

# Information Geometry and Minimum Description Length Networks

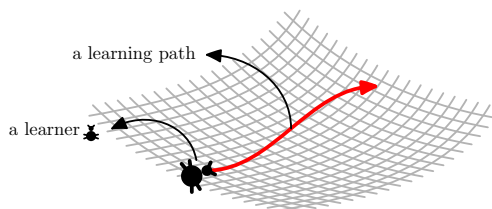
K. Sun<sup>\*</sup>, J. Wang<sup>†</sup>, A. Kalousis<sup>‡\*</sup>, S. Marchand-Maillet<sup>\*</sup>

<sup>\*</sup> University of Geneva

<sup>†</sup> Expedia

<sup>‡</sup> University of Applied Sciences Western Switzerland

# A Statistical Manifold $\mathcal{S}$



$$\mathcal{S} = \{\theta : p(\mathbf{x} | \theta) \text{ has certain structures}\}$$

**A point**  $\theta \in \mathcal{S}$  is a probability distribution

**Learning** forms a path  $\theta^0 \rightarrow \theta^1 \rightarrow \dots$

**Geometry** of  $\mathcal{S}$ , defined by

$$\left\{ \begin{array}{l} \text{Fisher Information Metric (FIM) [Rao45]} \\ \alpha\text{-connections [\check{C}encov82][Amari00]} \\ \text{Divergence [Csiszár63][Bregman67][Amari00]} \end{array} \right.$$

# Let $\mathcal{S}$ be an Exponential Family

In an exponential family  $\mathcal{S}$  with a base measure  $\sigma(\mathbf{x})$ ,

$$p(\mathbf{x} | \boldsymbol{\theta}) = \exp \left( \boldsymbol{\theta}^T \mathbf{t}(\mathbf{x}) - \psi(\boldsymbol{\theta}) \right) \quad (\psi: \text{convex})$$

$\mathcal{S}$  has two coordinate systems

- ▶ canonical parameters  $\boldsymbol{\theta}$
- ▶ expectation parameters  $\boldsymbol{\eta} = E_{p(\mathbf{x} | \boldsymbol{\theta})}(\mathbf{t}(\mathbf{x}))$

The coordinate transformation  $\boldsymbol{\theta} \leftrightarrow \boldsymbol{\eta}$  is given by [\[Amari00\]](#)

$$\boldsymbol{\eta} = \frac{\partial \psi}{\partial \boldsymbol{\theta}} \quad \text{and} \quad \boldsymbol{\theta} = \frac{\partial \psi^*}{\partial \boldsymbol{\eta}},$$

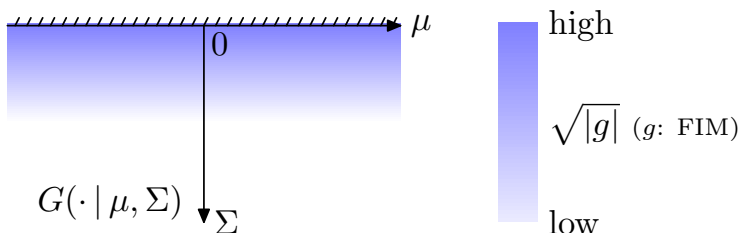
where  $\psi^* = \int p(\mathbf{x} | \boldsymbol{\theta}) \ln p(\mathbf{x} | \boldsymbol{\theta}) d\sigma(\mathbf{x})$  is the negative entropy.

# Let $\mathcal{S}$ be Gaussian

A multi-variate Gaussian distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$  has the form

$$p(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp \left( \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \text{tr} \left( \boldsymbol{\Sigma}^{-1} \mathbf{x} \mathbf{x}^T \right) \right)$$

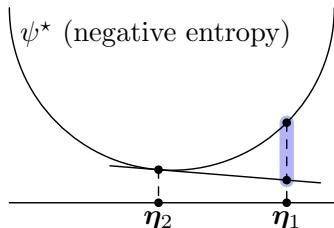
- ▶ Canonical parameters:  $\boldsymbol{\theta}^{(1)} = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$ ,  $\boldsymbol{\theta}^{(2)} = -\frac{1}{2} \boldsymbol{\Sigma}^{-1}$ ;
- ▶ Expectation parameters:  $\boldsymbol{\eta}^{(1)} = \boldsymbol{\mu}$ ,  $\boldsymbol{\eta}^{(2)} = \boldsymbol{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}^T$ .



# Information Divergence

**Geodesic Distance**  $\delta(\eta_1, \eta_2)$  is too complex

**Divergence**  $D(\eta_1 \parallel \eta_2)$  is a convenient dissimilarity measure



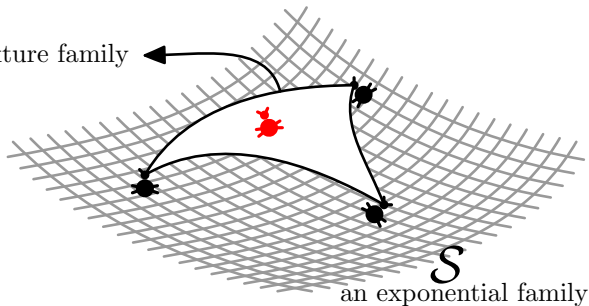
$$\begin{aligned} D(\eta_1 \parallel \eta_2) &= \psi^*(\eta_1) - \psi^*(\eta_2) - \left. \frac{\partial \psi^*}{\partial \eta} \right|_{\eta=\eta_2} (\eta_1 - \eta_2) \\ &= \dots = \psi^*(\eta_1) - \eta_1^T \theta_2 + \psi(\theta_2) \end{aligned}$$

This Kullback-Leibler (KL) divergence is among a large family of divergence measures [Csiszár63][Bregman67]

# A Mixture Model

$$p(\mathbf{x}) = \sum_{i=1}^m \alpha_i p(\mathbf{x} | \theta_i), \quad \sum_{i=1}^m \alpha_i = 1, \quad \forall i, \alpha_i \geq 0. \quad (1)$$

a mixture family



# A Geometric View of Mixture Learning [Amari95]

$\{\mathbf{x}_i\}_{i=1}^n$  i.i.d. samples in  $\mathfrak{R}^d$

$\{y_i\}_{i=1}^n$  corresponding mixture components (discrete hidden variable)

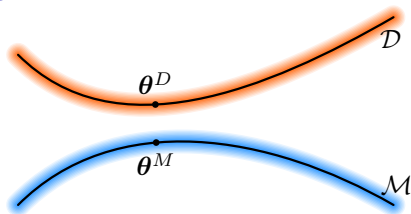
**Data sub-manifold**  $\mathcal{D} = \{p(\mathbf{x}, y)\}$

Spanned by all possible  $p(y_i)$  while fixing  $\mathbf{x}_i$

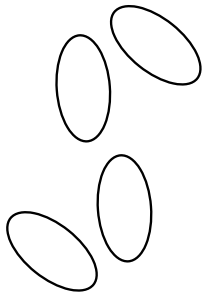
**Model sub-manifold**  $\mathcal{M} = \{p(\mathbf{x}, y)\}$

Spanned by all possible mixture models on  $\mathcal{S}$  with  $m$  components

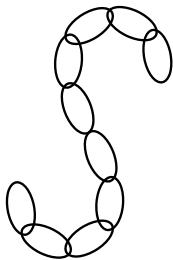
**Learning: finding  $\theta^D \in \mathcal{D}$  and  $\theta^M \in \mathcal{M}$  with minimal divergence**



proposed method



simple perception



complex perception



observations



# Divergence-induced Priors

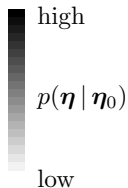
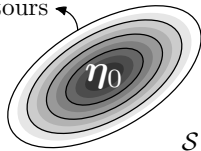
In  $\mathcal{R}^d$  distance  $\rightarrow$  probability:

$$p(\mathbf{x} | \mathbf{x}_0) \propto \exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}_0\|_2^2\right)$$

In  $\mathcal{S}$  divergence  $\rightarrow$  probability:

$$p(\boldsymbol{\eta} | \boldsymbol{\eta}_0) \propto \exp(-D(\boldsymbol{\eta} \| \boldsymbol{\eta}_0)) \quad (\text{over a compact region on } \mathcal{S})$$

equal-divergence contours

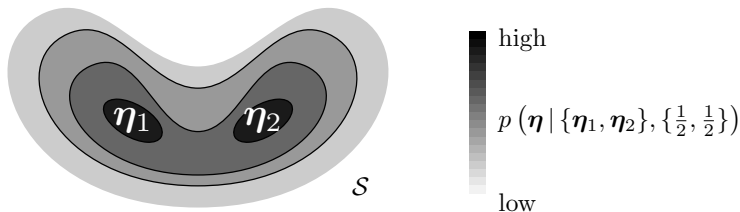


# Divergence-induced Priors

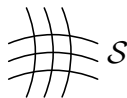
- ▶  $\mathcal{B} = \{\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_m\} \subset \mathcal{S}$
- ▶  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_m)$  s.t.  $\sum_{i=1}^m \alpha_i = 1, \forall i, \alpha_i > 0$

$$p(\boldsymbol{\eta} | \mathcal{B}, \boldsymbol{\alpha}) \propto \sum_{i=1}^m \alpha_i \exp(-D(\boldsymbol{\eta} \| \boldsymbol{\eta}_i)), \quad (2)$$

- ▶  $\ln(\sum_{i=1}^m \alpha_i \exp(-D(\boldsymbol{\eta} \| \boldsymbol{\eta}_i))) \geq \sum_{i=1}^m \alpha_i (-D(\boldsymbol{\eta} \| \boldsymbol{\eta}_i))$
- ▶ Like a kernel density estimator on  $\mathcal{S}$



# The Description Length



The code length [Shannon48] of  $\eta$  is

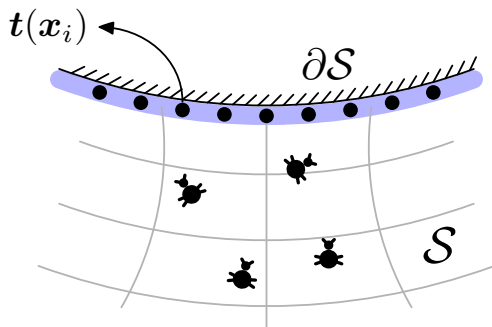
$$-\ln \left[ p(\eta | \mathcal{B}, \alpha) \sqrt{|g(\eta)|} \delta^{\dim \mathcal{S}} \right] = -\ln \left( \sum_{i=1}^m \alpha_i \exp(-D(\eta || \eta_i)) \right) \\ + \ln N(\mathcal{B}, \alpha) - \frac{1}{2} \ln |g(\eta)| - \dim \mathcal{S} \ln \delta$$

(  $N(\mathcal{B}, \alpha)$ : normalizer;  $g(\eta)$ : FIM;  $\delta$ : precision )

- ▶ Learning is based on the red term, which is
  - ▶ **novelty** of the new knowledge  $\eta$  w.r.t.  $\mathcal{B}$  and  $\alpha$
  - ▶ sensitive to parameter variations
- ▶ Other terms are useful in post-learning model assessment

## Samples on the Boundary $\partial\mathcal{S}$

Given a single observation  $\mathbf{x}$ ,  $\mathbf{t}(\mathbf{x})$  in the  $\eta$ -coordinates is on  $\partial\mathcal{S}$

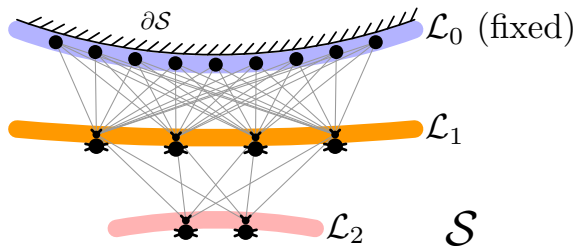


$$\sum_{j=1}^m \alpha_j \exp(-D(\mathbf{t}(\mathbf{x}) \parallel \boldsymbol{\eta}_j)) \quad \propto \quad \sum_{j=1}^m \alpha_j p(\mathbf{x} \mid \boldsymbol{\eta}_j)$$

divergence-induced prior                      sample likelihood

learning  $\min \left[ -\sum_{i=0}^n \ln \left( \sum_{j=1}^m \alpha_j \exp(-D(\mathbf{t}(\mathbf{x}_i) \parallel \boldsymbol{\eta}_j)) \right) \right]$

## Minimum Description Length [Rissanen78] Network



$\mathcal{N}$  a network of mixture components  $\{\eta_{li}\} \subset \mathcal{S}$

layer  $\mathcal{L}_0$  input samples  $\{t(\mathbf{x}_i)\}$

layer  $\mathcal{L}_l$  free points (distributions) on  $\mathcal{S}$

shape  $|\mathcal{L}_l| > |\mathcal{L}_{l+1}|, \quad \forall l$

cost function  $E = -\sum_l \sum_i \ln \left( \sum_j \alpha_{l+1,j} \exp(-D(\eta_{li} \parallel \eta_{l+1,j})) \right)$

## Implementations

$$\text{HARDN} \quad E \leq \sum_l \sum_i \min_j (-\ln \alpha_{l+1,j} + D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j}))$$

$$\text{SOFTN} \quad E \leq \sum_l \sum_i \sum_j \beta_{li}^j \left( \ln \frac{\beta_{li}^j}{\alpha_{l+1,j}} + D(\boldsymbol{\eta}_{li} \parallel \boldsymbol{\eta}_{l+1,j}) \right)$$

where  $\forall l, i, j, \beta_{li}^j \geq 0$  and  $\sum_{j=1}^{n_{l+1}} \beta_{li}^j = 1$

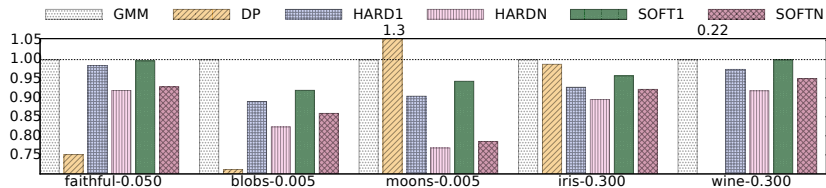
### Symmetrized Centroid [Nielsen & Nock09]

$$\min \sum_i w_i^L D(\boldsymbol{\eta}_i^L \parallel \boldsymbol{\eta}) + \sum_j w_j^R D(\boldsymbol{\eta} \parallel \boldsymbol{\eta}_j^R)$$

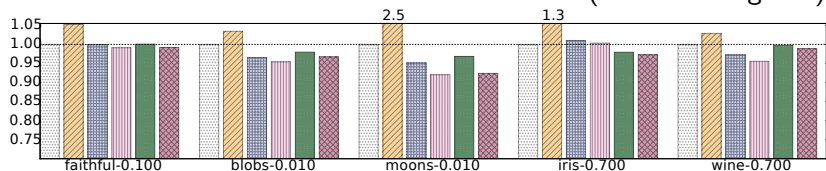
w. r. t. some given  $\{(\boldsymbol{\eta}_i^L, w_i^L)\}, \{(\boldsymbol{\eta}_j^R, w_j^R)\} \subset \mathcal{S} \times \mathbb{R}^+$ .

**Strategy:** Natural gradient [Amari98] descent  
(see paper for details)

# Testing Error on Toy Datasets



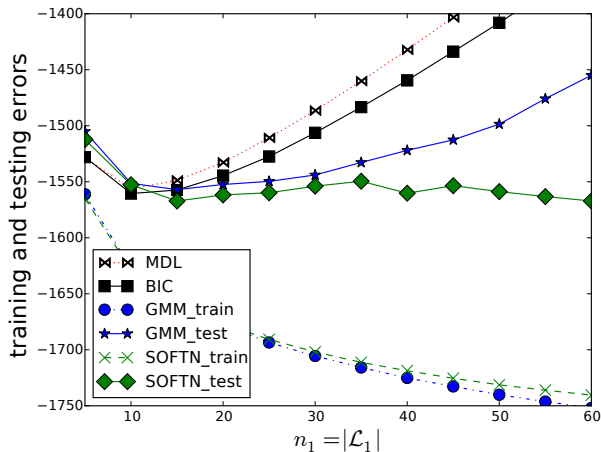
(small training size)



(large training size)

4 v.s. 3 (6 v.s. 5) shows the effectiveness of letting  $\mathcal{L}_2$  regularize  $\mathcal{L}_1$

# Testing Error on hand-written “1”s



The size of each layer is a hyper-parameter



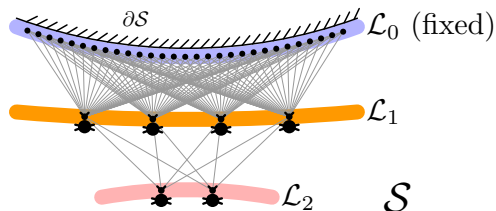
Effective regularization means relying less on model selection



# Consistency

## Theorem 1

If the truth is a finite mixture model w. r. t.  $\{\eta_i^t\}$ , then as the sample size  $n \rightarrow \infty$ ,  $\mathcal{L}_1^*$  in an optimal MDL network is exactly  $\{\eta_i^t\}$ .



(intuitively, the cost between  $\mathcal{L}_0$  and  $\mathcal{L}_1$  will dominate)

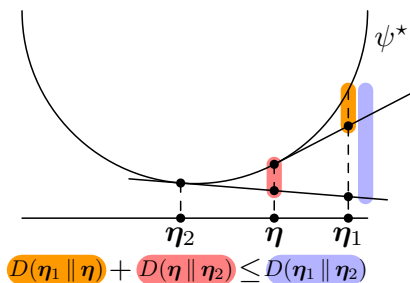
# The Gain

## Theorem 2

$\forall \eta_1, \eta_2 \in \mathcal{S}, \eta_1 \neq \eta_2$ , then  $\exists \eta \in \mathcal{S}$ , s.t.

$$D(\eta_1 \parallel \eta) + D(\eta \parallel \eta_2) \leq D(\eta_1 \parallel \eta_2) - \max\{D(\eta_{lc} \parallel \eta_1) + D(\eta_1 \parallel \eta_{lc}), \\ D(\eta_{rc} \parallel \eta_2) + D(\eta_2 \parallel \eta_{rc})\},$$

where  $\theta_{lc} = (\theta_1 + \theta_2)/2$ , and  $\eta_{rc} = (\eta_1 + \eta_2)/2$ .



**Surprise** the “path” becomes shorter after taking an intermediate stop

## v.s. Bayesian Mixture Models

- ▶ No heavy integration involved
- ▶ Geometric centroids instead of Bayes' rule

## v.s. Neural Networks

- ▶ Non-linearity is introduced by symmetrized centroid
- ▶ layers: +regularization instead of +flexibility

## v.s. Hierarchical Clustering

- ▶ The whole network is learned at once

# information geometric compactness of a learning network

$\mathcal{S}$  Gaussian  $\rightarrow$  Bernoulli, categorical,  $\dots$

$D(\cdot\|\cdot)$  other divergence measures

**Network Structure** sparsity, dropout,  $\dots$

## Thank You

Codes available at <https://git.unige.ch/gitweb/marchand/mdlnetworks>

- ★ Olivier Schwander helped proof-read the slides
  - ★ Tom SF Hanines helped with the Dirichlet Process codes
- 

Q & A