

Supervised learning of GRFsA. Maximum likelihood estimators

Given: A parametrised family of p.d.s $p_u(x)$, $x \in \mathcal{X}$, $u \in \mathcal{U}$ and a sample $\mathcal{T}^m = \{x^j \mid j=1, \dots, m\}$, where x^j are generated i.i.d. from $p_{u_0}(x)$ with unknown $u_0 \in \mathcal{U}$

Task: Estimate the unknown parameter u_0

Maximum likelihood estimator

$$u_* \in \operatorname{argmax}_{u \in \mathcal{U}} \frac{1}{m} \sum_{j=1}^m \log p_u(x^j) = \operatorname{argmax}_{u \in \mathcal{U}} L(u, \mathcal{T}^m)$$

Consistency of MLEs

$$\begin{array}{ccc} L(u, \mathcal{T}^m) & \xrightarrow[m \rightarrow \infty]{\mathbb{P}} & L(u) = \sum_{x \in \mathcal{X}} p_{u_0}(x) \log p_u(x) \\ \operatorname{argmax}_{u \in \mathcal{U}} \downarrow & & \downarrow \operatorname{argmax}_{u \in \mathcal{U}} \\ u_*(\mathcal{T}^m) & \xrightarrow[m \rightarrow \infty]{\mathbb{P}} & u_0 \end{array}$$

a) $L(u, \mathcal{T}^m) \xrightarrow[m \rightarrow \infty]{\mathbb{P}} L(u)$? Yes, LLN

$$\mathbb{P}(|L(u, \mathcal{T}^m) - L(u)| \geq \varepsilon) \xrightarrow[m \rightarrow \infty]{} 0 \quad \forall u \in \mathcal{U}, \forall \varepsilon > 0$$

b) $u_0 = \operatorname{argmax}_{u \in \mathcal{U}} L(u)$? Yes, if the model is identifiable,

i.e. $p_{u_0}(x) \neq p_u(x) \quad \forall u \neq u_0$. This follows from

$$D_{\text{KL}}(p \parallel q) \geq 0 \quad \text{with equality iff } p = q$$

c) $u_*(\tilde{T}^m) \xrightarrow[m \rightarrow \infty]{P} u_0$? Yes, if

1) \mathcal{U} is compact and $L(u, \tilde{T}^m)$ is continuous in u .

2) $L(u, \tilde{T}^m)$ converges uniformly, i.e.

$$P\left(\sup_{u \in \mathcal{U}} |L(u, \tilde{T}^m) - L(u)| \geq \varepsilon\right) \xrightarrow[m \rightarrow \infty]{} 0$$

Remark 1 Conditions 1) & 2) are sufficient but not necessary.

E.g. condition 1) can be replaced by $L(u, \tilde{T}^m)$ is concave in u and $u_0 \in \text{int}(\mathcal{U})$.

B. MLEs for Gibbs random fields

$S = \{S_i \mid i \in V\}$ is a GRF w.r.t. the graph (V, E) and has p.d.

$$p_u(s) = \frac{1}{Z(u)} \exp \sum_{ij \in E} u_{ij}(s_i, s_j)$$

The parameters u_{ij} are unknown. We are given an i.i.d sample of training data $\tilde{T}^m = \{s^j \in K^V \mid j=1, \dots, m\}$.

Since the GRF is an exponential family, we can write

$$\begin{aligned} L(u, \tilde{T}^m) &= \frac{1}{m} \sum_{s^j \in \tilde{T}^m} \log p_u(s^j) = \\ &= \frac{1}{m} \sum_{j=1}^m \langle \Phi(s^j), u \rangle - \log Z(u) \end{aligned}$$

- $L(u, \tilde{T}^m)$ is concave in u
- The model is identifiable up to re-parametrisations
- If \mathcal{U} is compact, then the MLE is consistent

How to solve the task

$$L(u, T^m) = \frac{1}{m} \sum_{j=1}^m \langle \Phi(s^j), u \rangle - \log Z(u) \rightarrow \max_{u \in \mathcal{U}}$$

for a GRF? Gradient ascend requires to compute

$$\nabla L(u, T^m) = \frac{1}{m} \sum_{j=1}^m \Phi(s^j) - \mathbb{E}_u[\Phi].$$

Computing $\mathbb{E}_u[\Phi]$ amounts to computing pairwise marginals for all edges of the graph. Notice, that the marginals must be (re-)estimated in each iteration of the gradient ascend.

CG Pseudo-likelihood estimators for GRFs

Can we do simpler? Besag, 1975 \Rightarrow Recall that a GRF on a graph (V, E) is defined by fixing the family of conditional distr.

$$p_u(s_i | S_{N_i}) = \frac{1}{Z_i(u, S_{N_i})} \exp \sum_{j \in N_i} u_{ij}(s_i, s_j)$$

(see Gibbs sampler). Hence, we may use pseudo-likelihood

$$\tilde{L}(u, T^m) = \frac{1}{m} \sum_{S \in T^m} \sum_{i \in V} \log p_u(s_i | S_{N_i}) \rightarrow \max_u$$

instead of MLE. We obtain

$$\begin{aligned} \tilde{L}(u, T^m) &= \frac{1}{m} \sum_{S \in T^m} \sum_{i \in V} \sum_{j \in N_i} u_{ij}(s_i, s_j) - \frac{1}{m} \sum_{S \in T^m} \sum_{i \in V} \log Z_i(u, S_{N_i}) \\ &= 2 \sum_{j \in E} \frac{1}{m} \sum_{S \in T^m} u_{ij}(s_i, s_j) - \sum_{i \in V} \frac{1}{m} \sum_{S \in T^m} \log Z_i(u, S_{N_i}) \end{aligned}$$

• Computing $\tilde{L}(u, T^m)$ and $\nabla \tilde{L}(u, T^m)$ has complexity $\mathcal{O}(m|E||K|^2)$

- $\tilde{L}(y, T^m)$ is concave. It can be proved that pseudo-likelihood estimators are consistent.
- Its variance is higher than the variance of MLE
- Pseudo-likelihood estimators can not be used for unsupervised learning, whereas MLE can (by EM-algorithm).