# Classification Semi-supervised learning based on network

#### Speakers: Hanwen Wang, Xinxin Huang, and Zeyu Li CS 249-2 2017 Winter

# Semi-Supervised Learning Using Gaussian Fields and Harmonic Functions

Xiaojin Zhu, Zoubin Ghahramani, John Lafferty

School of Computer Science, Carnegie Mellon University, Gatsby Computational Neuroscience Unit, University College London

# Introduction

#### **Supervised Learning**

labeled data is expensive

- Skilled human anotators
- Time consuming
- example: protein shape classficaton

#### Semi-supervised Learning

Exploit the manifold structure of data Assumption: similar unlabeled data should be under one category

# Frame Work

#### Annotation:

- Labeled points: L = {1,..,I}
- Unlabel points: U = {I+1,..,I+u}
- The similairty betwen point i and j: w(i,j)

#### **Objective:**

- Find a function:  $f: V \to \mathbb{R}$ such that the energy function  $E(f) = \frac{1}{2} \sum_{i,j} w_{xy} (f(i) - f(j))^2$  is minimized.
- Similar points have higher weight

How to find the minimum of a function? Ans: first derivation

$$E(f) = \frac{1}{2} \sum_{i,j} w_{xy} (f(i) - f(j))^2$$
Partial derivation
$$\frac{\partial E(f)}{\partial f(i)} = \frac{1}{2} \sum_j w_{ij} * 2(f(i) - f(j))$$

$$= \sum_j (w_{ij}f(i) - w_{ij}f(j))$$

Assign the right hand size to zero gives us:

$$f(i) = \frac{\sum_{j} w_{ij} f(j)}{\sum_{j} w_{ij}}$$

$$= \frac{\sum_j w_{ij} f(i)}{d_i}$$

Since  $f = \arg \min_{f|_L = f_l} E(f)$  is harmonic, f satisfy  $\Delta f = 0$ 

https://en.wikipedia.org/wiki/Harmonic\_function

(D-W)f = 0



If we pick a row and expand the matrix multiplication, we will get

$$f(i) = \frac{\sum_{j} w_{ij} f(i)}{d_i}$$

#### Now, we do the calculation in matrix form

$$\Delta f = 0 \longrightarrow Df = Wf$$

$$\begin{pmatrix} D_{ll} & 0\\ 0 & D_{uu} \end{pmatrix} \begin{pmatrix} f_l\\ f_u \end{pmatrix} = \begin{pmatrix} w_{uu} & w_{lu}\\ w_{ul} & w_{uu} \end{pmatrix} \begin{pmatrix} f_l\\ f_u \end{pmatrix}$$

Expanding the second row, we get:

$$D_{uu}f_u = w_{ul}f_l + w_{uu}f_u$$
$$f_u = (D_{uu} - w_{uu})^{-1}w_{ul}f_l$$

Since: 
$$P = D^{-1}w \rightarrow w_{uu} = D_{uu}P_{uu}$$

We get:

$$f_u = (I - P_{uu})^{-1} P_{ul} f_l$$

Further expand the equation

$$f_{u} = (I - P_{uu})^{-1} P_{ul} f_{l}$$

$$(I - A)^{-1} = (I + A + A^{2} + A^{3} + \dots A^{n})$$

$$f_{u} = p_{ul} f_{l} + p_{uu} p_{ul} f_{l} + p_{uu}^{2} p_{ul} f_{l} + \dots p_{uu}^{n} p_{ul} f_{l}$$

### Example



### Interpretation 1: Random Walk



		x1	x2	х3	x4
P =	x1	0.0	0.0	0.0	0.0
	x2	0.5	0.0	0.5	0.0
	х3	0.0	0.5	0.0	0.5
	x4	0.0	0.0	1.0	0.0

# Interpretation 2: Electric Network

Edges: resistor with conductance

Point Labels: voltage

 $P = V^2 / R$ 

Energy dissipation is minimized since the voltage difference between two neighbors are minimized

• Heat equation: a parabolic partial differential equation that describes the distribution of heat (or variation in temperature) in a given region over time

$$\frac{\partial u_t(x,y)}{\partial t} = \alpha \left(\frac{(\partial^2 u_t(x,y))}{\partial^2 x} + \frac{(\partial^2 u_t(x,y))}{\partial^2 y}\right) = \alpha \Delta u_t(x,y)$$

• Heat kernel, it is a solution:

$$K_t = e^{-t\Delta}$$

 $K_t(i, j)$ : the solution of heat equation with initial conditions being a point source at i.

If we use this kernel in kernel classifier:

$$\hat{f}_t(j) = \sum_{i \in L} \alpha_i y_i K_t(i, j)$$

The kernel classifier can be considered as the solution of the heat equation with initial heat resource at the labeled data.

If we don't consider the time and we only consider about the temprature relation between different points

$$G = \int_0^\infty K_t dt = \Delta^{-1}$$

Consider the green funtion on unlabel data

$$f_{u} = G_{uu} W_{ul} f_{l}$$
  $f(i) = \sum_{i=1}^{l} \sum_{k} y_{i} w_{ik} G(k, j)$ 

1

Our method can be interpreted as a kernel classifier with kernel G

- Spectrum of the G is the inverse of spectrum of  $\ \Delta$ 
  - This can indicate a connection to the work of Chapelle et al. 2002 about cluster kernels for semi-supervised learning.
  - By manipulating the eigenvalues of graph Laplacian, we can construct kernels which implement the cluster assumption: the induced distance depends on whether the points are in the same cluster or not.

# **Interpretation 4: Spectral Clustering**

• Normalized cutting problem: Minimize the cost function:

$$min_f \frac{f^T (D - W)f}{f^T D f}$$

The solution is the eigenvector corresponding to the second smallest eigenvalue of the generalized eigenvalue problem

$$\Delta f = \lambda D f$$
 or  $(I - D^{-1}W)f = \lambda f$ 

# Spectral Clustering with group constraints

- Yu and Shi (2001) added group bias into the normalized cutting problem to specify which points should be in the same group.
- They proposed some pairwise grouping constraints of the labeled data.
- Imply the intuition that the points tend to be in the same cluster(have the labels) as its neighbors.

# Label Propagation v.s.Constrainted Clustering

Semi-supervised learning on the graph can be interpreted in two ways

- In label propagation algorithms, the known labels are propagated to the unlabeled nodes.
- In constrained spectral clustering algorithms, known labels are first converted to pairwise constraints, then a constrained cut is computed as a tradeoff between minimizing the cut cost and maximizing the constraint satisfaction
- (Wang and Qian 2012)

# Incorporating Class Prior Knowledge

- Decision rule is  $f(i) > \frac{1}{2}$ , assign to label 1, otherwise assign to label 0.
  - it works only when the classes are well separated. However, in real datasets ,the situation is different. Using f tend to produce severely unbiased classification
- Reason: W may be poorly estimated and does not reflect the classification goal.

We can not fully trust the the graph structure.

We want to incorprate the class prior knowledge in our model

# Incorporating Class Prior Knowledge

- q: proportion for class 1; 1-q: proportion for class 0.
- To match this priors, we modified the decision rule by class mass normalization as

$$q\frac{f_u(i)}{\sum_i f_u(i))} > (1-q)\frac{1-f_u(i)}{\sum_i (1-f_u(i))}$$

Example: f = [0.1, 0.2, 0.3, 0.4] and q = 0.5

L.H.: [0.05,0.1,0.15,0.2]

The first 2 will be assigned with label1,

R.H.: [0.15,0.133,0.117,0.1]

while the last 2 will be assigned with label 0.

### **Incorporating Externel Classifier**

- Assume the external classifer produces label  $h_u$  on the unlabeled data.
  - it can either be 0/1 or soft label [0.1]



# Learning the Weight Matrix W

Recall the defination of Weight Matrix

$$w_{ij} = exp(-\sum_{d=1}^{m} \frac{(x_{jd} - x_{id})^2}{\sigma_d^2})$$

This will be a feature selection mechanism which better aligns the graph structure with the data.

Learn  $\sigma_d$  by minimizing the average label entropy

$$H(f) = \frac{1}{u} \sum_{i=l+1}^{l+u} H_i(f(i))$$

# Learning the Weight Matrix W

$$H_i(f(i)) = -f(i)log(f(i)) - (1 - f(i))log(1 - f(i))$$

Why we can get optimal  $\sigma_d$  by minimizing H?

- Small H(i) implies that f(i) is close to 0 or 1.
- This captures intuition that a good W (equivalently a good set of {  $\sigma_d$ }) should result in confident labeling.
- min H lead to a set of optimal  $\sigma_d$  which can result in confident labeling u.

# Learning the Weight Matrix W

- Important property of H is that H has a minimum at 0 as  $\sigma_d \rightarrow 0$
- Solution:

$$\widetilde{P} = \epsilon U + (1 - \epsilon)P$$

The label will not be dominated by its nearest neighbor. It can also be influenced by all the other nodes.

• Use the gradient decsent to get the hyperparameter  $\sigma_d$ 

# Conclusion

- Harmonic function is strong model to solve the semi-supervised learing problem.
- Label propagation and constrained spectral clustering algorithms can also be implemented to solve the semi-supervised learning tasks.
- This model is flexible and can be easily incorprated with external helpful information.

# Graph Regularized Transductive Classification on Heterogeneous Information Networks

Ming Ji, Yizhou Sun, Marina Danilevsky, Jiawei Han and Jing Gao Dept. of Computer Science, University of Illinois at Urbana-Champaign Semi-supervised Learning:

Classify the unlabeled data based on known information

Two groups of classification:

Transductive classification - to predict labels for the given unlabeled data

Inductive classification - construct decision function in whole data space

Homogeneous network & Heterogeneous network

Classifying multi-typed objects into classes.



Fig. 1. Knowledge propagation in a bibliographic information network

#### **Definition 1: Heterogeneous information network**

*m* types of data objects:  $\mathcal{X}_1 = \{x_{11}, \dots, x_{1n_1}\}, \dots, \mathcal{X}_m = \{x_{m1}, \dots, x_{mn_m}\}$ Graph  $G = \langle V, E, W \rangle$ 

$$V = \bigcup_{i=1}^{m} \mathcal{X}_i \quad m \ge 2$$

#### **Definition 2: Class**

Given 
$$G = \langle V, E, W \rangle$$
 and  $V = \bigcup_{i=1}^{m} \mathcal{X}_i$ ,  
Class:  $G' = \langle V', E', W' \rangle$ , where  $V' \subseteq V$ ,  $E' \subseteq E$ .  
 $W'_{x_{ip}x_{jq}} = W_{x_{ip}x_{jq}} \ \forall e = \langle x_{ip}, x_{jq} \rangle \in E'$ 

#### **Definition 2: Class**

Given  $G = \langle V, E, W \rangle$  and  $V = \bigcup_{i=1}^{m} \mathcal{X}_i$ , Class:  $G' = \langle V', E', W' \rangle$ , where  $V' \subseteq V$ ,  $E' \subseteq E$ .  $W'_{x_{ip}x_{jq}} = W_{x_{ip}x_{jq}} \quad \forall e = \langle x_{ip}, x_{jq} \rangle \in E'$ 



# Definition 3: Transductive classification on heterogeneous information networks

Given  $G = \langle V, E, W \rangle$  and  $V' \subseteq V$  which are labeled with value  $\mathcal{Y}$ ,

Predict the class labels for all unlabeld object V - V'



# Definition 3: Transductive classification on heterogeneous information networks

Given  $G = \langle V, E, W \rangle$  and  $V' \subseteq V$  which are labeled with value  $\mathcal{Y}$ ,

Predict the class labels for all unlabeld object V - V'



Suppose the number of classifiers is *K* 

Compute  $\boldsymbol{F}_i = [\boldsymbol{f}_i^{(1)}, \dots, \boldsymbol{f}_i^{(K)}] \in \mathbb{R}^{n_i \times K}$ , where each  $\boldsymbol{f}_i^k = [f_{i1}^k, \dots, f_{in_i}^k]$  measures the confidence of  $x_{ip} \in \mathcal{X}_i$  belongs to class k.

Class of 
$$x_{ip}$$
 is  $\arg \max_{1 \le k \le K} f_{ip}^{(k)}$ 

Use  $R_{ij}$  to denote the relation matrix of Type *i* and Type *j*,  $R_{ij} \in \mathbb{R}^{n_i \times n_j}$ .

 $R_{ij,pq}$  represents the weight on link  $\langle x_{ip}, x_{jq} \rangle$ 

$$R_{ij,pq} = \begin{cases} 1 & \text{if data objects } x_{ip} \text{ and } x_{jq} \text{ are linked together} \\ 0 & \text{otherwise} \end{cases}$$

Another vector to use:

$$y_{ip} = \begin{cases} 1 & \text{if } x_{ip} \text{ is labeled to the k-th class} \\ 0 & \text{otherwise} \end{cases}$$

The goal is to predict infer a set of  $f_i^k$  from  $R_{ij}$  and  $y_i^{(k)}$ .

#### Intuition:

Prior Knowledge: A1, P1 and C1 belong to "data mining" => Infer: A2, T1 are highly related to data mining.

Similarly: A3, C2, T2, and T3 highly related to "database".

Knowledge propagation.



 ${\bf Fig. 1.}\ {\bf Knowledge\ propagation\ in\ a\ bibliographic\ information\ network}$ 

Formulate Intuition as follows:

- (1) The estimated confidence measure of two objects  $x_{ip}$  and  $x_{jq}$  belonging to class k,  $f_{ip}^{(k)}$  and  $f_{jq}^{(k)}$ , should be similar if  $x_{ip}$  and  $x_{jq}$  are linked together, i.e., the weight value  $R_{ij,pq} > 0$ .
- (2) The confidence estimation  $f_i^k$  should be similar to the ground truth,  $y_i^{(k)}$ .

#### The Algorithm:

Define a diagonal matrix  $D_{ij}$  of size  $n_i \times n_i$ . The (p,p)-th of  $D_{ij}$  is the sum of the p-th row of  $R_{ij}$ .

#### **Objective function:**

$$J(\boldsymbol{f}_{i}^{(k)},\ldots,\boldsymbol{f}_{m}^{(k)}) = \sum_{i,j=1}^{m} \lambda_{i,j} \sum_{p=1}^{n_{i}} \sum_{q=1}^{n_{j}} R_{ij,pq} \left( \frac{1}{\sqrt{D_{ij,pp}}} f_{ip}^{(k)} - \frac{1}{\sqrt{D_{ji,qq}}} f_{jq}^{(k)} \right)^{2} + \sum_{i=1}^{m} \alpha_{i} (\boldsymbol{f}_{i}^{(k)} - \boldsymbol{y}_{i}^{(k)})^{T} (\boldsymbol{f}_{i}^{(k)} - \boldsymbol{y}_{i}^{(k)})$$

0

#### Trade-off:

Controlled by  $\lambda_{ij}$  and  $\alpha_i$  where

Larger  $\lambda_{ij}$ : more rely on relationship of  $\mathcal{X}_i$  and  $\mathcal{X}_j$ .

Larger  $\alpha_i$ : The label of *i* is more trustworthy.

#### Prior Knowledge.

Define normalized form  $S_{ij} = D_{ij}^{(-1/2)} R_{ij} D_{ji}^{-1/2}$ 

Rewrite the objective function as

$$J(\boldsymbol{f}_{1}^{(k)}, \dots, \boldsymbol{f}_{m}^{(k)}) = \sum_{i,j=1}^{m} \lambda_{ij} ((\boldsymbol{f}_{i}^{(k)})^{T} \boldsymbol{f}_{i}^{(k)} + (\boldsymbol{f}_{j}^{(k)})^{T} \boldsymbol{f}_{j}^{(k)} - 2(\boldsymbol{f}_{i}^{(k)})^{T} \boldsymbol{S}_{ij} \boldsymbol{f}_{j}^{(k)}) + \sum_{i=1}^{m} \alpha_{i} (\boldsymbol{f}_{i}^{(k)} - \boldsymbol{y}_{i}^{(k)})^{T} (\boldsymbol{f}_{i}^{(k)} - \boldsymbol{y}_{i}^{(k)})$$

Reduce to

$$J(\boldsymbol{f}_{1}^{(k)}) = 2\lambda_{11}(\boldsymbol{f}_{1}^{(k)})^{T} \mathbf{L}_{11} \boldsymbol{f}_{1}^{(k)} + \alpha_{1}(\boldsymbol{f}_{1}^{(k)} - \mathbf{y}_{1}^{(k)})^{T} (\boldsymbol{f}_{1}^{(k)} - \mathbf{y}_{1}^{(k)})$$

in homogeneous information networks.  $\mathbf{L}_{ii} = \mathbf{I}_i - \mathbf{S}_{ii}$  is the *normalized graph Laplacian*.

Given the following definition:

$$\boldsymbol{f}^{(k)} = [(\boldsymbol{f}_1^{(k)})^T, \dots, (\boldsymbol{f}_m^{(k)})^T]^T; \, \mathbf{y}^{(k)} = [(\mathbf{y}_1^{(k)})^T, \dots, (\mathbf{y}_m^{(k)})^T]^T$$
$$\boldsymbol{\alpha} = \operatorname{diag}(\alpha_1, \alpha_2, \dots, \alpha_m) \qquad \mathbf{L}_{ij} = \begin{bmatrix} \mathbf{I}_i & -\mathbf{S}_{ij} \\ -\mathbf{S}_{ji} & \mathbf{I}_j \end{bmatrix}, \, \text{where} \, i \neq j$$

Let  $H_{ij}$  be the  $n \times n$  matrix,  $\sum_{i=1}^{m} n_i = n$  And  $\mathbf{H} = \sum_{i \neq j} \lambda_{ij} \mathbf{H}_{ij} + 2 \sum_{i=1}^{m} \lambda_{ii} \mathbf{H}_{ii}$ 

We can rewrite the objective function as:

$$J(\boldsymbol{f}_{1}^{(k)},\ldots,\boldsymbol{f}_{m}^{(k)}) = (\boldsymbol{f}^{(k)})^{T} \mathbf{H} \boldsymbol{f}^{(k)} + (\boldsymbol{f}^{(k)} - \mathbf{y}^{(k)})^{T} \boldsymbol{\alpha} (\boldsymbol{f}^{(k)} - \mathbf{y}^{(k)})$$

#### Solution

Closed form solution

Hessian matrix of  $H(J(\boldsymbol{f}_{1}^{(k)}, \dots, \boldsymbol{f}_{m}^{(k)})) = 2\mathbf{H} + 2\boldsymbol{\alpha}$  is positive semi-definite.  $\frac{\partial J}{\partial (\boldsymbol{f}_{i}^{(k)})^{T}} = \sum_{j=1, j \neq i}^{m} \lambda_{ij} (2\boldsymbol{f}_{i}^{(k)} - 2\mathbf{S}_{ij}\boldsymbol{f}_{j}^{(k)}) + 4\lambda_{ii}\mathbf{L}_{ii}\boldsymbol{f}_{i}^{(k)} + 2\alpha_{i}(\boldsymbol{f}_{i}^{(k)} - \mathbf{y}_{i}^{(k)})$ And setting  $\frac{\partial J}{\partial (\boldsymbol{f}_{i}^{(k)})^{T}} = 0$  for all *i*.  $f_{i}^{(k)} = \left( \left(\sum_{j=1, j \neq i}^{m} \lambda_{ij} + \alpha_{i}\right)\mathbf{I}_{i} + 2\lambda_{ii}\mathbf{L}_{ii} \right)^{-1} \left(\alpha_{i}\mathbf{y}_{i}^{(k)} + \sum_{j=1, j \neq i}^{m} \lambda_{ij}\mathbf{S}_{ij}\boldsymbol{f}_{j}^{(k)} \right), i \in \{1, \dots, m\}$ 

• Iterative solution

#### Solution

- Closed form solution
- Iterative solution

Step 0: Initialization.  $f_i^{(k)}(0) = \mathbf{y}_i^{(k)}$ 

# Step 1: Based on current $f_i^{(k)}(t)$ , compute the $f_i^{(k)}(t+1) = \frac{\sum_{j=1, j \neq i}^m \lambda_{ij} \mathbf{S}_{ij} f_j^{(k)}(t) + 2\lambda_{ii} \mathbf{S}_{ii} f_i^{(k)}(t) + \alpha_i \mathbf{y}_i^{(k)}}{\sum_{j=1, j \neq i}^m \lambda_{ij} + 2\lambda_{ii} + \alpha_i}$

Step 2: Repeat Step 1 until converge. ( $f_i^{(k)}(t)$  change little over t-th iteration)

Step 3: Assign class label to p-th object of  $\mathcal{X}_i$  by  $c_{ip} = \arg \max_{1 \le k \le K} f_{ip}^{(k)*}$ .

Complexity

• Iterative solution O(NK(|E| + |V|))

*N*: **#** of iterations. *K*: **#** of classes.

• Closed form solution

Worse than the iterative solution since iterative solution bypass the matrix inversion operation.



Hanwen Wang: Xinxin Huang: Zeyu Li: hanwenwang@g.ucla.edu xinxinh@ucla.edu zyli@cs.ucla.edu