

Restricted Boltzmann Machines for Collaborative Filtering

Ruslan Salakhutdinov Andriy Mnih Geoffrey Hinton

November 29, 2016

- Wikipedia:
 - In the newer, narrower sense, collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating).

- Wikipedia:
 - In the newer, narrower sense, collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating).
- Fundamental ideas:
 - If two items get similar rating patterns then they are probably similar

- Wikipedia:
 - In the newer, narrower sense, collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating).
- Fundamental ideas:
 - If two items get similar rating patterns then they are probably similar
 - If two users rated items in a similar fashion, then they will probably give similar ratings to an unrated item

- Wikipedia:
 - In the newer, narrower sense, collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating).
- Fundamental ideas:
 - If two items get similar rating patterns then they are probably similar
 - If two users rated items in a similar fashion, then they will probably give similar ratings to an unrated item
 - Properties of items are unknown

- Wikipedia:
 - In the newer, narrower sense, collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating).
- Fundamental ideas:
 - If two items get similar rating patterns then they are probably similar
 - If two users rated items in a similar fashion, then they will probably give similar ratings to an unrated item
 - Properties of items are unknown
- Applications:
 - Amazon (Customers Who Bought This Item Also Bought)

- Wikipedia:
 - In the newer, narrower sense, collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating).
- Fundamental ideas:
 - If two items get similar rating patterns then they are probably similar
 - If two users rated items in a similar fashion, then they will probably give similar ratings to an unrated item
 - Properties of items are unknown
- Applications:
 - Amazon (Customers Who Bought This Item Also Bought)
 - Netflix

- Wikipedia:
 - In the newer, narrower sense, collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating).
- Fundamental ideas:
 - If two items get similar rating patterns then they are probably similar
 - If two users rated items in a similar fashion, then they will probably give similar ratings to an unrated item
 - Properties of items are unknown
- Applications:
 - Amazon (Customers Who Bought This Item Also Bought)
 - Netflix
 - Spotify

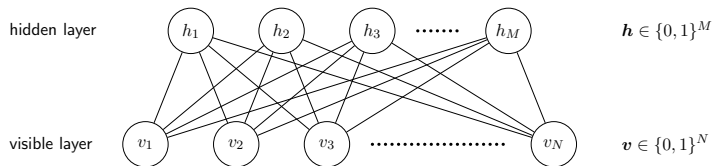
- Basic statistics:
 - 480,189 users, $\langle \text{CustomerID} \rangle$
 - 17,770 movies, $\langle \text{MovieID}, \text{YearOfRelease}, \text{Title} \rangle$

- Basic statistics:
 - 480,189 users, $\langle \text{CustomerID} \rangle$
 - 17,770 movies, $\langle \text{MovieID}, \text{YearOfRelease}, \text{Title} \rangle$
- Training set:
 - 100,480,507 ratings, $\langle \text{CustomerID}, \text{MovieID}, \text{Rating}, \text{Date} \rangle$
 - Rating is an integer that ranges from 1 to 5

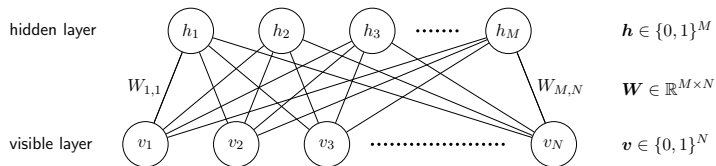
- Basic statistics:
 - 480,189 users, $\langle \text{CustomerID} \rangle$
 - 17,770 movies, $\langle \text{MovieID}, \text{YearOfRelease}, \text{Title} \rangle$
- Training set:
 - 100,480,507 ratings, $\langle \text{CustomerID}, \text{MovieID}, \text{Rating}, \text{Date} \rangle$
 - Rating is an integer that ranges from 1 to 5
- Qualifying set:
 - 2,817,131 ratings, $\langle \text{CustomerID}, \text{MovieID}, \text{Date} \rangle$

- Basic statistics:
 - 480,189 users, $\langle \text{CustomerID} \rangle$
 - 17,770 movies, $\langle \text{MovieID}, \text{YearOfRelease}, \text{Title} \rangle$
- Training set:
 - 100,480,507 ratings, $\langle \text{CustomerID}, \text{MovieID}, \text{Rating}, \text{Date} \rangle$
 - Rating is an integer that ranges from 1 to 5
- Qualifying set:
 - 2,817,131 ratings, $\langle \text{CustomerID}, \text{MovieID}, \text{Date} \rangle$
- Sparsity of ratings:
 - $$\frac{100,480,507 + 2,817,131}{480,189 \times 17,770} = 0.0121 = 1.21\%$$

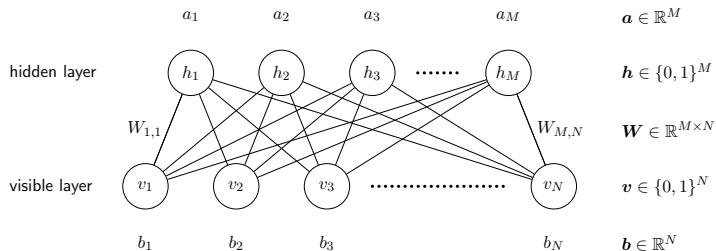
Binary RBM



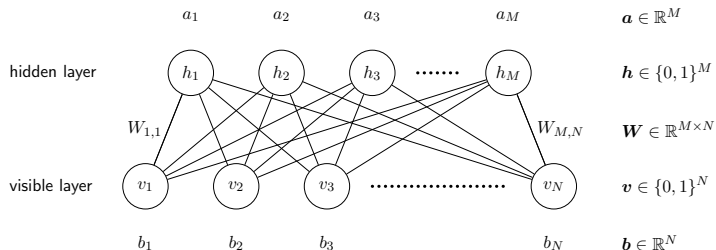
Binary RBM



Binary RBM

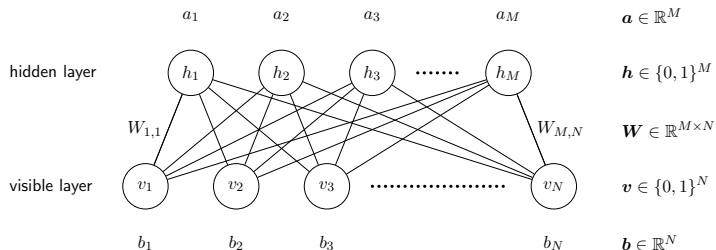


Binary RBM



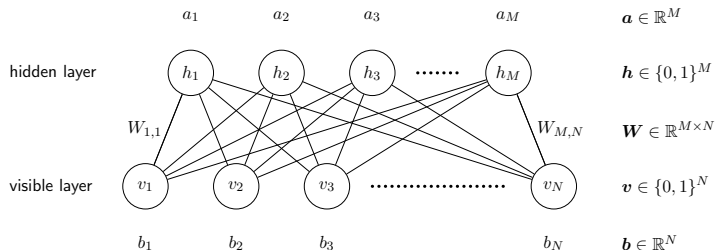
- $E(\mathbf{h}, \mathbf{v}) = -\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$

Binary RBM



- $E(\mathbf{h}, \mathbf{v}) = -\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$
- $p(\mathbf{h}, \mathbf{v}) = \frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$

Binary RBM



- $E(\mathbf{h}, \mathbf{v}) = -\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$
- $p(\mathbf{h}, \mathbf{v}) = \frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$
- $Z = \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$

$$p(\mathbf{h}|\mathbf{v})$$

$$p(\mathbf{h}|\mathbf{v}) = \frac{p(\mathbf{h}, \mathbf{v})}{p(\mathbf{v})}$$

$$\begin{aligned} p(\mathbf{h}|\mathbf{v}) &= \frac{p(\mathbf{h}, \mathbf{v})}{p(\mathbf{v})} \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{b}^T \mathbf{v} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \end{aligned}$$

$$\begin{aligned} p(\mathbf{h}|\mathbf{v}) &= \frac{p(\mathbf{h}, \mathbf{v})}{p(\mathbf{v})} \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{b}^T \mathbf{v} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{b}^T \mathbf{v}) \exp(\mathbf{a}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \end{aligned}$$

$$\begin{aligned} p(\mathbf{h}|\mathbf{v}) &= \frac{p(\mathbf{h}, \mathbf{v})}{p(\mathbf{v})} \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{b}^T \mathbf{v} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{b}^T \mathbf{v}) \exp(\mathbf{a}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{Z'} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \end{aligned}$$

$$\begin{aligned} p(\mathbf{h}|\mathbf{v}) &= \frac{p(\mathbf{h}, \mathbf{v})}{p(\mathbf{v})} \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{b}^T \mathbf{v} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{b}^T \mathbf{v}) \exp(\mathbf{a}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{Z'} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{Z'} \exp \left(\sum_{i=1}^M a_i h_i + \sum_{i=1}^M h_i \mathbf{W}_{i,:} \mathbf{v} \right) \end{aligned}$$

$$\begin{aligned} p(\mathbf{h}|\mathbf{v}) &= \frac{p(\mathbf{h}, \mathbf{v})}{p(\mathbf{v})} \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{b}^T \mathbf{v} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{b}^T \mathbf{v}) \exp(\mathbf{a}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{Z'} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{Z'} \exp \left(\sum_{i=1}^M a_i h_i + \sum_{i=1}^M h_i \mathbf{W}_{i,:} \mathbf{v} \right) \\ &= \frac{1}{Z'} \exp \left(\sum_{i=1}^M h_i (a_i + \mathbf{W}_{i,:} \mathbf{v}) \right) \end{aligned}$$

$$\begin{aligned} p(\mathbf{h}|\mathbf{v}) &= \frac{p(\mathbf{h}, \mathbf{v})}{p(\mathbf{v})} \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{b}^T \mathbf{v} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{p(\mathbf{v})} \frac{1}{Z} \exp(\mathbf{b}^T \mathbf{v}) \exp(\mathbf{a}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{Z'} \exp(\mathbf{a}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{v}) \\ &= \frac{1}{Z'} \exp \left(\sum_{i=1}^M a_i h_i + \sum_{i=1}^M h_i \mathbf{W}_{i,:} \mathbf{v} \right) \\ &= \frac{1}{Z'} \exp \left(\sum_{i=1}^M h_i (a_i + \mathbf{W}_{i,:} \mathbf{v}) \right) \\ &= \frac{1}{Z'} \prod_{i=1}^M \exp(h_i (a_i + \mathbf{W}_{i,:} \mathbf{v})) \end{aligned}$$

$$p(\mathbf{h}|\mathbf{v}) = \frac{1}{Z'} \prod_{i=1}^M \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$p(\mathbf{h}|\mathbf{v}) = \frac{1}{Z'} \prod_{i=1}^M \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$p(h_i|\mathbf{v}) \propto \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$p(\mathbf{h}|\mathbf{v}) = \frac{1}{Z'} \prod_{i=1}^M \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$p(h_i|\mathbf{v}) \propto \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$p(h_i = 1|\mathbf{v}) = \frac{p(h_i = 1|\mathbf{v})}{p(h_i = 0|\mathbf{v}) + p(h_i = 1|\mathbf{v})}$$

$$p(\mathbf{h}|\mathbf{v}) = \frac{1}{Z'} \prod_{i=1}^M \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$p(h_i|\mathbf{v}) \propto \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$\begin{aligned} p(h_i = 1|\mathbf{v}) &= \frac{p(h_i = 1|\mathbf{v})}{p(h_i = 0|\mathbf{v}) + p(h_i = 1|\mathbf{v})} \\ &= \frac{\exp(a_i + \mathbf{W}_{i,:}\mathbf{v})}{1 + \exp(a_i + \mathbf{W}_{i,:}\mathbf{v})} \end{aligned}$$

$$p(\mathbf{h}|\mathbf{v}) = \frac{1}{Z'} \prod_{i=1}^M \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$p(h_i|\mathbf{v}) \propto \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$\begin{aligned} p(h_i = 1|\mathbf{v}) &= \frac{p(h_i = 1|\mathbf{v})}{p(h_i = 0|\mathbf{v}) + p(h_i = 1|\mathbf{v})} \\ &= \frac{\exp(a_i + \mathbf{W}_{i,:}\mathbf{v})}{1 + \exp(a_i + \mathbf{W}_{i,:}\mathbf{v})} \\ &= \text{sigmoid}(a_i + \mathbf{W}_{i,:}\mathbf{v}) \end{aligned}$$

$$p(\mathbf{h}|\mathbf{v}) = \frac{1}{Z'} \prod_{i=1}^M \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$p(h_i|\mathbf{v}) \propto \exp(h_i(a_i + \mathbf{W}_{i,:}\mathbf{v}))$$

$$\begin{aligned} p(h_i = 1|\mathbf{v}) &= \frac{p(h_i = 1|\mathbf{v})}{p(h_i = 0|\mathbf{v}) + p(h_i = 1|\mathbf{v})} \\ &= \frac{\exp(a_i + \mathbf{W}_{i,:}\mathbf{v})}{1 + \exp(a_i + \mathbf{W}_{i,:}\mathbf{v})} \\ &= \text{sigmoid}(a_i + \mathbf{W}_{i,:}\mathbf{v}) \end{aligned}$$

$$p(v_j = 1|\mathbf{h}) = \text{sigmoid}(b_j + \mathbf{h}^T \mathbf{W}_{:,j})$$

- Binary RBM's parameters: $\theta = \{\mathbf{W}, \mathbf{a}, \mathbf{b}\}$

- Binary RBM's parameters: $\theta = \{\mathbf{W}, \mathbf{a}, \mathbf{b}\}$
- Maximum likelihood principle: $\operatorname{argmax}_{\theta} p(\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(T)} | \theta)$

- Binary RBM's parameters: $\theta = \{\mathbf{W}, \mathbf{a}, \mathbf{b}\}$
- Maximum likelihood principle: $\operatorname{argmax}_{\theta} p(\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(T)} | \theta)$
- No closed-form solution

- Binary RBM's parameters: $\theta = \{\mathbf{W}, \mathbf{a}, \mathbf{b}\}$
- Maximum likelihood principle: $\operatorname{argmax}_{\theta} p(\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(T)} | \theta)$
- No closed-form solution
- Resort to gradient ascent

$$\ell(\theta)$$

$$\ell(\theta) = \sum_{t=1}^T \log p(\mathbf{v}^{(t)})$$

$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log p(\mathbf{v}^{(t)}) \\ &= \sum_{t=1}^T \log \sum_{\mathbf{h}} p(\mathbf{h}, \mathbf{v}^{(t)})\end{aligned}$$

$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log p(\mathbf{v}^{(t)}) \\ &= \sum_{t=1}^T \log \sum_{\mathbf{h}} p(\mathbf{h}, \mathbf{v}^{(t)}) \\ &= \sum_{t=1}^T \log \sum_{\mathbf{h}} \frac{1}{Z} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)\end{aligned}$$

$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log p(\mathbf{v}^{(t)}) \\ &= \sum_{t=1}^T \log \sum_{\mathbf{h}} p(\mathbf{h}, \mathbf{v}^{(t)}) \\ &= \sum_{t=1}^T \log \sum_{\mathbf{h}} \frac{1}{Z} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) \\ &= \sum_{t=1}^T \log \frac{1}{Z} \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)\end{aligned}$$

$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log p(\mathbf{v}^{(t)}) \\&= \sum_{t=1}^T \log \sum_{\mathbf{h}} p(\mathbf{h}, \mathbf{v}^{(t)}) \\&= \sum_{t=1}^T \log \sum_{\mathbf{h}} \frac{1}{Z} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) \\&= \sum_{t=1}^T \log \frac{1}{Z} \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) \\&= \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log Z\end{aligned}$$

$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log p(\mathbf{v}^{(t)}) \\&= \sum_{t=1}^T \log \sum_{\mathbf{h}} p(\mathbf{h}, \mathbf{v}^{(t)}) \\&= \sum_{t=1}^T \log \sum_{\mathbf{h}} \frac{1}{Z} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) \\&= \sum_{t=1}^T \log \frac{1}{Z} \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) \\&= \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log Z \\&= \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))\end{aligned}$$

Binary RBM Learning, continued

$$\ell(\theta) = \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$$

$$\frac{\partial \ell(\theta)}{\partial \beta}$$

Binary RBM Learning, continued

$$\ell(\theta) = \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$$

$$\frac{\partial \ell(\theta)}{\partial \beta} = \sum_{t=1}^T \frac{\partial \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta} - T \frac{\partial \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}$$

Binary RBM Learning, continued

$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \\ \frac{\partial \ell(\theta)}{\partial \beta} &= \sum_{t=1}^T \frac{\partial \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta} - T \frac{\partial \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta} \\ &= \sum_{t=1}^T \frac{\frac{\partial \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\frac{\partial \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}\end{aligned}$$

Binary RBM Learning, continued

$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \\ \frac{\partial \ell(\theta)}{\partial \beta} &= \sum_{t=1}^T \frac{\partial \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta} - T \frac{\partial \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta} \\ &= \sum_{t=1}^T \frac{\frac{\partial \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\frac{\partial \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \frac{\sum_{\mathbf{h}} \frac{\partial \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\sum_{\mathbf{h}, \mathbf{v}} \frac{\partial \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}\end{aligned}$$

Binary RBM Learning, continued

$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \\ \frac{\partial \ell(\theta)}{\partial \beta} &= \sum_{t=1}^T \frac{\partial \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta} - T \frac{\partial \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta} \\ &= \sum_{t=1}^T \frac{\frac{\partial \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\frac{\partial \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \frac{\sum_{\mathbf{h}} \frac{\partial \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\sum_{\mathbf{h}, \mathbf{v}} \frac{\partial \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \frac{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) \frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}\end{aligned}$$

Binary RBM Learning, continued

$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \\ \frac{\partial \ell(\theta)}{\partial \beta} &= \sum_{t=1}^T \frac{\partial \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta} - T \frac{\partial \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta} \\ &= \sum_{t=1}^T \frac{\frac{\partial \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\frac{\partial \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \frac{\sum_{\mathbf{h}} \frac{\partial \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\sum_{\mathbf{h}, \mathbf{v}} \frac{\partial \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \frac{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) \frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \sum_{\mathbf{h}} \frac{\exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} \frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta} - T \sum_{\mathbf{h}, \mathbf{v}} \frac{\exp(-E(\mathbf{h}, \mathbf{v}))}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta}\end{aligned}$$

Binary RBM Learning, continued

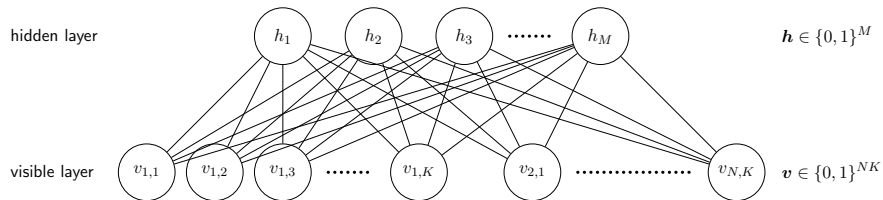
$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \\ \frac{\partial \ell(\theta)}{\partial \beta} &= \sum_{t=1}^T \frac{\partial \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta} - T \frac{\partial \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta} \\ &= \sum_{t=1}^T \frac{\frac{\partial \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\frac{\partial \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \frac{\sum_{\mathbf{h}} \frac{\partial \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\sum_{\mathbf{h}, \mathbf{v}} \frac{\partial \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \frac{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) \frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \sum_{\mathbf{h}} \frac{\exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} \frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta} - T \sum_{\mathbf{h}, \mathbf{v}} \frac{\exp(-E(\mathbf{h}, \mathbf{v}))}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta} \\ &= \sum_{t=1}^T \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(t)}) \frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta} - T \sum_{\mathbf{h}, \mathbf{v}} p(\mathbf{h}, \mathbf{v}) \frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta}\end{aligned}$$

Binary RBM Learning, continued

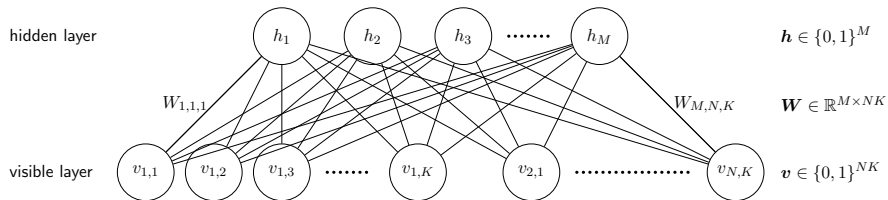
$$\begin{aligned}\ell(\theta) &= \sum_{t=1}^T \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) - T \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \\ \frac{\partial \ell(\theta)}{\partial \beta} &= \sum_{t=1}^T \frac{\partial \log \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta} - T \frac{\partial \log \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta} \\ &= \sum_{t=1}^T \frac{\frac{\partial \sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\frac{\partial \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \frac{\sum_{\mathbf{h}} \frac{\partial \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\sum_{\mathbf{h}, \mathbf{v}} \frac{\partial \exp(-E(\mathbf{h}, \mathbf{v}))}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \frac{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right) \frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta}}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} - T \frac{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v})) \frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta}}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \\ &= \sum_{t=1}^T \sum_{\mathbf{h}} \frac{\exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)}{\sum_{\mathbf{h}} \exp \left(-E(\mathbf{h}, \mathbf{v}^{(t)}) \right)} \frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta} - T \sum_{\mathbf{h}, \mathbf{v}} \frac{\exp(-E(\mathbf{h}, \mathbf{v}))}{\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))} \frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta} \\ &= \sum_{t=1}^T \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(t)}) \frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta} - T \sum_{\mathbf{h}, \mathbf{v}} p(\mathbf{h}, \mathbf{v}) \frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta} \\ &= \sum_{t=1}^T \mathbb{E}_{p(\mathbf{h}|\mathbf{v}^{(t)})} \left[\frac{\partial -E(\mathbf{h}, \mathbf{v}^{(t)})}{\partial \beta} \right] - T \mathbb{E}_{p(\mathbf{h}, \mathbf{v})} \left[\frac{\partial -E(\mathbf{h}, \mathbf{v})}{\partial \beta} \right]\end{aligned}$$

- A two layer RBM can be fully characterized by the following:
 - $E(\mathbf{h}, \mathbf{v})$
 - $p(\mathbf{h}, \mathbf{v})$
 - Z
 - $p(h_i = s | \mathbf{v})$
 - $p(v_j = t | \mathbf{h})$

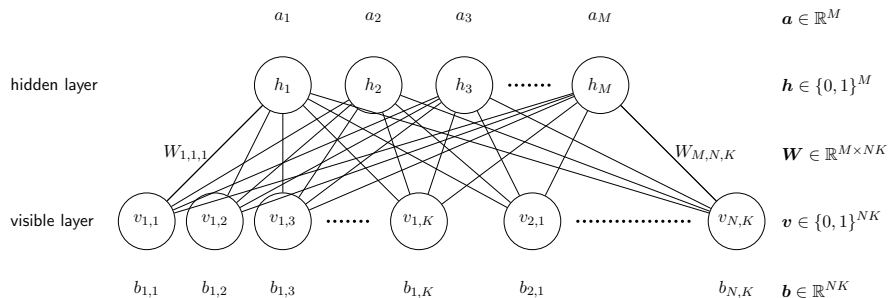
Towards RBM for CF, First Step: Make Visible Units K-nary



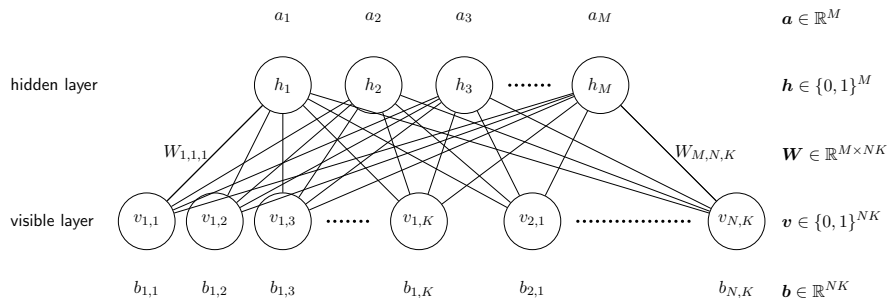
Towards RBM for CF, First Step: Make Visible Units K-nary



Towards RBM for CF, First Step: Make Visible Units K-nary



Towards RBM for CF, First Step: Make Visible Units K-nary



$v_{i,:}$	Valid Assignment?
00001	✓
01000	✓
11010	×
00000	×

Towards RBM for CF, First Step: Make Visible Units K-nary, continued

Binary

K-nary

Towards RBM for CF, First Step: Make Visible Units K-nary, continued

Binary

K-nary

$$E(\mathbf{h}, \mathbf{v}) \quad -\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v} \quad -\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v} \quad \text{K-nary requires valid } \mathbf{v}.$$

Towards RBM for CF, First Step: Make Visible Units K-nary, continued

	Binary	K-nary	
$E(\mathbf{h}, \mathbf{v})$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	K-nary requires valid \mathbf{v} .
$p(\mathbf{h}, \mathbf{v})$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	K-nary requires valid \mathbf{v} .

Towards RBM for CF, First Step: Make Visible Units K-nary, continued

	Binary	K-nary	
$E(\mathbf{h}, \mathbf{v})$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	K-nary requires valid \mathbf{v} .
$p(\mathbf{h}, \mathbf{v})$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	K-nary requires valid \mathbf{v} .
Z	$\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$	$\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$	K-nary requires valid \mathbf{v} .

Towards RBM for CF, First Step: Make Visible Units K-nary, continued

	Binary	K-nary	
$E(\mathbf{h}, \mathbf{v})$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	K-nary requires valid \mathbf{v} .
$p(\mathbf{h}, \mathbf{v})$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	K-nary requires valid \mathbf{v} .
Z	$\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$	$\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$	K-nary requires valid \mathbf{v} .
$p(h_i = 1 \mathbf{v})$	$\text{sigmoid}(a_i + \mathbf{W}_{i,:} \mathbf{v})$	$\text{sigmoid}(a_i + \mathbf{W}_{i,:} \mathbf{v})$	K-nary requires valid \mathbf{v} .

Towards RBM for CF, First Step: Make Visible Units K-nary, continued

	Binary	K-nary	
$E(\mathbf{h}, \mathbf{v})$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	K-nary requires valid \mathbf{v} .
$p(\mathbf{h}, \mathbf{v})$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	K-nary requires valid \mathbf{v} .
Z	$\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$	$\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$	K-nary requires valid \mathbf{v} .
$p(h_i = 1 \mathbf{v})$	$\text{sigmoid}(a_i + \mathbf{W}_{i,:} \mathbf{v})$	$\text{sigmoid}(a_i + \mathbf{W}_{i,:} \mathbf{v})$	K-nary requires valid \mathbf{v} .
$p(v_j = 1 \mathbf{h})$	$\text{sigmoid}(b_j + \mathbf{h}^T \mathbf{W}_{:,j})$		Undefined for K-nary.

Towards RBM for CF, First Step: Make Visible Units K-nary, continued

	Binary	K-nary	
$E(\mathbf{h}, \mathbf{v})$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	$-\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$	K-nary requires valid \mathbf{v} .
$p(\mathbf{h}, \mathbf{v})$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$	K-nary requires valid \mathbf{v} .
Z	$\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$	$\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$	K-nary requires valid \mathbf{v} .
$p(h_i = 1 \mathbf{v})$	$\text{sigmoid}(a_i + \mathbf{W}_{i,:} \mathbf{v})$	$\text{sigmoid}(a_i + \mathbf{W}_{i,:} \mathbf{v})$	K-nary requires valid \mathbf{v} .
$p(v_j = 1 \mathbf{h})$	$\text{sigmoid}(b_j + \mathbf{h}^T \mathbf{W}_{:,j})$		Undefined for K-nary.
$p(v_{j,k} = 1 \mathbf{h})$		$\frac{\exp(b_{j,k} + \mathbf{h}^T \mathbf{W}_{:,j,k})}{\sum_{l=1}^K \exp(b_{j,l} + \mathbf{h}^T \mathbf{W}_{:,j,l})}$	Undefined for Binary.

- How to handle missing values?

- How to handle missing values?
- Imputation
 - Fill in the blanks with some estimation based on available information.

- How to handle missing values?
- Imputation
 - Fill in the blanks with some estimation based on available information.
- Parameter Sharing
 - Use RBM with different numbers of visible units for different training/test cases.

Suppose $K = 3$ for the following case:

	m1	m2	m3
u1	1	?	?
u2	2	?	2
u3	1	3	?

Suppose $K = 3$ for the following case:

	m1	m2	m3
u1	1	?	?
u2	2	?	2
u3	1	3	?

	v	\tilde{v}	Visible Units
u1	100 ??? ???	100	$v_{1,1}, v_{1,2}, v_{1,3}$
u2	010 ??? 010	010 010	$v_{1,1}, v_{1,2}, v_{1,3}, v_{3,1}, v_{3,2}, v_{3,3}$
u3	100 001 ???	100 001	$v_{1,1}, v_{1,2}, v_{1,3}, v_{2,1}, v_{2,2}, v_{2,3}$

$$p(v_{q,k} = 1 | \tilde{\mathbf{v}})$$

$$p(v_{q,k} = 1 | \tilde{\mathbf{v}}) = \sum_{\mathbf{h}} p(\mathbf{h}, \tilde{\mathbf{v}}, v_{q,k})$$

$$\begin{aligned} p(v_{q,k} = 1 | \tilde{\mathbf{v}}) &= \sum_{\mathbf{h}} p(\mathbf{h}, \tilde{\mathbf{v}}, v_{q,k}) \\ &\propto \sum_{\mathbf{h}} \exp(-E(\mathbf{h}, \tilde{\mathbf{v}}, v_{q,k})) \end{aligned}$$

$$\begin{aligned} p(v_{q,k} = 1 | \tilde{\mathbf{v}}) &= \sum_{\mathbf{h}} p(\mathbf{h}, \tilde{\mathbf{v}}, v_{q,k}) \\ &\propto \sum_{\mathbf{h}} \exp(-E(\mathbf{h}, \tilde{\mathbf{v}}, v_{q,k})) \end{aligned}$$

Can be computed in polynomial time.

$$\begin{aligned} p(v_{q,k} = 1 | \tilde{\mathbf{v}}) &= \sum_{\mathbf{h}} p(\mathbf{h}, \tilde{\mathbf{v}}, v_{q,k}) \\ &\propto \sum_{\mathbf{h}} \exp(-E(\mathbf{h}, \tilde{\mathbf{v}}, v_{q,k})) \end{aligned}$$

Can be computed in polynomial time.

$$p(v_{q_1,k_1} = 1, v_{q_2,k_2} = 1, \dots, v_{q_S,k_S} = 1 | \tilde{\mathbf{v}})$$

$$\begin{aligned} p(v_{q,k} = 1 | \tilde{\mathbf{v}}) &= \sum_{\mathbf{h}} p(\mathbf{h}, \tilde{\mathbf{v}}, v_{q,k}) \\ &\propto \sum_{\mathbf{h}} \exp(-E(\mathbf{h}, \tilde{\mathbf{v}}, v_{q,k})) \end{aligned}$$

Can be computed in polynomial time.

$$p(v_{q_1,k_1} = 1, v_{q_2,k_2} = 1, \dots, v_{q_S,k_S} = 1 | \tilde{\mathbf{v}})$$

However, query of S movies on a user requires K^S evaluations.

$$\hat{h}_i = p(h_i = 1 | \tilde{\mathbf{v}}) = \text{sigmoid} \left(a_i + \widetilde{\mathbf{W}}_{i,:} \tilde{\mathbf{v}} \right)$$

$$\hat{h}_i = p(h_i = 1 | \tilde{\mathbf{v}}) = \text{sigmoid} \left(a_i + \widetilde{\mathbf{W}}_{i,:} \tilde{\mathbf{v}} \right)$$

$$p(v_{q,k} = 1 | \hat{\mathbf{h}}) = \frac{\exp \left(b_{q,k} + \hat{\mathbf{h}}^T \mathbf{W}_{:,q,k} \right)}{\sum_{l=1}^K \exp \left(b_{q,l} + \hat{\mathbf{h}}^T \mathbf{W}_{:,q,l} \right)}$$

$$\hat{h}_i = p(h_i = 1 | \tilde{\mathbf{v}}) = \text{sigmoid} \left(a_i + \widetilde{\mathbf{W}}_{i,:} \tilde{\mathbf{v}} \right)$$

$$p(v_{q,k} = 1 | \hat{\mathbf{h}}) = \frac{\exp \left(b_{q,k} + \hat{\mathbf{h}}^T \mathbf{W}_{:,q,k} \right)}{\sum_{l=1}^K \exp \left(b_{q,l} + \hat{\mathbf{h}}^T \mathbf{W}_{:,q,l} \right)}$$

A lot faster but slightly less accurate.

Extension 1: Gaussian Hidden Units

K-nary

$$E(\mathbf{h}, \mathbf{v}) \quad -\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$$

$$p(\mathbf{h}, \mathbf{v}) \quad \frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$$

$$Z \quad \sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$$

$$p(v_{j,k} = 1 | \mathbf{h}) \quad \frac{\exp(b_{j,k} + \mathbf{h}^T \mathbf{W}_{:,j,k})}{\sum_{l=1}^K \exp(b_{j,l} + \mathbf{h}^T \mathbf{W}_{:,j,l})}$$

$$p(h_i = 1 | \mathbf{v}) \quad \text{sigmoid}(a_i + \mathbf{W}_{i,:} \mathbf{v})$$

$$p(h_i = h | \mathbf{v}) \quad \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(h - a_i - \sigma_i \mathbf{W}_{i,:} \mathbf{v})^2}{2\sigma_i^2}\right)$$

K-nary with Gaussian

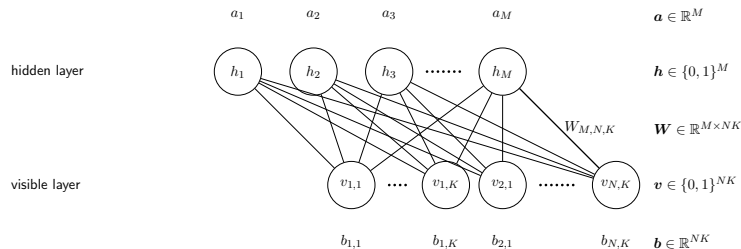
$$\sum_{i=1}^M \frac{(h_i - a_i)^2}{2\sigma_i^2} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v}$$

$$\frac{1}{Z} \exp(-E(\mathbf{h}, \mathbf{v}))$$

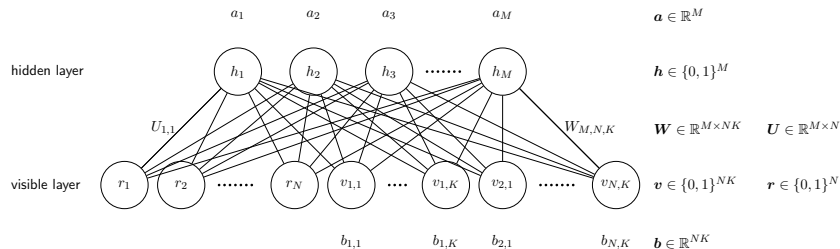
$$\sum_{\mathbf{h}, \mathbf{v}} \exp(-E(\mathbf{h}, \mathbf{v}))$$

$$\frac{\exp(b_{j,k} + \mathbf{h}^T \mathbf{W}_{:,j,k})}{\sum_{l=1}^K \exp(b_{j,l} + \mathbf{h}^T \mathbf{W}_{:,j,l})}$$

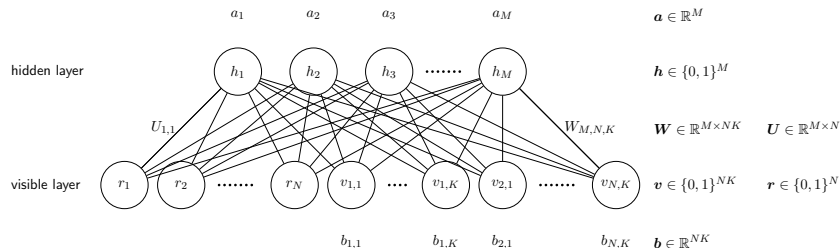
Extension 2: Condition On What Users Have Watched



Extension 2: Condition On What Users Have Watched



Extension 2: Condition On What Users Have Watched



- $E(\mathbf{h}, \mathbf{v}, \mathbf{r}) = -\mathbf{a}^T \mathbf{h} - \mathbf{b}^T \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{v} - \mathbf{h}^T \mathbf{W} \mathbf{r}$
- $p(v_{j,k} = 1 | \mathbf{h}) = \frac{\exp(b_{j,k} + \mathbf{h}^T \mathbf{W}_{:,j,k})}{\sum_{l=1}^K \exp(b_{j,l} + \mathbf{h}^T \mathbf{W}_{:,j,l})}$
- $p(h_i = 1 | \mathbf{v}, \mathbf{r}) = \text{sigmoid}(a_i + \mathbf{W}_{i,:} \mathbf{v} + \mathbf{U}_{i,:} \mathbf{r})$

Extension 3: Factorize \mathbf{W}

- $\mathbf{W} \in \mathbb{R}^{M \times NK}$

Extension 3: Factorize \mathbf{W}

- $\mathbf{W} \in \mathbb{R}^{M \times NK}$
- $100 \times 17,770 \times 5 = 8,885,000$

Extension 3: Factorize \mathbf{W}

- $\mathbf{W} \in \mathbb{R}^{M \times NK}$
- $100 \times 17,770 \times 5 = 8,885,000$
- Factorize \mathbf{W} into \mathbf{PQ}
- $\mathbf{P} \in \mathbb{R}^{M \times C}$, $\mathbf{Q} \in \mathbb{R}^{C \times NK}$
- $C \ll M$ and $C \ll NK$

Extension 3: Factorize \mathbf{W}

- $\mathbf{W} \in \mathbb{R}^{M \times NK}$
- $100 \times 17,770 \times 5 = 8,885,000$
- Factorize \mathbf{W} into \mathbf{PQ}
- $\mathbf{P} \in \mathbb{R}^{M \times C}, \mathbf{Q} \in \mathbb{R}^{C \times NK}$
- $C \ll M$ and $C \ll NK$
- $100 \times 30 + 30 \times 17,770 \times 5 = 2,668,500$

Experimental Results

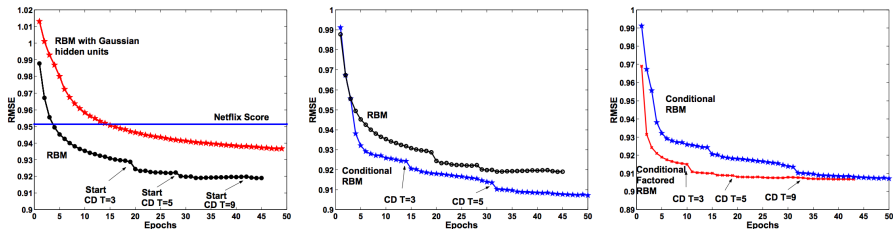


Figure 3. Performance of various models on the validation data. Left panel: RBM vs. RBM with Gaussian hidden units. Middle panel: RBM vs. conditional RBM. Right panel: conditional RBM vs. conditional factored RBM. The y-axis displays RMSE (root mean squared error), and the x-axis shows the number of epochs, or passes through the entire training dataset.

Questions?

Thank you!