

# Latent Labeled Graph Classification

Joshua T. Vogelstein, John C. Conroy, Lou Podrazik, Steve Kratzer, Glen A. Coppersmith, Mark Dredze, R. Jacob Vogelstein, and Carey E. Priebe

**Abstract**—Because graphs can encode more information in their structure than vectors, they are becoming increasingly popular data structures for representing information. While the last century has witnessed the development of a plethora of statistical tools for the analysis of data, the vast majority of these tools natively operate in vector spaces, not graph spaces. Thus, algorithms for even relatively simple statistical inference tasks, such as two-class classification, are essentially absent for graph data. In this work, we propose a number of complementary algorithms to classify graphs, with special attention to the possibility of unknown vertex labels. Since exactly solving the graph-matching problem is currently computational intractable, we consider several approximate approaches. We introduce a multiple-restart Frank-Wolfe approach to solving the graph matching problem by formulating it as a quadratic assignment problem. Although this approach has superior performance than previous state-of-the-art approaches to the graph matching problem, even when it “should” do well in classification problems, it is outperformed by a graph invariant strategy. This is just the beginning.

**Index Terms**—statistical inference, graph theory, network theory, structural pattern recognition, connectome.



## 1 INTRODUCTION

THE statistical analysis of collections of graphs is becoming an increasingly popular desideratum [1]. Specifically, we consider the following idealized scenario. Let  $\mathbb{G} : \Omega \mapsto \mathcal{G}$  be a graph-valued random variable taking values  $G \in \mathcal{G}$ . Let  $Y$  be a categorical random variable,  $Y : \Omega \mapsto \mathcal{Y} \subseteq \mathbb{Z}$ , such that each graph has an associated class. Given a collection of graphs and classes, we assume they were jointly sampled independently and identically from some true but unknown joint distribution,  $\{(\mathbb{G}_i, Y_i)\}_{i \in [n]} \stackrel{iid}{\sim} F_{\mathbb{G}, Y}(\cdot; \theta)$ , where  $\theta$  is a set of parameters. Note that  $F_{\mathbb{G}, Y}(\cdot; \theta)$  is but one of a (possibly infinite) set of distributions, collectively comprising the model:  $\mathcal{F}_{\mathbb{G}, Y} = \{F_{\mathbb{G}, Y}(\cdot; \theta) : \theta \in \Theta\}$ , where  $\Theta$  is the set of feasible parameters. The goal of such an analysis is to learn about the relationship between the random variables  $\mathbb{G}$  and  $Y$ . Most standard classification techniques fail in this domain as they typically require classifying objects that live in finite dimensional Euclidean space,  $\mathbb{R}^d$ , whereas the object of interest here live in graph space,  $\mathcal{G}$  (even finite graphs do not natively live in Euclidean space). In this work, therefore, we propose novel extensions of several classification algorithms appropriate for the graph domain.

## 2 GRAPH CLASSIFICATION

The graph classification problem may be stated thusly: given training data  $\mathcal{T}_n = \{(\mathbb{G}_i, Y_i)\}_{i \in [n]}$ , and a new graph,  $\mathbb{G}$ , estimate the new graph’s corresponding class,

$Y$ , assuming each graph/class pair was sampled identically and independently from some true but unknown distribution,  $\mathcal{T}_n, (\mathbb{G}, Y) \stackrel{iid}{\sim} F_{\mathbb{G}, Y}(\cdot; \theta)$ . Given an appropriately defined risk-function, such as misclassification rate,  $R_h = \mathbb{P}[h(\mathbb{G}) \neq Y]$ , one can then search for the function  $h^* \in \mathcal{H}$  that minimizes the risk function of interest:

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathbb{P}[h(\mathbb{G}) \neq Y]. \quad (1)$$

In general,  $h^*$  is unavailable and dependent on the true but unknown distribution,  $F = F_{\mathbb{G}, Y}$  (which includes the vertex labels). When  $h^*$  is unavailable, one can utilize training data,  $\mathcal{T}_n$ , to obtain  $\hat{h}$ , an approximation to  $h^*$ :

$$\hat{h} \approx \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathbb{P}[h(\mathbb{G}) \neq Y | \mathcal{T}_n], \quad (2)$$

where  $\approx$  indicates that in general, we will not be able to find the actual minimum in the set  $\mathcal{H}$ . Regardless, any approach necessarily estimates a decision boundary in the space of graphs separating them into two classes.

## 3 MODELS

Let a graph be a 4-tuple consisting of a set of (i) vertices (ii) edges, (iii) vertex labels, and (iv) edge attributes. Assume that each graph has the same set of vertices. Moreover, assume that edge attributes can be represented as scalar quantities, so the edges and edge attributes of a random graph are jointly encoded in an adjacency matrix representation. Thus, the joint distribution on graphs and classes can be expanded and factorized:

$$F_{\mathbb{G}, Y} = F_{\mathbb{L}, \mathbb{A}, Y} = F_{\mathbb{L}, \mathbb{A} | Y} F_Y \quad (3)$$

where  $\mathbb{L} : \Omega \mapsto \mathcal{L} \subseteq \{L_1, L_2, \dots, L_{n_v}\} = [L_{n_v}]$  is a random vertex labeling,  $\mathbb{A} : \Omega \mapsto \mathcal{A} \subseteq \mathbb{R}^{n_v \times n_v}$  is a random adjacency matrix, and  $n_v$  is the number of vertices in the graphs. The class-conditional signal therefore derives from two possible sources of variability,  $\mathbb{L}$  and  $\mathbb{A}$ . *Labeled*

- J.T. Vogelstein and C.E. Priebe are with the Department of Applied Mathematics and Statistics, Johns Hopkins University, Baltimore, MD 21218.  
E-mail: joshuaov@jhu.edu
- R.J. Vogelstein is with the Johns Hopkins University Applied Physics Laboratory, Laurel, MD, 20723.

graph classifiers therefore proceed accordingly, using the labels to facilitate classification (see [?]). *Unlabeled* graph classifiers, however, assume that the labels carry no class-conditional signal, that is, they assume that  $F_{\mathbb{L},\mathbb{A}|Y} = F_{\mathbb{A}|Y}$ . In *Latent labeled* graph classifiers, those under investigation here, the labels are not observed, yet they are not assumed to be wholly uninformative with respect to the task at hand.

Let  $h_{\mathbb{A}}^* : \mathcal{A} \mapsto \mathcal{Y}$  be the *edge attribute only* Bayes optimal classifier, and  $R_{\mathbb{A}}^*$  be the risk associated with this classifier. Similarly, let  $h_{\mathbb{L},\mathbb{A}}^* : \mathcal{A} \times \mathcal{L} \mapsto \mathcal{Y}$  be the edge attribute and vertex label Bayes optimal classifier, and let  $R_{\mathbb{L},\mathbb{A}}^*$  be its risk. We are interested in parsing model space into those models for which  $R_{\mathbb{L},\mathbb{A}}^* < R_{\mathbb{A}}^*$ . A vertices are *stochastically equivalent* whenever

Given such a scenario, one could consider at least three complementary approaches: (i) a graph *labeling* approach in which one effectively tries to estimate labels, or (ii) a graph *invariant* approach in which one does not, and (iii) a graph *dissimilarity* approach in which one could. The preferred approach will depend on the assumed model, the data, and the computational resources. Section 4 provides details for these strategies, and Section 5 provides some guidance with regard to when to use each.

## 4 LATENT LABELED GRAPH CLASSIFIERS

### 4.1 Graph Labeling Classifiers

We define a graph labeling function as any algorithm that assigns a label to each vertex of a graph:  $Q_n : \mathcal{A}, \Xi^n \mapsto \mathcal{L}$ , where  $\Xi \subseteq \mathcal{G}$ , for example,  $\Xi = \mathcal{A}$  (note that we have actually defined a sequence of graph labeling functions). Remember that we have defined  $\mathcal{L}$  as a subset of  $[n_v]$ , so each vertex need not have a unique label.

To use a graph labeling classifier, by definition, one first explicitly attempts to estimate the labels of each vertex in each graph. Given these label estimates, one can plug them in and apply a labeled graph classifier,  $h_{\mathbb{L},\mathbb{A}}^*$ , such as those proposed in [?] (although perhaps keeping the label uncertainty could help). We therefore proceed by providing a novel approach to solving the graph labeling problem, followed by some options of how one might use them in a labeled graph classifier.

#### 4.1.1 Graph Labeling

A common approach to (approximately) solving the graph labeling problem follows from an adjacency matrix representation of the graph. A labeled graph can be represented by its adjacency matrix,  $A$ , whenever its edge attributes are univariate. Unlabeled graphs, on the other hand, can be represented by a set of adjacency matrices,  $\{QAQ^T : Q \in \mathcal{Q}\}$ , where  $Q$  is any permutation matrix. Thus, one can define a graph labeling function that finds a permutation matrix that permutes the rows and columns of one graph to match another:

$$Q_{QAP} \triangleq Q_{QAP}(A, B) = \operatorname{argmin}_{Q \in \mathcal{Q}} \|QAQ^T - B\|_F^2, \quad (4)$$

where the permutation matrix  $Q_{QAP}$  induces a labeling of the vertices of  $A$  onto those of  $B$ . A bit of linear algebra simplifies Eq. (4):

$$\begin{aligned} \operatorname{argmin}_{Q \in \mathcal{Q}} \|QAQ^T - B\|_F^2 \\ = \operatorname{argmin}_{Q \in \mathcal{Q}} -\operatorname{tr}(B^T QAQ^T) - \operatorname{tr}(QAQ^T B), \end{aligned} \quad (5)$$

which is equivalent to the standard representation of the quadratic assignment problem (QAP) [2]:

$$\hat{\sigma} = \operatorname{argmin}_{\sigma} a_{\sigma(i), \sigma(j)} b_{ij} = \operatorname{argmin}_{q \in \mathcal{Q}} q_{ij} a_{ij}, q_{ji} b_{ij} \quad (6)$$

where  $\sigma$  is a permutation function,  $\sigma : [n] \mapsto [n]$ . Unfortunately, Eq. (4) is an NP-complete problem [3]. The primary difficulty in solving Eq. (4) is the discrete non-convex constraint set. Thus, it is natural to consider an approximation with the constraints relaxed. Since the convex hull of permutation matrices is the set of doubly stochastic matrices, we define the approximate quadratic assignment problem:

$$Q_{AQAP} \triangleq Q_{AQAP}(A, B) = \operatorname{argmin}_{Q \in \mathcal{D}} \|QAQ^T - B\|_F^2, \quad (7)$$

where  $\mathcal{D}$  is the set of doubly stochastic matrices. When the permutation matrix constraint is relaxed, the equivalence relation shown in Eq. (5) no longer holds. Nonetheless, we proceed by attempting to solve:

$$\hat{Q}_{AQAP} \approx \operatorname{argmin}_{Q \in \mathcal{D}} -\operatorname{tr}(B^T QAQ^T) - \operatorname{tr}(QAQ^T B), \quad (8)$$

considering it an auxiliary function for which we can compute gradients and ascend a likelihood, unlike the permutation constrained case.

The Frank-Wolfe (FW) algorithm is a successive linear programming algorithm for nonlinear programming problems; specifically, for quadratic problems with linear (equality and/or inequality) constraints. Let  $f(Q) = -\operatorname{tr}(B^T QAQ^T) - \operatorname{tr}(QAQ^T B)$ . With each iteration  $j$ , the FW algorithm takes the following steps:

**Step 1: Compute the gradient** The gradient of  $f$  with respect to  $Q$  is given by:

$$\nabla_Q^{(j)} = \partial f / \partial Q^{(j)} = A Q^{(j)} B^T + A^T Q^{(j)} B. \quad (9)$$

**Step 2: Find the closest doubly stochastic matrix** Instead of directly descending this gradient, we search for the direction of the doubly stochastic matrix closest to this gradient. Noting that that direction may be computed by the dot-product operator, we have:

$$W^{(j)} = \operatorname{argmin}_{W^{(j)} \in \mathcal{D}} \langle W^{(j)}, \nabla_Q^{(j)} \rangle. \quad (10)$$

Eq. (10) can be solved as a Linear Assignment Problem (LAP). More specifically, a LAP can be written as:

$$Q_{LAP} \triangleq Q_{LAP}(A, B) = \operatorname{argmin}_{Q \in \mathcal{Q}} \|QA - B\|_F^2, \quad (11)$$

which, when  $B = I$ , can be simplified:

$$\begin{aligned}
Q_{\text{LAP}}(A, I) &= \operatorname{argmin}_{Q \in \mathcal{Q}} \|QA - I\|_F^2 \\
&= \operatorname{argmin}_{Q \in \mathcal{Q}} (QA - I)^\top (QA - I) \\
&= \operatorname{argmin}_{Q \in \mathcal{Q}} A^\top Q^\top QA - 2QA - II \\
&= \operatorname{argmin}_{Q \in \mathcal{Q}} -\langle Q, A \rangle. \tag{12}
\end{aligned}$$

In other words, letting  $B = I$ , the projection of a matrix onto its nearest doubly stochastic matrix is a LAP problem. While Eq. (12) cannot be solved directly, as above, we can relax the permutation matrix constraint to the doubly stochastic matrix constraint:

$$Q_{\text{LAP}}(A, I) = \operatorname{argmin}_{Q \in \mathcal{D}} -\langle Q, A \rangle. \tag{13}$$

Since the permutation matrices are the vertices of the set of doubly stochastic matrices, finding the minimum of Eq. (13) is guaranteed to yield a permutation matrix (as minima are necessarily at the vertices). Thus, letting  $A = \nabla_Q^{(j)}$ , solving Eq. (13)—which is a linear problem with linear and non-negative constraints—is equivalent to solving Eq. (10). Fortunately, the Hungarian algorithm solves any LAP in  $\mathcal{O}(n^3)$  [4], thus this projection is relatively computationally efficient.<sup>1</sup>

**Step 3: Update the direction** Given  $W^{(j)}$ , the new direction is given by:

$$d^{(j)} = W^{(j)} - Q^{(j)}. \tag{14}$$

**Step 4: Line search** Given this direction, one can then perform a line search to find the doubly stochastic matrix that minimizes the objective function along that direction:

$$\alpha^{(j)} = \operatorname{argmin}_{\alpha \in [0,1]} f(Q^{(j)} + \alpha^{(j)}d^{(j)}). \tag{15}$$

This can be performed exactly, because  $f$  is a quadratic function.

**Step 5: Update  $Q$**  Finally, the new estimated doubly stochastic matrix is given by:

$$Q^{(j+1)} = Q^{(j)} + \alpha^{(j)}d^{(j)}. \tag{16}$$

**The grand finale** Steps 1–5 are iterated until convergence, computational budget limits, or some other stopping criterion is met. These 5 steps collective comprise the FW algorithm. Note that while  $Q^{(j)}$  will generally not be a permutation matrix, we do not project  $Q^{(j)}$  back onto the set of permutation matrices between each iteration, as that projection requires  $\mathcal{O}(n^3)$  time. After the final iteration, however, we have  $\hat{Q}_{AQAP}$ , which we project onto the set of permutation matrices:

$$\hat{Q}_{QAP} = \operatorname{argmin}_{Q \in \mathcal{Q}} \langle \hat{Q}_{AQAP}, Q \rangle, \tag{17}$$

1. More efficient algorithms are available for certain special cases, that is, whenever the matrix-vector multiplication operation is fast (for example, when both  $A$  and  $B$  are sparse).

which is a LAP, and yields an approximate solution to QAP. Let FW appended with a projection onto the permutation matrices be denoted by  $\text{QAP}$ .

**Multiple restarts** Note that FW will not generally achieve the global optimum even of Eq. (7), because  $f$  is not necessarily positive definite. This is clear upon computing the Hessian of  $f$  with respect to  $Q$ :

$$\nabla_Q^2 = B \otimes A + B^\top \otimes A^\top, \tag{18}$$

where  $\otimes$  indicates the Kronecker product. This means that the initialization,  $Q^{(0)}$ , could be important. While any doubly stochastic matrix would be a feasible initial point, two choices seem natural: (i) the “flat doubly stochastic matrix,”  $J = \mathbf{1}^\top \mathbf{1}/n_v$ , which is the middle of the feasible region, and (ii) the identity matrix, which is a permutation matrix. Therefore, if we run  $\text{QAP}$  once, we always start with one of those two. If we use multiple restarts, each initial point is “near” the flat matrix. Specifically, we sample  $J'$ , a random doubly stochastic matrix using 10 iterations of Sinkhorn balancing [5], and let  $Q^{(0)} = (J + J')/2$ . We refer to multiple re-starts of  $\text{QAP}$  with subscripts, that is, the performance of  $\text{QAP}_n$  is the best result of  $n$  pseudo-random re-starts of  $\text{QAP}$ . Note that  $\text{QAP}$  natively operates on matrices, which could correspond to either weighted or unweighted graphs.

#### 4.1.2 Using the estimated vertex labels

Given the above strategy to (approximately) solve the graph labeling problem, one can represent the graphs as adjacency matrices, and use many possible classifiers, some of which designed specifically for the labeled graph domain [?]. We focus here on a graph model based strategy, while acknowledging other possibilities below.

A labeled graph classification model based strategy begins by defining a model for the data. For example, [?] define an independent edge graph classification model,  $\mathcal{F}_{\mathbb{L}, \mathbb{A}, Y} = \{F_{\mathbb{L}, \mathbb{A}, Y}(\cdot; \theta) : \theta \in \Theta\}$ , where  $F_{\mathbb{L}, \mathbb{A}, Y}$  can be factorized,  $F_{\mathbb{L}, \mathbb{A}, Y} = F_{\mathbb{L}, \mathbb{A}} F_Y$ . The likelihood can be further expanded,  $F_{\mathbb{L}, \mathbb{A}}(\cdot; \theta) = \prod_{(u,v) \in \mathcal{E}} \text{Bern}(a_{uv}; p_{uv|y})$ , where  $\mathcal{E}$  is the set of edges,  $a_{uv}$  is the value of the edge from the vertex labeled  $v$  to the vertex labeled  $u$ , and  $p_{uv|y}$  is the probability of an edge from  $v$  to  $u$  in class  $y$ . A Bayes plug-in classifier proceeds by estimating the parameters,  $\mathbf{p} = \{p_{uv|y}\}$ ,  $\pi = \{\pi_y\}$ , and plugging them into the Bayes classifier:

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \prod_{(u,v) \in \mathcal{E}} \text{Bern}(a_{uv}; \hat{p}_{uv|y}) \hat{\pi}_y. \tag{19}$$

When the labels are latent, the extension of the Bayes plug-in classifier also plugs-in the estimated vertex labels. Therefore, to use a Bayes plug-in classifier, given a model, one can take the following steps. First, obtain  $s$  “prototpye” graphs that will be  $\text{QAP}$ ed to, and assign each training graph,  $A_i$ , to a prototype,  $A_{i*}$  (see Section 5 for details on the caveats of the various ways of choosing prototypes). Second,  $\text{QAP}$  each graph to its corresponding

prototype:  $\hat{Q}_i = \hat{Q}_{QAP}(A_i, A_{i*})$ , yielding  $\hat{A}_i = \hat{Q}_i A_i \hat{Q}_i^\top$ . Third, estimate the parameters of the model by plugging-in the vertex labels. For instance, one can obtain the maximum likelihood estimators:

$$\{\hat{p}, \hat{\pi}\} = \underset{p, \pi}{\operatorname{argmax}} \prod_{i \in [n] \mid y} p_{uv|y}^{a_{\hat{u}\hat{v}}^{(i)}} (1 - p_{uv|y})^{1 - a_{\hat{u}\hat{v}}^{(i)}} \pi_y, \quad (20)$$

where  $a_{\hat{u}\hat{v}}^{(i)}$  comprise the vertex-label plugin-in adjacency matrix,  $\hat{A}_i$ . Finally, one can then plug the estimates into a Bayes classifier, Eq. (19). Call this algorithm  $\text{BP I} \circ \text{QAP}$ .

## 4.2 Graph Invariant Classifiers

A graph invariant (GI) is any function that maps a graph to a scalar whose value is independent of the vertex labels,  $T : (\mathcal{V}, \mathcal{E}, \mathcal{A}) \mapsto \mathbb{R}$ . By defining a set of GIs, one can embed a collection of graphs into a quotient space invariant to vertex labels. Whenever this quotient space is a subset of finite dimensional Euclidean space, all standard machine learning classifiers may be implemented to solve the classification problem.

To use graph invariants to classify, one first must choose as set of graph invariants. Unfortunately, there is no known set of graph invariants that are collectively optimal for graph classification models. Fortunately, some recent theoretical work shows that certain graph invariants have greater discriminability with regard to certain graph inference tasks [9]. With that in mind, we consider the following invariants described fully in [?]:

- $T_{\text{weight}}$ : total weight of all the edges in the graph
- $T_{\text{maxweight}}$ , is the max over  $d(v)$  for all  $v \in \mathcal{V}$ , where  $d(v)$  is the sum of all weights incident to vertex  $v$
- $T_{MAW_g}$ : a greedy approximation of the maximum average weight (MAW), akin to maximum average degree
- $T_{MAW_e}$ : an eigenvalue approximation of MAW
- $T_{wS_1}$ : the maximum weighted locality statistic, akin to the typical scan statistic, but sums the weight of edges in each neighborhood rather than just the number of edges.
- $T_{Dijkstra}$ : the average Dijkstra path distance between each pair of vertices

For each graph  $G_i$  in the training set, we compute a graph invariant vector:  $T_i : \mathcal{G} \mapsto \mathbb{R}^d$ . We stack these  $n$   $d$ -dimensional vectors to form a matrix  $T \in \mathbb{R}^{n \times d}$ . Letting  $T_{ij}$  indicate the  $i^{th}$  graph invariant of the  $j^{th}$  graph, and normalize each element to be between zero and one according to:  $T_{ij} \leftarrow \frac{T_{ij} - \min_i(T_{ij})}{\max_i(T_{ij}) - \min_i(T_{ij})}$ .

Now, to estimate the class of a test graph, we first compute its invariant vector,  $t$ , and normalize it appropriately. We then apply a variety of machine learning algorithms, including  $k$ NN, linear classifiers, and SVMs. For the below connectome data, the best performing algorithm is the exact confidence weighted classifier [10]. Call this algorithm  $\text{CW} \circ \text{GI}$ .

## 4.3 Graph dissimilarity Classifiers

A *graph dissimilarity* is any non-negative function  $\delta : \mathcal{G} \times \mathcal{G} \mapsto \mathbb{R}_+$ , typically also satisfying several criteria for being a metric. It is becoming increasingly popular to use a *graph kernel* as the dissimilarity [1].<sup>2</sup> Graph kernels have a number of desirable properties, perhaps most notably, that one can then use standard *kernel machines* to classify [6]. Regardless, given any dissimilarity, one can apply at least two different kinds of classifiers.

First, one could implement a nearest neighbor style classifiers, such as the  $k_n$  nearest neighbor ( $k$ NN) classifier. In addition to being universally consistent<sup>3</sup>,  $k$ NN classifiers are computationally efficient, in that they only require  $n+1$  graph dissimilarity computations (assuming a single test graph).

Alternately, one could implement an interpoint-dissimilarity matrix based algorithm [8]. This strategy has the advantage of using all available information to generate a class prediction, but a disadvantage that it requires  $\binom{n+1}{2}$  graph dissimilarity computations. Moreover, it may be more sensitive to outliers.

For concreteness and simplicity, consider the  $k$ NN strategy using the QAP objective function as a dissimilarity. Given a test adjacency matrix,  $A$ , QAP it to all  $\{B_i\}_{i \in [n]}$  training adjacency matrices. Instead of using the argmin of the QAP objective function as above, here we use the min as a dissimilarity, that is, we obtain  $d_i = d(\hat{A}_i, B_i)$  for all  $i \in [n]$ , and sort them:  $d_{(1)} \leq d_{(2)} \leq \dots \leq d_{(n)}$ . Let the  $k_n$  nearest neighbors of  $A$  be the graphs with the  $k_n$  smallest distances,  $\{d_{(1)}, \dots, d_{(k_n)}\}$ . The estimated class of the training sample  $A$  is then the plurality class of the  $k_n$  nearest neighbors:  $\hat{y} = \operatorname{argmax}_y \mathbb{I}\{\sum_{i \in [k_n]} y(i) = y\}$ . Denote this approach by  $k\text{NN} \circ \text{QAP}$ .

## 5 RESULTS

### 5.1 Theoretical results

**Theorem 1.** *Under the graph-class model,  $\mathcal{F}_{\mathbb{G}, \mathbb{Y}}$ , whenever all vertices are stochastically equivalent,  $R_{\mathbb{L}, \mathbb{A}}^* = R_{\mathbb{A}}^*$ . Otherwise,  $R_{\mathbb{L}, \mathbb{A}}^* < R_{\mathbb{A}}^*$ .*

*Proof:* Assume without loss of generality that  $p_{iy} = \pi$  for all  $y \in \mathcal{Y}$ . Thus,  $h_{\mathbb{A}}^*$  and  $h_{\mathbb{L}, \mathbb{A}}^*$  can be written as follows:

$$\hat{y}_{\mathbb{A}} = \operatorname{argmax}_{y \in \mathcal{Y}} F_{\mathbb{A}}|Y \quad (21)$$

$$\hat{y}_{\mathbb{L}, \mathbb{A}} = \operatorname{argmax}_{y \in \mathcal{Y}} F_{\mathbb{A}, \mathbb{L}|Y} \quad (22)$$

□

**Theorem 2.** *Under the graph-class model,  $\mathcal{F}_{\mathbb{G}, \mathbb{Y}}$ , if all the graphs are passed through a vertex shuffle channel, then  $R_{\mathbb{L}, \mathbb{A}}^* = R_{\mathbb{A}}^*$ .*

2. A graph kernel is any function  $\kappa(G, G') = \langle \phi(G), \phi(G') \rangle$ , where  $\phi(\cdot)$  maps from graph space to finite dimensional Euclidean space,  $\phi : \mathcal{G} \mapsto \mathbb{E}^d$ .

3. A sequence of  $k$ NN classifiers is guaranteed to converge to the Bayes optimal classifier if as  $n \rightarrow \infty$ ,  $k \rightarrow \infty$  but  $k/n \rightarrow 0$  [7].

*Proof:* Information processing lemma....  $\square$

Specifically, if one has good reason to believe that the vertex labels lack much class-conditional signal, then one “should” ignore them and use a graph invariant approach. Alternately, if one believes that the vertex labels might contain some class-conditional signal, then one “could” try to use them. However, whether those efforts are fruitful will depend on a bias-variance trade-off. Certainly, the class-conditional entropy for the joint vertex labels and edge attributes is at least as great as the class-conditional entropy for the edge attributes alone,  $\mathcal{I}(\mathbb{L}, \mathbb{A}|Y) \geq \mathcal{I}(\mathbb{A}|Y)$ . But if the vertex have labels, ignoring them will induce more bias.

Another view of how to determine whether to try using the labels concerns the relative effective signal-to-noise-ratios. Specifically, if given the vertex labels, the adjacency matrices are all very similar, then  $F_{\mathbb{A}|\mathbb{L},Y}$  has relatively low entropy. Therefore, one can try to impute the missing labels with perhaps great success. Alternately, if  $F_{\mathbb{A}|\mathbb{L},Y}$  has relatively high entropy, then imputing the missing labels might be too difficult, and we are better off using only  $F_{\mathbb{A}|Y}$ .

It is our contention that many interesting graph classification problems can be cast as latent (or noisy) labeled graph classification problems. For instance, in the MR connectome situation, labels are assigned to vertices using (nonlinear) image registration, which works well, but struggles in certain cases (for example, the boundaries of the regions). Further, in the chemical compound classification setting, one could perhaps use the vertex labels (each vertex is a specific element), to enhance classification performance.

The sequel considers these two complementary approaches to solving latent labeled graph classification problems. Via synthetic data analysis, we demonstrate scenarios in which the above intuition regarding the bias-variance trade-off suggests which approach is likely more effective.

### 5.1.1 Graph labeling approaches

Which one to use depends on whether we first assume there is a “canonical labeling.” A canonical labeling is vertex labeling that introduces a vertex ordering for an equivalence class of graphs (that is, all graphs isomorphic to one another). In the absence of an assumed canonical labeling, a *graph dissimilarity* approach seems more natural to us. In the presence of an assumed canonical labeling, a *graph model* approach seems more natural to us.

When a canonical labeling is assumed to exist, we can pursue a complementary strategy. Note that in the graph dissimilarity approach described above, the vertex labels of adjacency matrix  $A$  is a pairwise property,  $Q_i^A$ , which changed for each  $B_i$ . When a canonical labeling is assumed, we search a simple permutation matrix,  $Q^A$ , invariant with respect to all the other training samples. Sometimes a canonical labeling is somehow known in

advance. Alternately, one could choose a set of “prototype” graphs (possibly with cardinality one), and graph match all other graphs to that set (one). While some efforts have been directed at choosing prototypes (see [1]), it remains a bit of a sticky wicket. Regardless, assuming a canonical labeling (somehow), a graph model approach proceeds as follows.

Like the above  $k$ NN strategy, a straightforward implementation of a *graph model* based approach depends on first choosing a particular graph onto which one implements the graph labeling function. One option would be to label each training graph onto the test graph. While simple, a disadvantage of this approach is that it might diminish the differences between the two classes. An alternative strategy would be to choose a “prototype” from each class,

## 5.2 QAP benchmarks vs. PATH algorithm

We first compare the performance of  $QAP_n$  with recent state-of-the-art approaches on the QAP benchmark library [11]. Specifically, [12] reported improved performance in all but two cases, in which the QPB method of Cremers et al. [13] achieved a lower minimum. We compare  $QAP_n$  with the previous best performing algorithm. In *all* cases,  $QAP_3$  outperforms the previous best result, often by orders of magnitude in terms of relative error. In three cases,  $QAP_{100}$  achieves the absolute minimum. In 12 out of 16 cases, 75%, the simple  $QAP_1$  algorithm outperforms the others (starting with the flat doubly stochastic matrix). See Figure 1 for quantitative comparisons.

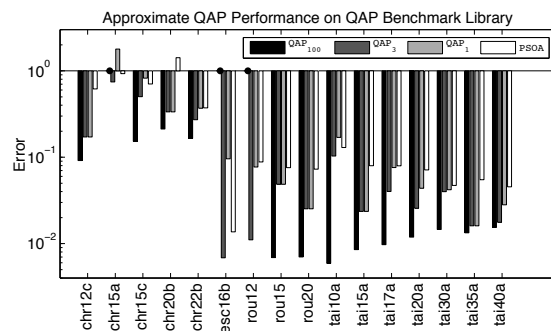


Fig. 1.  $QAP_3$  outperforms PSOA on all 16 benchmark graph matching problems. Moreover,  $QAP_1$  outperforms PSOA on 12 of 16 tests. For 3 of 16 tests,  $QAP_{100}$  achieves the minimum (none of the other algorithms ever find the absolute minimum), as indicated by a black dot. Let  $f_*$  be the minimum and  $\hat{f}_x$  be the minimum achieved by algorithm  $x$ . Error is  $f_*/\hat{f}_x - 1$ .

## 5.3 Unweighted Graph Simulation

$QAP_n$ ’s near perfect performance gave us hope for using  $QAP$  as part of an unlabeled graph classifier. To investigate further, we generated some simulations using

the following assumptions. First, assume an independent edge random graph model for each class:  $F_y = \prod_{(u,v) \in \mathcal{E}} F_{uv|y}$ . For this simple graph scenario, each edge is a Bernoulli random variable,  $F_{uv|y} = \text{Bern}(a_{uv}; p_{uv|y})$ , where  $\{p_{uv|y}\}$  are the likelihood parameters. Then, assume class prior probabilities are equal,  $\mathbb{P}[Y = 1] = \mathbb{P}[Y = 0] = 1/2$ . For simplicity, we sample one graph from each class, meaning  $n = 2$ , and a single training graph sampled according to the class priors. Thus, each simulation is defined by  $\mathcal{M} = (P_0, P_1, n_v)$ , where  $P_0$  and  $P_1$  are the class-conditional likelihoods, and  $n_v$  is the number of vertices per graph. Given a model,  $\mathcal{M}$ , we generate  $n_{MC}$  Monte Carlo trials. For each model,  $R_{\text{chance}} \approx 0.5$ , as estimated by using BPI to classify without first implementing QAP. We estimate Bayes error,  $R_*$ , by using the true parameters and a Bayes plug-in classifier (that is, we use all the labels). We then implement  $\text{BPI} \circ \text{QAP}_1$ , and plot both the misclassification rate and objective function  $f(Q^{(j)})$  as a function of iteration number (not number of restarts, which we hold fixed at one). Figure 2 shows a model (right panel) and results (left panel) of one such simulation. Intriguingly, the first iteration of  $\text{QAP}_1$  seems to basically do the trick. This led us to investigate the relationship between LAP and QAP.

#### 5.4 LAP vs. QAP

Much like the QAP objective function from Eq. (4) can be simplified to Eq. (5), the LAP objective function can be similarly simplified:

$$\underset{Q \in \mathcal{Q}}{\text{argmin}} \|QA - B\|_F^2 = \underset{Q \in \mathcal{Q}}{\text{argmin}} \text{tr}(QAB^T). \quad (23)$$

Letting  $f_{\text{LAP}}(Q) = \text{tr}(QAB^T)$ , the gradient is:

$$\nabla_{\text{LAP}} = 2AB^T. \quad (24)$$

Comparing this gradient to that of QAP—Eq. (9)—one can see that when  $Q^{(j)}$  is the identity matrix, the two gradients are identical. Thus, if QAP is initialized at the identity matrix, the first permutation matrix—output of Step 2—is identical to  $\hat{Q}_{\text{LAP}}$ ; although the line search will make  $Q^{(1)} \neq \hat{Q}_{\text{LAP}}$ , in general. In the above simulation, the first iteration of QAP is essentially the only useful one. Thus, we compare the performance of  $\text{BPI} \circ \text{LAP}$  (dark gray). The performance of LAP and QAP are not statistically different for this simulation. This suggests that for certain problems, LAP (which is  $\mathcal{O}(n^3)$ ) is both an efficient and useful approximation to solving NP-hard graph matching problems. We were unable to find a model for simple graphs in which multiple iterations of QAP improved performance over LAP.

#### 5.5 Weighted simulation

Because  $\text{QAP}_1$  works for both weighted and unweighted, we next simulated multi-graphs (that is, integer valued weights). The below simulation is identical to the

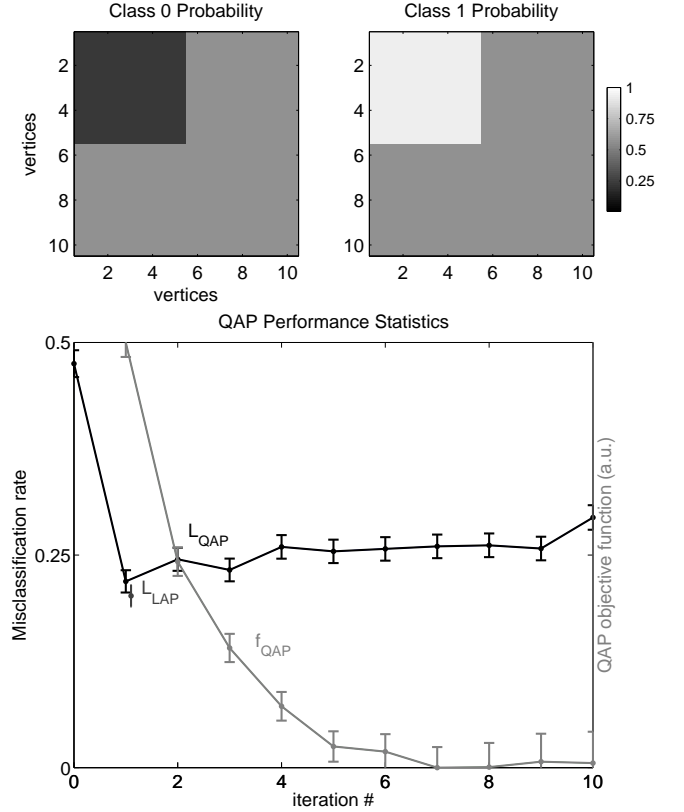


Fig. 2. Homogeneous-kidney-egg model simulation. The left and middle panels show the model parameters for class 0 and 1, respectively. Each edge in the “kidney” in both classes has probability 0.5; in the egg, class 0 edges are sampled with probability 0.25, and class 1 edges are sampled with probability 0.75. The right panel shows the QAP objective function (gray) and misclassification rate (black) as a function of iteration number. LAP does approximately as well as QAP on this (and other) simple graph simulation. Note that  $R_{\text{chance}} \approx 0.5$  and  $R_* \approx 0$  for this simulation. The units of the right-side ordinate are arbitrary. The QAP objective function evaluation prior to any QAP iterations is beyond the bounds of this figure.

above except for  $F_{uv|y}$ —the edge random variable—was Bernoulli before and is now Poisson:  $F_{uv|y} = \text{Poisson}(a_{uv}; \lambda_{uv|y})$ . Figure 3 shows misclassification rate steadily decreasing with each iteration.

#### 5.6 Connectome Classification

A “connectome” is a graph in which vertices correspond to biological neural units, and edges correspond to connections between the units. Diffusion Magnetic Resonance (MR) Imaging and related technologies are making the acquisition of MR connectomes routine [14]. We use 49 subjects from the Baltimore Longitudinal Study on Aging, with acquisition and connectome inference details as reported in [15]. For each connectome, we obtain a  $70 \times 70$  element adjacency matrix, where each element of the matrix encodes the number of streamlines

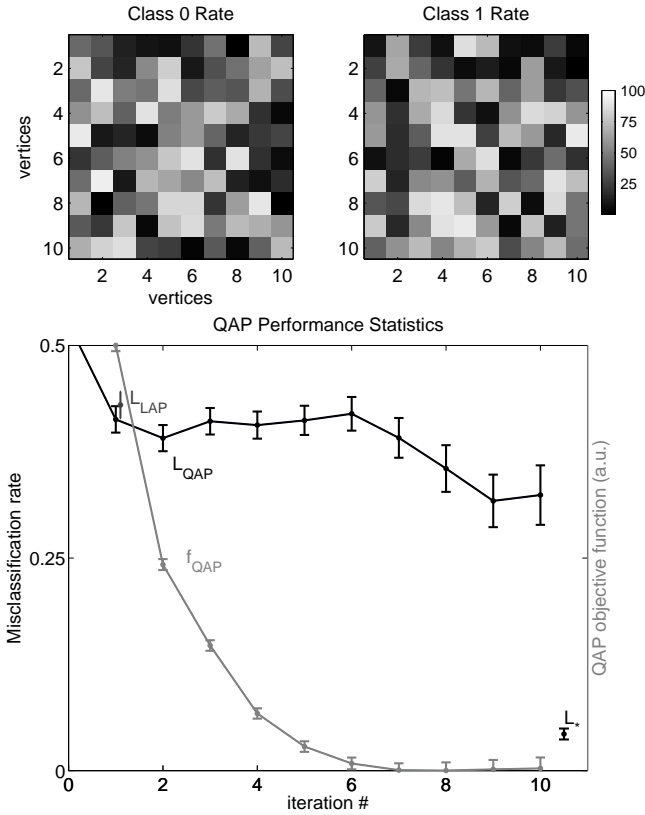


Fig. 3. Multigraph model simulation. The left and middle panels show the parameters for class 0 and 1, respectively. The right panel shows the QAP objective function (gray) and misclassification rate (black) as a function of iteration number; both clearly descending. For comparison, both  $\hat{R}_{LAP}$  and  $\hat{R}_*$  are shown.

between a pair of regions, ranging between 0 and about 65,000. Associated with each graph is class label based on the gender of the individual (24 males, 25 females). Because the vertices are labeled, we can compare the results of having the labels and not having the labels. As such, we implement the following classification strategies. In each case, we use a leave-one-out strategy to evaluate performance:

- N/A-QAP Using the vertex labels, implement a standard 1NN classifier, where distance is the norm of the difference between any pair of adjacency matrices.
- 1-QAP Permute only the vertex labels of the test graph, and then implement  $1NN \circ QAP_1$ .
- 48-QAP Permuting the vertex labels, then implement  $1NN \circ QAP_1$ .
- AVG-QAP Permuting the vertex labels,  $QAP_1$  each of the 48 training graphs to the test graph. Then, given those permuted adjacency matrices, compute the average, and then implement a standard 1NN classifier.
- 1NN-GI Use the graph invariant approach as described above. We provide the normalized graph in-

variants as inputs into a number of standard classifiers, including  $k$ NN, linear classifiers, support vector machines, random forests, and CW. On this data, the CW classifier performed best; we therefore only report its results.

Table 1 shows leave-one-out misclassification rates for the various strategies.

TABLE 1  
MR Connectome Leave-One-Out Misclassification Rates

N/A-QAP	1-QAP	48-QAP	AVG-QAP	1NN-GI
20%	31%	45%	??	25%

## 6 DISCUSSION

In this work, we have presented a number of approaches one could take to classifier graphs. Importantly, when the vertex labeling function is unavailable, one must deal with this uncertainty somehow. We compare a number of approaches on both simulated and connectome data. A multiple-restart Frank-Wolfe approach to approximating QAP outperforms previous state-of-the-art approaches in terms of approximating the graph matching problem. Simulations demonstrate that only the first iteration of such an iterative algorithm, starting from the identity matrix, yields classification performance better than chance. Moreover, the first iteration is identical to LAP, which is a linear problem with linear and non-negativity constraints, and therefore can be solved quite easily.

On a connectome dataset, we compare the performance of various QAP classification algorithms with several graph invariant (GI) based strategies. Of the algorithms that we tried, a graph invariant approach was most effective, even though, in theory, a QAP based approach could have done better (compare the first and last columns of Table 1).

These analyses leave many open questions. Perhaps most interestingly, when might one expect a QAP-based approach to outperform a GI-based approach? Resorting to a generative model, it should be clear that if the class conditional difference is independent of the vertex labels, then there is no reason to even try to implement graph matching. However, if one believes that the labeling function might convey some class-conditional signal (as in the connectome data), then QAP-based approaches could outperform any approach that ignores the labeling function. Which QAP-based approach to use in such a scenario, however, will depend on many factors, including the assumed model and computational resources.

## ACKNOWLEDGMENTS

## REFERENCES

- [1] H. Bunke and K. Riesen, "Towards the Unification of Structural and Statistical Pattern Recognition," *Pattern Recognition Letters*, vol. in press, 2011. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0167865511001309>



- [2] D. Conte, P. Foggia, C. Sansone, and M. Vento, "THIRTY YEARS OF GRAPH MATCHING," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 18, no. 3, pp. 265–298, 2004.
- [3] M. R. Garey and D. S. Johnson, *Computer and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman, 1979.
- [4] R. E. Burkard, M. Dell'Amico, and S. Martello, *Assignment Problems*. SIAM, 2009. [Online]. Available: <http://books.google.com/books?id=nHlzbApLor0C\&pgis=1>
- [5] R. Sinkhorn, "A relationship between arbitrary positive matrices and doubly stochastic matrices," *The Annals of Mathematical Statistics*, vol. 35, no. 2, pp. 876–879, 1964. [Online]. Available: <http://www.jstor.org/stable/2238545>
- [6] V. N. Vapnik, *Statistical Learning Theory*, ser. Wiley Series on Adaptive and Learning Systems for Signal Processing, Communications and Control, S. Haykin, Ed. Wiley, 1998, vol. 2, no. 4. [Online]. Available: <http://www.amazon.com/Statistical-Learning-Theory-Vladimir-Vapnik/dp/0471030031>
- [7] L. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition (Stochastic Modelling and Applied Probability)*. Springer, 1997. [Online]. Available: <http://www.amazon.ca/exec/obidos/redirect?tag=citeulike09-20\&path=ASIN/0387946187>
- [8] R. P. Duin and E. Pkalskab, "The dissimilarity space: bridging structural and statistical pattern recognition," *Pattern Recognition Letters*, vol. in press, 2011. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0167865511001322>
- [9] A. Rukhin and C. E. Priebe, "A Comparative Power Analysis of the Maximum Degree and Size Invariants for Random Graph Inference," *Journal of Statistical Planning and Inference*, vol. 141, no. 2, pp. 1041–1046, 2011.
- [10] K. Crammer, M. Dredze, and F. Pereira, "Exact Convex Confidence-Weighted Learning," *Journal of Multivariate Analysis*, vol. 99, no. 5, pp. 1–8, 2008. [Online]. Available: <http://linkinghub.elsevier.com/retrieve/pii/S0047259X07000784>
- [11] R. E. Burkard, S. E. Karisch, and F. Rendl, "QAPLIB A Quadratic Assignment Problem Library," *Journal of Global Optimization*, vol. 10, no. 4, pp. 391–403, 1997. [Online]. Available: <http://www.springerlink.com/content/n5468q3g33184443/fulltext.pdf>
- [12] M. Zaslavskiy, F. Bach, and J.-P. Vert, "A path following algorithm for the graph matching problem." *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 12, pp. 2227–2242, 2009. [Online]. Available: <http://eprints.pascal-network.org/archive/00004435/>
- [13] C. Schellewald, S. Roth, and C. Schnörr, "Evaluation of Convex Optimization Techniques for the Weighted Graph-Matching Problem in Computer Vision," in *Proceedings of the 23rd DAGMSymposium on Pattern Recognition*. Springer-Verlag, 2001, pp. 361–368.
- [14] P. Hagmann, L. Cammoun, X. Gigandet, S. Gerhard, P. Ellen Grant, V. Wedeen, R. Meuli, J. P. Thiran, C. J. Honey, and O. Sporns, "MR connectomics: Principles and challenges," *J Neurosci Methods*, vol. 194, no. 1, pp. 34–45, 2010. [Online]. Available: [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve\&db=PubMed\&dopt=Citation\&list\\_uids=20096730](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve\&db=PubMed\&dopt=Citation\&list_uids=20096730)
- [15] J. T. Vogelstein, W. R. Gray, J. L. Prince, L. Ferrucci, S. M. Resnick, C. E. Priebe, and R. J. Vogelstein, "Graph-Theoretical Methods for Statistical Inference on MR Connectome Data," *Organization Human Brain Mapping*, 2010.

**William R. Gray** William R. Gray graduated from Vanderbilt University in 2003 with a Bachelors degree in electrical engineering, and received his MS in electrical engineering in 2005 from the University of Southern California. Currently, Will is a PhD student in electrical engineering at The Johns Hopkins University, where he is conducting research in the areas of connectivity, signal and image processing, and machine learning. He is also a member of the technical staff at the Johns Hopkins University Applied Physics Laboratory, where he manages projects in the Biomedicine and Undersea Warfare business areas. Will is a member of IEEE, Eta Kappa Nu, and Tau Beta Pi

**R. Jacob Vogelstein** R. Jacob Vogelstein received the Sc.B. degree in neuroengineering from Brown University, Providence, RI, and the Ph.D. degree in biomedical engineering from the Johns Hopkins University School of Medicine, Baltimore, MD. He currently oversees the Applied Neuroscience programs at the Johns Hopkins University (JHU) Applied Physics Laboratory as an Assistant Program Manager, and has an appointment as an Assistant Research Professor at the JHU Whiting School of Engineering's Department of Electrical and Computer Engineering. He has worked on neuroscience technology for over a decade, focusing primarily on neuromorphic systems and closed-loop brain-machine interfaces. His research has been featured in a number of prominent scientific and engineering journals including the IEEE Transactions on Neural Systems and Rehabilitation Engineering, the IEEE Transactions on Biomedical Circuits and Systems, and the IEEE Transactions on Neural Networks.

**Carey E. Priebe** Buddha in training.

PLACE  
PHOTO  
HERE

PLACE  
PHOTO  
HERE

**Joshua T. Vogelstein** Joshua T. Vogelstein is a spritely young man, engorged in a novel post-buddhist metaphor.