

E9: 309 Advanced Deep Learning

11-11-2020



Housekeeping

✳ 1st mini-project

✓ Deadlines

- ★ Presentation on Nov19 and Nov20
- ★ Your date allocation has been finalized
- ★ Presentation and report template will be sent out this week.
- ★ Report 1 page + references and tools
- ★ Slides 4 slides for individual project and 6 slides for 2-member.



Recap of previous class



Representation learning/data-visualization

- ✳ Learning a lower dimensional representation

- ➡ Unsupervised dimensionality reductions

- ✓ Based on neighborhood preservation

- ➡ t-SNE embeddings

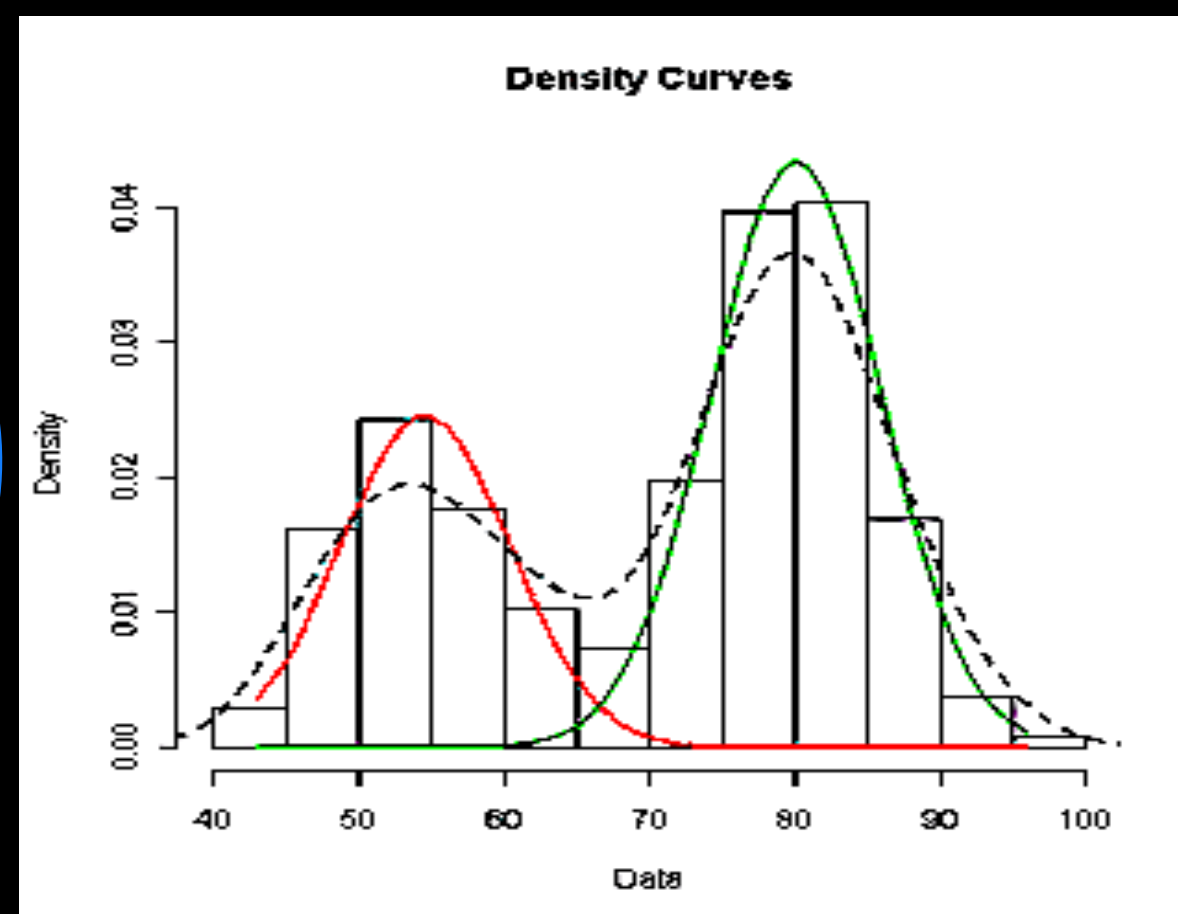
- ✓ visualization of neural network layers.



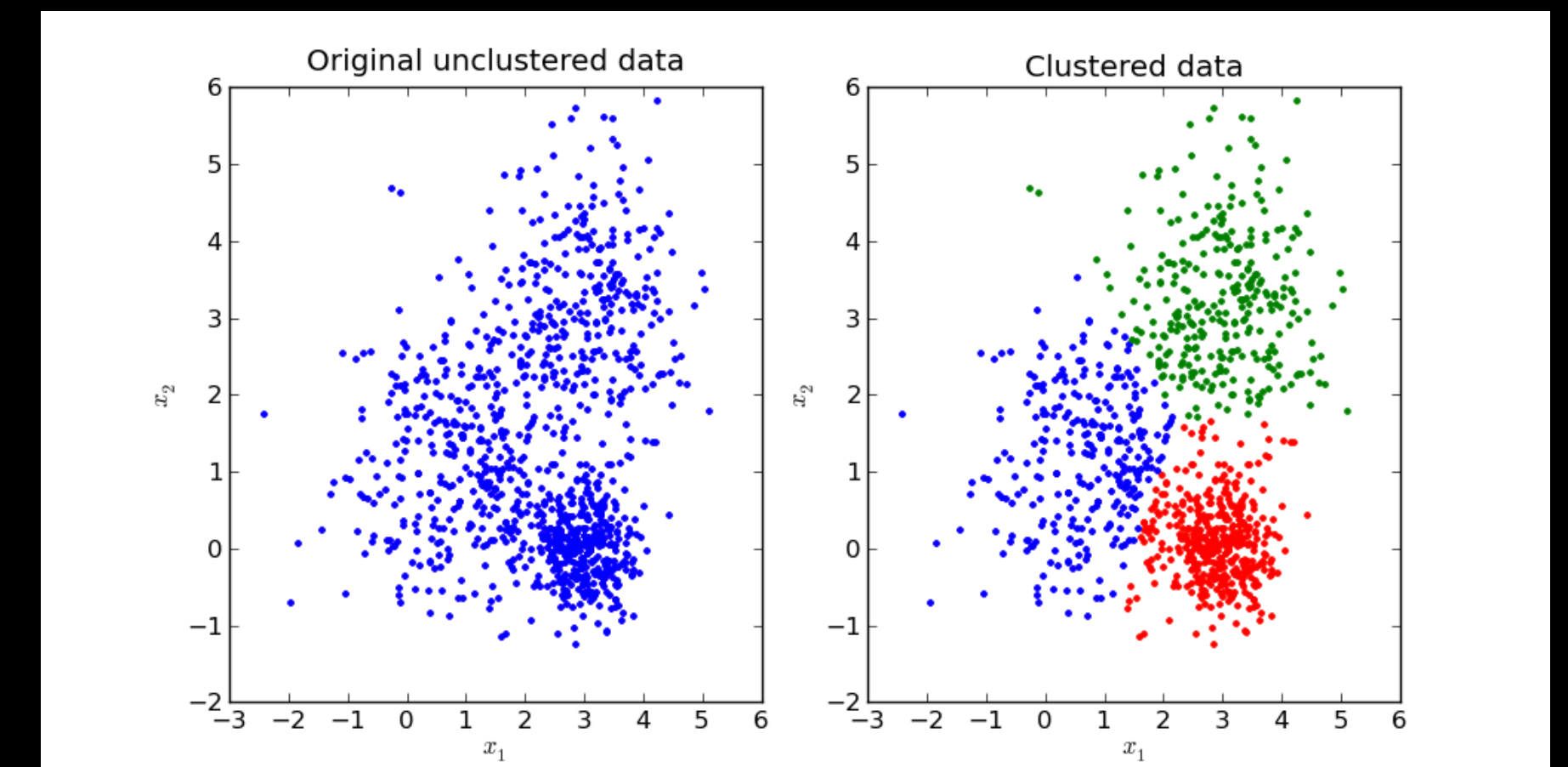
Unsupervised learning

- * Developing models that do not need labels
 - May model the generation of data.
 - May allow generation of new data samples
- * Broad strategies for unsupervised learning

Based on
maximizing
likelihood



Based on
clustering



Boltzmann machine [Deep Learning Book] (Chap 20) Binary data

- * Energy based distribution of the data

$$p(\mathbf{x}) = \frac{\exp(-E(\mathbf{x}))}{Z}$$

$$E(\mathbf{x}) = -\mathbf{x}^T \mathbf{U} \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

$$\mathbf{x} \in \mathcal{B}^{D \times 1}$$

Z - scalar
normalization
constant

- * The model parameters consists are learned by maximizing the likelihood.

- * The data can be partitioned as visible and hidden units as well

$$\mathbf{x} = [\mathbf{v}^T \mathbf{h}^T]^T$$

$$\mathbf{x} = \begin{bmatrix} \mathbf{v} \\ \mathbf{h} \end{bmatrix}$$



Restricted Boltzmann machine

- ✳ Removing the self-connections between visible and hidden units.

$$p([\mathbf{v} \ \mathbf{h}]) = \frac{\exp(-E([\mathbf{v} \ \mathbf{h}]))}{Z}$$
$$\underline{E([\mathbf{v} \ \mathbf{h}])} = \underline{-\mathbf{v}^T \mathbf{W} \mathbf{h} - b^T \mathbf{v} - \mathbf{c}^T \mathbf{h}}$$
$$Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \underline{\exp(-E([\mathbf{v} \ \mathbf{h}]))}$$

- ✳ The normalizing constant Z is called partition function.

➡ The partition function is intractable to compute explicitly.

$p(\mathbf{x})$ is not explicitly computed.

$$\mathbf{x} = \begin{bmatrix} \mathbf{v} \\ \mathbf{h} \end{bmatrix}$$

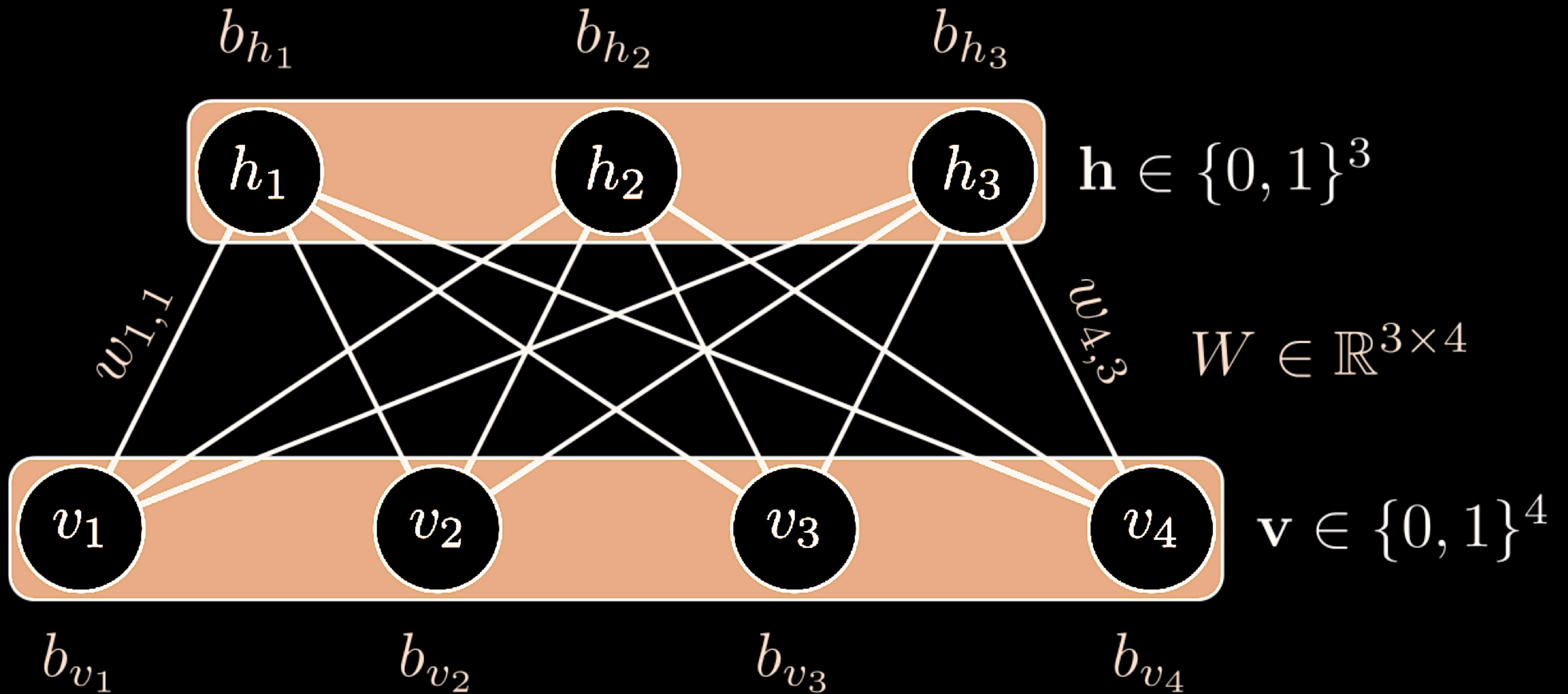
$$\mathbf{x}^T \mathbf{U} \mathbf{x}$$

$$\mathbf{U} = \begin{bmatrix} \mathbf{u}_1^T \mathbf{W} \\ \mathbf{W}^T \mathbf{u}_2 \end{bmatrix}$$

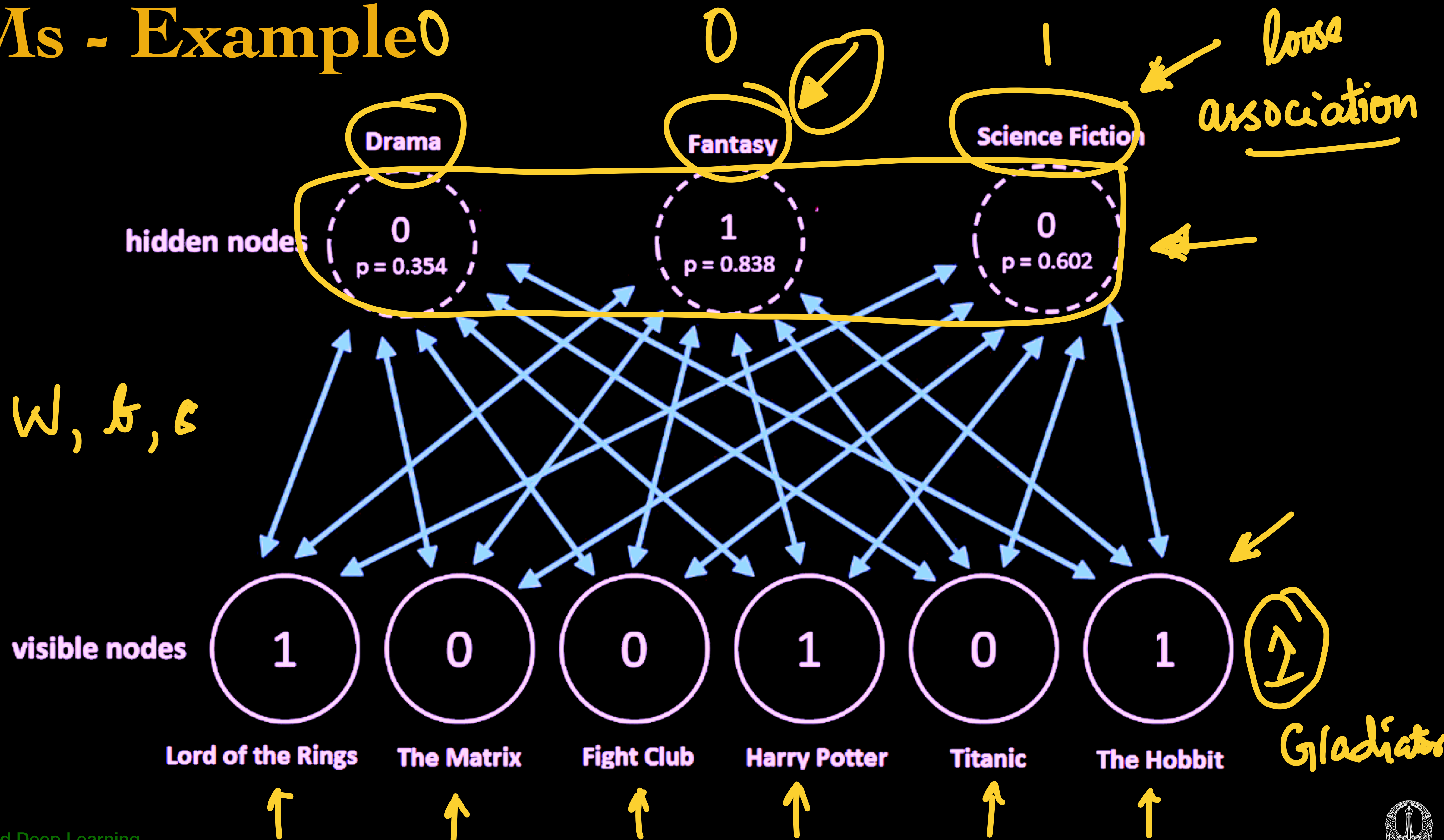
$$\cancel{\mathbf{v}^T \mathbf{u}_1 \mathbf{v}} + \cancel{\mathbf{h}^T \mathbf{u}_2 \mathbf{h}} + \underline{\underline{\mathbf{v}^T \mathbf{W} \mathbf{h}}}$$



RBM - Example



RBM - Example 0



Conditional independence

* Conditional probability of hidden given visible

visible - observations

v, h

$n_h \times 1$
 $h \in \mathbb{R}^{n_h \times 1}$

const.

constant

$$Z' = \frac{\exp(\bar{b}^T v)}{Z \cdot p(v)}$$

$$\begin{aligned} H &= \underline{h} \\ V &= \underline{v} \end{aligned}$$

$$\begin{aligned} \underline{p(h|v)} &= \frac{p(v, h)}{p(v)} \\ &= \frac{1}{Z} \frac{\exp(v^T W h + b^T v + c^T h)}{p(v)} \\ &= \frac{1}{Z'} \exp(v^T W h + c^T h) \\ &= \frac{1}{Z'} \exp\left(\sum_{j=1}^{n_h} (v^T \underline{w}_j) h_j + c_j h_j\right) \\ &= \frac{1}{Z'} \exp(\underline{v}^T \underline{w}_1 h_1 + c_1 h_1) \exp(v^T \underline{w}_2 h_2 + c_2 h_2) \times \dots \end{aligned}$$

$$p(h/v) = \frac{1}{Z_1} \exp \{ v^T w_1 h_1 + c_1 h_1 \} \times$$

$$\frac{1}{Z_2} \exp \{ \dots \} \times \dots$$

$$= \prod_{j=1}^{n_h} \underbrace{p(h_j^o/v)}$$

$$p(h_j^o/v) = \frac{\exp (v^T w_{[:j]} h_j + c_j h_j)}{Z_j}$$

$$\underline{p(h_j = 0 | v)} = \frac{1}{z_j'}$$

$$p(h_j = 1 | v) + p(h_j = 0 | v) = 1$$

$$\frac{1}{z_j'} \exp \{ v^T W_{:,j} + c_j \} + \frac{1}{z_j'} = 1.$$

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

$$z_j' = 1 + \exp \{ v^T W_{:,j} + c_j \}$$

Sigmoidal activation

✱ Conditional probability of hidden variable having a value of 1 given visible

$$p(h_j = 1 | \mathbf{v}) = \frac{1}{Z_j} \exp(c_j + \mathbf{v}^T \mathbf{W}_{:,j})$$

$$= \frac{\exp(c_j + \mathbf{v}^T \mathbf{W}_{:,j})}{1 + \exp(c_j + \mathbf{v}^T \mathbf{W}_{:,j})}$$

$$= \sigma(c_j + \mathbf{v}^T \mathbf{W}_{:,j})$$

$$p(\mathbf{h} | \mathbf{v}) = \prod_j \sigma((2h_j - 1)(c_j + \mathbf{v}^T \mathbf{W}_{:,j}))$$

✱ Note that these are probabilities.

$\sigma(x)$

$$\frac{1}{1 + e^{-x}}$$

$$\sigma(-z) = 1 - \sigma(z)$$

RBM - Training

✳ Model parameters

$$\Theta = \{W, b, c\}$$

visible bias

hidden bias

$$\mathbf{x} = [\mathbf{v}^T \ \mathbf{h}^T]^T$$

✓ Learnt by maximizing the log-likelihood $\log(p(\mathbf{x}; \Theta))$

✓ Non-convex optimization.

✳ Gradient ascent based optimization

$$\Theta^{n+1} = \Theta^n + \eta \frac{\partial \log p(\mathbf{x}; \Theta)}{\partial \Theta}$$

$$\frac{\partial \log(p(\mathbf{x}; \Theta))}{\partial \Theta} = \frac{\partial \log(\tilde{p}(\mathbf{x}; \Theta))}{\partial \Theta} - \frac{\partial \log(Z(\Theta))}{\partial \Theta}$$

$$\frac{\partial \log(Z(\Theta))}{\partial \Theta} = \frac{1}{Z} \frac{\partial Z(\Theta)}{\partial \Theta}$$

$$Z(\Theta) = \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}; \Theta)$$



$$p(x; \theta) =$$

↑

normalized

$$\tilde{p}(x; \theta)$$

$$Z(\theta)$$

$$x = \begin{bmatrix} v \\ h \end{bmatrix}$$

unnormalized

$$Z(\theta) =$$

$$\sum_x \tilde{p}(x; \theta) \exp(-E(x; \theta))$$

RBM - Training

* For exponential families

$$\tilde{p}(\mathbf{x}; \Theta) > 0 \quad \forall \mathbf{x}$$

$$Z(\Theta) = \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}; \Theta)$$

$$\frac{\partial \log(Z(\Theta))}{\partial \Theta} = \frac{1}{Z} \sum_{\mathbf{x}} \frac{\partial \tilde{p}(\mathbf{x}; \Theta)}{\partial \Theta}$$

$$\begin{aligned} \frac{\partial \log(Z(\Theta))}{\partial \Theta} &= \frac{1}{Z} \sum_{\mathbf{x}} \frac{\partial \exp(\log(\tilde{p}(\mathbf{x}; \Theta)))}{\partial \Theta} \\ \tilde{p} &= \exp(\log \tilde{p}) \\ &= \frac{1}{Z} \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}; \Theta) \left(\frac{\partial(\log(\tilde{p}(\mathbf{x}; \Theta)))}{\partial \Theta} \right) \\ \frac{\tilde{p}(\mathbf{x}; \Theta)}{Z} &= \underline{\underline{p(\mathbf{x})}} \\ &= \sum_{\mathbf{x}} \underline{\underline{p(\mathbf{x}; \Theta)}} \underline{\underline{\frac{\partial(\log(\tilde{p}(\mathbf{x}; \Theta)))}{\partial \Theta}}} \end{aligned}$$



$$\sum_x \frac{\frac{\partial}{\partial \theta} \log z(\theta)}{p(x; \theta)} \frac{\partial}{\partial \theta} \log \tilde{p}(x; \theta) = E_{x \sim p(x)} \left(\frac{\partial \log \tilde{p}(x)}{\partial \theta} \right)$$

$$\frac{\partial}{\partial \theta} \log \underline{p(x; \theta)} = \frac{\partial}{\partial \theta} \log \tilde{p}(x; \theta) - \underbrace{E_{x \sim p(x)} \left[\frac{\partial \log \tilde{p}(x; \theta)}{\partial \theta} \right]}_{\text{negative phase}}$$

positive phase

$p(x; \theta)$

RBM - Training

$$\frac{\partial}{\partial \theta} \log Z(\theta) = \underbrace{E_{x \sim p(x; \theta)} \left[\frac{\partial}{\partial \theta} \log \tilde{p}(x; \theta) \right]}$$

✱ Intractable to compute the exact gradient

➔ Using approximations to expectations *(computationally simple)*

✓ Based on sampling methods.

★ Monte-carlo Markov Chain (MCMC) based approximation

★ Resorting to Gibbs sampling.



Next lecture

→ Gibbs sampling

→ Gaussian RBM $U \in \mathbb{R}^{n_v \times 1}$

→ Gaussian mixture model

→ Application

→ Deep Belief Network

$$E[f(x)]$$

$$x \sim p(x)$$

$$\approx \frac{1}{n} \sum_{i=1}^n f(x_i)$$

$$x_i \sim p(x)$$