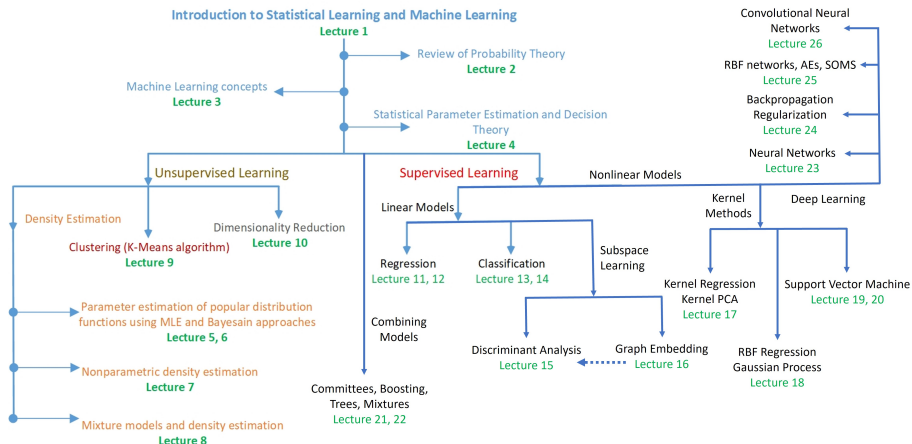


Statistical Learning and Machine Learning

Lecture 16 - Graph-based Learning

October 13, 2021

Course overview and where do we stand



Definitions

Let us consider a graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$:

- $\mathcal{V} = \{v_1, \dots, v_N\}$ is called the *vertex set* of the graph
- $\mathcal{E} = \{e_{ij}\}$ is called the *edge set* of the graph
 - e_{ij} is a edge connecting v_i and v_j
 - e_{ij} can indicate a link between the i -th and the j -th vertices (*unweighted graph*)
 - e_{ij} can be associated with a weight value w_{ij} encoding how strong is the link between the i -th and the j -th vertices (*weighted graph*)

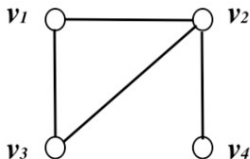
The graph can be:

- fully connected: there is an edge between every pair of vertices (and the associated weight)
- partially connected: each vertex is connected only with some other vertices (those which are more similar to it according to a similarity/distance metric)

Definitions

The *adjacency* matrix $A \in \mathbb{R}^{N \times N}$ of the graph is defined as:

$$A := \begin{cases} A_{ij} = 1, & \text{if there is an edge } e_{ij} \\ A_{ij} = 0, & \text{if there is no edge} \\ A_{ii} = 0 \end{cases} \quad (1)$$



$$A = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

The weight matrix $W \in \mathbb{R}^{N \times N}$ of the graph has the same form, but each element encodes the strength of connection (similarity) between the corresponding vertices.

Laplacian Embedding

Consider a value x_i associated with vertex v_i of an undirected weighted graph:

- We want the relation between the values associated with the graph vertices to be the same as those expressed by the corresponding weights w_{ij}
- Example of graph weights (Gaussian kernel):

$$w_{ij} = \exp \left(-\frac{(x_i - x_j)^2}{\sigma^2} \right) \quad (2)$$

$$0 \leq w_{min} \leq w_{ij} \leq w_{max} \leq 1. \quad (3)$$

This means that if w_{ij} is high, the representations of vertices v_i and v_j should be close using the Euclidean distance

- This means that the vertex representations should be those minimizing:

$$\mathcal{J} = \frac{1}{2} \sum_{e_{ij}} w_{ij} (x_i - x_j)^2 \quad (4)$$

Laplacian Embedding

The minimization criterion is:

$$\mathcal{J} = \frac{1}{2} \sum_{e_{ij}} w_{ij} (x_i - x_j)^2 = \mathbf{x}^T \mathbf{L} \mathbf{x} \quad (5)$$

where $\mathbf{L} = \mathbf{D} - \mathbf{W} \in \mathbb{R}^{N \times N}$ is the *Laplacian* matrix.

$\mathbf{D} \in \mathbb{R}^{N \times N}$ is a diagonal matrix (called *Degree* matrix) having elements:

$$D_{ii} = \sum_j W_{ij} \quad (6)$$

We introduce two constraints:

$$\mathbf{x}^T \mathbf{x} = 1 \quad \text{and} \quad \mathbf{x}^T \mathbf{1} = 0 \quad (7)$$

The solution is obtained by solving the eigenanalysis problem:

$$\mathbf{L} \mathbf{x} = \lambda \mathbf{x} \quad (8)$$

Graph Embedding

Graph Embedding uses a graph $\mathcal{G} = \{\mathcal{V}, W, X\}$:

- $\mathcal{V} = \{v_1, \dots, v_N\}$ is the vertex set of the graph
- $W \in \mathbb{R}^{N \times N}$ is the weight matrix the elements of which encode the similarity between vertex pairs
- $X = [x_1, \dots, x_N] \in \mathbb{R}^{D \times N}$ is a set of D - dimensional vectors with x_i representing v_i

Graph Embedding defines two graphs:

- an (*intrinsic*) graph \mathcal{G}^I expressing properties of the data that we want the embedding to enhance
- a *penalty* graph \mathcal{G}^P expressing properties of the data that we want the embedding to penalize

Graph Embedding

Suppose that:

- $L \in \mathbb{R}^{N \times N}$ is the Laplacian matrix of \mathcal{G}^I
- $B \in \mathbb{R}^{N \times N}$ is the Laplacian matrix of \mathcal{G}^P

Graph Embedding optimizes the following criterion:

$$y^* = \arg \min_{y^T B y = q} \sum_{i \neq j} (y_i - y_j)^2 W_{ij} = \arg \min_{y^T B y = q} y^T L y \quad (9)$$

where:

- q is a constant value
- $y_i \in \mathbb{R}$ is the value associated with vertex v_i
- $y = [y_1, \dots, y_N]^T$

Graph Embedding

If we define a linear mapping:

$$y_i = w^T x_i \quad (10)$$

then the criterion of obtaining the projection matrix $w \in \mathbb{R}^D$ is:

$$\begin{aligned} w^* &= \arg \min_{\substack{w^T X B X^T w = q \\ \text{or } w^T w = q}} \sum_{i \neq j} (w^T x_i - w^T x_j)^2 W_{ij} \\ &= \arg \min_{\substack{w^T X B X^T w = q \\ \text{or } w^T w = q}} w^T X L X^T w \end{aligned} \quad (11)$$

w^* is calculated by solving the *generalized eigenanalysis* problem:

$$\tilde{L} w = \lambda \tilde{B} w \quad (12)$$

where:

$$\tilde{L} = X L X^T \quad \text{and} \quad \tilde{B} = X B X^T \quad (13)$$

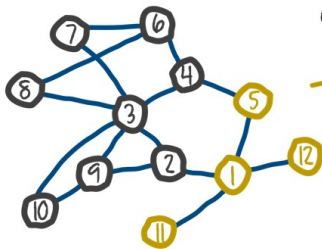
Graph Embedding

To define more than one dimensions $y \in \mathbb{R}^d$ we solve again the same generalized eigenanalysis problem:

$$\tilde{L}w = \lambda \tilde{B}w \quad (14)$$

and we keep the eigenvectors corresponding to the d smallest eigenvalues.

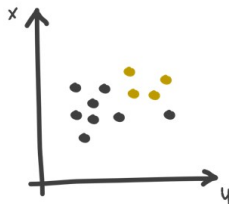
from a graph representation ...



embedding
algorithm



to real vector representation



Let us take again a look at the within-class scatter matrix of LDA:

$$S_W = \sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \quad (15)$$

Using:

- $\mathbf{1}_k \in \mathbb{R}^N$ a vector having elements:

$$1_k(i) := \begin{cases} 1, & \text{if } \mathbf{x}_i \in \mathcal{C}_k \\ 0, & \text{otherwise} \end{cases} \quad (16)$$

- $\mathbf{J}_k = \text{diag}(\mathbf{1}_k)$

Graph Embedding: LDA

Let us take again a look at the within-class scatter matrix of LDA:

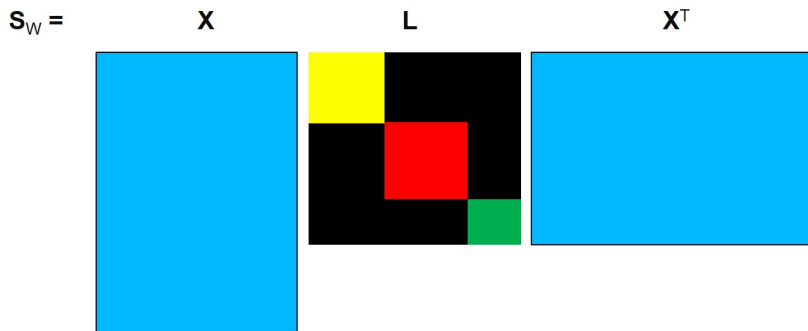
$$S_W = \sum_{k=1}^K \sum_{i, l_i=k} (x_i - \mu_k)(x_i - \mu_k)^T \quad (17)$$

$$\begin{aligned} S_W &= \sum_{k=1}^K \left(X J_k - \frac{1}{N_k} X 1_k 1_k^T \right) \left(X J_k - \frac{1}{N_k} X 1_k 1_k^T \right)^T \\ &= \sum_{k=1}^K X \left(J_k - \frac{1}{N_k} 1_k 1_k^T \right) \left(J_k - \frac{1}{N_k} 1_k 1_k^T \right)^T X^T \\ &= X \left(\sum_{k=1}^K \left(J_k - \frac{1}{N_k} 1_k 1_k^T \right) \left(J_k - \frac{1}{N_k} 1_k 1_k^T \right)^T \right) X^T \\ &= X L X^T \end{aligned} \quad (18)$$

Graph Embedding: LDA

Let us take again a look at the within-class scatter matrix of LDA:

$$S_W = \sum_{k=1}^K \sum_{i, l_j=k} (x_i - \mu_k)(x_i - \mu_k)^T \quad (19)$$



Black color corresponds to pairs of data points belonging to different classes (these elements take zero values).

Graph Embedding: LDA

Let us take again a look at the between-class scatter matrix of LDA:

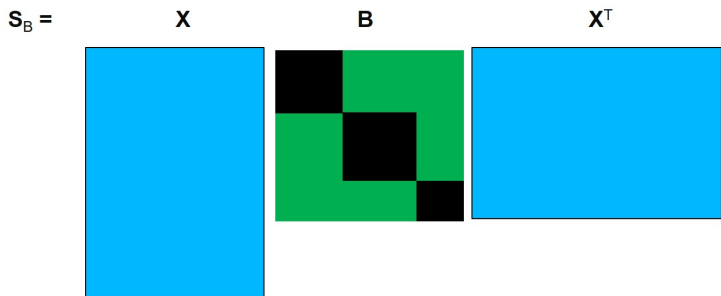
$$S_B = \sum_{k=1}^K N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T \quad (20)$$

$$\begin{aligned} S_B &= \sum_{k=1}^K N_k \left(\frac{1}{N_k} \mathbf{X} \mathbf{1}_k - \frac{1}{N} \mathbf{X} \mathbf{1} \right) \left(\frac{1}{N_k} \mathbf{X} \mathbf{1}_k - \frac{1}{N} \mathbf{X} \mathbf{1} \right)^T \\ &= \sum_{k=1}^K N_k \mathbf{X} \left(\frac{1}{N_k} \mathbf{1}_k - \frac{1}{N} \mathbf{1} \right) \left(\frac{1}{N_k} \mathbf{1}_k - \frac{1}{N} \mathbf{1} \right)^T \mathbf{X}^T \\ &= \mathbf{X} \left(\sum_{k=1}^K N_k \left(\frac{1}{N_k} \mathbf{1}_k - \frac{1}{N} \mathbf{1} \right) \left(\frac{1}{N_k} \mathbf{1}_k - \frac{1}{N} \mathbf{1} \right)^T \right) \mathbf{X}^T \\ &= \mathbf{X} \mathbf{B} \mathbf{X}^T \end{aligned} \quad (21)$$

Graph Embedding: LDA

Let us take again a look at the between-class scatter matrix of LDA:

$$S_B = \sum_{k=1}^K N_k (\mu_k - \mu)(\mu_k - \mu)^T \quad (22)$$



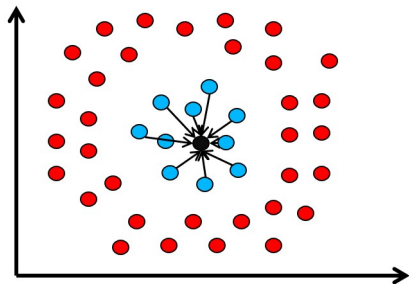
Black color corresponds pairs of data points belonging to the same class (these elements take values equal to $(1/N_k - 1/N)^2$).

Graph Embedding: Other types of graphs

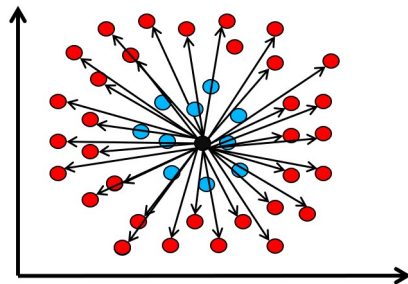
Using the Graph Embedding procedure, we can define subspace learning methods optimize several types of criteria.

Example: Class-Specific Discriminant Analysis using intra-class variance S_i and out-of-class variance matrix S_p .

Intra-class scatter S_i



Out-of-class scatter S_p



Graph Embedding: Generic Graph types

Marginal Discriminant Analysis:

- For each data point find its K_w nearest neighbors considering only the data points belonging to the same class (call the neighborhood of x_i as \mathcal{N}_{x_i}). Create the intrinsic graph weight matrix W_I having elements:

$$W_{I,ij} := \begin{cases} 1, & \text{if } x_j \in \mathcal{N}_{x_i} \text{ and } x_i \in \mathcal{N}_{x_j} \\ 0, & \text{otherwise} \end{cases} \quad (23)$$

- For each pair of classes \mathcal{C}_k and \mathcal{C}_l find the K_b data points that correspond to the K_w smallest (between-class) distances include them in the set \mathcal{N}_{kl} . Create the penalty graph weight matrix W_P having elements:

$$W_{P,ij} := \begin{cases} 1, & \text{if } x_i \text{ and } x_j \in \mathcal{N}_{kl} \\ 0, & \text{otherwise} \end{cases} \quad (24)$$

- Calculate the Laplacian matrices $\tilde{L} = D_I - W_I$ and $\tilde{B} = D_I - W_P$
- Solve for $\tilde{L}w = \lambda \tilde{B}w$ and we keep the eigenvectors corresponding to the d smallest eigenvalues

Spectral Clustering

Spectral Clustering exploits a data transformation based on the spectrum (eigenanalysis) of the graph Laplacian matrix for applying nonlinear clustering.

Algorithm:

- 1 Construct a fully connected graph $\mathcal{G} = \{\mathcal{V}, W, X\}$ using $x_i \in \mathbb{R}^D, i = 1, \dots, N$
- 2 Calculate the graph Laplacian matrix $L = D - W$
- 3 Apply eigenanalysis to L and keep the eigenvectors corresponding to the d smallest eigenvalues to form new d -dimensional data representations $y_i \in \mathbb{R}^d$
- 4 Apply a clustering algorithm (e.g. K -Means) using $y_i, i = 1, \dots, N$
- 5 Assign the cluster labels to the original data $x_i, i = 1, \dots, N$

Spectral Clustering

