

8. Supervised learning, MaxMin estimate

Given training data $\mathcal{T}_e = \{(x^j, s^j) \mid j=1, \dots, \ell\}$ solve

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

We write the task as a convex optimisation task

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

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

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

Construct the Lagrange function and the dual task

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

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

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

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

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

Iterate:

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

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

$$\text{set } \vec{\lambda}^{(t+1)} \sim \vec{\lambda}^{(t)} + \vec{e}_i$$

until

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

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 $\mathcal{T}_\ell = \{(x^j, s^j) \mid j = 1, \dots, \ell\}$

• loss function $C(s', s) = \mathbb{1}\{s' \neq s\}$, 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

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

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.

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

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 $j = 1, \dots, \ell$

This task can be solved by the perceptron algorithm:

Start from arbitrary \vec{u} and iterate:

- solve

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

$$\forall j=1, \dots, \ell$$

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

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

Let us reconsider the general task

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

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

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

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

The approx. task reads

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

Solve it

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