

# PRE-DRAFT: On unlabeled graph classification

Joshua T. Vogelstein<sup>1</sup>, John M. Conroy<sup>2</sup>, Glen M. Coppersmith<sup>3</sup>, Mark Dredze<sup>3</sup>, R. Jacob Vogelstein<sup>4</sup>,  
Carey E. Priebe<sup>1</sup>

1

May 24, 2011

## 1 Introduction

The statistical analysis of collections of graphs is becoming an increasingly popular desideratum [cite]. Specifically, we consider the following idealized and simplified scenario. Let  $\mathbb{G} : \Omega \mapsto \mathcal{G}$  be a graph-valued random variable taking values  $G \in \mathcal{G}$ . Each graph is a 4-tuple:  $G = (\mathcal{V}, \mathcal{E}, \alpha_V, \alpha_E)$ , where  $\mathcal{V}$  is a set of  $|\mathcal{V}| = V$  vertices,  $\mathcal{E}$  is a set of  $|\mathcal{E}| = E$  edges,  $\alpha_V : \mathcal{V} \mapsto \mathcal{A}_V$  is a vertex labeling function, and  $\alpha_E : \mathcal{E} \mapsto \mathcal{A}_E$  is an edge attributing function (for example, edge weights). Given a graph, one can construct an adjacency matrix representation,  $A \in \mathcal{A}_E^{V \times V}$ . When the edge attributing function is binary,  $\mathcal{A}_E = \{0, 1\}$ , resulting in a binary adjacency matrix (generalizations are straightforward). Let  $Y$  be a Bernoulli covariate:  $Y : \Omega \mapsto \{0, 1\}$ , yielding graph *classification* problem. Given a collection of graphs and associated covariates, we assume they were jointly sampled independently and identically from some true but unknown distribution,  $\{(\mathbb{G}_i, Y_i)\}_{i \in [n]} \sim F_{\mathbb{G}, Y}(\cdot; \theta)$ . Note that  $F_{\mathbb{G}, Y}(\cdot; \theta)$  is but one of a (possibly infinite) set of distributions, collectively comprising the model  $\mathcal{F}_{\mathbb{G}, Y} = \{F_{\mathbb{G}, Y}(\cdot; \theta) : \theta \in \Theta\}$ . The goal of such an analysis is to learn about the relationship between  $\mathbb{G}$  and  $Y$ . Standard classification techniques fail in this domain as they typically require classifying finite dimensional Euclidean objects ( $G \in \mathbb{E}^d$ ), whereas the object of interest here are graphs ( $G \in \mathcal{G}$ ). In this work, therefore, we propose a novel extension of classification algorithms appropriate for the graph domain.

## 2 Graph Classification

The graph classification problem may be stated thusly: given  $\mathcal{T}_n = \{(\mathbb{G}_i, Y_i)\}_{i \in [n]} \sim F_{\mathbb{G}, Y}(\cdot; \theta)$  and a new graph,  $\mathbb{G}$ , estimate the new graph's corresponding class,  $Y$ . Given an appropriately defined loss-function, such as misclassification rate:  $L_h = F[h(\mathbb{G}) \neq Y]$ , one can then search for the algorithm  $h^* \in \mathcal{H}$  that minimizes the misclassification rate:

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

In general,  $h^*$  is unavailable and dependent on the model,  $\mathcal{F}_{\mathbb{G}, Y}$ . Instead, one can therefore utilize training data,  $\mathcal{T}_n$ , to obtain  $\tilde{h}$ , an approximation to  $h^*$ :

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

where  $\approx$  indicates that in general, we might not be able to find the actual minimum in the set  $\mathcal{H}$ . Regardless, any approach requires estimating a decision boundary in the space of graphs separating them into two classes. We consider a few distinct such approaches to constructing such a decision boundary:

**Graph dissimilarity approach** Define a dissimilarity on graph spaces:  $d : \mathcal{G} \times \mathcal{G} \mapsto \mathbb{R}_+$ , in which one can compute the dissimilarity between any pair of graphs [cite]. Given an adjacency matrix representation, many such dissimilarities are possible (e.g., graph edit distance, Hamming distance, etc.). It is becoming increasingly popular to use a *graph kernel*,  $\kappa(G, G') = \langle \phi(G), \phi(G') \rangle$ , as the dissimilarity [cite]. Graph kernels have a number of desirable properties, perhaps most notably, that one can then use standard *kernel machines* to classify.

Whether or not one uses a graph kernel, given such a dissimilarity, standard classification algorithms, including  $k_n$  nearest neighbor ( $knn$ ) algorithms [cite] and interpoint-dissimilarity matrix based algorithms [cite], can be straightforwardly applied.

**Graph model approach** Define a random graph model:  $\mathcal{F}_G = \{F_G[\cdot; \theta] : \theta \in \Theta\}$ . Given such a model, one could then, for instance, estimate  $\theta$  and then use standard model-based classifiers (for example, the Bayes plugin classifier) [VogelsteinPriebe11].

**Graph embedding approach** Define an embedding of graphs into finite dimensional Euclidean space:  $\phi : \mathcal{G} \mapsto \mathbb{E}^d$ . Once in  $\mathbb{E}^d$ , one can apply one of many possible standard machine learning or other such approaches [Bunke]. Note that using a graph kernel is closely related to graph embedding.

### 3 Unlabeled sticky wicket

In certain graph classification problems the vertex labels,  $\alpha_V$ , are unobserved. In such scenarios, one must (either implicitly or explicitly) deal with the *graph matching* (GM) problem. In words, graph matching is the operation of finding a set of labels for a collection of vertices on multiple graphs. For the graph classification setting, we require an *approximate* GM approach [Conte04], as we do not expect, in general, for graphs to be isomorphic to one another; rather, we expect graphs within the same class to be “similar” to one another, in some sense. We consider two complimentary approaches to graph matching.

**Adjacency matrix space approach** In this representation, GM can be considered a special case of a quadratic assignment problem (QAP) [cite]. Unfortunately, no polynomial time algorithm is known to solve QAP, although much work has been devoted to this problem [cite]. Fortunately, efficient and approximate QAP solvers, such as the Frank-Wolfe (FW) algorithm [cite], are readily available. Given an approximate assignment of each graph, one can use any of the above classification approaches.

**Graph invariant (GI) approach** A graph invariant is any function that maps a graph to a scalar whose value is independent of the vertex labeling,  $T : (\mathcal{V}, \mathcal{E}) \mapsto \mathbb{R}$ . By defining a set of GIs, one can embed a collection of graphs into an invariant space. This could be considered a special case of the “Graph embedding approach” to classification described above.

## 4 Methods

### 4.1 QAP Approach

In the adjacency matrix approach, we first use an approximate QAP approach to approximate the GM problem. Then, given the estimated labels, we can implement a  $knn$  classifier.

The QAP is defined by the following, where we seek to find a permutation matrix,  $\hat{Q}$ , such that

$$Q_{QAP}(A, B) = Q_{QAP} = \underset{Q \in \mathcal{Q}}{\operatorname{argmin}} \|QAQ^\top - B\|_F^2 \quad (3)$$

where  $A$  and  $B$  are adjacency matrix representations of two different graphs. A bit of linear algebra [HornJohnson] shows that Eq (3) can be simplified:

$$Q_{QAP} = \underset{Q \in \mathcal{Q}}{\operatorname{argmin}} \|QAQ^\top - B\|_F^2 = \underset{Q \in \mathcal{Q}}{\operatorname{argmin}} -\operatorname{tr}(B^\top QAQ^\top) - \operatorname{tr}(QAQ^\top B), \quad (4)$$

which is equivalent to the standard representation of the quadratic assignment problem []:

$$\hat{\sigma} = \underset{\sigma}{\operatorname{argmin}} a_{\sigma(i), \sigma(j)} b_{ij} = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} q_{ij} a_{ij}, q_{ji} b_{ij} \quad (5)$$

where  $\sigma$  is a permutation, that is,  $\sigma : [n] \mapsto [n]$ . As hinted at above, solving Eq. (4) is NP-Incomplete (not known to belong either to P or NP). Because the primary difficulty is the discrete, non-convex constraint set, it is natural to

consider an approximate solution with the constraints relaxed. And because the set of permutation matrices is a subset of the doubly stochastic matrices, we define the approximate quadratic assignment problem:

$$Q_{AQAP} = \operatorname{argmin}_{Q \in \mathcal{D}} \|QAQ^T - B\|_F^2, \quad (6)$$

where  $\mathcal{D}$  is the set of doubly stochastic matrices. Note that when the permutation matrix constraint is relaxed, the equivalence relation shown in Eq. (4) no longer holds, that is:

$$\operatorname{argmin}_{Q \in \mathcal{D}} \|QAQ^T - B\|_F^2 \neq \operatorname{argmin}_{Q \in \mathcal{D}} -\operatorname{tr}(B^T QAQ^T) - \operatorname{tr}(QAQ^T B). \quad (7)$$

Nonetheless, we proceed by trying to solve Eq. (6), 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 (SLP) [] algorithm for nonlinear programming problems, specifically, for quadratic problems with linear (equality and/or inequality) constraints. Let  $f(Q) = \|QAQ^T - B\|_F^2$ . With each step  $k$ , the gradient of  $f$  with respect to  $Q$  is given by:

$$\nabla_Q^{(k)} = AQ^{(k)}B^T + A^T Q^{(k)}B, \quad (8)$$

see [?] pg. 168 for details. Instead of directly ascending this gradient, we traverse 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:

$$W^{(k)} = \operatorname{argmin}_{W^{(k)} \in \mathcal{D}} \langle \nabla_Q^{(k)}, W^{(k)} \rangle. \quad (9)$$

Although  $W^{(k)}$  is constrained only to be a doubly stochastic matrix, it is guaranteed to be a permutation, because the permutation matrices are the vertices of the set of doubly stochastic matrices. Note that Eq. (14) is a linear assignment problem (LAP) []. The Hungarian algorithm is an efficient algorithm for finding the global optimum of any LAP in  $\mathcal{O}(V^3)$  [].<sup>1</sup> Given this direction, one can then perform a line search to find the doubly stochastic matrix that minimizes the objective function along that direction:

$$\alpha^{(k)} = \operatorname{argmin}_{\alpha \in [0,1]} f(Q^{(k)} + \alpha^{(k)} W^{(k)}) \quad (10)$$

This can be performed exactly, because  $f$  is a quadratic function. Finally, the new estimated doubly stochastic matrix is given by:

$$Q^{(k+1)} = Q^{(k)} + \alpha^{(k)} W^{(k)} \quad (11)$$

Eqs. (8)–(11) are iterated until convergence, computational budget limits, or some other stopping criterion is met. Note that while  $Q^{(k+1)}$  is not a permutation matrix, we do not project  $Q^{(k+1)}$  back onto the set of permutation matrices between each iteration, as that projection is a LAP, and requires  $\mathcal{O}(n^3)$  time. At convergence, however, we have  $\hat{Q}_{AQAP}$ , which we project onto the set of permutation matrices:

$$\hat{Q}_{QAP} = \operatorname{argmin}_{Q \in \mathcal{Q}} \langle \hat{Q}_{AQAP}, Q \rangle, \quad (12)$$

which is our approximate solution to QAP. Note that FW will not generally achieve the global optimal even of Eq. (6), because  $f$  is not necessarily positive definite. This is clear upon computing the Hessian of  $f$  with respect to  $Q$ :

$$\nabla_Q^2 = B \otimes A + B^T \otimes A^T, \quad (13)$$

where  $\otimes$  indicates the Kronecker product. This means that the initialization,  $Q^{(0)}$ , will be important. While any doubly stochastic matrix would be a feasible initial point, the “flat doubly stochastic matrix,”  $J = \mathbf{1}^T \mathbf{1}/V$ , is the middle of the feasible region. Therefore, if we run the FW algorithm once, we always start with the flat matrix. 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 rebalancing, and let  $Q^{(0)} = (J + J')/2$ . We call this approach the Multiple Frank-Wolfe (MFW) approach to approximating QAP, versus 1-QAP which refers to a single FW run.

<sup>1</sup>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).

#### 4.1.1 knn $\circ$ QAP

We use the following algorithm to utilize the above approach within a knn classification framework. Given a test adjacency matrix,  $A$ , find  $\hat{Q}_i^A = \hat{Q}_{QAP}(A, A_i)$  for all  $n$  training samples,  $\{A_i\}_{i \in [n]}$ . Given these solutions, let  $\tilde{A}_i = \hat{Q}_i^A A \hat{Q}_i^{A^\top}$  for all  $i$ . Now, define a suitable dissimilarity  $d : \mathcal{A} \times \mathcal{A} \mapsto \mathbb{R}_+$ , one can compute  $d(\tilde{A}_i, A_i)$  for all  $i \in [n]$ , and then one can sort the distances,  $d_{(1)} \leq d_{(2)} \leq \dots \leq d_{(n)}$ . Let the  $k_n$  nearest neighbors of  $A$  be the graphs with the  $k_n$  smallest distances,  $\{d_{(1)}, \dots, d_{(k)}\}$ .<sup>2</sup> The estimated class of the training sample  $A$  is then the plurality class of the  $k_n$  nearest neighbors. Formally,  $\hat{y} = \operatorname{argmax}_y \mathbb{I}\{\sum_{i \in [k_n]} y_{(i)} = y\}$ .

#### 4.1.2 Bayes Plugin $\circ$ QAP

To utilize a Bayes plugin classifier with QAP, we take the following strategy. First, we assume an independent edge random graph model for each class:  $F_y = \prod_{(u,v) \in \mathcal{E}} p_{uv|y}^{a_{uv}} (1 - p_{uv|y})^{(1-a_{uv|y})}$ , where  $\{p_{uv|y}\}$  are the likelihood parameters. For simplicity, we assume class prior probabilities are equal,  $\mathbb{P}[Y = 1] = \mathbb{P}[Y = 0] = 1/2$ . Given a test adjacency matrix,  $A$ , do 1-QAP with respect to a single training graph from each class:  $\hat{Q}_i^A = \hat{Q}_{QAP}(A, A_i)$  for  $i = 0, 1$ . Then, compute the likelihood of  $A$  coming from each class, using the true parameters. Because class-prior probabilities are 1/2, the likelihood is equal to the posterior, so  $\hat{y} = \operatorname{argmax}_y \mathbb{P}[Y = y|A; \{p_{uv|y}\}]$ .

### 4.2 Graph invariant approach

In the graph invariant approach, we first define a set of  $d$  graph invariants: (i) , (ii) , (iii).... 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  $T \in \mathbb{R}^{n \times d}$ . We then whiten this matrix to control for the divergence means and scales of the various graph invariants,  $\tilde{T} = (T - \mu)\Sigma^{-1}$ , where  $\mu$  and  $\Sigma$  are the mean and variance of  $T$ , respectively. Now, to estimate the class of a test graph, we first compute its invariant vector,  $t$ , and normalize it appropriately. Then, we can apply a standard knn algorithm, given a suitably defined dissimilarity (and rule for  $k_n$ ).

## 5 Results

### 5.1 QAP benchmarks vs. PATH algorithm

Before comparing the QAP and GI approach, we first compare the performance of MFW-QAP with recent state-of-the-art approaches on the QAP benchmark library [Cela07 # 38 in PATH paper]. Specifically, in [PATH] this benchmark was used to compare the PATH algorithm against previous state-of-the-art approaches. In all but two cases the PATH algorithm achieved a lower minimum. For those two, the QPB method of Cremers et al. [18 in PATH] achieved a lower minimum. We compare MFW-QAP with the previous best performing algorithm. In *all* cases, MFW-QAP outperforms the previous best result. In 12 out of 16 cases 75%, the simple FW algorithm outperforms the others. See Table 2 and Figure 1 for quantitative results.

### 5.2 Simulations

For each of the 3 below simulations,  $V = 10$ ,  $n = 1000$ . The results consistent show that multiple iterations of 1-QAP improve the objective function, but do not improve the misclassification rate. This led us to investigate the relationship between QAP and Linear Assignment Problem (LAP).

<sup>2</sup>Note that  $k_n$  is a function of  $n$ , typically chosen so that as  $n \rightarrow \infty$ ,  $k \rightarrow \text{infinity}$  but  $k/n \rightarrow 0$ .

Table 1: Comparison of Frank-Wolfe with Minimum Solution and Path Algorithm

#	Problem	Min	FW <sub>100</sub>	FW <sub>3</sub>	FW <sub>2</sub>	FW <sub>1</sub>	min(PATH,QPB)
1	chr12c	11156	12176	13072	13072	13072	18048
2	chr15a	9896	9896	17272	17272	27584	19086
3	chr15c	9504	10960	14274	14274	17324	16206
4	chr20b	2298	2786	3068	3068	3068	5560
5	chr22b	6194	7218	7876	7876	8482	8500
6	esc16b	292	292	294	294	320	296
7	rou12	235528	235528	238134	253684	253684	256320
8	rou15	354210	356654	371458	371458	371458	381016
9	rou20	725522	730614	743884	743884	743884	778284
10	tai10a	135028	135828	148970	157954	157954	152534
11	tai15a	388214	391522	397376	397376	397376	419224
12	tai17a	491812	496598	511574	511574	529134	530978
13	tai20a	703482	711840	721540	721540	734276	753712
14	tai30a	1818146	1844636	1890738	1894640	1894640	1903872
15	tai35a	2422002	2454292	2460940	2460940	2460940	2555110
16	tai40a	3139370	3187738	3194826	3194826	3227612	3281830

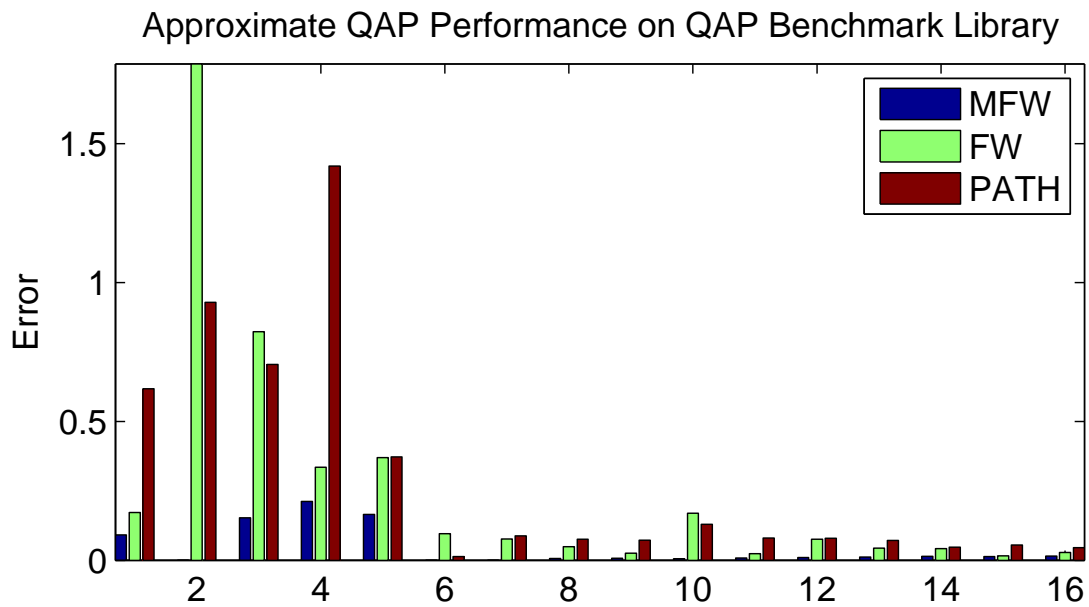


Figure 1: relative performance plot. XXX: would a log yscale be better?

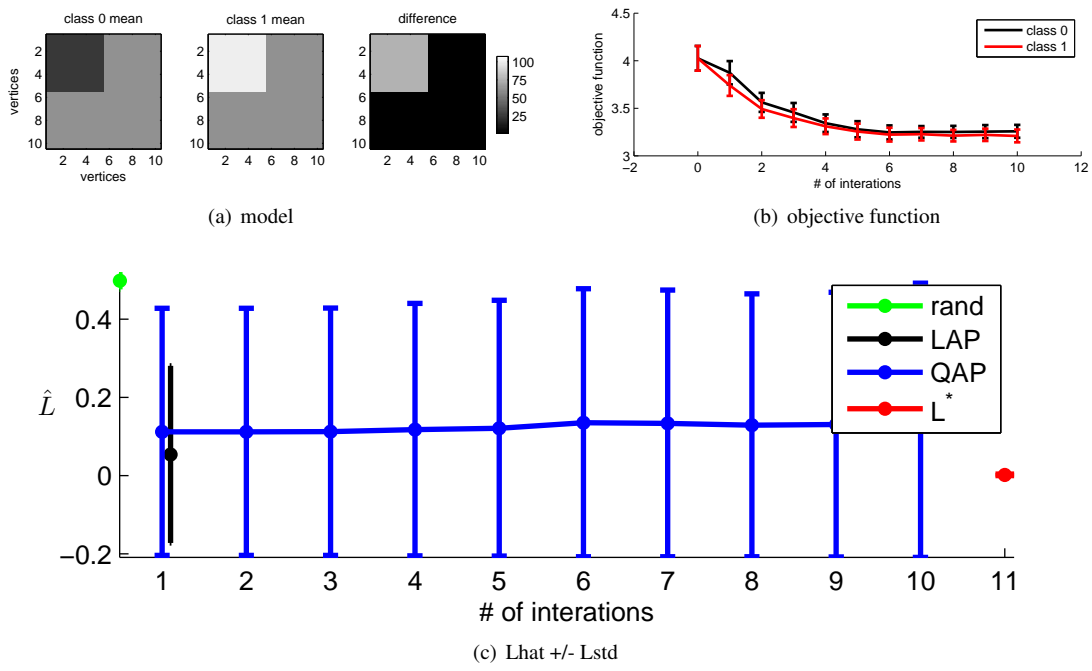


Figure 2: homo kidney egg model

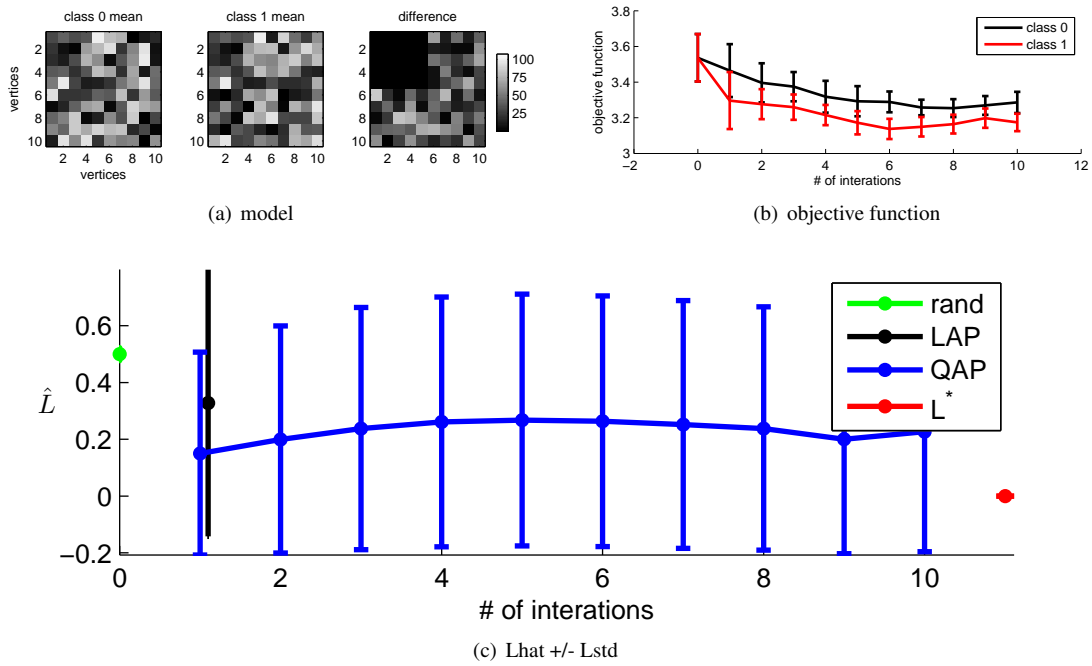


Figure 3: hetero kidney egg model

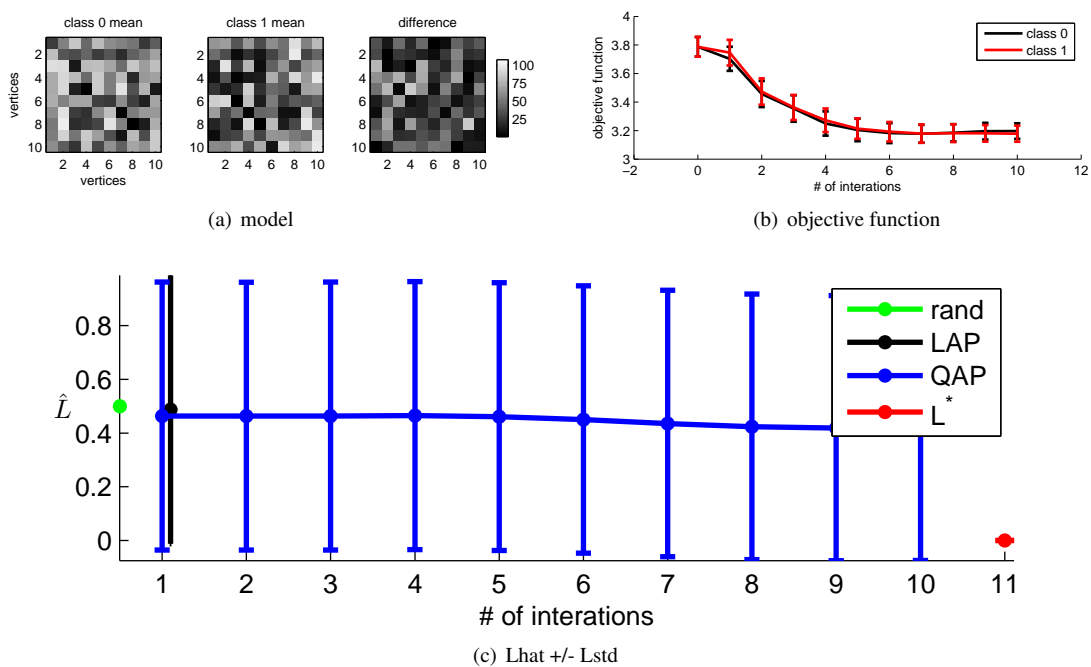


Figure 4: hetero model

### 5.3 LAP vs. QAP

LAPs can be cast as follows:

The LAP approximation for this class problem is given by the following equation:

$$\hat{Q}_{\text{LAP}} = \underset{Q \in \mathcal{Q}}{\operatorname{argmin}} \|AQ^T - B\|_F^2, \quad (14)$$

which is quite similar to QAP, except  $A$  is only post-multiplied by  $Q$ , instead of also being pre-multiplied by  $Q$ . To solve a LAP, we compute the gradient of  $f'(Q) = \|AQ^T - B\|_F^2$ :

$$\nabla_{Q'}^{(k)} = 2AB^T. \quad (15)$$

Interestingly, this is identical to the gradient of QAP, Eq. (8), when  $Q^{(k)}$  is the identical matrix. Thus, if QAP is initialized at the identity matrix, the first iteration is identical to LAP. This suggests that for certain problems, LAP is both an efficient and useful approximation to graph matching. We confirm this intuition by replacing QAP with LAP in the above simulations. As depicted in the above figures, this intuition is consistent with the numerical results. In other words, while naively one might implement an algorithm with exponential time complexity, LAP, which is only quadratic complexity, will often suffice.

### 5.4 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 [cite]. We use 49 subjects from the Baltimore Longitudinal Study on Aging, with acquisition and connectome inference details as reported in [cite]. For each connectome, we obtain a  $70 \times 70$  element adjacency matrix, where each element of the matrix encodes the number of streamlines between a pair of regions [FACT cite], 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:

1. Using the vertex labels, implement a standard 1nn classifier, where distance is the norm of the difference between any pair of adjacency matrices.
2. Permute only the vertex labels of the test graph, and then implement 1nn○1-QAP.
3. Permuting the vertex labels, then implement 1nn○1-QAP as described above.
4. Permuting the vertex labels, QAP 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.
5. Use the graph invariant approach as described above. We provide the normalized graph invariants 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 2: MR Connectome Classification Results

1	2	3	4	5
20%	31%	45%	??	25%

## 6 Discussion

In this work, we have presented a number of approaches one could take



## References