



دانشگاه صنعتی اصفهان

دانشکده مهندسی برق و کامپیوتر

دستورکار آزمایشگاه هوش محاسباتی

جلسه ۱۲

عامل بندی ماتریس و سیستم پیشنهاددهنده

استاد درس: دکتر مهران صفایانی

فصل ۱۲

عامل بندی ماتریس و سیستم پیشنهاددهنده

اهداف این جلسه

شما در این جلسه یاد خواهید گرفت که :

- یک سیستم پیشنهاددهنده^۱ بسازید.
- عملکرد این سیستم را ارزیابی کنید.
- عامل بندی ماتریس^۲ را با استفاده از SGD^۳ درک و پیاده سازی کنید.
- تعداد مناسبی از عامل ها، مانند عامل های منظم سازی، را انتخاب کنید.
- سیستم خود را با یکسری پایه^۴ مقایسه کنید.

در این جلسه، ما از مجموعه داده `movielens100k.csv` استفاده خواهیم کرد. همچنین شما می توانید از کدهای آماده ی مفیدی که برای شما آماده شده است، استفاده کنید.

System Recommender^۱
Factorization Matrix^۲
Descent Gradient Stochastic^۳
baseline^۴

۱.۱۲ ایجاد و مصورسازی تقسیم بندی داده های آموزشی و آزمون

از آنجایی که هدف ما، پیش بینی کردن رتبه بندی های دیده نشده است، می توانیم یک مجموعه ی آزمون را با پنهان کردن برخی عناصر ماتریس، ماتریس ایجاد کنیم. منظور ما، انتخاب کردن تعدادی رتبه بندی به صورت تصادفی از توی ماتریس است. این کار یک عمل بسیار معمولی در مجموعه داده های سیستم های پیشنهاد دهنده و کلان داده ها است.

تمرین اول

تابع `split_data` که درون فایل ژوپیتِر مربوط به این جلسه قرار دارد را به منظور تقسیم مجموعه داده به داده های آموزشی و آزمون، تکمیل کنید. ما فقط فیلمها و کاربرانی را در نظر می گیریم، که بیش از ۱۰ رتبه بندی داشته باشند و دیگر فیلمها و کاربران را نادیده می گیریم. در میان رتبه بندی های معتبر باقی مانده، شما باید بصورت تصادفی با احتمال ۱۰٪ یک رتبه بندی را برای داده ی آزمون و با احتمال ۹۰٪ یک رتبه بندی را برای داده ی آموزشی، انتخاب کنید. ما تعداد کمی از داده های آموزشی ای را برای کاربران و فیلم هایی که رتبه بندی کمی دارند، نگه می داریم. می توان نتیجه ی تقسیم بندی را مانند شکل دو، نشان داد.

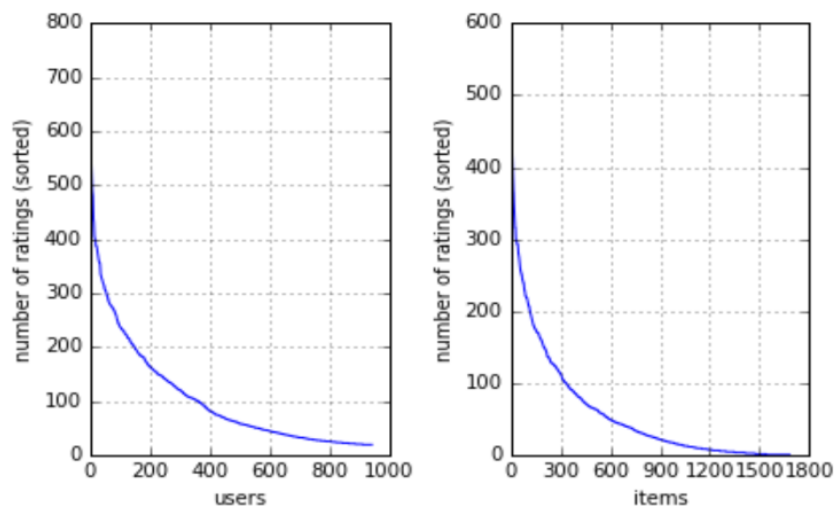
۲.۱۲ ارزیابی عملکرد

ما از مجذور میانگین مربعات خطا^۱ (RMSE) برای ارزیابی عملکرد استفاده خواهیم کرد. با داشتن رتبه بندی صحیح x_{dn} برای d - اُمین فیلم و n - اُمین کاربر و همچنین پیش بینی $(WZ^T)_{dn}$ ، می توانیم RMSE را بصورت زیر تعریف کنیم:

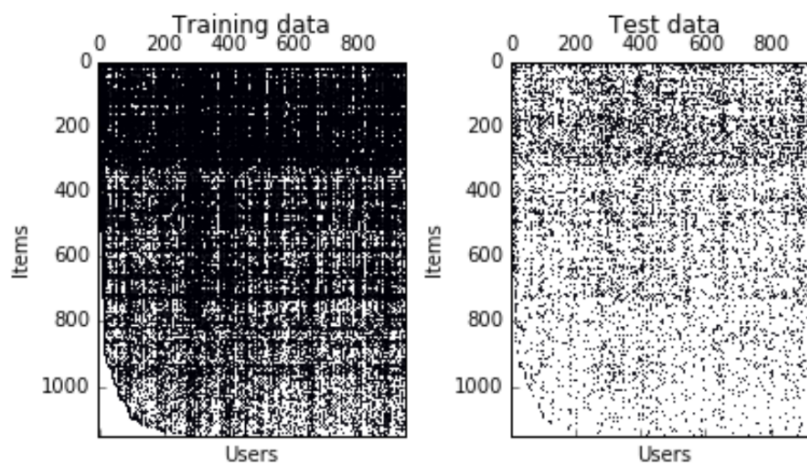
$$\text{RMSE}(W, Z) := \sqrt{\frac{1}{|\Omega|} \sum_{(d,n) \in \Omega} \frac{1}{2} [x_{dn} - (WZ^T)_{dn}]^2} \quad (۱.۱۲)$$

که $W \in \mathbb{R}^{D \times K}$, $Z \in \mathbb{R}^{N \times K}$ در اینجا $\Omega \subseteq [D] \times [N]$ مجموعه ی شماره های رتبه بندی های مشاهده شده ی ماتریس ورودی X می باشد. RMSE می تواند هم بر روی مجموعه آموزشی Ω محاسبه شود و هم بر روی یک مجموعه آزمون آماده به کار.

^۱root mean-square error



شکل ۱۰.۱۲: شکل سمت چپ، تعداد رتبه‌بندی (مرتب شده) هر کاربر را نشان می‌دهد و شکل سمت راست، تعداد رتبه‌بندی (مرتب شده) فیلم‌ها را نشان می‌دهد. در هر دو نمودار، اعداد به منظور وضوح بیشتر، به صورت نزولی نمایش داده شده‌اند.



شکل ۲۰.۱۲: این تصویر، تقسیم‌بندی آموزشی و آزمون داده‌ها را نشان می‌دهد. نمودار سمت چپ، نشان‌دهنده‌ی داده‌های آموزشی و نمودار سمت راست، نشان‌دهنده‌ی داده‌های آزمون، می‌باشند. در هر نمودار، هر نقطه بیانگر یک زوج کاربر-فیلم با یک رتبه‌بندی غیر صفر، می‌باشد.

۳.۱۲ مدل های پایه

از مدل های ذیل، که برای پیش بینی، از میانگین بهره می برند، به عنوان پایه، استفاده خواهیم کرد.

$$\text{Mean: Global } \hat{x} := \frac{1}{|\Omega|} \sum_{(d,n) \in \Omega} x_{dn} = \frac{1}{|\Omega|} \sum_{n=1}^N \sum_{d \in \Omega_n} x_{dn} = \frac{1}{|\Omega|} \sum_{d=1}^D \sum_{n \in \Omega_d} x_{dn}, \quad (2.12)$$

$$\text{Mean: User } \hat{x}_n := \frac{1}{|\Omega_{:,n}|} \sum_{d \in \Omega_{:,n}} x_{dn}, \quad (3.12)$$

$$\text{Mean: Movie } \hat{x}_d := \frac{1}{|\Omega_{d,:}|} \sum_{n \in \Omega_{d,:}} x_{dn}, \quad (4.12)$$

که Ω مجموعه ای از نشان گرهای غیر صفر رتبه بندی (d, n) در ماتریس داده های آموزشی و $\Omega_{:,n}$ مجموعه ای از فیلم های رتبه بندی شده توسط n - آمین کاربر و $\Omega_{d,:}$ مجموعه ای کاربرانی است که به فیلم d - ام، رتبه داده اند.

تمرین دوم

در ابتدا، سه پایه ی بالا را با هم مقایسه خواهیم کرد.

- به نظر شما، کدامیک از سه مدل، بهترین عملکرد را خواهد داشت ؟ چرا؟
- توابع `baseline_global_mean()`، `baseline_user_mean()` و `baseline_item_mean()` را پیاده سازی کنید.

- مدل های حاصله را با یکدیگر مقایسه کنید. کدام مدل، کمترین RMSE آموزشی را به ما می دهد؟ کدامیک کمترین RMSE آزمون را به ما می دهد؟
- راهنمایی: شما می توانید مقداردهی تصادفی تقسیم آموزشی یا آزمون خود را تغییر دهید و در نتیجه، تعداد زیادی تخمین، ایجاد کنید.

۴.۱۲ عامل بندی ماتریس با استفاده از SGD

تمرین سوم

هدف پیاد سازی نزول گرادیانی تصادفی است

- نزول گرادیانی نزولی را برای تابع هدف، RMSE همانطور که در معادله ی (۱) نشان داده شده است، پیاده سازی کنید.
- SGD را برای تابع هزینه ی منظم شده ی

$$\mathcal{L}(\mathbf{W}, \mathbf{Z}) := \frac{1}{2} \sum_{(d,n) \in \Omega} [x_{dn} - (\mathbf{W} \mathbf{Z}^T)_{dn}]^2 + \frac{\lambda_w}{2} \|\mathbf{W}\|_{\text{Frob}}^2 + \frac{\lambda_z}{2} \|\mathbf{Z}\|_{\text{Frob}}^2$$

پیاده سازی کنید. که در آن $\lambda_w, \lambda_z > 0$ اسکالر هستند.

- بصورت تجربی، بهترین اندازه گام γ را به منظور گرفتن کمترین خطای آموزشی، پیدا کنید.
- آیا با این کار، خطای آزمون نیز بصورت یکنواخت کاهش می یابد یا دوباره بعد از مدت زمانی، افزایش می یابد؟ در این مورد، انتخاب های دیگر K را نیز امتحان کنید.