

A Constant Factor Approximation Algorithm for a Class of Classification Problems

Anupam Gupta*

Éva Tardos†

ABSTRACT

In a traditional classification problem, we wish to assign labels from a set L to each of n objects so that the labeling is consistent with some observed data that includes pairwise relationships among the objects. Kleinberg and Tardos recently formulated a general classification problem of this type, the “metric labeling problem”, and gave an $O(\log |L| \log \log |L|)$ approximation algorithm for it. The algorithm is based on solving a linear programming relaxation of a natural integer program and then randomized rounding.

In this paper we consider an important case of the metric labeling problem, in which the metric is the truncated linear metric. This is a natural non-uniform and robust metric, and it arises in a number of applications. We give a combinatorial 4-approximation algorithm for this metric. Our algorithm is a natural local search method, where the local steps are based on minimum cut computations in an appropriately constructed graph. Our method extends previous work by Boykov, Veksler and Zabih on more restricted classes of metrics.

1. INTRODUCTION

In a traditional classification problem, we wish to assign one of l labels (or classes) to each of n objects, in a way that is consistent with some observed data that we have about the problem that includes pairwise relationships among the objects to

*Computer Science Division, UC Berkeley, Berkeley, CA 94720. Email: angup@cs.berkeley.edu. Supported by NSF grant CCR-9820951

†Department of Computer Science, Cornell University, Ithaca NY 14853. Email: eva@cs.cornell.edu. Research done in part while a visiting Miller Professor at the Computer Science Division, UC Berkeley, Berkeley, CA 94720. Partially supported by a Guggenheim Fellowship, NSF grant CCR-9700163 and ONR grant N00014-98-1-0589.

be classified [7; 12]. Kleinberg and Tardos [19] recently formulated a general classification problem which they called the *metric labeling problem*. They also gave an $O(\log l \log \log l)$ -approximation algorithm for the general metric labeling problem, where l is the number of labels. Their approximation algorithm is based on solving a linear programming relaxation of an integer program formulation and then rounding the solution randomly. In this paper, we give a fast combinatorial 4-approximation algorithm for the special case of the truncated linear metric.

Classification and the Metric Labeling Problem: The traditional classification problem can give rise to the following labeling problem. Consider a set P of n objects that we wish to classify, and a set L of l possible labels. A *labeling* of P over L is simply a function $f : P \rightarrow L$; we choose a label for each object. The quality of our labeling is based on the contribution of two sets of terms.

- For each object $p \in P$ and label $i \in L$, we have a non-negative *assignment cost* $c(p, i)$ associated with assigning the label i to the object p .
- We have a graph G over the vertex set P , with edge set E indicating the objects that are related; each edge $e = \{p, q\}$ will have a non-negative weight w_e , indicating the strength of the relation. In addition, we have a distance function $d(\cdot, \cdot)$ on the set L of labels. If we assign label i to object p and label j to object q , and $e = \{p, q\}$ is an edge of G , then we pay a *separation cost* $w_e d(i, j)$.

Thus, the *total cost* of a labeling f is given by

$$Q(f) = \sum_{p \in P} c(p, f(p)) + \sum_{e = \{p, q\} \in E} w_e d(f(p), f(q)).$$

The *labeling problem* asks for a discrete labeling of minimum total cost. Recall that a *distance* $d : L \times L \rightarrow \mathbb{R}^+$ is a function that is symmetric and $d(i, i) = 0$ for all $i \in L$. If d also satisfies the triangle inequality, then d is a *metric* [11, page 27]. The labeling problem is called *metric* if the distance function $d(\cdot, \cdot)$ is a metric on the label set L .

In this paper, we consider the important case when the labels are the positive integers and the metric distance between labels i and j is given by the *truncated linear metric* $d(i, j) =$

$\min\{M, |i - j|\}$, for which we give a fast flow-based 4-approximation algorithm. In the next subsection, we discuss why this is a natural metric to be singled out for study by considering two well-studied problems from vision: the problem of “restoring” an image that has been degraded by noise [3; 14], and the visual correspondence problem [20] that is the basis of determining depth and motion in an image.

The Truncated Linear Metric: Motivations and Applications: The labeling problem is closely related to the theory of *Markov Random Fields* (MRFs). Suppose that we have observed data $f'(p)$ for each object $p \in P$. We assume that f' was obtained from a labeling f by the introduction of independent random noise at each object. We wish to decide on the most probable “true” labeling f of P given this data. A Markov Random Field can be defined by a graph on the objects, with edges indicating dependencies between the objects. We need to assume that in the “true” labeling f the probability that an object has a certain label depends only on the labels of the neighboring objects. Now the labeling problem above can be restated as the problem of finding the f that maximizes the *a posteriori* probability $\Pr[f|f']$, if the underlying MRF model satisfies two standard assumptions of Homogeneity and Pairwise Interactions. (See [19] for more details on the connection to MRFs.) Next we give two applications from vision when the truncated linear metric is important. Both these problems can also be formulated using MRFs, but here we will give a direct description.

In both these applications, we are given a large grid of pixels as our graph. In the *image restoration problem* we want to label each pixel with its “true” intensity. We are given an “observed” intensity $f'(p)$ for each pixel, that is the result of corruption by noise. We would like to find the best way to label each pixel with a (true) intensity value, based on the observed intensities and the special structure of the image. Our determination of the “best” labeling is a labeling problem based on the trade-off between two competing influences: We would like to give each pixel an intensity close to what we have observed; furthermore, since real images are mainly smooth, with occasional boundary regions of sharp discontinuity, we would like spatially neighboring pixels to receive similar intensity values. The labels for this problem are the possible intensities, and we assign a penalty of $w_\epsilon d(i, j)$ if neighboring pixels are assigned labels i and j respectively.

Now consider the desired properties of a distance function $d(\cdot, \cdot)$. Small differences in intensity can reflect, for example, gradual changes in lighting, and hence should receive a small penalty. Large changes in intensity are likely to reflect object boundaries, and hence should receive a large penalty. In addition, we want the penalty to be *robust*: once the difference of the assigned labels is large enough to suggest an object boundary, the penalty should be constant. If the penalty function is not robust, there is a possibility of over-smoothing at the boundaries, i.e., object boundaries may become fuzzy. In summary, we want a metric on the labels with the following two properties.

- The metric should not be uniform, if two labels $i, j \in L$ are similar, than $d(i, j)$ should be small.
- The metric should be robust, once two labels $i, j \in L$

are sufficiently different to suggest object boundaries, the distance $d(i, j)$ should be constant.

In the image restoration problem for a grey-scale image, the possible intensities are integers, and the simplest and most natural metric function that satisfies our properties is the truncated linear metric, $d(i, j) = \min\{M, |i - j|\}$.

An application where non-uniform, robust metrics are even more important arises in the visual correspondence problem, which is the basis of determining motion and depth from the camera in an image. In the *visual correspondence problem*, we are given two images of the same scene, a pixel in one image corresponds to a pixel in the other if both pixels are projections along lines of sight of the same physical scene element. The problem is to determine this correspondence between pixels of two images. To apply the labeling problem to solve this problem we arbitrarily select one of the images to be the primary image. The quantity to be estimated for each pixel is the difference of its place on the image, and the place of the corresponding pixel on the other image (called the *disparity*). This disparity is in one-to-one correspondence with the depth of the corresponding point in the scene from the camera. Small differences in depth on neighboring pixels can reflect sloping objects (i.e., an object that starts close to the camera, but extends to some distance), while large differences in depth are likely to reflect object boundaries. As before, a non-uniform and robust metric seems natural to use. The possible depths values are integer multiples of some small shift, and the truncated linear metric is a extremely natural metric to use in this case.

Previous Results: Traditional local search methods (such as simulated annealing) with simple local steps (such as re-labeling one object) have previously been used extensively to solve such classification problems. However, the algorithms were known to perform poorly in practice, and no bounds were known for the performance.

Some special cases of the labeling problem are known to be solvable in polynomial time using max-flow computations. Besag [3] and Greig et al. [15] showed that the case of $l = 2$ labels is polynomially solvable as a two-terminal minimum cut problem. Another special case when the metric labeling problem reduces to two-terminal minimum cuts was discovered independently by Boykov et al. [4; 21] and Ishikawa and Geiger [16]. They considered metric labeling with the label set being $L = \{1, \dots, l\}$ and the distance between two labels i and j being $d(i, j) = |i - j|$. However, the linear metric is not robust and labeling with it can lead to over-smoothing. The work of Karzanov [18] focussed on some other polynomially solvable special cases of the labeling problem.

Kleinberg and Tardos [19] provided the first polynomial-time approximation algorithm for the metric labeling problem. They gave an $O(\log l \log \log l)$ -approximation algorithm with respect to an arbitrary graph on the set of objects, and an arbitrary metric on the set of labels. For the special case of the *uniform metric*, in which all distances are equal to 1, they gave a 2-approximation algorithm. For the general case, they use a result of Bartal [1; 2] which says that every metric can be approximated by hierarchically well-separated tree metric. However, the best approximation of the truncated linear metric with

a random hierarchically well-separated tree metric has distortion $O(\log M)$, and hence for the case of the truncated linear metric the analysis gives only an $O(\log M)$ -approximation. Furthermore, the approximation algorithm is based on solving a linear program (that corresponds to a fractional classification problem), and rounding its solution. The linear program involved is quite large, and this causes the method to be too slow to be practical for images.

The metric labeling problem is an extension of the well-studied *multiway cut problem* [8; 9; 10; 13; 17], in which one must find a partition of a graph into l sets in a way that is consistent with an initial assignment of certain vertices to these sets, and which cuts as few edges as possible. Boykov et al. developed the connection between *uniform* labelings, with $l > 2$ labels, and multiway cuts in a graph, showing a direct reduction from labelings to multiway cuts [4]. However, their reduction requires the use of edges of enormously large weight, and so it is not approximation-preserving. In [4; 6; 5; 21], they also developed some flow-based local search heuristics for labeling problems using the uniform metric. Their methods can be thought of as clever extensions of the multiway cut *Isolation Heuristic* of Dahlhaus et al. [10]. For one of their heuristics [6], they prove that any local optimum is a 2-approximation. In fact, it is possible to adapt their proof to show that the local improvement heuristic provides a $2 + \epsilon$ -approximation in polynomial time. In [4; 21], they show that the heuristics also perform very well in practice. They also consider other metrics, and in [5; 21], they show that a similar local improvement step can be described for any metric, but in general the locally optimal solution can be arbitrarily far from optimum. In particular, for the truncated linear metric, this method provides a $2M$ -approximation algorithm [5; 21], where M is the upper-bound used in defining the metric.

Our Results: In this paper we consider the metric labeling problem with the *truncated linear metric*, which is the simplest and arguably the most natural non-uniform and robust metric. We assume that the metric distance between labels i and j is given by the truncated linear norm $d(i, j) = \min\{M, |i - j|\}$.

Our first observation is that the truncated linear metric can be approximated by a random tree metric with expected stretch at most 3 (See figure 2.1). The trees used by this approximation are not hierarchically well-separated, and hence cannot be used directly in the algorithm of [19]. However, it is possible to combine the ideas of the linear programming and rounding with the flow algorithm for the (untruncated) linear metric, to obtain a constant factor approximation algorithm. Unfortunately, this constant is quite large, and furthermore, solving the large LP still remains practically infeasible. We will not describe the details of this linear programming and rounding method here as our main result is a simpler and better approximation algorithm.

Our main result is a flow based local improvement algorithm that gives a factor 4-approximation to the optimum. A key insight this algorithm comes from the above tree embedding. Any tree in the distribution consists of roughly l/M equal length edges leaving the root, and there is a linear segment of length at most M attached to each of them. (See figure 2.1.) Thus the tree metrics in the embedding can be thought of as a combination of the uniform metric and a set of (untruncated)

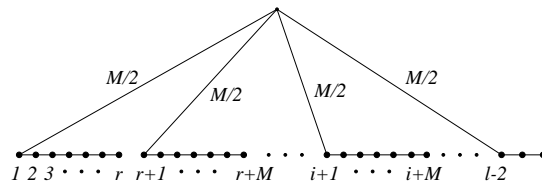


Figure 2.1: The embedding for the truncated linear metric

linear metrics. Now recall that the metric labeling problem on the (untruncated) linear metric can be solved optimally using flows and cuts, and further that the uniform metric can be approximated well using a flow based local search. Our algorithm is a local search method for the truncated linear metric that is obtained by a clever combination of these two techniques which gives us the claimed approximation guarantee.

2. TREE METRICS AND RANDOM PARTITIONING

Consider a path of length l , where the vertices are the labels $L = \{1, 2, \dots, l\}$ and edges $\{i, i + 1\}$ are of unit length. To this, we add edges of length M between all pairs of labels $\{i, j\}$ with $|i - j| > M$. The geodesic metric generated by this graph is exactly the truncated linear metric $\mathcal{L} = (L, d)$.

To embed this into a random tree, consider picking a random offset $r \in \{1, 2, \dots, M\}$, and for each position $1 \leq k \leq l$ which is equal to $r \pmod{M}$, deleting all edges that connect a pair of labels $i \leq k$ and $j > k$. This results in a set of *intervals*, each of length at most M . Let us call this set of intervals \mathcal{S}_r , the *partition* associated with r . Now we add the lowest numbered vertex of each interval in \mathcal{S}_r to a new *root* vertex with an edge of length $M/2$ to get a tree. (See figure 2.1 for an example.)

Theorem 2.1. *For any pair of labels $i, j \in L$, the distance between them in any of the (random) trees is at least $d(i, j)$, and the expected distance between them is at most $3d(i, j)$.*

We shall not use this tree embedding directly; however, it will influence the structure of both the algorithm and the analysis we develop in the following section. We will need the following easy fact about the probability that two labels lie in different intervals of a partition \mathcal{S}_r .

Lemma 2.2. *For any pair of labels $i, j \in L$, the probability that i and j lie in different intervals of the partition \mathcal{S}_r is exactly $d(i, j)/M$.*

3. THE LOCAL SEARCH ALGORITHM

The local search algorithm starts with an arbitrary labeling f of the vertices of the graph G . It then makes a simple “local” move whose goal is to decrease the cost of the labeling. If making the move causes it to reach a new labeling whose cost is lower than the cost of the current labeling, it moves to the new labeling and continues.

Boykov et al. [6] use the following local improvement step. Consider a label i . In a single “local” improvement step they allow relabeling any subset of the vertices with the label i , and they call this local move an i -expansion move. They prove that if the distance function on the labels is a metric, then the best i -expansion move can be found via a single min-cut computation in an appropriately constructed graph [6; 5]. Further, they prove that in the case of the uniform metric, any local optimum is a 2-approximation to the global optimum [6; 21]. For metrics other than the uniform metric, however, a locally optimal solution can be arbitrarily far from optimum. In particular, for the truncated linear metric, their methods give a $2M$ -approximation algorithm [5; 21], where M is the parameter used in defining the truncated metric.

In this section we develop an analogous, but much more general local improvement method. In a single local step we consider an interval I of labels of length at most M , and allow any subset of vertices to be relabeled by any of the labels in I . Given a labeling f , we call another labeling f' a *local relabeling* if it can be obtained from f by such a local move, i.e., if for all objects $f(p) \neq f'(p)$ implies that $f'(p) \in I$. Unfortunately, we will not be able to find the best possible such local move. However, we will later show that if the current labeling has cost sufficiently far above the minimum possible cost, then our method will find a move that significantly decreases the cost of the labeling.

In the algorithm, we repeatedly pick a random interval I (of length at most M) and try to relabel some subset of objects to the labels in I to decrease the cost of our labeling. After this local step, each object will either have its label unchanged, or will have a label in the interval I . To perform this relabeling, we will create a flow network and find a minimum $s - t$ cut in this network. This minimum cut can be associated with a new labeling f' , and if this new labeling f' has a lower cost than the cost of the original labeling f , we move to the new labeling. In summary, the algorithm is the following:

Algorithm Local-Search:
repeat
 pick a random interval I
 build the flow network N_I associated with I
 if labeling given by the minimum cut
 on N_I has lower cost
 then move to new labeling.
until some stopping rule.

More formally, we will pick the random intervals I in the following manner: we pick a random integer $-M < r < l$, and set I to be the part of the interval of length M starting from offset r that lies in L . I.e., $I = \{r + 1, r + 2, \dots, r + M\} \cap \{1, 2, \dots, l\}$. We had claimed that this algorithm was motivated by the tree-metrics of Section 2. To see this, note that the probability of any fixed interval J being picked is almost the same as by the following alternative procedure: first choose a random tree from the tree embedding of the previous section, and then randomly pick one of the intervals in it.

In the next subsection, we will describe how to build the flow network associated an interval I to which the labels can be changed. We then analyze this algorithm in subsection 3.2, showing, for instance, that running the algorithm until we reach

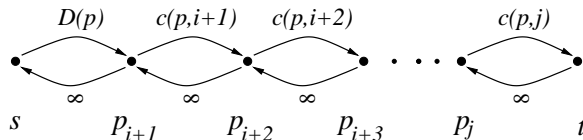


Figure 3.2: The chain for vertex p

a local optimum (i.e., a state where none of the intervals leads to an improvement) implies that we are close to the optimal value.

3.1 The flow network

Let us consider an interval I of labels. Assume that the labels in the interval I are $\{i + 1, i + 2, \dots, j\}$, where $(j - i)$, the number of labels in I , is at most M . (In fact, any interval we consider will have length exactly M except when it is some initial or final portion of the line.) In this move, each vertex in G will either retain its label or move to a label in I . We now describe a flow network $N_I = (V, A)$ associated with I : it will be a directed graph with a source s and a sink t , and with capacities on the edges. We describe the construction in two steps.

For each vertex p of the original graph G , we add $(j - i)$ nodes $\{p_{i+1}, \dots, p_j\}$ to N_I . (Figure 3.2 shows this construction.) We add directed edges (p_k, p_{k+1}) and (p_{k+1}, p_k) for all $k = i + 1, \dots, j - 1$. We then attach p_{i+1} to s and p_j to t by directed edges in both directions. For the following discussion, it will be convenient to think of the vertex s as being p_i , and also of t as being p_{j+1} .

Now we shall assign the capacities. For $i + 1 \leq k \leq j$, each edge (p_k, p_{k+1}) is assigned a capacity equal to the assignment cost $c(p, k)$, while the edge (p_{k+1}, p_k) is assigned an infinite capacity. Finally, the edge (p_{i+1}, s) is assigned an infinite capacity, while (s, p_{i+1}) is assigned a capacity $D(p)$ which is defined as follows: if $f(p) \in I$, then we set $D(p) = \infty$, else $D(p) = c(p, f(p))$.

To see the rationale for this construction, consider any minimum $s - t$ cut in N_I : the infinite capacity edges ensure that this cut will include exactly one edge from the chain corresponding to each vertex p . We interpret the edge (p_k, p_{k+1}) being cut to mean that the vertex p is assigned the label k (if $k \neq i$); if the edge (s, p_{i+1}) is cut and $f(p) \notin I$, then we retain the original label $f(p)$ for the vertex p . It is clear to see that the assignment cost is exactly the capacity of the edge in the cut, and thus the chain for vertex p captures the assignment costs.

We also need to model the separation costs; for this we will add edges between the chains. Consider an edge $e = \{p, q\}$ in the original graph. There will be different cases depending on the original labels of p and q . One of the cases, when neither p nor q have their original label in I is shown on Figure 3.3. For each of the corresponding nodes p_k and q_k , where $(i + 1) < k \leq j$, we add a pair of oppositely directed edges between them, each with capacity w_e . If both p and q are labeled with vertices in I , then we do nothing more. If $f(p) \notin I$ but $f(q) \in I$, then we add an edge (p_{i+1}, q_{i+1}) with capacity

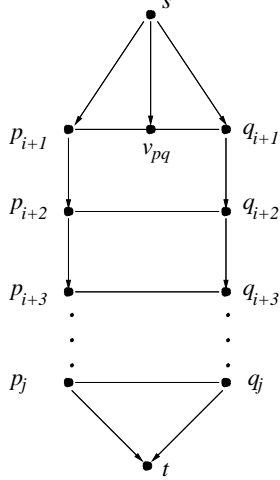


Figure 3.3: Construction for the edge $\{p, q\}$

$w_e d(f(p), i+1)$. Finally, if both the labels do not lie in I , then we add a new vertex v_{pq} . We then add two oppositely directed edges between p_{i+1} and v_{pq} having capacities $w_e d(f(p), i+1)$, and opposite edges between v_{pq} and q_{i+1} having capacities $w_e d(f(q), i+1)$. To complete the construction, we attach an edge (s, v_{pq}) with capacity $w_e d(f(p), f(q))$. (A view of this third case is given in Figure 3.3, where we have dropped some of the infinite capacity edges, and used undirected edges to indicate a pair of oppositely directed edges with identical costs.)

To see that this captures the separation costs, let f' be the labeling corresponding to a cut, and let us focus on the edge $e = \{p, q\}$. If both the vertices retain their original labels, then the cut will be minimized when it passes through (s, v_{pq}) , and incurs a cost of $w_e d(f(p), f(q))$, which is exactly the separation cost in this case. In case both the vertices are labeled with labels in I , then the cut value will be exactly $w_e |f'(p) - f'(q)|$, which is again the right separation cost.

However, if the cut indicates that one of the vertices (say p) retains its label $f(p) \notin I$, and q is labeled with some element $f'(q) = k \in I$, then the cut will incur a separation cost of $w_e [d(f(p), i+1) + (k - (i+1))]$. Note that this value possibly overestimates the actual separation cost in the new labeling.

Before we formalize the relation between cuts and labeling in N_I , let us consider two simpler special cases. The first case is when the interval I consists of all the labels. In this case, the construction is the same as the one discovered independently by Boykov et al. [4] and Ishikawa and Geiger [16] for finding the optimal labeling for the line metric.

Theorem 3.1 (Boykov et al. [4], Ishikawa and Geiger [16]). *The minimum cut of the flow network associated with an interval I consisting of all the labels gives the optimal labeling for the line metric on I .*

The other special case is when the interval consists of a single label $I = \{i+1\}$. In this case, the construction above is the same as the one by Boykov et al. [6] for finding the best

$(i+1)$ -expansion move, i.e., relabeling where labels can only change to the label $(i+1)$.

Theorem 3.2 (Boykov et al. [6; 5]). *The minimum cut of the flow network N_I associated with an interval $I = \{i+1\}$ consisting of a single label is the optimal $(i+1)$ -expansion move.*

We now summarize the important properties of the labeling obtained by the minimum cut in the flow network for the general case, when interval I has more than a single node, but does not contain all labels in L .

A minimum cut in N_I can contain at most one of the edges between the pairs of nodes p_{i+1} and v_{pq} , or between q_{i+1} and v_{pq} , or the edge (s, v_{pq}) for any edge $e = \{p, q\}$ in the graph G . To see this, note that only one edge in any pair of opposite edges can be in the cut. Further, in any set of up to three permissible edges, the cost of any two edges is more than the cost of the third one, since $d(\cdot, \cdot)$ satisfies the triangle inequality. We call a cut *simple* if it has finite capacity, and it does not cut more than one of the above five edges associated with any edge $e \in E$.

Theorem 3.3. *The simple cuts in the flow network N_I are in one-to-one correspondence with local relabelings f' . The cost of the relabeling $Q(f')$ is no more than the cost of the associated cut, and the cost of the cut overestimates the cost of the labeling by replacing the separation cost $w_e d(f'(p), f'(q))$ for edges $e = \{p, q\}$ where exactly one end receives a label in I by a possibly larger term $w_e [d(f'(p), i+1) + d(i+1, f'(q))]$. Further, we have that*

$$\begin{aligned} d(f'(p), f'(q)) &\leq d(f'(p), i+1) + d(i+1, f'(q)) \\ &\leq 2M. \end{aligned} \quad (3.1)$$

3.2 The analysis

In this subsection, we show that when the algorithm reaches a local minimum, then the cost of the labeling is no more than 4 times the minimal cost of any labeling. Further, we also show that the algorithm converges to a such solution in polynomial time.

Let us fix an optimal labeling f^* , and let the algorithm's current labeling be f . For any subset $X \subseteq P$, let $A^*(X)$ and $A(X)$ be the assignment cost that the optimum and the current labeling pay respectively for the vertices in X , and for a set of edges $Y \subseteq E$, let $S^*(Y)$ and $S(Y)$ be the separation cost for those edges paid by the optimum and the current solution respectively. Clearly, $Q(f) = A(P) + S(E)$ and the optimum $Q(f^*) = A^*(P) + S^*(E)$.

Consider the case when the algorithm chooses an interval I . Let P_I be the set of vertices of G to which f^* assigns labels from the interval I . Let E_I be the set of edges in $E(G)$ such that the f^* -labels of both endpoints lie in I . Let $\partial_I^- \subset E_I$ be the set of edges such that exactly one end of the edge has f^* -label in I , the end with higher f^* -label, and ∂_I^+ be the set of edges that only the end with lower f^* -label has an f^* -label in I .

In the proof we will consider a random partition \mathcal{S}_r of the labels, as defined in Section 2. Clearly, $P = \cup_{I \in \mathcal{S}_r} P_I$ and $\cup_{I \in \mathcal{S}_r} \partial_I^- = \cup_{I \in \mathcal{S}_r} \partial_I^+$. We will use ∂_r to denote this union of the *boundary edges* in the partition. Also note that $E = \partial_r \cup (\cup_{I \in \mathcal{S}_r} E_I)$. The edges in ∂_r will play a special role in the analysis. The following lemma will be useful later.

Lemma 3.4. *For a random partition \mathcal{S}_r , the expected value of $M \sum_{e \in \partial_r} w_e$ is $S^*(E)$.*

Proof. The probability that edge $e = \{p, q\}$ will be part of ∂_r is exactly $d(f^*(p), f^*(q))/M$, hence the edge e contributes $w_e d(f^*(p), f^*(q))$ to the expectation of the sum. This is exactly the same as its contribution to $S^*(E)$. \square

Consider the situation when the algorithm picks the interval I . One possible local move is changing the labels of vertices in the region P_I to their labels in the optimum labeling f^* . We will use this move to bound the improvement obtained in our algorithm.

Lemma 3.5. *For a labeling f , and an interval I , the local relabeling move that corresponds to the minimum cut in N_I decreases the cost of the solution by at least*

$$\left(A(P_I) + S(E_I \cup \partial_I^- \cup \partial_I^+) \right) - \left((A^*(P_I) + S^*(E_I \cup \partial_I^-)) + M \sum_{e \in \partial_I^-} w_e + 2M \sum_{e \in \partial_I^+} w_e \right). \quad (3.2)$$

Proof. The algorithm selects the minimum cut in N_I . The labeling corresponding to this cut has cost at most the capacity of the cut in N_I . To prove the bound we have to show a cut with small capacity. Consider the cut that corresponds to the local move of relabeling all nodes in P_I to their f^* -label. The capacity corresponding to the assignment cost for nodes in P_I will be exactly $A^*(P_I)$, and the capacity corresponding to separation cost for the edges in E_I will be $S^*(E_I)$. For nodes and edges that do not change their label in this move, the capacity corresponding to their assignment and separation cost will be unchanged.

Consider the edges in $\partial_I^- \cup \partial_I^+$. These are the edges that possibly change the label of exactly one of their ends. We estimate the capacity corresponding to the separation cost of such an edge $e = \{p, q\}$, assuming p is the end in P_I by

$$w_e [d(f^*(p), i+1) + d(i+1, f(q))]$$

using Theorem 3.3. If $e \in \partial_I^+$, we use (3.1) to get a bound of $2M w_e$. If $e \in \partial_I^-$ we have that $f^*(q) \leq i$, and so $d(f^*(p), i+1) \leq d(f^*(p), f^*(q))$. Bounding $d(i+1, f(q))$ by M , we get that the capacity corresponding to the separation of edge e is at most $w_e [d(f^*(p), f^*(q)) + M]$. Summing this over all edges in $\partial_I^- \cup \partial_I^+$, we get exactly the claimed bound. \square

Let us first assume that we have reached a local optimum, i.e., a local move with any interval does not decrease the cost of the labeling.

Theorem 3.6. *If the labeling f is a local optimum, its cost $Q(f)$ is at most 4 times the optimal cost $Q(f^*)$.*

Proof. The fact that f is a local optimum implies the improvement indicated by the expression (3.2) is non-positive for any interval I , i.e., for all I

$$\begin{aligned} A(P_I) + S(E_I \cup \partial_I^- \cup \partial_I^+) &\leq A^*(P_I) + S^*(E_I \cup \partial_I^-) \\ &\quad + M \sum_{e \in \partial_I^-} w_e + 2M \sum_{e \in \partial_I^+} w_e. \end{aligned} \quad (3.3)$$

Now consider a partition \mathcal{S}_r and sum these inequalities for each interval $I \in \mathcal{S}_r$. On the left hand side we get $A(P) + S(E) + S(\partial_r)$ as edges in ∂_r occur in the boundary of two intervals. This is at least $Q(f)$, the cost of the labeling f . Summing the right hand side, we get exactly $A^*(P) + S^*(E) + 3M \sum_{e \in \partial_r} w_e$. So we have that

$$Q(f) \leq A^*(P) + S^*(E) + 3M \sum_{e \in \partial_r} w_e$$

for any partition \mathcal{S}_r . Now we can take expectations. The left side is a constant, and by lemma 3.4, the expected value of the right hand side is at most $A^*(P) + 4S^*(E)$. Thus we get $Q(f) \leq A^*(P) + 4S^*(E) \leq 4Q(f^*)$, as claimed. \square

In fact, we can also bound the number of steps we have to run the algorithm to reach a near-optimal solution. Let Q_0 be the cost of the initial labeling.

Theorem 3.7. *If the main loop of the algorithm is repeated $O((l/M)(\log Q_0 + \log \epsilon^{-1}))$ times, then the expected cost of the resulting labeling is at most $(4 + \epsilon) Q(f^*)$.*

Proof. To estimate the expected decrease of the labeling cost in one iteration, consider the following alternate process for selecting a random interval. First select a partition \mathcal{S}_r at random, and then select an interval $I \in \mathcal{S}_r$. This process selects each interval with roughly the same probability as the one our algorithm used. The difference is that some partitions have one more interval than others, and intervals in such partitions have a slightly smaller chance of being selected.

Consider a labeling f and a partition \mathcal{S}_r . We use Lemma 3.5 to estimate the improvement obtained if interval I is used. As in the previous proof, we will sum the estimates over the set of intervals in the partition. There are at most $\lceil l/M \rceil + 1$ intervals in \mathcal{S}_r , hence the expected decrease in the labeling cost when a random interval is selected from \mathcal{S}_r is at least

$$\frac{1}{\lceil l/M \rceil + 1} \left(Q(f) - Q(f^*) - 3M \sum_{e \in \partial_r} w_e \right).$$

But, by Lemma 3.4, for a randomly chosen partition \mathcal{S}_r , the expected value of this decrease is at least

$$\Omega \left(\frac{M}{l} (Q(f) - 4Q(f^*)) \right),$$

so in $O(l/M)$ iterations the difference $Q(f) - 4Q(f^*)$ is expected to decrease by a constant factor. This implies the claimed bound on the running time. \square

It is easy to see that the algorithm can be trivially derandomized at a cost of a factor $(l + M)$ increase in the running time. This can be done by considering all possible $(l + M)$ intervals and making the moves corresponding to the best possible interval, or making local moves for each of the $(l + M)$ intervals in turn in a round robin fashion.

We note that in the special case of the uniform metric, our algorithm becomes the 2-approximation algorithm of Boykov et al. [6]. Since the uniform metric can be thought of as a truncated linear metric with $M = 1$, the flow construction finds the optimal $(i + 1)$ -expansion move (see Theorem 3.2). Using this instead of Theorem 3.3 in the proof of Theorem 3.6, we get the result of [6] that any local optimal solution is a 2-approximation. Further, we also get an analog of Theorem 3.7, saying that the expected cost of the labeling after $O(\log Q_0 + \log \epsilon^{-1})$ repetitions of the main loop of the algorithm have been performed is at most $(2 + \epsilon) Q(f^*)$.

Acknowledgements

We would like to thank Yuri Boykov, Jon Kleinberg, Olga Veksler and Ramin Zabih for numerous discussions on the topic of this paper.

4. REFERENCES

- [1] Yair Bartal. Probabilistic approximations of metric spaces and its algorithmic applications. In *Proceedings of the 37th Annual IEEE Symposium on Foundations of Computer Science*, pages 184–193, 1996.
- [2] Yair Bartal. On approximating arbitrary metrics by tree metrics. In *Proceedings of the 30th Annual ACM Symposium on Theory of Computing*, pages 161–168, 1998.
- [3] Julian Besag. On the statistical analysis of dirty pictures. *J. Roy. Statist. Soc. Ser. B*, 48(3):259–302, 1986.
- [4] Yuri Boykov, Olga Veksler, and Ramin Zabih. Markov Random Fields with efficient approximations. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 648–655, 1998.
- [5] Yuri Boykov, Olga Veksler, and Ramin Zabih. Fast approximate energy minimization via graph cuts. In *Proceedings of the 7th IEEE International Conference on Computer Vision*, volume 1, pages 377–384, 1999.
- [6] Yuri Boykov, Olga Veksler, and Ramin Zabih. A new algorithm for energy minimization with discontinuities. In *International Workshop on Energy Minimization Methods in Computer Vision and Pattern Recognition*, 1999.
- [7] Leo Breiman, Jerome H. Friedman, Richard A. Olshen, and Charles J. Stone. *Classification and regression trees*. Wadsworth Advanced Books and Software, Belmont, CA, 1984.
- [8] Gruiă Călinescu, Howard Karloff, and Yuval Rabani. Improved approximation algorithms for multiway cut. In *Proceedings of the 30th Annual ACM Symposium on Theory of Computing*, pages 48–52, 1998.
- [9] William H. Cunningham and Lawrence Tang. Optimal 3-terminal cuts and linear programming. In *Proceedings of the MPS Conference on Integer Programming and Combinatorial Optimization*, 1999.
- [10] Elias Dahlhaus, David S. Johnson, Christos H. Papadimitriou, Paul D. Seymour, and Mihalis Yannakakis. The complexity of multiterminal cuts. *SIAM J. Comput.*, 23(4):864–894, 1994.
- [11] Michel Marie Deza and Monique Laurent. *Geometry of Cuts and Metrics*. Springer Verlag, 1997.
- [12] Thomas G. Dietterich. Machine learning research: Four current directions. *AI Magazine*, 18(4):97–136, 1997.
- [13] Péter L. Erdős and László A. Székely. Evolutionary trees: an integer multicommodity max-flow–min-cut theorem. *Adv. in Appl. Math.*, 13(4):375–389, 1992.
- [14] Stuart Geman and Donald Geman. Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:721–741, 1984.
- [15] D. Greig, B.T. Porteous, and A. Seheult. Exact maximum a posteriori estimation for binary images. *J. Royal Statistical Society B*, 51(2):271–279, 1989.
- [16] Hiroshi Ishikawa and Davi Geiger. Segmentation by grouping junctions. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 125–131, 1998.
- [17] David R. Karger, Philip Klein, Cliff Stein, Mikkel Thorup, and Neal E. Young. Rounding algorithms for a geometric embedding of multiway cut. In *Proceedings of the 31th Annual ACM Symposium on Theory of Computing*, pages 668–678, 1999.
- [18] Alexander V. Karzanov. Minimum 0-extensions of graph metrics. *European J. Combin.*, 19(1):71–101, 1998.
- [19] Jon Kleinberg and Éva Tardos. Approximation algorithms for classification problems with pairwise relationships: Metric labeling and Markov Random Fields. In *Proceedings of the 40th Annual IEEE Symposium on Foundations of Computer Science*, 1999. To appear.
- [20] Sébastien Roy and Ingemar J. Cox. A maximum-flow formulation of the N -camera stereo correspondence problem. In *Proceedings of the 6th IEEE International Conference on Computer Vision*, pages 492–499, 1998.
- [21] Olga Veksler. *Efficient graph-based energy minimization methods in computer vision*. PhD thesis, Department of Computer Science, Cornell University, 1999.