

Problem Setting

non-normalized model = statistical model with an intractable normalization constant
 ▶ Markov random field, Boltzmann machine, overcomplete ICA, ...

$$p(x | \theta) = \frac{1}{Z(\theta)} \tilde{p}(x | \theta), \quad Z(\theta) = \int \tilde{p}(x | \theta) dx$$

↑
computationally intractable

non-normalized mixture model
 = finite mixture of non-normalized models

$$p(x | \theta, \pi) = \sum_{k=1}^K \pi_k \cdot p(x | \theta_k)$$

$$p(x | \theta_k) = \frac{1}{Z(\theta_k)} \tilde{p}(x | \theta_k), \quad \pi_k > 0, \quad \sum_{k=1}^K \pi_k = 1$$

$x_1, \dots, x_N \sim p(x | \theta, \pi)$

We develop a general method for estimating $\theta = (\theta_1, \dots, \theta_K)$ and $\pi = (\pi_1, \dots, \pi_K)$ **without** computing $Z(\theta_k)$

- ▶ extension of noise contrastive estimation (Gutmann and Hyvärinen, 2012)
- ▶ can even be used on deep image representations

Proposed Method

Reparametrization: $(\theta, \pi) \rightarrow (\theta, c)$

$$p(x | \theta, c) = \sum_{k=1}^K p(x | \theta_k, c_k)$$

$$\log p(x | \theta_k, c_k) = \log \tilde{p}(x | \theta_k) + c_k, \quad c_k = \log \pi_k - \log Z(\theta_k)$$

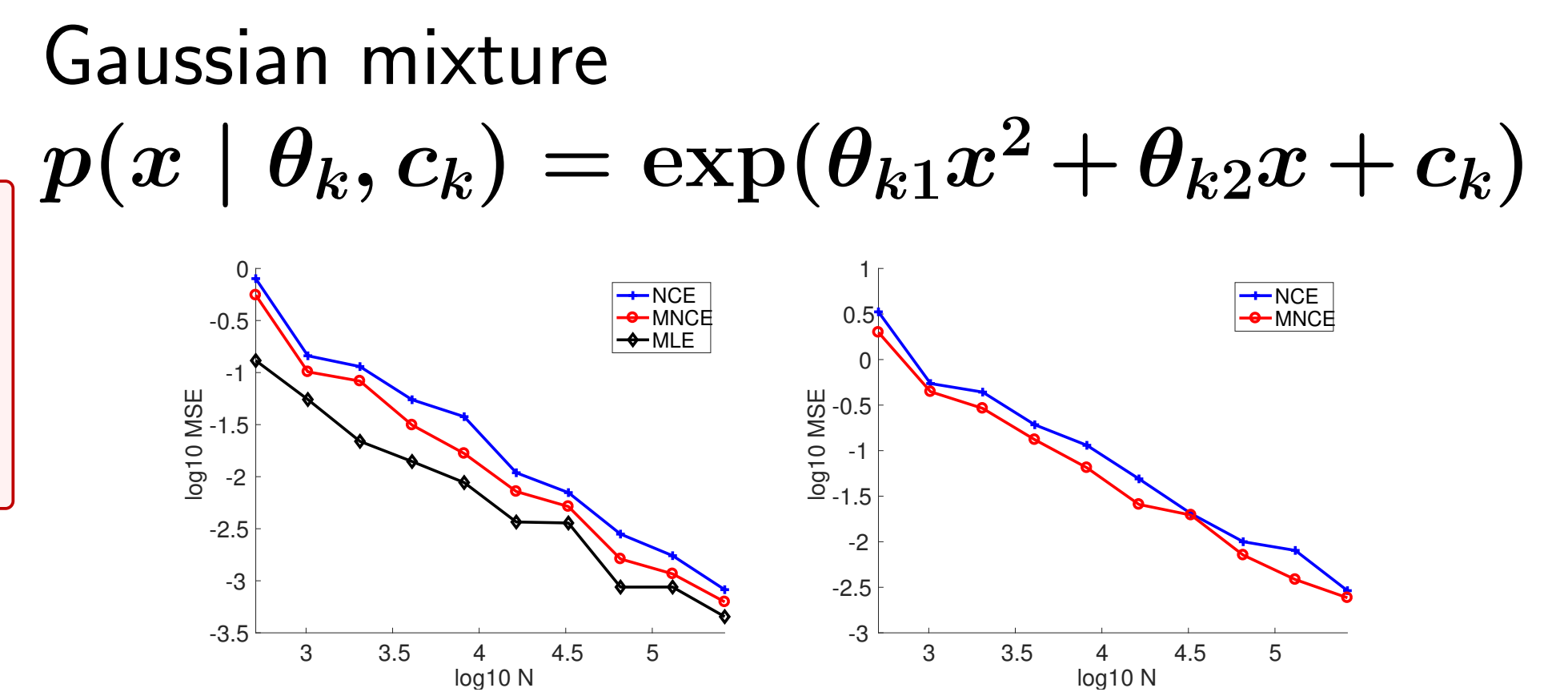
Noise generation

We generate artificial noise $y_1, \dots, y_M \sim n(y)$

- ▶ should be difficult to discriminate from data (cf. GAN)
- ▶ e.g., Gaussian with same mean and covariance as data

idea: estimate θ and c by discriminating between data and noise

$$(\hat{\theta}, \hat{c}) = \arg \max_{\theta, c} \sum_{t=1}^N \log \frac{Np(x_t | \theta, c)}{Np(x_t | \theta, c) + Mn(x_t)} + \sum_{t=1}^M \log \frac{Mn(y_t)}{Np(y_t | \theta, c) + Mn(y_t)}$$



This estimator has consistency under mild regularity conditions (Theorem 1)

We can improve estimation accuracy by using multiple noise distributions:

$$y_1^{(l)}, \dots, y_{M_l}^{(l)} \sim n_l(y) \quad (l = 1, \dots, L)$$

$$(\hat{\theta}, \hat{c}) = \arg \max_{\theta, c} \sum_{t=1}^N \log \frac{Np(x_t | \theta, c)}{Np(x_t | \theta, c) + \sum_{l=1}^L M_l n_l(x_t)} + \sum_{l=1}^L \sum_{t=1}^{M_l} \log \frac{M_l n_l(y_t^{(l)})}{Np(y_t^{(l)} | \theta, c) + \sum_{l=1}^L M_l n_l(y_t^{(l)})}$$

This estimator is equivalent to the original one with a mixture noise distribution (Theorem 2)

Clustering with Deep Representation

x : data (e.g., image), z : label (e.g., "dog")

Classification with neural network (softmax)

$$p(z = l | x) = \frac{\exp(\sum_{i=1}^d w_{li} f_i(x))}{\sum_{j=1}^L \exp(\sum_{i=1}^d w_{ji} f_i(x))}$$

$f(x) = (f_1(x), \dots, f_d(x))$: feature vector (activation of last hidden layer)

Non-normalized exponential family

$$p(x | z = l) = h(x) \exp \left(\sum_{i=1}^d w_{li} f_i(x) - A(w_l) \right)$$

↑ unknown ↑ unknown

We propose a probabilistically principled method for transferring the deep representation f to clustering of unlabeled data x_1, \dots, x_N .

$$p(x | \theta, c) = h(x) \sum_{k=1}^K \exp \left(\sum_{i=1}^d \theta_{ki} f_i(x) + c_k \right) : \text{target data}$$

$$n_l(x) = h(x) \exp \left(\sum_{i=1}^d w_{li} f_i(x) - A(w_l) \right) : \text{original training data}$$

($l = 1, \dots, L$)

$$\rightarrow (\hat{\theta}, \hat{c}) \rightarrow p(z_t = k | x_t; \hat{\theta}, \hat{c}) \rightarrow \text{clustering}$$

Note: the unknown function h **cancels out**

Image clustering

data: 12,500 dog images & 12,500 cat images
 deep representation: inception-v3 ($d = 2048, L = 149, K = 2$)
 clustering results (GMM = Gaussian Mixture Model; diagonal/isotropic covariance)

proposed	dog	cat	GMM1	dog	cat	GMM2	dog	cat
cluster 1	12400	145	cluster 1	12490	325	cluster 1	12490	792
cluster 2	100	12355	cluster 2	10	12175	cluster 2	10	11708

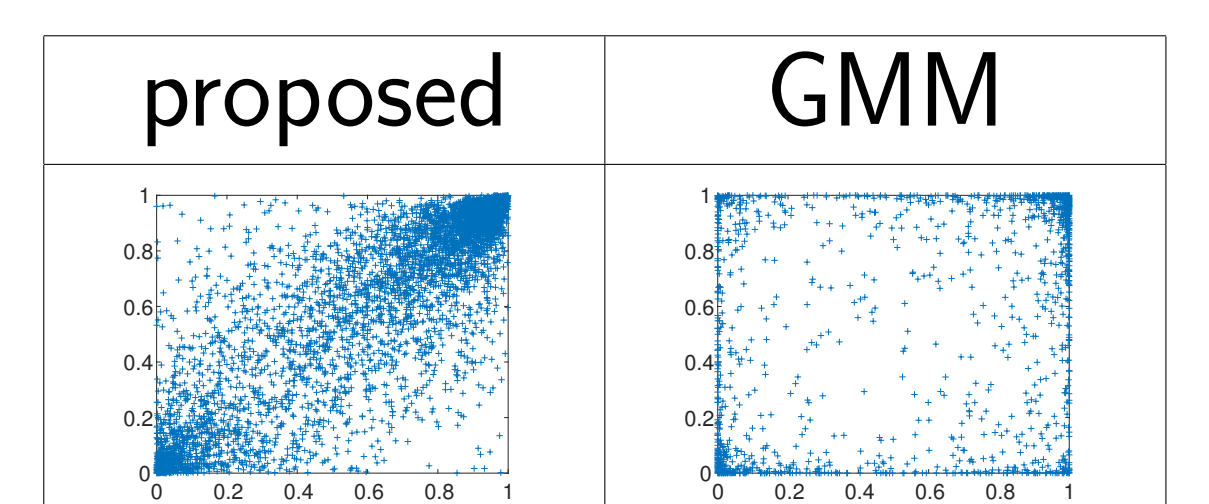
The proposed method has better classification accuracy
 estimated and actual numbers of misclassifications

	estimate	actual	estimate
proposed	169.98	245	$= \sum_x \min(p(z = 1 x), p(z = 2 x))$
GMM1	0.66	335	The proposed method quantifies clustering uncertainty more accurately
GMM2	5.58	802	

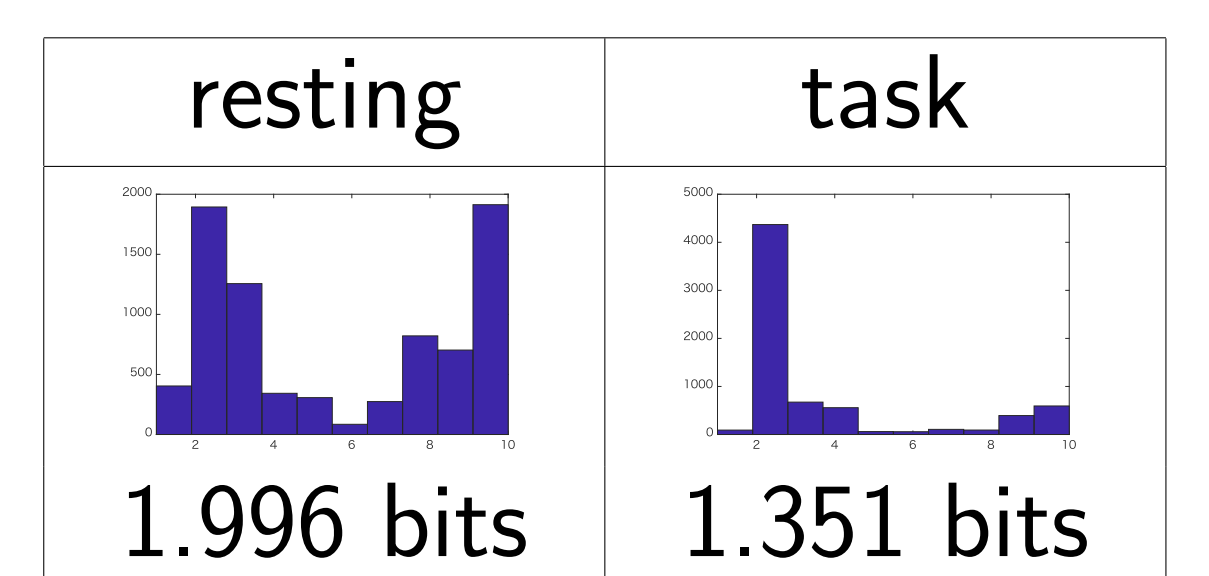
Brain state clustering

data: 306-ch magnetoencephalography (MEG) from CamCAN repository
 deep representation: obtained by nonlinear ICA with Time Contrastive Learning (Hyvärinen and Morioka, 2016)

scatter plots of $p(z_{t-1} = 1 | x_{t-1})$ and $p(z_t = 1 | x_t)$ for $K = 2$
 → proposed method extracts stable brain states



state histogram for $K = 10$
 → Resting MEG has more temporal variability of brain states



▶ consistent with previous findings in neuroscience