

**STRUCTURED MODEL LEARNING (SML2019)**  
**SEMINAR 4.**

**Assignment 1.** Let  $s: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a function defined as

$$s(x, x') = (\phi(x) - \phi(x'))^T \mathbf{W} (\phi(x) - \phi(x')),$$

which measures a dissimilarity between two images where  $\phi: \mathcal{X} \rightarrow \mathbb{R}^n$  is map extracting  $n$  features from image  $x \in \mathcal{X}$ , and  $\mathbf{W} \in \mathbb{R}^{n \times n}$  is a matrix. Consider a classifier  $h: \mathcal{X} \times \mathcal{X} \rightarrow \{-1, +1\}$  assigning a pair of images  $(x, x') \in \mathcal{X} \times \mathcal{X}$  into the positive class if their dissimilarity  $s(x, x')$  is not higher than a threshold  $b \in \mathbb{R}$  and to the negative class otherwise, i.e.

$$h(x, x'; \mathbf{W}, b) = \begin{cases} +1 & \text{if } s(x, x') \leq b, \\ -1 & \text{if } s(x, x') > b, \end{cases} \quad (1)$$

where  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  $b \in \mathbb{R}$  are parameters of the classifier.

**a)** Let  $\mathcal{T}^m = \{(x_A^j, x_b^j, y^j) \in (\mathcal{X} \times \mathcal{X} \times \{+1, -1\}) \mid j = 1, \dots, m\}$  be a set of training examples composed of a pair of images  $(x_A, x_b)$  and their label  $y$ . Describe a variant of the Perceptron algorithm which finds the parameters  $\mathbf{W} \in \mathbb{R}^{n \times n}$  and  $b \in \mathbb{R}$  such that the classifier (1) predicts all examples from  $\mathcal{T}^m$  correctly provided such parameters exists.

**b)** Extend the algorithm from assignment a) such that the found matrix  $\mathbf{W}$  is symmetric and positive definite, i.e.  $\mathbf{W}^T = \mathbf{W}$  and  $\langle \mathbf{u}, \mathbf{W} \mathbf{u} \rangle > 0, \forall \mathbf{u} \in \mathbb{R}^n$ .

**Assignment 2.** Consider a linear classifier  $h: \mathcal{X} \rightarrow \mathcal{Y}$  assigning inputs  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^n$  to classes  $\mathcal{Y} = \{1, \dots, Y\}$  based on the rule

$$h(\mathbf{x}; \mathbf{w}, b_1, \dots, b_{Y-1}) = 1 + \sum_{y=1}^{Y-1} \mathbb{I}[\langle \mathbf{x}, \mathbf{w} \rangle \geq b_y] \quad (2)$$

where  $\mathbf{w} \in \mathbb{R}^n$  and  $(b_1, \dots, b_{Y-1}) \in \mathbb{R}^{Y-1}$  are parameters. Let  $\mathcal{T}^m = \{(\mathbf{x}^j, y^j) \in (\mathcal{X} \times \mathcal{Y}) \mid j = 1, \dots, m\}$  be a training set of examples. Describe a variant of the Perceptron algorithm which finds the parameters such that the classifier (2) predicts all examples from  $\mathcal{T}^m$  correctly provided such parameters exist.

**Assignment 3.** Consider a linear max-sum classifier  $h: \mathcal{X} \rightarrow \mathcal{Y}^n$  for a sequence prediction

$$\mathbf{y}^* = h(\mathbf{x}; \mathbf{w}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} \left( \sum_{i=1}^n \langle \mathbf{w}, \phi^A(x, y_i) \rangle + \sum_{i=1}^{n-1} \langle \mathbf{w}, \phi^B(y_i, y_{i+1}) \rangle \right) \quad (3)$$

where

- $n$  is the length of the output sequence
- $\mathcal{X}$  is an arbitrary set of inputs
- $\mathcal{Y}$  is a finite set labels
- $\phi^A: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$  and  $\phi^B: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^d$  are the fixed feature maps
- $\mathbf{w} \in \mathbb{R}^d$  are the weights

Let  $\mathcal{T}^m = \{(x^j, y_1^j, \dots, y_n^j) \in (\mathcal{X} \times \mathcal{Y}^n) \mid j = 1, \dots, m\}$  be a set of training examples. Describe a variant of the Perceptron algorithm which finds the weights  $\mathbf{w} \in \mathbb{R}^d$  such that the classifier (3) predicts all examples from  $\mathcal{T}^m$  correctly provided such parameters exists. Describe two variants:

- Perceptron using the dynamic programming to implement the classification oracle.
- Perceptron which does not use the dynamic programming.

**Assignment 4.** Consider a linear classifier  $h: \mathcal{X} \rightarrow \mathcal{Y}$  assigning inputs  $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^n$  to classes  $\mathcal{Y} = \{1, \dots, Y\}$  based on the rule

$$h(\mathbf{x}; \mathbf{w}, b_1, \dots, b_Y) = \operatorname{argmax}_{y \in \mathcal{Y}} (y \langle \mathbf{x}, \mathbf{w} \rangle + b_y) \quad (4)$$

where  $\mathbf{w} \in \mathbb{R}^n$  and  $b_y \in \mathbb{R}$ ,  $y \in \mathcal{Y}$ , are parameters. Let  $\mathcal{T}^m = \{(\mathbf{x}^j, y^j) \in (\mathcal{X} \times \mathcal{Y}) \mid j = 1, \dots, m\}$  be a training set of examples. The goal is to learn parameters  $\mathbf{w}$  such that the predictor (4) has a small expectation of the Mean Absolute Deviation loss  $\ell(y, y') = |y - y'|$ . To this end, we employ the SO-SVM framework learning parameters of a generic linear classifier

$$h(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle \mathbf{w}, \phi(x, y) \rangle \quad (5)$$

by solving the convex problem  $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} [\frac{\lambda}{2} \|\mathbf{w}\|^2 + R(\mathbf{w})]$  where  $\lambda > 0$  is a regularization constant and

$$R(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max_{y \in \mathcal{Y}} \left( \ell(y^i, y) + \langle \mathbf{w}, \phi(x^i, y) - \phi(x^i, y^i) \rangle \right).$$

- Give an interpretation of the classification rule (4), i.e. for which type of prediction problems it is appropriate?
- Define the joint feature map  $\phi(x, y)$  so that (5) and (4) are equivalent.
- Write the risk  $R(\mathbf{w})$  instantiated for the classification rule (4) and  $\ell(y, y') = |y - y'|$ . Write a formula for a sub-gradient of  $R(\mathbf{w})$  at  $\mathbf{w}$ .

**Assignment 5.** Consider problem of learning a linear two-class SVM classifier

$$h(\mathbf{x}; \mathbf{w}, b) = \operatorname{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$$

from a training set  $\{(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^m, y^m)\} \in (\mathbb{R}^n \times \{+1, -1\})^m$  by solving

$$(\mathbf{w}^*, b^*) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \left[ \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right] \quad (6a)$$

subject to

$$\begin{aligned} y^i(\langle \mathbf{w}, \mathbf{x}^i \rangle + b) &\geq 1 - \xi_i, & i \in \{1, \dots, m\}, \\ \xi_i &\geq 0, & i \in \{1, \dots, m\}. \end{aligned} \tag{6b}$$

Note that the bias  $b$  is not contained in the quadratic regularizer. Convert the problem (6) to an unconstrained convex problem

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left[ \frac{\lambda}{2} \|\mathbf{w}\|^2 + R(\mathbf{w}) \right]$$

and derive an algorithm for evaluating  $R(\mathbf{w})$  and its sub-gradient.