

# SDS 385: Stat Models for Big Data Lecture 10a: Semisupervised learning

Purnamrita Sarkar Department of Statistics and Data Science The University of Texas at Austin

https://psarkar.github.io/teaching

# Semi-supervised learning

- You are given a lot of unlabeled data.
- Only a few points are labeled.
- Is this useful?



- Two broad ways
  - Label propagation:
    - Graph Based algorithm
    - Does not generalize to unseen data, i.e. Transductive
  - Manifold regularization
    - Graph Based regularization
    - Does generalize to unseen data, i.e. Inductive

- Input *n* data points  $x_1, \ldots, x_n$
- Define similarity matrix  $S \in \mathbb{R}^{n \times n}$

$$S_{ij} = \exp(-\|x_i - x_j\|^2/2\sigma^2)$$

• Since *S* is dense, often *k* nearest neighbor graphs are also used. (Your homework!)

# Label propagation [Zhu et al, 2003]

- Input:
  - $\ell$  labeled datapoints  $(x_1, y_1), \ldots, (x_\ell, y_\ell)$
  - *u* unlabeled datapoints  $x_{\ell+1}, \ldots, x_{\ell+u}$
- Predict: labels  $y_{\ell+1}, \ldots, y_{\ell+u}$



• Compute  $P \in [0,1]^{(\ell+u) \times (\ell+u)}$ 

• 
$$P_{ij} = \frac{S_{ij}}{\sum_j S_{ij}}$$

- Harmonic function:
  - Function value at an unlabeled node is an average of function values at its neighbors

• For 
$$j = \ell + 1 : \ell + u_i$$

$$f(j) = \frac{\sum_{i} S_{ji} f(i)}{\sum_{i} S_{ji}} = \sum_{i} P_{ji} f(i)$$

• In other words, for the unlabeled nodes, this fixed point equation is satisfied

$$f[U] = Pf[U] \qquad f[L] = y[L]$$

- For a vector v and set S, we denote by v[S] the subset of values in S
- U and L denote the set of unlabeled and labeled points respectively.

convex combination of values of neighbors

## **Closed form**

- Assume there are just two classes. Set  $y[L] \in \{0,1\}^{\ell}$  accordingly.
- We have:

$$\begin{bmatrix} P_{LL} & P_{LU} \\ P_{UL} & P_{UU} \end{bmatrix} \begin{bmatrix} \mathbf{Y}_{L} \\ \mathbf{Y}_{U} \end{bmatrix} = \begin{bmatrix} \mathbf{Y}_{L} \\ \mathbf{Y}_{U} \end{bmatrix}$$

• Expanding, we get:

$$P_{UL}\mathbf{Y}_{L} + P_{UU}\mathbf{Y}_{U} = \mathbf{Y}_{U}$$

• Moving things around:

$$Y_U = (I - P_{UU})^{-1} P_{UL} \mathbf{Y}_L$$

• Can use a linear system solver

# Label propagation - Random walk interpretation

- Think of the labeled nodes as absorbing states
- Use

$$Y_{U} = \sum_{t=0}^{\infty} P_{UU}^{t} P_{UL} Y_{L}$$
  
= 
$$\underbrace{P_{UL}Y_{L}}_{\text{Probability of reaching "1"s in one step}} + \underbrace{P_{UU}P_{UL}Y_{L}}_{\text{Probability of reaching in two steps}}$$
  
= Probability of reaching a label "1" in a long random walk

- Why is this useful?
  - If the labels are all reachable, a long walk must hit a "0" or a "1"
  - So if Y<sub>i</sub> > 1/2, that means from i, its more likely to reach a "1" than a "0"

+

- Graph Laplacian: L = D S, where  $D_{ii} = \sum_{i} S_{ij}$  is a diagonal matrix
- L is positive semi-definite (we are assuming  $S_{ij} > 0$ )
- Why?
  - For any vector v,

$$v^T L v = \sum_{ij} S_{ij} (v_i - v_j)^2$$

- So this measures how unsmooth v is w.r.t S
- L is also singular, why?

• We can also frame label propagation as

$$\arg\min_{v} v^{T} L v \qquad s.t.v[L] = y_{L}$$

• Why?

• We can also frame label propagation as

$$\arg\min_{v} v^{T}Lv \qquad s.t.v[L] = y_{L}$$

- Why? Write  $v = \begin{bmatrix} y_L & v_U \end{bmatrix}$

• We can also frame label propagation as

$$\arg\min_{v} v^{T}Lv \qquad s.t.v[L] = y_{L}$$

- Why?
- Write  $v = \begin{bmatrix} y_L & v_U \end{bmatrix}$
- Now set the derivative to zero.

• 
$$Lv = 0$$
 such that  $v[L] = y_L$ 

•

$$\begin{bmatrix} D_{LL} - S_{LL} & -S_{LU} \\ -S_{UL} & D_{UU} - S_{UU} \end{bmatrix} \begin{bmatrix} \mathbf{y}_L \\ \mathbf{v}_U \end{bmatrix} = \mathbf{0}$$

• Solving:

$$-S_{UL}\mathbf{y}_{L} + (D_{UU} - S_{UU})\mathbf{v}_{U} = 0$$

• We can also frame label propagation as

$$\arg\min_{v} v^{T}Lv \qquad s.t.v[L] = y_{L}$$

- Why?
- Write  $v = \begin{bmatrix} y_L & v_U \end{bmatrix}$
- Now set the derivative to zero.

• 
$$Lv = 0$$
 such that  $v[L] = y_L$ 

•

$$\begin{bmatrix} D_{LL} - S_{LL} & -S_{LU} \\ -S_{UL} & D_{UU} - S_{UU} \end{bmatrix} \begin{bmatrix} \mathbf{y}_L \\ \mathbf{v}_U \end{bmatrix} = \mathbf{0}$$

• Solving:

$$-S_{UL}\mathbf{y}_{L}+(D_{UU}-S_{UU})\mathbf{v}_{U}=0$$

• rearranging:  $v_U = (D_{UU} - S_{UU})^{-1} S_{UL} y_L = (I - P_{UU})^{-1} P_{UL} y_L$ 



Figure 3. Harmonic energy minimization on digits "1" vs. "2" (left) and on all 10 digits (middle) and combining voted-perceptron with harmonic energy minimization on odd vs. even digits (right)



Figure 4. Harmonic energy minimization on PC vs. MAC (left), baseball vs. hockey (middle), and MS-Windows vs. MAC (right)

- Input:
  - $\ell$  labeled datapoints  $(x_1, y_1), \ldots, (x_\ell, y_\ell)$
  - *u* unlabeled datapoints  $x_{\ell+1}, \ldots, x_{\ell+u}$
- Predict: labels  $y_{\ell+1}, \ldots, y_{\ell+u}$

- Input:
  - $\ell$  labeled datapoints  $(x_1, y_1), \ldots, (x_\ell, y_\ell)$
  - *u* unlabeled datapoints  $x_{\ell+1}, \ldots, x_{\ell+u}$
- Predict: labels  $y_{\ell+1}, \ldots, y_{\ell+u}$





Figure 1: Unlabeled data and prior beliefs

- Input:
  - $\ell$  labeled datapoints  $(x_1, y_1), \ldots, (x_\ell, y_\ell)$
  - *u* unlabeled datapoints  $x_{\ell+1}, \ldots, x_{\ell+u}$
- Predict: labels  $y_{\ell+1}, \ldots, y_{\ell+u}$

- Input:
  - $\ell$  labeled datapoints  $(x_1, y_1), \ldots, (x_\ell, y_\ell)$
  - *u* unlabeled datapoints  $x_{\ell+1}, \ldots, x_{\ell+u}$
- Predict: labels  $y_{\ell+1}, \ldots, y_{\ell+u}$
- Before the constraint was  $v[L] = y_L$ , now instead we will use a loss function and learn a classifier on the labeled data

$$\min_{w} \underbrace{\sum_{i=1}^{\ell} \mathsf{loss}(y_i, w^T x_i)}_{\mathsf{loss}} + \lambda \underbrace{R(w)}_{\mathsf{regularization}}$$

- Input:
  - $\ell$  labeled datapoints  $(x_1, y_1), \ldots, (x_\ell, y_\ell)$
  - *u* unlabeled datapoints  $x_{\ell+1}, \ldots, x_{\ell+u}$
- Predict: labels  $y_{\ell+1}, \ldots, y_{\ell+u}$
- Before the constraint was  $v[L] = y_L$ , now instead we will use a loss function and learn a classifier on the labeled data

$$\min_{w} \underbrace{\sum_{i=1}^{\ell} \operatorname{loss}(y_{i}, w^{T} x_{i})}_{loss} + \lambda \underbrace{R(w)}_{regularization}$$

• How about the unlabeled data?

#### Manifold regularization: Belkin et al 2006

$$\min_{w} \underbrace{\sum_{i=1}^{\ell} \operatorname{loss}(y_{i}, w^{T}x_{i}) + \lambda}_{loss} + \underbrace{R(w)}_{regularization} + \beta(Xw)^{T} L(Xw)$$

- Assume a linear predictor  $w^T x$
- Idea: close/similar points have similar predicted labels.
- LapSVM:

$$\min_{w} \sum_{i=1}^{\ell} (1 - y_i f(x_i))_+ + \lambda \underbrace{\|f\|_{K}^2}_{regularization} + \beta f^{\mathsf{T}} Lf$$

$$\min_{w,y_{\ell+1},\dots,y_{\ell+u}} \sum_{i=1}^{\ell} (1-y_i f(x_i))_+ + C' \sum_{i=\ell+1}^n (1-y_i f(x_i))_+ + \lambda \underbrace{\|f\|_K^2}_{regularization}$$

- Iteratively solves SVM quadratic programs
- Switches labels to improve objective function
- Suffers from local optima, inherently combinatorial problem

# Transductive SVM VS LapSVM



- Cho-Jui Hsieh's lecture notes from UC Davis
- Zhu et al's paper in ICML "Semi-Supervised Learning Using Gaussian Fields and Harmonic Functions"
- Ng et al's paper on ranking Stability "Link Analysis, Eigenvectors and Stability", IJCAI 2001
- Belkin, Niyogi and Sindhwani's paper "Manifold Regularization: A Geometric Framework for Learning from Labeled and Unlabeled Examples" in JMLR 2006
- My old talk on Random walks.