



الگوریتم‌های خلاصه‌سازی برای مه‌داده

محمد هادی فروغ‌منداعرابی
پاییز ۱۳۹۹

مروری بر تجزیه SVD و نشانیدن زیرفضا، حل مسئله‌ی رگرسیون

جلسه شانزدهم

نگارنده: مائده حشمتی

۱ مروری بر مباحث گذشته

در جلسه گذشته مروری بر خلاصه‌سازی خطی غافل داشتیم و با نشانیدن زیرفضا جلسه را به پایان رساندیم. در این جلسه با کمک ابزارهایی همچون نشانیدن زیرفضا و تجزیه‌ی SVD به حل مسئله رگرسیون می‌پردازیم.

۲ تجزیه SVD :

قضیه ۱: هر ماتریس حقیقی $A \in \mathbb{R}^{n \times d}$ با رتبه‌ی r می‌تواند به شکل زیر تجزیه شود:

$$A = U \Sigma V^T$$

به طوریکه $U \in \mathbb{R}^n \times r$, $V \in \mathbb{R}^{d \times r}$, $U^T U = I$, $V^T V = I$, $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$, $\sigma_i > 0$ و σ_i ها را مقادیر تکی ماتریس A می‌نامند و برابر رادیکال مقادیر ویژه‌ی $A^T A$ می‌باشند. براساس این تجزیه، برای هر ماتریس حقیقی، شبه وارون نیز تعریف می‌شود.

تعریف ۱: برای یک ماتریس حقیقی داده شده با تجزیه‌ی $M = U \Sigma V^T$ ، شبه وارون این ماتریس (pseudoinverse) برابر است با:
 $M^+ = V \Sigma^{-1} U^T$ (مقادیر تکی در حالت معکوس در Σ^{-1} قرار می‌گیرند: $\Sigma^{-1} = \text{diag}(1/\sigma_1, 1/\sigma_2, \dots, 1/\sigma_r)$)

ادعا ۱: ماتریس $A(A^T A)^+ A^T$ تصویر در زیرفضای ستون‌های A می‌باشد.

اثبات ادعا ۱: با کمک تجزیه SVD اثبات می‌کنیم:

$$A = U \Sigma V^T A^T A = V \Sigma U^T U \Sigma V^T = V \Sigma^\dagger V^T (A^T A)^+ = V \Sigma^{-\dagger} V^T$$

در نهایت خواهیم داشت:

$$A(A^T A)^+ A^T = U \Sigma V^T (V \Sigma^{-\dagger} V^T) V \Sigma U^T = U U^T$$

که یعنی ضرب داخلی در پایه‌ها و نشان‌دهنده‌ی تصویر از A می‌باشد و حکم ثابت می‌شود.

۳ نشانندن زیرفضا:

تعریف ۲: زیرفضای خطی $E \subset \mathbb{R}^n$ را در نظر بگیرید. ماتریس $\Pi \in \mathbb{R}^{m \times n}$ یک نشاننده برای این زیرفضاست (ε -subspace embedding) اگر داشته باشیم:

$$\forall x \in E : (1 - \varepsilon) \|x\|_2^2 \leq \|\Pi x\|_2^2 \leq (1 + \varepsilon) \|x\|_2^2$$

یعنی ماتریسی که تقریباً طول بردارهای درون این زیرفضا را ثابت نگه می‌دارد.

تعریف ۳: یک نشاننده‌ی زیرفضای غافل (OSE) oblivious subspace embedding (ε, d, δ) ، یک توزیعی D روی ماتریس‌های $\Pi \in \mathbb{R}^{m \times n}$ است به طوریکه به ازای هر ماتریس $U \in \mathbb{R}^{n \times d}$ با ستون‌های متعام داریم:

$$\mathbb{P}_{\Pi \sim D} (\|(\Pi U)^T (\Pi U) - I\| > \varepsilon) < \delta$$

۴ رگرسیون:

هدف این است که با توجه به داده‌هایی که تا به حال داشتیم و قضایا، یک مسئله‌ی رگرسیون را در زمان کمتری حل کنیم. ابتدا به تعریف مسئله‌ی رگرسیون می‌پردازیم.

تعریف ۴: در یک مسئله‌ی رگرسیون، یک ماتریس ویژگی‌ها داریم که سطرها، نمونه‌های ما هستند و تعداد آنها بسیار زیاد است و ستون‌های آن، ویژگی‌های مورد نظر ما می‌باشد که تعدادش از سطرها بسیار کمتر است. هدف یافتن یک ماتریسی است که در صورت ضرب شدن در ماتریس ما، بتوانیم جواب نهایی را پیش‌بینی کنیم. یعنی با یک ترکیب خطی خوب از سطرهای ماتریس داده‌هایمان، به جواب نهایی با صحت بالا برسیم:

$$X\beta \simeq y$$

X همان ماتریس ویژگی‌ها و داده‌های اولیه‌ی ما، y هدف پیش‌بینی و β مجهول ما است. ($X \in \mathbb{R}^n \times d : n \gg d$) پیشنهاد اولیه برای حل این مسئله در زمان کمتر، استفاده از کمترین مربعات می‌باشد.

$$\beta^{LS} = \operatorname{argmin}_{\beta} \|X\beta - y\|_2, \quad X \in \mathbb{R}^{n \times d}, y \in \mathbb{R}^n$$

می‌دانیم هر بردار، نمایشی از برداری درون زیرفضای ساخته شده از ستون‌های X و برداری درون زیرفضای عمود بر آن دارد، ($y = y^\perp + y^\parallel$)، خواهیم داشت:

$$\begin{aligned} \|X\beta - y\|_2^2 &= \|X\beta - y^\perp - y^\parallel\|_2^2 \\ &= \|X\beta - y^\parallel\|_2^2 - 2\langle X\beta - y^\parallel, y^\perp \rangle + \|y^\perp\|_2^2 \end{aligned}$$

جمله دوم به وضوح صفر است و چون y^\perp معین می‌باشد، آنچه دست ماست تنها جمله‌ی اول یعنی $\|X\beta - y^\parallel\|_2^2$ می‌باشد. می‌توانیم β را به نحوی اتخاذ کنیم که داشته باشیم: $\beta^\parallel = y^\parallel$ و برای تصویر y داریم: $X(X^T X)^+ X^T y$ در نهایت برای β^{LS} داریم: $\beta^{LS} = (X^T X)^+ X^T y$.

در این حالت زمان اجرای ما بسیار زیاد بود و در ادامه تلاش خواهیم کرد که تعداد سطرها را به نحوی کاهش دهیم. در این صورت به یک رگرسیون جدید نیاز داریم به طوری که به جای X از ΠX و به جای y از Πy استفاده کنیم. ماتریس Π باید به نحوی باشد که تعداد صف‌هایش زیاد و به تعداد کمی ۱ یا ۱- داشته باشد تا زمان اجرا را کاهش دهد و نیز جواب رگرسیون اصلی را زیاد تغییر ندهد. به طور کلی چنین ایده‌ای داریم:

$$\hat{\beta}^{LS} = \operatorname{argmin} \|\Pi X \beta - \Pi y\|_2$$

و تصویر y برابر خواهد بود با: $\Pi X (\Pi X^T \Pi X)^+ \Pi X^T \Pi y$

لم ۱: زیرفضای خطی $E := \operatorname{span}(\operatorname{cols}(X), y)$ را در نظر بگیرید. ماتریس Π یک نشاننده ی این زیرفضاست. در اینصورت خواهیم داشت:

$$\|X \hat{\beta}^{LS} - y\|_2 \leq \frac{1+\varepsilon}{1-\varepsilon} \|X \beta^{LS} - y\|_2$$

اثبات لم ۱: با توجه به تعریف نشانند و تعریف $\hat{\beta}$ خواهیم داشت:

$$(1-\varepsilon) \|X \hat{\beta}^{LS} - y\|_2 \leq \|\Pi X \hat{\beta}^{LS} - \Pi y\|_2 \leq \|\Pi X \beta^{LS} - \Pi y\|_2 \leq (1+\varepsilon) \|X \beta^{LS} - y\|_2$$

پس حالا نیاز داریم تا یک Π بیابیم تا در شروط ما صدق کند و بتواند جوابی برای مسئله رگرسیون بدهد.

قضیه ۲: برای ماتریس $X \in \mathbb{R}^{n \times d}$ و بردار $y \in \mathbb{R}^n$ در زمان $O(\operatorname{nnz}(X) + n) + \operatorname{poly}(d/\varepsilon)$ میتوان با احتمال $1 - \varepsilon$ به گونه‌ای یافت که در شرایط زیر صدق کند:

$$\|X \hat{\beta}^{LS} - y\|_2 \leq \frac{1+\varepsilon}{1-\varepsilon} \|X \beta^{LS} - y\|_2$$

یکی از انواع ماتریس‌های مطلوب برای ما، ماتریسی است که خاصیت OSE را دارد و در شروط زیر صدق می‌کند:

$$\mathbb{P}_{\Pi \sim D}(\|(\Pi U)^T (\Pi U) - I\| > \varepsilon) < \delta$$

یکی از این دسته، ماتریس‌های *countsketch* می‌باشد که در هر ستون با یک تابع هش، مکان تنها درایه‌ی ناصفر مشخص می‌شود و با احتمال $1/2$ برابر $+1$ یا -1 خواهد شد. این ماتریس با تعداد سطرهای $m = d^2/\varepsilon^2$ خاصیت OSE را دارد و در شروط مورد نظر ما صدق می‌کند.

۵ احساس فشردگی:

صورت کلی مسئله‌ی *Compressed Sensing* یا احساس فشردگی، بدین صورت می‌باشد که یک بردار بزرگ و k تنک داریم و چون برای ذخیره‌سازی این بردار، به فضای زیادی نیاز داریم، میخواهیم به جای این بردار، Πx را ذخیره کنیم تا فضای کمتری اشغال کند و در صورت نیاز به x ، عملیات بازیابی را به نحوی انجام دهیم که با تقریب خوبی به x اولیه برسیم. ماتریس Π باید به نحوی باشد که تعداد ستون‌هایش برابر با تعداد سطرهای x باشد، اما تعداد سطرهایش به قدر کافی کوچک باشد تا در نهایت ضرب این دو، به برداری با طول کمتر رسیده تا فضای کمی برای ذخیره‌سازی اشغال کند.

$$x \rightarrow y = \Pi x$$

و برای بازیابی x از روی y و Π :

$$\begin{aligned} \min \|z\|_1 \\ \text{s.t. } \Pi z = y \end{aligned}$$

در واقع با این الگوریتم، به دنبال تنک ترین z می‌گردیم که حاصلضرب Πz برابر با مقدار ذخیره شده یعنی y شود. علت قرار دادن نرم یک در تابع هدف: قرار دادن نرم صفر که دقیقاً تعداد درایه‌های ناصفر بردار را می‌دهد در اینجا موجب سخت شدن مسئله خواهد شد و مسئله به $NP - Complete$ تبدیل می‌شود. بقیه‌ی نرم‌ها مانند دو از نظر هندسی برای رسیدن به جواب، مطلوب نیستند و ساده ترین و بهترین نرم برای یافتن تنک ترین بردار، نرم یک می‌شود.

این الگوریتم برای حالتی که x حقیقی باشد و به جای دقیقاً تنک، تقریباً تنک نیز باشد کار می‌کند. در واقع با توجه به خواسته‌ی ما و ورودی مسئله، این الگوریتم می‌تواند برای ما مطلوب باشد. به طور مثال، می‌توانیم با این روش، k تا بزرگترین درایه‌های x را بازیابی کنیم و می‌دانیم این روش، تقریباً جواب خوبی به ما خواهد داد.

تعریف ۵: ماتریس $\Pi \in \mathbb{R}^{m \times n}$ دارای خاصیت $(\varepsilon, k) - \text{restricted isometry property}$ است، اگر به ازای تمام بردارهای k تنک با طول واحد در شرط زیر صدق کند:

$$1 - \delta \leq \|\Pi x\|_2^2 \leq 1 + \delta$$

یعنی ماتریسی که برای بردارهای k تنک، طول بردار را تقریباً حفظ کند. به صورت دیگر همانطور که در جلسات گذشته نیز گفتیم یعنی داشته باشیم:

$$\sup_{T \subset [n], |T|=k} \|I_k - (\Pi^{(T)})^* \Pi^{(T)}\| < \delta$$

تعریف ۶: یک ماتریس A $m \times n$ دارای خاصیت k *null - space property of order* k می‌باشد اگر به ازای هر $\eta \in \mathbb{R}^n$ که $A\eta = 0$ و به ازای هر مجموعه k عضوی $T \subset \{1, 2, \dots, n\}$ داشته باشیم:

$$\|\eta\|_1 \leq C \|\eta_{-T}\|_1$$

$-T$ به معنای مجموعه‌ی مکمل T است.

این ویژگی بدین معنی است که به ازای بردار η دریه‌هایش بسیار بزرگ نباشند و بتوان k تا از بزرگترین دریه‌هایش را با بقیه‌ی دریه‌ها کران کرد و طبق تعریف، زیرفضای پوچ A دارای این خاصیت می‌باشد.

لم ۲: اگر ماتریس A دارای خاصیت RIP از *order* برابر با $(c+2)k$ و ثابت برابر با $1 < c, \delta$ باشد، آنگاه A دارای خاصیت $null - space property$ از *order* برابر با $2k$ و ثابت $C = 1 + \sqrt{2/c}(1+\delta)(1-\delta)$ خواهد بود. یعنی خاصیت RIP خاصیت $null - space$ را نتیجه می‌دهد.

اثبات لم بالا در جلسه آینده گفته خواهد شد.