

## 8. Empirical risk minimisation

Given: training data (i.i.d.)  $T = \{(x_i^j, s_i^j) \mid x_i^j \in F, s_i^j \in K, j=1, \dots, \ell\}$   
 loss function  $C(s, s') = \mathbb{1}\{s \neq s'\}$

Recall: risk (average loss) for an inference strategy  $g: F \rightarrow K$

$$\begin{aligned} \sum_{x, s} p(x, s) \mathbb{1}\{s \neq g(x)\} &= \sum_{x, s} p(x, s) [1 - \delta(s, g(x))] \\ &= 1 - \sum_x p(x, g(x)) \rightarrow \min_g \end{aligned}$$

i.e., if the true model  $p_{\vec{u}}(x, s)$  is known  $\rightarrow$  optimal inference strategy is

$$g_{\vec{u}}(x) \in \operatorname{argmax}_{s \in K} p_{\vec{u}}(x, s)$$

But, the true model is not known  $\rightarrow$  replace risk by empirical risk

$$\sum_{x, s} p(x, s) \mathbb{1}\{s \neq g_{\vec{u}}(x)\} \approx \frac{1}{\ell} \sum_{(x, s) \in T} \mathbb{1}\{s \neq g_{\vec{u}}(x)\} \rightarrow \min_{\vec{u}}$$

This task is not tractable in general, the objective function is neither convex nor differentiable.

However, it is simple to solve for HMMs if there exists a  $\vec{u}^*$  such that

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

i.e. there  $\exists$  a  $\vec{u}^*$  for which the empirical risk is zero.

The above equations are equivalent to the following system of inequalities

$$\begin{aligned} \langle \vec{\varphi}(x_i^j, s_i^j), \vec{u}^* \rangle &> \langle \vec{\varphi}(x_i^j, s), \vec{u}^* \rangle && \forall s \neq s_i^j \\ &&& \forall j=1, 2, \dots, \ell \end{aligned}$$

Such a  $\vec{u}^*$  can be found by the perceptron algorithm (if it exists):

Start with arbitrary  $\vec{u}$  and iterate

- solve  $\tilde{s}^j = \operatorname{argmax}_{S \in K^n} \langle \vec{\varphi}(x^j, s), \vec{u} \rangle \quad \forall j=1, \dots, \ell$

for an HMM this can be done by the algorithm given in sec. 4 (dynamic programming)

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

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

Let us reconsider the general task of empirical risk minimisation

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

$$\frac{1}{\ell} \sum_{(x, s) \in \mathcal{T}} \mathbb{1}\{s \neq q_{\vec{u}}(x)\} \rightarrow \min_{\vec{u}}$$

Approximate the loss (as a function of  $\vec{u}$ ) by a convex upper bound. E.g. "margin rescaling" loss surrogate

$$\mathbb{1}\{s \neq q_{\vec{u}}(x)\} \leq \max_{S' \in K^n} \left[ \mathbb{1}\{s \neq S'\} + \langle \vec{\varphi}(x, S') - \vec{\varphi}(x, s), \vec{u} \rangle \right]$$

The approximation task reads

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

- Solve it by subgradient descent, cutting plane algorithm, ...
- The inner optimisation tasks  $\max_{S \in K^n} [\dots]$  are solved by the algorithm given in sec. 4

## 9. A sidestep: non-convex optimisation, DC-algorithm

Let  $\mathbb{E}$  denote a finite dimensional Euclidian space

- A function  $f: \mathbb{E} \rightarrow (-\infty, +\infty]$  is convex if its epigraph

$$\text{epi}(f) = \{(x, r) \in \mathbb{E} \times \mathbb{R} \mid f(x) \leq r\}$$

is a convex set.

- The domain of a function  $f: \mathbb{E} \rightarrow (-\infty, +\infty]$  is the set

$$\text{dom } f = \{x \in \mathbb{E} \mid f(x) < +\infty\}$$

- An element  $\varphi$  of  $\mathbb{E}$  is a subgradient of  $f$  in  $\bar{x}$  if it satisfies

$$\langle \varphi, x - \bar{x} \rangle \leq f(x) - f(\bar{x}) \quad \forall x \in \mathbb{E}$$

The set of subgradients of  $f$  in  $\bar{x}$  (called subdifferential) is denoted by  $\partial f(\bar{x})$ . ( $\partial f(\bar{x}) = \emptyset$  if  $\bar{x} \notin \text{dom } f$ )

- The Fenchel conjugate of a function  $f: \mathbb{E} \rightarrow [-\infty, +\infty]$  is the function  $f^*: \mathbb{E} \rightarrow [-\infty, +\infty]$  defined by

$$f^*(\varphi) = \sup_{x \in \mathbb{E}} \{ \langle \varphi, x \rangle - f(x) \}$$

The function  $f^*$  is convex. If  $f$  is convex and closed then

$$\varphi \in \partial f(x) \Leftrightarrow x \in \partial f^*(\varphi)$$

Let  $g, h: \mathbb{E} \rightarrow (-\infty, +\infty]$  be convex functions and consider the problem

$$g(x) - h(x) \rightarrow \min_{x \in \mathbb{E}}$$

We call such a problem DC-problem (difference of convex functions).

- Each DC-program has a DC-dual program

$$g(x) - h(x) \rightarrow \min_{x \in \mathbb{E}}$$

$$h^*(y) - g^*(y) \rightarrow \min_{y \in \mathbb{E}}$$

The optimal values of both tasks coincide

- The DC-algorithm aims at solving both task simultaneously by constructing a pair of decreasing sequences  $x^{(t)}, y^{(t)}$ :

$$a) \quad y^{(t)} \in \partial h(x^{(t)})$$

$$b) \quad x^{(t+1)} \in \partial g^*(y^{(t)})$$