

# Information Geometric view of Belief Propagation

*Yunshu Liu*

2013-10-17

## References:

- [1]. Shiro Ikeda, Toshiyuki Tanaka and Shun-ichi Amari, *Stochastic reasoning, Free energy and Information Geometry*, Neural Computation, 16, 1779-1810, 2004.
- [2]. Shiro Ikeda, Toshiyuki Tanaka and Shun-ichi Amari, *Information Geometry of Turbo and Low-Density Parity-Check Codes*, IEEE Transaction on Information Theory, VOL. 50, NO.6, June 2004.
- [3]. Kevin P. Murphy, *Machine Learning - A Probabilistic Perspective*, The MIT Press, 2012.

# Outline

## Information Geometry

- Examples of probabilistic manifold
- Important Submanifold
- Information projection

## Information Geometry of Belief Propagation

- Belief Propagation/Sum-product algorithms
- Information Geometry of Graphical structure
- Information-Geometrical view of Belief Propagation

# Probability distribution viewed as manifold

## Outline for Information Geometry

Examples of probabilistic manifold

Important Submanifold

Information projection

# Manifold

## Manifold $S$

**Manifold:** a set with a coordinate system, a one-to-one mapping from  $S$  to  $\mathbb{R}^n$ , supposed to be "locally" looks like an open subset of  $\mathbb{R}^n$ "

**Elements of the set(points):** points in  $\mathbb{R}^n$ , probability distribution, linear system.

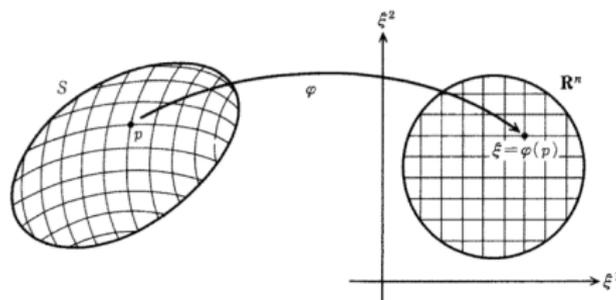


Figure: A coordinate system  $\xi$  for a manifold  $S$

# Example manifold

## Parametrization of color models

Parametrization: map of color into  $\mathbb{R}^3$

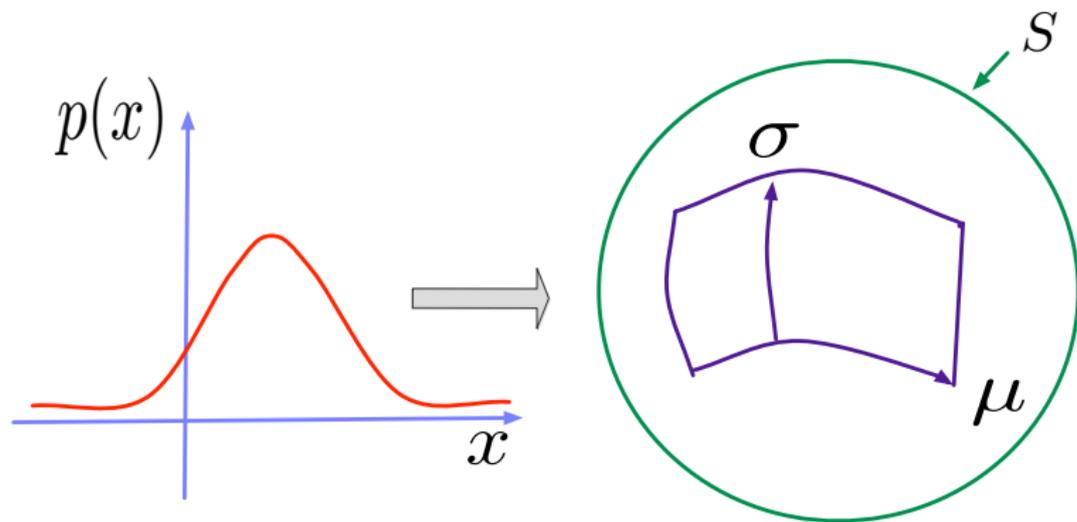
Examples:

<b>RGB 0÷255</b>	<b>RGB 0÷FF</b>	<b>RGB 0÷1</b>	<b>XYZ</b>	<b>CMY 0÷1</b>	<b>CMYK %</b>
32.00 R	20 R	0.12549 R	39.300 X	0.87451 C	87.451 C
200.00 G	C8 G	0.78431 G	48.836 Y	0.21569 M	21.569 M
255.00 B	FF B	1.00000 B	101.963 Z	0.00000 Y	0.000 Y
					0.000 K
<b>CIE-L*ab</b>	<b>CIE-L*CH</b>	<b>CIE-L*uv</b>	<b>Yxy (Y=LRV)</b>	<b>Hunter-Lab</b>	
75.349 L*	75.349 L*	75.349 L*	48.836 Y	69.882 L	
-21.250 a*	43.688 C*	-50.913 u*	0.20673 x	-21.911 a	
-38.172 b*	240.896 H°	-59.274 v*	0.25690 y	-37.590 b	
<b>HTML</b>					
#20C8FF					
<b>Web-Safe</b>					
#33CCFF					
	→ Get commercial tints				
	→ Get color harmonies				

# Example of manifold

## Parametrization of Gaussian distribution

$$S = \{p(x; \mu, \sigma)\} \quad p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$



# Example of manifold

## Exponential family

$$p(x; \theta) = \exp\left[C(x) + \sum_{i=1}^n \theta^i F_i(x) - \psi(\theta)\right]$$

$[\theta^i]$  are called the natural parameters (coordinates), and  $\psi$  is the potential function for  $[\theta^i]$ , which can be calculated as

$$\psi(\theta) = \log \int \exp\left[C(x) + \sum_{i=1}^n \theta^i F_i(x)\right] dx$$

The exponential families include many of the most common distributions, including the normal, exponential, gamma, beta, Dirichlet, Bernoulli, binomial, multinomial, Poisson, and so on.

## Example of two binary discrete random variables as Exponential families

Let  $x = \{x_1, x_2\}$ ,  $x_i = 1$  or  $-1$ ,  
consider the family of all the  
probability distributions  $S =$   
 $\{p(x) \mid p(x) > 0, \sum_i p(x = C_i) = 1\}$

Outcomes:

$C_0 = \{-1, -1\}$ ,  $C_1 = \{-1, 1\}$

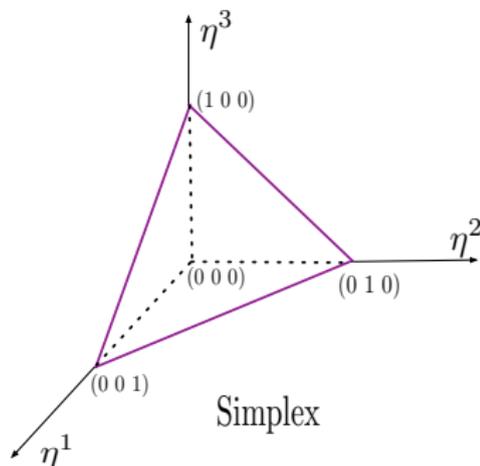
$C_2 = \{1, -1\}$ ,  $C_3 = \{1, 1\}$

Parameterization:

$p(x = C_i) = \eta^i$  for  $i = 1, 2, 3$

$p(x = C_0) = 1 - \sum_{j=1}^3 \eta^j$

$\eta = (\eta^1 \ \eta^2 \ \eta^3)$  is called the  
**m-coordinate** of  $S$ .

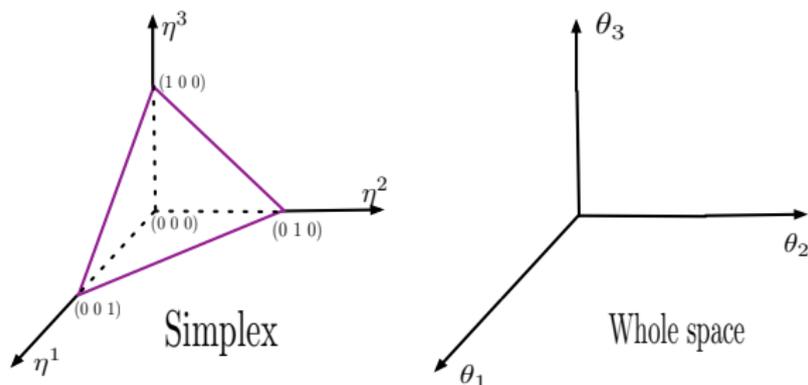


# Example of two binary discrete random variables as Exponential families

A new coordinate  $\{\theta_i\}$ , which is defined as

$$\theta_i = \ln \frac{\eta^i}{1 - \sum_{j=1}^3 \eta^j} \text{ for } i = 1, 2, 3. \quad (1)$$

$\theta = (\theta_1 \theta_2 \theta_3)$  is called the **e-coordinate** of S.



# Example of two binary discrete random variables as Exponential families

Introduce new random variables:

$$\delta_{C_i}(x) = \begin{cases} 1 & \text{if } x = C_i \\ 0 & \text{otherwise} \end{cases}$$

and let

$$\psi = -\ln p(x = C_0) \quad (2)$$

Then  $S$  is an **exponential family** with natural parameter  $\{\theta_i\}$ :

$$p(x, \theta) = \exp\left(\sum_{i=1}^3 \theta_i \delta_{C_i}(x) - \psi(\theta)\right) \quad (3)$$

# Probability distribution viewed as manifold

## Outline for Information Geometry

Examples of probabilistic manifold

**Important Submanifolds**

Information projection

# Submanifolds

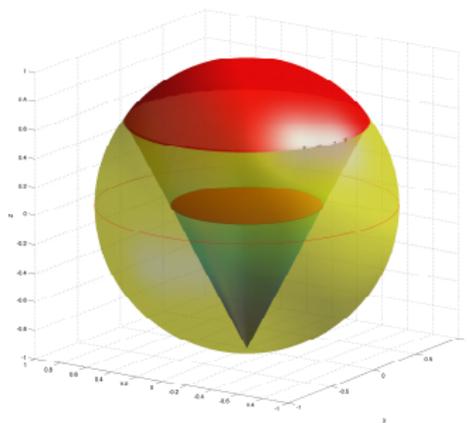
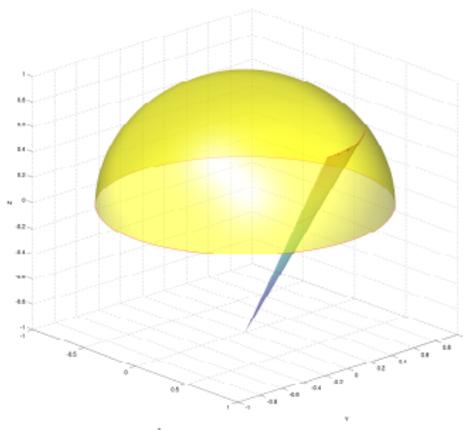
## Submanifolds

Definition: a submanifold  $M$  of a manifold  $S$  is a subset of  $S$  which itself has the structure of a manifold

An open subset of  $n$ -dimensional manifold forms an  $n$ -dimensional submanifold.

One way to construct  $m (< n)$  dimensional manifold: fix  $n-m$  coordinates.

Examples:



# e-autoparallel submanifold and m-autoparallel submanifold

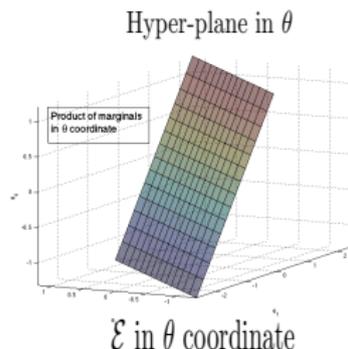
## e-autoparallel submanifold

A subfamily  $E$  is said to be e-autoparallel submanifold of  $S$  if for some coordinate  $\lambda$  there exist a matrix  $A$  and a vector  $b$  such that for the e-coordinate  $\theta$  of  $S$

$$\theta(p) = \mathbf{A}\lambda(p) + \mathbf{b} \quad (\forall p \in E) \quad (4)$$

If  $\lambda$  is one-dimensional, it is called a *e-geodesic*.

An e-autoparallel submanifold of  $S$  is a **hyperplane** in e-coordinate;  
A e-geodesic of  $S$  is a **straight line** in e-coordinate.



# e-autoparallel submanifold and m-autoparallel submanifold

## m-autoparallel submanifold

A subfamily  $M$  is said to be  $m$ -autoparallel submanifold of  $S$  if for some coordinate  $\gamma$  there exist a matrix  $\mathbf{A}$  and a vector  $\mathbf{b}$  such that for the  $m$ -coordinate  $\eta$  of  $S$

$$\eta(p) = \mathbf{A}\gamma(p) + \mathbf{b} \quad (\forall p \in M) \quad (5)$$

If  $\gamma$  is one-dimensional, it is called a *m-geodesic*.

# e-autoparallel submanifold and m-autoparallel submanifold

## Two independent bits

The set of all product distributions, which is defined as

$$\mathcal{E} = \{p(x) \mid p(x) = p_{X_1}(x_1)p_{X_2}(x_2)\} \quad (6)$$

$$p_{X_1}(x_1 = 1) = p_1, p_{X_2}(x_2 = 1) = p_2$$

$$\text{then } p_{X_1}(x_1 = -1) = 1 - p_1, p_{X_2}(x_2 = -1) = 1 - p_2$$

## Parameterization

Define  $\lambda_i = \frac{1}{2} \ln \frac{p_i}{1-p_i}$  for  $i = 1, 2$ , and

$$\psi(\lambda) = \sum_{i=1,2} \ln(e^{\lambda_i} + e^{-\lambda_i}),$$

$$p(x, \lambda) = e^{\lambda_1 x_1 - \psi(\lambda_1)} e^{\lambda_2 x_2 - \psi(\lambda_2)} = e^{\lambda x - \psi(\lambda)} \quad (7)$$

Thus  $\mathcal{E}$  is an exponential family with nature parameter  $\{\lambda_i\}$

# e-autoparallel submanifold and m-autoparallel submanifold

## Dual parameters: expectations

$$\gamma = \begin{pmatrix} \gamma^1 \\ \gamma^2 \end{pmatrix} = \begin{pmatrix} \partial_1 \psi(\lambda) \\ \partial_2 \psi(\lambda) \end{pmatrix} = \begin{pmatrix} 2p_1 - 1 \\ 2p_2 - 1 \end{pmatrix} = E_p(x)$$

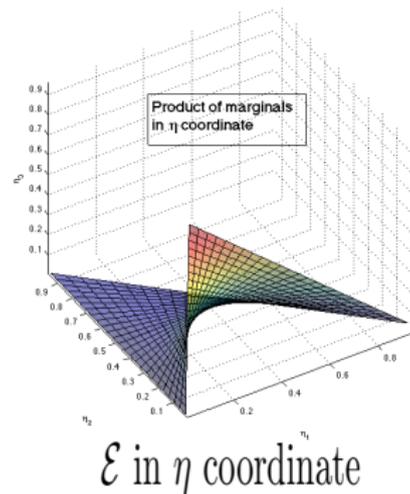
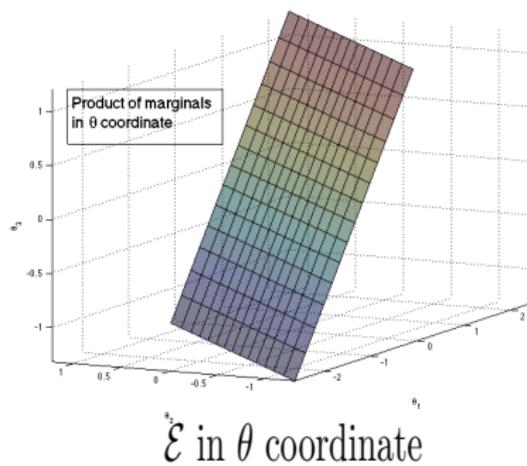
## Properties

$\mathcal{E}$  is an e-autoparallel submanifold(plane in  $\theta$  coordinate) of  $S$ , but not a m-autoparallel submanifold(not a plane in  $\eta$  coordinate):

$$\theta(p) = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{pmatrix} = \begin{pmatrix} 2 & 2 \\ 2 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} \lambda_1 \\ \lambda_2 \end{pmatrix} = \begin{pmatrix} 2 & 2 \\ 2 & 0 \\ 0 & 2 \end{pmatrix} \lambda(p)$$

# e-autoparallel submanifold and m-autoparallel submanifold

Two independent bits: e-autoparallel submanifold but not m-autoparallel submanifold



# Probability distribution viewed as manifold

## Outline for Information Geometry

Examples of probabilistic manifold

Important Submanifolds

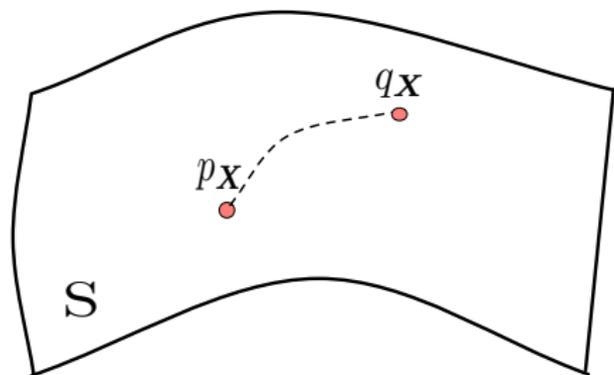
Information projection

# Projections and Pythagorean Theorem

## KL-divergence

On the manifold of probability mass functions for random variables taking values in the set  $\mathcal{X}$ , we can also define the Kullback Leibler divergence, or relative entropy, measured in bits, according to

$$D(p_{\mathbf{X}}||q_{\mathbf{X}}) = \sum_{\mathbf{x} \in \mathcal{X}} p_{\mathbf{X}}(\mathbf{x}) \log \left( \frac{p_{\mathbf{X}}(\mathbf{x})}{q_{\mathbf{X}}(\mathbf{x})} \right) \quad (8)$$



# Projections and Pythagorean Theorem

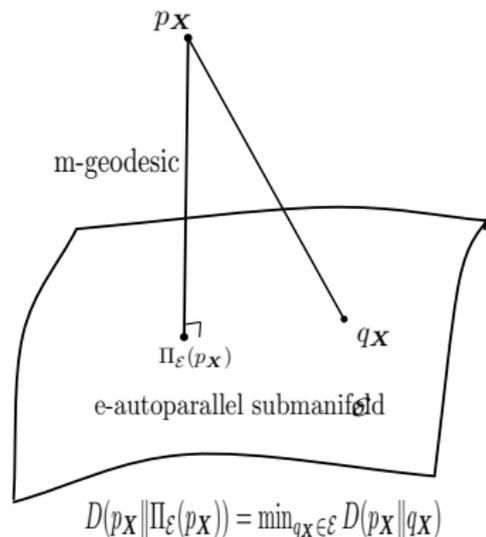
## Information Projection

Let  $p$  be a point in  $S$  and let  $\mathcal{E}$  be a  $e$ -autoparallel submanifold of  $S$ . We want to find point  $\Pi_{\mathcal{E}}(p_X)$  in  $\mathcal{E}$  that minimize the KL-divergence

$$D(p_X \| \Pi_{\mathcal{E}}(p_X)) = \min_{q_X \in \mathcal{E}} D(p_X \| q_X) \quad (9)$$

The  $m$ -geodesic connecting  $p_X$  and  $\Pi_{\mathcal{E}}(p_X)$  is orthogonal to  $\mathcal{E}$  at  $\Pi_{\mathcal{E}}(p_X)$ .

- ▶ In geometry, orthogonal usually means right angle
- ▶ In information geometry, orthogonal usually means certain values are uncorrelated

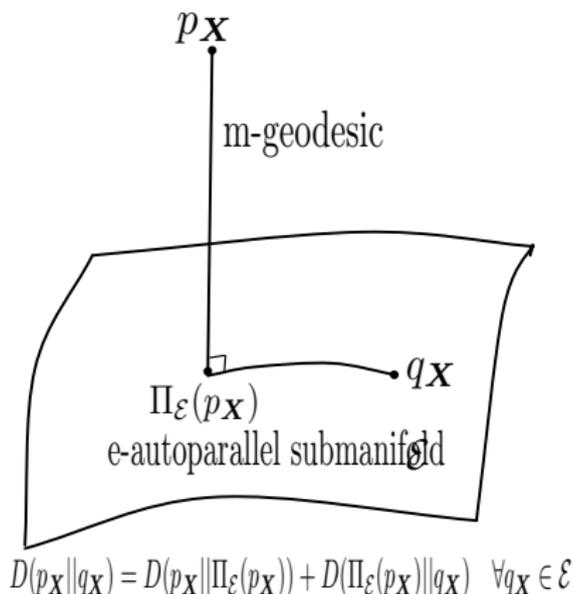


# Projections and Pythagorean Theorem

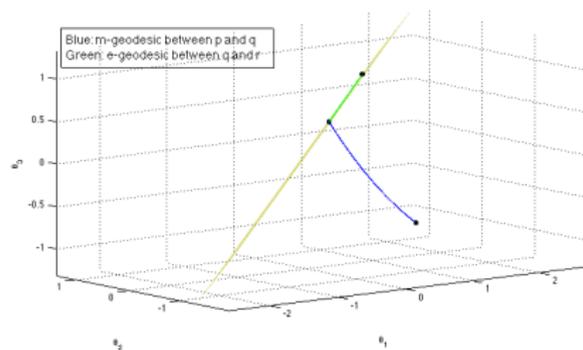
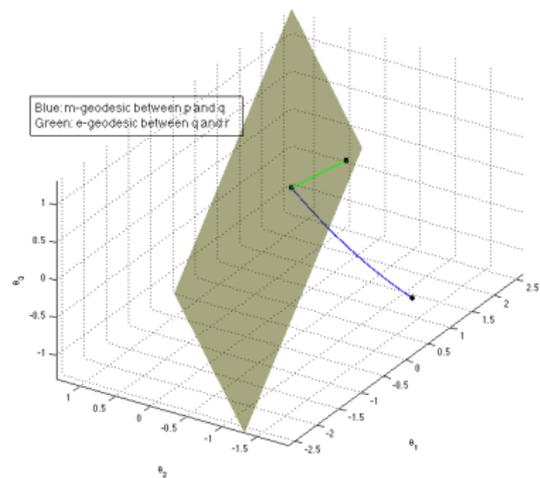
## Pythagorean relation

For  $\forall q_X \in \mathcal{E}$ , we have the following Pythagorean relation.

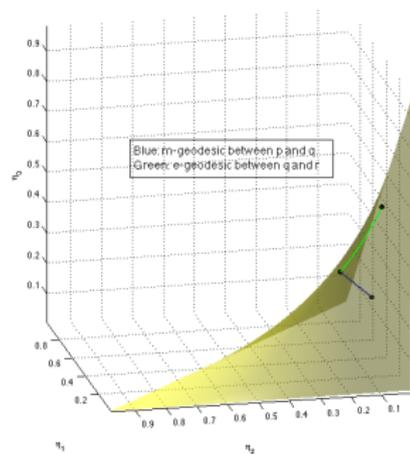
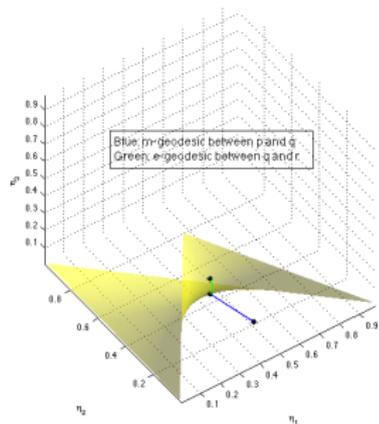
$$D(p_X || q_X) = D(p_X || \Pi_{\mathcal{E}}(p_X)) + D(\Pi_{\mathcal{E}}(p_X) || q_X) \quad \forall q_X \in \mathcal{E} \quad (10)$$



# Projection in $\theta$ coordinate



# Projection in $\eta$ coordinate



# Information Geometry of Belief Propagation

## Outline for Information Geometry of Belief Propagation

Belief Propagation/Sum-product algorithms

Information Geometry of Graphical structure

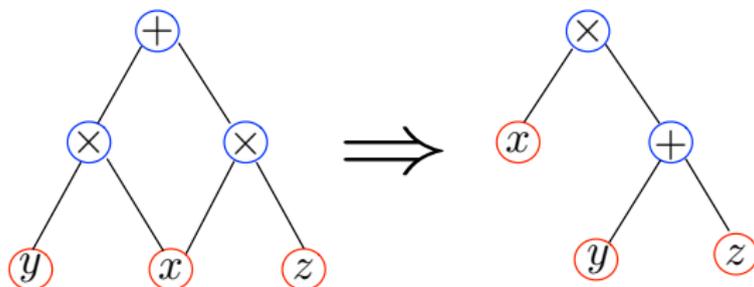
Information-Geometrical view of Belief Propagation

# Belief Propagation

## Belief Propagation/Sum-product algorithm

The Belief Propagation algorithm involves using a message passing scheme to explore the conditional independence and use these properties to change the order of sum and product.

$$xy + xz = x(y + z)$$



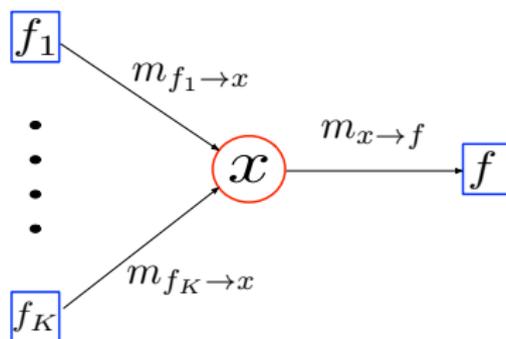
Expression tree of  $xy + xz$  and  $x(y+z)$ , where the leaf nodes represent variables and internal nodes represent arithmetic operators (addition, multiplication, negation, etc.).

# Belief Propagation

## Message passing on a factor graph:

Variable  $x$  to factor  $f$ :

$$m_{x \rightarrow f} = \prod_{k \in nb(x) \setminus f} m_{f_k \rightarrow x}(x)$$



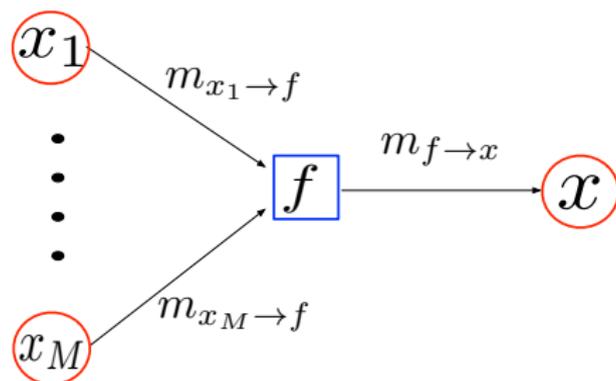
$nb(x)$  denote the neighboring factor nodes of  $x$   
If  $nb(x) \setminus f = \emptyset$  then  $m_{x \rightarrow f} = 1$ .

# Belief Propagation

Message passing on a factor graph:

Factor  $f$  to variable  $x$ :

$$m_{f \rightarrow x} = \sum_{x_1, \dots, x_M} f(x, x_1, \dots, x_M) \prod_{m \in nb(f) \setminus x} m_{x_m \rightarrow f}(x_m)$$



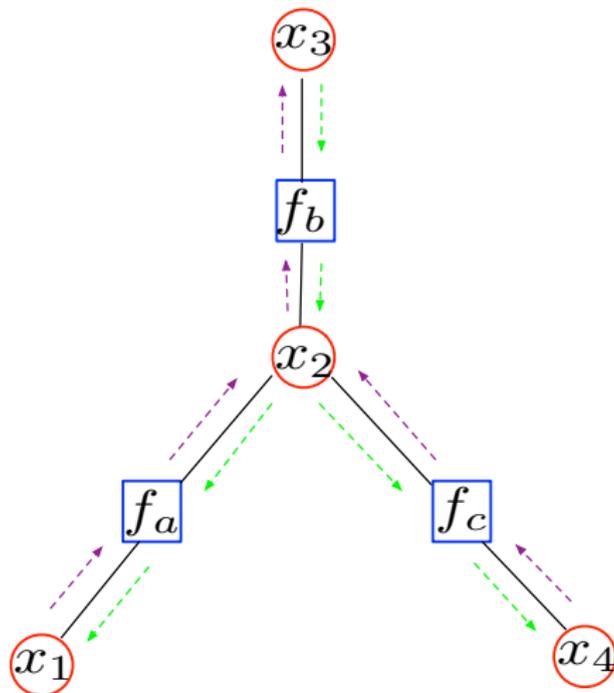
$nb(f)$  denote the neighboring variable nodes of  $f$

If  $nb(f) \setminus x = \emptyset$  then  $m_{f \rightarrow x} = f(x)$ .

# Serial protocol of Belief Propagation

Message passing on a tree structured factor graph

Serial protocol: Send messages up to the root and back



# Serial protocol of Belief Propagation

## Message passing on a tree structured factor graph

Down - up:

$$m_{x_1 \rightarrow f_a}(x_1) = 1$$

$$m_{f_a \rightarrow x_2}(x_2) = \sum_{x_1} f_a(x_1, x_2)$$

$$m_{x_4 \rightarrow f_c}(x_4) = 1$$

$$m_{f_c \rightarrow x_2}(x_2) = \sum_{x_4} f_c(x_2, x_4)$$

$$m_{x_2 \rightarrow f_b}(x_2) = m_{f_a \rightarrow x_2}(x_2) m_{f_c \rightarrow x_2}(x_2)$$

$$m_{f_b \rightarrow x_3}(x_3) = \sum_{x_2} f_b(x_2, x_3) m_{x_2 \rightarrow f_b}$$

# Serial protocol of Belief Propagation

## Message passing on a tree structured factor graph

Up - down:

$$m_{x_3 \rightarrow f_b}(x_3) = 1$$

$$m_{f_b \rightarrow x_2}(x_2) = \sum_{x_3} f_b(x_2, x_3)$$

$$m_{x_2 \rightarrow f_a}(x_2) = m_{f_b \rightarrow x_2}(x_2) m_{f_c \rightarrow x_2}(x_2)$$

$$m_{f_a \rightarrow x_1}(x_1) = \sum_{x_2} f_a(x_1, x_2) m_{x_2 \rightarrow f_a}(x_2)$$

$$m_{x_2 \rightarrow f_c}(x_2) = m_{f_a \rightarrow x_2}(x_2) m_{f_b \rightarrow x_2}(x_2)$$

$$m_{f_c \rightarrow x_4}(x_4) = \sum_{x_2} f_c(x_2, x_4) m_{x_2 \rightarrow f_c}(x_2)$$

# Serial protocol of Belief Propagation

## Message passing on a tree structured factor graph

Calculating  $p(x_2)$

$$\begin{aligned} p(x_2) &= \frac{1}{Z} m_{f_a \rightarrow x_2}(x_2) m_{f_b \rightarrow x_2}(x_2) m_{f_c \rightarrow x_2}(x_2) \\ &= \frac{1}{Z} \left\{ \sum_{x_1} f_a(x_1, x_2) \right\} \left\{ \sum_{x_3} f_b(x_2, x_3) \right\} \left\{ \sum_{x_4} f_c(x_2, x_4) \right\} \\ &= \frac{1}{Z} \sum_{x_1} \sum_{x_3} \sum_{x_4} f_a(x_1, x_2) f_b(x_2, x_3) f_c(x_2, x_4) \\ &= \sum_{x_1} \sum_{x_3} \sum_{x_4} p(\mathbf{x}) \end{aligned}$$

# Parallel protocol of Belief Propagation

## Message passing on general graphs with loops

Parallel protocol:

- ▶ Initialize all messages
- ▶ All nodes receive message from their neighbors in parallel and update their belief states
- ▶ Send new messages back out to their neighbors
- ▶ Repeats the above two process until convergence

# Information Geometry of Belief Propagation

## Outline for Information Geometry of Belief Propagation

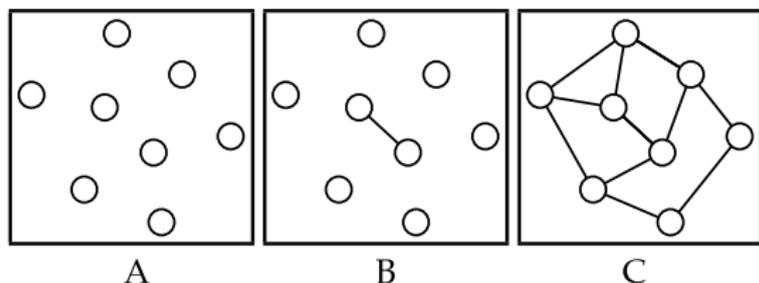
Belief Propagation/Sum-product algorithms

**Information Geometry of Graphical structure**

Information-Geometrical view of Belief Propagation

# Information Geometry of Graphical structure

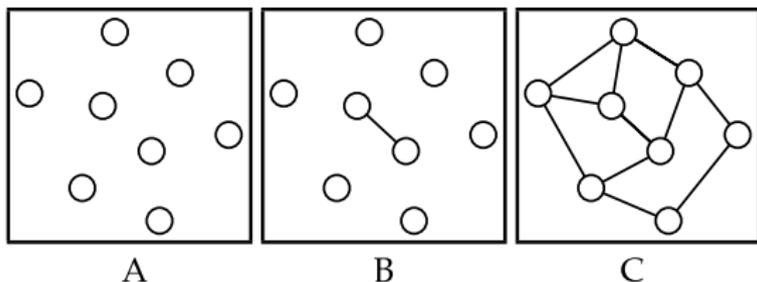
## Important Graphs(Undirected) for Belief Propagation



- A: Belief graph with no edges(all  $n$  nodes are independent)
- B: Graph with a single edge  $r$
- C: Graph of the true distribution with  $n$  nodes and  $L$  edges

# Information Geometry of Graphical structure

## Important Graphs(Undirected) for Belief Propagation



$$A: M_0 = \{p_0(\mathbf{x}, \theta) = \exp[\sum_{i=1}^n h_i x_i + \sum_{i=1}^n \theta_i x_i - \psi_0(\theta)] | \theta \in \mathcal{R}^n\}$$

$$B: M_r = \{p_r(\mathbf{x}, \zeta_r) = \exp[\sum_{i=1}^n h_i x_i + s_r c_r(\mathbf{x}) + \sum_{i=1}^n \zeta_{ri} x_i - \psi_r(\zeta_r)] | \zeta_r \in \mathcal{R}^n\}$$

$$C: S = \{p(\mathbf{x}, \theta, \mathbf{s}) = \exp[\sum_{i=1}^n h_i x_i + \sum_{r=1}^L s_r c_r(\mathbf{x}) - \psi(\theta, \mathbf{s})] | \theta \in \mathcal{R}^n, \mathbf{s} \in \mathcal{R}^L\}$$

# Information Geometry of Graphical structure

## Properties of $M_0, M_r$

- ▶  $M_0$  is an e-autoparallel submanifold of  $S$ , with e-affine coordinate  $\theta$
- ▶  $M_r$  for  $r = 1, \dots, L$  are all e-autoparallel submanifold of  $S$ , with e-affine coordinate  $\zeta_r$
- ▶ A distribution in  $M_r$  include only one nonlinear interaction term  $c_r(\mathbf{x})$  corresponding to the single edge in  $B$

# Information Geometry of Belief Propagation

## Outline for Information Geometry of Belief Propagation

Belief Propagation/Sum-product algorithms

Information Geometry of Graphical structure

Information-Geometrical view of Belief Propagation

# Information-Geometrical view of Belief Propagation

## Problem setup for belief propagation

Suppose we have the true distribution  $q(\mathbf{x}) \in \mathcal{S}$ , and want to calculate  $q_0(\mathbf{x}) = \prod_{M_0} q(\mathbf{x})$ , the m-projection of  $q(\mathbf{x})$  to  $M_0$ , the goal of belief propagation is to find a good approximation  $\hat{q}(\mathbf{x}) \in M_0$  of  $q(\mathbf{x})$ .

$$M_0 = \{p_0(\mathbf{x}, \theta) = \exp[\sum_{i=1}^n h_i x_i + \sum_{i=1}^n \theta_i x_i - \psi_0(\theta)] \mid \theta \in \mathcal{R}^n\}$$

$$M_r = \{p_r(\mathbf{x}, \zeta_r) = \exp[\sum_{i=1}^n h_i x_i + \mathbf{s}_r \mathbf{c}_r(\mathbf{x}) + \sum_{i=1}^n \zeta_{ri} x_i - \psi_r(\zeta_r)] \mid \zeta_r \in \mathcal{R}^n\}$$

$$\mathcal{S} = \{p(\mathbf{x}, \theta, \mathbf{s}) = \exp[\sum_{i=1}^n h_i x_i + \sum_{r=1}^L \mathbf{s}_r \mathbf{c}_r(\mathbf{x}) - \psi(\theta, \mathbf{s})] \mid \theta \in \mathcal{R}^n, \mathbf{s} \in \mathcal{R}^L\}$$

# Geometrical Belief Propagation Algorithm

- ▶ 1) Put  $t = 0$ , and start with initial guesses  $\zeta_r^0$ , for example,  $\zeta_r^0 = 0$
- ▶ 2) For  $t = 0, 1, 2, \dots$ , m-project  $p_r(\mathbf{x}, \zeta_r^t)$  to  $M_0$  and obtain the linearized version of  $s_r c_r(\mathbf{x})$

$$\xi_r^{t+1} = \prod_{M_0} p_r(\mathbf{x}, \zeta_r^t) - \zeta_r^t$$

- ▶ 3) Summarize all the effects of  $s_r c_r(\mathbf{x})$ , to give

$$\theta^{t+1} = \sum_{r=1}^L \xi_r^{t+1}$$

- ▶ 4) Update  $\zeta_r$  by

$$\zeta_r^{t+1} = \sum_{r' \neq r} \xi_{r'}^{t+1} = \theta^{t+1} - \xi_r^{t+1}$$

- ▶ 5) Repeat 2-4 until convergence

# Geometrical Belief Propagation Algorithm

## Analysis

Given  $q(\mathbf{x})$ , let  $\zeta_r$  be the current solution which  $M_r$  believes to give a good approximation of  $q(\mathbf{x})$ . then we project it to  $M_0$ , giving the solution  $\theta_r$  in  $M_0$  having the same expectation:

$$p(\mathbf{x}, \theta_r) = \prod_{M_0} p_r(\mathbf{x}, \zeta_r).$$

Since  $\theta_r$  includes both the effect of the single nonlinear term  $s_r c_r(\mathbf{x})$  specific to  $M_r$  and that due to the linear term  $\zeta_r$ , we can use  $\xi_r = \theta_r - \zeta_r$  to represent the effect of the single nonlinear term  $s_r c_r(\mathbf{x})$ , we consider it as the linearized version of  $s_r c_r(\mathbf{x})$ .

# Geometrical Belief Propagation Algorithm

## Analysis:continue

Once  $M_r$  knows the linearized version of  $s_r c_r(\mathbf{x})$  is  $\xi_r$ , it broadcasts  $\xi_r$  to all the other models  $M_0$  and  $M_{r'} (r' \neq r)$ . Receiving these messages from all  $M_r$ ,  $M_0$  guesses that the equivalent linear term of  $\sum_{r=1}^L s_r c_r(\mathbf{x})$  is  $\theta = \sum_{r=1}^L \xi_r$ .

Since  $M_r$  in turn receives messages  $\xi_{r'}$  from all other  $M_{r'} (r' \neq r)$ ,  $M_r$  uses them to form a new  $\zeta_r = \sum_{r' \neq r} \xi_{r'} = \theta - \xi_r$ . The whole process is repeated until converge.

# Geometrical Belief Propagation Algorithm

## Analysis:continue

Assume the algorithm has converged to  $\{\zeta_r^*\}$  and  $\{\theta^*\}$ . We find the estimation of  $q(\mathbf{x}) \in \mathcal{S}$  in  $M_0$  is  $p_0(\mathbf{x}, \theta^*)$  (which may not be the exact solution). We write  $q(\mathbf{x})$ ,  $p_0(\mathbf{x}, \theta^*)$  and  $p_r(\mathbf{x}, \zeta_r^*)$  as:

$$p_0(\mathbf{x}, \theta^*) = \exp\left[\sum_{i=1}^n h_i x_i + \sum_{r=1}^L \xi_r^* \cdot \mathbf{x} - \psi_0\right]$$

$$p_r(\mathbf{x}, \zeta_r^*) = \exp\left[\sum_{i=1}^n h_i x_i + s_r c_r(\mathbf{x}) + \sum_{r' \neq r} \xi_{r'}^* \cdot \mathbf{x} - \psi_r\right]$$

$$q(\mathbf{x}) = \exp\left[\sum_{i=1}^n h_i x_i + \sum_{r=1}^L s_r c_r(\mathbf{x}) - \psi\right]$$

The whole idea of the algorithm is to approximate  $s_r c_r(\mathbf{x})$  by  $\xi_r^* \cdot \mathbf{x}$  in  $M_r$ , then the independent distribution  $p_0(\mathbf{x}, \theta^*) \in M_0$  integrates all the information.

# Information-Geometrical view of Belief Propagation

## Convergence Analysis

The convergence point satisfies the two conditions:

- ▶ 1. m-condition:  $\theta^* = \prod_{M_0} p_r(\mathbf{x}, \zeta_r^*)$
- ▶ 2. e-condition:  $\theta^* = \frac{1}{L-1} \sum_{r=1}^L \zeta_r^*$

The m-condition and e-condition of the algorithm

- ▶ e-condition is satisfied at each iteration of BP algorithm
- ▶ m-condition is satisfied at the equilibrium of the BP algorithm

# Information-Geometrical view of Belief Propagation

## Convergence Analysis

Let us consider the  $m$ -autoparallel submanifold  $M^*$  connecting  $p_0(\mathbf{x}, \theta^*)$  and  $p_r(\mathbf{x}, \zeta_r^*)$ ,  $r = 1, \dots, L$ :

$$M^* = \{p(\mathbf{x}) | p(\mathbf{x}) = d_0 p_0(\mathbf{x}, \theta^*) + \sum_{r=1}^L d_r p_r(\mathbf{x}, \zeta_r^*); d_0 + \sum_{r=1}^L d_r = 1\}$$

We also consider the  $e$ -autoparallel submanifold  $E^*$  connecting  $p_0(\mathbf{x}, \theta^*)$  and  $p_r(\mathbf{x}, \zeta_r^*)$ :

$$E^* = \{p(\mathbf{x}) | \log p(\mathbf{x}) = v_0 \log p_0(\mathbf{x}, \theta^*) + \sum_{r=1}^L v_r \log p_r(\mathbf{x}, \zeta_r^*) - \psi; \\ v_0 + \sum_{r=1}^L v_r = 1\}$$

# Information-Geometrical view of Belief Propagation

## Convergence Analysis

$M^*$  hold the expectation of  $\mathbf{x}$  unchange, thus a equivalent definition is given by:

$$M^* = \{p(\mathbf{x}) | p(\mathbf{x}) \in \mathcal{S}, \sum_{\mathbf{x}} \mathbf{x}p(\mathbf{x}) = \sum_{\mathbf{x}} \mathbf{x}p_0(\mathbf{x}, \theta^*) = \eta_0(\theta^*)\}$$

Thus m-condition implies that  $M^*$  intersects all of  $M_0$  and  $M_r$  orthogonally.

The e-condition implies that  $E^*$  include the true distribution  $q(\mathbf{x})$  since if we set  $v_0 = -(L-1)$ ,  $v_r = 1$  for  $\forall r$ , using e-condition, we will get

$$q(\mathbf{x}) = Cp_0(\mathbf{x}, \theta^*)^{-(L-1)} \prod_{r=1}^L p_r(\mathbf{x}, \zeta_r^*) \in E^*$$

where  $C$  is the normalization factor

# Information-Geometrical view of Belief Propagation

## Convergence Analysis

The proposed BP algorithm search for  $\theta^*$  and  $\zeta^*$  until both the m-condition and e-condition are satisfied. We know  $E^*$  includes  $q(\mathbf{x})$ , if  $M^*$  include  $q(\mathbf{x})$ , its m-projection to  $M_0$  is  $p_0(\mathbf{x}, \theta^*)$ ,  $\theta^*$  gives the true solution; if  $M^*$  doesn't include  $q(\mathbf{x})$ , the algorithm only give a approximation of the true distribution. Thus we have the following theorem:

**Theorem:** When the underlying graph is a tree, both  $E^*$  and  $M^*$  include  $q(\mathbf{x})$ , and the algorithm gives the exact solution.

# Information-Geometrical view of Belief Propagation

Example of  $L = 2$  where  $\theta = \zeta_1 + \zeta_2 = \xi_1 + \xi_2$ ,  $\zeta_1 = \xi_2$  and  $\zeta_2 = \xi_1$

- ▶ 1) Put  $t = 0$ , and start with initial guesses  $\zeta_2^0$ , for example,  $\zeta_2^0 = 0$
- ▶ 2) For  $t = 0, 1, 2, \dots$ , m-project  $p_r(\mathbf{x}, \zeta_r^t)$  to  $M_0$  and obtain the linearized version of  $s_r c_r(\mathbf{x})$

$$\zeta_1^{t+1} = \prod_{M_0} p_2(\mathbf{x}, \zeta_2^t) - \zeta_2^t, \quad \zeta_2^{t+1} = \prod_{M_0} p_1(\mathbf{x}, \zeta_1^{t+1}) - \zeta_1^{t+1}$$

- ▶ 3) Summarize all the effects of  $s_r c_r(\mathbf{x})$ , to give

$$\theta^{t+1} = \zeta_1^{t+1} + \zeta_2^{t+1}$$

- ▶ 4) Repeat 2)-3) until convergence, and we get  $\theta^* = \zeta_1^* + \zeta_2^*$

# Information-Geometrical view of Belief Propagation

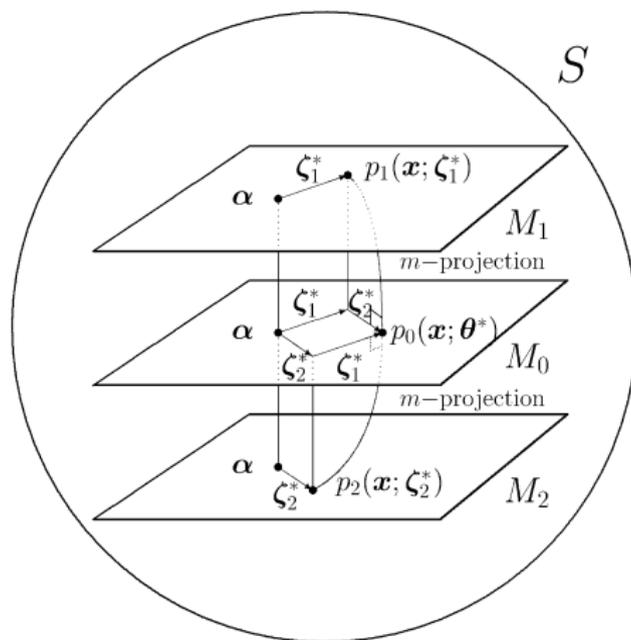


Figure: Information geometrical view of the BP algorithm with  $L = 2$

# Information-Geometrical view of Belief Propagation

## Analysis of the example with $L = 2$

The convergence point satisfies the two conditions:

- ▶ 1. m-condition:  $\theta^* = \prod_{M_0} p_1(\mathbf{x}, \zeta_1^*) = \prod_{M_0} p_2(\mathbf{x}, \zeta_2^*)$
- ▶ 2. e-condition:  $\theta^* = \zeta_1^* + \zeta_2^* = \xi_1^* + \xi_2^*$

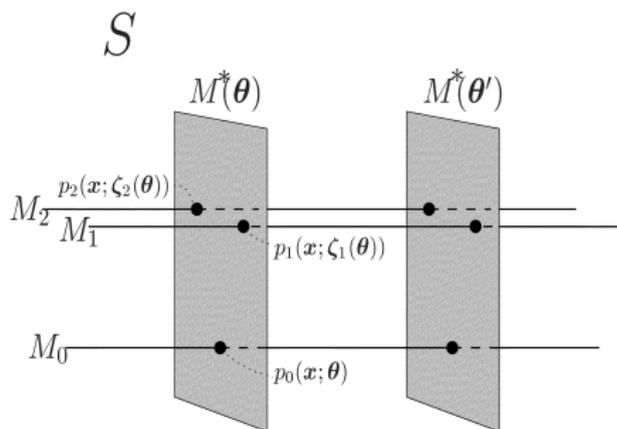
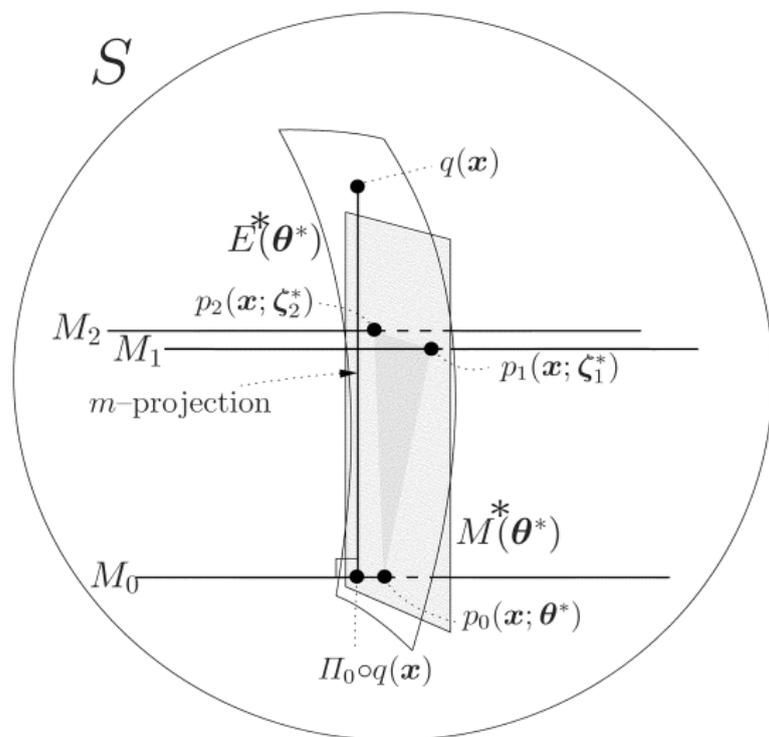


Figure: Equimarinial submanifold  $M^*$

# Information-Geometrical view of Belief Propagation



# Information-Geometrical view of CCCP-Bethe

## Information Geometric Convex Concave Computational Procedure(CCCP-Bethe)

**Inner loop:** Given  $\theta^t$ , calculate  $\{\zeta_r^{t+1}\}$  by solving

$$\prod_{M_0} p_r(\mathbf{x}, \zeta_r^{t+1}) = L\theta^t - \sum_{r=1}^L \zeta_r^{t+1} = p_0(\mathbf{x}, \theta^t), \quad r = 1, \dots, L$$

**outer loop:** Given a set of  $\{\zeta_r^{t+1}\}$  as the result of the inner loop, update

$$\theta^{t+1} = L\theta^t - \sum_{r=1}^L \zeta_r^{t+1}, \quad r = 1, \dots, L$$

**Notes:** CCCP enforces the m-condition at each iteration, the e-condition is satisfied only at the convergent point, which means we search for new  $\theta^{t+1}$  until the e-condition is satisfied.

# A New e-Constraint Algorithm

Define a cost function

$$\mathcal{F}_e(\{\zeta_r\}) = \sum_r \|\eta_0(\theta) - \eta_r(\zeta_r)\|^2$$

under the e-constraint  $\theta = \sum_r \zeta_r / (L - 1)$ . We minimize  $\mathcal{F}_e$  by the gradient descent algorithm. The gradient is

$$\frac{\partial \mathcal{F}_e}{\partial \zeta_r} = -2I_r(\zeta_r)[\eta_0(\theta) - \eta_r(\zeta_r)] + \frac{2}{L-1}I_0(\theta) \sum_r [\eta_0(\theta) - \eta_r(\zeta_r)]$$

Here,  $I_0(\theta)$  and  $I_r(\zeta_r)$  are the Fisher information matrices of  $p_0(\mathbf{x}, \theta)$  and  $p_r(\mathbf{x}, \zeta_r)$ , respectively, which are defined as

$$I_0(\zeta_r) = \partial_\theta \eta_0(\theta) = \partial_\theta^2 \psi_0(\theta), \quad I_r(\zeta_r) = \partial_{\zeta_r} \eta_r(\zeta_r) = \partial_{\zeta_r}^2 \psi_r(\zeta_r), \\ r = 1, \dots, L.$$

# A New e-Constraint Algorithm

Under the gradient descent algorithm,  
 $\zeta_r$  and  $\theta$  are updated as

$$\zeta_r^{t+1} = \zeta_r^t - \delta \frac{\partial \mathcal{F}_e}{\partial \zeta_r^t}$$

$$\theta^{t+1} = \frac{1}{L-1} \sum_r \zeta_r^{t+1}$$

where  $\delta$  is a small positive learning rate.

# A New e-Constraint Algorithm

## *A New e-Constraint Algorithm*

1. Set  $t = 0$ ,  $\theta^t = \mathbf{0}$ ,  $\zeta_r^t = \mathbf{0}$ ,  $r = 1, \dots, L$ .
2. Calculate  $\eta_0(\theta^t)$ ,  $I_0(\theta^t)$ , and  $\eta_r(\zeta_r^t)$ ,  $r = 1, \dots, L$ .
3. Let  $h_r = \eta_0(\theta^t) - \eta_r(\zeta_r^t)$  and calculate  $\eta_r(\zeta_r^t + \alpha h_r)$  for  $r = 1, \dots, L$ , where  $\alpha > 0$  is small. Then calculate

$$g_r = \frac{\eta_r(\zeta_r^t + \alpha h) - \eta_r(\zeta_r^t)}{\alpha}.$$

4. For  $t = 1, 2, \dots$ , update  $\zeta_r^{t+1}$  as follows:

$$\zeta_r^{t+1} = \zeta_r^t - \delta \left[ -2g_r + \frac{2}{L-1} I_0(\theta^t) \sum_r h_r \right],$$
$$\theta^{t+1} = \frac{1}{L-1} \sum_r \zeta_r^{t+1}.$$

5. If  $\mathcal{F}_e(\{\zeta_r\}) = \sum_r \|\eta_0(\theta) - \eta_r(\zeta_r)\|^2 > \epsilon$  ( $\epsilon$  is a threshold) holds,  $t+1 \rightarrow t$  and go to 2.

Thanks!

Thanks!