

Exploiting Cluster Structure to Predict the Labeling of a Graph

ALT 2008

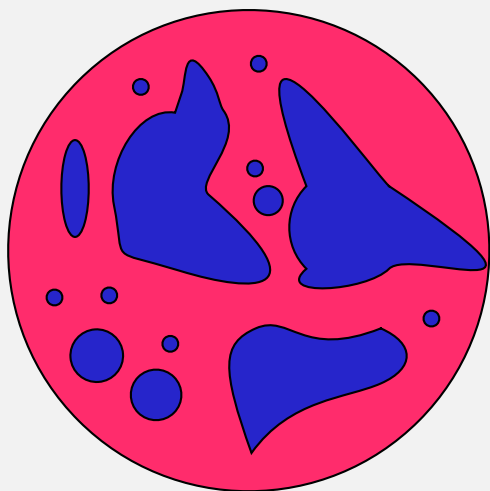
Mark Herbster

University College London
Department of Computer Science

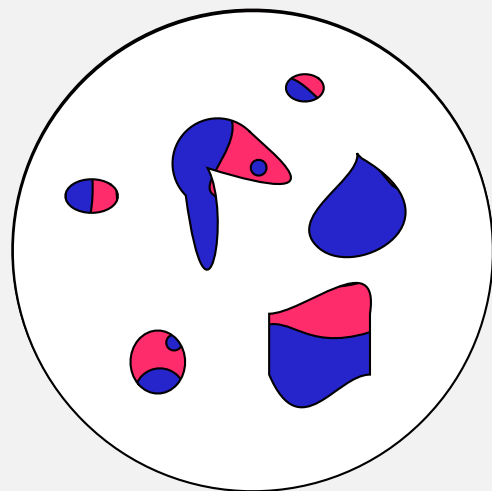
16 October, 2008

Overview

1. Give perceptron-like algorithm for graph label prediction
2. Improve on Perceptron bound when **cluster-structure**



“Default” Assumption



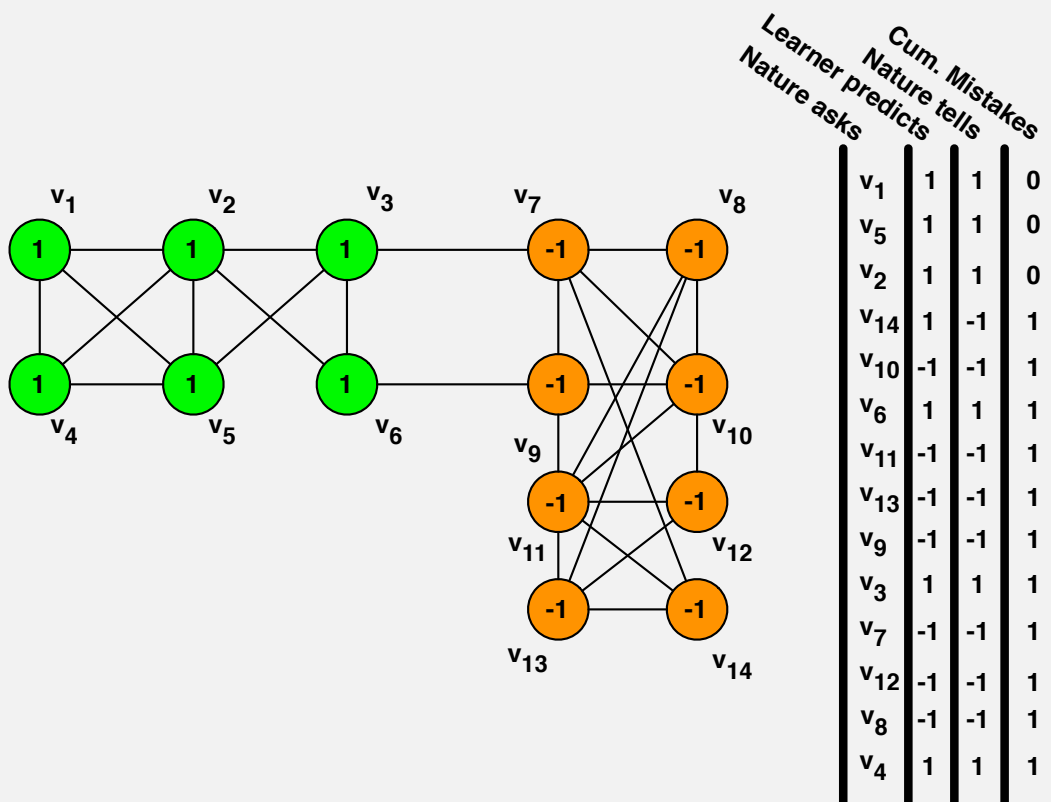
“Cluster” Assumption

–

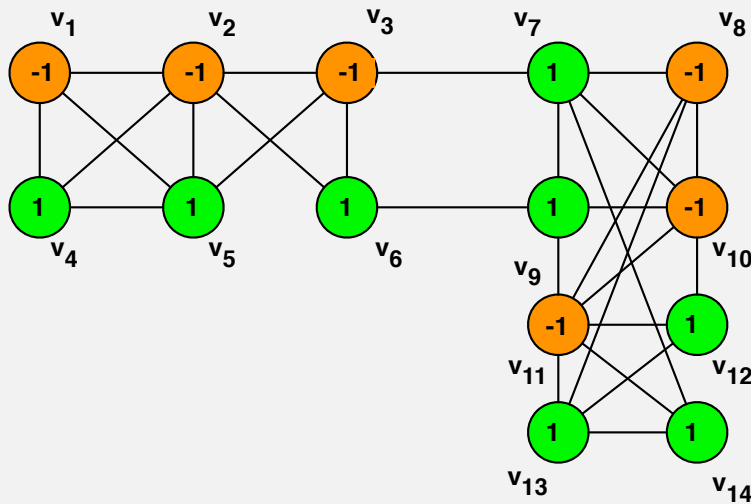
Outline

- ▶ Review: Online graph label prediction
- ▶ Review: Predicting labeling of a graph with a perceptron
- ▶ Problem: Multiple clusters (**Perceptron fails**)
- ▶ Solution: Pounce algorithm
- ▶ Bound: Via cover of input space
- ▶ Application: Heat kernel
- ▶ Moral: Value of unlabeled data

A prediction game



Let's try again



	Learner predicts	Nature tells	Cum. Mistakes
v1	1	-1	1
v5	-1	1	2
v2	1	-1	3
v14	-1	1	4
v10	1	-1	5
v6	-1	1	6
v11	1	-1	7
v13	-1	1	8
v9	-1	1	9
v3	1	-1	10
v7	-1	1	11
v12	-1	1	12
v8	1	-1	13
v4	-1	1	14

Online Learning Model

- ▶ Aim: learn a function $\mathbf{u} : V \rightarrow \{-1, +1\}$ corresponding to a labeling of a graph $G = (V, E)$ and $V = \{1, \dots, n\}$.
- ▶ Learning proceeds in trials
 - for** $t = 1, \dots, \ell$ **do**
 1. Nature selects $v_t \in V$
 2. Learner predicts $\hat{y}_t \in \{-1, +1\}$
 3. Nature selects $y_t \in \{-1, +1\}$
 4. If $\hat{y}_t \neq y_t$ then mistakes = mistakes + 1
- ▶ Learner's goal: *minimize* mistakes
- ▶ Bound: mistakes $\leq f(\text{complexity}(\mathbf{u}))$

Perceptron Bound (Novikoff)

Theorem [Novikoff]:

Given a sequence $\{(\mathbf{x}_t, y_t)\}_{t=1}^{\ell} \subseteq \mathcal{H} \times \{-1, 1\}$ then the mistakes of the perceptron are bounded by

$$M \leq \|\mathbf{u}\|^2 R$$

with $R = \max_t(\|\mathbf{x}_t\|^2)$ for all $\mathbf{u} \in \mathcal{H}$ such that

$$\langle \mathbf{u}, \mathbf{x}_t \rangle y_t \geq 1$$

for $t = 1, \dots, \ell$.

Signed Laplacian

Matrix-Edge Decomposition

If \mathbf{G} is a symmetric matrix then

$$\mathbf{G} = \sum_{(i,j) \in E^+} a_{ij}(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top + \sum_{(i,j) \in E^-} b_{ij}(\mathbf{e}_i + \mathbf{e}_j)(\mathbf{e}_i + \mathbf{e}_j)^\top + \sum_{i=1}^n c_i \mathbf{e}_i \mathbf{e}_i^\top$$

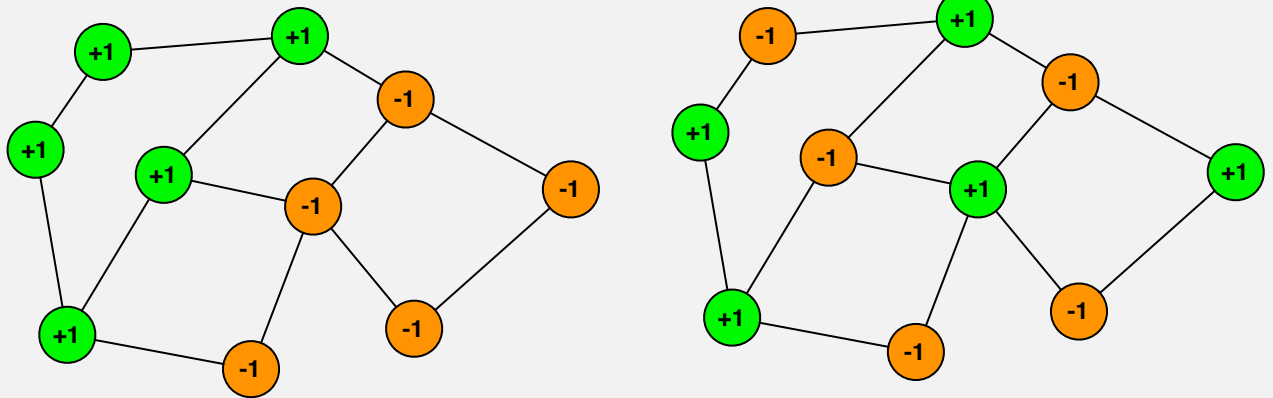
where $(\mathbf{e}_i)_j = \delta_{ij}$ with $a_{ij}, b_{ij} > 0$ and $c_i \in \mathbb{R}$ and thus

$$\|\mathbf{u}\|_{\mathbf{G}}^2 = \mathbf{u}^\top \mathbf{G} \mathbf{u} = \sum_{(i,j) \in E^+} a_{ij}(u_i - u_j)^2 + \sum_{(i,j) \in E^-} b_{ij}(u_i + u_j)^2 + \sum_{i=1}^n c_i u_i^2$$

Definitions

- ▶ \mathbf{G} is a graph Laplacian if $E^- = \emptyset$ and $c_1 = \dots = c_n = 0$
- ▶ \mathbf{G} is a signed graph Laplacian if $c_1 \geq 0, \dots, c_n \geq 0$ (p.s.d.)

Examples



$$\begin{array}{ll} \|\mathbf{u}\|^2 = 3 \times 4 & \text{positive unit-edges} \\ \|\mathbf{u}\|^2 = 10 \times 4 & \text{negative unit-edges} \end{array} \quad \begin{array}{ll} \|\mathbf{u}\|^2 = 12 \times 4 & \\ \|\mathbf{u}\|^2 = 1 \times 4 & \end{array}$$

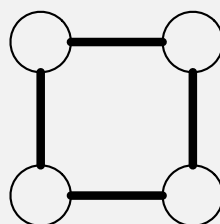
Resistance distance

The *resistance distance* [KR93] between vertex \mathbf{v}_p and \mathbf{v}_q ,

$$\|\mathbf{v}_p - \mathbf{v}_q\|_{\mathbf{G}}^2 = (\mathbf{e}_p - \mathbf{e}_q)^\top \mathbf{G}^+ (\mathbf{e}_p - \mathbf{e}_q)$$

is the *effective resistance* between \mathbf{v}_p and \mathbf{v}_q . The graph is the circuit and edge (i, j) has resistance a_{ij}^{-1} .

- ▶ Resistance Diameter : $R_{\mathbf{G}} := \max_{p, q \in V} \|\mathbf{v}_p - \mathbf{v}_q\|_{\mathbf{G}}^2$
- ▶ Geodesic distance upper bounds the resistance distance



Resistance Diameter: $R_{\mathbf{G}} = 1$

Predicting the Labeling of a Graph with the Perceptron

Theorem[HP06]:

The mistakes of perceptron are bounded by

$$M \leq 2\|\mathbf{u}\|^2 R_{\mathbf{G}} + 2$$

for all consistent labelings $\mathbf{u} \in [-1, 1]^n$.

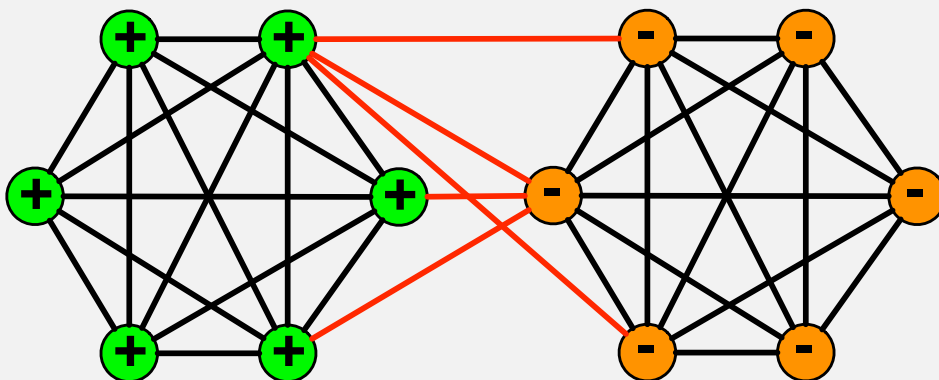
Proof.

If \mathbf{G} is a Laplacian use kernel $\mathbf{G}^+ + \mathbf{1}\mathbf{1}^\top R_{\mathbf{G}}$.

Observations

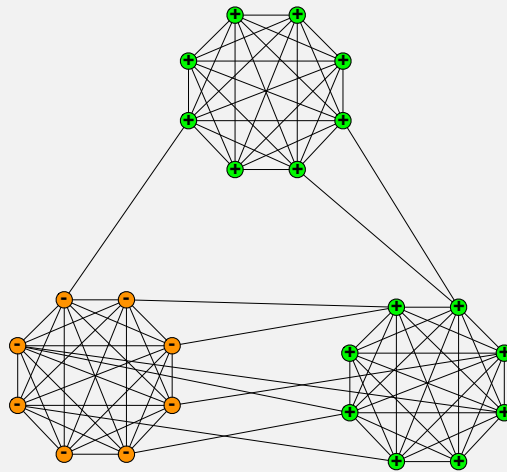
- ▶ Optimal \mathbf{u}^* are the voltages which minimizes energy
- ▶ Optimal $\|\mathbf{u}^*\|^2$ is bounded by each label-separating cut.

Bound: 2 Prototype Clusters



- ▶ Two m -cliques with ℓ edges ($\ell < m$) between cliques
- ▶ Norm: $\|\mathbf{u}\|^2 = 4\ell$
- ▶ Resistance diameter: $R_{\mathbf{G}} \leq 5/\ell$
- ▶ Perceptron: $M \leq 42$ (independent of ℓ and m)
- ▶ Does this generalize to multiple clusters? No!

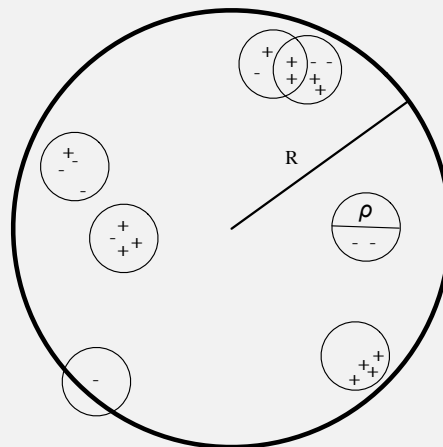
Problem: 3 Clusters



3 clusters (one in isolation)

- ▶ Three m -cliques
- ▶ Two m -cliques with ℓ edges ($cm < \ell < m$) between cliques
- ▶ An “isolated” clique with $\Theta(1)$ outgoing edges.
- ▶ Norm: $\|\mathbf{u}\|^2 = \Theta(\ell)$, Resistance diameter: $R_{\mathbf{G}} = \Theta(1)$
- ▶ **Problem: perceptron: $M \leq \Theta(\ell)$ (dependent on ℓ)**

Pounce Bound Motivation



Input space X of radius R with cover number $\mathcal{N}(X, \rho) = 7$.

- ▶ Bounds to be dependent on structure of input space X .
- ▶ Novikoff is only dependent on X through radius R .
- ▶ Expectation is that a typical ambient input space is only sparsely populated (cf manifold/cluster hypotheses).
- ▶ Pounce will depend on the **cover** of X .
- ▶ In particular the number of balls $\mathcal{N}(X, \rho)$ of diameter ρ .

Pounce Algorithm

Notation: $\mathcal{V}_{\mathbf{G}} := \{\mathbf{v}_i := \mathbf{e}_i^\top \mathbf{G}^+ : i \in \{1, 2, \dots, n\}\}$

Input: $\{(\mathbf{v}_{i_t}, y_t)\}_{t=1}^\ell \subseteq \mathcal{V}_{\mathbf{G}} \times \{-1, 1\}$.

Initialization: $\mathbf{w}_2 = \mathbf{0}$; $\mathcal{M} = \{1\}$.

for $t = 2, \dots, \ell$ **do**

Receive: $i_t \in \{1, \dots, n\}$

$\eta_t = \arg \min_{j \in \mathcal{M}} \|\mathbf{v}_{i_t} - \mathbf{v}_{i_j}\|$

Predict: $\hat{y}_t = \text{sign}(y_{\eta_t} + \mathbf{w}_t(i_t) - \mathbf{w}_t(i_{\eta_t}))$

Receive: y_t

if $\hat{y}_t = y_t$ **then**

$\mathbf{w}_{t+1} = \mathbf{w}_t$

else

$\mathbf{w}_{t+1} = \mathbf{w}_t + \frac{y_t - y_{\eta_t} - (\mathbf{w}_t(i_t) - \mathbf{w}_t(i_{\eta_t}))}{\|\mathbf{v}_{i_t} - \mathbf{v}_{i_{\eta_t}}\|^2} (\mathbf{v}_{i_t} - \mathbf{v}_{i_{\eta_t}})$

$\mathcal{M} = \mathcal{M} \cup \{t\}$

end

Pounce Bound

Theorem

The mistakes M of POUNCE are bounded by

$$M \leq \mathcal{N}(X, \rho) + \|\mathbf{u}\|^2 \rho + 1,$$

for all $0 < \rho$, and for all $\mathbf{u} \in \mathbb{R}^n$ such that

$$\mathbf{u}(i_t) y_t \geq 1$$

for all $t = 1, \dots, \ell$.

- ▶ **Definition:** $\mathcal{N}(X, \rho)$ is the minimum number of balls of squared diameter ρ that cover X .
- ▶ **Three Clusters:** $|\mathcal{M}| \leq 20$ with $(\rho = \frac{2}{m-1}, \mathcal{N}(X, \rho) = 3,$
and $\|\mathbf{u}^*\|^2 < 8\ell)$

Application: Heat Kernel

The Heat Kernel is used to build Laplacian from data.

Discrepancy Function

- ▶ $d : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$
- ▶ $d(x, y) = d(y, x)$
- ▶ $(d(x, y) = 0) \iff (x = y)$

Heat Kernel Laplacian

Given $X = \{x_1, \dots, x_n\} \subseteq \mathcal{X}$, discrepancy d , and $a > 0$ then

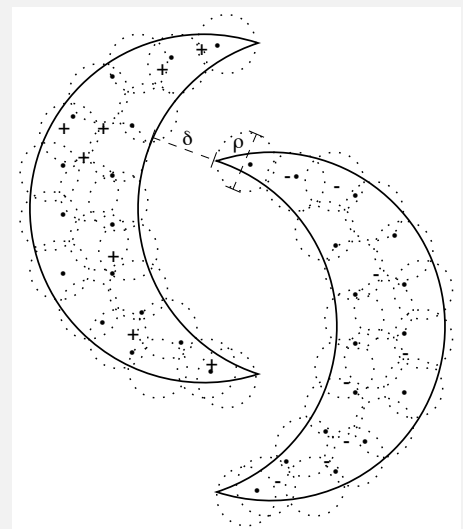
$$G_{ij}^a := \begin{cases} -e^{-ad(x_i, x_j)} & i \neq j \\ \sum_{k \neq i}^n e^{-ad(x_i, x_k)} & i = j \end{cases}$$

a is a scale parameter

(Component) Separating Cover

Definitions

- ▶ $\mathcal{N}^\circ(\mathcal{X}, \mathcal{Y}, d)$ is the *separating cover number* the minimal number of **balls** of diameter ρ covering \mathcal{X} with $\mathcal{Y}(x) \neq \mathcal{Y}(x') \rightarrow d(x, x') > \rho$
- ▶ $x, x' \in X$ are ρ -*path-connected* if $d(x, x') \leq \rho$ or $\exists x''$ such that $d(x, x'') \leq \rho$ and x'', x' are ρ -path-connected.
- ▶ $\mathcal{C}^\circ(X, \mathcal{Y}, d)$ is the *component-separating cover number* ... **ρ -path-connected sets** ... $\mathcal{Y}(x) \neq \mathcal{Y}(x') \rightarrow d(x, x') > \rho$.



$$\mathcal{N}^\circ(\mathcal{X}, \mathcal{Y}, d) = 42$$

$$\mathcal{C}^\circ(X, \mathcal{Y}, d) = 2$$

Heat Kernel Bound

Theorem

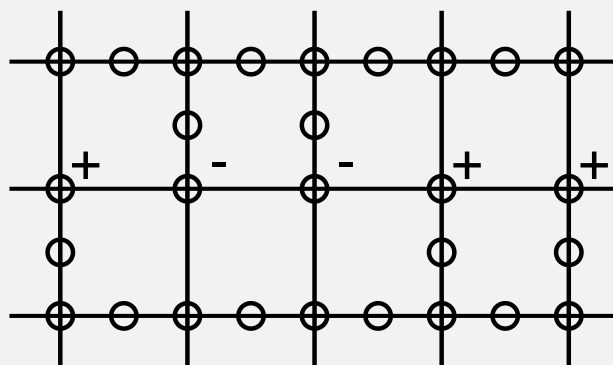
There exists an $a' > 0$ such that for all $a > a'$ the mistakes M of the POUNCE algorithm with Laplacian \mathbf{G}^a are bounded by

$$M \leq \mathcal{C}^\circ(X, \mathcal{Y}, d) + 1 \leq \mathcal{N}^\circ(\mathcal{X}, \mathcal{Y}, d) + 1 \quad (*)$$

Observe

- ▶ As X increases potentially $\mathcal{C}^\circ(X, \mathcal{Y}, d) \ll \mathcal{N}^\circ(\mathcal{X}, \mathcal{Y}, d)$

The Value of Unlabeled Data



Two Path connected components

- ▶ Center row of m labels is the task.
- ▶ Center row is labeled randomly.
- ▶ Expected mistakes of any algorithm not using unlabeled data is $M = \frac{m}{2}$.
- ▶ **Aligned** unlabeled data \rightarrow POUNCE's bound is $M \leq 3$.
- ▶ Example implies a limitation in the type bounds provable for non-transductive algorithms.

Thanks

Thank You!