Self-Paced Learning for Latent Variable Models

Kumar M P, Packer B, Koller D. //NIPS 2010

In ICML 2009

 \bigstar Learning from easier aspects of the task, and gradually increase the difficulty level

- ★ Expected two advantages:
- Help find a better local minima (as a regularizer)
- Speed the convergence of training towards the global minimum (for convex problem)

★ Basic steps:

- Sort samples according to certain "easiness" measure
- Gradually add samples into training from easy to complex





Φ(x,y,h)

Latent SSVM

$$\min_{\mathbf{w},\xi_i \ge 0} \frac{1}{2} ||\mathbf{w}||^2 + \frac{C}{n} \sum_{i=1}^n \xi_i, \quad \text{For any given w, the value of } \xi_i \\ \text{can be shown to be an upper bound on the risk } \Delta(y_i, y^{-i}(w))$$
s.t.
$$\max_{h_i \in \mathcal{H}} \mathbf{w}^\top \left(\Phi(\mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i) - \Phi(\mathbf{x}_i, \hat{\mathbf{y}}_i, \hat{\mathbf{h}}_i) \right) \ge \Delta(\mathbf{y}_i, \hat{\mathbf{y}}_i) - \xi_i, \\ \forall \hat{\mathbf{y}}_i \in \mathcal{Y}, \forall \hat{h}_i \in \mathcal{H}, i = 1, \cdots, n.$$

Algorithm 2 The CCCP algorithm for parameter estimation of latent SSVM.

input $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), \cdots, (\mathbf{x}_n, \mathbf{y}_n)\}, \mathbf{w}_0, \epsilon.$ 1: $t \leftarrow 0$

- 2: repeat
- 3: Update $\mathbf{h}_i^* = \operatorname{argmax}_{h_i \in \mathcal{H}} \mathbf{w}_t^\top \Phi(\mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i).$
- 4: Update \mathbf{w}_{t+1} by fixing the hidden variables for output \mathbf{y}_i to \mathbf{h}_i^* and solving the corresponding SSVM problem. Specifically,

 $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} ||\mathbf{w}||^2 + \frac{C}{n} \sum_i \max\{0, \Delta(\mathbf{y}_i, \hat{\mathbf{y}}_i) + \mathbf{w}^\top (\Phi(\mathbf{x}_i, \hat{\mathbf{y}}_i, \hat{\mathbf{h}}_i) - \Phi(\mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i^*))\}.$ 5: $t \leftarrow t+1.$

6: **until** Objective function cannot be decreased below tolerance ϵ .

$$\mathbf{w}_{t+1} = \operatorname*{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left(r(\mathbf{w}) + \sum_{i=1}^n f(\mathbf{x}_i, \mathbf{y}_i; \mathbf{w}) \right)$$

f(.) is the negative log-likelihood for EM or an upper bound on the risk for latent SSVM (or any other criteria for parameter learning).

binary variables vi that indicate whether the ith sample is easy or not. Only easy samples contribute to the objective function.

$$(\mathbf{w}_{t+1}, \mathbf{v}_{t+1}) = \operatorname*{argmin}_{\mathbf{w} \in \mathbb{R}^d, \mathbf{v} \in \{0,1\}^n} \left(r(\mathbf{w}) + \sum_{i=1}^n v_i f(\mathbf{x}_i, \mathbf{y}_i; \mathbf{w}) - \frac{1}{K} \sum_{i=1}^n v_i \right)$$

$$(\mathbf{w}_{t+1}, \mathbf{v}_{t+1}) = \operatorname*{argmin}_{\mathbf{w} \in \mathbb{R}^d, \mathbf{v} \in \{0,1\}^n} \left(r(\mathbf{w}) + \sum_{i=1}^n v_i f(\mathbf{x}_i, \mathbf{y}_i; \mathbf{w}) - \frac{1}{K} \sum_{i=1}^n v_i \right)$$

Algorithm 3 *The self-paced learning algorithm for parameter estimation of latent* SSVM.

input
$$\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), \cdots, (\mathbf{x}_n, \mathbf{y}_n)\}, \mathbf{w}_0, K_0, \epsilon.$$

1: $t \leftarrow 0, K \leftarrow K_0.$

2: repeat

- 3: Update $h_i^* = \operatorname{argmax}_{h_i \in \mathcal{H}} \mathbf{w}_t^\top \Phi(\mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i).$
- 4: Update \mathbf{w}_{t+1} by using ACS to minimize the objective $\frac{1}{2} ||\mathbf{w}||^2 + \frac{C}{n} \sum_{i=1}^n v_i \xi_i \frac{1}{K} \sum_{i=1}^n v_i$ subject to the constraints of problem (2) as well as $\mathbf{v} \in \{0, 1\}^n$.

5:
$$t \leftarrow t+1, K \leftarrow K/\mu$$
.

6: **until** $v_i = 1, \forall i$ and the objective function cannot be decreased below tolerance ϵ .



Noun Phrase Coreference

Given the occurrence of all the nouns in a document, the goal of noun phrase coreference is to provide a clustering of the nouns such that each cluster refers to a single object.



 $(obj_{cccp} - obj_{spl})/obj_{cccp}$

Motif Finding

binary classification of DNA sequences

For all experiments we use C = 150 and ϵ = 0.001 (the large size of the dataset made cross-validation highly time consuming)

(a) Objective function value

CCCP	92.77 ± 0.99	106.50 ± 0.38	94.00 ± 0.53	116.63 ± 18.78	75.51 ± 1.97
SPL	92.37 ± 0.65	106.60 ± 0.30	93.51 ± 0.29	107.18 ± 1.48	74.23 ± 0.59
(b) Training error (%)					
CCCP	27.10 ± 0.44	32.03 ± 0.31	26.90 ± 0.28	34.89 ± 8.53	20.09 ± 0.81
SPL	26.94 ± 0.26	32.04 ± 0.23	26.81 ± 0.19	30.31 ± 1.14	19.52 ± 0.34
(c) Test error (%)					
CCCP	27.10 ± 0.36	32.15 ± 0.31	27.10 ± 0.37	35.42 ± 8.19	20.25 ± 0.65
SPL	27.08 ± 0.38	32.24 ± 0.25	27.03 ± 0.13	30.84 ± 1.38	19.65 ± 0.39

one could argue that difficult examples can be more informative than easy examples. Here the difficult examples are probably not useful because they confuse the learner rather than help it establish the right location of the decision surface.



