



## گزارش پروژه

دانشکده مهندسی کامپیوتر

نام و نام خانوادگی دانشجو: سید احسان حسن بیگی

استاد: دکتر نجفی

پاییز ۱۴۰۳

## فهرست مطالب

|  |    |
|--|----|
| ۱. مقدمه   | 3  |
| ۲. پیش‌نیازها و تعاریف   | 4  |
| ۳. روش‌های پیشین   | 5  |
| ۴. روش ارائه شده   | 6  |
| ۴.۱. مفهوم sample compressibility                              | 6  |
| ۴.۲. محاسبه‌ی sample compressibility برای چند توزیع نمونه      | 6  |
| ۴.۳. قابل یادگیری بودن توزیع‌های دارای sample compressibility  | 7  |
| ۴.۴. بسته بودن sample compressibility نسبت به mixture, product | 7  |
| ۴.۵. کران یادگیری برای توزیع‌های با پایه‌ی گاوسی               | 8  |
| ۴.۶. کران یادگیری برای توزیع‌های با پایه‌ی گاوسی axis-aligned  | 9  |
| ۵. نتیجه‌گیری  | 11 |
| ۶. مراجع   | 12 |

## ۱. مقدمه

در این گزارش به بررسی مقاله‌ی زیر می‌پردازیم:

### Nearly Tight Sample Complexity Bounds for Learning Mixtures of Gaussians via Sample Compression Schemes [1]

به طور کلی یادگیری توزیع (distribution learning) یک مسئله‌ی معروف در حوزه‌ی یادگیری است و هدف آن تخمین توزیع از روی داده‌ای مشاهده شده است. یک جنبه‌ی مهم یادگیری توزیع، ارائه‌ی sample complexity مورد نیاز برای یادگیری آن خانواده از توزیع است. پیش از این کار، کران‌های ارائه شده برای خانواده‌هایی از توزیع‌ها که به صورت mixture هستند، به اندازه‌ی کافی tight نبودند. این مقاله کران بالا و پایین tight تری برای خانواده‌های توزیع که از جنس mixture هستند، به ویژه mixtures of gaussians و mixtures of axis-aligned gaussians، ارائه می‌دهد.

برای به دست آوردن کران بالا، مفهوم جدیدی به نام sample compressibility ارائه می‌شود که بیان می‌کند اگر بتوان پارامترهای توزیع را توسط تعداد اندکی از نمونه‌های دست‌چین شده encode کرد، یادگیری آن توزیع نیز با تعداد کمی از نمونه‌ها امکان پذیر است. همچنین نشان داده می‌شود که اگر توزیع‌های پایه، ویژگی sample compressibility را داشته باشند، خانواده‌ی mixture و product آن‌ها نیز این خاصیت را دارند و می‌توان برای یادگیری آن‌ها نیز کران tight ارائه کرد.

به طور خاص برای یادگیری توزیع mixture of  $k$  gaussians در فضای  $R^d$ ، کران‌های بالا و پایین به یک order می‌رسند و می‌توان گفت که تعداد  $\tilde{\Theta}(kd^2/\epsilon^2)$  سمپل برای یادگیری این خانواده از توزیع‌ها لازم و کافی است. همچنین شایان ذکر است که این کار بر روی یافتن الگوریتم بهینه تمرکز ندارد. با اینکه sample complexity ارائه شده بهینه است، اما time complexity الگوریتم، نسبت به بُعد فضا و تعداد mixture‌ها نمایشی است.

## ۲. پیش‌نیازها و تعاریف

**$\varepsilon$  - approximation**: توزیع  $\hat{g}$  یک  $\varepsilon$  - approximation از توزیع  $g$  است، اگر  $\|\hat{g} - g\|_1 \leq \varepsilon$ .

در حقیقت در این مقاله هر وقت از فاصله بین دو توزیع صحبت می‌شود، منظور فاصله‌ی (TV) total variation می‌باشد که به صورت زیر تعریف می‌شود:

$$\text{TV}(f_1, f_2) := \sup_{B \in \mathcal{L}} \int_B (f_1(x) - f_2(x)) dx = \frac{1}{2} \|f_1 - f_2\|_1$$

دلیل استفاده نکردن از معیاری مانند KL divergence آن است که این معیار نسبت به ساپورت توزیع‌ها از TV distance حساس‌تر است و برای توزیع‌هایی که مشابه‌اند اما ساپورت آن‌ها در نقاطی اشتراک ندارد، بی‌نهایت می‌شود. همچنین معیار TV از جنس فاصله بوده و نسبت به ترتیب توزیع‌ها متقارن است.

### PAC-learning distributions

#### Realizable case:

$$X_1, \dots, X_m \sim g \quad (i.i.d \text{ sample})$$

$$A(X_1, \dots, X_m) \rightarrow \hat{g} \in \mathcal{G}$$

$$P(\|g - \hat{g}\|_1 \leq \varepsilon) \geq 1 - \delta$$

#### Non-realizable case:

$$X_1, \dots, X_m \sim f \notin \mathcal{G} \quad (i.i.d \text{ samples})$$

$$A(X_1, \dots, X_m) \rightarrow \hat{g} \in \mathcal{G}$$

$$P\left(\|f - \hat{g}\|_1 \leq c \inf_{g^* \in \mathcal{G}} \|g^* - f\|_1 + \varepsilon\right) \geq 1 - \delta$$

### $k - \text{mix}(F)$

اگر  $F$  خانواده‌ی توزیع پایه باشد، آنگاه  $k - \text{mix}(F)$  خانواده‌ی توزیع‌های mixture آن است که به صورت زیر تعریف می‌شود:

$$k - \text{mix}(F) := \left\{ \sum_{i=1}^k w_i f_i : w_i \geq 0, \quad \sum_i w_i = 1, \quad f_1, \dots, f_k \in F \right\}$$

### ۳. روش های پیشین

روش های متنوعی برای یادگیری توزیع وجود دارد که در ادامه به شرح مختصر برخی از آن ها می پردازیم.

- روش maximum likelihood estimation یک روش پارامتریک برای تخمین پارامتر های توزیع از روی داده های مشاهده شده است. مشکل این روش آن است که برای خانواده ای مانند mixture of gaussians، به دست آوردن میانگین ها، ماتریس های کوواریانس و وزن های مربوط به هر توزیع پایه، در شرایطی که میانگین ها و کوواریانس ها بسیار به هم نزدیک هستند کار دشواری است. بنابراین در این روش به قیدی از جنس separability پارامترها نیاز است که لزوما برقرار نیست.
- همچنین به طور کلی می توان نشان داد که entry-wise approximation نیز روش مناسبی نیست و به condition number وابسته است. به عنوان مثال می توان دو ماتریس کوواریانس  $\Sigma, \hat{\Sigma}$  را در نظر گرفت که فاصله ی پارامترهای آن ها کم است:

$$|\Sigma - \hat{\Sigma}| < \varepsilon, \quad \varepsilon = \frac{2}{\kappa(\Sigma) + 1}$$

اما فاصله ی توزیع های  $\mathcal{N}(0, \Sigma), \mathcal{N}(0, \hat{\Sigma})$  از حیث KL و یا TV بیشینه است.

- روش دیگری که در عمل زیاد استفاده می شود kernel density estimation است که نیازمند قیدهایی از جنس smoothness و boundedness می باشد. این در حالی است که خانواده ی گاوسی ها به طور کلی Lipschitz و یا bounded نیست.
- روش دیگری که برای یادگیری توزیع استفاده می شود histogram estimation است اما استفاده از این روش به sample complexity ای منجر می شود که نسبت به بُعد فضا نمایی است.
- روش دیگری که می توان استفاده کرد minimum distance estimation است که بر اساس uniform convergence و به دست آوردن VC dimension برای Yatracos class می باشد. این روش برای توزیع یک گاوسی d بُعدی به کران بهینه ی  $O(d^2/\varepsilon^2)$  می رسد. اما برای خانواده k-mix(gaussian) در فضای d بُعدی به کران loose برابر با  $O(k^4 d^4/\varepsilon^2)$  می رسد.
- همچنین تخمینگرهای بر پایه ی piecewise polynomials نیز نمی توانند به خوبی گاوسی های d بُعدی را تخمین بزنند و کران هایی که به دست می دهند به صورت نمایی به d وابسته است.

بهترین کرانی که قبل از این کار برای خانواده ی mixture of gaussians و mixture of axis-aligned gaussians وجود داشت به ترتیب  $\tilde{O}(kd^2/\varepsilon^4)$  و  $\tilde{O}(kd/\varepsilon^4)$  بودند که به واسه ی وجود  $\varepsilon^4$  در مخرج، بهینه نبودند. بنابراین هیچ کدام از این روش ها برای یادگیری mixture ها بهینه نیستند.

## ۴. روش ارائه شده

### ۴.۱. مفهوم sample compressibility

ابتدا به معرفی مفهوم sample compressibility می‌پردازیم.

دیکودر  $J$  برای خانواده‌ی توزیع  $F$  تابعی است deterministic که تعداد متناهی سمپل ( $Z$ ) و تعداد متناهی بیت دریافت کرده و عضوی از خانواده توزیع  $F$  را به دست می‌دهد:

$$J: \bigcup_{n=0}^{\infty} Z^n \times \bigcup_{n=0}^{\infty} \{0,1\}^n \rightarrow F$$

حال می‌گوییم  $F$  دارای  $(\tau, t, m)$  sample compressibility است اگر برای  $F$  دیکودر  $J$  ای وجود داشته باشد که به ازای هر  $g \in F$  داشته باشیم:

به ازای هر  $\varepsilon \in (0,1)$  اگر توالی  $S$  شامل  $m(\varepsilon)$  تا سمپل  $i.i.d$  از  $g$  داشته باشیم، وجود داشته باشد توالی  $L$  به طول حداکثر  $\tau(\varepsilon)$  که از المان‌های  $S$  ایجاد شده است و توالی  $B$  که شامل حداکثر  $t(\varepsilon)$  بیت است، به طوری که:

$$P(\|J(L, B) - g\|_1 \leq \varepsilon) \geq 1 - \delta$$

بنابراین طبق این مفهوم، توزیع  $g$  توسط  $L, B$  انکود می‌شود و هرچه  $m$  کوچکتر باشد به این معناست که توزیع توسط تعداد نمونه‌های کمتری قابل انکود شدن است و بنابراین توزیع compressibility بالاتری دارد. همچنین دقت شود که  $S, L$  توالی هستند نه مجموعه و بنابراین می‌توانند شامل عضو تکراری باشند.

### ۴.۲. محاسبه‌ی sample compressibility برای چند توزیع نمونه

- $F = \{Unif(a, b)\}$   
 $m \in O\left(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon}\right), \quad t = 0, \quad \tau = 2$
- $F = \{\mathcal{N}(\mu, \sigma)\}$   
 $m \in O\left(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon}\right) \log \frac{1}{\varepsilon}, \quad t = 0, \quad \tau = 2$

### ۴.۳. قابل یادگیری بودن توزیع های دارای sample compressibility

حال می‌خواهیم نشان دهیم که sample compressibility برای توزیع نتیجه می‌دهد که آن توزیع learnable است. طبق قضیه‌ی ۳.۴ مقاله داریم:

الگوریتم deterministic ای وجود دارد که اگر  $M$  تا توزیع کاندید مانند  $\{f_1, \dots, f_M\}$  داشته باشیم و  $\varepsilon > 0$  و تعداد  $\log(3M^2/\delta)/2\varepsilon^2$  تا سمپل  $i.i.d$  از توزیع ناشناخته  $g$  داشته باشیم، این الگوریتم  $j \in [M]$  ای خروجی می‌دهد که به احتمال حداقل  $1 - \delta/3$  داشته باشیم:

$$\|f_j - g\|_1 \leq 3 \min_{i \in [M]} \|f_i - g\|_1 + 4\varepsilon$$

این نامساوی با پهن کردن grid در فضای پارامترها به دست می‌آید. به این صورت که ابتدا توسط confidence interval فضای پارامترهای توزیع را محدود می‌کنیم و سپس در فضای باقی مانده grid پهن کرده و عبارت اول سمت راست تساوی را  $O(\varepsilon)$  می‌کنیم.

حال می‌توان به جای این کار از مفهوم sample compressibility ارائه شده استفاده کرد. می‌دانیم که اگر این ویژگی را داشته باشیم و  $m$  سمپل از توزیع  $g$  داشته باشیم، دیکودر یک توزیع در  $\varepsilon$  همسایگی  $g$  به ما می‌دهد. همچنین  $M$  نیز در این حوزه معادل حداکثر حالات خروجی دیکودر یعنی  $2^{t(\varepsilon)} \times m(\varepsilon)^{\tau(\varepsilon)}$  می‌شود.

حال طبق قضیه‌ی ۳.۵ مقاله می‌توان گفت:

اگر  $F$  دارای sample compressibility  $(\tau, t, m)$  باشد و داشته باشیم:

$$\tau'(\varepsilon) := \tau(\varepsilon/6) + t(\varepsilon/6)$$

آنگاه  $F$  توسط  $\tilde{O}(m(\varepsilon/6) + \tau'(\varepsilon)/\varepsilon^2)$  سمپل قابل یادگیری است.

### ۴.۴. بسته بودن sample compressibility نسبت به mixture, product

حال می‌دانیم که اگر خانواده‌ی توزیع پایه و یا خانواده‌ی توزیع mixture دارای sample compressibility باشند، آنگاه مطابق sample complexity ذکر شده، قابل یادگیری اند. این در حالی است که بررسی sample compressible بودن توزیع های mixture نسبت به توزیع های پایه، پیچیده تر است. به این منظور لم های زیر در مقاله آورده شده اند که نشان می‌دهد sample compressibility نسبت به mixture و product توزیع ها بسته است:

لم ۳.۶ مقاله، درباره‌ی product توزیع پایه:

اگر  $F$  دارای sample compressibility  $(\tau(\varepsilon), t(\varepsilon), m(\varepsilon))$  باشد، آنگاه  $F^d$  دارای

$$(\tau(\varepsilon/d), t(\varepsilon/d), m(\varepsilon/d) \log(3d)) \text{ sample compressibility}$$

می‌باشد.

لم ۳.۷ مقاله، درباره‌ی mixture توزیع پایه:

اگر  $F$  دارای

$(\tau(\varepsilon), t(\varepsilon), m(\varepsilon))$  sample compressibility

باشد، آنگاه  $k - \text{mix}(F)$  دارای

(

$$k\tau(\varepsilon/3),$$

$$kt(\varepsilon/3) + k \log_2(4k/\varepsilon),$$

$$48m(\varepsilon/3) k \log(6k)/\varepsilon$$

) sample compressibility

می‌باشد.

در حالت کلی با فرض قابل یادگیری بودن توزیع پایه نمی‌توانستیم نتیجه بگیریم که توزیع mixture نیز قابل یادگیری است. حال اما با فرض sample compressible بودن توزیع پایه می‌توان sample compressible بودن توزیع mixture را نتیجه گرفت و با توجه به قضایای پیشین، به قابل یادگیری بودن این توزیع ها نیز رسید.

#### ۴.۵. کران یادگیری برای توزیع های با پایه‌ی گاوسی

تا این جا خواصی که ذکر شد، به طور کلی و برای هر توزیع پایه‌ای صادق بود. حال اما می‌خواهیم به طور دقیق تر، خواص ذکر شده را برای توزیع پایه‌ی گاوسی بررسی کنیم.

قضیه‌ی ۱.۱ مقاله: خانواده‌ی k-mixture از گاوسی های  $d$  بُعدی، توسط  $\tilde{O}(kd^2/\varepsilon^2)$  سمپل قابل یادگیری اند.

اثبات: طبق لم ۴.۱ مقاله می‌دانیم که خانواده‌ی گاوسی های  $d$  بُعدی دارای

(

$$O(d \log(2d)),$$

$$O(d^2 \log(2d) \log(d/\varepsilon)),$$

$$O(d \log(2d))$$

) sample compressibility



می‌باشد. حال اگر این لم را با لم ۳.۷ مقاله که پیش تر ذکر شد ترکیب کنیم خواهیم داشت که خانواده‌ی  $k$ -mixture از گاوسی های  $d$  بُعدی دارای

(

$$O(kd \log(2d)),$$

$$O(kd^2 \log(2d) \log(d/\varepsilon) + k \log(k/\varepsilon)),$$

$$O(kd \log k \log(2d)/\varepsilon)$$

) sample compressibility

می‌باشد.

حال اگر مقادیر  $\tau, t, m$  به دست آمده را در قضیه‌ی ۳.۵ مقاله جایگذاری کنیم، طبق

$$m(\varepsilon) = \tilde{O}(dk/\varepsilon), \quad \tau'(\varepsilon) = \tilde{O}(d^2k)$$

خواهیم داشت که sample complexity برای یادگیری  $k$ -mixture گاوسی های  $d$  بُعدی برابر  $\tilde{O}(kd^2/\varepsilon^2)$  می‌باشد.

قضیه‌ی ۱.۲ مقاله: هر روشی که برای یادگیری خانواده‌ی  $k$ -mixture از گاوسی های  $d$  بُعدی استفاده شود، حداقل به  $\tilde{\Omega}(kd^2/\varepsilon^2)$  سمپل نیاز دارد.

ترکیب قضیه‌ی ۱.۱ و ۱.۲ مقاله نشان دهنده‌ی بهینه بودن کران بالای ارائه شده است و می‌توان گفت که تعداد  $\tilde{O}(kd^2/\varepsilon^2)$  سمپل برای یادگیری توزیع mixture of  $k$  gaussians در فضای  $R^d$  لازم و کافی است.

#### ۴.۶. کران یادگیری برای توزیع های با پایه‌ی گاوسی axis-aligned

برای حالت mixture of gaussians نیز کران جدیدی به دست می‌آید.

قضیه‌ی ۱.۳ مقاله: خانواده‌ی  $k$ -mixture از گاوسی های axis-aligned که  $d$  بُعدی هستند، توسط  $\tilde{O}(kd/\varepsilon^2)$  سمپل قابل یادگیری اند.

اثبات: به طور مشابه با ترکیب کردن لم های ۴.۱ و ۳.۶ و ۳.۷ مقاله می‌توان نشان داد که خانواده‌ی  $k$ -mixture از گاوسی های axis-aligned که  $d$  بُعدی هستند، دارای

(

$$O(kd),$$

$$O(kd \log(d/\varepsilon) + k \log(k/\varepsilon)),$$

$$O(k \log k \log(3d)/\varepsilon)$$

) *sample compressibility*

می‌باشد.

حال اگر مقادیر  $\tau, t, m$  به دست آمده را در قضیه ۳.۵ مقاله جاگذاری کنیم، خواهیم داشت که sample complexity برای یادگیری k-mixture گاوسی های axis-aligned که d بُعدی هستند برابر  $\tilde{O}(kd/\varepsilon^2)$  می‌باشد.

## ۵. نتیجه گیری

یکی از مسائل حل نشده در یادگیری توزیع، به دست آوردن sample complexity مورد نیاز برای یادگیری توزیع است. یک نگاه شهودی از حوزه یادگیری نظارتی بیان می‌کند که sample complexity مورد نیاز برای یادگیری هر concept ای، اعم از توابع یا توزیع ها، با بُعد فضای فرضیه تقسیم بر  $\epsilon^2$  نسبت دارد. در حالتی که درباره‌ی دسته بندی دوتایی صحبت کنیم، بُعد فضای فرضیه معادل VC dimension خواهد بود. برای یادگیری توزیع نیز، این بُعد معادل تعداد پارامتر های توزیع می‌باشد. قبل از این کار، برای توزیع هایی مانند گاوسی و گاوسی های axis-aligned که به ترتیب  $O(d^2)$  و  $O(d)$  پارامتر دارند، این شهود برقرار بود. اما کران هایی که برای mixture این توزیع ها وجود داشت، loose تر از این مقدار بودند. حال در این کار اثبات شد که این شهود برای mixture های این توزیع ها که به ترتیب  $O(kd^2)$  و  $O(kd)$  پارامتر دارند نیز برقرار است.

مفهوم جدید sample compressibility که در این کار ارائه شد، در مورد sample complexity مورد نیاز برای یادگیری توزیع، یک شرط کافی ارائه می‌دهد. با این که اثبات شد برای mixture of gaussians این شرط لازم نیز می‌باشد (یعنی در این حالت کران ارائه شده بهینه است)، اما در حالات دیگر لزوم این شرط اثبات نشده و به صورت یک مسئله‌ی حل نشده باقی می‌ماند. یک خاصیت جالب مفهوم sample compressibility آن است که نسبت به product و همچنین mixture توزیع های پایه بسته است. به طور کلی اگر بخواهیم از قابل یادگیری بودن توزیع پایه به قابل یادگیری بودن توزیع mixture برسیم، مسیر سختی را پیش رو خواهیم داشت. حال اما با معرفی مفهوم sample compressibility تنها کافی است وجود این ویژگی را برای خانواده‌ی توزیع پایه نشان دهیم تا وجود آن برای خانواده‌ی توزیع mixture نیز به دست آید و قابل یادگیری بودن آن اثبات شود.

## ٦. مراجع

- [1] Ashtiani, H., Ben-David, S., Harvey, N., Liaw, C., Mehrabian, A., & Plan, Y. (2018). Nearly tight sample complexity bounds for learning mixtures of gaussians via sample compression schemes. *Advances in Neural Information Processing Systems*, 31.