

کنترل کردن سویدگی در استفاده کردن از اطلاعات

ارائه‌ی پروژه‌ی درس تئوری اطلاعات

بهراد منیری محمد رضا رحمانی

استاد درس
دکتر میرمحسنی

دستیار آموزشی
امیرحسین بساره

دانشکده‌ی مهندسی برق
دانشگاه صنعتی شریف

فهرست

۱ معرفی

۲ انگیزه

۳ مدل ریاضی

۴ کاربردهای چهارچوب معرفی شده

۵ تصادفی سازی

۶ پیشنهادها

How much does your data exploration overfit? Controlling bias via information usage

Daniel Russo¹, James Zou²

IEEE Transactions on Information Theory, 2019

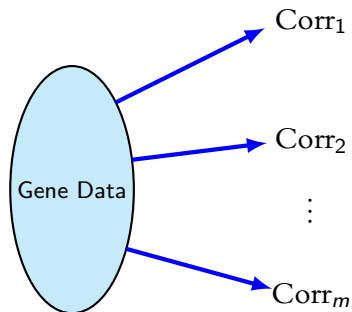
¹The Division of Decision, Risk and Operations, Columbia University

²Biomedical Data Science, Computer Science and Electrical
Engineering at Stanford University

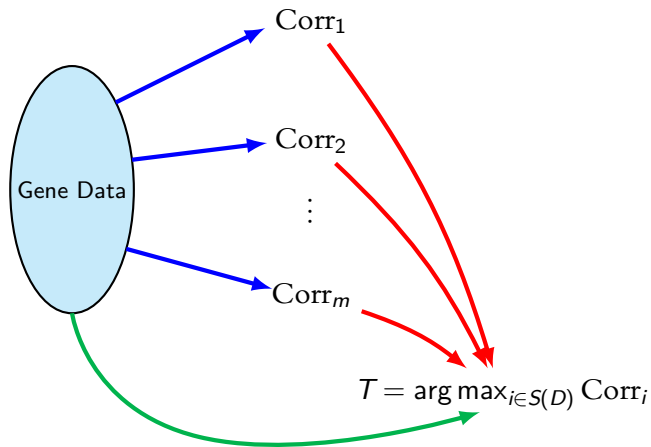
انگیزه



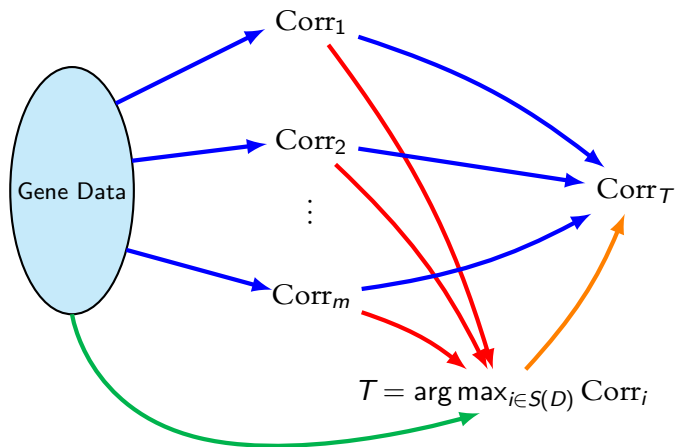
انگیزه



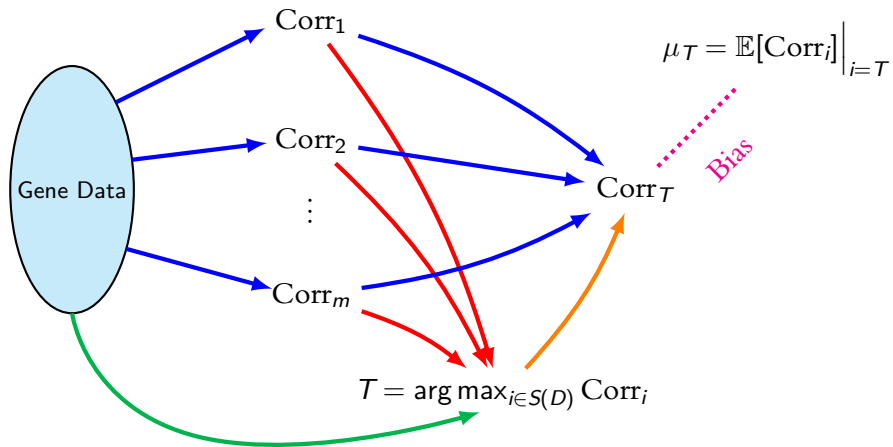
انگیزه



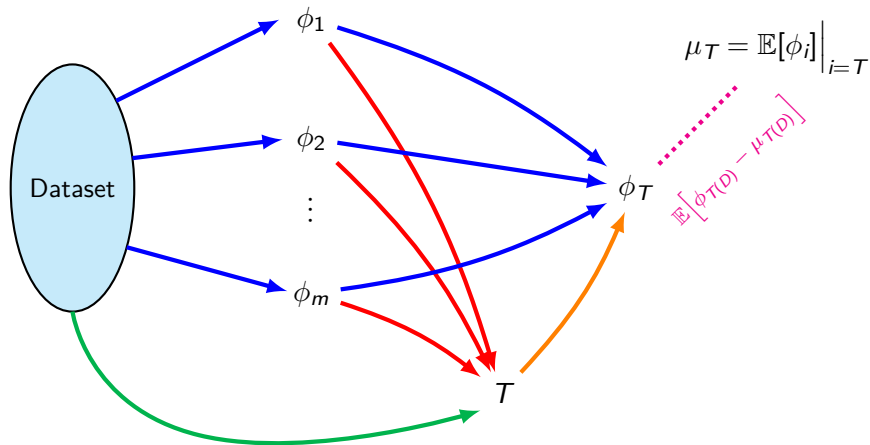
انگیزه



انگیزه



انگیزه



قضیه

متغیر تصادفی برداری $\phi = (\phi_1, \dots, \phi_m)$ را در نظر بگیرید، اگر تعریف کنیم

$$\mu = (\mu_1, \mu_2, \dots, \mu_m) = \mathbb{E}[\phi]$$

و اگر برای هر $\phi_i - \mu_i, i \in \{1, \dots, m\}$ یک متغیر تصادفی زیرگوسی با پارامتر σ باشد، آنگاه:

$$\left| \mathbb{E}[\phi_T - \mu_T] \right| \leq \sigma \sqrt{2I(T; \phi)} \quad (1)$$

$$\mathbb{E} \left[|\phi_T - \mu_T| \right] \leq \sigma + 36\sigma \sqrt{2I(T; \phi)} \quad (2)$$

$$\mathbb{E} [(\phi_T - \mu_T)^2] \leq 1.25\sigma^2 + 10\sigma^2 I(T; \phi) \quad (3)$$

رتبه‌بندی با وجود سیگنال در داده‌ها

◀ مدل تولید داده‌ها:

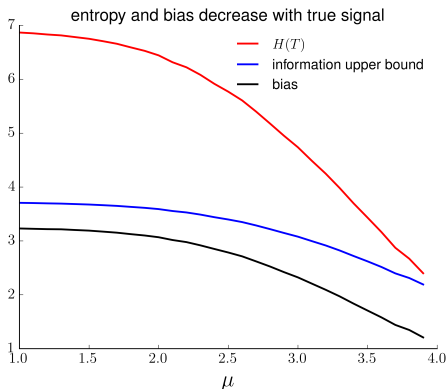
$$\phi_i \sim \begin{cases} \mathcal{N}(\mu, \sigma^2) & i = I^* \\ \mathcal{N}(0, \sigma^2) & i \neq I^* \end{cases}$$

◀ انتخاب:

$$T = \arg \max_i \phi_i$$

◀ کران مبتنی بر آنتروپی:

$$I(T; \phi) = H(T) \leq \log(m)$$



تفکیک تقریباً مستقل داده‌های غیر i.i.d.

$$\underbrace{s_1 - \dots - s_{n_1} - s_{n_1+1} - \dots - s_{n_2} - s_{n_2+1} - \dots - s_t}_{\text{انتخاب تحلیل}} \quad \underbrace{\hspace{10em}}_{\text{تخمین}}$$

$$\forall \tau \in \mathbb{N} \quad \max_s D\left(\mathbf{P}(s_\tau = \cdot | s_1 = s) || \pi\right) \leq c_0 e^{-c_1 \tau}$$

$$T \leftarrow \{s_1, \dots, s_{n_1}\} - s_{n_1} - s_{n_2+1} - \{s_{n_2+1}, \dots, s_t\} \rightarrow \phi$$

$$\begin{aligned} I(T; \phi) &\leq I(s_{n_1}; s_{n_2+1}) \\ &\leq c_0 e^{-c_1(n_2 - n_1)} \end{aligned}$$

انتخاب تصادفی اندیس

◀ چرا تصادفی سازی خوب است؟

$$I(T; \phi) = H(T) - H(T|\phi) = H(T) - H(\pi)$$

◀ در مسئله‌ی انتخاب ماکزیمم چه می توان کرد؟

$$\underset{\pi \in \mathbb{R}_+^m}{\text{maximize}} \quad H(\pi)$$

$$\text{subject to} \quad \sum_{i=1}^k \pi_i \phi_i \geq b \text{ and } \sum_{i=1}^k \pi_i = 1.$$

پاسخ این مسئله‌ی بهینه سازی:

$$\pi^* = A e^{\beta \phi_i}$$

انتخاب تصادفی اندیس

◀ نحوه تولید داده‌ها:

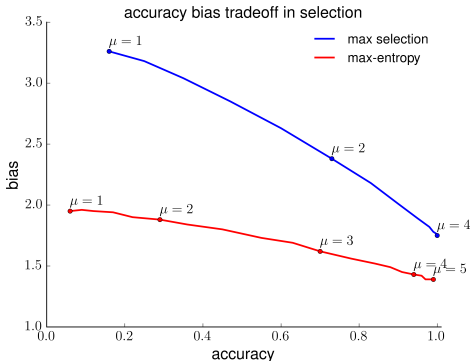
$$\begin{cases} \mu_i = \mu & i \leq N_1 \\ \mu_i = 0 & N_1 \leq i \leq N \end{cases}$$

سویدگی:

$$\frac{1}{K} \sum_{i=1}^K \phi_{T_i} - \mu_{T_i}$$

دقت:

$$\frac{|\{T_i : T_i \leq N_1\}|}{K}$$



مدل تحلیل تطبیقی داده‌ها

- در قدم اول، تحلیل ϕ_{T_1} انتخاب می‌شود. نتیجه‌ی این تحلیل $Y_{T_1} \in \mathbb{R}$ است.
- در تکرار k ام:

$$T_k = f(Y_{T_1}, Y_{T_2}, \dots, Y_{T_{k-1}}, T_1, \dots, T_{k-1}) \Rightarrow \phi_{T_k} \Rightarrow Y_{T_k}$$

گزاره

سویدگی Y_{T_k} دو جزء دارد:

$$\begin{aligned} & \mathbb{E}[|Y_{T_k} - \mu_{T_k}|] - \sigma \\ & \leq \mathbb{E}[|Y_{T_k} - \phi_{T_k}|] + \mathbb{E}[|\phi_{T_k} - \mu_{T_k}| - \sigma] \\ & \leq \underbrace{\mathbb{E}[|Y_{T_k} - \phi_{T_k}|]}_{\text{اعوجاج}} + \underbrace{c\sigma\sqrt{2l(T_k; \phi)}}_{\text{سویدگی ناشی از انتخاب}}. \end{aligned}$$

مدل تحلیل تطبیقی داده‌ها

$$T_{k+1} \leftarrow H_k = \{T_1, Y_{T_1}, T_2, Y_{T_2}, \dots, T_k, Y_{T_k}\} \leftarrow D \rightarrow \phi$$

لم

$$I(T_{k+1}; \phi) \leq I(H_k; \phi) = \sum_{i=1}^k I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i)$$

مدل تحلیل تطبیقی داده‌ها

$$T_{k+1} \leftarrow H_k = \{T_1, Y_{T_1}, T_2, Y_{T_2}, \dots, T_k, Y_{T_k}\} \leftarrow D \rightarrow \phi$$

لم

$$I(T_{k+1}; \phi) \leq I(H_k; \phi) = \sum_{i=1}^k I(Y_{T_i}; \phi_{T_i} | H_{i-1}, T_i)$$

ایده

از آنجایی که باید اطلاعات متقابل Y_{T_k} و ϕ_{T_k} را کنترل کنیم:

$$Y_{T_i} = \phi_{T_i} + \text{نویز}$$

پیش نهادها

- انتخاب تعدادی تصادفی از تحلیل ها به جای تنها یک تحلیل
- یافتن توزیع نوین بهینه در مسئله‌ی تحلیل تطبیقی برای کمترین سویدگی
- بررسی الگوریتم‌های آنلاین دیگر در چهارچوب اطلاعات استفاده شده
- مدل کردن دیتاست به صورت یک فرآیند تصادفی

مراجع



Russo, D. and J. Zou (2016).

Controlling bias in adaptive data analysis using information theory.
In *Artificial Intelligence and Statistics*, pp. 1232–1240.



Russo, D. and J. Zou (2019).

How much does your data exploration overfit? controlling bias via
information usage.
IEEE Transactions on Information Theory 66(1), 302–323.



Wainwright, M. J. (2019).

High-dimensional statistics: A non-asymptotic viewpoint, Volume 48.
Cambridge University Press.

با تشکر