

# **Reduced Products of Abstract Domains for Fairness Certification of Neural Networks**

Denis Mazzucato and Caterina Urban

**SAS 2021 - June 16th, 2022**

The image is a collage of various news articles and a screenshot of the Google Translate interface, all centered around the theme of machine learning's social impact.

**Top News Headlines:**

- WIRED** - [In 2019, predictive algorithms will start to make banking fair for all](#) (October 10, 2018)
- WIRED** - [Amazon scraps secret AI recruiting tool that showed bias against women](#) (March 25, 2019)
- The Telegraph** - [AI used for first time in job interviews in UK to find best applicants](#) (September 27, 2019)
- WIRED** - [The AI Doctor Will See You Now](#) (December 21, 2019)
- Google Translate** (Screenshot showing English to French translation of medical terms like "nurse" and "doctor")
- nature** - [Millions of black people affected by racial bias in health-care algorithms](#) (October 24, 2019)
- WIRED** - [Can AI Be a Fair Judge in Court? Estonia Thinks So](#) (March 25, 2019)
- WIRED** - [AUTOMATED BACKGROUND CHECKS ARE DECIDING WHO'S FIT FOR A HOME](#) (February 1, 2019)
- Machine Bias** (Text overlay: "There's software used across the country to predict future criminals. And it's biased against blacks.")

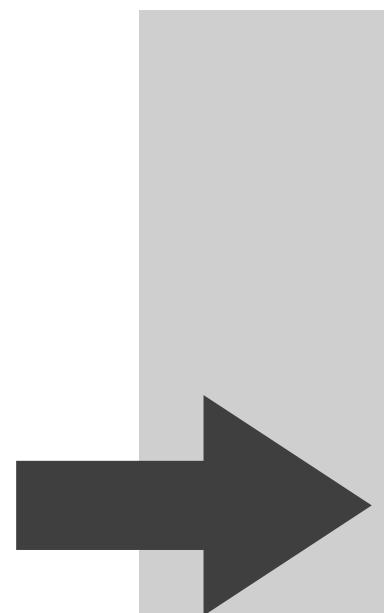
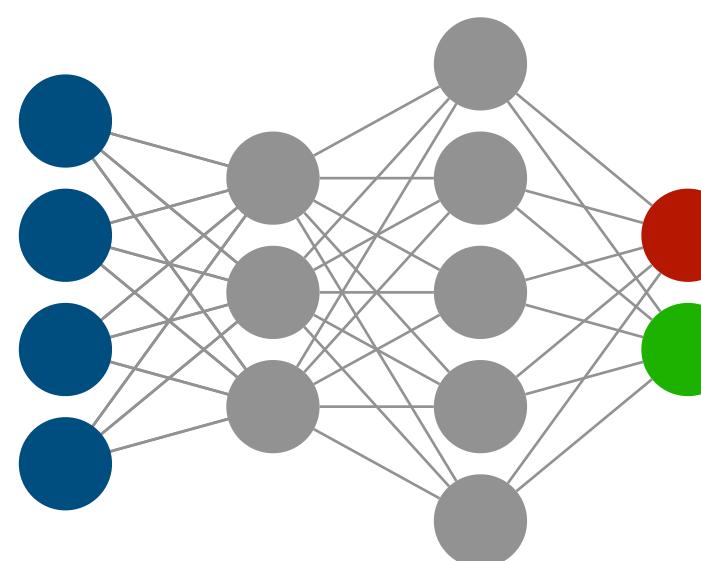
**Text Labels:**

- 8 MIN READ
- TOM SIMONITE
- BUSINESS 12.21.2019 08:00 AM
- By Charles Hymas 27 SEPTEMBER 2019 • 10:00 PM
- By Colin Lecher | @colinlecher | Feb 1, 2019, 8:00am EST
- by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica May 23, 2016

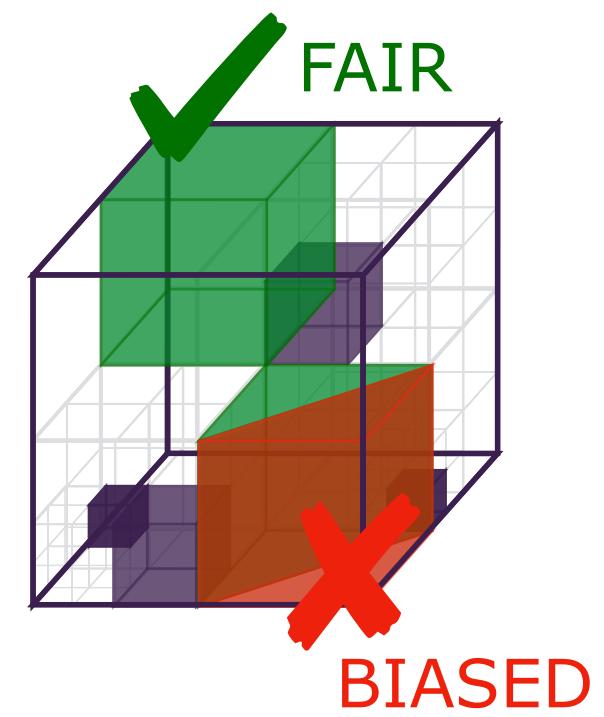
A black silhouette map of Europe on a white background. There are twelve yellow five-pointed stars arranged in a curve from the bottom left towards the top right, representing the European Union.

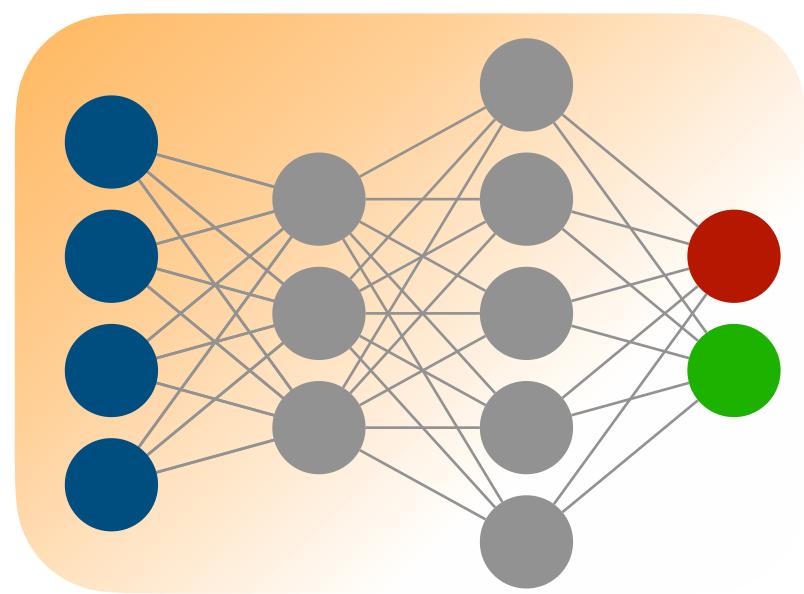
# Artificial Intelligence Act

April 2021

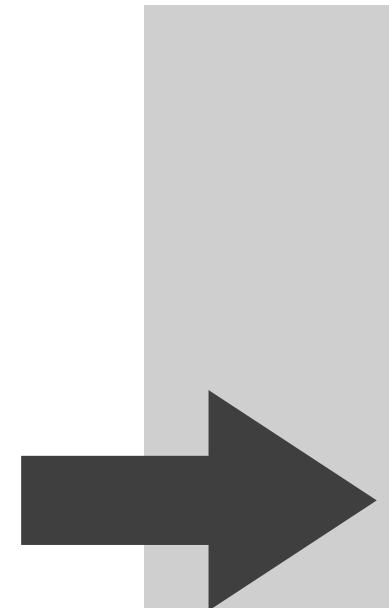


# Libra

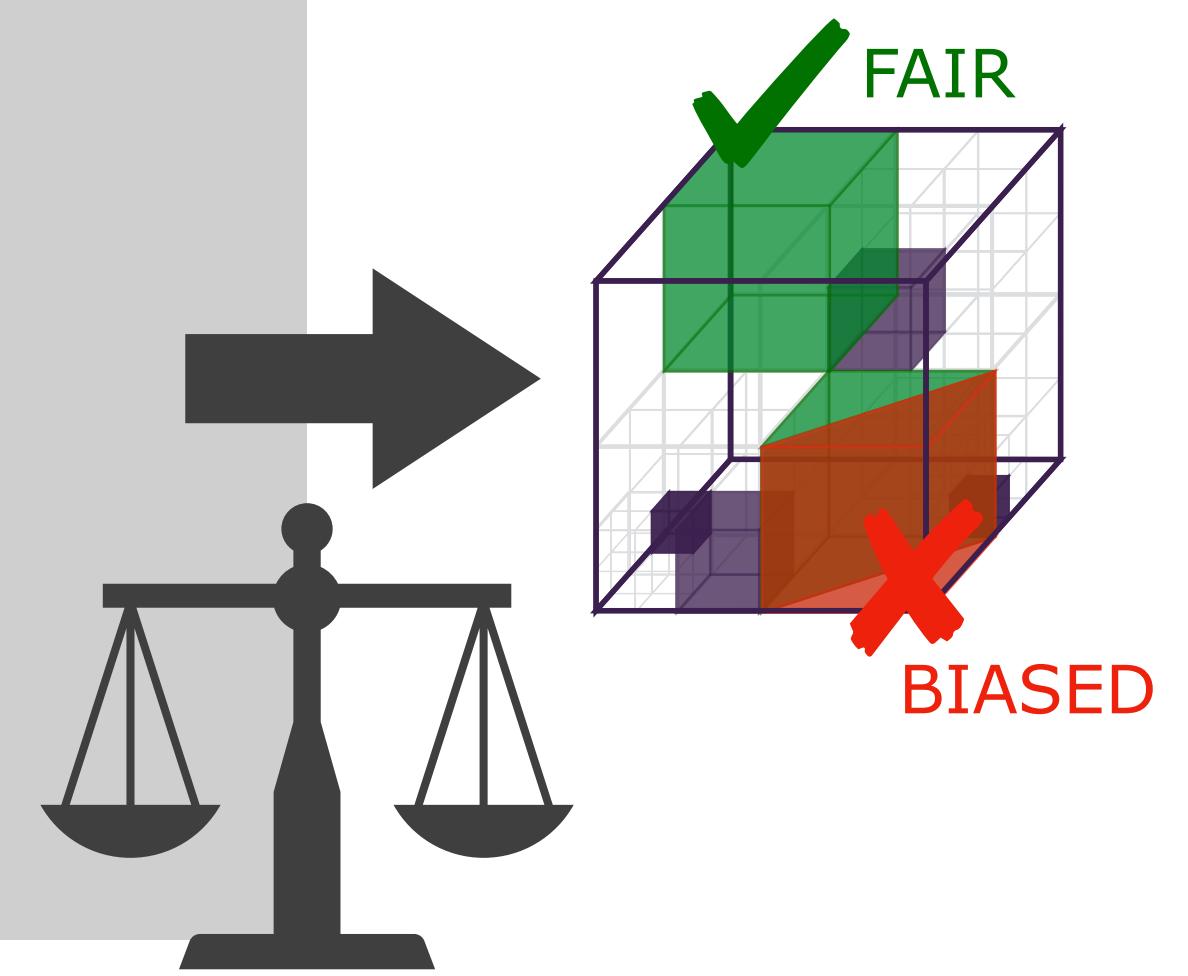




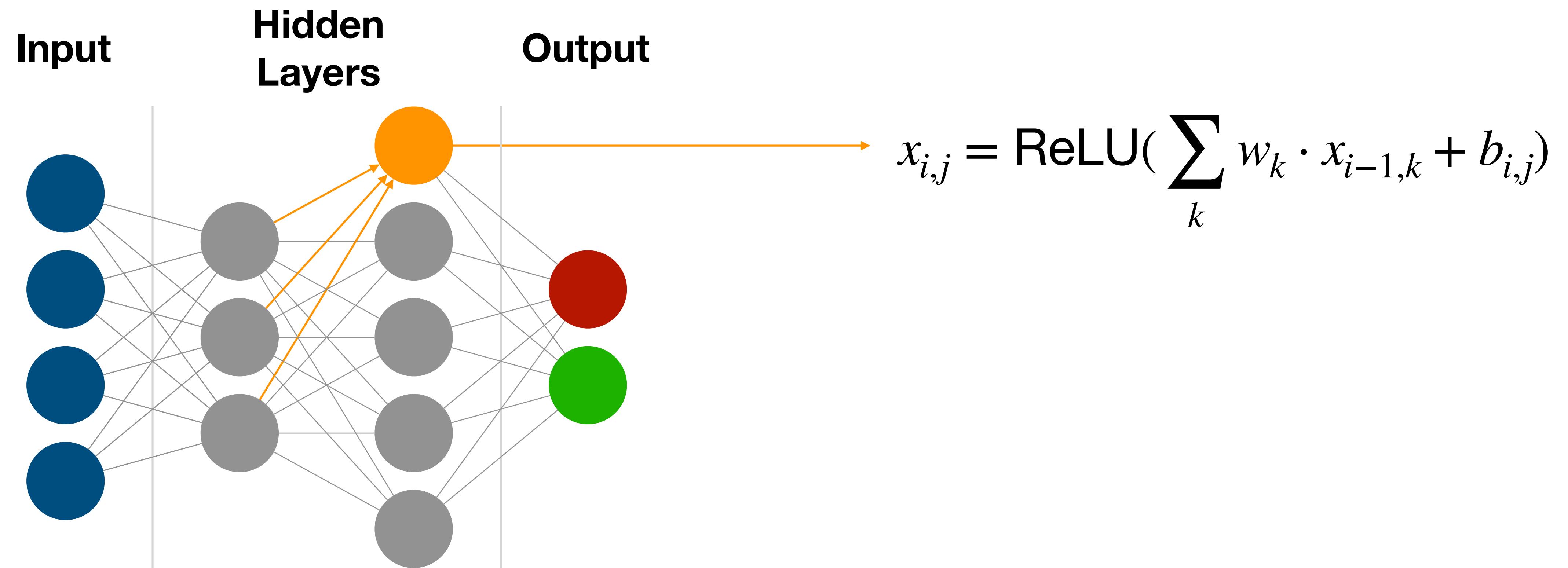
Neural Network



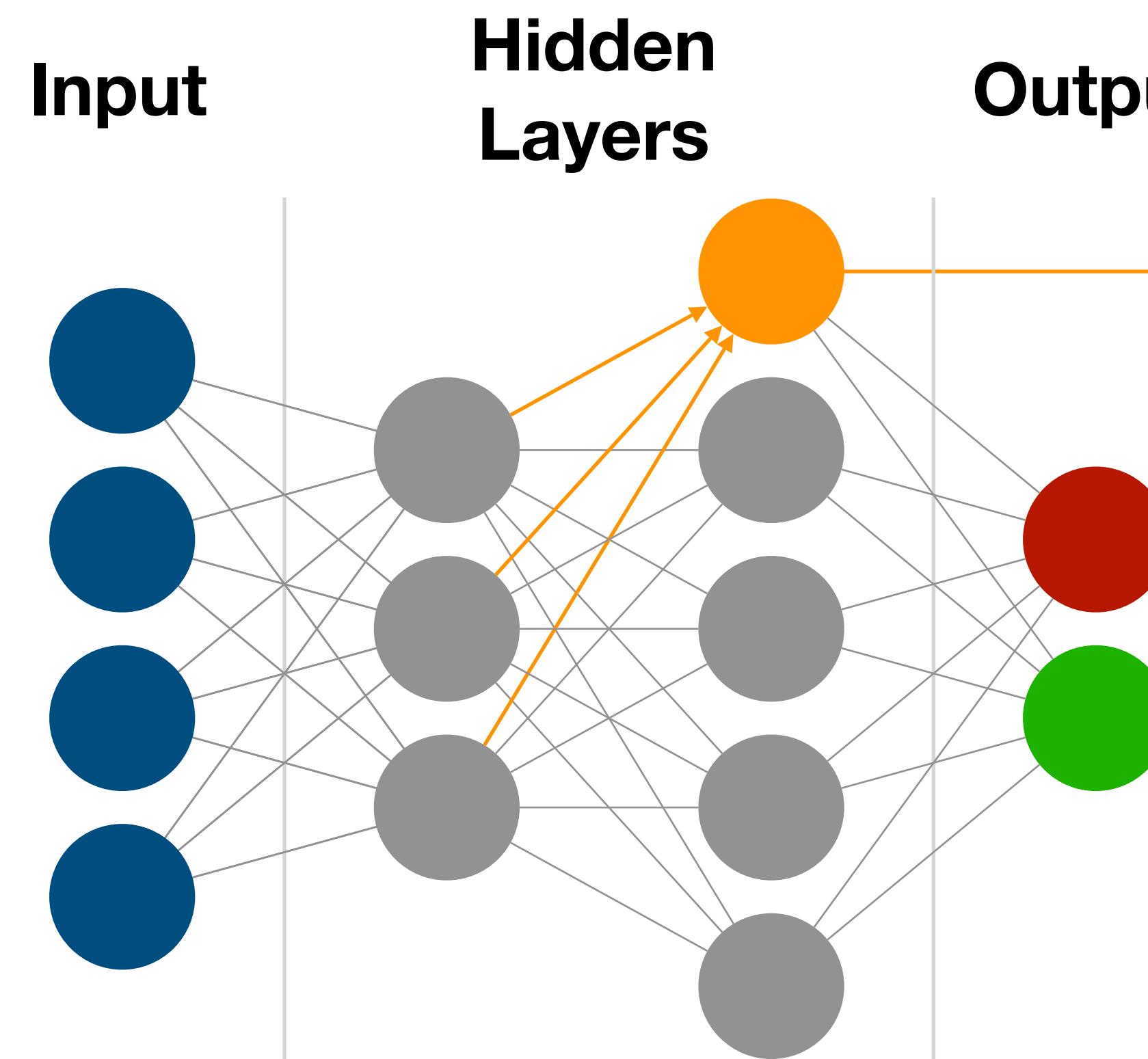
# Libra



# Feed-Forward Neural Networks with ReLU Activations

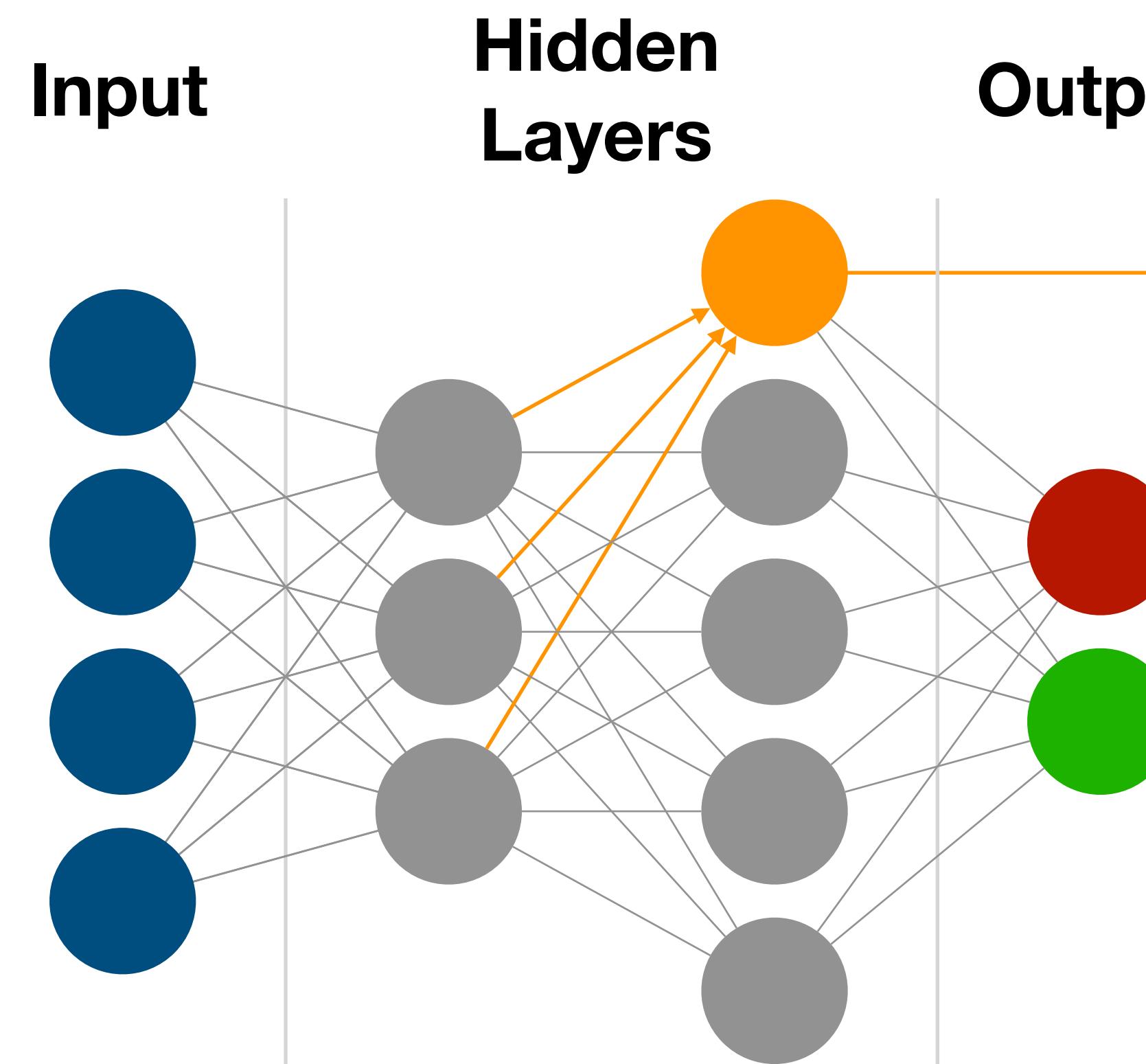


# Feed-Forward Neural Networks with ReLU Activations



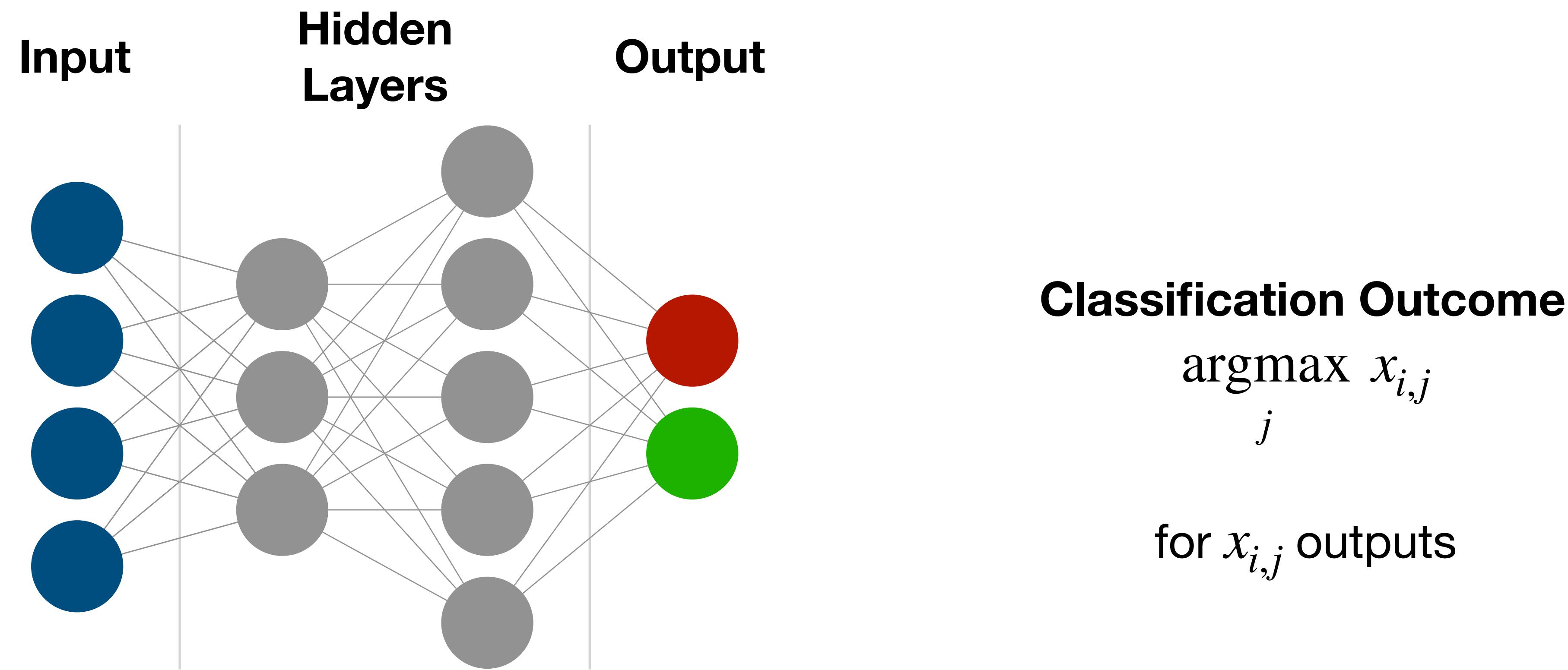
$$\begin{aligned}x_{i,j} &= \text{ReLU}\left(\sum_k w_k \cdot x_{i-1,k} + b_{i,j}\right) \\&= \max\left(\sum_k w_k \cdot x_{i-1,k} + b_{i,j}, 0\right)\end{aligned}$$

# Feed-Forward Neural Networks with ReLU Activations

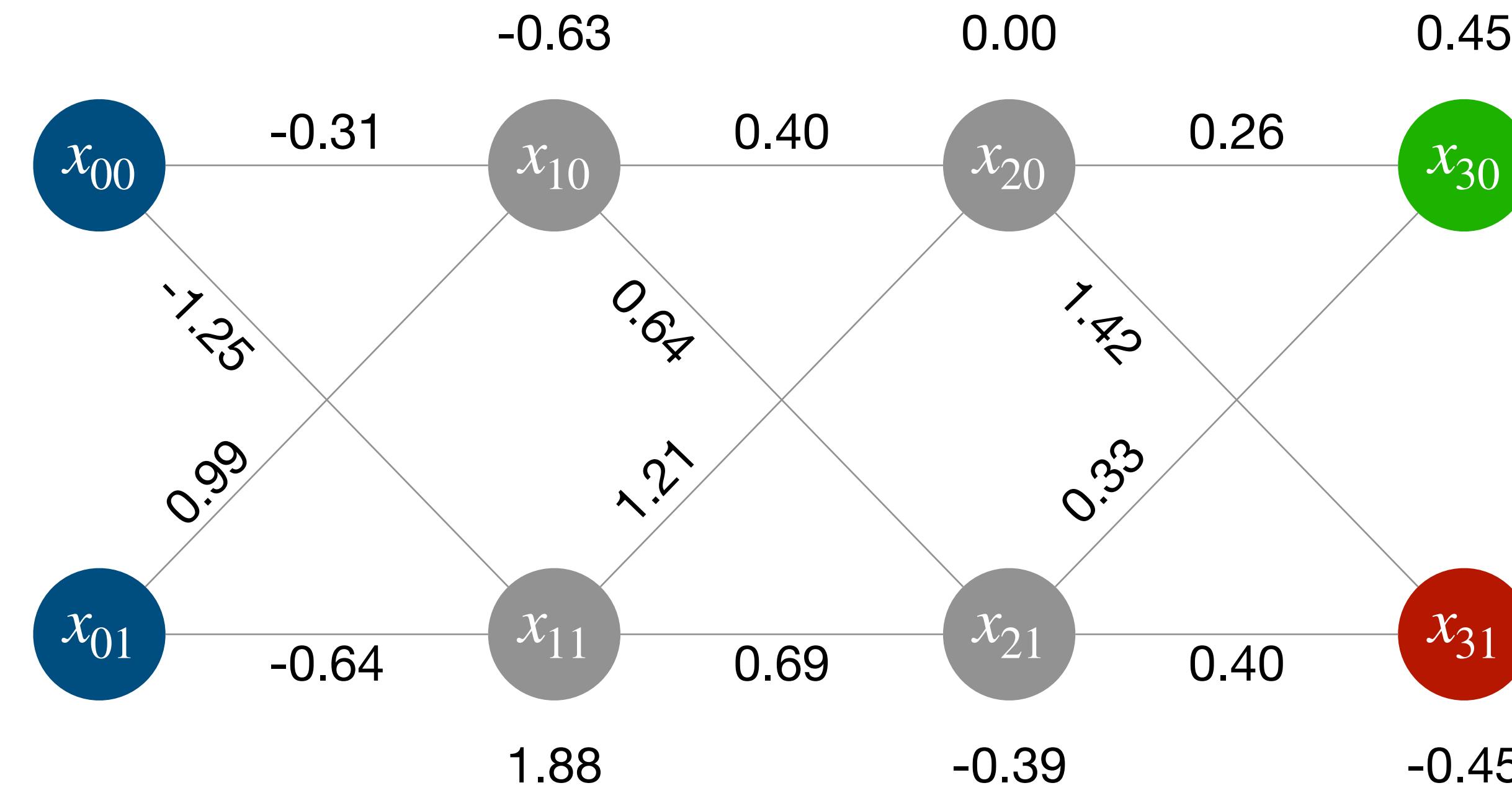


$$\begin{aligned}x_{i,j} &= \text{ReLU}\left(\sum_k w_k \cdot x_{i-1,k} + b_{i,j}\right) \\&= \max\left(\underbrace{\sum_k w_k \cdot x_{i-1,k} + b_{i,j}}_k, 0\right) \\ \hat{x}_{i,j} &= \sum_k w_k \cdot x_{i-1,k} + b_{i,j}\end{aligned}$$

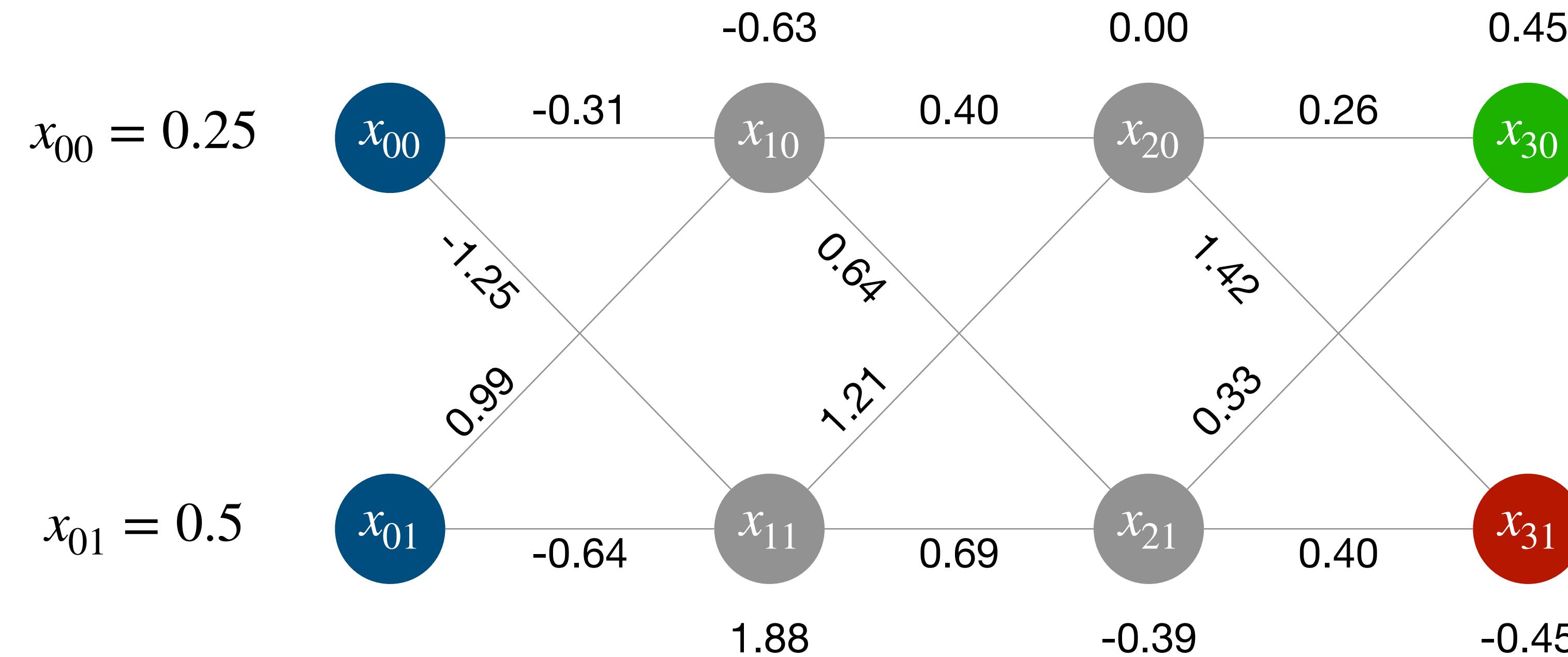
# Feed-Forward Neural Networks with ReLU Activations



# Feed-Forward Neural Networks with ReLU Activations

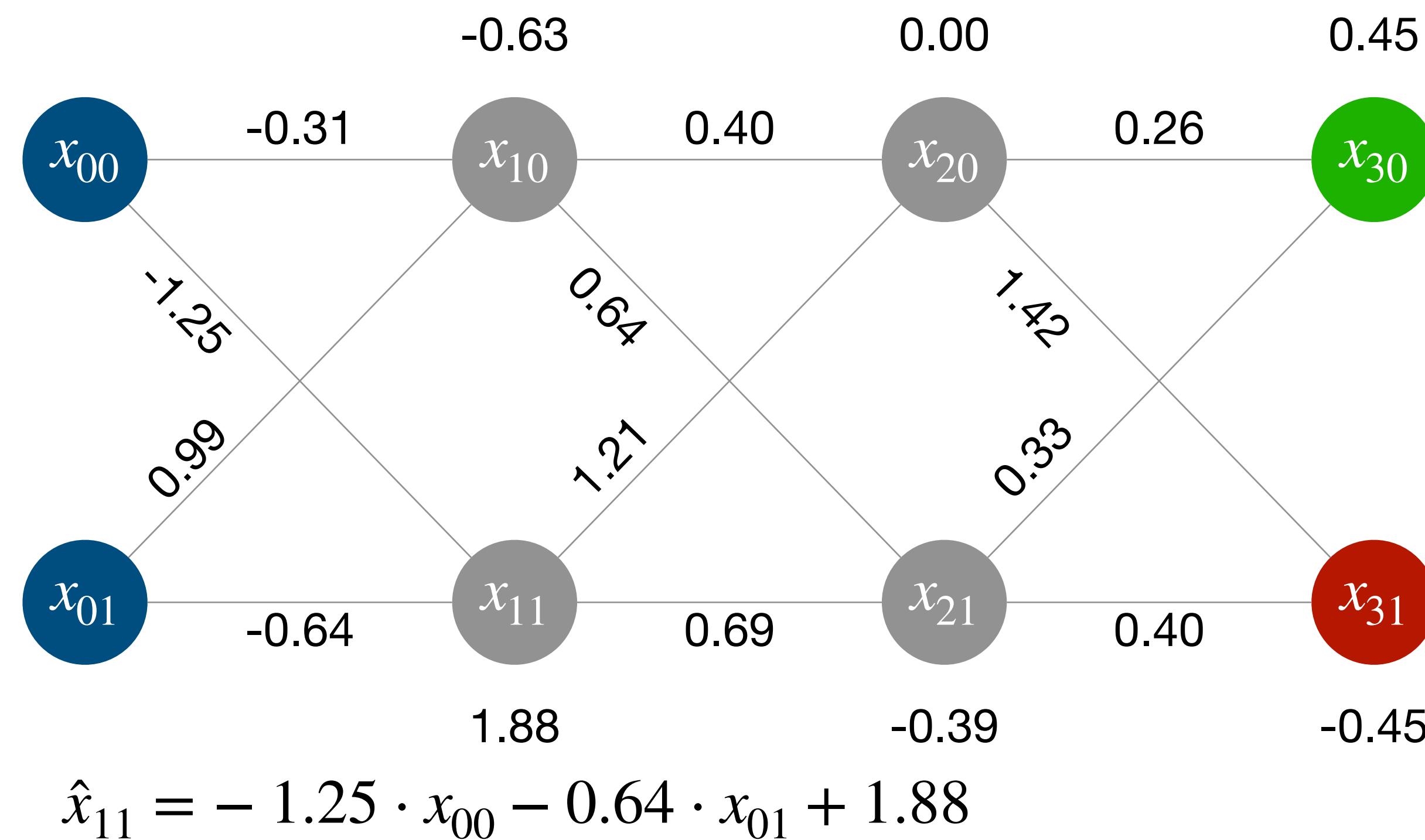


# Feed-Forward Neural Networks with ReLU Activations

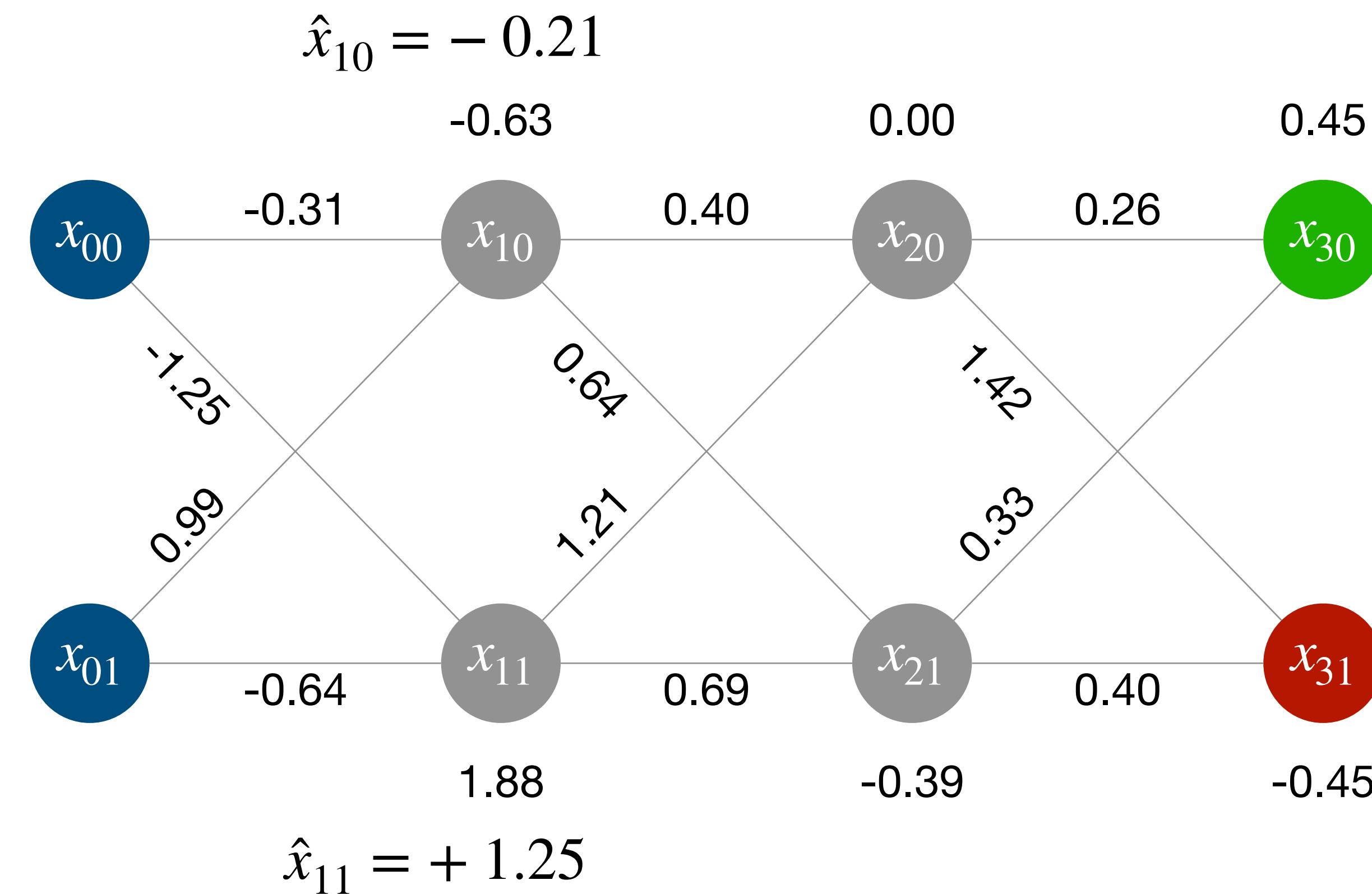


# Feed-Forward Neural Networks with ReLU Activations

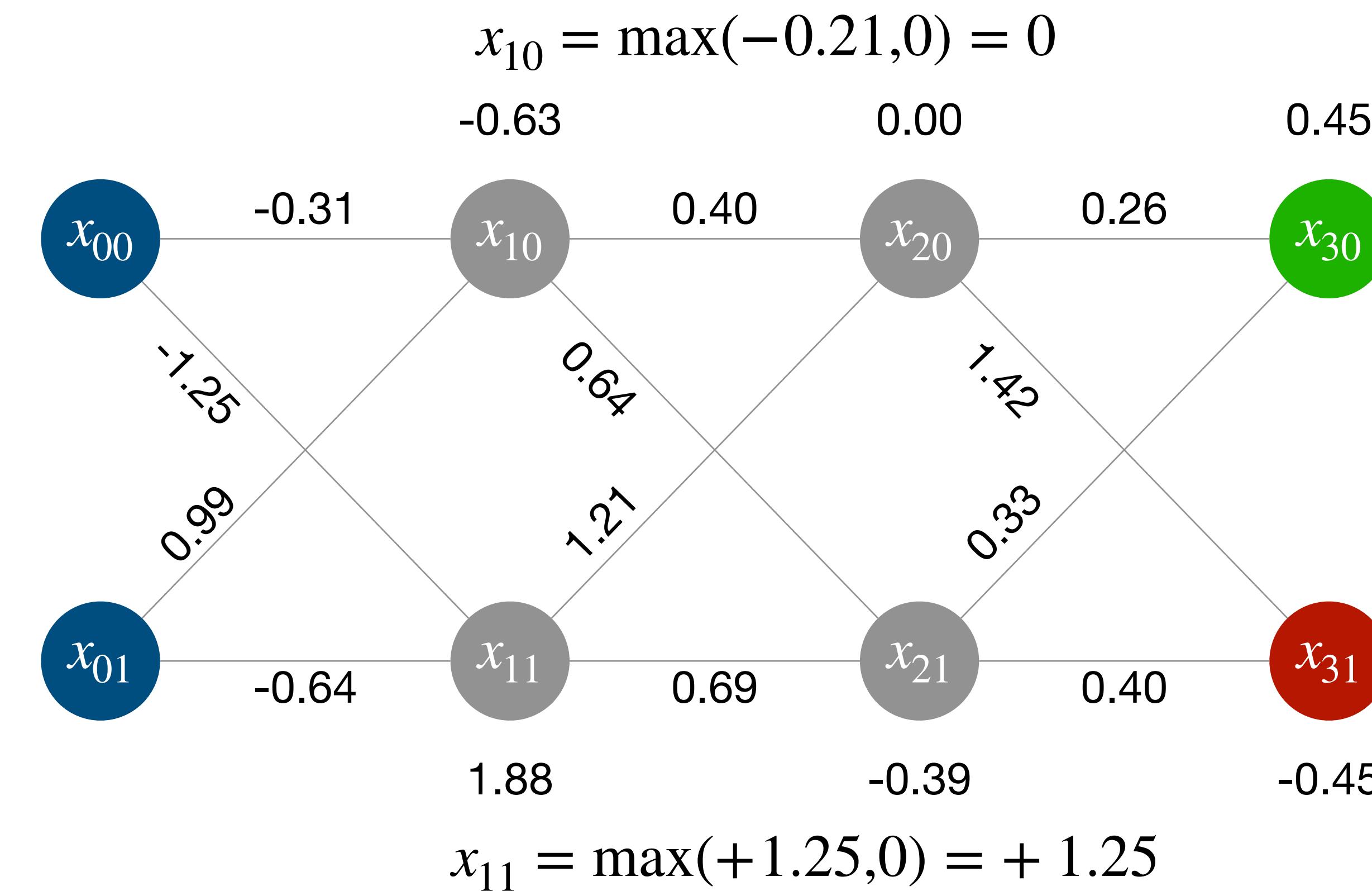
$$\hat{x}_{10} = -0.31 \cdot x_{00} + 0.99 \cdot x_{01} - 0.63$$



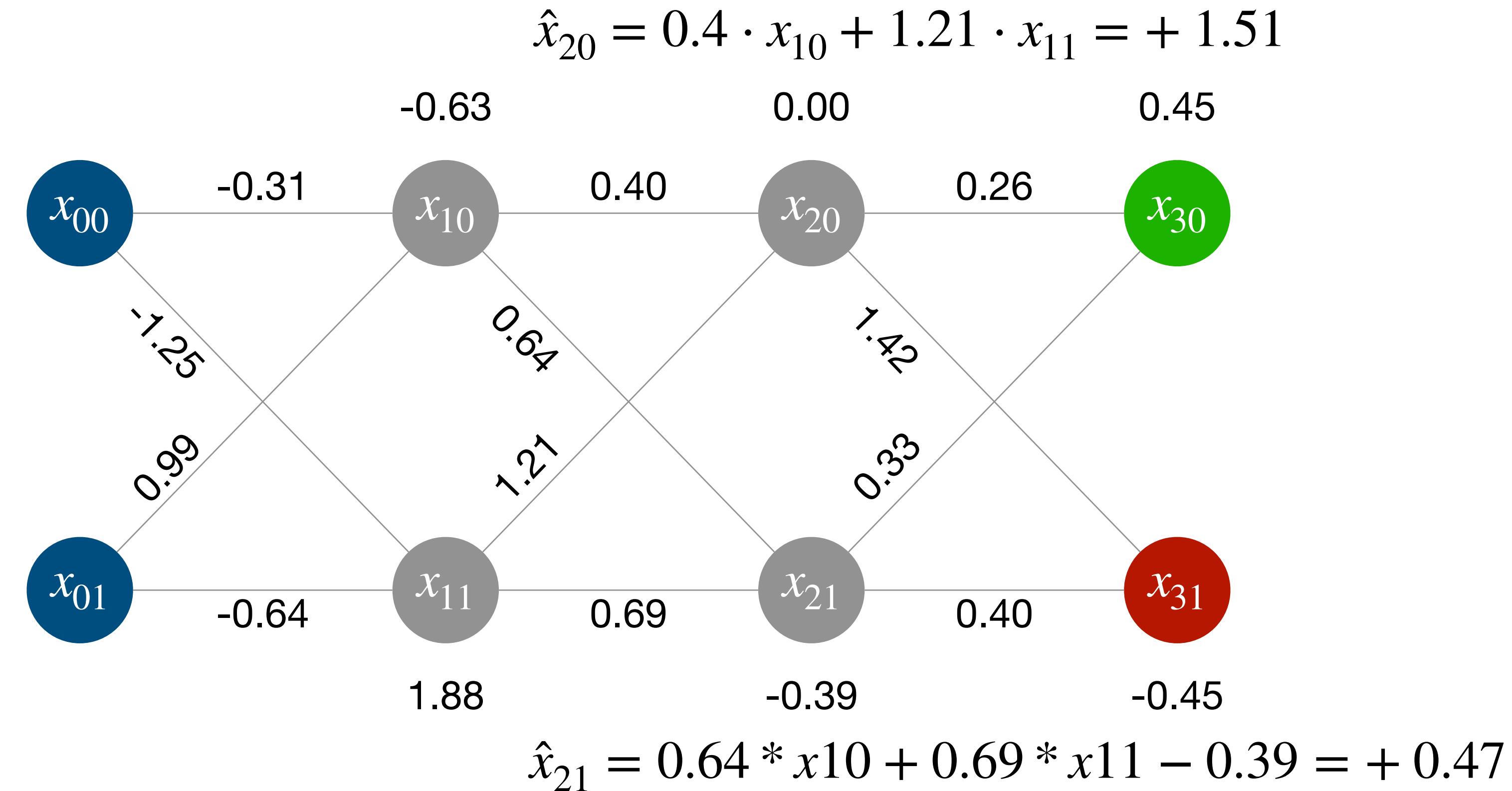
# Feed-Forward Neural Networks with ReLU Activations



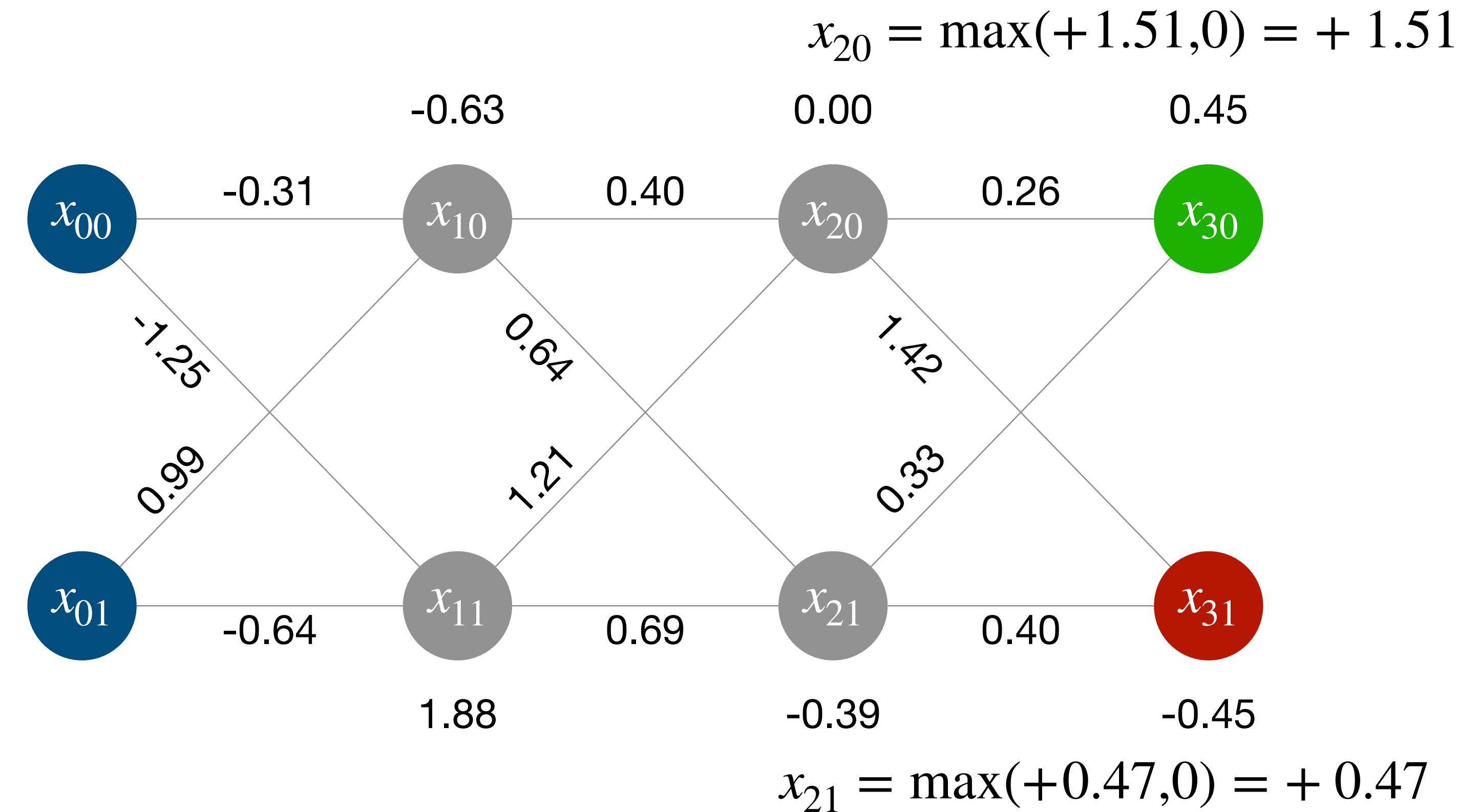
# Feed-Forward Neural Networks with ReLU Activations



# Feed-Forward Neural Networks with ReLU Activations

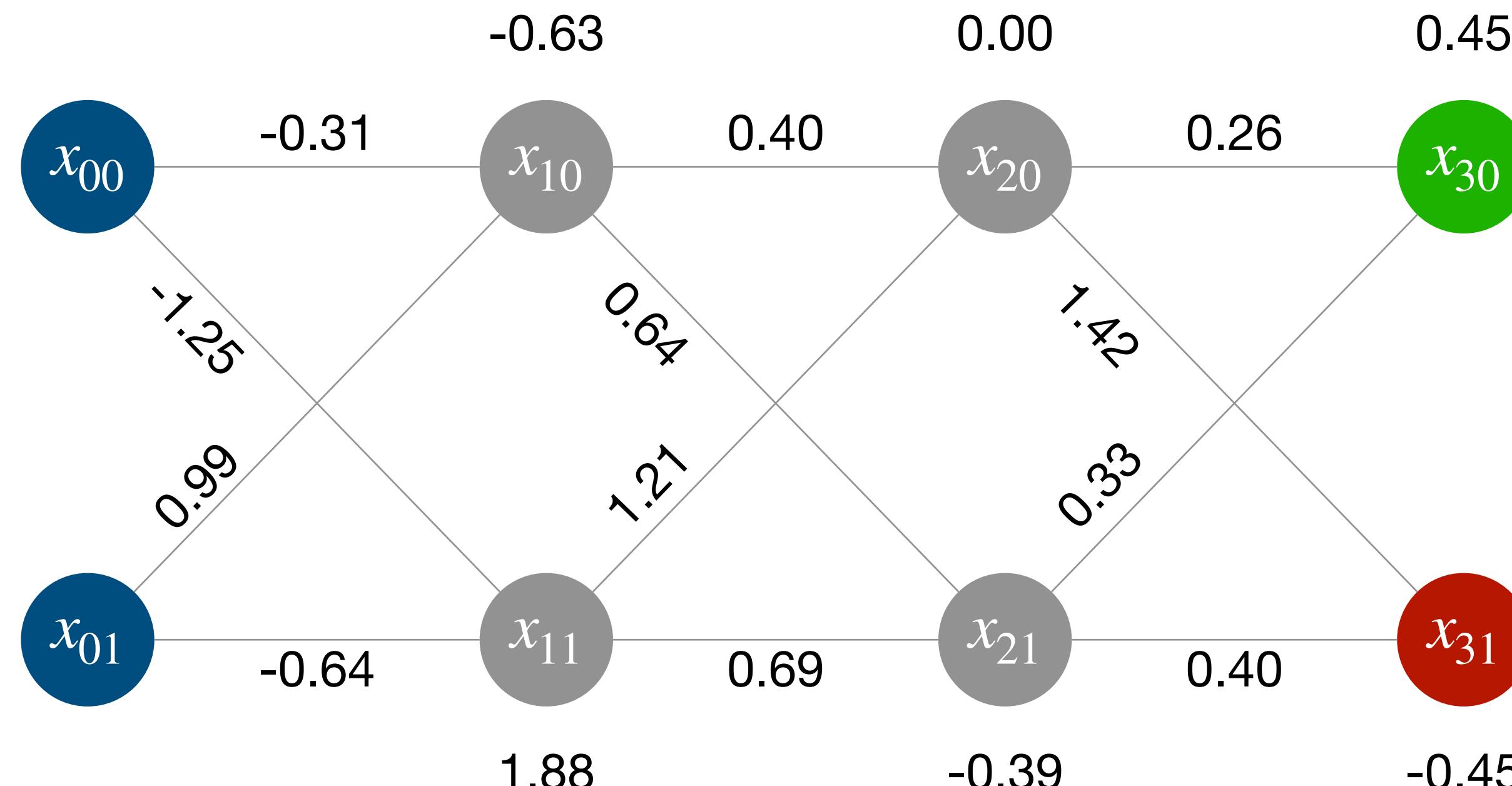


# Feed-Forward Neural Networks with ReLU Activations



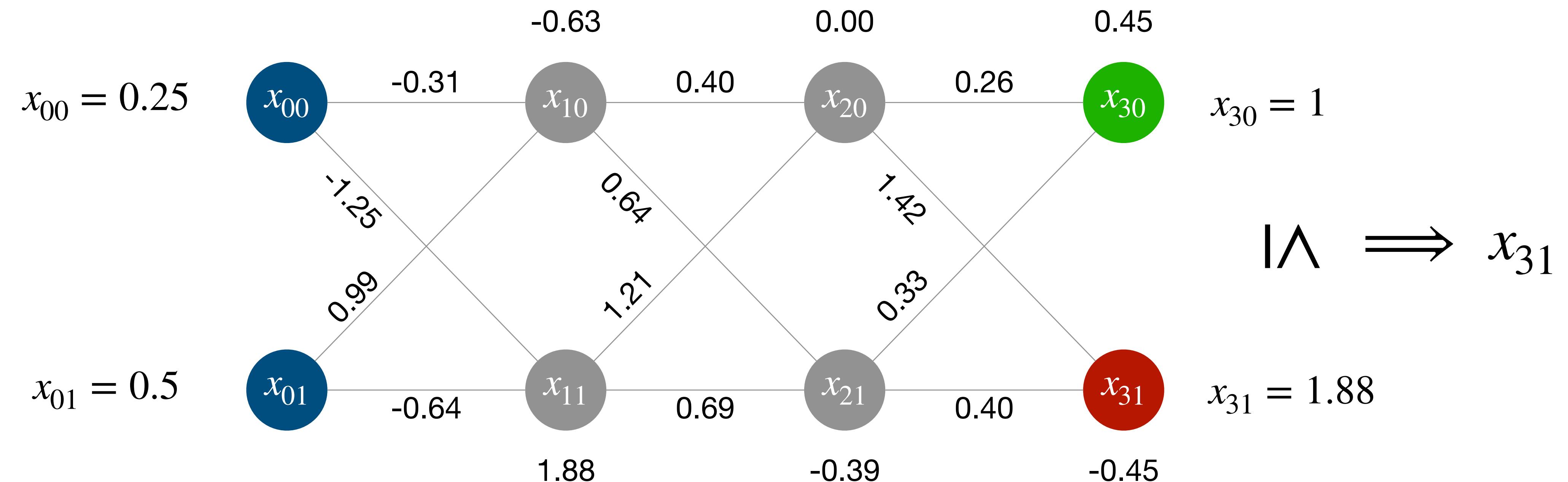
# Feed-Forward Neural Networks with ReLU Activations

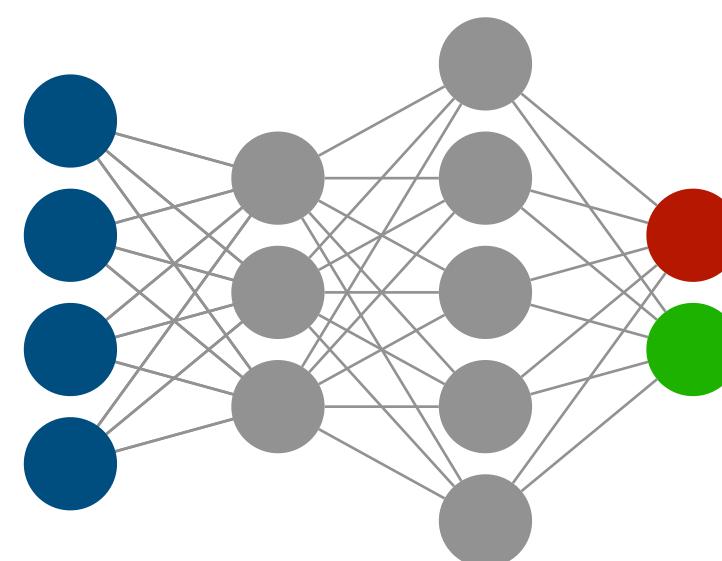
$$\begin{aligned}x_{30} &= 0.26 \cdot x_{20} + 0.33 \cdot x_{21} + 0.45 \\&= 1\end{aligned}$$



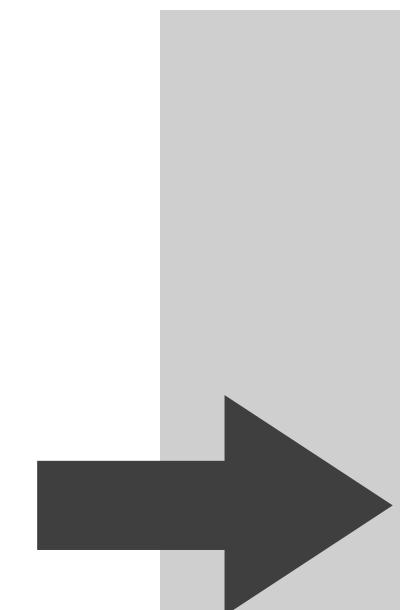
$$\begin{aligned}x_{31} &= 1.42 \cdot x_{20} + 0.4 \cdot x_{21} - 0.45 \\&= 1.88\end{aligned}$$

# Feed-Forward Neural Networks with ReLU Activations

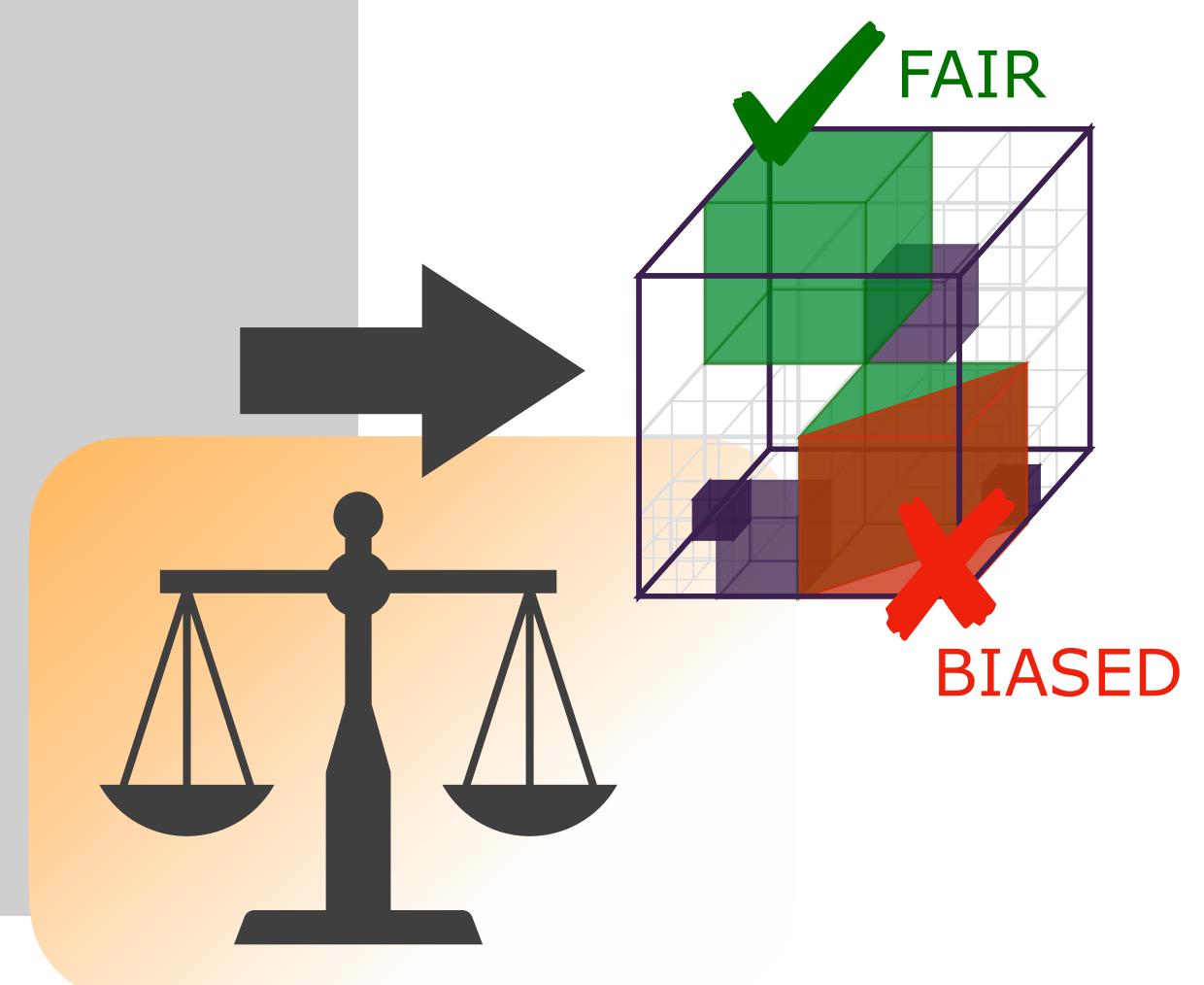




Neural Network



# Libra



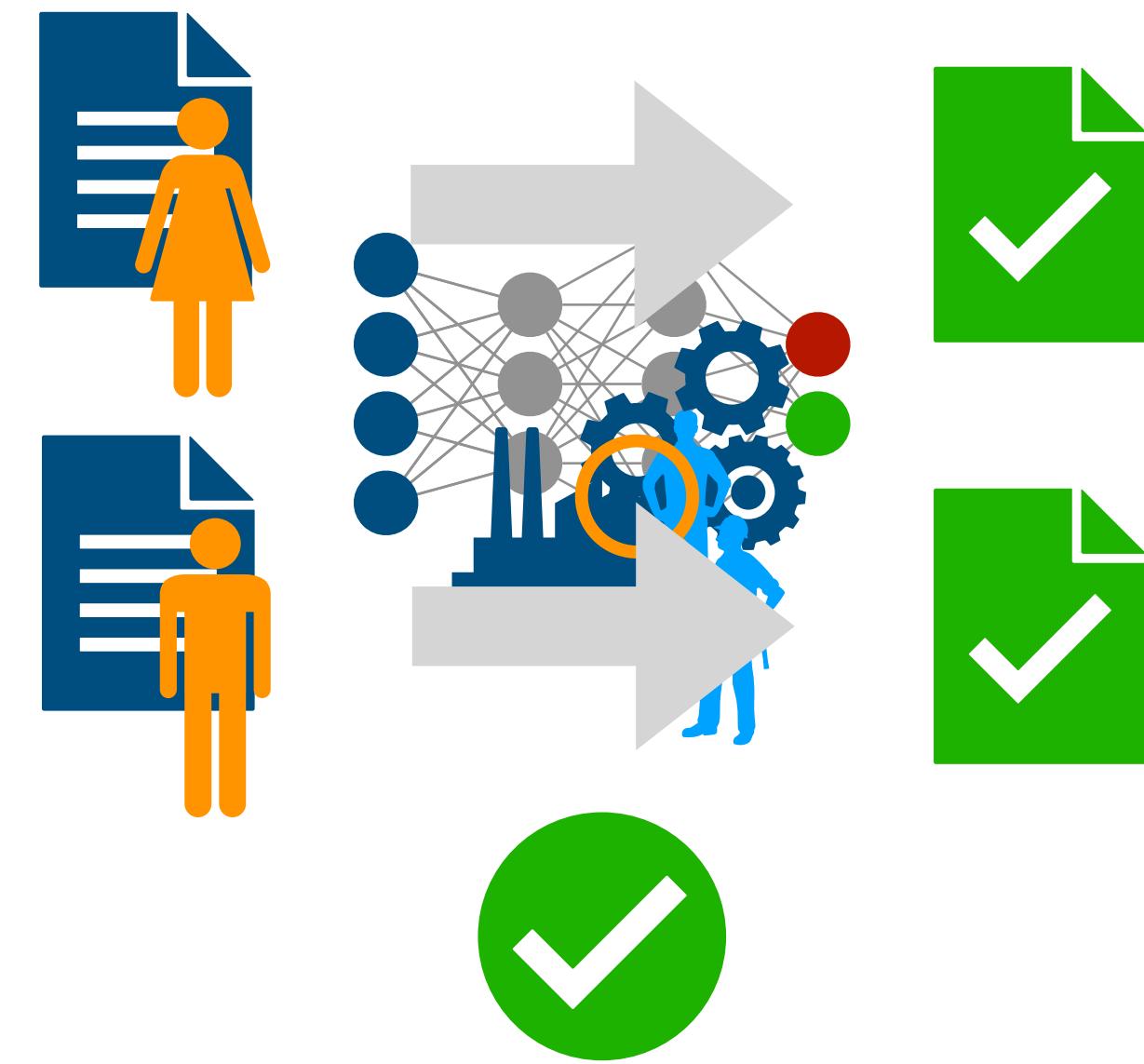
# Dependency **Fairness**

The classification outcome is  
**Independent** on the  
**Sensitive Features**

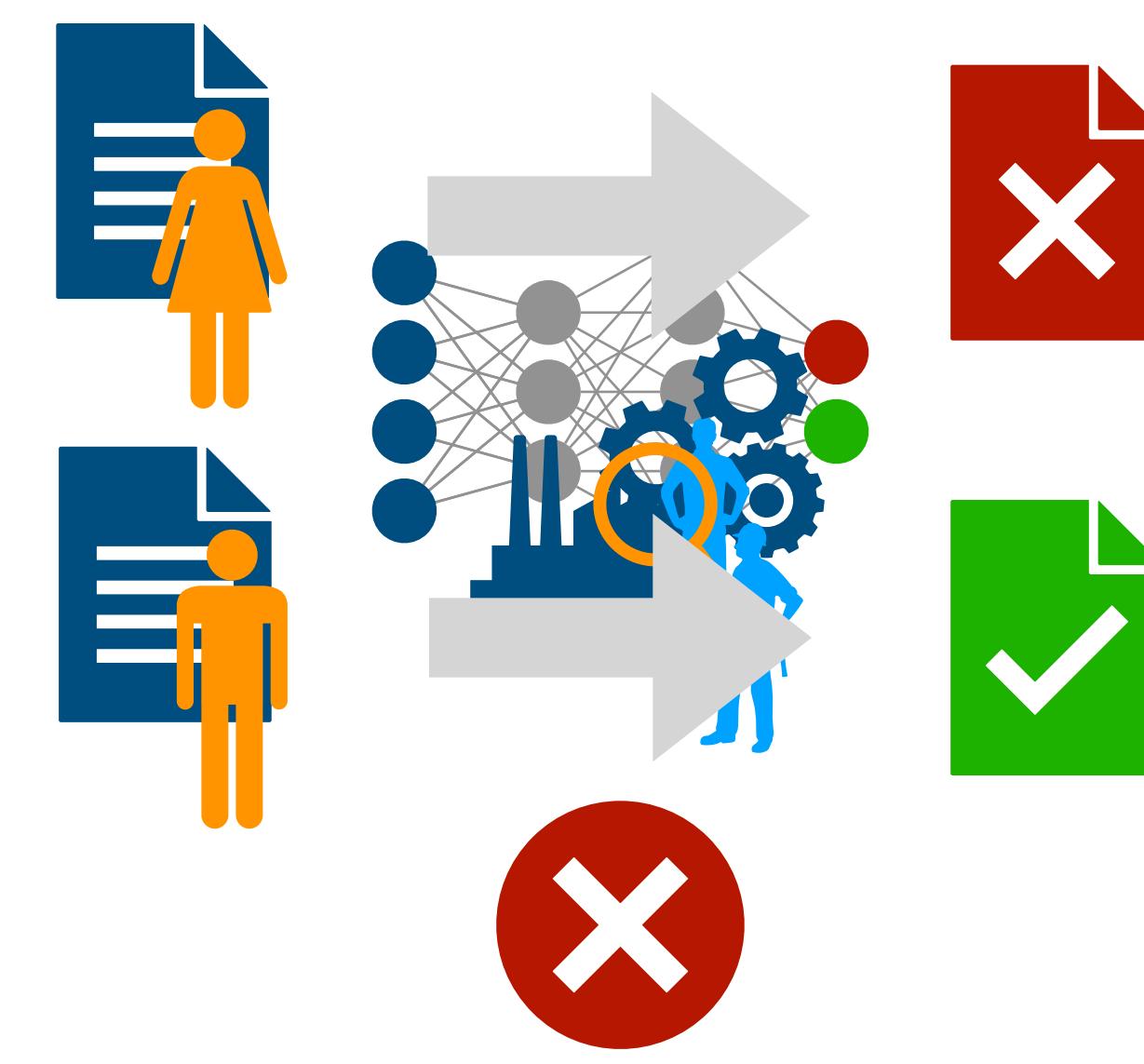
# Recruiting Process



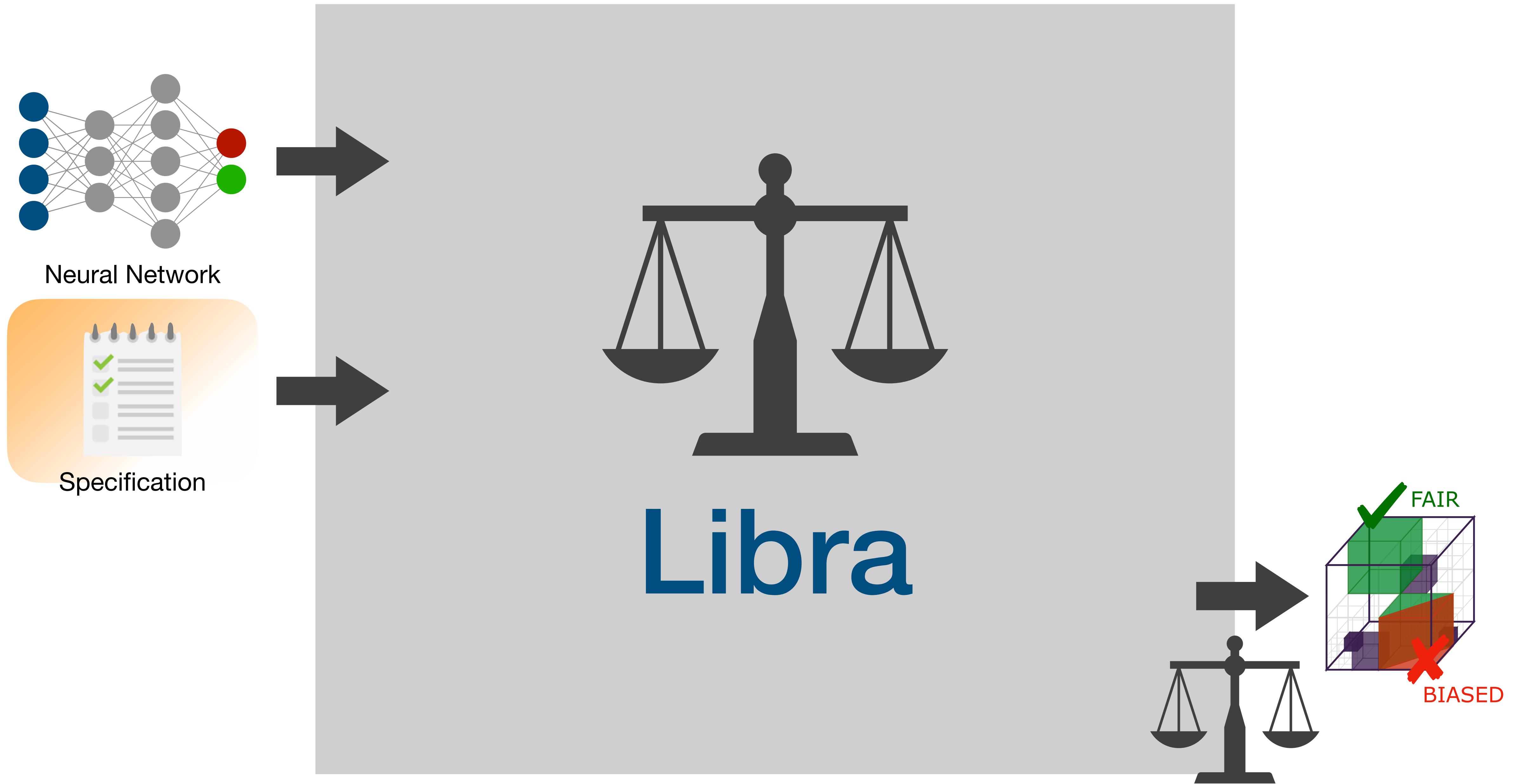
# Recruiting Process



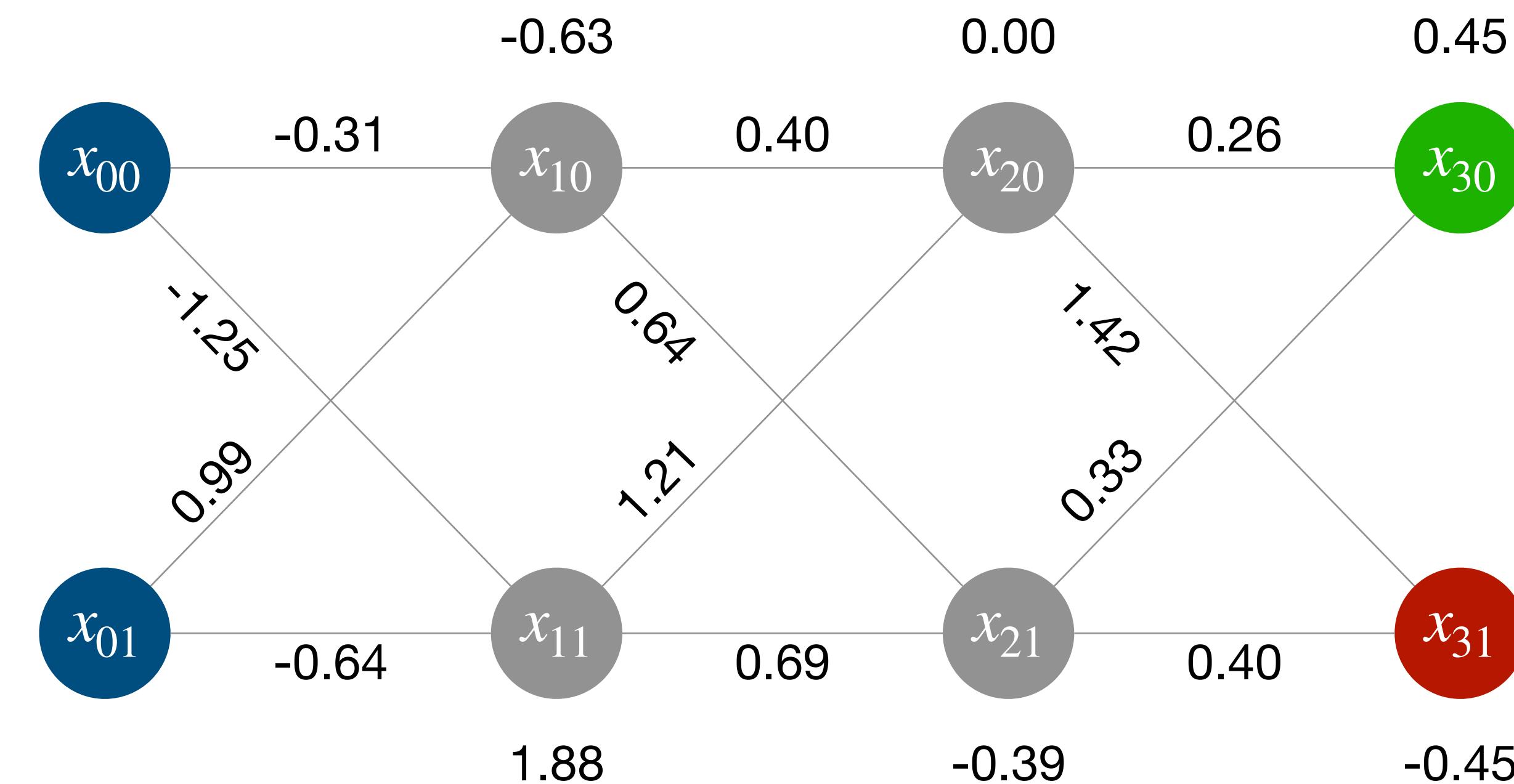
Fair



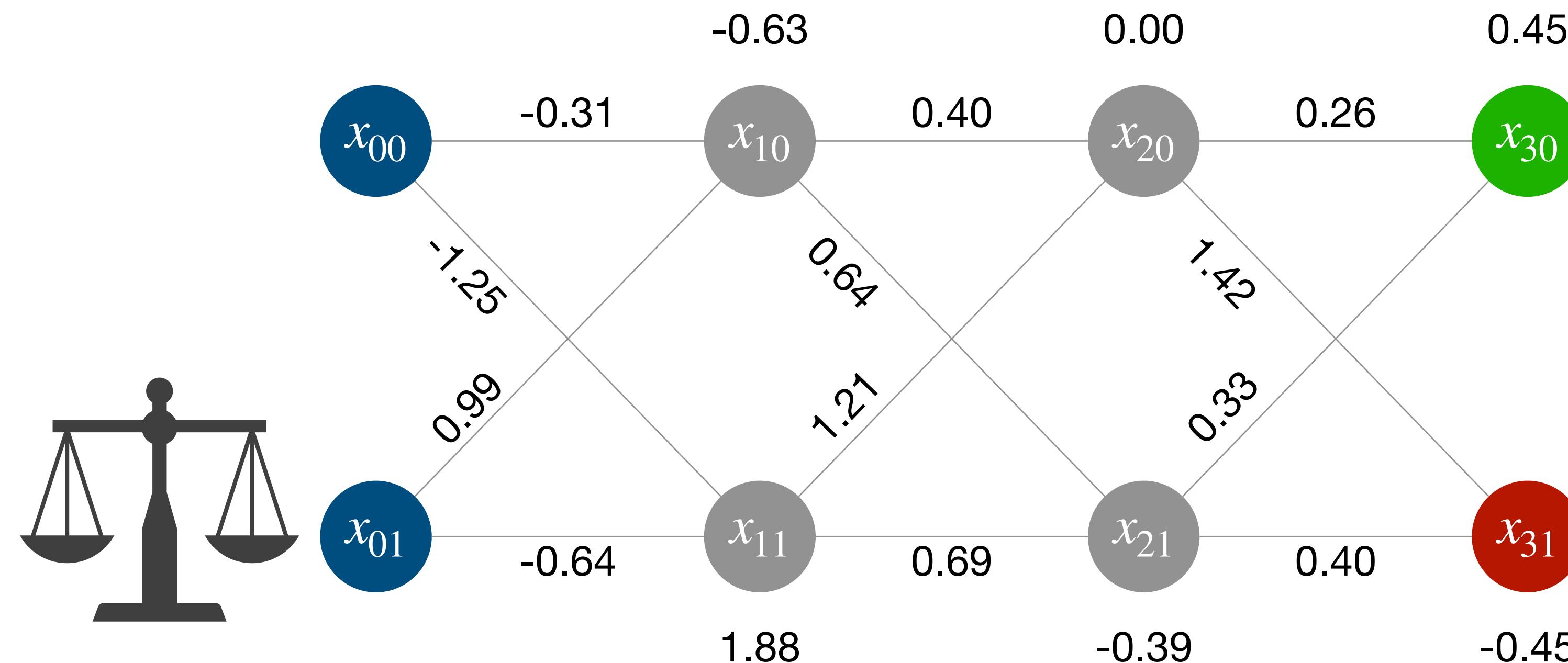
Unfair

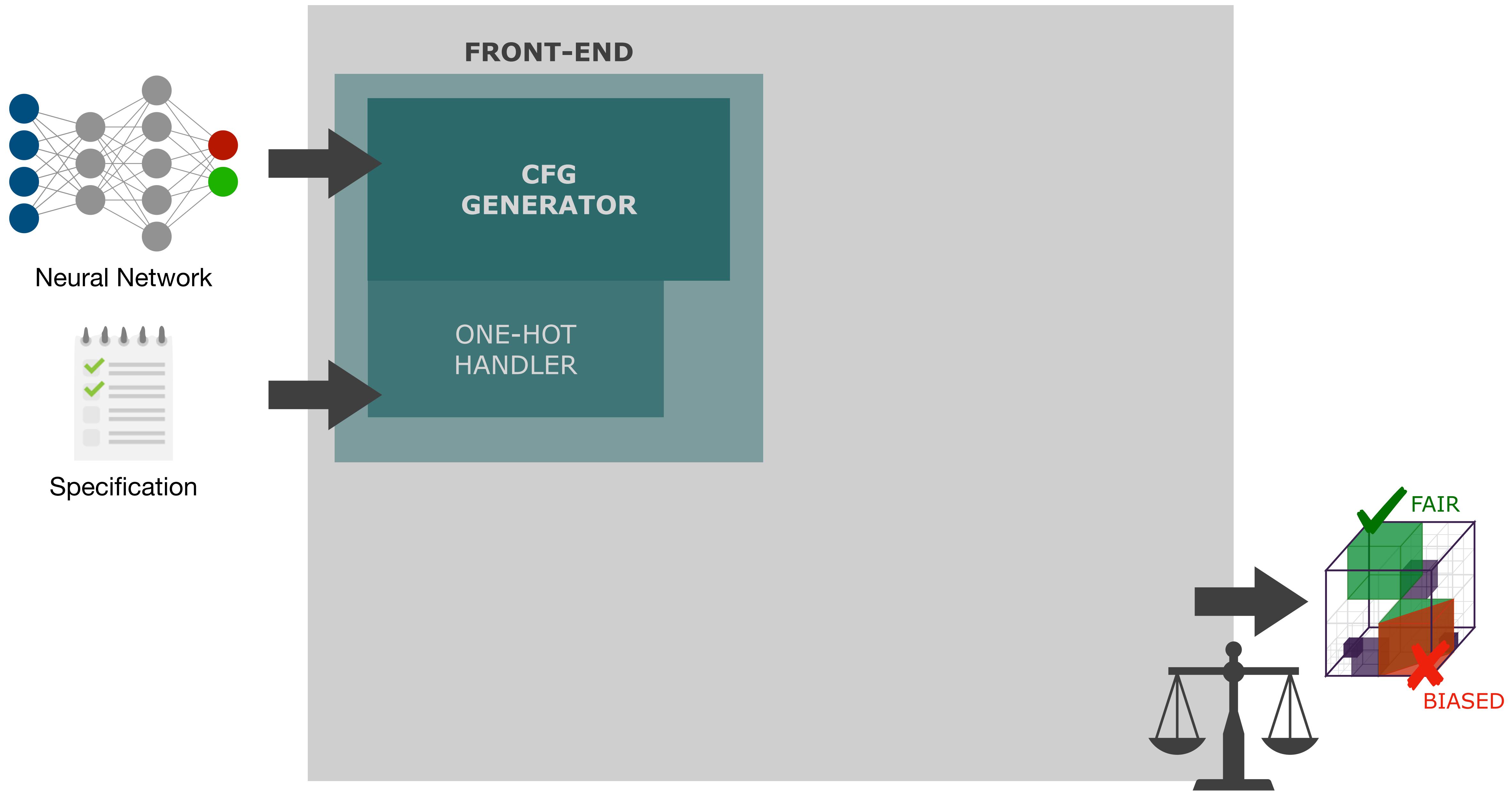


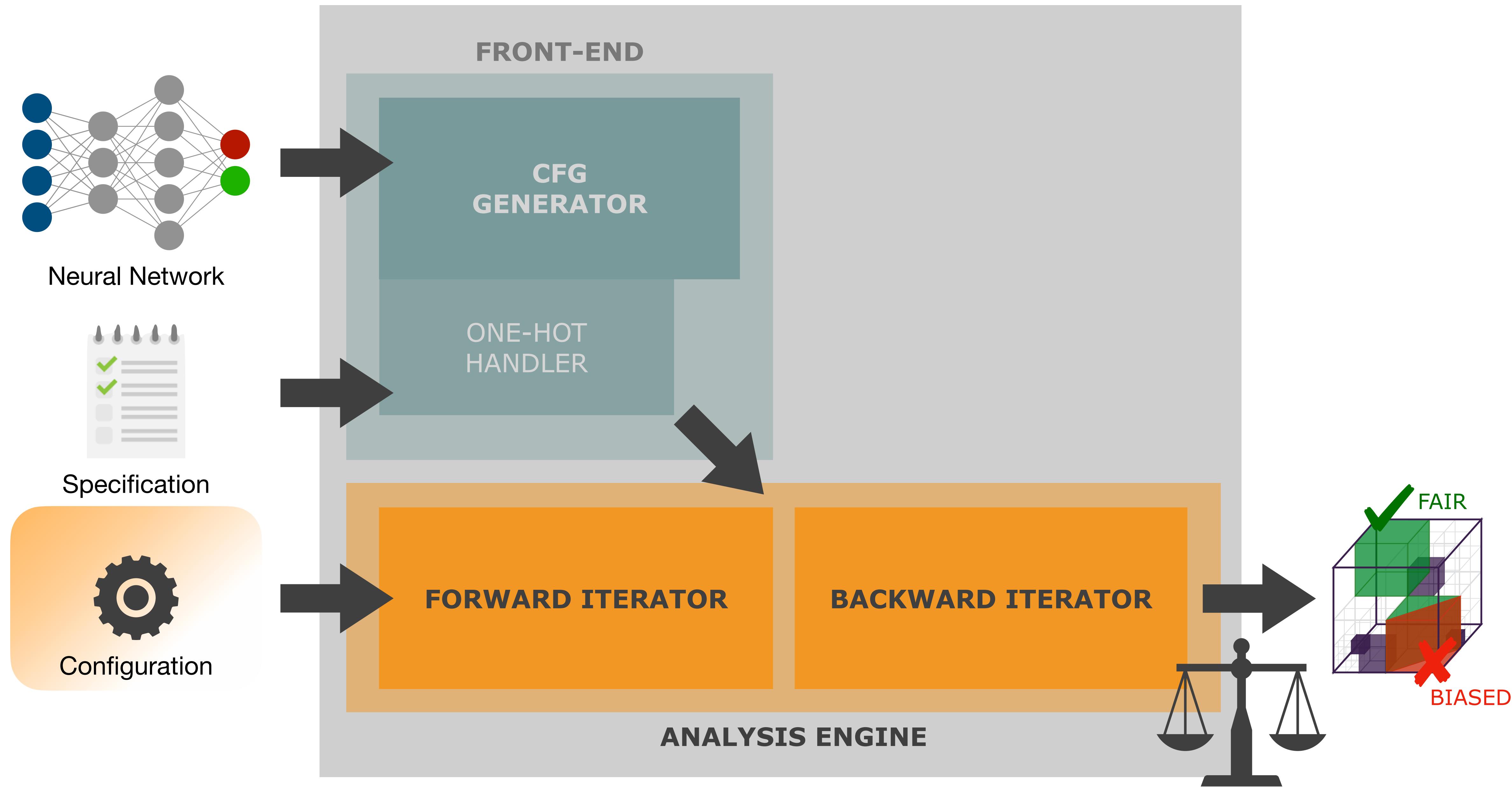
# Specification



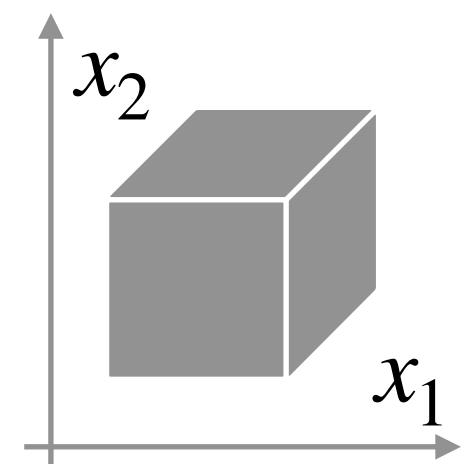
# Specification



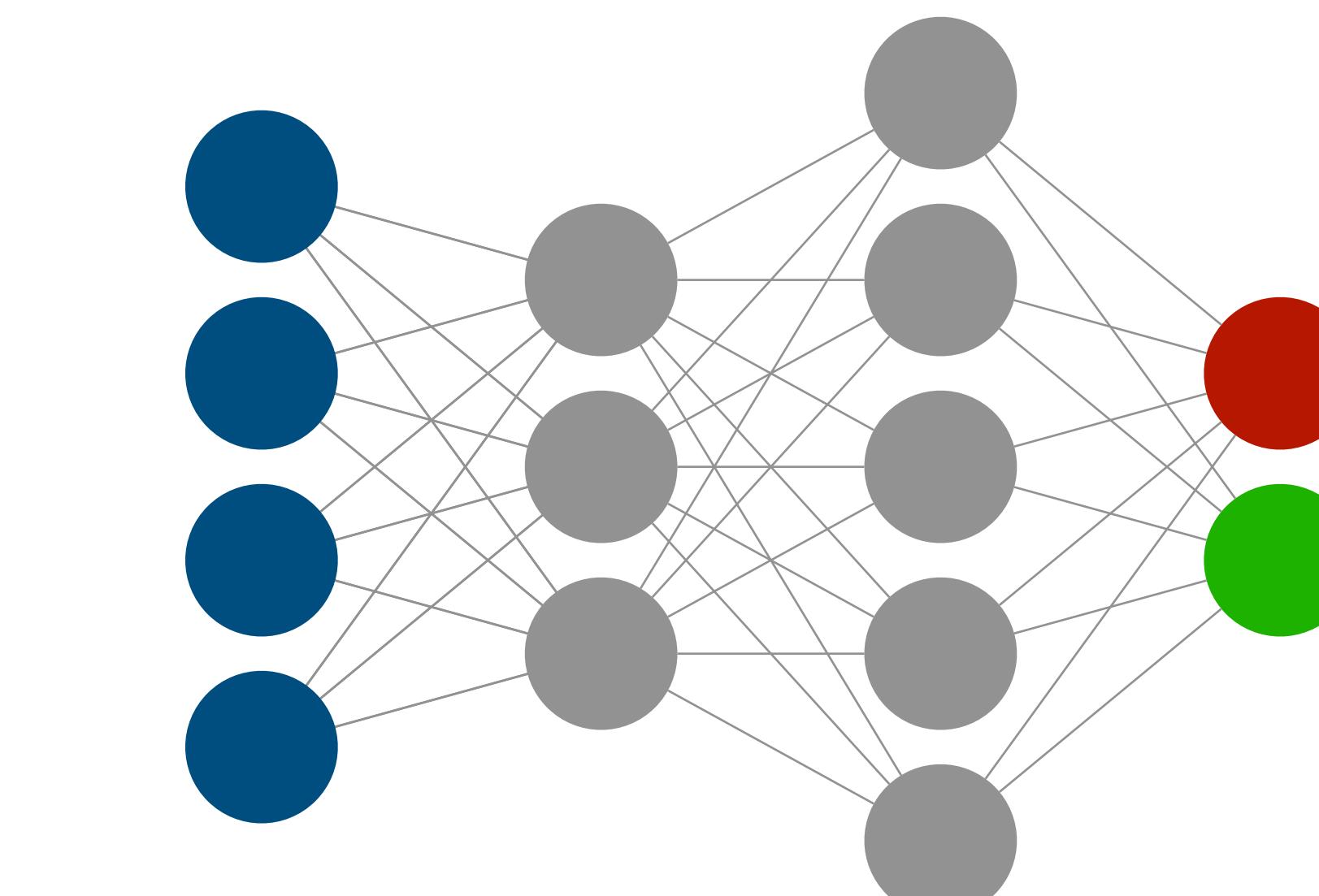
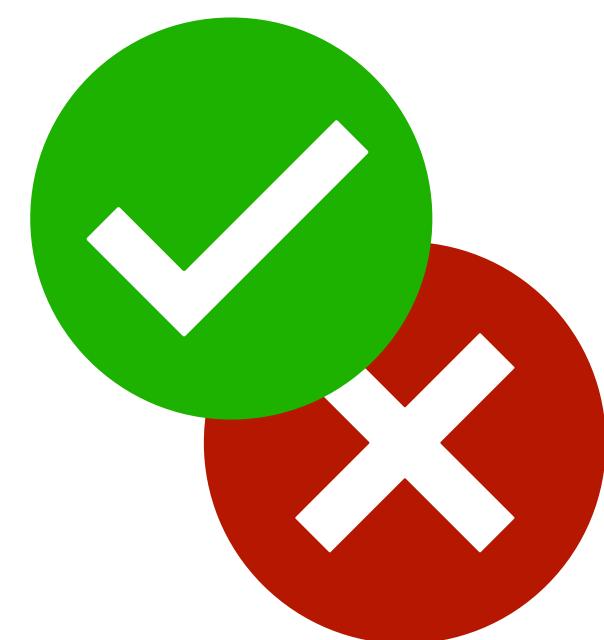




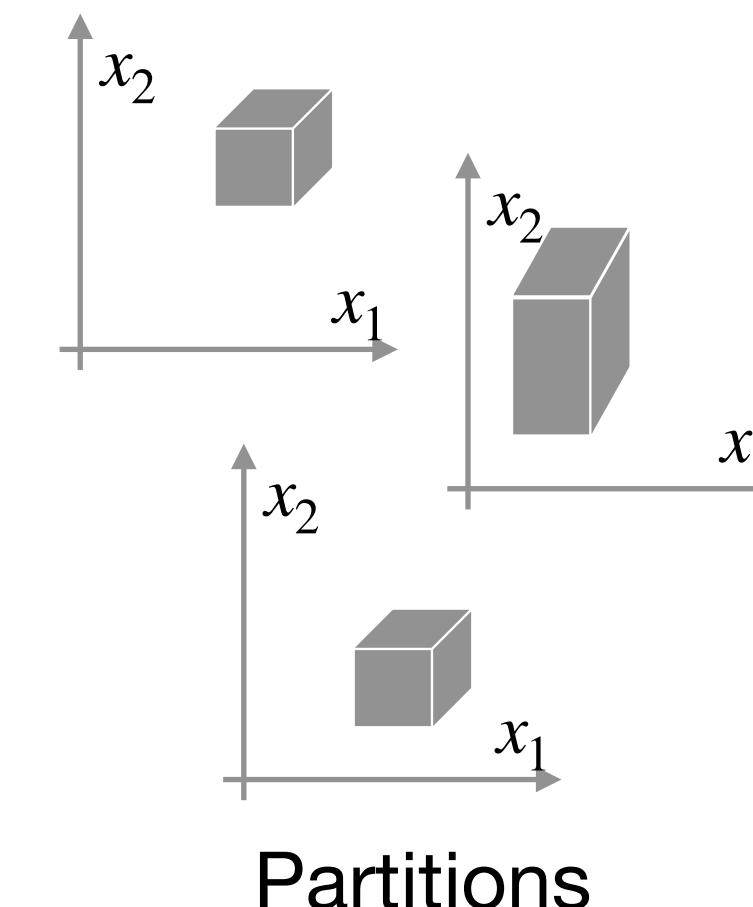
# Cheap Forward Pre-Analysis



Input Space

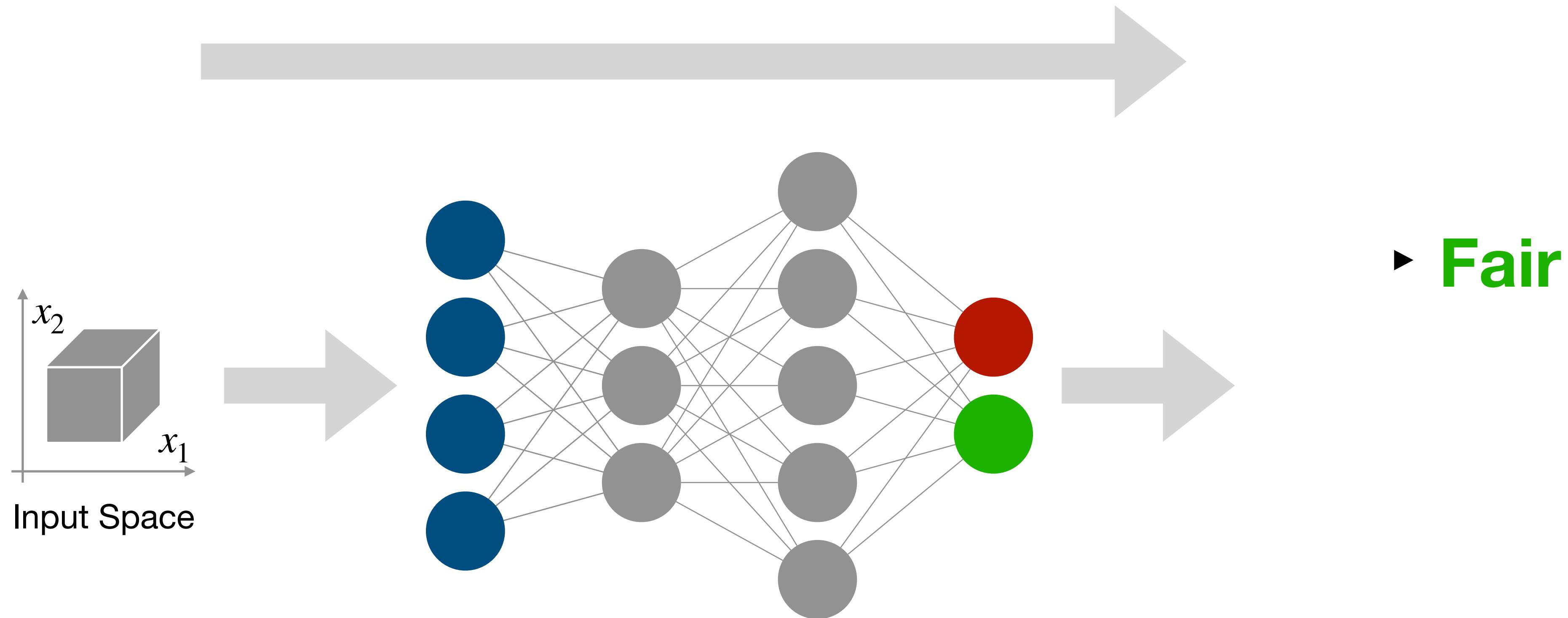


**Exact Backward Analysis  
using Polyhedra**



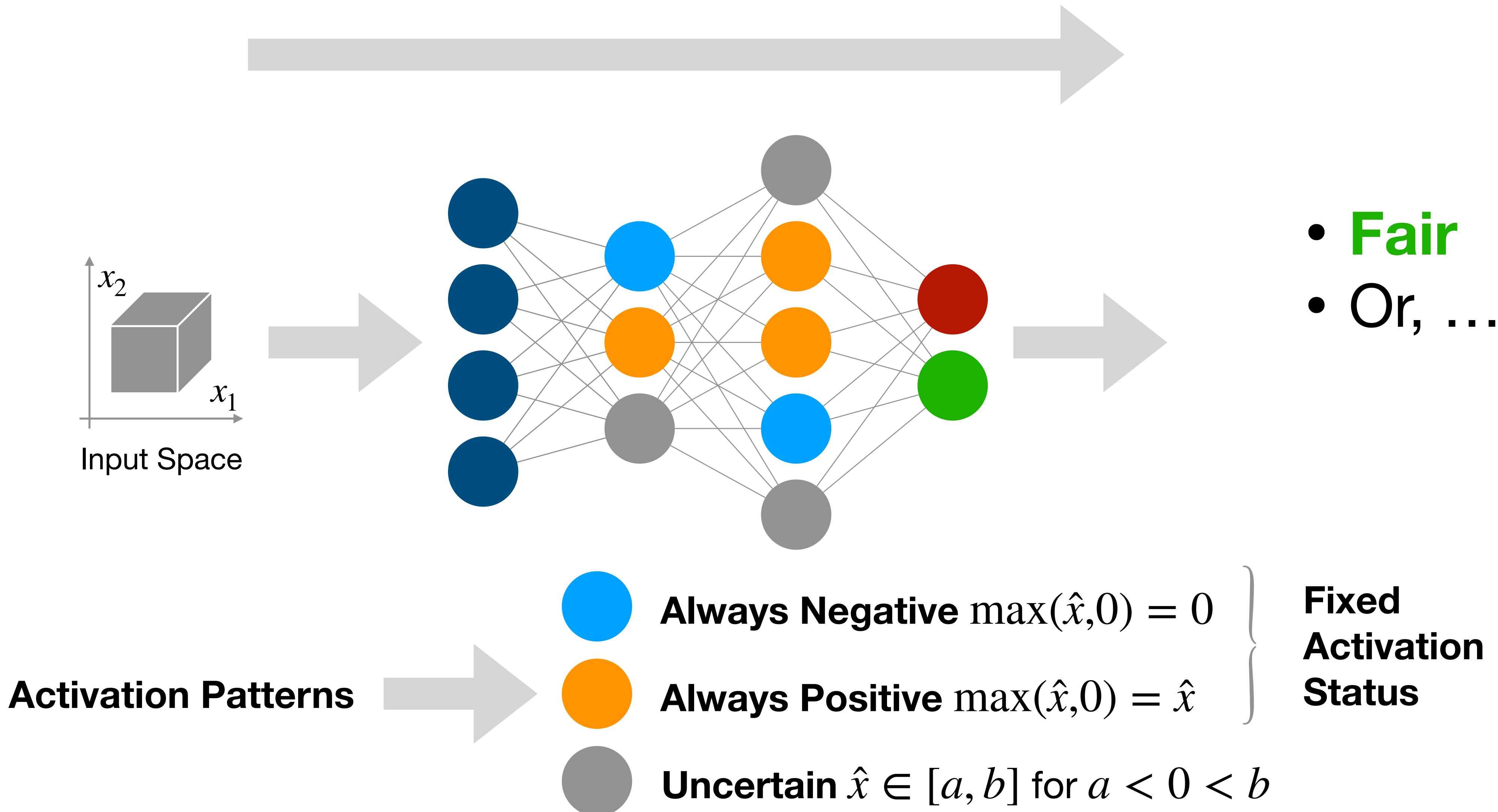
Partitions

# Cheap Forward Pre-Analysis

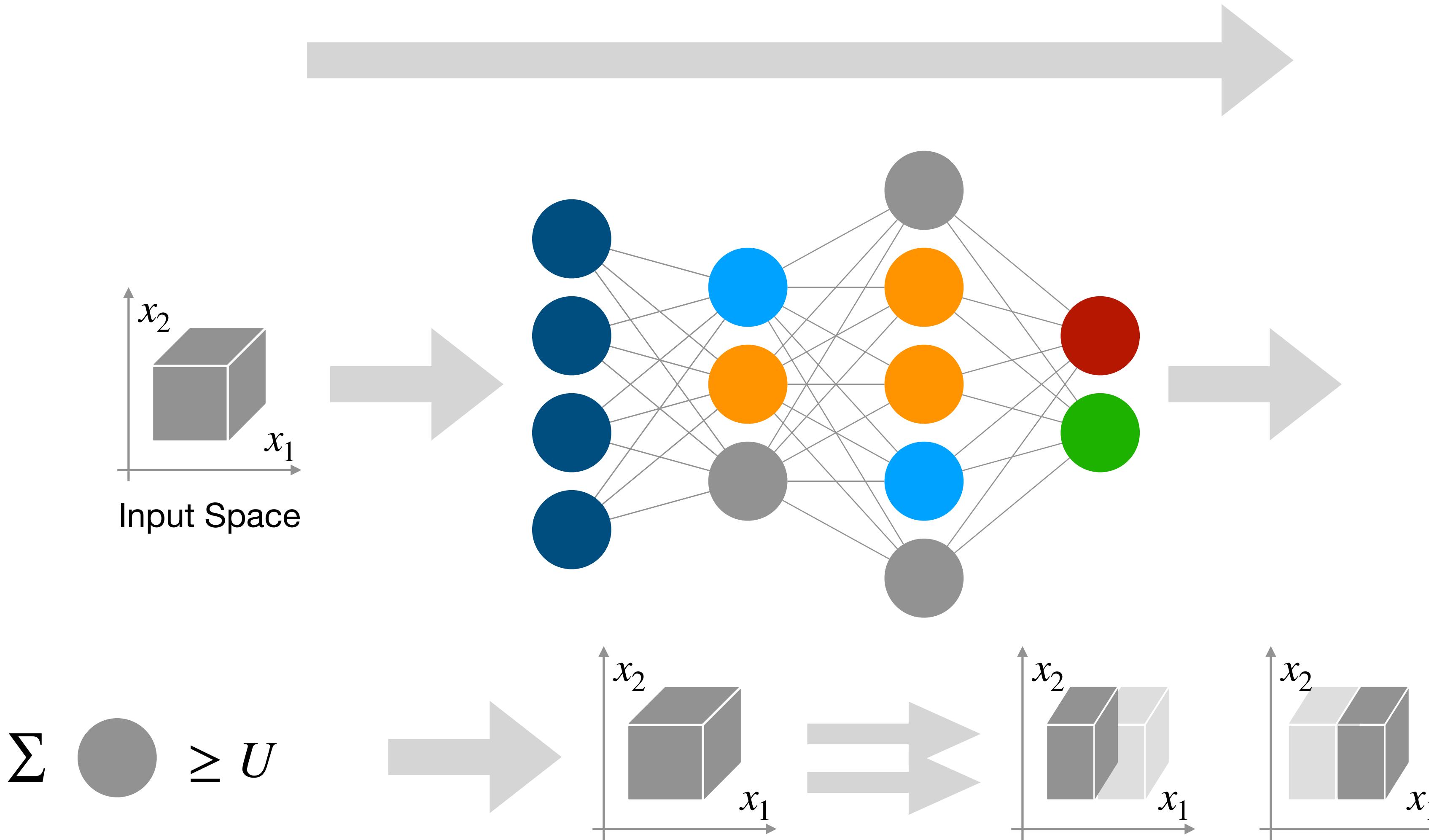


Propagate the partition through the network via **abstract domains**

# Cheap Forward Pre-Analysis



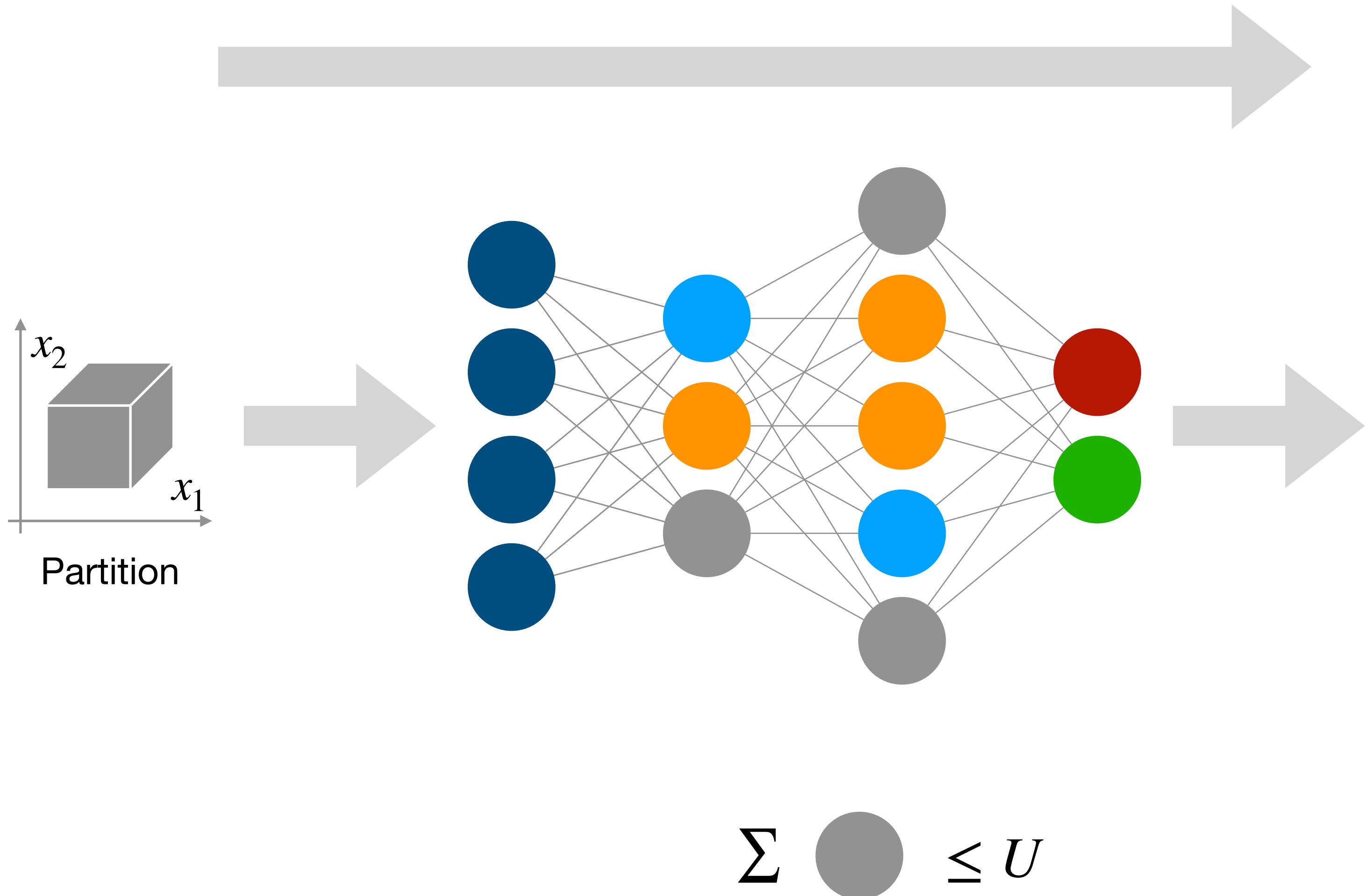
# Cheap Forward Pre-Analysis



- Fair
- Partitioned

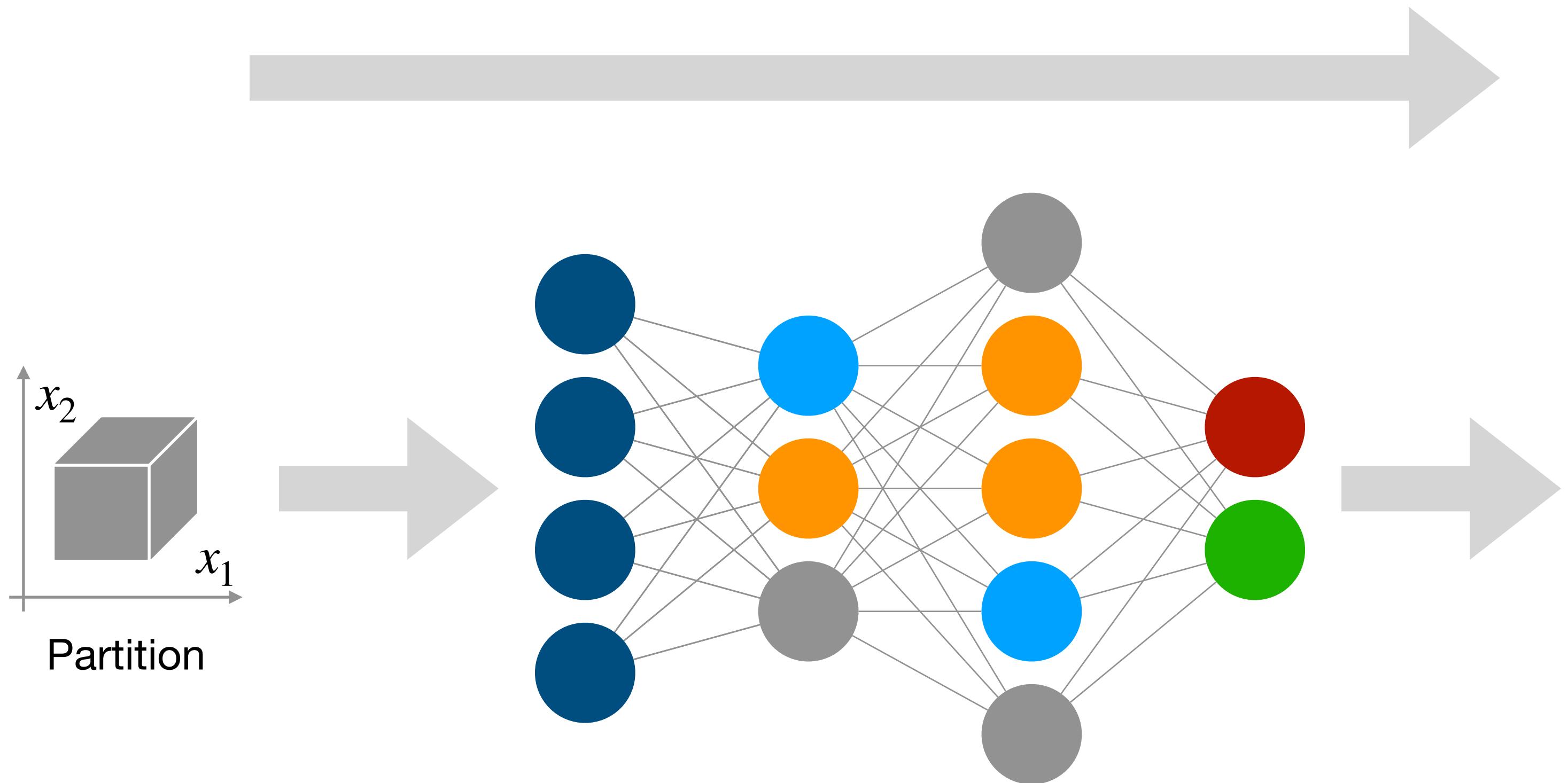
Along non-sensitive  
features only

# Cheap Forward Pre-Analysis



- Fair
- Partitioned
- Feasible

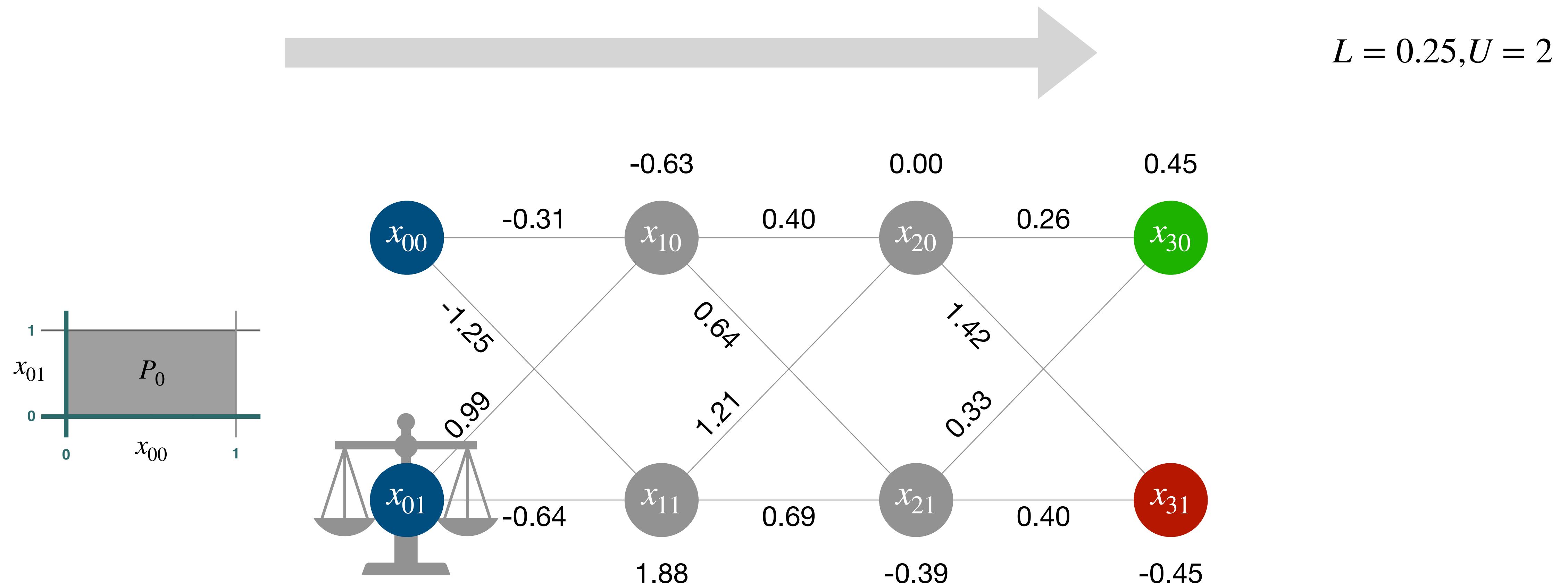
# Cheap Forward Pre-Analysis



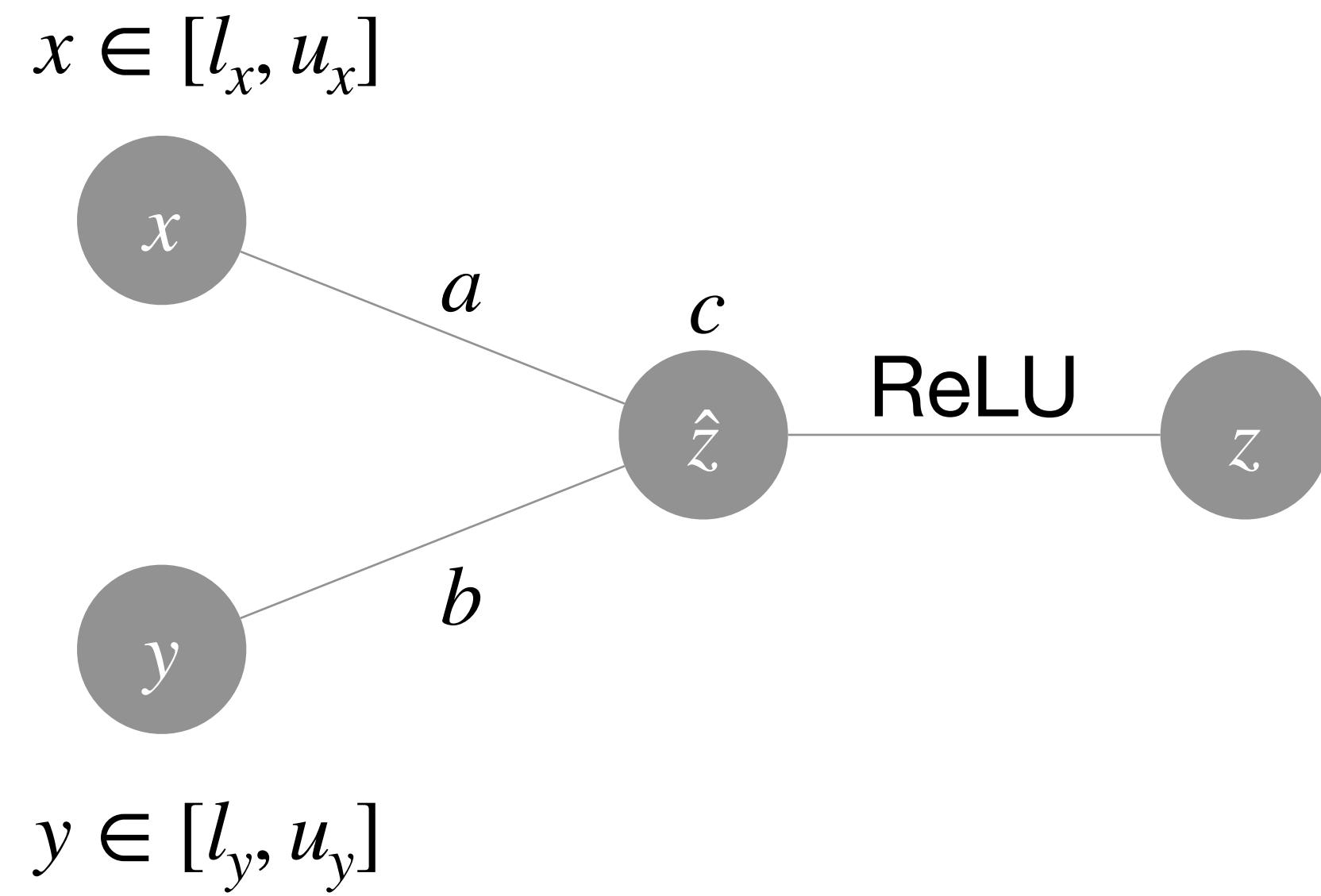
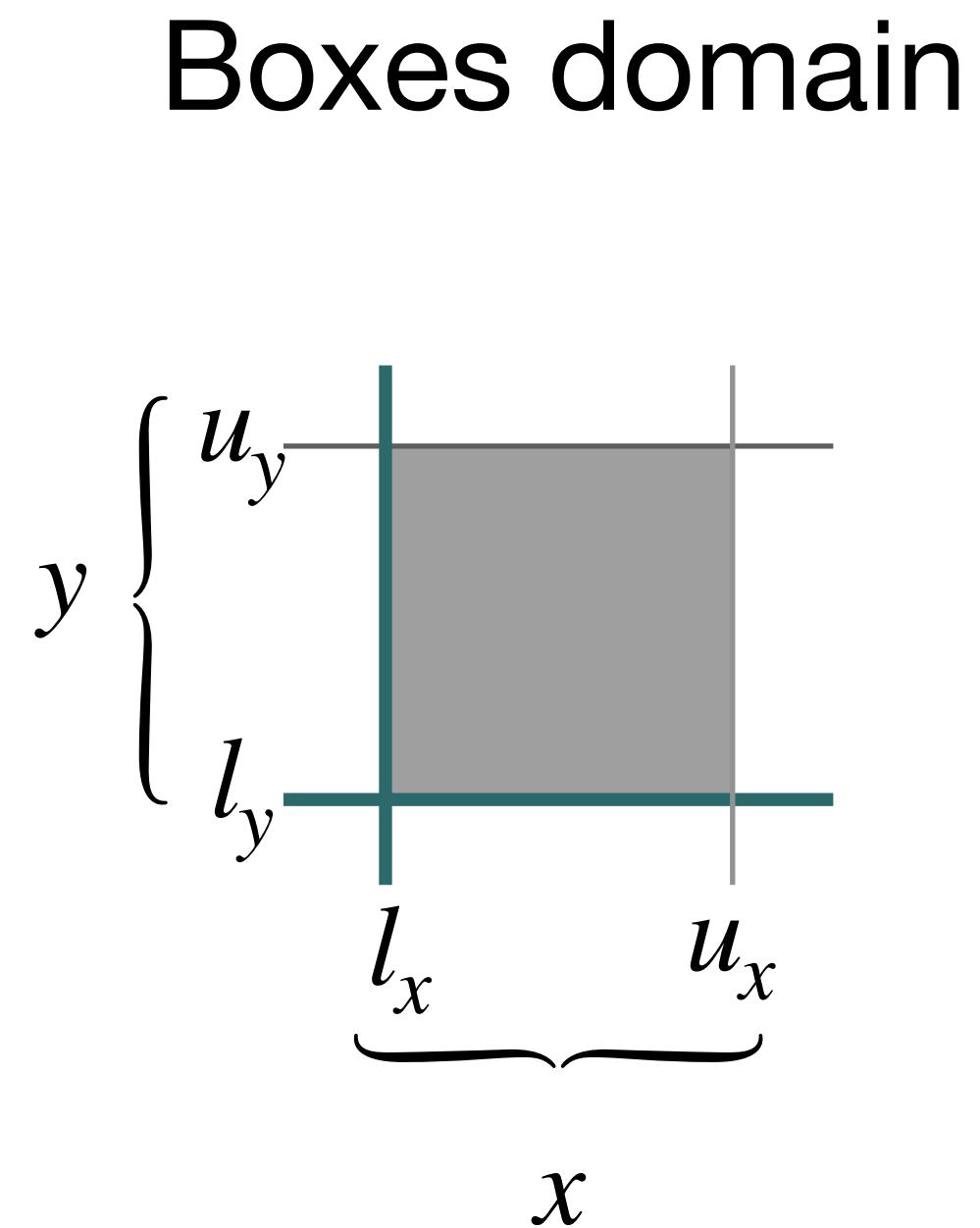
- Fair
- Partitioned
- Feasible
- Excluded

$\sum \bullet \geq U$ , and the partition becomes smaller than  $L$

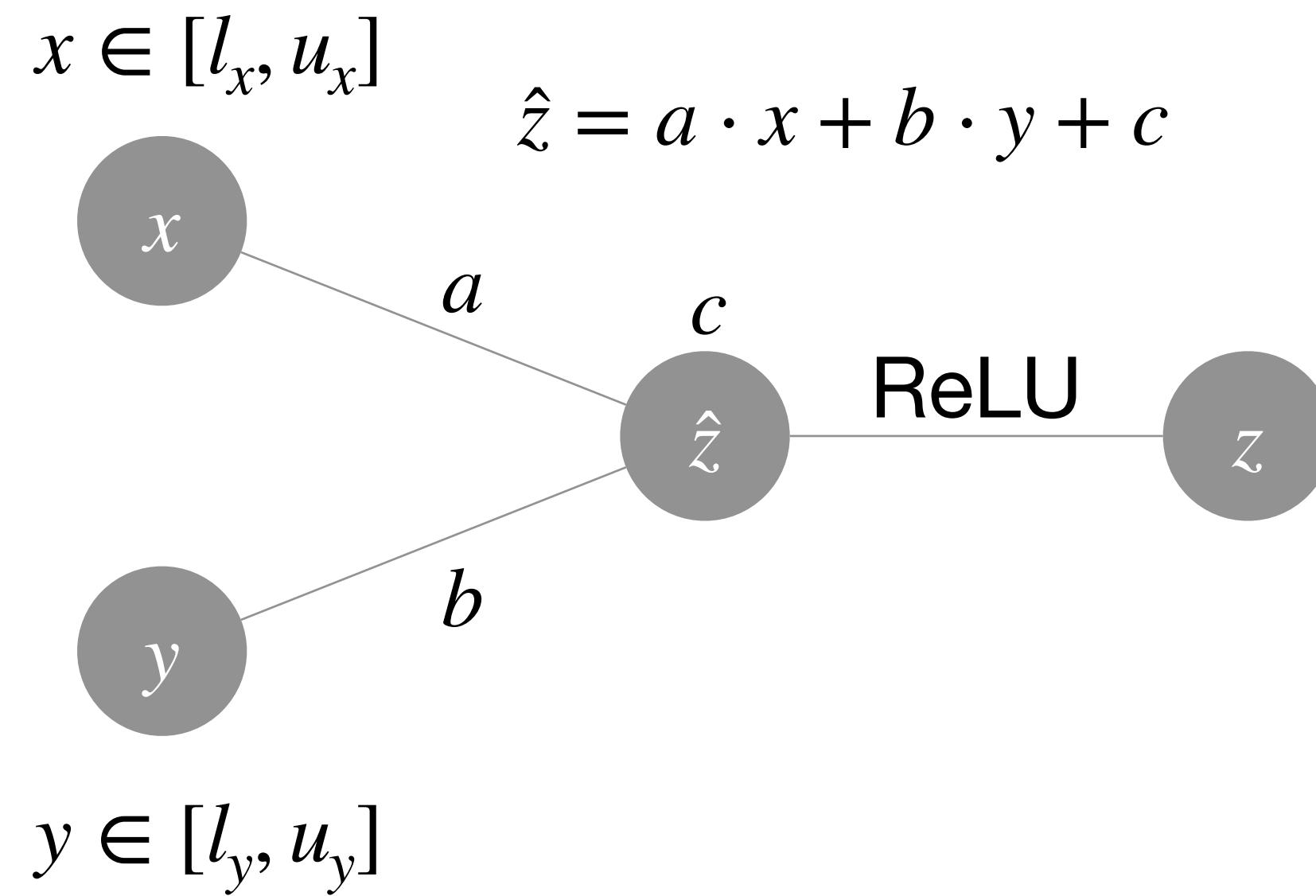
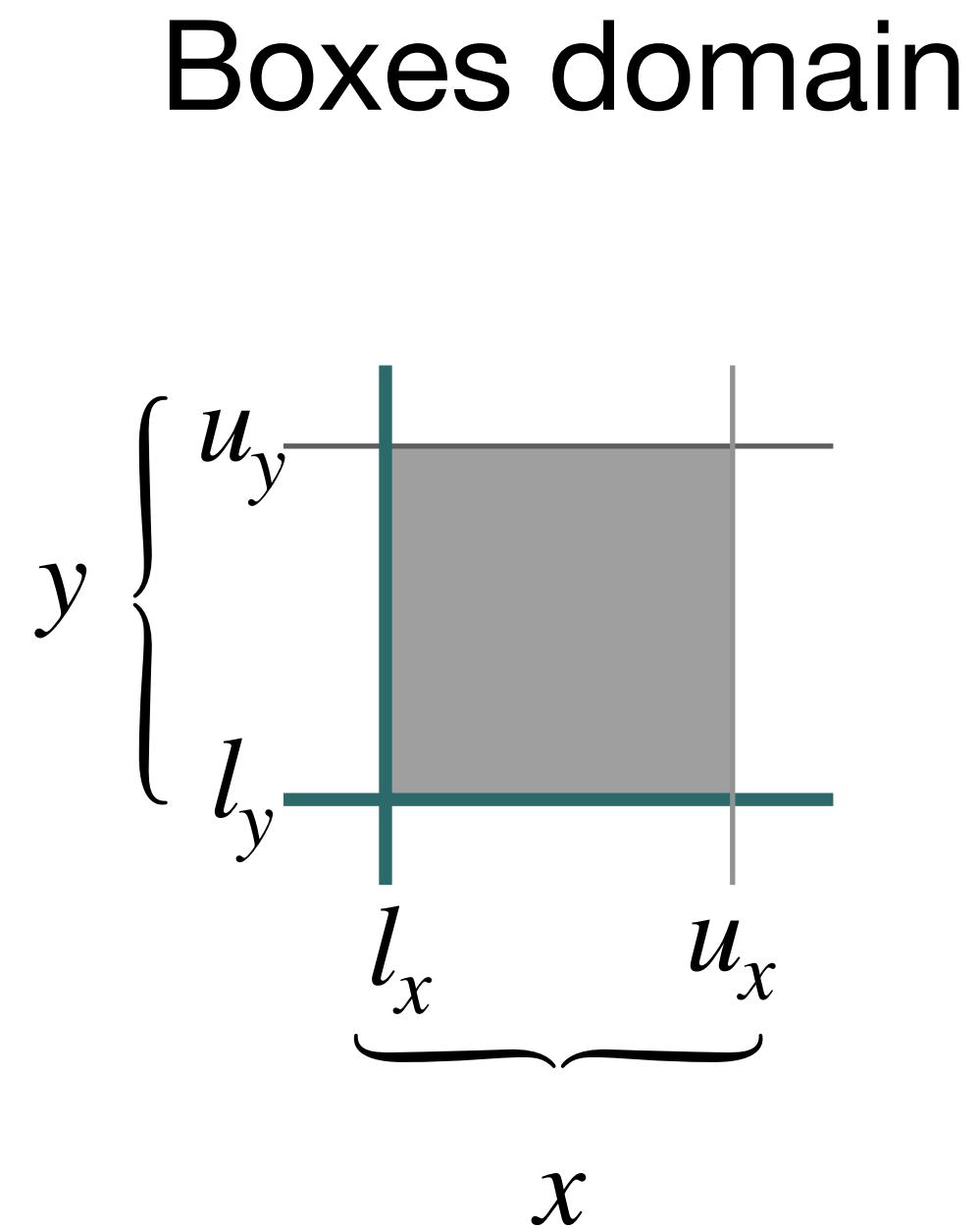
# Forward Analysis



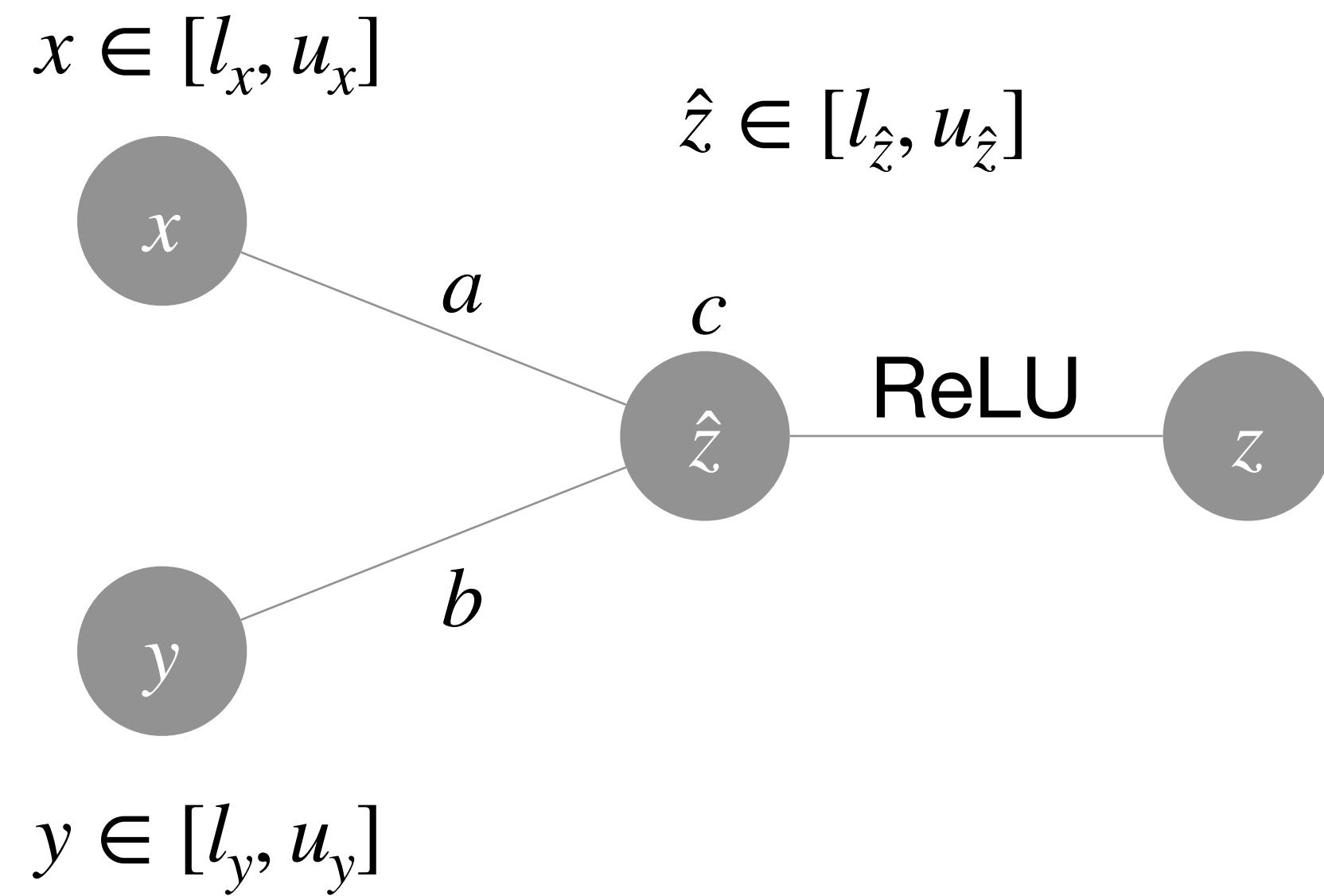
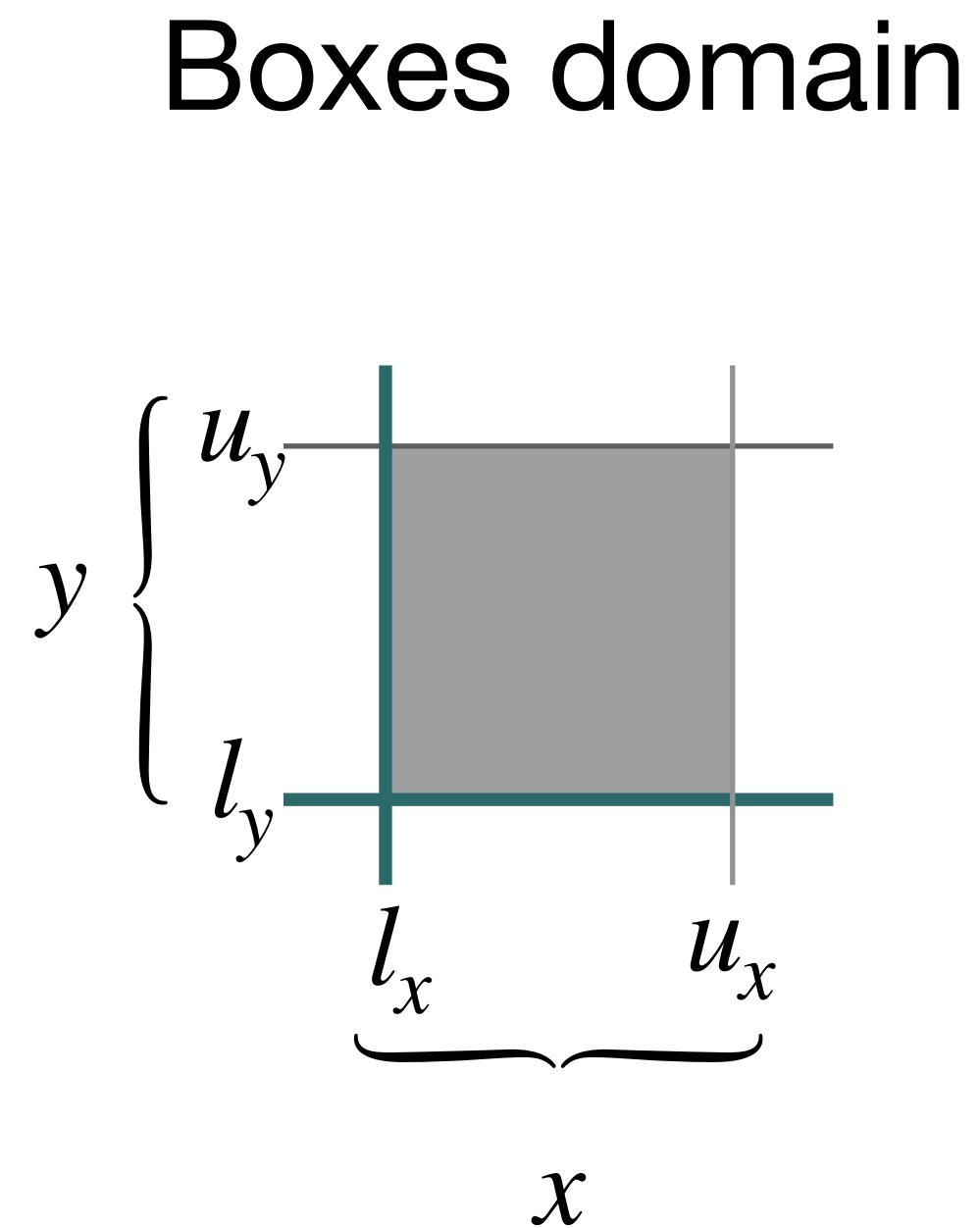
# Forward Analysis



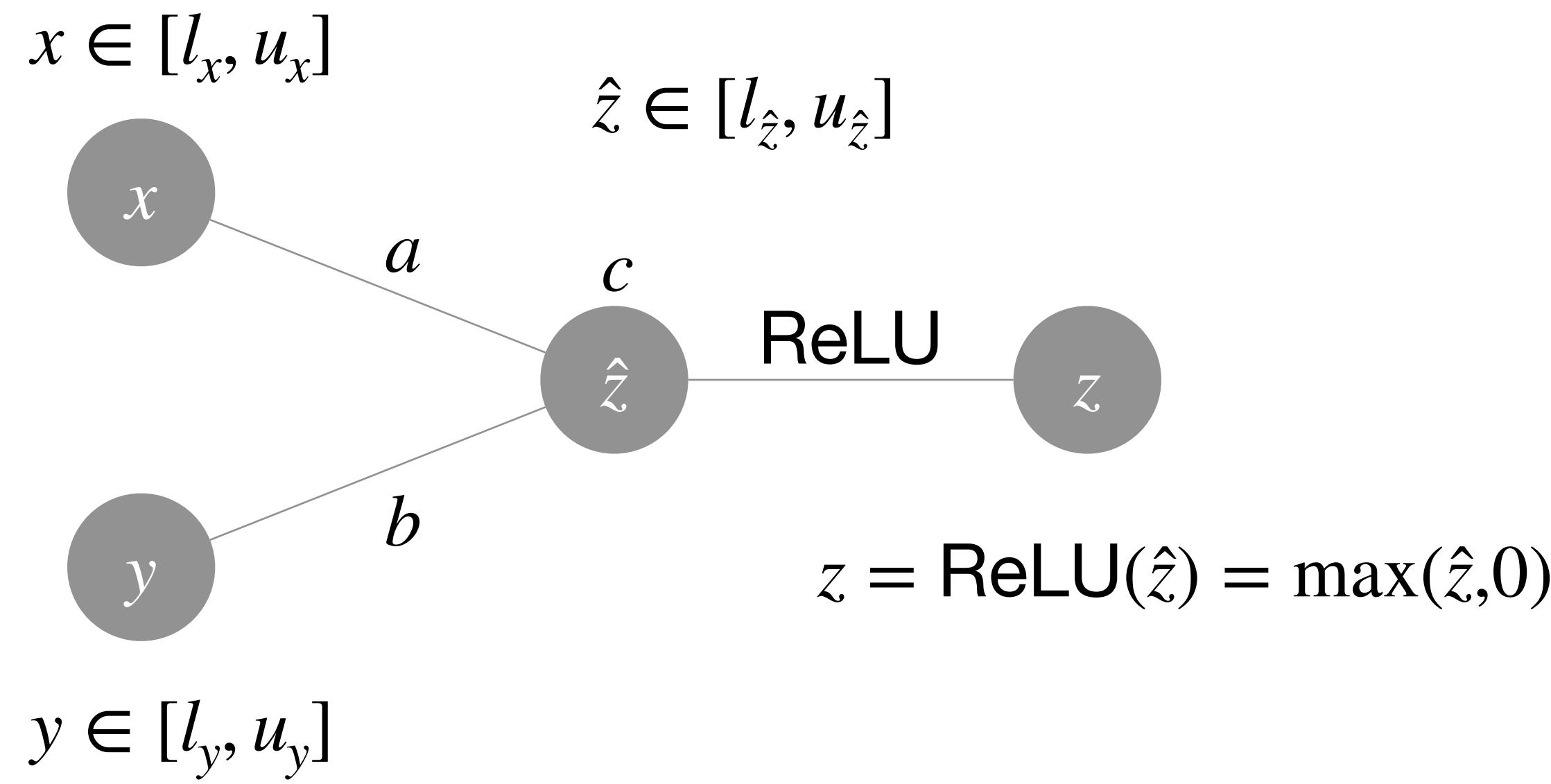
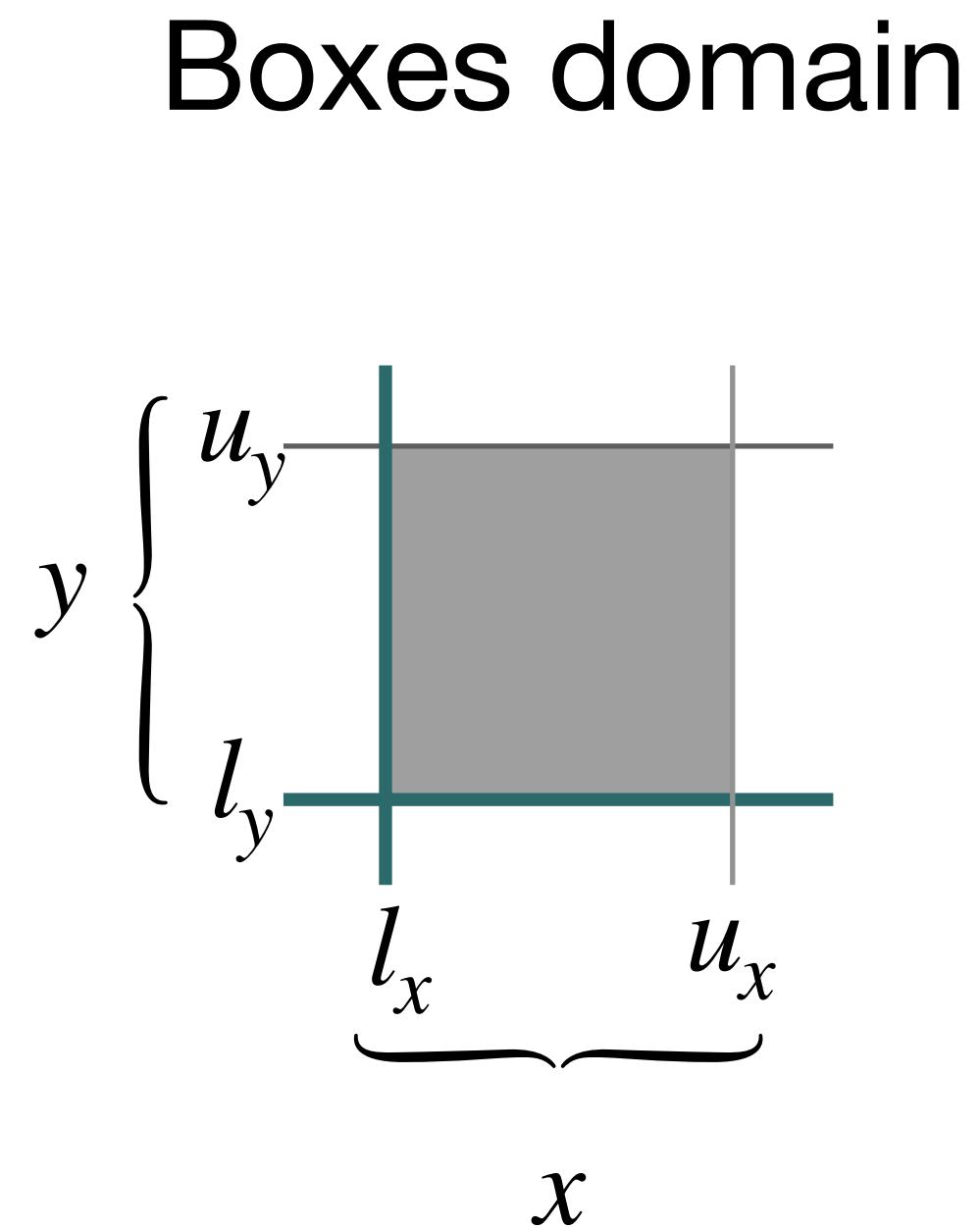
# Forward Analysis



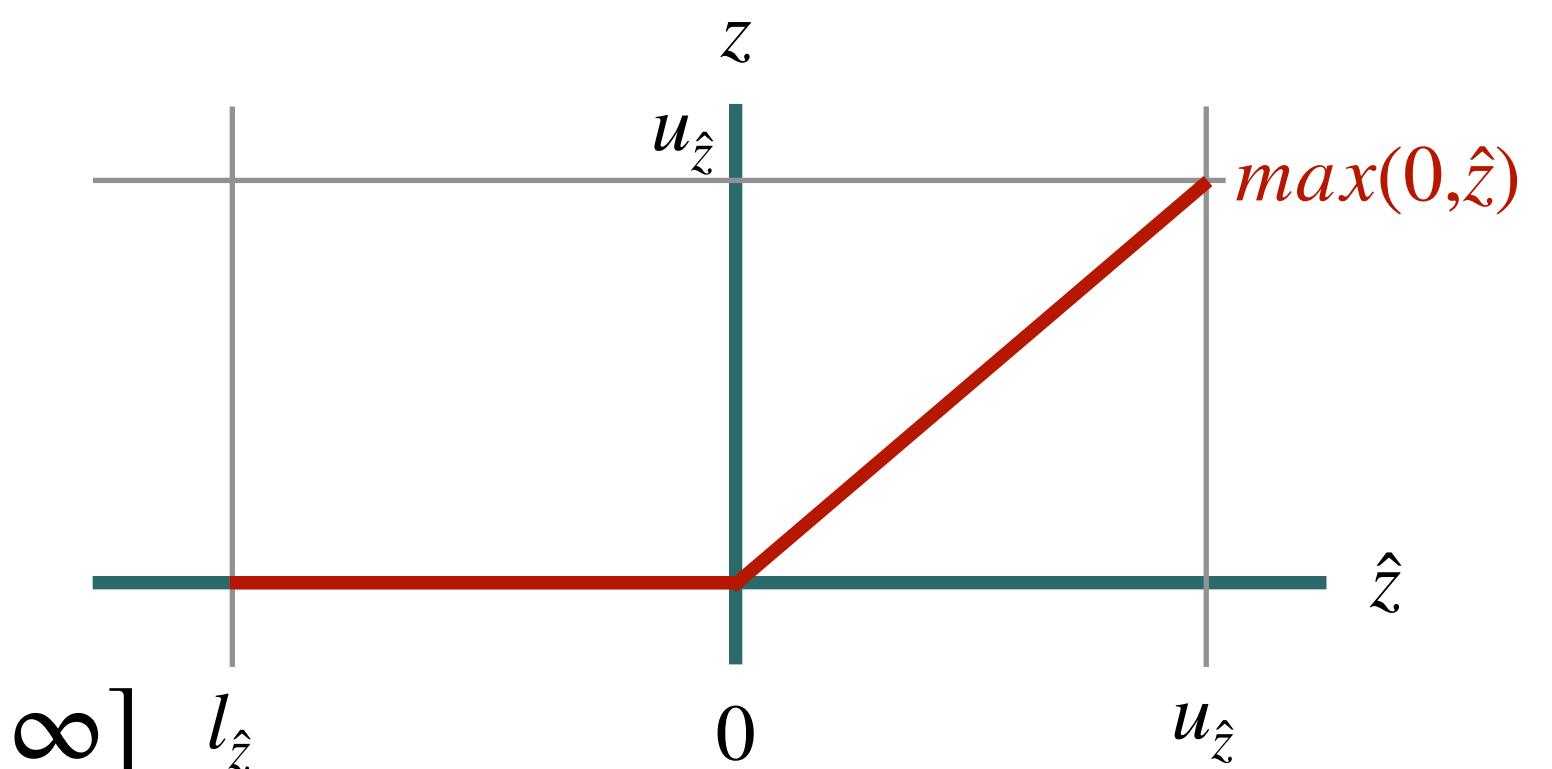
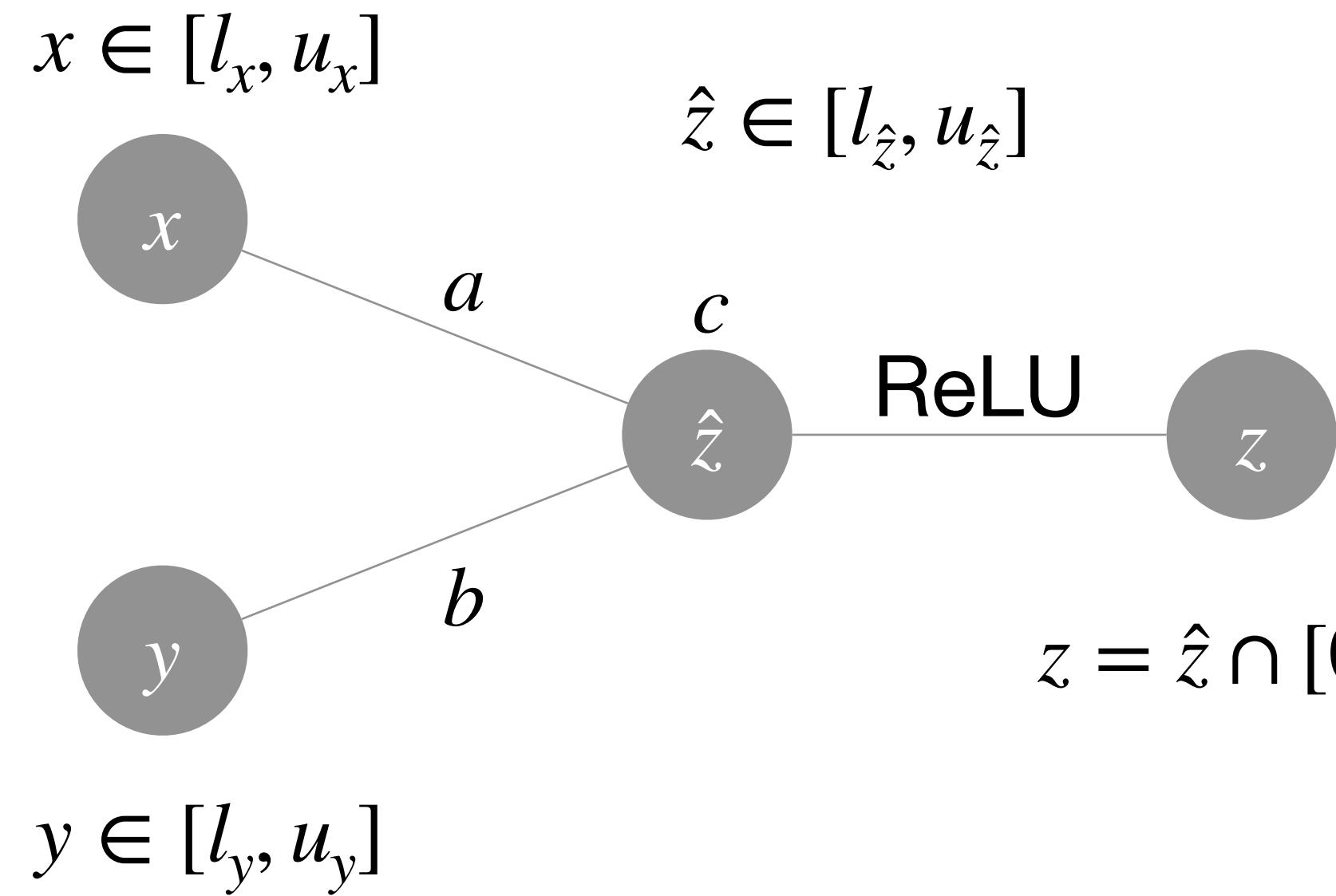
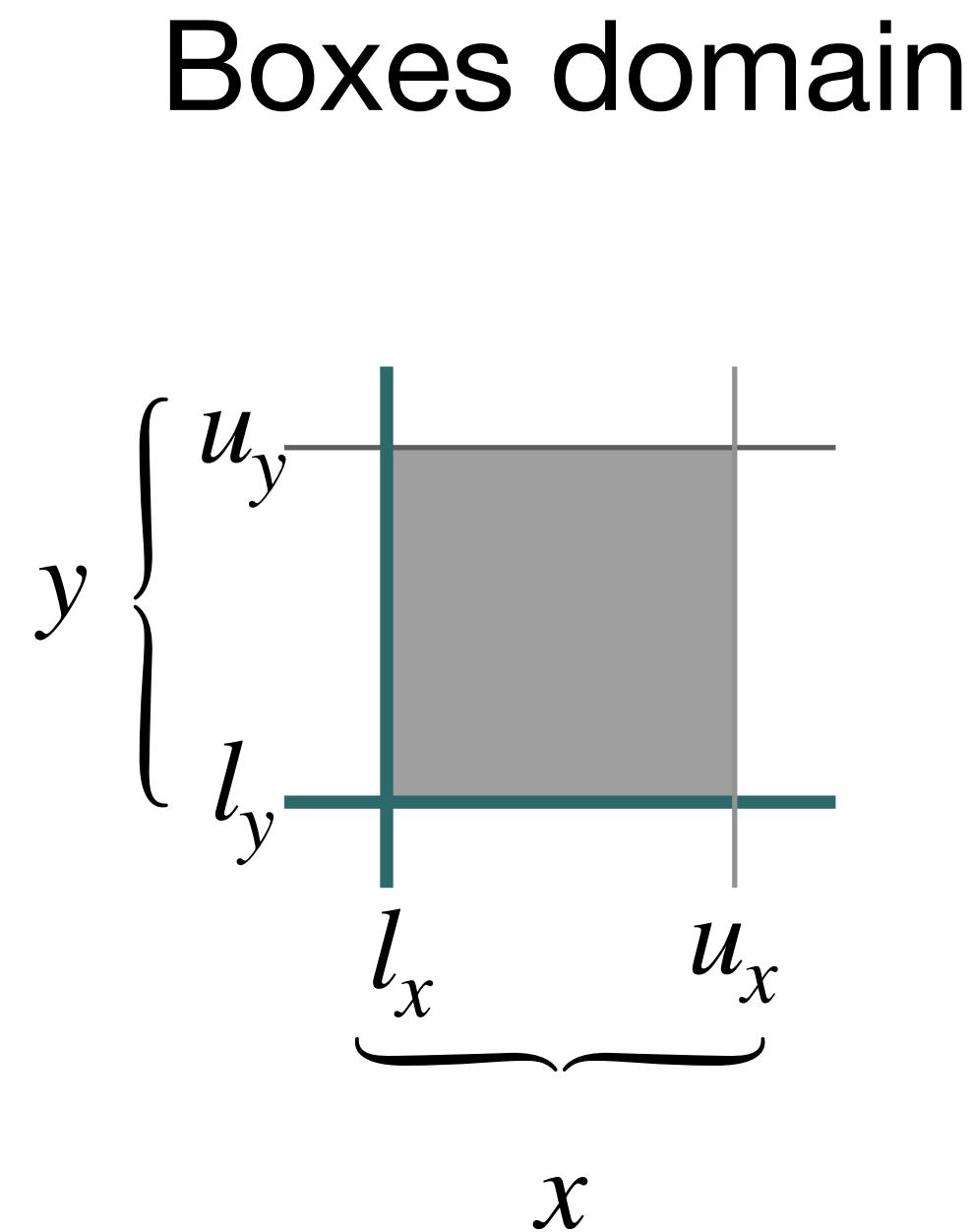
# Forward Analysis



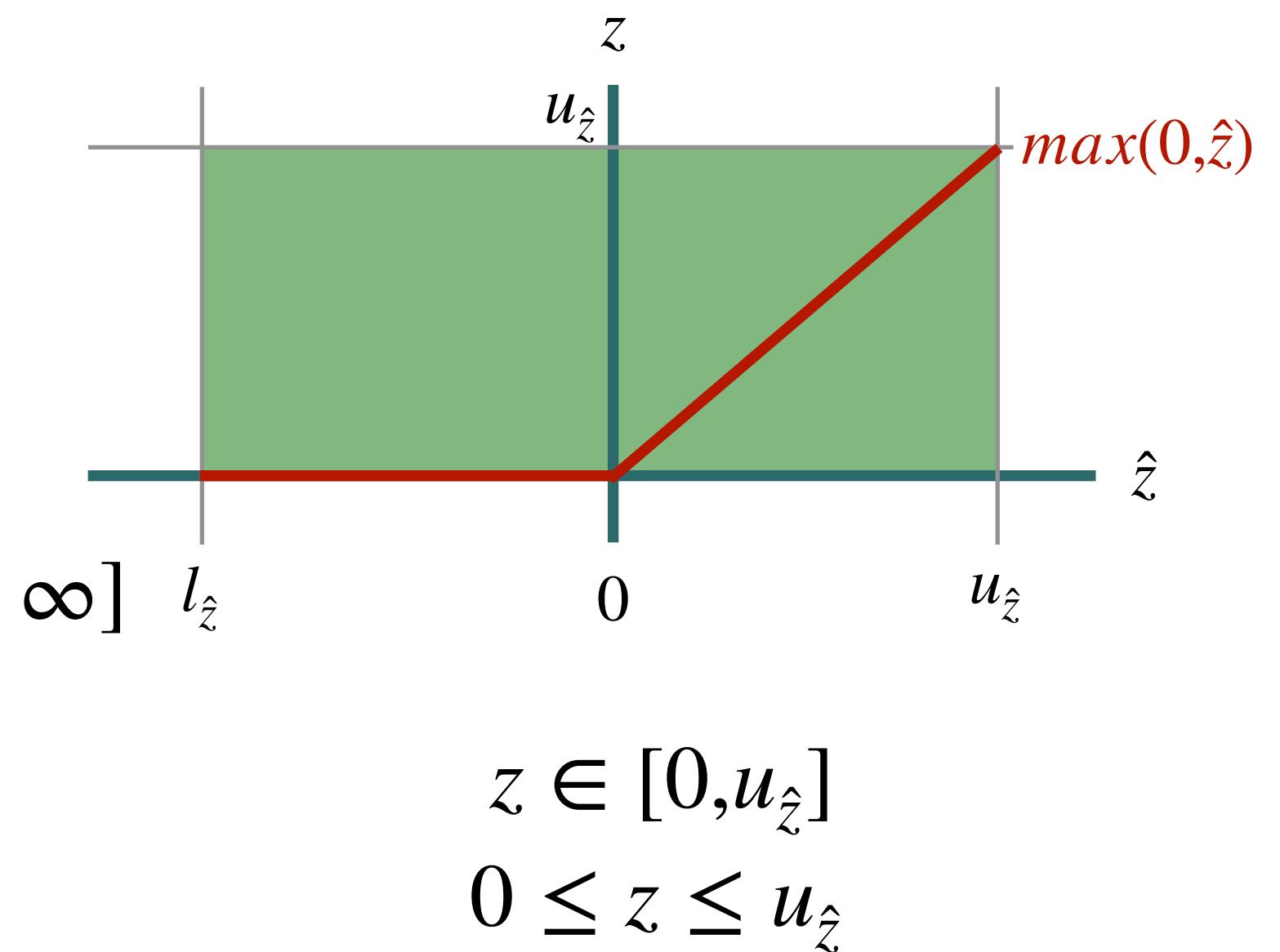
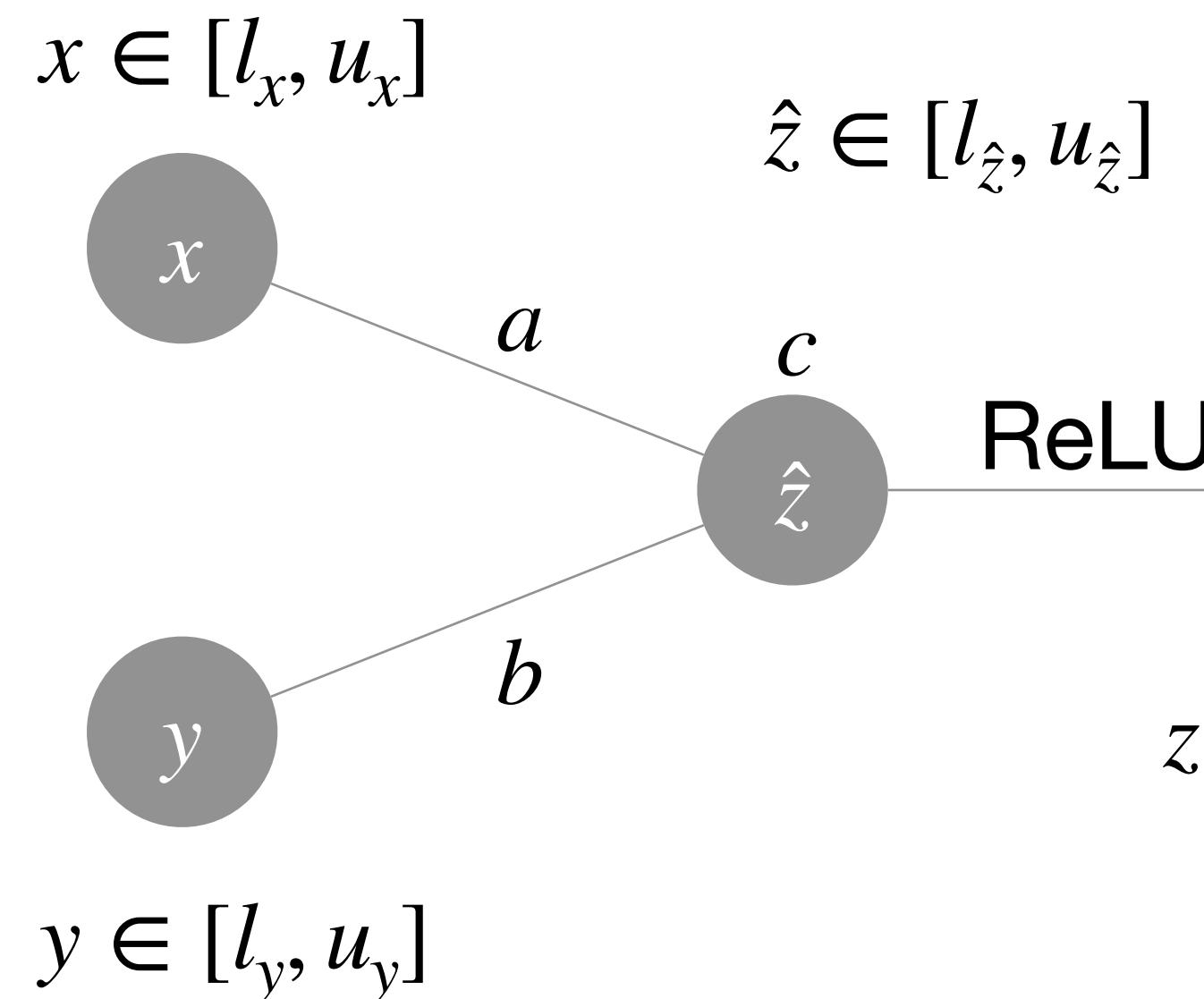
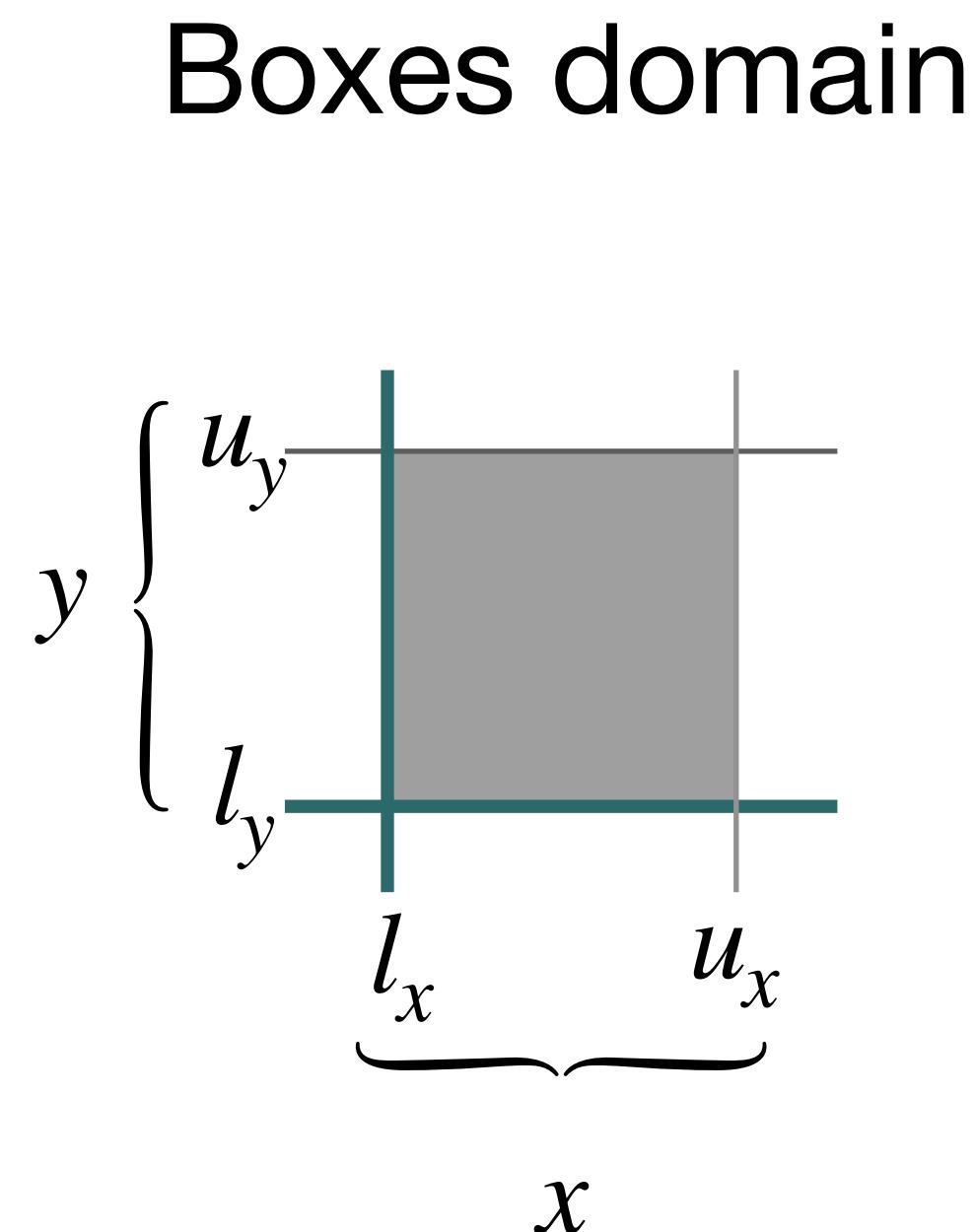
# Forward Analysis



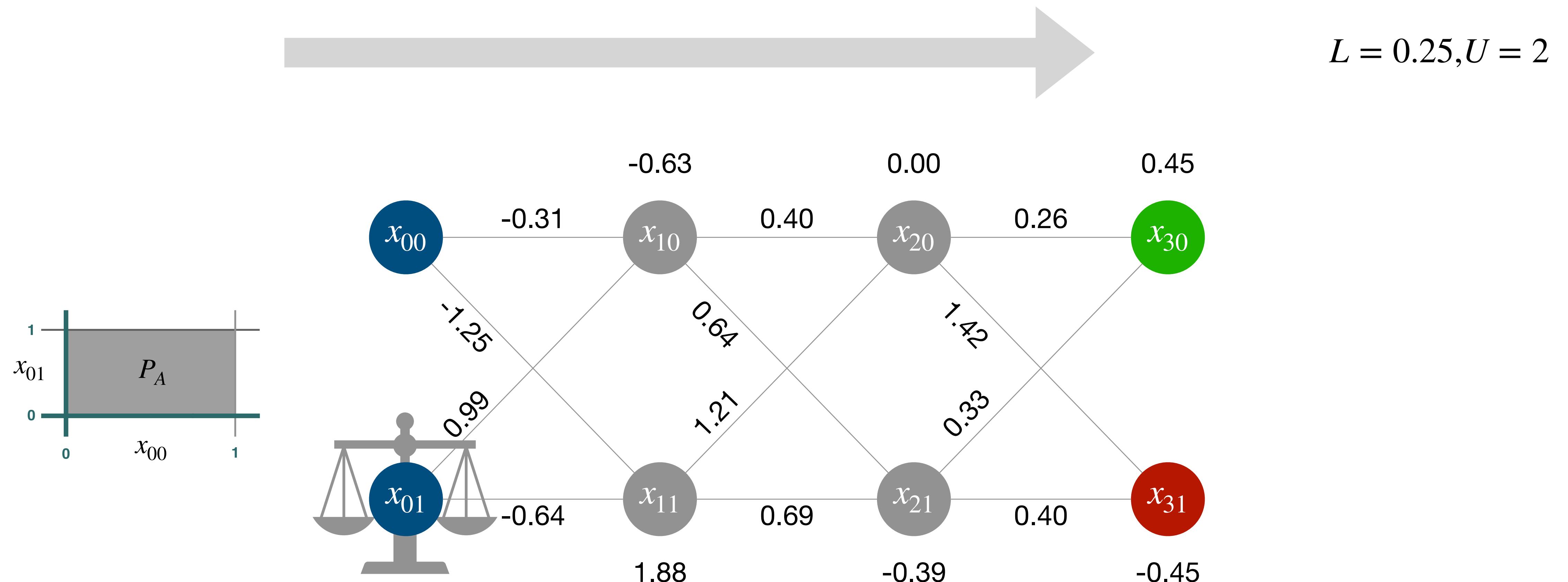
# Forward Analysis



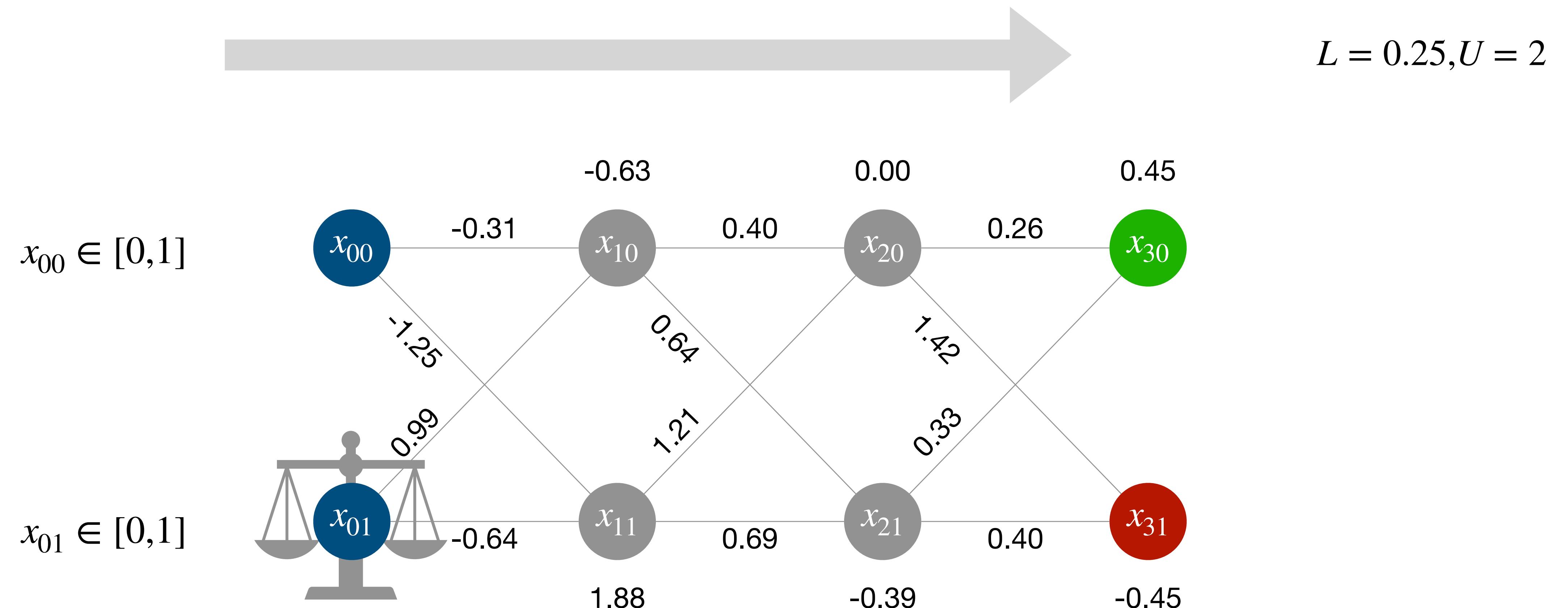
# Forward Analysis



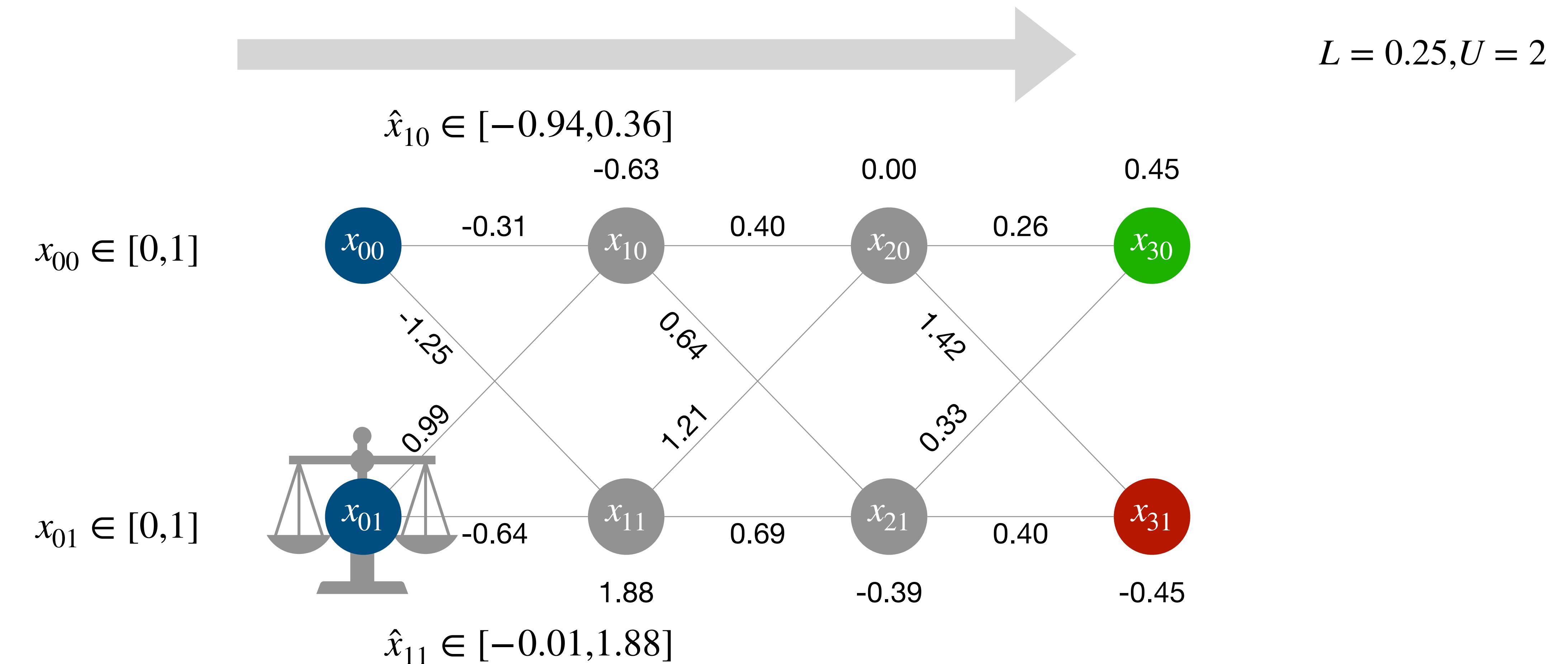
# Forward Analysis



# Forward Analysis

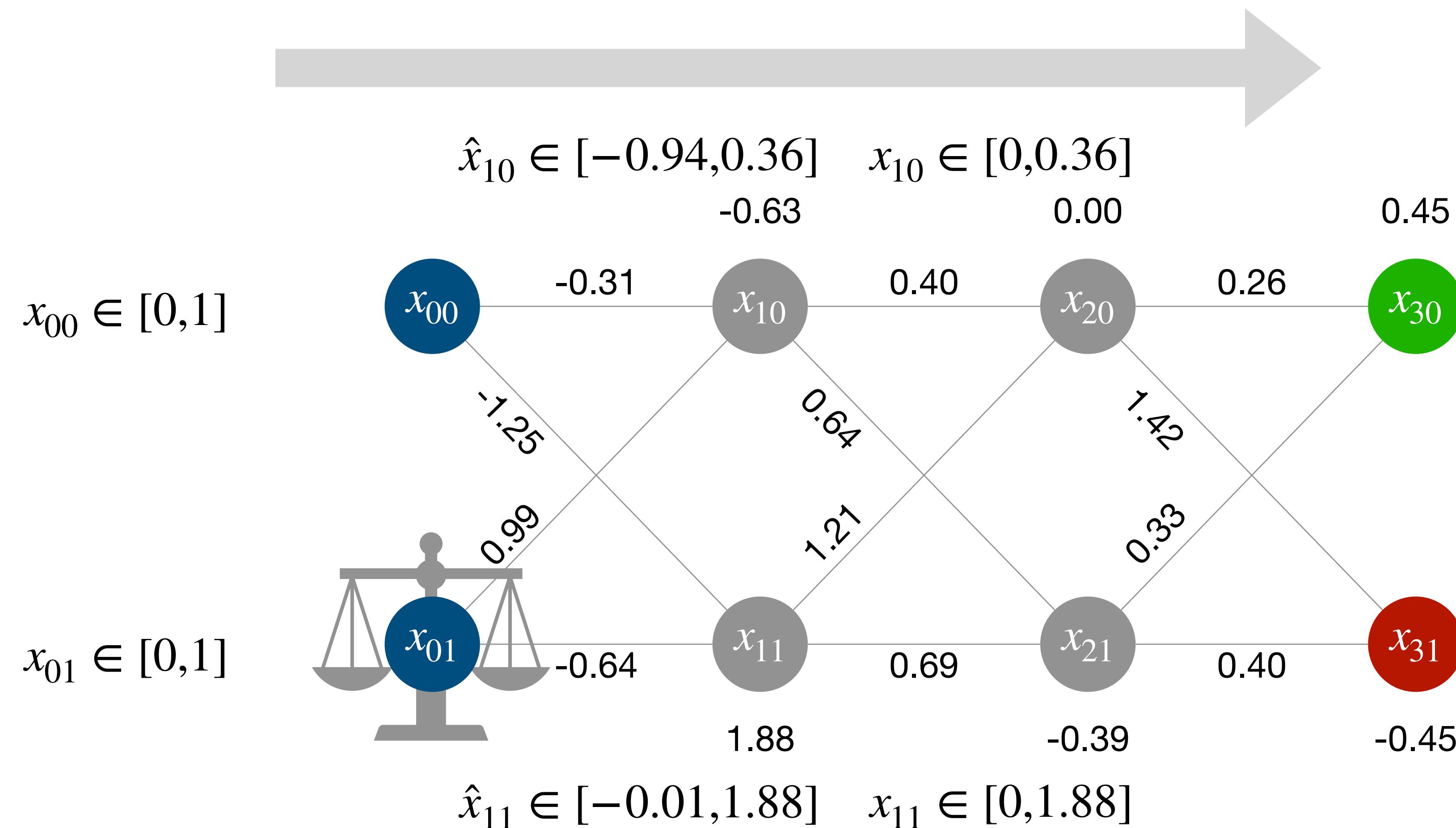


# Forward Analysis

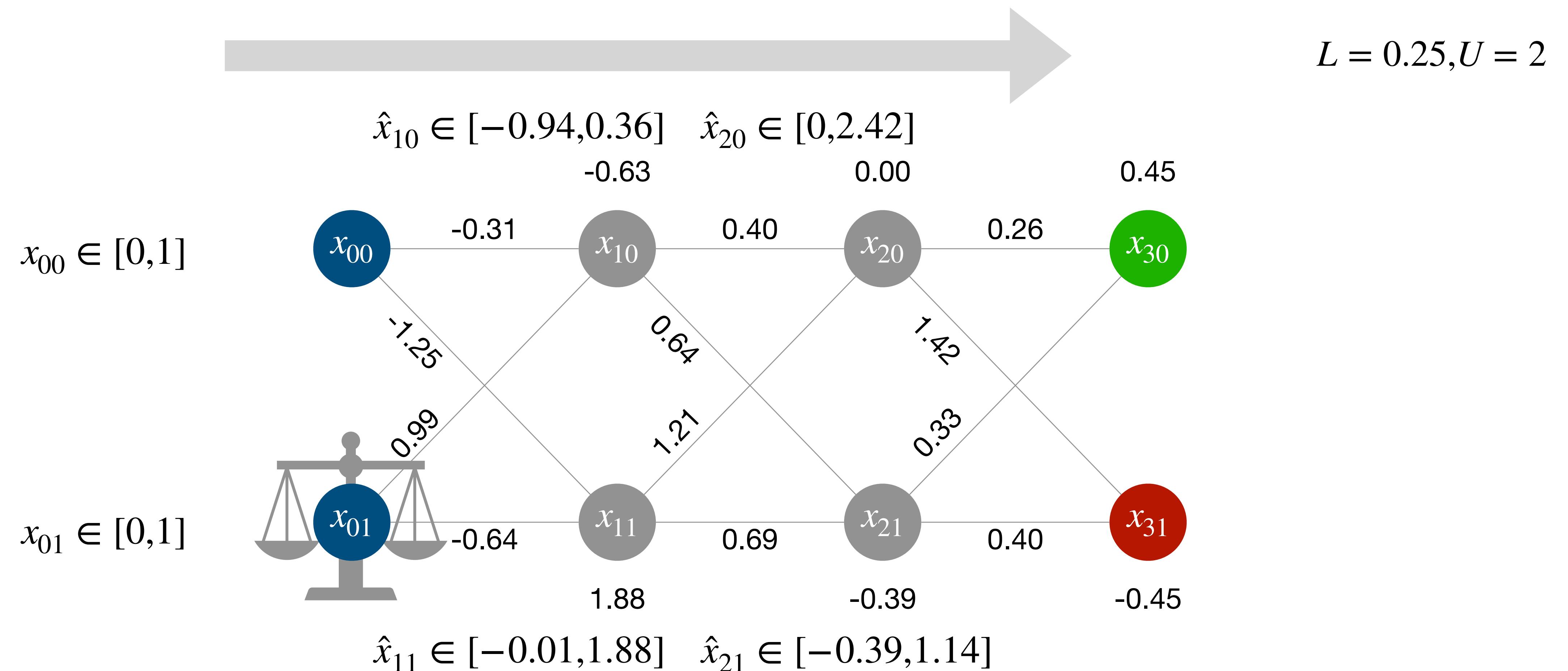


# Forward Analysis

$L = 0.25, U = 2$

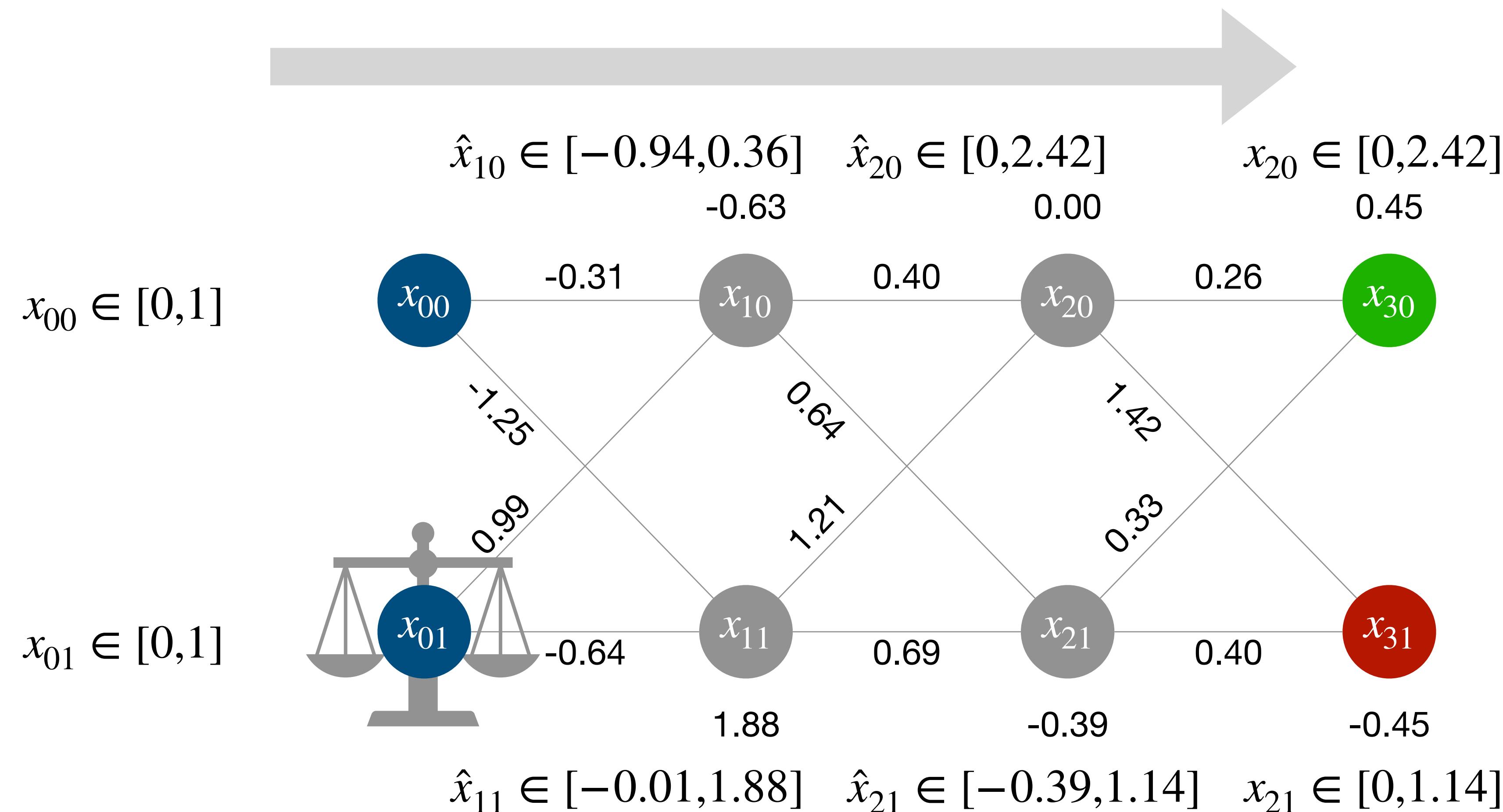


# Forward Analysis

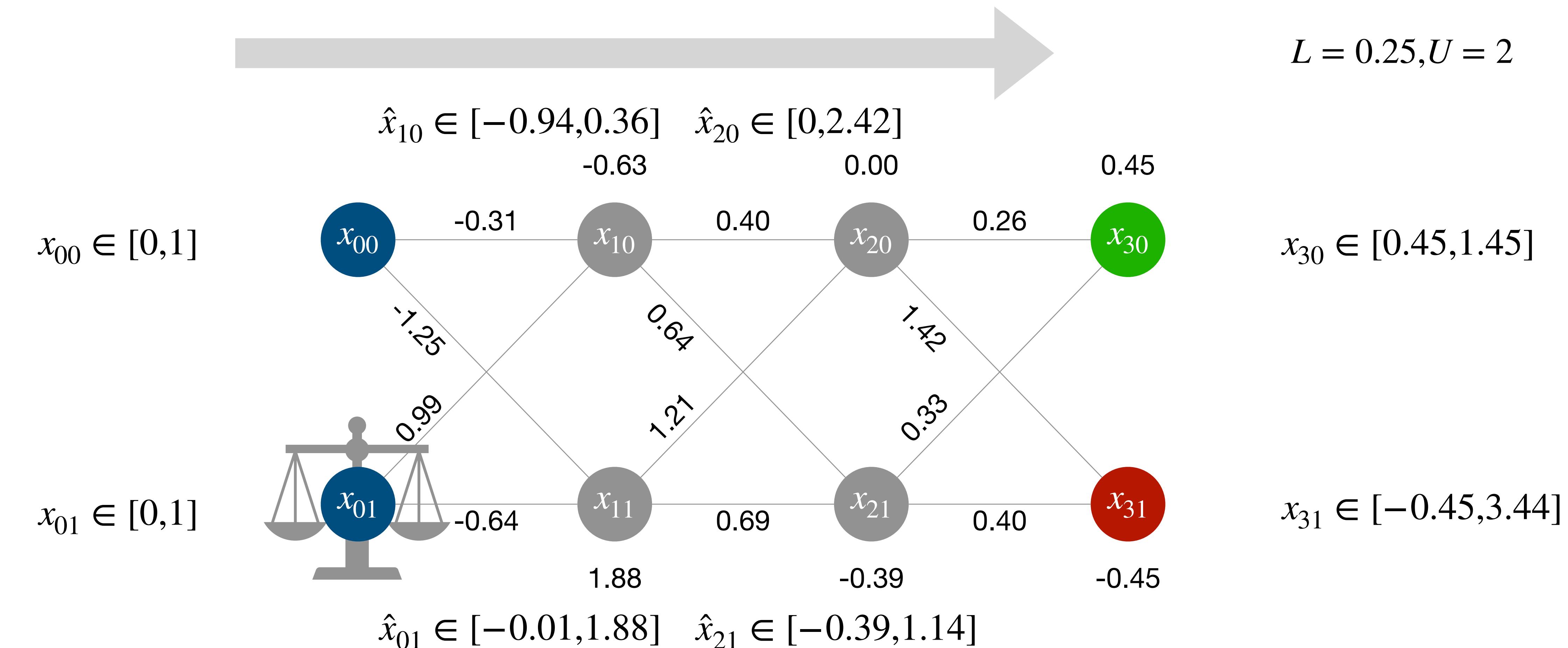


# Forward Analysis

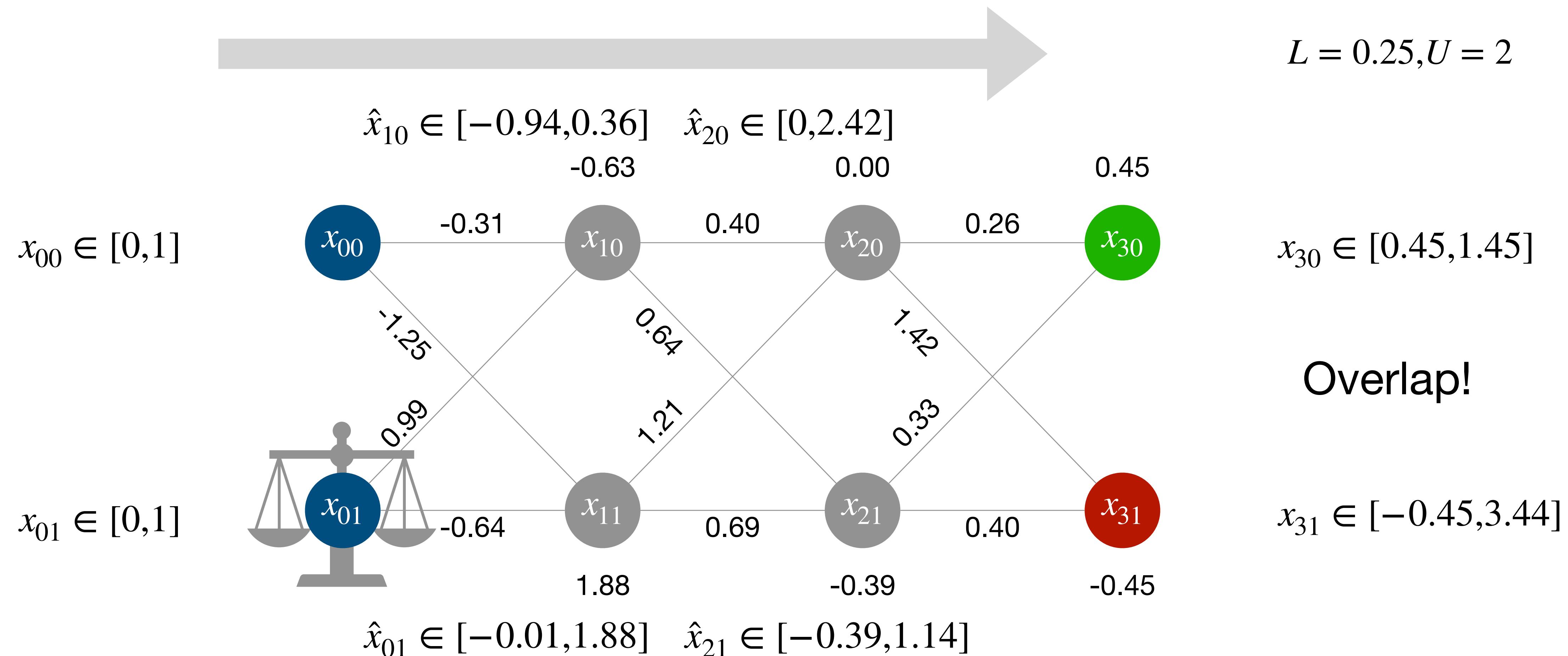
$$L = 0.25, U = 2$$



# Forward Analysis

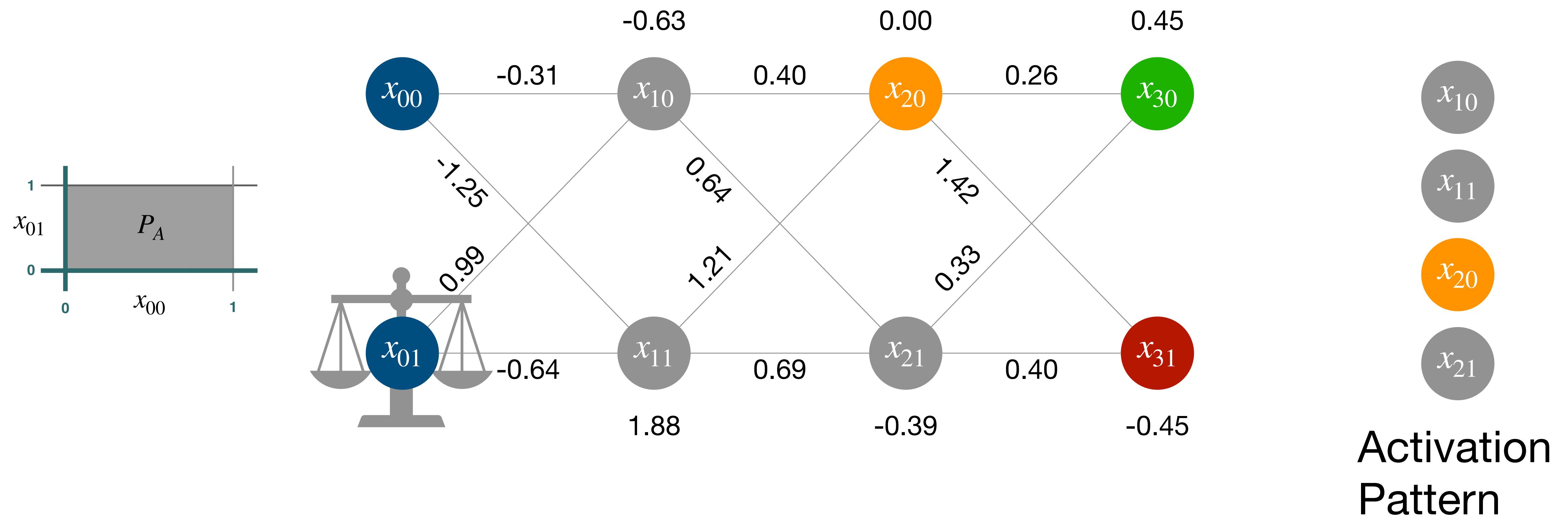


# Forward Analysis

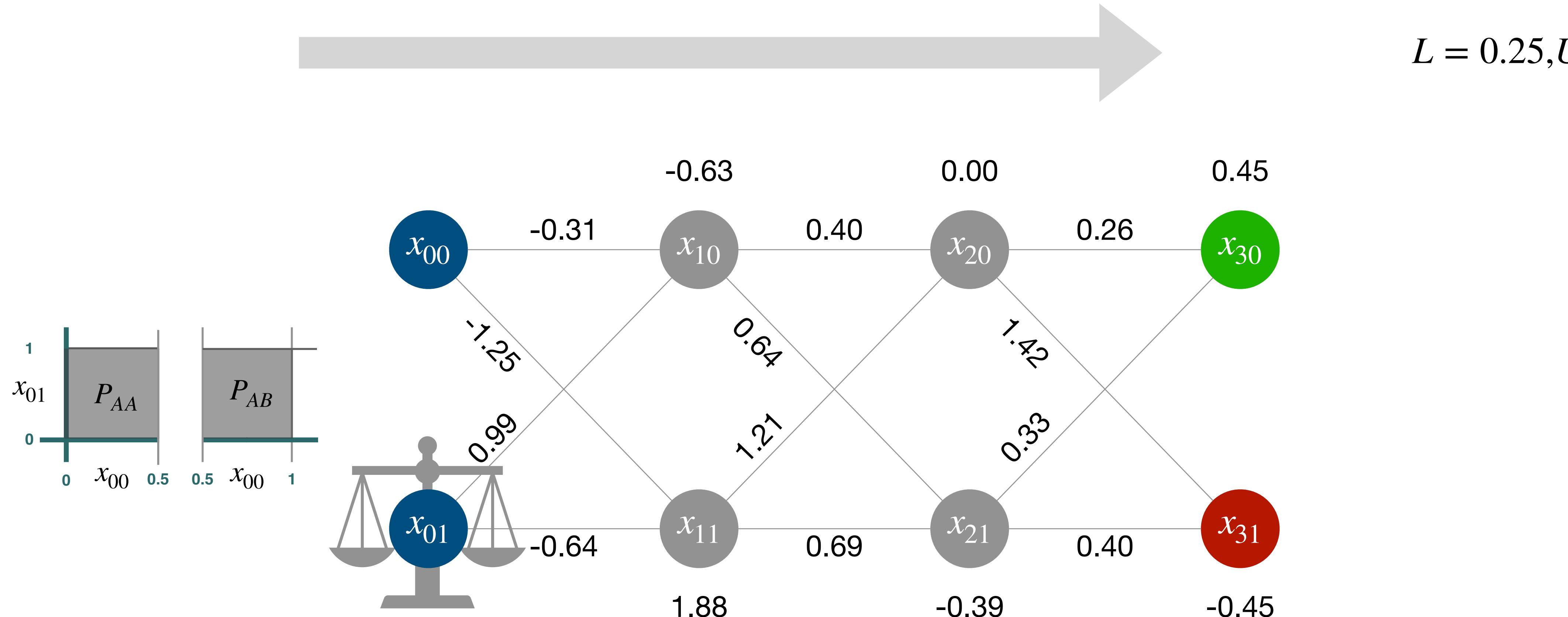


# Forward Analysis

$$L = 0.25, U = 2$$



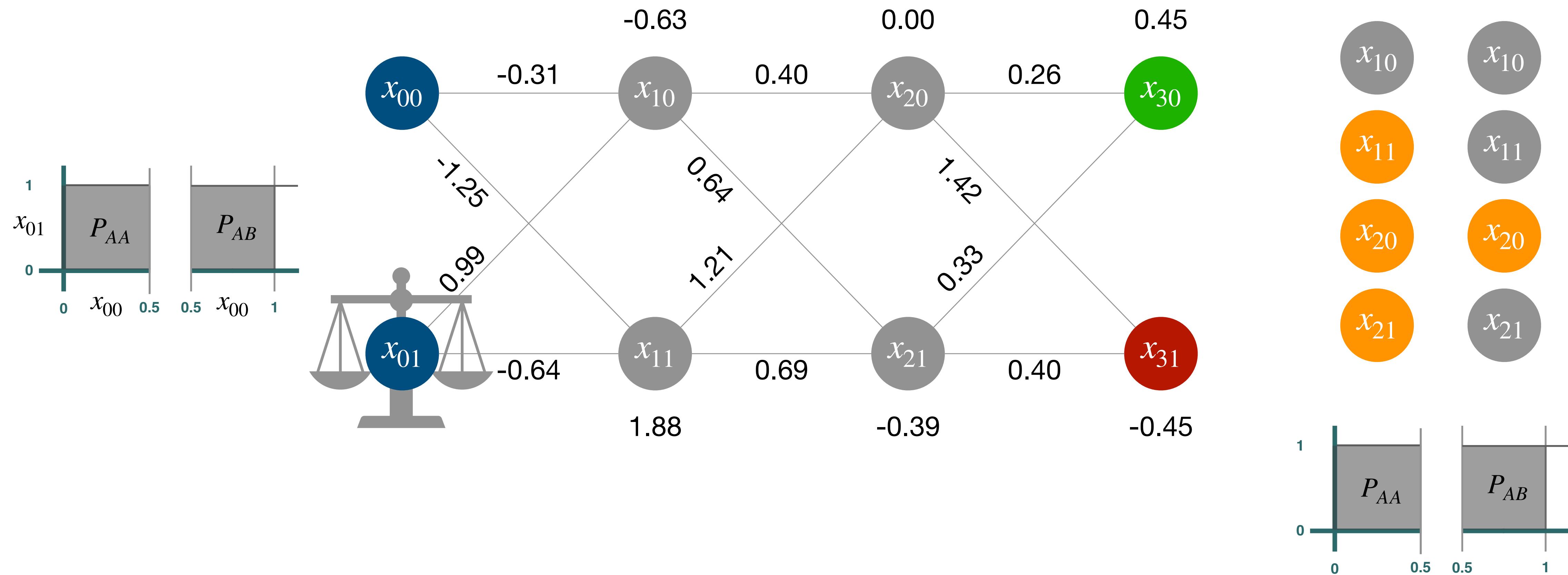
# Forward Analysis



# Forward Analysis



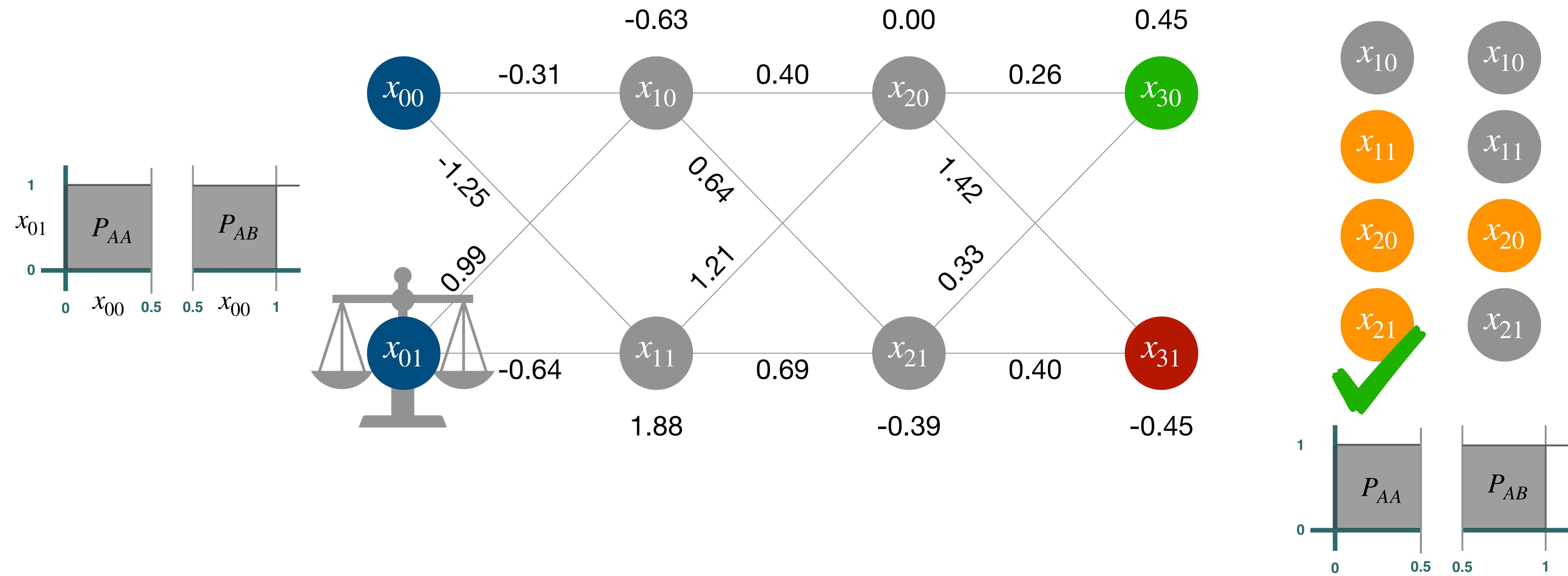
$$L = 0.25, U = 2$$



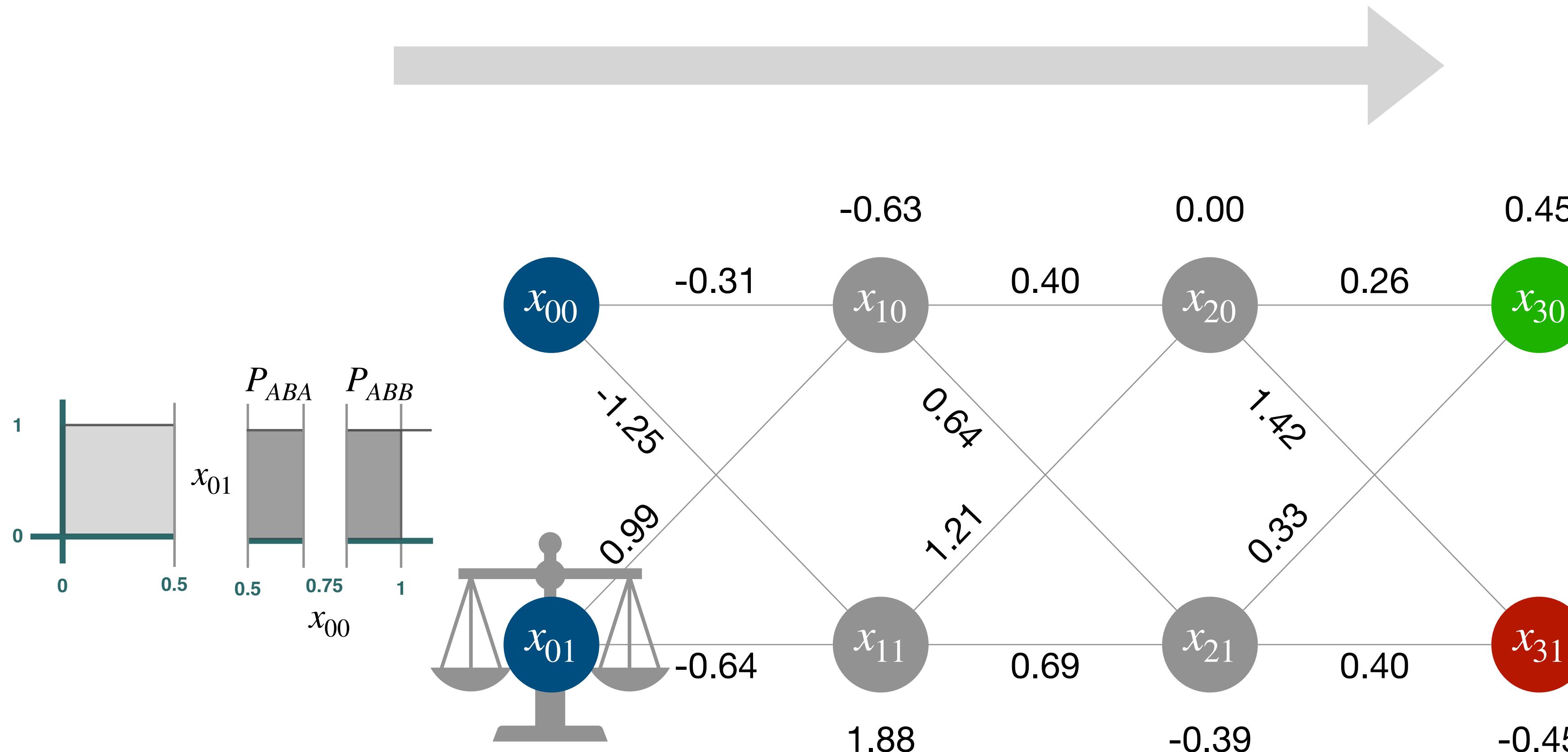
# Forward Analysis



$$L = 0.25, U = 2$$

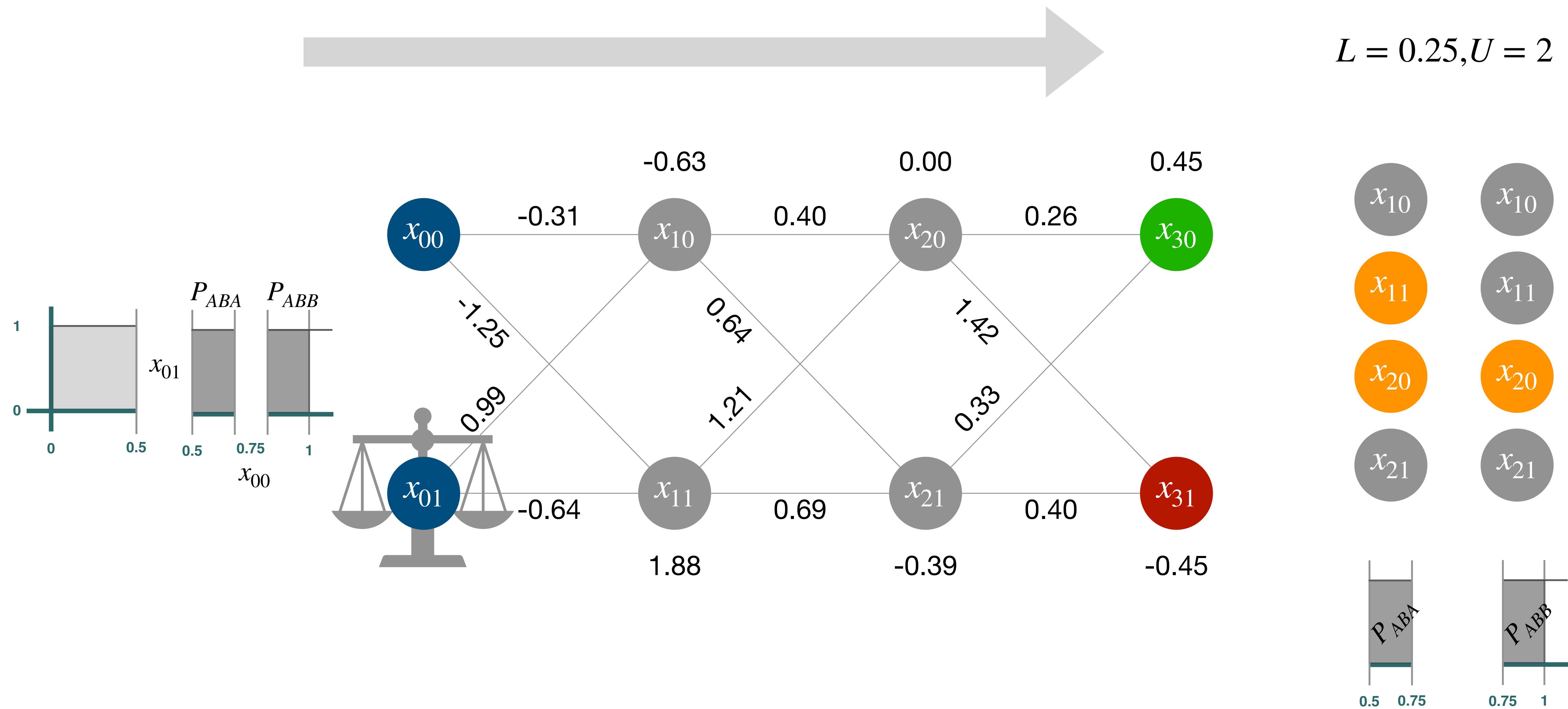


# Forward Analysis



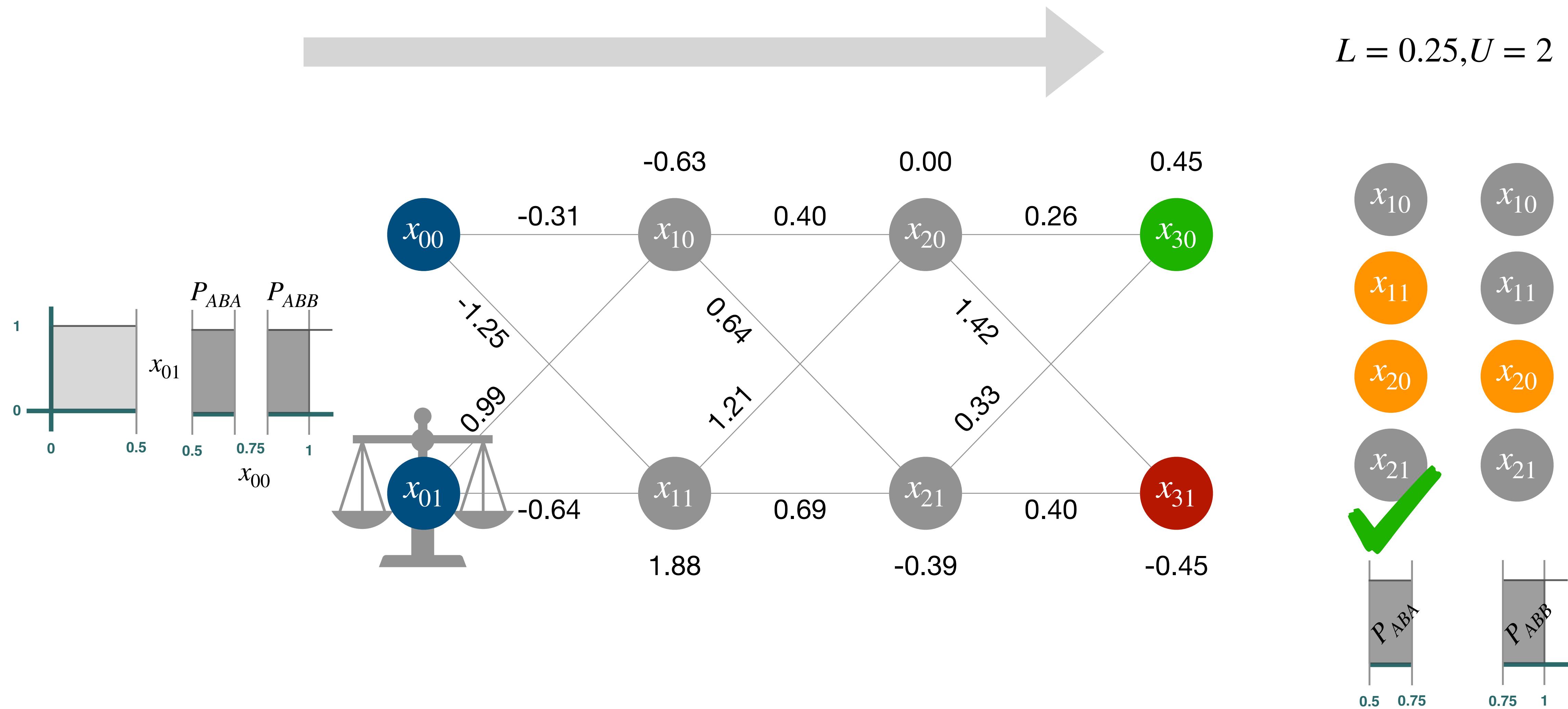
# Forward Analysis

$$L = 0.25, U = 2$$



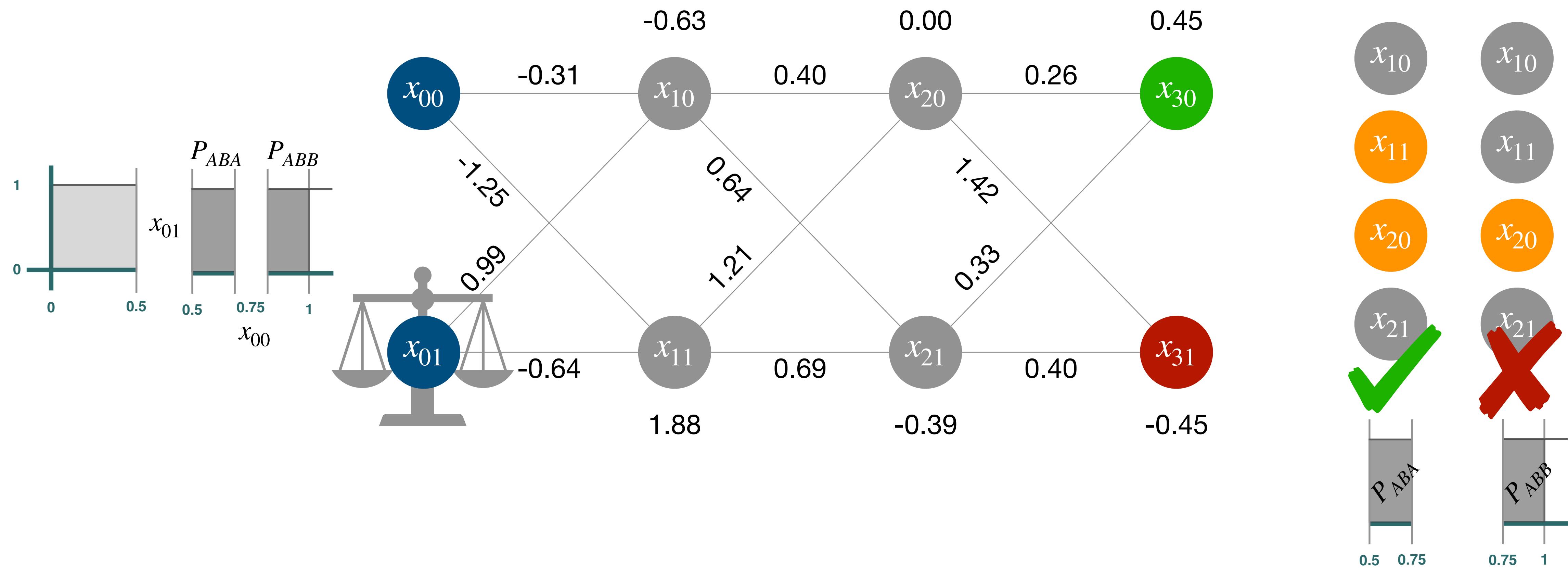
# Forward Analysis

$$L = 0.25, U = 2$$

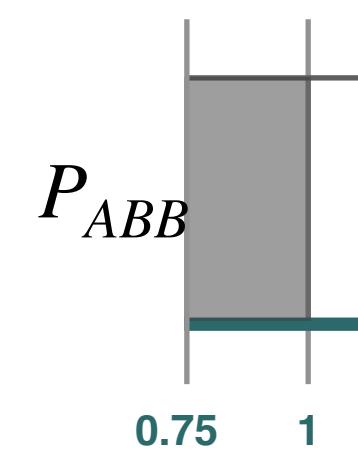
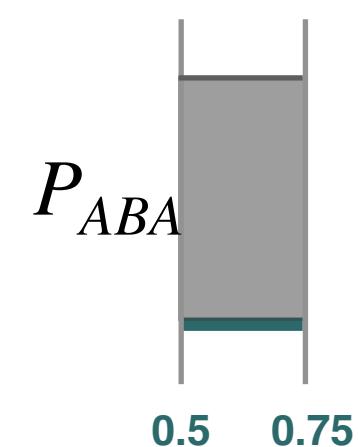
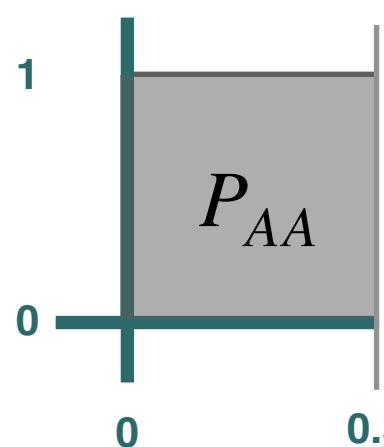
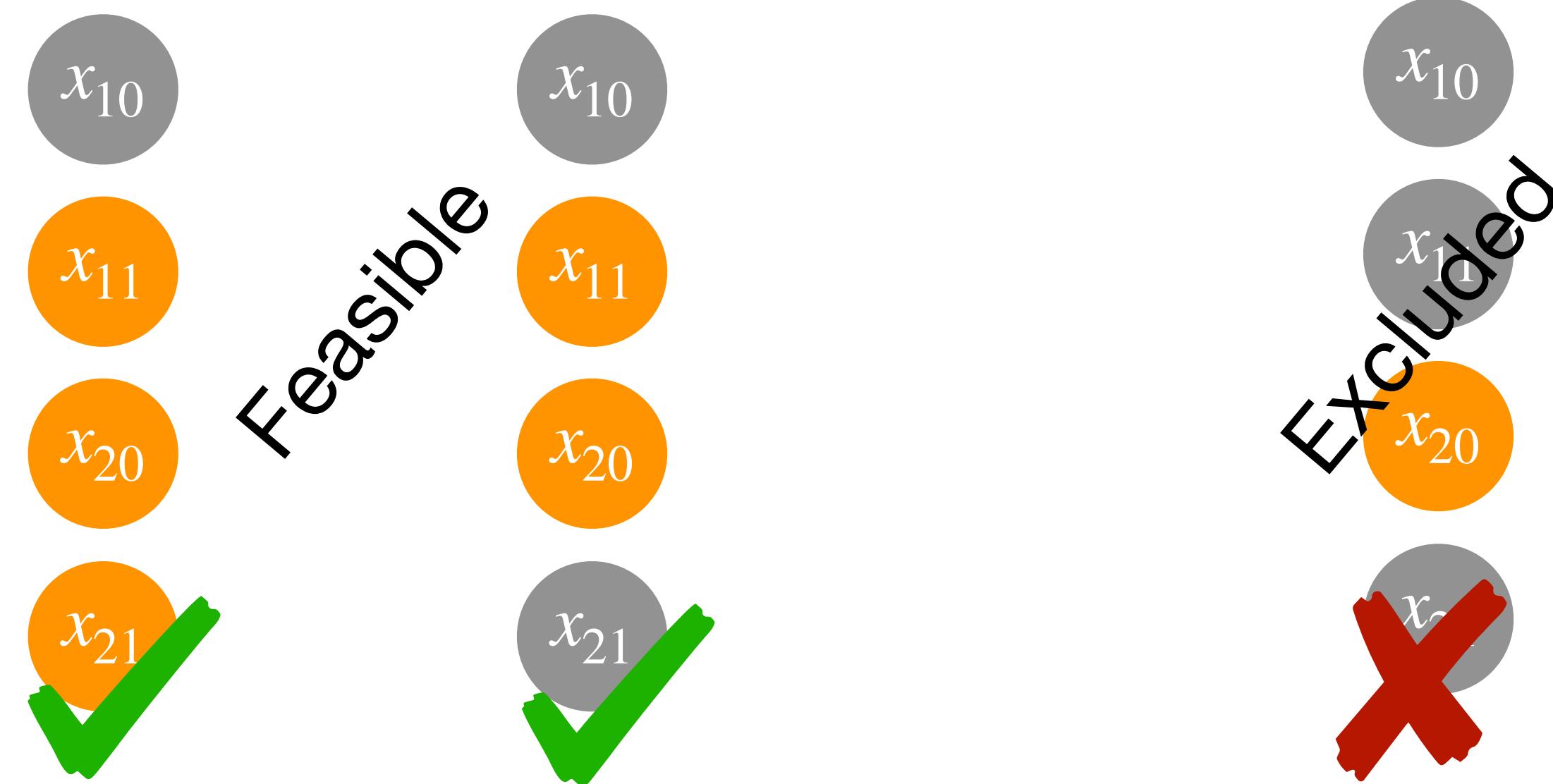


# Forward Analysis

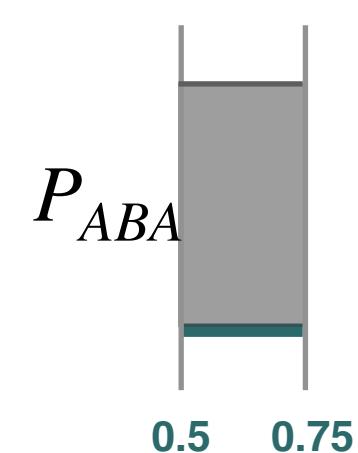
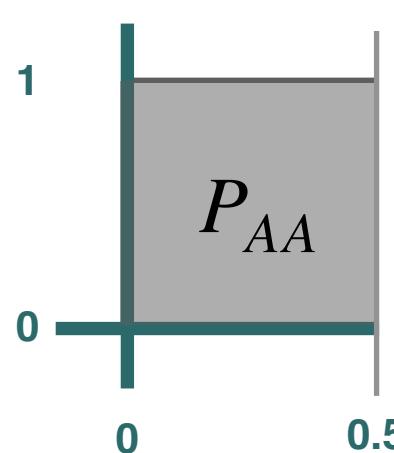
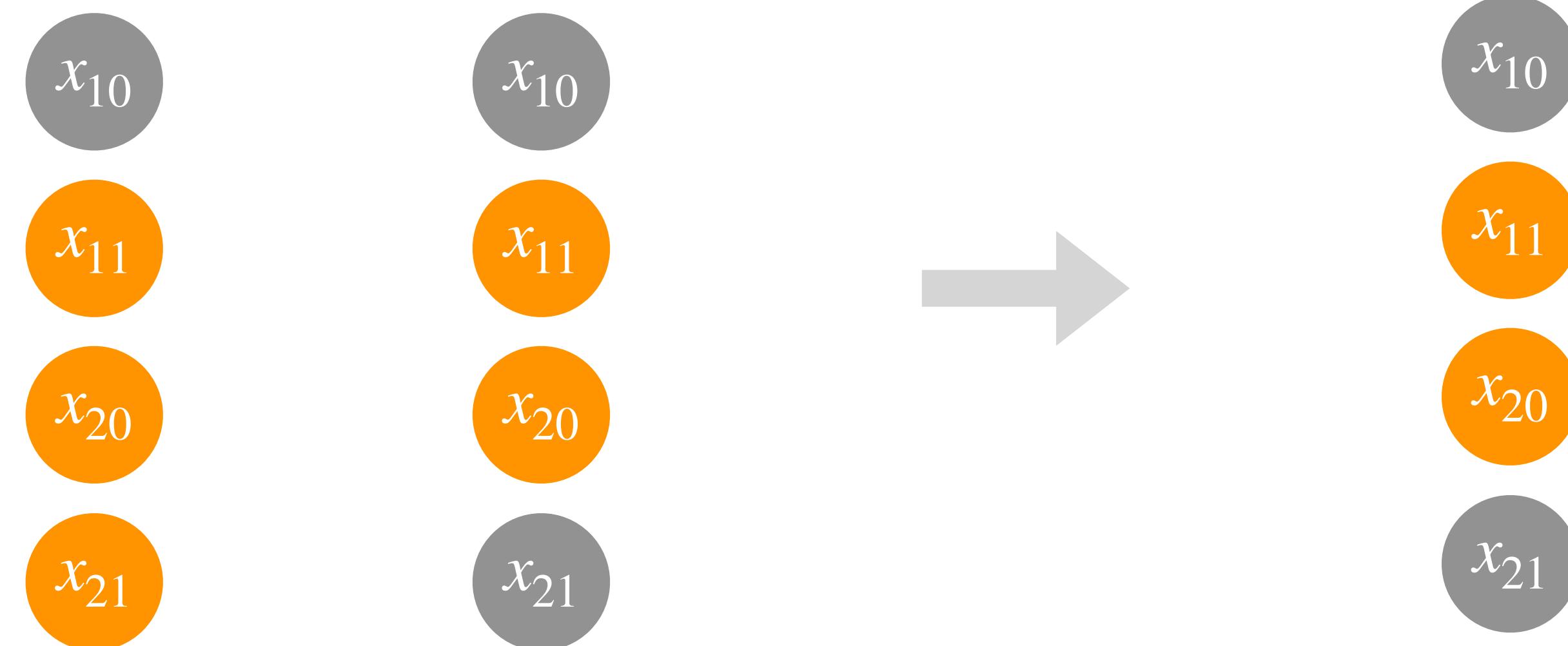
$$L = 0.25, U = 2$$



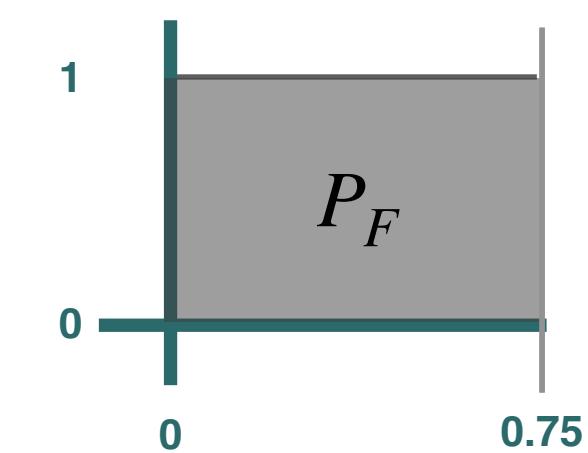
# Forward Analysis



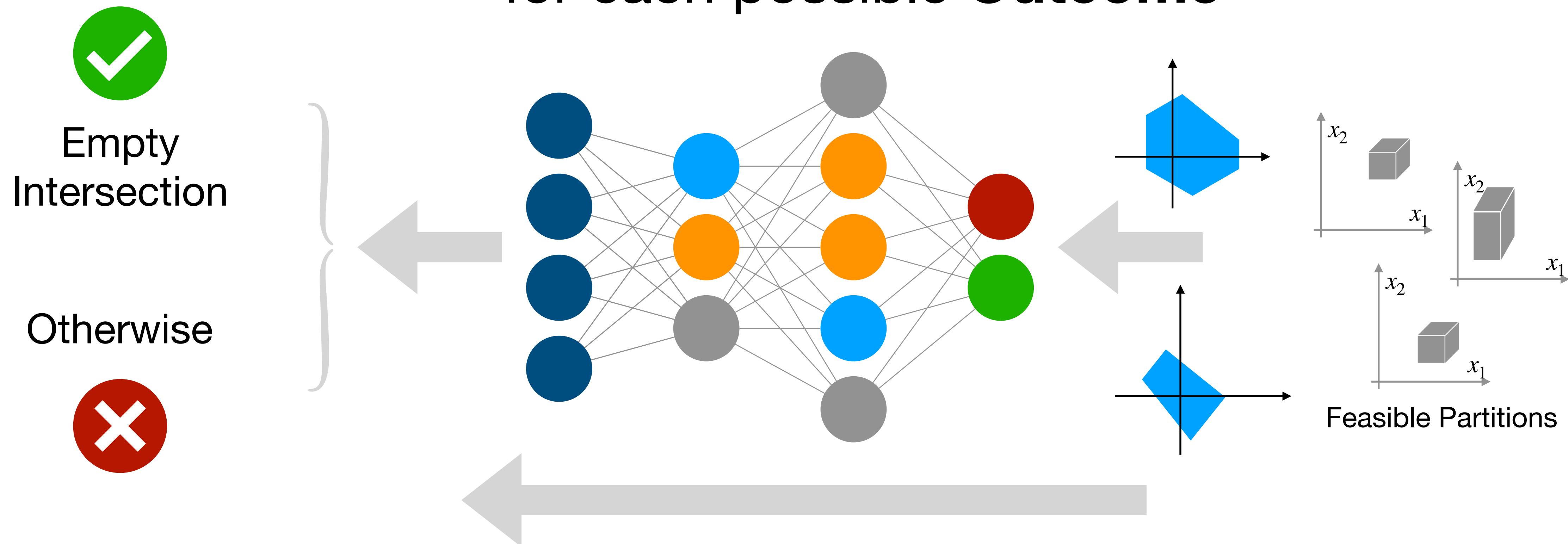
# Forward Analysis



$\Rightarrow$



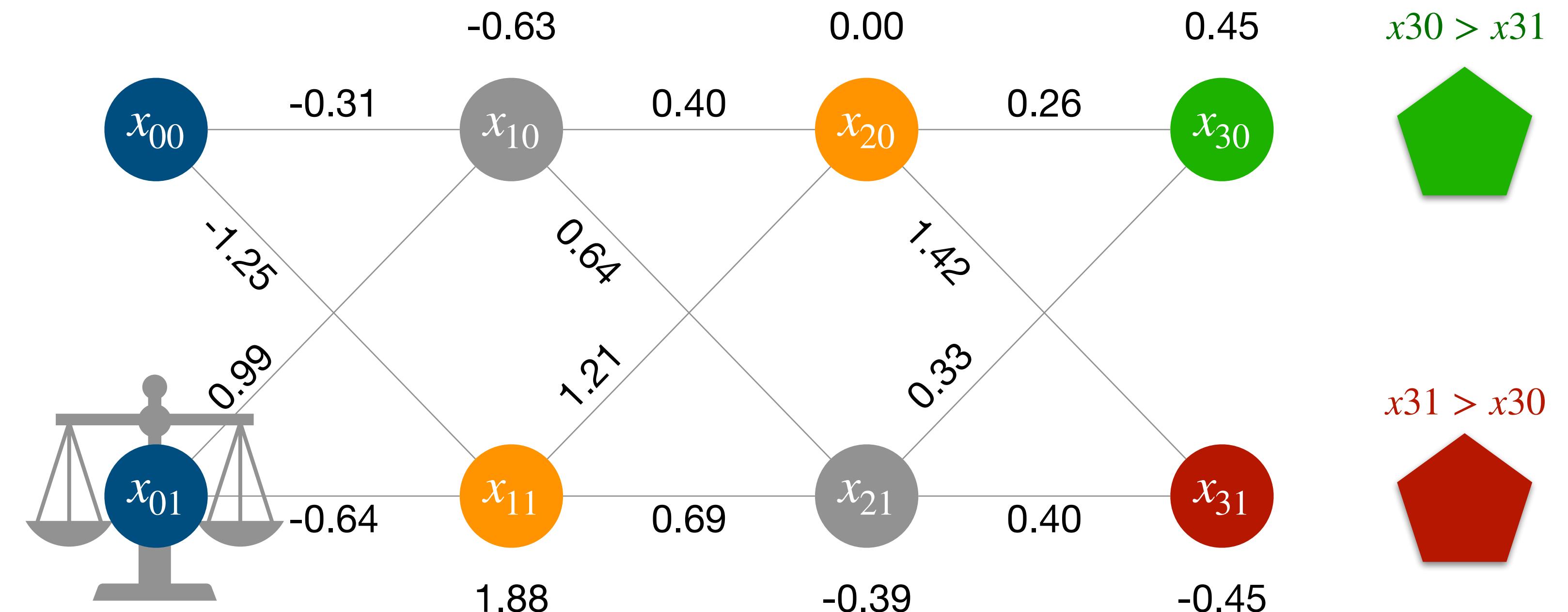
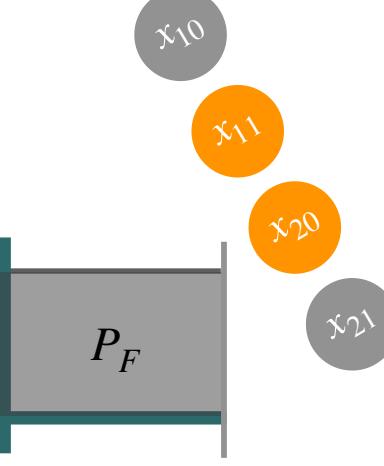
**Proceed Backwards**  
for each **Feasible** partitions  
for each possible **Outcome**



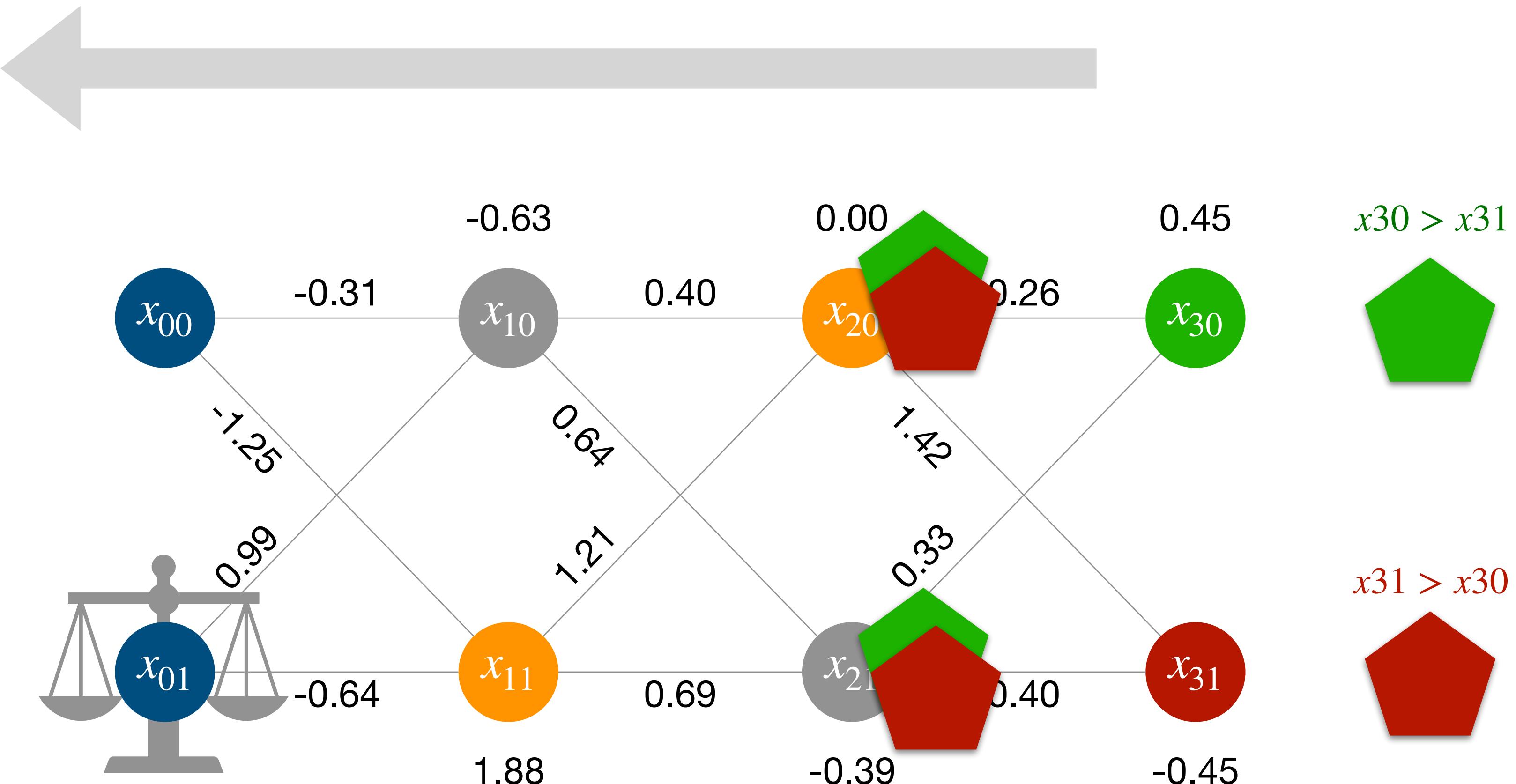
**Exact Backward Analysis**  
using **Polyhedra**

# Backward Analysis

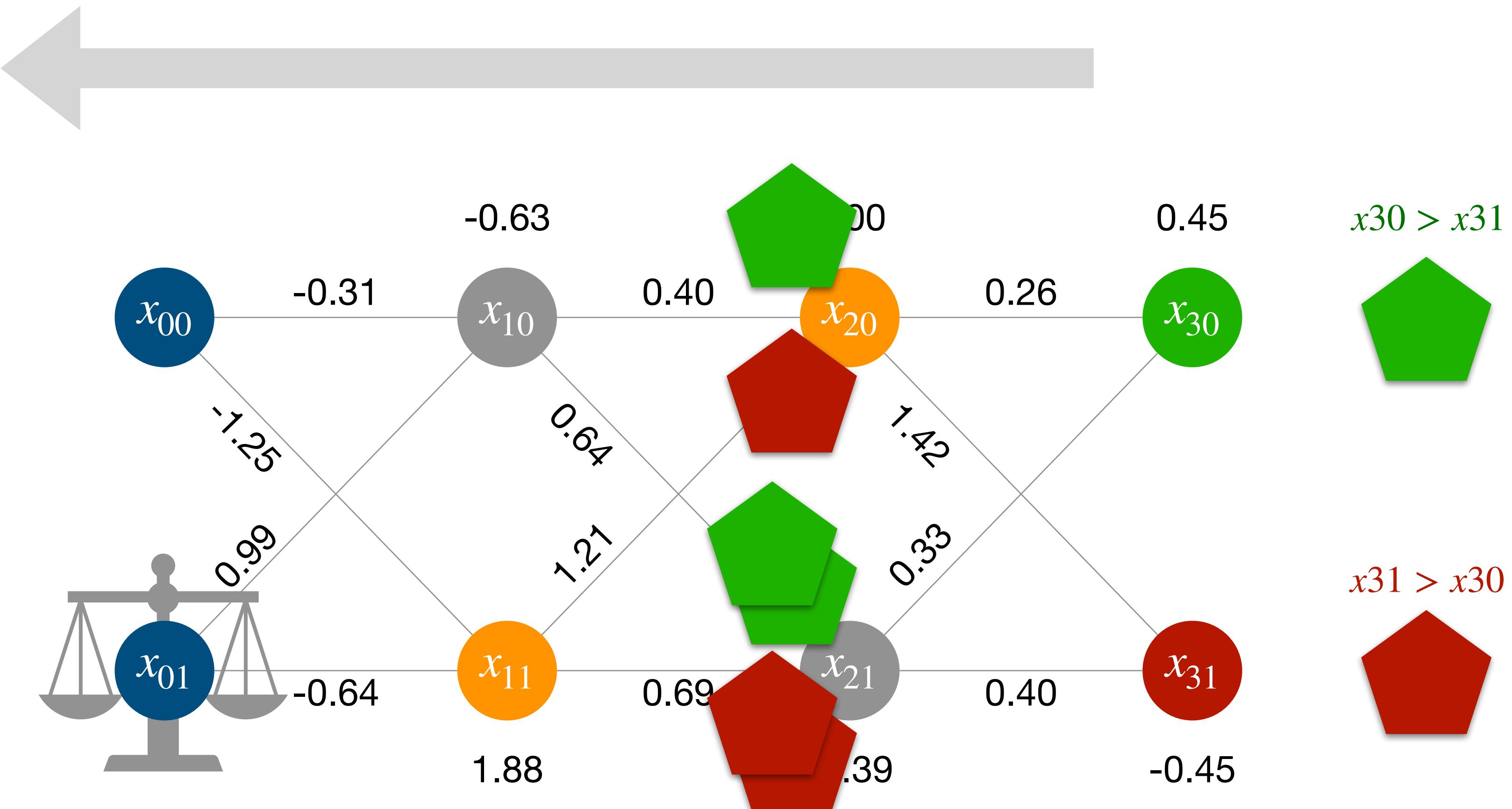
From the forward analysis



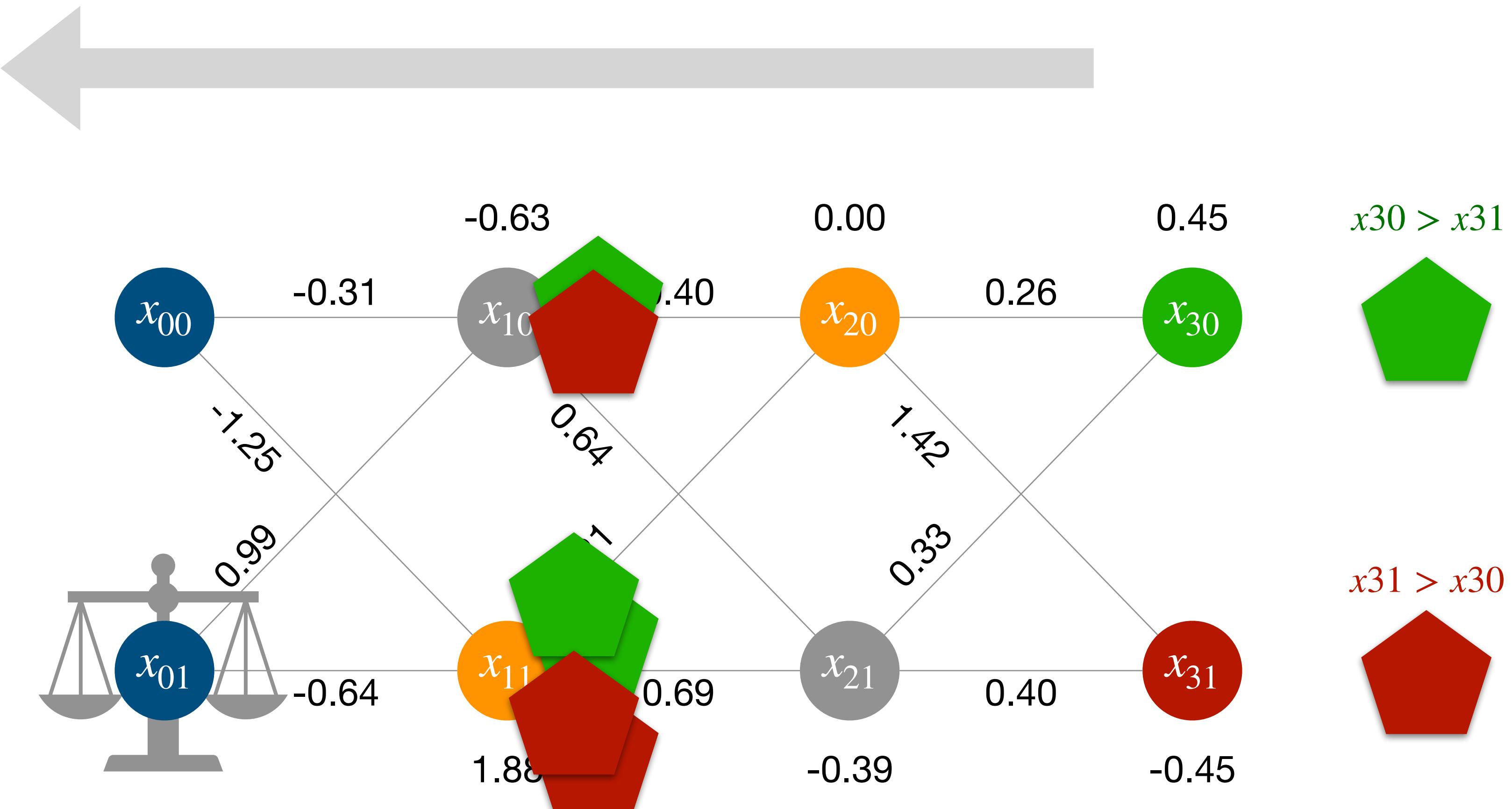
# Backward Analysis



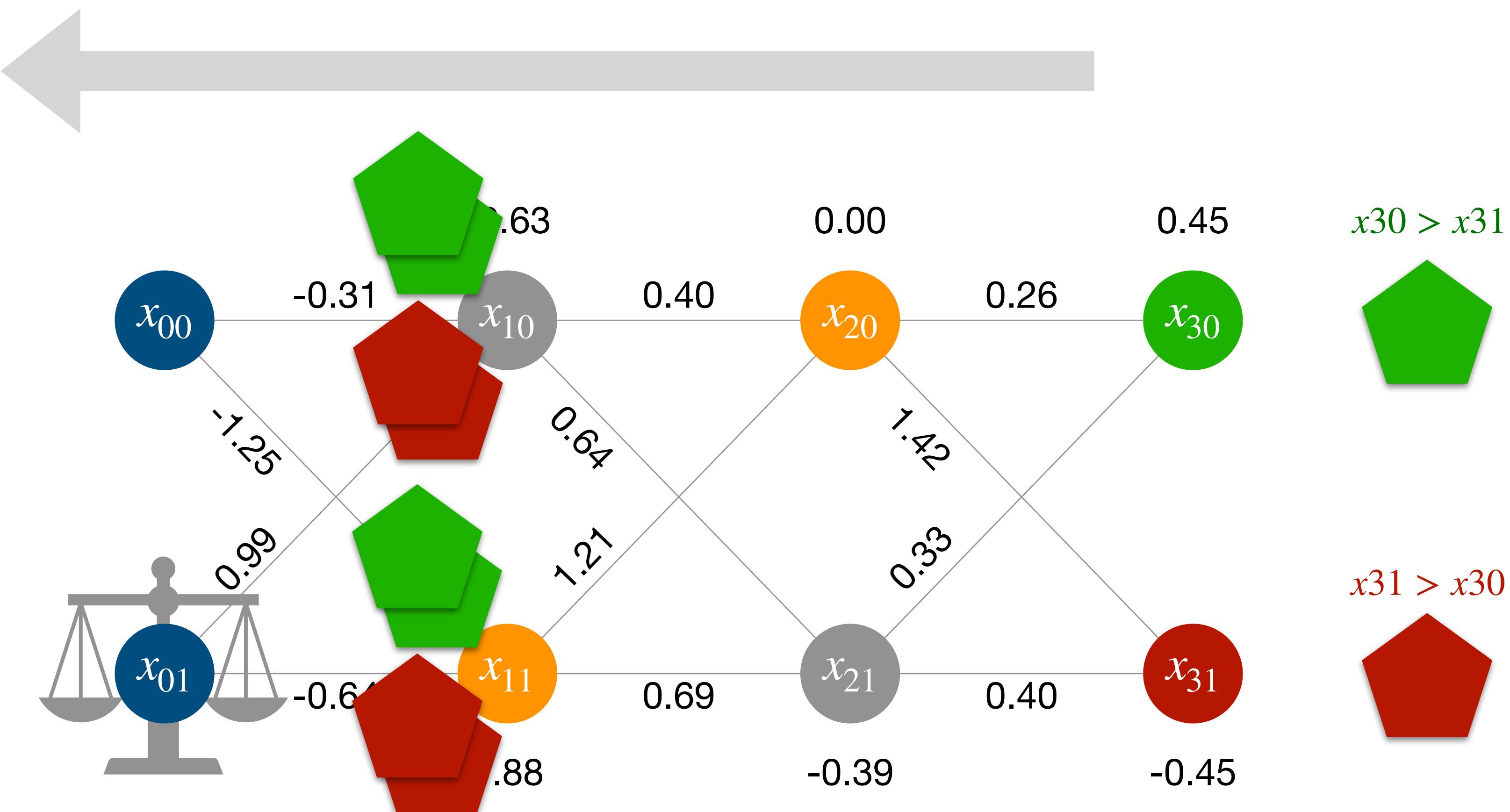
# Backward Analysis



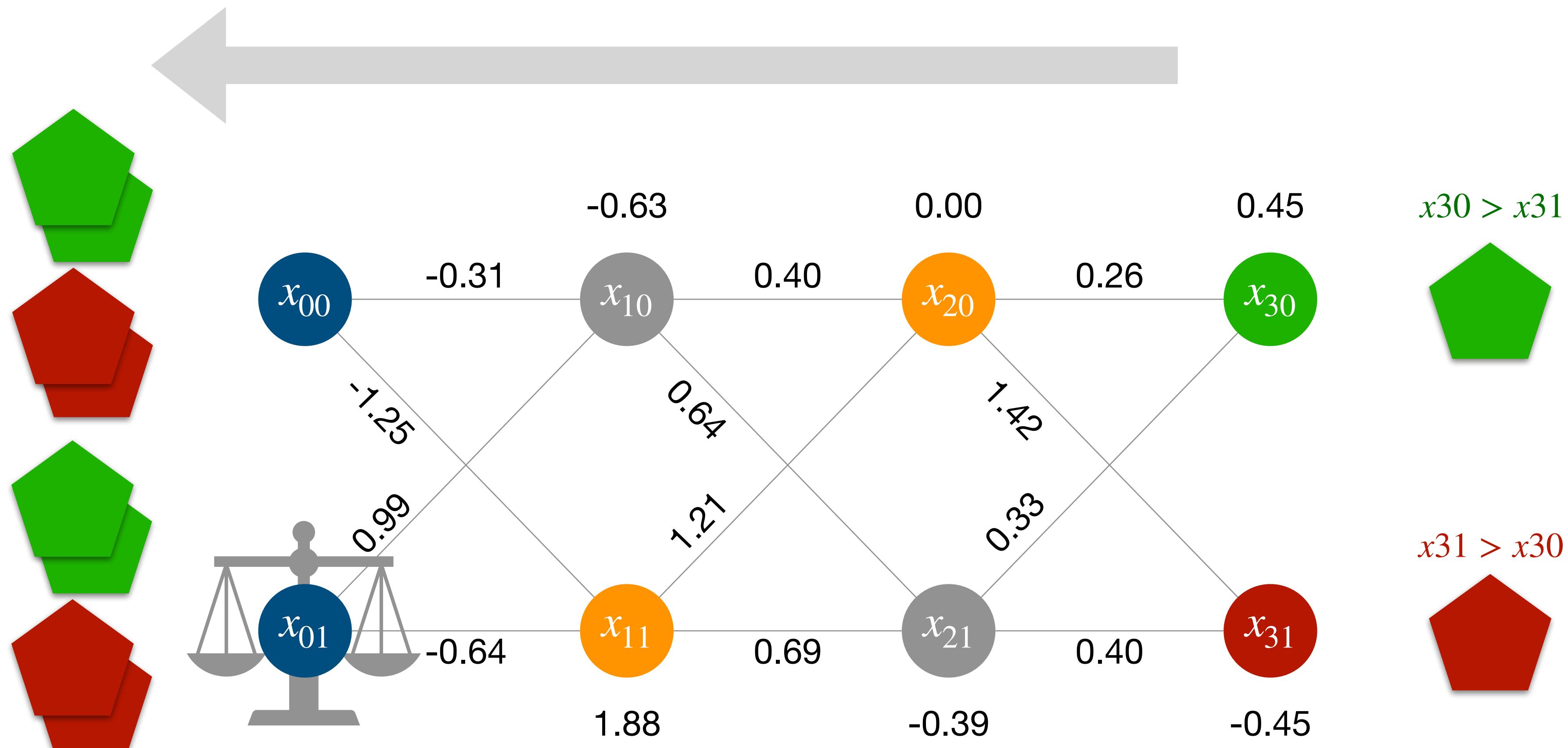
# Backward Analysis



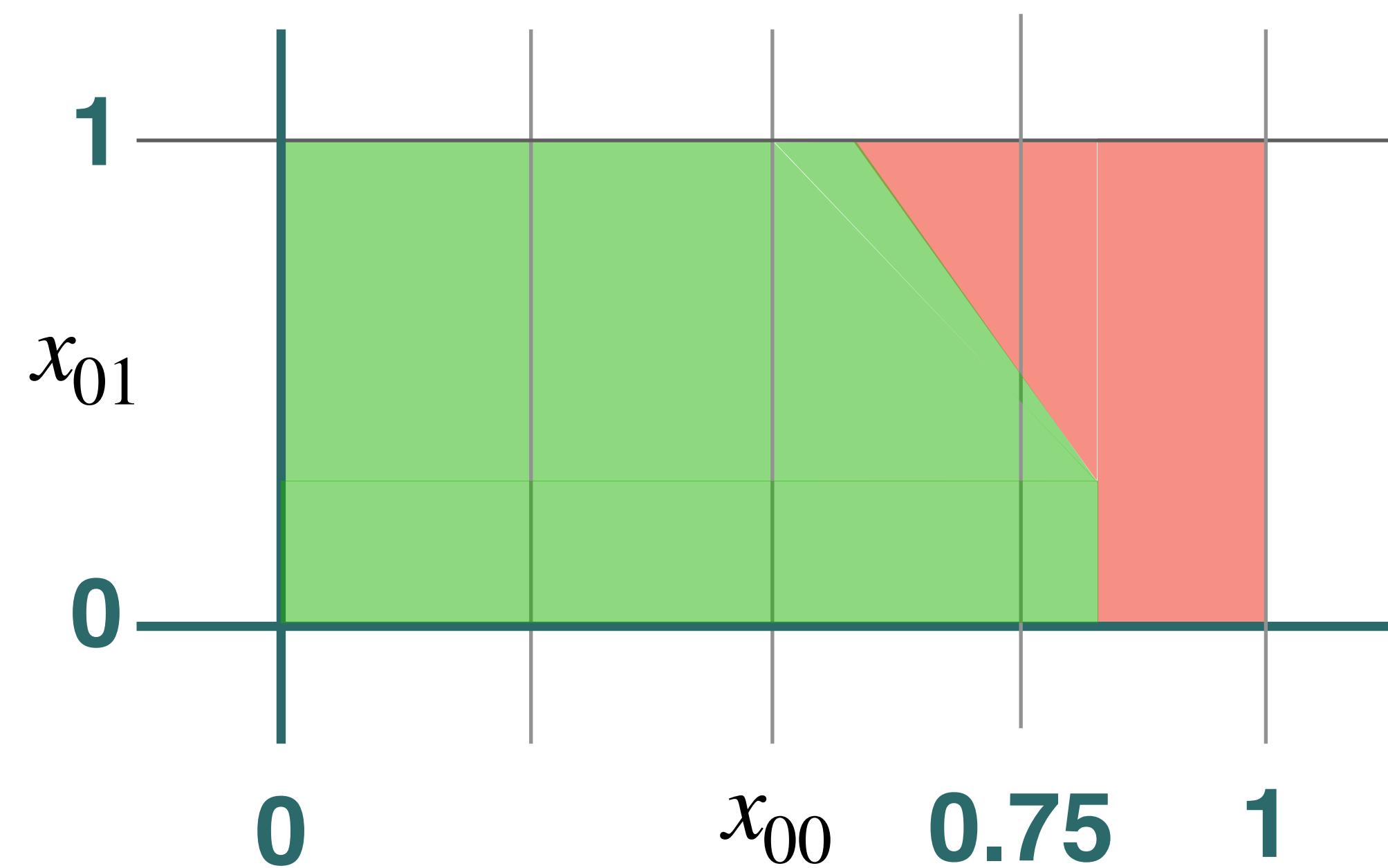
# Backward Analysis



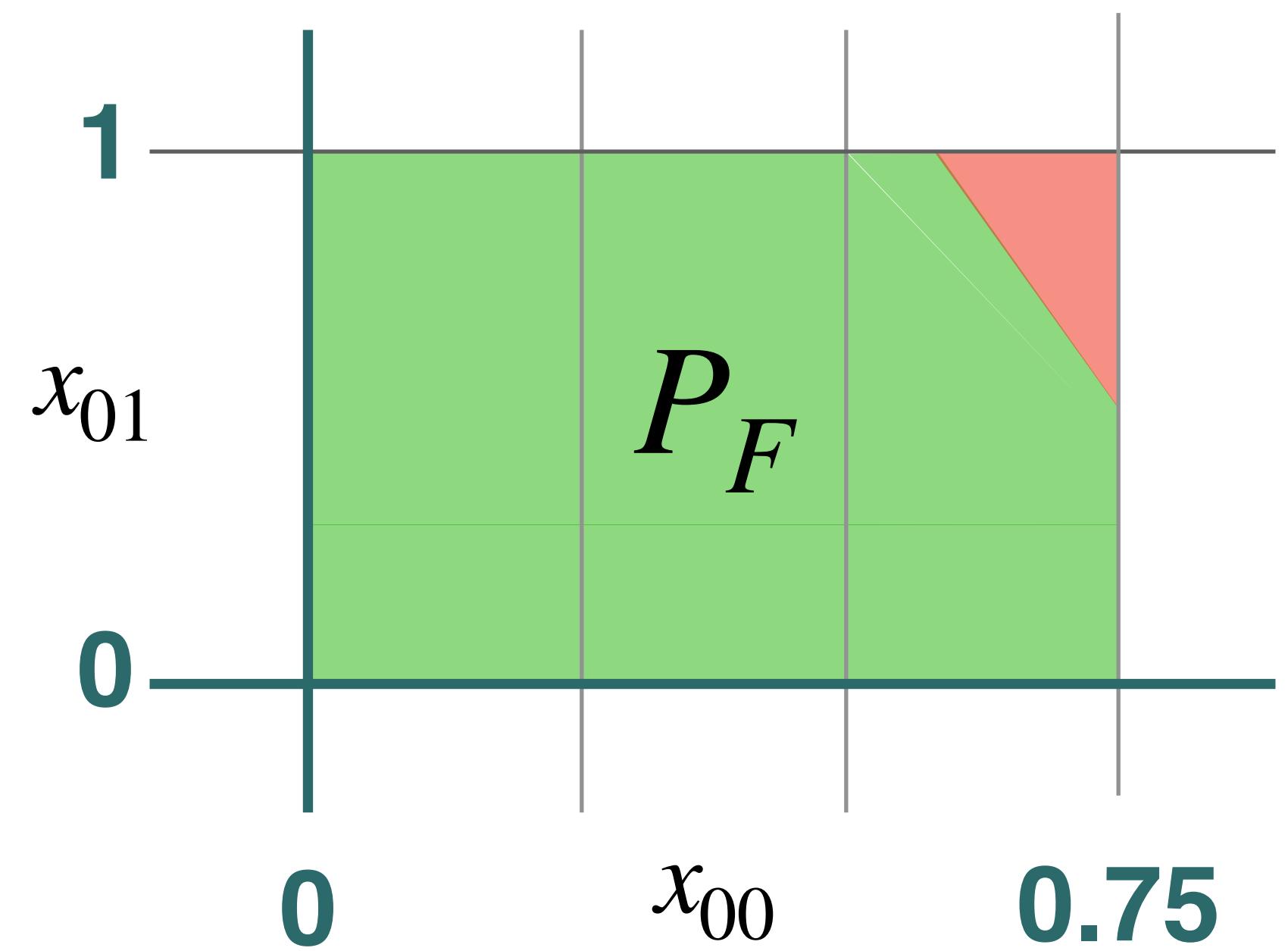
# Backward Analysis



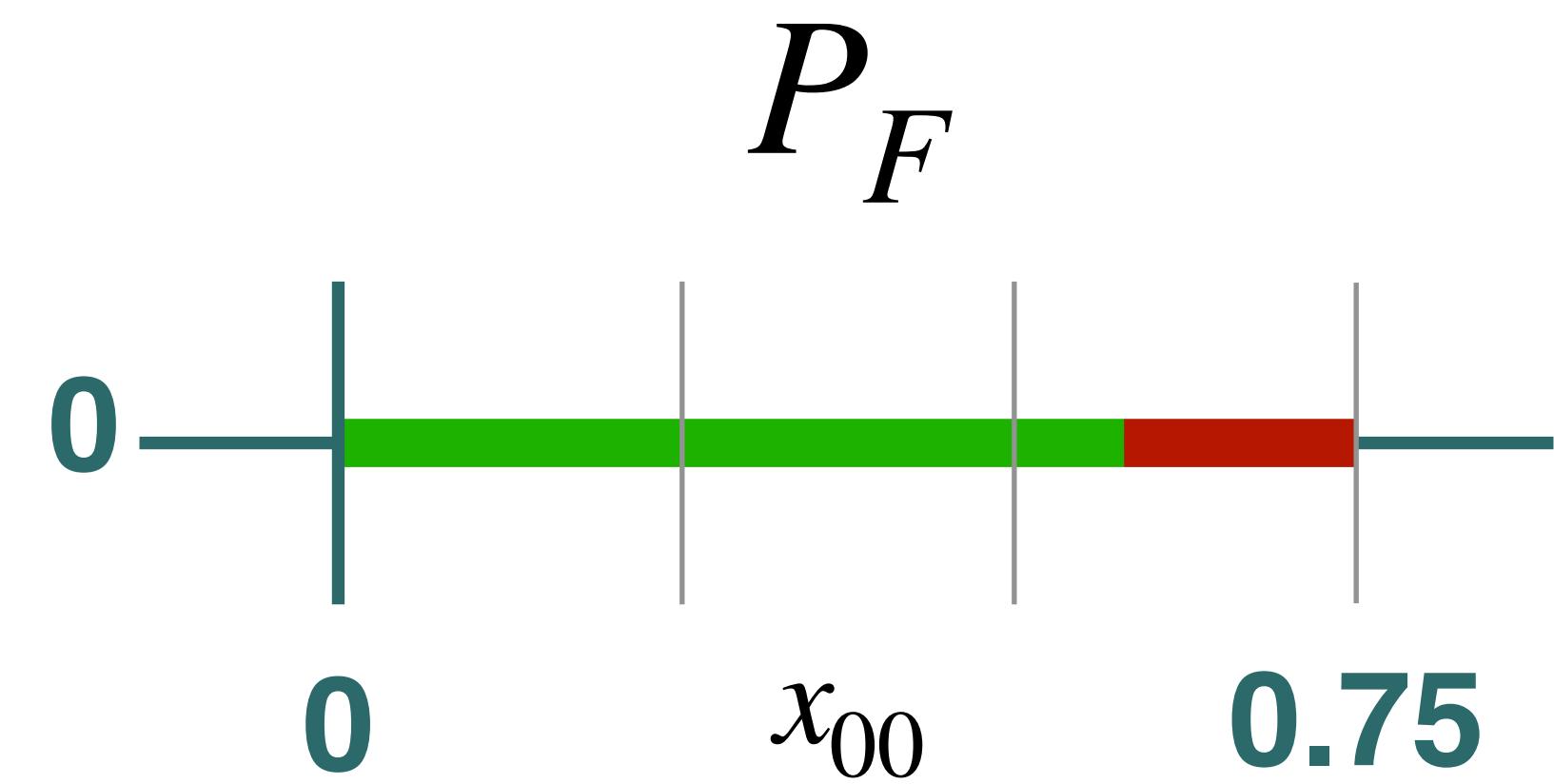
# Analysis Result



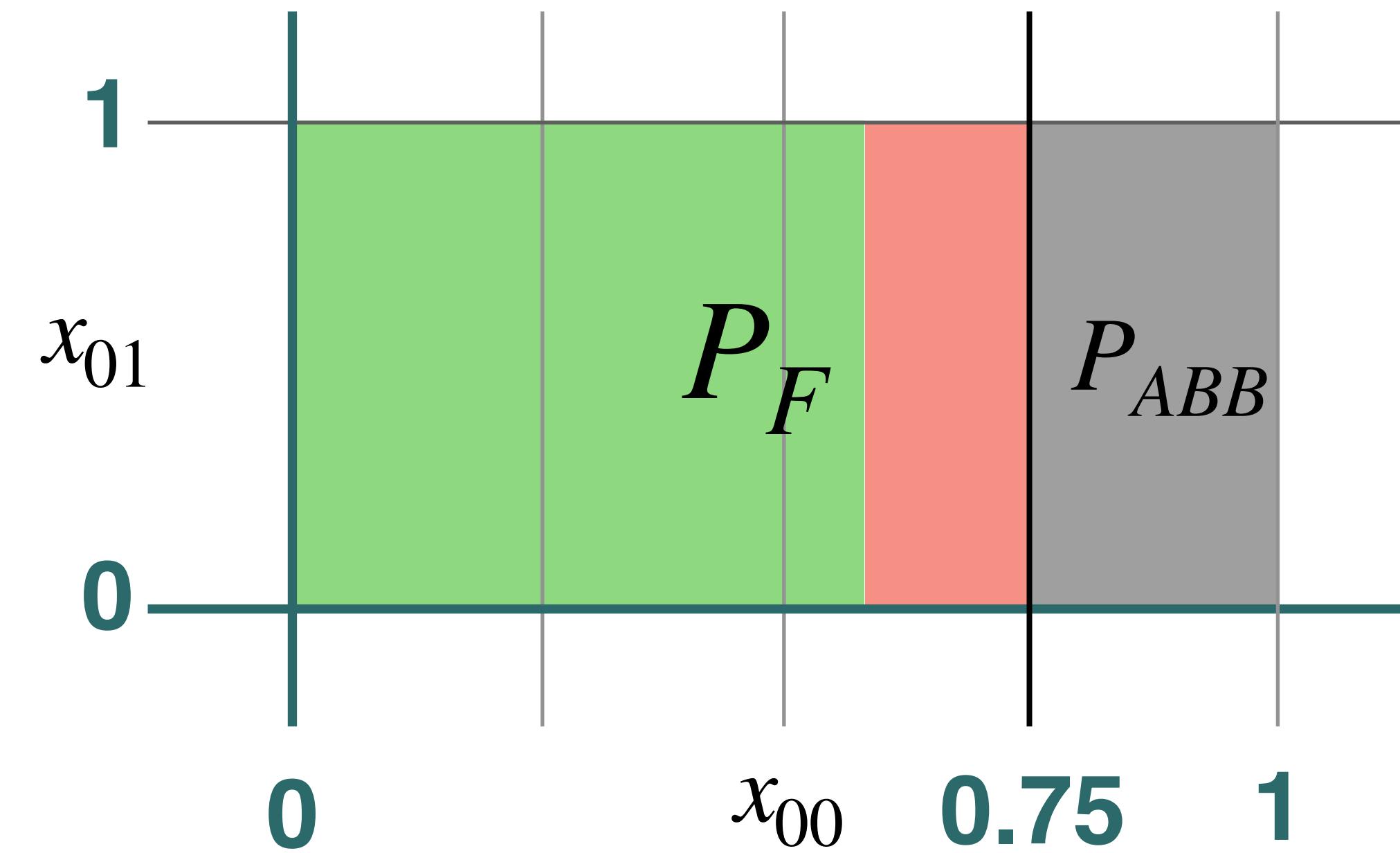
# Analysis Result

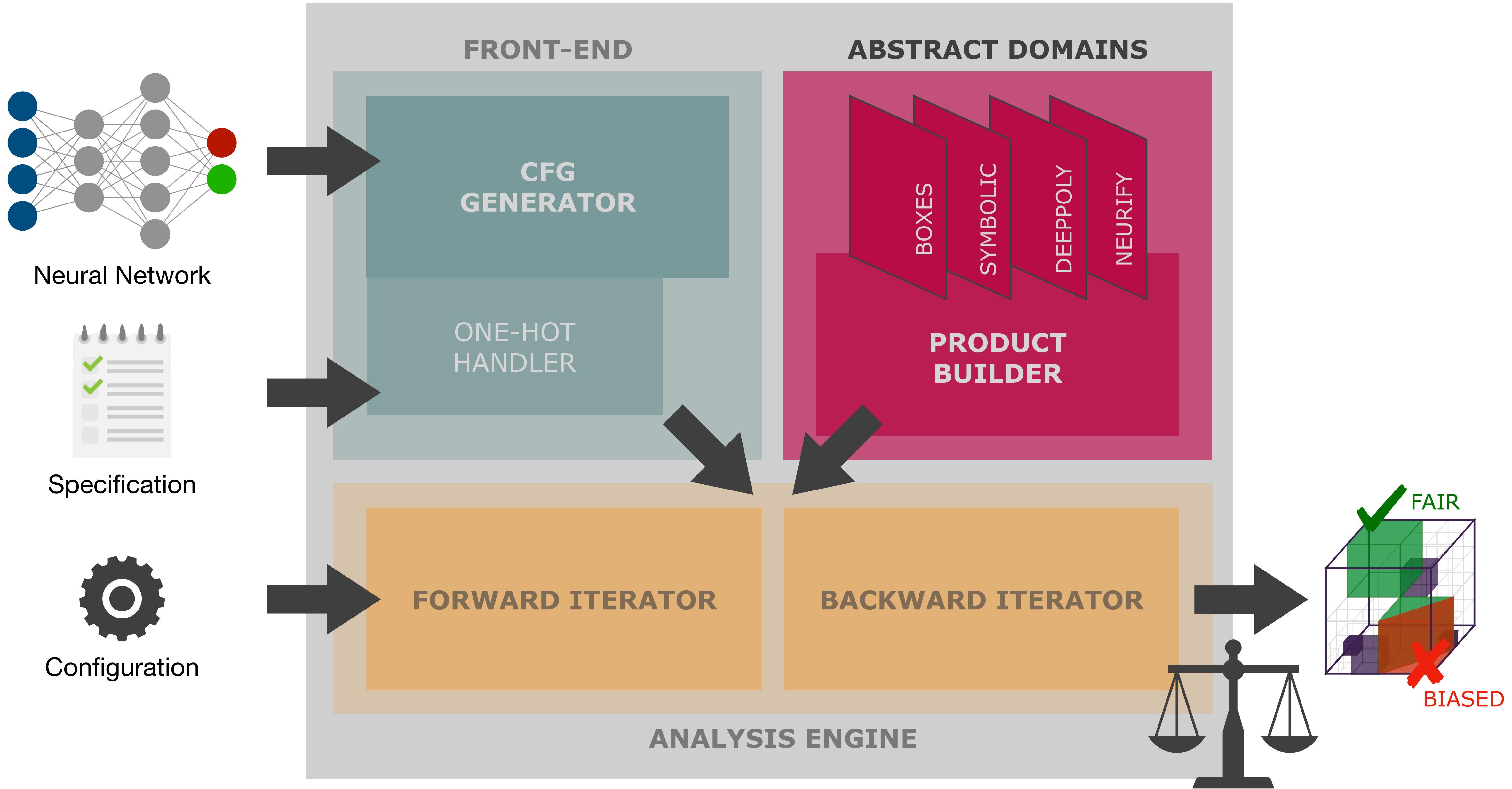


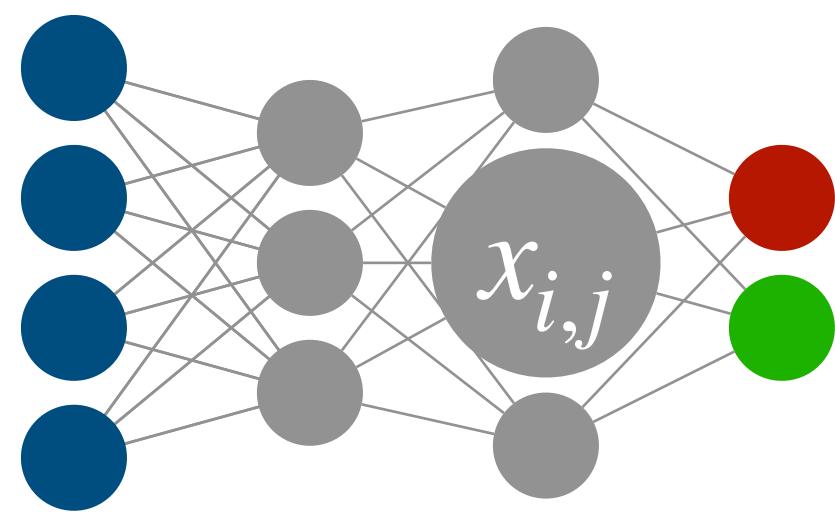
# Analysis Result



# Analysis Result





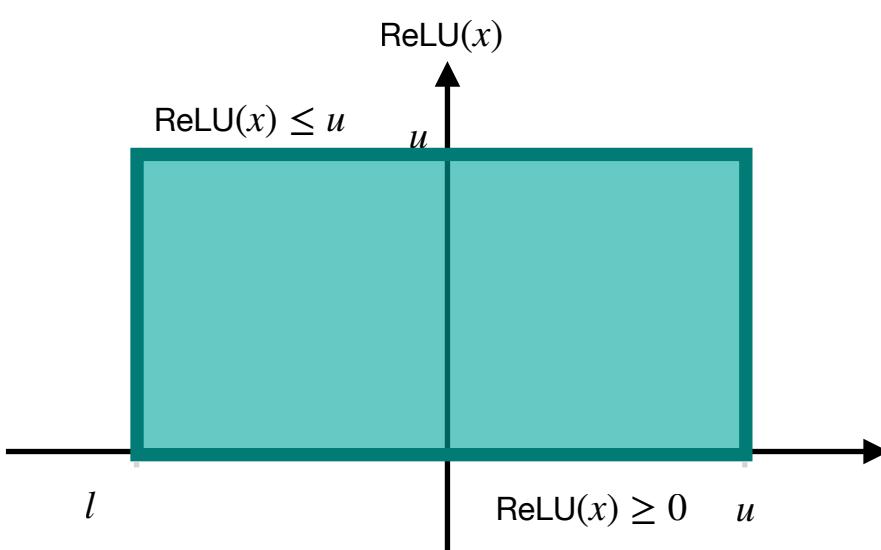
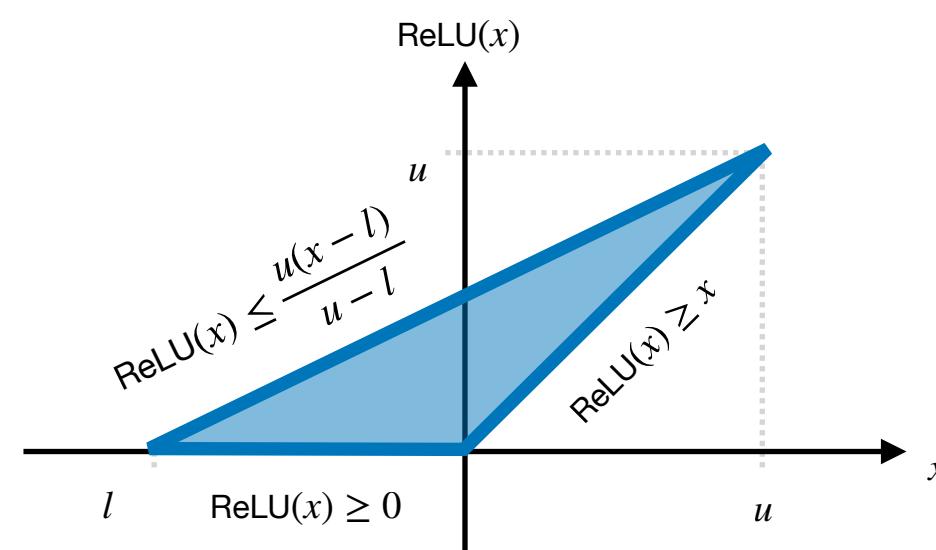


## Symbolic

Li et al. @ SAS 2019

$$[l, u]$$

$$\sum_k m_k \cdot x_k + q$$

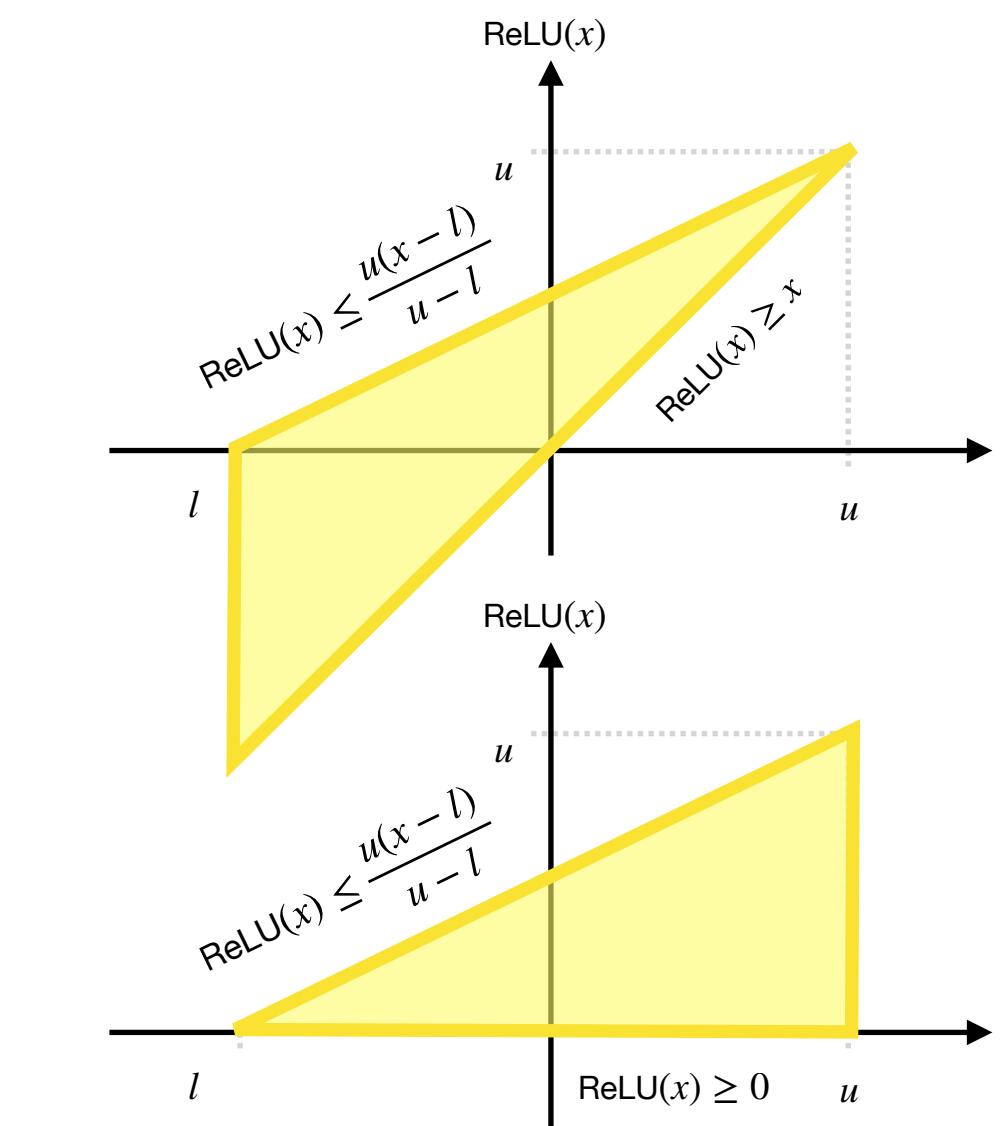


## DeepPoly

Singh et al. @ POPL 2019

$$[l, u]$$

$$[\text{eq}_{\text{low}}, \text{eq}_{\text{up}}]$$

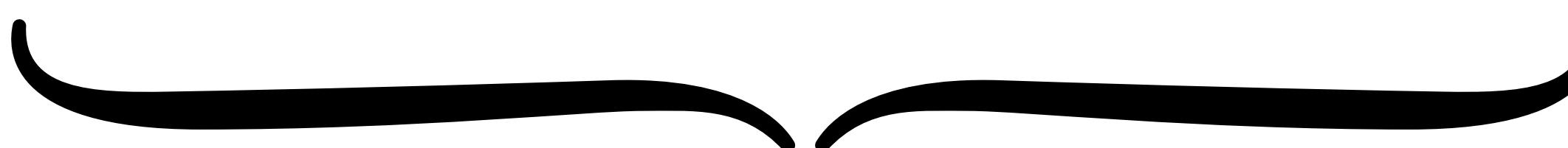
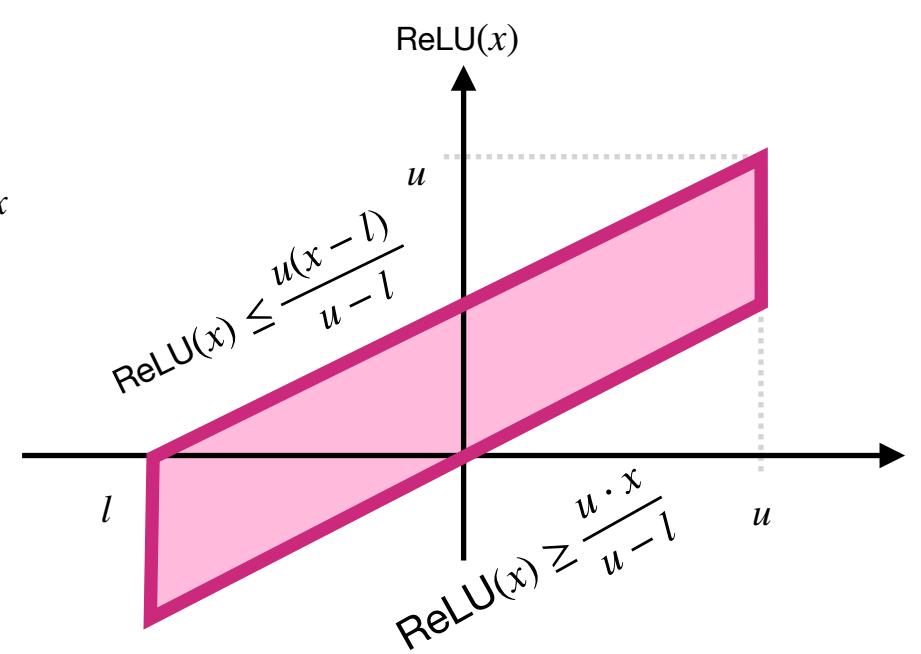


## Neurify

Wang et al. @ NeurIPS 2018

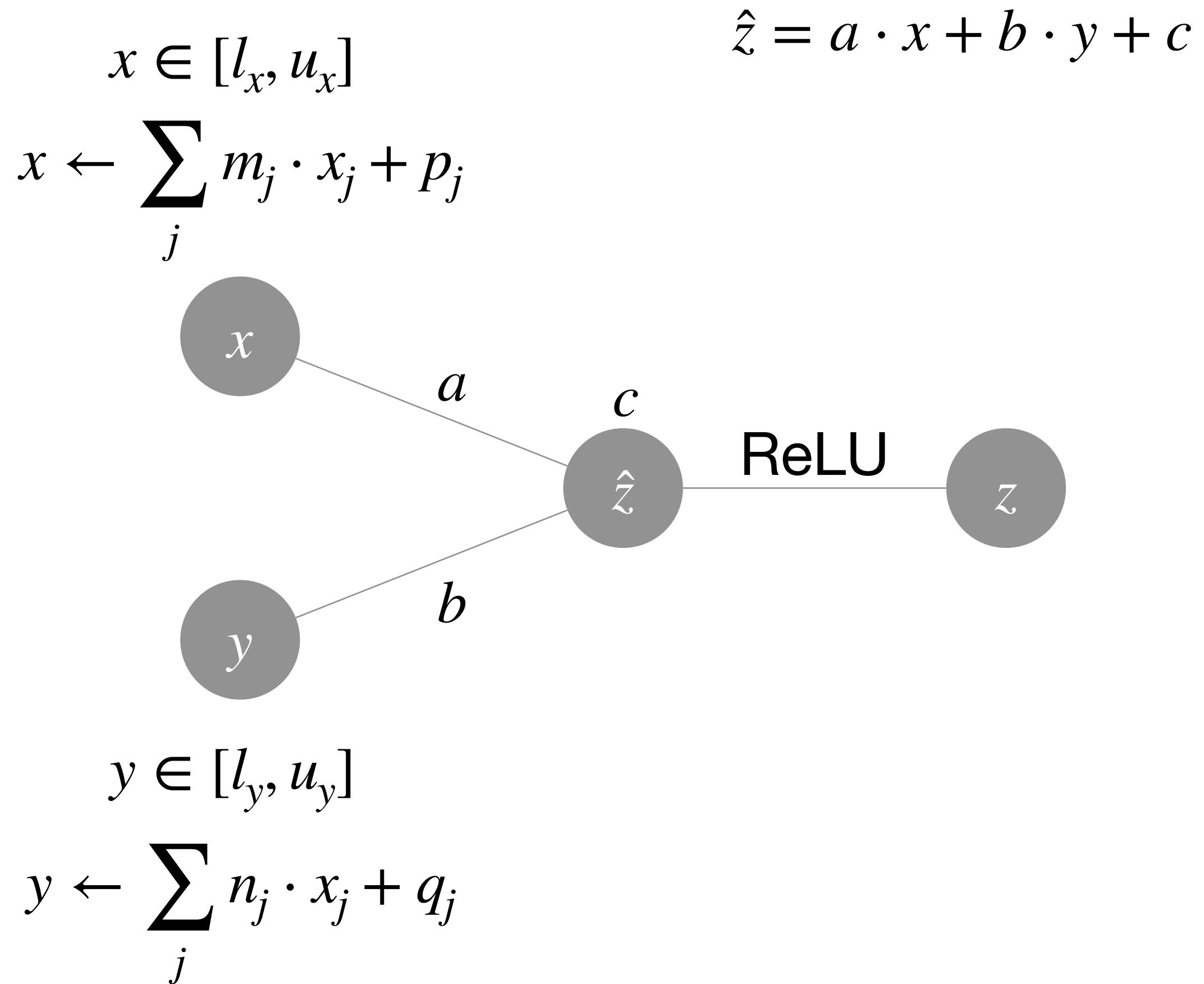
$$[l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$

$$[\text{eq}_{\text{low}}, \text{eq}_{\text{up}}]$$



## Reduced Product

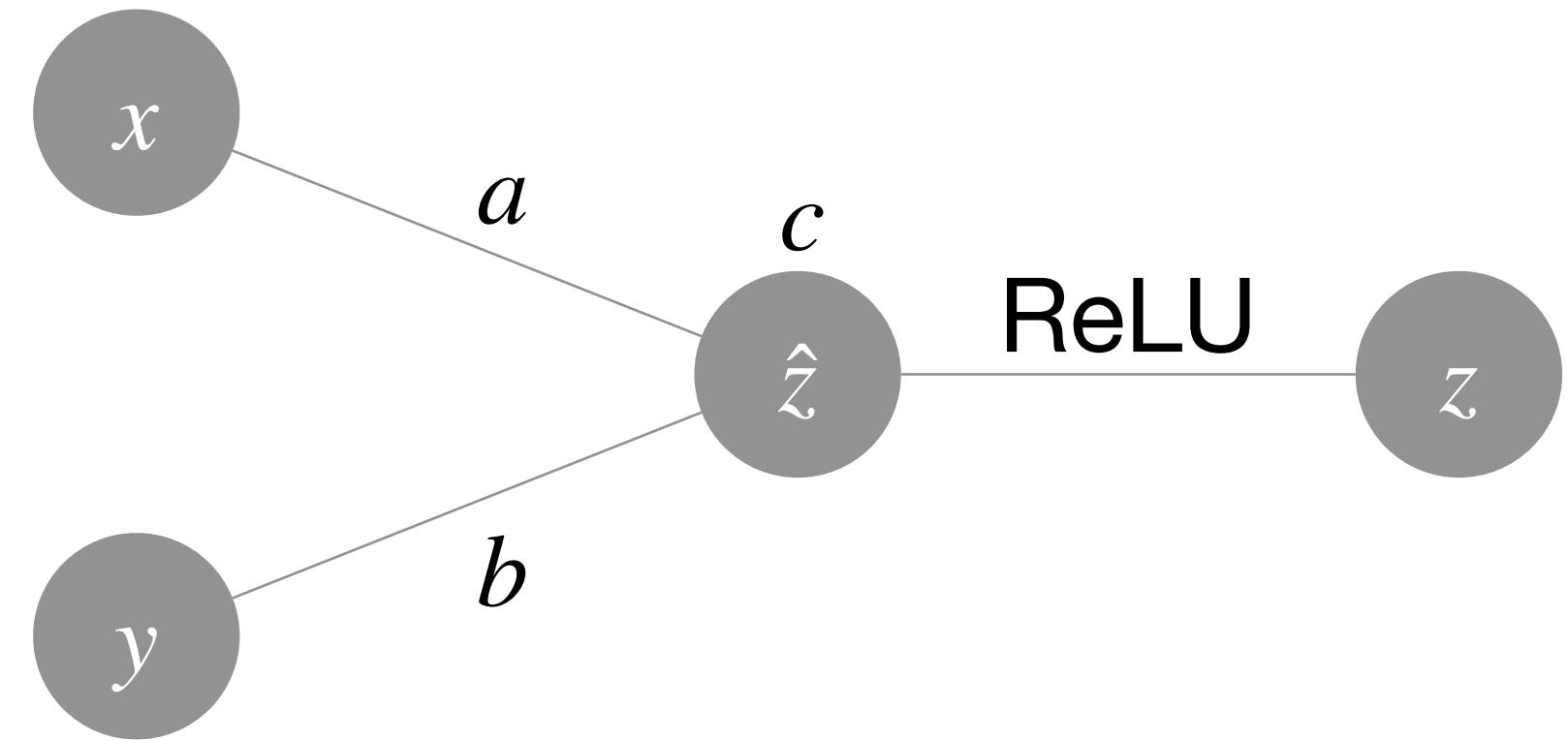
Mazzucato, Urban @ SAS 2021



# Symbolic

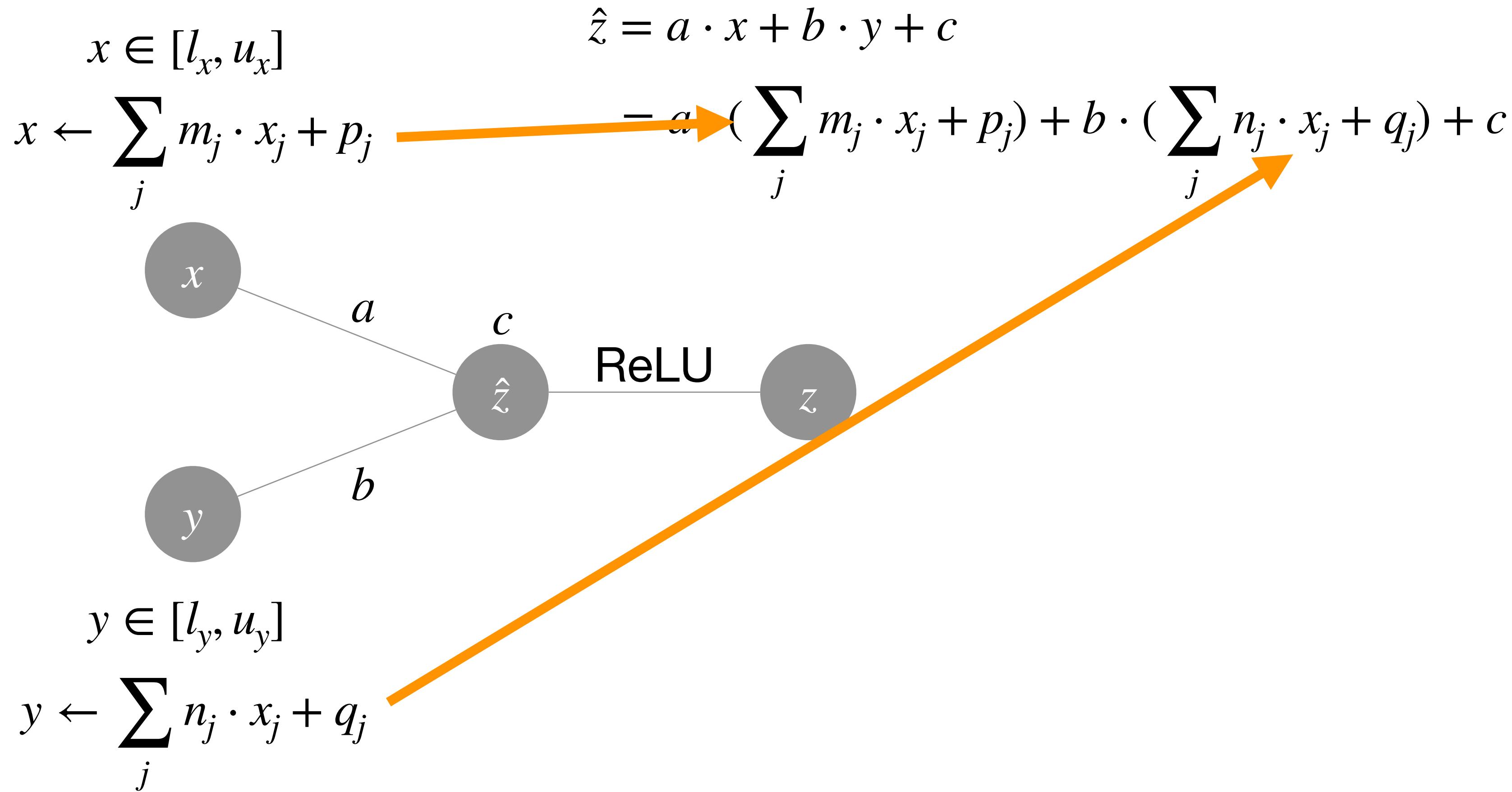
Li et al. @ SAS 2019

$$x \in [l_x, u_x]$$
$$x \leftarrow \sum_j m_j \cdot x_j + p_j$$



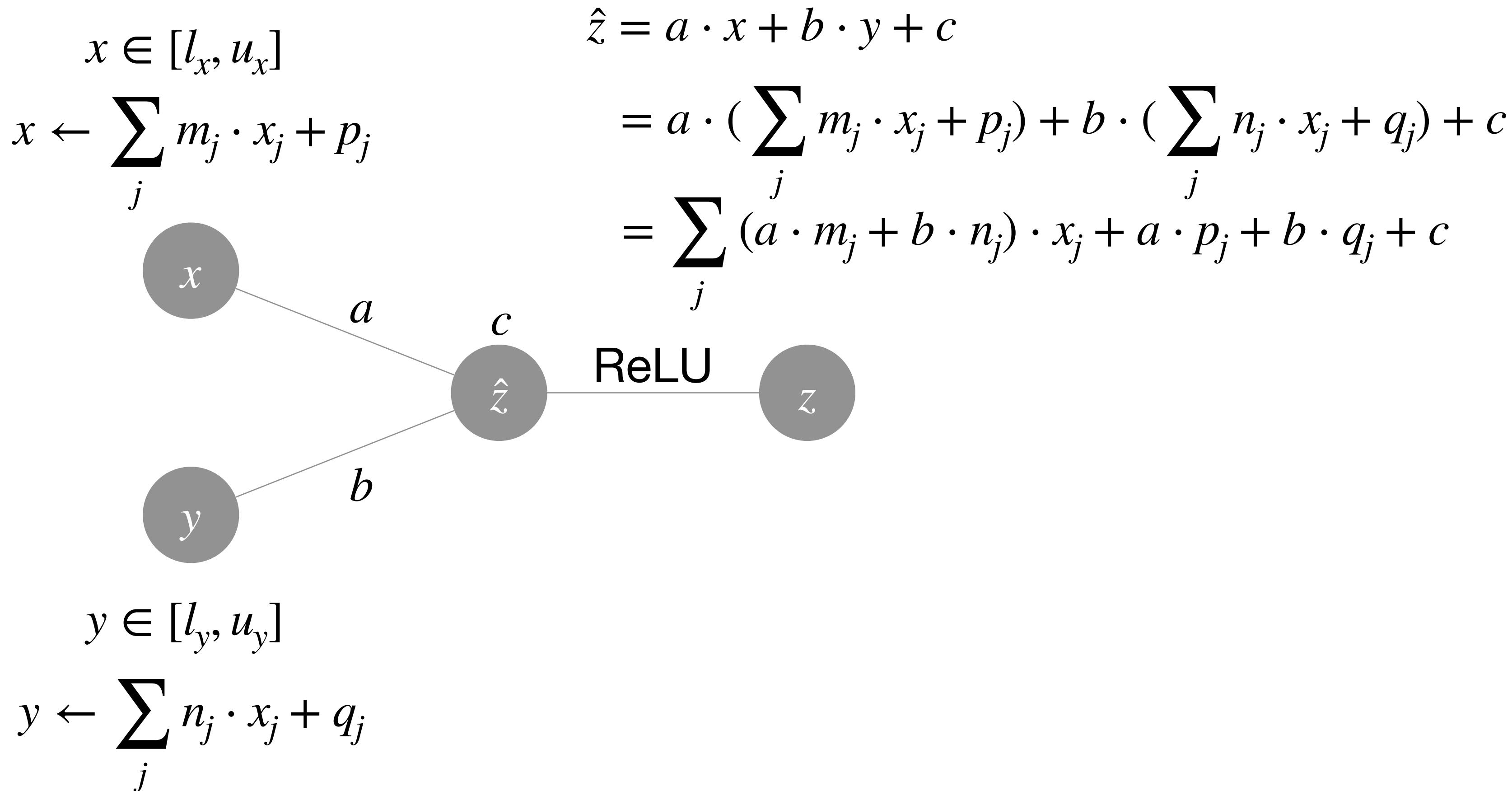
$$y \in [l_y, u_y]$$
$$y \leftarrow \sum_j n_j \cdot x_j + q_j$$

$$\begin{aligned}\hat{z} &= a \cdot x + b \cdot y + c \\ &= a \cdot \left( \sum_j m_j \cdot x_j + p_j \right) + b \cdot \left( \sum_j n_j \cdot x_j + q_j \right) + c\end{aligned}$$



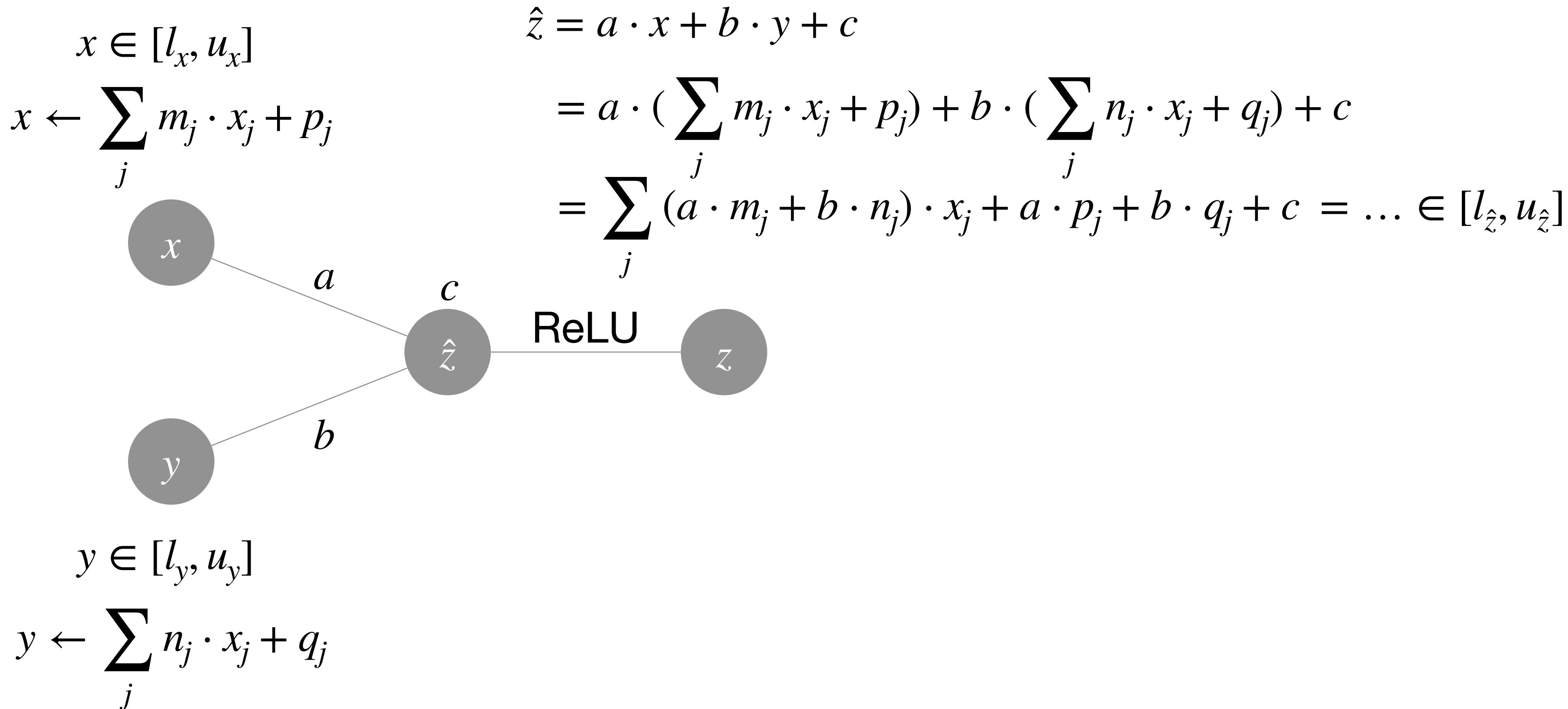
# Symbolic

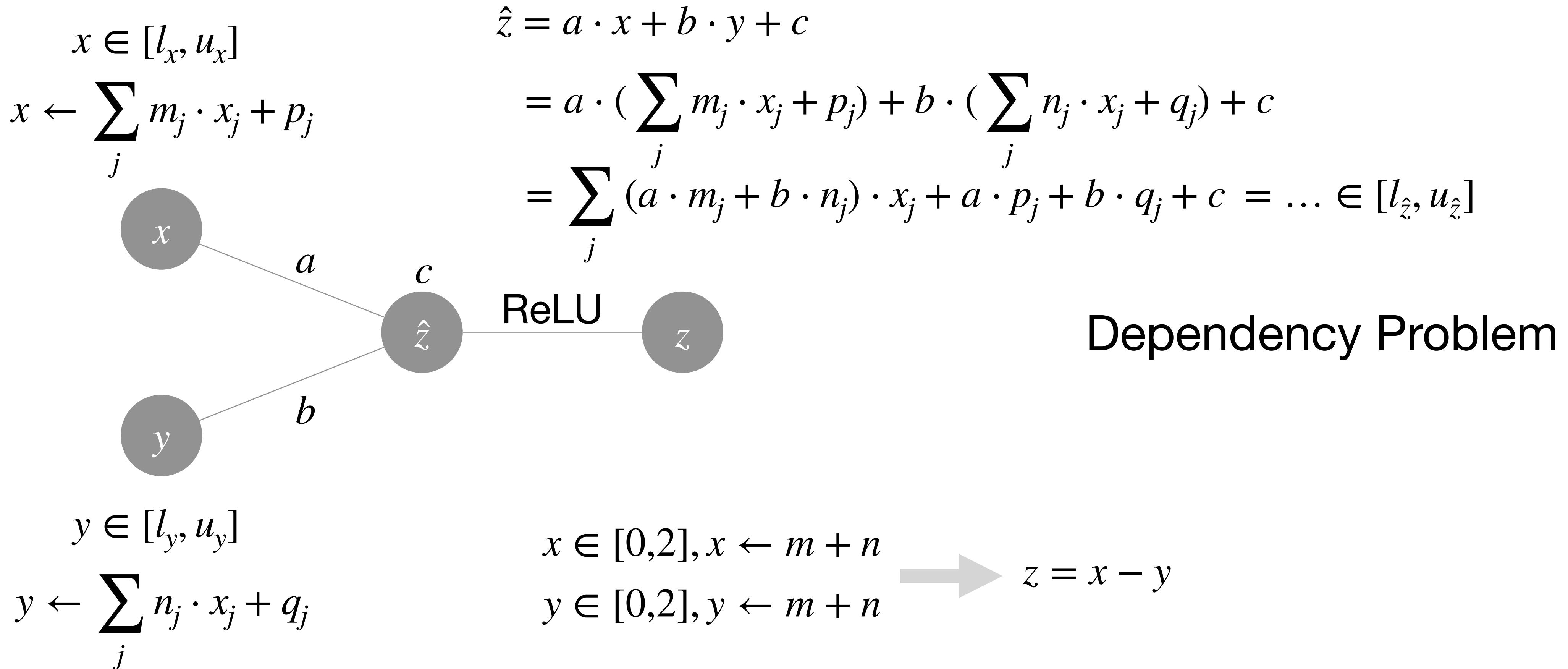
Li et al. @ SAS 2019



# Symbolic

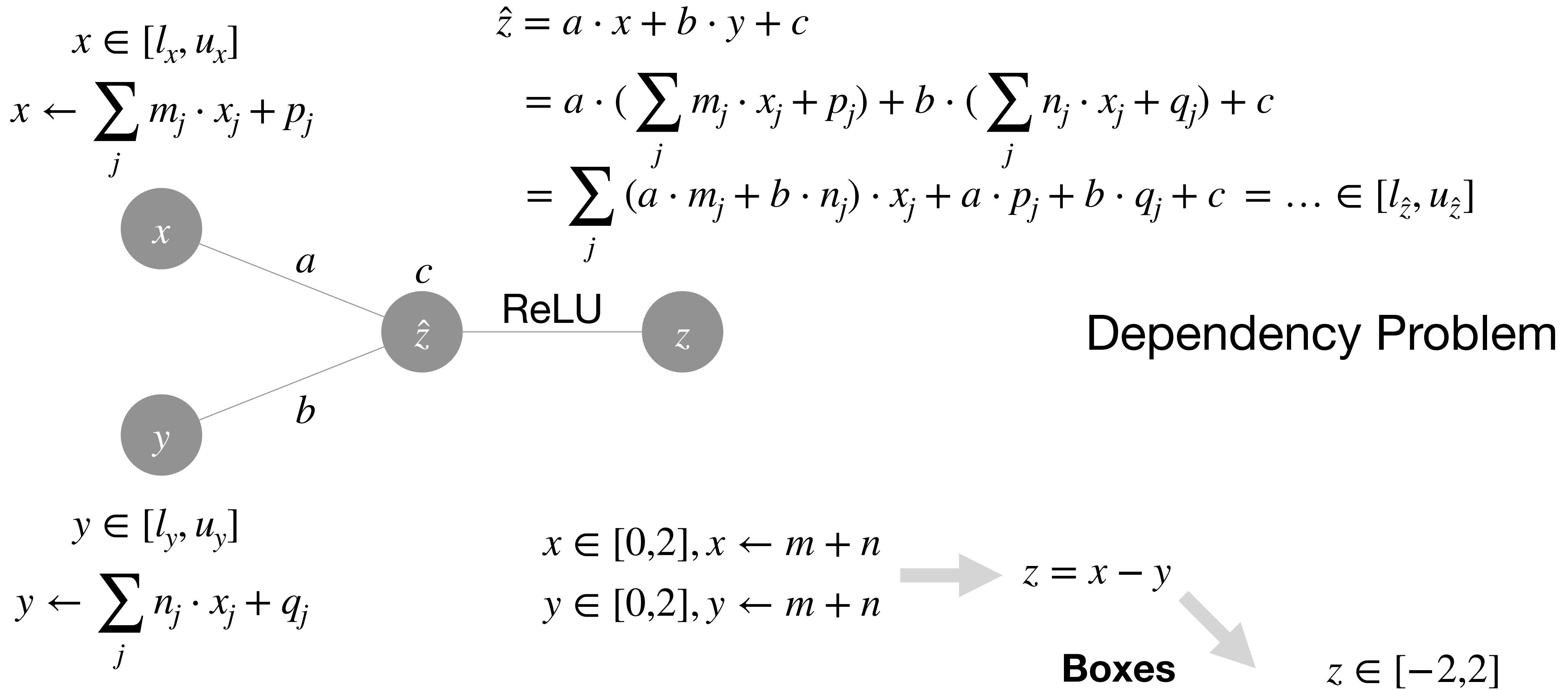
Li et al. @ SAS 2019





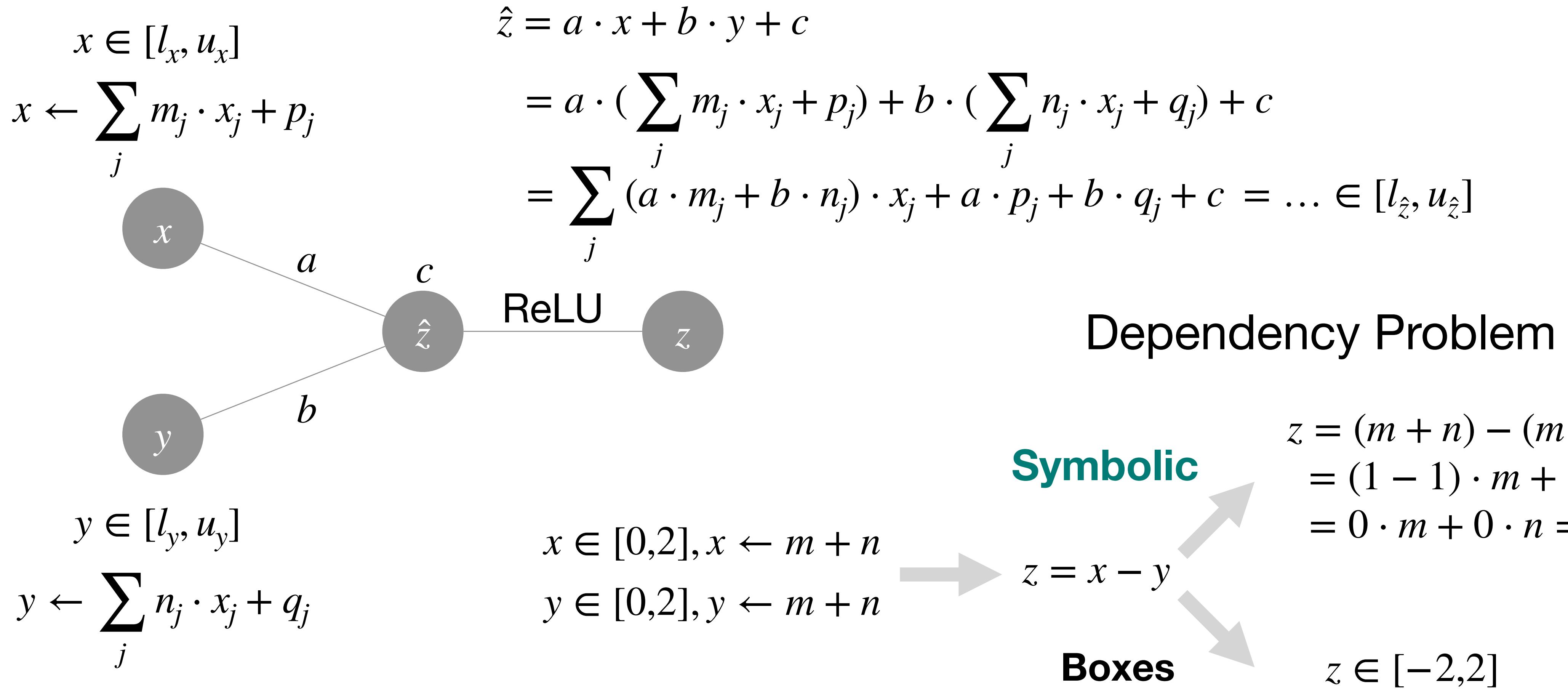
# Symbolic

Li et al. @ SAS 2019



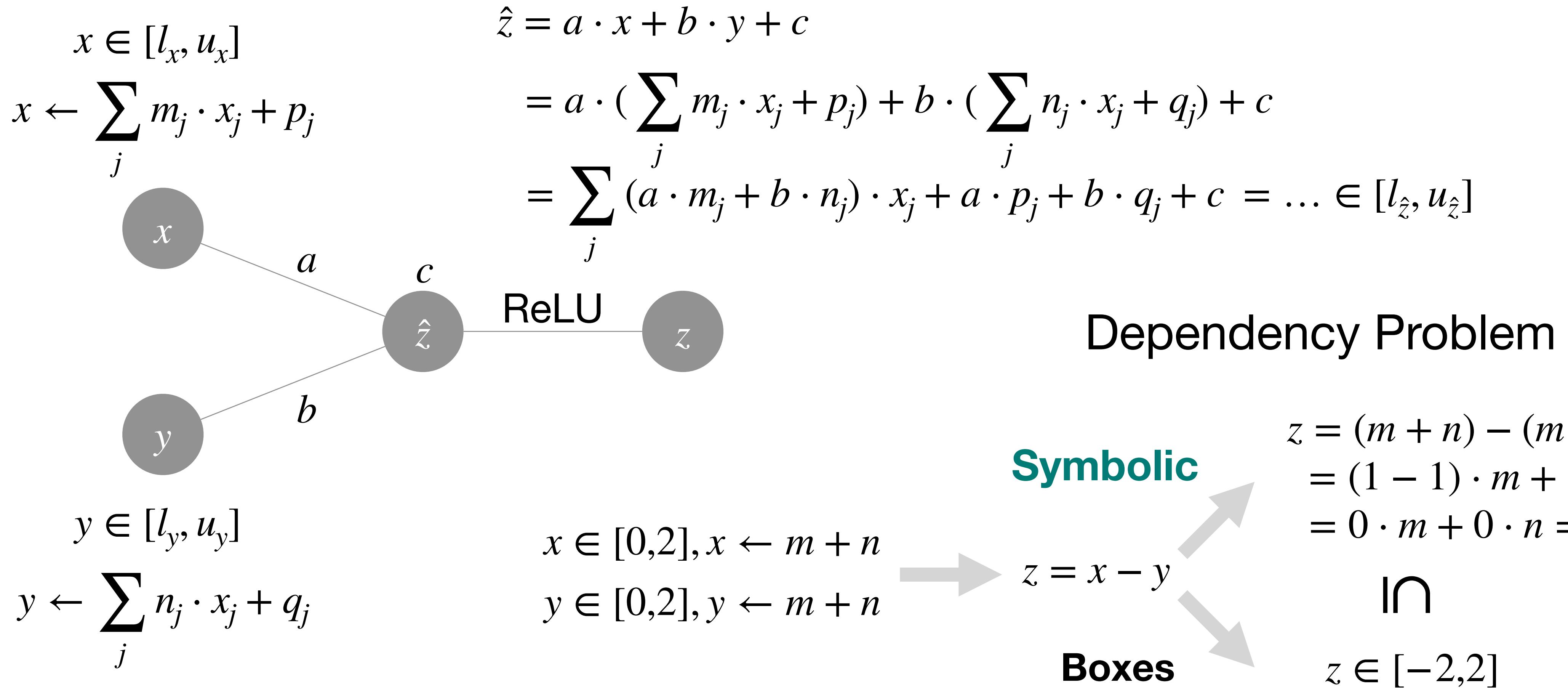
# Symbolic

Li et al. @ SAS 2019



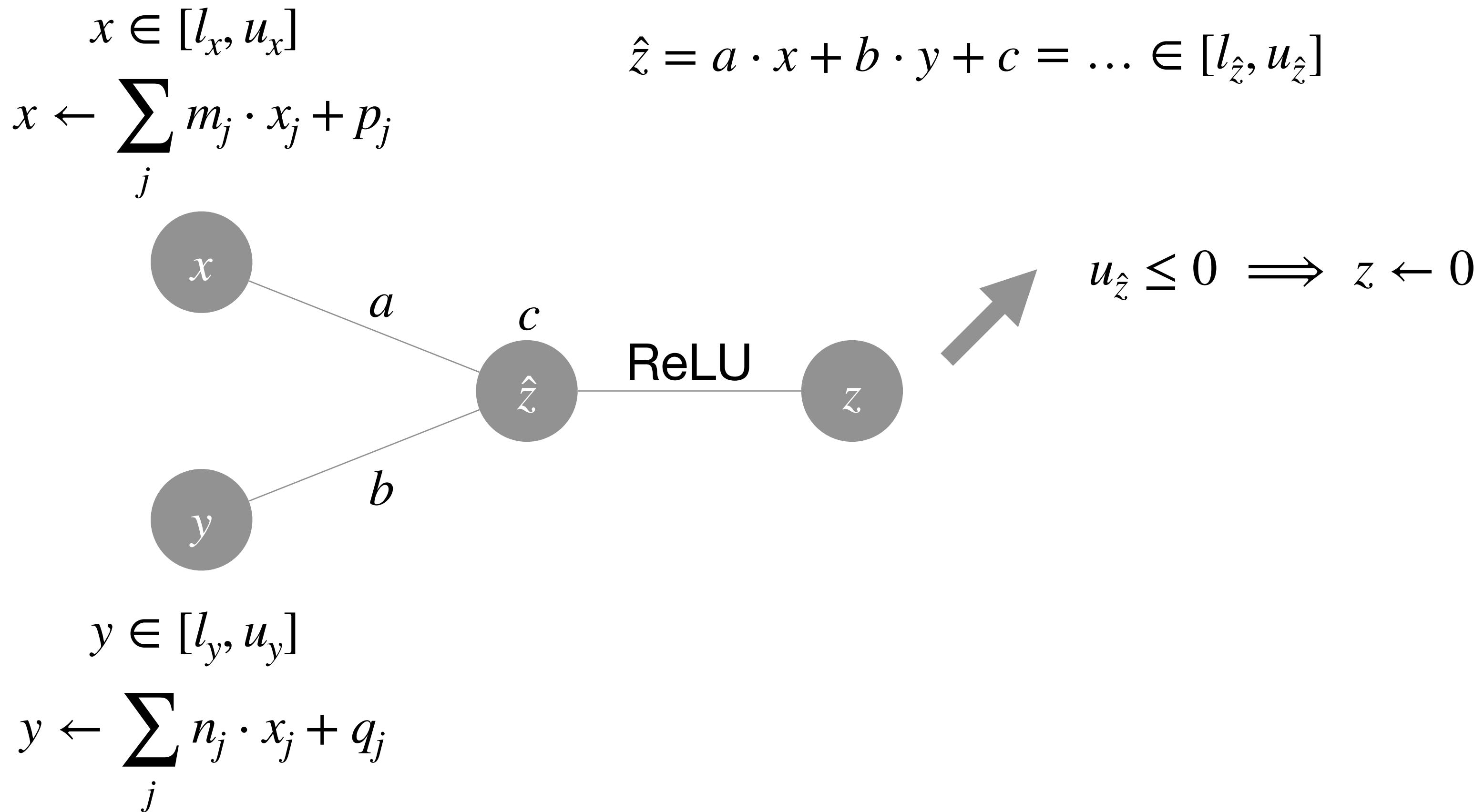
# Symbolic

Li et al. @ SAS 2019



# Symbolic

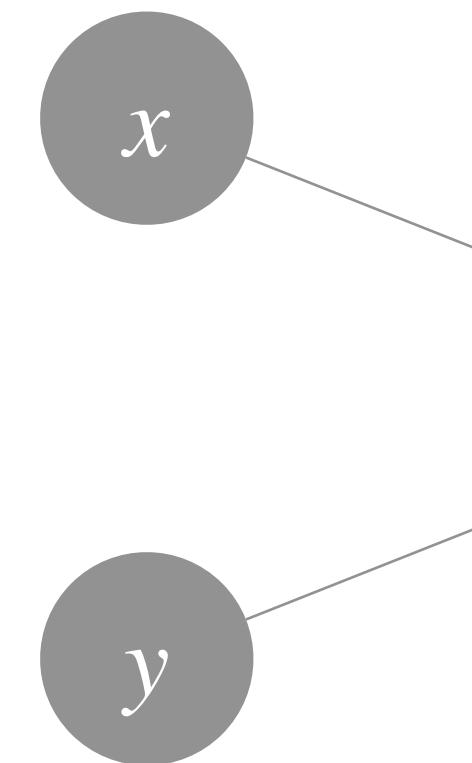
Li et al. @ SAS 2019



# Symbolic

Li et al. @ SAS 2019

$$x \in [l_x, u_x]$$
$$x \leftarrow \sum_j m_j \cdot x_j + p_j$$



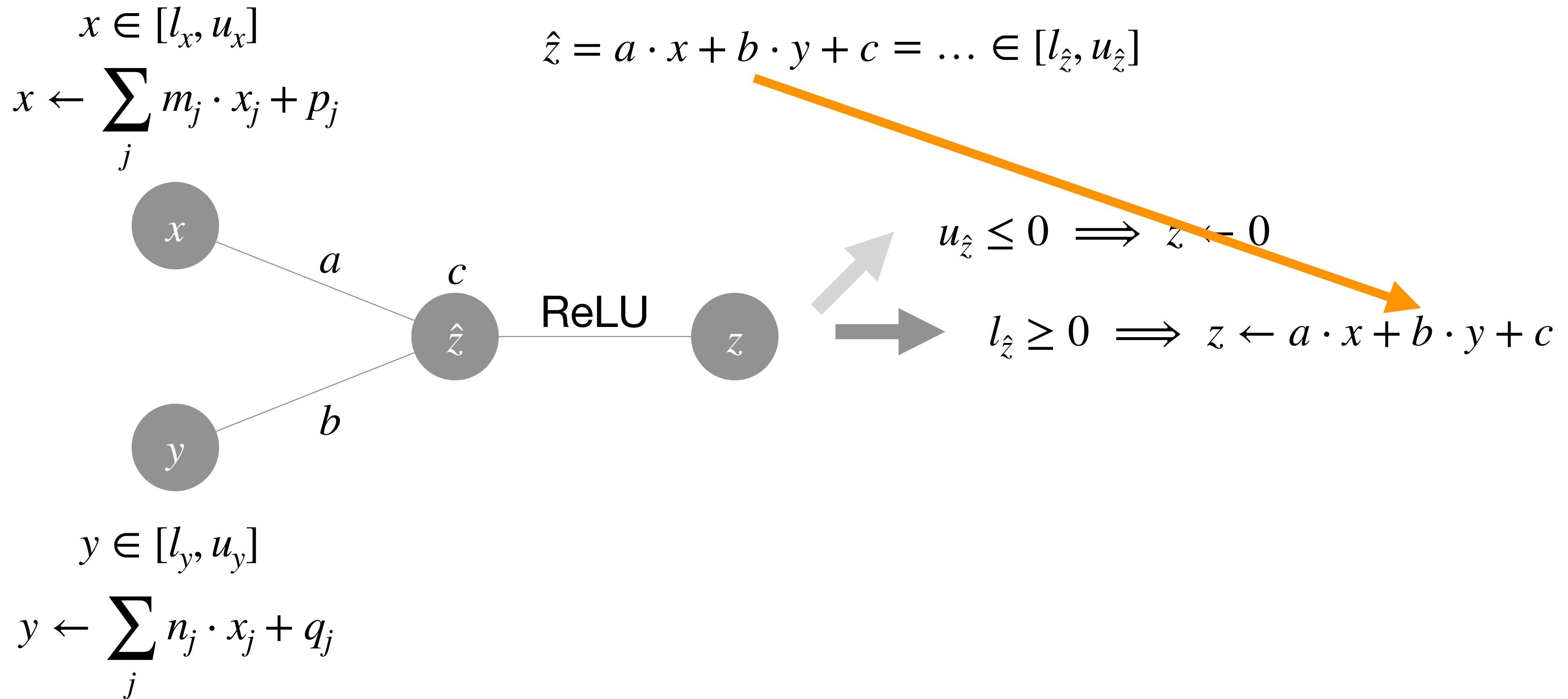
$$\hat{z} = a \cdot x + b \cdot y + c = \dots \in [l_{\hat{z}}, u_{\hat{z}}]$$

$$y \in [l_y, u_y]$$
$$y \leftarrow \sum_j n_j \cdot x_j + q_j$$

$$u_{\hat{z}} \leq 0 \implies z \leftarrow 0$$
$$l_{\hat{z}} \geq 0 \implies z \leftarrow a \cdot x + b \cdot y + c$$

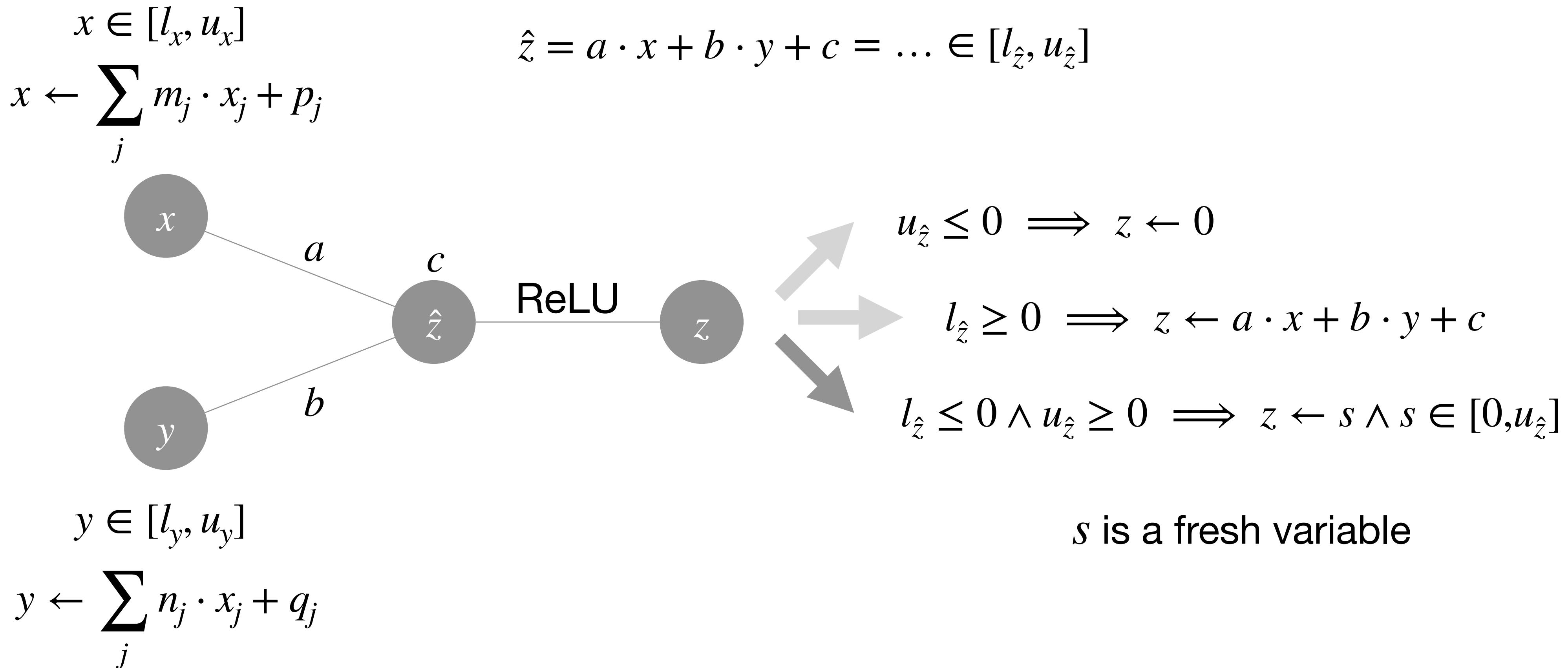
# Symbolic

Li et al. @ SAS 2019



# Symbolic

Li et al. @ SAS 2019



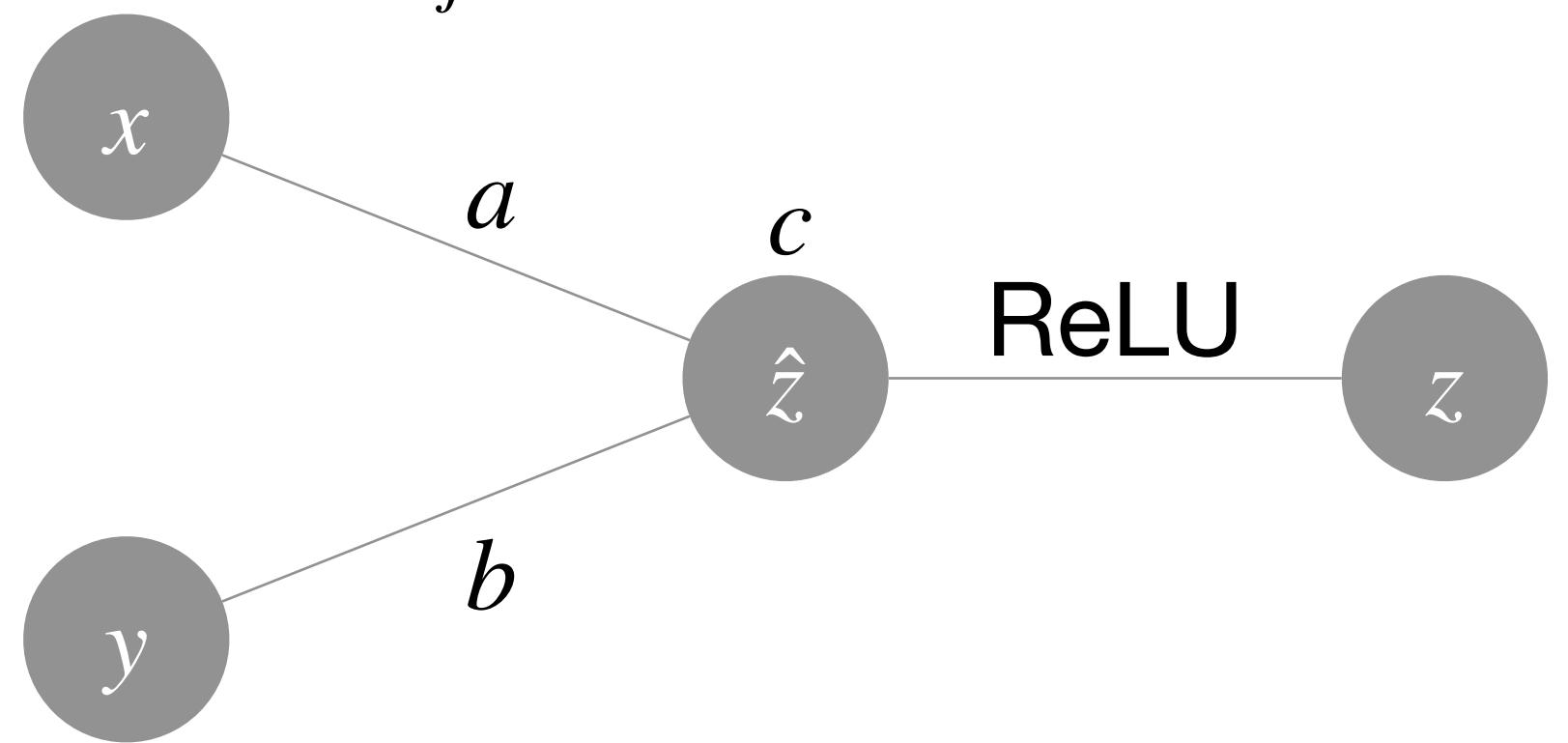
# DeepPoly

Singh et al. @ POPL 2019

$$x \in [l_x, u_x]$$

$$\sum_j \underline{m}_j \cdot x_j + \underline{p}_j \leq x \leq \sum_j \bar{m}_j \cdot x_j + \bar{p}_j$$

$$\hat{z} = a \cdot x + b \cdot y + c$$



$$y \in [l_y, u_y]$$

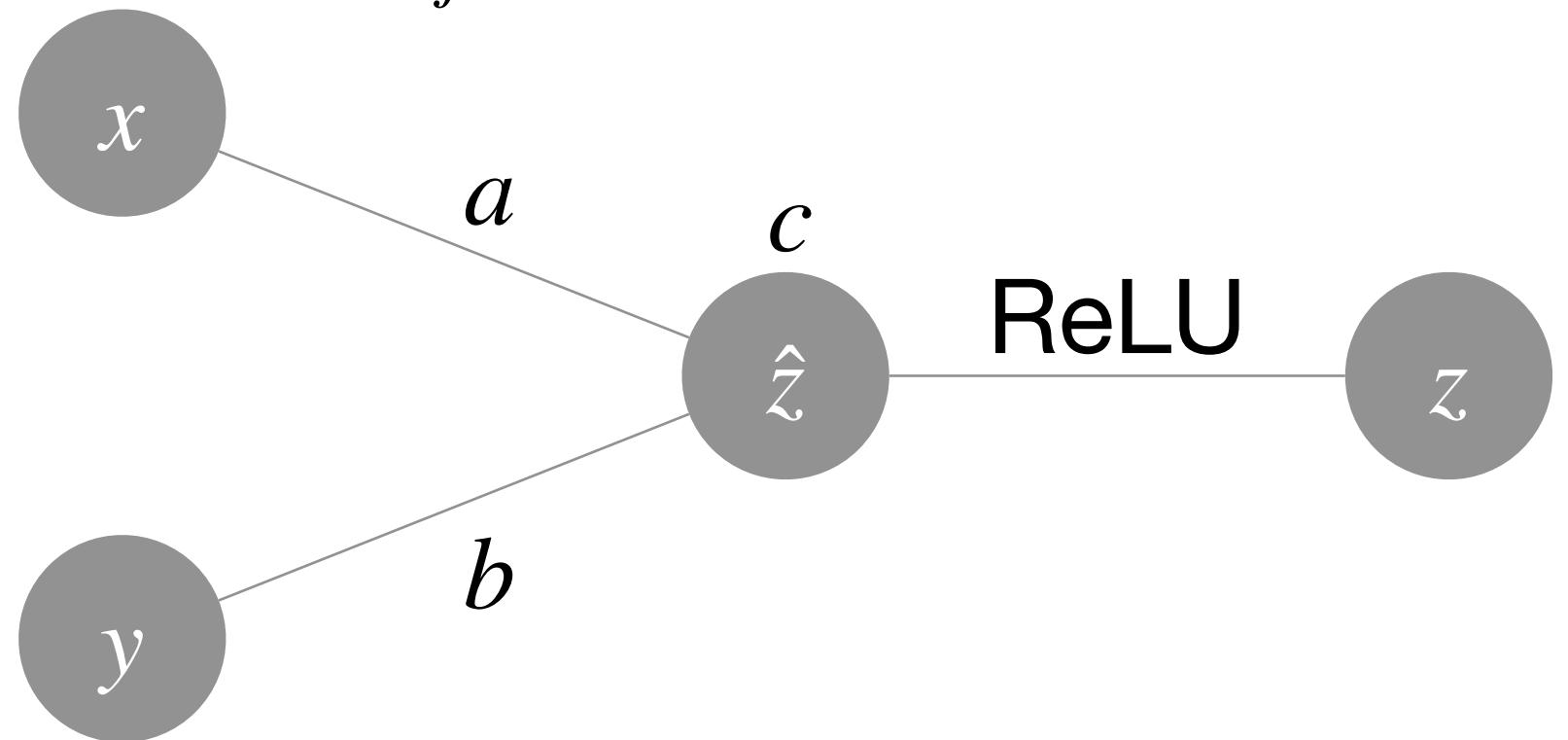
$$\sum_j \underline{n}_j \cdot y_j + \underline{q}_j \leq y \leq \sum_j \bar{n}_j \cdot y_j + \bar{q}_j$$

# DeepPoly

Singh et al. @ POPL 2019

$$x \in [l_x, u_x]$$

$$\sum_j \underline{m}_j \cdot x_j + \underline{p}_j \leq x \leq \sum_j \bar{m}_j \cdot x_j + \bar{p}_j$$



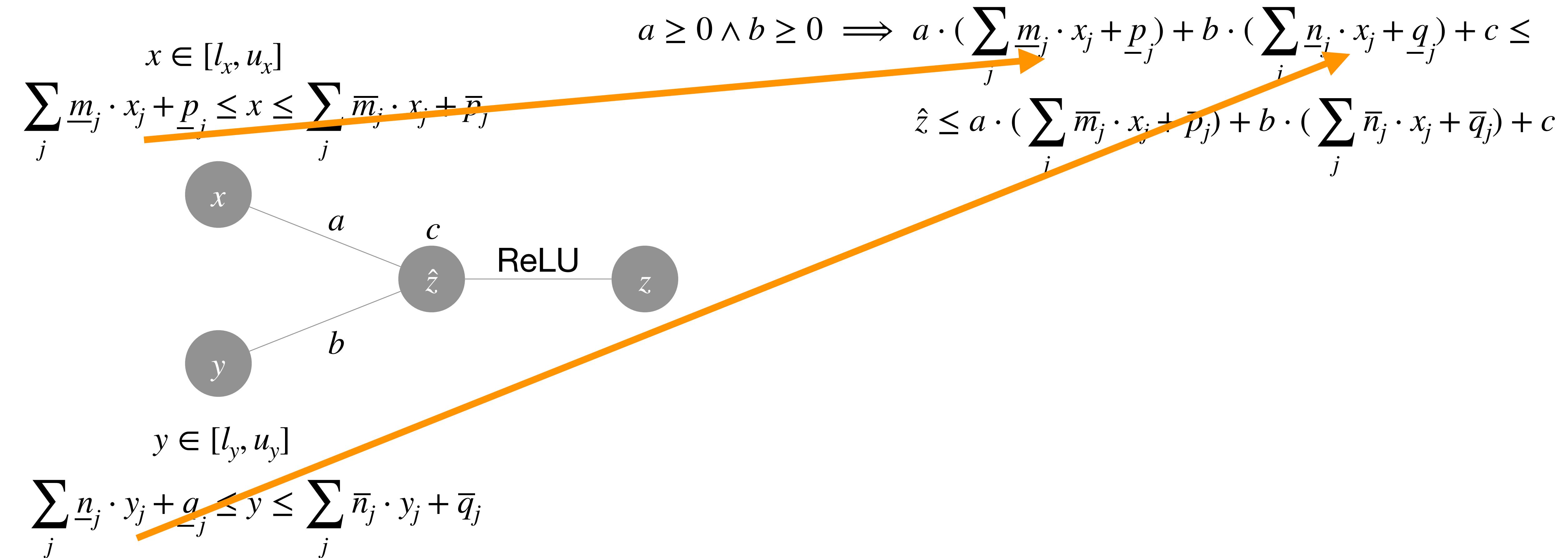
$$y \in [l_y, u_y]$$

$$\sum_j \underline{n}_j \cdot y_j + \underline{q}_j \leq y \leq \sum_j \bar{n}_j \cdot y_j + \bar{q}_j$$

$$a \geq 0 \wedge b \geq 0 \implies a \cdot \left( \sum_j \underline{m}_j \cdot x_j + \underline{p}_j \right) + b \cdot \left( \sum_j \bar{m}_j \cdot x_j + \bar{p}_j \right) + c \leq \hat{z} \leq a \cdot \left( \sum_j \bar{m}_j \cdot x_j + \bar{p}_j \right) + b \cdot \left( \sum_j \bar{n}_j \cdot y_j + \bar{q}_j \right) + c$$

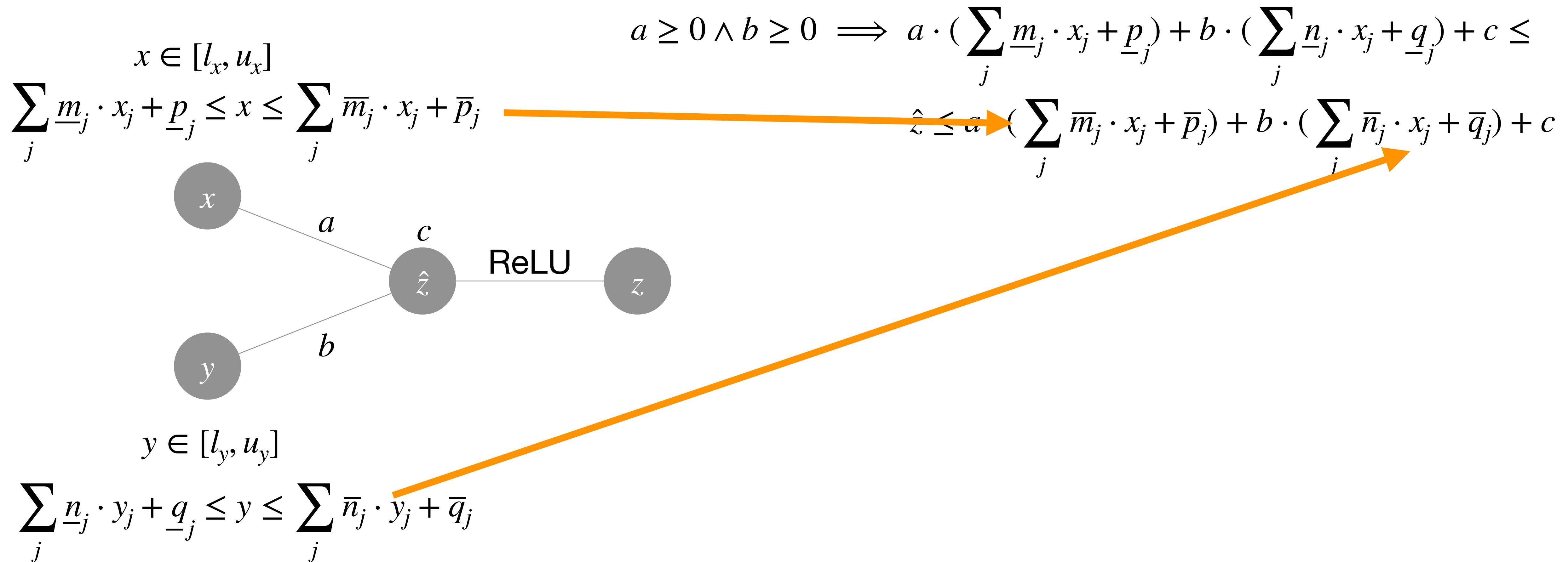
# DeepPoly

Singh et al. @ POPL 2019



# DeepPoly

Singh et al. @ POPL 2019

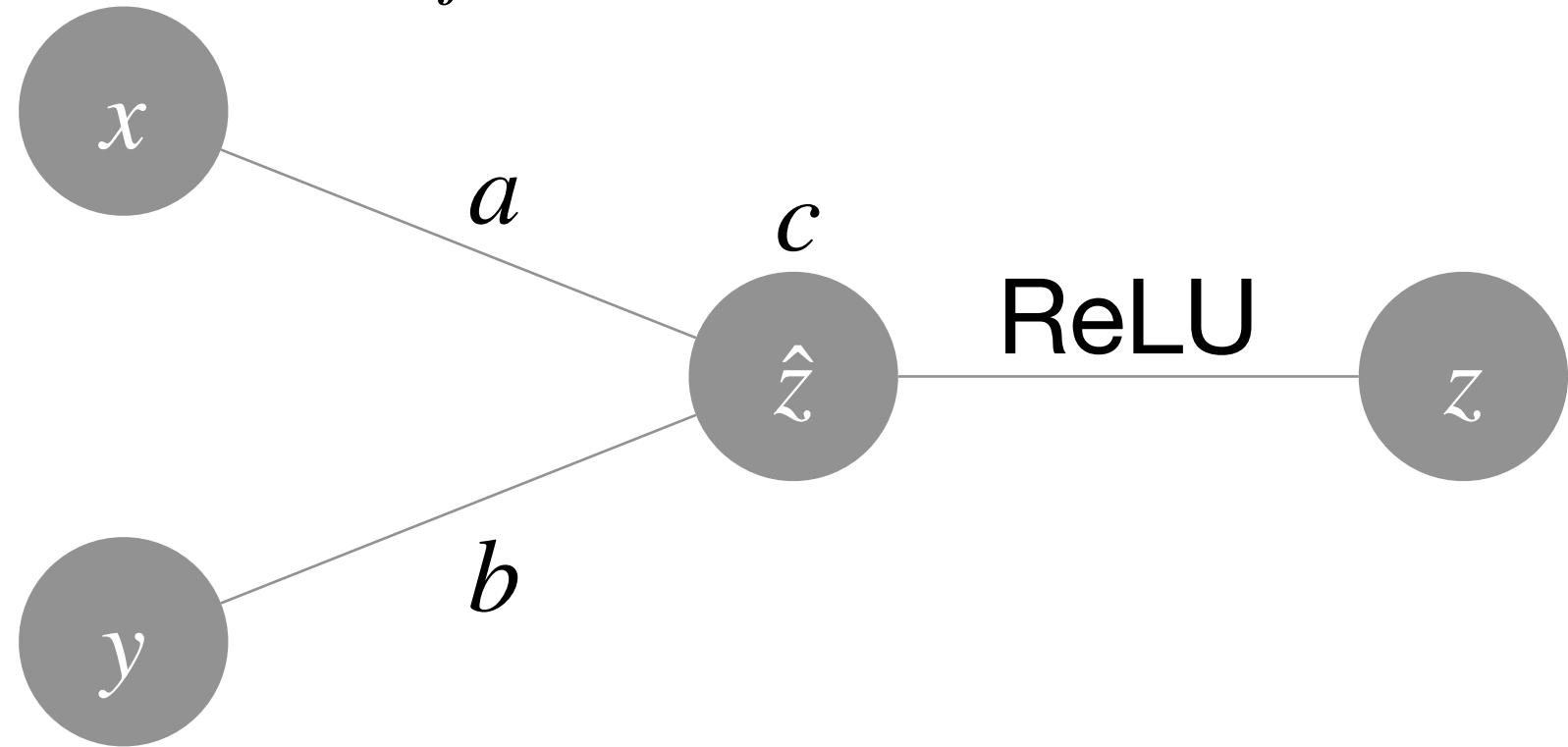


# DeepPoly

Singh et al. @ POPL 2019

$$x \in [l_x, u_x]$$

$$\sum_j \underline{m}_j \cdot x_j + \underline{p}_j \leq x \leq \sum_j \bar{m}_j \cdot x_j + \bar{p}_j$$



$$y \in [l_y, u_y]$$

$$\sum_j \underline{n}_j \cdot y_j + \underline{q}_j \leq y \leq \sum_j \bar{n}_j \cdot y_j + \bar{q}_j$$

$$a \geq 0 \wedge b \geq 0 \implies a \cdot \left( \sum_j \underline{m}_j \cdot x_j + \underline{p}_j \right) + b \cdot \left( \sum_j \bar{n}_j \cdot x_j + \bar{q}_j \right) + c \leq$$

$$\hat{z} \leq a \cdot \left( \sum_j \bar{m}_j \cdot x_j + \bar{p}_j \right) + b \cdot \left( \sum_j \bar{n}_j \cdot x_j + \bar{q}_j \right) + c$$

$$\iff \sum_j (a \cdot \underline{m}_j + b \cdot \bar{n}_j) \cdot x_j + a \cdot \underline{p}_j + b \cdot \bar{q}_j + c \leq$$

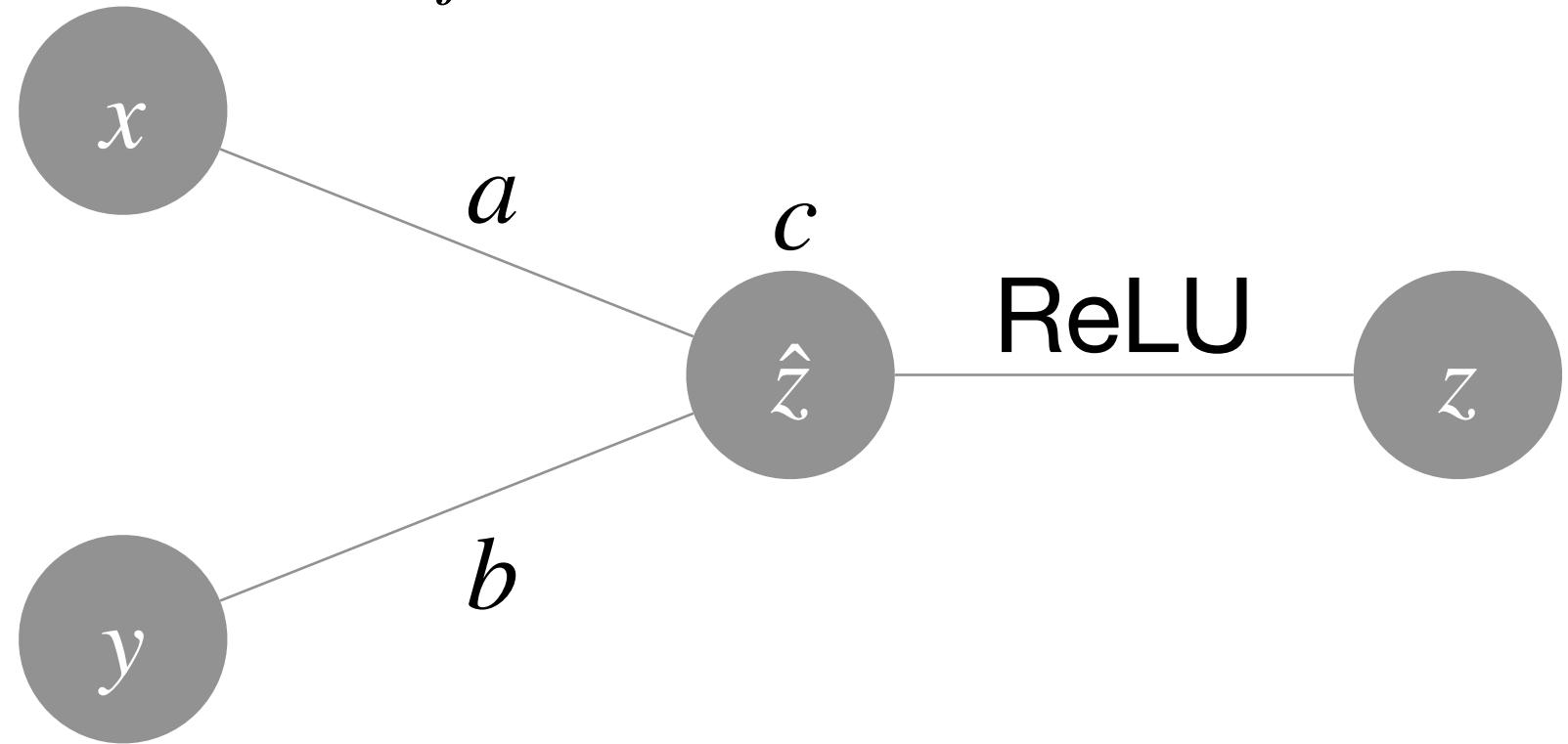
$$\hat{z} \leq \sum_j (a \cdot \bar{m}_j + b \cdot \bar{n}_j) \cdot x_j + a \cdot \bar{p}_j + b \cdot \bar{q}_j + c$$

# DeepPoly

Singh et al. @ POPL 2019

$$x \in [l_x, u_x]$$

$$\sum_j \underline{m}_j \cdot x_j + \underline{p}_j \leq x \leq \sum_j \bar{m}_j \cdot x_j + \bar{p}_j$$



$$y \in [l_y, u_y]$$

$$\sum_j \underline{n}_j \cdot y_j + \underline{q}_j \leq y \leq \sum_j \bar{n}_j \cdot y_j + \bar{q}_j$$

$$a \geq 0 \wedge b \geq 0 \implies a \cdot \left( \sum_j \underline{m}_j \cdot x_j + \underline{p}_j \right) + b \cdot \left( \sum_j \bar{n}_j \cdot x_j + \bar{q}_j \right) + c \leq$$

$$\hat{z} \leq a \cdot \left( \sum_j \bar{m}_j \cdot x_j + \bar{p}_j \right) + b \cdot \left( \sum_j \bar{n}_j \cdot x_j + \bar{q}_j \right) + c$$

$$\iff \sum_j (a \cdot \underline{m}_j + b \cdot \bar{n}_j) \cdot x_j + a \cdot \underline{p}_j + b \cdot \bar{q}_j + c \leq$$

$$\hat{z} \leq \sum_j (a \cdot \bar{m}_j + b \cdot \bar{n}_j) \cdot x_j + a \cdot \bar{p}_j + b \cdot \bar{q}_j + c$$

$$\iff \dots \in [l_{\hat{z}}, u_{\hat{z}}]$$

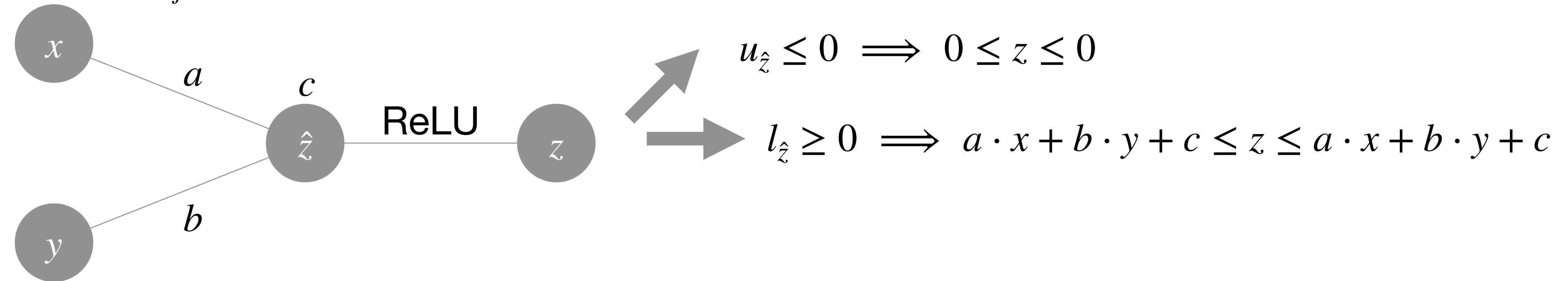
# DeepPoly

Singh et al. @ POPL 2019

$$a \cdot x + b \cdot y + c \leq \hat{z} \leq a \cdot x + b \cdot y + c$$

$$x \in [l_x, u_x]$$

$$\sum_j \underline{m}_j \cdot x_j + \underline{p}_j \leq x \leq \sum_j \bar{m}_j \cdot x_j + \bar{p}_j$$



$$y \in [l_y, u_y]$$

$$\sum_j \underline{n}_j \cdot y_j + \underline{q}_j \leq y \leq \sum_j \bar{n}_j \cdot y_j + \bar{q}_j$$

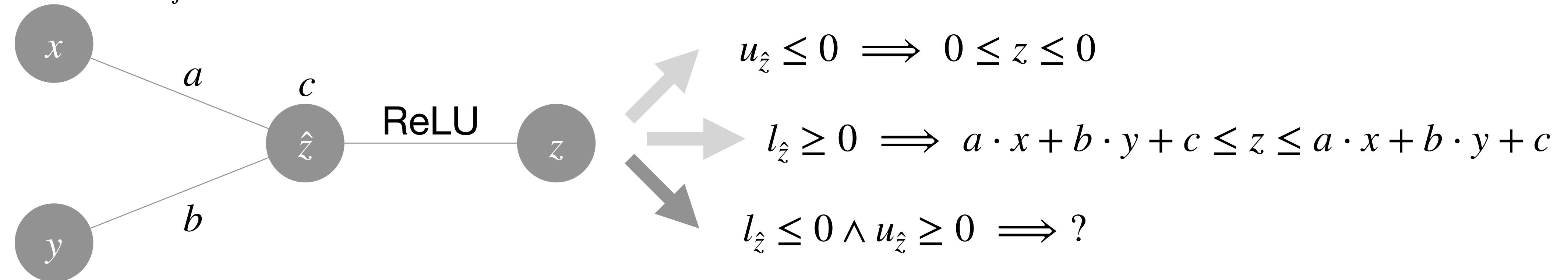
# DeepPoly

Singh et al. @ POPL 2019

$$a \cdot x + b \cdot y + c \leq \hat{z} \leq a \cdot x + b \cdot y + c$$

$$x \in [l_x, u_x]$$

$$\sum_j \underline{m}_j \cdot x_j + \underline{p}_j \leq x \leq \sum_j \bar{m}_j \cdot x_j + \bar{p}_j$$

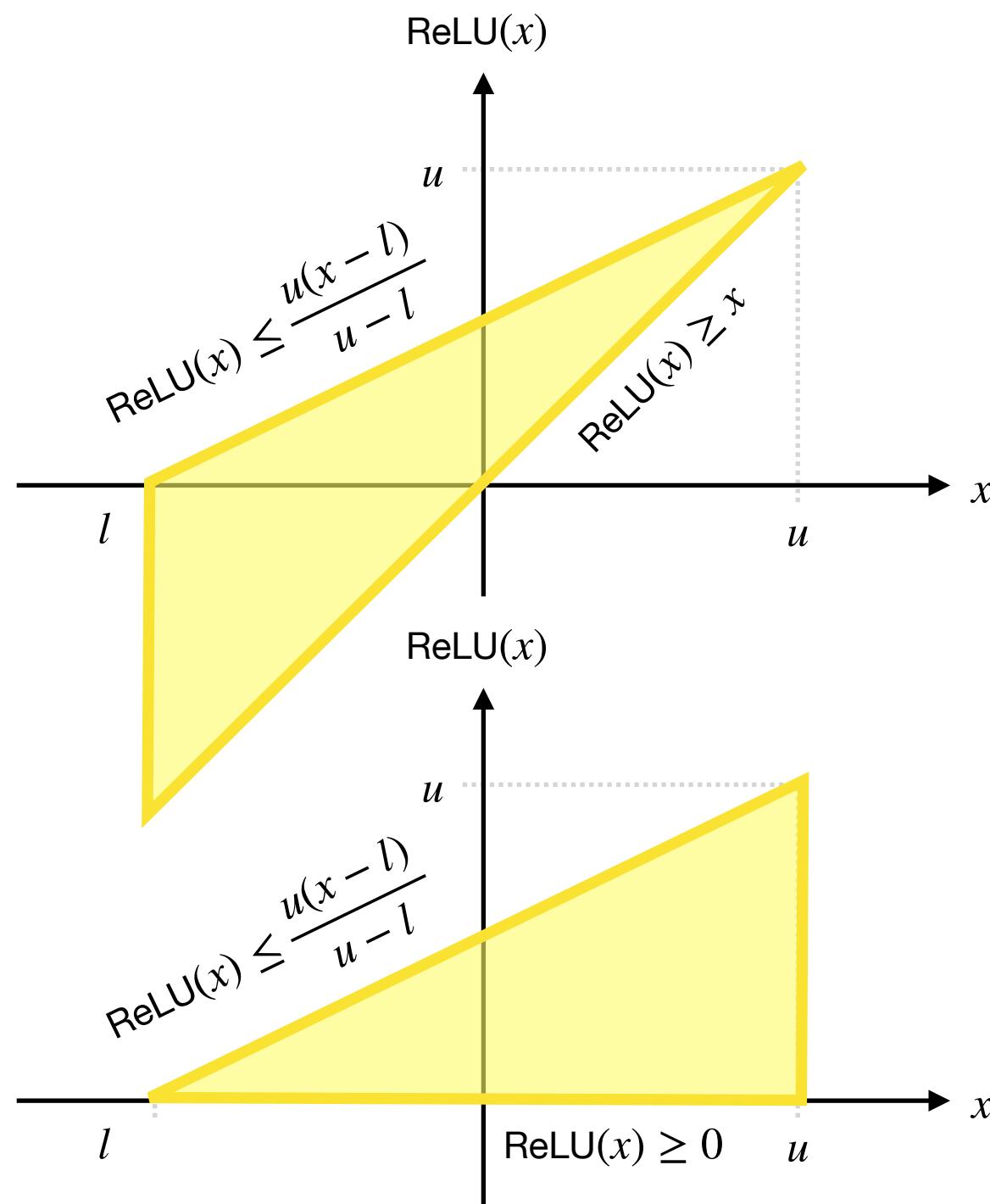


$$y \in [l_y, u_y]$$

$$\sum_j \underline{n}_j \cdot y_j + \underline{q}_j \leq y \leq \sum_j \bar{n}_j \cdot y_j + \bar{q}_j$$

# DeepPoly

Singh et al. @ POPL 2019



$$l_{\hat{z}} \leq 0 \wedge u_{\hat{z}} \geq 0 \wedge -l_{\hat{z}} \leq u_{\hat{z}} \implies \begin{cases} z \leq \frac{u_{\hat{z}}(\hat{z} - l_{\hat{z}})}{u_{\hat{z}} - l_{\hat{z}}} \\ z \geq \hat{z} \end{cases}$$

$$l_{\hat{z}} \leq 0 \wedge u_{\hat{z}} \geq 0 \wedge -l_{\hat{z}} > u_{\hat{z}} \implies \begin{cases} z \leq \frac{u_{\hat{z}}(\hat{z} - l_{\hat{z}})}{u_{\hat{z}} - l_{\hat{z}}} \\ z \geq 0 \end{cases}$$

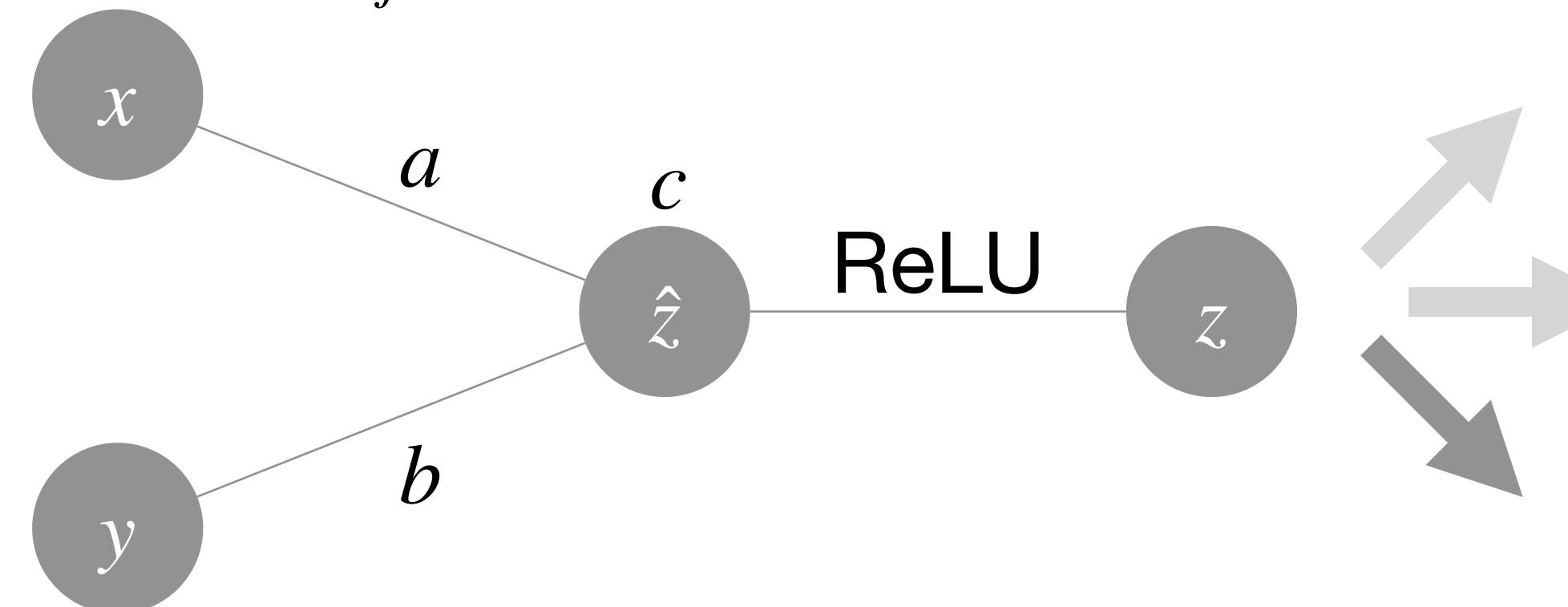
# DeepPoly

Singh et al. @ POPL 2019

$$a \cdot x + b \cdot y + c \leq \hat{z} \leq a \cdot x + b \cdot y + c$$

$$x \in [l_x, u_x]$$

$$\sum_j \underline{m}_j \cdot x_j + \underline{p}_j \leq x \leq \sum_j \bar{m}_j \cdot x_j + \bar{p}_j$$



$$y \in [l_y, u_y]$$

$$\sum_j \underline{n}_j \cdot y_j + \underline{q}_j \leq y \leq \sum_j \bar{n}_j \cdot y_j + \bar{q}_j$$

$$u_{\hat{z}} \leq 0 \implies 0 \leq z \leq 0$$

$$l_{\hat{z}} \geq 0 \implies a \cdot x + b \cdot y + c \leq z \leq a \cdot x + b \cdot y + c$$

$$l_{\hat{z}} \leq 0 \wedge u_{\hat{z}} \geq 0 \implies \begin{cases} \hat{z} \leq z \leq \frac{u_{\hat{z}}(\hat{z} - l_{\hat{z}})}{u_{\hat{z}} - l_{\hat{z}}} & -l_{\hat{z}} \leq u_{\hat{z}} \\ 0 \leq z \leq \frac{u_{\hat{z}}(\hat{z} - l_{\hat{z}})}{u_{\hat{z}} - l_{\hat{z}}} & -l_{\hat{z}} > u_{\hat{z}} \end{cases}$$

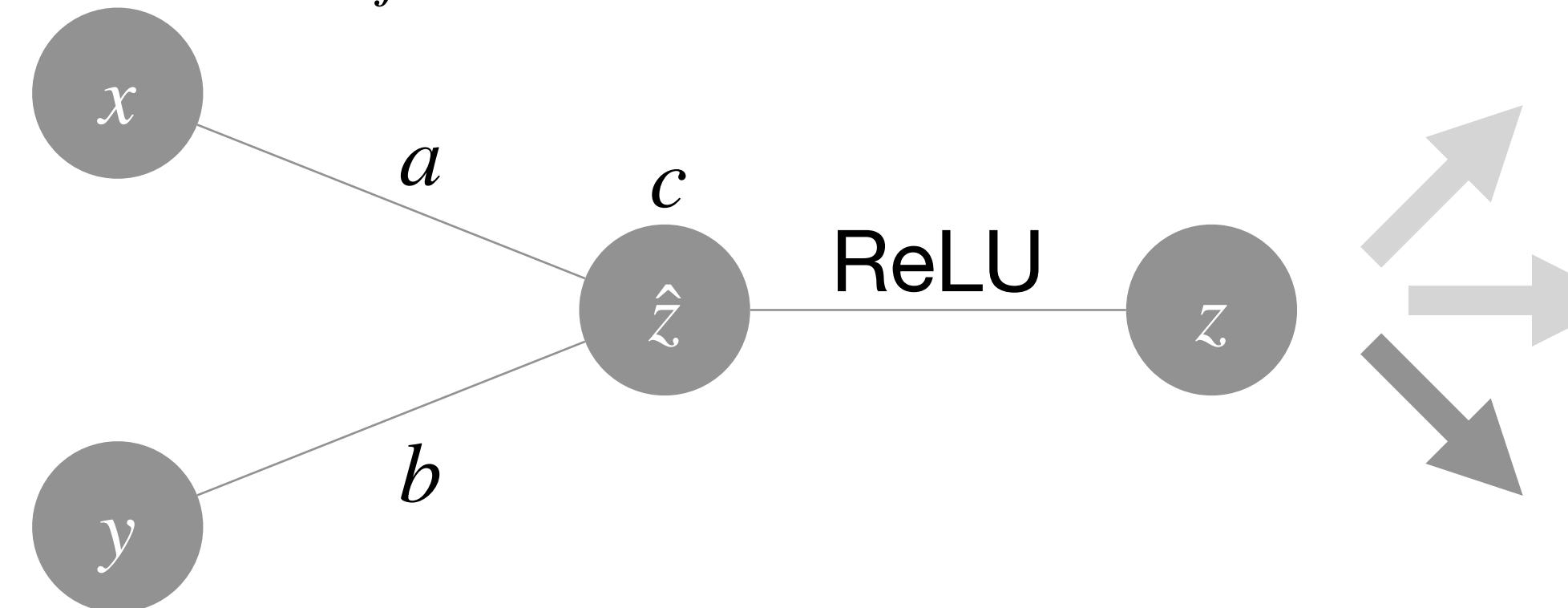
# DeepPoly

Singh et al. @ POPL 2019

$$a \cdot x + b \cdot y + c \leq \hat{z} \leq a \cdot x + b \cdot y + c$$

$$x \in [l_x, u_x]$$

$$\sum_j \underline{m}_j \cdot x_j + \underline{p}_j \leq x \leq \sum_j \bar{m}_j \cdot x_j + \bar{p}_j$$



$$y \in [l_y, u_y]$$

$$\sum_j \underline{n}_j \cdot y_j + \underline{q}_j \leq y \leq \sum_j \bar{n}_j \cdot y_j + \bar{q}_j$$

$$u_{\hat{z}} \leq 0 \implies 0 \leq z \leq 0$$

$$l_{\hat{z}} \geq 0 \implies a \cdot x + b \cdot y + c \leq z \leq a \cdot x + b \cdot y + c$$

$$l_{\hat{z}} \leq 0 \wedge u_{\hat{z}} \geq 0 \implies \begin{cases} \hat{z} \leq z \leq \frac{u_{\hat{z}}(\hat{z} - l_{\hat{z}})}{u_{\hat{z}} - l_{\hat{z}}} & -l_{\hat{z}} \leq u_{\hat{z}} \\ 0 \leq z \leq \frac{u_{\hat{z}}(\hat{z} - l_{\hat{z}})}{u_{\hat{z}} - l_{\hat{z}}} & -l_{\hat{z}} > u_{\hat{z}} \end{cases}$$

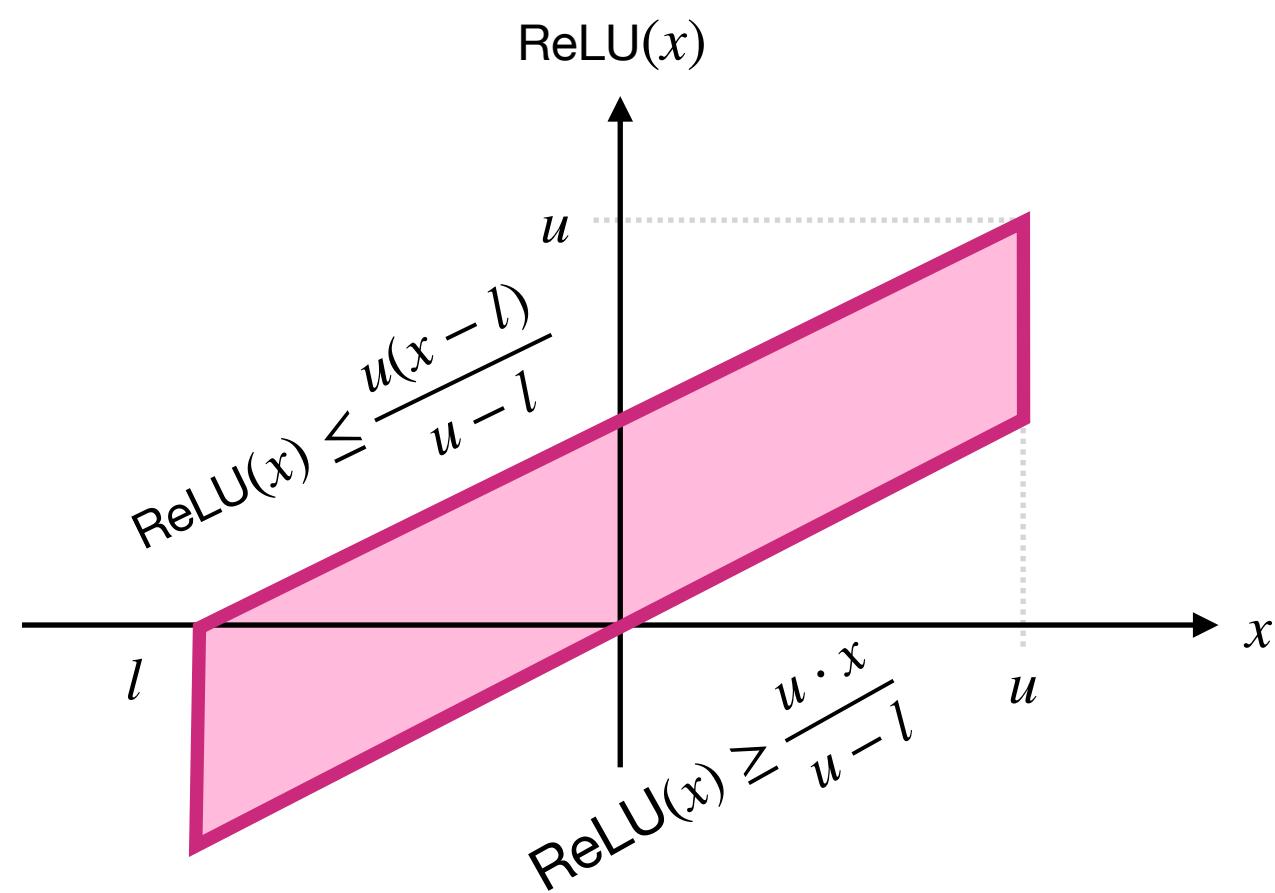
$\cap$

**Symbolic**  $\implies z \leftarrow s$

# Neurify

Wang et al. @ NeurIPS 2018

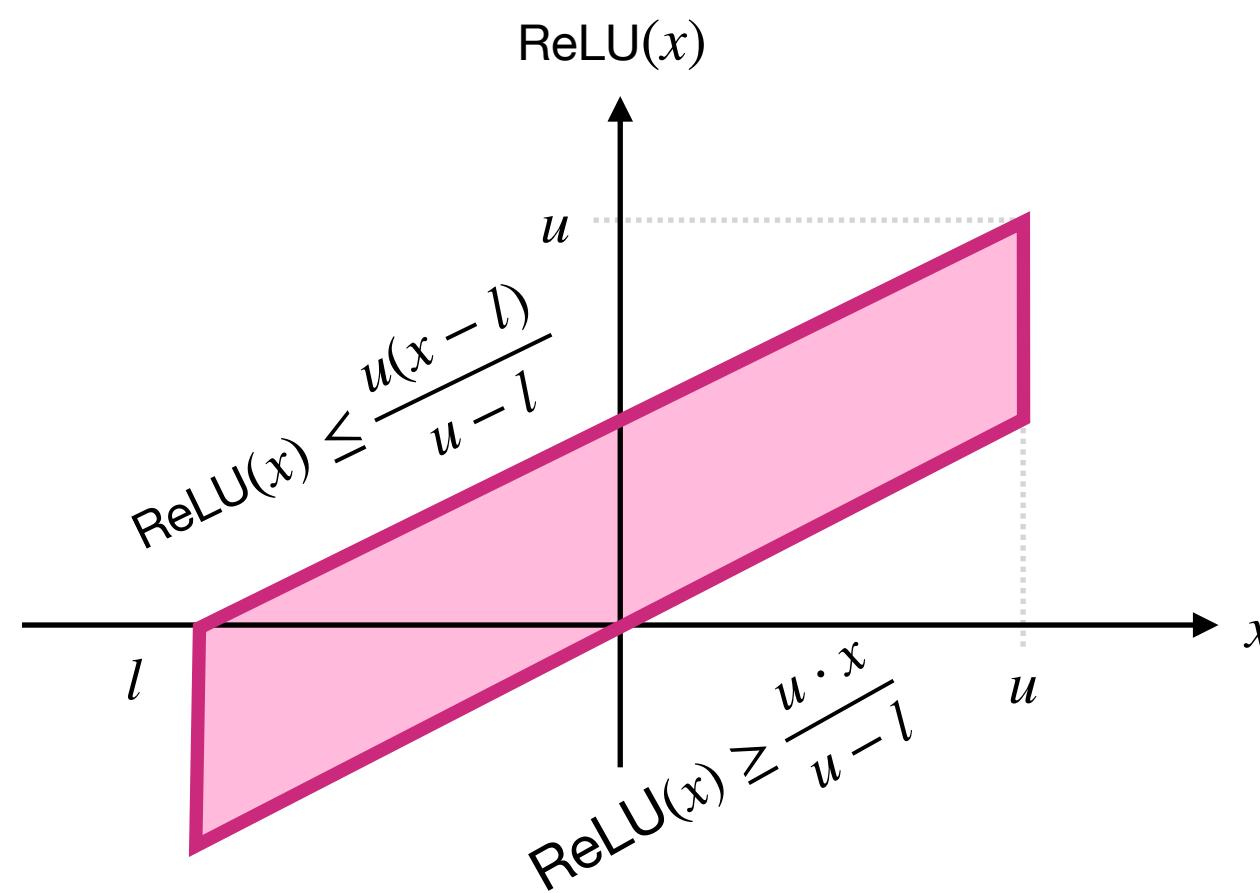
$$z \leftarrow \text{ReLU}(\hat{z}) \quad \hat{z} \in [l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$



# Neurify

Wang et al. @ NeurIPS 2018

$$z \leftarrow \text{ReLU}(\hat{z}) \quad \hat{z} \in [l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$

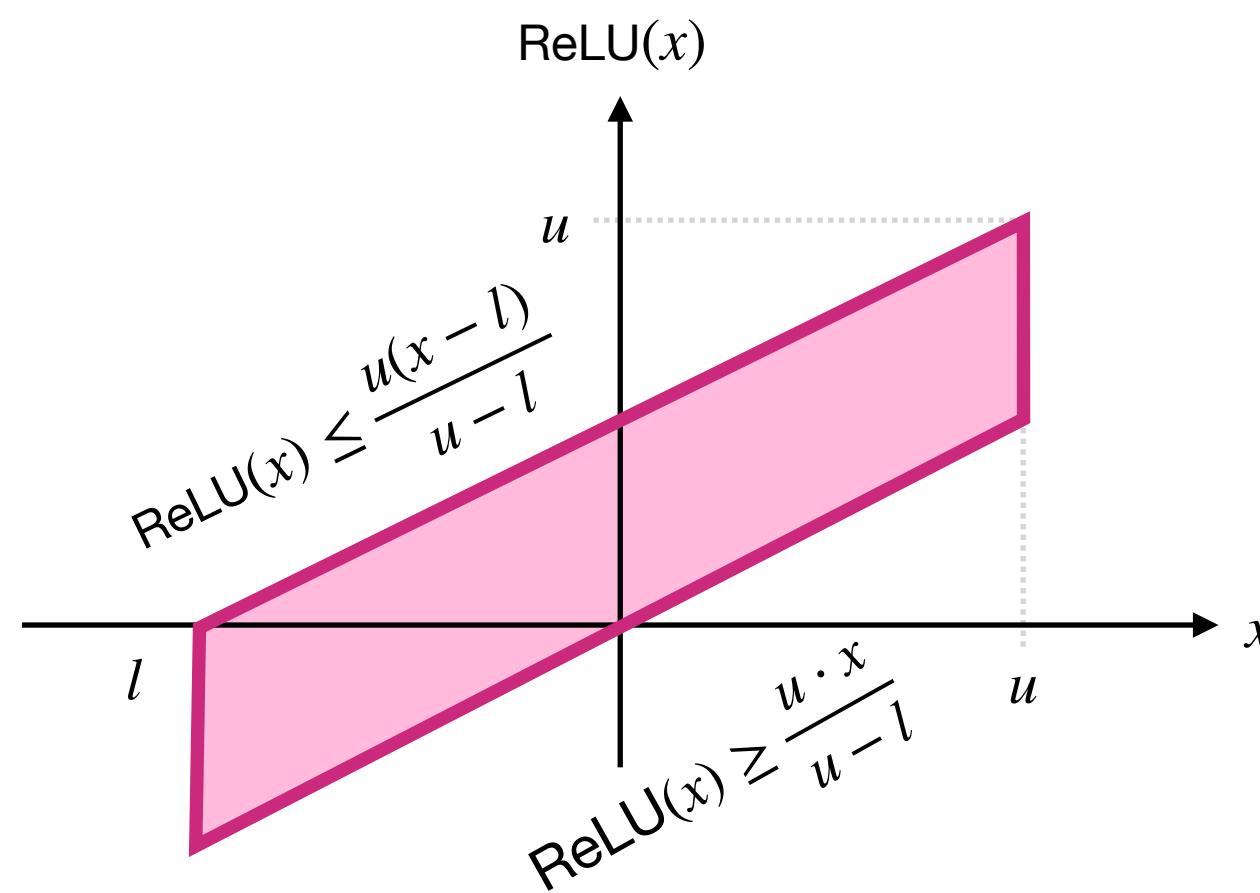


$$\underline{\text{ReLU}}(\hat{z}) \leq z \leq \overline{\text{ReLU}}(\hat{z})$$

# Neurify

Wang et al. @ NeurIPS 2018

$$z \leftarrow \text{ReLU}(\hat{z}) \quad \hat{z} \in [l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$



$$\underline{\text{ReLU}}(\hat{z}) \leq z \leq \overline{\text{ReLU}}(\hat{z})$$

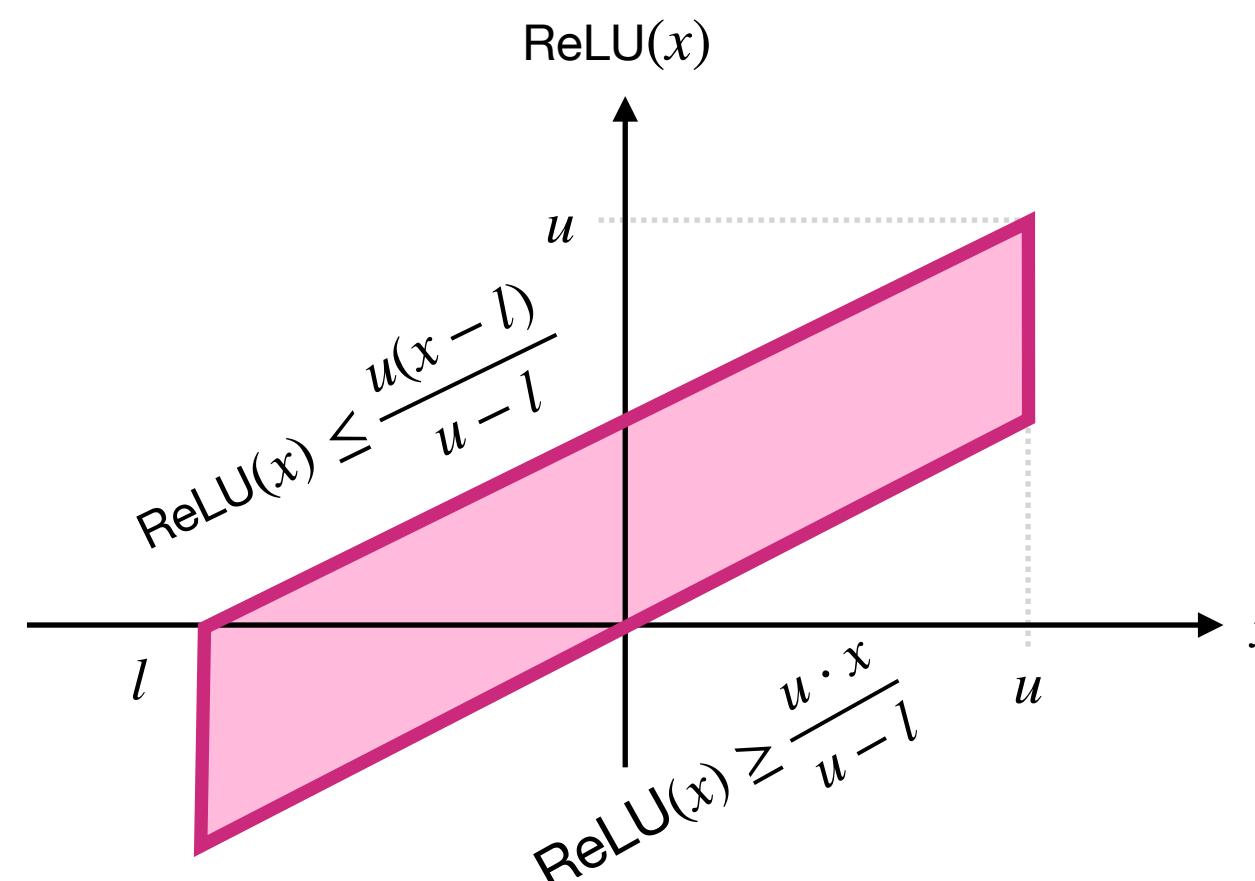
$$l_{\text{low}} \leq 0 \wedge l_{\text{up}} \geq 0 \implies \underline{\text{ReLU}}(\hat{z}) \leq z$$

$$l_{\text{low}} \leq 0 \wedge u_{\text{up}} \geq 0 \wedge$$

# Neurify

Wang et al. @ NeurIPS 2018

$$z \leftarrow \text{ReLU}(\hat{z}) \quad \hat{z} \in [l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$



$$l_{\text{low}} \leq 0 \wedge u_{\text{up}} \geq 0 \wedge$$

$$\underline{\text{ReLU}}(\hat{z}) \leq z \leq \overline{\text{ReLU}}(\hat{z})$$

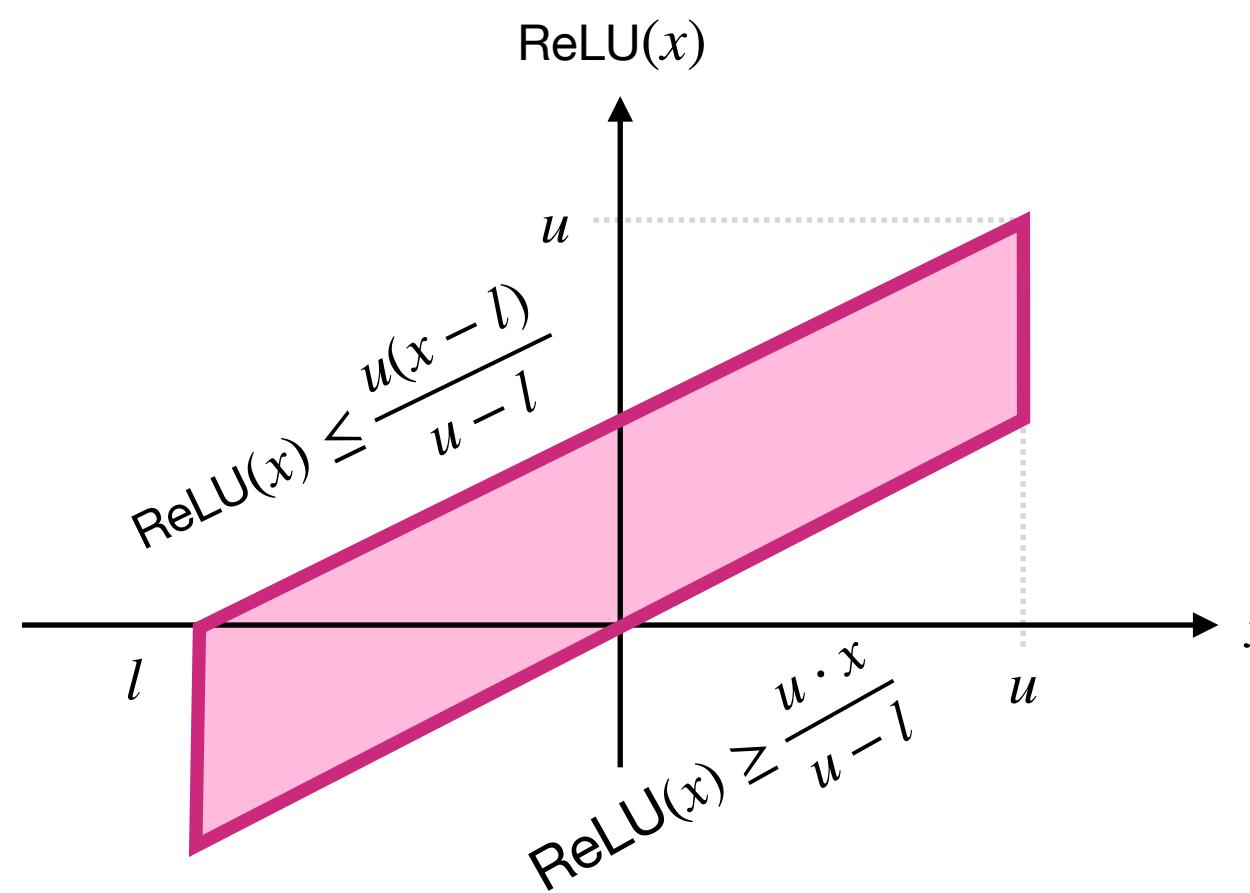
$$l_{\text{low}} \leq 0 \wedge l_{\text{up}} \geq 0 \implies \underline{\text{ReLU}}(\hat{z}) \leq z$$

$$l_{\text{low}} \leq 0 \wedge l_{\text{up}} \leq 0 \implies 0 \leq z$$

# Neurify

Wang et al. @ NeurIPS 2018

$$z \leftarrow \text{ReLU}(\hat{z}) \quad \hat{z} \in [l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$



$$l_{\text{low}} \leq 0 \wedge u_{\text{up}} \geq 0 \wedge$$

$$\underline{\text{ReLU}}(\hat{z}) \leq z \leq \overline{\text{ReLU}}(\hat{z})$$

$$l_{\text{low}} \leq 0 \wedge l_{\text{up}} \geq 0 \implies \underline{\text{ReLU}}(\hat{z}) \leq z$$

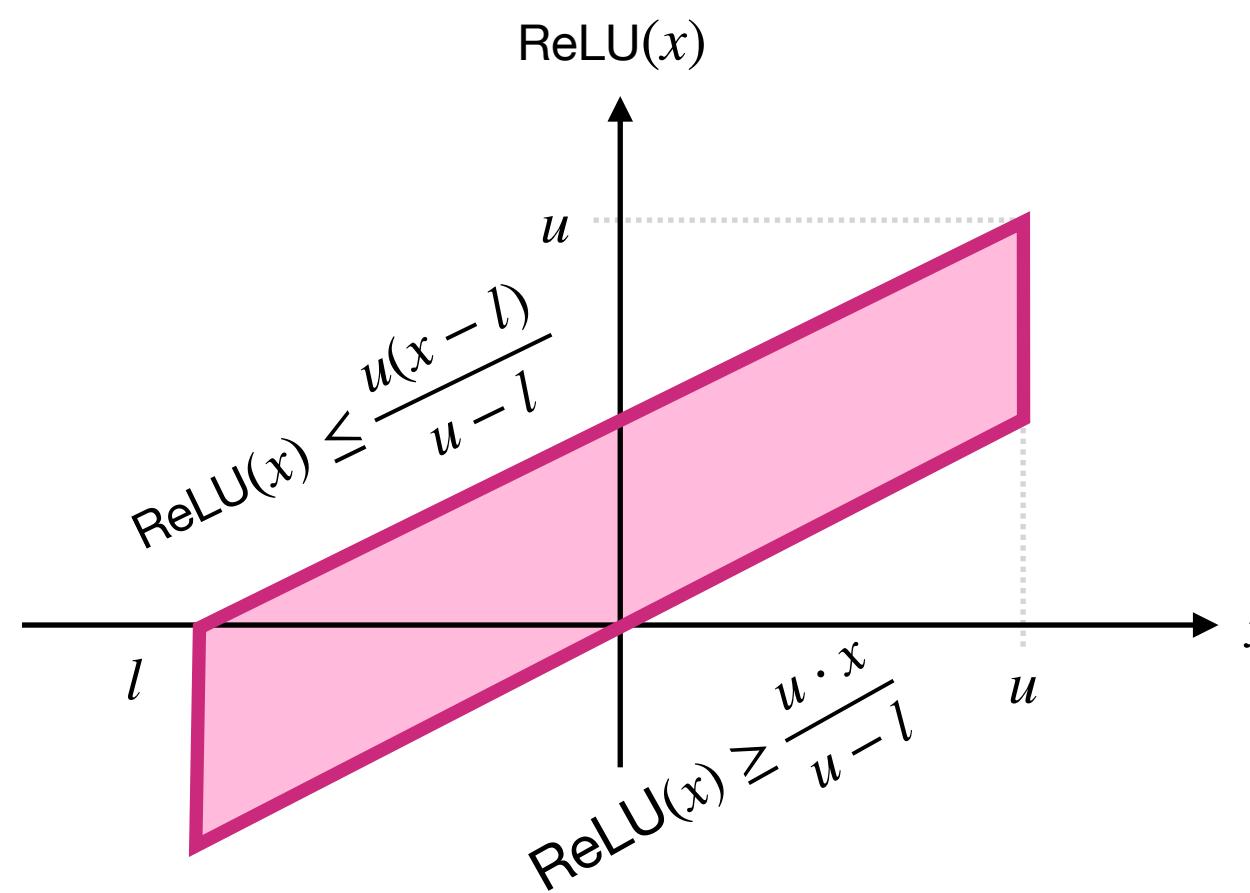
$$l_{\text{low}} \leq 0 \wedge l_{\text{up}} \leq 0 \implies 0 \leq z$$

$$u_{\text{low}} \leq 0 \wedge u_{\text{up}} \geq 0 \implies z \leq \overline{\text{ReLU}}(\hat{z})$$

# Neurify

Wang et al. @ NeurIPS 2018

$$z \leftarrow \text{ReLU}(\hat{z}) \quad \hat{z} \in [l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$



$$l_{\text{low}} \leq 0 \wedge u_{\text{up}} \geq 0 \wedge$$

$$\underline{\text{ReLU}}(\hat{z}) \leq z \leq \overline{\text{ReLU}}(\hat{z})$$

$$l_{\text{low}} \leq 0 \wedge l_{\text{up}} \geq 0 \implies \underline{\text{ReLU}}(\hat{z}) \leq z$$

$$l_{\text{low}} \leq 0 \wedge l_{\text{up}} \leq 0 \implies 0 \leq z$$

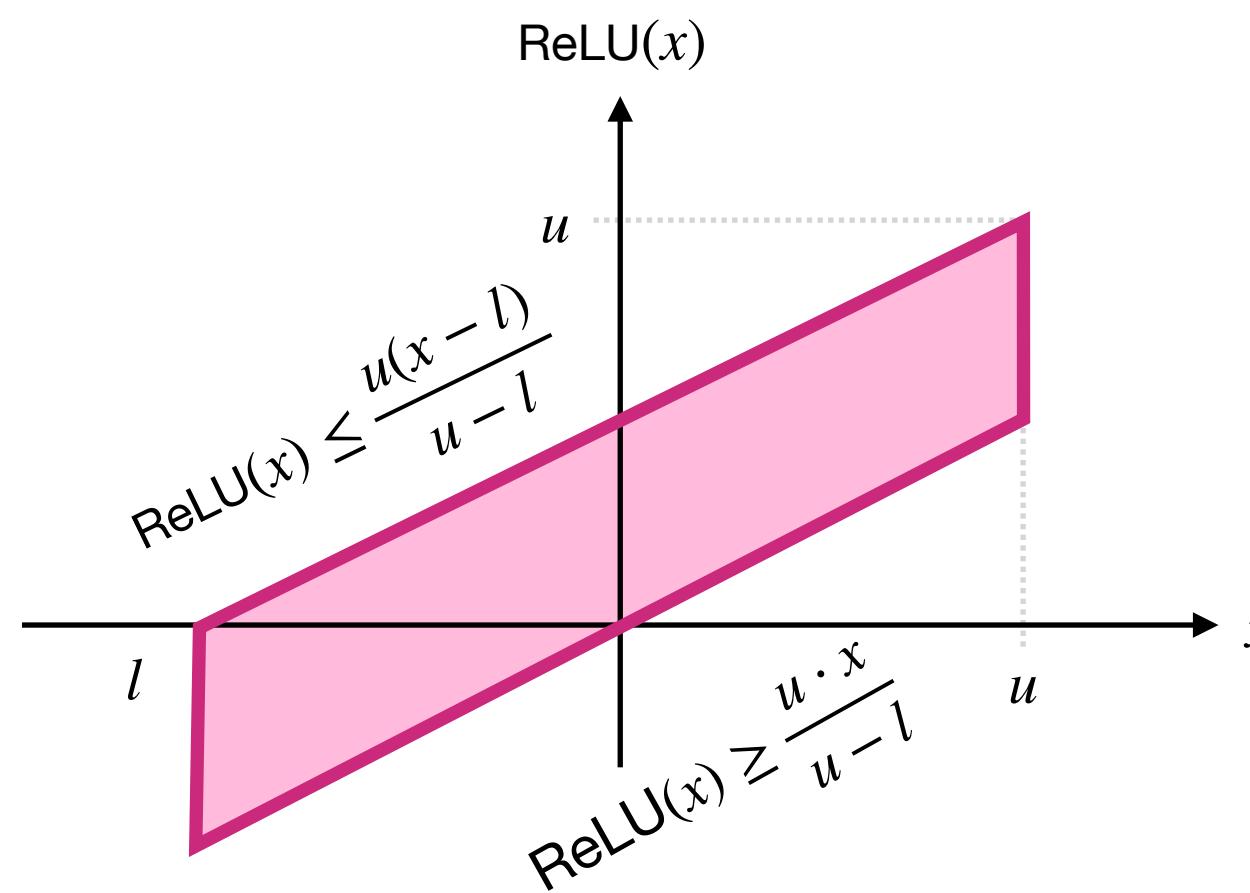
$$u_{\text{low}} \leq 0 \wedge u_{\text{up}} \geq 0 \implies z \leq \overline{\text{ReLU}}(\hat{z})$$

$$u_{\text{low}} \geq 0 \wedge u_{\text{up}} \geq 0 \implies z \leq \hat{z}$$

# Neurify

Wang et al. @ NeurIPS 2018

$$z \leftarrow \text{ReLU}(\hat{z}) \quad \hat{z} \in [l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$



$$l_{\text{low}} \leq 0 \wedge u_{\text{up}} \geq 0 \wedge$$

$$\underline{\text{ReLU}}(\hat{z}) \leq z \leq \overline{\text{ReLU}}(\hat{z})$$

$$l_{\text{low}} \leq 0 \wedge l_{\text{up}} \geq 0 \implies \underline{\text{ReLU}}(\hat{z}) \leq z$$

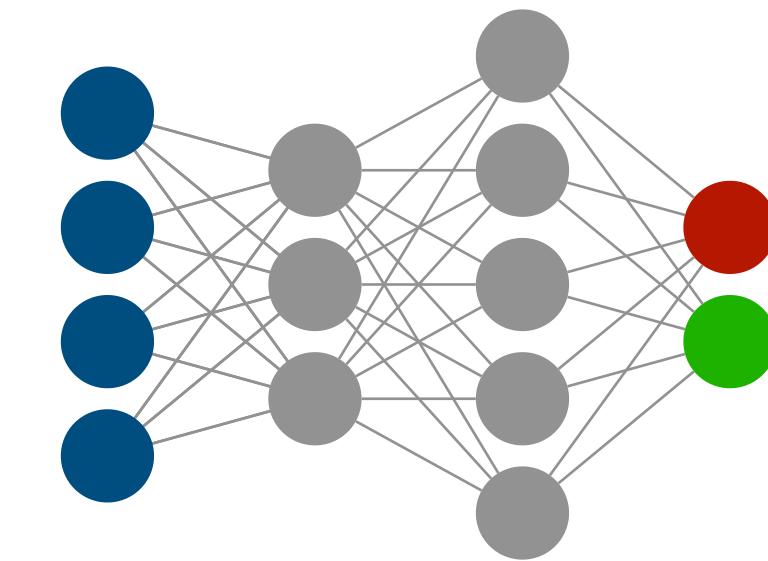
$$l_{\text{low}} \leq 0 \wedge l_{\text{up}} \leq 0 \implies 0 \leq z \quad \xleftarrow{\text{No relaxation}}$$

$$u_{\text{low}} \leq 0 \wedge u_{\text{up}} \geq 0 \implies z \leq \overline{\text{ReLU}}(\hat{z})$$

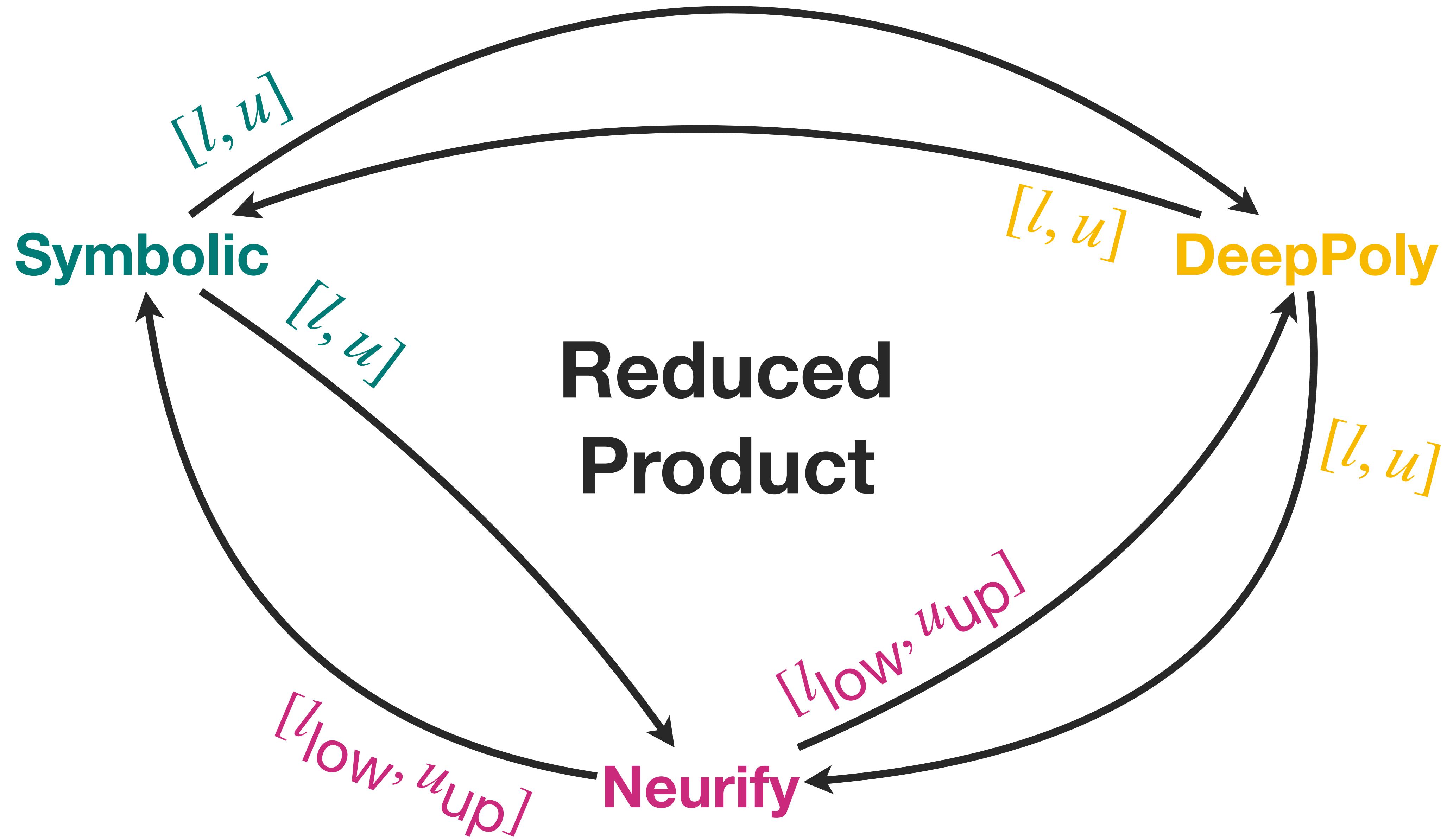
$$u_{\text{low}} \geq 0 \wedge u_{\text{up}} \geq 0 \implies z \leq \hat{z} \quad \xleftarrow{\text{No relaxation}}$$

# Precision-vs-Scalability

L	U	Symbolic	DeepPoly	Neurify
0.5	3	48.78%	49.01%	46.49%
	5	56.11%	56.15%	53.06%
0.25	3	83.63%	81.82%	81.40%
	5	91.67%	91.58%	92.33%



- 4 Hidden Layers
- 5 Neuron per Layer
- 23 inputs  $\in [0,1]$
- 2 Output classes

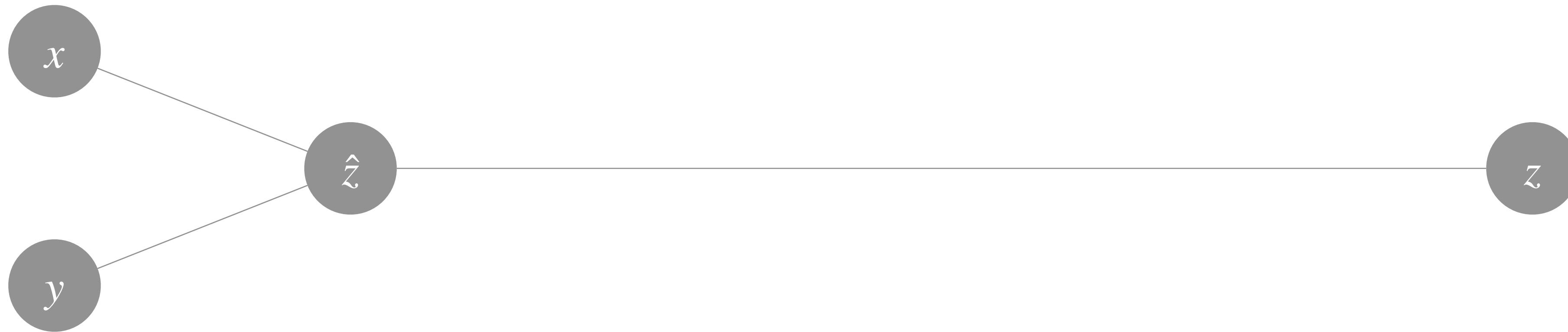


# Reduced Product

Mazzucato, Urban @ SAS 2021

$$x \in \begin{cases} [l_x, u_x] \\ [l_x, u_x] \\ [l_{x\text{low}}, l_{x\text{up}}, u_{x\text{low}}, u_{x\text{up}}] \end{cases}$$

**Symbolic**  
**DeepPoly**  
**Neurify**

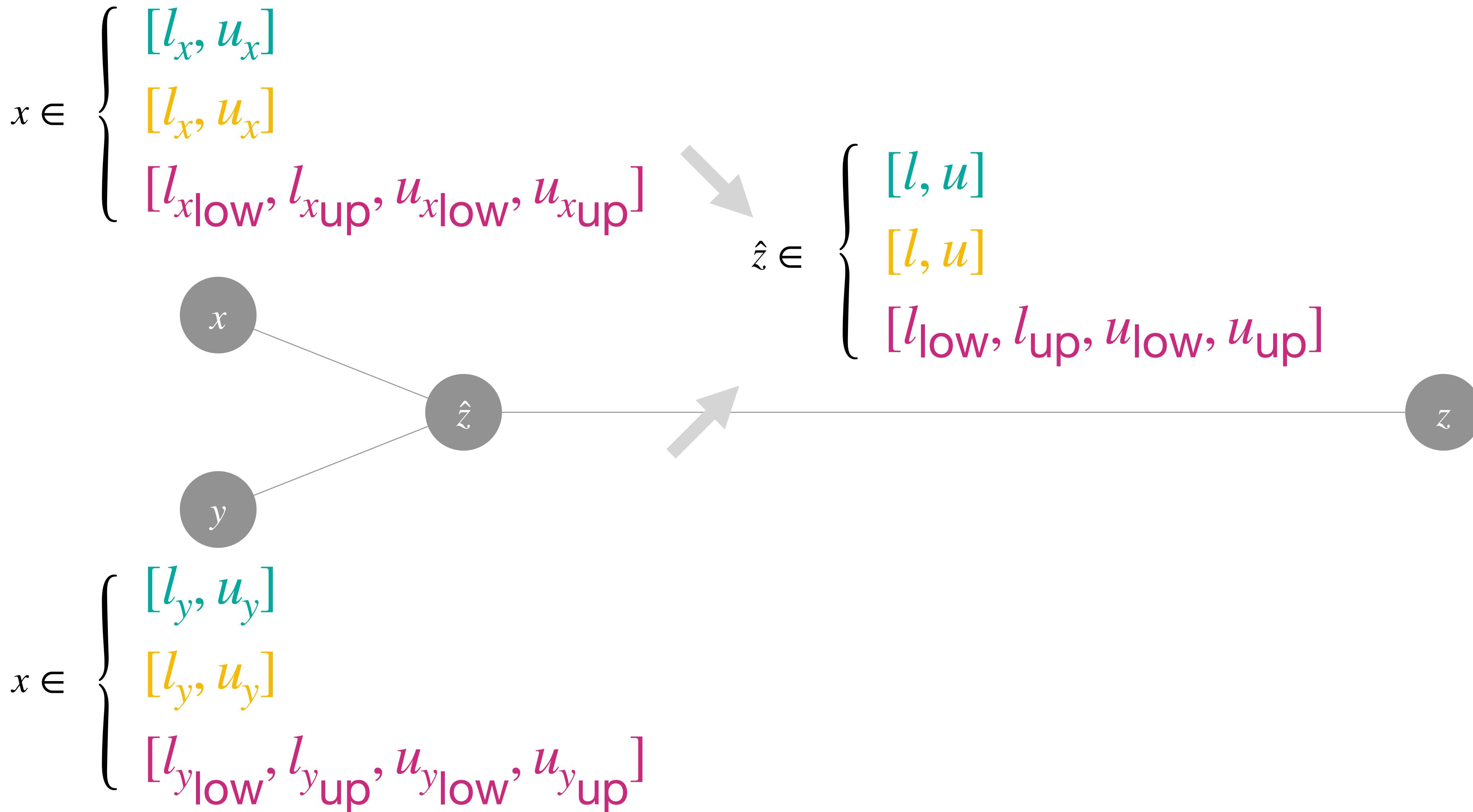


$$x \in \begin{cases} [l_y, u_y] \\ [l_y, u_y] \\ [l_{y\text{low}}, l_{y\text{up}}, u_{y\text{low}}, u_{y\text{up}}] \end{cases}$$

**Symbolic**  
**DeepPoly**  
**Neurify**

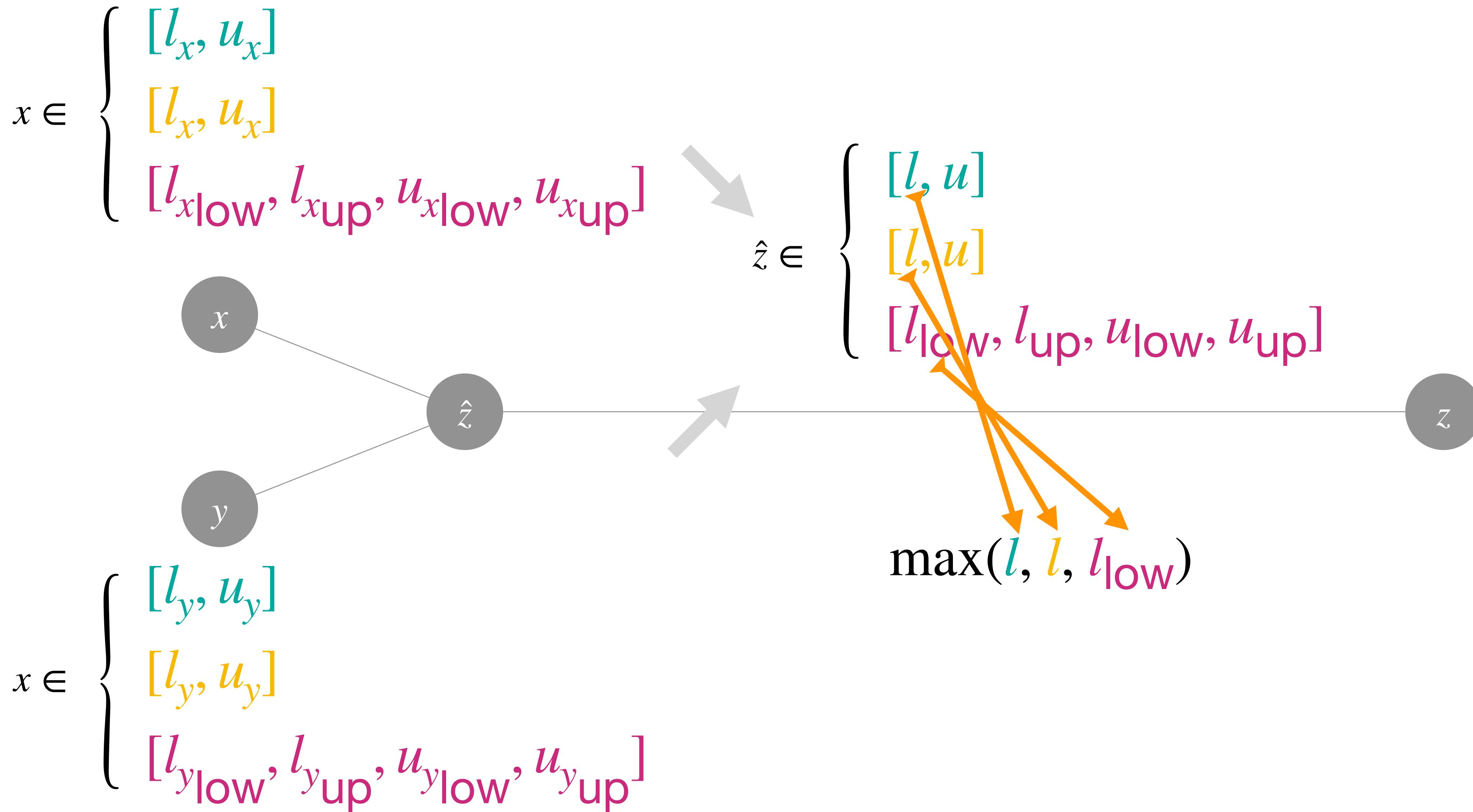
# Reduced Product

Mazzucato, Urban @ SAS 2021



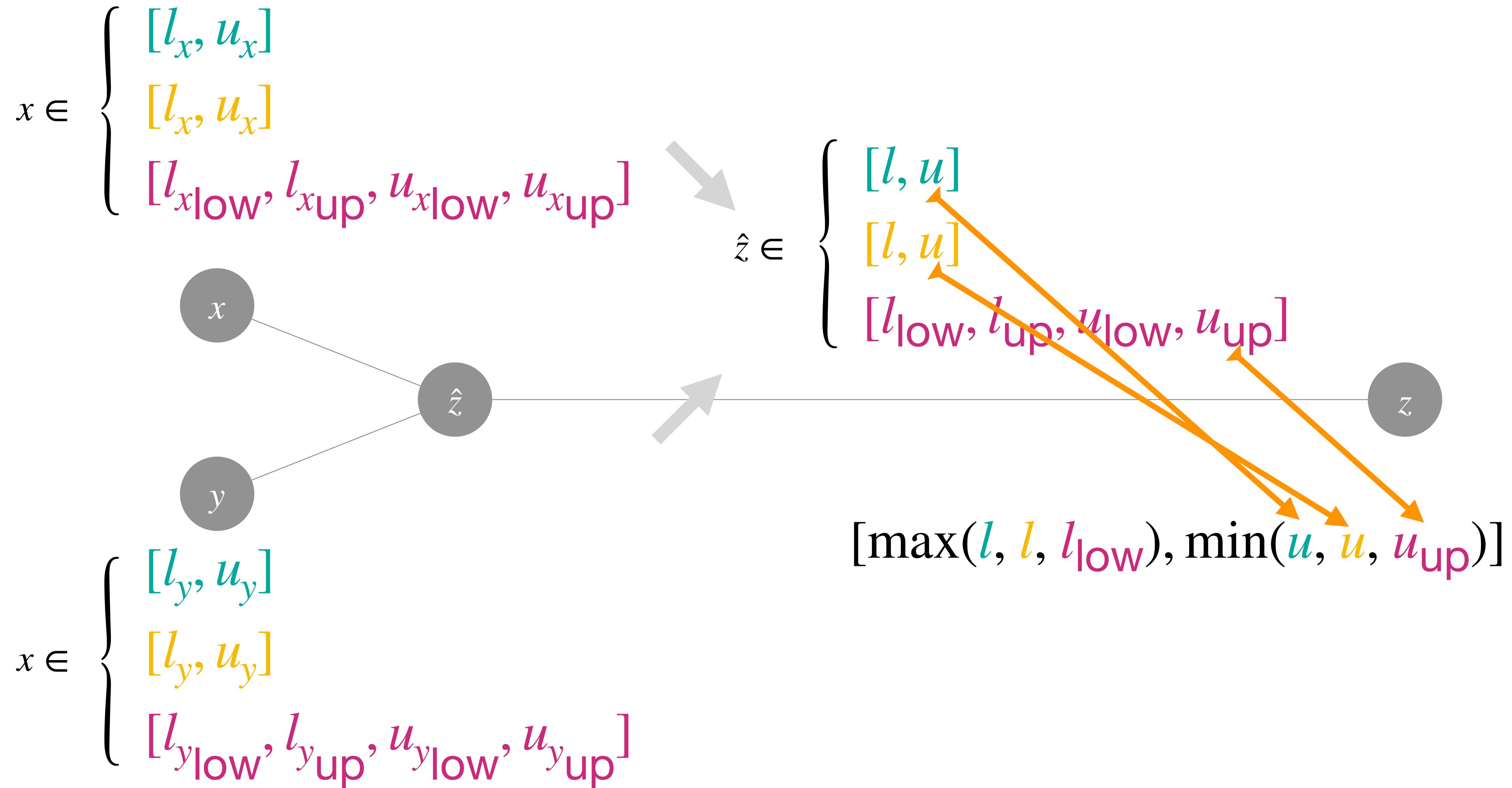
# Reduced Product

Mazzucato, Urban @ SAS 2021



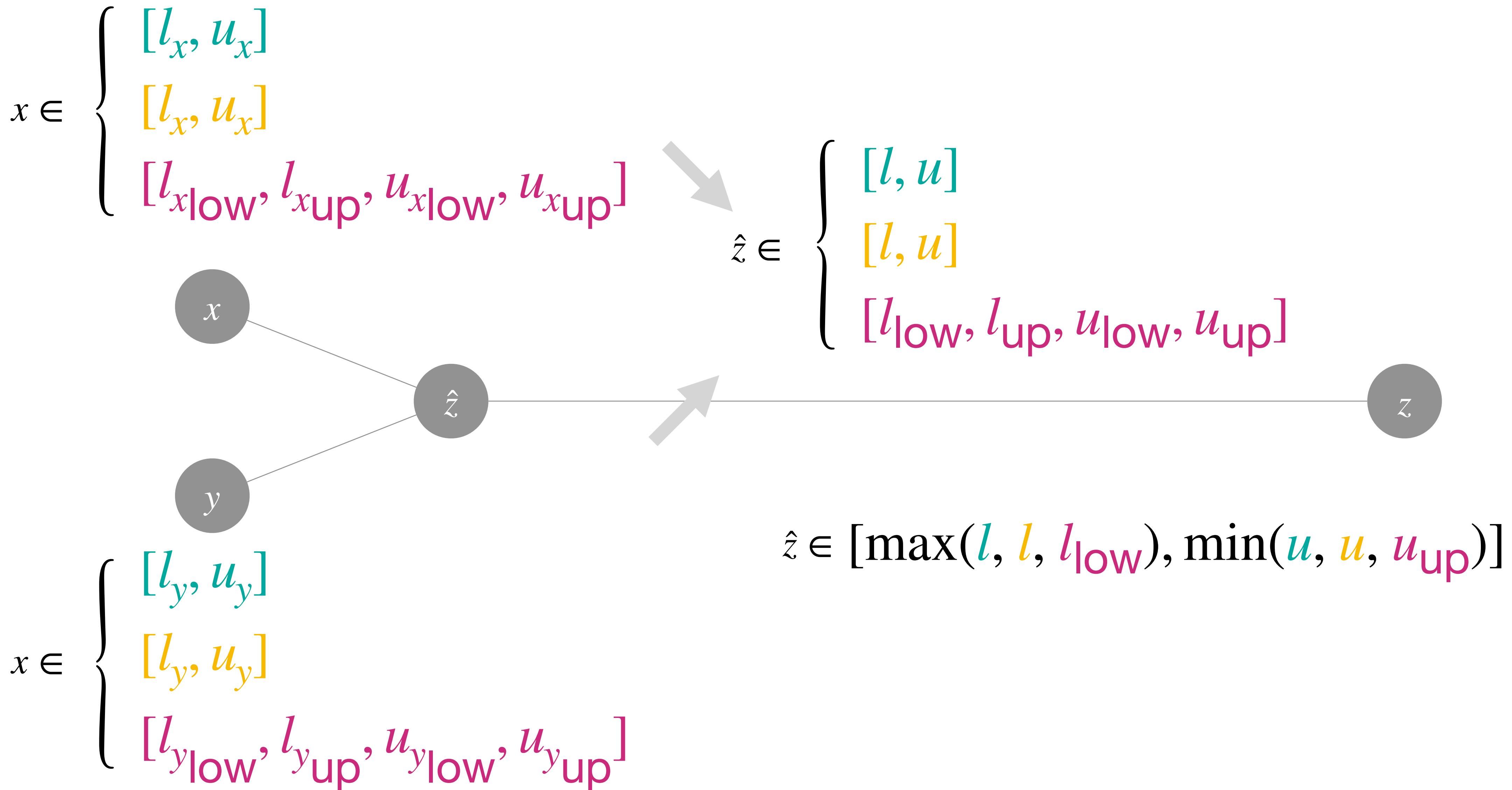
# Reduced Product

Mazzucato, Urban @ SAS 2021



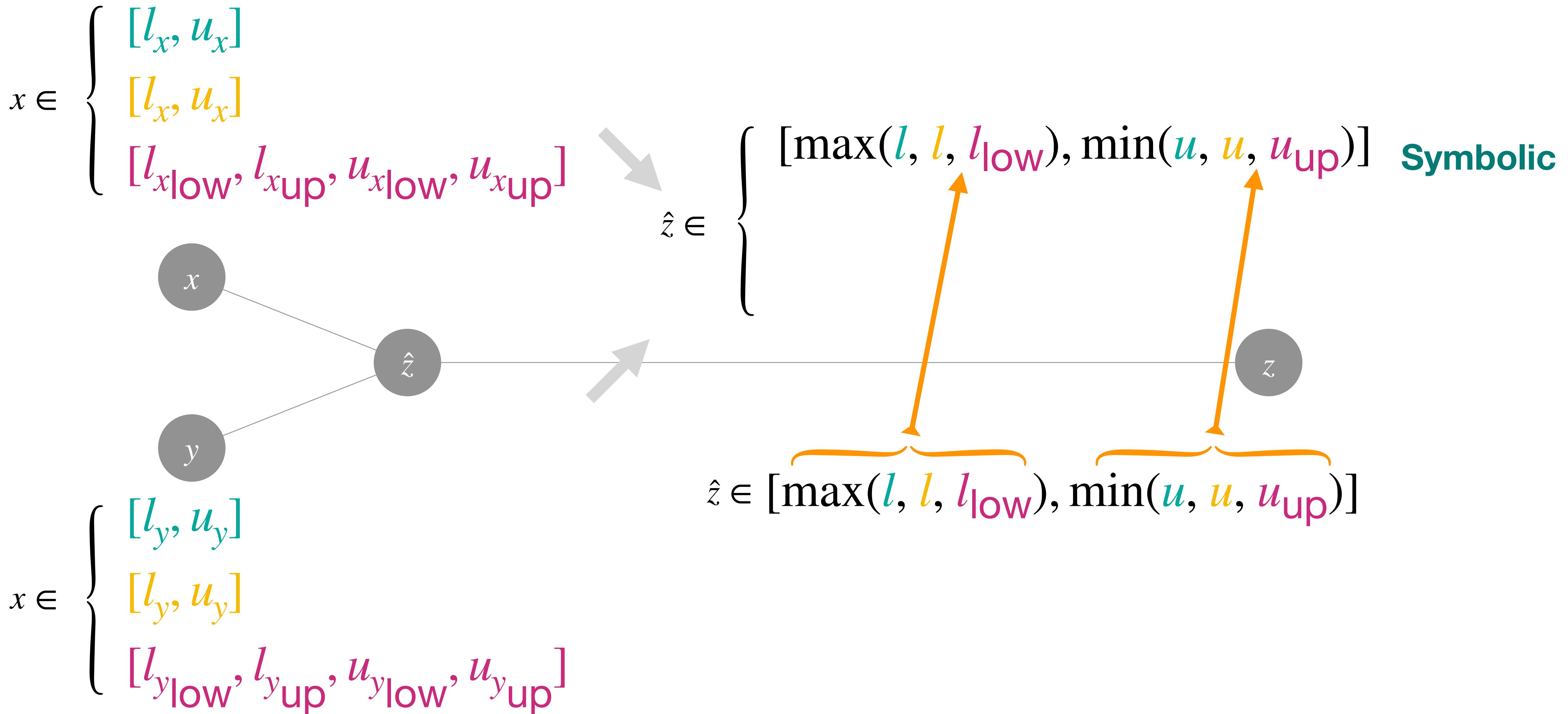
# Reduced Product

Mazzucato, Urban @ SAS 2021



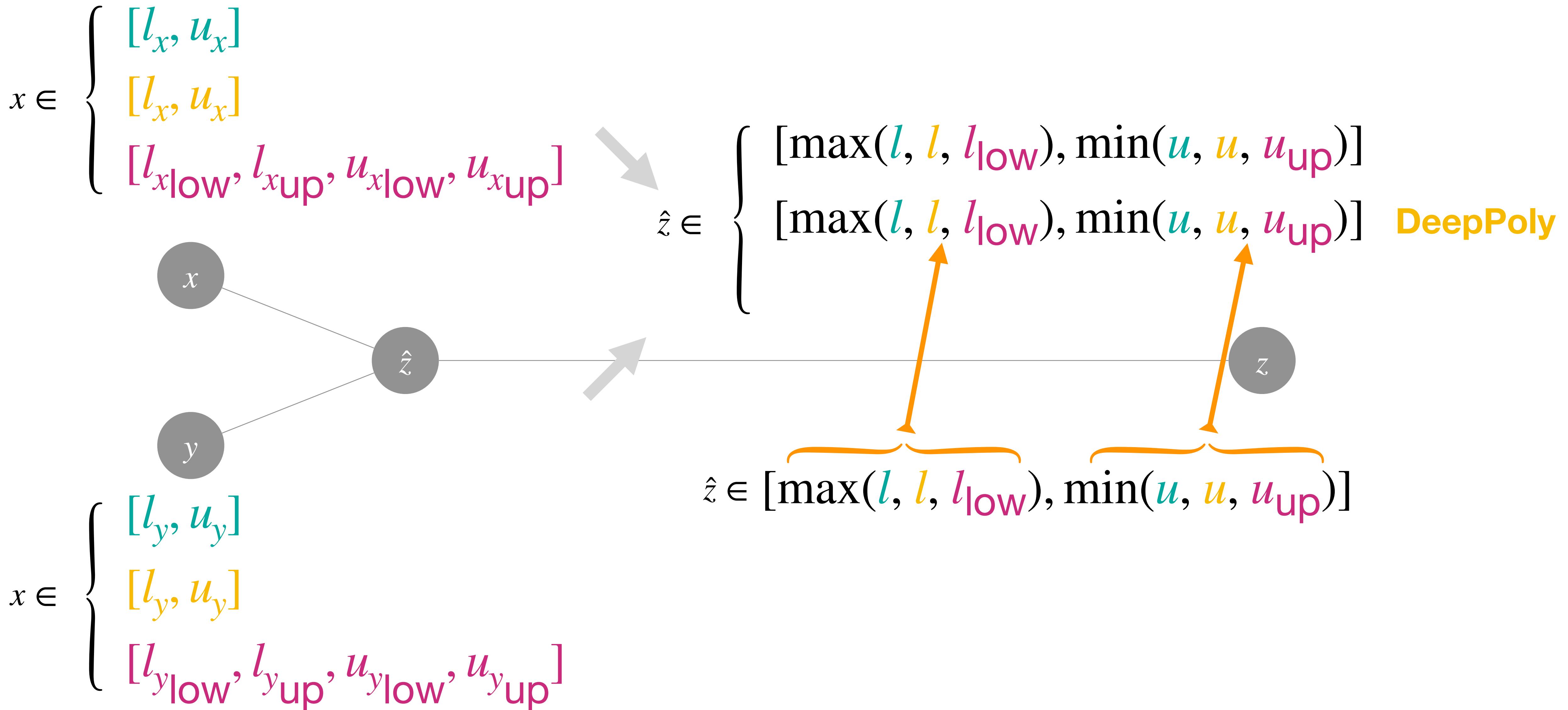
# Reduced Product

Mazzucato, Urban @ SAS 2021



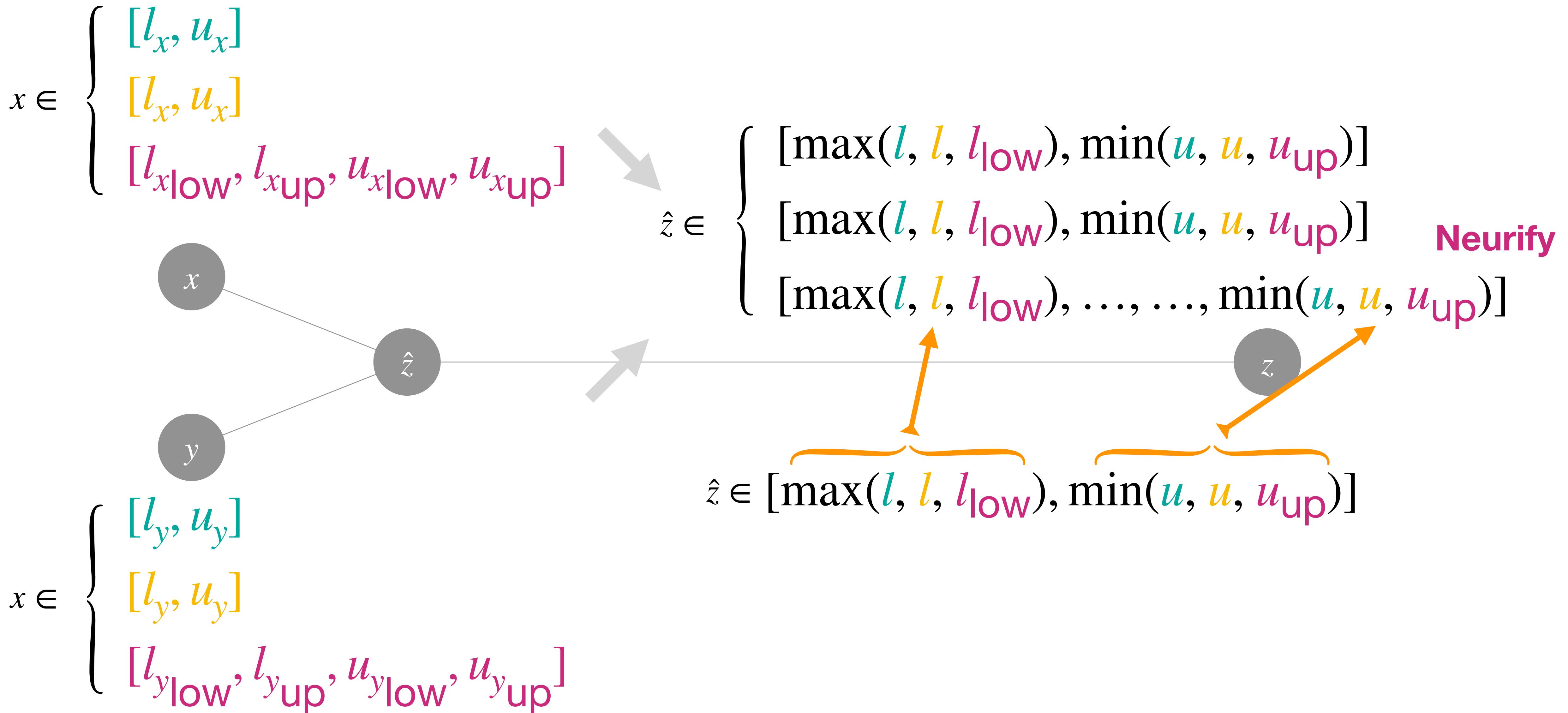
# Reduced Product

Mazzucato, Urban @ SAS 2021



# Reduced Product

Mazzucato, Urban @ SAS 2021



# Reduced Product

Mazzucato, Urban @ SAS 2021

$$\hat{z} \in [l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$

$$\hat{z} \in \left\{ \begin{array}{l} [\max(l, l, l_{\text{low}}), \min(u, u, u_{\text{up}})] \\ [\max(l, l, l_{\text{low}}), \min(u, u, u_{\text{up}})] \\ [\max(l, l, l_{\text{low}}), \dots, \dots, \min(u, u, u_{\text{up}})] \end{array} \right.$$

**Neurify**

$$\max(\max(l, l, l_{\text{low}}), l_{\text{up}})$$
$$\hat{z} \in [\max(l, l, l_{\text{low}}), \min(u, u, u_{\text{up}})]$$

# Reduced Product

Mazzucato, Urban @ SAS 2021

$$\hat{z} \in [l_{\text{low}}, l_{\text{up}}, u_{\text{low}}, u_{\text{up}}]$$

$$\hat{z} \in \left\{ \begin{array}{l} [\max(l, \underline{l}, l_{\text{low}}), \min(u, \underline{u}, u_{\text{up}})] \\ [\max(l, \underline{l}, l_{\text{low}}), \min(u, \underline{u}, u_{\text{up}})] \\ [\max(l, \underline{l}, l_{\text{low}}), \dots, \dots, \min(u, \underline{u}, u_{\text{up}})] \end{array} \right.$$

**Neurify**

$$\min(\min(u, \underline{u}, u_{\text{up}}), \underline{u}_{\text{low}})$$
$$\hat{z} \in [\max(l, \underline{l}, l_{\text{low}}), \min(u, \underline{u}, u_{\text{up}})]$$

# Reduced Product

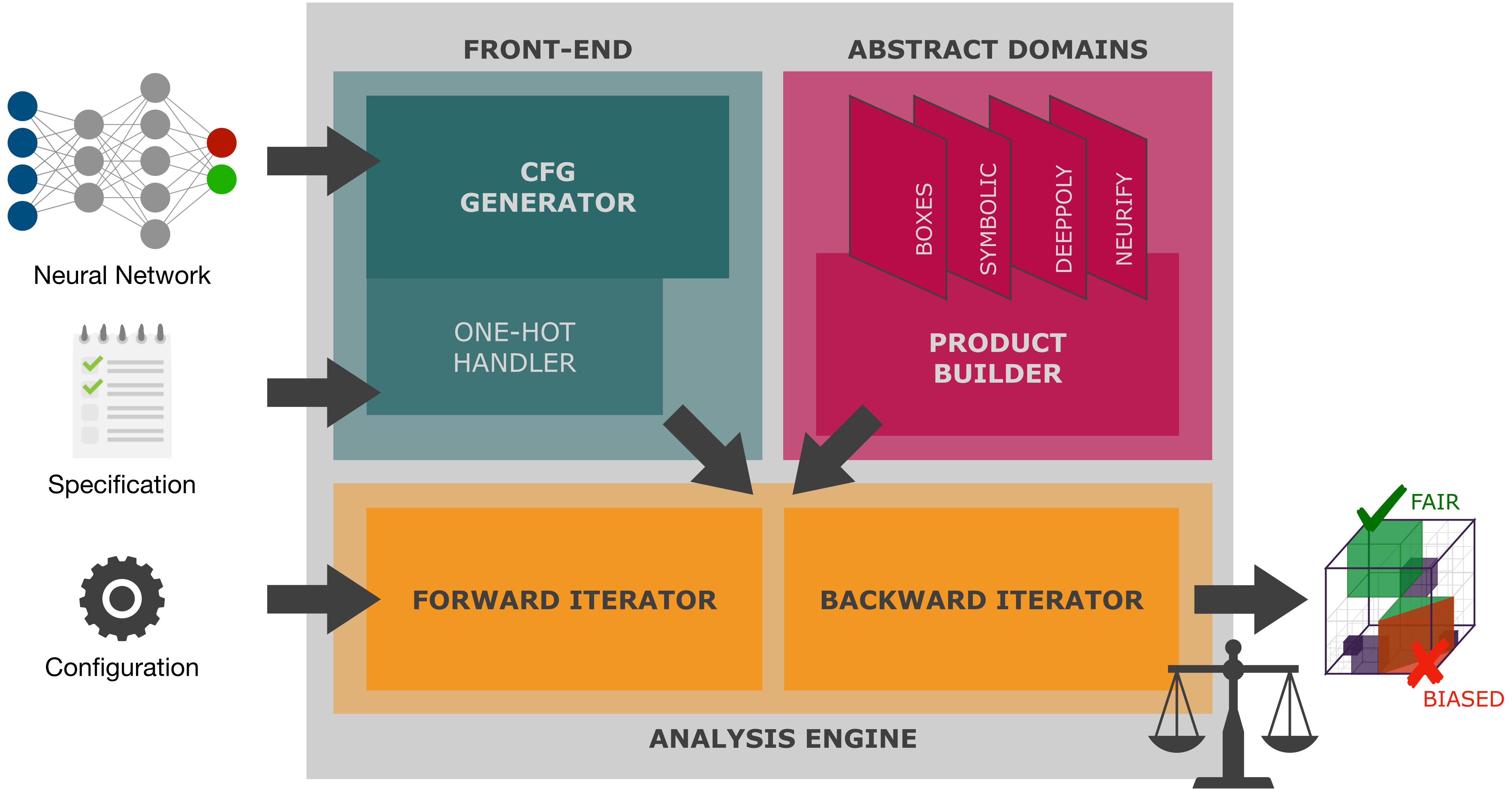
Mazzucato, Urban @ SAS 2021

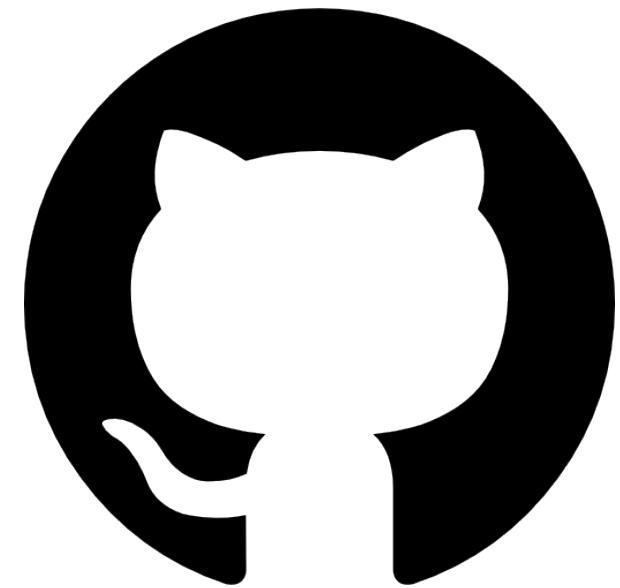
$$\hat{z} \in \left\{ \begin{array}{ll} [\max(\textcolor{teal}{l}, \textcolor{orange}{l}, l_{\text{low}}), \min(\textcolor{teal}{u}, \textcolor{orange}{u}, u_{\text{up}})] & \textbf{Symbolic} \\ [\max(\textcolor{teal}{l}, \textcolor{orange}{l}, l_{\text{low}}), \min(\textcolor{teal}{u}, \textcolor{orange}{u}, u_{\text{up}})] & \textbf{DeepPoly} \\ [\max(\textcolor{teal}{l}, \textcolor{orange}{l}, l_{\text{low}}), \max(\max(\textcolor{teal}{l}, \textcolor{orange}{l}, l_{\text{low}}), l_{\text{up}}), \\ \min(\min(\textcolor{teal}{u}, \textcolor{orange}{u}, u_{\text{up}}), u_{\text{low}}), \min(\textcolor{teal}{u}, \textcolor{orange}{u}, u_{\text{up}})] & \textbf{Neurify} \end{array} \right.$$

# Precision-vs-Scalability

L	U	Symbolic	DeepPoly	Neurify	Product	
0.5	3	48.78%	49.01%	46.49%	59.20%	+10.3%
	5	56.11%	56.15%	53.06%	68.23%	+11.9%
0.25	3	83.63%	81.82%	81.40%	87.04%	+3.4%
	5	91.67%	91.58%	92.33%	95.48%	+3.2%







Check it out on **GitHub!**

<https://github.com/caterinaurban/libra>

Ready-to-go **Docker image\*** at

<https://doi.org/10.5281/zenodo.4737450>

\* no installation needed!