

Online Prediction on Large Diameter Graphs

Authored by:

Mark Herbster
Massimiliano Pontil
Guy Lever

Abstract

Current on-line learning algorithms for predicting the labelling of a graph have an important limitation in the case of large diameter graphs; the number of mistakes made by such algorithms may be proportional to the square root of the number of vertices, even when tackling simple problems. We overcome this problem with an efficient algorithm which achieves a logarithmic mistake bound. Furthermore, current algorithms are optimised for data which exhibits cluster-structure; we give an additional algorithm which performs well locally in the presence of cluster structure and on large diameter graphs.

1 Paper Body

We continue our study of online prediction of the labelling of a graph. We show a fundamental limitation of Laplacian-based algorithms: if the graph has a large diameter then the number of mistakes made by such algorithms may be proportional to the square root of the number of vertices, even when tackling simple problems. We overcome this drawback by means of an efficient algorithm which achieves a logarithmic mistake bound. It is based on the notion of a spine, a path graph which provides a linear embedding of the original graph. In practice, graphs may exhibit cluster structure; thus in the last part, we present a modified algorithm which achieves the “best of both worlds”: it performs well locally in the presence of cluster structure, and globally on large diameter graphs.

1

Introduction

We study the problem of predicting the labelling of a graph in the online learning framework. Consider the following game for predicting the labelling of a graph: Nature presents a graph; nature queries a vertex v_1 ; the learner predicts $y_1 \in \{1, 1\}$, the label of the vertex; nature presents a label y_1 ; nature queries a vertex v_2 ; the learner predicts y_2 ; and so forth. The learner’s goal is to minimise the total number of mistakes $M = \sum_{t=1}^T \mathbb{1}_{y_t \neq y_t}$. If nature is

adversarial, the learner will always mispredict, but if nature is regular or simple, there is hope that a learner may make only a few mispredictions. Thus, a central goal of online learning is to design algorithms whose total mispredictions can be bounded relative to the complexity of nature's labelling. In [9, 8, 7], the cut size (the number of edges between disagreeing labels) was used as a measure of the complexity of a graph's labelling, and mistake bounds relative to this and the graph diameter were derived. The strength of the methods in [8, 7] is in the case when the graph exhibits 'cluster structure'. The apparent deficiency of these methods is that they have poor bounds when the graph diameter is large relative to the number of vertices. We observe that this weakness is not due to insufficiently tight bounds, but is a problem in their performance. In particular, we discuss an example of a n -vertex labelled graph with a single edge between disagreeing label sets. On this graph, sequential prediction using the common method based upon minimising the Laplacian semi-norm of a labelling, subject to ℓ_1 constraints, incurs $\Omega(n)$ mistakes (see Theorem 3). The expectation is that the number of mistakes incurred by an optimal online algorithm is bounded by $O(\ln n)$. We solve this problem by observing that there exists an approximate structure-preserving embedding of any graph into a path graph. In particular the cut-size of any labelling is increased by no more than a factor of two. We call this embedding a spine of the graph. The spine is the foundation on which we build two algorithms. Firstly we predict directly on the spine with the 1-nearest-neighbor algorithm. We demonstrate that this equivalent to the Bayes-optimal classifier for a particular Markov random field. A logarithmic mistake bound for learning on a path graph follows by the Halving algorithm analysis. Secondly, we use the spine of the graph as a foundation to add a binary support tree to the original graph. This enables us to prove a bound which is the 'best of both worlds' – if the predicted set of vertices has cluster-structure we will obtain a bound appropriate for that case, but if instead, the predicted set exhibits a large diameter we will obtain a polylogarithmic bound.

Previous work. The seminal approach to semi-supervised learning over graphs in [3] is to predict with a labelling which is consistent with a minimum label-separating cut. More recently, the graph Laplacian has emerged as a key object in semi-supervised learning, for example the semi-norm induced by the Laplacian is commonly either directly minimised subject to constraints, or used as a regulariser [14, 2]. In [8, 7] the online graph labelling problem was studied. An aim of those papers was to provide a natural interpretation of the bound on the cumulative mistakes of the kernel perceptron when the kernel is the pseudoinverse of the graph Laplacian – bounds in this case being relative to the cut and (resistance) diameter of the graph. In this paper we necessarily build directly on the very recent results in [7] as those results depend on the resistance diameter of the predicted vertex set as opposed to the whole graph [8]. The online graph labelling problem is also studied in [13], and here the graph structure is not given initially. A slightly weaker logarithmic bound for the online graph labelling problem has also been independently derived via a connection to an online routing problem in the very recent [5].

Preliminaries

We study the process of predicting a labelling defined on the vertices of a graph. Following the classical online learning framework, a sequence of labelled vertices $\{(v_1, y_1), (v_2, y_2), \dots\}$, the trial sequence, is presented to a learning algorithm such that, on sight of each vertex v_t , the learner makes a prediction \hat{y}_t for the label value, after which the correct label is revealed. This feedback information is then used by the learning algorithm to improve its performance on further examples. We analyse the performance of a learning algorithm in the mistake bound framework [12] – the aim is to minimise the maximum possible cumulative number of mistakes made on the training sequence. A graph $G = (V, E)$ is a collection of vertices $V = \{v_1, \dots, v_n\}$ joined by connecting (possibly weighted) edges. Denote $i \sim j$ whenever v_i and v_j are connected so that $E = \{(i, j) : i \sim j\}$ is the set of unordered pairs of connected vertex indices. Associated with each edge $(i, j) \in E$ is a weight A_{ij} , so that A is the $n \times n$ symmetric adjacency matrix. We say that G is unweighted if $A_{ij} = 1$ for every $(i, j) \in E$ and is 0 otherwise. In this paper, we consider only connected graphs – that is, graphs such that there exists a path between any two vertices. The Laplacian L of a graph G is the $n \times n$ matrix $L = D - A$, where D is the diagonal degree matrix such that $D_{ii} = \sum_j A_{ij}$. The quadratic form associated with the Laplacian relates to the cut size of graph labellings. Definition 1. Given a labelling $u \in \mathbb{R}^n$ of $G = (V, E)$ we define the cut size of u by $\frac{1}{2} \sum_{(i,j) \in E} (u_i - u_j)^2$.

n

In particular, if $u \in \{-1, 1\}^n$ we say that a cut occurs on edge (i, j) if $u_i \neq u_j$ and $\frac{1}{2} \sum_{(i,j) \in E} (u_i - u_j)^2$ measures the number of cuts. We evaluate the performance of prediction algorithms in terms of the cut size and the resistance diameter of the graph. There is an established natural connection between graphs and resistive networks where each edge $(i, j) \in E$ is viewed as a resistor with resistance $1/A_{ij}$ [4]. Thus the effective resistance $r_G(v_i, v_j)$ between vertex v_i and v_j is the potential difference needed to induce a unit current flow between v_i and v_j . The effective resistance may be computed by the formula [11] $r_G(v_i, v_j) = (e_i - e_j)^T G^+ (e_i - e_j)$, where G^+ denotes the pseudoinverse and e_1, \dots, e_n are the canonical basis vectors of \mathbb{R}^n . The resistance diameter of a graph $R_G := \max_{v_i, v_j \in V} r_G(v_i, v_j)$ is the maximum effective resistance between any pair of vertices on the graph.

3

Limitations of online minimum semi-norm interpolation

As we will show, it is possible to develop online algorithms for predicting the labelling of a graph which have a mistake bound that is a logarithmic function of the number of vertices. Conversely, we first highlight a deficiency in a standard Laplacian based method for predicting a graph labelling. Given a partially labelled graph $G = (V, E)$ with $|V| = n$ that is, such that for some $\ell \leq n$, $y' \in \{-1, 1\}^\ell$ is a labelling defined on the ℓ vertices $V' = \{v_1, v_2, \dots, v_\ell\}$ the minimum semi-norm interpolant is defined by $y' = \operatorname{argmin}\{u^T G u : u \in \mathbb{R}^n, u_{ik} = y'_k, k = 1, \dots, \ell\}$.

We then predict using $\hat{y}_i = \operatorname{sgn}(y'_i)$, for $i = 1, \dots, n$. The common

justification behind the above learning paradigm [14, 2] is that minimizing the cut (1) encourages neighbouring vertices to be similarly labelled. However, we now demonstrate that in the online setting such a regime will perform poorly on ? certain graph constructions ? there exists a trial sequence on which the method will make at least $\Omega(n)$ mistakes. Definition 2. An octopus graph of size d is defined to be d path graphs (the tentacles) of length d (that is, with $d + 1$ vertices) all adjoined at a common end vertex, to which a further single head vertex is attached, so that $n = |V| = d^2 + 2$. This corresponds to the graph $O(1, d, d)$ discussed in [8]. Theorem 3. Let $G = (V, E)$ be an octopus graph of size d and $y = (y_1, \dots, y_{|V|})$ the labelling such that $y_i = 1$ if v_i is the head vertex and $y_i = 0$ otherwise. There exists a trial sequence for p which online minimum semi-norm interpolation makes $\Omega(d)$ mistakes. Proof. Let the first query vertex be the head vertex, and let the end vertex of a tentacle be queried at each subsequent trial. We show that this strategy forces at least d mistakes. The solution to the minimum semi-norm interpolation with boundary P_n values problem is precisely the harmonic solution [4] y^* (that is, for every unlabeled vertex v_j , $i=1 \dots d$ $A_{ij} (y_i - y_j) = 0$). If the graph is connected y^* is unique and the graph labelling problem is identical to that of identifying the potential at each vertex of a resistive network defined on the graph where each edge corresponds to a resistor of 1 unit; the harmonic principle corresponds to Kirchoff's current law in this case. Using this analogy, suppose that the end points of k tentacles are labelled and that the end vertex v_q of an unlabelled tentacle is queried. Suppose a current of k flows from the head to the body of the graph. By Kirchoff's law, a current of k flows along each labelled tentacle (in order to obey the harmonic principle at every vertex it is clear that no current flows along the unlabelled tentacles). By Ohm's law $k = d + k$. Minimum semi-norm interpolation therefore results in the solution $2k y^*_q = 1$ iff $k \leq d$. Hence the minimum semi-norm solution predicts incorrectly whenever $k > d$ and the algorithm makes at least d mistakes. The above demonstrates a limitation in the method of online Laplacian minimum semi-norm interpolation for predicting a graph labelling ? the mistake bound can be proportional to the square root of the number of data points. We solve these problems in the following section.

4

A linear graph embedding

We demonstrate a method of embedding data represented as a connected graph G into a path graph, we call it a spine of G , which partially preserves the structure of G . Let P_n be the set of path graphs with n vertices. We would like to find a path graph with the same vertex set as G , which solves $\min_{P \in P_n} \max_{u \in \{1, \dots, n\}} |P(u) - G(u)|$. If a Hamiltonian path H of G (a path on G which visits each vertex precisely once) exists, then (u) the approximation ratio is $|H(u) - G(u)| \leq 1$. The problem of finding a Hamiltonian path is NP-complete $G(u)$ however, and such a path is not guaranteed to exist. As we shall see, a spine S of G may be found $S(u)$ efficiently and satisfies $|S(u) - G(u)| \leq 2$. We now detail the construction of a spine of a graph $G = (V, E)$, with $|V| = n$. Starting from any node, G is traversed in the manner of a depth-first search (that is, each

vertex is fully explored before backtracking to the last unexplored vertex), and an ordered list $VL = \{v_1, v_2, \dots, v_{2m+1}\}$ of the vertices ($m \geq 0$) in the order that they are visited is formed, allowing repetitions when a vertex is visited more than once. Note that each edge in EG is traversed no more than twice when forming VL . Define an edge multiset $EL = \{(v_1, v_2), (v_2, v_3), \dots, (v_{2m}, v_{2m+1})\}$ the set of pairs of consecutive vertices in VL . Let u be an arbitrary labelling of G and denote, as usual, $P(u) = \sum_{(i,j) \in EG} (u_i - u_j)^2$ and $L(u) = \sum_{(i,j) \in EL} (u_i - u_j)^2$. Since the multiset EL contains every element of EG no more than twice, $L(u) \leq 2P(u)$. We then take any subsequence VL_0 of VL containing every vertex in V exactly once. A spine $S = (V, ES)$ is a graph formed by connecting each vertex in V to its immediate neighbours in

the subsequence VL_0 with an edge. Since a cut occurs between connected vertices v_i and v_j in S only if a cut occurs on some edge in EL located between the corresponding vertices in the list VL we have $S(u) \leq L(u) \leq 2P(u)$.

(3)

Thus we have reduced the problem of learning the cut on a generic graph to that of learning the cut on a path graph. In the following we see that 1-nearest neighbour (1-NN) algorithm is a Bayes optimal algorithm for this problem. Note that the 1-NN algorithm does not perform well on general graphs; on the octopus graph discussed above, for example, it can make at least $\Omega(n)$ mistakes, and even $\Omega(n)$ mistakes on a related graph construction [8].

5

Predicting with a spine

We consider implementing the 1-NN algorithm on a path graph and demonstrate that it achieves a mistake bound which is logarithmic in the length of the line. Let $G = (V, E)$ be a path graph, where $V = \{v_1, v_2, \dots, v_n\}$ is the set of vertices and $E = \{(1, 2), (2, 3), \dots, (n-1, n)\}$. The nearest neighbour algorithm, in the standard online learning framework described above, attempts to predict a graph labelling by producing, for each query vertex v_i , the prediction y_i^t which is consistent with the label of the closest labelled vertex (and predicts randomly in the case of a tie). Theorem 4. Given the task of predicting the labelling of any unweighted, n -vertex path graph P in the online framework, the number of mistakes, M , incurred by the 1-NN algorithm satisfies

$$M \leq P(u) \log_2 n + 1, \quad (4)$$
 where $u \in \{0, 1\}^n$ is any labelling consistent with the trial sequence. Proof. We shall prove the result by noting that the Halving algorithm [1] (under certain conditions on the probabilities assigned to each hypothesis) implements the nearest neighbour algorithm on a path graph. Given any input space X and finite binary concept class $C \subseteq \{0, 1\}^X$, the Halving algorithm learns any target concept $c \in C$ as follows. Each hypothesis $c \in C$ is given an associated probability $p(c)$. A sequence of labelled examples $\{(x_1, y_1), \dots, (x_t, y_t)\} \subseteq X \times \{0, 1\}$, is revealed in accordance with the usual online framework. Let F_t be the set of feasible hypotheses at trial t ; $F_t = \{c : c(x_s) = y_s \text{ for } s \leq t\}$. Given an unlabelled example $x_t \in X$ at trial t the predicted label y_t^t is that which agrees with the majority vote — that is, such that it predicts randomly if this is equal to most

MH mistakes with

$$\begin{aligned} & \frac{1}{2} \cdot \\ & \frac{c(F_t, c(x_t))}{P} = y_t \\ & \frac{c(F_t)}{p(c)} \\ & \frac{p(c)}{p(c)} \\ & \frac{1}{2} \\ & \text{(and} \\ & \text{It is well known [1] that the Halving algorithm makes at} \\ & \text{MH } \log 2 \\ & \frac{1}{2} p(c) \\ & \cdot \\ & (5) \end{aligned}$$

We now define a probability distribution over the space of all labellings $u \in \{0, 1\}^n$ of P such that the Halving algorithm with these probabilities implements the nearest neighbour algorithm. Let a cut occur on any given edge with probability $\frac{1}{2}$, independently of all other cuts; $\text{Prob}(u_{i+1} \neq u_i) = \frac{1}{2}$ $\forall i \leq n$. The position of all cuts fixes the labelling up to flipping every label, and each of these two resulting possible arrangements are equally likely. This recipe associates with each possible labelling $u \in \{0, 1\}^n$ a probability $p(u)$ which is a function of the labelling's cut size $1/2 P(u) = (1/2)^n P(u)$. (6) This induces a full joint probability distribution on the space of vertex labels. In fact (6) is a Gibbs measure and as such defines a Markov random field over the space of vertex labels [10]. The mass function p therefore satisfies the Markov property $p(u) =$

$$p(u_i = ? \mid u_j = ?j \neq i) = p(u_i = ? \mid u_j = ?j \in N_i),$$

(7) where here N_i is the set of vertices neighbouring v_i — those connected to v_i by an edge. We will give an equivalent Markov property which allows a more general conditioning to reduce to that over boundary vertices.

Definition 5. Given a path graph $P = (V, E)$, a set of vertices $V_0 \subseteq V$ and a vertex $v_i \in V$, we define the boundary vertices v^l, v^r (either of which may be vacuous) to be the two vertices in V_0 that are closest to v_i in each direction along the path; its nearest neighbours in each direction. The distribution induced by (6) satisfies the following Markov property; given a partial labelling of P defined on a subset $V_0 \subseteq V$, the label of any vertex v_i is independent of all labels on V_0 except those on the vertices v^l, v^r (either of which could be vacuous) $p(u_i = ? \mid u_j = ?j, ?j : v_j \in V_0) = p(u_i = ? \mid u^l = ?l, u^r = ?r)$.

(8) Given the construction of the probability distribution formed by independent cuts on graph edges, we can evaluate conditional probabilities. For example, $p(u_j = ? \mid u_k = ?)$ is the probability of an even number of cuts between vertex v_j and vertex v_k . Since cuts occur with probability $\frac{1}{2}$ and there

are k possible arrangements of s cuts we have s

$$p(u_j = ? \mid u_k = ?) =$$

$$X - k \mid j - 1 \mid s \mid (1 \mid ? \mid ?) - k \mid j - s = (1 + (1 \mid ? \mid 2?) - k \mid j - s) \mid s \mid s \text{ even}$$

(9)

$$X - k \mid j - 1 \mid s \mid (1 \mid ? \mid ?) - k \mid j - s = (1 \mid ? \mid (1 \mid ? \mid 2?) - k \mid j - s) \mid s \mid 2$$

(10)

Likewise we have that $p(u_j = ? \mid u_k = ?) =$

s odd

Note also that for any single vertex we have $p(u_i = ?) = 1/2$ for $? \in \{1, 1\}$.

Lemma 6. Given the task of predicting the labelling of an n -vertex path graph online, the Halving algorithm, with a probability distribution over the labellings defined as in (6) and such that $0 \leq i \leq 1/2$, implements the nearest neighbour algorithm. Proof. Suppose that $t \geq 1$ trials have been performed so that we have a partial labelling of a subset $V_0 \subseteq V$, $\{(v_1, y_1), (v_2, y_2), \dots, (v_t, y_t)\}$. Suppose the label of vertex v_{t+1} is queried so that the Halving algorithm makes the following prediction y_{t+1} for vertex v_{t+1} : $y_{t+1} = y$ if $p(u_{t+1} = y \mid u_{ij} = y_j \mid 1 \leq j \leq t) \geq 1/2$, $y_{t+1} = ?$ if $p(u_{t+1} = y \mid u_{ij} = y_j \mid 1 \leq j \leq t) < 1/2$ (and predicts randomly if this probability is equal to $1/2$). We first consider the case where the conditional labelling includes vertices on both sides of v_{t+1} . We have, by (8), that $p(u_{t+1} = y \mid u_{ij} = y_j \mid 1 \leq j \leq t)$

$$= p(u_{t+1} = y \mid u' = y' \mid (r), u_r = y' \mid (r)) =$$

$$p(u' = y' \mid (r) \mid u_r = y' \mid (r), u_{t+1} = y) p(u_r = y' \mid (r), u_{t+1} = y) p(u' = y' \mid (r) \mid u_r = y' \mid (r))$$

=

$$p(u' = y' \mid (r) \mid u_{t+1} = y) p(u_r = y' \mid (r) \mid u_{t+1} = y) p(u' = y' \mid (r) \mid u_r = y' \mid (r))$$

(11)

where v' and v_r are the boundary vertices and (r) and (r) are trials at which vertices v' and v_r are queried, respectively. We can evaluate the right hand side of this expression using (9, 10). To show equivalence with the nearest neighbour method whenever $i \leq 1/2$, we have from (9, 10, 11) $p(u_{t+1} = y \mid u' = y, u_r = y) =$

=

$$(1 + (1 \mid ? \mid 2?) - (u_{t+1} \mid u' \mid u_r)) (1 \mid ? \mid (1 \mid ? \mid 2?) - (u_{t+1} \mid u' \mid u_r)) \mid 2 \mid (1 \mid ? \mid (1 \mid ? \mid 2?) - (u_{t+1} \mid u' \mid u_r))$$

which is greater than $1/2$ if $u_{t+1} \geq u' \geq u_r$ and less than $1/2$ if $u_{t+1} < u' < u_r$. Hence, this produces predictions exactly in accordance with the nearest neighbour scheme. We also have more simply that for all $i \leq 1/2$, u' and u_r and $i \leq 1/2$ $p(u_{t+1} = y \mid u' = y, u_r = y) \geq 1/2$, and $p(u_{t+1} = y \mid u' = y) \geq 1/2$.

This proves the lemma for all cases. A direct application of the Halving algorithm mistake bound (5) now gives

$$M \leq \log_2 \frac{1}{p(u)} = \log_2 \frac{1}{p(u)} \mid (1 \mid ? \mid ?) \mid n \mid 1 \mid ? \mid P(u)$$

$P(u) \geq 1$ where u is any labelling consistent with the trial sequence. We choose $P(u) = \min(1/n, 1/2)$ (note $P(u)$ that the bound is vacuous when $1/n \leq 1/2$ since M is necessarily upper bounded by n) giving

is vacuous while (13) is $M \leq 8c + 3 \cdot 2^d$ (with $c = m$, $N(X, \epsilon, rG) = 2$, and $\epsilon G(u) = cm$). An input space V may have both local cluster structure yet have a large diameter. Imagine a "universe" such that points are distributed into many dense clusters such that some sets of clusters are tightly packed but overall the distribution is quite diffuse. A given "problem" $X \subseteq V$ may then be centered on a few clusters or alternatively encompass the entire space. Thus, for practical purposes, we would like a prediction algorithm

which achieves the "best of both worlds", that is a mistake bound which is no greater, in order of magnitude, than the maximum of (12) and (13). The rest of this paper is directed toward this goal. We now introduce the notion of binary support tree, formalise the Pounce method in the support tree setting and then prove the desired result. Definition 8. Given a graph $G = (V, E)$, with $|V| = n$, and spine S , we define a binary support tree of G to be any binary tree $T = (VT, ET)$ of least possible depth, D , whose leaves are the vertices of S , in order. Note that $D \leq \log_2(n) + 1$. We show that there is a weighting of the support tree which ensures that the resistance diameter of the support tree is small, but also such that any labelling of the leaf vertices can be extended to the support tree such that its cut size remains small. This enables effective learning via the support tree. A related construction has been used to build preconditioners for solving linear systems [6]. Lemma 9. Given any spine graph $S = (V, E)$ with $|V| = n$, and labelling $u \in \{1, 1\}^n$, with $\epsilon \in [1, 1]$ — VT —support tree $T = (VT, ET)$, there exists a weighting A of T , and a labelling u' and u are identical on V , $\epsilon T(u) \leq \epsilon S(u)$ and $RT \leq (\log_2 n + 1)(\log_2 n + 4)(\log_2(\log_2 n + 2))^2$. Proof. Let v_r be the root vertex of T . Suppose each edge $(i, j) \in ET$ has a weight A_{ij} , which is a function of the edge's depth $d = \max\{dT(v_i, v_r), dT(v_j, v_r)\}$, $A_{ij} = W(d)$ where $dT(v, v_0) \leq$ such is the number of edges in the shortest path from v to v_0 . Consider the unique labelling u that, for $1 \leq i \leq n$ we have $u_i = u_i$ and such that for every other vertex $vp \in VT$, with child $u' \rightarrow u$ vertices vc_1, vc_2 , we have $u'_p = c_1 \cdot 2^{c_2}$, or $u'_p = u'_c$ in the case where vp has only one child, vc . Suppose the edges $(p, c_1), (p, c_2) \in ET$ are at some depth d in T , and let $V_0 \subseteq V$ correspond to the leaf vertices of T descended from vp . Define $\epsilon S(u|V_0)$ to be the cut of u restricted to vertices in V_0 . If $u'_c = u'_c$ then $(\epsilon u'_p - u'_c)^2 + (\epsilon u'_p - u'_c)^2 = 0 \leq 2\epsilon S(u|V_0)$, and if $u'_c = u'_c$ then $2^2(\epsilon u'_p - u'_c)^2 + (\epsilon u'_p - u'_c)^2 \leq 2\epsilon S(u|V_0)$. Hence $W(d)(\epsilon u'_p - u'_c)^2 + (\epsilon u'_p - u'_c)^2 \leq 2W(d)\epsilon S(u|V_0)$ (14) (a similar inequality is trivial in the case that vp has only one child). Since the sets of leaf descendants of all vertices at depth d form a partition of V , summing (14) first over all parent nodes at a given depth and then over all integers $d \in [1, D]$ gives $\epsilon \leq 4\epsilon T(u)$

$D \cdot X$

$W(d)\epsilon S(u)$.

$d=1$

(15) We then choose $1/(d+1)(\log_2(d+1))^2 R \leq 21 + \ln 2 \cdot 2 \cdot \ln 2 \cdot x$
 $dx =$

$W(d) =$ and note that

P?

$$1 \leq d \leq (d+1)(\log_2 (d+1))^2$$

$$(16) \quad 1 \leq 2$$

$$+ \ln 2 \leq 2.$$

PD Further, $RT = 2 \sum_{d=1}^{\infty} (d+1)(\log_2 (d+1))^2 \leq D(D+3)(\log_2 (D+1))^2$ and so $D \leq \log_2 n + 1$ gives the resistance bound. Definition 10. Given the task of predicting the labelling of an unweighted graph $G = (V, E)$ the Pounce algorithm proceeds as follows: An augmented graph $G' = (V', E')$ by attaching a binary support tree of G , with weights defined as in (16), to G ; formally let $T = (V_T, E_T)$ be such a binary support tree of G , then $G' = (V \cup V_T, E \cup E_T)$. The Pounce algorithm is then used to predict the (partial) labelling defined on G . Theorem 11. Given the task of predicting the labelling of any unweighted, connected, n -vertex graph $G = (V, E)$ in the online framework, the number of mistakes, M , incurred by the augmented Pounce algorithm satisfies $M \leq \min\{N(X, r_G) + 12 \sum_{u \in V} G(u) + 1, (17) \leq 0$

where $N(X, r_G)$ is the covering number of the input set $X = \{v_1, v_2, \dots, v_n\} \subseteq V$ relative to the resistance distance r_G of G and $u \in \mathbb{R}^n$ is any labelling consistent with the trial sequence. Furthermore, $M \leq 12 \sum_{u \in V} G(u)(\log_2 n + 1)(\log_2 n + 4)(\log_2 (\log_2 n + 2))^2 + 2$. (18)

Proof. Let u be some labelling consistent with the trial sequence. By (3) we have that $\sum_{u \in V} G(u) \leq 2 \sum_{u \in V} G(u)$ for any spine S of G . Moreover, by the arguments in Lemma 9 there exists some labelling \tilde{u} of the weighted support tree T of G , consistent with u on V , such that $\sum_{u \in V} T(u) \leq \sum_{u \in V} S(u)$. We then have $\sum_{u \in V} T(u) \leq \sum_{u \in V} G(u) + 3 \sum_{u \in V} G(u)$. (19)

By Rayleigh's monotonicity law the addition of the support tree does not increase the resistance between any vertices on G , hence $N(X, r_{G'}) \leq N(X, r_G)$. (20)

(19) yields (17) on G . Combining inequalities (19) and (20) with the pounce bound (13) for predicting $u \leq 1 \leq N(X, r_{G'}) + 12 \sum_{u \in V} G(u) + 1$. $M \leq N(X, r_{G'}) + 4 \sum_{u \in V} G(u) \leq G + 2 \leq$ which proves (17). We prove (18) by covering G with single ball so that $M \leq 4 \sum_{u \in V} G(u) R 12 \sum_{u \in V} G(u) RT + 2$ and the result follows from the bound on RT in Lemma 9.

7

Conclusion

We have explored a deficiency with existing online techniques for predicting the labelling of a graph. As a solution, we have presented an approximate cut-preserving embedding of any graph $G = (V, E)$ into a simple path graph, which we call a spine, such that an implementation of the 1nearest-neighbours algorithm is an efficient realisation of a Bayes optimal classifier. This therefore achieves a mistake bound which is logarithmic in the size of the vertex set for any graph, and the complexity of our algorithm is of $O(\sum_{u \in V} \deg(u) + \sum_{u \in V} \ln \deg(u))$. We further applied the insights gained to a second algorithm - an augmentation of the Pounce algorithm, which achieves a polylogarithmic

performance guarantee, but can further take advantage of clustered data, in which case its bound is relative to any cover of the graph.

2 References

- [1] J. M. Barzdin and R. V. Frievald. On the prediction of general recursive functions. *Soviet Math. Doklady*, 13:1224–1228, 1972.
- [2] M. Belkin and P. Niyogi. Semi-supervised learning on riemannian manifolds. *Machine Learning*, 56:209–239, 2004.
- [3] A. Blum and S. Chawla. Learning from labeled and unlabeled data using graph mincuts. In *Proc. 18th International Conf. on Machine Learning*, pages 19–26. Morgan Kaufmann, San Francisco, CA, 2001.
- [4] P. Doyle and J. Snell. Random walks and electric networks. *Mathematical Association of America*, 1984.
- [5] J. Fakcharoenphol and B. Kijssirikul. Low congestion online routing and an improved mistake bound for online prediction of graph labeling. *CoRR*, abs/0809.2075, 2008.
- [6] K. Gremban, G. Miller, and M. Zagha. Performance evaluation of a new parallel preconditioner. *Parallel Processing Symposium, International*, 0:65, 1995.
- [7] M. Herbster. Exploiting cluster-structure to predict the labeling of a graph. In *The 19th International Conference on Algorithmic Learning Theory*, pages 54–69, 2008.
- [8] M. Herbster and M. Pontil. Prediction on a graph with a perceptron. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information Processing Systems 19*, pages 577–584. MIT Press, Cambridge, MA, 2007.
- [9] M. Herbster, M. Pontil, and L. Wainer. Online learning over graphs. In *ICML ’05: Proceedings of the 22nd international conference on Machine learning*, pages 305–312, New York, NY, USA, 2005. ACM.
- [10] R. Kinderman and J. L. Snell. *Markov Random Fields and Their Applications*. Amer. Math. Soc., Providence, RI, 1980.
- [11] D. Klein and M. Randić. Resistance distance. *Journal of Mathematical Chemistry*, 12(1):81–95, 1993.
- [12] N. Littlestone. Learning when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1988.
- [13] K. Pelckmans and J. A. Suykens. An online algorithm for learning a labeling of a graph. In *Proceedings of the 6th International Workshop on Mining and Learning with Graphs*, 2008.
- [14] X. Zhu, Z. Ghahramani, and J. Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *20-th International Conference on Machine Learning (ICML-2003)*, pages 912–919, 2003.