

8. Supervised learning, MaxMin estimate

Given training data $\mathcal{T}_e = \{(x^j, s^j) / j=1,..,l\}$ solve

$$\vec{u}_* = \underset{\vec{u}}{\operatorname{argmax}} \min_j \log p_{\vec{u}}(x^j, s^j)$$

We write the task as a convex optimisation task

$$\log Z(\vec{u}) - c \rightarrow \min_{\vec{u}, c}$$

$$\text{s.t. } c - \langle \vec{\varphi}^j, \vec{u} \rangle \leq 0 \quad \forall j=1,..,l \quad (1)$$

$$\text{where } \vec{\varphi}^j = \vec{\varphi}(x^j, s^j)$$

Construct the Lagrange function and the dual task

$$L(c, \vec{u}, \vec{\lambda}) = \log Z(\vec{u}) - c + \sum_{j=1}^l \lambda_j [c - \langle \vec{\varphi}^j, \vec{u} \rangle]$$

$$\max_{\vec{\lambda} \geq 0} \min_{\vec{u}, c} L(c, \vec{u}, \vec{\lambda}) =$$

$$\max_{\vec{\lambda} \in \Delta^{l-1}} \underbrace{\min_{\vec{u}} \left[\log Z(\vec{u}) - \sum_j \lambda_j \langle \vec{\varphi}^j, \vec{u} \rangle \right]}_{\text{MLE}}$$

Where Δ^{l-1} denotes the simplex $\Delta^{l-1} = \{\vec{\lambda} \in \mathbb{R}_+^l \mid \sum_j \lambda_j = 1\}$

Algorithm Choose arbitrary $\vec{\lambda}^{(0)} \in \Delta^{l-1}$ e.g. $\lambda_j^{(0)} = \frac{1}{l}$

Iterate:

1) Solve the MLE task (see sec. 7.) for $\vec{\lambda}^{(H)} \rightarrow \vec{u}^{(H)}$

2) find $i \in \operatorname{argmin}_j p_{\vec{u}^{(H)}}(x^j, s^j) = \operatorname{argmin}_j \langle \vec{\varphi}(x^j, s^j), \vec{u}^{(H)} \rangle$

Set $\vec{\lambda}^{(H+1)} \sim \cancel{\vec{\lambda}^{(H+1)}}(t+1) \vec{\lambda}^{(H)} + \vec{e}_i$

until

$$\sum_{j=1}^l \lambda_j \langle \bar{\varphi}^j, \bar{u} \rangle - \min_j \langle \bar{\varphi}^j, \bar{u} \rangle \leq \varepsilon$$

Theorem 1 (w/o proof)

The algorithm stops after a finite number of iterations

The obtained model \bar{u}_* is ε -optimal, i.e.

$$\min_j \log p_{\bar{u}_0}(x^j, s^j) - \min_j \log p_{\bar{u}_*}(x^j, s^j) \leq \varepsilon,$$

where \bar{u}_0 is the optimal model.

9. Empirical risk minimisation

Given:

- training data $\tilde{T}_c = \{(x^j, s^j) \mid j = 1, \dots, l\}$

- loss function $C(s', s) = \mathbb{I}\{s' \neq s\}$, i.e.

Bayes optimal decision $g_{\bar{u}}(x) \in \operatorname{argmax}_{s \in K^n} p_{\bar{u}}(x, s)$

The task of empirical risk minimisation reads

$$\frac{1}{l} \sum_{j=1}^l \mathbb{I}\{s^j \neq g_{\bar{u}}(x^j)\} \rightarrow \min_{\bar{u}}$$

It is not tractable in general; the objective function is neither convex nor differentiable.

It is simple to solve if $\exists \bar{u}_*$ s.t.

$$s^j = \operatorname{argmax}_{s \in K^n} p_{\bar{u}_*}(x^j, s) \quad \forall j = 1, \dots, l$$

i.e. we must find a \bar{u} s.t. inequalities

$$\langle \bar{\varphi}(x^j, s^j), \bar{u} \rangle > \langle \bar{\varphi}(x^j, s), \bar{u} \rangle \quad \forall s \neq s^j$$

hold for all $j = 1, \dots, l$

This task can be solved by the perceptron algorithm:

Start from arbitrary \vec{u} and iterate:

- solve

$$\tilde{s}^j = \operatorname{argmax}_{s \in K^n} \langle \vec{\varphi}(x^j, s), \vec{u} \rangle \quad (\text{see sec. 4})$$

$$+ j=1, \dots, l$$

- if for some j $\tilde{s}^j \neq s^j$, update \vec{u} by

$$\vec{u} \rightarrow \vec{u} + \vec{\varphi}(x^j, s^j) - \vec{\varphi}(x^j, \tilde{s}^j)$$

Let us reconsider the general task

$$q_{\vec{u}}(x) \in \operatorname{argmax}_{s \in K^n} \langle \vec{\varphi}(x, s), \vec{u} \rangle$$

$$\frac{1}{l} \sum_{j=1}^l \mathbb{1}\{s^j \neq q_{\vec{u}}(x^j)\} \rightarrow \min_{\vec{u}}$$

Approximate the loss (as a function of \vec{u}) by a convex upper bound, e.g., as follows

$$\mathbb{1}\{s' \neq q_{\vec{u}}(x)\} \leq \max_s \left\{ \mathbb{1}\{s' \neq s\} - \langle \vec{\varphi}(x, s'), \vec{u} \rangle + \langle \vec{\varphi}(x, s), \vec{u} \rangle \right\}$$

The approx. task reads

$$\frac{1}{l} \sum_{j=1}^l \max_{s \in K^n} \left[\mathbb{1}\{s' \neq s\} + \langle \vec{\varphi}(x^j, s) - \vec{\varphi}(x^j, s^j), \vec{u} \rangle \right] \rightarrow \min_{\vec{u}}$$

Solve it

- by subgradient descent or
- by cutting plane algorithm or
- ...

[GMW WS 14 / L 6]

10. Unsupervised learning, EM algorithm

Given: i.i.d. training data $T_e = \{x^j \in F^n \mid j=1 \dots e\}$

Task: $\vec{u}_* \in \operatorname{argmax}_{\vec{u}} \sum_{x \in F^n} \beta(x) \log \sum_{s \in K^n} p_{\vec{u}}(x, s)$

Substituting $p_{\vec{u}}$ we get

$$\underbrace{\log Z(\vec{u})}_{g(\vec{u})} - \underbrace{\sum_{x \in F^n} \beta(x) \log \sum_{s \in K^n} \exp \langle \vec{\varphi}(x, s), \vec{u} \rangle}_{h(\vec{u})} \rightarrow \min_{\vec{u}} \quad (1)$$

The objective function is a difference of convex functions
i.e. we have to solve a DC-program

A. DC-programs, DC-duality, DC-algorithm

Definition 1 The Fenchel conjugate of a function $g: \mathbb{R}^n \rightarrow [-\infty, +\infty]$ is the function $g^*: \mathbb{R}^n \rightarrow [-\infty, +\infty]$ defined by

$$g^*(\vec{v}) = \sup_{\vec{u} \in \mathbb{R}^n} \{ \langle \vec{v}, \vec{u} \rangle - g(\vec{u}) \}$$

- The function g^* is convex
- If g is convex and closed then

$$\vec{v} \in \partial g(\vec{u}) \Leftrightarrow \vec{u} \in \partial g^*(\vec{v}) \quad (2)$$

- Each DC program has a DC-dual program

$$g(\vec{u}) - h(\vec{u}) \rightarrow \min_{\vec{u}}$$

$$h^*(\vec{v}) - g^*(\vec{v}) \rightarrow \min_{\vec{v}}$$

The optimal values of both tasks coincide

- The DC-algorithm aims at solving both tasks simultaneously by constructing a pair of sequences $\vec{u}^{(t)}, \vec{v}^{(t)}$ in an alternating way
 - $\vec{v}^{(t)} \in \partial h(\vec{u}^{(t)})$
 - $\vec{u}^{(t+1)} \in \partial g^*(\vec{v}^{(t)})$

|B. DC-algorithm $\hat{=}$ EM algorithm for (1)

Applying DCA for (1) is a "reincarnation" of the EM-algorithm:

Choose an arbitrary $\vec{u}^{(0)}$. Iterate

- a) E-step: compute

$$\begin{aligned}\vec{v}^{(t)} &= \nabla h(\vec{u}^{(t)}) \\ &= \sum_x^1 \beta(x) \sum_s^1 \exp \langle \vec{\varphi}(x,s), \vec{u}^{(t)} \rangle \vec{\varphi}(x,s) / \sum_s^1 \exp \langle \vec{\varphi}(x,s), \vec{u}^{(t)} \rangle \\ &= \sum_{x,s}^1 \beta(x) p_{\vec{u}^{(t)}}(s|x) \vec{\varphi}(x,s)\end{aligned}$$

i.e. $\vec{v}^{(t)}$ is the vector of marginal statistics of the distr. $\beta(x) p_{\vec{u}^{(t)}}(s|x)$

For this it is necessary to compute pairwise posterior marginals $P_{\vec{U}^{(t)}}(s_{i-1}, s_i | x) \neq i=2, \dots, n$ for each $x \in \mathcal{X}_c$ (see sec. 5).

b) M-step: compute $\vec{u}^{(t+1)} \in \partial g^*(\vec{v}^{(t)})$

From Def. 1 and (2) we have

$$\vec{u}_* \in \arg \max_{\vec{u}} \underbrace{\left[\langle \vec{v}^{(t)}, \vec{u} \rangle - \log Z(\vec{u}) \right]}_{\text{MLE}} \Rightarrow \vec{u}_* \in \partial g^*(\vec{v}^{(t)})$$

But this task is easy to solve (see sec. 7)!

■

Theorem 1 (no proof)

- The sequences $g(\vec{u}^{(t)}) - h(\vec{u}^{(t)})$ and $h^*(\vec{v}^{(t)}) - g^*(\vec{v}^{(t)})$ are non-increasing
- The sequence $\vec{v}^{(t)}$ is convergent. Its fixpoint is a local minimum of $h^*(\vec{v}) - g^*(\vec{v})$. ■