMP to Gr-AMP



The Gr-AMP Algorithm

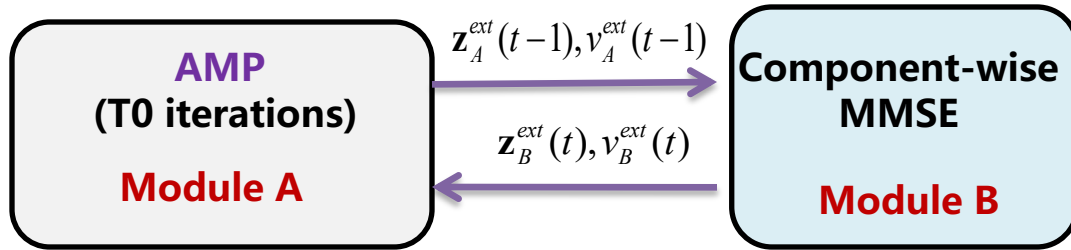
- Initialization $\mathbf{z}_A^{ext}(0), v_A^{ext}(0)$
- For $t = 1: T$, Do
 1. Perform component-wise MMSE
 2. Update $\mathbf{z}_B^{ext}(t), v_B^{ext}(t)$
 3. **Perform AMP for T₀ iterations**
 4. Compute $\mathbf{z}_A^{post}(t), v_A^{post}(t)$ and then update $\mathbf{z}_A^{ext}(t), v_A^{ext}(t)$

• Relation of Gr-AMP to GAMP

- ✓ Gr-AMP is precisely GAMP when $T_0 = 1$ and thus **provides an EP perspective on GAMP** [MWZ18]
In essence, GAMP first transforms nonlinear model to linear model using EP and then directly apply AMP on the linear model in each iteration.
- ✓ This perspective provides **a concise derivation of GAMP using EP as in** [MWZ18]
- ✓ A more flexible message passing schedule: double-loop implementation.

A Unified Inference Framework for GLM

□ From AMP to Gr-AMP

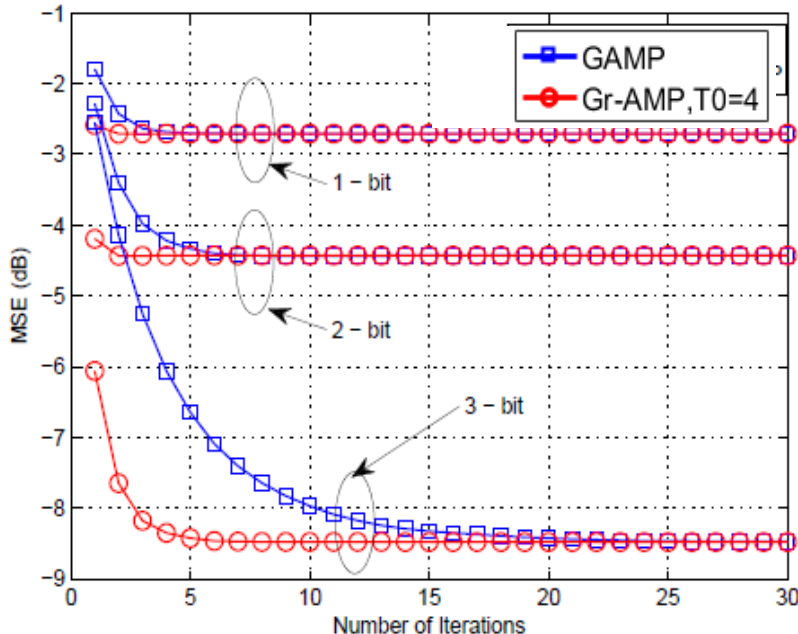


The Gr-AMP Algorithm

- Initialization $\mathbf{z}_A^{ext}(0), v_A^{ext}(0)$
- For $t = 1: T$, Do
 1. Perform component-wise MMSE
 2. Update $\mathbf{z}_B^{ext}(t), v_B^{ext}(t)$
 3. Perform AMP for T_0 iterations
 4. Compute $\mathbf{z}_A^{post}(t), v_A^{post}(t)$ and then update $\mathbf{z}_A^{ext}(t), v_A^{ext}(t)$

• Relation of Gr-AMP to GAMP

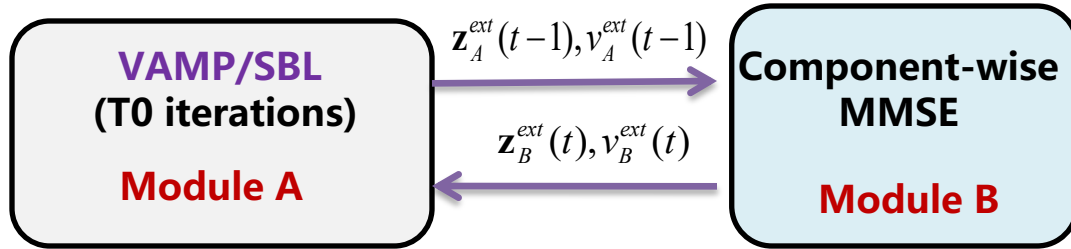
- ✓ Gr-AMP is precisely GAMP when $T_0 = 1$ and thus **provides an EP perspective on GAMP** [MWZ18]. In essence, GAMP first transforms nonlinear model to linear model using EP and then directly apply AMP on the linear model in each iteration.
- ✓ This perspective provides **a concise derivation of GAMP using EP as in** [MWZ18]
- ✓ A more flexible message passing schedule: double-loop implementation.



- Quantized CS for 1,2,3-bit cases: $N=1024, M=512, \text{SNR}=50\text{dB}$
- Gr-AMP and GAMP converge to the same performance for i.i.d. Gaussian A
- Total number iterations of AMP are about the same while **the number of MMSE operations is reduced** for Gr-AMP.

A Unified Inference Framework for GLM

□ From VAMP/SBL to Gr-AMP/Gr-SBL

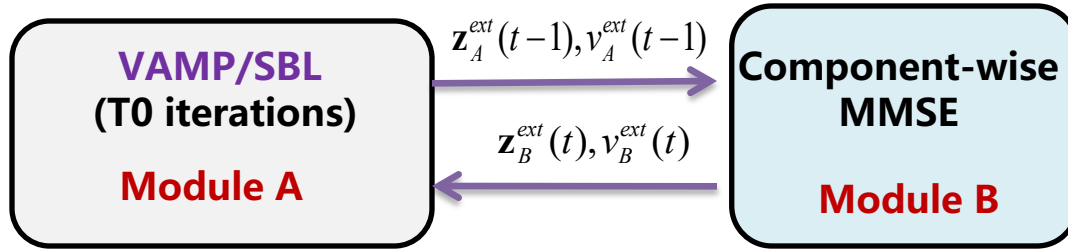


The Gr-VAMP/Gr-SBL Algorithm

- Initialization $\mathbf{z}_A^{ext}(0), v_A^{ext}(0)$
- For $t = 1: T$, Do
 1. Perform component-wise MMSE
 2. Update $\mathbf{z}_B^{ext}(t), v_B^{ext}(t)$
 3. Perform VAMP/SBL for T_0 iterations
 4. Compute $\mathbf{z}_A^{post}(t), v_A^{post}(t)$ and then update $\mathbf{z}_A^{ext}(t), v_A^{ext}(t)$

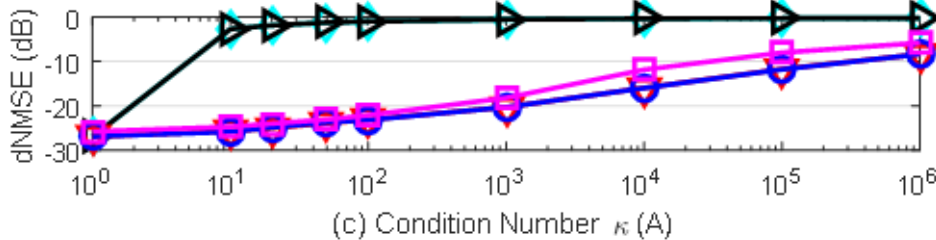
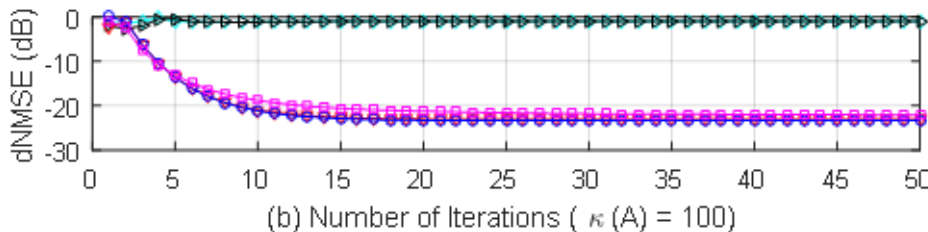
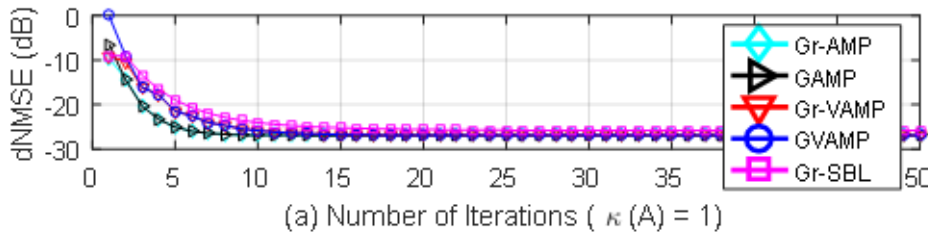
A Unified Inference Framework for GLM

From VAMP/SBL to Gr-AMP/Gr-SBL



The Gr-VAMP/Gr-SBL Algorithm

- Initialization $\mathbf{z}_A^{ext}(0), \mathbf{v}_A^{ext}(0)$
- For $t = 1: T$, Do
 1. Perform component-wise MMSE
 2. Update $\mathbf{z}_B^{ext}(t), \mathbf{v}_B^{ext}(t)$
 3. Perform VAMP/SBL for T_0 iterations
 4. Compute $\mathbf{z}_A^{post}(t), \mathbf{v}_A^{post}(t)$ and then update $\mathbf{z}_A^{ext}(t), \mathbf{v}_A^{ext}(t)$



Performance of de-biased NMSE for **1-bit CS**

- ✓ $N = 512, M = 2048, \text{SNR} = 50\text{dB}$, sparse ratio 0.1
- ✓ $T_0 = 1$ for both Gr-VAMP and Gr-SBL
- ✓ When conditional number is 1, all kinds of algorithms performs nearly the same.
- ✓ As the condition number increases, the recovery performances degrade smoothly for Gr-VAMP/GVAMP/Gr-SBL while both Gr-AMP and GAMP diverge for even mild condition number, which show the robustness of Gr-VAMP/Gr-SBL/GVAMP for general matrices.

X. Meng, S. Wu and J. Zhu, "A unified Bayesian inference framework for generalized linear model," IEEE Signal Processing Letters., vol. 25, no. 3, Mar. 2018.

Conclusions

- **A high-bias low-variance introduction to approximate Bayesian inference**
- **An overview of variational inference framework**
- **A tutorial introduction of expectation propagation**
- **A unified EP perspective on AMP and its extensions.**

References

- [DMM09] Donoho, Maleki, Montanari. "Message-passing algorithms for compressed sensing." Proceedings of the National Academy of Sciences 106.45 (2009): 18914-18919.
- [DMM10] Donoho, Maleki, Montanari, "Message passing algorithms for compressed sensing: I. motivation and construction." Proc IEEE ITW, 20110
- [BM11] Bayati, Montanari. "The dynamics of message passing on dense graphs, with applications to compressed sensing." IEEE Transactions on Information Theory 57.2 (2011): 764-785.
- [Rangan10] Rangan, Sundeep. "Estimation with random linear mixing, belief propagation and compressed sensing." Information Sciences and Systems (CISS), 2010 44th Annual Conference on. IEEE, 2010.
- [Tanaka 02] Tanaka, Toshiyuki. "A statistical-mechanics approach to large-system analysis of CDMA multiuser detectors." IEEE Transactions on Information theory 48.11 (2002): 2888-2910.
- [Minka01] Minka, Thomas P. "Expectation propagation for approximate Bayesian inference." 2001
- [Minka01b] Minka, Thomas P. "A family of algorithms for approximate bayesian inference", (thesis) 2001.
- [Opper05] Opper, Manfred, and Ole Winther. "Expectation consistent approximate inference." Journal of Machine Learning Research 6.Dec (2005): 2177-2204.
- [Rangan11] Rangan, "Generalized approximate message passing for estimation with random linear mixing." Proc IEEE ISIT 2011
- [RSF16] Rangan, Schniter, Fletcher, "Vector approximate message passing", 2016
- [SRF16] P. Schniter, S. Rangan, and A. K. Fletcher, "Vector approximate message passing for the generalized linear model," in Proc. 50th Asilomar Conf. Signals, Syst. Comput., Nov. 2016, pp. 1525-1529.
- [Kabashima 03] Kabashima Y. A , " CDMA multiuser detection algorithm on the basis of belief propagation" , Journal of Physics A: Mathematical and General, 2003, 36(43)
- [KMSSZ12] F. Krzakala, M. Mézard, F. Sausset, Y. Sun, and L. Zdeborová, "Probabilistic reconstruction in compressed sensing: algorithms, phase diagrams, and threshold achieving matrices," Journal of Statistical Mechanics: Theory and Experiment, vol. 2012, no. 08, p. P08009, 2012.
- [MWKL15a] X. Meng, S. Wu, L. Kuang, and J. Lu, "An expectation propagation perspective on approximate message passing," IEEE Signal Process. Lett., vol. 22, no. 8, pp. 1194-1197, Aug. 2015.
- [MWKL15b] X. Meng, S. Wu, L. Kuang, and J. Lu, " Concise derivation of complex Bayesian approximate message passing via expectation propagation," arXiv preprint arXiv:1509.08658, 2015.
- [WKNLHDQ14] S. Wu, L. Kuang, Z. Ni, J. Lu, D. Huang, and Q. Guo, "Low-complexity iterative detection for large-scale multiuser MIMO-OFDM systems using approximate message passing," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 902–915, Oct. 2014.

References

- [MWZ18] X. Meng, S. Wu and J. Zhu, "A unified Bayesian inference framework for generalized linear model," IEEE Signal Process. Lett., vol. 25, no. 3, Mar. 2018.
- [MZ18] X. Meng, and J. Zhu, "Bilinear Adaptive Generalized Vector Approximate Message Passing," arXiv preprint arXiv:1810.08129, 2018
- [PSC14a] J. T. Parker, P. Schniter, and V. Cevher, "Bilinear generalized approximate message passing-Part I: Derivation," IEEE Trans. Signal Process., vol. 62, no. 22, pp. 5839-5853, Nov. 2014.
- [PSC14b] J. T. Parker, P. Schniter, and V. Cevher, "Bilinear generalized approximate message passing-Part II: Applications," IEEE Trans. Signal Process., vol. 62, no. 22, pp. 5854-5867, Nov. 2014.
- [PS16] J. T. Parker and P. Schniter, "Parametric bilinear generalized approximate message passing," IEEE J. Sel. Topics Signal Process., vol. 10, no. 4, pp. 795-808, 2016.
- [KKMSZ16] Y. Kabashima, F. Krzakala, M. Mézard, A. Sakata and L. Zdeborov' a, "Phase transitions and sample complexity in bayesoptimal matrix factorization," IEEE Trans. Inf. Theory, vol. 62, no. 7, pp. 4228-4265, 2016.
- [SRF] P. Schniter, S. Rangan, and A. K. Fletcher, "Vector approximate message passing for the generalized linear model," in Proc. 50th Asilomar Conf. Signals, Syst. Comput., Nov. 2016, pp. 1525-1529.
- [ML17] J. Ma and L. Ping, "Orthogonal AMP," IEEE Access, vol. 5, pp. 2020-2033, 2017.
- [HWJ17] H. He, C. K. Wen, and S. Jin, "Generalized expectation consistent signal recovery for nonlinear measurements," in Proc. IEEE Int. Symp. Inf. Theory, Jun. 2017, pp. 2333-2337.
- [Yedidia et al 05] Yedidia J S, Freeman W T, Weiss Y. Constructing free-energy approximations and generalized belief propagation algorithms[J]. IEEE Transactions on information theory, 2005, 51(7): 2282-2312.
- [Yedidia et al 02] J. S. Yedidia, W. T. Freeman, and Y. Weiss. Understanding belief propagation and its generalizations. Technical Report TR2001-22, Mitsubishi Electric Research Labs, January 2002.
- [Opper&Saad 01] M. Opper and D. Saad. Adaptive TAP equations. In M. Opper and D. Saad, editors, Advanced mean field methods: Theory and practice, pages 85–98. MIT Press, 2001.
- [Hoffman et al 2013] Hoffman, Matthew D., et al. "Stochastic variational inference." The Journal of Machine Learning Research 14.1 (2013): 1303-1347.
- [Mézard et al 87] M. Mézard, G. Parisi and M. A. Virasoro, Spin Glass Theory and Beyond (World Scientific, Singapore, 1987).
- [Bishop06] Bishop C M. Pattern recognition and machine learning[M]. springer, 2006.
- [Mehta et al 19] Mehta P, Bukov M, Wang C H, et al. A high-bias, low-variance introduction to machine learning for physicists[J]. Physics reports, 2019.

Thank You

ありがとうございます

Q&A