

→ Unsupervised DL

- Three approaches:

- ML based:

DBN, RBM, DBM

- Reconstruction error:

Autoencoders

- Adversarial training:

GAN

→ RBM:

- Visible units

v_1
0
⋮
0
⋮
0
 v_d

- Hidden units

h_1
0
⋮
0
⋮
0
 h_k

$W \in \mathbb{R}^{d \times k}$

$b \in \mathbb{R}^{1 \times k}$

$a \in \mathbb{R}^{1 \times d}$

- Classical RBM's units are binary,

$\Theta: \{W, b, a\} \Rightarrow$ it has 2^{b+k} possibility of v & h

- Energy function, $E(v, h)$:

$$- \sum_{i=1}^d a_i v_i - \sum_{j=1}^k b_j h_j - \sum_{i=1}^d \sum_{j=1}^k w_{ij} v_i h_j$$

$$= -a^T v - b^T h - v^T w h$$

• Joint prob dist:

$$P(v, h; \theta) = \frac{1}{Z} e^{-E(v, h)}, \text{ normalizer: } Z = \sum_{v, h} e^{-E(v, h)}$$

• Marginal P:

$$P(v; \theta) = \frac{1}{Z} \sum_h e^{-E(v, h)}$$

• Cond P:

$$P(h|v; \theta) = \frac{\frac{1}{Z} e^{-E(v, h)}}{\frac{1}{Z} \sum_h e^{-E(v, h)}} = \frac{e^{-a^T v - b^T h - v^T w h}}{\sum_h e^{-a^T v - b^T h - v^T w h}}$$

$$\propto e^{\frac{b}{L} h_j + \frac{1}{L} \sum_i w_{ij} v_i h_j} = \frac{1}{L} e^{b_j h_j + (w_j^T v) h_j}$$

$$\Rightarrow P(h_j = 1 | v; \theta) = \frac{\tilde{P}(h_j = 1 | v; \theta)}{\tilde{P}(h_j = 1 | v; \theta) + \tilde{P}(h_j = 0 | v; \theta)}$$

$$= \sigma(w_j^T v + b_j)$$

- Different energy func gives different cond P

$$- P(v|h; \theta) = \prod_{i=1}^d P(v_i | h; \theta) = \prod_{i=1}^d \sigma(\tilde{h}^T w_i + a_i)$$

you can go backwards with this \uparrow

- Generative model: since we know joint prob.

- ML based method to learn θ

- $P(v; \theta) = \frac{1}{Z} \sum_h e^{-E(v, h)}$

- observe v_1, \dots, v_n :

- $P(v; \theta) = \prod_{i=1}^n P(v_i; \theta) \Rightarrow \log P(v; \theta) =$

$$\sum_{i=1}^n \left(\log \left(\sum_h e^{-E(v_i, h)} \right) - \log \left(\sum_v \sum_h e^{-E(v, h)} \right) \right)$$

$$\frac{\partial}{\partial w_{ij}} E(v, h) = -v_i h_j$$

$$\frac{\partial}{\partial w_{ij}} \log P(v; \theta) = \sum_{i=1}^n \frac{\sum_h e^{-E(v_i, h)} \cdot v_i h_j}{\sum_h e^{-E(v_i, h)}} - n \sum_{v, h} P(v, h) v_i h_j$$

$$= \underbrace{\sum_{i=1}^n \sum_h P(h | v_i^{(P)}; \theta) \cdot v_i h_j}_{\bar{E}_{\text{data}}} - n \underbrace{\sum_{v, h} P(v, h) \cdot v_i h_j}_{\bar{E}_{\text{model}}}$$

$$\Rightarrow \log P(v; \theta) = \sum_{i=1}^n \left[\bar{E}_{\text{data}}(v_i^{(P)} h_j) - \bar{E}_{\text{model}}(v_i \cdot h_j) \right]$$

Too hard to calculate mean, 2^{d+k} combinations

so we estimate: first term

sample h_j for each $v_i \Rightarrow$ unbiased estimate.

second term:

we need to sample $P(v, h)$: Gibbs sampling

Repeat:

sample $\tilde{h} \sim P(h|v)$

sample $\tilde{v} \sim P(v|h)$

set $v = \tilde{v}, h = \tilde{h}$

return $P(v, h)$

- This way of estimate gradient with sampling is called Contrastive Divergence.