

Eraser Machine (ϵ -machine)

I. Introduction

- System Inference

: Given observed configurations, σ , obtain their underlying prob. dist.
 $p(\sigma)$

↳ ① reconstructing missing parts

$$p(\sigma_i | \sigma_1 \sigma_2 \dots \sigma_n) = \frac{p(\sigma)}{\sum_{\sigma_i} p(\sigma)}$$

② generating samples

$$p(\sigma) \rightarrow \sigma$$

③ underlying model, e.g., Hopfield model, Boltzmann machine,

$$P(\sigma) = \frac{e^{-E(\sigma)}}{Z} \quad E(\sigma) = -\sum_i b_i \sigma_i - \sum_{j > k} w_{jk} \sigma_j \sigma_k$$

$$Z = \sum_{\{\sigma\}} e^{-E(\sigma)}$$

- Issue for computing partition function.

For example, if $n=20$, # of $\sigma = 2^{20} \approx 10^6$

How to overcome this issue? literature search

- Our idea

by reweighting data configuration, trivialize
to compute partition function.

For example, at high temperature, $Z = \sum_{\{\sigma\}} e^{-\beta E(\sigma)} \approx 2^n$

"Attain complete vacuity" by Lao Tzu

⇒ develop a very efficient algorithm to obtain
underlying model distributions from general observations.

II. Maximum likelihood estimation (MLE)

- Data: $f(\sigma)$ relative observed # of σ configuration

- Model: $p(\sigma)$

- Likelihood: $L = \prod_{\{\sigma\}} p(\sigma)^{f(\sigma)}$

$$p^*(\sigma) = f(\sigma), \text{ solution maximizing } \log L$$

- Specific model

$$p(\sigma) = \frac{e^{-E(\sigma)}}{Z}, \quad E(\sigma) = -\sum_i b_i \sigma_i - \sum_{j < k} w_{jk} \sigma_j \sigma_k$$

$$Z = \sum_{\{\sigma\}} e^{-E(\sigma)} = \omega \cdot \sigma$$

$$\left\{ \begin{array}{l} \omega = (b_i, w_{jk}) \\ \sigma = (\sigma_i, \sigma_j \sigma_k) \end{array} \right.$$

for example, $n=2$

$$\omega = (b_1, b_2, w_{12})$$

$$\sigma = (\sigma_1, \sigma_2, \sigma_1 \sigma_2)$$

Then,

$$L = L(\omega)$$

$$\omega^* = ?$$

$$\textcircled{1} \text{ Hopfield solution } \omega_I^* = \sum_{\{\sigma\}} \sigma_I \cdot f(\sigma) \equiv \langle \sigma_I \rangle_f$$

$$\hookrightarrow b_i = \langle \sigma_i \rangle_f, \quad w_{jk} = \langle \sigma_j \sigma_k \rangle_f$$

$$\textcircled{2} \text{ MLE (Boltzmann machine)}$$

$$\left. \frac{\partial \log L}{\partial \omega_I} \right|_{\omega=\omega_I^*} = 0, \quad I \in \{1, 2, \dots, n + \frac{n(n-1)}{2}\}$$

$$\frac{\partial \log L}{\partial w_I} = -\langle \sigma_I \rangle_f + \langle \sigma_I \rangle_p$$

$$\langle \sigma_I \rangle_f = \sum_{\{\sigma\}} \sigma_I \cdot f(\sigma)$$

$$\langle \sigma_I \rangle_p = \sum_{\{\sigma\}} \sigma_I \cdot \frac{e^{-E(\sigma)}}{Z}$$

$$w_I^{\text{new}} = w_I + d \cdot \frac{\partial \log L}{\partial w_I}$$

↖ learning rate

Problem:

- { Hopfield solution \rightarrow not accurate
- Boltzmann machine \rightarrow computationally heavy
to obtain Z
even intractable
for large $n > 20$

③ erasure machine (ϵ -machine)

* reweighting with an arbitrary distribution $g(\sigma)$

$$f(\sigma) \rightarrow \tilde{f}(\sigma) \propto f(\sigma) \cdot g^{\epsilon-1}(\sigma) \quad \tilde{f}(\sigma) = \frac{f(\sigma) \cdot g(\sigma)}{\sum_{\sigma'} f(\sigma') \cdot g(\sigma')}^{\epsilon-1}$$

with a small parameter $\epsilon \in [0, 1]$.

$$p^*(\sigma) = f(\sigma)$$

↙

$$\tilde{p}^*(\sigma) = \tilde{f}(\sigma)$$

$$\text{Here, if } g(\sigma) = p^*(\sigma), \quad \tilde{p}^*(\sigma) \propto f(\sigma) \cdot f(\sigma)^{\epsilon-1}$$

$$\propto f^\epsilon(\sigma)$$

$$\propto [p^*(\sigma)]^\epsilon$$

$$\tilde{p}(\sigma) = \frac{e^{-\epsilon E(\sigma)}}{\tilde{Z}}, \quad \tilde{Z} = \sum_{\{\sigma\}} e^{-\epsilon E(\sigma)}$$

tempting to interpret ϵ as the inverse temperature β in statistical mechanics.

But it should be cautious !!

$$P^*(\sigma) = f(\sigma)$$

$$\tilde{P}^*(\sigma) \propto [P^*(\sigma)]^\beta \quad \text{not by simply exponentiating with } \beta$$

$$\tilde{P}^*(\sigma) \propto P^*(\sigma) \cdot [P^*(\sigma)]^{\epsilon-1} = [P^*(\sigma)]^\epsilon$$

Now let's work with reweighted data distribution

$$\tilde{L} = \prod_{\{\sigma\}} \tilde{P}(\sigma)^{\tilde{f}(\sigma)} \quad \text{where } \tilde{P}(\sigma) = \frac{e^{-\epsilon E(\sigma)}}{\tilde{Z}}$$

$$\begin{aligned} \log \tilde{L} &= \sum_{\{\sigma\}} \tilde{f}(\sigma) \cdot \log \tilde{P}(\sigma) \\ &= \sum_{\{\sigma\}} \tilde{f}(\sigma) \cdot [-\epsilon \cdot E(\sigma) - \log \tilde{Z}] \end{aligned}$$

$$\frac{\partial \log \tilde{L}}{\partial w_I} = \epsilon \langle \sigma_I \rangle_{\tilde{f}} - \epsilon^2 w_I$$

$$\text{where } \langle \sigma_I \rangle_{\tilde{f}} = \sum_{\{\sigma\}} \sigma_I \cdot \tilde{f}(\sigma)$$

① Second term
contributes as a regularizer
to prevent w_I from blowing up

$$\text{② } \epsilon=1, \frac{\partial \log \tilde{L}}{\partial w_I} = 0 \rightarrow \text{Hopfield solution}$$

To obtain the second term, we use a small ϵ expansion.

(Supplementary material)

$$\begin{aligned} \tilde{Z} &= \sum_{\{\sigma\}} e^{-\epsilon E(\sigma)} \\ &= \sum_{\{\sigma\}} e^{\sum_i \epsilon b_i \sigma_i + \sum_{j < k} \epsilon w_{jk} \sigma_j \sigma_k} \\ &= \sum_{\{\sigma\}} \prod_i \cosh(\epsilon b_i) \cdot [1 + \sigma_i \tanh(\epsilon b_i)] \prod_{j < k} \cosh(\epsilon w_{jk}) \\ &\quad \times [1 + \sigma_j \sigma_k \tanh(\epsilon w_{jk})] \end{aligned}$$

$$= 2^n \prod_i \cosh(\epsilon b_i) \prod_{j < k} \cosh(\epsilon w_{jk}) \cdot \left[1 + \sum_{\alpha < \beta} \tanh(\epsilon b_\alpha) \cdot \tanh(\epsilon b_\beta) \right. \\ \left. + \sum_{\substack{\alpha < \beta \\ \alpha < r \\ \beta < r}} \tanh(\epsilon w_{\alpha\beta}) \cdot \tanh(\epsilon w_{\alpha r}) \cdot \tanh(\epsilon w_{\beta r}) + O(\epsilon^4) \right]$$

- $\log \tilde{Z} = n \cdot \log 2 + \sum_i \log \cosh(\epsilon b_i) + \sum_{j < k} \log \cosh(\epsilon w_{jk}) + O(\epsilon^3)$
- $\delta \log \tilde{Z} = \sum_i \tanh(\epsilon b_i) \cdot \epsilon \cdot \delta b_i + \sum_{j < k} \tanh(\epsilon w_{jk}) \cdot \epsilon \cdot \delta w_{jk} + O(\epsilon^3)$
 $\approx \sum_i \epsilon^2 b_i \delta b_i + \sum_{j < k} \epsilon^2 w_{jk} \delta w_{jk} = \sum_I \epsilon^2 w_I \cdot \delta w_I$
- $\frac{\partial \log \tilde{Z}}{\partial w_I} = \epsilon^2 w_I$

Algorithm

① reweight $\tilde{f}(\sigma) = \frac{f(\sigma) \cdot q(\sigma)^{\epsilon-1}}{\sum_{\{\sigma'\}} f(\sigma') \cdot q(\sigma')^{\epsilon-1}}$ for $q(\sigma)$ initially
 $q(\sigma) = \frac{e^{-E(\sigma)}}{Z}$

$$E(\sigma) = -w \cdot \sigma$$

\uparrow
random w

② $\langle \sigma_I \rangle_{\tilde{f}} = \sum_{\{\sigma\}} \sigma_I \cdot \tilde{f}(\sigma)$

③ $\frac{\partial \log L}{\partial w_I} = \epsilon \langle \sigma_I \rangle_{\tilde{f}} - \epsilon^2 w_I$

④ Update $w_I^{\text{new}} = w_I + \alpha \cdot \frac{\partial \log L}{\partial w_I} = w_I + \alpha' (\langle \sigma_I \rangle_{\tilde{f}} - \epsilon w_I)$

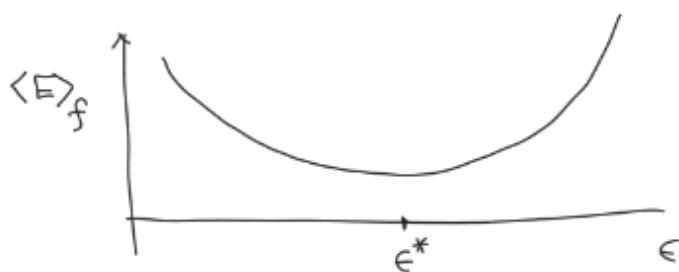
$$\alpha' = \alpha \cdot \epsilon$$

⑤ Iterate this until $D_{KL}(\tilde{f} || \tilde{p})$ minimizes.

$$\tilde{p}(\sigma) = \frac{e^{-E(\sigma)}}{\tilde{Z}}$$

$$\begin{aligned}
D_{KL}(\tilde{f} \parallel \tilde{p}) &= \sum_{\{\sigma\}} \tilde{f}(\sigma) \cdot \log \frac{\tilde{f}(\sigma)}{\tilde{p}(\sigma)} \\
&= \sum_{\{\sigma\}} \tilde{f}(\sigma) \cdot \log \tilde{f}(\sigma) - \tilde{f}(\sigma) \cdot \log \tilde{p}(\sigma) \\
&= \sum_{\{\sigma\}} \tilde{f}(\sigma) \cdot \log \tilde{f}(\sigma) - \tilde{f}(\sigma) \cdot \left[-\epsilon E(\sigma) - \log \tilde{Z} \right] \\
&= -S + \epsilon \langle E \rangle_{\tilde{f}} + \log \tilde{Z} \\
&\approx -S + \epsilon \langle E \rangle_{\tilde{f}} + n \cdot \log 2 + \sum_I \log \cosh(\epsilon w_I)
\end{aligned}$$

Optimal ϵ



$$\langle E \rangle_p = \sum_{\{\sigma\}} E(\sigma) \cdot p(\sigma), \quad p(\sigma) = \frac{e^{-E(\sigma)}}{Z}, \quad E(\sigma) = w \cdot \sigma$$

if $w = w^{\text{true}}$

$$\Rightarrow E(\sigma) = E^{\text{true}}(\sigma) \quad \text{and} \quad p(\sigma) = f(\sigma)$$

$$\text{Then, } \langle E \rangle_p = \langle E \rangle_f$$

if the inferred values of $\langle E \rangle_f$ do not vary for similar values of ϵ , this means that

$$\langle E \rangle_f = \langle E^{\text{true}} \rangle_f \text{ is inferred correctly.}$$

Elaborate this argument!

random ϵ

$$\epsilon = \frac{\text{random}}{\max[1, w]}$$

Fig.1. Performance of ϵ -machine

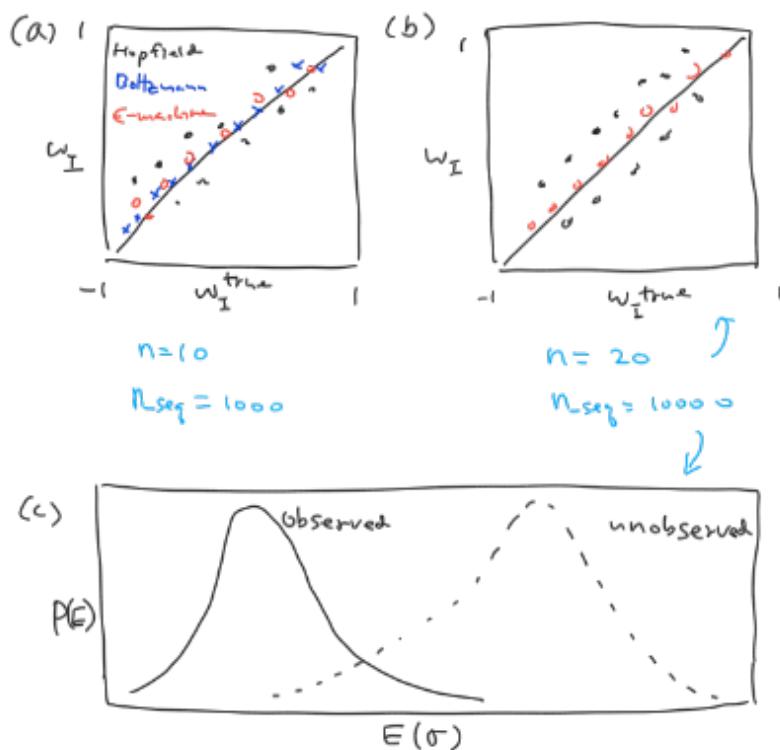
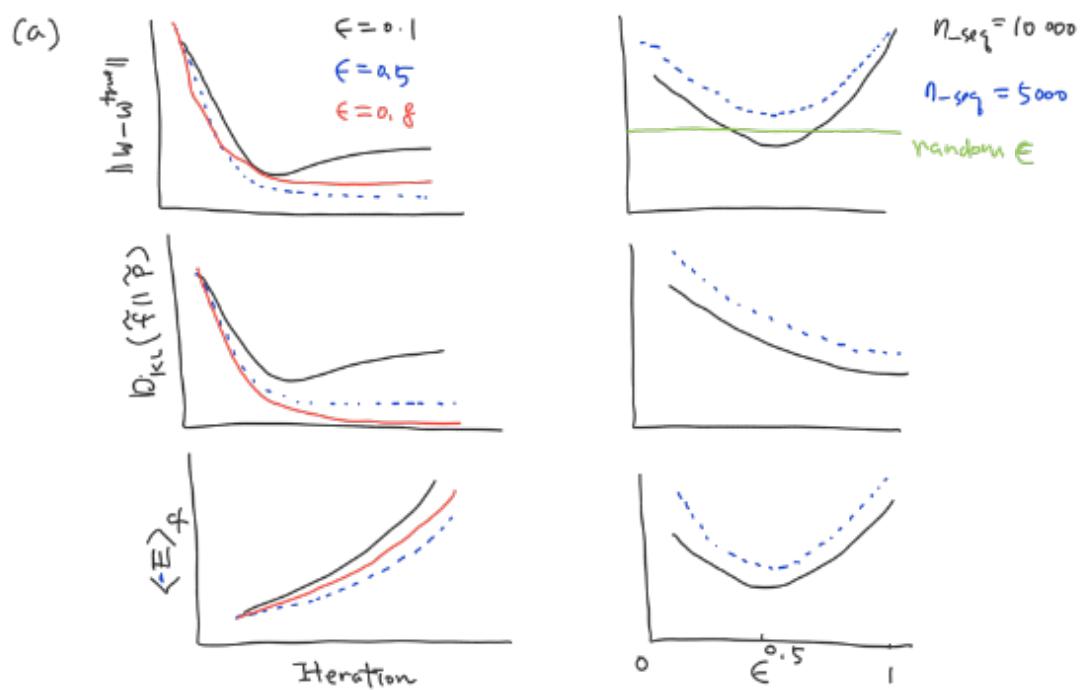


Fig.1 Learning of ϵ -machine

$n=20,$



III. Discussions

- Applications (e.g. equilibrium inference of neural networks)
- Other ways to avoid computing exact partition function.

마지막 수정: 오후 12:34