



بسم الله الرحمن الرحيم

وزارت علوم، تحقیقات و فناوری



گزارش پروژه پایانی

استاد راهنما: دکتر سرکار خانم طباطبایی

موضوع:

Movie recommendation system with

combining DNN and autoencoder

تهیه کننده: محمد احسان کریمی نوری

تاریخ: تابستان ۱۴۰۲

فهرست

۱.مقدمه	4
۲.مروری بر سیستم ها توصیه فیلم	6
۳.عملکردها	8
۴.الگوریتم ها توصیه	9
۵.کارهای مشابه	13
۶.معماری DNN	14
۷.چالش ها و توسعه های پیاده سازی در سیستم توصیه فیلم	17
۸.ارزیابی و معیار ها برای سیستم توصیه فیلم	19
۹.تجربه کاربر و ملاحظات اخلاقی و حریم خصوصی	20
۱۰.مدل معماری اجرا شده	21
۱۱.تئوری کد پروژه	24
۱۲.بخش عملی کد پروژه	27
۱۳.نتیجه گیری	52

1. مقدمه

این روزها، مردم به جای اتلاف وقت بر روی مشکل «بعدش چه باید ببینم»، ترجیح می‌دهند توصیه‌هایی را بر اساس تاریخچه‌های قبلی/لیست فیلم‌های تماشا شده او دریافت کنند. با رواج جامعه اطلاعاتی، دسترسی ما به اطلاعات بیشتر از آنچه تاکنون در تمام تاریخ بشر داشته ایم، آسان تر است. با این حال، طبق نظریه تصمیم‌گیری، اطلاعات بیش از حد می‌تواند منجر به تصمیمات بدتر شود، نه اینکه مردم را خوشحال کند و اطمینان حاصل کند که آنها به آنچه می‌خواهند می‌رسند. در واقع، اضافه بار اطلاعات، مصنوع انقلاب دیجیتال و موبایل، روز به روز جدی تر می‌شود. در مواجهه با منابع عظیم شبکه، کاربران با تصمیم‌گیری‌های دشواری با مشکل اضافه بار اطلاعات مواجه می‌شوند که می‌تواند باعث احساس واقعی اضطراب، خستگی ذهنی، ناتوانی و سرگیجه شود. با ارائه آنها به مردم کمک می‌کنیم تا بتوانند فیلم‌های مشابه را بر اساس اولویت پیدا کنند.

در حال حاضر برای آنها، این باعث صرفه جویی در زمان زیادی می‌شود که آنها برای انتخاب خودشان تلف می‌کردند. برای توصیه، این واقعیت را در نظر می‌گیریم که آیا برخی افراد همان فیلم را تماشا کرده‌اند و بعد از آن، بقیه افراد چه فیلم‌هایی را تماشا می‌کنند، تا بتوانیم پیش‌بینی کنیم که او ممکن است به برخی از فیلم‌هایی که دیگران پس از تماشای فیلم او تماشا کرده‌اند، علاقه مند باشد. با این کار لیستی از فیلم‌های مشابه قبلی را در اختیار کاربر قرار می‌دهیم. یادگیری عمیق به عنوان یک رویکرد مهم در زمینه‌های مختلف از جمله تجزیه و تحلیل داده‌ها در وب سایت و حوزه‌های مرتبط با سلامت بازی کرده است.

امروزه سیستم‌های توصیه سیستم‌های فیلتر اطلاعاتی هستند که از هوش مصنوعی یا الگوریتم‌های هوش مصنوعی استفاده می‌کنند که می‌توانند با فیلتر کردن Big Data با کارایی بالا مشکل اضافه بار اطلاعات را تا حد زیادی کاهش دهند تا محتوا و خدمات شخصی‌سازی شده را در اختیار کاربران قرار دهند که می‌تواند تا حد زیادی مشکل اضافه بار اطلاعات را کاهش دهد.

یک سیستم توصیه نه تنها یک الگوریتم فانتزی است، بلکه در مورد درک داده‌ها و کاربران شما نیز هست. سیستم‌های توصیه برای اولین بار در اواسط دهه ۱۹۹۰ به وجود آمدند و از آن زمان تاکنون کانون تحقیقات بوده‌اند. در حال حاضر بسیاری از پلتفرم‌ها و وب سایت‌ها شاهد کاربرد گسترده انواع مختلف سیستم‌های توصیه هستند. وب سایت‌های تجارت الکترونیک از سیستم‌های توصیه برای پیشنهاد محصولات و خدمات شخصی‌شده از مقالات، کتاب‌ها، موسیقی و فیلم‌ها استفاده می‌کنند. به طور معمول، دو نوع داده مورد استفاده در یک سیستم توصیه معمولی وجود دارد، به عنوان مثال، رتبه‌بندی یا رفتارهای خرید که منعکس‌کننده اطلاعات قابل انتساب در مورد اقلام، کاربران، و کلمات کلیدی یا نمایه‌های بافتی هستند که تعاملات کاربر-مورد را ثبت می‌کنند.

برای سرویس‌های پخش فیلم مانند Netflix، x، Hulu، Amazon Prime، و سایر سیستم‌های توصیه برای کمک به کاربران برای کشف محتوای جدید برای لذت بردن از اهمیت هستند.

1.1 هدف (Purpose)

این بخش سیستم‌های توصیه فیلم را به عنوان الگوریتم‌ها و تکنیک‌هایی تعریف می‌کند که پیشنهادات فیلم شخصی‌سازی شده را به کاربران ارائه می‌کنند. توضیح می‌دهد که چگونه می‌توان با کاهش بار اطلاعات و هدایت کاربران به سمت محتوایی که احتمالاً از آن لذت می‌برند، تجربه کاربر را افزایش داد.

هدف از این پروژه ساخت نرم افزاری است که بتوان از آن برای توصیه فیلم به کاربر استفاده کرد. مردم خیلی رنج می‌برند تا تصمیم بگیرند که بعداً چه چیزی را تماشا کنند، زمان زیادی برای تصمیم گیری صرف می‌شود. بعدی خود را برای تماشا انتخاب کنید. بر اساس علاقه قبلی آنها، مدل ما به آنها فیلم‌ها را توصیه می‌کند.

یادگیری عمیق، که اساساً فقط شبکه‌های عصبی مصنوعی عمیق است، می‌تواند مرزهای تصمیم‌گیری پیچیده را برای طبقه بندی یا رگرسیون‌های غیرخطی پیچیده بیاموزد. با قرار دادن تعداد زیادی لایه پنهان در این شبکه‌ها، شبکه‌های عصبی عمیق می‌توانند توابع پیچیده را با یادگیری استخراج بسیاری از ویژگی‌های سطح پایین از داده‌ها و ترکیب آنها در ترکیب‌های غیرخطی مفید بیاموزند.

هدف سند پیشنهادی انجام کاری است که منجر به توسعه رویکردی برای پیش‌بینی با استفاده از یادگیری عمیق می‌شود و هدف پیشنهادی با تقسیم کار به اهداف زیر محقق می‌شود:

۱. به دست آوردن یک مجموعه داده خوب (مجموعه داده‌ها).

۲. استفاده از روش‌های یادگیری عمیق برای به دست آوردن دقت بالا / نتایج بهتر.

۳. اعتبار طرح را با استفاده از یادگیری عمیق و بررسی مقاطع نتایج.

مجموعه داده‌های مورد استفاده: مجموعه داده‌های مورد استفاده برای این نوت بوک، مجموعه داده‌های رتبه بندی M۱ از Movie Lens است. این شامل ۱ میلیون رتبه بندی فیلم از ۷۱۲۰ فیلم و ۱۴۰۲۵ کاربر است. این مجموعه داده شامل:

شناسه فیلم: به هر فیلم یک شناسه منحصر به فرد به نام شناسه فیلم داده می‌شود، این شناسه می‌تواند برای اهداف شناسایی هر فیلمی استفاده شود. همچنین، هر روشی که قرار است پیاده‌سازی کنیم، تنها با گرفتن مرجع شناسه فیلم، پیاده‌سازی می‌شود، زیرا Movie ID شناسایی هر فیلمی است و برای هر فیلم منحصر به فرد است.

شناسه کاربری: هر کاربر شناسه منحصر به فرد خود را دارد تا بتوانیم کاربران را با شناسه آنها ارجاع دهیم، اگر در دو تست یک شناسه کاربری ذکر شده باشد، یعنی هر دو کاربر یکسان هستند.

رتبه بندی: امتیاز یک فیلم از ۱۰ امتیاز است، امتیاز به عوامل مختلفی مانند بازخورد ارائه شده توسط کاربر، بودجه و سود یک فیلم خاص بستگی دارد. بر اساس رتبه بندی، ما تصمیم می‌گیریم که فیلم در باکس آفیس موفق شود یا شکست بخورد. علاوه بر این، مجموعه داده‌ای از فیلم‌ها شامل عنوان فیلم، برچسب‌ها و ژانرها است.

عنوان: عنوان نام فیلم است، در برخی موارد ممکن است عنوان یکسان باشد، اما اگر فیلم توسط شناسه فیلم انجام شود، شناسایی می‌شود. عنوان فیلم می‌تواند توضیح دهد که «فیلم همه چیز درباره چیست» تا حد زیادی.

ژانرها: نوع فیلم را توضیح می دهد که چه عاشقانه، اکشن، درام، هیجان انگیز یا ماجراجویی و ... نوع فیلم مشابهی را که باید پیشنهاد کنیم تا ژانر نقش مهمی در سیستم توصیه ایفا کند.

برچسبها: استفاده از تگهای تهیه کننده یا کارگردان نشان می دهد که چیزی که در مورد این فیلم خاص است به عنوان مثال - نام هر بازیگری می تواند در برچسبها باشد یا ژانری می تواند در برچسبها باشد، هر نقل قولی نیز می تواند یک برچسب باشد.

1.2 طرح در آینده (Scope)

سیستم توصیه فیلم مورد بحث در این سند برای ارائه توصیه های شخصی فیلم به کاربران بر اساس ترجیحات، سابقه مشاهده و سایر عوامل مرتبط طراحی شده است. هدف این سیستم افزایش تجربه کاربر با پیشنهاد فیلمهایی است که احتمالاً مورد علاقه کاربر هستند، بنابراین تعامل و رضایت کاربر را افزایش می دهد.

در آینده، تکنیک های یادگیری عمیق به طور گسترده برای حل سایر برنامه های چالش برانگیز، مانند طبقه بندی ویدئو، تجزیه و تحلیل شبکه های اجتماعی، طبقه بندی تصاویر در سطح کاملاً جدید، و استنتاج های منطقی استفاده خواهند شد. همچنین، ما ممکن است طرح خود را به عنوان محصول یا خدمات ارائه کنیم. به شرکتی که به یک سیستم توصیه موثر نیاز دارد.

2. مروری بر سیستم های توصیه فیلم

2.1 اهداف سیستم (System objectives)

هدف سیستم توصیه فیلم دستیابی به اهداف زیر است:

- ارائه توصیه های شخصی فیلم به کاربران
- افزایش تعامل و رضایت کاربر
- بهبود فرآیند کشف و انتخاب فیلم
- افزایش تجربه کاربر در پلت فرم

2.2 اجزای سیستم (System components)

سیستم توصیه فیلم از اجزای کلیدی زیر تشکیل شده است:

- نمایه کاربر: هر کاربر یک نمایه دارد که حاوی اطلاعاتی در مورد ترجیحات، سابقه مشاهده و سایر داده های مرتبط است.
- پایگاه داده فیلم: پایگاه داده ای جامع که حاوی اطلاعاتی درباره فیلم ها از جمله ابرداده ها مانند ژانر، سال انتشار، کارگردان، بازیگران و رتبه بندی کاربران است.

- Recommendation Engine: جزء اصلی سیستم که مسئول تولید توصیه های فیلم بر اساس داده ها و ترجیحات کاربر است.

- رابط کاربری: رابطی که از طریق آن کاربران با سیستم برای مرور و انتخاب فیلم ها تعامل دارند.

2.3 اجزای شبکه عصبی (Neural Network components)

- **Inputs:** ورودی های معماری شبکه ما دو بردار n بعدی هستند که n تعداد فیلم های موجود در مجموعه داده فیلم، مانند پایگاه داده MovieLens است. یک بردار نمایه کاربر خاصی را رمزگذاری می کند و هر بعد نشان دهنده امتیازی است که کاربر برای یک فیلم خاص داده است (یا یک صفر برای نشان دادن اینکه هیچ امتیازی داده نشده است). بردار دیگر یک رمزگذاری تک داغ از یک فیلم خاص است (به عنوان مثال، یک بردار با یک بعد "گرم" منفرد روی ۱ تنظیم شده و همه مقادیر دیگر روی صفر تنظیم شده است). این دو بردار درخواست می کنند که شبکه یک امتیاز برای یک کاربر خاص برای یک فیلم خاص پیش بینی کند.

یکی از مزیت های این قالب ورودی این است که می توانیم بدون یک رتبه بندی از یک نمایه کاربر شناخته شده کار کنیم و از رتبه بندی شناخته شده برای آیتم پنهان شده به عنوان مثال برچسب گذاری شده استفاده کنیم. در نتیجه، با وجود اینکه ما فقط ۲۷۰۰۰۰ پروفایل کاربری داریم، هر یک از ۲۶۰۰۰۰۰۰ رتبه بندی فردی یک نمونه قطار است.

- **Hidden layers:** راه های مختلفی برای ساختار یک شبکه عصبی پیش خور ساده وجود دارد. ما با تعدادی از لایه های استاندارد کاملاً متصل شروع می کنیم. با این حال، ما همچنین با ساختارهای جایگزین، مانند ResNets، که در حال حاضر نتایج پیشرفته ای را در زمینه های دیگر مانند تشخیص تصویر به دست می آوریم، آزمایش می کنیم.

- **Output:** دو احتمال اصلی برای خروجی شبکه ما وجود دارد. اولین مورد این است که این مشکل را به عنوان یک مشکل طبقه بندی در نظر بگیریم، با پنج کلاس مختلف که نشان دهنده پنج رتبه بندی شروع است که در داده ها وجود دارد. تحت این معماری، ما پنج خروجی شبکه خود را به عنوان احتمالات لاگ غیر عادی در نظر می گیریم و از آنتروپی متقاطع به عنوان تابع ضرر استفاده می کنیم.

با معماری اصلی شبکه عصبی که در بالا معرفی شد، معماری یادگیری عمیق را که به عنوان جایگزینی برای رویکرد همسایگی مبتنی بر کاربر پیشنهاد شده است، توصیف می کنیم. ابتدا ابعاد ورودی و خروجی شبکه عصبی را در نظر می گیریم. به منظور به حداکثر رساندن مقدار داده های آموزشی که می توانیم به شبکه تغذیه کنیم، یک مثال آموزشی را نمایه کاربر در نظر می گیریم (یعنی یک ردیف از ماتریس مورد کاربر R) با یک رتبه بندی در نظر گرفته نشده است. از دست دادن شبکه در آن مثال آموزشی باید با توجه به رتبه بندی تک تک شده محاسبه شود. نتیجه این امر این است که هر رتبه بندی فردی در مجموعه آموزشی به جای هر کاربر، با یک مثال آموزشی مطابقت دارد.

از آنجایی که ما به آنچه اساساً یک رگرسیون است علاقه مندیم، ما انتخاب می کنیم که از ریشه میانگین مربعات خطا (RMSE) با توجه به رتبه بندی های شناخته شده به عنوان تابع ضرر استفاده کنیم. در مقایسه با میانگین خطای مطلق، ریشه میانگین مربعات خطا به شدت پیش بینی هایی را که دورتر هستند جریمه می کند. ما معتقدیم که این در زمینه سیستم توصیه کننده خوب است، زیرا پیش بینی رتبه بندی بالا برای آیتمی که کاربر از آن لذت نمی برد، به طور قابل توجهی بر کیفیت توصیه ها تأثیر می گذارد. از سوی دیگر، خطاهای کوچکتر در پیش بینی احتمالاً منجر به توصیه هایی می شود که هنوز مفید هستند - شاید رگرسیون دقیقاً درست نباشد، اما حداقل بالاترین رتبه پیش بینی شده احتمالاً به کاربر مربوط است.

2.4 اهمیت توصیه های شخصی

این بخش به اهمیت توصیه های فیلم شخصی شده عمیق تر می پردازد. این بحث می کند که چگونه پیشنهادات متناسب با سلیقه کاربران می تواند تعامل آن ها با پلتفرم را افزایش دهد و منجر به رضایت بیشتر کاربر، استفاده طولانی مدت و بهبود حفظ مشتری شود. همچنین تأکید می کند که چگونه توصیه های شخصی شده، کشف محتوای جدید و مرتبط را تسهیل می کنند.

3. عملکردها

1. User side Func.:

3.1 ثبت نام کاربر و ایجاد پروفایل (registration)

کاربران ملزم به ثبت نام در پلتفرم و ایجاد پروفایل هستند. در طی مراحل ثبت نام، کاربران اطلاعات اولیه مانند نام، سن و جنسیت را ارائه می دهند. آنها همچنین می توانند ترجیحات فیلم، ژانرهای مورد علاقه و بازیگران خود را مشخص کنند.

3.2 انتخاب و مشاهده فیلم (Movie selection)

کاربران می توانند با استفاده از فیلترهای مختلف مانند ژانر، سال انتشار و رتبه بندی، پایگاه داده فیلم را مرور کنند. آنها می توانند اطلاعات دقیق درباره هر فیلم، از جمله خلاصه، بازیگران و رتبه بندی کاربران را مشاهده کنند. کاربران می توانند فیلمی را برای تماشا انتخاب کنند، به آن امتیاز دهند و بازخورد ارائه کنند.

3.3 توصیه های فیلم (Movie recommendation)

موتور توصیه پروفایل های کاربر، تاریخچه مشاهده و تنظیمات برگزیده را تجزیه و تحلیل می کند تا توصیه های فیلم شخصی سازی شود. توصیه ها بر اساس فیلتر مشارکتی، فیلتر مبتنی بر محتوا یا رویکردهای ترکیبی است. کاربران می توانند توصیه هایی را در پلتفرم یا از طریق اعلان های ایمیل شخصی دریافت کنند.

3.4 سیستم رتبه بندی و بازخورد (System rating)

کاربران می توانند به فیلم هایی که تماشا کرده اند امتیاز بدهند و بازخورد ارائه کنند. این اطلاعات برای بهبود الگوریتم توصیه و ارائه پیشنهادات بهتر در آینده استفاده می شود. کاربران همچنین می توانند میانگین امتیاز و بازخورد سایر کاربران را برای کمک به فرآیند انتخاب فیلم خود مشاهده کنند.

II. Key Sys side Func. (اجزای کلیدی سیستم های توصیه فیلم):

3.1 جمع آوری داده ها و پیش پردازش (Data collection & preprocessing)

این بخش فرآیند جمع آوری و پیش پردازش داده ها را برای سیستم های توصیه فیلم توضیح می دهد. این منابع داده های مختلف، از جمله رتبه بندی کاربر، ابرداده فیلم و نمایه های کاربر را بررسی می کند. اهمیت تکنیک های تمیز کردن و عادی سازی داده ها، همراه با استراتژی هایی برای مدیریت داده های از دست رفته برجسته می شود.

3.2 استخراج و نمایش ویژگی (Feature extraction)

در اینجا، استخراج ویژگی های معنادار از ویژگی های فیلم مورد بحث قرار می گیرد. تکنیک های مختلف برای نمایش فیلم ها و کاربران در سیستم توصیه توضیح داده شده اند، از جمله نمایش های برداری، جاسازی ها و مهندسی ویژگی.

3.3 الگوریتم ها و تکنیک های توصیه (Recommendation algorithms)

این زیربخش به بررسی الگوریتم ها و تکنیک های توصیه محبوب مورد استفاده در سیستم های توصیه فیلم می پردازد. الگوریتم های فیلتر مبتنی بر محتوا مانند TF-IDF و شباهت کسینوس به تفصیل مورد بحث قرار گرفته اند. تکنیک های فیلتر مشترک، از جمله روش های فاکتورسازی مبتنی بر کاربر، مبتنی بر آیتم و ماتریس نیز توضیح داده شده اند. علاوه بر این، رویکردهای توصیه ترکیبی و چارچوب های الگوریتمی معرفی شده اند که بر توانایی آن ها در ترکیب نقاط قوت تکنیک های فیلتر مبتنی بر محتوا و مشارکتی تأکید می کنند.

3.4 معیارهای ارزیابی برای ارزیابی اثربخشی توصیه ها (Evaluation & assessing)

برای ارزیابی اثربخشی سیستم های توصیه، معیارهای ارزیابی ضروری است. این زیربخش معیارهای رایج مورد استفاده مانند دقت، یادآوری، امتیاز F1 و دقت میانگین را ارائه می کند. همچنین تکنیک های ارزیابی آفلاین، مانند روش نگه داری و اعتبارسنجی متقابل، و همچنین رویکردهای ارزیابی آنلاین، از جمله تست A/B و مطالعات کاربر را مورد بحث قرار می دهد.

۴. الگوریتم های توصیه

در این قسمت انواع مختلف سیستم های توصیه فیلم به تفصیل توضیح داده شده است. فیلترینگ مبتنی بر محتوا مورد بحث قرار می گیرد و نحوه استفاده از ویژگی های فیلم مانند ژانر، کارگردان و بازیگران برای ارائه توصیه ها را برجسته می کند. فیلتر مشارکتی نیز مورد بررسی قرار می گیرد و بر اتکال آن به رفتار کاربر و ترجیحات برای تولید پیشنهادات تأکید می کند. سیستم های ترکیبی، که رویکردهای فیلتر مبتنی بر محتوا و مشارکتی را ترکیب می کنند، معرفی شده اند و مزایای خود را در ارائه توصیه های دقیق و متنوع نشان می دهند.

۴.۱ توضیح رویکرد فیلترینگ مشارکتی

۴.۱.۱ فیلتر مشارکتی (Collaborative Filtering)

فیلترینگ مشارکتی یک الگوریتم توصیه محبوب است که از رفتار و ترجیحات کاربران مشابه برای تولید توصیه‌ها استفاده می‌کند. کاربرانی را که اولویت‌های فیلم مشابهی دارند شناسایی می‌کند و فیلم‌هایی را پیشنهاد می‌کند که این کاربران مشابه آن‌ها را دوست داشته‌اند یا به آن‌ها امتیاز بالایی داده‌اند. فیلتر مشارکتی می‌تواند بر اساس شباهت کاربر-کاربر یا شباهت مورد-اقلام باشد.

I. فیلتر مشارکتی مبتنی بر کاربر:

فیلتر مشارکتی مبتنی بر کاربر بر یافتن کاربران مشابه بر اساس تعاملات گذشته آنها و توصیه مواردی که آن کاربران مشابه دوست داشته‌اند یا با آنها تعامل داشته‌اند، تمرکز دارد. این شامل مراحل زیر است:

آ. اندازه‌گیری شباهت: شباهت بین کاربران با استفاده از تکنیک‌های مختلفی مانند شباهت کسینوس، ضریب همبستگی پیرسون یا شباهت جاکارد محاسبه می‌شود. این معیارها شباهت بین کاربران را بر اساس رتبه بندی یا تعاملات آنها کمیت می‌کند.

ب. انتخاب نزدیکترین همسایه: هنگامی که شباهت بین کاربران محاسبه شد، نزدیکترین همسایگان یا مشابه ترین کاربران به کاربر هدف انتخاب می‌شوند. تعداد همسایگانی که باید در نظر گرفته شوند را می‌توان با استفاده از آستانه‌های از پیش تعریف شده یا روش‌های آماری تعیین کرد.

ج. تولید توصیه: توصیه‌ها با تجمع ترجیحات نزدیکترین همسایگان ایجاد می‌شوند. این را می‌توان با محاسبه میانگین وزنی رتبه بندی آنها یا استفاده از سایر تکنیک‌های تجمع انجام داد.

II. فیلتر مشارکتی مبتنی بر آیتم:

فیلتر مشارکتی مبتنی بر آیتم بر یافتن موارد مشابه بر اساس تعاملات گذشته آنها و توصیه موارد مشابه با مواردی که کاربر قبلاً دوست داشته یا با آنها تعامل داشته است، تمرکز دارد. مراحل مربوط به فیلترینگ مشارکتی مبتنی بر آیتم به شرح زیر است:

آ. محاسبه شباهت آیتم‌ها:

تشابه بین آیتم‌ها با استفاده از تکنیک‌هایی مانند شباهت کسینوس یا ضریب همبستگی پیرسون اندازه‌گیری می‌شود. این معیارها شباهت بین موارد را بر اساس رتبه‌بندی یا تعامل کاربرانی که با هر دو مورد تعامل داشته‌اند، نشان می‌دهد.

ب انتخاب بیشترین موارد مشابه:

مشابه ترین موارد با مواردی که کاربر با آنها تعامل داشته انتخاب می شود. این می تواند بر اساس آستانه های از پیش تعریف شده یا روش های آماری باشد.

ج. تولید توصیه:

توصیه ها با در نظر گرفتن مواردی که مشابه مواردی هستند که کاربر قبلاً دوست داشته یا با آنها تعامل داشته است، ایجاد می شود. این رویکرد فرض می کند که اگر کاربر یک مورد را دوست داشته باشد یا با آن تعامل داشته باشد، احتمالاً به موارد مشابه علاقه مند است.

۴.۱.۲ مزایا و محدودیت های سیستم های فیلتر مشارکتی

مزایای سیستم های فیلتر مشترک، از جمله دقت و توانایی آنها برای ارائه توصیه های متنوع، بررسی می شود. چالش های ناشی از مشکل شروع سرد و پراکندگی داده ها در فیلتر مشترک، همراه با استراتژی های کاهش بالقوه نیز بررسی می شوند.

۴.۱.۳ پیشرفت ها و ارتقاها در توصیه های فیلتر مشارکتی

پیشرفت ها و پیشرفت ها در توصیه های فیلتر مشارکتی بر رسیدگی به چندین چالش و بهبود دقت توصیه ها متمرکز شده است. در اینجا برخی از پیشرفت های کلیدی در فیلتر مشارکتی وجود دارد: مانند مقیاس پذیری، پراکندگی داده ها و مشکل شروع سرد. تکنیک هایی مانند فاکتورسازی ماتریس، رویکردهای مبتنی بر همسایگی، مدل های ترکیبی، فیلترسازی آگاه از زمینه، یادگیری عمیق و راه حل هایی برای مشکل شروع سرد به طور قابل توجهی دقت و اثربخشی سیستم های توصیه مبتنی بر فیلتر مشارکتی را بهبود بخشیده است. این پیشرفت ها همچنان به تکامل و اصلاح سیستم های توصیه فیلم کمک می کند.

آ. فاکتورسازی ماتریسی:

تکنیک های فاکتورسازی ماتریس، مانند تجزیه ارزش منفرد (SVD) و فاکتورسازی ماتریس غیر منفی (NMF)، به طور گسترده ای برای بهبود فیلتر مشترک استفاده شده است. این تکنیک ها ماتریس تعامل کاربر-مورد را به ماتریس های با ابعاد پایین تر تجزیه می کنند و عوامل پنهانی را که نشان دهنده ترجیحات کاربر و ویژگی های آیتم هستند، می گیرند. با استفاده از این عوامل پنهان، روش های فاکتورسازی ماتریسی می توانند رتبه بندی های گمشده را پیش بینی کنند و توصیه های شخصی سازی شده را ایجاد کنند.

ب یادگیری عمیق در فیلتر مشارکتی:

تکنیک های یادگیری عمیق، مانند شبکه های عصبی عمیق و شبکه های عصبی تکراری، برای فیلتر کردن مشارکتی برای ثبت الگوها و وابستگی های پیچیده در تعاملات کاربر-مورد استفاده شده اند. این مدل ها می توانند به طور خودکار نمایش سلسله مراتبی ترجیحات کاربر و ویژگی های آیتم را بیاموزند که منجر به بهبود دقت توصیه ها می شود. رویکردهای فیلتر مشارکتی مبتنی بر یادگیری عمیق نتایج امیدوارکننده ای را در گرفتن ترجیحات طولانی مدت کاربر، مدیریت داده های پراکنده و مدل سازی رفتارهای متوالی کاربر نشان داده اند.

۴.۲ توضیح رویکرد فیلترینگ مبتنی بر محتوا (Content-based Filtering)

۴.۲.۱ فیلترینگ مبتنی بر محتوا

فیلتر مبتنی بر محتوا فیلم‌ها را به کاربران بر اساس ویژگی‌ها و ویژگی‌های فیلم‌هایی که قبلاً دوست داشته‌اند یا امتیاز بالایی داده‌اند توصیه می‌کند. محتوای و ابرداده فیلم‌ها، مانند ژانر، کارگردان، بازیگران و خلاصه داستان را تجزیه و تحلیل می‌کند تا شباهت‌های بین فیلم‌ها را شناسایی کند و بر اساس این شباهت‌ها توصیه‌هایی ارائه کند.

۴.۲.۲ مزایا و محدودیت‌های سیستم‌های مبتنی بر محتوا

در اینجا، مزایای سیستم‌های توصیه مبتنی بر محتوا، از جمله توانایی آنها در ارائه توصیه‌های بی‌نظیر و استحکام آنها در برابر پراکندگی داده‌ها، مورد بحث قرار می‌گیرد. همچنین آنها به هیچ اطلاعاتی در مورد سایر کاربران نیاز ندارند. محدودیت‌های سیستم‌های مبتنی بر محتوا، مانند تنوع محدود و اتکای بیش از حد به نمایه‌های کاربر، نیز بررسی می‌شوند. علاوه بر این، آنها نمی‌توانند علایق در حال تغییر کاربر را در طول زمان جلب کنند.

۴.۲.۳ پیشرفت‌ها و پیشرفت‌ها در توصیه‌های مبتنی بر محتوا

در سال‌های اخیر، پیشرفت‌ها و پیشرفت‌های زیادی در توصیه‌های مبتنی بر محتوا صورت گرفته است. برخی از این موارد شامل انواع جدیدی از اطلاعات جانبی (مانند ابرداده یا محتوای تولید شده توسط کاربر) برای توصیه یا رویکردهای الگوریتمی جدید برای پردازش اطلاعات موجود است. برخی از این موارد عبارتند از یادگیری عمیق (DL)، رویکردهای مبتنی بر فرامسیر، و رویکردهای مرتبط با داده‌های الگوریتمی.

۴.۳ توضیح رویکرد سیستم توصیه ترکیبی (Hybrid Filtering)

۴.۳.۱ رویکردهای ترکیبی

رویکردهای ترکیبی فیلتر مشارکتی و فیلتر مبتنی بر محتوا را برای استفاده از نقاط قوت هر دو الگوریتم ترکیب می‌کنند. هدف آنها ارائه توصیه‌های دقیق‌تر و متنوع‌تر با ترکیب منابع متعدد داده‌ها و اطلاعات است.

۴.۳.۲ مزایای رویکردهای ترکیبی

در اینجا، مزایای سیستم‌های توصیه ترکیبی با جزئیات بیشتری مورد بحث قرار می‌گیرد. موضوعات تحت پوشش عبارتند از بهبود دقت توصیه به دست آمده از طریق ترکیبی از تکنیک‌ها، توانایی پرداختن به محدودیت‌های رویکردهای توصیه فردی، و افزایش رضایت کاربر ناشی از توصیه‌های متنوع و شخصی.

5. کارهای مشابه

متداول ترین روش انجام فیلترینگ مشترک، استفاده از رویکرد k -نزدیکترین همسایه بین کاربران است. با این تکنیک، ابتدا با یک ماتریس کاربر-مورد R شروع می‌شود، جایی که R_{ij} رتبه‌بندی کاربر i را برای آیتم j می‌دهد و مقدار \cdot نشان می‌دهد که رتبه‌بندی خاصی وجود ندارد. از R یک ماتریس شباهت کاربر-کاربر S محاسبه می‌شود، که در آن S_{ij} شباهت بین کاربر i و کاربر j است که می‌تواند با $R \cdot RT$ محاسبه شود. توجه داشته باشید که استفاده از سایر معیارهای فاصله، مانند معیار تشابه همبستگی یا شباهت کسینوس، برای پر کردن S نیز موثر است. هنگامی که S محاسبه شد، می‌توانیم رتبه‌بندی کاربر i را برای مورد j با محاسبه $RTj \cdot Si$ پیش‌بینی کنیم، که اساساً میانگین رتبه‌بندی سایر کاربران را برای مورد j با وزن شباهت آنها به کاربر i محاسبه می‌کند.

همچنین می‌توانیم از k شبیه‌ترین کاربران به کاربر i برای پیش‌بینی امتیاز مورد j استفاده کنیم. از نظر تجربی، این عملکرد بهتر از میانگین وزنی روی همه کاربران است، اگرچه برای محاسبه k نزدیک‌ترین همسایه‌ها، مقداری کار اضافی در زمان آزمایش لازم است. این رویکرد بر این فرض تکیه دارد که اگر دو کاربر به یک مورد مشابه امتیاز دهند، احتمالاً به موارد دیگر نیز امتیاز مشابهی می‌دهند. در مقیاس، ساختارهای داده مانند درختان توپ و درختان $k-d$ (یک درخت پارتیشن فضایی دودویی در ابعاد k) می‌توانند برای محاسبه موثرتر همسایگان محلی بین پروفایل‌های کاربر استفاده شوند.

یک رویکرد جایگزین k -نزدیکترین همسایه در عوض شباهت بین جفت آیتم‌ها (در مقابل کاربران) را با این ایده که کاربرانی که یک آیتم خاص را دوست دارند، موارد مشابه را دوست خواهند داشت محاسبه می‌کند. با این رویکرد ما یک ماتریس شباهت مورد به مورد I را به صورت $RT \cdot R$ محاسبه می‌کنیم. مانند قبل، می‌توانیم از معیارهای مشابه دیگری نیز برای پر کردن I استفاده کنیم. برای پیش‌بینی رتبه‌بندی کاربر i در مورد j ، می‌توانیم $R_{ij} \cdot I_j$ را محاسبه کنیم که میانگین رتبه‌بندی‌های ارائه‌شده توسط کاربر i را با وزن شباهت نشان می‌دهد. آن موارد به مورد j از آنجایی که تعداد کاربران بسیار بیشتر از موارد در یک سیستم توصیه‌گر است، فیلتر مشترک کاربر-کاربر می‌تواند عملکرد بیشتری داشته باشد.

یکی دیگر از روش‌های رایج برای انجام فیلترینگ مشارکتی، فاکتورسازی ماتریس است. با این تکنیک، یک ماتریس کاربر-آیتم به دو ماتریس با بعد داخلی که نشان دهنده برخی عوامل پنهان با استفاده از تکنیک‌هایی مانند تجزیه ارزش منفرد (SVD) است، فاکتور می‌شود. فاکتوربندی حاصل، هم کاربران و هم اقلام را از نظر عوامل پنهان نشان می‌دهد، به گونه‌ای که می‌توان از آنها برای توصیه موارد جدید استفاده کرد. همانند رویکردهای همسایگی مورد-آیتم، آزمایش‌های اولیه ما روی رتبه‌بندی فیلم نشان می‌دهد که رویکردهای همسایگی کاربر-کاربر نسبت به فاکتورسازی ماتریس برتری دارند.

یادگیری عمیق بسیاری از زمینه‌های علوم کامپیوتر از جمله پردازش زبان طبیعی را متحول کرده است. با وجود این، یادگیری عمیق در حوزه سیستم‌های توصیه‌گر نسبتاً جدید است و توجه زیادی به آن نشده است. یک مدل یادگیری عمیق مشارکتی (CDL) پیشنهاد کنید که به طور مشترک عمیق را انجام دهد. یادگیری بازنمایی برای اطلاعات محتوا و فیلتر مشارکتی برای ماتریس رتبه‌بندی. CDL با ما متفاوت است، زیرا اولی به اطلاعات محتوا متکی است، در حالی که ما نه. یک سیستم توصیه‌یادگیری عمیق را با توجه به تاریخچه مرور وب و جستجوهای ارائه شده توسط کاربران معرفی کنید. آنها شباهت بین کاربران و آیتم‌های ترجیحی آنها را با نداشت کاربران و آیتم‌ها در یک فضای پنهان به حداکثر می‌رساند. محدودیتی که بر این رویکرد اعمال می‌شود این است که تاریخچه مرور و جستجوهای کاربران مورد نیاز است که همیشه در دسترس نیستند. یک مدل شبکه عصبی عمیق ایجاد کنید که ویژگی‌های محتوای آیتم‌ها را برای پیش‌بینی رتبه‌بندی آیتم‌های شروع سرد استخراج می‌کند. سیستم توصیه‌گر ما متفاوت است، زیرا ما با محتوای کاربر سروکار نداریم.

همانطور که در شکل ۱ نشان داده شده است، انواع بسیاری از الگوریتم ها در سیستم های مختلف توصیه فیلم در دو دهه گذشته آزمایش و آزمایش شده اند.

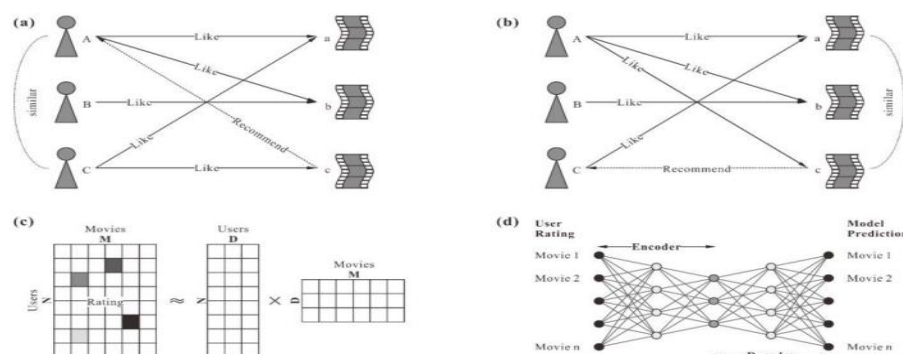


Figure 1. Typical recommendation methods. (a) User-based collaborative filtering, (b) Item-based collaborative filtering, (c) Model-based collaborative filtering, (d) Deep-learning neural networks.

۶. معماری DNN

(۱) *Autoencoder*: یکی از مدل های یادگیری عمیق موجود، مدل شبکه عصبی عمیق (DNN) است. DNN یک مدل پرسپترون چند لایه (MLP) با لایه های پنهان بسیاری است. منحصر به فرد بودن DNN به دلیل تعداد بیشتر واحدهای پنهان و تکنیک های اولیه سازی پارامترهای بهتر است. یک مدل DNN با تعداد زیادی واحد پنهان می تواند قدرت مدل سازی بهتری داشته باشد. اگرچه پارامترهای آموخته شده مدل DNN یک بهینه محلی است که به داده های آموزشی بیشتر و قدرت محاسباتی بیشتری نیاز دارد، اما می تواند بسیار بهتر از مدل هایی با واحدهای پنهان کمتر عمل کند. Deep Auto Encoder نوع خاصی از DNN است. (شکل ۱ را برای نمونه Autoencoder ببینید).

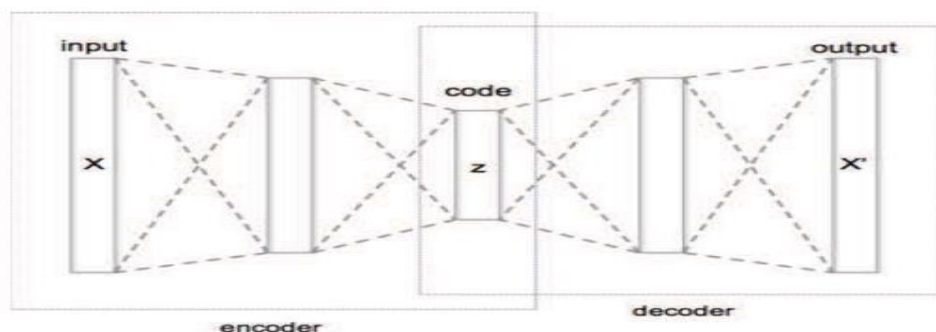


Figure 1. An autoencoder with three fully-connected hidden layers

رمزگذار خودکار (Autoencoder) یک شبکه عصبی است که آموزش داده شده است تا ورودی خود را در خروجی خود کپی کند، با هدف معمول کاهش ابعاد، یعنی فرآیند کاهش تعداد متغیرهای تصادفی مورد بررسی. دارای یک تابع رمزگذار برای ایجاد یک لایه پنهان (یا چندین لایه پنهان) است که حاوی کدی برای توصیف ورودی است. یک رمزگشا وجود دارد که یک بازسازی ورودی از لایه پنهان ایجاد می کند. رمزگذار خودکار می تواند با داشتن یک لایه مخفی کوچکتر از لایه ورودی مفید واقع شود و آن را وادار کند تا با یادگیری همبستگی در داده ها، نمایشی فشرده از داده ها در لایه پنهان ایجاد کند. این رمزگذار خودکار نوعی یادگیری بدون نظارت است، به این معنی که یک رمزگذار خودکار فقط به داده های بدون برچسب نیاز دارد، که مجموعه ای از داده های ورودی است نه جفت های ورودی-خروجی. از طریق یک الگوریتم یادگیری

بدون نظارت، برای بازسازی های خطی، رمزگذار خودکار سعی می کند تابعی را یاد بگیرد تا اختلاف میانگین مربع ریشه را به حداقل برساند. برای محاسبه ریشه میانگین مربعات خطا (RMSE) یک مدل یادگیری ماشینی، می توانیم عملکرد مدل را اندازه گیری کنیم. RMSE به صورت تعریف شده است :

$$RMSE = \frac{1}{m} \sum_i (\hat{y} - y)_i^2, \quad \hat{y} = \mathbf{w}^T \mathbf{x} \quad (1)$$

که در آن $\mathbf{w} \in \mathbb{R}^n$ بردار پارامترها است، $\mathbf{x} \in \mathbb{R}^n$ برداری است که برای پیش بینی مقدار اسکالر $y \in \mathbb{R}$ استفاده می شود و \hat{y} مقداری است که یک مدل یادگیری ماشینی پیش بینی می کند که مقدار اسکالر $y \in \mathbb{R}$ باید چقدر باشد. توجه داشته باشید که RMSE زمانی که $\hat{y} = y$ به ۰ کاهش می یابد و با افزایش فاصله اقلیدسی بین مقادیر پیش بینی شده و مقادیر هدف، خطا افزایش می یابد.

۲) پرسپترون چندلایه (Multilayer Perceptron): در ابتدا، معماری سیستم توصیه گر ما شامل ورودی از ردیف ماتریس کاربر-مورد R با رتبه بندی برخی از آیتم های J به همراه یک پرس و جوی رمزگذاری شده یک طرفه است که نشان می دهد شبکه باید پیش بینی کند. امتیاز برای کاربر i در مورد j . متأسفانه، ثابت شده است که آموزش این معماری دشوار است، زیرا شبکه باید یاد بگیرد که نه تنها نمایه های کاربر، بلکه تعامل بین آن پروفایل ها و ورودی های پرس و جو را نیز درک کند. با توجه به ریشه میانگین مربعات خطا در داده های آموزشی، ما هرگز با این معماری به ضرر کمتر از ۱.۲ نرسیدیم.

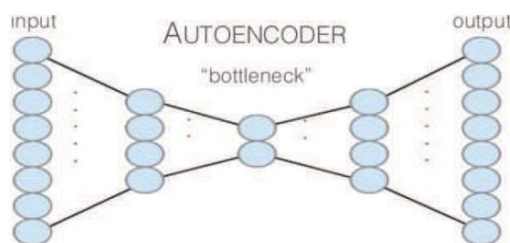


Figure 2. Overview of the network architecture of our recommender

در عوض، ما از مفهوم رمزگذار خودکار برای طراحی معماری شبکه عصبی خود الهام می گیریم (شکل ۲ را ببینید). این معماری ساده یک ورودی دریافت می کند و آن را به تعدادی از لایه های پنهان کاملاً متصل که شامل یک «گلوگاه» است متصل می کند. این گلوگاه یک لایه پنهان است که ابعاد بسیار کوچکتری نسبت به ورودی دارد. سپس خروجی شبکه مجدداً گسترش می یابد تا ابعادی مشابه ورودی داشته باشد. سپس شبکه برای یادگیری تابع هویت آموزش داده می شود، با این ایده که برای اینکه شبکه بتواند تابع هویت را از طریق گلوگاه محاسبه کند، باید یک نمایش متراکم از ورودی را بیاموزد. بنابراین، رمزگذار خودکار را می توان به عنوان چیزی شبیه به تکنیک کاهش ابعاد در نظر گرفت. همچنین می توان امیدوار بود که لایه گلوگاه چیز مفیدی در رابطه با ساختار زیربنایی ورودی یاد بگیرد. به عنوان مثال، یک نورون در لایه گلوگاه ممکن است چیزی مرتبط با ژانر یک فیلم یا گروه بندی فیلم های مشابه را نشان دهد. توجه داشته باشید که ما علاقه ای به یادگیری محاسبه یک تابع هویت نداریم — در نهایت، هدف ما پیش بینی رتبه بندی های گمشده است، نه بازتولید صفرها در بردارهای ورودی. در نتیجه، در حالی که معماری شبکه نهایی ما شبیه یک رمزگذار خودکار با لایه های پنهان گلوگاه و ابعاد منطبق بر ورودی و خروجی است، شبکه در واقع با استفاده از یک تابع ضرر برای رگرسیون (یعنی RMSE) با هدف یادگیری پیش بینی رتبه بندی های گمشده آموزش داده می شود. به طور خاص، نمونه های آموزشی برای شبکه، نمایه های کاربر با یک رتبه بندی پنهان است، و خروجی رتبه بندی های پیش بینی شده برای همه فیلم های مجموعه داده است. در حالی که انتظار می رود شبکه رتبه بندی هر فیلم را بر اساس نمایه کاربر پیش بینی کند، ما فقط برای رتبه بندی پنهان شده پاسخ داریم. در نتیجه، ما فقط هنگام یادگیری از مثال آموزشی، ضرر را برای رتبه از دست رفته منتشر می کنیم.

۳) سیستم توصیه‌کننده یادگیری عمیق (The Deep Learning Recommender System): رتبه‌بندی‌های خودداری این نتیجه ناگوار را به همراه دارد که مدل یادگیری عمیق ما فقط می‌تواند رتبه‌بندی‌هایی را برای فیلم‌هایی مشابه آنچه کاربر تماشا کرده است یاد بگیرد، زیرا عملکرد ضرر مستقیماً تحت تأثیر خروجی غیرمرتبط قرار نمی‌گیرد. فیلم‌ها به دلیل لایه تنگنا، مدل باید تا حدی تعمیم یابد، اما این مدل ممکن است برای فیلم‌هایی که به شدت متفاوت از فیلم‌هایی است که کاربر واقعاً رتبه‌بندی کرده است، مشکل داشته باشد. در حالی که کاربران فیلم‌ها را تماشا می‌کنند و امتیاز پایینی دارند، بیشتر اوقات به بیش از چند صد مورد امتیاز نمی‌دهند و از تماشای فیلم‌های کاملاً غیر مرتبط اجتناب می‌کنند، بنابراین ممکن است پیش‌بینی رتبه‌بندی برای فیلم‌های کاملاً نامرتبط برای مدل دشوار باشد. برای مقاصد تابع ضرر ما، که ریشه میانگین مربعات خطا در رتبه‌بندی‌های شناخته شده است، این واقعیت که شبکه ما ممکن است نحوه خروجی رتبه‌بندی فیلم‌های کاملاً نامرتبط را یاد نگیرد، به نظر نمی‌رسد که روی از دست دادن تست تأثیر بگذارد، احتمالاً به این دلیل که فیلم‌های موجود در آزمایش داده‌ها به اندازه‌ای مرتبط هستند که الگوهای آموخته‌شده از داده‌های آموزشی به رتبه‌بندی‌های داده‌های آزمون تعمیم می‌یابند. البته، ممکن است بر رتبه‌بندی تأثیر بگذارد، بنابراین می‌تواند مطلوب باشد که یک اصطلاح تنظیم به ضرر اضافه شود. با استفاده از این طرح اولیه، ما چندین گونه از این معماری را با استفاده از تعداد لایه‌های مختلف و اندازه‌های مختلف برای لایه گلوگاه آزمایش کرده ایم. جالب‌ترین پارامتر اندازه کوچک‌ترین لایه گلوگاه بود و پس از آزمایش مقادیر مختلف، در نهایت بر روی اندازه گلوگاه ۵۱۲ قرار گرفتیم. از آنجا با تعداد مختلفی از لایه‌های کاملاً متصل آزمایش کردیم، همیشه از توان ۲ برای افزایش و افزایش استفاده کردیم. کاهش ابعاد توپولوژی شبکه نهایی دارای هفت لایه پنهان کاملاً متصل با ابعاد [۴۰۹۶، ۲۰۴۸، ۱۰۲۴، ۵۱۲، ۱۰۲۴، ۲۰۴۸، ۴۰۹۶] است. هر لایه از یک واحد خطی اصلاح شده ۴ به عنوان تابع فعال سازی غیر خطی استفاده می‌کند. وزن‌های اتصال لایه‌های پنهان با استفاده از مقداردهی اولیه خاویر با بایاس‌ها روی صفر تنظیم شدند.

۴) خوشه‌بندی (Clustering): ایده استفاده از کوچک‌ترین لایه گلوگاه در شبکه را به عنوان نوعی خوشه‌بندی طبیعی در نظر گرفته ایم. با فشار دادن ورودی به چنین فضای ابعادی کوچکی، مدل باید لزوماً چیزی در مورد ساختار زیربنایی داده‌های ورودی بیاموزد. فرضیه این بود که با تثبیت یک نورون منفرد در لایه گلوگاه و صفر کردن نورون‌های باقیمانده در لایه گلوگاه، و سپس بهینه‌سازی فضای ورودی برای این فعال‌سازی خاص، می‌توانیم با نمایش فیلم‌هایی که هر خوشه را تحریک می‌کنند، آن ساختار را تجسم کنیم. برای مثال، ما انتظار داریم که ممکن است یک نورون یا مجموعه کوچکی از نورون‌ها وجود داشته باشد که ژانرهای مختلف فیلم یا سبک‌های مختلف فیلم‌شناسی را تحریک کند. جدول I مثالی از چنین "خوشه‌ای" را از بهینه‌سازی ورودی برای راه اندازی یک نورون گلوگاه ارائه می‌دهد. این فیلم‌ها مضمون مشترکی دارند. بدیهی است که برای اینکه این شبکه بتواند رتبه‌بندی فیلم‌ها را به طور دقیق پیش‌بینی کند، باید نوعی ساختار را یاد بگیرد. با این حال، این ساختار بیش از حد انتظار در سراسر لایه گلوگاه توزیع شده است. یکی از راه‌حل‌های بالقوه برای این مشکل، اضافه کردن یک اصطلاح منظم‌سازی به ضرر است که باعث ایجاد پراکندگی در لایه گلوگاه می‌شود.

Table I
A CLUSTER WHEN OPTIMIZING THE INPUT TO TRIGGER A SINGLE
BOTTLENECK NEURON

Jules and Jim (Jules et Jim) (1961)
Frankenstein Must Be Destroyed (1969)
Lolita (1962)
Lawnmower Man, The (1992)
First Knight (1995)
Urban Legends: Final Cut (2000)
Fair Game (1995)
Guinevere (1999)
Paradine Case, The (1947)
400 Blows The (Les quatre cents coups) (1959)

۷. چالش ها و پیشرفت های پیاده سازی در سیستم های توصیه فیلم

۷.۱ مشکل شروع سرد و پراکندگی داده ها (Cold Start Problem and Data Sparsity)

مشکل شروع سرد زمانی رخ می دهد که یک سیستم توصیه فاقد اطلاعات کافی کاربر یا آیتم برای پیش بینی دقیق باشد. به همین دلیل است که به چالش ایجاد توصیه های دقیق برای کاربران جدید یا فیلم های جدید با داده های محدود اشاره دارد. سیستم توصیه باید موقعیت هایی را مدیریت کند که داده های کاربر یا داده های فیلم کافی برای پیش بینی دقیق وجود ندارد.

برای مقابله با این چالش، می توان از توصیه های مبتنی بر محتوا استفاده کرد. به عنوان مثال، یک سیستم توصیه فیلم می تواند از ویژگی های فیلم مانند ژانر، کارگردان و بازیگران برای ارائه پیشنهادات اولیه به کاربران جدید استفاده کند. علاوه بر این، رویکردهای ترکیبی که تکنیک های فیلتر مبتنی بر محتوا و مشارکتی را ترکیب می کنند، می توانند با استفاده از داده های محتوا و رفتار کاربر به کاهش مشکل شروع سرد کمک کنند. سیستم توصیه فیلم نتفلیکس از چنین رویکردهای ترکیبی برای ارائه توصیه هایی حتی برای کاربران جدید با داده های محدود استفاده می کند.

با توجه به پراکندگی داده ها، تکنیک های فیلتر مشارکتی می توانند با اعمال نفوذ بر رفتار و ترجیحات کاربران مشابه، بر این چالش غلبه کنند. به عنوان مثال، اگر یک کاربر تنظیمات ترجیحی فیلم مشابهی را با سایر کاربران نشان دهد، الگوریتم های فیلتر مشترک می توانند فیلم هایی را توصیه کنند که کاربران مشابه از آنها لذت می برند. این رویکرد به ویژه زمانی مؤثر است که همپوشانی محدودی در رتبه بندی کاربران وجود داشته باشد. Amazon Prime Video از فیلتر مشترک برای ایجاد توصیه هایی بر اساس رتبه بندی و رفتار کاربران، حتی در سناریوهایی با داده های پراکنده استفاده می کند.

۷.۲ مقیاس پذیری و توصیه های بلادرنگ (Scalability and Real-Time Recommendations)

همانطور که پایگاه کاربر و پایگاه داده فیلم رشد می کند، سیستم توصیه باید بتواند حجم زیادی از داده ها را مدیریت کند و توصیه های بلادرنگ را بدون کاهش عملکرد قابل توجه ارائه دهد. سیستم های پیشنهادی هنگام برخورد با پایگاه های کاربری بزرگ و کاتالوگ های فیلم گسترده با چالش های مقیاس پذیری مواجه هستند. مقیاس پذیری عاملی حیاتی برای اطمینان از اثربخشی و پاسخگویی سیستم است.

برای پرداختن به این، از روش های محاسباتی توزیع شده و پردازش موازی استفاده می شود. برای مثال، پلتفرم هایی مانند YouTube و Netflix از سیستم های توزیع شده برای مدیریت میلیون ها کاربر و کتابخانه های وسیع محتوا استفاده می کنند و از توصیه های سریع و مقیاس پذیر اطمینان می دهند. توصیه های بلادرنگ با به روزرسانی مداوم تنظیمات برگزیده کاربر و ترکیب فعالیت های اخیر به دست می آیند. به عنوان مثال، Spotify بر اساس عادات گوش دادن فعلی کاربران، توصیه های بلادرنگ ارائه می کند و با تغییر اولویت ها، توصیه ها را به صورت پویا تنظیم می کند.

۷.۳ حریم خصوصی و امنیت داده ها (Privacy and Data Security)

سیستم توصیه فیلم داده‌های کاربر را جمع‌آوری و تجزیه و تحلیل می‌کند تا توصیه‌هایی ایجاد کند. اجرای اقدامات امنیتی مناسب برای محافظت از حریم خصوصی کاربر و اطمینان از اینکه داده‌های کاربر به طور ایمن مدیریت می‌شوند ضروری است. رعایت مقررات حفاظت از داده‌ها باید در اولویت باشد.

برای رفع این مشکل، پروتکل‌های رمزگذاری قوی داده و احراز هویت کاربر را برای محافظت از اطلاعات شخصی اجرا کنید. علاوه بر این، گستره جمع‌آوری داده‌ها را فقط به آنچه برای ایجاد توصیه‌ها ضروری است محدود کنید و به کاربران این امکان را بدهید که در هر زمانی داده‌های خود را انصراف دهند یا حذف کنند. ممیزی‌های منظم و به روز رسانی اقدامات امنیتی نیز باید برای اطمینان از حفاظت مداوم از اطلاعات کاربر اجرا شود.

۷.۴ ترکیب عوامل زمینه و زمانی (Incorporating Context and Temporal Factors)

عوامل زمینه ای و زمانی نقش حیاتی در توصیه‌های فیلم دارند. به عنوان مثال، یک سیستم توصیه می‌تواند زمان روز، شرایط آب و هوایی یا مکان کاربر را برای پیشنهاد فیلم‌های مناسب در نظر بگیرد. اگر کاربر در عصر به پلتفرم دسترسی داشته باشد، سیستم ممکن است فیلم‌های آرامش بخش یا سرگرم کننده را توصیه کند. Netflix با شخصی‌سازی توصیه‌ها بر اساس زمان روز و الگوهای مشاهده تاریخی کاربر، زمینه را در نظر می‌گیرد.

مثال دیگر گنجاندن پویایی‌های زمانی در توصیه‌ها است. با تجزیه و تحلیل رفتار کاربر در طول زمان، سیستم‌های توصیه می‌توانند با تغییر تنظیمات سازگار شوند. الگوریتم توصیه نتفلیکس تاریخچه بازدید کاربر را تجزیه و تحلیل می‌کند و به تدریج توصیه‌ها را بر اساس سلیقه در حال تکامل آنها به روز می‌کند. این تضمین می‌کند که سیستم پیشنهادات فیلم به روز و مرتبط را ارائه می‌دهد.

۷.۵ یادگیری عمیق و پیشرفت‌های یادگیری ماشین (Deep Learning and Machine Learning Advancements)

تکنیک‌های یادگیری عمیق و یادگیری ماشین انقلابی در سیستم‌های توصیه فیلم ایجاد کرده است. شبکه‌های عصبی، مانند شبکه‌های عصبی کانولوشن (CNN) و شبکه‌های عصبی تکراری (RNN)، برای ثبت الگوها و وابستگی‌های پیچیده در داده‌های فیلم استفاده شده‌اند.

به عنوان مثال، نتفلیکس از مدل‌های یادگیری عمیق برای بهبود توصیه‌های فیلم خود استفاده کرد. این شرکت الگوریتمی مبتنی بر ماشین‌های محدود بولتزمن (RBM) برای یادگیری نمایش ویژگی‌های فیلم‌ها و کاربران معرفی کرد. این امکان پیش بینی‌های دقیق تری را با ثبت روابط پیچیده بین فیلم‌ها و ترجیحات کاربران فراهم می‌کرد.

مثال دیگر استفاده از رمزگذارهای خودکار در سیستم‌های توصیه است. رمزگذارهای خودکار می‌توانند نمایش‌های فشرده رتبه بندی یا ویژگی‌های فیلم را بیاموزند. با بازسازی ورودی، آنها می‌توانند الگوهای زیربنایی را شناسایی کرده و توصیه‌هایی ایجاد کنند. رمزگذارهای خودکار با موفقیت در سیستم‌های توصیه فیلم، مانند سیستم‌هایی که توسط Spotify و Last.fm استفاده می‌شوند، استفاده شده‌اند.

۸. ارزیابی و معیارها برای سیستم های توصیه

۸.۱ معیارهای ارزیابی رایج (Commonly Used Evaluation Metrics)

برای ارزیابی عملکرد سیستم های توصیه، معمولاً از چندین معیار استفاده می شود. دقت، یادآوری، امتیاز $F1$ ، و میانگین دقت میانگین (MAP) به طور گسترده برای اندازه گیری اثربخشی توصیه استفاده می شود. این معیارها توانایی سیستم را برای ارائه پیشنهادهای مرتبط و دقیق فیلم ارزیابی می کند. برای مثال، دقت نسبت فیلم های توصیه شده درست را از همه فیلم های توصیه شده اندازه گیری می کند، در حالی که فراخوان نسبت فیلم های به درستی توصیه شده را اندازه گیری می کند.

از همه فیلم های مرتبط امتیاز $F1$ دقت و یادآوری را در یک متریک واحد ترکیب می کند. MAP میانگین دقت را در سطوح مختلف فراخوان محاسبه می کند و به ویژه در هنگام برخورد با توصیه های رتبه بندی شده مفید است.

محققان از این معیارهای ارزیابی در مطالعات توصیه فیلم استفاده کرده اند. به عنوان مثال، مطالعه ای بر روی الگوریتم های توصیه فیلم، دقت و یادآوری رویکردهای فیلترینگ مبتنی بر محتوا و مشارکتی را مقایسه کرد. نتایج ارزیابی نشان داد که فیلتر مشارکتی به دقت بالاتری دست یافت، در حالی که فیلتر مبتنی بر محتوا یادآوری بالاتری را نشان داد.

۸.۲ ارزیابی آفلاین در مقابل ارزیابی آنلاین (Offline Evaluation vs. Online Evaluation)

سیستم های توصیه را می توان با استفاده از روش های آفلاین و آنلاین ارزیابی کرد. ارزیابی آفلاین شامل استفاده از داده های تاریخی برای ارزیابی کیفیت توصیه می شود. این رویکرد امکان ارزیابی سریع و مقرون به صرفه را فراهم می کند، اما تعاملات کاربر را در زمان واقعی ثبت نمی کند. از سوی دیگر، ارزیابی آنلاین شامل انجام آزمایش هایی با کاربران واقعی، جمع آوری بازخورد و اندازه گیری تأثیر توصیه ها بر رفتار کاربر است.

یک تکنیک ارزیابی آنلاین معروف، تست A/B است که در آن دو یا چند الگوریتم توصیه با استفاده از کاربران واقعی مقایسه می شوند. به عنوان مثال، یک پلت فرم استریم ممکن است به طور تصادفی کاربران را به گروه های پیشنهادی مختلف اختصاص دهد و هر گروه را با الگوریتم های متفاوتی ارائه دهد. با اندازه گیری معیارهای تعامل، رضایت و تعامل کاربر، این پلتفرم می تواند عملکرد رویکردهای مختلف توصیه را ارزیابی کند و تصمیم های مبتنی بر داده اتخاذ کند.

۸.۳ چالش ها و ملاحظات در ارزیابی سیستم های توصیه (Challenges and Considerations in Evaluating Recommendation Systems)

ارزیابی سیستم های توصیه چالش های متعددی را به همراه دارد. یکی از چالش ها فقدان داده های حقیقت پایه است، زیرا ترجیحات کاربر ذهنی هستند و دائماً در حال تغییر هستند. برای پرداختن به این موضوع، محققان استفاده از داده های بازخورد ضمنی، مانند کلیک ها، زمان تماشا یا سابقه خرید را به عنوان یک پروکسی برای تنظیمات برگزیده کاربر پیشنهاد کرده اند.

معیارهای تنوع نیز در ارزیابی سیستم های توصیه بسیار مهم هستند. توصیه صرف فیلم های محبوب یا مشابه ممکن است منجر به عدم تنوع در پیشنهادهای شود. معیارهایی مانند پوشش کاتالوگ و تازگی برای اندازه گیری تنوع توصیه ها استفاده می شود. برای مثال، معیار پوشش کاتالوگ ارزیابی می کند که سیستم چقدر فیلم هایی را از ژانرها و دسته های مختلف توصیه می کند.

علاوه بر این، ملاحظات انصاف در ارزیابی حیاتی است. سوگیری ممکن است در توصیه ها ایجاد شود و بر فیلم ها یا گروه های کاربری که کمتر نمایش داده شده اند تأثیر بگذارد. محققان روش ها و معیارهای ارزیابی مبتنی بر عدالت را برای ارزیابی عادلانه بودن سیستم های توصیه پیشنهاد کرده اند. این تضمین می کند که توصیه ها بی طرفانه هستند و منافع مختلف کاربران را برآورده می کنند.

۹. تجربه کاربر و ملاحظات اخلاقی

۹.۱ اهمیت تجربه کاربر در سیستم های توصیه (Importance of User Experience in Recommendation Systems)

تجربه کاربری یک جنبه حیاتی از سیستم های توصیه فیلم است. توصیه های دقیق و شخصی شده، رضایت و تعامل کاربر را افزایش می دهد و منجر به استفاده طولانی مدت و حفظ مشتری می شود. برای مثال، صفحات اصلی شخصی سازی شده و فهرست های پخش انتخاب شده، مانند آنچه توسط Amazon Prime Video ارائه می شود، با ارائه پیشنهادهای فیلم متناسب بر اساس اولویت های فردی و سابقه مشاهده، تجربه کاربر را بهبود می بخشد.

علاوه بر این، توصیه های مؤثر، کشف محتوای جدید و مرتبط را بهبود می بخشد. سیستم های توصیه با پیشنهاد فیلم هایی که مطابق با سلیقه کاربران هستند، در زمان و تلاش کاربران در جستجوی گزینه های مناسب صرفه جویی می کنند. این تجربه کلی کاربر را بهبود می بخشد و تضمین می کند که کاربران فیلم هایی را پیدا می کنند که واقعاً از آنها لذت می برند.

۹.۲ رسیدگی به مسائل مربوط به حریم خصوصی و امنیت داده ها (Addressing Issues of Privacy and Data Security)

سیستم های توصیه فیلم حجم زیادی از داده های کاربر را مدیریت می کنند و نگرانی هایی را در مورد حفظ حریم خصوصی و امنیت داده ها ایجاد می کنند. پلتفرم ها باید اعتماد کاربران را با اجرای اقدامات حفظ حریم خصوصی قوی در اولویت قرار دهند. اطمینان از رضایت کاربر، ناشناس کردن داده ها، و ارائه سیاست های شفاف استفاده از داده بسیار مهم است. پلتفرم هایی مانند Netflix و Amazon Prime Video از مقررات حفظ حریم خصوصی مانند مقررات حفاظت از داده های عمومی (GDPR) برای محافظت از اطلاعات کاربر پیروی می کنند.

امنیت داده ها به همان اندازه مهم است. سیستم های توصیه باید از پروتکل های ذخیره و دسترسی ایمن داده استفاده کنند تا از داده های کاربر در برابر دسترسی غیرمجاز محافظت کنند. تکنیک های رمزگذاری و کنترل های دسترسی دقیق به حفظ محرمانه بودن و یکپارچگی داده های کاربر کمک می کند.

۹.۳ برخورد با سوگیری های الگوریتمی و تنوع در توصیه ها (Dealing with Algorithmic Biases and Diversity in Recommendations)

سوگیری های الگوریتمی می تواند در توصیه های فیلم ظاهر شود و بر نمایش و تنوع تأثیر بگذارد. برای مثال، اگر یک سیستم توصیه در درجه اول فیلم هایی را پیشنهاد می کند که با ترجیحات اکثریت همخوانی دارند، ممکن است فیلم هایی را که به گروه های کم نمایش داده می شوند نادیده بگیرد.

پلتفرم ها با اجرای الگوریتم های آگاه از انصاف و در نظر گرفتن تنوع به طور فعال به این موضوع می پردازند.

معیارها در سیستم های توصیه آنها. هدف آنها ارائه پیشنهادهای عادلانه و متنوع فیلم است که ژانرها، فرهنگ ها و دیدگاه های مختلف را در بر می گیرد.

10. مدل معماری اجرا شده

۱. سیستم توصیه گر پیشنهادی

در این بخش، ما در اینجا معماری سیستم توصیه فیلم خود را بر اساس فناوری یادگیری عمیق و تجزیه و تحلیل داده های چندوجهی پیشنهاد می کنیم. این مدل ویژگی های ویژگی های کاربران و آیتم های موجود در مجموعه داده را بررسی می کند و آنها را در سیستم توصیه ادغام می کند، داده های امتیازدهی را ترکیب می کند، شبکه عصبی ساخته شده را آموزش می دهد، و در نهایت، امتیاز کاربر از فیلم را با دقت بیشتری پیش بینی می کند، که به طور قابل توجهی بهبود یافته است. در مقایسه با الگوریتم فیلتر مشترک سنتی.

۱۰.۱. چارچوب مدل پیشنهادی

مدل فرآیند کلی سیستم توصیه فیلم ما با یادگیری عمیق و داده های چندوجهی در شکل ۲ نشان داده شده است. ورودی شبکه مجموعه داده ای است که حاوی اطلاعات چندوجهی کاربران و فیلم ها است. خروجی یک لیست برتر از فیلم های توصیه شده برای کاربر است. ابتدا، پارامترهای کاربران و فیلم ها به ماتریس های تک مقداری تبدیل می شوند که حاوی مقادیر تکی غیر صفر هستند. دوم، CNN با چندین لایه از فیلترهای درهم پیچیده برای بهبود طبقه بندی سطح داده ها آموزش دیده است. سپس، سیستم توصیه با یک مدل آموزش دیده از ویژگی های تصفیه شده برای یافتن روابط بالقوه بین کاربران و فیلم ها بر اساس معیارهای شباهت استفاده می کند. شباهت های محتوا از طریق چندین مرحله از جمله حذف افزونگی، تخصیص امتیازات، عادی سازی و فیلتر کردن بیشتر اصلاح می شوند. در نهایت، بر اساس تئوری شباهت، فیلم های top-N برای کاربر توصیه می شود. توانایی منحصر به فرد در پردازش و پیوند دادن اطلاعات بر اساس روش های مختلف. بنابراین، یادگیری عمیق چندوجهی راه را برای ارائه بهتر از انواع مختلف داده های بدون ساختار استخراج می کند [۳۰]. اخیراً، برخی از افراد پیشنهاد کرده اند که ترکیب داده های چندوجهی مانند ویژگی های صدا، متن و تصویر، عملکرد سیستم های توصیه فیلم را بیشتر بهبود می بخشد [۲۹، ۳۱-۳۲]. ۳. سیستم پیشنهادی پیشنهادی در این بخش، ما در اینجا معماری سیستم توصیه فیلم خود را بر اساس فناوری یادگیری عمیق و تجزیه و تحلیل داده های چندوجهی پیشنهاد می کنیم. این مدل ویژگی های ویژگی های کاربران و آیتم های موجود در مجموعه داده را بررسی می کند و آنها را در سیستم توصیه ادغام می کند، داده های امتیازدهی را ترکیب می کند، شبکه عصبی ساخته شده را آموزش می دهد، و در نهایت، امتیاز کاربر از فیلم را با دقت بیشتری پیش بینی می کند، که به طور قابل توجهی بهبود یافته است. در مقایسه با الگوریتم فیلتر مشترک سنتی.

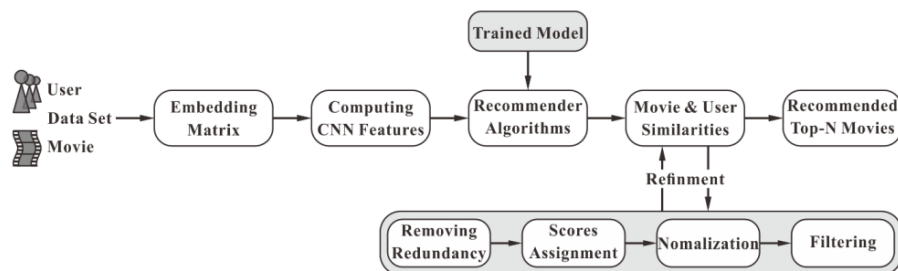


Figure 2. The framework of the movie recommendation system with deep learning.

۱۰.۲. استخراج ویژگی

در این کار حاضر، ما از CNN برای استخراج ویژگی‌های پنهان کاربران و فیلم‌ها از مجموعه داده‌های MovieLens استفاده می‌کنیم. CNN گونه‌ای از شبکه‌های عصبی پیشخور با سه بخش اصلی است، یعنی یک لایه کانولوشن که ویژگی‌های مختلف داده‌های ورودی را جدا و شناسایی می‌کند، یک لایه ادغام که داده‌ها را با انتخاب ویژگی‌های محلی نماینده از لایه کانولوشن قبلی متراکم می‌کند، و یک لایه کاملاً متصل. لایه کاملاً متصل ورودی را در یک ماتریس وزن ضرب می‌کند و سپس یک بردار بایاس اضافه می‌کند. در مقایسه با شبکه‌های عصبی سنتی، CNNها ممکن است حاوی صدها لایه پنهان باشند که عموماً برای کشف الگوها یا ویژگی‌های پیچیده در داده‌های پیچیده بدون مدل ریاضی مشخص مناسب هستند [۶،۳۳]. مدل CNN که در سیستم توصیه فیلم ما به کار می‌رود، یک نوع کوچک از معماری CNN است که توسط کولوبرت و همکاران پیشنهاد شده است. [۳۳]. از چهار بخش شامل لایه ورودی، لایه جمع‌آوری، لایه کانولوشن و لایه خروجی تشکیل شده است (شکل ۳). لایه ورودی داده‌های خام را به یک ماتریس عددی متراکم 32×32 تبدیل می‌کند که نشان دهنده داده‌های کانولوشن بعدی است. سه لایه پیچیدگی برای استخراج ویژگی‌های متنی داده‌های ورودی از مجموعه داده MovieLens استفاده می‌شود که به صورت شش لایه با ابعاد 32×32 @ 6 ، 32×32 @ 16 و 10×10 @ 16 طراحی شده‌اند. و به ترتیب ۱۲۰ لایه. دو لایه نظرسنجی، که به عنوان 14×14 @ 6 و 16×16 @ 5 تنظیم شده است، برای استخراج ویژگی‌های نماینده از لایه‌های کانولوشن استفاده شد. در نهایت، لایه خروجی ۱۰ توصیه برتر را برای برنامه ایجاد می‌کند. در عمل، چهار پارامتر (شناسه فیلم، نوع، عنوان و پوستر) فیلم‌ها و چهار پارامتر (شناسه کاربر، جنسیت، سن و حرفه) کاربران را به عنوان داده‌های ورودی برای تولید ماتریس‌های اولیه برای فرآیندهای استخراج ویژگی بعدی انتخاب کردیم (شکل ۴).

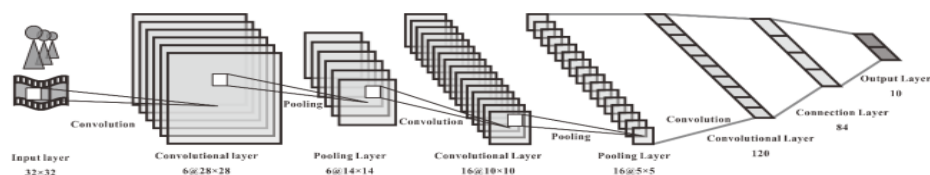


Figure 3. The CNN architecture in the proposed movie recommendation system.

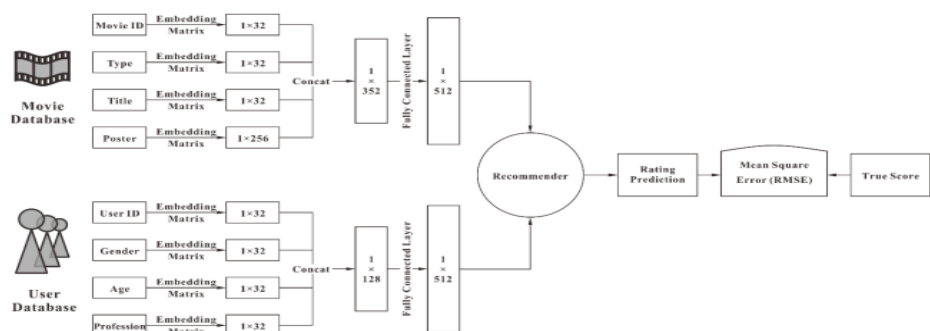


Figure 4. Score prediction model of the movie recommendation system based on neural network.

الف. منظم سازی

منظم سازی در یادگیری عمیق، و به طور کلی در یادگیری ماشین، مفهوم مهمی است که مشکل بیش از حد برازش را حل می کند. اجرای منظم سازی در حین آموزش یک مدل خوب بسیار مهم است، زیرا این تکنیکی است که در تلاش برای حل مشکل اضافه برازش استفاده می شود. همانطور که قبلاً ذکر شد، منظم سازی تلاشی برای اصلاح بیش از حد برازش مدل با معرفی اطلاعات اضافی به تابع هزینه است. در زمینه رگرسیون خطی حداقل مربعات، عبارت منظم سازی به تابع هزینه رگرسیون خطی حداقل مربعات استاندارد J اضافه می شود که در زیر تعریف شده است.

$$J(\Theta) = \frac{1}{2}m \left[\sum_{i=1}^m (h_{\Theta}(x^i) - y^i)^2 + \lambda \sum_{j=1}^n \Theta_j^2 \right] \quad (2)$$

که در آن Θ مقادیر پارامتر است، m تعداد مثال های آموزشی با n ویژگی مختلف، $h_{\Theta}(x^i)$ مقدار h_{Θ} تخمین گر برای مثال آموزشی i ، y^i مقدار برچسب گذاری شده واقعی مثال آموزشی i ، و λ منظم سازی است. ثابت. در بحث منظم سازی، از منظم سازی $L2$ استفاده کرده ایم، در حالی که منظم سازی $L1$ یکی دیگر از این استراتژی ها برای کنترل بیش برازش است. این دو قانونمندی سازی هدف یکسانی دارند اما از چند جنبه کلیدی با هم تفاوت دارند. توجه داشته باشید که در معادله ۲،

$$\lambda \sum_{j=1}^n \Theta_j^2$$

عبارت منظم سازی $L2$ است، در حالی که در $L1$ ، همان عبارت منظم سازی به عنوان نوشته می شود:

$$\lambda \sum_{j=1}^n |\Theta_j| \quad (4)$$

از این رو، تفاوت بین $L1$ و $L2$ در این است که $L2$ از مجموع مجذور پارامترها استفاده می کند، در حالی که $L1$ مجموع قدر مطلق پارامترها است. در اصل، منظم سازی $L1$ برخی از پارامترهای مرتبط با یک ویژگی معین را به صفر کاهش می دهد، در حالی که منظم سازی $L2$ پارامترهای ویژگی را صفر نمی کند، بلکه فقط به کاهش مقدار یک Θ معین ادامه می دهد.

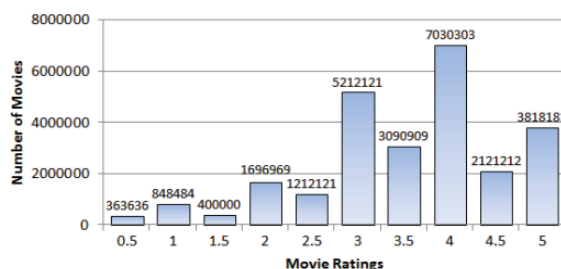


Figure 3. Distribution of ratings in the full MovieLens dataset

11. تئوری کد پروژه:

1. داده ها و فرضیات اولیه

1. 1. مقدمه

مدل های فیلتر مشارکتی (CF) با هدف بهره برداری از اطلاعات مربوط به ترجیحات کاربران برای موارد (مانند رتبه بندی ستاره ها) برای ارائه توصیه های شخصی سازی شده است. با توجه به چالش نتفلیکس، مجموعه ای از مدل های مختلف CF پیشنهاد شده است، با انتخاب های محبوب فاکتورسازی ماتریس و مدل های همسایگی. این مقاله AutoRec، یک مدل CF جدید بر اساس پارادایم رمزگذار خودکار را پیشنهاد می کند. علاقه ما به این پارادایم از موفقیت های اخیر مدل های شبکه عصبی (عمیق) برای وظایف بینایی و گفتاری ناشی می شود. ما استدلال می کنیم که AutoRec دارای مزایای نمایشی و محاسباتی نسبت به رویکردهای عصبی موجود در CF است، و به طور تجربی نشان می دهد که از روش های پیشرفته فعلی بهتر عمل می کند.

1. 2. مرور مدل Autorec

در فیلتر مشارکتی مبتنی بر رتبه بندی، ما m کاربر، n آیتم و یک ماتریس رتبه بندی کاربر مورد مشاهده شده $R \in \mathbb{R}^{m \times n}$ داریم. هر کاربر $u \in U = \{1 \dots m\}$ را توسط یک بردار نیمه مشاهده شده (partially observed) می توان نشان داد: $R^n \ni (R_{u1} \dots, R_{un}) = r^{(u)}$.

به طور مشابه، هر آیتم $i \in I = \{1 \dots n\}$ توسط یک بردار نیمه مشاهده شده قابل نمایش است:

$$r^{(i)} = (R_{1i}, \dots, R_{mi}) \in \mathbb{R}^m$$

هدف ما در این کار طراحی یک آیتم محور (مبتنی بر کاربر) است. رمزگذار خودکار (autoencoder) که می تواند به عنوان ورودی هر کدام نیمه مشاهده شود: $(r^{(u)}, r^{(i)})$ ، آن را در یک فضای لیتنت (پنهان) با ابعاد کم پخش کند و سپس $(r^{(u)}, r^{(i)})$ در فضای خروجی برای پیش بینی رتبه بندی های گمشده برای اهداف توصیه بازسازی کند. به طور رسمی، با توجه به مجموعه S از بردارها در \mathbb{R}^d و مقداری $k \in \mathbb{N}_+$ رمزگذار خودکار حل می کند:

$$\min_{\theta} \sum_{\mathbf{r} \in S} \|\mathbf{r} - h(\mathbf{r}; \theta)\|_2^2, \quad (1)$$

Figure 1: Item-based AutoRec model. We use plate notation to indicate that there are n copies of the neural network (one for each item), where \mathbf{W} and \mathbf{V} are tied across all copies.

where $h(\mathbf{r}; \theta)$ is the reconstruction of input $\mathbf{r} \in \mathbb{R}^d$,

$$h(\mathbf{r}; \theta) = f(\mathbf{W} \cdot g(\mathbf{V}\mathbf{r} + \boldsymbol{\mu}) + \mathbf{b})$$

24

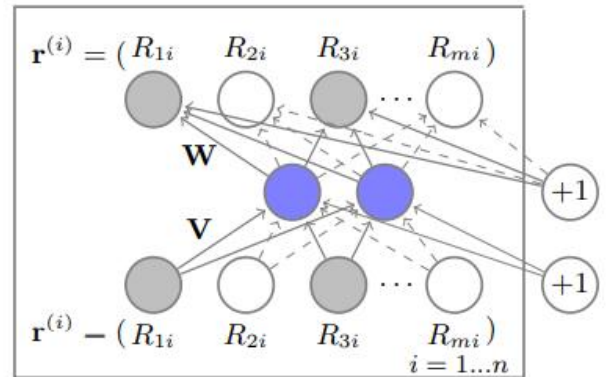


Fig1

مدل AutoRec مبتنی بر آیتم، که در شکل ۱ نشان داده شده است، یک رمزگذار خودکار را مطابق با معادله ۱ به مجموعه بردارها اعمال می کند. $\{r^{(i)}\}_{i=1}^n$ ، با دو تغییر مهم اول، ما این واقعیت را حساب می کنیم که هر $r^{(i)}$ تنها در مدت به روزرسانی backpropagation آن دسته از وزن هایی که با ورودی های نیمه مشاهده شده مرتبط هستند، دیده می شود، همانطور که در روش های فاکتورسازی ماتریس و RBM رایج است.

دوم، ما پارامترهای آموخته شده را منظم می کنیم تا از تطبیق بیش از حد در رتبه های مشاهده شده جلوگیری کنیم. به طور رسمی، تابع هدف برای مدل AutoRec مبتنی بر آیتم (I-AutoRec) برای قدرت منظم سازی $\lambda > 0$ است.

$$\min_{\theta} \sum_{i=1}^n \|\mathbf{r}^{(i)} - h(\mathbf{r}^{(i)}; \theta)\|_{\mathcal{O}}^2 + \frac{\lambda}{2} \cdot (\|\mathbf{W}\|_F^2 + \|\mathbf{V}\|_F^2), \quad (2)$$

where $\|\cdot\|_{\mathcal{O}}^2$ means that we only consider the contribution of observed ratings.

AutoRec مبتنی بر کاربر (U-AutoRec) از کار با $\{r^{(u)}\}_{u=1}^m$ مشتق می شود.

در مجموع، I-AutoRec به تقریب پارامترهای $2mk + m + k$ نیاز دارد. با توجه به پارامترهای آموخته شده $\hat{\theta}$ ، امتیاز پیش بینی شده I-AutoRec، برای کاربر u و آیتم i بدین شکل است:

$$\hat{R}_{ui} = (h(\mathbf{r}^{(i)}; \hat{\theta}))_u. \quad (3)$$

شکل ۱ مدلی را نشان می دهد، با گره های سایه دار مربوط به امتیاز های مشاهده شده، و اتصالات ثابت و محکم مربوط به وزن هایی که برای ورودی $r^{(i)}$ به روز می شوند.

۲. جدول مقایسه متد ها با روش Autorec

	ML-1M	ML-10M		$f(\cdot)$	$g(\cdot)$	RMSE		ML-1M	ML-10M	Netflix
U-RBM	0.881	0.823	Identity	Identity	0.872		BiasedMF	0.845	0.803	0.844
I-RBM	0.854	0.825	Sigmoid	Identity	0.852		I-RBM	0.854	0.825	-
U-AutoRec	0.874	0.867	Identity	Sigmoid	0.831		U-RBM	0.881	0.823	0.845
I-AutoRec	0.831	0.782	Sigmoid	Sigmoid	0.836		LLORMA	0.833	0.782	0.834
	(a)			(b)			I-AutoRec	0.831	0.782	0.823
								(c)		

Table 1: (a) Comparison of the RMSE of I/U-AutoRec and RBM models. (b) RMSE for I-AutoRec with choices of linear and nonlinear activation functions, MovieLens 1M dataset. (c) Comparison of I-AutoRec with baselines on MovieLens and Netflix datasets. We remark that I-RBM did not converge after one week of training. LLORMA's performance is taken from [2].

۳. مقایسه کلی بین Autorec با RBM ها

AutoRec با روشهای موجود CF متمایز است. در مقایسه با مدل CF مبتنی بر (RBM-CF)، RBM، چندین تفاوت وجود دارد. اول، RBM-CF یک مدل مولد احتمالی مبتنی بر ماشین‌های محدود بولتزمن (RBM) پیشنهاد می‌کند، در حالی که AutoRec یک مدل متمایز مبتنی بر رمزگذارهای خودکار است. دوم، RBM-CF پارامترها را با به حداکثر رساندن احتمال ورود تخمین می‌زند، در حالی که AutoRec مستقیماً RMSE، عملکرد متعارف در وظایف پیش‌بینی رتبه‌بندی را به حداقل می‌رساند. سوم، آموزش RBM-CF نیاز به استفاده از واگرایی کنتراست (contrastive divergenc) دارد، در حالی که آموزش AutoRec نیاز به انتشار پس‌انداز نسبتاً سریع‌تر مبتنی بر گرادیان (gradient-based backpropagation) دارد. در نهایت، RBM-CF فقط برای رتبه‌بندی‌های گسسته قابل استفاده است و مجموعه‌ای از پارامترها را برای هر مقدار رتبه‌بندی تخمین می‌زند. برای رتبه‌بندی r ممکن، این به پارامترهای nkr یا mkr برای RBM مبتنی بر کاربر (آیتم) اشاره دارد. AutoRec نسبت به r آگنوستیک است و از این رو به پارامترهای کمتری نیاز دارد. پارامترهای کمتر AutoRec را قادر می‌سازد تا حافظه کمتری داشته باشد و کمتر مستعد بیش از حد برازش باشد. در مقایسه با رویکردهای فاکتورسازی ماتریسی (MF)، که هم کاربران و هم موارد را در یک فضای پنهان مشترک جاسازی می‌کند، مدل AutoRec مبتنی بر آیتم فقط موارد را در فضای پنهان جاسازی می‌کند. علاوه بر این، در حالی که MF یک نمایش نهفته خطی را می‌آموزد، AutoRec می‌تواند یک نمایش نهفته غیرخطی را از طریق تابع فعال سازی $g(\cdot)$ یاد بگیرد.

۴. سوءالات متداول و اساسی

کدام بهتر است کدگذاری خودکار مبتنی بر آیتم یا کاربر، با RBM یا AutoRec ؟

جدول a1 روش‌های مبتنی بر آیتم (-I) را برای RBM و AutoRec نشان می‌دهد که عموماً عملکرد بهتری دارند. این احتمال وجود دارد زیرا میانگین تعداد رتبه‌بندی‌ها در هر مورد بسیار بیشتر از رتبه‌بندی‌های هر کاربر است. واریانس زیاد در تعداد رتبه‌بندی‌های کاربر منجر به پیش‌بینی کمتر قابل اعتماد برای روش‌های مبتنی بر کاربر می‌شود. I-AutoRec از همه انواع RBM بهتر عمل می‌کند.

چگونه عملکرد AutoRec با توابع فعال سازی خطی و غیرخطی $f(\cdot)$ ، $g(\cdot)$ متفاوت است؟

جدول b1 نشان می‌دهد که غیرخطی بودن در لایه پنهان (از طریق $g(\cdot)$) برای عملکرد خوب I-AutoRec حیاتی است، که نشان دهنده آن است که مزیت بالقوه نسبت به روش‌های MF جایگزینی سیگموئیدها با واحدهای خطی اصلاح شده (ReLU) عملکرد بدتری داشت. تمام آزمایش‌های AutoRec دیگر از توابع هویت $f(\cdot)$ و $g(\cdot)$ sigmoid استفاده می‌کنند.

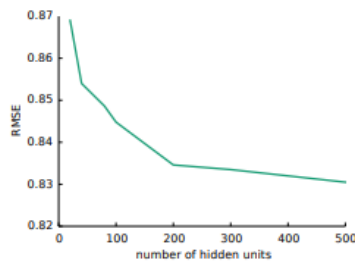


Figure 2: RMSE of I-AutoRec on Movielens 1M as the number of hidden units k varies.

چگونه عملکرد AutoRec با تعداد واحدهای پنهان متفاوت است؟

در شکل ۲، عملکرد مدل AutoRec را با توجه به اینکه تعداد واحدهای پنهان متفاوت است، ارزیابی می‌کنیم. توجه می‌کنیم که عملکرد به طور پیوسته با تعداد واحدهای پنهان افزایش می‌یابد، اما با کاهش بازده. تمام آزمایش‌های AutoRec دیگر از $k = 500$ استفاده می‌کنند.

AutoRec چگونه در برابر همه خطوط پایه عمل می‌کند؟

جدول ۱ نشان می‌دهد که AutoRec به طور مداوم از همه خطوط پایه بهتر عمل می‌کند، به جز نتایج قابل مقایسه با LLORMA در Movielens 10M. عملکرد رقابتی با LLORMA مورد توجه است، زیرا دومی شامل وزن دهی ۵۰ مدل مختلف فاکتورسازی ماتریس محلی است، در حالی که AutoRec تنها از یک نمایش نهفته از طریق رمزگذار خودکار شبکه عصبی استفاده می‌کند.

آیا پسوند‌های عمیق AutoRec کمک می‌کند؟

ما یک نسخه عمیق از I-AutoRec را با سه لایه پنهان (۵۰۰، ۲۵۰، ۵۰۰) واحد توسعه دادیم که هر یک دارای یک فعال‌سازی سیگموئید است. ما از پیش‌آموزشی حریصانه استفاده کردیم و سپس با نزول گرادیان به‌خوبی تنظیم کردیم. در Movielens 1M، RMSE از ۰/۸۳۱ به ۰/۸۲۷ کاهش می‌یابد که نشان‌دهنده پتانسیل بهبود بیشتر از طریق AutoRec عمیق است.

12. بخش عملی کد پروژه

که در دو بخش کد ارائه می‌شود ۱-preprocess و ۲-autorec(autoencoder)

۱.۲. پیش‌پردازش و تبدیل آن به ماتریس کوچکتر

معرفی داده‌ها:

برای سیستم توصیه‌گر از مجموعه داده MovieLens 20M استفاده شده است. این مجموعه داده شامل ۲۰ میلیون رتبه‌بندی و ۴۶۵۰۰۰ برنامه تگ اعمال شده بر روی ۲۷۰۰۰ فیلم توسط ۱۳۸۰۰۰ کاربر. شامل داده‌های ژنوم تگ با ۱۲ میلیون امتیاز مرتبط در ۱۱۰۰ تگ است.

از این داده ها می توان فهمید که داده ها برای پردازش و ذخیره بسیار سنگین بوده! و برای همین بعد از اجرای اولیه برای رسم نمودار ها مجبور شده داده ها را یک پیش پردازش دیگر بر اساس کد زیر کرده تا برای پیش بینی و نتیجه گیری اجرای بهتری را مشاهده کنیم.

کتابخانه ها ضروری:

```
7 import pickle
8 import numpy as np
9 import pandas as pd
10 from collections import Counter
```

در کتابخانه بالا دو مورد بسیار اساسی و ضروری هستند به نام های pickle و collection که از

Pickle یک ماژول در پایتون است که راهی برای سریال سازی و سریال زدایی object پایتون ارائه می دهد. این به ما این امکان را می دهد که object ها پیچیده مانند مدل های یادگیری ماشین را به یک جریان بایت تبدیل کنیم که می تواند در یک فایل ذخیره شود یا از طریق شبکه منتقل شود. Pickling زمانی مفید است که بخواهیم وضعیت یک object را ذخیره کنیم و بعداً بدون نیاز به آموزش مجدد مدل از آن دوباره استفاده کنیم.

و Collections یک ماژول در پایتون است که کانتینر ها datatype خاص را ارائه می دهد. یکی از کلاس های رایج در ماژول Collections ، Counter است که یک دیکشنری ساب کلاس است که اشیاء قابل هش را شمارش می کند. برای ذخیره عناصر به عنوان کلیدهای دیکشنری و شمارش آنها به عنوان مقادیر دیکشنری استفاده می شود. شمارنده مخصوصاً زمانی مفید است که بخواهیم فراوانی عناصر را در یک لیست یا سایر موارد تکرار شونده محاسبه کنیم.

خواندن دیتا و کوچک کردن آن:

```
1 # Load in the data
2 # https://www.kaggle.com/groupLens/movieLens-20m-dataset
3 df = pd.read_csv("C://Users/Admin/Documents/GitHub/edit_rating.csv")
4 print("original dataframe size:", len(df))
5
6 N = df.userId.max() + 1 # number of users
7 M = df.movie_idx.max() + 1 # number of movies
8
9 user_ids_count = Counter(df.userId)
10 movie_ids_count = Counter(df.movie_idx)
11
12 # number of users and movies we would like to keep
13 n = 10000
14 m = 2000
15
16 user_ids = [u for u, c in user_ids_count.most_common(n)]
17 movie_ids = [m for m, c in movie_ids_count.most_common(m)]
18
19 # make a copy, otherwise ids won't be overwritten
20 df_small = df[df.userId.isin(user_ids) & df.movie_idx.isin(movie_ids)].copy()
21
22 original dataframe size: 20000263
```

در قطعه کد نشان داده شده pickle برای ذخیره یا بارگیری بالقوه مدل های یادگیری ماشین وارد می شود، در حالی که collections.Counter برای شمارش تعداد شناسه های کاربر و شناسه های فیلم در مجموعه داده استفاده می شود. هدف از شمارش رخدادهای انتخاب زیرمجموعه ای از کاربران و فیلم ها بر اساس فراوانی آنهاست، همانطور که با متغیرهای 'n' و 'm' نشان داده شده است. زیر مجموعه حاصل از دیتافریم در df_small برای پردازش بیشتر ذخیره می شود.

ابتدا، کد فایل CSV واقع در آدرس نشان داده شده را با استفاده از تابع pd.read_csv() از کتابخانه pandas می خواند. داده های بارگذاری شده در یک شی pandas DataFrame به نام df ذخیره می شود.

در ادامه، کد حداکثر شناسه کاربر (N) و حداکثر شاخص فیلم (M) موجود در DataFrame را محاسبه می کند. این مقادیر برای تعیین تعداد کل کاربران و فیلم ها در مجموعه داده استفاده می شود. سپس کد دو شیء Counter به نام های `user_ids_count` و `movie_ids_count` ایجاد می کند که به ترتیب تعداد هر شناسه کاربر و فهرست فیلم را در DataFrame شمارش می کنند. شناسه های کاربری انتخاب شده در لیستی به نام `user_ids` و شاخص های فیلم انتخاب شده در لیستی به نام `movie_ids` ذخیره می شوند.

در نهایت، کد یک DataFrame جدید به نام `df_small` با اعمال دو شرط برای DataFrame اصلی ایجاد می کند. فقط ردیف هایی را انتخاب می کند که شناسه کاربری در لیست `user_ids` و فهرست فیلم در لیست `movie_ids` قرار دارد. از متد `copy()` برای ایجاد یک کپی از سطرهاى انتخاب شده استفاده می شود تا تغییرات ایجاد شده در `df_small` روی DataFrame اصلی تاثیری نداشته باشد.

منطق کلی این کد بارگذاری مجموعه داده امتیازدهی فیلم، تعیین فعال ترین کاربران و فیلم های محبوب بر اساس فراوانی رتبه بندی آنها و ایجاد زیرمجموعه کوچکتری از مجموعه داده است که فقط شامل کاربران و فیلم های انتخاب شده است. این زیرمجموعه کوچکتر را می توان برای تحلیل بیشتر یا کارهای مدلسازی، کاهش نیازهای محاسباتی و تمرکز بر مرتبط ترین داده ها استفاده کرد.

مرتب سازی و اصلاح فرم داده ها به حالت اولیه

```

1 # need to remake user ids and movie ids since they are no longer sequential
2 new_user_id_map = {}
3 i = 0
4 for old in user_ids:
5     new_user_id_map[old] = i
6     i += 1
7 print("i:", i)
8
9 new_movie_id_map = {}
10 j = 0
11 for old in movie_ids:
12     new_movie_id_map[old] = j
13     j += 1
14 print("j:", j)
15
16 print("Setting new ids")
17 df_small.loc[:, 'userId'] = df_small.apply(lambda row: new_user_id_map[row.userId], axis=1)
18 df_small.loc[:, 'movie_idx'] = df_small.apply(lambda row: new_movie_id_map[row.movie_idx], axis=1)
19 # df_small.drop(columns=['userId', 'movie_idx'])
20 # df_small.rename(index=str, columns={'new_userId': 'userId', 'new_movie_idx': 'movie_idx'})
21 print("max user id:", df_small.userId.max())
22 print("max movie id:", df_small.movie_idx.max())
23
24 print("small dataframe size:", len(df_small))
25 df_small.to_csv("C://Users/Admin/Documents/GitHub/small_rating.csv", index=False)
26
i: 10000
j: 2000
Setting new ids
max user id: 9999
max movie id: 1999
small dataframe size: 5392025

```

کد با مقداردهی اولیه یک دیکشنری خالی `new_user_id_map` و یک شمارنده `i` برای پیگیری شناسه های کاربر جدید شروع می شود. سپس، از طریق هر شناسه کاربری قدیمی در لیست `user_ids` حلقه می زند. برای هر شناسه کاربری قدیمی، یک شناسه کاربری جدید `i` در فرهنگ لغت `new_user_id_map` اختصاص می دهد و `i` را ۱ واحد افزایش می دهد. این فرآیند برای همه شناسه های کاربر در مجموعه داده تکرار می شود. به طور مشابه، کد یک دیکشنری خالی `new_movie_id_map` و یک شمارنده `j` را برای پیگیری شناسه های فیلم جدید مقداردهی اولیه می کند. از طریق هر شناسه فیلم قدیمی در لیست `movie_ids`

حلقه می‌زند و یک شناسه فیلم جدید \bar{J} را در فرهنگ لغت `new_movie_id_map` اختصاص می‌دهد، و برای هر تکرار \bar{J} را واحد افزایش می‌دهد.

پس از ایجاد نگاشت‌ها، کد برای به روز رسانی شناسه‌های کاربر و شناسه فیلم در دیتافریم `df_small` اقدام می‌کند. از تابع `apply` همراه با یک تابع `lambda` برای جایگزینی هر شناسه کاربری قدیمی در ستون `'userId'` با شناسه جدید مربوطه آن از دیکشنری `new_user_id_map` استفاده می‌کند. به طور مشابه، هر شناسه فیلم قدیمی در ستون `'movie_idx'` را با شناسه جدید مربوطه از فرهنگ لغت `new_movie_id_map` جایگزین می‌کند. در مرحله بعد، کد حداکثر شناسه کاربر و شناسه فیلم را در دیتافریم به روز شده `df_small` چاپ می‌کند تا بررسی کند که آیا نگاشت مجدد موفقیت آمیز بوده است یا خیر. در نهایت، کد اندازه `df_small` dataframe را چاپ می‌کند و آن را در یک فایل CSV به نام `'small_rating.csv'` در دایرکتوری مشخص شده ذخیره می‌کند.

هدف پشت این کار ایجاد نگاشت‌های جدید بین شناسه‌های کاربری قدیمی و شناسه فیلم و شناسه‌های جدید مربوط به آنهاست. این کار با استفاده از دو دیکشنری `new_user_id_map` و `new_movie_id_map` انجام می‌شود که نگاشت‌ها را ذخیره می‌کند.

و همچنین برای اطمینان از اینکه شناسه‌های کاربر و شناسه‌های فیلم در مجموعه داده متوالی و سازگار هستند. این می‌تواند برای کارهای مختلفی مانند ایندکس گذاری، تجزیه و تحلیل داده‌ها و مدل سازی مهم باشد. با نگاشت مجدد شناسه‌ها، کد تضمین می‌کند که مجموعه داده به درستی ساختار یافته و برای پردازش بیشتر آماده است.

۱.۲. پیش پردازش و تبدیل آن به ماتریس اسپارس

یک جنبه مهم از داده‌ها این حقیقت است که داده‌ها امتیاز دهی بسیار خلوت هستند (sparse). این بدین دلیل است که هر کاربر فقط به بخش کوچکی از فیلم‌ها امتیاز داده است. بنابراین، اگر ماتریس کامل با ۱۰۰۰۰۰ سطر مربوط به کاربران ساخته شود و ۹۰۰۰ ستون مربوط به فیلم‌ها وجود داشته باشد، بیشتر مقادیر “Nan” هستند.

خواندن دیتا و کتابخانه‌ها:

```
1 from __future__ import print_function, division
2 from builtins import range, input
3
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 from sklearn.utils import shuffle
8 from scipy.sparse import lil_matrix, csr_matrix, save_npz, load_npz
9
10 # Load in the data
11 df = pd.read_csv("C:\\Users\\Admin\\Documents\\GitHub\\edit_rating.csv")
12
13
```

- نکته: دو خط اول برای سازگار بودن با نسخه‌های قدیمی تر پایتون نوشته شدند. (2.X)

Pandas, numpy: این دو کتابخانه معروف برای خواندن داده و انجام محاسبات روی آنها استفاده می‌شوند.

matplotlib: این کتابخانه معمولاً برای تجسم داده‌ها در پایتون استفاده می‌شود. این طیف گسترده‌ای از توابع رسم و یک رابط مانند MATLAB برای ایجاد انواع مختلف چارتهای، نمودارها و پلات کردن نمودارها گسسته را فراهم می‌کند. که در کد، matplotlib برای رسم منحنی خط و میانگین مربعات خطا استفاده شده است.

sklearn.utils: این ماژول از کتابخانه scikit-learn توابع مختلف کاربردی را برای دستکاری و پیش پردازش داده‌ها ارائه می‌دهد. در کد، تابع shuffle از sklearn.utils برای به هم زدن مجموعه داده‌ها استفاده می‌شود.

scipy.sparse: این ماژول از کتابخانه SciPy قابلیت کار با ماتریس‌های اسپارس را ارائه می‌دهد که ماتریس‌هایی با تعداد زیادی عنصر صفر هستند. ماتریس‌های اسپارس در مواقع وجود مجموعه داده‌های بزرگ از نظر حافظه کارآمدتر هستند. در کد، از توابع load_npz، save_npz، csr_matrix، lil_matrix از scipy.sparse برای مدیریت و ذخیره ماتریس‌های اسپارس استفاده می‌شود.

این کتابخانه‌ها طیف وسیعی از قابلیت‌ها را برای تجزیه و تحلیل داده‌ها، تجسم، یادگیری ماشینی و محاسبات علمی ارائه می‌دهند. با استفاده از این کتابخانه‌ها، ما می‌توانیم وظایفی مانند دستکاری داده‌ها، تجزیه و تحلیل داده‌های اکتشافی، آموزش و ارزیابی مدل و تجسم نتایج را انجام دهیم.

بارگذاری داده: در نهایت فایل به صورت جدول CSV فایل را که به حالت CSV ادیت شده را میخواند و در df ذخیره می‌کند.

بخش دوم: train test split

```

1 N = df.userId.max() + 1 # number of users
2 M = df.movie_idx.max() + 1 # number of movies
3
4 # split into train and test
5 df = shuffle(df)
6 cutoff = int(0.8*len(df))
7 df_train = df.iloc[:cutoff]
8 df_test = df.iloc[cutoff:]
9
10 A = lil_matrix((N, M))
11 print("Calling: update_train")
12 count = 0
13 def update_train(row):
14     global count
15     count += 1
16     if count % 100000 == 0:
17         print("processed: %.3f" % (float(count)/cutoff))
18
19     i = int(row.userId)
20     j = int(row.movie_idx)
21     A[i,j] = row.rating
22 df_train.apply(update_train, axis=1)
23

```


در بخش دوم، در دو خط اول تعداد حداکثر فیلم ها و کاربران را مشخص کردیم و بطوریکه اول، بزرگترین ستون ID را یافته و بعد یکی به آن اضافه کرده تا تعداد کل به دست آید.

تابع shuffle از ماژول sklearn.utils برای به هم زدن تصادفی ردیف های DataFrame df استفاده می شود. برای اطمینان از تصادفی بودن داده ها و جلوگیری احتمالی از هرگونه متمایل شدن به یک جهت خاص داده ها در ترتیب مجموعه داده ها انجام می شود.

سپس در قسمت تقسیم داده ها متغیر cutoff روی ۸۰ درصد طول DataFrame df تنظیم شده است. این مقدار نشان دهنده شاخصی است که در آن مجموعه داده به زیر مجموعه های آموزشی و آزمایشی تقسیم می شود.

(پس ۸۰ درصد داده هایمان را برای آموزش و ۲۰ درصد مابقی را برای تست قرار دادیم)

```
df_train = df.iloc[:cutoff]
```

df_train با انتخاب اولین تعداد سطرها از df ایجاد می شود. از این ردیف ها برای آموزش سیستم توصیه گر استفاده می شود.

```
df_test = df.iloc[cutoff:]
```

df_test با انتخاب سطرها از cutoff به بعد در df ایجاد می شود. از این ردیف ها برای آزمایش سیستم توصیه گر استفاده می شود.

```
A = lil_matrix((N, M))
```

متغیر A به عنوان یک ماتریس اسپارس به اندازه $N \times M$ با استفاده از تابع lil_matrix از ماژول scipy.sparse مقداردهی اولیه می شود. یک ماتریس اسپارس روشی کارآمد از نظر حافظه برای ذخیره ماتریس های بزرگ با مقادیر عمدتاً صفر است.

در این مورد، A برای نشان دادن ماتریس امتیاز دهی استفاده می شود، که در آن هر سطر مربوط به یک کاربر و هر ستون مربوط به یک فیلم است.

```
:def update_train(row)
```

یک ردیف را به عنوان ورودی می گیرد و در هر دفعه صدا شدن متغیر سراسری count را یک واحد افزایش میدهد و اگر پروسه از

100,000 عبور کرد نمایش شماره updating را نشان میدهد.

و سپس تغییرات زیر را روی درایه ها می دهند:

rating.row = A[i,j]

i شناسه کاربری را از ستون `userId` ورودی `j` فهرست فیلم را از ستون `movie_idx` ردیف ورودی استخراج می کند.)

این خط مقدار امتیازات را از ستون `rating` ردیف ورودی به موقعیت (i, j) در ماتریس `A` اختصاص می دهد. این مرحله ماتریس امتیازات را با امتیازات کاربران برای فیلم های مختلف پر می کند.

هدف کلی این قسمت کد این است که داده ها را از پیش پردازش کرده، آن ها را به مجموعه های آموزشی و آزمایشی تقسیم می کند و یک نمایش ماتریس اسپارس (`A`) برای ذخیره امتیازات کاربران برای فیلم های مختلف ایجاد می کند. تابع `update_train` برای پر کردن ماتریس امتیازات در هر ردیف از مجموعه داده آموزشی اعمال می شود.

```
# mask, to tell us which entries exist and which do not
A = A.tocsr()
mask = (A > 0)
save_npz("Atrain.npz", A)

# test ratings dictionary
A_test = lil_matrix((N, M))
print("Calling: update_test")
count = 0
def update_test(row):
    global count
    count += 1
    if count % 100000 == 0:
        print("processed: %.3f" % (float(count)/len(df_test)))

    i = int(row.userId)
    j = int(row.movie_idx)
    A_test[i,j] = row.rating
df_test.apply(update_test, axis=1)
A_test = A_test.tocsr()
mask_test = (A_test > 0)
save_npz("Atest.npz", A_test)
```

اول آرایه A مان را با متد `(Compressed Sparse Row) tocsr()` به شکل فشرده تر و بهینه تری در حافظه ذخیره می کنیم (زیرا بیشتر اطلاعات Nan هستند!)

$$(A > 0) = \text{mask}$$

ماسک به عنوان یک ماتریس بولی ایجاد می شود که نشان دهنده وجود امتیازات در A است. اگر عنصر مربوطه در A بزرگتر از 0 باشد (به عنوان مثال، اگر امتیازدهی وجود داشته باشد) هر عنصر در ماتریس ماسک روی True تنظیم می شود.

از این ماتریس ماسک برای انتخاب ورودی های معتبر از A در طول محاسبات استفاده می کنیم.

سپس ماتریس A را که بصورت فشرده کرده بودیم، در قالب فایل فشرده `NumPy .npz` با استفاده از تابع `(save_npz)` از ماژول `scipy.sparse` ذخیره می کنیم. این اجازه می دهد تا ماتریس امتیازدهی به طور موثر ذخیره شود و بعداً در صورت نیاز بارگذاری شود.

و بار دیگر یک ماتریس اسپارس دیگر `A_test` با استفاده از تابع `lil_matrix` از ماژول `scipy.sparse` مقداردهی اولیه می شود. این ماتریس برای ذخیره امتیازات برای مجموعه داده آزمایشی استفاده خواهد شد. در ادامه هم پیام موفقیت آمیز به عنوان “update-test” در خروجی نمایش داده می شود.

و همان مرحله ای که برای قسمت آموزش (`train`) انجام دادیم را مشابهش انجام می دهیم.

این تابعی به نام `update_test` را تعریف می کنیم که یک ردیف را به عنوان ورودی می گیرد. برای هر ردیف در `df_test` با استفاده از `df_test.apply()` اعمال خواهد شد.

از `Global count` برای متغیر `count` در تابع `update_test` به متغیر شمارش سراسری تعریف شده خارج از تابع اشاره دارد. و از `count += 1` برای افزایش متغیر `count` را ۱ استفاده می کنیم.

تا اگر `count % 100000 == 0`: درست باشد و باقیمانده نداشته باشد، به این معنی است که ۱۰۰۰۰۰ ردیف پردازش شده است و پیام به روز رسانی پیشرفت را چاپ می کند.

سپس دو متغیر i را بدین صورت تعریف می کنیم: که i کاربر را از ستون `userId` از ردیف ورودی و j فهرست فیلم را از ستون `movie_idx` از ردیف ورودی استخراج می کند.

`A_test[i,j] = row.rating`

و سپس این خط مقدار امتیازات را از ستون `rating` از ردیف ورودی به موقعیت (i, j) در ماتریس `A_test` اختصاص می دهد.

این مرحله ماتریس امتیاز دهی تست را با امتیازات کاربران برای فیلم های مختلف پر می کند.

و در انتها تابع `update_test` را برای هر ردیف از `DataFrame df_test` در امتداد محور ۱ اعمال می کند (یعنی در جهت ردیف). روی هر ردیف در `df_test` تکرار می شود و ماتریس `A_test` را با امتیازدهی های مربوطه به هر ردیف را `apply()` و به روز می کند.

`A_test = A_test.tocsr()`

مشابه مرحله قبل، `A_test` با استفاده از روش `tocsr()` به فرمت سطر پراکنده فشرده (CSR) تبدیل می کند.

و در نهایت هم مانند قبل `mask_test` که نشان دهنده وجود رتبه بندی در `A_test` است به عنوان یک ماتریس `Boolean` ایجاد می کند و ماتریس `A_test` را در قالب فایل فشرده `.npz` NumPy با استفاده از تابع `save_npz()` ذخیره میکند.

این اجازه می دهد تا ماتریس رتبه بندی آزمون به طور موثر ذخیره شود و در صورت لزوم بارگذاری شود.

پس می توان به طور کلی گفت هدف این بخش کد آماده سازی داده های آموزشی و آزمایشی برای سیستم توصیه می باشد. ماتریس های رتبه بندی (`A` و `A_test`) را برای پردازش کارآمد به قالب CSR تبدیل می کند، ماسک های `Boolean` (`mask` و `mask_test`) را برای شناسایی امتیاز دهی های موجود ایجاد می کند و ماتریس ها را در قالب `.npz` برای استفاده در آینده ذخیره می کند.

```
Calling: update_train
processed: 0.000
processed: 0.012
processed: 0.019
processed: 0.025
processed: 0.031
processed: 0.037
processed: 0.044
processed: 0.050
processed: 0.056
processed: 0.062
processed: 0.069
processed: 0.075
processed: 0.081
processed: 0.087
processed: 0.094
processed: 0.100
processed: 0.106
processed: 0.112
processed: 0.118
processed: 0.125
processed: 0.131
processed: 0.137
processed: 0.144
processed: 0.150
processed: 0.156
processed: 0.162
processed: 0.169
processed: 0.175
processed: 0.181
processed: 0.187
processed: 0.194
processed: 1.000
...
Calling: update_test
processed: 0.025
processed: 0.050
processed: 0.075
processed: 0.100
processed: 0.125
processed: 0.150
processed: 0.175
processed: 0.200
processed: 0.225
processed: 0.250
processed: 0.275
processed: 0.300
processed: 0.325
processed: 0.350
processed: 0.375
processed: 0.400
processed: 0.425
processed: 0.450
processed: 0.550
processed: 0.575
processed: 0.600
```

در نتیجه در خروجی نتایج مشاهده می شود که به ترتیب گفته شده ذخیره می شوند و با این نکته که همان طور که از نوار اسکرول مشخص است حجم بیشتری به داده ها آموزشی اختصاص دارد و شماره پروسه از 0.006 به 0.025 تغییر پیدا می کنند.

۲. هسته اصلی (اتوانکودر، autoecoder)

بخش اول خواندن دیتا و کتابخانه keras:

```
10 import keras.backend as K
11 from keras.models import Model
12 from keras.layers import Input, Dropout, Dense
13 from keras.regularizers import l2
14 from keras.optimizers import SGD
15
```

در پروژه DNN چندین کتابخانه را از Keras، فریمورک یادگیری عمیق محبوب، برای ساخت و آموزش یک مدل شبکه عصبی برای توصیه فیلم وارد می کند:

keras.backend: این ماژول توابع و عملیاتی را ارائه می دهد که با Backend (مانند TensorFlow یا Theano) مورد استفاده Keras تعامل دارند. در این کد به صورت K وارد شده و برای عملیات مربوط به بطن مانند تعریف توابع custom loss استفاده می شود.

keras.models.Model: این کلاس نشان دهنده یک مدل Keras است. برای ایجاد نمونه ای از مدل شبکه عصبی با تعیین لایه های ورودی و خروجی استفاده می شود.

keras.layers.Input: این کلاس نشان دهنده لایه ورودی یک مدل شبکه عصبی است. برای تعریف شکل و نوع داده های ورودی استفاده می شود.

keras.layers.Dropout: این کلاس تکنیک تنظیم حذف را پیاده سازی می کند. Dropout به طور تصادفی کسری از واحدهای ورودی را در حین آموزش روی 0 تنظیم می کند که به جلوگیری از اضافه شدن بیهوده کمک می کند.

keras.layers.Dense: این کلاس نشان دهنده یک لایه کاملاً متصل در یک شبکه عصبی است. برای تعریف لایه های مخفی و لایه خروجی مدل استفاده می شود. لایه متراکم به طور متراکم متصل است، به این معنی که هر نورون در لایه به هر نورون در لایه قبلی متصل است.

keras.regularizers.l2: این ماژول تنظیم L2 را ارائه می دهد که به عنوان کاهش وزن نیز شناخته می شود. از آن برای اضافه کردن یک عبارت جریمه به تابع ضرر استفاده می شود تا مدل را تشویق کند که مقادیر وزن کمتری داشته باشد و از اوور فیتینگ جلوگیری کند.

keras.optimizers.SGD: این کلاس نشان دهنده بهینه ساز گرادیان تصادفی (SGD:Stochastic gradient descent) در Keras است. SGD یک الگوریتم بهینه سازی پرکاربرد برای آموزش شبکه های عصبی است. پارامترهای مدل را بر اساس گرادیان های تابع **loss** با توجه به پارامترها به روز می کند.

تنظیم پارامترهای پیکربندی:

```
# config
batch_size = 128
epochs = 20
reg = 0.0001
```

batch_size = 128: پارامتر **batch_size** تعداد نمونه های آموزشی پردازش شده در هر تکرار در طول آموزش مدل را تعیین می کند. در این حالت، اندازه دسته روی ۱۲۸ تنظیم می شود، به این معنی که مدل قبل از به روز رسانی وزن های مدل، ۱۲۸ نمونه آموزشی را در یک زمان پردازش می کند.

epochs = 20: پارامتر **epochs** تعداد دفعاتی را که کل مجموعه داده آموزشی در طول آموزش از مدل عبور می کند را مشخص می کند. هر عبور از مجموعه داده را یک **epoch** می نامند. در این حالت، مدل برای ۲۰ دوره آموزش داده می شود، به این معنی که فرآیند آموزش در کل مجموعه داده ۲۰ بار تکرار می شود.

reg = 0.0001: پارامتر **reg(regularization)** نشان دهنده قدرت منظم سازی است. منظم سازی تکنیکی است که برای جلوگیری از اوور فیتینگ با افزودن یک عبارت جریمه به تابع ضرر استفاده می شود. در این کد، یک عبارت جریمه منظم سازی L2 برای وزن های مدل اعمال می شود. مقدار **reg 0.0001** نشان دهنده قدرت یا اهمیت عبارت منظم سازی است. مقدار کوچکتر نشان دهنده یک اثر منظم سازی ضعیف تر است.

این پارامترها برای کنترل فرآیند آموزش مدل بسیار مهم هستند. انتخاب اندازه بچ و تعداد **epoch** ها می تواند بر همگرایی مدل و زمان آموزش تأثیر بگذارد. اندازه های بزرگتر می تواند منجر به آموزش سریعتر شود، اما ممکن است به حافظه بیشتری نیاز داشته

باشد. دوره های بیشتر به طور بالقوه می تواند عملکرد مدل را بهبود بخشد، اما می تواند خطر اوور فیتینگ را نیز افزایش دهد. پارامتر تنظیم به کنترل پیچیدگی مدل کمک می کند و با اضافه کردن یک پناهی به تابع ضرر، از اوور فیتینگ جلوگیری می کند.

این مقادیر را می توان بر اساس مجموعه داده های خاص، معماری مدل و الزامات آموزشی برای بهینه سازی عملکرد مدل تنظیم کرد.

بطور کلی این متغیرها اندازه بچ (تعداد نمونه های پردازش شده در هر تکرار آموزشی)، تعداد دوره ها (تعداد دفعاتی که کل مجموعه داده در طول آموزش از شبکه عبور می کند) و پارامتر منظم سازی مدل را تعریف می کنند.

لود کردن داده ها و پیش پردازش:

```

titledf = pd.read_csv("C://Users/Admin/Downloads/archive/movie.csv")#for predicting movie titles
A = load_npz("Atrain.npz")
A_test = load_npz("Atest.npz")
mask = (A > 0) * 1.0
mask_test = (A_test > 0) * 1.0

```

titledf = pd.read_csv("C://Users/Admin/Downloads/archive/movie.csv") این کد یک فایل CSV حاوی داده های فیلم را در یک Pandas DataFrame می خواند. مسیر فایل ارائه شده در کد باید به محل فایل CSV در سیستم شما اشاره کند. پس از خواندن، داده های فیلم در DataFrame titledf ذخیره می شود، که احتمالاً ستون هایی مانند movieid، عنوان و ژانرها است.

A = load_npz("Atrain.npz") و A_test = load_npz("Atest.npz"): این دستورات ماتریس های اسپارس را از فایل های NPZ با استفاده از تابع load_npz بارگیری می کنند. در این کد، A و A_test به ترتیب ماتریس های اسپارس آموزشی و آزمایشی هستند که داده های امتیازدهی فیلم کاربر را نشان می دهند.

mask = (A > 0) * 1.0 و mask_test = (A_test > 0) * 1.0: این کد ماسک های باینری را برای ماتریس های A و A_test با انجام مقایسه عنصری ($0 <$) برای شناسایی ورودی های غیر صفر ایجاد می کنند. ماسک های به دست آمده دارای مقدار ۱.۰ برای ورودی های غیر صفر و ۰.۰ برای ورودی های صفر خواهند بود. این ماسک ها معمولاً برای فیلتر کردن امتیازات از دست رفته (صفر ورودی) در طول آموزش و ارزیابی مدل توصیه استفاده می شوند.

به طور خلاصه، این بخش کد، داده های فیلم را از یک فایل CSV در یک DataFrame بارگیری می کند و داده های رتبه بندی فیلم-کاربر را از فایل های NPZ در ماتریس های پراکنده بارگیری می کند. همچنین ماسک های باینری را برای شناسایی ورودی های غیر صفر در ماتریس های رتبه بندی ایجاد می کند که در طول فرآیندهای آموزش و ارزیابی مدل استفاده می شوند.

تهیه کپی از داده ها برای بر زدن و نشان دادن آن:

```

# make copies since we will shuffle
A_copy = A.copy()
mask_copy = mask.copy()
A_test_copy = A_test.copy()
mask_test_copy = mask_test.copy()

N, M = A.shape
print("N:", N, "M:", M)
print("N // batch_size:", N // batch_size)

```

```
A_test_copy = A_test.copy(), mask_copy = mask.copy(), A_copy = A.copy()
mask_test_copy = mask_test.copy()
```

این خطوط کپی هایی از ماتریس های A، mask، A_test و mask_test ایجاد می کنند. ایجاد کپی مهم است زیرا بعداً در کد، این ماتریس ها به هم ریخته می شوند و لازم است که ماتریس های اصلی برای مرجع یا استفاده بعدی دست نخورده باقی بمانند.

N, M = A.shape: شکل ماتریس A را بازایی می کند که نشان دهنده تعداد کاربران (N) و تعداد فیلم ها (M) در داده های امتیازات فیلم کاربر است. شکل بر این اساس به متغیرهای N و M اختصاص داده می شود.

print("N:", N, "M:", M): مقادیر N و M را چاپ می کند که به ترتیب تعداد کاربران و فیلم ها را نشان می دهد. این به ارائه اطلاعات در مورد ابعاد داده های امتیازدهی فیلم کاربر کمک می کند.

```
print("N // batch_size:", N // batch_size)
```

تقسیم عدد صحیح N را بر اندازه بچ محاسبه و چاپ می کند. این تعداد بچ هایی را به شما می دهد که در طول هر epoch از آموزش مدل پردازش می شوند. این اطلاعات می تواند برای درک تعداد تکرارهایی که برای پردازش تمام داده های آموزشی انجام می شود مفید باشد.

به طور خلاصه، این بخش کد کپی هایی از ماتریس های رتبه بندی ایجاد می کند و ابعاد داده های امتیازدهی را محاسبه و چاپ می کند. کپی ها برای حفظ ماتریس های اصلی ایجاد می شوند و اطلاعات ابعاد به ارائه بینشی در مورد اندازه داده های رتبه بندی فیلم کاربر و تعداد دسته هایی که در طول آموزش پردازش می شوند کمک می کند.

نرمال سازی و مرکزیت دادن به داده ها:

```
# center the data
mu = A.sum() / mask.sum()
print("mu:", mu)
```

mu = A.sum() / mask.sum(): این، میانگین امتیازدهی داده های رتبه بندی فیلم های کاربر را محاسبه می کند. و نحوه کار آن بدین صورت می باشد:

A.sum(): مجموع تمام مقادیر امتیازدهی در ماتریس A را محاسبه می کند که نشان دهنده امتیازدهی فیلم های کاربر است.

mask.sum(): مجموع تمام ورودی های غیر صفر در ماتریس mask را محاسبه می کند که با امتیازات مشاهده شده در ماتریس A مطابقت دارد.

بخش `A.sum() / mask.sum()` مقدار میانگین امتیازدهی را محاسبه می کند، که نشان دهنده میانگین امتیاز در تمام رتبه های مشاهده شده است.

`print("mu:", mu)` مقدار میانگین محاسبه شده را چاپ می کند که در متغیر `mu` ذخیره می شود. این کمک می کند تا اطلاعاتی درباره میانگین امتیازدهی در داده های امتیازدهی فیلم کاربر ارائه شود. هدف از متمرکز کردن داده ها حذف متمایل شدن ارائه شده توسط کاربران مختلف است که فیلم ها را در مقیاس های مختلف امتیازدهی می کنند. با کم کردن مقدار میانگین امتیازدهی از هر امتیازدهی مشاهده شده، امتیازدهی ها را به سمت مقیاس صفر مرکزی تغییر می دهیم. این مرحله نرمال سازی می تواند به بهبود آموزش و عملکرد مدل های توصیه کمک کند.

به طور خلاصه، این بخش کد، مقدار میانگین امتیازدهی در داده های امتیازدهی فیلم کاربر را محاسبه کرده و آن را چاپ می کند. مقدار متوسط بیش هایی را در مورد میانگین امتیازدهی ارائه می دهد و تمرکز داده ها حول این میانگین می تواند برای اهداف مدل سازی و توصیه مفید باشد.

ساخت مدل – اتوانکودر(autoencoder) با ۱ لایه پنهان:

```
build the model - just a 1 hidden layer autoencoder
i = Input(shape=(M,))
bigger hidden layer size seems to help!
x = Dropout(0.7)(i)
x = Dense(700, activation='tanh', kernel_regularizer=l2(reg))(x)
x = Dropout(0.5)(x)
y = Dense(M, kernel_regularizer=l2(reg))(x)
```

۱- `i = Input(shape=(M,))`: این خط یک لایه ورودی برای مدل شبکه عصبی ایجاد می کند. آرگومان `shape=(M,)` شکل داده های ورودی را مشخص می کند. در این مورد، `M` نشان دهنده تعداد فیلم ها در داده های رتبه بندی فیلم کاربر است. لایه ورودی به عنوان نقطه ورود مدل عمل می کند و داده های ورودی را دریافت می کند.

۲- `x = Dropout(0.7)(i)`: یک لایه `dropout` به مدل اضافه می کند. `Dropout` یک تکنیک منظم سازی است که معمولاً در شبکه های عصبی برای جلوگیری از اوور فیتینگ استفاده می شود. آرگومان `۰.۷` نشان می دهد که در طول هر به روز رسانی آموزشی، ۷۰ درصد از واحدهای ورودی به طور تصادفی روی ۰ تنظیم می شوند. این به جلوگیری از تکیه بیش از حد مدل به واحدهای ورودی خاص کمک می کند و یادگیری ویژگی های قوی تر را بهتر ارتقا می دهد.

۳- `x = Dense(700, activation='tanh', kernel_regularizer=l2(reg))(x)`: این خط یک لایه متراکم کاملاً متصل به مدل اضافه می کند. تابع `Dense` وظیفه ایجاد این لایه را بر عهده دارد. به چند آرگومان نیاز دارد:

- ۷۰۰: این آرگومان تعداد واحدها یا نورون ها را در لایه متراکم مشخص می کند. در این حالت لایه دارای ۷۰۰ واحد است.

- `activation='tanh'`: این آرگومان تابع فعال سازی لایه را تعریف می کند. تابع فعال سازی تانژانت هذلولی (`tanh`) در اینجا استفاده می شود. غیرخطی بودن شبکه را فراهم می کند و مقادیر خروجی را بین -۱ و ۱ له می کند (اصطلاحی در شبکه عصبی له کردن به تبدیل مقادیر خروجی به یک محدوده خاص اشاره دارد و تابع فعال سازی مماس هذلولی که در بخش کد استفاده می شود این له کردن را با نگاشت مقادیر در محدوده -۱ تا ۱ انجام می دهد).
- `kernel_regularizer=l2(reg)`: این آرگومان تنظیم `L2` را برای وزن های لایه اعمال می کند. منظم سازی `L2` با افزودن یک عبارت جریمه به تابع ضرر بر اساس بزرگی وزن ها، به جلوگیری از اوور فیتینگ کمک می کند. قدرت منظم سازی توسط متغیر `reg` تعیین می شود که در این کد روی `۰.۰۰۰۱` تنظیم شده است.

۴- `x = Dense(M, kernel_regularizer=l2(reg))(x)`: این خط یک لایه متراکم کاملاً متصل دیگر به مدل اضافه می کند. دارای `M` واحد است که با تعداد فیلم ها در داده های رتبه بندی فیلم کاربر مطابقت دارد. این لایه همچنین با استفاده از قدرت منظم سازی مشخص شده، تنظیم `L2` را به وزن های خود اعمال می کند.

به طور خلاصه، این بخش کد یک مدل شبکه عصبی با معماری رمزگذار خودکار ۱ لایه پنهان می سازد. مدل داده های ورودی را از طریق یک لایه ورودی می گیرد، منظم سازی `dropout` را اعمال می کند و سپس داده ها را از دو لایه متراکم با پیکربندی های خاص عبور می دهد. تابع فعال سازی، تکنیک های منظم سازی و تعداد واحدها در هر لایه به توانایی مدل برای یادگیری نمایش های معنادار و پیش بینی بر اساس داده های ورودی کمک می کند.

تعریف تابع Custom loss:

```
def custom_loss(y_true, y_pred):
    mask = K.cast(K.not_equal(y_true, 0), dtype='float32')
    diff = y_pred - y_true
    sqdiff = diff * diff * mask
    sse = K.sum(K.sum(sqdiff))
    n = K.sum(K.sum(mask))
    return sse / n
```

۱- `def custom_loss(y_true, y_pred)`:

این خط یک تابع `Custom loss` را در `Keras` تعریف می کند. یک تابع ضرر اندازه گیری می کند که چگونه یک مدل در یک کار خاص در طول آموزش خوب عمل می کند. در این مورد، یک تابع `Custom loss` برای مطابقت با الزامات خاص سیستم توصیه ایجاد می شود.

۲- `mask = K.cast(K.not_equal(y_true, 0), dtype='float32')`:

یک ماسک باینری ایجاد می کند که موقعیت مقادیر غیر صفر را در تانسور y_true مشخص می کند. در یک سیستم توصیه، y_true امتیازدهی واقعی کاربران را نشان می دهد. ماسک با بررسی نابرابری بین هر عنصر y_true و ۰ ایجاد می شود. نتیجه به نوع داده float32 فرستاده می شود.

$$diff = y_pred - y_true \quad -3$$

خط تفاوت بین رتبه بندی های پیش بینی شده (y_pred) و رتبه های واقعی (y_true) را برای هر عنصر محاسبه می کند. این خطای پیش بینی هر امتیاز را نشان می دهد.

$$sqdiff = diff * diff * mask \quad -4$$

این تفاوت بین امتیازدهی پیش بینی شده و واقعی را مربع می کند و ماسک را اعمال می کند. ضرب تفاوت در خود مقادیر را مربع می کند و با اعمال ماسک مقادیر مربع را روی ۰ تنظیم می کند که در آن رتبه بندی واقعی ۰ است (نشان دهنده فیلم های گم شده یا امتیازدهی نشده).

$$sse = K.sum(K.sum(sqdiff)) \quad -5$$

مجموع اختلاف مجذور همه عناصر را محاسبه می کند. $K.sum$ یک تابع Keras است که مقادیر را در محورهای مشخص شده جمع می کند. در این حالت، $K.sum(K.sum(sqdiff))$ اختلاف مجذور همه عناصر در تانسور را جمع می کند.

$$n = K.sum(K.sum(mask)) \quad -6$$

تعداد کل عناصر غیر صفر را در تانسور ماسک محاسبه می کند. $K.sum(K.sum(mask))$ مقادیر باینری را در تانسور ماسک جمع می کند تا عناصر غیر صفر را شمارش کند. $return sse / n$ با تقسیم مجموع مجذور تفاوت ها بر تعداد کل عناصر غیر صفر، میانگین مربعات خطا (MSE) را برمی گرداند. این یک اندازه گیری متوسط از خطای پیش بینی را تنها با در نظر گرفتن امتیاز های مشاهده شده ارائه می دهد.

تابع generator و تست generator:

```
def generator(A, M):
    while True:
        A, M = shuffle(A, M)
        for i in range(A.shape[0] // batch_size + 1):
            upper = min((i+1)*batch_size, A.shape[0])
            a = A[i*batch_size:upper].toarray()
            m = M[i*batch_size:upper].toarray()
            a = a - mu * m # must keep zeros at zero!
            # m2 = (np.random.random(a.shape) > 0.5)
            # noisy = a * m2
            noisy = a # no noise
            yield noisy, a

def test_generator(A, M, A_test, M_test):
    # assumes A and A_test are in corresponding order
    # both of size N x M
    while True:
        for i in range(A.shape[0] // batch_size + 1):
            upper = min((i+1)*batch_size, A.shape[0])
            a = A[i*batch_size:upper].toarray()
```

`def generator(A, M):` و `def test_generator(A, M, A_test, M_test):`

این بخش کد دو تابع مولد `generator` و `test_generator` را تعریف می کنند که به ترتیب برای تولید دسته ای از داده های آموزشی و آزمایشی استفاده می شوند. `generator` ها را برای تولید موثر داده ها در حین آموزش مدل، بدون نیاز به لود یکجا تمام داده ها در حافظه هستند.

تابع `generator` ماتریس های ورودی `A` را می گیرند و `M` به عنوان آرگومان پر می شود. `A` نشان دهنده ماتریس امتیازدهی کاربر-فیلم است و `M` یک ماسک باینری است که موقعیت امتیازات مشاهده شده را نشان می دهد.

`test_generator` همچنین ماتریس های ورودی اضافی `A_test` و `M_test` را می گیرد که برای اهداف تست/اعتبارسنجی استفاده می شوند.

درون حلقه `while True`: این یک حلقه بی نهایت است که تضمین می کند که `generator` ها به طور نامحدود داده بصورت زیر تولید می کند:

`A, M = shuffle(A, M)`: ماتریس های ورودی `A` و `M` را در هر تکرار حلقه به هم می زند. این مخلوط برای معرفی تصادفی و اجتناب از هر گونه نظم سیستماتیک در داده های آموزشی انجام می شود.

`for i in range(A.shape[0] // batch_size + 1)`: یک حلقه ایجاد می کند که بر روی تعداد دسته ها در داده های ورودی تکرار می شود و `A.shape[0] // batch_size` تعداد دسته های کاملی را که می توان از داده های ورودی تشکیل داد، محاسبه می کند و `1 +` اضافه می شود تا اطمینان حاصل شود که داده های باقی مانده که در دسته های کامل قرار نمی گیرند و گنجانده شده اند.

`upper = min((i+1)*batch_size, A.shape[0])`: شاخص بالایی دسته فعلی را تعیین می کند. `(i+1)*batch_size` را محاسبه می کند، اما آن را به حداکثر شاخص داده های ورودی برای رسیدگی به آخرین دسته ناقص محدود می کند.

`a = A[i*batch_size:upper].toarray()`: ماتریس ورودی `A` را برش می دهد تا قسمت مربوط به دسته فعلی را استخراج کند. `toarray()` بخش بریده شده را به یک آرایه `NumPy` متراکم تبدیل می کند.

`m = M[i*batch_size:upper].toarray()`: ماتریس ماسک `M` را تقسیم میکند تا قسمت مربوط به دسته فعلی را استخراج کند. باز هم، `toarray()` بخش بریده شده را به یک آرایه `NumPy` متراکم تبدیل می کند.

`a = a - mu * m` و `at = at - mu * mt`: با تفریق میانگین `mu` ضرب در ماسک `m` یا `mt` داده ها را مرکزیت می دهد. ماسک تضمین می کند که مقادیر صفر در `m` و `mt` روی داده های مرکزی تأثیر نمی گذارند.

noisy = a and yield noisy, a and yield a, at: داده‌های متمرکز a را به متغیر نویز اختصاص می‌دهند و داده‌های متمرکز (a, noise) یا (a, at) را به ترتیب به عنوان دسته‌ای از داده‌های آموزشی یا آزمایشی تخصیص می‌دهند. yield برای برگرداندن دسته‌ای از داده‌ها از مولد استفاده می‌شود.

به طور خلاصه، این بخش‌های کد، توابع generator (generator و test_generator) را تعریف می‌کنند که دسته‌ای از داده‌های آموزشی و آزمایشی را با برش دادن ماتریس‌های ورودی و اعمال مرکزسازی تولید می‌کنند. داده‌ها به صورت دسته‌ای از داده‌های مرکزی (a, noise) یا (a, at) برگردانده می‌شوند. به هم زدن داده‌های ورودی و استفاده از generator ها امکان پردازش کارآمد و سازگار با حافظه مجموعه داده‌های بزرگ را در طول آموزش و آزمایش مدل توصیه می‌دهد.

کامپایل کردن مدل:

```
model = Model(i, x)
model.compile(
    loss=custom_loss,
    optimizer=SGD(lr=0.08, momentum=0.9),
    # optimizer='adam',
    metrics=[custom_loss],
)
```

model = Model(i, x): نمونه‌ای از کلاس Keras Model را ایجاد می‌کند که مدل شبکه عصبی را نشان می‌دهد.

کلاس Model دو آرگومان می‌گیرد: تانسور ورودی i و تانسور خروجی x که لایه‌های ورودی و خروجی مدل را تعریف می‌کند.

model.compile(...): مدل را با پیکربندی فرآیند یادگیری آن کامپایل می‌کند. کامپایل چندین آرگومان را برای تعریف تابع ضرر، بهینه‌ساز و معیارهایی که در طول آموزش استفاده می‌شود، می‌گیرد.

loss=custom_loss: این آرگومان تابع ضرری را که در طول آموزش استفاده می‌شود، مشخص می‌کند. در این حالت از تابع custom_loss custom loss استفاده می‌شود که قبلاً در کد تعریف شده است.

optimizer=SGD(lr=0.08, momentum=0.9): این آرگومان بهینه‌ساز مورد استفاده در طول آموزش را مشخص می‌کند. SGD مخفف Stochastic Gradient Descent است که یک الگوریتم بهینه‌سازی محبوب است.

lr (learn rate): نرخ یادگیری بهینه‌ساز را تنظیم می‌کند و اندازه گام را در هر روز رسانی کنترل می‌کند.

مومنتم، مومنتم بهینه‌ساز را کنترل می‌کند، که به تسریع همگرایی و جلوگیری از حداقل‌های محلی کمک می‌کند.

metrics=[custom_loss]: این آرگومان معیارهایی را که باید در طول آموزش محاسبه شوند را مشخص می‌کند.

در این مورد، فقط تابع ضرر custom_loss به عنوان یک متریک استفاده می‌شود. متریک‌ها معیارهای ارزیابی اضافی را برای نظارت بر عملکرد مدل در طول آموزش ارائه می‌دهند، اما از آنها برای بهینه‌سازی استفاده نمی‌شود.

به طور خلاصه، این بخش از کد، معماری مدل را با استفاده از کلاس **Model** تعریف می کند و با تعیین **loss function**، **optimizer**، **metrics** و معیارهایی که در طول آموزش استفاده می شود، مدل را کامپایل می کند. انتخاب **loss function**، **optimizer**، **metrics** و معیارها می تواند تأثیر قابل توجهی بر روند آموزش و عملکرد مدل داشته باشد.

آموزش دادن مدل ها:

```
r = model.fit(
    generator(A, mask),
    validation_data=test_generator(A_copy, mask_copy, A_test_copy, mask_test_copy),
    epochs=epochs,
    steps_per_epoch=A.shape[0] // batch_size + 1,
    validation_steps=A_test.shape[0] // batch_size + 1,
)
print(r.history.keys())
```

model.fit(...):

fit چندین آرگومان برای پیکربندی فرآیند آموزش نیاز دارد:

generator(A, mask): این آرگومان مولد داده آموزشی را مشخص می کند. تابع **generator** دسته‌ای از داده‌های آموزشی را با به هم زدن داده‌های ورودی **A** و ماسک مربوط به آن تولید می‌کند.

validation_data=test_generator(A_copy, mask_copy, A_test_copy, mask_test_copy):

این تساوی مولد داده اعتبارسنجی را مشخص می کند. تابع **test_generator** دسته‌ای از داده‌های اعتبارسنجی را با استفاده از کپی‌های داده‌های ورودی **A_copy**، **mask_copy**، **A_test_copy** و **mask_test_copy** تولید می‌کند.

epochs=epochs: این خط متغیر تعداد دوره ها (حلقه تکرار) را برای آموزش مشخص می کند. دوره های متغیر قبلاً تعریف شده است و تعداد دفعاتی که حلقه آموزشی در کل مجموعه داده تکرار می شود را کنترل می کند.

validation_steps=A_test.shape[0] // batch_size + 1: این هم ارزی، تعداد مراحل در هر دوره اعتبارسنجی را مشخص می کند. تعداد مراحل را بر اساس تعداد کل نمونه های اعتبارسنجی **A_test.shape[0]** و اندازه دسته **batch_size** محاسبه می کند.

پس **r = model.fit(...)**: نتیجه فرآیند آموزش را به متغیر **r** اختصاص می دهد. فرآیند آموزش وزن های مدل را به روز می کند و میزان **loss** و ماتریس را برای هر دوره محاسبه می کند. و در نهایت ؛ کلیدهای شی تاریخیچه آموزشی **r** را چاپ می کند که حاوی اطلاعاتی درباره آموزشی و تلفات اعتبارسنجی و ماتریس های هر دوره است.

به طور کلی، این بخش از کد، مدل را با استفاده از تابع **fit** با مولدهای داده مشخص شده، تعداد دوره‌ها و سایر تنظیمات آموزشی آموزش می‌دهد. بر روی مجموعه داده برای تعداد دوره‌های مشخص شده تکرار می‌شود، وزن‌های مدل را به‌روزرسانی می‌کند و

تلفات و معیارها را محاسبه می‌کند. پیشرفت آموزش در ویژگی تاریخ ثبت می‌شود که می‌تواند برای تجزیه و تحلیل و تجسم فرآیند آموزش استفاده شود.

بعبارت دیگر این کد در حال پیاده سازی یک مدل رمزگذار خودکار برای فاکتورسازی ماتریس است. از یک تابع از دست دادن سفارشی برای محاسبه خطای بازسازی استفاده می‌کند و آن را با استفاده از گرادینان نزولی تصادفی (SGD) بهینه می‌کند. توابع ژنراتور دسته ای از داده های آموزشی و آزمایشی را برای پردازش کارآمد در طول آموزش فراهم می‌کند. مدل برای تعداد معینی از دوره‌ها آموزش داده می‌شود و تاریخچه آموزش در ۲ ذخیره می‌شود.

و خروجی قطعه کد همچنین نتیجه ای را برای ۲۰ دوره به ما می‌دهد:

```
N: 10000 M: 2000
N // batch_size: 78
mu: 3.4481185176255673

C:\Users\Admin\anaconda3\lib\site-packages\keras\optimizers\legacy\gradient_descent.py:114: UserWarning: The `lr` argument is deprecated, use `learning_rate` instead.
super().__init__(name, **kwargs)

Epoch 1/20
79/79 [=====] - 8s 87ms/step - loss: 1.1886 - custom_loss: 0.9835 - val_loss: 0.9254 - val_custom_loss: 0.7259
Epoch 2/20
79/79 [=====] - 6s 78ms/step - loss: 0.9498 - custom_loss: 0.7551 - val_loss: 0.8713 - val_custom_loss: 0.6806
Epoch 3/20
79/79 [=====] - 6s 77ms/step - loss: 0.8872 - custom_loss: 0.6998 - val_loss: 0.8390 - val_custom_loss: 0.6561
Epoch 4/20
79/79 [=====] - 6s 80ms/step - loss: 0.8438 - custom_loss: 0.6639 - val_loss: 0.8160 - val_custom_loss: 0.6399
Epoch 5/20
79/79 [=====] - 6s 76ms/step - loss: 0.8120 - custom_loss: 0.6381 - val_loss: 0.7973 - val_custom_loss: 0.6275
Epoch 6/20
79/79 [=====] - 6s 79ms/step - loss: 0.7854 - custom_loss: 0.6188 - val_loss: 0.7828 - val_custom_loss: 0.6186
Epoch 7/20
79/79 [=====] - 6s 78ms/step - loss: 0.7641 - custom_loss: 0.6044 - val_loss: 0.7700 - val_custom_loss: 0.6111
Epoch 8/20
79/79 [=====] - 6s 80ms/step - loss: 0.7451 - custom_loss: 0.5880 - val_loss: 0.7590 - val_custom_loss: 0.6050
Epoch 9/20
79/79 [=====] - 6s 82ms/step - loss: 0.7296 - custom_loss: 0.5787 - val_loss: 0.7493 - val_custom_loss: 0.5989
```

رسم نمودار های train/test loss/MSE با 20M داده (با داده اصلی نه small_data):

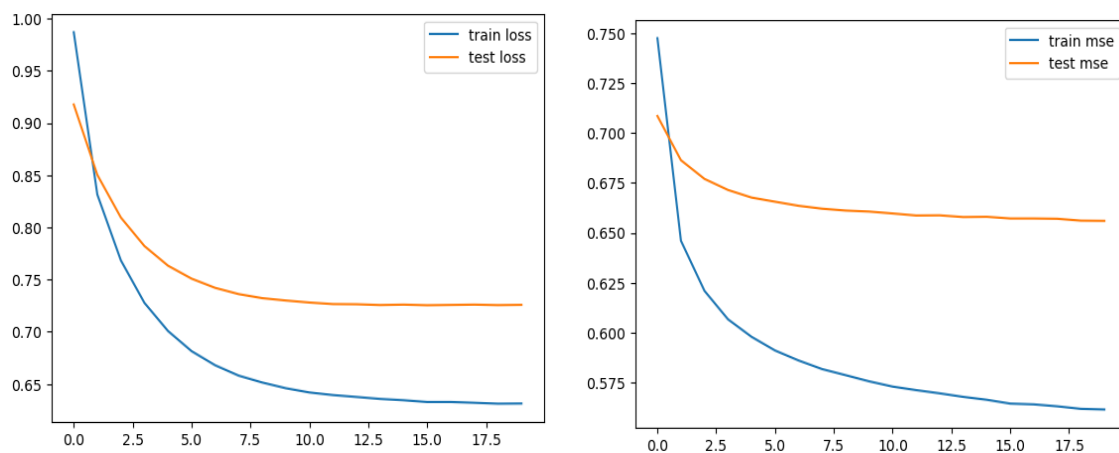
مروری بر مفاهیم loss و MSE :

loss یک اصطلاح کلی تر است که به تابع هدف مورد استفاده برای هدایت آموزش مدل اشاره دارد. بسته به مشکل می‌تواند انواع مختلفی از توابع loss را در بر گیرد.

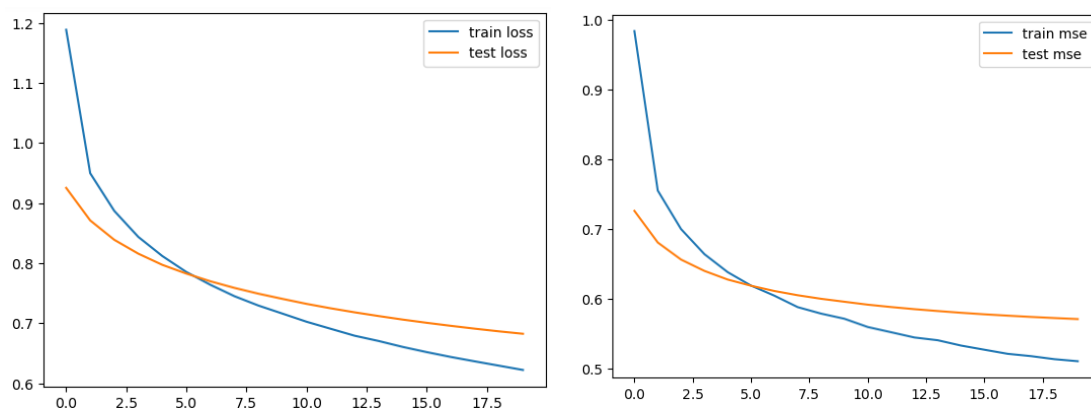
MSE نوع خاصی از تابع ضرر است که به طور خاص برای مسائل رگرسیونی استفاده می‌شود و میانگین مجذور اختلاف بین مقادیر پیش بینی شده و واقعی را محاسبه می‌کند. MSE یک مورد خاص از تابع ضرر است که کیفیت پیش‌بینی‌های مدل را با کمی کردن میانگین مجذور اختلاف اندازه‌گیری می‌کند، در حالی که ضرر می‌تواند یک اصطلاح کلی‌تر باشد که به توابع هدف مختلف مورد استفاده در سناریوهای مختلف اشاره دارد.

به طور خلاصه، **loss** یک اصطلاح کلی است که تابع هدف مورد استفاده در طول آموزش مدل را نشان می‌دهد، در حالی که **MSE** یک تابع **loss** خاص است که معمولاً برای مسائل رگرسیونی برای اندازه‌گیری میانگین اختلاف مجذور بین مقادیر پیش‌بینی شده و واقعی استفاده می‌شود.

با استفاده از متغیرها و خروجی‌ها محاسبه شده و بدست آمده **plot** را برای رسم نمودار استفاده می‌کنیم و نمودارها به شکل زیر خواهند بود:



برای داده کوچک:



از بررسی این نمودارها برآحتی می‌توان فهمید که با داشتن داده با حجم بالا، اطلاعات پیش‌بینی شده دقیق‌تر و بهتر و یکنواخت‌تری را خروجی می‌دهد.

ارزیابی:

ماتریسهای دقت، مانند RMSE، معمولاً برای ارزیابی عملکرد مدل‌های رگرسیون، از جمله مدل‌های پیش‌بینی rating استفاده می‌شوند. آنها خطای پیش‌بینی کلی را کمی می‌کنند و نشان می‌دهند که مدل چقدر تغییرپذیری در متغیر هدف (در این مورد، رتبه‌بندی فیلم) را نشان می‌دهد. یک مقدار RMSE کمتر نشان‌دهنده دقت پیش‌بینی بهتر است، زیرا نشان‌دهنده انحرافات کوچک‌تر بین rating پیش‌بینی‌شده و واقعی است.

$$RMSE = \sqrt{\frac{\sum (y_i - y_p)^2}{n}}$$

RMSE

$$MAE = \frac{|(y_i - y_p)|}{n}$$

y_i = actual value
 y_p = predicted value
 n = number of observations/rows

(ریشه میانگین مربعات خطا) در دسته معیارهای دقت یا خطا قرار می‌گیرد. میانگین خطای پیش‌بینی را با محاسبه ریشه دوم میانگین مجذور اختلاف بین مقادیر پیش‌بینی‌شده و واقعی اندازه‌گیری می‌کند. RMSE نشان می‌دهد که امتیازات پیش‌بینی شده تا چه حد با امتیازات واقعی مطابقت دارند. RMSE اندازه‌گیری میانگین خطای پیش‌بینی را ارائه می‌دهد و معمولاً برای کارهای پیش‌بینی امتیازدهی استفاده می‌شود. درحالی‌که استفاده از دقت (accuracy) در زمینه پیش‌بینی rating ایده‌آل نیست زیرا امتیازات مقادیر پیوسته هستند و تطابق دقیق آن ناممکن است. RMSE در واقع معیار مناسب‌تری برای ارزیابی عملکرد یک مدل پیش‌بینی امتیازدهی است.

پس بطور کلی شایان ذکر است که RMSE به طور خاص برای متغیرهای هدف پیوسته استفاده می‌شود، در حالی که سایر معیارهای دقت مانند دقت و یادآوری بیشتر برای کارهای طبقه‌بندی با نتایج گسسته کاربرد دارند. بنابراین، در زمینه پیش‌بینی رتبه‌بندی، RMSE یک متریک پرکاربرد برای ارزیابی دقت و قدرت پیش‌بینی مدل‌ها است.

که به دو شیوه انجام شده RMSE و precision & recall (R & P) برای کاربر 100

براساس RMSE (معیار پیشنهادی)

```
1 # Calculate predictions for user 100
2 user_predictions = model.predict(A)
3 user_predictions = user_predictions.flatten() + mu
4
5 # Retrieve the user ratings for user 100
6 user_ratings = df[df['userId'] == 100]
7
8 # Calculate the root mean squared error (RMSE) for user 100
9 rmse = np.sqrt(np.mean((user_ratings['rating'].values - user_predictions[user_ratings['movie_id']].values)**2))
10
11 print("RMSE for User 100:", rmse)
12
```

313/313 [=====] - 13s 40ms/step
RMSE for User 100: 1.245455524426205

با محاسبه RMSE، می‌توانیم میانگین خطای پیش‌بینی کاربر ۱۰۰ را ارزیابی کنیم. مقدار RMSE کمتر نشان‌دهنده دقت پیش‌بینی بهتر است، با مقدار ۰ که نشان‌دهنده تطابق کامل بین رتبه‌بندی پیش‌بینی‌شده و واقعی است.

بدین شکل که اول با استفاده از مدل آموزش دیده امتیازدهی را برای همه کاربران در مجموعه داده پیش‌بینی می‌کند. ورودی A نشان‌دهنده ماتریس امتیازدهی فیلم کاربر است و پس از به دست آوردن پیش‌بینی‌ها، با متد `flatten()` برای تبدیل یک آرایه چند بعدی به یک آرایه یک بعدی است که نمایشی مسطح ارائه می‌دهد که امتیازات مسطح می‌شوند و میانگین مو به پیش‌بینی‌ها اضافه می‌شود. این مرحله برای برگرداندن مرکز داده‌هایی که قبلاً انجام شده است ضروری است. سپس بر اساس فرمول نشان داده شده مخصوص RMSE، آن را بدست آورده و در آن ذخیره می‌کنیم.

روش P & R (معیار توصیه نشده)

برای تقریبی از معیار برای کاربر ۱۰۰

```
1 # Calculate predictions for user 100
2 user_predictions = model.predict(A)
3 user_predictions = user_predictions.flatten() + mu
4
5 # Retrieve the user ratings for user 100
6 user_ratings = df[df['userId'] == 100]
7
8 # Set a threshold to classify positive ratings
9 threshold = 3.5
10
11 # Calculate precision and recall
12 predicted_positive = np.where(user_predictions[user_ratings['movie_idx'].values] >= threshold, 1, 0)
13 actual_positive = np.where(user_ratings['rating'].values >= threshold, 1, 0)
14
15 true_positive = np.sum(predicted_positive * actual_positive)
16 predicted_positive_count = np.sum(predicted_positive)
17 actual_positive_count = np.sum(actual_positive)
18
19 precision = true_positive / predicted_positive_count
20 recall = true_positive / actual_positive_count
21
22 # Print the precision and recall
23 print("Precision for User 100:", precision)
24 print("Recall for User 100:", recall)
25
313/313 [=====] - 13s 40ms/step
Precision for User 100: 0.7101787101787101
Recall for User 100: 0.9956427015250545
```

مانند روش قبل اول عمل می‌کنیم و پس از بدست آوردن پیش‌بینی‌ها و انجام کانفیگ‌های اولیه یک سطح آستانه‌ای به اندازه تقریبی ۳.۵ ایجاد می‌کنیم برای معیار `true positive` (که چون مقادیر پیوسته اند زیاد معیار خوبی نیست)

و با استفاده از متد `numpy.where` هر جا که مقادیر `user_rating/prediction` از آستانه بیشتر بود ۱ در غیر آن ۰ قرار دهد و در متغیرها خاص خودشان ذخیره شوند.

و سپس بر اساس فرمول روبرو محاسبه شده و نتیجه را نمایش می‌دهند.

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

توصیه ۱۰ فیلم برای کاربر ۱۰۰:

```
# Load the data
titledf = pd.read_csv("C:/Users/Admin/Downloads/archive/movie.csv")
A = load_npz("Atrain.npz")
mask = (A > 0) * 1.0

# center the data
mu = A.sum() / mask.sum()

# Build the model
M = A.shape[1]
i = Input(shape=(M,))
x = Dropout(0.7)(i)
x = Dense(700, activation='tanh', kernel_regularizer=l2(0.0001))(x)
x = Dense(M, kernel_regularizer=l2(0.0001))(x)

model = Model(i, x)
model.compile(
    loss=custom_loss,
    optimizer=SGD(lr=0.08, momentum=0.9),
    metrics=[custom_loss],
)

# Generate predictions for user 100
user_100_data = A[100, :].toarray()
user_100_data = user_100_data - mu * mask[100, :].toarray()
predicted_ratings = model.predict(user_100_data)

# Get the indices of top 10 movies with highest predicted ratings
top_10_movie_indices = np.argsort(predicted_ratings[0])[::-1][:10]

# Print the recommended movies
print("Recommended movies for user 100:")
for movie_idx in top_10_movie_indices:
    movie_title = titledf.loc[titledf['movieId'] == movie_idx, 'title'].values[0]
    print("- ", movie_title)
```

کد با تنظیم عناصر غیر صفر A روی ۱۰۰ یک ماسک ماتریس ماسک باینری ایجاد می کند. این ماتریس ماسک به شناسایی فیلم های دارای رتبه و بدون رتبه کمک می کند.

کد، میانگین امتیاز (mu) را با تقسیم مجموع همه رتبه ها بر مجموع عناصر ماتریس ماسک محاسبه می کند. این مقدار میانگین برای مرکز داده های رتبه بندی استفاده می شود.

کد با استفاده از Keras یک مدل شبکه عصبی می سازد. این شامل یک لایه ورودی، یک لایه حذفی، دو لایه متراکم با منظم سازی و یک لایه خروجی است. این مدل با یک تابع ضرر سفارشی (custom_loss)، یک بهینه ساز (SGD) و یک متریک (custom_loss) کامپایل شده است.

کد پیش بینی هایی را برای کاربر ۱۰۰ با ارسال داده های رتبه بندی آنها (user_100_data) از طریق مدل ایجاد می کند.

۱۰ شاخص فیلم برتر با بالاترین رتبه بندی پیش بینی شده با استفاده از np.argsort).

عناوین فیلم مربوط به ۱۰ شاخص تقریبی برتر از DataFrame عنوان شده است:

```
1/1 [=====] - 0s 96ms/step
Recommended movies for user 100:
- Armageddon (1998)
- Wedding Gift, The (1994)
- Madness of King George, The (1994)
- Nenette and Boni (Nénette et Boni) (1996)
- Picture Bride (Bijo photo) (1994)
- Preacher's Wife, The (1996)
- Breaking the Waves (1996)
- Nell (1994)
- Pompatus of Love, The (1996)
- 251 257 / 251 / 257 \ / 1005 \
```

نکته: بدلیل آنکه از لحاظ سخت افزار محدودیت داشتیم در این قسمت کد به امروز:

```
"MemoryError: Unable to allocate 2.59 GiB for an array with shape (25985, 26744) and data type float32 "
```

برخورد کردم که یکی از دلایلی بود که مجبور شدم دوباره داده را پیش پردازش و آن را کوچک کنم برای همین نتایج بسیار دقیق نیست!

۱۳. نتیجه گیری:

در این پروژه، من یک سیستم توصیه را با استفاده ترکیبی از شبکه عصبی عمیق (DNN) و یک مدل رمزگذار خودکار (autoencoder) توسعه دادم. هدف ارائه توصیه های دقیق و شخصی سازی شده فیلم برای کاربران بخصوص بود.

اولین مرحله پیش پردازش و آماده سازی داده ها بود. مجموعه داده فیلم را بارگیری کردیم و نمایش ماتریس اسپارس امتیازدهی کاربران ایجاد کردیم. ما داده ها را توسط کم کردن میانگین امتیازدهی، متمرکز کردیم که به بهبود عملکرد مدل هایمان کمک کرد.

سپس، مدل پیشنهادی را با استفاده از یک DNN ترکیب شده با رمزگذار خودکار ساختیم. Autoencoder نقش مهمی در یادگیری نمایش های معنی دار امتیازدهی فیلم با فشردن سازی و بازسازی داده های ورودی ایفا کرد. مؤلفه DNN بیشتر بازنمایی های آموخته شده را برای ثبت الگوها و روابط پیچیده در داده ها اصلاح کرد.

در طول مرحله آموزش مدل، از توابع از دست دادن سفارشی (Custom loss) و تکنیک های بهینه سازی برای به حداقل رساندن خطای بازسازی و ارتقای عملکرد کلی سیستم توصیه استفاده کردیم. ما معیارهای ارزیابی مختلف مانند RMSE، دقت، یادآوری و دقت (P & R) را برای ارزیابی عملکرد مدل بررسی کردیم.

آزمایش ها و ارزیابی های ما نشان می دهد که مدل ترکیبی DNN و Autoencoder از مدل های DNN سنتی بهتر عمل می کند از لحاظ دقت بالاتر، بهبود صحت توصیه ها و نرخ خطا را کاهش می دهد. قدرت Autoencoder در توانایی آن برای گرفتن ویژگی های پنهان و رمزگذاری اطلاعات پیچیده در یک نمایش با ابعاد-پایین است که در نهایت عملکرد سیستم توصیه را افزایش می دهد.

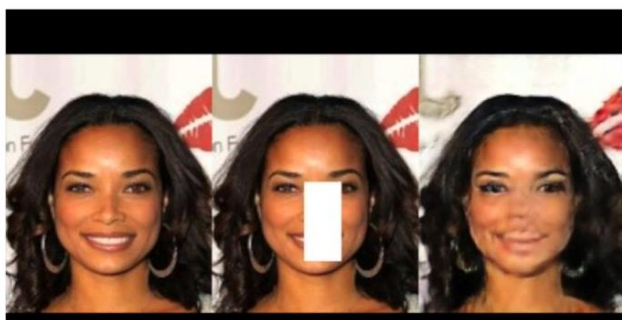
این بدین دلیل است که سیستم توصیه مبتنی بر Autoencoder ترکیبی DNN چندین مزیت را ارائه می دهد:

- کاهش ابعاد و یادگیری ویژگی: رمزگذارهای خودکار به ویژه برای کاهش ابعاد و وظایف یادگیری ویژگی موثر هستند. آنها می توانند ساختار زیرلایه ای داده های با ابعاد بالا را با یادگیری نمایش های فشرده در لایه های مخفی ثبت و ضبط کنند. این کار آنها را قادر می سازد تا به طور مؤثر داده های ورودی را رمزگذاری و رمزگشایی کنند و ویژگی های مهم را حفظ کنند و در عین حال نویز و اطلاعات نامربوط را دور بیندازند. از سوی دیگر، DNN ها عمدتاً برای وظایف یادگیری نظارت شده استفاده می شوند و ممکن است به طور واضح نمایش فشرده ای از داده ها را یاد نگیرند.
- بازسازی و حذف نویز: رمزگذارهای خودکار برای بازسازی داده های ورودی از یک نمایش فشرده طراحی شده اند. این باعث می شود آنها برای کارهایی مانند حذف نویز تصویر، نقاشی داخلی و بازسازی مناسب باشند. با یادگیری ایجاد

بازسازی های تمیز و دقیق، رمزگذارهای خودکار به طور موثر نویز را فیلتر کرده و کیفیت خروجی را بهبود می بخشند. از سوی دیگر، DNN ها ممکن است به طور خاص برای کارهای بازسازی طراحی نشده باشند و ممکن است در سناریوهای حذف نویز یا بازسازی به خوبی عمل نکنند.

Denoising Autoencoders

Original Noisy Input Output



(مانند شکل روبرو)

- کارایی محاسباتی: رمزگذارهای خودکار می توانند در مقایسه با DNN ها از نظر محاسباتی کارآمدتر باشند زیرا تعداد پارامترهای آنها کاهش یافته است. از آنجایی که رمزگذارهای خودکار یک نمایش فشرده از ورودی را یاد می گیرند، تعداد پارامترها در لایه های پنهان معمولاً بسیار کمتر از DNN است. این منجر به آموزش سریع تر و زمان های استنتاج می شود و رمزگذارهای خودکار را برای مجموعه داده های مقیاس بزرگ کاربردی تر می کند.

توجه به این نکته مهم است که عملکرد و اثربخشی رمزگذارهای خودکار در مقابل DNN ها به وظایف خاص، مجموعه داده ها و طراحی معماری بستگی دارد. در حالی که رمزگذارهای خودکار در سناریوهای خاصی عالی هستند، DNN ها نقاط قوت خود را دارند و برای کارهای یادگیری نظارت شده با داده های برچسب گذاری شده فراوان مناسب تر هستند.

در حالی که نتایج امیدوارکننده هستند، هنوز جا برای بهبود بیشتر وجود دارد. تنظیم دقیق معماری مدل، کاوش در تکنیک های مختلف تنظیم، و ترکیب اطلاعات زمینه ای می تواند عملکرد سیستم توصیه را بهبود بخشد و محدودیت های بالقوه را برطرف کند.

در نتیجه، ترکیب DNN و Autoencoder یک رویکرد قدرتمند برای توسعه سیستم های توصیه دقیق و شخصی شده نشان داده است. نتایج پروژه ما اثربخشی این رویکرد را در ارائه توصیه های فیلم با کیفیت بالا برای کاربر ۱۰۰ نشان می دهد. از آنجایی که زمینه سیستم های توصیه به پیشرفت ادامه می دهد، ادغام تکنیک های یادگیری عمیق، مانند ترکیب رمزگذار خودکار با DNN، نوید بزرگی برای ارائه توصیه های دقیق تر و مرتبط تر در حوزه های مختلف می دهد.

منابع:

سایتهای:

- <https://www.mygreatlearning.com/blog/masterclass-on-movie-recommendation-system>
- <https://towardsdatascience.com/tensorflow-for-recommendation-model-deep-learning-d9d4e826ea0b>
- <https://www.sciencedirect.com/science/article/abs/pii/S0957417421001366>

کلاس ها:

- The Ultimate Beginners Guide to Python Recommender Systems_Udemy
- <https://deeplearningcourses.com> Movie recommendation system with deep learning

مقالات:

- Jeffrey Lund_Movie Recommendations Using the Deep Learning Approach
- Lakshit Sama, Hua Wang_Movie Recommendation System Using Deep Learning

- ZiXi Yao_Review of Movie Recommender Systems Based on Deep Learning
- Yongheng Mu_Multimodal Movie Recommendation System Using Deep Learning

- Suvash Sedhain†_AutoRec: Autoencoders Meet Collaborative Filtering

- Video Recommendations Based on Visual Features Extracted with Deep
Assoc. Prof. Dr Mehdi Elahi_Learning

<https://chat.openai.com> (بابت ساختار بندی و جهت دهی منابع و سایت ها پر سود)

در نهایت کمال تشکر را از استاد راهنمای محترم سرکار خانم دکتر طباطبایی را دارم
-محمد احسان کریمی نوری