

## 二次錘計画法による多タスク学習算法

加藤 毅<sup>†1</sup> 鹿島 久嗣<sup>†3</sup>  
杉山 将<sup>†4</sup> 浅井 潔<sup>†1,†2</sup>

関連する複数のタスクがあるとき、これらを別々に解くより同時に解いたほうが効果的なことがある。このようはアプローチは多タスク学習 (MTL) と呼ばれる。この論文では新しい多タスク学習法を提案する。提案法はすべてのタスクの解が必ず近くなることが保証されるように、タスク間の関係を局所的に制御する。この考えをサポートベクトルマシンに適用すると、その最適化問題は二次錘計画問題で表現できることを示す。多クラス識別、順序回帰、リンク予測の問題が多タスク学習問題として扱えることを示し、実験により提案法の有用性を示す。

## An SOCP Formulation for Multi-Task Learning

TSUYOSHI KATO,<sup>†1</sup> HISASHI KASHIMA,<sup>†3</sup>  
MASASHI SUGIYAMA<sup>†4</sup> and KIYOSHI ASAI<sup>†1,†2</sup>

When we have several related tasks, solving them simultaneously has been shown to be more effective than solving them individually. This approach is called *multi-task learning* (MTL). In this paper, we propose a novel MTL algorithm. Our method controls the relatedness among the tasks *locally*, so all pairs of related tasks are guaranteed to have similar solutions. We apply the above idea to support vector machines and show that the optimization problem can be cast as a *second-order cone program*, which is convex and can be solved efficiently. The usefulness of our approach is demonstrated in ordinal regression, and link prediction, each of which can be formulated as a structured multi-task problem.

## 1. はじめに

多くの応用において、いくつかの関連する学習タスクがある。関連するタスクは共通する要素を共有していることがしばしばあり、これらを同時に解いたほうが、別々に解くより利点があると期待できる。このようなアプローチは、多タスク学習 (Multi-Task Learning, MTL) と呼ばれており、理論的にも実験的にも有用性が証明されている<sup>2),3),5),7),9),10),20),23),26),32)</sup>。

典型的には、タスク間の関係は関連タスク間の解の近さを要求することで実現される。しかし、多タスク学習法はこれまで次のような欠点があった。従来法では、タスク間の関連性を、すべてのタスク間の解の距離の合計の上限を与えることで実現させていた<sup>23)</sup>。このような制約は大域的制約 (Global Constraint)<sup>27)</sup> と呼ばれる。この制約では、関連するタスクの解は必ずしも近ならず、いずれかはかなり離れた解になりうる。

この論文では、上述の短所を克服する新しい多タスク学習法を提案する。提案法では、関連タスクの各ペアに対して、それぞれ上限距離を設定する。この制約を局所的制約 (Local Constraint) と呼ぶ。さらに、提案法では、タスク関連ネットワークによってあらわされるタスクの関連性を扱うことができる。我々は、この考えをサポートベクトルマシン (SVM) の枠組みに適用すると、学習算法は二次錘計画問題 (Second-Order Cone Program)<sup>6)</sup> に帰着できることを見出した。二次錘計画問題は、凸計画問題であり、大域解を効率的に計算できる。実験を通して、提案法は既存の多タスク学習法より予測性能が良いことを示す。

多タスク学習の考え方は様々な応用において有用である：

- 多クラス識別: 複数のクラスからクラスラベルを予測する<sup>1)</sup>,
- 順序回帰: ユーザの好み (「好き」「普通」「嫌い」) など順序のあるクラスラベルを予測する<sup>23),25)</sup>,
- リンク予測: 部分的にリンクの有無が分かっているネットワークに対して、残りのリンクの有無を予測する<sup>4),17),29),31)</sup>,

†1 東京大学大学院新領域創成科学研究科

Graduate School of Frontier Sciences, University of Tokyo

†2 産総研生命情報工学研究センター

AIST Computational Biology Research Center

†3 東京大学情報理工学系研究科

Department of Mathematical Informatics, University of Tokyo

†4 東京工業大学大学院情報理工学系研究科

Graduate School of Information Science and Engineering, Tokyo Institute of Technology

● 協調フィルタリング: 多くのユーザの好みの情報を使って商品の推薦などを行う<sup>8)</sup>. 提案法は多クラス識別問題, 順序回帰, リンク予測において従来法を上回る予測性能が有することを実験を通して示す. 協調フィルタリングに関しても良好な実験結果は得られているが, 紙面の制約により割愛する. 定理の証明も割愛せざるを得なかったが, 文献<sup>18)</sup>には掲載したので, そちらを参照されたい.

#### 問題設定

$M$  個の 2 クラス識別タスクを考える. それらのタスクは入出力空間  $\mathcal{X} \times \{\pm 1\}$  を共有している<sup>30)</sup>. 当分の間, 簡単のため,  $\mathcal{X} \subset \mathbb{R}^d$  を仮定する. 第 3 節以降,  $\mathbb{R}^d$  を再生カーネルヒルベルト空間に拡張する.  $M$  タスクを学習するとし, 第  $i$  タスクは  $n_i$  個の学習用例題  $(\mathbf{x}_{t,i}, y_{t,i}) \in \mathcal{X} \times \{\pm 1\}$  ( $t = 1, \dots, n_i$ ) をもつ. 例題の総数を  $\ell$  とする:  $\ell \equiv \sum_{i=1}^M n_i$ . 最終的に得たいものは各識別タスクのスコア関数

$$f_i(\mathbf{x}; \mathbf{w}_i, b_i) = \mathbf{w}_i^\top \mathbf{x} + b_i, \quad \text{for } i = 1, \dots, M,$$

である. ただし,  $\mathbf{w}_i \in \mathbb{R}^d$  および  $b_i \in \mathbb{R}$  は第  $i$  タスクのモデルパラメータである.

## 2. タスクネットワークを使った局所的 multitask learning: 線形版

本節では, 新しい multitask learning 法を提案する.

### 2.1 基本的な考え方

タスク間に関連がない時,  $M$  個の学習タスクを別々に解くことができる:

$$\forall i = 1, \dots, M: \quad \frac{1}{2} \|\mathbf{w}_i\|^2 + C_\alpha \sum_{t=1}^{n_i} \text{Hinge}(f_i(\mathbf{x}_{t,i}; \mathbf{w}_i, b_i), y_{t,i}). \quad (1)$$

ただし,  $C_\alpha \in \mathbb{R}_+$  は正則化定数であり,  $\text{Hinge}(\cdot, \cdot)$  はヒンジ損失関数と呼ばれ:

$$\text{Hinge}(f, y) \equiv \max(1 - fy, 0)$$

のように定義されている. 式 (1) はサポートベクトルマシン (SVM)<sup>28)</sup> として知られている. 第 1 項目は 2 つのクラスを分けるマージンの単調減少関数である. SVM は 2 つのクラスを分ける超平面のうちマージンが最大となるものを見つける.

もしそれぞれのタスクの学習用例題が少ないと, この個別に解く方法の性能は悪くなる. 性能を良くするにはより多くの例題が必要となる.

この問題を克服するために, 我々は関連する解が互いに近くなるように細工を施す. つまり, 最適化問題 (1) 上で次の制約を付け加えることとする:

$$\forall i, \forall j: \quad \frac{1}{2} \|\mathbf{w}_i - \mathbf{w}_j\|^2 \leq \rho. \quad (2)$$

つまり, タスクの解の差に上限  $\rho \in \mathbb{R}_+$  を与えるのである. この制約を局所的制約と呼ぶ<sup>27)</sup>. しかし,  $b_i$  には制約を加えないことにする. なぜなら, これはタスク間で値が大きく異なることがあるからである. 本論文では, 一つの上限  $\rho$  のみを使うことにするが, タスクのペアごとに異なる上限  $\rho_{i,j}$  を与えるように一般化するのは簡単である. この制約 (2) によって, タスク間で訓練用例題を共有することになり, 結果的に訓練用例題が増えるような効果を持つ.

式 (1) と (2) を組み合わせると

$$\frac{1}{2M} \sum_{i=1}^M \|\mathbf{w}_i\|^2 + C_\alpha \sum_{i=1}^M \sum_{t=1}^{n_i} \text{Hinge}(f_i(\mathbf{x}_{t,i}; \boldsymbol{\theta}), y_{t,i}) + C_\rho \rho,$$

を得る. ただし,  $C_\rho \in \mathbb{R}_+$  は非負の定数である.

すると, 学習のための最適化問題は次のようになる:

$$\min \quad \frac{1}{2M} \sum_{i=1}^M \|\mathbf{w}_i\|^2 + C_\alpha \|\boldsymbol{\xi}\|_1 + C_\rho \rho, \quad (3)$$

$$\text{wrt } \mathbf{w} \in \mathbb{R}^{Md}, \mathbf{b} \in \mathbb{R}^M, \boldsymbol{\xi}_\alpha \in \mathbb{R}_+^\ell, \rho \in \mathbb{R}_+,$$

$$\text{subj. to } \forall i, \forall j \in \mathbb{N}_M: \quad \frac{1}{2} \|\mathbf{w}_i - \mathbf{w}_j\|^2 \leq \rho,$$

$$\forall i \in \mathbb{N}_M, \forall t \in \mathbb{N}_{n_i}: \quad y_{t,i} (\mathbf{w}_i^\top \mathbf{x}_{t,i} + b_i) \geq 1 - \xi_{t,i}^\alpha,$$

$$\text{where } \mathbf{w} \equiv [\mathbf{w}_1^\top, \dots, \mathbf{w}_M^\top]^\top, \quad \boldsymbol{\xi}_\alpha \equiv [\xi_{1,1}^\alpha, \dots, \xi_{n_1,1}^\alpha, \xi_{1,2}^\alpha, \dots, \xi_{n_M,M}^\alpha]^\top.$$

一般に, 制約が少なければ凸問題は高速に解くことができる. 問題 (3) の欠点は  $\mathbf{w}_i$  のすべてのペアに制約があることである. したがって, この問題は多大な計算コストがかかる. タスクをノードとする全結合のネットワークでこの制約をあらわすとする (図 1(a) 参照). 後の節で, エッジをいくつか削っても性能が下がらないことを実験的に示す. 残ったエッジの集合を  $\mathcal{E} \equiv \{i_k, j_k\}_{k=1}^K$  とあらわすことにする. すると, 最適化問題は次のようになる:

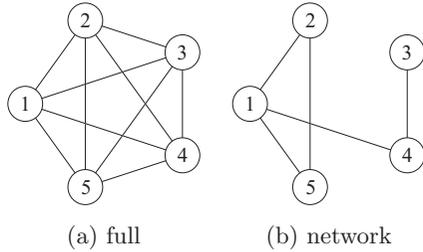


図 1 MTL-SVM の主要な特徴はそれぞれのタスクペアのモデルパラメータの差に上限を設けることである．MTL-SVM (local/full) の上限は (a) に示すようにすべてのタスクペアに与えられる．MTL-SVM (local/network) の上限は (b) に示すように一部のペアだけに設けられる．MTL-SVM (global/full) および MTL-SVM (global/network) の上限は、それぞれ、すべてのタスクペアの差の合計、および、一部のタスクペアの差の合計に与えられる．

$$\begin{aligned}
 \min \quad & \frac{1}{2M} \sum_{i=1}^M \|\mathbf{w}_i\|^2 + C_\alpha \|\boldsymbol{\xi}\|_1 + C_\rho \rho, \\
 \text{wrt} \quad & \mathbf{w} \in \mathbb{R}^{Md}, \mathbf{b} \in \mathbb{R}^M, \boldsymbol{\xi}_\alpha \in \mathbb{R}_+^\ell, \rho \in \mathbb{R}_+, \\
 \text{subj. to} \quad & \forall k \in \mathbb{N}_K : \frac{1}{2} \|\mathbf{w}_{i_k} - \mathbf{w}_{j_k}\|^2 \leq \rho, \\
 & \forall i \in \mathbb{N}_M, \forall t \in \mathbb{N}_{n_i} : y_{t,i} (\mathbf{w}_i^\top \mathbf{x}_{t,i} + b_i) \geq 1 - \xi_{t,i}^\alpha, \\
 \text{where} \quad & \mathbf{w} \equiv [\mathbf{w}_1^\top, \dots, \mathbf{w}_M^\top]^\top, \quad \boldsymbol{\xi}_\alpha \equiv [\xi_{1,1}^\alpha, \dots, \xi_{n_1,1}^\alpha, \xi_{1,2}^\alpha, \dots, \xi_{n_M,M}^\alpha]^\top.
 \end{aligned} \tag{4}$$

今後、エッジ集合を  $\mathcal{E}$  とするネットワークをタスクネットワークと呼ぶことにする (図 1(b) 参照)．タスクネットワークの自動獲得も挑戦的な課題ではあるが、本研究では、タスクネットワークは事前に与えられていると仮定する．

### 3. タスクネットワークを使った局所的多タスク学習：カーネル版

本節では、非線形な識別面も学習できるようカーネルトリックを提案法に適用する．

#### 3.1 双対形式

$K_{\text{fea}}$  を半正定値行列とし、その  $(s, t)$  要素は特徴ベクトル  $\mathbf{x}_s$  と  $\mathbf{x}_t$  の内積とする：

$$K_{s,t}^{\text{fea}} \equiv \langle \mathbf{x}_s, \mathbf{x}_t \rangle.$$

これは特徴ベクトルのカーネル行列である．タスク間のカーネルを

$$K_{\text{net}}(\boldsymbol{\lambda}) \equiv \left( \frac{1}{M} \mathbf{I}_M + \mathcal{U} \boldsymbol{\lambda} \right)^{-1},$$

と定義する．ただし、 $\boldsymbol{\lambda} \in \mathbb{R}_+^K$  は  $K$  次元のパラメータベクトルである．ここで、

$$\mathcal{U} \boldsymbol{\lambda} \equiv \sum_{k=1}^K \lambda_k \mathbf{U}_k, \quad \mathbf{U}_k \equiv \mathbf{E}^{i_k i_k} + \mathbf{E}^{j_k j_k} - \mathbf{E}^{i_k j_k} - \mathbf{E}^{j_k i_k}$$

と定義した  $\mathbf{E}^{(i,j)} \in \mathbb{R}^{M \times M}$  は、 $(i, j)$  要素のみが 1 で、それ以外は 0 の行列である．これは第  $k$  エッジの重みが  $\lambda_k$  のグラフラプラシアンカーネル<sup>33)</sup> として知られている．行列  $\mathbf{Z} \in \mathbb{N}^{M \times \ell}$  をそれぞれの例題がどのタスクに対応するか示すように次のように定義する：

$$\mathbf{Z}^\top \equiv \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{0}_{n_1} & \cdots & \mathbf{0}_{n_1} \\ \mathbf{0}_{n_2} & \mathbf{1}_{n_2} & \cdots & \mathbf{0}_{n_2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_{n_M} & \mathbf{0}_{n_M} & \cdots & \mathbf{1}_{n_M} \end{bmatrix}.$$

すると、タスクの情報は  $\ell \times \ell$  のカーネル行列  $\mathbf{Z}^\top K_{\text{net}}(\boldsymbol{\lambda}) \mathbf{Z}$  によってあらわされる．アダマール積  $\circ$  を用いて、この 2 つのカーネル行列  $K_{\text{fea}}$  および  $\mathbf{Z}^\top K_{\text{net}}(\boldsymbol{\lambda}) \mathbf{Z}$  は

$$K_{\text{int}}(\boldsymbol{\lambda}) \equiv K_{\text{fea}} \circ (\mathbf{Z}^\top K_{\text{net}}(\boldsymbol{\lambda}) \mathbf{Z}), \tag{5}$$

と統合できる．このパラメータつき行列  $K_{\text{int}}(\boldsymbol{\lambda})$  は、 $\boldsymbol{\lambda} \geq \mathbf{0}_K$  である限り、半正定値であることが保証されている<sup>13)</sup>．

次の定理に示すように、上記の記法を用いると、双対形式を簡潔に記述できる：

**Theorem 3.1.** 主問題 (4) の双対問題はパラメータつき統合カーネル行列  $K_{\text{int}}(\boldsymbol{\lambda})$  を使って次のようにあらわされる．

$$\begin{aligned}
 \min \quad & \frac{1}{2} \boldsymbol{\alpha}^\top \text{diag}(\mathbf{y}) K_{\text{int}}(\boldsymbol{\lambda}) \text{diag}(\mathbf{y}) \boldsymbol{\alpha} - \|\boldsymbol{\alpha}\|_1, \\
 \text{wrt} \quad & \boldsymbol{\alpha} \in \mathbb{R}_+^\ell, \boldsymbol{\lambda} \in \mathbb{R}_+^K, \\
 \text{subj. to} \quad & \boldsymbol{\alpha} \leq C_\alpha \mathbf{1}_\ell, \quad \mathbf{Z} \text{diag}(\mathbf{y}) \boldsymbol{\alpha} = \mathbf{0}_M, \quad \|\boldsymbol{\lambda}\|_1 \leq C_\rho.
 \end{aligned} \tag{6}$$

□

$\forall i, \mathbb{N}_M \forall t \in \mathbb{N}_{n_i}$  に対して、 $\alpha_{t,i}$  は

$$\boldsymbol{\alpha} = [\alpha_{1,1}, \dots, \alpha_{n_1,1}, \alpha_{1,2}, \dots, \alpha_{n_{M-1},M-1}, \alpha_{1,M}, \dots, \alpha_{n_M,M}].$$

を満たすような双対変数  $\alpha \in \mathbb{R}_+^{\ell}$  の要素である。

証明は、文献<sup>18)</sup>を参照されたい。 $\alpha$  および  $\lambda$  の最適解は、 $\ell_1$  ノルムの影響から、疎になる傾向にある。次の定理により、提案法は二次錘計画法の既製のエンジンを使って効率的に学習することができる。

**Theorem 3.2.** 双対問題 (6) は二次錘計画問題に帰着される。 □

カーネル行列  $K_{\text{fea}}$  の定義を線形カーネルから非線形カーネルに変えることによって、提案する多タスク学習法は非線形な識別面も得られるようになる。さらに、文字列カーネルやグラフカーネル<sup>11),12),14)–16),19),21),22)</sup> などを使うことによって、ベクトルではない構造データも扱うことができる。

予測の段階では、第  $j$  タスクに属する未知の例題  $x$  は

$$f_j(x) = \sum_{i=1}^M \sum_{t=1}^{n_i} \alpha_{t,i} y_{t,i} k_{\text{fea}}(x_{t,i}, x) k_{\text{net}}(i, j) + b_j,$$

によって識別できる。ただし、 $k_{\text{fea}}(\cdot, \cdot)$  および  $k_{\text{net}}(\cdot, \cdot)$  は、それぞれ、特徴のカーネル関数とタスクのカーネル関数である。

## 4. 議 論

### 4.1 標準的な SVM との関係

提案する多タスク学習法は標準的な SVM を特殊ケースとして含む。実際に、タスク数が 1 のとき、問題 (6) は標準的な SVM 最適化問題に帰着される。このように、提案法は SVM の自然な拡張とみることができる。

タスクネットワークにエッジが全くなかったとき、提案法は、個別に SVM を訓練する場合と全く等価になる。タスクネットワークのエッジをなくした SVM を **Individually Learned SVM (IL-SVM)** と呼ぶことにする。

### 4.2 大域的制約・局所的制約

本研究に先んじて、多タスク学習法のための異なる算法がすでに提案されている<sup>9),10)</sup>。その手法は関連するタスク間のすべてのペアに対して、解の距離の合計に上限  $\frac{1}{2} \sum_{i,j=1}^M \|w_i - w_j\|^2 \leq \rho$  を与えるものであった。合計に上限を与える制約を大域的制約<sup>27)</sup>と呼ぶ。タスク間のすべてのペアに対する合計に上限を与えることから、タスクネットワークは完全結合しているとみることができる。このような理由からこのアプローチを **MTL-SVM (global/full)** と呼ぶことにする。大域的制約は距離の合計に上限を与えてい

るだけなので、一部の距離が大きくなることを許してしまう。実際に、これが顕著に性能を低下させてしまうことを第 5 節にて実験的に示す。対照的に、我々は各タスクペアの距離に上限を与えている。これによって、すべてのタスクペアに対して解が近くなることが保証される。各ペアに上限をかけることを局所的制約<sup>27)</sup>と呼び、我々のアプローチを **MTL-SVM (local/full)** と呼ぶことにする。

Micchelli & Pontil<sup>23)</sup> は強い関連のあるペアだけに適用するほかの定式化も与えている： $\frac{1}{2} \sum_{k=1}^K \|w_{i_k} - w_{j_k}\|^2 \leq \rho$ 。ただし、各々の  $k = 1, \dots, K$  に対して、第  $i_k$  タスクと第  $j_k$  タスクは強く関連しているが、それ以外はそうではないとしている。この定式化はタスクネットワークの情報を利用するので、このアプローチを **MTL-SVM (global/network)** と呼ぶ。我々の局所的制約による方法もタスクネットワークを  $\forall k \in \mathbb{N}_K : \frac{1}{2} \|w_{i_k} - w_{j_k}\|^2 \leq \rho$  のように利用できる。このアプローチを **MTL-SVM (local/network)** と呼ぶ。この 4 つの方法は図 1 にまとめた。

### 4.3 リンク予測

多タスク学習はリンク予測にも応用できる。 $n$  ノードの無向グラフ  $G$  が与えられた。ノードの集合を  $\mathcal{V} = \{1, \dots, n\}$  であらわし、すべてのノードのペアを  $\mathcal{P} \equiv \{(i, j) \mid 1 \leq i < j \leq n\}$  であらわすことにする。 $\mathcal{P}$  の要素数は  $n(n-1)/2$  である。リンクの集合  $\mathcal{E}$  は  $\mathcal{P}$  の部分集合であり、その補集合  $\bar{\mathcal{E}} \equiv \mathcal{P} \setminus \mathcal{E}$  にあるペアはリンクしていない。リンク情報を

$$y_i^j = \begin{cases} +1 & \text{if } (i, j) \in \mathcal{E}, \\ -1 & \text{if } (i, j) \notin \mathcal{E} \end{cases}$$

であらわす。それぞれのノードには特徴ベクトル  $x \in \mathbb{R}^d$  が与えられているとする。 $\mathcal{P}$  の部分集合  $\mathcal{P}_{\text{tra}}$  にあるノードペアはリンクの有無が既知であるとする。つまり、 $\mathcal{P}_{\text{tra}}$  にあるノードペアはクラスラベルが分かっている。リンク予測は、 $\mathcal{P}_{\text{tra}}$  から、残りの集合  $\mathcal{P} \setminus \mathcal{P}_{\text{tra}}$  にあるノードペアのリンクの有無を予測する問題といえることができる。

一つの標準的な方法は大域的モデル<sup>17),29),31)</sup>を使うことである。それは、ネットワーク全体で一つのモデルを構築するものである。しかし、大域的モデルは関連のない情報を取り込みやすい。この問題を解決するために、Bleakley ら<sup>4)</sup> は局所的モデルを使うことを提案している。局所的モデルとは一つのノードに対して一つのモデルを作るものである。この方法では、一つの局所的モデルは局所的な情報のみから訓練されるので、無関係な情報の影響を受けにくい。その代わりに、訓練に必要な情報は限られるので、この方法から得られる予測値の精度は信頼しがたい。

これに対して、我々は多タスク学習を使うことを提案する。ある指定された局所的モデルを訓練するために、リンクしているほかのノードに対応する局所的モデルを利用する。これは複雑ネットワークの性質<sup>24)</sup>を生かしたものになっている。一方、Bleakley らの方法<sup>4)</sup>は IL-SVM に対応する。

## 5. 実 験

提案法を多クラス識別問題、順序回帰、リンク予測に適用した結果をそれぞれ表 1, 3, 4 に示す。また、計算時間の比較を表 2 に示す。詳細な実験条件は文献<sup>18)</sup>を参照されたい。多くの場合、提案法がほかの方法より性能が上回っている。MTL-SVM(local/full) と MTL-SVM(local/network) の間には統計的有意な差はあまり見られなかったが、タスクネットワークを用いることにより計算時間が激減していることが分かる。協調フィルタリングへの応用でも提案法が優れた性能を示したが、その実験結果は紙面の制約から割愛する。

## 参 考 文 献

- 1) Amit, Y., Fink, M., Srebro, N. and Ullman, S.: Uncovering shared structures in multiclass classification, *Proceedings of the 24th International Conference on Machine Learning*, pp.17–24 (2007).
- 2) Bakker, B. and Heskes, T.: Task clustering and gating for Bayesian multitask learning, *Journal of Machine Learning Research*, Vol.4, pp.83–99 (2003).
- 3) Baxter, J.: A Model of Inductive Bias Learning, *Journal of Artificial Intelligence Research*, Vol.12, pp.149–198 (2000).
- 4) Bleakley, K., Biau, G. and Vert, J.-P.: Supervised Reconstruction of Biological Networks with Local Models, *Bioinformatics*, Vol.23, No.13, pp.i57–i65 (2007).
- 5) Bonilla, E.V., Agakov, F.V. and Williams, C.K.I.: Kernel Multi-task Learning using Task-specific Features, *Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics*, pp.43–50 (2007).
- 6) Boyd, S. and Vandenberghe, L.: *Convex Optimization*, Cambridge University Press (2004).
- 7) Caruana, R.: Multitask Learning, *Machine Learning*, Vol. 28, No. 1, pp. 41–75 (1997).
- 8) Deerwester, S., Dumais, S.T., Furnas, G.W., Landauer, T.K. and Harshman, R.: Indexing by latent semantic analysis., *Journal of the American Society for Information Science*, Vol.41, No.6, pp.391–407 (1990).
- 9) Evgeniou, T., Micchelli, C.A. and Pontil, M.: Learning Multiple Tasks with Kernel Methods, *Journal of Machine Learning Research*, Vol.6, pp.615–637 (2005).

- 10) Evgeniou, T. and Pontil, M.: Regularized Multitask Learning, *Proceedings of the 17th SIGKDD Conference on Knowledge Discovery and Data Mining*, pp.109–117 (2004).
- 11) Gärtner, T.: A Survey of Kernels for Structured Data, *SIGKDD Explorations*, Vol.5, No.1, pp.S268–S275 (2003).
- 12) Gärtner, T., Flach, P. and Wrobel, S.: On Graph Kernels: Hardness Results and Efficient Alternatives, *Proceedings of the Sixteenth Annual Conference on Computational Learning Theory*, pp.129–143 (2003).
- 13) Haussler, D.: Convolution kernels on discrete structures, Technical Report UCSC-CRL-99-10, UC Santa Cruz (1999).
- 14) Jaakkola, T. and Haussler, D.: Exploiting Generative Models in Discriminative Classifiers, *Advances in Neural Information Processing Systems 11* (Kearns, M.S., Solla, S.A. and Cohn, D.A., eds.), Cambridge, MA., MIT Press, pp.487–493 (1999).
- 15) Kashima, H. and Koyanagi, T.: Kernels for Semi-Structured Data, *Proceedings of the Nineteenth International Conference on Machine Learning*, pp.291–298 (2002).
- 16) Kashima, H., Tsuda, K. and Inokuchi, A.: Marginalized Kernels between Labeled Graphs, *Proceedings of the Twentieth International Conference on Machine Learning*, pp.321–328 (2003).
- 17) Kato, T., Tsuda, K. and Asai, K.: Selective integration of multiple biological data for supervised network inference, *Bioinformatics*, Vol.21, pp.2488–2495 (2005).
- 18) Kato, T., Kashima, H., Sugiyama, M. and Asai, K.: Conic Programming for Multi-Task Learning, *IEEE Transactions on Knowledge and Data Engineering*.
- 19) Kondor, R.I. and Lafferty, J.: Diffusion Kernels on Graphs and Other Discrete Input Spaces, *Proceedings of the Nineteenth International Conference on Machine Learning*, pp.315–322 (2002).
- 20) Lawrence, N.D. and Platt, J.C.: Learning to learn with the informative vector machine, *Proceedings of the Twenty First International Conference on Machine Learning*, pp.512–519 (2004).
- 21) Leslie, C., Eskin, E. and Noble, W.S.: The Spectrum Kernel: A String Kernel for SVM Protein Classification, *Proceedings of the Pacific Symposium on Biocomputing*, pp.566–575 (2002).
- 22) Lodhi, H., Saunders, C., Shawe-Taylor, J., Cristianini, N. and Watkins, C.: Text Classification Using String Kernels, *Journal of Machine Learning Research*, Vol.2, pp.419–444 (2002).
- 23) Micchelli, C.A. and Pontil, M.: Kernels for Multi-task Learning, *Advances in Neural Information Processing Systems 17*, Cambridge, MA, MIT Press, pp.921–928 (2005).
- 24) Milo, R., Shen-Orr, S., Itzkovitz, S., Kashtan, N., Chklovskii, D. and Alon, U.:

表 1 タンパク質スーパーファミリー識別．タスクネットワークは無作為に生成した木である．太字は最高性能，下線は最高性能と統計的有意差がないことを示す．

Class	IL-SVM	MTL-SVM (global/full)	MTL-SVM (global/network)	MTL-SVM (local/full)	MTL-SVM (local/network)
1	<u>0.893</u> (0.033)	<b>0.896</b> (0.032)	<u>0.885</u> (0.053)	0.887 (0.040)	0.884 (0.043)
2	0.803 (0.066)	<u>0.827</u> (0.040)	0.820 (0.043)	<b>0.833</b> (0.035)	<u>0.831</u> (0.038)
3	0.916 (0.013)	0.919 (0.016)	0.918 (0.018)	<b>0.926</b> (0.012)	<u>0.925</u> (0.012)
4	0.702 (0.103)	<u>0.762</u> (0.046)	<u>0.769</u> (0.041)	<b>0.776</b> (0.055)	<u>0.771</u> (0.067)
5	0.938 (0.030)	<b>0.948</b> (0.024)	<u>0.947</u> (0.027)	0.942 (0.023)	<u>0.942</u> (0.026)
6	0.755 (0.085)	0.763 (0.057)	<u>0.782</u> (0.041)	<b>0.786</b> (0.048)	<u>0.785</u> (0.061)
7	0.591 (0.061)	0.612 (0.067)	<u>0.616</u> (0.073)	<b>0.629</b> (0.055)	<u>0.627</u> (0.060)
ave	0.800 (0.024)	0.818 (0.019)	0.819 (0.021)	<b>0.825</b> (0.017)	<u>0.823</u> (0.020)

Network Motifs: Simple Building Blocks of Complex Networks, *Science*, Vol.298, pp.824–827 (2002).

- 25) Shashua, A. and Levin, A.: Ranking with large margin principle: two approaches, *Advances in Neural Information Processing Systems 15*, Cambridge, MA, MIT Press, pp.937–944 (2003).
- 26) Thrun, S. and Pratt, L.: *Learning to Learn*, Springer (1997).
- 27) Tsuda, K. and Noble, W.S.: Learning kernels from biological networks by maximizing entropy, *Bioinformatics*, Vol.20, No.Suppl. 1, pp.i326–i333 (2004).
- 28) Vapnik, V.N.: *Statistical Learning Theory*, Wiley, New York (1998).
- 29) Vert, J.-P. and Yamanishi, Y.: Supervised graph inference, *Advances in Neural Information Processing Systems 17*, Cambridge, MA, MIT Press (2005).
- 30) Xue, Y., Liao, X., Carin, L. and Krishnapuram, B.: Multi-Task Learning for Classification with Dirichlet Process Priors, *Journal of Machine Learning Research*, Vol.8, pp.35–63 (2007).
- 31) Yamanishi, Y., Vert, J.P. and Kanehisa, M.: Supervised enzyme network inference from the integration of genomic data and chemical information, *Bioinformatics*, Vol.21 Suppl.1, pp.i468–i477 (2005).
- 32) Yu, K., Tresp, V. and Schwaighofer, A.: Learning Gaussian Processes from Multiple Tasks, *Proceedings of the 22nd International Conference on Machine Learning*, pp.1012–1019 (2005).
- 33) Zhu, X., Kandola, J., Ghahramani, Z. and Lafferty, J.: Nonparametric transforms of graph kernels for semi-supervised learning, *Advances in Neural Information Processing Systems 17*, Cambridge, MA, MIT Press, pp.1641–1648 (2004).

表 2 計算時間．

Class	IL-SVM	MTL-SVM (global/full)	MTL-SVM (global/network)	MTL-SVM (local/full)	MTL-SVM (local/network)
Time (sec)	0.069 (0.006)	0.269 (0.007)	0.288 (0.016)	<b>2.339</b> (0.089)	0.884 (0.057)

表 3 順序回帰における正解率．太字は最高性能，下線は最高性能と統計的有意差がないことを示す．

Dataset	SVOR	IL-SVM	MTL-SVM (global/full)	MTL-SVM (global/network)	MTL-SVM (local/full)	MTL-SVM (local/network)
abalone	0.965 (0.016)	0.965 (0.016)	0.965 (0.016)	0.965 (0.016)	<b>0.972</b> (0.012)	0.966 (0.016)
bodyfat	0.957 (0.016)	0.958 (0.013)	0.958 (0.016)	0.958 (0.015)	<b>0.962</b> (0.013)	<b>0.962</b> (0.013)
cadata	0.974 (0.010)	0.974 (0.010)	0.974 (0.010)	0.974 (0.010)	<b>0.975</b> (0.009)	<b>0.975</b> (0.010)
housing	0.969 (0.012)	0.970 (0.012)	0.969 (0.012)	0.969 (0.012)	<b>0.973</b> (0.011)	<u>0.972</u> (0.011)
mg	0.969 (0.010)	0.969 (0.010)	0.969 (0.010)	0.969 (0.010)	<u>0.969</u> (0.009)	<b>0.970</b> (0.010)
mpg	0.964 (0.012)	0.966 (0.012)	0.963 (0.011)	0.963 (0.011)	<b>0.969</b> (0.011)	<b>0.969</b> (0.011)

表 4 酵素ネットワーク予測における ROC カーブの AUC．太字は最高性能，下線は最高性能と統計的有意差がないことを示す．

Dataset	IL-SVM	MTL-SVM (global/full)	MTL-SVM (global/network)	MTL-SVM (local/full)	MTL-SVM (local/network)
ady	0.733 (0.113)	0.744 (0.122)	0.746 (0.120)	n/a	<b>0.752</b> (0.118)
blast	0.786 (0.102)	0.792 (0.104)	<u>0.800</u> (0.097)	n/a	<b>0.806</b> (0.098)
diff	0.620 (0.130)	0.630 (0.124)	0.639 (0.122)	n/a	<b>0.654</b> (0.109)
expr	0.630 (0.104)	0.635 (0.109)	<u>0.636</u> (0.107)	n/a	<b>0.645</b> (0.100)
fft	0.652 (0.111)	0.663 (0.112)	<u>0.668</u> (0.109)	n/a	<b>0.675</b> (0.103)
lin_int	0.588 (0.111)	<u>0.609</u> (0.129)	<u>0.611</u> (0.127)	n/a	<b>0.614</b> (0.119)
pfam_hmm	0.740 (0.123)	0.748 (0.123)	0.749 (0.125)	n/a	<b>0.760</b> (0.124)
sw	0.732 (0.126)	0.731 (0.135)	0.725 (0.146)	n/a	<b>0.753</b> (0.122)