

Predicting Switching Graph Labelings with Cluster Specialists

MoN18: Eighteenth Mathematics of Networks Meeting

James Robinson (joint work with Mark Herbster)

8 April 2019

Department of Computer Science
University College London

Outline

Introduction

Predicting Switching Graph Labelings

Cluster Specialists

Experiments

Conclusion

Introduction

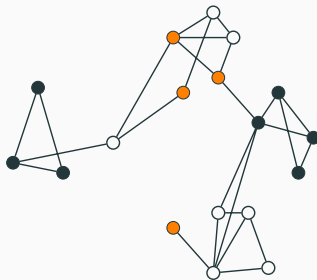
Introduction

- Graph prediction is a foundational problem in machine learning
- Many flavours/settings (node classification, edge classification, clustering)
- Today: Node classification in the *online learning* setting (sequential prediction)
- Want to develop algorithms with *performance guarantees*

Predicting Switching Graph Labelings

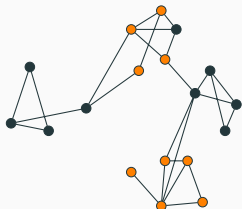
Predicting Graph Labelings Online

- n -vertex Graph $\mathcal{G} = (V, E)$,
 $V = \{1, \dots, n\}$
- A labeling is a function $\mathbf{u} : V \mapsto \{-1, 1\}$
- Online learning protocol:
For $t = 1, \dots, T$ **do**:
 1. Nature selects a vertex $i_t \in V$
 2. Learner predicts $\hat{y}_t \in \{-1, 1\}$
 3. Nature reveals label $\mathbf{u}_t(i_t) \in \{-1, 1\}$
 4. Learner incurs loss $m_t = [\mathbf{u}_t(i_t) \neq \hat{y}_t]$
- No statistical assumptions are made!
Nature could be adversarial
- Performance guarantees hold in the worst case

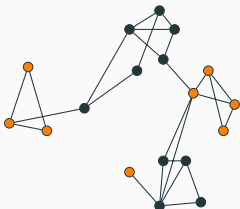


Switching Graph Labelings

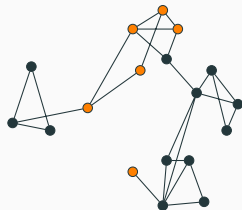
Sequence of labelings $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_T$ s.t. $|\{t : \mathbf{u}_t \neq \mathbf{u}_{t+1}\}| = K$



$t \rightarrow$ $t = 1 \dots$



$\dots t = 7 \dots$



$\dots t = 20 \dots$

The learner **doesn't** know when switches occur

Assume K is 'small'

Objectives

- Minimize the number of mistakes

$$M = \sum_{t=1}^T m_t = \sum_{t=1}^T [\mathbf{u}_t(i_t) \neq \hat{y}_t]$$

- Provide *good* mistake bound guarantees for *switching*:

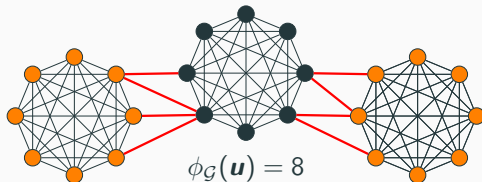
$$M \leq f(\text{complexity}(\mathbf{u}_1, \dots, \mathbf{u}_T), K, \text{structure}(\mathcal{G}))$$

- Algorithms should be *fast* (online predictions)

complexity($\mathbf{u}_1, \dots, \mathbf{u}_T$) - Cut-size ϕ

We assume that a graph \mathcal{G} consists of tightly-connected clusters, with loose inter-cluster connections. Nodes in a cluster (mostly) share the same label.

A labeling $\mathbf{u} : V \mapsto \{-1, 1\}$ induces a *cut* $\phi_{\mathcal{G}}(\mathbf{u}) = \sum_{(i,j) \in E} [\mathbf{u}(i) \neq \mathbf{u}(j)]$



Static mistake bounds typically scale *linearly* with $\phi_{\mathcal{G}}(\mathbf{u})$ - **sensitive!**

- [HLP08] - $\mathcal{O}\left(\phi_{\mathcal{G}}(\mathbf{u}) \log \frac{n}{\phi_{\mathcal{G}}(\mathbf{u})} + \phi_{\mathcal{G}}(\mathbf{u})\right)$
- [HP06] - $\mathcal{O}(\phi_{\mathcal{G}}(\mathbf{u})R_{\mathcal{G}})$, $R_{\mathcal{G}} = f(\text{structure}(\mathcal{G}))$

complexity($\mathbf{u}_1, \dots, \mathbf{u}_T$) - Effective Resistance $r_{i,j}$

Define $r_{i,j}$ to be the *effective resistance* between nodes i and j when \mathcal{G} is a network of *unit* resistors (edges)



$$r_{a,b} = 1$$

$$r_{c,d} = \frac{1}{1 + \frac{1}{2}} = \frac{2}{3}$$

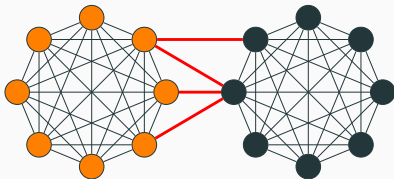
(Kirchoff's laws for resistors in series and parallel)

$r_{i,j}$ is a measure of connectivity - effective resistance between two nodes *decreases* with increased connectivity

Definition

Define the *resistance-weighted cut-size* to be:

$$\phi^r(\mathbf{u}) = \sum_{(i,j) \in E} r_{i,j} [\mathbf{u}(i) \neq \mathbf{u}(j)]$$



- Two m -cliques with $\ell < m$ edges between them
- For all vertices $i, j \in V$, we have $r_{i,j} \leq \Theta(\frac{1}{\ell})$
- Hence,

$$\phi(\mathbf{u}) = \sum_{(i,j) \in E} [\mathbf{u}(i) \neq \mathbf{u}(j)] = \ell$$

$$\phi^r(\mathbf{u}) = \sum_{(i,j) \in E} r_{i,j} [\mathbf{u}(i) \neq \mathbf{u}(j)] \leq \Theta(1)$$

- $\phi^r(\mathbf{u})$ is robust!

Random Spanning Tree - Resistance weighted cut-size

How to exploit $\phi^r(\mathbf{u})$?



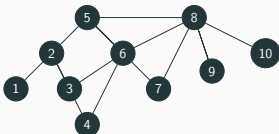
Expected cut-size of a random spanning tree **generated uniformly at random** ([CBGV09]):

$$\begin{aligned}\mathbb{E}[\phi_{\mathcal{T}}(\mathbf{u})] &= \sum_{(i,j) \in E} \mathbb{P}((i,j) \in E_{\mathcal{T}}) [\mathbf{u}(i) \neq \mathbf{u}(j)] \\ &= \sum_{(i,j) \in E} r_{i,j} [\mathbf{u}(i) \neq \mathbf{u}(j)] \\ &= \phi^r(\mathbf{u})\end{aligned}$$

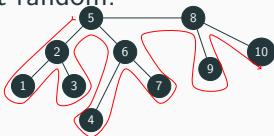
Mistake bounds in terms of $\phi(\mathbf{u})$ become *expected* mistake bounds in terms of $\phi^r(\mathbf{u})$!

Two Transformations - Trees and Linear Embeddings

Original graph \mathcal{G} :



Sample \mathcal{T} uniformly at random:



Compute *spine* \mathcal{S} from \mathcal{T} (depth-first search):



Properties: ([HLP08])

$$\phi_{\mathcal{S}}(\mathbf{u}) \leq 2\phi_{\mathcal{T}}(\mathbf{u}) \leq 2\phi_{\mathcal{G}}(\mathbf{u})$$

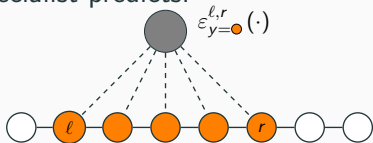
$$\mathbb{E}[\phi_{\mathcal{S}}(\mathbf{u})] \leq 2\mathbb{E}[\phi_{\mathcal{T}}(\mathbf{u})] = 2\phi^r(\mathbf{u})$$

Cluster Specialists

Cluster Specialists

- A specialist is a *basis function* $\varepsilon : V \rightarrow \{-1, 1, \square\}$
- “ \square ” - on some inputs a specialist can offer no prediction
- Given \mathcal{S} denote the vertices $\{1, \dots, n\}$ in linear order
- For a vertex $v \in V$ a cluster specialist predicts:

$$\varepsilon_y^{\ell, r}(v) := \begin{cases} y & \ell \leq v \leq r \\ \square & \text{otherwise} \end{cases}$$

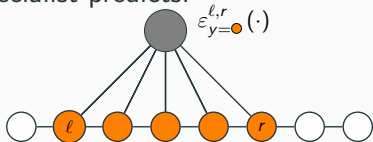


- How to construct a specialist set?
 - Needs to be *complete* (any labeling $\mathbf{u} \in \{-1, 1\}^{|V|}$ is covered)
 - The ‘covering set’ of a labeling should **not** be *too large*

Cluster Specialists

- A specialist is a *basis function* $\varepsilon : V \rightarrow \{-1, 1, \square\}$
- “ \square ” - on some inputs a specialist can offer no prediction
- Given \mathcal{S} denote the vertices $\{1, \dots, n\}$ in linear order
- For a vertex $v \in V$ a cluster specialist predicts:

$$\varepsilon_y^{\ell,r}(v) := \begin{cases} y & \ell \leq v \leq r \\ \square & \text{otherwise} \end{cases}$$

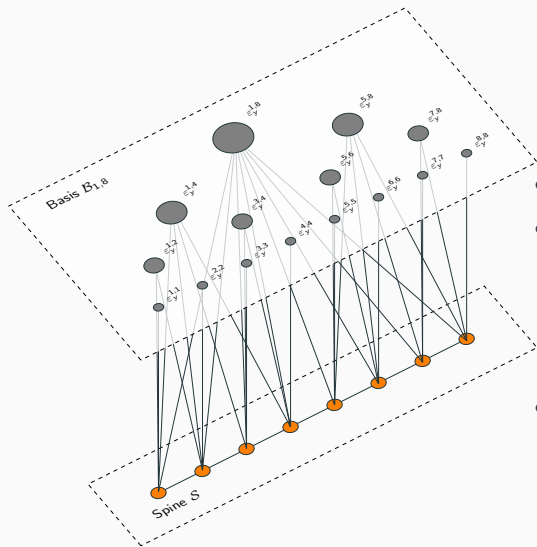


- Two specialist sets:

$$\mathcal{F}_n := \{\varepsilon_y^{\ell,r} : \ell, r \in [n], \ell \leq r; y \in \{-1, 1\}\}, \quad |\mathcal{F}_n| = \mathcal{O}(n^2)$$

$$\mathcal{B}_{m,n} := \begin{cases} \{\varepsilon_{-1}^{m,n}, \varepsilon_1^{m,n}\} & m = n \\ \{\varepsilon_{-1}^{m,n}, \varepsilon_1^{m,n}\} \cup \mathcal{B}_{m, \lfloor \frac{m+n}{2} \rfloor} \cup \mathcal{B}_{\lceil \frac{m+n}{2} \rceil, n} & m \neq n \end{cases}, \quad |\mathcal{B}_{1,n}| = \mathcal{O}(n)$$

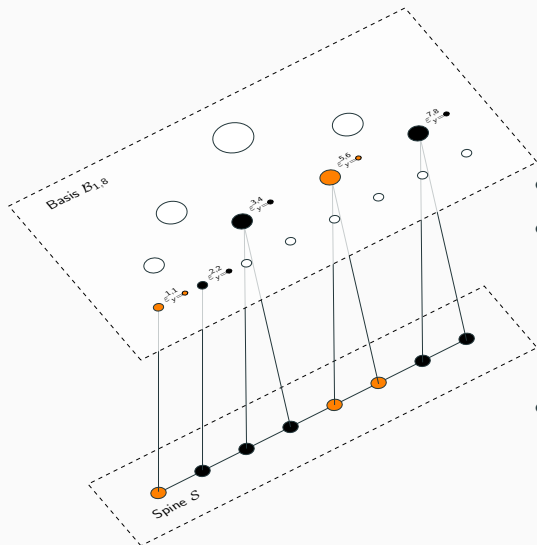
Basis Set $\mathcal{B}_{1,n}$



$$\varepsilon_y^{\ell,r}(v) := \begin{cases} y & \ell \leq v \leq r \\ \square & \text{otherwise} \end{cases}$$

- $\mathcal{B}_{1,n}$ is **complete**
- Maximum number of specialists required to cover a labeling $\mathbf{u} \in \{-1, 1\}^{|\mathcal{V}|}$ is bounded above by $2(\phi_S(\mathbf{u}) + 1) \lceil \log_2 \frac{n}{2} \rceil$
- Only $\Theta(\log n)$ specialists are 'active' at any given time ($\mathcal{O}(n^2)$ for basis set \mathcal{F}_n)

Basis Set $\mathcal{B}_{1,n}$



$$\varepsilon_y^{\ell,r}(v) := \begin{cases} y & \ell \leq v \leq r \\ \square & \text{otherwise} \end{cases}$$

- $\mathcal{B}_{1,n}$ is **complete**
- Maximum number of specialists required to cover a labeling $\mathbf{u} \in \{-1, 1\}^{|V|}$ is bounded above by $2(\phi_S(\mathbf{u}) + 1) \lceil \log_2 \frac{n}{2} \rceil$
- Only $\Theta(\log n)$ specialists are 'active' at any given time ($\mathcal{O}(n^2)$ for basis set \mathcal{F}_n)

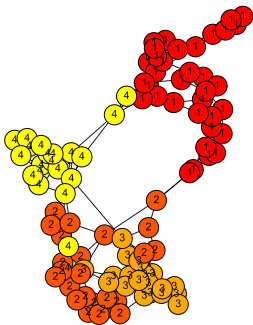
Specialists Example - USPS



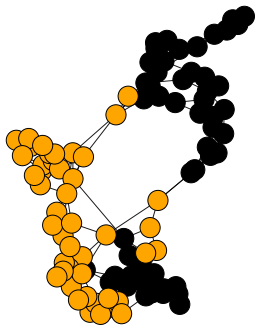
Image Source: [YHW18]

- USPS Dataset (hand-written digits)
- 16×16 pixels \rightarrow points in \mathbb{R}^{256}
- Build graph by connecting each point with its 3 nearest neighbors

Specialists Example - USPS

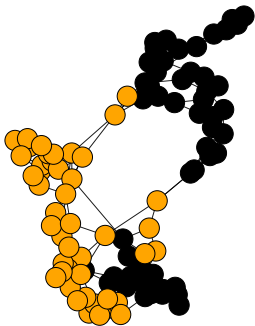


Original Graph \mathcal{G}

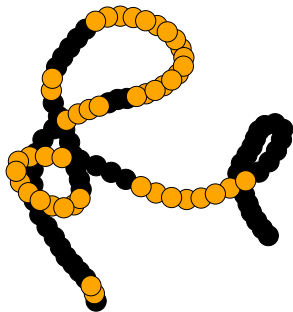


Simulated binary labeling

Specialists Example - USPS

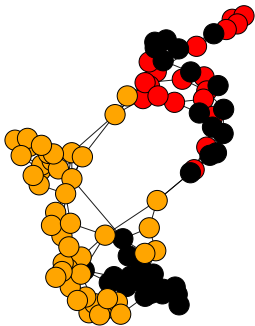


Original Graph \mathcal{G}

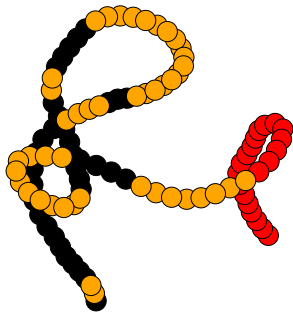


Linear Embedding (Spine) \mathcal{S}

Specialists Example - USPS



Original Graph \mathcal{G}



Linear Embedding (Spine) \mathcal{S}

Algorithm 1: SWITCHING CLUSTER SPECIALISTS

```

input      : Specialists set  $\mathcal{E}$ 
parameter :  $\alpha \in [0, 1]$ 
initialize :  $\omega_1 \leftarrow \frac{1}{|\mathcal{E}|}\mathbf{1}, \dot{\omega}_0 \leftarrow \frac{1}{|\mathcal{E}|}\mathbf{1}, \mathbf{p} \leftarrow \mathbf{0}, m \leftarrow 0$ 

for  $t = 1$  to  $T$  do
  receive  $i_t \in V$ 
  set  $\mathcal{A}_t := \{\varepsilon \in \mathcal{E} : \varepsilon(i_t) \neq \square\}$ 
  foreach  $\varepsilon \in \mathcal{A}_t$  do
    
$$\omega_{t,\varepsilon} \leftarrow (1 - \alpha)^{m-p_\varepsilon} \dot{\omega}_{t-1,\varepsilon} + \frac{1 - (1 - \alpha)^{m-p_\varepsilon}}{|\mathcal{E}|}$$

    // delayed share update
  predict  $\hat{y}_t \leftarrow \text{sign}(\sum_{\varepsilon \in \mathcal{A}_t} \omega_{t,\varepsilon} \varepsilon(i_t))$ 
  receive  $y_t \in \{-1, 1\}$ 
  set  $\mathcal{Y}_t := \{\varepsilon \in \mathcal{E} : \varepsilon(i_t) = y_t\}$ 
  if  $\hat{y}_t \neq y_t$  then
    // loss update
    
$$\dot{\omega}_{t,\varepsilon} \leftarrow \begin{cases} 0 & \varepsilon \in \mathcal{A}_t \cap \bar{\mathcal{Y}}_t \\ \dot{\omega}_{t-1,\varepsilon} & \varepsilon \notin \mathcal{A}_t \\ \omega_{t,\varepsilon} \frac{\omega_t(\mathcal{A}_t)}{\omega_t(\mathcal{Y}_t)} & \varepsilon \in \mathcal{Y}_t \end{cases}$$

    (2)
    foreach  $\varepsilon \in \mathcal{A}_t$  do
       $p_\varepsilon \leftarrow m$ 
     $m \leftarrow m + 1$ 
  else
     $\dot{\omega}_t \leftarrow \dot{\omega}_{t-1}$ 

```

Algorithm Intuition

- Weight vector $\omega_t \in [0, 1]^{|\mathcal{E}|}$ maintained
- Weight $\omega_{t,\varepsilon}$ corresponds to our ‘confidence’ in specialist ε
- On each trial set “active” specialists
 $\mathcal{A}_t := \{\varepsilon \in \mathcal{E} : \varepsilon(i_t) \neq \square\}$
- Take the weighted-majority vote of specialists in \mathcal{A}_t
- Decrease weight of *incorrect* specialists
- Increase weight of *correct* specialists
- Share some of the weight among all specialists after each update (can be done efficiently)

Mistake Bound Guarantees

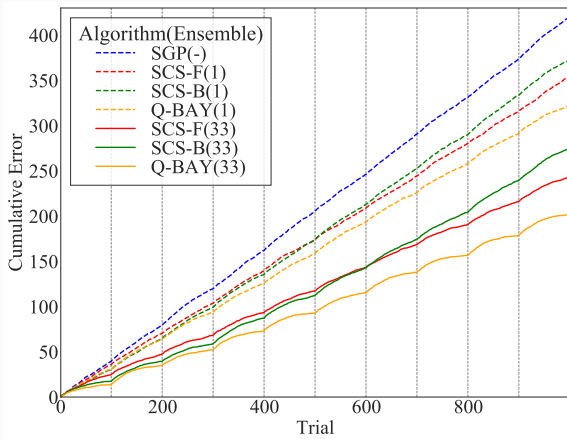
For a sequence of K distinct labelings $\mathbf{u}^1, \dots, \mathbf{u}^K$, let

$$H_k := \sum_{(i,j) \in E_S} \left[\left[[\mathbf{u}^k(i) \neq \mathbf{u}^k(j)] \vee [\mathbf{u}^{k+1}(i) \neq \mathbf{u}^{k+1}(j)] \right] \wedge \right. \\ \left. \left[[\mathbf{u}^k(i) \neq \mathbf{u}^{k+1}(i)] \vee [\mathbf{u}^k(j) \neq \mathbf{u}^{k+1}(j)] \right] \right]$$

	Static Bounds
[HLP08]	$\mathcal{O} \left(\phi_{\mathcal{G}}(\mathbf{u}) \log \frac{n}{\phi_{\mathcal{G}}(\mathbf{u})} + \phi_{\mathcal{G}}(\mathbf{u}) \right)$
[HP06]	$\mathcal{O}(\phi_{\mathcal{G}}(\mathbf{u}) R_{\mathcal{G}})$
Switching Mistake Bounds	
\mathcal{F}_n	$\mathcal{O} \left(\phi_{\mathcal{G}}(\mathbf{u}_1) \log n + \sum_{k=1}^{K-1} H_k (\log n + \log K + \log \log T) \right)$
$\mathcal{B}_{1,n}$	$\mathcal{O} \left(\left(\phi_{\mathcal{G}}(\mathbf{u}_1) \log n + \sum_{k=1}^{K-1} H_k (\log n + \log K + \log \log T) \right) \log n \right)$
Time Complexity (per trial)	
\mathcal{F}_n	$\mathcal{O}(n^2)$
$\mathcal{B}_{1,n}$	$\mathcal{O}(\log n)$

Experiments

Experiments



Mean cumulative error over 12 iterations of 10 switches every 100 trials on an 4096-vertex graph. Solid lines SCS-F and SCS-B show the mean cumulative error of an ensemble size of 33, dashed lines show the average cumulative error of a single instance (ensemble size 1).

Conclusion

Conclusion

- Solved the problem of efficient online prediction of switching graph labelings
- Described the machinery of Cluster Specialists
- Proved *smooth* mistake bounds
- Exponential speed up with $\mathcal{B}_{1,n}$
- Future work:
 - New methods of constructing specialist sets (e.g., hierarchical clustering)
 - Further experiments

Thank you!

(Thank you to Fabio Vitale for some slides)

References



N. Cesa-Bianchi, C. Gentile, and F. Vitale, *Fast and optimal prediction on a labeled tree*, Proceedings of the 22nd Annual Conference on Learning Theory, Omnipress, 2009, pp. 145–156.



M. Herbster, G. Lever, and M. Pontil, *Online prediction on large diameter graphs*, Proceedings of the 21st International Conference on Neural Information Processing Systems, NIPS '08, 2008, pp. 649–656.



M. Herbster and M. Pontil, *Prediction on a graph with a perceptron*, Proceedings of the 19th International Conference on Neural Information Processing Systems, NIPS'06, 2006, pp. 577–584.



Y. Yang, Y. Hu, and F. Wu, *Sparse and low-rank subspace data clustering with manifold regularization learned by local linear embedding*, Applied Sciences **8** (2018), 2175.