

STRUCTURED
CLASSIFICATION

29

CHAPTER CONTENTS

Sequence Classification	329
Probabilistic Classification for Sequences	330
Conditional Random Field	330
MLE	333
Recursive Computation	333
Prediction for New Sample	336
Deterministic Classification for Sequences	337

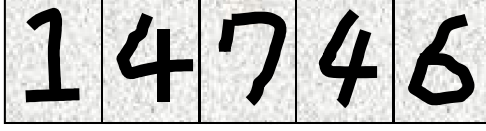
When classifying a set of input patterns having certain data structure such as a sequence of letters and a parsing tree of a sentence, the classification performance is expected to be improved if such structural information is utilized. In this chapter, taking a sequence of letters as an example of such structure, methods of structured classification are introduced.

29.1 SEQUENCE CLASSIFICATION

Let us consider the problem of classifying a sequence of m patterns where each pattern belongs to one of the c classes. Let $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ be sequences of m patterns and m classes, and let $\mathbf{x}^{(k)}$ and $y^{(k)}$ be their k th elements:

$$\begin{aligned}\bar{\mathbf{x}} &= (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}), \\ \bar{\mathbf{y}} &= (y^{(1)}, \dots, y^{(m)}).\end{aligned}$$

For example, when recognizing a sequence of five hand-written digits (see Fig. 29.1), $m = 5$ and $c = 10$. If such a sequence is decomposed into five digits, ordinary classification algorithms introduced in the previous chapters can be used for classifying each digit. However, this decomposition approach loses useful contextual information that can drastically improve the classification accuracy, for example, the same number may rarely appear consecutively and some number appears frequently after a certain number. On the other hand, if the entire sequence of digits is directly recognized, the number of classes grows exponentially with respect to the length of

**FIGURE 29.1**

Classification of sequence of hand-written digits.

the sequence. In the example of Fig. 29.1, $(m, c) = (5, 10)$ and thus the number of classes is 10^5 , which is hard to handle.

In this chapter, intermediate approaches are introduced, which take into account some contextual information with the computational complexity kept moderately. Note that, although only sequence data are considered, methods introduced in this chapter can be applied to more general structured data.

29.2 PROBABILISTIC CLASSIFICATION FOR SEQUENCES

In this section, the *logistic model* introduced in Section 28.1 is extended to the *conditional random field* [65] for sequence classification.

29.2.1 CONDITIONAL RANDOM FIELD

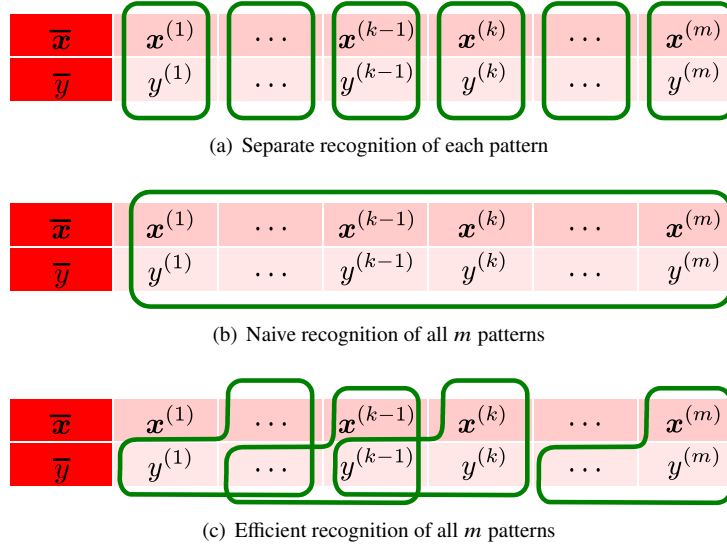
To introduce the conditional random field, let us first consider the logistic model applied to sequence classification, i.e. each pattern $\mathbf{x}^{(k)}$ is independently classified into one of the c classes (Fig. 29.2(a)):

$$q(y|\mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\theta}_y^\top \boldsymbol{\phi}(\mathbf{x}))}{\sum_{y'=1}^c \exp(\boldsymbol{\theta}_{y'}^\top \boldsymbol{\phi}(\mathbf{x}))},$$

where $\boldsymbol{\phi}(\mathbf{x}) \in \mathbb{R}^b$ is the vector of basis functions for a single pattern \mathbf{x} , $\boldsymbol{\theta}_y \in \mathbb{R}^b$ is the parameter vector for class y , and

$$\boldsymbol{\theta} = (\boldsymbol{\theta}_1^\top, \dots, \boldsymbol{\theta}_c^\top)^\top \in \mathbb{R}^{bc}$$

is the vector of all parameters. However, this decomposition approach cannot utilize the contextual information contained in sequence data.

**FIGURE 29.2**

Sequence classification.

On the other hand, if the entire sequence is classified at once, a classification problem with $\bar{c} = c^m$ classes needs to be solved (Fig. 29.2(b)):

$$q(\bar{y}|\bar{\mathbf{x}}; \bar{\boldsymbol{\theta}}) = \frac{\exp\left(\bar{\boldsymbol{\theta}}_{\bar{y}}^{\top} \bar{\boldsymbol{\phi}}(\bar{\mathbf{x}})\right)}{\sum_{y^{(1)}, \dots, y^{(m)}=1}^c \exp\left(\bar{\boldsymbol{\theta}}_{\bar{y}}^{\top} \bar{\boldsymbol{\phi}}(\bar{\mathbf{x}})\right)}, \quad (29.1)$$

where

$$\bar{\boldsymbol{\phi}}(\bar{\mathbf{x}}) = (\boldsymbol{\phi}(\mathbf{x}^{(1)})^{\top}, \dots, \boldsymbol{\phi}(\mathbf{x}^{(m)})^{\top})^{\top} \in \mathbb{R}^{bm}$$

is the vector of basis functions for pattern sequence $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$,

$$\bar{\boldsymbol{\theta}}_{\bar{y}} = (\boldsymbol{\theta}_{y^{(1)}}^{\top}, \dots, \boldsymbol{\theta}_{y^{(m)}}^{\top})^{\top} \in \mathbb{R}^{bm}$$

is the parameter vector for class $\bar{y} = (y^{(1)}, \dots, y^{(m)})$, and

$$\bar{\boldsymbol{\theta}} = (\bar{\boldsymbol{\theta}}_1^{\top}, \dots, \bar{\boldsymbol{\theta}}_{\bar{c}}^{\top})^{\top} \in \mathbb{R}^{bm\bar{c}}$$

is the vector of all parameters. However, since the number of classes $\bar{c} = c^m$ grows exponentially with respect to the sequence length m , handling the above logistic model may be intractable.

To avoid the exponential growth in the number of classes, let us assume that $y^{(k)}$ is determined only by $\mathbf{x}^{(k)}$ and $y^{(k-1)}$ when classifying the entire sequence $\bar{\mathbf{x}}$. Note that this approach is different from merely classifying two consecutive patterns into one of the c^2 classes, because of the overlap (see Fig. 29.2(c)).

To implement the above idea, let us consider the following model, called the *conditional random field*:

$$q(\bar{y}|\bar{\mathbf{x}}; \zeta) = \frac{\exp(\zeta^\top \boldsymbol{\varphi}(\bar{\mathbf{x}}, \bar{y}))}{\sum_{y^{(1)}, \dots, y^{(m)}=1}^c \exp(\zeta^\top \boldsymbol{\varphi}(\bar{\mathbf{x}}, \bar{y}'))}, \quad (29.2)$$

where the basis function $\boldsymbol{\varphi}(\bar{\mathbf{x}}, \bar{y})$ depend not only on input $\bar{\mathbf{x}}$ but also on output \bar{y} . Conditional random field (29.2) is reduced to the logistic model for sequence $\bar{\mathbf{x}}$ given by Eq. (29.1), if

$$\zeta = \bar{\boldsymbol{\theta}} \quad \text{and} \quad \boldsymbol{\varphi}(\bar{\mathbf{x}}, \bar{y}) = \mathbf{e}_{\bar{y}}^{(\bar{c})} \otimes \bar{\boldsymbol{\phi}}(\bar{\mathbf{x}}),$$

where $\mathbf{e}_{\bar{y}}^{(\bar{c})}$ is the \bar{c} -dimensional indicator vector of class \bar{y} (i.e. the element corresponding class \bar{y} is 1 and all other elements are 0), and \otimes denotes the *Kronecker product* (see Fig. 6.5):

$$\begin{aligned} \text{For } \mathbf{f} &= (f_1, \dots, f_u)^\top \in \mathbb{R}^u, \quad \mathbf{g} = (g_1, \dots, g_v)^\top \in \mathbb{R}^v, \\ \mathbf{f} \otimes \mathbf{g} &= (f_1 g_1, f_1 g_2, \dots, f_1 g_v, f_2 g_1, f_2 g_2, \dots, f_2 g_v, \dots, \\ &\quad f_u g_1, f_u g_2, \dots, f_u g_v)^\top \in \mathbb{R}^{uv}. \end{aligned}$$

Thus, naively using conditional random field (29.2) is just an overly redundant expression of Eq. (29.1).

Let us simplify conditional random field (29.2) by taking into account the structure described in Fig. 29.2(c). More specifically, basis function $\boldsymbol{\varphi}(\bar{\mathbf{x}}, \bar{y})$ is defined as the sum of basis functions depending only on two consecutive patterns:

$$\boldsymbol{\varphi}(\bar{\mathbf{x}}, \bar{y}) = \sum_{k=1}^m \boldsymbol{\varphi}(\mathbf{x}^{(k)}, y^{(k)}, y^{(k-1)}), \quad (29.3)$$

where $y^{(0)} = y^{(1)}$. As two-pattern basis function $\boldsymbol{\varphi}(\mathbf{x}^{(k)}, y^{(k)}, y^{(k-1)})$, for example,

$$\boldsymbol{\varphi}(\mathbf{x}^{(k)}, y^{(k)}, y^{(k-1)}) = \begin{pmatrix} \mathbf{e}_{y^{(k)}}^{(c)} \otimes \boldsymbol{\phi}(\mathbf{x}^{(k)}) \\ \mathbf{e}_{y^{(k)}}^{(c)} \otimes \mathbf{e}_{y^{(k-1)}}^{(c)} \end{pmatrix} \in \mathbb{R}^{cb+c^2}, \quad (29.4)$$

may be used. Then the dimensionality of parameter vector ζ is $cb + c^2$, which is independent of the sequence length m .

1. Initialize ζ .
2. Choose the i th training sample randomly $(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}_i)$.
3. Update parameter ζ to go up the gradient:

$$\zeta \leftarrow \zeta + \varepsilon \left(\varphi_i(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}_i) - \frac{\sum_{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m_i)}=1}^c \exp(\zeta^\top \varphi_i(\bar{\mathbf{x}}_i, \bar{\mathbf{y}})) \varphi_i(\bar{\mathbf{x}}_i, \bar{\mathbf{y}})}{\sum_{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m_i)}=1}^c \exp(\zeta^\top \varphi_i(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}))} \right),$$

where $\varepsilon > 0$ is the step size and $\bar{\mathbf{y}} = (y^{(1)}, \dots, y^{(m_i)})$.

4. Iterate 2 and 3 until convergence.

FIGURE 29.3

Stochastic gradient algorithm for conditional random field.

29.2.2 MLE

Suppose that labeled sequences,

$$\{(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}_i) \mid \bar{\mathbf{x}}_i = (\mathbf{x}_i^{(1)}, \dots, \mathbf{x}_i^{(m_i)}), \bar{\mathbf{y}}_i = (\mathbf{y}_i^{(1)}, \dots, \mathbf{y}_i^{(m_i)})\}_{i=1}^n,$$

are provided as training samples, where the length m_i of sequence $\bar{\mathbf{x}}_i$ can be different for each sample. With these training samples, the parameter ζ in conditional random field (29.2) is learned by *MLE* (see Chapter 12):

$$\max_{\zeta} \sum_{i=1}^n \log \frac{\exp(\zeta^\top \varphi_i(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}_i))}{\sum_{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m_i)}=1}^c \exp(\zeta^\top \varphi_i(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}))},$$

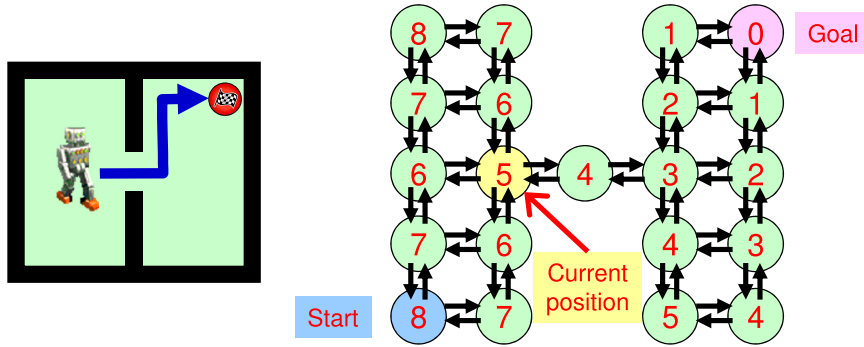
where

$$\varphi_i(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}_i) = \sum_{k=1}^{m_i} \varphi(\mathbf{x}_i^{(k)}, \mathbf{y}_i^{(k)}, \mathbf{y}_i^{(k-1)}).$$

The maximum likelihood solution may be obtained by a *stochastic gradient* algorithm (Section 15.3), as described in Fig. 29.3.

29.2.3 RECURSIVE COMPUTATION

However, the computational complexity of the second term of the gradient (see Fig. 29.3),

**FIGURE 29.4**

Dynamic programming, which solves a complex optimization problem by breaking it down into simpler subproblems recursively. When the number of steps to the goal is counted, dynamic programming trace back the steps from the goal. In this case, many subproblems of counting the number of steps from other positions are actually shared and thus dynamic programming can efficiently reuse the solutions to reduce the computation costs.

$$\frac{\sum_{y^{(1)}, \dots, y^{(m_i)}=1}^c \exp(\zeta^\top \varphi_i(\bar{x}_i, \bar{y})) \varphi_i(\bar{x}_i, \bar{y})}{\sum_{y^{(1)}, \dots, y^{(m_i)}=1}^c \exp(\zeta^\top \varphi_i(\bar{x}_i, \bar{y}))}, \quad (29.5)$$

grows exponentially with respect to the sequence length m_i . To cope with this problem, let us utilize the decomposition given by (29.3) based on *dynamic programming* (Fig. 29.4).

First, let us split the summation in the denominator of Eq. (29.5) into $(y^{(1)}, \dots, y^{(m_i-1)})$ and $y^{(m_i)}$:

$$\sum_{y^{(1)}, \dots, y^{(m_i)}=1}^c \exp\left(\sum_{k=1}^{m_i} \zeta^\top \varphi(x_i^{(k)}, y^{(k)}, y^{(k-1)})\right) = \sum_{y^{(m_i)}=1}^c A_{m_i}(y^{(m_i)}), \quad (29.6)$$

where

$$A_\tau(y) = \sum_{y^{(1)}, \dots, y^{(\tau-1)}=1}^c \exp\left(\sum_{k=1}^{\tau-1} \zeta^\top \varphi(x_i^{(k)}, y^{(k)}, y^{(k-1)}) + \zeta^\top \varphi(x_i^{(\tau)}, y, y^{(\tau-1)})\right).$$

$A_\tau(y^{(\tau)})$ can be expressed recursively using $A_{\tau-1}(y^{(\tau-1)})$ as

$$\begin{aligned} A_\tau(y^{(\tau)}) &= \sum_{y^{(1)}, \dots, y^{(\tau-1)}=1}^c \exp \left(\sum_{k=1}^{\tau} \zeta^\top \varphi(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)}) \right) \\ &= \sum_{y^{(\tau-1)}=1}^c A_{\tau-1}(y^{(\tau-1)}) \exp \left(\zeta^\top \varphi(\mathbf{x}_i^{(\tau)}, y^{(\tau)}, y^{(\tau-1)}) \right). \end{aligned}$$

The computational complexity of naively computing $A_{m_i}(y^{(m_i)})$ is $O(c^{m_i})$, while that of computing $A_{m_i}(y^{(m_i)})$ using the above recursive expression from $A_1(y^{(1)})$ is only $O(c^2 m_i)$. Note that the complexity for computing $\zeta^\top \varphi(\mathbf{x}, y, y')$ is not included here.

Next, let us split the summation in the numerator of Eq. (29.5) into $(y^{(1)}, \dots, y^{(k'-2)})$, $(y^{(k'-1)}, y^{(k')})$, and $(y^{(k'+1)}, \dots, y^{(m_i)})$:

$$\begin{aligned} & \sum_{y^{(1)}, \dots, y^{(m_i)}=1}^c \exp \left(\sum_{k=1}^{m_i} \zeta^\top \varphi(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)}) \right) \left(\sum_{k'=1}^{m_i} \varphi(\mathbf{x}_i^{(k')}, y^{(k')}, y^{(k'-1)}) \right) \\ &= \sum_{k'=1}^{m_i} \sum_{y^{(1)}, \dots, y^{(k'-2)}=1}^c \sum_{y^{(k'-1)}, y^{(k')}=1}^c \sum_{y^{(k'+1)}, \dots, y^{(m_i)}=1}^c \\ & \quad \exp \left(\sum_{k=1}^{m_i} \zeta^\top \varphi(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)}) \right) \varphi(\mathbf{x}_i^{(k')}, y^{(k')}, y^{(k'-1)}) \\ &= \sum_{k'=1}^{m_i} \sum_{y^{(k'-1)}, y^{(k')}=1}^c \varphi(\mathbf{x}_i^{(k')}, y^{(k')}, y^{(k'-1)}) \\ & \quad \times \exp \left(\zeta^\top \varphi(\mathbf{x}_i^{(k')}, y^{(k')}, y^{(k'-1)}) \right) A_{k'-1}(y^{(k'-1)}) B_{k'}(y^{(k')}), \end{aligned} \quad (29.7)$$

where

$$\begin{aligned} B_\tau(y) &= \sum_{y^{(\tau+1)}, \dots, y^{(m_i)}=1}^c \exp \left(\sum_{k=\tau+2}^{m_i} \zeta^\top \varphi(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)}) \right. \\ & \quad \left. + \zeta^\top \varphi(\mathbf{x}_i^{(\tau+1)}, y^{(\tau+1)}, y) \right). \end{aligned}$$

$B_\tau(y^{(\tau)})$ can be expressed recursively using $B_{\tau+1}(y^{(\tau+1)})$ as

$$\begin{aligned} B_\tau(y^{(\tau)}) &= \sum_{y^{(\tau+1)}, \dots, y^{(m_i)}=1}^c \exp \left(\sum_{k=\tau+1}^{m_i} \zeta^\top \varphi(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)}) \right) \\ &= \sum_{y^{(\tau+1)}=1}^c B_{\tau+1}(y^{(\tau+1)}) \exp \left(\zeta^\top \varphi(\mathbf{x}_i^{(\tau+1)}, y^{(\tau+1)}, y^{(\tau)}) \right). \end{aligned}$$

The computational complexity of naively computing $B_1(y^{(1)})$ is $O(c^{m_i})$, while that of computing $B_1(y^{(1)})$ using the above recursive expression from $B_{m_i}(y^{(m_i)})$ is only $O(c^2 m_i)$. Note that the complexity for computing $\zeta^\top \varphi(\mathbf{x}, y, y')$ is not included here.

Summarizing the above computations, stochastic gradient ascent for the conditional random field can be performed efficiently by computing $\{A_k(y^{(k)})\}_{k=1}^{m_i}$ and $\{B_k(y^{(k)})\}_{k=1}^{m_i}$ in advance using the recursive expressions.

29.2.4 PREDICTION FOR NEW SAMPLE

Trained conditional random field $q(\bar{y}|\bar{\mathbf{x}}; \hat{\zeta})$ allows us to predict the most probable label sequence $\bar{y} = (y^{(1)}, \dots, y^{(m)})$ for a test pattern sequence $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$ as

$$\operatorname{argmax}_{y^{(1)}, \dots, y^{(m)} \in \{1, \dots, c\}} q(\bar{y}|\bar{\mathbf{x}}; \hat{\zeta}).$$

However, the computational complexity for naive maximization over $y^{(1)}, \dots, y^{(m)}$ grows exponentially with respect to the sequence length m . Here, dynamic programming is used again to speed up the maximization.

Let us simplify the maximization problem as

$$\operatorname{argmax}_{y^{(1)}, \dots, y^{(m)} \in \{1, \dots, c\}} \frac{\exp(\hat{\zeta}^\top \varphi(\bar{\mathbf{x}}, \bar{y}))}{\sum_{\bar{y}=1}^{\bar{c}} \exp(\hat{\zeta}^\top \varphi(\bar{\mathbf{x}}, \bar{y}))} = \operatorname{argmax}_{y^{(1)}, \dots, y^{(m)} \in \{1, \dots, c\}} \hat{\zeta}^\top \varphi(\bar{\mathbf{x}}, \bar{y})$$

and decompose the maximization in the right-hand side into $(y^{(1)}, \dots, y^{(m-1)})$ and $y^{(m)}$:

$$\max_{y^{(1)}, \dots, y^{(m)} \in \{1, \dots, c\}} \hat{\zeta}^\top \varphi(\bar{\mathbf{x}}, \bar{y}) = \max_{y^{(m)} \in \{1, \dots, c\}} P_m(y^{(m)}),$$

where

$$P_\tau(y) = \max_{y^{(1)}, \dots, y^{(\tau-1)} \in \{1, \dots, c\}} \left[\sum_{k=1}^{\tau-1} \hat{\zeta}^\top \varphi(\mathbf{x}^{(k)}, y^{(k)}, y^{(k-1)}) + \hat{\zeta}^\top \varphi(\mathbf{x}^{(\tau)}, y, y^{(\tau-1)}) \right].$$

$P_\tau(y^{(\tau)})$ can be expressed by recursively using $P_{\tau-1}(y^{(\tau-1)})$ as

$$\begin{aligned} P_\tau(y^{(\tau)}) &= \max_{y^{(1)}, \dots, y^{(\tau-1)} \in \{1, \dots, c\}} \left[\sum_{k=1}^{\tau} \hat{\zeta}^\top \varphi(\mathbf{x}^{(k)}, y^{(k)}, y^{(k-1)}) \right] \\ &= \max_{y^{(\tau-1)} \in \{1, \dots, c\}} \left[P_{\tau-1}(y^{(\tau-1)}) + \hat{\zeta}^\top \varphi(\mathbf{x}^{(\tau)}, y^{(\tau)}, y^{(\tau-1)}) \right]. \end{aligned}$$

The computational complexity of naively computing $P_m(y^{(m)})$ is $O(c^m)$, while that of computing $P_m(y^{(m)})$ using the above recursive expression from $P_1(y^{(1)})$ is only $O(c^2 m)$. Note that the complexity for computing $\zeta^\top \varphi(\mathbf{x}, y, y')$ is not included here.

29.3 DETERMINISTIC CLASSIFICATION FOR SEQUENCES

The conditional random field models the class-posterior probability, which corresponds to a probabilistic classification method. On the other hand, in deterministic classification, only the class label is learned. In this section, based on the formulation of *multiclass support vector classification* given in Section 27.5, a deterministic sequence classification method is introduced. More specifically,

$$f_{\mathbf{w}_y}(\mathbf{x}) = \mathbf{w}_y^\top \mathbf{x}$$

is used as a score function for class y , and a test sample \mathbf{x} is classified into class \hat{y} as

$$\hat{y} = \operatorname{argmax}_{y=1, \dots, c} f_{\mathbf{w}_y}(\mathbf{x}).$$

The parameter \mathbf{w}_y is learned as

$$\min_{\mathbf{w}_1, \dots, \mathbf{w}_c, \xi} \left[\frac{1}{2} \sum_{y=1}^c \|\mathbf{w}_y\|^2 + C \sum_{i=1}^n \xi_i \right]$$

$$\text{subject to } \mathbf{w}_{y_i}^\top \mathbf{x}_i - \mathbf{w}_y^\top \mathbf{x}_i \geq t_{i,y} - \xi_i, \quad \forall i = 1, \dots, n, \forall y = 1, \dots, c,$$

where

$$t_{i,y} = \begin{cases} 0 & (y = y_i), \\ 1 & (y \neq y_i). \end{cases}$$

Note that, when $y = y_i$, the constraint $\mathbf{w}_{y_i}^\top \mathbf{x}_i - \mathbf{w}_y^\top \mathbf{x}_i \geq t_{i,y} - \xi_i$ is reduced to $\xi_i \geq 0$.

Applying the feature expression of the conditional random field, $\zeta^\top \boldsymbol{\varphi}(\bar{\mathbf{x}}, \bar{\mathbf{y}})$, to the above optimization yields

$$\begin{aligned} \min_{\zeta, \xi} & \left[\frac{1}{2} \|\zeta\|^2 + C \sum_{i=1}^n \xi_i \right] \\ \text{subject to } & \zeta^\top \Delta \boldsymbol{\varphi}_i(\bar{\mathbf{y}}) \geq t_{i,\bar{\mathbf{y}}} - \xi_i, \\ & \forall i = 1, \dots, n, \forall \bar{\mathbf{y}} = (\underbrace{1, \dots, 1}_m), \dots, (\underbrace{c, \dots, c}_m), \end{aligned}$$

where

$$\Delta \boldsymbol{\varphi}_i(\bar{\mathbf{y}}) = \boldsymbol{\varphi}_i(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}_i) - \boldsymbol{\varphi}_i(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}).$$

The *Lagrange dual* (Fig. 23.5) of the above optimization problem is given by

$$\begin{aligned} \max_{\{\alpha_{i,\bar{\mathbf{y}}}\}_{i,\bar{\mathbf{y}}}} & J(\{\alpha_{i,\bar{\mathbf{y}}}\}_{i,\bar{\mathbf{y}}}) \\ \text{subject to } & \sum_{y^{(1)}, \dots, y^{(m)}=1}^c \alpha_{i,\bar{\mathbf{y}}} = C, \quad \forall i = 1, \dots, n, \\ & \alpha_{i,\bar{\mathbf{y}}} \geq 0, \quad \forall i = 1, \dots, n, \forall y^{(1)}, \dots, y^{(m)} = 1, \dots, c, \end{aligned} \quad (29.8)$$

where

$$J(\{\alpha_{i,\bar{y}}\}_{i,\bar{y}}) = \sum_{i=1}^n \sum_{y^{(1)}, \dots, y^{(m)}=1}^c \alpha_{i,\bar{y}} t_{i,\bar{y}} \\ - \frac{1}{2} \sum_{i,i'=1}^n \sum_{y^{(1)}, \dots, y^{(m)}=1}^c \sum_{y^{(1)}, \dots, y^{(m)}=1}^c \alpha_{i,\bar{y}} \alpha_{i',\bar{y}'} \Delta \boldsymbol{\varphi}_i(\bar{y})^\top \Delta \boldsymbol{\varphi}_{i'}(\bar{y}').$$

This is a quadratic programming problem and thus the solution $\{\hat{\alpha}_{i,\bar{y}}\}_{i,\bar{y}}$ can be obtained by standard optimization software in principle. However, the number of optimization variables is exponential with respect to the sequence length, m , and thus it is not computationally tractable.

Here, let us utilize the same simplification trick as the conditional random field:

$$\boldsymbol{\varphi}(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \sum_{k=1}^m \boldsymbol{\varphi}(\mathbf{x}^{(k)}, y^{(k)}, y^{(k-1)}),$$

which yields

$$\Delta \boldsymbol{\varphi}_i(\bar{\mathbf{y}}) = \sum_{k=1}^{m_i} \Delta \boldsymbol{\varphi}_i(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)}),$$

where

$$\Delta \boldsymbol{\varphi}_i(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)}) = \boldsymbol{\varphi}_i(\mathbf{x}_i^{(k)}, y_i^{(k)}, y_i^{(k-1)}) - \boldsymbol{\varphi}_i(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)}).$$

Then the objective function in optimization problem (29.8) can be expressed as

$$J(\{\alpha_{i,\bar{y}}\}_{i,\bar{y}}) = \sum_{i=1}^n \sum_{k=1}^m \sum_{y^{(k)}=1}^c \mu_i(y^{(k)}) t_{i,\bar{y}} \\ - \frac{1}{2} \sum_{i,i'=1}^n \sum_{k,k'=1}^m \sum_{y^{(k-1)}, y^{(k)}=1}^c \sum_{y^{(k'-1)}, y^{(k')}=1}^c \Delta \boldsymbol{\varphi}_i(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)})^\top \\ \times \Delta \boldsymbol{\varphi}_{i'}(\mathbf{x}_{i'}^{(k')}, y^{(k')}, y^{(k'-1)}) \mu_i(y^{(k-1)}, y^{(k)}) \mu_{i'}(y^{(k'-1)}, y^{(k')}),$$

where

$$\mu_i(y^{(k)}) = \sum_{y^{(1)}, \dots, y^{(k-1)}, y^{(k+1)}, \dots, y^{(m_i)}=1}^c \alpha_{i,\bar{y}}, \\ \mu_i(y^{(k-1)}, y^{(k)}) = \sum_{y^{(1)}, \dots, y^{(k-2)}, y^{(k+1)}, \dots, y^{(m_i)}=1}^c \alpha_{i,\bar{y}}.$$

Thus, exponentially many variables $\{\alpha_{i,\bar{y}}\}_{i,\bar{y}}$ do not have to be optimized, but only optimizing $\{\mu_i(y^{(k)})\}_{k=1}^{m_i}$ and $\{\mu_i(y^{(k-1)}, y^{(k)})\}_{k=1}^{m_i}$ for $i = 1, \dots, n$ is sufficient.

However, since $\{\mu_i(y^{(k)})\}_{k=1}^{m_i}$ and $\{\mu_i(y^{(k-1)}, y^{(k)})\}_{k=1}^{m_i}$ are related to each other, the additional constraint,

$$\sum_{y^{(k-1)}=1}^c \mu_i(y^{(k-1)}, y^{(k)}) = \mu_i(y^{(k)}),$$

is necessary, together with the original constraints:

$$\sum_{y^{(k-1)}=1}^c \mu_i(y^{(k)}) = C \quad \text{and} \quad \mu_i(y^{(k-1)}, y^{(k)}) \geq 0.$$

From the dual solution $\{\widehat{\mu}_i(y^{(k)})\}_{k=1}^{m_i}$ and $\{\widehat{\mu}_i(y^{(k-1)}, y^{(k)})\}_{k=1}^{m_i}$, the primal solution $\widehat{\zeta}$ can be obtained as

$$\widehat{\zeta} = \sum_{i=1}^n \sum_{k=1}^{m_i} \sum_{y^{(k-1)}, y^{(k)}=1}^c \widehat{\mu}_i(y^{(k-1)}, y^{(k)}) \Delta \varphi_i(\mathbf{x}_i^{(k)}, y^{(k)}, y^{(k-1)}).$$

Once the solution $\widehat{\zeta}$ is obtained, classification of a test sample $\bar{\mathbf{x}} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$,

$$\operatorname{argmax}_{y^{(1)}, \dots, y^{(m)} \in \{1, \dots, c\}} \widehat{\zeta}^\top \varphi(\bar{\mathbf{x}}, \bar{\mathbf{y}}),$$

can be performed exactly in the same way as that described in Section 29.2.4.

This method is referred to as the *structured support vector machine* [114].

