

## Exercise Sheet 9

### Exercise 1: Structured Prediction for Classification (20 P)

While structured output learning is typically used for predicting complex output signals such as sequences or trees, the same framework can also be used to address the more standard problem of classification. Let  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$  be two data points and  $y, y' \in \{1, \dots, C\}$  their respective classes. Consider the structured output kernel

$$k_{\text{struct}}((\mathbf{x}, y), (\mathbf{x}', y')) = k(\mathbf{x}, \mathbf{x}') \cdot I(y = y'),$$

where  $k(\mathbf{x}, \mathbf{x}')$  is a positive semi-definite kernel with associated feature map  $\phi(\mathbf{x})$ , and where  $I(\cdot)$  is an indicator function that is 1 when the argument is true and 0 otherwise.

- (a) Show that the kernel  $k_{\text{struct}}((\mathbf{x}, y), (\mathbf{x}', y'))$  is positive semi-definite, that is, show that

$$\sum_{i=1}^N \sum_{j=1}^N c_i c_j k_{\text{struct}}((\mathbf{x}_i, y_i), (\mathbf{x}_j, y_j)) \geq 0$$

for all input/output pairs  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  and choice of real numbers  $c_1, \dots, c_N$ .

- (b) Find a feature map  $\phi_{\text{struct}}(\mathbf{x}, y)$  associated this kernel, i.e. satisfying

$$\langle \phi_{\text{struct}}(\mathbf{x}, y), \phi_{\text{struct}}(\mathbf{x}', y') \rangle = k_{\text{struct}}((\mathbf{x}, y), (\mathbf{x}', y'))$$

for all pairs  $(\mathbf{x}, y)$  and  $(\mathbf{x}', y')$ .

### Exercise 2: Dual Formulation of Structured Output Learning (20 P)

In structured output learning, we look for a linear model in joint feature space that produces a large margin between the correct prediction and all other possible predictions. The primal formulation of this problem can be expressed as:

$$\min_{\mathbf{w}, \xi} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n$$

subject to the constraints:

$$\forall_{n=1}^N \forall_{\mathbf{y} \neq \mathbf{y}_n} : \mathbf{w}^\top \Psi_{n, \mathbf{y}} \geq 1 - \xi_n \quad \text{and} \quad \forall_{n=1}^N : \xi_n \geq 0$$

where we have used the shortcut notation  $\Psi_{n, \mathbf{y}} = \phi(\mathbf{x}, \mathbf{y}_n) - \phi(\mathbf{x}, \mathbf{y})$ .

- (a) Show that the associated dual optimization problem is given by:

$$\max_{\alpha} \quad \sum_{n=1}^N \sum_{\mathbf{y} \neq \mathbf{y}_n} \alpha_{n, \mathbf{y}} - \frac{1}{2} \sum_{n, n'} \sum_{\substack{\mathbf{y} \neq \mathbf{y}_n \\ \mathbf{y}' \neq \mathbf{y}_{n'}}} \alpha_{n, \mathbf{y}} \alpha_{n', \mathbf{y}'} \langle \Psi_{n, \mathbf{y}}, \Psi_{n', \mathbf{y}'} \rangle$$

subject to the constraints:

$$\forall_{n=1}^N \forall_{\mathbf{y} \neq \mathbf{y}_n} : 0 \leq \alpha_{n, \mathbf{y}} \quad \text{and} \quad \forall_{n=1}^N : \sum_{\mathbf{y} \neq \mathbf{y}_n} \alpha_{n, \mathbf{y}} \leq C$$

- (b) Assuming  $k((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}'))$  is the kernel that induces the feature map  $\phi(\mathbf{x}, \mathbf{y})$ , express the dot product  $\langle \Psi_{n, \mathbf{y}}, \Psi_{n', \mathbf{y}'} \rangle$  in terms of this kernel.

**Exercise 3: Prediction of Output Sequences (20 P)**

Consider output sequences to predict to be of the type  $\mathbf{y} \in \{-1, 1\}^L$ , and the feature map:

$$\phi(\mathbf{x}, \mathbf{y}) = [\mathbf{x} \odot \mathbf{y}, 2 \cdot (\mathbf{y}_{1 \dots L-1} \odot \mathbf{y}_{2 \dots L})]$$

where  $\odot$  denotes the element-wise product between two vectors. The structured output model looks for the output  $\mathbf{y}$  that maximizes the matching function, i.e.

$$\max_{\mathbf{y}} \mathbf{w}^\top \phi(\mathbf{x}, \mathbf{y})$$

with  $\mathbf{w} \in \mathbb{R}^{2L-1}$ . In the following, we assume that  $L = 3$ , that the current parameter is  $\mathbf{w} = (1, 1, 1, 1, 1)$  and that we receive the input  $\mathbf{x} = (1, -1, 1)$ .

(a) *Show* that the problem of maximizing the matching function simplifies to:

$$\max_{y_1 \in \{-1, 1\}} \left\{ y_1 + \max_{y_2 \in \{-1, 1\}} \left\{ 2y_1y_2 - y_2 + \max_{y_3 \in \{-1, 1\}} \left\{ 2y_2y_3 + y_3 \right\} \right\} \right\}$$

(b) *Find* using the Viterbi procedure the best output  $(y_1, y_2, y_3)$ , that is, solve  $\max_{y_3} \{\}$  for every  $y_2$ , then solve  $\max_{y_2} \{\}$  for every  $y_1$ , and then solve  $\max_{y_1} \{\}$ . While doing so, keep track of the values in the sequence that have produced the respective maximums so that the optimal sequence can be reconstructed.

**Exercise 4: Programming (40 P)**

Download the programming files on ISIS and follow the instructions.