



Motivation

Principal matrix feature (PMaF): a single vector summarising a data matrix.

- Useful for feature representation with a lower dimension (a vector)
- Better and possibly faster solutions that are constrained on a sphere
- Differentiable for end-to-end learning
- Implicit differentiation [1] with higher running speed and lower hardware memory requirements than unrolling the forward optimization iteration [2] or without exploited structures [3]

Two Deep Declarative Layers

1. Least Squares on Sphere (LESS)

$$\text{minimize}_{\mathbf{u} \in \mathbb{R}^n} f(\mathbf{A}, \mathbf{b}, \mathbf{u}) \triangleq \frac{1}{2} \|\mathbf{A}\mathbf{u} - \mathbf{b}\|^2$$

$$\text{subject to } \|\mathbf{u}\|^2 = 1.$$

Iterative optimization:

1. Projected gradient descent (PGD)
2. + Direction weight (DW)

$$\mathbf{w}_t = 1 - S_c(\mathbf{d}_t, \mathbf{u}_0)$$

$$S_c(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a}^T \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|}, \mathbf{d}_t = \begin{cases} -\nabla f(\mathbf{u}_t) & \text{if } \|\mathbf{u}_0\| \geq 1 \\ \nabla f(\mathbf{u}_t) & \text{otherwise} \end{cases}$$

3. + Riemannian manifold (RM)

$$\text{Proj}_{\text{RM}}(-\nabla f(\mathbf{u}_t)) = (\mathbf{I}_n - \mathbf{u}_t \mathbf{u}_t^T)(-\nabla f(\mathbf{u}_t))$$

$$\mathbf{u}_{t+1} = \text{Proj}_{\text{sph}}(\mathbf{u}_t + \eta \text{Proj}_{\text{RM}}(-\nabla f(\mathbf{u}_t)))$$

$$\eta \leftarrow \beta \eta$$

4. + Backtracking line search (BLS)

$$f(\mathbf{u}_t + \eta \Delta \mathbf{u}_t) \leq f(\mathbf{u}_t) + \alpha \eta \nabla^T f(\mathbf{u}_t) \Delta \mathbf{u}_t$$

5. + Tangent weight decay (TWD)

$$S_c(\text{Proj}_{\text{RM}}(-\nabla f(\mathbf{u}_t)), \text{Proj}_{\text{RM}}(-\nabla f(\mathbf{u}_{t+1}))) < 0$$

2. Implicit Eigen Decomposition (IED)

$$\text{minimize}_{\mathbf{u} \in \mathbb{R}^{m \times n}} f(\mathbf{A}, \mathbf{u}) \triangleq -\text{tr}(\mathbf{u}^T \mathbf{A} \mathbf{u})$$

$$\text{subject to } \mathbf{h}(\mathbf{u}) \triangleq \mathbf{u}^T \mathbf{u} = \mathbf{I}_n.$$

Iterative optimization:

1. Power iteration (PI)

$$\mathbf{u}_{t+1} = \mathbf{A} \mathbf{u}_t / \|\mathbf{A} \mathbf{u}_t\|$$

$$\mathbf{y} = \mathbf{u}_K \text{ and } \lambda = \mathbf{y}^T \mathbf{A} \mathbf{y}$$

2. Simultaneous iteration (SI)

$$\{\mathbf{Q}_t, \mathbf{R}_t\} = \text{QR}(\mathbf{x}_t) \text{ and } \mathbf{x}_{t+1} = \mathbf{x}_t \mathbf{Q}_t$$

$$\mathbf{y} = \text{the component of } \mathbf{Q}_K \text{ to } \lambda = \max(\mathbf{R}_K)$$

➤ Solution consistency in iterations

$$\text{Historical: } \mathbf{u}_t \leftarrow \nabla(\mathbf{u}_t, \mathbf{u}_{t-1}) \mathbf{u}_t$$

$$\text{Hard-coded: } \mathbf{y} \leftarrow \nabla(\mathbf{y}, \mathbf{r}) \mathbf{y}$$

$$\nabla(\mathbf{a}, \mathbf{b}) = \text{Sign}(\mathbf{a}^T \mathbf{b}) \text{ if } (\mathbf{a} \not\perp \mathbf{b}) \text{ and otherwise } 1$$

Deep Declarative Networks as Backward

Implicit differentiation of the learning loss over inputs using DDN [1] and exploited structures [3]:

$$\nabla_X L = \nabla_Y L (\mathcal{H}^{-1} \mathbf{A}^T (\mathcal{A} \mathcal{H}^{-1} \mathbf{A}^T)^{-1} \mathbf{A} - \mathbf{I}_n) \mathcal{H}^{-1} \mathbf{B}$$

Vanilla

$$\mathbf{A} = 2\mathbf{y}^T \in \mathbb{R}^{1 \times n}$$

$$\mathbf{B} = \nabla_{XY}^2 f(\mathbf{A}, \mathbf{b}, \mathbf{y}) \in \mathbb{R}^{n \times (m \times n)}$$

$$\mathcal{H} = \mathbf{A}^T \mathbf{A} - 2\beta \mathbf{I}_n \in \mathbb{R}^{n \times n}$$

$$\beta = \frac{1}{2} \mathbf{y}^T \mathbf{A}^T (\mathbf{A} \mathbf{y} - \mathbf{b}) \in \mathbb{R}$$

With exploited structures

$$\mathbf{B}_{ij} = \mathbf{A}_i \mathbf{y}_j^T, \forall i, j \in \mathcal{N},$$

$$\mathbf{B}_{ii} \leftarrow \mathbf{B}_{ii} + (\mathbf{A} \mathbf{y} - \mathbf{b}), \forall i \in \mathcal{N}.$$

Vanilla

$$\mathbf{A} = 2\mathbf{y}^T \in \mathbb{R}^{1 \times m}$$

$$\mathbf{B} = \nabla_{XY}^2 f(\mathbf{A}, \mathbf{y}) \in \mathbb{R}^{m \times (m \times m)}$$

$$\mathcal{H} = -(\mathbf{A} + \mathbf{A}^T) - 2\beta \mathbf{I}_m \in \mathbb{R}^{m \times m}$$

$$\beta = -\frac{1}{2} \mathbf{y}^T (\mathbf{A} \mathbf{y} + \mathbf{A}^T \mathbf{y}) \in \mathbb{R}$$

With exploited structures

$$\mathbf{B}_{ijk} = 0, \forall i, j, k \in \mathcal{M},$$

$$\mathbf{B}_{iji} \leftarrow \mathbf{B}_{iji} - \mathbf{y}_j, \forall i, j \in \mathcal{M},$$

$$\mathbf{B}_{iij} \leftarrow \mathbf{B}_{iij} - \mathbf{y}_j, \forall i, j \in \mathcal{M}.$$

$$\text{Further } \nabla_X L = -\mathcal{K} \mathbf{y}^T - \mathbf{y} \mathcal{K}^T.$$

IED-implicit function theorem (IFT)

Objective function

$$f(\mathbf{A}, \mathbf{y}) = \mathbf{y} - \mathbf{A} \mathbf{y} / \|\mathbf{A} \mathbf{y}\|$$

Implicit gradients

$$\nabla_X \mathbf{y} = -(\nabla_Y f(\mathbf{A}, \mathbf{y}))^{-1} \nabla_X f(\mathbf{A}, \mathbf{y})$$

$$\nabla_X L = \nabla_Y L \nabla_{XY} = -\nabla_Y L \mathcal{H}^{-1} \mathbf{B}$$

$$\mathcal{H} = \mathbf{I}_m - (\mathbf{I}_m - \mathbf{y} \mathbf{y}^T) \mathbf{A} / \lambda,$$

$$\mathbf{B} = -(\mathbf{I}_m - \mathbf{y} \mathbf{y}^T) \mathbf{y}^T / \lambda.$$

Improvements on LESS Solvers

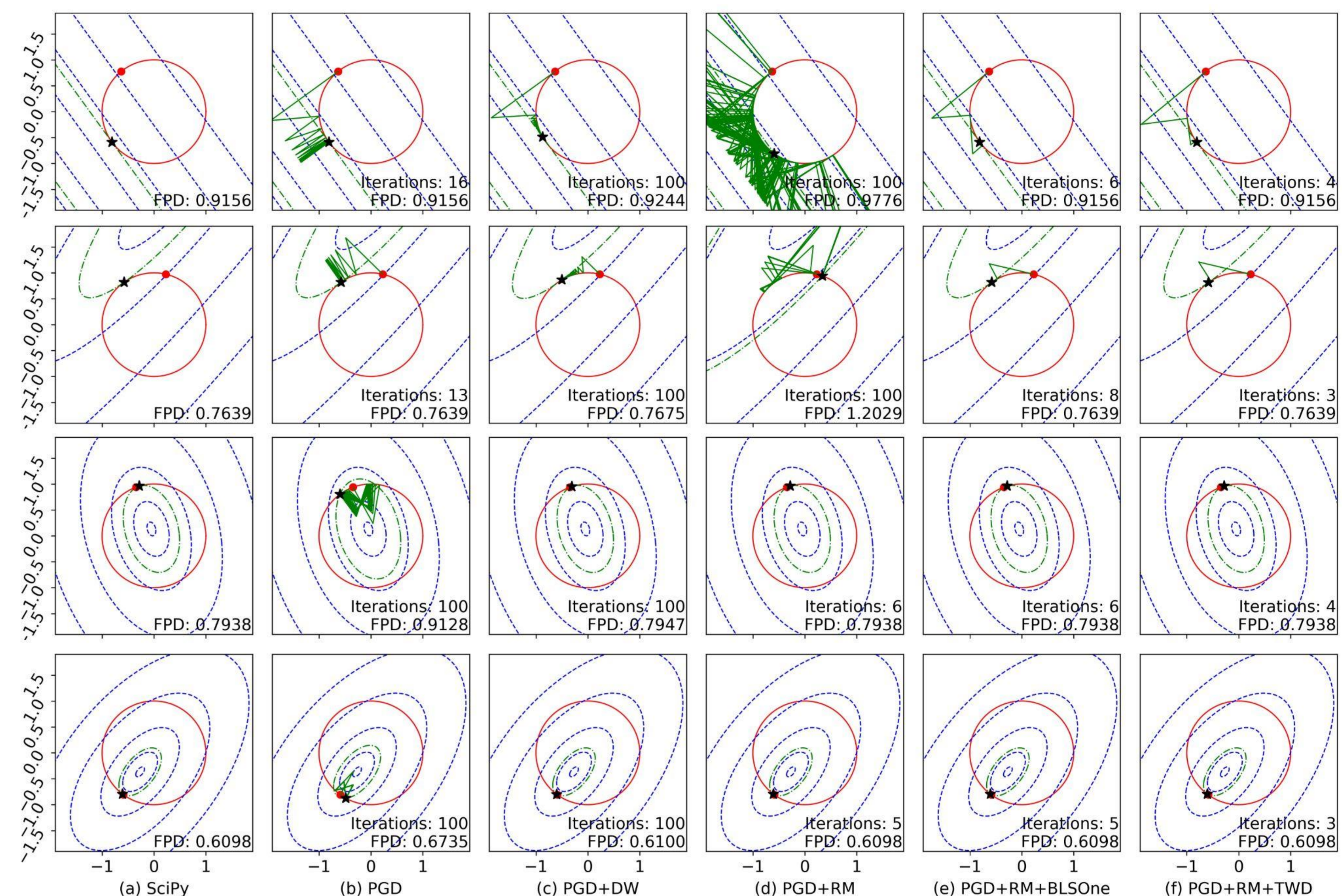


Figure 1. Iterative optimization in LESS. Each row is a sample with 6 methods. “blue dash”: the least squares function; “red solid”: the sphere constraints; “green solid”: solution updates before and after Riemannian projection; “green dot-dash”: the least squares function with the final solution; “red dot”: the initial solution; “black star”: the final solution, not always the optimal. *Ours (the last two columns) require much fewer iterations than the others for the optimal solution with comparable fixed point distance (FPD).*

Table 1. Effectiveness on 1,000 random Gaussian samples. “In” and “Out”: failed cases inner and outer of the sphere respectively. A case is failed when the solution update reaches 100 iterations, “Imp.”: the number of cases with energy no greater than SciPy. Problem sizes, $m-n$ for the input matrix, are 2-2 (282 inner and 718 outer), 64-32 (555 inner and 445 outer), and 1024-256 (all inner).

	2x2				64x32				1024x256			
	In↓	Out↓	Imp.↑	MRE↓	In↓	Out↓	Imp.↑	MRE↓	In↓	Out↓	Imp.↑	MRE↓
SciPy	✓											
PGD		✓										
RM			✓									
BLS $\eta = 1$				✓								
$\eta = \nabla$								✓				
TWD												
DW												
2x2	172	353	389	10.84	548	440	52	0.15	1000	-	0	1.35
64x32	179	438	37	4.89	474	363	63	0.07	1000	-	0	0.47
1024x256	0	118	844	4.30	30	0	980	0.00	4	-	1000	-0.03
Imp.	0	14	938	0.14	6	0	994	0.00	4	-	1000	-0.03
Energy	175	421	40	3.51	285	142	286	0.01	984	-	999	-0.03
Outer	0	108	853	4.21	7	0	994	0.00	2	-	1000	-0.03
MRE	204	536	17	1.19	472	355	92	0.04	1000	-	0	0.42
Outer	0	0	806	0.00	29	0	954	0.00	4	-	1000	-0.03
MRE	204	533	15	3.12	466	352	90	0.04	1000	-	0	0.42

Improvements on IED Solvers

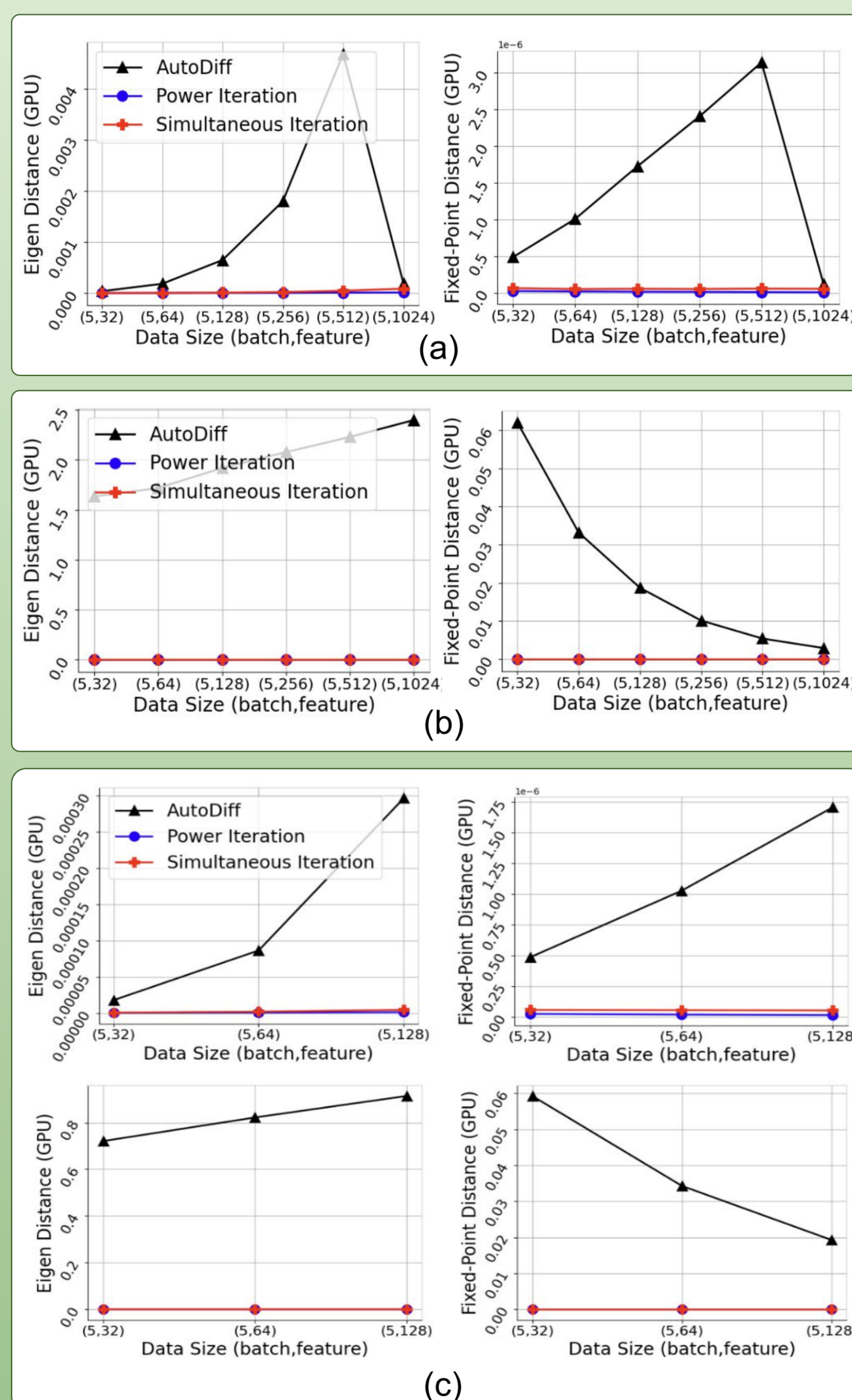


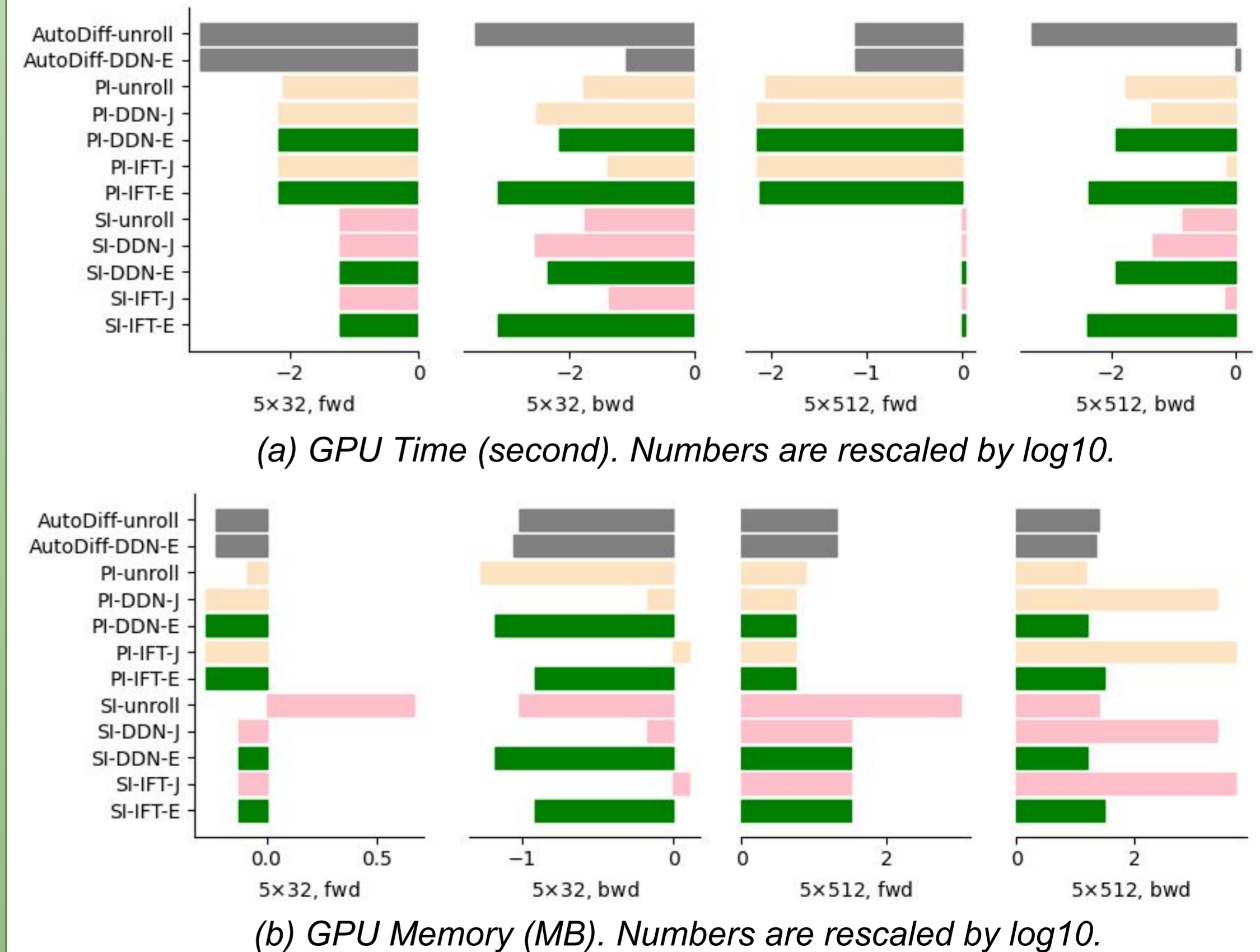
Figure 2. Accuracy evaluation using eigen distance and FPD (float32). (a) Symmetric and (b) non-symmetric matrices, sampled from the Normal distribution with activation. (c) Matrices from a pretrained ResNet50 with symmetric (top two rows) and non-symmetric (bottom two rows) matrices. Metrics are the less the better.

Improvements on Backward Efficiency

Table 2. Backward speedup of LESS with exploited structures. Numbers are averaged over 100 samples. “PGD+RM+TWD” is used for LESS. “Speedup” is (“AutoDiff”-“LESS”)/“LESS”, “AutoDiff” from SciPy solver.

	Size 256 × 8 × 64		Size 256 × 16 × 128	
	Time (s)	Memory (MB)	Time (s)	Memory (MB)
AutoDiff	4.44	65.22	8.64	519.32
LESS	0.01	49.20	0.02	325.70
Speedup	×444↑	×0.33↑	×432↑	×0.59↑

Figure 3. IED evaluation on symmetric Gaussian matrices with absolute activation. “fwd”: forward pass; “bwd”: backward pass; “AutoDiff”: PyTorch eigh(); “PI”: power iteration; “SI”: simultaneous iteration; “unroll”: unrolling the forward iteration via PyTorch autodiff mechanism; “J”: autodiff Jacobian without exploited structure; “E”: ours with exploited structure. Best suggestions are highlighted with green color considering the overall precision in Fig. 2 and computational requirements. More results on non-symmetric matrices and different data types are in the Appendix.



Conclusion

- Two deep declarative layers for PMaF, namely LESS and IED, for end-to-end learning
- Overall better and faster (for regularized matrices in IED) solutions using iterative optimization
- Efficient and implicit differentiation with exploited matrix structures
- Code available <https://github.com/anucvml/ddn.git>

References

- [1] Gould, S., Hartley, R., and Campbell, D. Deep declarative networks. *TPAMI*, 2022.
- [2] Boyd, S. and Vandenberghe, L. Convex optimization. Cambridge University Press, 2004.
- [3] Gould, S., Campbell, D., Ben-Shabat, I., Koneputugodage, C. H., and Xu, Z. Exploiting problem structure in deep declarative networks: Two case studies. *AAAI Workshop on OTSDM*, 2022.