

Lecture 17: Sem.-supervised learning on graphs

We introduced the idea of diffusion maps last lecture. The mathematical background & analysis, specially in connection to diffusion processes & elliptic operators is quite deep & contemporary. (see ref. on canvas)

Ignoring the mathematical details, the construction of diffusion maps is quite simple:

Step 1) Construct Graph $G = \{X, W\}$

Step 2) Compute eigen pairs $\{\lambda_j, q_j\}$ of $D^{-1}W$

(equivalently eigen pairs $\{\alpha_j, q_j\}$ of L_{RW} set $\lambda_j = 1 - \alpha_j$)

Step 3) Define diff. map

$$f_t: X \rightarrow \mathbb{R}^n, \quad f_t(x_i) = (\lambda_1^t q_{1(i)}, \dots, \lambda_n^t q_{n(i)})^T$$

Of course, thinking about diff. maps this way we immediately see that the choice of L_{RW} is not crucial & we can define such maps for other graph Laplacians.

Let $L = I - A$ be a graph Laplacian with an arbitrary adjacency matrix (normalized, RW, etc) with eigen pairs $\{\lambda_j, q_j\}$, set $\tilde{\lambda}_j = 1 - \lambda_j$. Define

$$\mathcal{F}_t(x_i) = (\tilde{\lambda}_1^{2t} q_1(i), \dots, \tilde{\lambda}_n^{2t} q_n(i))$$

* Understanding whether your choice of A leads to a diff. process is an interesting theoretical question.

Now observe that given a diffusion map $\tilde{\mathcal{F}}_t$ we can define a kernel

$$K_t: X \times X \rightarrow \mathbb{R}, \quad K_t(x_i, x_l) = \sum_{j=1}^n \tilde{\lambda}_j^{2t} q_j(i) q_j(l)$$

Recall $|X| = n$ Under some technical assumptions, we can show that if $x_i \stackrel{iid}{\sim} U(M)$ (some compact smooth manifold) then $\lim_{n \rightarrow \infty} K_t^n \rightarrow K_t$ where K_t is a "Heat kernel" ie, it describes a diffusion process on our manifold M .

This suggests, the RKHS of kernels such as above, are good candidates for designing alg. that can automatically take advantage of an existing manifold M on which the data lies.

(2)

17.1 Semi-supervised learning with manifold regularization

regularization

We will now use the above intuition to introduce another approach to exploiting spectral geometric information of graph Laplacians in the context of semi-supervised Learning (SSL)

Supervised learning

Training data $\{x_i, y_i\}_{i=1}^N$

Predict $y_{N+1}(x_{N+1})$



Unsupervised learning

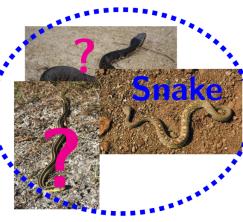
Given unlabeled data

$\{x_i\}_{i=1}^N$, discover structure
in X .



Given $\{x_i, y_i\}_{i=1}^N$ & $\{x_i\}_{N+1}^M$

Predict $y_i(x_i)$ for $i=N+1, \dots, M$



(3)

explore clusters \uparrow
& labels simultaneously

We will introduce an approach for SSL based on the so called manifold regularization technique of Belkin et al (2006) which is indeed nothing more than RKHS regression.

Let $X = \{\underline{x}_1, \dots, \underline{x}_M\} \subset M \subset \mathbb{R}^d$ & suppose we are given (possibly noisy) labels $\mathbf{y} = (y_1, \dots, y_N) \in \{-1, +1\}^N$ where $N < M$ (typically we would like $N \ll M!$)

Goal: Predict $f(\underline{x}_i)$ for $i=N+1, \dots, M$

First, we need a model for the labels & a data misfit term as we did before in classification. For this we will use probit.

Recall:

$$(\text{model}) \quad y_i = \text{Sign}(f(\underline{x}_i) + \varepsilon_i), \quad \varepsilon_i \stackrel{iid}{\sim} \psi$$

$$(\text{Probit misfit}) \quad L(f; \mathbf{y}) = - \sum_{j=1}^N \log \Psi(y_j f(\underline{x}_j))$$

Since our goal is to exploit possible clustering/geometric info in X to get better

labelling it is natural for us to build a graph on X .

Step 1) Build graph $G = \{X, W\}$

(prox. or k -NN graph with weight function
 $\eta: [0, \infty) \rightarrow [0, \infty)$ defined by user)

e.g. $\eta(t) = \exp\left(\frac{-t^2}{2r^2}\right)$ & set $w_{ij} = \eta(\|\underline{x}_i - \underline{x}_j\|_2)$

Note, the choice of the radius $r > 0$ is crucial in getting a good graph.

A surprisingly good heuristic:

In most applications we simply compute the pairwise distances $p_{ij} = \|\underline{x}_i - \underline{x}_j\|_2$ & simply set r to be the $1/25$ quantile of the p_{ij} .

This can be further fine tuned by CV if needed.

Step 2) Define a graph Laplacian L on G
(your fav. normalization, e.g.

$$L = D - W$$

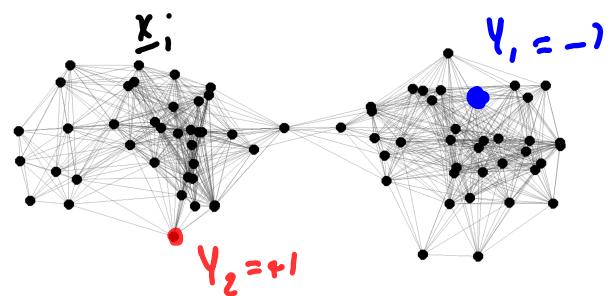
⑤

Step 3) Since we are only interested in labeling x_{N+1}, \dots, x_m that are additional nodes on the graph we can think of the latent function f in the probit model as a function on the vertices of the graph

$$f: V \rightarrow \mathbb{R}, \quad f \equiv \underline{f} \in \mathbb{R}^M$$

That is $f(i) = f_i$.

Then we define the probit optimization problem



$$\underline{f}^* = \arg \min_{\underline{f} \in \mathbb{R}^M} - \sum_{j=1}^N \log \Psi(f_j y_j) + \lambda \underline{f}^T C \underline{f}$$

where we take $C = (L + \gamma^2 I)^{-\alpha}$ for parameters $\gamma^2, \alpha > 0$.

• First, observe this is completely analogous to our formulation of Probit on general state spaces

⑥

$$f^* = \underset{f \in H}{\operatorname{argmin}} - \sum_{j=1}^M \log \Psi(f(x_j) y_j) + \gamma \|f\|_K^2$$

RKHS Norm

By this observation we readily infer that

$\underline{f}^T C \underline{f}$ must be some kind of an RKHS norm!

- In fact our kernel here turns out to be $K: \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$, $K(i, j) = C_{ij}$.

Then we can write down a rep. thm. (Next lecture).

- Suppose L is symm so we can write

$$L = U \Sigma U^T \text{ then } C = U (\Sigma + \gamma^2 I)^{-\alpha} U^T$$

that is $C = \sum_{j=1}^M \frac{1}{(\sigma_j + \gamma^2)^\alpha} \underline{u}_j \underline{u}_j^T$

Suppose there are P clusters in X then $\sigma_1 = \sigma_2 = \dots = \sigma_P = 0$ & $\sigma_{P+1} = G(1)$ then

$$C = \underbrace{\sum_{j=1}^P \frac{1}{\gamma^{2\alpha}} \underline{u}_j \underline{u}_j^T}_{G(\gamma^{-2\alpha})} + \underbrace{\sum_{j=P+1}^M \frac{1}{(\sigma_j + \gamma^2)^\alpha} \underline{u}_j \underline{u}_j^T}_{G(1)}$$

(7)

