

# Data science- Recommender system for Amazon

Sadaf Fatollahy

۱۳ دی ۱۴۰۲



نام استاد: دکتر خردپیشه  
نام درس: مبانی علوم داده

# فهرست مطالب

۳	..... مقدمه	۱.۰
۳	..... درباره دیتاست	۲.۰
۳	..... EDA and feature engineering	۳.۰
۵	..... Data sparsity	۴.۰
۵	..... Colaborative filtering	۵.۰

## ۱.۰ مقدمه

یک سیستم توصیه توسعه یافته به کسب و کارها کمک می کند تا تجربه خریدار خود را در وب سایت بهبود بخشند و منجر به جذب و حفظ مشتری بهتر شود. Amazon یکی از بزرگترین شرکت های تجارت الکترونیک و رایانش ابری است. آنها ۸.۴ میلیون دلار در آگوست ۲۰۱۳ از دست دادند، زمانی که وب سایت آنها برای ۴۰ دقیقه از کار افتاد. آمازون به شدت به موتور توصیه ای متکی است که رتبه بندی مشتریان و تاریخچه خرید را برای توصیه اقلام و بهبود فروش بررسی می کند.

## ۲.۰ درباره دیتاست

این مجموعه داده مربوط به بیش از ۲ میلیون بررسی و رتبه بندی مشتری از محصولات مرتبط با زیبایی است که در وب سایت آنها فروخته شده است.

۱. UserId : شناسه کاربری منحصر به فرد مشتری.

۲. ProductId : کد شناسایی منحصر به فرد محصول آمازون برای هر محصول

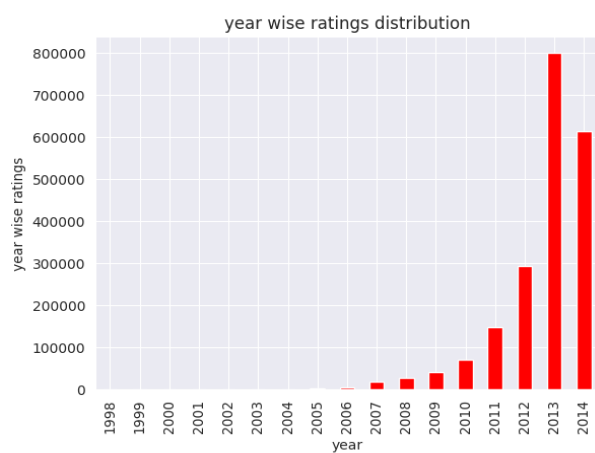
۳. Ratings : رتبه بندی (از ۱ تا ۵ بر اساس رضایت مشتری)

۴. Timestamp : مهر زمانی رتبه بندی (در زمان یونیکس)

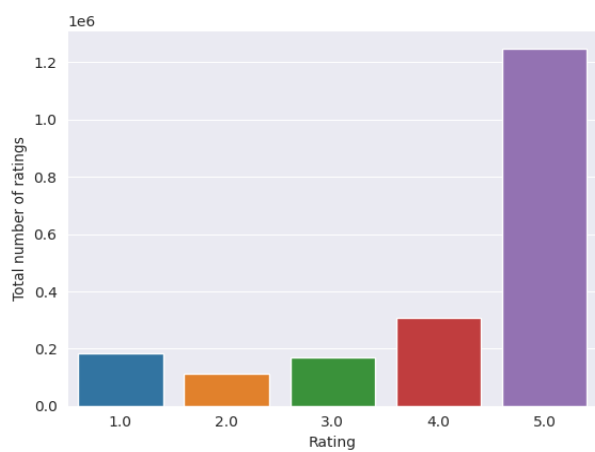
## ۳.۰ EDA and feature engineering

: در ابتدا در دیتاست کمی گشت و گذار کرده و به بررسی فیچر ها میپردازیم.

- این دیتاست شامل ۲۰۲۳۰۷۰ سطر و ۴ ستون است
- میانگین نمرات داده شده به محصولات ۴ بوده است که بیانگر کیفیت بالای محصولات است که آمازون میفروشد.
- دیتاست missing value ندارد.
- دیتاست داده تکراری ندارد.
- ستون Timestamp بیانگر تعداد ثانیه هایی است که از تاریخ ۱ ژانویه ۱۹۷۰ گذشته است. ما آن را به فرمت قابل فهم تاریخ و ساعت تبدیل میکنیم. تاریخ این دیتا ست از سال ۱۹۹۸ است تا ۲۰۱۴.
- میخواهیم بررسی کنیم میزان نمره دهی به محصولات در طی این سال ها به چه صورت بوده است: داده ها از سال ۱۹۹۸ تا ۲۰۱۴ جمع آوری شده اند. و ما می توانیم هر ساله ببینیم که رتبه بندی محصولات ارایشی به طور مداوم در حال افزایش است، به جز افزایش غیرمعمول در سال ۲۰۱۳ به دلایل ناشناخته.

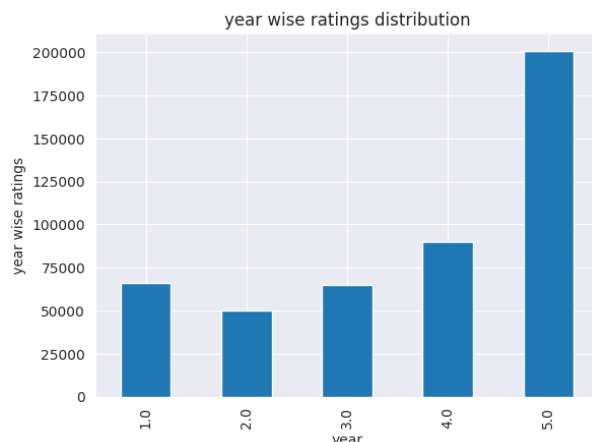


• سپس توزیع rating را بررسی کردیم:



که همانطور که مشخص است بیشتر نمره دهی ۵ بوده که بیانگر کیفیت محصولات عرضه شده می باشد.

- نمودار زیر را در نظر بگیرید:



نمودار بالا نشان می‌دهد کاربران معمولاً برای تجربه متوسط به محصول نمره نمی‌دهند، به همین دلیل است که رتبه‌بندی‌های ۲ و ۳ در مقایسه با سایر محصولات، فرکانس کمتری دارند. نمودار شمارش و تعداد محصول منحصر به فرد در هر دسته رتبه‌بندی نشان می‌دهد که در رتبه‌بندی ۵ فراوانی بیشتری وجود دارد و پس از آن رتبه‌بندی ۴ و رتبه‌بندی ۱ وجود دارد. این به این معنی است زمانی که کاربران بسیار راضی یا بسیار ناراضی هستند، رتبه‌بندی‌های بیشتری خواهند داد.

## ۴.۰ Data sparsity

بسیاری از سیستم‌های توصیه‌کننده با مشکلی به نام مشکل شروع سرد cold start مواجه می‌شوند. اساساً به یک کاربر نمی‌توان چیزی توصیه کرد زیرا به چیزی امتیاز نداده است! علاوه بر این، اگر محصول جدیدی را معرفی شود، کسی به آن امتیاز نداده و نمی‌توان آن را توصیه کرد. به این دلیل، فضای ریاضی نیز بسیار sparse می‌شود.

به همین دلیل، ما قصد داریم بررسی کنیم که چند رتبه برای هر کاربر ارسال شده است. به این ترتیب می‌توان ده کاربری که به محصولات رای داده بودند را پیدا کرد. سپس بررسی کردیم density ماتریس 0.00067 است سپس کاربرانی را که کمتر از ۵۰ بار به محصولات رتبه داده بودند را حذف کردیم به این ترتیب تعداد سطرها دیتاست به ۲۹۵۵۹ تا کاهش پیدا کرد. دوباره میزان density را بررسی کردیم که برابر با 0.48 شد. ماتریس user-item به صورت ۳۶۱ سطر در ۱۷۲۲۸ ستون تبدیل شد.

## ۵.۰ Collaborative filtering

User based collaborative filtering و Item based collaborative filtering دو رویکرد متداول هستند که در سیستم‌های توصیه‌گر برای ارائه توصیه‌های شخصی به کاربران استفاده می‌شوند. آنها از رفتار جمعی و ترجیحات کاربران برای ارائه توصیه‌ها استفاده می‌کنند.

• **User based collaborative filtering**: این روش بر یافتن شباهت‌های بین کاربران بر اساس تعاملات و ترجیحات گذشته آنها تمرکز دارد. فرض بر این است که کاربرانی که سلیقه و ترجیحات مشابهی در گذشته دارند، در آینده نیز ترجیحات مشابهی خواهند داشت.

• **Item based collaborative filtering**: این روش بر روی یافتن شباهت‌های بین آیتم‌ها بر اساس تعاملات کاربر تمرکز دارد. فرض بر این است که اگر دو مورد به طور مکرر توسط کاربران یکسان رتبه بندی یا با آنها تعامل داشته باشند، احتمالاً از نظر ترجیحات کاربر مشابه هستند.

به این ترتیب ابتدا ماتریس‌های **user-item** و **item-user** را حساب کردیم سپس **similarity** بین **user** ها و **similarity** بین **item** ها را بررسی کردیم. پس از محاسبه میزان شباهت، زیرمجموعه ای از کاربران مشابه که به "همسایگی" معروف است، انتخاب می شود. که متشکل از کاربرانی است که بیشترین شباهت را به کاربر هدف دارند. برای آیتم ها نیز به همین صورت عمل کردیم. سپس رتبه‌بندی‌های پیش‌بینی‌شده برای مواردی که کاربر هدف هنوز با آن‌ها تعامل نداشته است، بر اساس رتبه‌بندی کاربران همسایه محاسبه می‌شود. برای آیتم ها نیز به همین صورت است.

نتیجه حاصل برای را برای سه کاربر بر مبنای **User based collaborative filtering** بررسی کردیم و ۵ آیتم پیشنهادی برای هر یک از آنها به صورت زیر است:

user_ratings user_predictions			user_ratings user_predictions			user_ratings user_predictions		
Recommended Items			Recommended Items			Recommended Items		
B00AE0790U	0.0	2.078306	B000142FVW	0.0	0.783491	B0056VEYMS	0.0	0.336677
B00AQJ084Y	0.0	1.965524	B000PQDKU6	0.0	0.675815	B0092MCO88	0.0	0.526412
B00AWLB9I4	0.0	1.887036	B000ZVOD9W	0.0	0.636783	B0030HKJ8I	0.0	0.508133
B00AQ4EBOI	0.0	1.783238	B0068FDR96	0.0	0.618689	B0000YUX4O	0.0	0.499395
B008U2Y9BQ	0.0	1.770657	B006L1DNWY	0.0	0.609835	B000ELP5KA	0.0	0.497339