

Problem Set 2  
6.S091: Topics in Causality  
IAP 2024

**Problem 1: Implementing DAGs with NO TEARS [30 points]**

**Preliminaries [10 point]**

(a) We have

$$\begin{aligned}\text{tr}(\mathcal{A}(W)^k) &= \sum_{i=1}^d (\mathcal{A}(W)^k)_{ii} \\ &= \sum_{i=1}^d \sum_{t_1, \dots, t_{k-2}=1}^d \mathcal{A}(W)_{it_1} \mathcal{A}(W)_{t_1 t_2} \dots \mathcal{A}(W)_{t_{k-2} i}.\end{aligned}$$

Note that if  $i - t_1 - \dots - t_{k-2} - i$  is a (length- $k$ ) cycle in  $\mathcal{G}(W)$ ,

$$\mathcal{A}(W)_{it_1} \mathcal{A}(W)_{t_1 t_2} \dots \mathcal{A}(W)_{t_{k-2} i} = 1,$$

otherwise

$$\mathcal{A}(W)_{it_1} \mathcal{A}(W)_{t_1 t_2} \dots \mathcal{A}(W)_{t_{k-2} i} = 0.$$

As the sums in  $\text{tr}(\mathcal{A}(W)^k)$  consist of different and all sequences of  $i, t_1, \dots, t_{k-2}, i$ , it equals to the number of length- $k$  cycles in  $\mathcal{G}(W)$ .

(b) By the definition of matrix exponential and the linearity of trace, we have

$$\text{tr}(e^{W \circ W}) = \text{tr}(I) + \sum_{k=1}^{+\infty} \frac{1}{k!} \cdot \text{tr}((W \circ W)^k),$$

where  $I$  is the  $d \times d$  identity matrix. Since all entries of  $W \circ W$  are non-negative and  $\text{tr}(I) = d$ , we immediately have  $\text{tr}(e^{W \circ W}) \geq d$ , i.e.,  $h(W) \geq 0$ .

When  $h(W) = 0$ , there is  $\text{tr}((W \circ W)^k) = 0$  for all positive integer  $k$ . Note that using a similar proof as (a), we have that  $\text{tr}((W \circ W)^k)$  equals to the weighted sum of length- $k$  circles in  $\mathcal{G}(W)$ , where the weight of each circle is the square of the product of weights of the edges on this circle. When a circle is present, this weight is greater than 0. Therefore  $\text{tr}((W \circ W)^k) = 0$  if and only if there is no length- $k$  circle in  $\mathcal{G}(W)$ . As this holds for all positive integer  $k$ , we know that  $h(W) = 0$  if and only if there is no cycles in  $\mathcal{G}(W)$ .

The gradient of  $h(W)$  is  $\nabla h(W) = e^{W \circ W} \circ 2W$ .

**Implementing the constraint [10 point]**

(c) See the following snippet.

```
1 import numpy as np
2 import scipy.linalg as slin
3
4 def h(W):
5     E = slin.expm(W * W)
```

```

6     h = np.trace(E) - len(W)
7     return h

```

(d)  $h(W_1) = 0$ ,  $h(W_2) = 0.2826$ .

## Optimizing with augmented Lagrangian [10 point]

(e) See the following snippet.

```

1 def notears_linear(X, max_iter=100, h_tol=1e-6, rho_max=1e+16):
2     """Solve min_W F(W) s.t. h(W) = 0 using augmented Lagrangian.
3
4     Args:
5         X (np.ndarray): [n, d] sample matrix
6         max_iter (int): max num of dual ascent steps
7         h_tol (float): exit if |h(W_est)| <= htol
8         rho_max (float): exit if rho >= rho_max
9
10    Returns:
11        W_est (np.ndarray): [d, d] estimated weights
12    """
13    def _loss(W):
14        """Evaluate value and gradient of loss."""
15        W = W.reshape((d, d))
16        M = X @ W
17        R = X - M
18        loss = 0.5 / X.shape[0] * (R ** 2).sum()
19        G_loss = - 1.0 / X.shape[0] * X.T @ R
20        return loss, G_loss
21
22    def _h(W):
23        """Evaluate value and gradient of acyclicity constraint."""
24        W = W.reshape((d, d))
25        E = slin.expm(W * W) # (Zheng et al. 2018)
26        h = np.trace(E) - d
27        G_h = E.T * W * 2
28        return h, G_h
29
30    def _func(W):
31        """Evaluate value and gradient of augmented Lagrangian."""
32        loss, G_loss = _loss(W)
33        h, G_h = _h(W)
34        obj = loss + 0.5 * rho * h * h + alpha * h
35        g_obj = np.concatenate((G_loss + (rho * h + alpha) * G_h), axis=None)
36
37        return obj, g_obj
38
39    n, d = X.shape
40    W_est, rho, alpha, h = np.zeros((d*d)), 1.0, 0.0, np.inf
41
42    for _ in range(max_iter):
43        W_new, h_new = None, None
44
45        W_new = sopt.minimize(_func, W_est, method='L-BFGS-B', jac=True).x
46        h_new, _ = _h(W_new)
47
48        alpha += rho * h_new
49
50        if h_new > h / 2:

```

```

51     rho *= 10
52
53     W_est, h = W_new, h_new
54
55     if h <= h_tol or rho >= rho_max:
56         break
57
58     return W_est.reshape((d, d))

```

(f) See the following for one run of the above algorithm, compared with the ground-truth  $W$ .



Figure 1: One run of the above algorithm v.s. the ground-truth  $W$ .

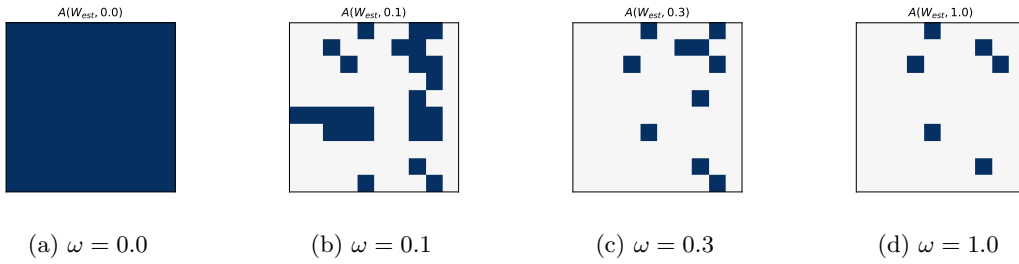


Figure 2: Adjacency matrix  $\mathcal{A}(W_{est}, \omega)$  with different threshold values.