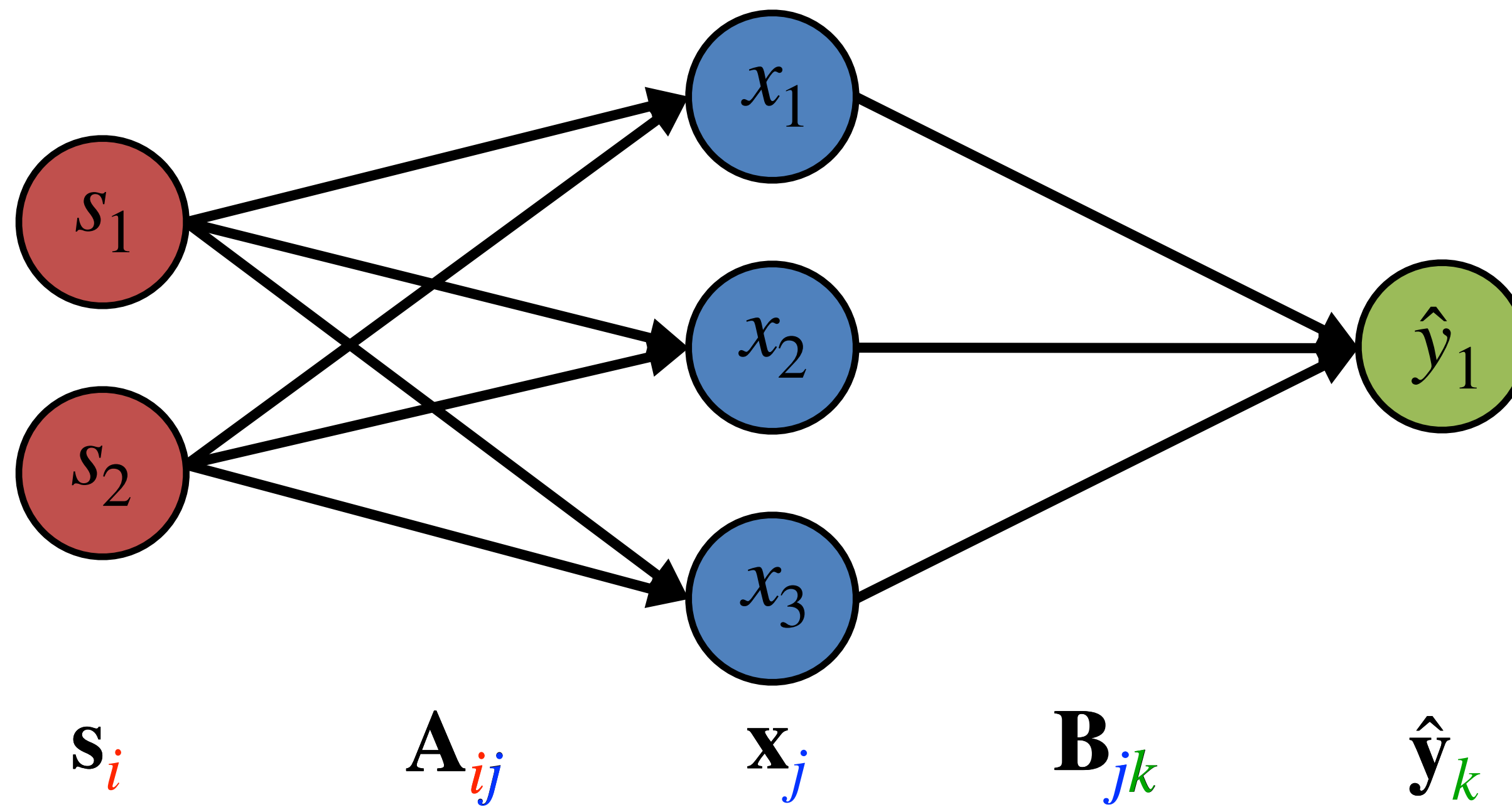


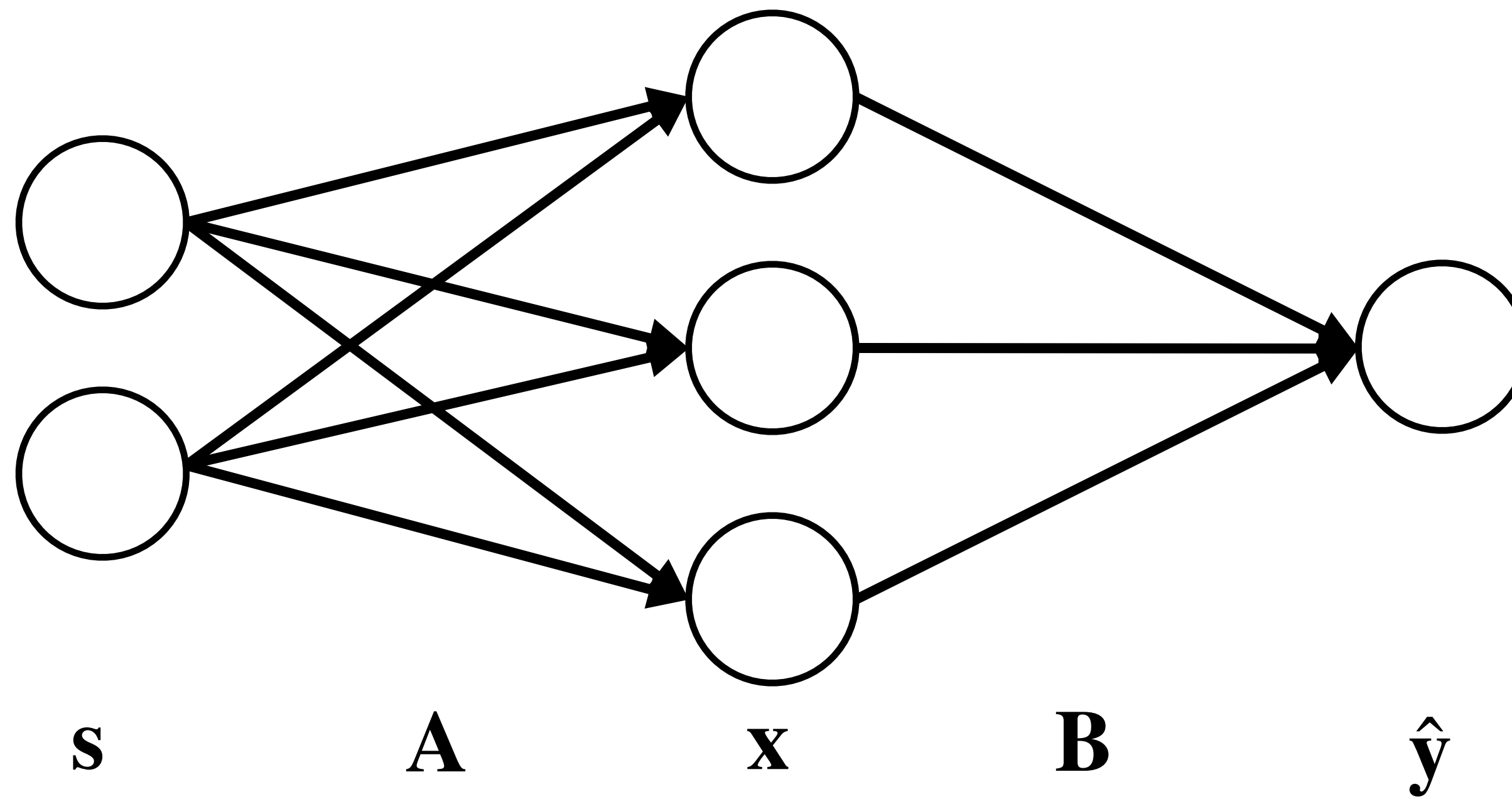
Prediction and Control with Approximation

Gradient Descent for Training Neural Networks

Notation

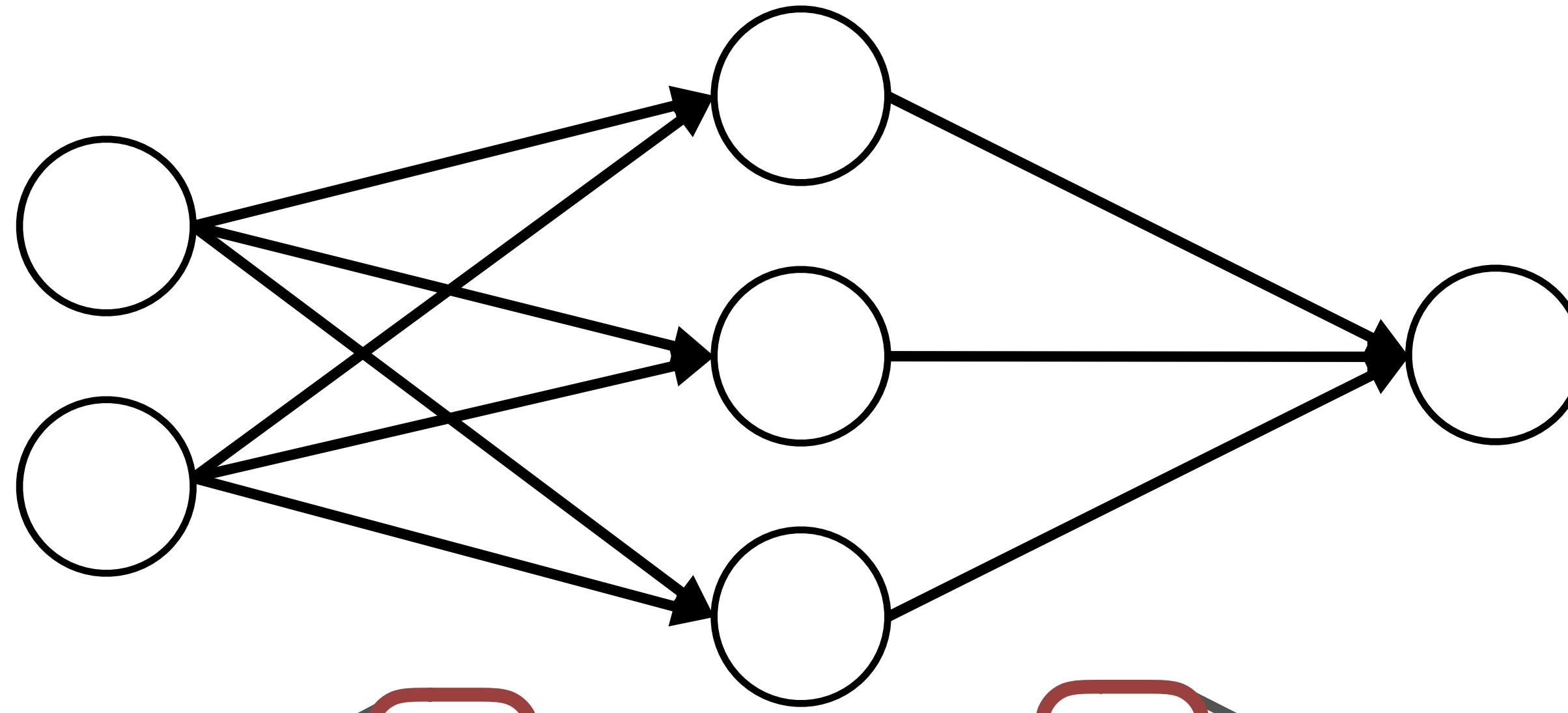


Notation



$$L(\hat{y}_k, y_k) = (\hat{y}_k - y_k)^2$$

Goal



s

A

x

B

\hat{y}

$$\mathbf{A} = \mathbf{A} - \alpha \delta^{\mathbf{A}} \mathbf{s}$$

$$\mathbf{B} = \mathbf{B} - \alpha \delta^{\mathbf{B}} \mathbf{x}$$

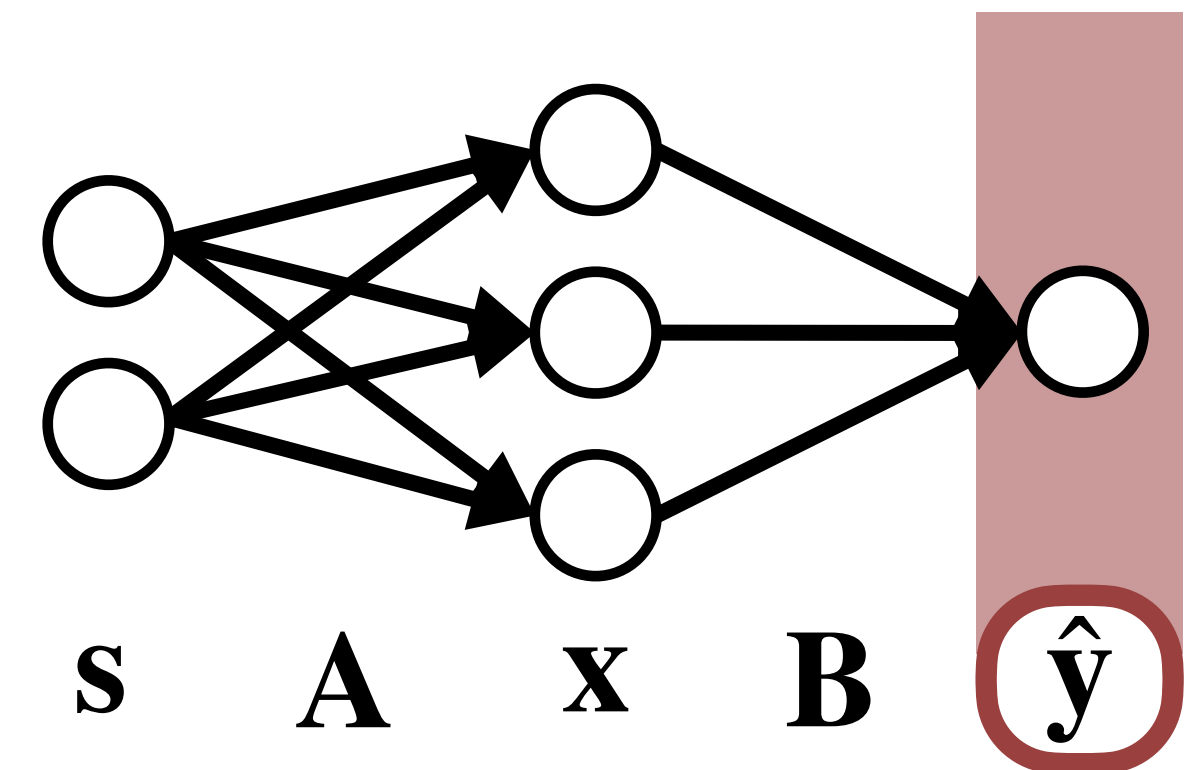
Deriving the gradient

$$\frac{\partial L(\hat{y}_k, y_k)}{\partial \mathbf{B}_{jk}} = \frac{\partial L(\hat{y}_k, y_k)}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial \mathbf{B}_{jk}}$$

$$\theta \doteq \mathbf{x}\mathbf{B}$$

$$\mathbf{x} \doteq f_{\mathbf{A}}(\mathbf{s}\mathbf{A})$$

$$\hat{y} \doteq f_{\mathbf{B}}(\mathbf{x}\mathbf{B})$$

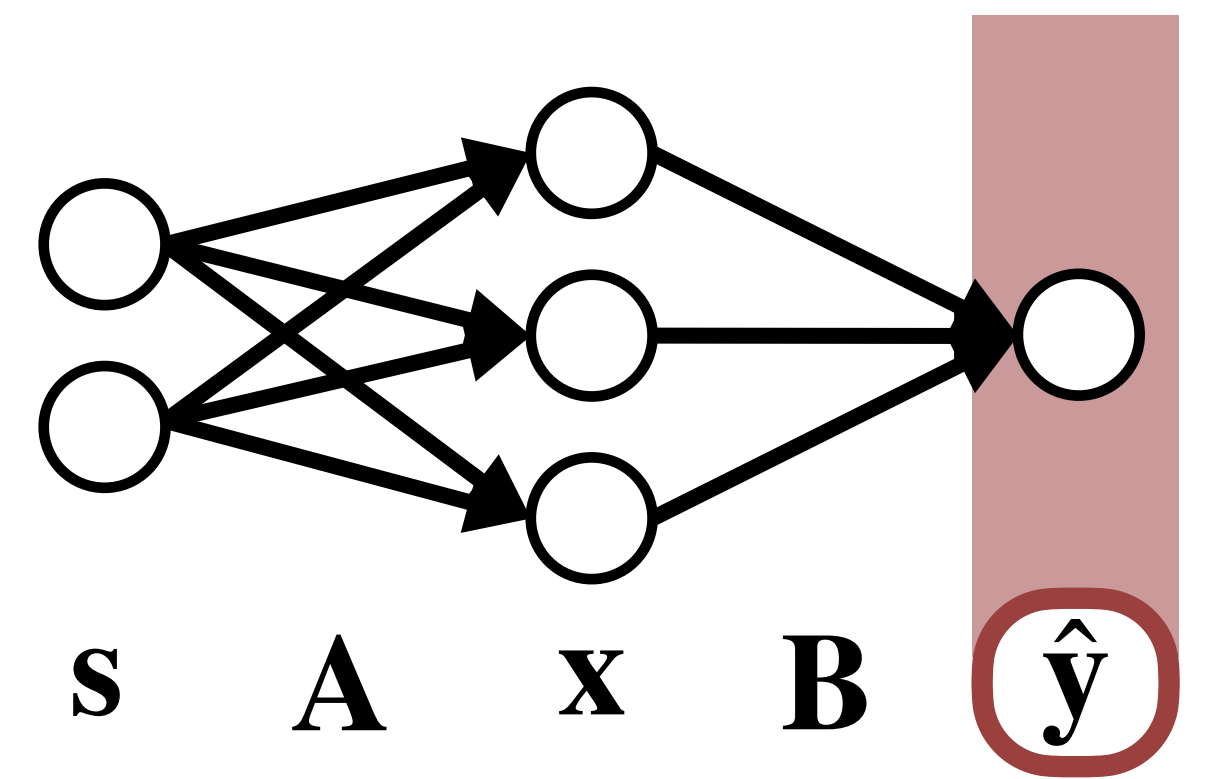


Deriving the gradient

$$\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \mathbf{B}_{jk}} = \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial \hat{\mathbf{y}}_k}{\partial \mathbf{B}_{jk}}$$
$$= \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k} \frac{\partial \theta_k}{\partial \mathbf{B}_{jk}}$$



$$\mathbf{x} \doteq f_{\mathbf{A}}(\mathbf{s}\mathbf{A})$$
$$\theta \doteq \mathbf{x}\mathbf{B}$$
$$\hat{\mathbf{y}} \doteq f_{\mathbf{B}}(\theta)$$



Deriving the gradient

$$\begin{aligned}\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \mathbf{B}_{jk}} &= \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial \hat{\mathbf{y}}_k}{\partial \mathbf{B}_{jk}} \\ &= \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k} \frac{\partial \theta_k}{\partial \mathbf{B}_{jk}}\end{aligned}$$

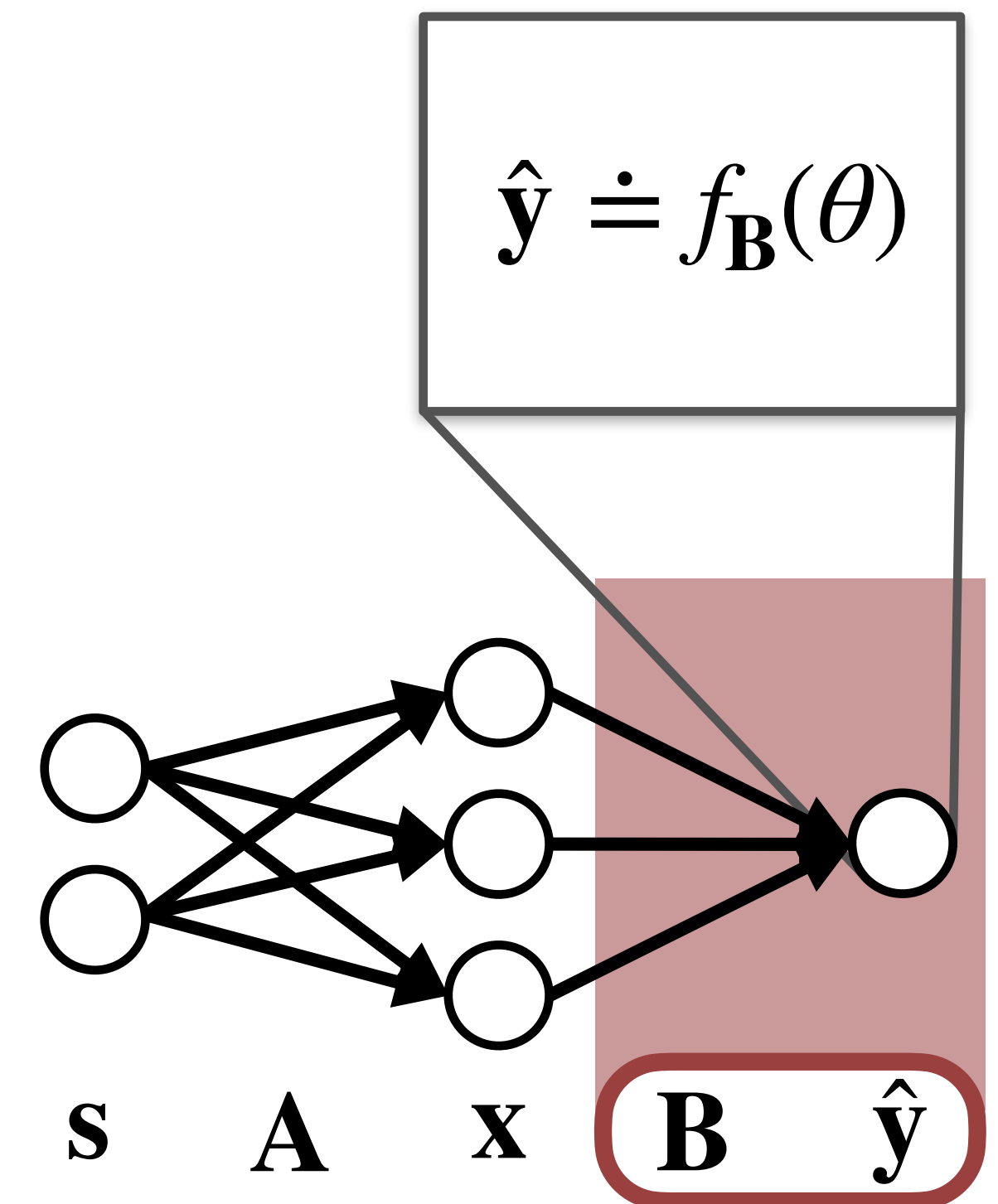
$$= \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k} \mathbf{x}_j$$

$$\frac{\partial \theta_k}{\partial \mathbf{B}_{jk}} = \mathbf{x}_j$$

$$\mathbf{x} \doteq f_{\mathbf{A}}(\mathbf{s}\mathbf{A})$$

$$\theta \doteq \mathbf{x}\mathbf{B}$$

$$\hat{\mathbf{y}} \doteq f_{\mathbf{B}}(\theta)$$



An example of the gradient

$$\mathbf{x} \doteq f_{\mathbf{A}}(\mathbf{s}\mathbf{A})$$

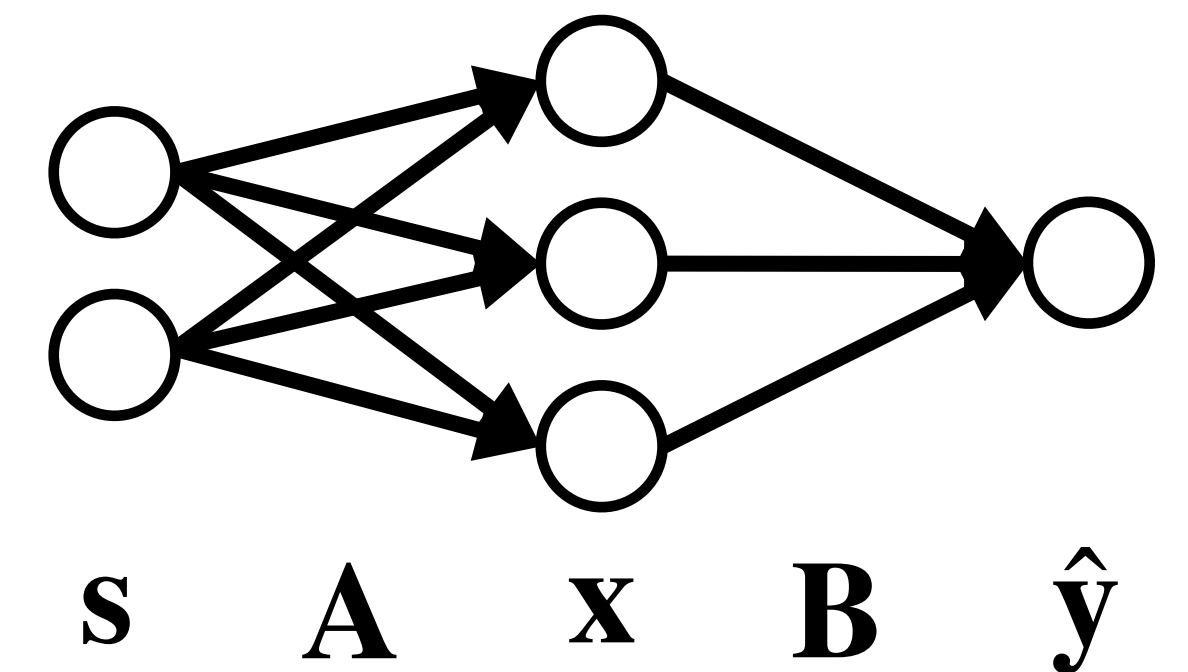
$$\theta \doteq \mathbf{x}\mathbf{B}$$

$$\hat{\mathbf{y}} \doteq f_{\mathbf{B}}(\theta)$$

Loss: $L = \frac{1}{2}(\hat{\mathbf{y}}_k - \mathbf{y}_k)^2$ $\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} = (\hat{\mathbf{y}}_k - \mathbf{y}_k)$

Activation: $f_{\mathbf{B}}(\theta_k) = \theta_k$ $\frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k} = \frac{\partial \theta_k}{\partial \theta_k} = 1$

$$\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \mathbf{B}_{jk}} = \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k} \mathbf{x}_j = (\hat{\mathbf{y}}_k - \mathbf{y}_k) \mathbf{x}_j$$



Deriving the gradient

$$\mathbf{x} \doteq f_{\mathbf{A}}(\mathbf{s}\mathbf{A})$$

$$\theta \doteq \mathbf{x}\mathbf{B}$$

$$\hat{\mathbf{y}} \doteq f_{\mathbf{B}}(\theta)$$

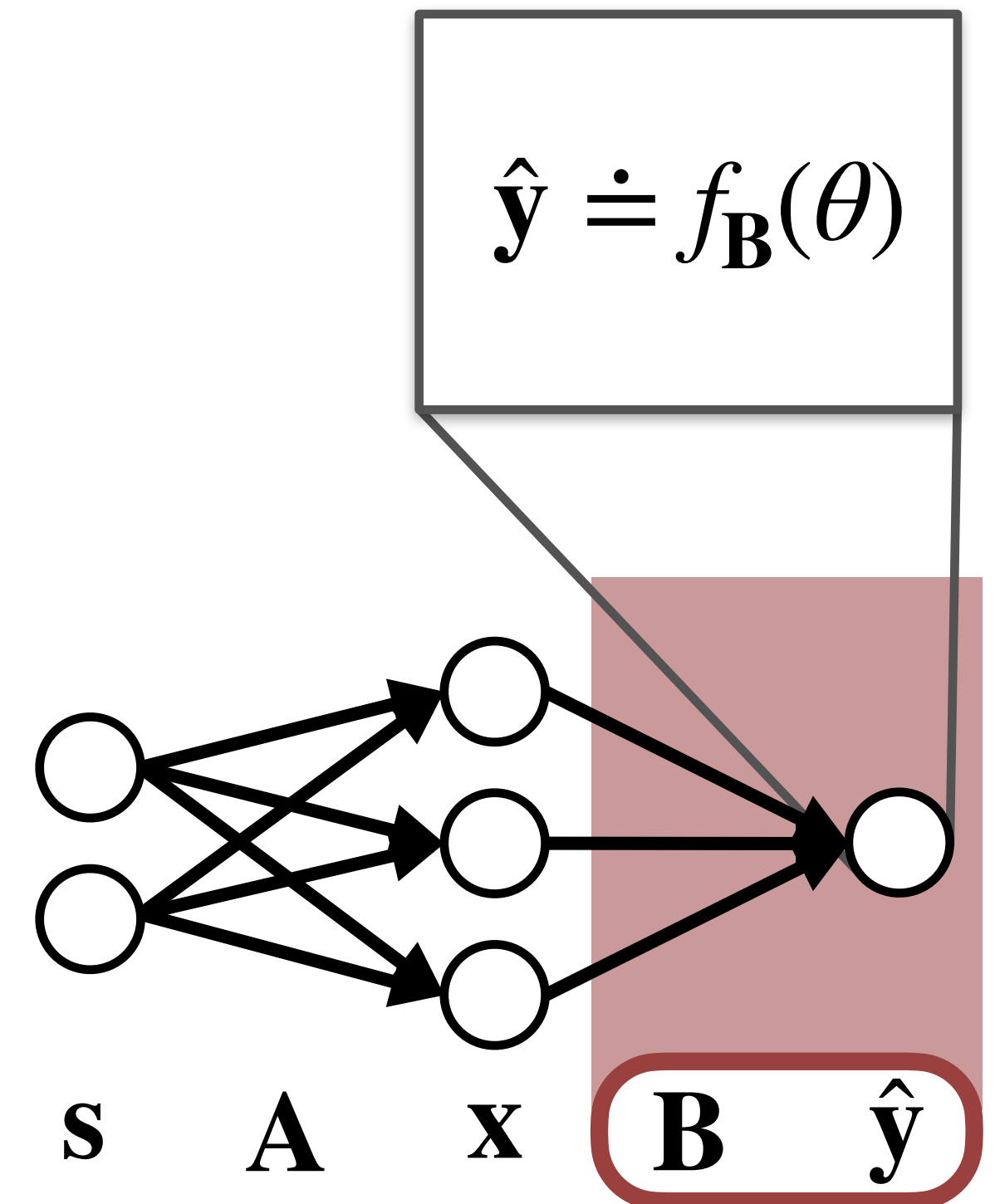
$$\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \mathbf{B}_{jk}} = \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial \hat{\mathbf{y}}_k}{\partial \mathbf{B}_{jk}}$$

$$= \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k} \frac{\partial \theta_k}{\partial \mathbf{B}_{jk}}$$

$$= \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k} \mathbf{x}_j$$

$$= \delta_k^{\mathbf{B}} \mathbf{x}_j$$

$$\delta_k^{\mathbf{B}} = \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k}$$



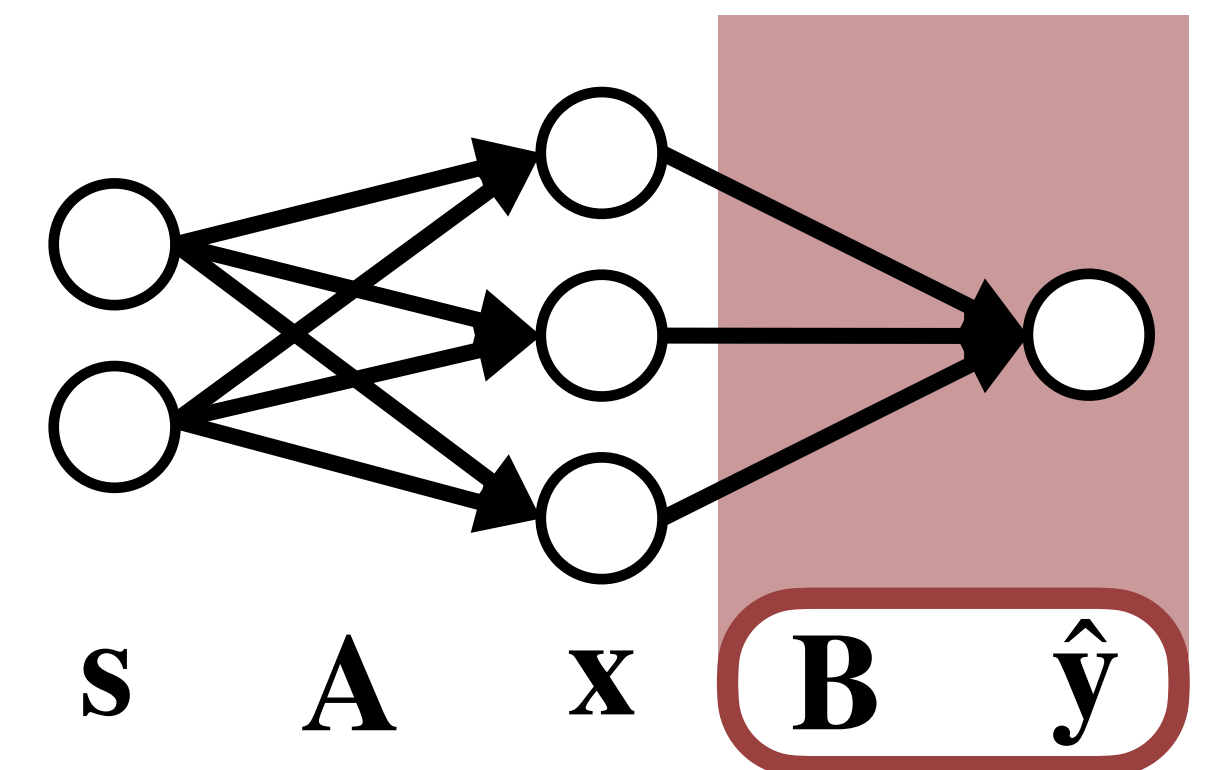
Deriving the gradient

$$\frac{\partial L(\hat{y}_k, y_k)}{\partial \mathbf{B}_{jk}} = \delta_k^{\mathbf{B}} \frac{\partial \theta_k}{\partial \mathbf{B}_{jk}}$$

$$\mathbf{x} \doteq f_{\mathbf{A}}(\mathbf{s}\mathbf{A})$$

$$\theta \doteq \mathbf{x}\mathbf{B}$$

$$\hat{y} \doteq f_{\mathbf{B}}(\theta)$$



Deriving the gradient

$$\frac{\partial L(\hat{y}_k, y_k)}{\partial \mathbf{A}_{ij}} = \delta_k^{\mathbf{B}} \frac{\partial \theta_k}{\partial \mathbf{A}_{ij}}$$

$$= \delta_k^{\mathbf{B}} \mathbf{B}_{jk} \frac{\partial \mathbf{x}_j}{\partial \mathbf{A}_{ij}}$$

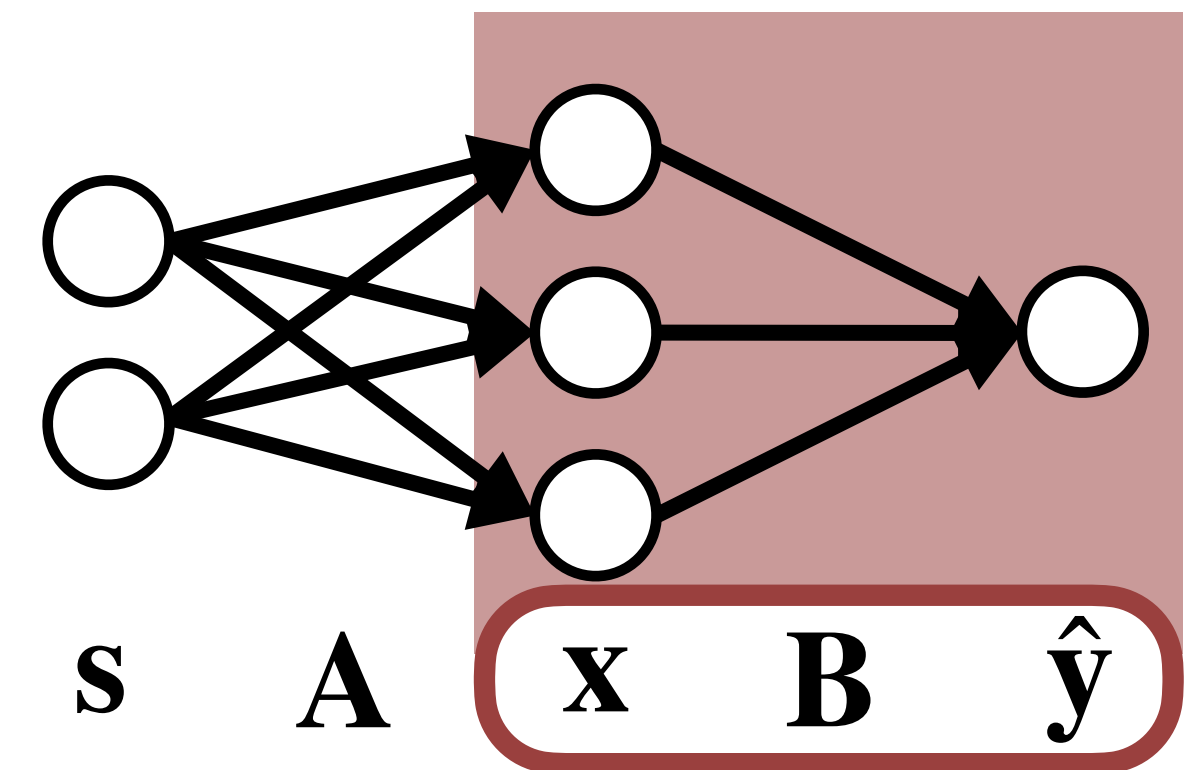
$$\frac{\partial \theta_k}{\partial \mathbf{A}_{ij}} = \mathbf{B}_{jk} \frac{\partial \mathbf{x}_j}{\partial \mathbf{A}_{ij}}$$

$$\psi \doteq \mathbf{sA}$$

$$\mathbf{x} \doteq f_{\mathbf{A}}(\mathbf{sA})$$

$$\theta \doteq \mathbf{xB}$$

$$\hat{y} \doteq f_{\mathbf{B}}(\theta)$$



Deriving the gradient

$$\frac{\partial L(\hat{y}_k, y_k)}{\partial \mathbf{A}_{ij}} = \delta_k^{\mathbf{B}} \frac{\partial \theta_k}{\partial \mathbf{A}_{ij}}$$

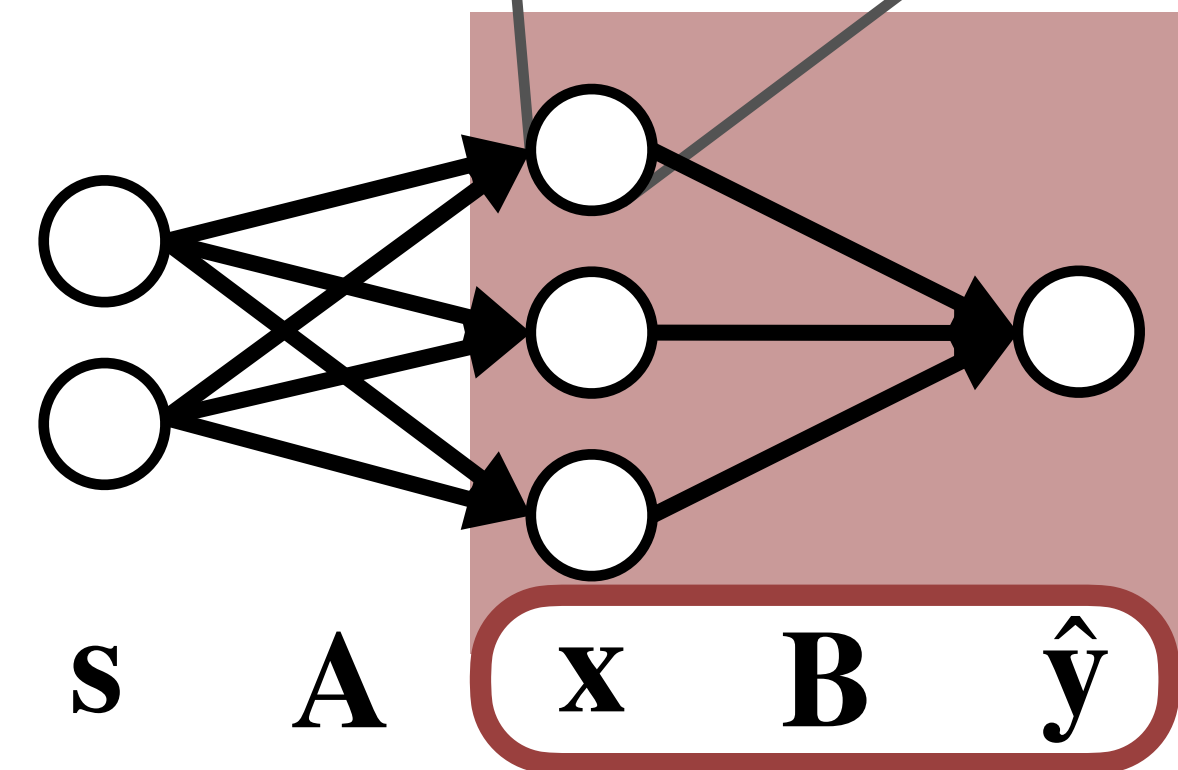
$$= \delta_k^{\mathbf{B}} \mathbf{B}_{jk} \frac{\partial \mathbf{x}_j}{\partial \mathbf{A}_{ij}}$$

$$\frac{\partial \mathbf{x}_j}{\partial \mathbf{A}_{ij}} = \frac{\partial f_{\mathbf{A}}(\psi_j)}{\partial \psi_j} \frac{\partial \psi_j}{\partial \mathbf{A}_{ij}}$$



$$\begin{aligned} \psi &\doteq \mathbf{sA} \\ \mathbf{x} &\doteq f_{\mathbf{A}}(\mathbf{sA}) \\ \theta &\doteq \mathbf{xB} \\ \hat{y} &\doteq f_{\mathbf{B}}(\theta) \end{aligned}$$

$$\mathbf{x} \doteq f_{\mathbf{A}}(\psi)$$



Deriving the gradient

$$\begin{aligned}\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \mathbf{A}_{ij}} &= \delta_k^{\mathbf{B}} \frac{\partial \theta_k}{\partial \mathbf{A}_{ij}} \\ &= \delta_k^{\mathbf{B}} \mathbf{B}_{jk} \frac{\partial \mathbf{x}_j}{\partial \mathbf{A}_{ij}} \\ &= \delta_k^{\mathbf{B}} \mathbf{B}_{jk} \frac{\partial f_{\mathbf{A}}(\psi_j)}{\partial \psi_j} \frac{\partial \psi_j}{\partial \mathbf{A}_{ij}} \\ &= \delta_k^{\mathbf{B}} \mathbf{B}_{jk} \frac{\partial f_{\mathbf{A}}(\psi_j)}{\partial \psi_j} \mathbf{s}_i\end{aligned}$$

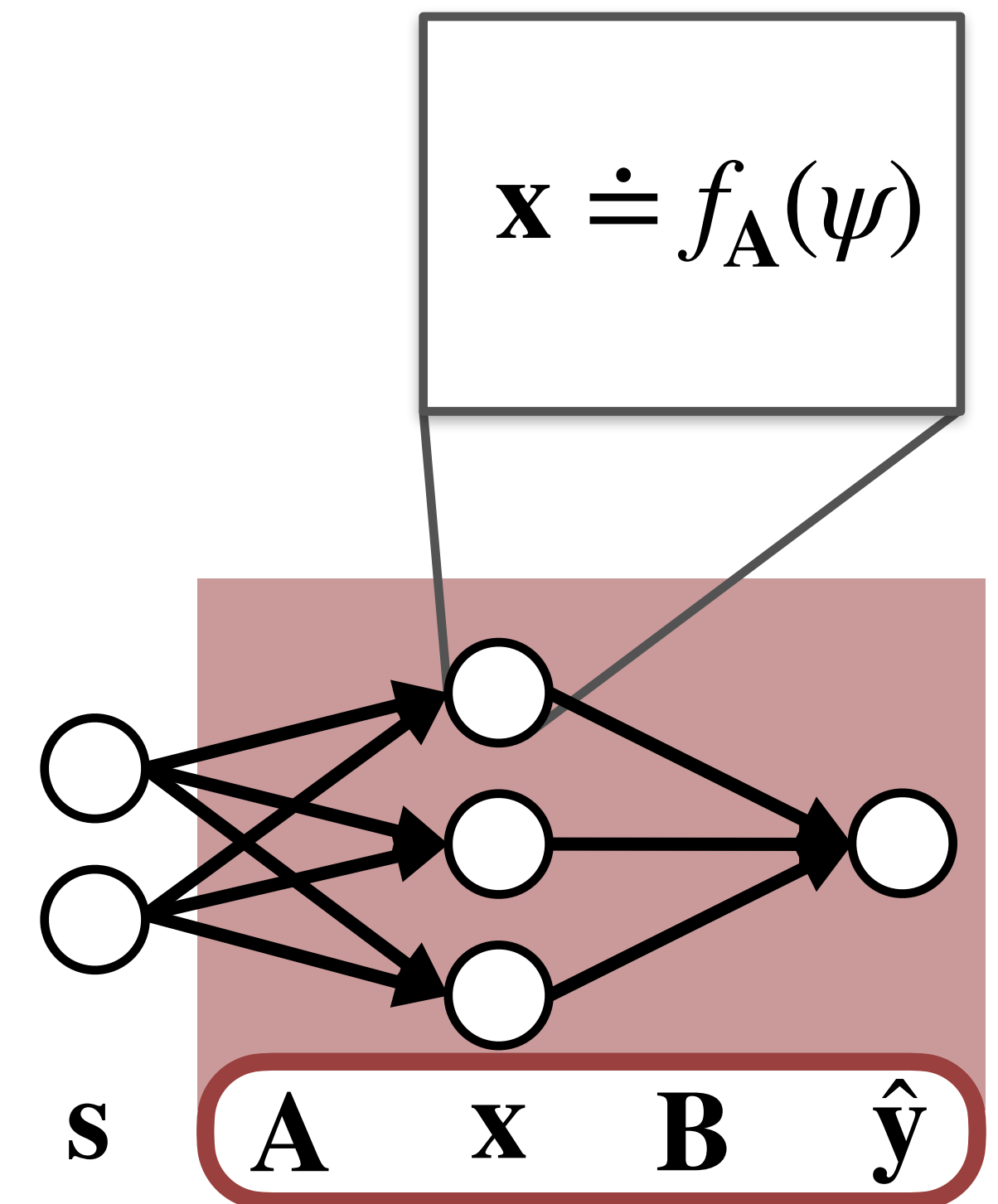
$$\frac{\partial \psi_j}{\partial \mathbf{A}_{ij}} = \mathbf{s}_i$$

$$\psi \doteq \mathbf{s}\mathbf{A}$$

$$\mathbf{x} \doteq f_{\mathbf{A}}(\psi)$$

$$\theta \doteq \mathbf{x}\mathbf{B}$$

$$\hat{\mathbf{y}} \doteq f_{\mathbf{B}}(\theta)$$



Deriving the gradient

$$\frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \mathbf{A}_{ij}} = \delta_k^{\mathbf{B}} \mathbf{B}_{jk} \frac{\partial f_{\mathbf{A}}(\psi_j)}{\partial \psi_j} \mathbf{s}_i$$
$$= \delta_j^{\mathbf{A}} \mathbf{s}_i$$

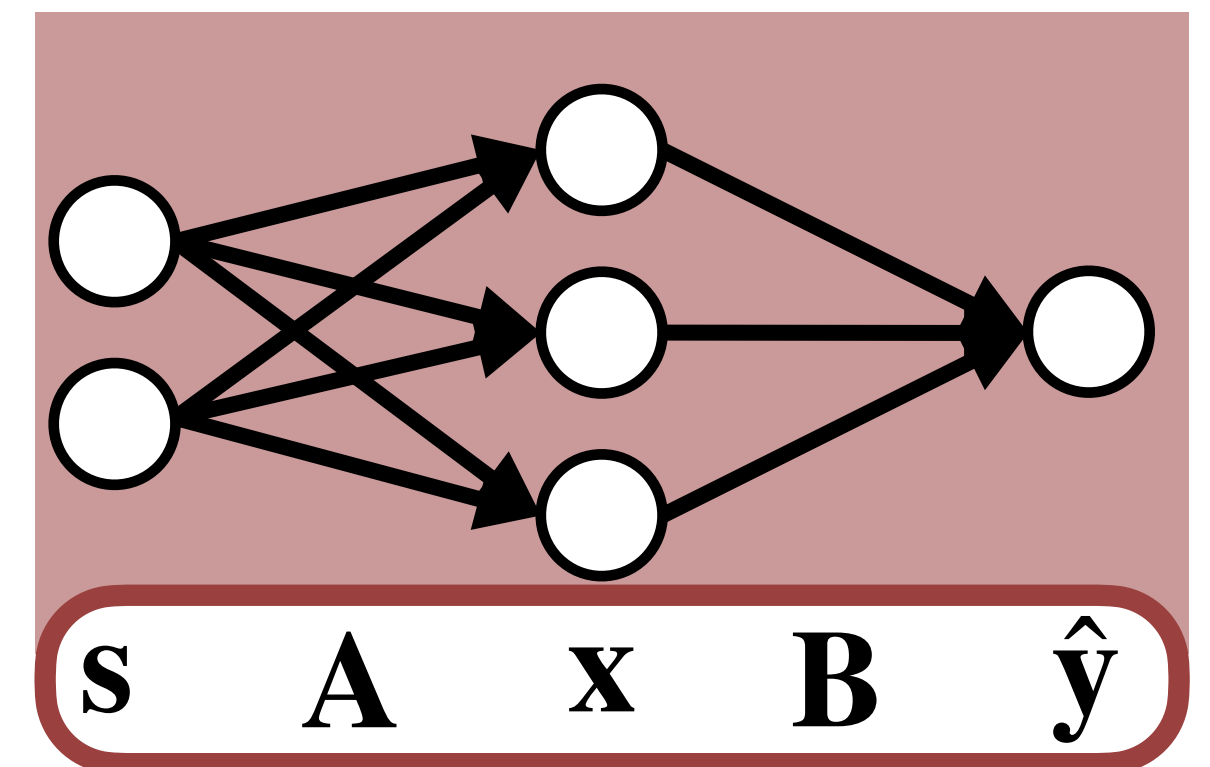
$$\delta_j^{\mathbf{A}} = (\mathbf{B}_{jk} \delta_k^{\mathbf{B}}) \frac{\partial f_{\mathbf{A}}(\psi_j)}{\partial \psi_j}$$

$$\psi \doteq \mathbf{sA}$$

$$\mathbf{x} \doteq f_{\mathbf{A}}(\psi)$$

$$\theta \doteq \mathbf{xB}$$

$$\hat{\mathbf{y}} \doteq f_{\mathbf{B}}(\theta)$$



Deriving the gradient

$$\frac{\partial L(\hat{y}_k, y_k)}{\partial \mathbf{A}_{ij}} = \delta_j^{\mathbf{A}} \mathbf{s}_i$$

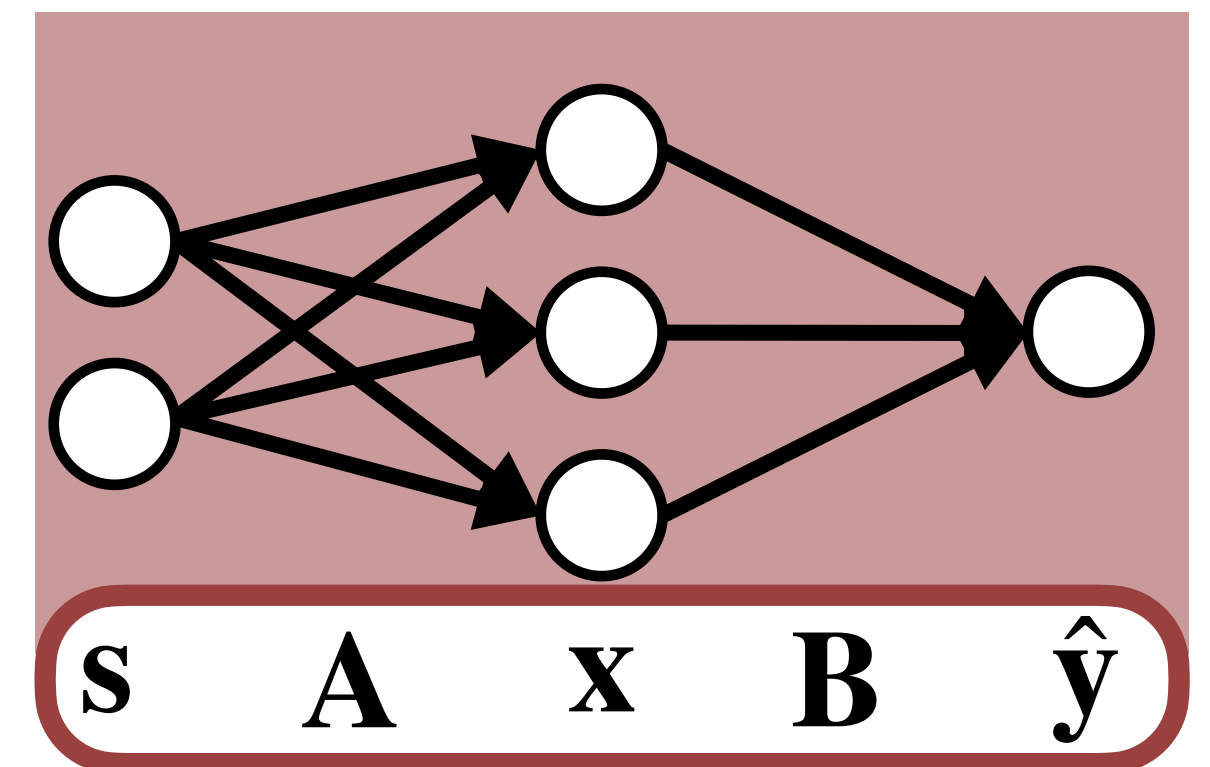
$$\frac{\partial L(\hat{y}_k, y_k)}{\partial \mathbf{B}_{jk}} = \delta_k^{\mathbf{B}} \mathbf{x}_j$$

$$\psi \doteq \mathbf{s}\mathbf{A}$$

$$\mathbf{x} \doteq f_{\mathbf{A}}(\psi)$$

$$\theta \doteq \mathbf{x}\mathbf{B}$$

$$\hat{y} \doteq f_{\mathbf{B}}(\theta)$$



The backprop algorithm

for each (s, y) in D :

$$\delta_k^{\mathbf{B}} = \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k}$$

$$\nabla_{\mathbf{B}}^{jk} = \delta_k^{\mathbf{B}} \mathbf{x}_j$$

$$\mathbf{B} = \mathbf{B} - \alpha_{\mathbf{B}} \nabla_{\mathbf{B}}$$

$$\delta_j^{\mathbf{A}} = \left(\mathbf{B}_{jk} \delta_k^{\mathbf{B}} \right) \frac{\partial f_{\mathbf{A}}(\psi_j)}{\partial \psi_j}$$

$$\nabla_{\mathbf{A}}^{ij} = \delta_j^{\mathbf{A}} \mathbf{s}_i$$

$$\mathbf{A} = \mathbf{A} - \alpha_{\mathbf{A}} \nabla_{\mathbf{A}}$$

The backprop algorithm

for each (s, y) in D :

$$\delta_k^{\mathbf{B}} = (\hat{\mathbf{y}}_k - \mathbf{y}_k) \mathbf{x}_j$$

$$\delta_k^{\mathbf{B}} = \frac{\partial L(\hat{\mathbf{y}}_k, \mathbf{y}_k)}{\partial \hat{\mathbf{y}}_k} \frac{\partial f_{\mathbf{B}}(\theta_k)}{\partial \theta_k}$$

$$\nabla_{\mathbf{B}}^{jk} = \delta_k^{\mathbf{B}} \mathbf{x}_j$$

$$\mathbf{B} = \mathbf{B} - \alpha_{\mathbf{B}} \nabla_{\mathbf{B}}$$

$$u = \begin{cases} \psi & \psi > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$u = \frac{\partial f_{\mathbf{A}}(\psi_j)}{\partial \psi_j}$$

$$\delta_j^{\mathbf{A}} = (\mathbf{B}_{jk} \delta_k^{\mathbf{B}}) u$$

$$\delta_j^{\mathbf{A}} = (\mathbf{B}_{jk} \delta_k^{\mathbf{B}}) \frac{\partial f_{\mathbf{A}}(\psi_j)}{\partial \psi_j}$$

$$\nabla_{\mathbf{A}}^{ij} = \delta_j^{\mathbf{A}} \mathbf{s}_i$$

$$\mathbf{A} = \mathbf{A} - \alpha_{\mathbf{A}} \nabla_{\mathbf{A}}$$

Summary

- The **gradient** can be used to update the parameters of a neural network with **stochastic gradient descent**
- **Backprop** can save **computation** by computing gradients starting at the output of the network.