

FYS 5429, APRIL 19, 2023

- Markov chain Monte Carlo

$p_i$  (model)

$$P_i(t) = \sum_j W(j \rightarrow i) P_j(t-1)$$

steady state

$$P(t=\infty) = W P(t=\infty)$$

$$P = W P \quad \lambda = 1$$

$W$  is a stochastic matrix

$$\lambda_0 = 1 > \lambda_1 > \lambda_2 \dots \lambda_{n-1}$$

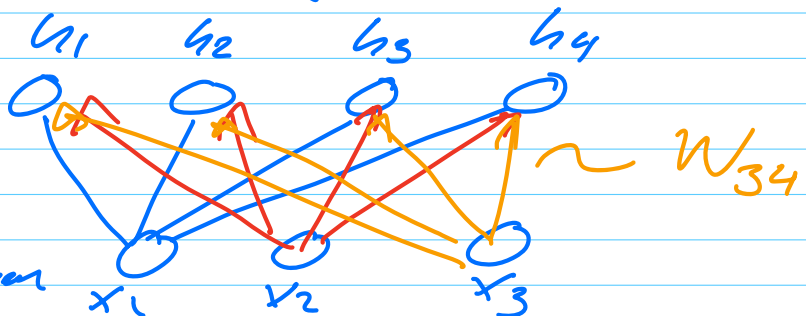
$$P_i = \frac{e^{-\beta E_i}}{Z}$$

$$E(x, h) = -b^T x - c^T h - x^T W h$$

↑  
visible  
nodes

↑  
hidden  
nodes

RBM =  
reduced  
Boltzmann



machine

$Z$  = partition function

$$Z = \sum_i e^{-\beta E_i}$$

$$P_i \rightarrow p(x, h) = e^{-\beta E(x, h)}$$

marginal probabilities

$$p(x|h) = \int_{h \in \mathcal{D}} p(x, h) dh$$

$$\left( \sum_{h_i} p(x, h_i) \right)$$

$$p(h|x) = \int_{x \in \mathcal{B}} p(x, h) dx$$

Sampling rules

- used when accepting new move/state/config

- (MC)<sup>2</sup>

- Metropolis-Hastings

- Gibbs

$$W(j \rightarrow i) = W_{i,j}$$

$$\sum_j W_{ij} = 1$$

$$W_{ij} = \underset{\substack{\uparrow \\ \text{Transition} \\ \text{probability}}}{T(j \rightarrow i)} A(j \rightarrow i) \underset{\substack{\uparrow \\ \text{acceptance} \\ \text{probability}}}{A(j \rightarrow i)}$$

$$\frac{P_i}{P_j} = \frac{T(j \rightarrow i) A(j \rightarrow i)}{T(i \rightarrow j) A(i \rightarrow j)}$$

$$T(j \rightarrow i) = T(i \rightarrow j)$$

$$\frac{P_i}{P_j} = \frac{A(j \rightarrow i)}{A(i \rightarrow j)} > 1 \quad 0 \leq A(j \rightarrow i) \leq 1$$

Metropolis's - algo

$$A(j \rightarrow i) = \min\left(1, \frac{P_i}{P_j}\right)$$

Metropolis's - Hastings

$$A(j \rightarrow i) = \min\left(1, \frac{P_i T(i \rightarrow j)}{P_j T(j \rightarrow i)}\right)$$

Gibbs-sampling.

We want k-sampler of

$$X = \{x_1, x_2, \dots, x_n\}$$

from a distribution (Model)

$$p(x_1, x_2, \dots, x_n) \mid p(x|h)$$

We label the  $i$ 'th sample

$$\text{as } X^{(i)} = \{x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}\}$$

1) Begin with  $k=0$ ,  $X^{(0)}$   
(random)

2) We want the next sample  
$$X^{(i+1)} = \{x_1^{(i+1)}, x_2^{(i+1)}, \dots, x_n^{(i+1)}\}$$

$x_j^{(i+1)}$  is updated using  
the distribution

$$p(x_j^{(i+1)} \mid x_1^{(i+1)}, x_2^{(i+1)}, \dots, x_{j-1}^{(i+1)}, x_{j+1}^{(i+1)}, \dots, x_n^{(i)})$$

3) Repeat k-times

Partition function (chap 18 of Goodfellow)

$$p(x, \theta) = \frac{e^{-\beta E(x, \theta)}}{Z}$$

$$p(x; \theta) = \int p(x, \theta) d\theta$$

$$p(x; \epsilon)$$

$$\epsilon = \{b, c, w, \dots\}$$

$$p(x; \epsilon) = \frac{\tilde{p}(x; \epsilon)}{Z}$$

$$\int \tilde{p}(x) dx = Z = Z(\epsilon) \\ = \sum_x \tilde{p}(x)$$

optimization of  $\log p(x; \epsilon)$

wrt  $\epsilon$

$$\hat{\epsilon} = \arg \min_{\epsilon} \log p(x; \epsilon)$$

$$\nabla_{\epsilon} \log p(x; \epsilon) = \nabla_{\epsilon} [\log \tilde{p}(x; \epsilon) - \log Z(\epsilon)]$$

$$\nabla_G \log \tilde{p}(x; \theta) - \nabla_G \log Z(\theta)$$

$\nwarrow e^{-\beta E(x; \theta)} = \text{known}$   
 $\nearrow$  positive phase of learning       $\nearrow$  negative phase

$$\nabla_G \log Z = \frac{\nabla_G Z}{Z}$$

$$= \frac{\nabla_G \sum_x \tilde{p}(x)}{Z} = \frac{\sum_x \nabla_G \tilde{p}(x)}{Z}$$

$$\tilde{p}(x) > 0 \text{ for all } x$$

$$\frac{\sum_x \nabla_G \exp(\log \tilde{p}(x))}{Z}$$

$$= \frac{\sum_x \exp(\log \tilde{p}(x)) \nabla_G \log \tilde{p}(x)}{Z}$$

$$\frac{\sum_x \tilde{p}(x) \nabla_G \log \tilde{p}(x)}{Z} \sim p(x)$$

=

$$= \sum_x p(x) \nabla_{\theta} \log \tilde{p}(x)$$

$$= \mathbb{E}[\nabla_{\theta} \log \tilde{p}(x)]$$

$$\approx \frac{1}{N} \sum_i \nabla_{\theta} \log \tilde{p}(x_i)$$

$$\nabla_{\theta} \log p(x; \theta) =$$

$$\nabla_{\theta} \log \tilde{p}(x; \theta) -$$

$$\mathbb{E}[\nabla_{\theta} \log \tilde{p}(x; \theta)]$$

Algorithm MCMC for maximizing the log-likelihood  
use gradient ascent

Minimize  $\rightarrow$  gradient descent

- set stepsize  $\epsilon \sim 10^{-3}$   
(learning rate)

- set of Gibbs steps

while not converged

- sample a minibatch

$\{x_1, x_2, \dots, x_m\}$

from training set

$$g \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\Theta} \log \tilde{p}(\tilde{x}_i; \Theta)$$

- initialize a set of  $m$  samples  $\{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_m\}$

for  $i = 1$  to  $K$

for  $j = 1$  to  $m$

$\tilde{x}_j \leftarrow \text{gibbs\_update}(\tilde{x}_j)$

end

end

$$g \leftarrow g - \frac{1}{m} \sum_{i=1}^m \nabla_{\Theta} \log \tilde{p}(\tilde{x}_i; \Theta)$$

$$\Theta \leftarrow \Theta + \epsilon \cdot g \quad / \quad (\epsilon - \epsilon g)$$

gradient descent.