

Leveraging Sequence Classification by Taxonomy-based Multitask Learning

Christian Widmer,¹ Jose Leiva,^{1,2} Yasemin Altun,² and Gunnar Rätsch¹

¹ Friedrich Miescher Laboratory, Max Planck Society, Spemannstr. 39, 72076 Tübingen, Germany

² Max Planck Institute for Biological Cybernetics, Spemanstr. 38, 72076 Tübingen, Germany

Abstract. In this work we consider an inference task that biologists are very good at: bringing together the knowledge that has been obtained in experiments on various organisms in order to understand the differences and commonalities of molecular processes in these related organisms. We look at this problem from an sequence analysis point of view, where we aim at solving the same classification task in different organisms. We investigate the challenge of combining information from related organisms. Here, we consider relation between the organisms that are defined by a tree or graph implied by their taxonomy or phylogeny. Multitask learning, a machine learning technique that recently received considerable attention, considers the problem of learning across tasks that are related to each other. We treat each organism as one task and present three novel multitask learning methods to handle situations in which the relationships among tasks can be described by a hierarchy or by a graph. These algorithms are designed for large-scale applications and scale to problems with a large number of training examples, which frequently appear in sequence analysis. We perform experimental analyses on related synthetic datasets in order to illustrate the properties of our algorithms. Moreover, we consider a problem from genomic sequence analysis, namely splice site recognition. We are able show that using data from 15 eukaryotic organisms one can indeed significantly improve the prediction performance compared to traditional approaches. On a broader perspective, we expect that approaches like the ones presented in this work have the potential to complement and enrich the strategy of homology-based sequence analysis that are currently the quasi-standard in biological sequence analysis.

1 Introduction

Ten years ago, an eight-year lasting collaborative effort resulted in the first completely sequenced genome of a multi-cellular organism, the free-living nematode *Caenorhabditis elegans*. Today, about a decade after this accomplishment, more than 50 eukaryotic genomes have been sequenced and several hundreds more are underway. The genome sequences are the basis for much of the research on the molecular processes in these organisms. Typically, the more closely related the organisms are, the more similar are these processes. For some organisms, certain biochemical experiments facilitating the analysis of particular processes can be performed more readily than for others. This understanding can then be transferred to other organisms, for instance by verifying or refining models of the processes—often at a fraction of the original cost. This is but one example of a situation where transfer of knowledge across organisms is very fruitful.

In computational biology we often study the problem of building statistical models from data in order to predict, analyze, and ultimately understand biological systems. Regardless of the problem at hand, be it the recognition of sequence signals such as splice sites, the prediction of protein-protein interactions, or the modeling of metabolic networks, we frequently have access to datasets for multiple organisms. In this paper, our goal is to develop methods that aim at taking advantage of the data from different organisms in order to improve the performance of the statistical models built for all organisms. We argue that, when building a predictor for a given organism, data from other organisms should be incorporated to the extent of the relation between the organisms.

Since it is assumed that all life can be traced back to an ancient common ancestor, all organisms can ultimately be related by phylogeny. Furthermore, if two organisms share a sufficiently long evolutionary history before divergence, it can be expected that certain biological mechanisms (e.g., splicing) are conserved to some degree. Thus, it is reasonable to assume that we can leverage data from other organisms to enhance model quality for the organism of interest. In Bioinformatics, this is traditionally done by considering sequence homology. This approach, however, is limited to almost exact correspondences of sequences between one or several biological sequences, while it fails to capture other features such as sequence composition that can be used to build an accurate model.

A family of machine learning methods, commonly referred as *domain adaptation* or *transfer learning*, investigates applying a predictor trained with data from a given domain to data from a different one (see e.g., [3, 9]). On the other hand, the so-called *multitask learning* techniques consider the problem of simultaneously obtaining predictors from different domains, by exploiting the fact that the domains are related (see e.g., [4]). Most of these methods assume uniform relations across domains/tasks.¹ However, it is conceivable that from sharing information between closely related domains one can derive more benefit than between domains that are only distantly related according to a given criterion. Hence, it is important to take into account the degree of relatedness among the domains when obtaining the set of models. In this paper, we investigate multitask learning scenarios where we are given *a priori* information about a hierarchy that relates the domains at hand, which is often the case for biological problems. In particular, we treat each organism as a domain and employ the hierarchy given by the phylogeny. The fact that the availability of data describing the same biological mechanism in several organisms is a reoccurring theme makes hierarchical multitask learning approach particularly well suited for many applications in computational biology.

¹ We use the terms *task* and *domain* interchangeably.

Building upon previous work [9], we propose a general framework for leveraging information from related organisms, by ensuring correspondence on model basis, rather than directly comparing sequences. We consider two principal approaches of incorporating relations across domains. In the first approach, models related to task t serve as prior information to the model of t , such that the parameters \mathbf{w}_t of task t are close to the parameters \mathbf{w}_o of the other models. This can be achieved by minimizing the norm of the differences of the parameter vectors, $\|\mathbf{w}_t - \mathbf{w}_o\|$, along with the original loss function. A convenient way of implementing this approach is training models in a top-down manner, where a model is learned for each node in the hierarchy over the datasets of tasks spanned by the node and the parent nodes are used as the prior information, $\|\mathbf{w}_t - \mathbf{w}_{parent(t)}\|$. Here, one can readily use existing inference techniques with slight changes in the implementation. We describe this method in Section 2.2. Alternatively, one can use the models of all tasks as prior information, $\|\mathbf{w}_t - \mathbf{w}_{t'}\|$ for related tasks t and t' , and train the parameters of all tasks jointly. This method is outlined in Section 2.3. Compared to the top-down approach, this formulation involves a larger set of parameters to be jointly optimized during training. However, it can be decomposed into sub-problems which in turn are solved in an iterative manner. Its advantage is that each problem can be trained on smaller datasets compared to the top-down approach where the model of the root node is trained on the union of all datasets.

An alternative to the latter pair-wise approach has been suggested in the context of support vector machines (SVMs) in the special case of two tasks [3]. We extend this idea and design an appropriate kernel function that not only considers the data of the task t , but also the data of all other tasks according to their similarities. We show that this kernel design can be derived by defining predictor functions over the task parameters as well as the parameters of the ancestors of the task. This leads to an essentially effortless multitask approach, since one can use standard SVM implementations by simply implementing the new kernel. We describe the new kernel and its derivation in Section 2.4.

In Section 3 we evaluate the proposed algorithms on simulated data and illustrate some of their properties. Moreover, we consider the problem of splice-site recognition in 15 different eukaryotic genomes and show that the proposed methods can significantly improve the prediction accuracy by combining the information available for all organisms. We conclude the paper with a discussion in Section 4.

2 Hierarchical Multitask Learning

2.1 Preliminaries

We are interested in the problem of learning M functions $f_t : \mathcal{X} \rightarrow \mathcal{Y}$ between the input space \mathcal{X} and discrete output space \mathcal{Y} , where $t = 1, \dots, M$ corresponds to a task. We assume that we are given the relations across tasks via a hierarchy \mathcal{T} , where the tasks are the leaves of \mathcal{T} . Our goal is to make use of the training samples of task t , S_t , while training the other tasks by exploiting the hierarchical information. In the more general case, the relationships do not necessarily need have tree structure. In this case it can be described by a graph with positive weights indicating their relatedness. This graph can be represented by a “relatedness matrix” $\mathbf{\Gamma}$, which can, in the special case of hierarchies, be inferred from the hierarchy \mathcal{T} .² In the following, we propose three

² In the simplest case, $\mathbf{\Gamma}$ is a binary matrix where $\gamma_{vu} = 1$ if node v is the parent of node u and 0 otherwise. Another approach is suggested in Section 2.3.

methods, of which the first assumes a hierarchical structure, while the other two work with the general weighted relatedness graph Γ .

We investigate two principal approaches for exploiting hierarchical information about task relations in a multitask learning (MTL) framework, namely incorporating prior knowledge via *regularization* and via *kernel design*. The first approach is used in two of the proposed methods. Hence, we describe the underlying idea before we go to the technical details. A regularization term is typically used to introduce a penalty for complex solutions into the optimization problem of a learning algorithm. In the existence of prior knowledge, the Empirical Risk Minimization framework [10] incorporates the prior knowledge \bar{f} by

$$\hat{f} = \min_f \left[R(f - \bar{f}) + \sum_{(\mathbf{x}, y) \in S} \ell(f(\mathbf{x}), y) \right], \quad (1)$$

where R is the regularization term that penalizes the deviation of the current model f from the previously obtained (fixed) model \bar{f} , and ℓ is a loss function (such as the squared loss, or the hinge loss) that penalizes the error on the training sample $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. In the multitask approach, we use the models of related nodes as \bar{f} , where relations are given by \mathcal{T} or Γ . Following this scheme, any regularized Machine Learning framework (e.g. regularized least squares, regularized logistic regression) can be extended to include prior information.

As one of the main goals of this work is to provide learning algorithms that scale to large amounts of data, we instantiate the concepts for Support Vector Machines (SVM) [10].³ It has been shown in previous work that SVMs using string kernels such as the Spectrum [6] or the Weighted Degree Kernel (WDK) [7], is well suited for nucleic and protein sequence analysis [1].

2.2 Top-Down Domain Adaptive Support Vector Machines

Our first approach uses a regularized multitask approach, where a model is learned for each node of the hierarchy \mathcal{T} and the parent of node v serves as prior information to v , such that the final model of v is close to the model of its parent.⁴ The idea is to train the models in a top-down fashion, where the most general model is obtained at the root node and more domain specific models are obtained by moving down towards the leaves.

In SVMs, a model f is a linear function parametrized by \mathbf{w} , $f = \langle \mathbf{x}, \mathbf{w} \rangle + b$. Since each model can be trained independently, we drop the indices for tasks and simply use \mathbf{w} for parameters of the current node, and \mathbf{w}_0 for the parameters of the parent node. The primal of the extended SVM formulation is

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|^2 + C \sum_{(\mathbf{x}, y) \in S} \ell(\langle \mathbf{x}, \mathbf{w} \rangle + b, y), \quad (2)$$

where ℓ is the hinge loss, $\ell(z, y) = \max\{1 - yz, 0\}$. We refer to this methodology as Domain Adaptive SVM (DA-SVM), since it describes an extension of standard SVMs to domain adaptation, where the model of the parent node is fixed. From the biological perspective, the penalty term $\|\mathbf{w} - \mathbf{w}_0\|^2$ enforces the models of an organism and the organism it is derived from to be similar, based on

³ Other options, however, may also be suitable to implement the ideas of this work.

⁴ The hierarchy \mathcal{T} corresponds to a binary relatedness matrix Γ where $\gamma_{ij} = 1$ if node i is the parent of node j and 0 otherwise. Note that this similarity is non-symmetric and it renders a node independent of all other nodes given its parent.

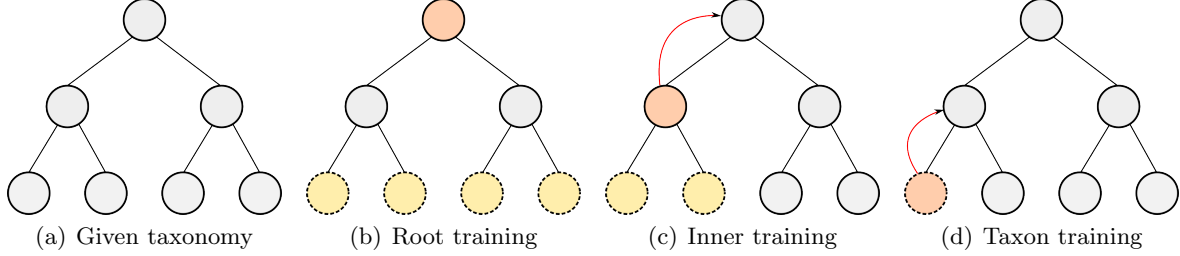


Fig. 1. Illustration of the hierarchical top-down MTL training procedure. In this example, we consider four tasks related via the tree structure in 1(a). Each leaf is associated with a task t and its training sample $S_t = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. In 1(b), training begins by obtaining the predictor at the root node, for which data from all leaves are taken into account in the loss term. Next, we move down one level to train a classifier at an inner node, as shown in 1(c). Here, the loss is measured w.r.t. $S_1 \cup S_2$ and the classifier is forced to be similar to the parent solution via the regularization term, as indicated by the red arrow. Finally, in 1(d), we obtain the final classifier of task 1 by only taking into account S_1 to measure the loss, while again regularizing versus the parent predictor. The procedure is applied in a top-down manner to the remaining nodes, until we obtain a predictor for each leaf.

the assumption that mutations take place in small steps. Hence, most properties of the model is conserved through evolution.

In order to employ kernels, we derive the dual of the above formulation:

$$\max_{\alpha} -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i \underbrace{\left(\sum_{j=1}^m \alpha'_j y_i y'_j k(\mathbf{x}_i, \mathbf{x}'_j) \right)}_{p_i} - 1,$$

$$\text{s.t. } \alpha^T \mathbf{y} = 0, \quad 0 \leq \alpha_i \leq C \quad \forall i \in \{1, n\},$$

where n and m are the number of training samples for \mathbf{w} and \mathbf{w}_0 respectively; α_i represent the dual variable of the current learning problem for sample (\mathbf{x}_i, y_i) , whereas the α'_j are the dual variables obtained from the prior model; in the case of the linear kernel it is described as $\mathbf{w}_0 = \sum_{j=1}^m \alpha'_j y'_j \mathbf{x}'_j$. The resulting prediction is performed by

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{j=1}^m \alpha'_j y'_j k(\mathbf{x}, \mathbf{x}_j) + b.$$

In the standard SVM formulation, we have $p_i = -1$, for $i = 1, \dots, n$. In this extended formulation, the p_i can be pre-computed and passed to the underlying SVM-solver as the linear term of the corresponding quadratic program (QP). To provide implementations that readily deal with large-scale learning problems, we have extended the SVM implementations *LibSVM* [2] and *SVMLight* [5] to handle prior information (which is going to be publicly available).

Top-Down Hierarchical Learning An illustration of the training procedure is given in Figure 1. In a top-down manner a predictor f is obtained at each node v , where the loss L is evaluated on the union of training data $S = \cup\{S_t \preccurlyeq_{\mathcal{T}} v\}$ at the leaves underneath the current node v , while the current model f is regularized versus the parent predictor f_p . Training is complete, once we have obtained a predictor f_t for each task t . The algorithm is described in Appendix 5.1.

2.3 Support Vector Machines with Pairwise Task Regularization

In the previous section, we described a method that learns models for internal nodes of the hierarchy and imposes regularization between a node and its parent. In this section, we describe a method where the relation across tasks are modeled directly via pairwise regularization. We refer to this method as *Pairwise* learning. To incorporate pairwise task regularization, we consider the relatedness matrix $\mathbf{\Gamma}$ between tasks. If only the hierarchy \mathcal{T} is given, one may, for instance, first compute the distance of the two tasks in the hierarchy by counting the number (or adding the weights) of edges separating the two tasks (“hop-distance”). This distance may then be transformed into a task relatedness score.⁵

The *prior* information in this method comes from the models of similar tasks. All the models are trained jointly by optimizing the regularization term along with the loss ℓ which is measured separately for each \mathbf{w}_t on the corresponding dataset S_t ,

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_M} \frac{1}{2} \sum_{t=1}^M \sum_{s=1}^M \gamma_{ts} \|\mathbf{w}_t - \mathbf{w}_s\|^2 + \sum_{t=1}^M C_t \sum_{(\mathbf{x}, y) \in S_t} \ell(\langle \mathbf{x}, \mathbf{w}_t \rangle, y).$$

The parameter C_t is used to trade off the loss of task t on the training sample versus its regularization terms in order to control generalization performance. In our experiments, we set $C_t = C$, $t = 1, \dots, M$, for simplicity. The biological interpretation of this formulation is that if two organisms share a sufficiently long evolutionary history before divergence (reflected by γ_{ts}), it can be expected that certain aspects of the model describing the biological mechanism are conserved. Thus, we use the task-similarity matrix $\mathbf{\Gamma}$ to control how strongly we regularize each $(\mathbf{w}_t, \mathbf{w}_s)$ pair to be close to each other.

Decomposition The pairwise formulation learns models only for the leaf nodes of the hierarchy, as opposed to the DA-SVM approach, where the number of models to learn is given by all the nodes of the hierarchy. Its comparative disadvantage, on the other hand, is that it couples the parameters of all models which leads to a large optimization problem. In order to overcome this limitation, we developed a decomposition of the optimization problem that allows the global solution to be obtained by solving a series of SVM-like quadratic programs iteratively until convergence. It can be shown that the above optimization problem has a fixed point that coincides with the optimization problem of

$$\min_{\mathbf{w}_t} \frac{1}{2} \lambda_t \|\mathbf{w}_t - \mathbf{r}_t\|^2 + C \sum_{(\mathbf{x}, y) \in S_t} \ell(\langle \mathbf{x}, \mathbf{w}_t \rangle, y), \quad \forall t \tag{3}$$

$$\mathbf{r}_t = \sum_{s \neq t} \beta_{ts} \mathbf{w}_s, \tag{4}$$

$$\lambda_t = \sum_{s \neq t} \gamma_{ts}, \quad \beta_{ts} = \gamma_{ts} / \lambda_t.$$

This formulation decouples the optimization problem into individual tasks and, hence, retains scalability. It states that the regularization prior \mathbf{r}_t of each task should be in the convex hull of $\{\mathbf{w}_s\}_{s \neq t}$. It can be solved iteratively by finding the optimal \mathbf{r}_t for the current \mathbf{w}_t (cf. (4)) and finding the optimal \mathbf{w}_t for the current \mathbf{r}_t (cf. (3)) for each task until convergence. Note that, the difference

⁵ Although it may be practically important for accurate hierarchical multitask learning, how to choose the transformation is not in the main focus of this work. In our experiments, we used $\gamma = 1 - d/2$, where d is the normalized distance and γ is the computed relatedness between the tasks.

between (3) and (2) is the prior model, where in the first case it is the weight average of related tasks and in the latter it is the parent model. Hence, the SVM implementations in Section 2.2 can be used to solve (3). Furthermore, kernelization follows similarly. Investigating the relation between the dual and primal parameters,

$$\mathbf{w}_t = \sum_{i:\mathbf{x}_i \in S_t} \alpha_i y_i \mathbf{x}_i + \sum_{s \neq t} \beta_{ts} \sum_{j:\mathbf{x}_j \in S_s} \alpha_j y_j \mathbf{x}_j,$$

reveals that the pairwise regularization method results in “borrowing” support vectors of related models with respect to the task similarities.

2.4 Multitask-Kernel Learning

In Sections 2.2 and 2.3, we presented two hierarchical multitask approaches based on regularization with respect to related models. In this section, we propose an alternative approach, *MultiKernel*, by defining a kernel function that incorporates data from related tasks. We first define a similarity matrix $\mathbf{\Gamma}$ that includes ancestry relationships: In particular, γ_{tr} , the similarity of a task t and a node r . As before, it is inversely related to the distance of t and r , if r is an ancestor of t with respect to \mathcal{T} , $t \preceq_{\mathcal{T}} r$, and 0 otherwise. We then define the predictor function of task t as

$$f_t(\mathbf{x}) = \mathbf{w}_t^T \mathbf{x} + b_t = (\mathbf{u}_t + \sum_{r \prec t} \gamma_{tr} \mathbf{v}_r)^T \mathbf{x} + b_t,$$

where $\{\mathbf{v}\}_r$ are the internal node parameters and \mathbf{u}_t are the leaf node (task) parameters. Here, the parameters of the node at each level in the hierarchy represent the corresponding level of abstraction. More precisely, the parameters of the root node capture the most common structure across all tasks and as we descend on the hierarchy, the parameters capture the deviation from previous level. Our goal is to achieve the proper specialization level for each task by combining these parameters in the prediction function.

We propose to obtain $\{\mathbf{u}_t\}$ and $\{\mathbf{v}_r\}$ by solving the regularized loss function for all tasks,

$$\min_{\{\mathbf{u}_t\}, \{\mathbf{v}_r\}} \frac{1}{2} \sum_{t=1}^T \|\mathbf{u}_t\|^2 + \frac{1}{2} \sum_{r=1}^R \|\mathbf{v}_r\|^2 + C \sum_{t=1}^M \sum_{(\mathbf{x}, y) \in S_t} \ell(\langle \mathbf{x}, \mathbf{w}_t \rangle, y),$$

where ℓ is the hinge loss. Note that the tasks are related to each other through the internal node parameters \mathbf{v} as in the case of DA-SVM. However, the loss term ℓ is evaluated only on the leaf nodes, as opposed to DA-SVM where ℓ is evaluated at all nodes in the hierarchy by combining the relevant datasets. Instead of learning the internal node models by error minimization and then enforcing the models of the parent-child nodes to be similar, the goal here is to learn the internal node models directly by minimizing the error of the leaf nodes, the tasks.

It can be shown that the dual formulation of the problem above is equivalent to that of standard SVMs

$$\begin{aligned} \max_{\alpha} & -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \tilde{k}(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i \\ \text{s.t.} & \quad \alpha^T \mathbf{y} = 0, \quad 0 \leq \alpha_i \leq C \quad \forall i \in \{1, n\}, \end{aligned}$$

where the kernel is defined over all related data,

$$\tilde{k}(\mathbf{x}_i, \mathbf{x}_j) = \tilde{\gamma}_{t(i),t(j)}k(\mathbf{x}_i, \mathbf{x}_j).$$

Here, $t(i)$ denotes the task that data point \mathbf{x}_i belongs to and $\tilde{\gamma}_{ts}$ are the entries of $\mathbf{I} + \mathbf{\Gamma}\mathbf{\Gamma}^T$. Hence, the kernel \tilde{k} incorporates the interaction among tasks in addition to the interaction among data points. This formulation is a generalization of the domain adaptation method of [3], where the original kernel k is reweighted with the task-similarity matrix $\mathbf{\Gamma}$ derived from the hierarchy. If $\mathbf{\Gamma}$ is positive semi-definite, the resulting kernel \tilde{k} is still a positive semi-definite kernel. One can readily use existing kernel methods by simply implementing \tilde{k} .

3 Results

Experimental Setup We performed experiments on two types of data sets. First, we considered synthetic sequence data, which was created by applying mutations to a Position Specific Scoring Matrix (PSSM) according to a pre-defined binary tree structure (details can be found in Appendix 5.2, the code for data generation will be made available). Balanced, equally sized training data sets, each with 100 examples, were sampled for each of the leaves. As test set, additional 5,000 examples were sampled for each task. We used the area under the ROC curve to evaluate the prediction performance. (An evaluation using the the area under the precision recall curve yields similar results.)

As a second dataset we considered the problem of splice site recognition. We generated and used labeled sequences acceptor splice sites of 15 different eukaryotic genomes which were similarly generated in [8, 9] (see Figure 3 for the taxonomic relation between the organisms). For each task, we obtained 10,000 training examples and an additional test set of 6,000 examples. We normalized the datasets such that there are 100 negative examples for each positive example. We report the area under the precision recall curve (auPRC), which is an appropriate measure for unbalanced detection problems [7].

Experiments were performed for each of the presented MTL methods (*DA-SVM*, *Pairwise*, *MultiKernel*) and the following two baseline methods. In *Union*, all data are combined into one dataset $S = \cup_{t=1}^M S_t$ and a single global model is obtained that is used to predict on all domains. On the other extreme, we consider the baseline method *Plain*, where an individual SVM is trained on the data S_i of each domain separately, without taking data from the other domains into account. For each method, the regularization parameter C was optimized via cross-validation with 4 and 3 splits for toy and splice data, respectively. After obtaining the optimal C , we retrained the model on all available training data. The performance was measured on the separate test sets, which are considered large enough to obtain reliable estimates of the predictors’ performances.

Results on the Toy Datasets We consider two different settings of generating the toy data sets: one with relatively high mutation rate (larger differences between neighbouring tasks in the hierarchy) and one with small mutation rate (all tasks are closely related). For this study we analyze how the different methods perform in these two settings with respect to the number of tasks (8, 16, 32, 64 tasks). The results are shown in Figure 2. We can observe that the two baseline methods perform very differently: While the *Plain* method performs essentially indifferent for different numbers of tasks and mutation rates, the *Union* method performs very well when the mutation rate is low, but quite poorly for large mutation rates. Moreover, the performance of *Union* degrades for an increasing number and, hence, diversity of tasks. This is particularly pronounced for high mutation

rates. The three proposed methods all perform better than the two baseline methods, indicating that it pays-off to take the additional data and the relation between the tasks into account. However, among the three proposed methods there is no method that is consistently performing better.

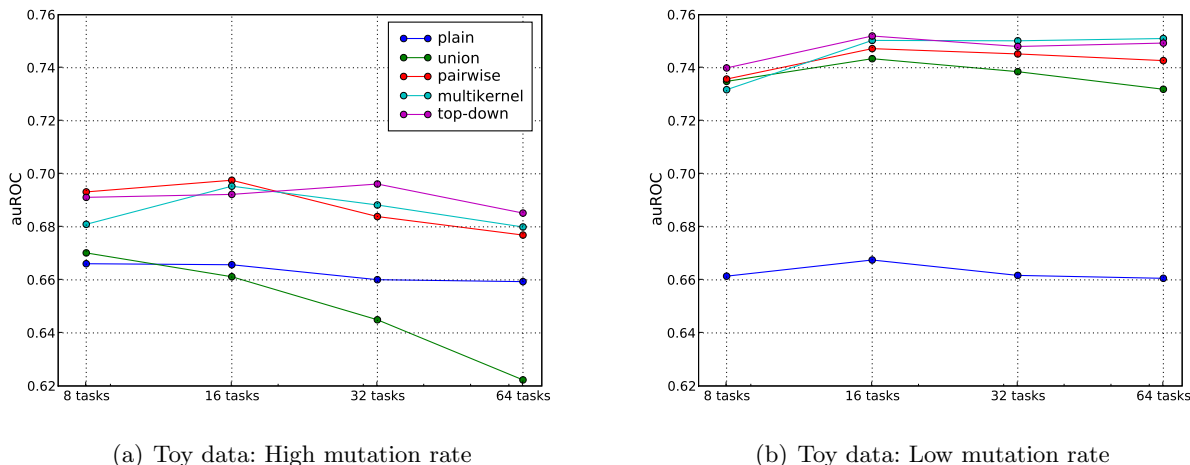


Fig. 2. Results for the five considered methods on the toy data sets with high and low mutation rate. Shown are the average areas under the ROC curves for different numbers of tasks that have been generated from full binary trees.

Results on the Splice Dataset On the toy dataset, all three MTL methods perform similarly well and clearly outperform the two baselines on all eight tasks. For the splice site dataset, examining the average performance across all organisms (the last column in Figure 4), we also observe superior performance of the MTL methods over the baseline methods. The *Plain* method is again outperformed by all other methods, which emphasizes the importance of using data from related organisms. Comparing the hierarchical MTL methods, we observe that the top-down approach performs slightly better than the other MTL methods. We conjecture that this is due to a suboptimal choice of the similarity matrix, where tree-hop-distance was used for the pairwise and multikernel approach.

Zooming in the results for individual organisms, we observe the same trend, where MTL methods outperform the baselines in most cases and the *Plain* method yields the worst performance. For 11 out of 15 organisms, all MTL methods perform better than the baseline. The *Top-Down* method performs always better than the baseline methods, except for *A. nidulans*, which is the only fungal organism in our set. The gain from hierarchy is more pronounced in the lower levels of the hierarchy, e.g., for *A. thaliana*, *O. sativa*, *O. latipes*, and *D. rerio*, where data from closely related organisms are used to leverage the learning process. An exception to this behaviour is seen on the mammals branch, where the performance gain from MTL methods get smaller for *H. sapiens* and essentially diminishes for *M. musculus*. Our conjecture is that the hierarchy for this branch is not deep enough in order to represent closely related organisms. Including more similar organisms, and hence extending the hierarchy to capture more evolutionary steps, can improve the performance gain for these organisms.

It is worthwhile investigating the performance of the methods on *A. nidulans*, which is the child of the root node and, hence, most distantly related to the other organisms. We observe that the

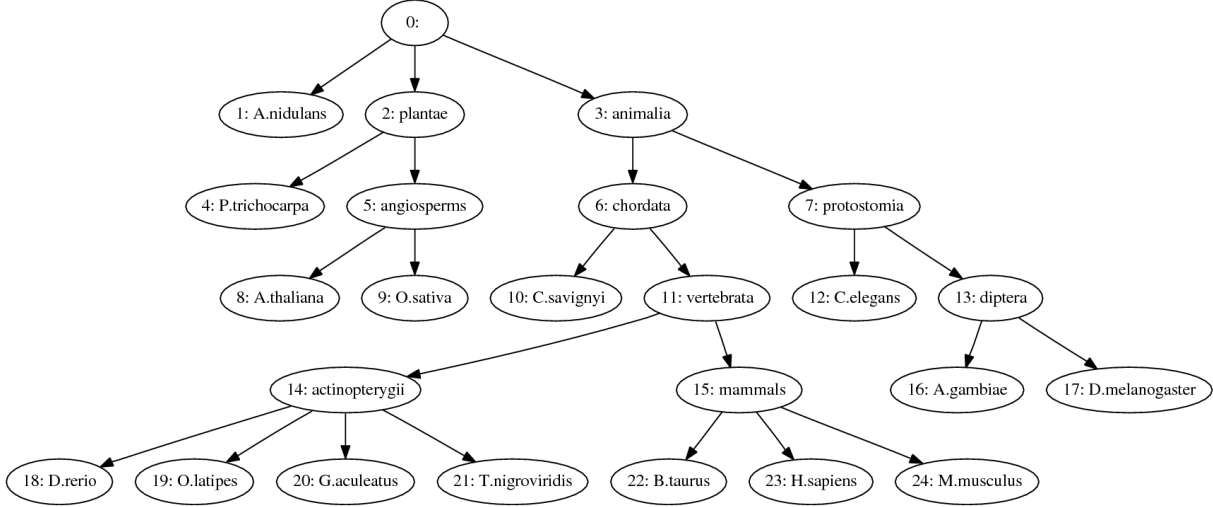


Fig. 3. Taxonomy used to relate organisms for the splicing experiment

Union method performs the worst on this organisms and the MTL methods perform on the same level as the *Plain* method. Hence, MTL methods manage to improve the performance for closely related organisms, while causing (essentially) no performance loss on the distantly related ones.

Moreover, it is interesting to observe that *A. nidulans* and *C. elegans* are the only organisms/tasks, for which the *Union* method is considerably worse than the *Plain* method. This hints at major differences between the recognition of splice sites in these two organisms. This is less surprising for *A. nidulans*, as it is for *C. elegans*. It appears worth investigating this property also for other nematode genomes to understand when these differences have been acquired.

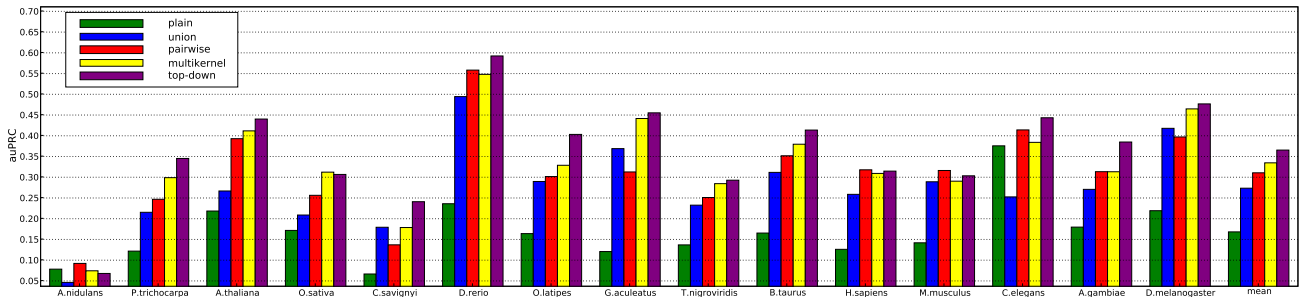


Fig. 4. Results for the splice site datasets from 15 eukaryotic genomes: Shown are auPRC performances of the five considered methods for each organism (two baseline methods: *Plain*, *Union*; three proposed methods: *Top-Down*, *Pairwise*, and *Multikernel*). We can observe that the two of the MTL methods consistently outperform the baseline methods. In 14/15 cases, the hierarchy information lead to improved prediction results.

4 Discussion and Conclusion

We outlined two principle ways of leveraging information from related organisms where the relation across organisms are defined with respect to a given hierarchy. We presented three algorithms that readily deal with large scale problems such as those frequently encountered in genomic sequence analysis. We have demonstrated that our methods outperform baseline methods on synthetic data, and data from splice site prediction. On the one hand, the poor performance of *Plain* relative to the MTL methods shows that exploiting information from other tasks is in fact beneficial. On the other hand, the poor result of *Union* demonstrates that there is no single model that fits all tasks equally. Clearly, methods that carefully combine the data from different tasks according to their relatedness perform best.

We are encouraged by the good performance of the *DA-SVM* method, as it provides a fast, simple and non-parametric way of exploiting hierarchical information. Inferring an accurate task-similarity matrix \mathbf{I} proves to be non-trivial, therefore one should think of additional ways of using the hierarchy to ease that task. Experiments on the splicing data shows that a simple task-similarity matrix based on tree-hop-distance can be suboptimal particularly for cases when edge lengths are unequal. Our immediate future work involves experiments where edge lengths are incorporated into the similarity matrix. For the phylogeny trees, the edge lengths can be given by evolutionary years.

To cope with unequal edge lengths, the *DA-SVM* method can be extended to locally estimate the optimal degree of regularization towards the parent model by performing a cross-validation at each node v to determine an optimal C_v (which trades off loss versus the similarity of the model to the parent node). With this extension, the hierarchy essentially restricts the hyper-parameter search space by decoupling hyper-parameters and removes the necessity to search over a grid of all combinations of parameters.

For all of the presented methods, we plan to provide publicly available scalable implementations based on modified versions of *SVMLight* [2] and *LibSVM* [5]. Lastly, we would like to emphasize that in computational biology there are a great number of problems for which there exists datasets for multiple organisms. We expect that hierarchical MTL methods can indeed make a big difference for these problems. Therefore, the methods and implementations that we presented can be employed for a wide range of problems.

Acknowledgments We would like to thank Sören Sonnenburg for help with the implementation of algorithms presented in this work. Also, we acknowledge Cheng Soon Ong for providing the raw data sets from which we sampled data for the splicing experiment. Furthermore, we would like to thank Gabrielle Schweikert, Georg Zeller and Klaus-Robert Mueller for helpful discussions. This work was supported by the DFG grant RA1894/1-1.

References

1. A. Ben-Hur, C.S. Ong, S. Sonnenburg, B. Schölkopf, and G. Rätsch. Support vector machines and kernels for computational biology. *PLoS Comput Biol*, 4(10):e1000173, Oct 2008.
2. C.C. Chang and C.J. Lin. LIBSVM: a library for support vector machines, 2001.
3. H. Daumé. Frustratingly easy domain adaptation. In *ACL*. The Association for Computer Linguistics, 2007.
4. T. Evgeniou, C.A. Micchelli, and M. Pontil. Learning multiple tasks with kernel methods. *Journal of Machine Learning Research*, 6:615–637, 2005.
5. T. Joachims. SVMlight: Support Vector Machine. *SVM-Light Support Vector Machine* <http://svmlight.joachims.org/>, University of Dortmund, 1999.
6. C. Leslie, E. Eskin, and W. S. Noble. The spectrum kernel: A string kernel for SVM protein classification. In *Proceedings of the Pacific Symposium on Biocomputing*, pages 564–575, 2002.
7. G. Rätsch and S. Sonnenburg. *Accurate Splice Site Detection for Caenorhabditis elegans*. MIT Press, 2004.
8. G. Rätsch, S. Sonnenburg, J. Srinivasan, H. Witte, K.-R. Müller, R. Sommer, and B. Schölkopf. Improving the *C. elegans* genome annotation using machine learning. *PLoS Computational Biology*, 3(2):e20, 2007.
9. G. Schweikert, C. Widmer, B. Schölkopf, and G. Rätsch. An empirical analysis of domain adaptation algorithms. In *Advances in Neural Information Processing System, NIPS*, volume 22, Vancouver, B.C., 2008.
10. V.N. Vapnik. *The nature of statistical learning theory*. Springer Verlag, New York, 1995.

5 Appendix

5.1 Top-Down Domain Adaptive SVM Algorithm

Input: Training Samples S_t for $t \in \{1, M\}$, Taxonomy \mathcal{T}
Output: Trained Predictors $F = \{f_t : t \in \{1, M\}\}$
 $F = \{\}$;
 $V = \{v_{root}\}$;
while $|V| > 0$ **do**
 $v = pop(V)$;
 $S = \cup_{\{S_t \preccurlyeq_{\mathcal{T}} v\}}$;
 $f_p = parentpred(v)$;
 $f = \min_f R(f - f_p) + \sum_{(x,y) \in S} L(y_i - f(x_i))$;
 if $isleaf(v)$ **then**
 $F = append(F, f)$
 end
 else
 $V = append(V, children(v))$;
 end
end
return F

Algorithm 1: Top-Down Domain Adaptive SVM Training Procedure

5.2 Toy Data Generation

With the application of our methods to computational biology in mind, a first task was to generate toy data for which we are able to tightly control the properties, such as task relatedness and difficulty of sub-problems. Since we are particular interested in learning problems in the domain of sequence biology, we generated data that are sequences $X \in \Sigma^L$ of a fixed length L , where Σ is an alphabet (a finite, non-empty set) and $L \in \mathbb{N}_+$. For DNA, the alphabet consists of four characters $\Sigma = \{A, G, C, T\}$, corresponding to the four nucleotides present in genomic sequences.

We define generative models $p(x|\Theta)$ which are used to generate sequence data according to the parameter vector Θ . In particular, we restrict the model space to be a zero-th order Markov Chains, where $\Theta_{i,a}$ defines the probability of observing symbol $a \in \Sigma$ at position i . Then, the probability of a sequence X of length L is given by

$$P(X|\Theta) = \prod_{i=1}^L \Theta_{i,X_i} (1 - \Theta_{i,X_i}).$$

To generate data for a binary classification task, we use two generative models, one for the positive class and one for the negative class. For simplicity, we set the distributions for the negative class to be uniform (e.g. $\Theta_{i,a} = 1/|\Sigma| \forall i \in [1, L], a \in \Sigma$) and therefore non-informative. Thus, the information contained for discrimination is solely in the positive class. This process outlines the data generation with respect to one model.

We are interested in a hierarchy of models that represent the evolution process. We achieve this by defining a model for the root node and mutate the parameter vector Θ at each edge on

the path from the root to the leaves, according to a previously fixed tree structure. At each step a new $\Theta_{child} = mut(\Theta_p)$, is obtained by applying a mutation operator mut . Data is generated at the leaves of the tree using the corresponding model.

When designing the mutation operator mut , we would like to control the magnitude of the mutation and the difficulty (e.g. separability) of the classification tasks at the leaves. To do this, we turn to the concept of relative Entropy or Kullback-Leibler divergence

$$D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}.$$

The KL-divergence is a non-symmetric measure of the difference between two probability distributions. We assess the separability of a classification problem, by taking the KL-divergence between the generative model for the positive class $P_{pos}(X|\Theta)$ and the generative model of negative class $P_{neg}(X|\Theta_{unif})$. Note that this is equivalent to measuring the entropy of the generative model of the positive class directly, as the negative class is non-informative. As the entropy of the class gets larger, the difficulty of the task increases. To assess the magnitude of mutation based on the generative models, we compute KL-divergence between a model at a node $P_{child}(X|\Theta_{child})$ and the model at the parent node $P_{parent}(X|\Theta_{parent})$. The mutation rate is large, if this difference is large.

For a given separability s and mutation rate r , we define a mutation for each edge as

$$\begin{aligned} \Theta_{child} = mut(\Theta_{parent}, s, r) &:= \max_{\Theta} const \\ \text{s.t. } D_{KL}(P_{\Theta}||P_{\Theta_{parent}}) &= r \\ E(P_{\Theta}) &= s. \end{aligned}$$

A mut operator respecting these constraints can thus be obtained by a at each edge of the tree.

The above discussion gives rise to the following recursive algorithm that generates datasets according to a given tree structure mutation rate r and entropy s .

Input: Current node v , Model parameters Θ , Tree structure T
Output: Multitask Dataset $D = \{(D_i, i) : \text{for each leaf } i\}$
if $isleaf(v)$ **then**
 | **return** $(sample(\Theta), v)$
end
else
 | $D = \{\}$;
 | **for** $w \in children(v, T)$ **do**
 | $\Theta_w = mut(\Theta, s, r)$;
 | $D_w = toy(w, T, \Theta_w)$;
 | $D = D \cup D_w$;
 | **return** D ;
 | **end**
end

Algorithm 2: Hierarchical Toy Sequence Generation