

Multiclass Perceptron

Key Idea: If the prediction is incorrect ($\hat{y} \neq y^*$) then increase the score of y^* and decrease the score of \hat{y}

Algo: (Multiclass Perceptron) $\Theta \in \mathbb{R}^M$ $f(\cdot, \cdot) \in \mathbb{R}^M$

$\Theta \leftarrow 0$

while not converged:

for (x, y^*) in D :

$$\hat{y} = \arg\max_{j \in \{1, \dots, K\}} \Theta^T f(\vec{x}, j)$$

if $\hat{y} \neq y^*$:

$$\Theta \leftarrow \Theta + \frac{f(\vec{x}, y^*) - f(\vec{x}, \hat{y})}{\text{Q1: fill in the blank}}$$

$$(\Theta + f(\vec{x}, y^*))^T f(\vec{x}, y^*) \geq \Theta^T f(\vec{x}, y^*)$$

Algo. (Structured Per.)

"

"

"

$$\hat{y} = \arg\max_{y \in \mathcal{Y}(\vec{x})} \Theta^T f(\vec{x}, y)$$

"

"

Structured Perceptron

- Algo:
- Same as Multiclass Perceptron w/ each possible output structure $\vec{y} \in \mathcal{Y}(\vec{x})$ as a "class"
- Training examples are (\vec{x}, \vec{y}) where $\vec{x} \in \mathcal{X}$ and $\vec{y} \in \mathcal{Y}(\vec{x})$
- Feature functions $f(\vec{x}, \vec{y}) \in \mathbb{R}^M$ ← repr. of (\vec{x}, \vec{y}) pair
- Predict $\hat{y} = \arg\max_{\vec{y} \in \mathcal{Y}(\vec{x})} \Theta^T f(\vec{x}, \vec{y})$

Key Question: How to compute?

If $\exp(\Theta^T f(\vec{x}, \vec{y}))$ decomposes multiplicatively according to some factor graph, then solve this "MAP inference problem" w/ MILP, e.g.,

$$\begin{aligned} \exp(\Theta^T f(\vec{x}, \vec{y})) &= \exp\left(\sum_c \Theta^T f_c(\vec{x}, \vec{y}_c)\right) \\ &= \prod_c \exp(\Theta^T f_c(\vec{x}, y_c)) \end{aligned}$$

$$= \prod_c \exp(\theta^T f_c(\vec{x}, y_c))$$

$$= \prod_c \psi_c(\vec{x}, y_c)$$

Structured SVM

Data: $\{(\vec{x}^{(i)}, y^{(i)})\}_{i=1}^N$ $\vec{x}^{(i)} \in \mathcal{X}$ $y^{(i)} \in \mathcal{Y}(\vec{x})$

Model: Linear $\hat{y} = h(\vec{x}) = \arg \max_{\tilde{y} \in \mathcal{Y}(\vec{x})} \underbrace{w^T f(\vec{x}, \tilde{y})}_{s_w(\vec{x}, \tilde{y})}$

Q.P.: $\min_{w, e} \frac{1}{2} (\|w\|_2)^2 + C \sum_{i=1}^N \sum_{\hat{y} \in \mathcal{Y}(\vec{x}^{(i)})} e_{i, \hat{y}}$

s.t. $s_w(\vec{x}^{(i)}, y^{(i)}) - s_w(\vec{x}^{(i)}, \hat{y}) \geq \underbrace{\ell(y^{(i)}, \hat{y})}_{\text{Hamming loss}} - e_{i, \hat{y}} \quad \forall i$
 $e_{i, \hat{y}} \geq 0 \quad \forall \hat{y} \in \mathcal{Y}(\vec{x}^{(i)}) \setminus y^{(i)}$

Key Idea:

- Score of true $y^{(i)}$ should be larger than the score of the predicted \hat{y} by some margin
- margin to be scaled by the loss

Ex: $\vec{x}^{(i)} = \text{time flies fast}$

	$s_w(\vec{x}^{(i)}, \cdot)$	$\ell(y^{(i)}, \cdot)$
$y^{(i)} = N \quad V \quad A$	7	0
$\hat{y}_1 = V \quad V \quad A$	5	1
$\hat{y}_2 = N \quad N \quad A$	4	1
$\hat{y}_3 = V \quad N \quad A$	3.5	2
\vdots	\vdots	\vdots
$\hat{y}_{\text{LAST}} = \dots$	-7	

Sorted by score

Annotations: $+2$ (between 7 and 5), $+3$ (between 5 and 4), $+3.5$ (between 4 and 3.5).

Q: How many constraints are in this Q.P.?

A: $O(N \cdot \max_i |\mathcal{Y}(\vec{x}^{(i)})|)$

Q.P. w/ fewer constraints

Key Idea: fold a maximization problem into constraint

Q.P.: $\min_{w, e} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N e_i$

s.t. $s_w(\vec{x}^{(i)}, \vec{y}^{(i)}) \geq \left[\max_{\hat{y} \in \mathcal{Y}(\vec{x})} s_w(\vec{x}^{(i)}, \hat{y}) + \ell(y^{(i)}, \hat{y}) \right] - e_i$

$$\text{s.t. } s_w(\vec{x}^{(i)}, \vec{y}^{(i)}) \geq \left[\max_{\hat{y} \in \mathcal{Y}(\vec{x})} s_w(\vec{x}^{(i)}, \hat{y}) + \ell(y^{(i)}, \hat{y}) \right] - e_i \quad \forall i$$

$$e_i \geq 0$$

these N constraints and N slack vars replaced
the $O(N \cdot \max_i |\mathcal{Y}(\vec{x}^{(i)})|)$ constraints and slack vars above

Q.P. w/ Hinge Loss

$$\min_w \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \left[\max(0, \left[\max_{\hat{y} \in \mathcal{Y}(\vec{x})} s_w(\vec{x}^{(i)}, \hat{y}) + \ell(y^{(i)}, \hat{y}) \right] - s_w(\vec{x}^{(i)}, y^{(i)}) \right]$$

$$\rightarrow \text{structured hinge loss} = \ell_w^{\text{s.H.}}(\vec{x}^{(i)}, y^{(i)})$$

Sub-Gradient

$$\nabla \ell_w^{\text{s.H.}}(\vec{x}^{(i)}, y^{(i)}) = \begin{cases} 0 & \text{if } \ell_w^{\text{s.H.}}(\vec{x}^{(i)}, y^{(i)}) = 0 \\ f(\vec{x}^{(i)}, y^{(i)}) - f(\vec{x}^{(i)}, \hat{y}) & \text{otherwise} \end{cases}$$

$$\text{where } \hat{y} = \underset{\hat{y} \in \mathcal{Y}(\vec{x}^{(i)})}{\operatorname{argmax}} s_w(\vec{x}^{(i)}, \hat{y}) + \ell(y^{(i)}, \hat{y})$$

not your standard
MAP inference problem

Train w/ Stochastic Subgradient Gradient

OH

$$p(z) = \phi_z$$

ϕ

0.2	0.1	0.6	0.1
-----	-----	-----	-----

$$p(z) = \phi_z$$

⁰⁾

0.2	0.1	0.6	0.1
-----	-----	-----	-----

$$\sum_z p(z) = 1$$

z discrete

$$\int_z p(z) dz = 1 \quad z \in \mathbb{R}$$

$$\int_z p(z|x) dz = 1$$