

Rozdělení vzdálenosti mezi body z multinomického rozdělení

Šárka Horská

11.3.2020

Struktura prezentace

- motivace
- základní pojmy
- odvozené rozdělení
- odvozené vlastnosti tohoto rozdělení
- nalezení nestranného odhadu $\|P_x - P_y\|^2$

- porovnání dvou multinomických rozdělení
- dvouvýběrový χ^2 - test je vhodný jen tehdy, když máme dostatečně velký počet pozorování v každé složce vektoru
- tento problém nemá např. dvouvýběrový ED (Euclidean distance) test
- ED test má za testovou statistikou odhad $\|\mathbf{P}_x - \mathbf{P}_y\|^2$
- chceme nalézt nestranný odhad

Základní pojmy

Multinomické rozdělení

Ať $\mathbf{p} = (p_1, \dots, p_K)^T$ je vektor konstant splňující $\sum_{k=1}^K p_k = 1$ a $p_k > 0$ pro $k \in \{1, \dots, K\}$. Pak $\mathbf{X} \sim \text{Mult}_K(n, \mathbf{p}) \iff$ jeho hustota je dána vztahem

$$P[X_1 = x_1, X_2 = x_2, \dots, X_K = x_K] = \begin{cases} \frac{n!}{x_1! \cdots x_K!} & \text{pokud } \sum_{k=1}^K x_k = n \\ & \text{a } x_k \in \mathbb{N}_0 \text{ pro } \forall k \\ 0 & \text{jinak.} \end{cases}$$

- Vzdáleností bodů \mathbf{X}_i a \mathbf{X}_j rozumíme eukleidovskou metriku

$$D_{ij}^2 = \sum_{k=1}^K (X_{ik} - X_{jk})^2 = \sum_{k=1}^K T_{ij}^2(k)$$

Odvození rozdělení vzdáleností

Víme, že pro $\mathbf{X}_i, \mathbf{X}_j \sim \text{Mult}_K(n, \mathbf{p})$, nezávislé platí:

- $X_{ik}, X_{jk} \sim \text{Bi}(n, p_k)$
- $X_{ik}, X_{jk} \in (0, \dots, n)$
- $P(X_{ik} = t) = \binom{n}{t} p_k^t (1 - p_k)^{n-t}$
- X_{ik}, X_{jk} jsou nezávislé
- $T_{ij}(k) = X_{ik} - X_{jk} \in (-n, \dots, n)$

Odvození rozdělení vzdáleností

Odvodíme postupně následující rozdělení:

$$P(T_{ij}(k) = t_k) = \sum_{m=0}^n P(X_{ik} = t_k + m)P(X_{jk} = m)$$

$$P(T_{ij}^2(k) = t_k^2) = \begin{cases} P(T_{ij}(k) = t_k) + P(T_{ij}(k) = -t_k) = 2P(T_{ij}(k) = t_k) \\ P(T_{ij}(k) = 0) \end{cases}$$

Odvození rozdělení vzdáleností

- Definujeme množinu

$$I = \{t = t_1^2 + \dots + t_K^2 : t_i \in \{0, \dots, n\} \text{ pro } i \in \{1, \dots, K\}\}$$

- Hledané rozdělení je tvaru

$$P(d_{ij}^2 = t) = \sum_{t \in I} \prod_{k=1}^K (1 + \delta_k) \sum_{m=0}^n P(X_{ik} = t_k + m) P(X_{jk} = m),$$

kde $\delta_k = 1$ pro $t_k > 0$ a $\delta_k = 0$ jinak.

Základní vlastnosti rozdělení

- $E[d_{ij}^2] = E[\sum_{k=1}^K (X_{ik} - X_{jk})^2] =$
 $\sum_{k=1}^K (E(X_{ik}^2) - 2E(X_{ik}X_{jk}) + E(X_{jk}^2)) =$
 $2 \sum_{k=1}^K (E(X_{ik}^2) - E(X_{ik})^2) = 2 \sum_{k=1}^K \text{Var}(X_{ik})$
- vyjádření rozptylu vzdálenosti je velice nepřehledné

Odhad $\|\mathbf{P}_x - \mathbf{P}_y\|^2$

Ať $\mathbf{X} = \{\mathbf{X}_i\} \sim \text{Mult}_K(n, \mathbf{P}_x)$ pro $i = 1, \dots, N_x$ jsou nezávislé,
 $\mathbf{Y} = \{\mathbf{Y}_i\} \sim \text{Mult}_K(n, \mathbf{P}_y)$ pro $i = 1, \dots, N_y$ jsou nezávislé a \mathbf{X}, \mathbf{Y} jsou
vzájemně nezávislé.

- $\hat{P}_{(x)k} = \frac{1}{nN_x} \sum_{i=1}^{N_x} X_{ik}$
- tento odhad je nestranný a konzistentní
- při výpočtu $E[\|\mathbf{P}_x - \mathbf{P}_y\|^2]$ zjistíme, že není nestranný

Odhad $\|\mathbf{P}_x - \mathbf{P}_y\|^2$

- hledané c splňující $E[\|\hat{\mathbf{P}}_x - \hat{\mathbf{P}}_y\|^2 - c] = \|\mathbf{P}_x - \mathbf{P}_y\|^2$ je tvaru

$$c = \frac{1}{n^2 N_x} \sum_{k=1}^K \hat{V}ar X_{1k} + \frac{1}{n^2 N_y} \sum_{k=1}^K \hat{V}ar Y_{1k}$$

- testová statistika je tvaru

$$ED = \|\hat{\mathbf{P}}_x - \hat{\mathbf{P}}_y\|^2 - \frac{1}{n^2 N_x} \sum_{k=1}^K \hat{V}ar X_{1k} - \frac{1}{n^2 N_y} \sum_{k=1}^K \hat{V}ar Y_{1k}$$

Svou prezentaci bych časem ráda doplnila o:

- porovnání ED testových statistik s neupravenou statistikou a nestrannou statistikou
- nalézt více důvodů, proč se o rozdělení vzdálenosti zajímat
- doplnit některé výpočty
- doplnit vlastnosti rozdělení

Děkuji za pozornost!