

MULTITASK LEARNING 34

CHAPTER CONTENTS

Task Similarity Regularization	391
Formulation	391
Analytic Solution	392
Efficient Computation for Many Tasks	393
Multidimensional Function Learning	394
Formulation	394
Efficient Analytic Solution	397
Matrix Regularization	397
Parameter Matrix Regularization	397
Proximal Gradient for Trace Norm Regularization	400

When solving multiple related learning tasks, solving them together simultaneously by sharing information is expected to give a better solution than solving them independently. This is the basic idea of *multitask learning* [24]. In this chapter, practically useful multitask learning methods are introduced.

Let us consider T tasks indexed by $t = 1, \dots, T$ and assume that each input-output paired training sample (\mathbf{x}_i, y_i) is accompanied with task index t_i :

$$\{(\mathbf{x}_i, y_i, t_i) \mid t_i \in \{1, \dots, T\}\}_{i=1}^n.$$

34.1 TASK SIMILARITY REGULARIZATION

In this section, a multitask learning method with *task similarity regularization* [40] is introduced.

34.1.1 FORMULATION

For the t th learning task, let us employ a linear-in-parameter model,

$$\sum_{j=1}^b \theta_{t,j} \phi_j(\mathbf{x}) = \boldsymbol{\theta}_t^\top \boldsymbol{\phi}(\mathbf{x}),$$

where $\theta_t = (\theta_{t,1}, \dots, \theta_{t,b})^\top$ is the parameter vector for the t th task and the basis functions $\phi(\mathbf{x})$ are *common* to all tasks.

The idea of task similarity regularization is to impose the parameters $\theta_1, \dots, \theta_T$ to take similar values and learn all parameters,

$$\theta = (\theta_1^\top, \dots, \theta_T^\top)^\top \in \mathbb{R}^{bT},$$

simultaneously. Let us employ the ℓ_2 -regularized LS (see Section 23.2) and learn θ so that the following $J(\theta)$ is minimized:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (y_i - \phi(\mathbf{x}_i)^\top \theta_i)^2 + \frac{1}{2} \sum_{t=1}^T \lambda_t \|\theta_t\|^2 + \frac{1}{4} \sum_{t,t'=1}^T \gamma_{t,t'} \|\theta_t - \theta_{t'}\|^2,$$

where $\lambda_t \geq 0$ is the ℓ_2 -regularization parameter for the t th task, and $\gamma_{t,t'} \geq 0$ is the similarity between the t th task and the t' th task. If $\gamma_{t,t'} = 0$ for all $t, t' = 1, \dots, T$, the third term in $J(\theta)$ disappears. Then there is no interaction between tasks and thus this corresponds to merely learning T tasks separately. On the other hand, if $\gamma_{t,t'} > 0$, θ_t and $\theta_{t'}$ are imposed to be close to each other and thus information can be implicitly shared between the t th task and the t' th task. If all $\gamma_{t,t'}$ are large enough, all solutions are imposed to be the same and thus a single solution that is common to all tasks is obtained from all training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$.

34.1.2 ANALYTIC SOLUTION

$J(\theta)$ can be compactly expressed as

$$J(\theta) = \frac{1}{2} \|\mathbf{y} - \Psi\theta\|^2 + \frac{1}{2} \theta^\top (\mathbf{C} \otimes \mathbf{I}_b) \theta, \quad (34.1)$$

where

$$\begin{aligned} \mathbf{y} &= (y_1, \dots, y_n)^\top \in \mathbb{R}^n, \\ \Psi &= (\psi_{t_1}(\mathbf{x}_1), \dots, \psi_{t_n}(\mathbf{x}_n))^\top \in \mathbb{R}^{n \times bT}, \\ \psi_t(\mathbf{x}) &= (\mathbf{0}_{b(t-1)}^\top, \phi(\mathbf{x})^\top, \mathbf{0}_{b(T-t)}^\top)^\top \in \mathbb{R}^{bT}. \end{aligned}$$

\mathbf{C} is the $T \times T$ matrix defined as

$$C_{t,t'} = \begin{cases} \lambda_t + \sum_{t''=1}^T \gamma_{t,t''} - \gamma_{t,t} & (t = t'), \\ -\gamma_{t,t'} & (t \neq t'), \end{cases}$$

and \otimes denotes the *Kronecker product*, i.e. for $\mathbf{E} \in \mathbb{R}^{m \times n}$ and $\mathbf{F} \in \mathbb{R}^{p \times q}$,

$$\mathbf{E} \otimes \mathbf{F} = \begin{pmatrix} E_{1,1}\mathbf{F} & \cdots & E_{1,n}\mathbf{F} \\ \vdots & \ddots & \vdots \\ E_{m,1}\mathbf{F} & \cdots & E_{m,n}\mathbf{F} \end{pmatrix} \in \mathbb{R}^{mp \times nq}.$$

Then the minimizer $\widehat{\theta}$ of $J(\theta)$ satisfies

$$(\Psi^\top \Psi + C \otimes I_b) \theta = \Psi^\top y, \quad (34.2)$$

and $\widehat{\theta}$ can be obtained analytically as

$$\widehat{\theta} = (\Psi^\top \Psi + C \otimes I_b)^{-1} \Psi^\top y. \quad (34.3)$$

34.1.3 EFFICIENT COMPUTATION FOR MANY TASKS

The size of matrix $(\Psi^\top \Psi + C \otimes I_b)$ is $bT \times bT$, and thus directly computing the solution by the above analytic form is not tractable if the number of tasks, T , is large. However, since the rank of matrix $\Psi^\top \Psi$ is at most n , the solution can be computed efficiently if $n < bT$.

More specifically, the solution $\widehat{\theta}_t^\top \phi(x)$ for the t th task can be expressed as

$$\widehat{\theta}_t^\top \phi(x) = \widehat{\theta}^\top \psi_t(x) = y^\top A^{-1} b_t. \quad (34.4)$$

Here, A and b_t are the $n \times n$ matrix and n -dimensional vector defined as

$$\begin{aligned} A_{i,i'} &= [\Psi(C^{-1} \otimes I_b) \Psi^\top + I_n]_{i,i'} \\ &= [C^{-1}]_{t_i, t_{i'}} \phi(x_i)^\top \phi(x_{i'}) + \begin{cases} 1 & (i = i'), \\ 0 & (i \neq i'), \end{cases} \\ b_{t,i} &= [\Psi(C^{-1} \otimes I_b) \psi_t(x)]_i \\ &= [C^{-1}]_{t_i, t_i} \phi(x_i)^\top \phi(x), \end{aligned}$$

where $[C^{-1}]_{t_i, t_{i'}}$ denotes the (t, t') th element of matrix C^{-1} . Since the size of matrix A and the dimensionality of vector b_t are independent of the number of tasks, T , Eq. (34.4) would be computationally more efficient than Eq. (34.3) if $n < bT$. Moreover, since A^{-1} is independent of task index t , it needs to be computed only once across all tasks. Note that the following formulas regarding the Kronecker product were utilized in the derivation of Eq. (34.4):

$$\begin{aligned} \Psi(\Psi^\top \Psi + C \otimes I_b)^{-1} &= (\Psi(C \otimes I_b)^{-1} \Psi^\top + I_n)^{-1} \Psi(C \otimes I_b)^{-1}, \\ (C \otimes I_b)^{-1} &= C^{-1} \otimes I_b. \end{aligned}$$

A MATLAB code for multitask LS is provided in Fig. 34.1, and its behavior is illustrated in Fig. 34.2. This shows that multitask classification outperforms single-task classification.

In the above multitask method, the task similarity $\gamma_{t,t'}$ was assumed to be known. When $\gamma_{t,t'}$ is unknown, it may be alternately learned as described in Fig. 34.3.

```

n=2; T=6; y=[ones(n/2,T); -ones(n/2,T)];
x=[randn(2,n,T); ones(1,n,T)]; r(1,1,:)=pi*[1:T]/T/10;
c= repmat(cos(r),[1 n/2]); x(1, :, :)=x(1, :, :)+[c -c];
s= repmat(sin(r),[1 n/2]); x(2, :, :)=x(2, :, :)+[s -s];
Ci=inv(-ones(T,T)+diag(T*ones(T,1)+0.01));
a= repmat([1:T], [n 1]); a=a(:); m=20; X=linspace(-4,4,m);
b= repmat(X,[m 1]); bt=b'; XX=[b(:)'; bt(:)'; ones(1,m^2)];
yAi=y(:)'*inv(Ci(a,a).*(x(:, :)'*x(:, :))+eye(n*T));

figure(1); clf; colormap([1 0.7 1; 0.7 1 1]);
for k=1:T
% Y=yAi*(repmat(Ci(a,k),[1 3]).*x(:, :)'*XX;
q=x(:, :, k); Y=((q*q'+0.01*eye(3))\ (q*y(:, k)))'*XX;
subplot(2,3,k); hold on; contourf(X,X,reshape(Y,m,m));
plot(x(1,y(:,k)==1,k),x(2,y(:,k)==1,k), 'bo');
plot(x(1,y(:,k)==-1,k),x(2,y(:,k)==-1,k), 'rx');
plot(99*sin(r(k))*[1 -1],99*cos(r(k))*[-1 1], 'k--');
axis([-4 4 -4 4]);
end

```

FIGURE 34.1

MATLAB code for multitask LS.

34.2 MULTIDIMENSIONAL FUNCTION LEARNING

In this section, the problem of *multidimensional function learning* is discussed, which can be regarded as a special case of multitask learning.

34.2.1 FORMULATION

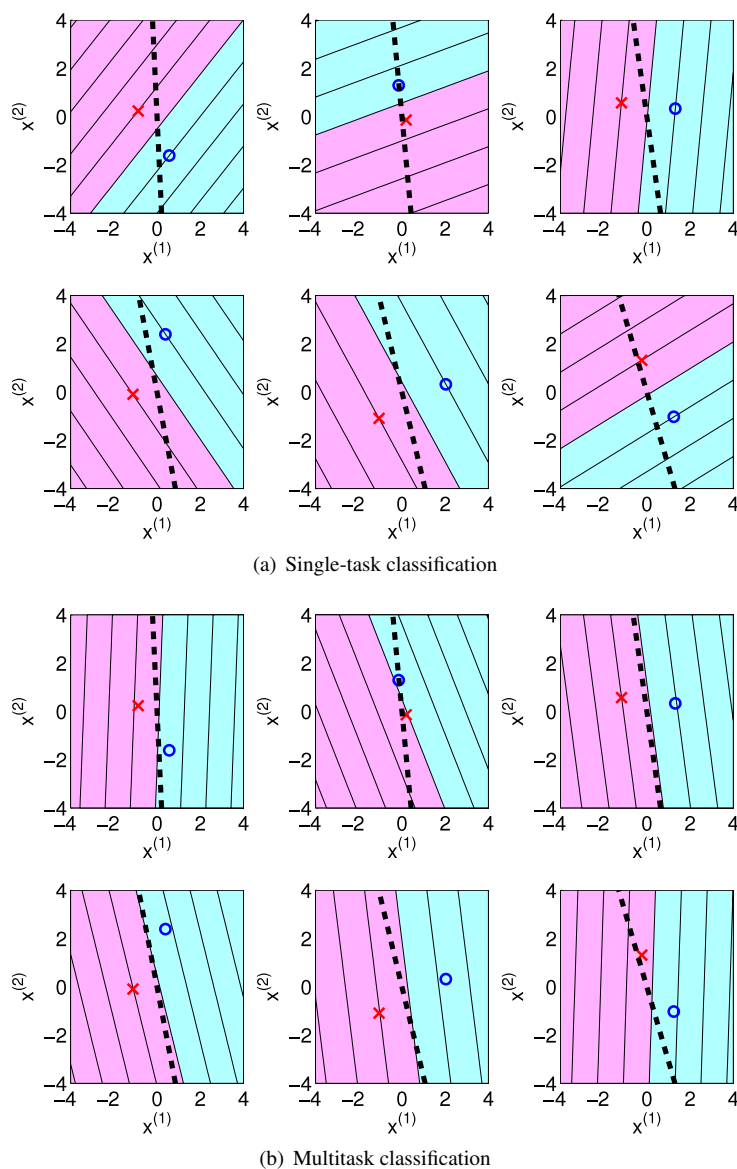
Let us consider the problem of learning a T -dimensional function,

$$\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_T(\mathbf{x}))^\top,$$

from input-output paired training samples:

$$\{(\mathbf{x}_i, \mathbf{y}_i) \mid \mathbf{x}_i \in \mathbb{R}^d, \mathbf{y}_i = (y_i^{(1)}, \dots, y_i^{(T)})^\top \in \mathbb{R}^T\}_{i=1}^n.$$

If each dimension of output \mathbf{y} is regarded as a task, the problem of multidimensional function learning can be regarded as multitask learning. The difference is that multidimensional function learning shares input points across all tasks, while input points are generally different in multitask learning (Fig. 34.4). Thus, the number of training samples in multidimensional function learning is actually nT in the context of multitask learning.

**FIGURE 34.2**

Examples of multitask LS. The dashed lines denote true decision boundaries and the contour lines denote learned results.

Multidimensional function learning for classification is specifically called *multilabel classification*, which can be regarded as a generalization of *multiclass*

1. Initialize task similarity, e.g. $\gamma_{t,t'} = \gamma > 0, \forall t, t'$.
2. Learn parameter θ based on the current task similarity $\gamma_{t,t'}$.
3. Update task similarity $\gamma_{t,t'}$ based on the similarity between θ_t and $\theta_{t'}$, e.g.

$$\gamma_{t,t'} = \eta \exp(-\kappa \|\theta_t - \theta_{t'}\|^2),$$

where $\eta \geq 0$ and $\kappa \geq 0$ are tuning parameters and may be determined by cross validation.

4. Iterate 2–3 until convergence.

FIGURE 34.3

Alternate learning of task similarity $\gamma_{t,t'}$ and solution θ .

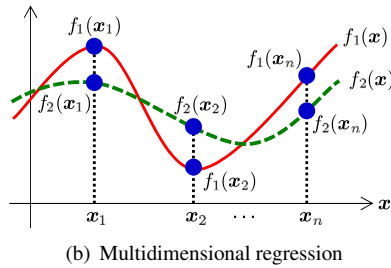
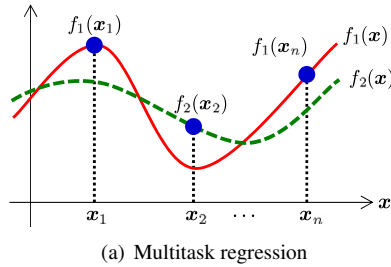


FIGURE 34.4

Multidimensional function learning.

classification (see Section 26.3). A pattern x belongs to one of the classes exclusively in multiclass classification, while x can belong to multiple classes simultaneously in multilabel classification. For example, an image can belong to the classes of “human,” “dog,” and “building” at the same time in image recognition, and a sound clip can belong to the classes of “conversation,” “noise,” and “background music” at the

same time in audio recognition. In multilabel classification, co-occurrence of multiple classes can be utilized by multitask learning. For example, an image containing a dog often contains a human at the same time.

However, in multidimensional function learning, the number of parameters is bT and the *actual* number of samples is nT , both of which are proportional to the number of tasks, T . Thus, both the solutions explained in Section 34.1.2 and Section 34.1.3 are not computationally efficient in multidimensional function learning.

34.2.2 EFFICIENT ANALYTIC SOLUTION

Let us arrange the parameter vectors $\theta_1, \dots, \theta_T$, the output vectors y_1, \dots, y_n , and the basis vectors $\phi(x_1), \dots, \phi(x_n)$ into matrices as

$$\begin{aligned}\Theta &= (\theta_1, \dots, \theta_T) \in \mathbb{R}^{b \times T}, \\ Y &= (y_1, \dots, y_n)^\top \in \mathbb{R}^{n \times T}, \\ \Phi &= (\phi(x_1), \dots, \phi(x_n))^\top \in \mathbb{R}^{n \times b}.\end{aligned}$$

Then, Eq. (34.2), which multitask solutions should satisfy, can be rewritten as

$$\Phi^\top \Phi \Theta + \Theta C = \Phi^\top Y.$$

This form is known as the *continuous Sylvester equation*, which often arises in control theory [94], and can be solved efficiently, as explained in Fig. 34.5.

A MATLAB code for multidimensional regression is provided in Fig. 34.6, and its behavior is illustrated in Fig. 34.7. This shows that multitask regression outperforms single-task regression.

34.3 MATRIX REGULARIZATION

The multitask learning method introduced in Section 34.1 explicitly used task similarity $\gamma_{t,t'}$ to control the amount of information sharing. In this section, another multitask learning approach is introduced, which does not involve task similarity.

34.3.1 PARAMETER MATRIX REGULARIZATION

The basic idea is to share information across multiple tasks by regularizing the *parameter matrix*:

$$\Theta = (\theta_1, \dots, \theta_T) \in \mathbb{R}^{b \times T}.$$

More specifically, for a squared loss function, the multitask learning criterion to be minimized with respect to Θ is given by

$$\frac{1}{2} \sum_{i=1}^n (y_i - \phi(x_i)^\top \theta_{t_i})^2 + \lambda R(\Theta),$$

where $R(\Theta)$ denotes some regularization functional for Θ .

The *continuous Sylvester equation* for some matrices $\mathbf{A} \in \mathbb{R}^{b \times b}$, $\mathbf{B} \in \mathbb{R}^{T \times T}$, and $\mathbf{Z} \in \mathbb{R}^{b \times T}$ with respect to $\mathbf{\Theta} \in \mathbb{R}^{b \times T}$ is given by

$$\mathbf{A}\mathbf{\Theta} + \mathbf{\Theta}\mathbf{B} = \mathbf{Z}.$$

The use of the *Kronecker product* and the *vectorization operator* (see Fig. 6.5) allows us to rewrite the above equation as

$$(\mathbf{I}_T \otimes \mathbf{A} + \mathbf{B} \otimes \mathbf{I}_b) \text{vec}(\mathbf{\Theta}) = \text{vec}(\mathbf{Z}),$$

where \mathbf{I}_T denotes the $T \times T$ identity matrix. This shows that the continuous Sylvester equation with respect to a $b \times T$ matrix can be seen as a linear equation with respect to a bT -dimensional vector, which is computationally expensive to solve naively. However, the continuous Sylvester equation can be solved more efficiently. For example, let u_1, \dots, u_b and $\mathbf{u}_1, \dots, \mathbf{u}_b$ be *eigenvalues* and *eigenvectors* of \mathbf{A} (see Fig. 6.2), and let v_1, \dots, v_T and $\mathbf{v}_1, \dots, \mathbf{v}_T$ be eigenvalues and eigenvectors of \mathbf{B} , respectively. Then, when $u_j + v_t \neq 0$ for all $j = 1, \dots, b$ and $t = 1, \dots, T$, the solution $\hat{\mathbf{\Theta}}$ of the above continuous Sylvester equation with respect to $\mathbf{\Theta}$ can be obtained analytically as

$$\hat{\mathbf{\Theta}} = (\mathbf{u}_1, \dots, \mathbf{u}_b) \mathbf{Q} (\mathbf{v}_1, \dots, \mathbf{v}_T)^\top,$$

where \mathbf{Q} is the $b \times T$ matrix defined as

$$Q_{j,t} = \frac{\mathbf{u}_j^\top \mathbf{Z} \mathbf{v}_t}{u_j + v_t}.$$

FIGURE 34.5

Continuous Sylvester equation.

If the squared *Frobenius norm* $\|\mathbf{\Theta}\|_{\text{Frob}}^2$ is used for regularization, no information is shared across different tasks since this is equivalent to the sum of ℓ_2 -norms of parameter vectors $\boldsymbol{\theta}_t$:

$$\|\mathbf{\Theta}\|_{\text{Frob}}^2 = \sum_{t=1}^T \sum_{j=1}^b \Theta_{t,j}^2 = \sum_{t=1}^T \|\boldsymbol{\theta}_t\|^2.$$

Thus, a more intricate norm should be used for matrix regularization.

For example, the *trace norm* (see Fig. 24.10) tends to produce a *low-rank* solution:


```

n=30; x=linspace(-3,3,n)'; pix=pi*x; T=3;
y=repmat(sin(pix)./(pix)+0.1*x,1,T)+0.1*repmat([1:T],n,1);
y=y+0.5*randn(n,T); N=1000; X=linspace(-3,3,N)'; piX=pi*X;
Y=repmat(sin(piX)./(piX)+0.1*X,1,T)+0.1*repmat([1:T],N,1);
G=10*ones(T,T); %G=zeros(T,T);
l=0.1; C=l*eye(T)+diag(sum(G))-G; hh=1; x2=x.^2;
k=exp(-(repmat(x2,1,n)+repmat(x2',n,1)-2*x*x')/hh);
K=exp(-(repmat(X.^2,1,n)+repmat(X.^2',N,1)-2*X*X')/hh);

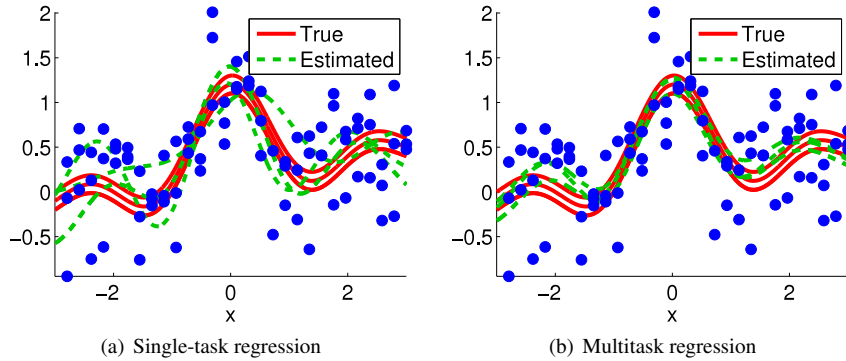
[U,u]=eig(k'*k); [V,v]=eig(C);
Q=U'*k'*y*V./(repmat(diag(u),1,T)+repmat(diag(v)',n,1));
S=U*Q*V'; F=K*S;

figure(1); clf; hold on; axis([-inf inf -inf inf])
plot(X,Y,'r-'); plot(X,F,'g--'); plot(x,y,'bo');

```

FIGURE 34.6

MATLAB code for multidimensional regression.

**FIGURE 34.7**

Examples of multidimensional regression.

$$\|\Theta\|_{\text{tr}} = \sum_{k=1}^{\min(b,T)} \sigma_k,$$

where σ_k is a *singular value* of Θ (see Fig. 22.2). Since a low-rank solution in multitask learning confines parameters of each task in a common subspace, information can be shared across different tasks [6]. In Section 23.1, the method

of *subspace-constrained LS* was introduced, which confines the LS solution in a given subspace. Multitask learning based on the trace norm can be regarded as automatically learning the subspace with the help of other learning tasks.

Another possibility of matrix regularization for multitask learning is to use the $\ell_{2,1}$ -norm $\|\Theta\|_{2,1}$:

$$\|\Theta\|_{2,1} = \sum_{t=1}^T \sqrt{\sum_{j=1}^b \Theta_{t,j}^2},$$

which tends to give a *row-wise* sparse solution (see Section 24.4.4). This means that $\theta_{1,j}, \dots, \theta_{T,j}$ tend to simultaneously vanish for several $j \in \{1, \dots, b\}$, and thus *feature selection* can be performed across different tasks, if linear basis function $\phi(\mathbf{x}) = \mathbf{x}$ is used [8] (see also Section 24.3).

Below, a practical optimization method for multitask LS with trace norm regularization is introduced.

34.3.2 PROXIMAL GRADIENT FOR TRACE NORM REGULARIZATION

Since the trace norm $\|\Theta\|_{\text{tr}}$ is a convex function, the global optimal solution to multitask LS with trace norm regularization can be obtained, for example, by the *proximal gradient method* described in Fig. 34.8.

More specifically, the proximal gradient method updates the parameter matrix Θ from some initial value as

$$\Theta \leftarrow \text{prox}(\Theta - \varepsilon \nabla L(\Theta)), \quad (34.5)$$

where $\varepsilon > 0$ is the step size. In Eq. (34.5), $L(\Theta)$ is the loss function for all training samples:

$$L(\Theta) = \frac{1}{2} \sum_{i=1}^n \left(y_i - \phi(\mathbf{x}_i)^\top \Theta \mathbf{e}_{t_i} \right)^2,$$

where \mathbf{e}_t denotes the T -dimensional vector with all zeros but the t th element being one, i.e. $\Theta \mathbf{e}_t = \theta_t$. $\nabla L(\Theta)$ is the gradient of L given by

$$\nabla L(\Theta) = \sum_{i=1}^n \left(\phi(\mathbf{x}_i)^\top \Theta \mathbf{e}_{t_i} - y_i \right) \phi(\mathbf{x}_i) \mathbf{e}_{t_i}^\top,$$

where the following matrix derivative formula is used:

$$\frac{\partial \phi(\mathbf{x}_i)^\top \Theta \mathbf{e}_{t_i}}{\partial \Theta} = \phi(\mathbf{x}_i) \mathbf{e}_{t_i}^\top.$$

The *proximal operator* $\text{prox}(\Theta)$ in Eq. (34.5) is given by

$$\text{prox}(\Theta) = \underset{U}{\text{argmin}} \left(\|U\|_{\text{tr}} + \frac{1}{2\varepsilon\lambda} \|U - \Theta\|_{\text{Frob}}^2 \right),$$

Let us consider the following optimization problem:

$$\min_{\theta} L(\theta) + R(\theta),$$

where $L(\theta)$ is a convex differentiable function and $R(\theta)$ is a closed convex function. The *proximal gradient method* finds the minimizer of $L(\theta) + R(\theta)$ by updating θ from some initial value as

$$\theta \leftarrow \operatorname{argmin}_{\mathbf{u}} \left(R(\mathbf{u}) + L(\theta) + \nabla L(\theta)^\top (\mathbf{u} - \theta) + \frac{1}{2\varepsilon} \|\mathbf{u} - \theta\|^2 \right),$$

where $\nabla L(\theta)$ is the gradient of L with respect to θ and $\varepsilon > 0$ corresponds to the step size. If the above minimization with respect to \mathbf{u} can be solved efficiently (e.g. analytically), the proximal gradient method is computationally efficient. The above update rule can be equivalently expressed as

$$\theta \leftarrow \operatorname{prox}_{\varepsilon R}(\theta - \varepsilon \nabla L(\theta)),$$

where $\operatorname{prox}_R(\theta)$ is called the *proximal operator* defined as

$$\operatorname{prox}_R(\theta) = \operatorname{argmin}_{\mathbf{u}} \left(R(\mathbf{u}) + \frac{1}{2} \|\mathbf{u} - \theta\|^2 \right).$$

This implies that the proximal gradient method can be regarded as a generalization of the projected gradient method. Indeed, when $R(\theta)$ is the *indicator function* for some set \mathcal{S} ,

$$R(\theta) = \begin{cases} 0 & (\theta \in \mathcal{S}), \\ \infty & (\theta \notin \mathcal{S}), \end{cases}$$

the proximal operator is reduced to an ordinary projection operator and thus the proximal gradient method is a projected gradient method.

FIGURE 34.8

Proximal gradient method.

where $\|\Theta\|_{\text{Frob}}$ denotes the *Frobenius norm*:

$$\|\Theta\|_{\text{Frob}} = \sqrt{\sum_{k_1, k_2=1}^{b_1, b_2} \Theta_{k_1, k_2}^2} = \sqrt{\operatorname{tr}(\Theta \Theta^\top)}.$$

```

n=2; T=6; y=[ones(n/2,T); -ones(n/2,T)];
x=[randn(2,n,T); ones(1,n,T)]; r(1,1,:)=pi*[1:T]/T/10;
c= repmat(cos(r),[1 n/2]); x(1, :, :)=x(1, :, :)+[c -c];
s= repmat(sin(r),[1 n/2]); x(2, :, :)=x(2, :, :)+[s -s];
t0=randn(3,T); e=0.1; l=4;
for o=1:1000
    for k=1:T
        gt(:,k)=x(:, :, k)*(x(:, :, k)'*t0(:,k)-y(:,k));
    end
    [U,S,V]=svd(t0-e*gt, 'econ');
    S=diag(max(0,diag(S)-e*l)); t=U*S*V';
    if norm(t-t0)<0.001, break, end
    t0=t;
end

figure(1); clf; colormap([1 0.7 1; 0.7 1 1]);
m=20; X=linspace(-4,4,m); b=repmat(X,[m 1]); bt=b';
for k=1:T
    Y=t(:,k)'*[b(:)'; bt(:)'; ones(1,m^2)];
    subplot(2,3,k); hold on; contourf(X,X,reshape(Y,m,m));
    plot(x(1,y(:,k)==1,k),x(2,y(:,k)==1,k), 'bo');
    plot(x(1,y(:,k)==-1,k),x(2,y(:,k)==-1,k), 'rx');
    plot(99*sin(r(k))*[1 -1],99*cos(r(k))*[-1 1], 'k--');
    axis([-4 4 -4 4]);
end

```

FIGURE 34.9

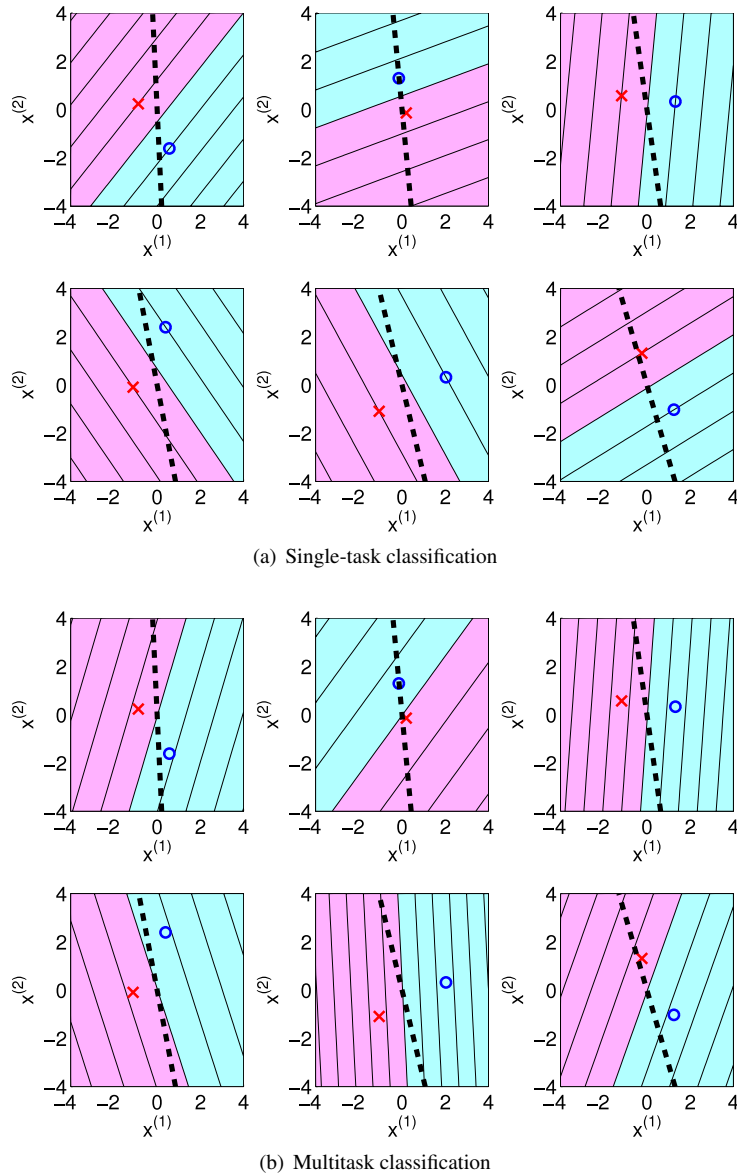
MATLAB code for multitask learning with trace norm regularization.

A notable fact is that the above proximal operator can be analytically expressed as follows [78]:

$$\text{prox}(\Theta) = \sum_{k=1}^{\min(b_1, b_2)} \max(0, \sigma_k - \varepsilon \lambda) \psi_k \phi_k^\top,$$

where ψ_k , ϕ_k , and σ_k are a left singular vector, a right singular vector, and a singular value of Θ , respectively:

$$\Theta = \sum_{k=1}^{\min(b_1, b_2)} \sigma_k \psi_k \phi_k^\top.$$

**FIGURE 34.10**

Examples of multitask LS with trace norm regularization. The data set is the same as [Fig. 34.2](#). The dashed lines denote true decision boundaries and the contour lines denote learned results.

This means that singular values less than $\varepsilon\lambda$ are rounded off to zero and other singular values are reduced by $\varepsilon\lambda$. This operation is called *soft thresholding*, which can be computed very efficiently.

A MATLAB code for multitask LS with trace norm regularization is provided in Fig. 34.9, and its behavior using the same data set as in Fig. 34.2 is illustrated in Fig. 34.10. This shows that multitask classification with trace norm regularization works reasonably well.