## 15. Supervised parameter estimation for GRFs

### A. Generative learning

- $S = \{S_i \mid i \in V\}$ is a $K$-valued GRF on a graph $(V,E)$ with joint p.d.

$$P_u(s) = \frac{1}{Z(u)} \exp\left[ \sum_{i \in V} u_i(s_i) + \sum_{ij \in E} u_{ij}(s_i, s_j) \right]$$

- $T = \{s^j \in K^V \mid j = 1..,m\}$ is an i.i.d. training sample

<u>Task</u>: Estimate unary and pairwise potentials (i.e. model parameters) $u_i, u_{ij}$ from training data

Maximum likelihood estimator

$$L(u) = \frac{1}{m} \sum_{s \in T} \log \frac{1}{Z(u)} \exp U(s) \longrightarrow \max_u$$

Using the exponential family representation (Sec. 6), we get

$$L(u) = \frac{1}{m} \sum_{s \in T} \log \frac{1}{Z(u)} e^{\langle \Phi(s), u \rangle}$$

$$= \frac{1}{m} \sum_{s \in T} \langle \Phi(s), u \rangle - \log \sum_{s \in K^V} e^{\langle \Phi(s), u \rangle} \longrightarrow \max_u$$

The task has the structure $\langle \Psi, u \rangle - g(u) \longrightarrow \max_u$ with a convex function $g(u)$. Can we solve it by gradient ascent? Computing $\nabla g(u)$ requires to compute statistics of $\Phi(u)$, i.e. computing unary and pairwise marginal probabilities.
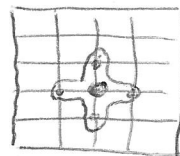
<u>Remark 1</u> The learning task is easy to solve for acyclic graphs $(V,E)$. See Sec. 7.

B. Pseudo-likelihood estimator (Besag, 1975)

Recall the Gibbs sampler (Sec. 14), which is defined by the conditional distributions

$$p(s_i \mid s_{N_i}), \quad i \in V, \quad s_i \in K$$

and in turn defines the joint p.d. $p(s)$

Idea    Use the pseudo-likelihood estimator defined by

$$L_p(u) = \frac{1}{m} \sum_{s \in T} \sum_{i \in V} \log p_u(s_i \mid s_{N_i}) \to \max_u$$

where

$$\log p_u(s_i \mid s_{N_i}) = \log \frac{\exp\left[u_i(s_i) + \sum_{j \in N_i} u_{ij}(s_i, s_j)\right]}{\sum_{s_i \in K} \exp\left[\ \_\_\ \text{''} \ \_\_\ \right]}$$

$$= u_i(s_i) + \sum_{j \in N_i} u_{ij}(s_i, s_j) - \log \sum_{k \in K} \exp\left[u_i(k) + \sum_{j \in N_i} u_{ij}(k, s_j)\right]$$

$\Rightarrow$ $L_p(u)$ is a concave function of $u$ and its gradient is easy to compute.

Theorem 1 (w/o proof)

The pseudo-likelihood estimator is asymptotically consistent for GRFs. However, it has a higher variance than MLE.  ▪

Remark 2

The pseudo-likelihood estimator can be easily generalised for GRF models as in C.

## C. Discriminative learning

- $X, S$ is a pair of $\mathfrak{X}$-valued and $K$-valued random fields on a graph $(V, E)$ with joint p.d.

$$P_u(x,s) = \frac{1}{Z[u]} \exp\left[\sum_{i \in V} u_i(x_i, s_i) + \sum_{ij \in E} u_{ij}(s_i, s_j)\right]$$

- loss function $l(s, s') = \sum_{i \in V} [\![s_i \neq s'_i]\!]$

- i.i.d. training data $T = \{(x^j, s^j) \mid x^j \in \mathfrak{X}^V, s^j \in K^V, j=1,..,m\}$

<u>Task</u> Estimate unary and pairwise potentials by minimising the empirical risk on training data.

$$R(u, T) = \frac{1}{m} \sum_{(x,s) \in \tilde{T}} l\left(s, \underset{s' \in K^V}{\operatorname{argmax}} \, p_u(x, s')\right)$$

$$= \frac{1}{m} \sum_{(x,s) \in T} l\left(s, \underset{s' \in K^V}{\operatorname{argmax}} \langle \varphi(x, s'), u \rangle\right) \to \min_u$$

The objective function is piecewise constant $\Rightarrow$ replace the true loss by a surrogate loss, e.g. margin rescaling loss

$$\tilde{R}(u, T) = \frac{1}{m} \sum_{(x,s) \in T} \max_{s' \in K^V} \left[l(s, s') - \langle \varphi(x, s), u \rangle + \langle \varphi(x, s'), u \rangle\right]$$

The objective upper bounds $R(u, T)$ and is convex in $u$. Computing its subgradient amounts to solve a $(\max, +)$ problem for each training example.

<u>Remark 3</u> The same approach can be applied for conditional random fields, where $p(s|x)$ is modelled as a GRF.