

A Deep Architecture for Log-Linear Models

Simon Luo¹², Sally Cripps¹² and Mahito Sugiyama³⁴

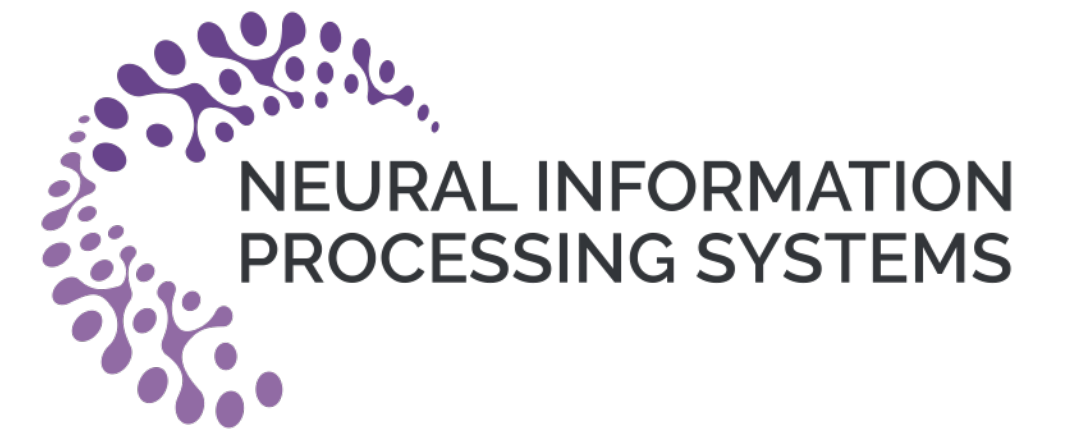
¹ School of Mathematics and Statistics, The University of Sydney

² Data Analytics for Resources and Environments (DARE), Australian Research Council

³ National Institute of Informatics

⁴ JST, PRESTO

Thirty-fourth Annual Conference on Neural Information Processing Systems (NeurIPS 2020) Workshop on Deep Learning through Information Geometry



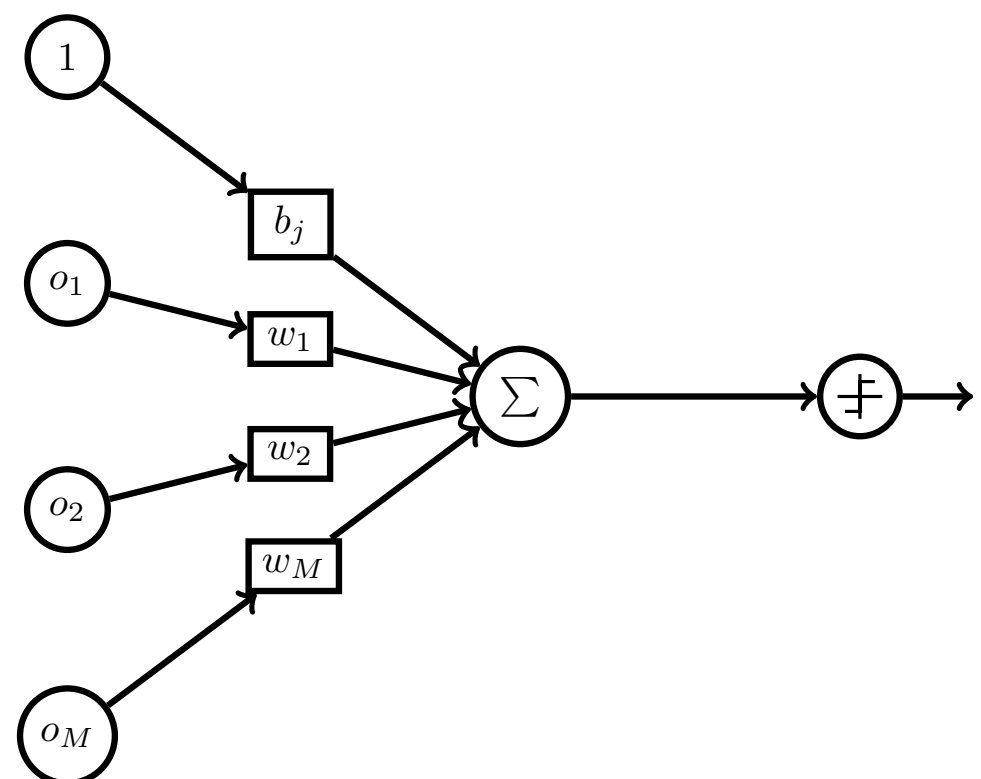
Summary

- Formulated a deep learning architecture using a partial order structure
- Bias and edge weights are realized on different layers
- Does not require gradients to update the parameters
- Minimizes Kullback-Leibler Divergence from a set of samples to our poset
- Uses statistical EM-Algorithm (Expectation-Maximization) for optimization
- Closed form formulae for both E-step and M-step

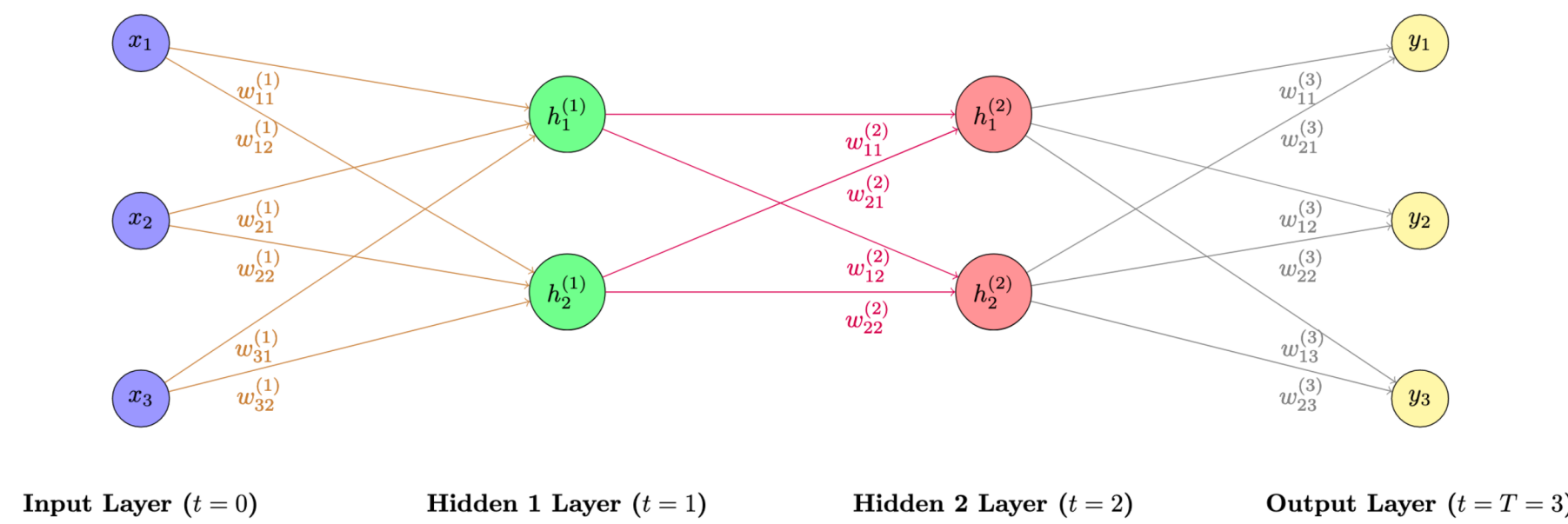
Perceptron

$$u_k^{(t)} = \sum_{i=1}^{M^{(t-1)}} w_{ij}^{(t-1)} o_i^{(t-1)} - b_j^{(t)} = \sum_{i=1}^{M^{(t-1)}} w_{ij}^{(t-1)} \sigma(u_i^{(t-1)}) - b_j^{(t)}$$

Inputs ($t-1$) Integrated Input (t) Output (t)
Edge Weights ($t-1$)



Neural Network



Log-Linear Model on a Partially Ordered Set (poset)

The log-linear model is defined over a partial order set (poset) (S, \preceq)

Dual coordinate system (θ, η) of a statistical manifold

• η Expectation parameter

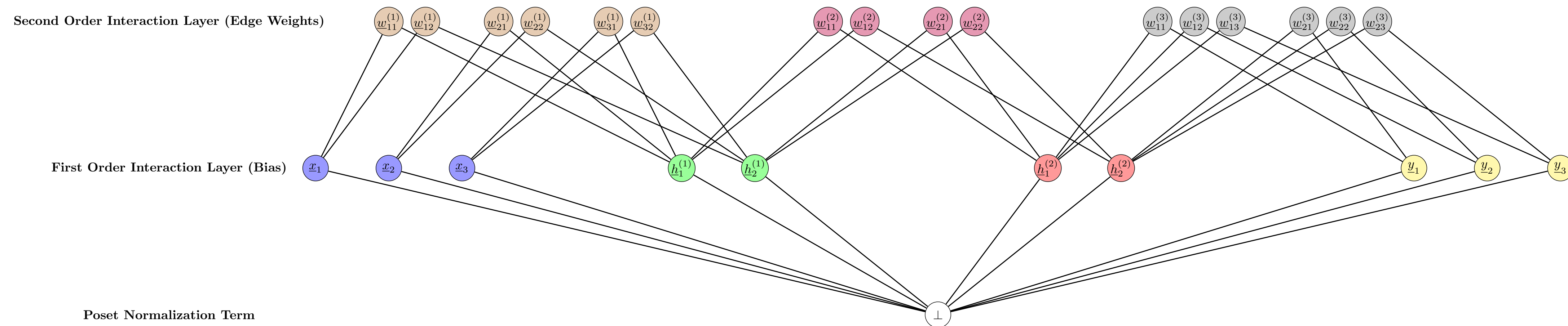
• θ Natural parameter in the exponential family

$$\eta_\omega = \sum_{s \in \Omega} \mathbf{1}_{s \preceq \omega} p(s), \quad \log p(\omega) = \sum_{s \in \Omega^+} \mathbf{1}_{s \preceq \omega} \theta_s - \psi(\theta).$$

All parameters in the model are connected via the partition function given by

$$\psi(\theta) = \log \sum_{\omega' \in \Omega^+} \exp \left[\sum_{\omega \in \Omega^+} \mathbf{1}_{\omega \preceq \omega'} \theta(\omega) \right] = -\theta(\perp).$$

Neural Network as a Partial Orders Structure



Representing Parameters as Partial Orders

Inputs: $\underline{x} \in \mathcal{X}$, Outputs: $\underline{y} \in \mathcal{Y}$

Hidden Nodes: $\underline{h} \in \mathcal{H}$, Edge Weights: $\underline{w} \in \mathcal{W}$

Integrated input: $u_k^{(t)} = \theta(\underline{n}_k)$
Edge Weight: $w_{ij}^{(t)} = \theta(\underline{n}_i^{(t-1)}) + \theta(\underline{w}_{ij}^{(t)}) + \theta(\underline{n}_j^{(t)})$

Representing Inputs and Outputs

$$\hat{\eta}(\underline{x}_k) = \frac{\exp(\mathbb{E}[x_k])}{\sum_i \exp(\mathbb{E}[x_i])}, \quad \hat{\eta}(\underline{y}_k) = \frac{\exp(\mathbb{E}[y_k])}{\sum_i \exp(\mathbb{E}[y_i])}.$$

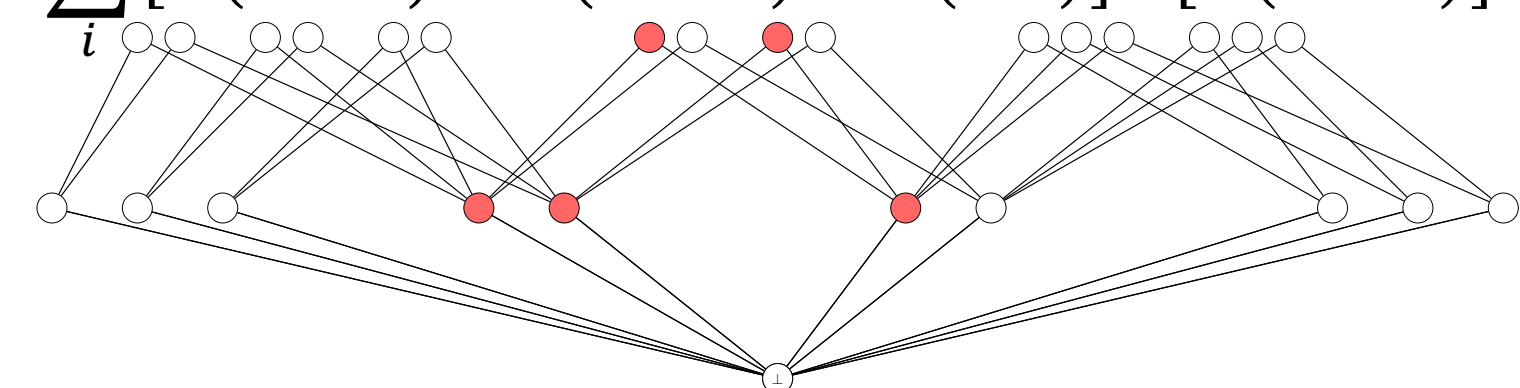
Expectation Step

Update probabilities for the input and output

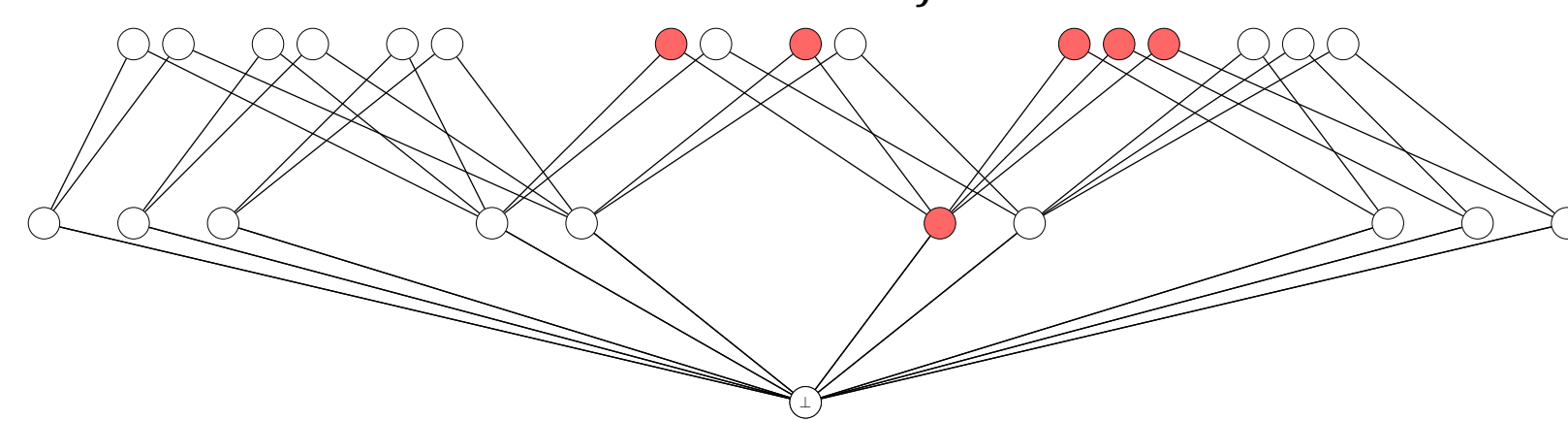
$$\hat{p}(\underline{x}_k) = \frac{\hat{\eta}(\underline{x}_k)}{[1 + \sum_j \exp[\hat{\theta}(\underline{w}_{kj}) + \hat{\theta}(\underline{h}_j)]]}, \quad \hat{p}(\underline{y}_k) = \frac{\hat{\eta}(\underline{y}_k)}{[1 + \sum_j \exp[\hat{\theta}(\underline{w}_{kj}) + \hat{\theta}(\underline{h}_j)]]}$$

Forward Propagation on a Poset

$$u_k^{(t)} = \sum_i \left[\hat{\theta}(\underline{n}_i^{(t-1)}) + \hat{\theta}(\underline{w}_{ik}^{(t-1)}) + \hat{\theta}(\underline{n}_k^{(t)}) \right] \sigma \left[\hat{\theta}(\underline{n}_i^{(t-1)}) \right] = \hat{\theta}(\underline{n}_k^{(t)}).$$



$$\hat{\eta}(\underline{h}_k^{(t)}) = \sum_i \hat{p}(\underline{w}_{ik}^{(t-1)}; \hat{\theta}) + \sum_j \hat{p}(\underline{w}_{kj}^{(t)}; \hat{\theta}) + \hat{p}(\underline{h}_k^{(t)}; \theta).$$



$$\hat{\eta}(\underline{w}_{ij}) = \hat{p}(\underline{w}_{ij}; \hat{\theta})$$

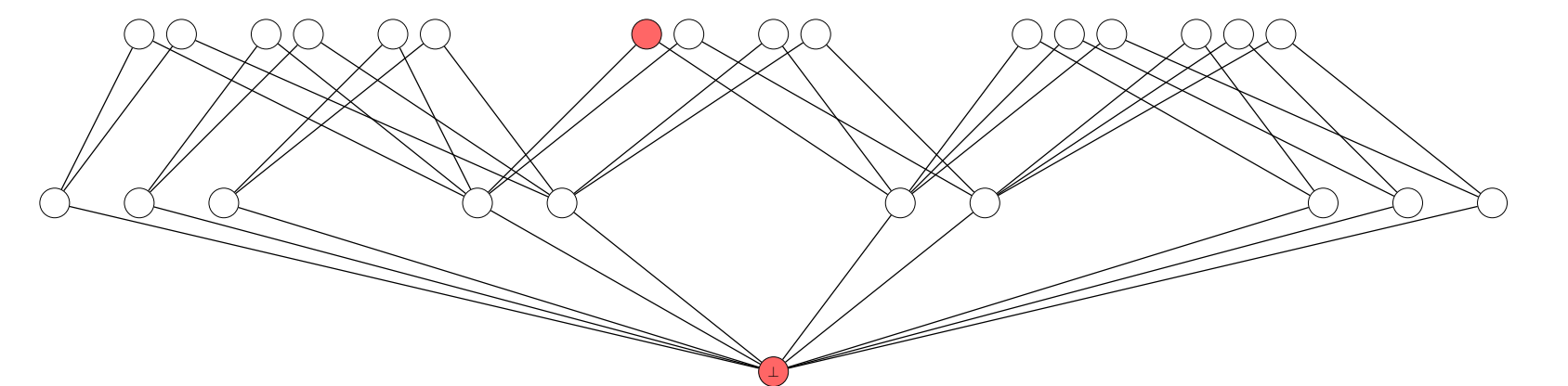
Optimization via EM-Algorithm

$$J(\theta, \eta) = -\mathbb{E}_{p(x,y)}[\log \hat{p}(\mathcal{H}, \mathcal{W} | \mathcal{X}, \mathcal{Y}; \theta)]$$

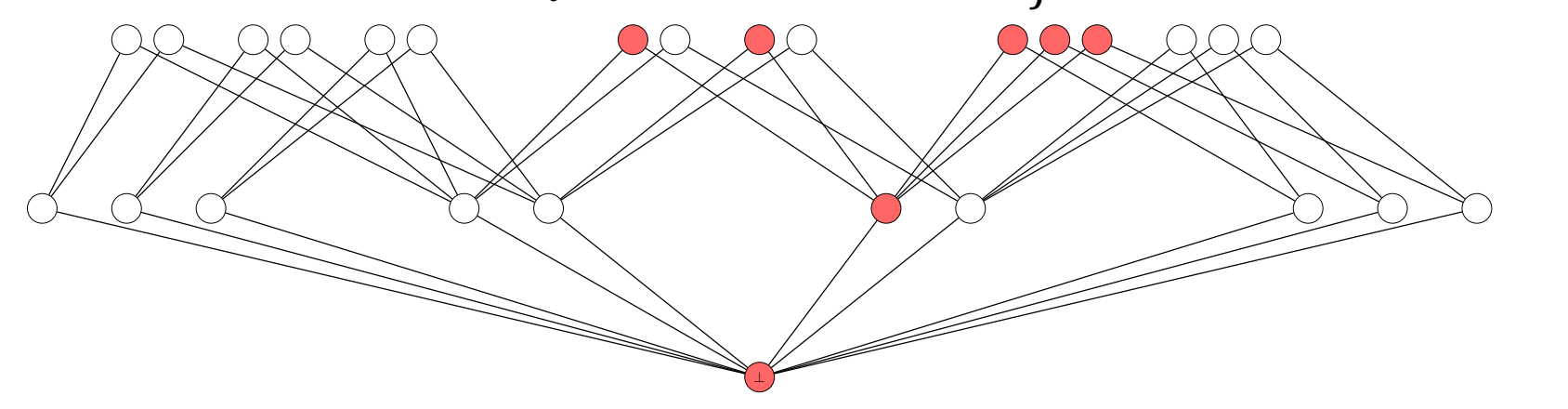
$$\eta_{\text{next}} = \arg \min_{\eta} J(\theta, \eta), \quad \theta_{\text{next}} = \arg \min_{\theta} J(\theta, \eta).$$

Maximization Step

$$\hat{\theta}(\underline{w}_{ij}) = \log[\hat{\eta}(\underline{w}_{ij})] + \psi(\hat{\theta})$$



$$\hat{\theta}(\underline{h}_k^{(t)}) = \log \left[\hat{\eta}(\underline{h}_k^{(t)}) - \sum_i \hat{p}(\underline{w}_{ik}^{(t-1)}; \hat{\theta}) - \sum_j \hat{p}(\underline{w}_{kj}^{(t)}; \hat{\theta}) \right] + \psi(\hat{\theta})$$



References

- Mahito Sugiyama, Hiroyuki Nakahara, and Koji Tsuda. **Tensor balancing on statistical manifold**, ICML 2017.
- Simon Luo and Mahito Sugiyama, **Bias-variance trade-off in hierarchical probabilistic models using higher-order feature interactions**, AAAI 2019.

