

Lecture 10: Structured Output SVM

Vojtěch Franc

May 3, 2015

8.B: Structured output SVM

- ◆ Convex surrogates of the empirical risk
- ◆ Risk surrogates for the max-sum classifier

8.B: Structured Output Support Vector Machines

- Given training examples $\mathcal{T} = \{(\mathbf{x}^i, \mathbf{y}^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \mathcal{I}\}$, the goal is to learn parameters $\mathbf{w} \in \mathbb{R}^n$ of a general linear classifier

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

where $\Psi: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^n$ is fixed feature map.

- Regularized empirical risk minimization based learning leads to solving

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left(\lambda \underbrace{\Omega(\mathbf{w})}_{\text{regularizer}} + \underbrace{\hat{R}_{\mathcal{T}}(\mathbf{w})}_{\substack{\text{surrogate of} \\ \text{empirical risk}}} \right)$$

where $\Omega: \mathbb{R}^n \rightarrow \mathbb{R}$ is a (convex) regularizer and $\hat{R}_{\mathcal{T}}: \mathbb{R}^n \rightarrow \mathbb{R}$ is a surrogate of the empirical risk

$$R_{\mathcal{T}}(h(\cdot; \mathbf{w})) = \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{y}^i, h(\mathbf{x}^i; \mathbf{w}))$$

and $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$ such that $\ell(\mathbf{y}, \mathbf{y}') = 0$ iff $\mathbf{y} = \mathbf{y}'$ is the loss.

Questions:

- How to construct the surrogate $\hat{R}_{\mathcal{T}}$ for a generic linear classifier and loss ?
- How to solve the risk minimization problem ?

8.B: Slack-rescaling loss

A loss of the linear classifier $h(\mathbf{x}; \mathbf{w})$ on $(\mathbf{x}^i, \mathbf{y}^i)$ can be upper-bounded by a convex function:

$$\begin{aligned}
 \ell(\mathbf{y}^i, h(\mathbf{x}^i, \mathbf{w})) &= \max_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}^i, \mathbf{y}) [\max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}') \rangle - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \leq 0] \\
 &\leq \max_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}^i, \mathbf{y}) [\langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \leq 0] \\
 &\leq \max_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}^i, \mathbf{y}) \max\{0, 1 - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle\} \\
 &= \max_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}^i, \mathbf{y}) \left[1 - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right] \\
 &= \hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) .
 \end{aligned}$$

Remark: The slack-rescaling loss is non-zero if the margin of incorrect label $\gamma = \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle$ is less than 1.

Learning: RRM leads to a convex problem $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} [\lambda \Omega(\mathbf{w}) + \hat{R}_{\mathcal{T}}(\mathbf{w})]$ where

$$\hat{R}_{\mathcal{T}}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}^i, \mathbf{y}) \left[1 - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right]$$

8.B: Slack-rescaling loss

RRM with the slack-rescaling loss is scaling invariant. Let $\ell^{sc} = K\ell$, $K > 0$ be a scaled version of the loss ℓ . Then

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \left[\lambda \Omega(\mathbf{w}) + \sum_{i=1}^m \hat{\ell}(\mathbf{x}^i, \mathbf{y}^i; \mathbf{w}) \right] = \operatorname{argmin}_{\mathbf{w}} \left[\frac{\lambda}{K} \Omega(\mathbf{w}) + \sum_{i=1}^m \hat{\ell}^{sc}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) \right]$$

First order oracle: compute the value $\hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$ and its sub-gradient $\hat{\ell}'(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$ at \mathbf{w}

1. Solve $\hat{\mathbf{y}}^i = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}^i, \mathbf{y}) \left[1 - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right]$
2. Compute sub-gradient $\hat{\ell}'(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \ell(\mathbf{y}^i, \hat{\mathbf{y}}^i) \left[\Psi(\mathbf{x}^i, \hat{\mathbf{y}}^i) - \Psi(\mathbf{x}^i, \mathbf{y}^i) \right]$.

Remarks:

- ◆ An efficient oracle can be constructed in special cases, e.g. when $|\mathcal{Y}|$ is small or the loss is simple like the 0/1-loss $\ell(\mathbf{y}, \mathbf{y}') = \llbracket \mathbf{y} \neq \mathbf{y}' \rrbracket$.
- ◆ It remains an open problem how to do it for a general loss.

8.B: Margin-rescaling loss

A loss of the linear classifier $h(\mathbf{x}; \mathbf{w})$ on $(\mathbf{x}^i, \mathbf{y}^i)$ can be upper-bounded by a convex function:

$$\begin{aligned}
 \ell(\mathbf{y}^i, h(\mathbf{x}^i, \mathbf{w})) &= \max_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}^i, \mathbf{y}) [\max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}') \rangle - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \leq 0] \\
 &\leq \max_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}^i, \mathbf{y}) [\langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \leq 0] \\
 &\leq \max_{\mathbf{y} \in \mathcal{Y}} \max \left\{ 0, \ell(\mathbf{y}^i, \mathbf{y}) - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right\} \\
 &= \max_{\mathbf{y} \in \mathcal{Y}} \left[\ell(\mathbf{y}^i, \mathbf{y}) - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right] \\
 &= \hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})
 \end{aligned}$$

Remark: The margin-rescaling loss is non-zero if the margin of incorrect label $\gamma = \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle$ is less than $\ell(\mathbf{y}^i, \mathbf{y})$.

Learning: RRM leads to a convex problem $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} [\lambda \Omega(\mathbf{w}) + \hat{R}_{\mathcal{T}}(\mathbf{w})]$ where

$$\hat{R}_{\mathcal{T}}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max_{\mathbf{y} \in \mathcal{Y}} \left[\ell(\mathbf{y}^i, \mathbf{y}) - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right]$$

8.B: Margin-rescaling loss

The margin-rescaling loss

$$\hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \max_{\mathbf{y} \in \mathcal{Y}} \left[\ell(\mathbf{y}^i, \mathbf{y}) - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right]$$

Note that the slack-rescaling and the margin-rescaling surrogates coincide for the 0/1-loss
 $\ell(\mathbf{y}, \mathbf{y}') = [\mathbf{y} \neq \mathbf{y}']$

First order oracle: compute the value $\hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$ and its sub-gradient $\hat{\ell}'(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$ at \mathbf{w}

1. Solve augmented classification problem (ACP)

$$\hat{\mathbf{y}}^i \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \left[\ell(\mathbf{y}^i, \mathbf{y}) + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right]$$

2. Compute sub-gradient $\hat{\ell}'(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \Psi(\mathbf{x}^i, \hat{\mathbf{y}}^i) - \Psi(\mathbf{x}^i, \mathbf{y}^i)$.

Notice the analogy with Perceptron algorithm: in order to learn you have to be able to classify.

8.B: Margin-rescaling loss for max-sum classifier

Consider the max-sum classifier

$$\hat{\mathbf{y}} \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \langle \mathbf{w}, \Psi(\mathbf{x}, \mathbf{y}) \rangle = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \left\langle \mathbf{w}, \sum_{v \in \mathcal{V}} \Psi_v(\mathbf{x}, y_v) + \sum_{\{v, v'\} \in \mathcal{E}} \Psi_{vv'}(\mathbf{x}, y_v, y_{v'}) \right\rangle$$

and an additive loss $\ell(\mathbf{y}, \mathbf{y}') = \sum_{v \in \mathcal{V}} \ell_v(y_v, y'_{v'})$ where $\ell_v: \mathcal{Y} \rightarrow [0, \infty)$, $v \in \mathcal{V}$
 (e.g. the Hamming loss)

Then the augmented classification problem boils down to solving

$$\begin{aligned} \hat{\mathbf{y}}^i &\in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \left[\ell(\mathbf{y}^i, \mathbf{y}) + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right] \\ &= \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \left[\sum_{v \in \mathcal{V}} \left(\ell_v(y_v^i, y_v) + \langle \mathbf{w}, \Psi_v(\mathbf{x}^i, y_v) \rangle \right) + \sum_{\{v, v'\} \in \mathcal{E}} \langle \mathbf{w}, \Psi_{vv'}(\mathbf{x}^i, y_v, y_{v'}) \rangle \right] \end{aligned}$$

The ACP is tractable if the graph $(\mathcal{V}, \mathcal{E})$ is acyclic.

8.B: Margin-rescaling loss for sub-modular max-sum classifier

Consider a class of the max-sum classifier

$$\hat{\mathbf{y}} \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \langle \mathbf{w}, \Psi(\mathbf{x}, \mathbf{y}) \rangle = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \left\langle \mathbf{w}, \sum_{v \in \mathcal{V}} \Psi_v(\mathbf{x}, y_v) + \sum_{\{v, v'\} \in \mathcal{E}} \Psi_{vv'}(y_v, y_{v'}) \right\rangle$$

where $\mathbf{w} \in \mathcal{W} \subset \mathbb{R}^n$ such that $-g_{vv'}(y, y') = \langle \mathbf{w}, \Psi_{vv'}(y, y') \rangle$ is sub-modular w.r.t. an ordering of $\mathcal{Y} = \{1, \dots, K\}$.

Learning: RRM with additive loss $\ell(\mathbf{y}, \mathbf{y}') = \sum_{v \in \mathcal{V}} \ell_v(y_v, y'_{v'})$ leads to a constrained convex problem

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left[\lambda \Omega(\mathbf{w}) + \sum_{i=1}^m \hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) \right]$$

subject to

$$\left\langle \mathbf{w}, \Psi_{vv'}(y, y') + \Psi_{vv'}(y+1, y'+1) - \Psi_{vv'}(y, y'+1) - \Psi_{vv'}(y+1, y') \right\rangle \geq 0,$$

$$\{v, v'\} \in \mathcal{E}, y, y' \in \{1, \dots, K-1\}$$

Provided the solver maintains the intermediate solution \mathbf{w} feasible the max-sum problems associated with the ACP are sub-modular and thus tractable.

8.B: LP relaxation of margin-rescaling loss

The ACP associated to a general max-sum classifier and additive loss

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{Y}} \left[\ell(\mathbf{y}^i, \mathbf{y}) + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right] \\ &= \max_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \underbrace{\left[\sum_{v \in \mathcal{V}} \left(\ell_v(y_v^i, y_v) + \langle \mathbf{w}, \Psi_v(\mathbf{x}^i, y_v) \rangle \right) + \sum_{\{v, v'\} \in \mathcal{E}} \langle \mathbf{w}, \Psi_{vv'}(\mathbf{x}^i, y_v, y_{v'}) \rangle \right]}_{f(\mathbf{y}; \mathbf{x}^i, \mathbf{y}^i, \mathbf{w})} \end{aligned}$$

can be upper bounded via the LP relaxation:

$$\max_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} f(\mathbf{y}; \mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) \leq \min_{\varphi} U(\varphi; \mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$$

where

$$U(\varphi; \mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \left[\sum_{v \in \mathcal{V}} \max_{y \in \mathcal{Y}} q_v^\varphi(y; \mathbf{x}^i, y_v^i, \mathbf{w}) + \sum_{\{v, v'\} \in \mathcal{E}} \max_{(y, y') \in \mathcal{Y}^2} g_{vv'}^\varphi(y, y'; \mathbf{x}^i, \mathbf{w}) \right]$$

and $\varphi \in \mathbb{R}^{2|\mathcal{E}||\mathcal{Y}|}$ is a vector of dual variables $\varphi_{vv'}: \mathcal{Y} \rightarrow \mathbb{R}$, $\varphi_{v'v}: \mathcal{Y} \rightarrow \mathbb{R}$, $\{v, v'\} \in \mathcal{E}$ and

$$\begin{aligned} q_v^\varphi(y; \mathbf{x}^i, y_v^i, \mathbf{w}) &= \ell_v(y_v^i, y_v) + \langle \mathbf{w}, \Psi_v(\mathbf{x}^i, y) \rangle - \sum_{v' \in \mathcal{N}(v)} \varphi_{vv'}(y), \quad v \in \mathcal{V}, y \in \mathcal{Y} \\ g_{vv'}^\varphi(y, y'; \mathbf{x}^i, \mathbf{w}) &= \langle \mathbf{w}, \Psi_{vv'}(\mathbf{x}^i, y, y') \rangle + \varphi_{vv'}(y) + \varphi_{v'v}(y'), \quad \{v, v'\} \in \mathcal{E}, y, y' \in \mathcal{Y} \end{aligned}$$

8.B: LP relaxation of margin-rescaling loss

10/13

The LP-relaxed margin-rescaling surrogate is obtained by substituting the LP dual:

$$\begin{aligned}
 \hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) &= \max_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \left[\ell(\mathbf{y}^i, \mathbf{y}) + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right] - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle \\
 &\leq \min_{\varphi} U(\varphi; \mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle \\
 &= \hat{\ell}_{\text{LP}}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})
 \end{aligned}$$

Learning: RRM leads to a convex problem $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} [\lambda \Omega(\mathbf{w}) + \hat{R}_{\mathcal{T}}(\mathbf{w})]$ where

$$\hat{R}_{\mathcal{T}}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \hat{\ell}_{\text{LP}}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \min_{\varphi} \left[U(\varphi; \mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle \right]$$

First order oracle: compute the value $\hat{\ell}_{\text{LP}}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$ and its sub-gradient $\hat{\ell}'_{\text{LP}}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$

1. Solve the LP dual

$$\varphi^i = \operatorname{argmin}_{\varphi} U(\varphi; \mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$$

2. Compute the sub-gradient

$$\hat{\ell}'_{\text{LP}}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) \in \partial_{\mathbf{w}} [U(\varphi^i; \mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle]$$

8.B: PosLearn: surrogate of additive loss

A linear classifier $h(\mathbf{x}; \mathbf{w}) = (h_v(\mathbf{x}; \mathbf{w}) \in \mathcal{Y} \mid v \in \mathcal{V}) \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \langle \mathbf{w}, \Psi(\mathbf{x}, \mathbf{y}) \rangle$ can be evaluated component-wise

$$h_v(\mathbf{x}; \mathbf{w}) \in \operatorname{argmax}_{y_v \in \mathcal{Y}} \max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}, \mathbf{y}) \rangle \quad \text{where} \quad \mathcal{Y}(y_v) = \{ \mathbf{y}' \in \mathcal{Y}^{\mathcal{V}} \mid y'_v = y_v \}$$

An additive loss $\ell(\mathbf{y}^i, h(\mathbf{x}^i; \mathbf{w})) = \sum_{v \in \mathcal{V}} \ell_v(y_v^i, h_v(\mathbf{x}^i; \mathbf{w}))$ can be bound by

$$\hat{\ell}(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \sum_{v \in \mathcal{V}} \hat{\ell}_v(\mathbf{x}^i, y_v^i, \mathbf{w}) \text{ where}$$

$$\begin{aligned} \ell_v(y_v^i, h_v(\mathbf{x}^i, \mathbf{w})) &= \max_{y_v \in \mathcal{Y}} \ell_v(y_v^i, y_v) [\max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}') \rangle - \max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \leq 0] \\ &\leq \max_{y_v \in \mathcal{Y}} \ell_v(y_v^i, y_v) [\langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle - \max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \leq 0] \\ &\leq \max_{y_v \in \mathcal{Y}} \ell_v(y_v^i, y_v) \max \left\{ 0, 1 - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right\} \\ &= \max_{y_v \in \mathcal{Y}} \ell_v(y_v^i, y_v) \left[1 - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right] \\ &= \hat{\ell}_v(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) \end{aligned}$$

8.B: PosLearn: surrogate of additive loss

Learning: RRM leads to a convex problem $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left[\lambda \Omega(\mathbf{w}) + \hat{R}_{\mathcal{T}}(\mathbf{w}) \right]$ where

$$\begin{aligned}\hat{R}_{\mathcal{T}}(\mathbf{w}) &= \frac{1}{m} \sum_{i=1}^m \sum_{v \in \mathcal{V}} \hat{\ell}_v(\mathbf{x}^i, y_v^i, \mathbf{w}) \\ &= \frac{1}{m} \sum_{i=1}^m \sum_{v \in \mathcal{V}} \max_{y_v \in \mathcal{Y}} \ell_v(y_v^i, y_v) \left[1 - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right]\end{aligned}$$

For example, in the case of the Hamming loss $\ell_v(y, y') = \llbracket y \neq y' \rrbracket$ the surrogate simplifies to

$$\hat{\ell}_v(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \max \left\{ 0, 1 - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \max_{y_v \in \mathcal{Y} \setminus \{y_v^i\}} \max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right\}$$

8.B: PosLearn: surrogate of additive loss

Consider surrogate of the Hamming loss $\ell_v(y, y') = \llbracket y \neq y' \rrbracket$ which reads

$$\hat{\ell}_v(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \max \left\{ 0, 1 - \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle + \max_{y_v \in \mathcal{Y} \setminus \{y_v^i\}} \max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle \right\}$$

First order oracle: compute the value $\hat{\ell}_v(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$ and its sub-gradient $\hat{\ell}'_v(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w})$ at \mathbf{w}

1. Compute

$$\hat{\mathbf{y}}_v^i \in \operatorname{argmax}_{y_v \in \mathcal{Y} \setminus \{y_v^i\}} \max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle$$

which is the best labeling not containing y_v^i at position v .

2. Compute sub-gradient

$$\hat{\ell}'_v(\mathbf{x}^i, \mathbf{y}^i, \mathbf{w}) = \begin{cases} \mathbf{0} & \text{if } \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle \geq 1 + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \hat{\mathbf{y}}_v^i) \rangle \\ \Psi(\mathbf{x}^i, \hat{\mathbf{y}}_v^i) - \Psi(\mathbf{x}^i, \mathbf{y}^i) & \text{otherwise} \end{cases}$$

For example, the max-marginals $\max_{\mathbf{y} \in \mathcal{Y}(y_v)} \langle \mathbf{w}, \Psi(\mathbf{x}, \mathbf{y}) \rangle$ can be computed efficiently for max-sum classifier with acyclic graph $(\mathcal{V}, \mathcal{E})$.