# Structured Output Support Vector Machines

Czech Technical University in Prague

V. Franc

February 24, 2020

◆ Margin-rescaling loss

◆ Structured Output Support Vector Machines

**XEP33SML – Structured Model Learning, Summer 2020**

♦ Learning $h(x; \boldsymbol{w}) = \text{Argmax}_{y \in \mathcal{Y}} \langle \boldsymbol{w}, \boldsymbol{\phi}(x, y) \rangle$ from examples $\mathcal{T}^m = \{(x^i, y^i) \in (\mathcal{X} \times \mathcal{Y}) \mid i = 1, \ldots, m\}$ by ERM leads to

$$\boldsymbol{w}^* \in \underset{\boldsymbol{w} \in \mathbb{R}^n}{\text{Argmin}} \, R_{\mathcal{T}^m}(\boldsymbol{w}) \quad \text{where} \quad R_{\mathcal{T}^m}(\boldsymbol{w}) = \frac{1}{m} \sum_{i=1}^{m} \ell(y^i, h(x^i; \boldsymbol{w}))$$

♦ The SO-SVM approximates the ERM by a convex problem

$$\boldsymbol{w}^* \in \underset{\boldsymbol{w} \in \mathbb{R}^n}{\text{Argmin}} \left( \frac{\lambda}{2} \|\boldsymbol{w}\|^2 + R^{\psi}(\boldsymbol{w}) \right) \quad \text{where} \quad R^{\psi}(\boldsymbol{w}) = \frac{1}{m} \sum_{i=1}^{m} \psi(x^i, y^i, \boldsymbol{w})$$

♦ The surrograte loss $\psi \colon \mathcal{X} \times \mathcal{Y} \times \mathbb{R}^n \to \mathbb{R}$ is an upper bound:

$$\ell(y, h(x; \boldsymbol{w})) \leq \psi(x, y, \boldsymbol{w}), \quad \forall (x, y, \boldsymbol{w}) \in (\mathcal{X} \times \mathcal{Y} \times \mathbb{R}^n)$$

which is convex in $\boldsymbol{w}$ for any $(x, y)$.

◆ We require the score of the correct label $y^i$ to be higher than the score of any incorrect label $y$ by margin proportional to the loss $\ell(y^i, y)$:

$$\langle \boldsymbol{w}, \phi(x^i, y^i) \rangle \geq \langle \boldsymbol{w}, \phi(x^i, y) \rangle + \ell(y^i, y), \qquad \forall y \in \mathcal{Y} \setminus \{y^i\}$$

◆ Example: Sequencial OCR, Hamming distance $\ell(\boldsymbol{y}, \boldsymbol{y}') = \sum_{i=1}^{L} [\![ y_i \neq y_i' ]\!]$

$$\left\langle \phi\!\left( \text{JOHN}, JOHN \right), \boldsymbol{w} \right\rangle \geq \left\langle \phi\!\left( \text{JOHN}, AAAA \right), \boldsymbol{w} \right\rangle + 4$$

$$\left\langle \phi\!\left( \text{JOHN}, JOHN \right), \boldsymbol{w} \right\rangle \geq \left\langle \phi\!\left( \text{JOHN}, JAAA \right), \boldsymbol{w} \right\rangle + 3$$

$$\left\langle \phi\!\left( \text{JOHN}, JOHN \right), \boldsymbol{w} \right\rangle \geq \left\langle \phi\!\left( \text{JOHN}, JOAA \right), \boldsymbol{w} \right\rangle + 2$$

$$\left\langle \phi\!\left( \text{JOHN}, JOHN \right), \boldsymbol{w} \right\rangle \geq \left\langle \phi\!\left( \text{JOHN}, JOHA \right), \boldsymbol{w} \right\rangle + 1$$

$$\vdots$$

◆ We require the score of the correct label $y^i$ to be higher than the score of any incorrect label $y$ by margin proportional to the loss $\ell(y^i, y)$:

$$\langle \boldsymbol{w}, \phi(x^i, y^i) \rangle \geq \langle \boldsymbol{w}, \phi(x^i, y) \rangle + \ell(y^i, y), \qquad \forall y \in \mathcal{Y} \setminus \{y^i\}$$

◆ The margin rescaling loss

$$\psi(x^i, y^i, \boldsymbol{w}) = \max \left\{ 0, \max_{y \in \mathcal{Y} \setminus \{y^i\}} \left( \ell(y^i, y) + \langle \boldsymbol{w}, \boldsymbol{\phi}(x^i, y) \rangle - \langle \boldsymbol{w}, \boldsymbol{\phi}(x^i, y^i) \rangle \right) \right\}$$

◆ Upper bound of the true loss:

$$y^i \neq \hat{y} = h(x^i; \boldsymbol{w}) = \operatorname*{Argmax}_{y \in \mathcal{Y}} \langle \boldsymbol{w}, \boldsymbol{\phi}(x^i, y) \rangle$$

implies $\langle \boldsymbol{w}, \boldsymbol{\phi}(x^i, \hat{y}) \rangle - \langle \boldsymbol{w}, \boldsymbol{\phi}(x^i, y^i) \rangle \geq 0$ and hence

$$\psi(x^i, y^i, \boldsymbol{w}) \geq \ell(y^i, h(x^i, \boldsymbol{w})), \qquad \forall \boldsymbol{w} \in \mathbb{R}^n$$

◆ Using shortcuts $\ell_i(y) = \ell(y^i, y)$ and $\boldsymbol{\phi}_i(y) = \boldsymbol{\phi}(x^i, y) - \boldsymbol{\phi}(x^i, y^i)$ we can simplify the margin rescaling loss:

$$
\begin{aligned}
\psi(x^i, y^i, \boldsymbol{w}) &= \max\{0, \max_{y \in \mathcal{Y} \setminus \{y^i\}} \left( \ell(y^i, y) + \langle \boldsymbol{w}, \boldsymbol{\phi}(x^i, y) \rangle - \langle \boldsymbol{w}, \boldsymbol{\phi}(x^i, y^i) \rangle \right) \} \\
&= \max_{y \in \mathcal{Y}} \left( \ell(y^i, y) + \langle \boldsymbol{w}, \boldsymbol{\phi}(x^i, y) \rangle - \langle \boldsymbol{w}, \boldsymbol{\phi}(x^i, y^i) \rangle \right) \\
&= \max_{y \in \mathcal{Y}} \left( \ell_i(y) + \langle \boldsymbol{w}, \boldsymbol{\phi}_i(y) \rangle \right)
\end{aligned}
$$

◆ The margin-rescaling loss is a point-wise maximum over $|\mathcal{Y}|$ linear terms, hence, it is convex.

♦ The SO-SVM with margin-rescaling loss:

$$\boldsymbol{w}^* \in \operatorname*{Argmin}_{\boldsymbol{w} \in \mathbb{R}^n} \left( \frac{\lambda}{2} \|\boldsymbol{w}\|^2 + \underbrace{\frac{1}{m} \sum_{i=1}^{m} \max_{y \in \mathcal{Y}} \{\ell_i(y) + \langle \boldsymbol{w}, \boldsymbol{\phi}_i(y) \rangle\}}_{R^\psi(\boldsymbol{w})} \right)$$

♦ By using slack variables it can be rewritten as a Quadratic Program:

$$\boldsymbol{w}^* = \operatorname*{argmin}_{\boldsymbol{w} \in \mathbb{R}^n, \boldsymbol{\xi} \in \mathbb{R}^m} \left( \frac{\lambda}{2} \|\boldsymbol{w}\|^2 + \frac{1}{m} \sum_{i=1}^{m} \xi_i \right)$$

subject to

$$\xi_i \geq \ell_i(y) + \langle \boldsymbol{w}, \boldsymbol{\phi}_i(y) \rangle, \quad \forall i \in \{1, \ldots, m\}, \forall y \in \mathcal{Y}$$

♦ Note that the QP has $m|\mathcal{Y}|$ linear constaints !

**Theorem 1.** *Let $\ell\colon \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$ be a loss such that*
$\ell(y, y') = 0 \iff y = y'$, *and* $h(x) = \operatorname{argmax}_{y \in \mathcal{Y}} f(x, y)$ *a classifier*
$h\colon \mathcal{X} \to \mathcal{Y}$. *Then*

$$\ell(h(x), y) \leq \max_{y' \in \mathcal{Y} \setminus \{y\}} \psi\Big(f(x, y) - f(x, y'), \ell(y, y')\Big), \qquad \forall x \in \mathcal{X}, y \in \mathcal{Y}$$

*where* $\psi\colon \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ *is a function such that* $\psi(t, u) \geq u \, [\![ t \leq 0 ]\!]$.

PROOF:

$$
\begin{aligned}
\ell(h(x), y) &= \ell(h(x), y) \, [\![ f(x, y) - \max_{y' \neq y} f(x, y') \leq 0 ]\!] \\[2mm]
&\leq \psi\Big( f(x, y) - \max_{y' \neq y} f(x, y') \,,\, \ell(h(x), y) \Big) \\[2mm]
&= \psi\Big( f(x, y) - f(x, h(x)) \,,\, \ell(h(x), y) \Big) \\[2mm]
&\leq \max_{y' \neq y} \psi\Big( f(x, y) - f(x, y') \,,\, \ell(y', y) \Big)
\end{aligned}
$$

∎

◆ **Margin re-scaling loss:** $\psi(t, u) = \max\{0, u - t\}$

$$
\begin{aligned}
\ell(h(x), y) &\leq \max_{y' \neq y} \max \left\{ 0, \Delta(y', y) - f(x, y) + f(x, y') \right\} \\
&= \max_{y'} \left( \Delta(y', y) - f(x, y) + f(x, y') \right)
\end{aligned}
$$

◆ **Slack re-scaling loss:** $\psi(t, u) = \max\{0, u\,(1 - t)\}$

$$
\begin{aligned}
\ell(h(x), y) &\leq \max_{y' \neq y} \max \left\{ 0, \Delta(y', y)\,(1 - f(x, y) + f(x, y')) \right\} \\
&= \max_{y'} \Delta(y', y) \left( 1 - f(x, y) + f(x, y') \right)
\end{aligned}
$$