

Локальная информационная геометрия двоичных цепей Маркова высокого порядка и ее приложения

В.А. Волошко

Пусть $\mathcal{MC}(s)$ – семейство Марковских стационарных вероятностных мер (далее, ВМ) порядка $s \in \mathbb{N}_0 = \{0, 1, \dots\}$ на множестве двоичных последовательностей $\{0, 1\}^{\mathbb{Z}}$, $\mathcal{MC} ::= \cup_{s \in \mathbb{N}_0} \mathcal{MC}(s)$ – семейство Марковских стационарных ВМ конечного порядка, \mathbf{U} – равномерная ВМ на $\{0, 1\}^{\mathbb{Z}}$. Очевидно:

$$\mathbf{U} \in \mathcal{MC}(0) \subset \mathcal{MC}(1) \subset \mathcal{MC}(2) \subset \dots \subset \mathcal{MC}.$$

Поэтому касательные пространства $\mathcal{TMC}(s)$ точки \mathbf{U} в многообразиях $\mathcal{MC}(s)$ ($\dim(\mathcal{TMC}(s)) = 2^s$) также образуют вложенную систему:

$$\mathcal{TMC}(0) \subset \mathcal{TMC}(1) \subset \mathcal{TMC}(2) \subset \dots \subset \mathcal{TMC}, \quad \mathcal{TMC} ::= \cup_{s \in \mathbb{N}_0} \mathcal{TMC}(s).$$

Локальную информационную геометрию двоичных цепей Маркова высокого порядка в окрестности равномерной ВМ \mathbf{U} определяет тензор Фишера-Римана – скалярное произведение на бесконечномерном пространстве \mathcal{TMC} :

$$\langle \tau_1, \tau_2 \rangle_{\mathcal{TMC}} \in \mathbb{R}, \quad \tau_1, \tau_2 \in \mathcal{TMC}.$$

Пусть теперь задана функция (информативный признак) $f : \{0, 1\}^r \rightarrow \mathbb{R}^d$ от двоичных r -слов. Будем рассматривать статистики сумм информативных признаков по пересекающимся r -граммам наблюдаемой двоичной последовательности $x_1, \dots, x_T \in \{0, 1\}$:

$$\langle f \rangle_x = \sum_{t=1}^{T-r+1} f(x_t, x_{t+1}, \dots, x_{t+r-1}). \quad (1)$$

Статистика (1) для гипотезы “ $H_0: \{x_t\}$ – РРСП” имеет асимптотическое нормальное распределение с матрицей ковариаций $T \cdot (\Sigma_f + o(1))$, в случае невырожденности которой может быть построена хи-квадрат статистика:

$$S_f(x) = T^{-1} (\langle f \rangle_x - \mathbf{E}_{\mathbf{U}}\{f\})' \Sigma_f^{-1} (\langle f \rangle_x - \mathbf{E}_{\mathbf{U}}\{f\}), \quad \mathbf{E}_{\mathbf{U}}\{f\} = 2^{-r} \sum_{q \in \{0,1\}^r} f(q). \quad (2)$$

Статистики вида (2) неоднозначно задаются информативными признаками f , но однозначно кодируются конечномерными подпространствами $\mathcal{T} \subset \mathcal{TMC}$, а взаимные свойства таких подпространств отображаются на асимптотические свойства взаимных распределений статистик вида (2). В частности, ортогональность подпространств влечет асимптотическую независимость отвечающих им статистик (2) при истинной гипотезе H_0 .