

۱. برای مجموعه $T \subseteq \mathbb{R}^n$ و بردارهای تصادفی مستقل گوسی استاندارد $G_j, j = 1, \dots, n$ ، نشان دهید

$$\mathbb{E} \left[\sup_{v \in T} \sum_{j=1}^n G_j |v_j| \right] \leq \mathbb{E} \left[\sup_{v \in T} \sum_{j=1}^n G_j v_j \right].$$

تفسیری هندسی از نامساوی فوق ارائه دهید.

۲. (آ) بعد VC خانواده همه کره‌های سه بعدی را بیابید.

(ب) مجموعه نقاط

$$\{(x_1, x_2, \dots, x_n) : x_i \in \{-1, +1\}\}$$

رئوس یک مکعب n بعدی را تشکیل می‌دهند. طول ضلع مربع برابر ۲ می‌باشد. حجم این مکعب n بعدی را محاسبه کنید. حال فرض کنید که به مرکز هر کدام از رئوس یک کره n بعدی به شعاع یک بزنیم. این کره‌ها به همدیگر مماس خواهند شد. بنابراین 2^n کره خواهیم داشت. حال کره‌ای جدید به مرکز مبدأ و با شعاع R در نظر بگیرید که به تمامی کره قبلی مماس باشد. شعاع این کره چیست؟ آیا این کره کاملاً داخل مکعب قرار می‌گیرد؟! نسبت حجم این کره به حجم مکعب برای ابعاد بالا چیست؟ برای سادگی فرض کنید n زوج است. در این حالت حجم کره به شعاع r برابر با $r^n \frac{\pi^{n/2}}{(n/2)!}$ است. ضمناً برای تخمین فاکتوریل می‌توانید از تقریب استرلینگ یا نامساوی ساده زیر استفاده کنید

$$\log n! < (n+1) \log(n+1) - n.$$

۳. در این مسئله به دنبال فشردن داده‌ها هستیم به گونه‌ای که عملیات آماری یا یادگیری روی داده‌های فشرده شده تقریباً مشابه حالت داده‌های خام باشد.

(آ) ابتدا مسئله زیر را در نظر بگیرید. فرض کنید Z متغیر تصادفی با دامنه محدود باشد یعنی $|Z| \leq B$. همچنین فرض کنید

$$U \text{ و } \gamma > B \text{ یک متغیر تصادفی یکنواخت روی } [-\gamma, \gamma] \text{ باشد. قرار دهید}$$

$$\hat{Z} = \text{sign}(U + Z)$$

در حقیقت \hat{Z} فشرده شده Z توسط یک بیت است. نشان دهید

$$\mathbb{E}[\gamma \hat{Z}] = \mathbb{E}[Z]$$

(ب) مجموعه داده زیر در اختیار ما گذاشته شده است:

$$\mathcal{D} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$$

که در آن $\mathbf{X}_k \in \mathbb{R}^d$ و $Y_k \in \mathbb{R}$ هستند. فرض کنید که رابطه بین بردار ویژگی \mathbf{X} و متغیر Y (که مجموعه داده بالا هم از روی این توزیع تولید شده‌اند) توسط رابطه خطی نویزی شده زیر توصیف می‌شود:

$$Y = \Theta^\top \mathbf{X} + \epsilon$$

فرض کنید که متغیر Y و نویز ϵ هر دو زیرگوسی با پارامتر محدود $O(1)$ باشند. به جای ذخیره کل داده علاقه مند به فشردن آن هستیم. فعلاً برای سادگی فقط Y_k را فشرده می‌کنیم. بدین منظور از فرآیند فشردن سازی قسمت قبل استفاده می‌کنیم. بدین منظور برای هر $k \in \{1, 2, \dots, n\}$ ، $f(n)$ متغیر یکنواخت $(U_{k,1}, \dots, U_{k,f(n)})$ از $[-\gamma, \gamma]$ انتخاب و از روی آن‌ها

$$\hat{Y}_{k,j} = \text{sign}(Y_k + U_{k,j}), \quad j \in \{1, \dots, f(n)\}$$

می‌سازیم. دقت کنید که این فرآیند معادل با فشردن کردن Y_k با $f(n)$ بیت است. بدین ترتیب به مجموعه داده فشرده شده زیر می‌رسیم:

$$\mathcal{D}_c = \{(\mathbf{X}_k, \hat{Y}_{k,j}), k \in \{1, \dots, n\}, j \in \{1, \dots, f(n)\}\}$$

با توجه به قسمت قبل، سعی کنید که ابتدا تقریبی از داده اصلی \mathcal{D} را از روی داده فشرده شده \mathcal{D}_c بدست آورید. داده تقریبی شما باید به فرم زیر باشد

$$\mathcal{D}_{appr} = \{(\mathbf{X}_1, \tilde{Y}_1), \dots, (\mathbf{X}_n, \tilde{Y}_n)\}$$

که در آن \tilde{Y}_k باید از روی داده‌های فشرده شده محاسبه شود. سپس مقادیر γ و $f(n)$ را به گونه‌ای انتخاب کنید به طوریکه رابطه زیر برقرار باشد:

$$\tilde{Y}_k = \Theta^\top \mathbf{X}_k + \xi_k$$

به طوری که نویز ξ_k هم مشابه مدل اصلی زیرگوسی از مرتبه $O(1)$ باشد. در نهایت چند بیت برای فشردن بدون از دست دادن خواص آماری مسئله اصلی لازم است.

۴. سبد سهام را در کارهای سرمایه‌گذاری و اقتصاد و مدیریت حتماً تاکنون شنیده‌اید. در این مساله سعی می‌کنیم مساله انتخاب سبد سهام را به صورت ریاضی مدل و در مورد آن تحقیق کنیم. فرض کنید یک سبد سهام با یک بردار $u \in \mathbb{R}^d$ که $\sum_{i=1}^d u_i = 1$ و $u_j \geq 0$ برای هر $i \in [1 : d]$ مدل می‌شود که در آن u_j میزان سرمایه‌گذاری فرد در دارایی j ام را نشان می‌دهد. بردار (تصادفی) بازدهی سرمایه‌گذاری فرد را با $X \in \mathbb{R}^d$ نشان می‌دهیم و فرض می‌کنیم که X یک بردار زیرگوسی با پارامتر واریانس Σ و ماتریس کواریانس نامشخص Σ است (اطلاعاتی در مورد آن نداریم!).

پاداش و ریسک را به این شکل مدل می‌کنیم: پاداش: $\mu(u) = \mathbb{E}[u^\top X]$ و ریسک: $R(u) = \text{Var}(X^\top u)$. مساله زیر را

در نظر بگیرید:

$$u^* = \arg \min_{u: \mu(u) \geq \lambda} R(u)$$

ابتدا بگویید در این مساله چه استراتژی برای سوددهی در نظر گرفته شده است؟ در ادامه فرض کنید این مساله جواب دارد، در عمل توزیع بردار X نامشخص است. فرض کنید توانسته‌ایم نمونه‌های مستقل و هم توزیع X_1, \dots, X_n از روی X را بدست بیاوریم و تقریبهای زیر را برای $\mu(u), R(u)$ زده‌ایم:

$$\hat{\mu}(u) = \bar{X}^\top u = \frac{1}{n} \sum_{i=1}^n X_i^\top u,$$

$$\hat{R}(u) = u^\top \hat{\Sigma} u, \quad \hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{X})(X_i - \hat{X})^\top$$

ما از سبد سهام تقریبی زیر استفاده می‌کنیم:

$$\hat{u} = \arg \min_{u: \hat{\mu}(u) \geq \lambda} \hat{R}(u)$$

همچنین فرض کنید تعداد نمونه‌ها در $d \ll n \ll \log d$ صدق می‌کند.

(آ) نشان دهید برای هر u ,

$$|\hat{\mu}(u) - \mu(u)| \lesssim \frac{1}{\sqrt{n}}$$

و همچنین

$$|\hat{R}(u) - R(u)| \lesssim \frac{1}{\sqrt{n}}$$

(ب) نشان دهید

$$\hat{R}(\hat{u}) - R(\hat{u}) \lesssim \sqrt{\frac{\log d}{n}}$$

(ج) تخمینگر زیر را در نظر بگیرید:

$$\tilde{u} = \arg \min_{u: \hat{\mu}(u) \geq \lambda - \epsilon} \hat{R}(u)$$

کوچکترین ϵ مثبت را پیدا کنید (ضریب ثابت را در نظر نگیرید!) که با احتمال 0.99

$$R(\tilde{u}) \leq R(u^*).$$

۵. خانواده زیر از ماتریسهای مرتبه یک را در نظر بگیرید:

$$\mathbb{M}^{(n,d)}(1) := \{A \in \mathbb{R}^{n \times d} | \text{rank}(A) = 1, \|A\|_F = 1\}$$

کران پایین زیر را برای عدد پکینگ این مجموعه اثبات نمایید:

$$\log \mathcal{P} \left(\mathbb{M}^{(n,d)}(\mathfrak{V}), \|\cdot\|_F, \delta \right) \geq c(n+d) \log \frac{1}{\delta}$$