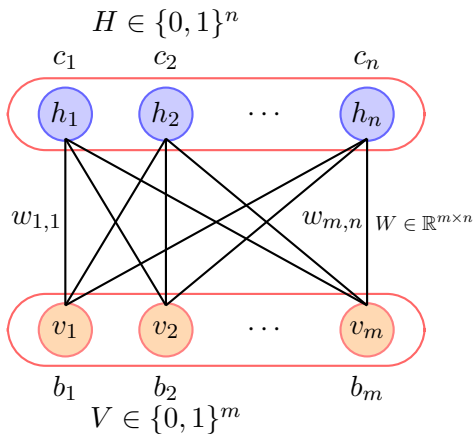


Module 19.6: Computing the gradient of the log likelihood

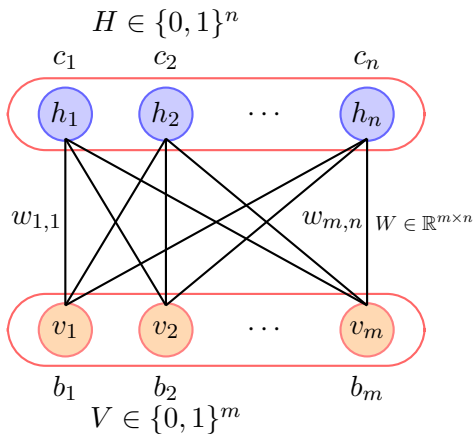
- We will just consider the loss for a single training example

$$\ln \mathcal{L}(\theta)$$



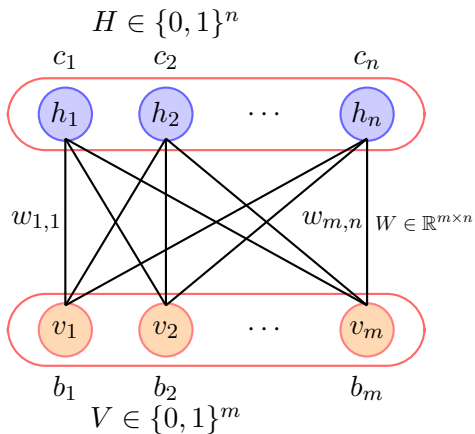
- We will just consider the loss for a single training example

$$\ln \mathcal{L}(\theta) = \ln p(V|\theta)$$



- We will just consider the loss for a single training example

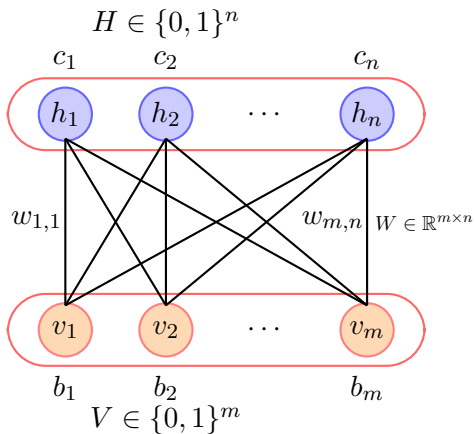
$$\ln \mathcal{L}(\theta) = \ln p(V|\theta) = \ln \frac{1}{Z} \sum_H e^{-E(V,H)}$$



- We will just consider the loss for a single training example

$$\ln \mathcal{L}(\theta) = \ln p(V|\theta) = \ln \frac{1}{Z} \sum_H e^{-E(V,H)}$$

$$= \ln \sum_H e^{-E(V,H)} - \ln \sum_{V,H} e^{-E(V,H)}$$

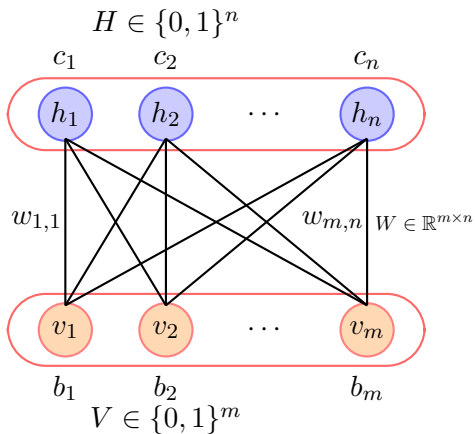


- We will just consider the loss for a single training example

$$\ln \mathcal{L}(\theta) = \ln p(V|\theta) = \ln \frac{1}{Z} \sum_H e^{-E(V,H)}$$

$$= \ln \sum_H e^{-E(V,H)} - \ln \sum_{V,H} e^{-E(V,H)}$$

$$\frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta}$$

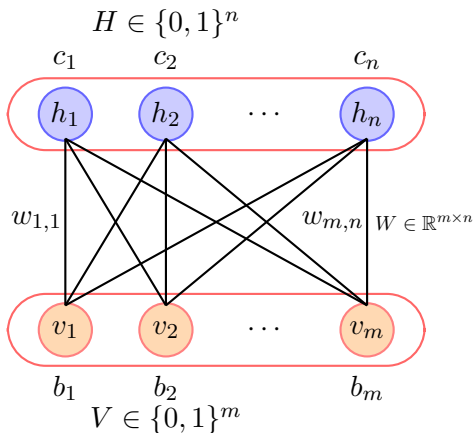


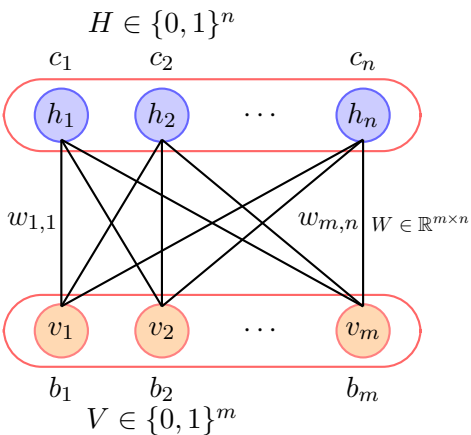
- We will just consider the loss for a single training example

$$\ln \mathcal{L}(\theta) = \ln p(V|\theta) = \ln \frac{1}{Z} \sum_H e^{-E(V,H)}$$

$$= \ln \sum_H e^{-E(V,H)} - \ln \sum_{V,H} e^{-E(V,H)}$$

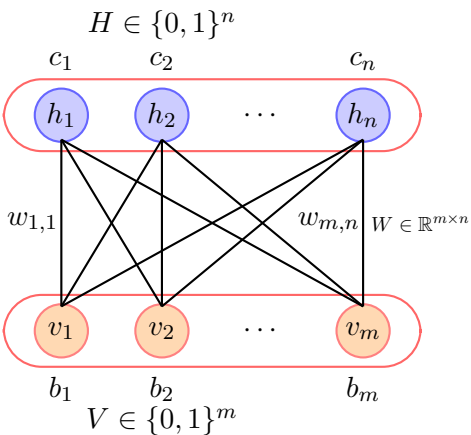
$$\frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} \left(\ln \sum_H e^{-E(V,H)} - \ln \sum_{V,H} e^{-E(V,H)} \right)$$





- We will just consider the loss for a single training example

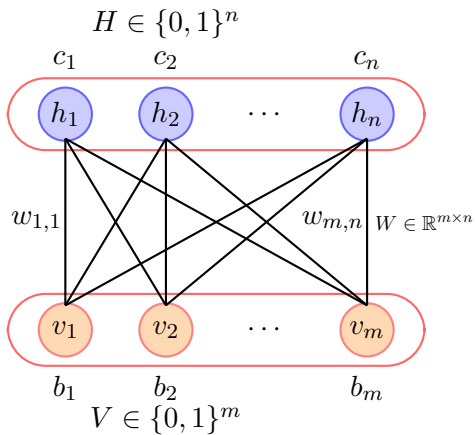
$$\begin{aligned}
 \ln \mathcal{L}(\theta) &= \ln p(V|\theta) = \ln \frac{1}{Z} \sum_H e^{-E(V,H)} \\
 &= \ln \sum_H e^{-E(V,H)} - \ln \sum_{V,H} e^{-E(V,H)} \\
 \frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \left(\ln \sum_H e^{-E(V,H)} - \ln \sum_{V,H} e^{-E(V,H)} \right) \\
 &= -\frac{1}{\sum_H e^{-E(V,H)}} \sum_H e^{-E(V,H)} \frac{\partial E(V,H)}{\partial \theta} \\
 &\quad + \frac{1}{\sum_{V,H} e^{-E(V,H)}} \sum_{V,H} e^{-E(V,H)} \frac{\partial E(V,H)}{\partial \theta}
 \end{aligned}$$



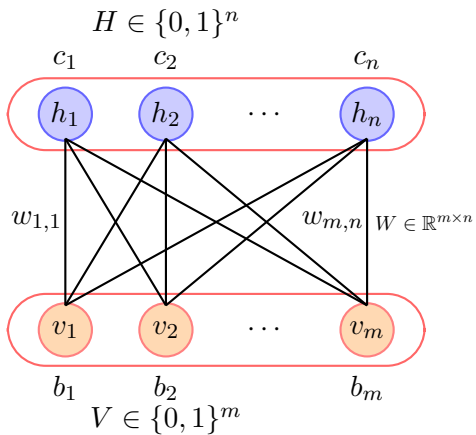
- We will just consider the loss for a single training example

$$\begin{aligned}
 \ln \mathcal{L}(\theta) &= \ln p(V|\theta) = \ln \frac{1}{Z} \sum_H e^{-E(V,H)} \\
 &= \ln \sum_H e^{-E(V,H)} - \ln \sum_{V,H} e^{-E(V,H)} \\
 \frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \left(\ln \sum_H e^{-E(V,H)} - \ln \sum_{V,H} e^{-E(V,H)} \right) \\
 &= -\frac{1}{\sum_H e^{-E(V,H)}} \sum_H e^{-E(V,H)} \frac{\partial E(V,H)}{\partial \theta} \\
 &\quad + \frac{1}{\sum_{V,H} e^{-E(V,H)}} \sum_{V,H} e^{-E(V,H)} \frac{\partial E(V,H)}{\partial \theta} \\
 &= -\sum_H \frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}} \frac{\partial E(V,H)}{\partial \theta} \\
 &\quad + \sum_{V,H} \frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} \frac{\partial E(V,H)}{\partial \theta}
 \end{aligned}$$

- Now,

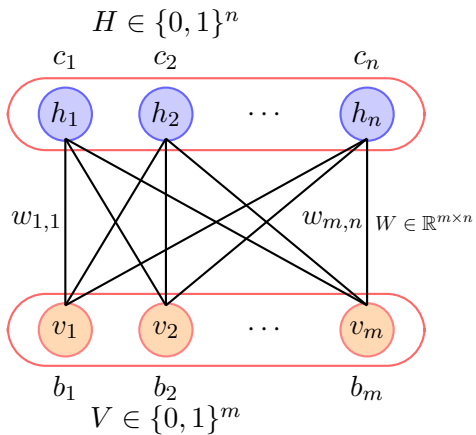


• Now,



$$\frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} = p(V,H)$$

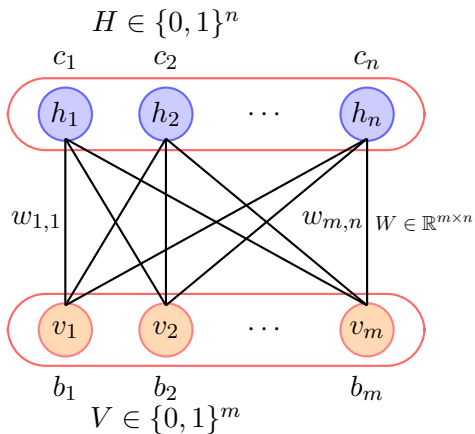
• Now,



$$\frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} = p(V,H)$$

$$\frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}}$$

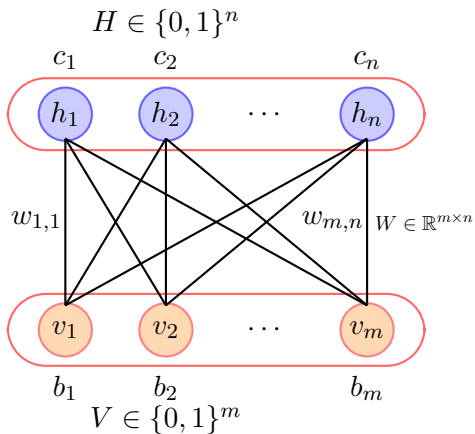
• Now,



$$\frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} = p(V,H)$$

$$\frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}} = \frac{\frac{1}{Z} e^{-E(V,H)}}{\frac{1}{Z} \sum_H e^{-E(V,H)}}$$

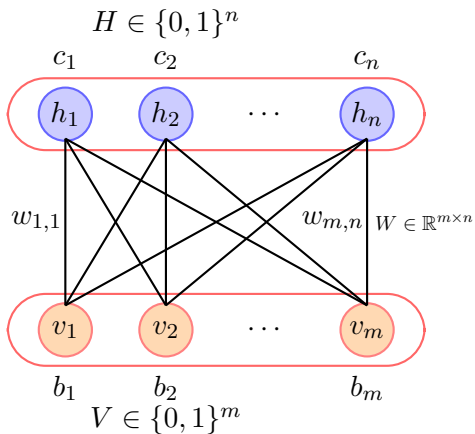
• Now,



$$\frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} = p(V, H)$$

$$\begin{aligned} \frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}} &= \frac{\frac{1}{Z} e^{-E(V,H)}}{\frac{1}{Z} \sum_H e^{-E(V,H)}} \\ &= \frac{p(V, H)}{p(V)} \end{aligned}$$

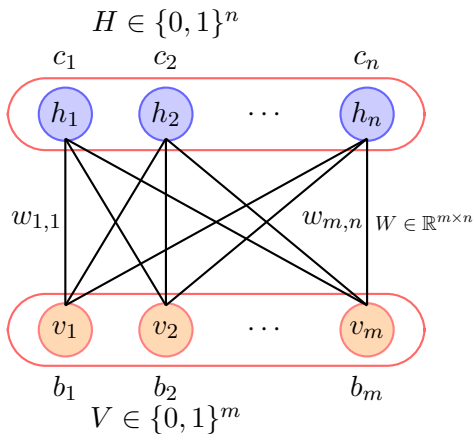
• Now,



$$\frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} = p(V,H)$$

$$\begin{aligned} \frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}} &= \frac{\frac{1}{Z} e^{-E(V,H)}}{\frac{1}{Z} \sum_H e^{-E(V,H)}} \\ &= \frac{p(V,H)}{p(V)} = p(H|V) \end{aligned}$$

• Now,

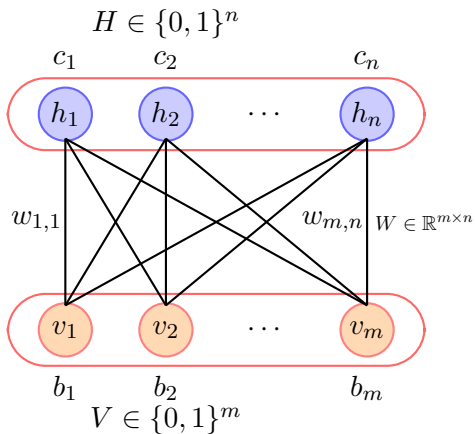


$$\frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} = p(V, H)$$

$$\begin{aligned} \frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}} &= \frac{\frac{1}{Z} e^{-E(V,H)}}{\frac{1}{Z} \sum_H e^{-E(V,H)}} \\ &= \frac{p(V, H)}{p(V)} = p(H|V) \end{aligned}$$

$$\frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta}$$

• Now,

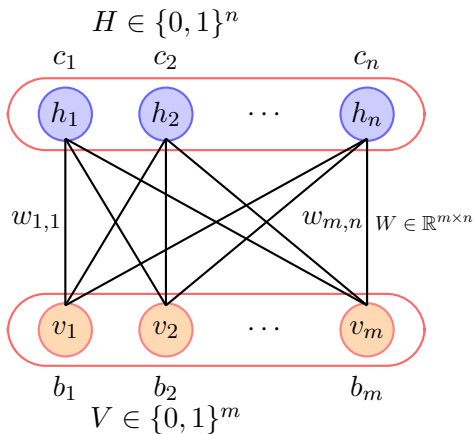


$$\frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} = p(V,H)$$

$$\begin{aligned} \frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}} &= \frac{\frac{1}{Z} e^{-E(V,H)}}{\frac{1}{Z} \sum_H e^{-E(V,H)}} \\ &= \frac{p(V,H)}{p(V)} = p(H|V) \end{aligned}$$

$$\begin{aligned} \frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta} &= - \sum_H \frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}} \frac{\partial E(V,H)}{\partial \theta} \\ &\quad + \sum_{V,H} \frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} \frac{\partial E(V,H)}{\partial \theta} \end{aligned}$$

• Now,

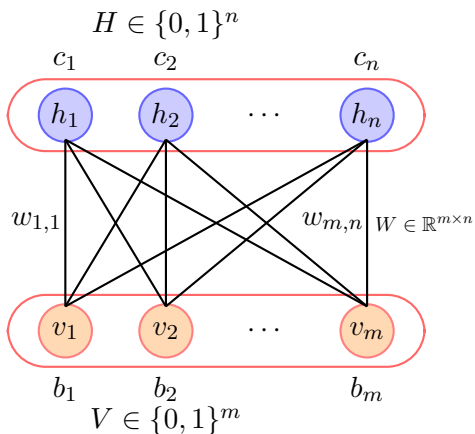


$$\frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} = p(V,H)$$

$$\begin{aligned} \frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}} &= \frac{\frac{1}{Z} e^{-E(V,H)}}{\frac{1}{Z} \sum_H e^{-E(V,H)}} \\ &= \frac{p(V,H)}{p(V)} = p(H|V) \end{aligned}$$

$$\begin{aligned} \frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta} &= - \sum_H \frac{e^{-E(V,H)}}{\sum_H e^{-E(V,H)}} \frac{\partial E(V,H)}{\partial \theta} \\ &\quad + \sum_{V,H} \frac{e^{-E(V,H)}}{\sum_{V,H} e^{-E(V,H)}} \frac{\partial E(V,H)}{\partial \theta} \\ &= - \sum_H p(H|V) \frac{\partial E(V,H)}{\partial \theta} + \sum_{V,H} p(V,H) \frac{\partial E(V,H)}{\partial \theta} \end{aligned}$$

- Okay, so we have,

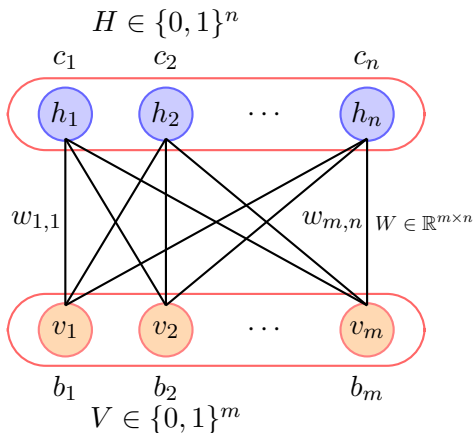


$$\begin{aligned} \frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta} = & - \sum_H p(H|V) \frac{\partial E(V, H)}{\partial \theta} \\ & + \sum_{V, H} p(V, H) \frac{\partial E(V, H)}{\partial \theta} \end{aligned}$$

- Okay, so we have,

$$\frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta} = - \sum_H p(H|V) \frac{\partial E(V, H)}{\partial \theta} + \sum_{V, H} p(V, H) \frac{\partial E(V, H)}{\partial \theta}$$

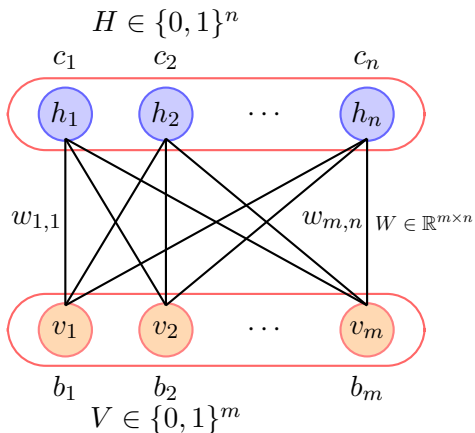
- Remember that θ is a collection of all the parameters in our model, i.e., $W_{ij}, b_i, c_j \forall i \in \{1, \dots, m\}$ and $\forall j \in \{1, \dots, n\}$

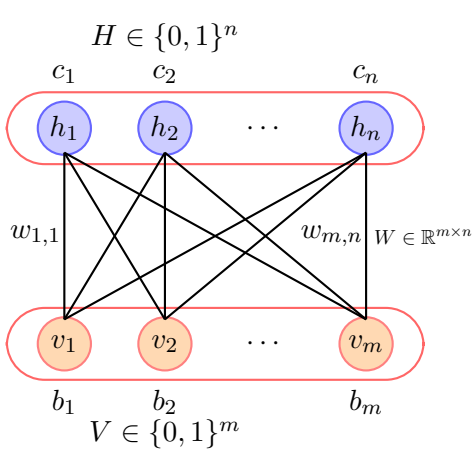


- Okay, so we have,

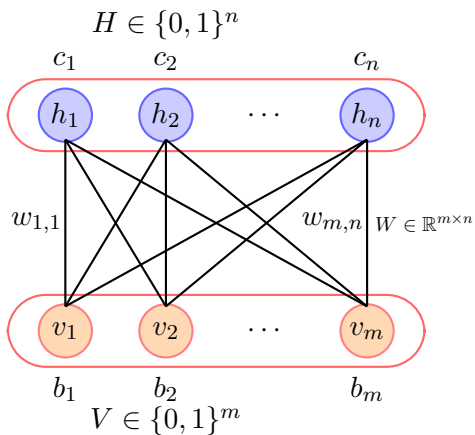
$$\frac{\partial \ln \mathcal{L}(\theta)}{\partial \theta} = - \sum_H p(H|V) \frac{\partial E(V, H)}{\partial \theta} + \sum_{V, H} p(V, H) \frac{\partial E(V, H)}{\partial \theta}$$

- Remember that θ is a collection of all the parameters in our model, i.e., $W_{ij}, b_i, c_j \forall i \in \{1, \dots, m\}$ and $\forall j \in \{1, \dots, n\}$
- We will follow our usual recipe of computing the partial derivative w.r.t. one weight w_{ij} and then generalize to the gradient w.r.t. the entire weight matrix W

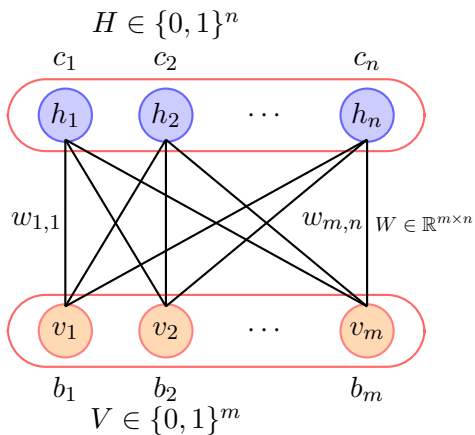




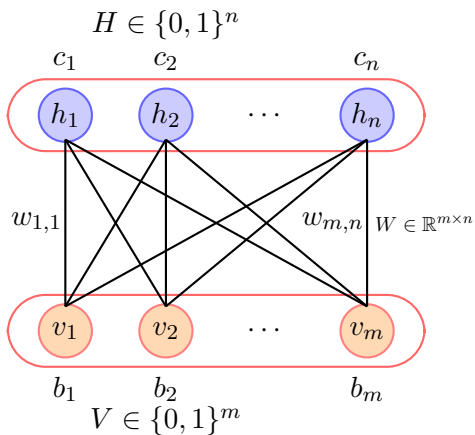
$$\frac{\partial \mathcal{L}(\theta)}{\partial w_{ij}}$$



$$\frac{\partial \mathcal{L}(\theta)}{\partial w_{ij}} = - \sum_H p(H|V) \frac{\partial E(V, H)}{\partial w_{ij}} + \sum_{V, H} p(V, H) \frac{\partial E(V, H)}{\partial w_{ij}}$$

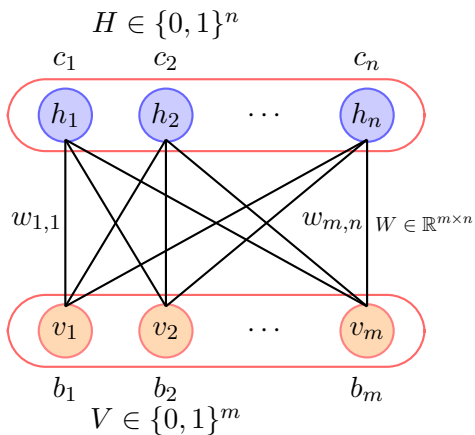


$$\begin{aligned}
 & \frac{\partial \mathcal{L}(\theta)}{\partial w_{ij}} \\
 &= - \sum_H p(H|V) \frac{\partial E(V, H)}{\partial w_{ij}} + \sum_{V, H} p(V, H) \frac{\partial E(V, H)}{\partial w_{ij}} \\
 &= \sum_H p(H|V) h_i v_j - \sum_{V, H} p(V, H) h_i v_j
 \end{aligned}$$



$$\begin{aligned}
 & \frac{\partial \mathcal{L}(\theta)}{\partial w_{ij}} \\
 &= - \sum_H p(H|V) \frac{\partial E(V, H)}{\partial w_{ij}} + \sum_{V, H} p(V, H) \frac{\partial E(V, H)}{\partial w_{ij}} \\
 &= \sum_H p(H|V) h_i v_j - \sum_{V, H} p(V, H) h_i v_j
 \end{aligned}$$

- We can write the above as a sum of two expectations



$$\begin{aligned}
 & \frac{\partial \mathcal{L}(\theta)}{\partial w_{ij}} \\
 &= - \sum_H p(H|V) \frac{\partial E(V, H)}{\partial w_{ij}} + \sum_{V, H} p(V, H) \frac{\partial E(V, H)}{\partial w_{ij}} \\
 &= \sum_H p(H|V) h_i v_j - \sum_{V, H} p(V, H) h_i v_j \\
 &= \mathbb{E}_{p(H|V)}[v_i h_j] - \mathbb{E}_{p(V, H)}[v_i h_j]
 \end{aligned}$$

- We can write the above as a sum of two expectations

- How do we compute these expectations?

$$\frac{\partial \mathcal{L}(\theta)}{\partial w_{ij}} = \mathbb{E}_{p(H|V)}[v_i h_j] - \mathbb{E}_{p(V,H)}[v_i h_j]$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial w_{ij}} = \mathbb{E}_{p(H|V)}[v_i h_j] - \mathbb{E}_{p(V,H)}[v_i h_j]$$

- How do we compute these expectations?
- The first summation can actually be simplified (we will come back and simplify it later)

$$\frac{\partial \mathcal{L}(\theta)}{\partial w_{ij}} = \mathbb{E}_{p(H|V)}[v_i h_j] - \mathbb{E}_{p(V,H)}[v_i h_j]$$

- How do we compute these expectations?
- The first summation can actually be simplified (we will come back and simplify it later)
- However, the second summation contains an exponential number of terms and hence intractable in practice

$$\frac{\partial \mathcal{L}(\theta)}{\partial w_{ij}} = \mathbb{E}_{p(H|V)}[v_i h_j] - \mathbb{E}_{p(V,H)}[v_i h_j]$$

- How do we compute these expectations?
- The first summation can actually be simplified (we will come back and simplify it later)
- However, the second summation contains an exponential number of terms and hence intractable in practice
- So how do we deal with this ?