

Unsupervised Learning Methods

Problem Set IV –

MDS, Isomap, Laplacian-Eigenmaps, and T-SNE



Or Livne - 203972922
Daniel Levi - 302506712

Classical MDS

Consider the following inner product:

$$\langle \mathbf{x}, \mathbf{y} \rangle_\phi := \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$$

for some suitable $\phi: \mathbb{R}^D \rightarrow \mathbb{R}^M$, where $\langle \cdot, \cdot \rangle$ is the standard Euclidean inner product, i.e. $\langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^T \mathbf{b}$. Consider:

1. The induced norm:

$$\|\mathbf{x}\|_\phi := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_\phi}$$

2. The induced metric:

$$d_\phi(\mathbf{x}, \mathbf{y}) := \|\mathbf{x} - \mathbf{y}\|_\phi$$

Consider a training set $\{\mathbf{x}_i\}_{i=1}^N$ and let $\mathbf{D}_\phi \in \mathbb{R}^{N \times N}$ where $\mathbf{D}_\phi[i, j] = d_\phi^2(\mathbf{x}_i, \mathbf{x}_j)$.

1.1

Show that

$$-\frac{1}{2} \mathbf{J} \mathbf{D}_\phi \mathbf{J} = \mathbf{J} \mathbf{K}_\phi \mathbf{J}$$

where:

1. $\mathbf{J} = \mathbf{I} - \frac{1}{N} \mathbf{1} \mathbf{1}^T \in \mathbb{R}^{N \times N}$
2. $\mathbf{K}_\phi := \Phi^T \Phi$, Φ is given by:

$$\Phi = \begin{bmatrix} | & | & & | \\ \phi_1 & \phi_2 & \dots & \phi_N \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{M \times N}$$

and $\phi_i = \phi(\mathbf{x}_i)$.

Steps:

1. Show that ϕ must be linear, namely:

$$\phi(\alpha \mathbf{x} + \beta \mathbf{y}) = \alpha \phi(\mathbf{x}) + \beta \phi(\mathbf{y})$$

Hint: Consider $\langle \alpha \mathbf{x} + \beta \mathbf{y}, \mathbf{z} \rangle_\phi$ and recall that this is true for all \mathbf{z} .

2. Show that:

$$\begin{aligned} d_\phi^2(\mathbf{x}, \mathbf{y}) &= \|\phi(\mathbf{x}) - \phi(\mathbf{y})\|_2^2 \\ &= \|\phi(\mathbf{x})\|_2^2 - 2 \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle + \|\phi(\mathbf{y})\|_2^2 \end{aligned}$$

3. Repeat\use the lecture notes to conclude that $-\frac{1}{2} \mathbf{J} \mathbf{D}_\phi \mathbf{J} = \mathbf{J} \mathbf{K}_\phi \mathbf{J}$.

1. MDS:

1.1

Considering $\langle \alpha \mathbf{x} + \beta \mathbf{y}, \mathbf{z} \rangle_\phi = \langle \phi(\alpha \mathbf{x} + \beta \mathbf{y}), \phi(\mathbf{z}) \rangle$

And on the other hand,

$$\begin{aligned} \langle \alpha \mathbf{x} + \beta \mathbf{y}, \mathbf{z} \rangle_\phi &\stackrel{\text{linearity in the first argument}}{=} \langle \alpha \mathbf{x}, \mathbf{z} \rangle_\phi + \langle \beta \mathbf{y}, \mathbf{z} \rangle_\phi = \\ &= \langle \alpha \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle + \langle \beta \phi(\mathbf{y}), \phi(\mathbf{z}) \rangle \end{aligned}$$

Then, $\langle \alpha \phi(x), \phi(z) \rangle + \langle \beta \phi(y), \phi(z) \rangle = \langle \phi(\alpha x + \beta y), \phi(z) \rangle$
 $\langle \alpha \phi(x) + \beta \phi(y) - \phi(\alpha x + \beta y), \phi(z) \rangle = 0$, true for any z , hence, ϕ is linear.

$$d_{\phi}^2(x, y) = \|x - y\|_{\phi}^2 = \langle \phi(x - y), \phi(x - y) \rangle = \langle \phi(x) - \phi(y), \phi(x) - \phi(y) \rangle = \|\phi(x) - \phi(y)\|_2^2 = \|\phi(x)\|_2^2 - 2\phi(x)^T \phi(y) + \|\phi(y)\|_2^2 = \|\phi(x)\|_2^2 - 2\langle \phi(x), \phi(y) \rangle + \|\phi(y)\|_2^2$$

$$D_{\phi}[i, j] = d_{\phi}^2(x_i, x_j) = \|\phi(x_i)\|_2^2 - 2\langle \phi(x_i), \phi(x_j) \rangle + \|\phi(x_j)\|_2^2$$

$$D_{\phi} = P - 2\Phi^T \Phi + P^T, \text{ where } P =$$

$$\begin{bmatrix} \|\phi(x_1)\|_2^2 \\ \|\phi(x_2)\|_2^2 \\ \vdots \\ \|\phi(x_N)\|_2^2 \end{bmatrix} \mathbf{1}_N^T$$

Hence, $-JD_{\phi}J = -J(P - 2\Phi^T \Phi + P^T)J$, from the lecture, $PJ = JP^T = 0$

Therefore, $-JD_{\phi}J = 2\Phi^T \Phi$, $-\frac{1}{2}JD_{\phi}J = JK_{\phi}J$

Consider a training set $\{\mathbf{x}_i\}_{i=1}^N$ and let $\mathbf{D} \in \mathbb{R}^{N \times N}$ where $\mathbf{D}[i, j] = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$.

1.2

Show that $\mathbf{v}^T \mathbf{D} \mathbf{v} < 0$ for any \mathbf{v} such that $\langle \mathbf{v}, \mathbf{1} \rangle = 0$.

.

1.2

$\langle \mathbf{v}, \mathbf{1} \rangle = 0$, and $\langle \mathbf{v}, \mathbf{1} \rangle = \mathbf{v}^T \mathbf{1}_N = \mathbf{1}_N^T \mathbf{v} = \langle \mathbf{1}, \mathbf{v} \rangle = 0$

Defining $\hat{\Lambda}$ as -

$$\begin{bmatrix} \|\mathbf{x}_1\|_2^2 \\ \|\mathbf{x}_2\|_2^2 \\ \vdots \\ \|\mathbf{x}_N\|_2^2 \end{bmatrix}$$

$$\mathbf{v}^T \mathbf{D} \mathbf{v} = \mathbf{v}^T (\Lambda - 2\mathbf{X}^T \mathbf{X} + \Lambda^T) \mathbf{v} = \mathbf{v}^T \Lambda \mathbf{v} - 2\mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} + \mathbf{v}^T \Lambda^T \mathbf{v} =$$

$$\langle \mathbf{v}, \Lambda \mathbf{v} \rangle - 2(\mathbf{X} \mathbf{v})^T \mathbf{X} \mathbf{v} + \langle \Lambda \mathbf{v}, \mathbf{v} \rangle = \langle \hat{\Lambda} \mathbf{1}_N^T \mathbf{v} \rangle - 2\|\mathbf{X} \mathbf{v}\|_2^2 - \langle \hat{\Lambda}_N^T \mathbf{v}, \mathbf{v} \rangle$$

$$\widehat{\Lambda}_N^T v = 0, \text{ hence, } -2\|Xv\|_2^2 \leq 0$$

MM (Majorization Minimization\Maximization)

Consider:

$$Y[i, j] = \begin{cases} X[i, j] & M[i, j] = 1 \\ 0 & M[i, j] = 0 \end{cases}$$

In other words:

$$Y = M \odot X$$

where $M \in \{0, 1\}^{M \times N}$ is a binary mask matrix.

Given $Y \in \mathbb{R}^{M \times N}$, the low-rank matrix completion objective is given by:

$$\begin{cases} \min_X \|M \odot (Y - X)\|_F^2 \\ \text{s.t.} \\ \text{rank}(X) \leq d \end{cases}$$

Consider the following function:

$$g(X, Z) := \|X - Z + M \odot (Z - Y)\|_F^2$$

1.3

Show that g surrogates the objective $f(X) := \|M \odot (Y - X)\|_F^2$.

Hint:

Show that $g(X, Z) = \|M \odot (X - Y) + \widetilde{M} \odot (X - Z)\|_F^2$ where $\widetilde{M} := 11^T - M$ is the complement of M .

1.3

$$g(X, Z) = \|X - Z + M \odot (Z - Y)\|_F^2 =$$

$$\|11^T \odot X - 11^T \odot Z + M \odot Z - M \odot Y + M \odot Y + M \odot X - M \odot X\|_F^2 =$$

$$\|M \odot (X - Y) - (11^T - M) \odot Z + (11^T - M) \odot X\|_F^2 =$$

$$\|M \odot (X - Y) - \overline{M} \odot Z + \overline{M} \odot X\|_F^2 =$$

$$\|M \odot (X - Y) + \overline{M} \odot (X - Z)\|_F^2$$

$$a. g(X, X) = \|M \odot (X - Y) + 0\|_F^2 = f(X)$$

$$b. g(X, X) = \|M \odot (X - Y) + 0\|_F^2 = \sum_{M[i,j]=1} \sqrt{(X - Y)_{ij}^2}$$

$$\text{Moreover, } Z = \overline{X}, g(X, \overline{X}) = \|M \odot (X - Y) + \overline{M} \odot (X - \overline{X})\|_F^2 =$$

$$\sum_{M[i,j]=1} \sqrt{(X - Y)_{ij}^2} + \sum_{M[i,j]=0} \sqrt{(X - \overline{X})_{ij}^2}$$

Now, \overline{M} is the complement of M , it is guaranteed that $f(X) \leq g(X, \overline{X})$ for all X and \overline{X} . Thus, $g(X, Z)$ is surrogate.

Metric MDS

The metric MDS objective is given by:

$$\min_{Z \in \mathbb{R}^{d \times N}} \|\Delta_x - D_z\|_F^2$$

where:

- $\Delta_x[i, j] = d(x_i, x_j)$ is a given distance matrix.
- $D_z[i, j] = \|z_i - z_j\|_2$.

Consider the surrogate function:

$$g(Z, \tilde{Z}) = \|\Delta_x\|_F^2 + 2N \text{Tr}\{ZJZ^T\} - 4\langle Z^T \tilde{Z}, B \rangle$$

where:

- $J = I - \frac{1}{N}\mathbf{1}\mathbf{1}^T$ is the centering matrix.
- $B = C - \text{diag}(C\mathbf{1})$
- $C[i, j] = \begin{cases} 0 & i = j \\ -\frac{\Delta_x[i, j]}{D_z[i, j]} & i \neq j \end{cases}$
- $\tilde{D}_z[i, j] = \|\tilde{z}_i - \tilde{z}_j\|_2$

1.4

Show that:

1.

$$BJ = B$$

2.

$$g(Z, Z) = \|\Delta_x - D_z\|_F^2$$

Notes: (See lecture slides)

1. $\|\Delta_x - D_z\|_F^2 = \|\Delta_x\|_F^2 + \|D_z\|_F^2 - 2\langle \Delta_x, D_z \rangle$
2. $\|D_z\|_F^2 = 2N \text{Tr}\{ZJZ^T\}$

Hint:

For $\tilde{Z} = Z$ we have:

$$\langle \Delta_x, D_z \rangle = -\langle C, D_z^{\circ 2} \rangle$$

$$\text{where } D_z^{\circ 2}[i, j] = p\mathbf{1}^T - 2Z^T Z + \mathbf{1}p^T \text{ and } p = \begin{bmatrix} \|z_1\|_2^2 \\ \vdots \\ \|z_N\|_2^2 \end{bmatrix}.$$

1.1) First let show that $BJ = B$

$$\begin{aligned} BJ &= (C - \text{diag}(C\mathbf{1}))J = (C - \text{diag}(C\mathbf{1}))\left(I_N - \frac{1}{N}\mathbf{1}_N\mathbf{1}_N^T\right) = \\ &= (C - \text{diag}(C\mathbf{1}))I_N - (C - \text{diag}(C\mathbf{1}))\frac{1}{N}\mathbf{1}_N\mathbf{1}_N^T \end{aligned}$$

We will assign:

$$\begin{aligned} B &= (C - \text{diag}(C\mathbf{1}))I_N \\ B' &= (C - \text{diag}(C\mathbf{1}))\frac{1}{N}\mathbf{1}_N\mathbf{1}_N^T \\ C' &= C - \text{diag}(C\mathbf{1}) \end{aligned}$$

Lets examine a general element ij:

$$(C' \mathbf{1}_N \mathbf{1}_N^T)_{ij} = \left(c'_{i1}, c'_{i2}, \dots, c'_{iN}\right)(1 \ 1 \ \dots \ 1) = c'_{i1} + c'_{i2} + \dots + c'_{iN} + c_{ii}$$

Note that

$$c'_{i1} + c'_{i2} + \dots + c'_{iN} + c_{ii} - \sum_{j=1}^N c_{ij} = 0$$

Because of that:

$$B' = \frac{1}{N}(C - \text{diag}(C_1))(\mathbf{1}_N \mathbf{1}_N^T)$$

Where $C - \text{diag}(C_1) = 0$ $B' = 0$

From the last point we can derive that:

$$BJ = B$$

Let show now that: $g(Z, Z) = \|\Delta_x - D_z\|_F^2$

Let's examine $\langle \Delta_x, D_z \rangle$

$$-\langle \Delta_x, D_z \rangle = \langle C, D_z^{\circ 2} \rangle = \langle C, \mathbf{p} \mathbf{1}_N^T - 2Z^T Z + \mathbf{1}_N \mathbf{p}^T \rangle$$

Similarly, to what we have saw in lecture:

$$\mathbf{p} = \text{diag}(Z^T Z) = \langle C, \text{diag}(Z^T Z) \mathbf{1}_N^T - 2Z^T Z + \mathbf{1}_N \text{diag}^T(Z^T Z) \rangle$$

From symmetric matrix W we know that:

$$\langle W, Y \rangle = \langle W, Y^T \rangle$$

Therefore:

$$\begin{aligned} \langle C, 2\text{diag}(Z^T Z) \mathbf{1}_N^T - 2Z^T Z \rangle &= \\ &= 2(\langle C, \text{diag}(Z^T Z) \mathbf{1}_N^T \rangle - \langle C, Z^T Z \rangle) \\ &= 2(\langle C \mathbf{1}_N, \text{diag}(Z^T Z) \rangle - \langle C, Z^T Z \rangle) \end{aligned}$$

From equation HW we know that: $\langle a, \text{diag}(X) \rangle = \langle \text{diag}(a), X \rangle$

$$= 2(\langle \text{diag}(C \mathbf{1}_N), Z^T Z \rangle - \langle C, Z^T Z \rangle) = -2 \langle B, Z^T Z \rangle$$

Now we can open the original equation:

$$\begin{aligned} \|\Delta_x - D_z\|_F^2 &= \|\Delta_x\|_F^2 + \|D_z\|_F^2 - 2 \langle \Delta_x, D_z \rangle = \\ &= \|\Delta_x\|_F^2 + 2N * \text{Tr}\{ZJZ^T\} + 2 \langle C, D_z^{\circ 2} \rangle = \\ &= \|\Delta_x\|_F^2 + 2N * \text{Tr}\{ZJZ^T\} - 4 \langle B, Z^T Z \rangle = g(Z, Z) \end{aligned}$$

2 Isomap

Let $G = (V, E, W)$ be a simple, undirected, and weighted graph, and assume no negative weights\edges. Let $\mathbf{D} \in \mathbb{R}^{N \times N}$ be the shortest path distance matrix, where $N = |V|$.

2.1

Prove or disprove:

Necessarily exists an embedding $\{\mathbf{z}_i \in \mathbb{R}^d\}_{i=1}^N$ (for some $d \in \mathbb{N}$) such that (for all i, j):

$$\mathbf{D}[i, j] = \|\mathbf{z}_i - \mathbf{z}_j\|_2$$

2.1

Assuming that \mathbf{D} , the shortest path distance matrix is computed based on graph G , where G is a graph constructed from training data $X \in \mathbb{R}^{D \times N}$.

Using MDS to solve the embeddings space Z , and for any training dataset $X \in \mathbb{R}^{D \times N}$ we can choose the embeddings vector size - d , to get the optimal solution Z^* that minimizes the problem.

For the optimal solution $Z^* = X$ -

$$D_x[i, j] = \|\mathbf{x}_i - \mathbf{x}_j\|_2 = \|\mathbf{z}_i^* - \mathbf{z}_j^*\|_2 = D_z^*[i, j]$$

- Let $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^N$ be the training set.
- Let $\mathcal{Z} = \{\mathbf{z}_i\}_{i=1}^N$ be the representation obtained by Isomap (training encoding).
- Consider a new point \mathbf{x}^* where $\mathbf{x}^* = \mathbf{x}_k$ for some $k \leq N$.
- Let \mathbf{z}^* be the out of sample encoding applied to \mathbf{x}^* .

2.2

Prove or disprove:

$$\mathbf{z}^* = \mathbf{z}_k$$

2.2

$$\mathbf{x}^* = \mathbf{x}_k \Rightarrow D_{xx} = D_{xz}$$

$$\begin{aligned} \bar{K}_{xz} &= -\frac{1}{2}J(D_{xz} - \frac{1}{N_x}D_{xx}\mathbf{1}_{N_x}\mathbf{1}_{N_z}^T) = -\frac{1}{2}J(D_{xx} - \frac{1}{N_x}D_{xx}\mathbf{1}_{N_x}\mathbf{1}_{N_x}^T) = -\frac{1}{2}JD_{xx}(\mathbf{I} - \frac{1}{N_x}\mathbf{1}_{N_x}\mathbf{1}_{N_x}^T) = \\ &= -\frac{1}{2}D_{xx}J = \bar{K}_{xx} \end{aligned}$$

$$\text{Hence, } Z_y = \sum_d^{-1} V_d^T \bar{K}_{xz} = \sum_d^{-1} V_d^T \bar{K}_{xx} = \sum_d^{-1} V_d^T V \Sigma^2 V^T = \Sigma_d V^T = Z$$

Z is the training encoding, hence, $\mathbf{z}^* = \mathbf{z}_k$

3 Laplacian Eigenmaps

- Consider $\mathcal{X} = \{\mathbf{x}_i \in \mathbb{R}^D\}_{i=1}^N$.
- Let $G = (V, E, W)$ be a weighted graph with $V = \mathcal{X}$ and:

$$W[i, j] = \begin{cases} \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right) & \mathbf{x}_i \in \mathcal{N}_j \text{ or } \mathbf{x}_j \in \mathcal{N}_i \\ 0 & \text{else} \end{cases}$$

- $e_{ij} \in E$ if $W[i, j] \neq 0$.
- Let $\mathbf{Z} \in \mathbb{R}^{d \times N}$ and $\mathbf{D}_z \in \mathbb{R}^{N \times N}$ such that $\mathbf{D}_z[i, j] = \|\mathbf{z}_i - \mathbf{z}_j\|_2^2$ where \mathbf{z}_i is the i th column of \mathbf{Z} .

3.1

Show that:

$$\frac{1}{2} \langle \mathbf{W}, \mathbf{D}_z \rangle = \text{Tr}\{\mathbf{Z}\mathbf{L}\mathbf{Z}^T\}$$

where:

- $\mathbf{L} = \mathbf{D} - \mathbf{W}$ is the graph-Laplacian.
- $\mathbf{D} = \text{diag}(\mathbf{W}\mathbf{1})$ is the degree matrix.

3.1)

Let's examine $\text{Tr}(\mathbf{Z}\mathbf{L}\mathbf{Z}^T)$

$$\text{Tr}(\mathbf{Z}\mathbf{L}\mathbf{Z}^T) = \sum_{i=1}^N (\mathbf{Z}\mathbf{L}\mathbf{Z}^T)_{ii} = \sum_{i=1}^N \mathbf{z}_i^T \mathbf{L} \mathbf{z}_i = (I)$$

Notice that \mathbf{z}_i is vector in \mathbb{R}^N and saw in lecture

For $v \in \mathbb{R}^N$, $v^T \mathbf{L} v = 0.5 \langle \mathbf{W}, \mathbf{D}_v \rangle$ and therefore:

$$(I) = 0.5 \sum_{i=1}^N \langle \mathbf{W}, \mathbf{D}_{\mathbf{z}_i} \rangle = 0.5 \langle \mathbf{W}, \mathbf{D}_z \rangle$$

Assume that G has two connected components, i.e. $V = V_1 \cup V_2$ such that:

$$\left\{ e_{ij} \mid i \in V_1, j \in V_2 \right\} = \emptyset$$

3.2

Show that the graph-Laplacian \mathbf{L} has two **orthogonal** eigenvectors corresponding to the zero eigenvalue. That is, exist $\mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^N$ such that:

1. $\mathbf{L}\mathbf{u}_1 = \mathbf{L}\mathbf{u}_2 = \mathbf{0}$
2. $\langle \mathbf{u}_1, \mathbf{u}_2 \rangle = 0$

3.2

Assuming that there are N points such that $|V| = N$, and $|V_1| = N_1$, $|V_2| = N_2$

D is diagonal, V_1 and V_2 are connected components, thus, $W_{ij} = 0$, for every $i \in V_1, j \in V_2$

Forming W as having first N_1 elements from V_1 , we get a block of size $N_1 \times N_1$ at the top

left corner that describes the weights of V_1 edges. Similarly, we get a block of size $N_2 \times N_2$

at the bottom right corner of W that describes the weights of V_2 edges.

The rest of W elements are 0. D is diagonal, hence, L has the same properties mentioned above.

Let's note the blocks on L as L_1, L_2 .

Now, we can find a u_1 that is constructed from N_1 1s and N_2 0s, and similarly, u_2 that is constructed by N_1 0s, and N_2 1s.

$\Rightarrow \langle u_1, u_2 \rangle = 0$, moreover, $L1 = 0$, thus

for each $0 \leq i \leq N_1$ $\sum_{j=0}^{N-1} L_{ij} = \sum_{j=0}^{N_1-1} L_{ij} u_1[j] = 0$, the multiplication for each of the first N_1 rows with u_1 is 0, the multiplication for each of the last N_2 rows is 0 by definition.

for each $N_1 \leq i \leq N$ $\sum_{j=0}^{N-1} L_{ij} = \sum_{j=N_1}^{N-1} L_{ij} u_2[j] = 0$, the multiplication for each of the first N_1 rows with u_1 is 0, the multiplication for each of the first N_1 rows is 0 by definition.

Thus, u_1, u_2 are eigenvectors with 0 eigenvalues.

4 t-SNE

The t-SNE objective is given by:

$$\min_{\mathbf{Z} \in \mathbb{R}^{d \times N}} \underbrace{D_{\text{KL}}(\mathbf{P}||\mathbf{Q})}_{:=f(\mathbf{Z})} = \min_{\mathbf{Z} \in \mathbb{R}^{d \times N}} \sum_{i=1}^N \sum_{j=1}^N p_{ij} \log \left(\frac{p_{ij}}{q_{ij}} \right)$$

See the definitions of \mathbf{P} and \mathbf{Q} in the lecture notes (do not get confused by the SNE definitions).
The goal of this question is to compute the gradient of the objective:

$$\nabla f(\mathbf{Z}) = ?$$

Let us break this task into several **smaller** steps.

4.1

Show that $f(\mathbf{Z}) = D_{\text{KL}}(\mathbf{P}||\mathbf{Q})$ can be written as:

$$f(\mathbf{Z}) = C - \langle \mathbf{P}, \log[\mathbf{Q}] \rangle$$

where C is some constant (the entropy of \mathbf{P}).

4.1) let's look at $f(z)$:

$$\begin{aligned} f(z) &= D_{\text{KL}}(P||Q) = \sum_{i=1}^N \sum_{j=1}^N p_{j|i} \log \left(\frac{p_{j|i}}{q_{j|i}} \right) \\ &= \sum_{i=1}^N \sum_{j=1}^N p_{j|i} \left(\log(p_{j|i}) - \log(q_{j|i}) \right) \\ &= \sum_{i=1}^N \sum_{j=1}^N p_{j|i} \log(p_{j|i}) - \sum_{i=1}^N \sum_{j=1}^N p_{j|i} \log(q_{j|i}) \\ &= C - \langle \mathbf{P}, \log[\mathbf{Q}] \rangle, \text{ where } C \text{ is the entropy of } \mathbf{P} \end{aligned}$$

4.2

Show that:

1.

$$B = \mathbf{1}^T (S - I) \mathbf{1} \in \mathbb{R}$$

2.

$$Q = B^{-1} (S - I) \in \mathbb{R}^{N \times N}$$

• Reminder:

$$Q[i, j] = \frac{1}{B} \begin{cases} 0 & i = j \\ (1 + \|z_i - z_j\|_2^2)^{-1} & i \neq j \end{cases}$$

• Let $D_z \in \mathbb{R}^{N \times N}$ such that $D_z[i, j] = \|z_i - z_j\|_2^2$.

• Let $S = (\mathbf{1}\mathbf{1}^T + D_z)^{0=1} \in \mathbb{R}^{N \times N}$, that is:

$$S[i, j] = (1 + D_z[i, j])^{-1}$$

4.2.1)

- in the lecture we seen the formula for B

$$\begin{aligned} B &= \sum_{i=1}^N \sum_{i \neq j}^N (1 + \|z_i - z_j\|)^{-1} \\ &= \sum_{i=1}^N \sum_{i \neq j}^N (1 + \|z_i - z_j\|)^{-1} \pm \sum_{i=1}^N \sum_{i=j}^N (1 + \|z_i - z_j\|)^{-1} \\ &\quad \sum_{i=1}^N \sum_{j=1}^N (1 + \|z_i - z_j\|)^{-1} - N \\ &= \mathbf{1}^T S \mathbf{1} \mathbf{1} - \mathbf{1}^T I \mathbf{1} = \mathbf{1}^T (S - I) \mathbf{1} \end{aligned}$$

- notice that:

- $\mathbf{1} \in \mathbb{R}^{N \times 1}$
- $S \in \mathbb{R}^{N \times N}$
- $B \in \mathbb{R}^{1 \times 1}$

4.2.2)

- Let look about W , where $W \in \mathbb{R}^{N \times N}$

$$W = S - I$$

- Let look on 2 cases

- $i = j$:

$$\begin{aligned} W[i, j] &= S - I[i, i] = S[i, i] - I[i, i] \\ &= (1 + \|z_i - z_j\|)^{-1} - 1 = 1 - 1 = 0 \end{aligned}$$

- therefore, we can infer that:

$$W = BQ \rightarrow Q = \frac{1}{B} W = \frac{1}{B} (S - I)$$

- because $W \in \mathbb{R}^{N \times N}$ multiplying by scalar does not change its dimension,

so $Q \in \mathbb{R}^{N \times N}$.

- $i \neq j$

$$\begin{aligned} W[i, j] &= S - I[i, j] = S[i, j] - I[i, j] \\ &= (1 + \|z_i - z_j\|)^{-1} - 0 = (1 + \|z_i - z_j\|)^{-1} \end{aligned}$$

4.3

Show that:

$$-\langle \mathbf{P}, \log [\mathbf{Q}] \rangle = \log (B) + \langle \mathbf{P}, \log [\mathbf{11}^T + \mathbf{D}_z] \rangle$$

Hints:

- $\mathbf{P}[i, i] = ?$
- $\mathbf{1}^T \mathbf{P} \mathbf{1} = ?$

4.3)

- Let look at $-\langle \mathbf{P}, \log [\mathbf{Q}] \rangle$:

$$\begin{aligned} -(\langle \mathbf{P}, \log \log [\mathbf{Q}] \rangle &= -\langle \mathbf{P}, \log \log [\mathbf{B}^{-1}(\mathbf{S} - \mathbf{I})] \rangle) \\ &= -(\langle \mathbf{P}, \log \log [(\mathbf{S} - \mathbf{I})] - \log \log [\mathbf{B} \mathbf{11}^T] \rangle) \\ &= -(\langle \mathbf{P}, \log \log [(\mathbf{S} - \mathbf{I})] \rangle - \langle \mathbf{P}, \log \log [\mathbf{B} \mathbf{11}^T] \rangle) \\ &= -(\langle \mathbf{P}, \log \log [(\mathbf{S} - \mathbf{I})] \rangle - \sum_{i=1}^N \sum_{j=1}^N p_{ji} \log(B)) \\ &= -(\langle \mathbf{P}, \log \log [(\mathbf{S} - \mathbf{I})] \rangle - \log(B) \sum_{i=1}^N \sum_{j=1}^N p_{ji}) \\ &= -(\langle \mathbf{P}, \log \log [(\mathbf{S} - \mathbf{I})] \rangle - \log \log (B) \mathbf{1}) \\ &= -\langle \mathbf{P}, \log \log [(\mathbf{S} - \mathbf{I})] \rangle + \log \log (B) = (*) \end{aligned}$$

- Now let examine $P_{i,i} \log[S - I]_{i,j}$, for 2 cases:

- $i \neq j$

$$P_{i,j} \log[S - I]_{i,j} = P_{i,j} \log[S]_{i,j}$$

- We can conclude that:

$$\langle \mathbf{P}, \log \log [\mathbf{S} - \mathbf{I}] \rangle = \langle \mathbf{P}, \log [\mathbf{S}] \rangle$$

- Now let return to the (*)

$$\begin{aligned} (*) &= -\langle \mathbf{P}, \log \log [(\mathbf{S} - \mathbf{I})] \rangle + \log \log (B) \\ &= -\langle \mathbf{P}, \log \log [(\mathbf{S})] \rangle + \log \log (B) \\ &= -\langle \mathbf{P}, \log \log \left[\left(\mathbf{11}^T + \mathbf{D}_z \right)^{\circ-1} \right] \rangle + \log \log (B) \end{aligned}$$

- Now let use the following rule:

$$\log \log (a^b) = b \log(a)$$

- Then we yield:

$$= \langle \mathbf{P}, \log \log \left[\left(\mathbf{11}^T + \mathbf{D}_z \right)^{\circ-1} \right] \rangle + \log \log (B)$$

- $i = j$

$$P_{i,i} \log[S - I]_{ii} = 0 * \log \log [S - I]_{ii} = 0 = P_{i,i} \log[S]_{i,i}$$

Let:

$$f(\mathbf{Z}) = C + \underbrace{\log(B)}_{(*)} + \underbrace{\langle \mathbf{P}, \log[\mathbf{1}\mathbf{1}^T + \mathbf{D}_z] \rangle}_{(**)}$$

4.4

Show that:

1.

$$\nabla_{\mathbf{Z}} \underbrace{\langle \mathbf{P}, \log[\mathbf{1}\mathbf{1}^T + \mathbf{D}_z] \rangle}_{(**)}[\mathbf{H}] = \langle \mathbf{S} \circ \mathbf{P}, \nabla \mathbf{D}_z[\mathbf{H}] \rangle$$

2.

$$\nabla_{\mathbf{Z}} \underbrace{\log(B)}_{(*)}[\mathbf{H}] = -\langle \mathbf{S} \circ \mathbf{Q}, \nabla \mathbf{D}_z[\mathbf{H}] \rangle$$

Hints:

- $\nabla \mathbf{S}[\mathbf{H}] = \nabla (\mathbf{1}\mathbf{1}^T + \mathbf{D}_z)^{\circ-1}[\mathbf{H}] = -(\mathbf{1}\mathbf{1}^T + \mathbf{D}_z)^{\circ-2} \circ \nabla(\mathbf{D}_z)[\mathbf{H}] = -\mathbf{S} \circ \mathbf{S} \circ \nabla(\mathbf{D}_z)[\mathbf{H}]$
- $\mathbf{Q} = \mathbf{B}^{-1}(\mathbf{S} - \mathbf{I})$

4.4.0.1) help prove for 4.4.1

- If D is diagonal matrix, and we want to calculate the product $\langle \mathbf{D}, \mathbf{D}_z \rangle = f(z) = ?$

$$L(z) \langle \mathbf{D}, \mathbf{D}_z \rangle = \text{Tr}(\mathbf{D}^T \mathbf{D}_z) = \sum_1^N D_{ii} D_z[i, i] = \sum_1^N D_{ii} D_z[i, i] = \sum_1^N D_{ii} 0 = 0$$

- Moreover, we know from that that $\nabla \langle \mathbf{D}, \mathbf{D}_z \rangle = 0$

4.4.0.2) Let express $\nabla L(z)$

$$\nabla L(z) = \nabla \langle \mathbf{D}, \text{diag}(\mathbf{Z}^T \mathbf{Z}) \mathbf{1}^T - 2\mathbf{Z}^T \mathbf{Z} + \mathbf{1}(\text{diag}^T(\mathbf{Z}^T \mathbf{Z})) \rangle$$

○ We saw in last homework that $\mathbf{Z} \in \mathbb{R}^{d \times N}$, $\mathbf{D} \in \mathbb{R}^{N \times N}$

$$\nabla \langle \mathbf{D}, \text{diag}(\mathbf{Z}^T \mathbf{Z}) \mathbf{1}^T - 2\mathbf{Z}^T \mathbf{Z} + \mathbf{1}(\text{diag}^T(\mathbf{Z}^T \mathbf{Z})) \rangle = -2 \langle \mathbf{D} - \text{diag}(\mathbf{D}\mathbf{1}), \mathbf{Z}^T \mathbf{Z} \rangle$$

- Let calculate the gradient of $\mathbf{Z}^T \mathbf{Z}$:

$$\nabla_g(\mathbf{Z})[\mathbf{H}] = \nabla(\mathbf{Z}^T \mathbf{Z})[\mathbf{H}] = \frac{\mathbf{Z}^T \mathbf{Z} + t\mathbf{H}^T \mathbf{Z} + t\mathbf{Z}^T \mathbf{H} + t^2 \mathbf{H}^T \mathbf{H} - \mathbf{Z}^T \mathbf{Z}}{t}$$

$$\mathbf{H}^T \mathbf{Z} + \mathbf{Z}^T \mathbf{H} + t\mathbf{H}^T \mathbf{H} = t\mathbf{H}^T \mathbf{H} + \mathbf{H}^T \mathbf{Z} + \mathbf{Z}^T \mathbf{H} = \mathbf{H}^T \mathbf{Z} + \mathbf{Z}^T \mathbf{H}$$

- $\nabla \langle \mathbf{D}, \text{diag}(\mathbf{Z}^T \mathbf{Z}) \mathbf{1}^T - 2\mathbf{Z}^T \mathbf{Z} + \mathbf{1}(\text{diag}^T(\mathbf{Z}^T \mathbf{Z})) \rangle = -2 \langle \mathbf{D} - \text{diag}(\mathbf{D}\mathbf{1}), \mathbf{Z}^T \mathbf{Z} \rangle$

$$\nabla L(z)[\mathbf{H}] = -2 \langle \mathbf{0}, \mathbf{Z}^T \mathbf{Z} \rangle + \langle \mathbf{D} - \text{diag}(\mathbf{D}\mathbf{1}), \nabla(\mathbf{Z}^T \mathbf{Z})[\mathbf{H}] \rangle$$

$$= -2 \langle \mathbf{D} - \text{diag}(\mathbf{D}\mathbf{1}), \mathbf{H}^T \mathbf{Z} + \mathbf{Z}^T \mathbf{H} \rangle$$

$$= -2 \langle \mathbf{D} - \text{diag}(\mathbf{D}\mathbf{1}), \mathbf{H}^T \mathbf{Z} \rangle + \langle \mathbf{D} - \text{diag}(\mathbf{D}\mathbf{1}), \mathbf{Z}^T \mathbf{H} \rangle$$

- For 2 NXN matrices A, B we know that:

$$\langle \mathbf{A}, \mathbf{B} \rangle = \langle \mathbf{A}^T, \mathbf{B}^T \rangle$$

- Now let's use this fact

$$-2 \langle \mathbf{D}^T - \text{diag}^T(\mathbf{D}\mathbf{1}), \mathbf{Z}^T \mathbf{H} \rangle + \langle \mathbf{D} - \text{diag}(\mathbf{D}\mathbf{1}), \mathbf{Z}^T \mathbf{H} \rangle$$

- Because D is symmetric $\mathbf{D} = \mathbf{D}^T$, and $\text{diag}^T(\mathbf{D}\mathbf{1}) = \text{diag}(\mathbf{D}\mathbf{1})$

$$= -4 \langle \mathbf{D} - \text{diag}(\mathbf{D}\mathbf{1}), \mathbf{Z}^T \mathbf{H} \rangle = 4\mathbf{Z} \langle \mathbf{D}\mathbf{1} - \mathbf{D}, \mathbf{H} \rangle$$

- Therefore $\nabla L(z) =$

$$\nabla L(z) = 4Z(\text{diag}(D1) - D)$$

4.4.0.2) help prove for 4.4.2

- If D is diagonal matrix, and we want to calculate the product $\langle D, D_z \rangle = ?$

$$\langle D, D_z \rangle = \text{Tr}(D^T D_z) = \sum_1^N D_{ii} D_z[i, i] = \sum_1^N D_{ii} D_z[i, i] = \sum_1^N D_{ii} 0 = 0$$

- Moreover, we know from that that $\nabla \langle D, D_z \rangle = 0$

4.4.1)

$$\nabla \langle P, \log \log [11^T + D_z] \rangle [H] = ?$$

- Let's use the product rule first. And then the chain rule

- product rule:

$$\begin{aligned} \langle P[H], \log[11^T + D_z] \rangle + \langle P, \nabla(\log[11^T + D_z])[H] \rangle \\ = 0 + \langle P, \nabla(\log[11^T + D_z])[H] \rangle \end{aligned}$$

- chain rule:

$$\langle P, \langle [11^T + D_z]^{-1}, \nabla D_z[H] \rangle \rangle = \langle P, S \circ \nabla D_z[H] \rangle$$

- now let use the first hint:

$$(P \circ S, \nabla(D_z)[H])$$

- missing part 1.4 ...

- now let use the second hint + the knowledge that S, P are symmetric + 4.4.0.2:

$$4Z \langle \text{diag}(P \circ S) 1 \rightarrow - (P \circ S), H \rangle$$

- therefore we know that:

$$\nabla \langle P, \log \log [11^T + D_z] \rangle [H] = 4Z(\text{diag}(P \circ S) 1 \rightarrow - (P \circ S))$$

4.4.2)

$$\nabla \log \log (B) = \log \log (1^T (S - I) 1) = ?$$

- Let's look on $\nabla B[H]$:

$$\begin{aligned} \nabla B[H] &= \nabla (1^T (S - I) 1)[H] = \nabla (1^T S 1 - 1^T I 1)[H] \\ &= \nabla (1^T S 1)[H] - \nabla (1^T I 1)[H] = \nabla (1^T S 1)[H] - (0) = \nabla (1^T S 1)[H] \end{aligned}$$

- know by using multiple time the product rule:

$$\begin{aligned} (1^T [H] S 1 + 1^T (\nabla S 1)[H]) &= 0 - 1^T (\nabla S 1)[H] = 1^T (\nabla S 1)[H] \\ &= 1^T ((\nabla S[H]) 1 + S(\nabla 1[H])) = 1^T ((\nabla S[H]) 1 + 0) = 1^T \nabla S[H] 1 \end{aligned}$$

- now lets use the the first hint:

$$\begin{aligned} -1^T (S \circ S \circ \nabla (D_z)[H] 1) &= -\langle 1 1^T, S \circ S \circ \nabla (D_z)[H] \rangle \\ &= -\langle S \circ S, \nabla (D_z)[G] \rangle \end{aligned}$$

- Let's do trick of adding and subtracting I from S:

$$-\langle (S - I + I) \circ S, \nabla (D_z)[H] \rangle = -\langle ((S - I) \circ S, \nabla (D_z)[H] \rangle - \langle (I \circ S, \nabla (D_z)[H] \rangle$$

- $-\langle (I \circ S, \nabla (D_z)[H] \rangle = 0$ follow(**) from 4.4.0.1

$$-\langle ((S - I) \circ S, \nabla (D_z)[H] \rangle - 0 = -\langle ((S - I) \circ S, \nabla (D_z)[H] \rangle$$

- Finally, we can show the results of $\nabla \log \log (B)[H] = ?$

$$\begin{aligned} \nabla \log \log (B) [H] &= \frac{1}{B} \nabla (B)[H] = -\frac{1}{B} \langle (S - I) \circ S, \nabla (D_z)[H] \rangle \\ &= -\langle Q \circ S, \nabla (D_z)[H] \rangle \\ &= -4Z(\text{diag}(Q \circ S 1) - Q \circ S) \end{aligned}$$

4.5

- Combine all previous results and write the gradient of the objective:

$$\nabla f(\mathbf{Z}) = ?$$

- Use $\mathbf{A} := (\mathbf{P} - \mathbf{Q}) \circ \mathbf{S}$ to simplify your answer.
- What can you say about the gradient $\nabla f(\mathbf{Z})$ when $\mathbf{P} = \mathbf{Q}$?

Hint: Use the lecture notes.

4.5)

$$\begin{aligned} \nabla f(\mathbf{Z}) &= \nabla (C - \langle \mathbf{P}, \log \log [\mathbf{Q}] \rangle) = \nabla \left(\log \log (\mathbf{B}) + \langle \mathbf{P}, \log \log \left[\mathbf{1}\mathbf{1}^T + \mathbf{D}_z \right] \rangle \right) \\ &= -4\mathbf{z}(\text{diag}(\mathbf{Q} \circ \mathbf{S}\mathbf{1}) - \mathbf{Q} \circ \mathbf{S}) + 4\mathbf{Z}(\text{diag}(\mathbf{P} \circ \mathbf{S}\mathbf{1}) - \mathbf{P} \circ \mathbf{S}) \\ 4\mathbf{z}(\text{diag}(\mathbf{P} - \mathbf{Q}) \circ \mathbf{S}\mathbf{1} - (\mathbf{P} - \mathbf{Q}) \circ \mathbf{S}\mathbf{1}) &= 4\mathbf{Z}(\text{diag}(\mathbf{A}\mathbf{1}) - \mathbf{A}) \end{aligned}$$

- notice that in this case $\mathbf{P} = \mathbf{Q}$
 - gradient will be 0 because \mathbf{A} will be 0
 - the same for $\text{diag}(\mathbf{A}\mathbf{1})$