

Permutation Invariance and Combinatorial Optimizations with Graph Deep Learning

Radu Balan

Department of Mathematics, CSCAMM and NWC
University of Maryland, College Park, MD

April 1-3, 2019

Workshop on "Dimension reduction in physical and data sciences"
Duke University, Durham NC 27708

Table of Contents:

- 1 Permutation Invariant Representations
 - Theory
 - Numerical Results
- 2 Optimizations using Deep Learning
 - Constructions
 - DNN as UA
 - Numerical Results

Linear Assignment Problems using GCN

The GCN design: Consider the GCN with $N + R$ nodes, adjacency/weight matrix $\mathbf{A} = \begin{bmatrix} 0 & C \\ C^T & 0 \end{bmatrix}$ and data matrix $X = \begin{bmatrix} \nu(C(i, :)) \\ \nu(C^T(j, :)) \end{bmatrix}$.

Linear Assignment Problems using GCN

The GCN design: Consider the GCN with $N + R$ nodes, adjacency/weight matrix $\mathbf{A} = \begin{bmatrix} 0 & C \\ C^T & 0 \end{bmatrix}$ and data matrix $X = \begin{bmatrix} \nu(C(i, :)) \\ \nu(C^T(j, :)) \end{bmatrix}$.

Key observation: When $C = uv^T$, that is, when the cost matrix is rank one then:

- ① Objective Function: $J(\Pi; C) = u^T \Pi v = \langle \Pi v, u \rangle$
- ② GCN output when no bias ($B_j = 0$): $\Gamma = \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix}$ satisfies $\Gamma_1 \Gamma_2^T = \alpha C$.

Consequence: the "greedy" algorithm produces the optimal solution.

Linear Assignment Problems using GCN

The GCN design: Consider the GCN with $N + R$ nodes, adjacency/weight matrix $\mathbf{A} = \begin{bmatrix} 0 & C \\ C^T & 0 \end{bmatrix}$ and data matrix $X = \begin{bmatrix} \nu(C(i, :)) \\ \nu(C^T(j, :)) \end{bmatrix}$.

Key observation: When $C = uv^T$, that is, when the cost matrix is rank one then:

- ① Objective Function: $J(\Pi; C) = u^T \Pi v = \langle \Pi v, u \rangle$
- ② GCN output when no bias ($B_j = 0$): $\Gamma = \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix}$ satisfies $\Gamma_1 \Gamma_2^T = \alpha C$.

Consequence: the "greedy" algorithm produces the optimal solution.

Network Objective: Once trained, the GCN produces a latent representation $Z = \Gamma_1 \Gamma_2^T$ close to the input cost matrix C so that the greedy algorithm applied on Z produces the optimal solution.

Quadratic Assignment Problem using GCN

Preliminary result

The GCN Design: Consider the GCN with n nodes, adjacency/weight

$$\text{matrix } \mathbf{A} = \begin{bmatrix} 0 & AB \\ BA & 0 \end{bmatrix} \text{ and data matrix } X = \begin{bmatrix} A \\ B \end{bmatrix}.$$

Quadratic Assignment Problem using GCN

Preliminary result

The GCN Design: Consider the GCN with n nodes, adjacency/weight

matrix $\mathbf{A} = \begin{bmatrix} 0 & AB \\ BA & 0 \end{bmatrix}$ and data matrix $X = \begin{bmatrix} A \\ B \end{bmatrix}$.

Key observation: When $A = uu^T$ and $B = vv^T$, that is, when the matrices are rank one then:

- Objective function: $J(\Pi; A, B) = (u^T \Pi v)^2 = (\langle \Pi v, u \rangle)^2$
- GCN output when no bias ($B_j = 0$): $\Gamma = \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix}$ satisfies

$$\Gamma_1 \Gamma_2^T \sim uv^T.$$

Consequence: the "greedy" algorithm or the solution to the linear assignment problem associated to uv^T produces the optimal solution.

Quadratic Assignment Problem using GCN

Preliminary result

The GCN Design: Consider the GCN with n nodes, adjacency/weight

matrix $\mathbf{A} = \begin{bmatrix} 0 & AB \\ BA & 0 \end{bmatrix}$ and data matrix $X = \begin{bmatrix} A \\ B \end{bmatrix}$.

Key observation: When $A = uu^T$ and $B = vv^T$, that is, when the matrices are rank one then:

- Objective function: $J(\Pi; A, B) = (u^T \Pi v)^2 = (\langle \Pi v, u \rangle)^2$
- GCN output when no bias ($B_j = 0$): $\Gamma = \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix}$ satisfies

$$\Gamma_1 \Gamma_2^T \sim uv^T.$$

Consequence: the "greedy" algorithm or the solution to the linear assignment problem associated to uv^T produces the optimal solution.

Network Objective: Once trained, the GCN produces a latent representation $Z = \Gamma_1 \Gamma_2^T$ so that the linear assignment problem associated to Z produces the same optimal permutation.

Deep Neural Networks as Universal Approximators

$$\begin{aligned}
 & \text{minimize/maximize} && \sum_{i=1}^N \sum_{j=1}^R \pi_{i,j} C_{i,j} \\
 & \text{subject to:} \\
 & \pi_{i,j} \in \{0, 1\}, \forall i, j \\
 & \sum_{i=1}^N \pi_{i,j} = 1, \forall 1 \leq j \leq R \\
 & \sum_{j=1}^R \pi_{i,j} \leq 1, \forall 1 \leq i \leq N
 \end{aligned}$$

Luckily, the convex relaxation (Linear Program) produces the same optimal solution:

$$\begin{aligned}
 & \text{minimize} && \sum_{i=1}^N \sum_{j=1}^R \pi_{i,j} C_{i,j} \\
 & \text{subject to:} \\
 & 0 \leq \pi_{i,j} \leq 1, \forall i, j \\
 & \sum_{i=1}^N \pi_{i,j} = 1, \forall 1 \leq j \leq R \\
 & \sum_{j=1}^R \pi_{i,j} \leq 1, \forall 1 \leq i \leq N
 \end{aligned}$$

Deep Neural Networks as Universal Approximators

Architectures

The overall system must output feasible solutions $\hat{\pi}$. Our architecture compose two components: (1) a deep neural network (DNN) that outputs a (generally) unfeasible estimate $\bar{\pi}$; (2) an enforcer (P) of the feasibility conditions that outputs the estimate $\hat{\pi}$:



Issues:

- 1 DNN architecture: how many layers; how many neurons per layer?
- 2 P , the feasibility enforcer

Deep Neural Networks as Universal Approximators

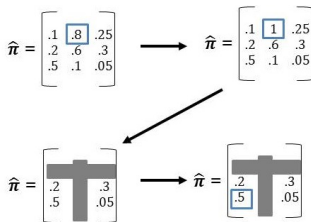
Feasibility Enforcer P

An "optimal" feasibility condition enforcer would minimize some "distance" to the feasibility set. However this may be a very computationally expensive component. An intermediate solution is to alternate between different feasibility conditions (equalities and inequalities) until convergence.

Instead we opt for a simpler and "greedier" approach:

Repeat R times:

1. Find (i, j) the largest entry in $\bar{\pi}$
2. Set $\hat{\pi}_{i,j} = 1$; set to 0 other entries in row i and column j ;
3. Remove row i and column j from both $\bar{\pi}$ and $\hat{\pi}$.



Deep Neural Networks as Universal Approximators

Baseline solution: The Greedy Algorithm

The "greedy" enforcer can be modified into a "greedy" optimization algorithm:

- 1 Initialize $E = C$ and $\hat{\pi} = 0_{N \times R}$
- 2 Repeat R times:
 - Find $(i, j) = \operatorname{argmin}_{(a,b)} E_{a,b}$;
 - Set $\hat{\pi}_{i,j} = 1$, $\hat{\pi}_{i,l} = 0 \forall l \neq j$, $\hat{\pi}_{l,j} = 0 \forall l \neq i$;
 - Set $E_{i,:} = \infty$, $E_{:,j} = \infty$.

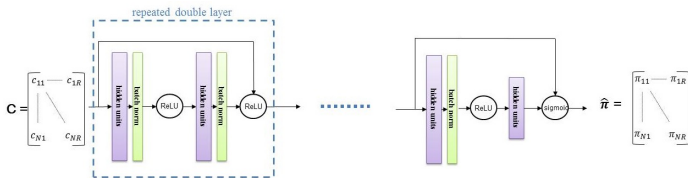
Proposition

The greedy algorithm produces the optimal solution if there is a positive number $\lambda > 0$ and two nonnegative vectors u, v such that

$$C = \lambda \mathbf{1} \cdot \mathbf{1}^T - u \cdot v^T.$$

Exp.1 : $N = 5$, $R = 4$ with ReLU activation

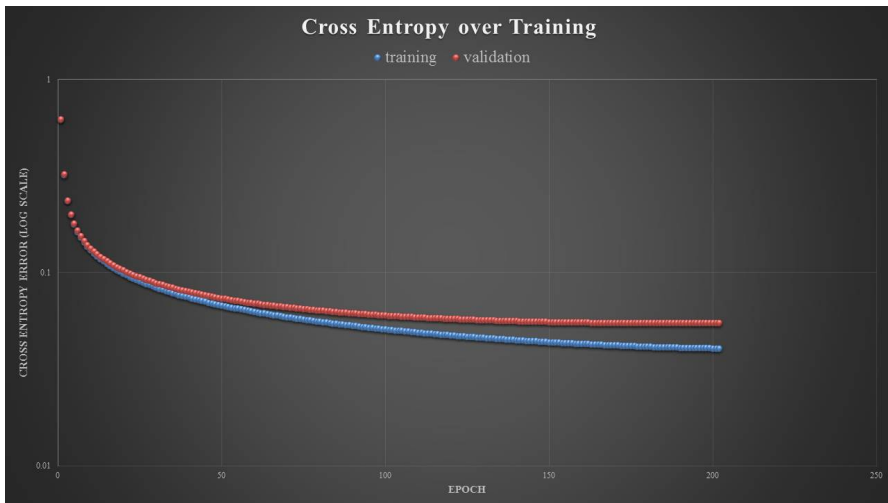
First architecture:



- Number of internal layers: 9
- Number of hidden units per layer: 250
- Batch size: 200; ADAM optimizer
- Loss function: cross-entropy:

$$\sum_{i,j} \pi_{i,j} (-\log(\hat{\pi}_{i,j})) + (1 - \pi_{i,j}) (-\log(1 - \hat{\pi}_{i,j}))$$
- Training data set: 1 million random instances $U(0, 1)$ i.i.d.
- Validation set: 20,000 random instances.

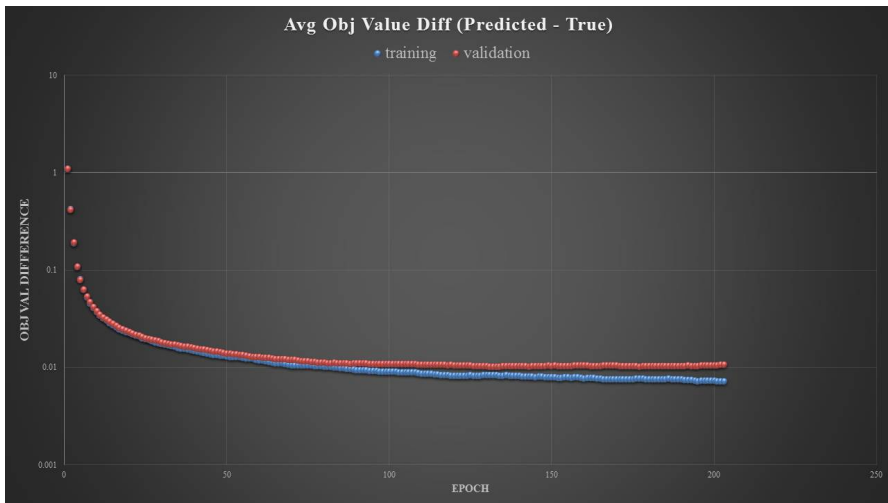
Numerical Results

Exp.1 : $N = 5$, $R = 4$ with ReLU activation

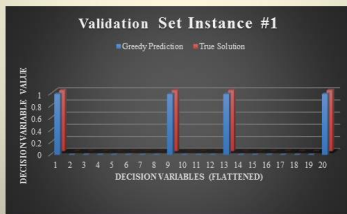
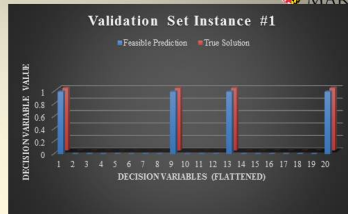
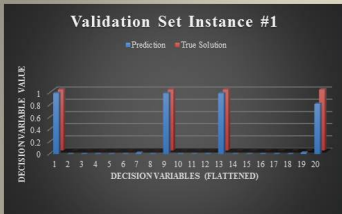
Numerical Results

Exp.1 : $N = 5$, $R = 4$ with ReLU activation

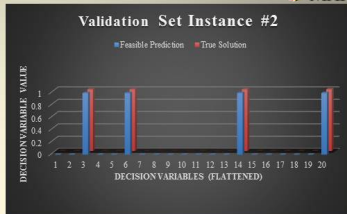
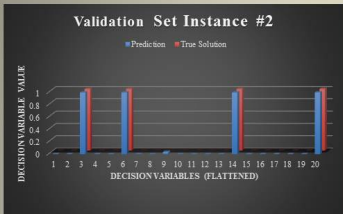
Numerical Results

Exp.1 : $N = 5$, $R = 4$ with ReLU activation

Numerical Results

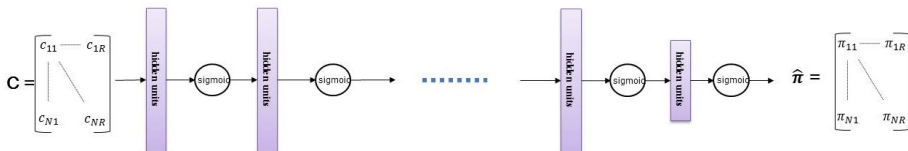
Exp.1 : $N = 5, R = 4$ with ReLU activation

Numerical Results

Exp.1 : $N = 5, R = 4$ with ReLU activation

Exp.2 : $N = 10, R = 8$ with sigmoid activation

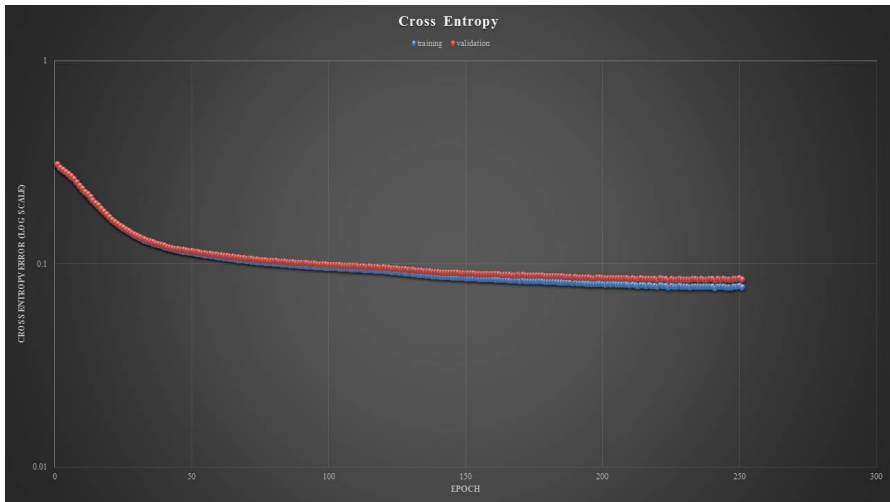
Second architecture:



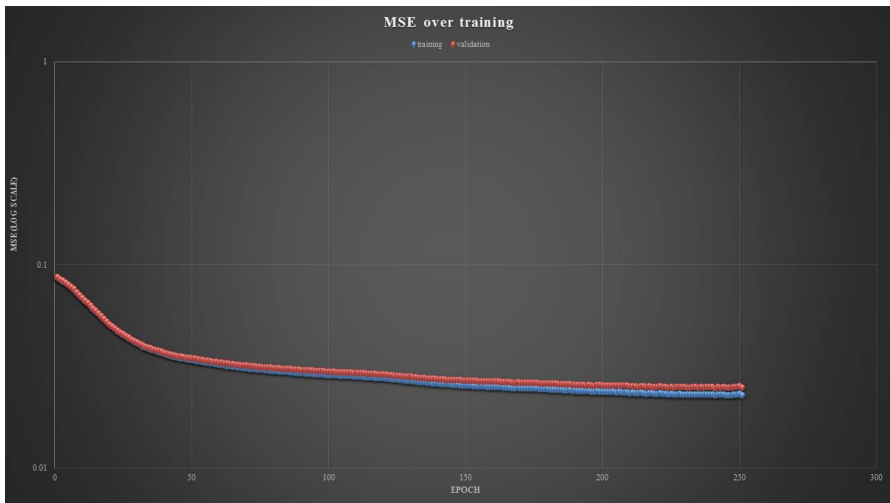
- Number of internal layers: 10
- Number of hidden units per layer: 250
- No Batch; ADAM optimizer
- Loss function: cross-entropy:

$$\sum_{i,j} \pi_{i,j} (-\log(\hat{\pi}_{i,j})) + (1 - \pi_{i,j}) (-\log(1 - \hat{\pi}_{i,j}))$$
- Training data set: 1 million random instances $U(0, 1)$ i.i.d.
- Validation set: 20,000 random instances.

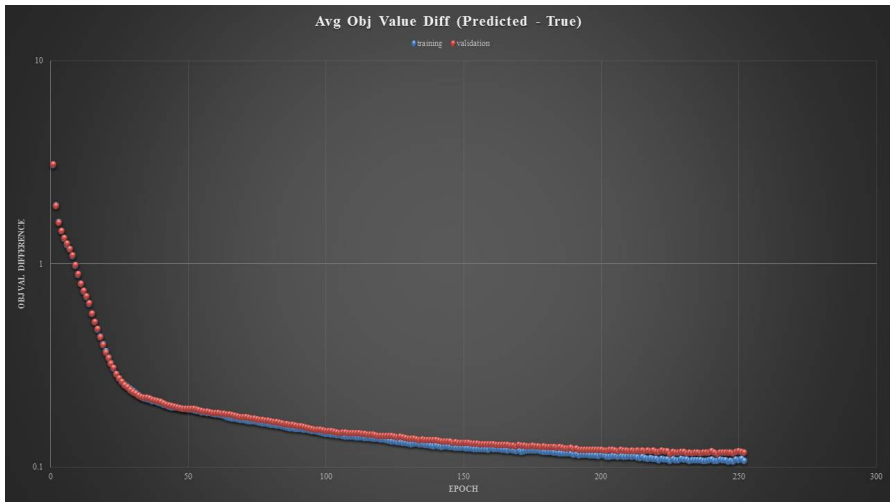
Numerical Results

Exp.2 : $N = 10$, $R = 8$ with sigmoid activation

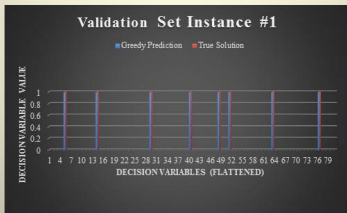
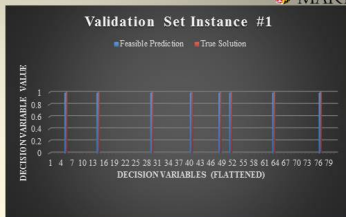
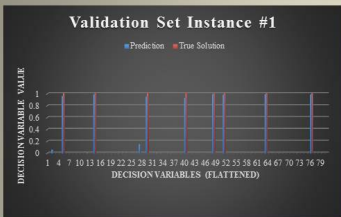
Numerical Results

Exp.2 : $N = 10$, $R = 8$ with sigmoid activation

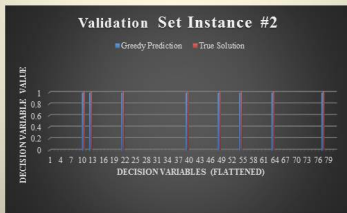
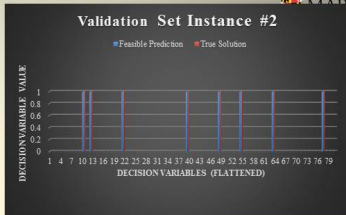
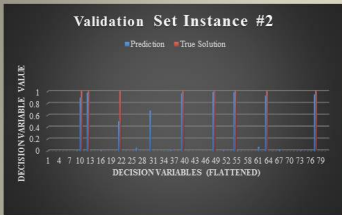
Numerical Results

Exp.2 : $N = 10$, $R = 8$ with sigmoid activation

Numerical Results

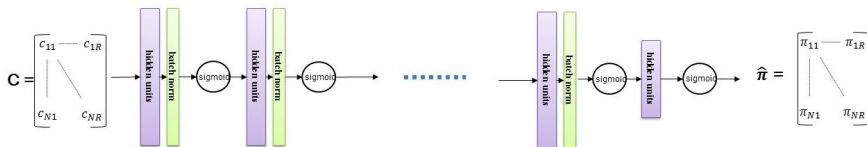
Exp.2 : $N = 10$, $R = 8$ with sigmoid activation

Numerical Results

Exp.2 : $N = 10, R = 8$ with sigmoid activation

Exp.3 : $N = 5, R = 4$ with sigmoid activation

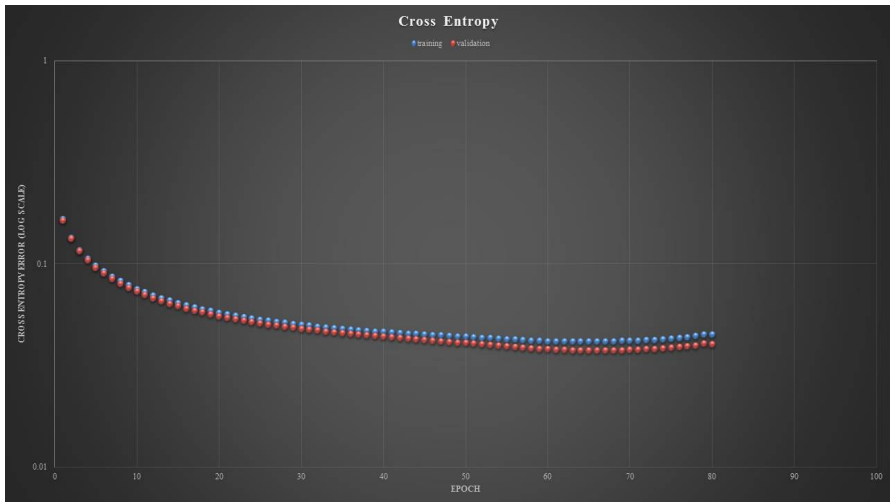
Second architecture:



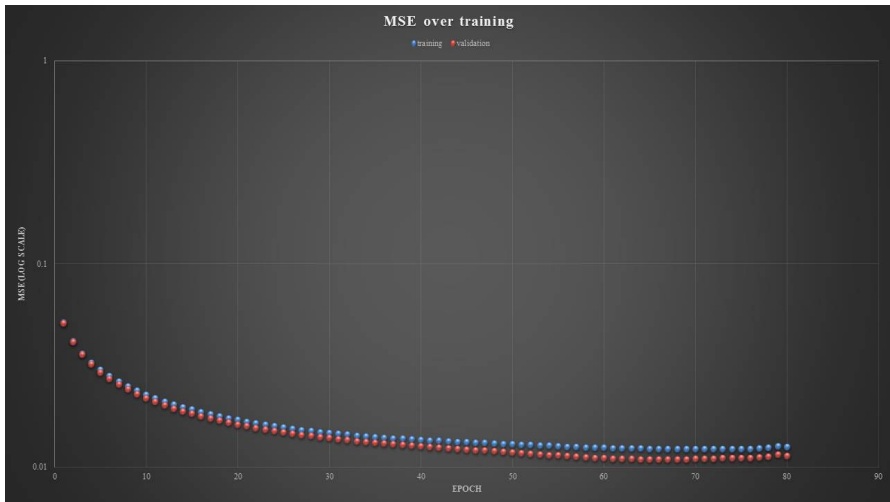
- Number of internal layers: 10
- Number of hidden units per layer: 250
- Batch size 200; ADAM optimizer
- Loss function: cross-entropy:

$$\sum_{i,j} \pi_{i,j} (-\log(\hat{\pi}_{i,j})) + (1 - \pi_{i,j}) (-\log(1 - \hat{\pi}_{i,j}))$$
- Training data set: 500,000 random instances $U(0, 1)$ i.i.d.
- Validation set: 20,000 random instances.

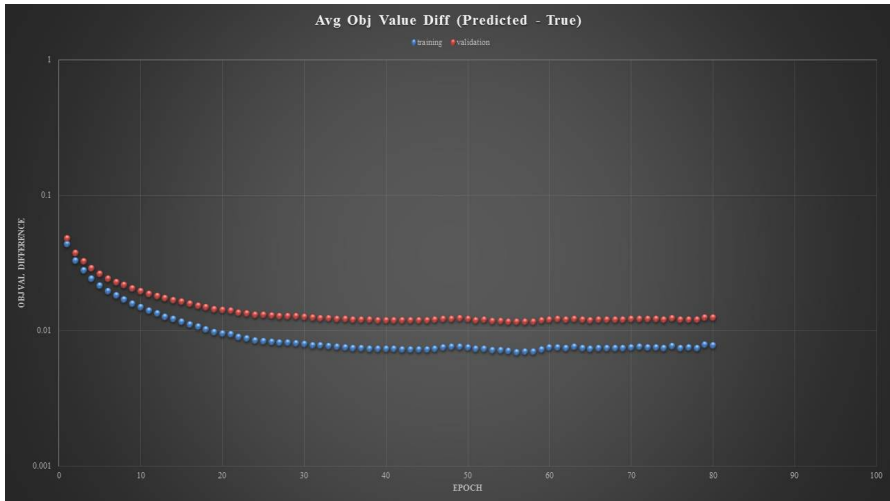
Numerical Results

Exp.3 : $N = 5$, $R = 4$ with sigmoid activation

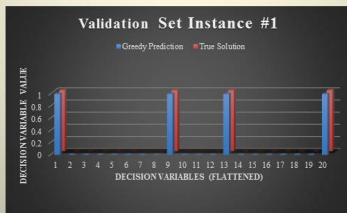
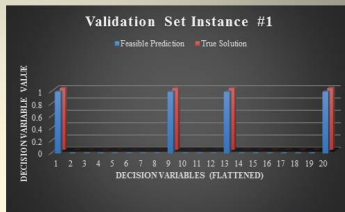
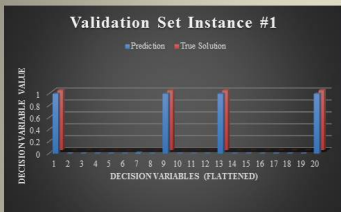
Numerical Results

Exp.3 : $N = 5$, $R = 4$ with sigmoid activation

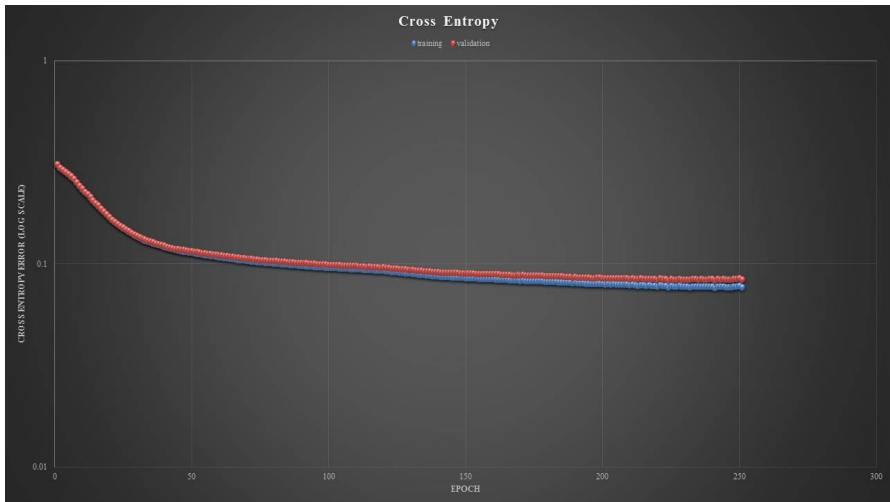
Numerical Results

Exp.3 : $N = 5$, $R = 4$ with sigmoid activation

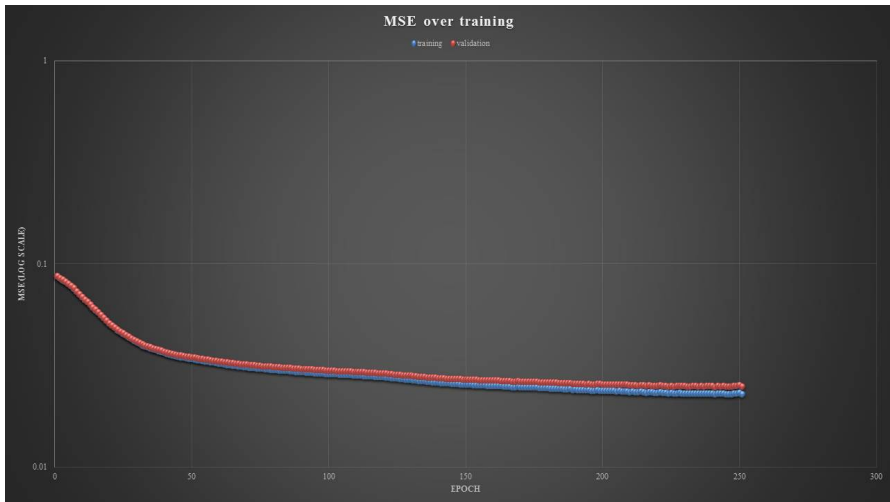
Numerical Results

Exp.3 : $N = 5, R = 4$ with sigmoid activation

Numerical Results

Exp.4 : $N = 10$, $R = 8$ with sigmoid activation

Numerical Results

Exp.4 : $N = 10$, $R = 8$ with sigmoid activation

Bibliography

1. M. Andrychowicz, M. Denil, S.G. Colmenarejo, M.W. Hoffman, D. Pfau, T. Schaul, B. Shillingford, N.de Freitas, *Learning to learn by gradient descent by gradient descent*, arXiv:1606.04474v2 [cs.NE]
2. T.N. Kipf, M. Welling, *Variational Graph Auto-Encoder*, arXiv:1611.07308 [stat.ML]
3. A. Nowak, S. Villar, A. Bandeira, J. Bruna, *Revised Note on Learning Quadratic Assignment with Graph Neural Network*, arXiv: 1706.07450 [stat.ML]