



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

## درس تحلیل کلان داده ها

استاد درس جناب آقای دکتر چهرقانی

(تمرین سوم)

محسن عبادپور | ۴۰۰۱۳۱۰۸۰ | [m.ebadpour@aut.ac.ir](mailto:m.ebadpour@aut.ac.ir)

نیمسال دوم سال تحصیلی ۱۴۰۱-۱۴۰۲



## فهرست مطالب

بخش اول: Data Stream	۳
سوال اول: عملکرد Adaptive Random Forest classifier (جنگل تصادفی تطبیقی)	۳
سوال دوم: بارگذاری مجموعه داده mnist	۴
سوال سوم: ایجاد جریان داده با استفاده از DataStream	۴
سوال چهارم: آموزش مدل Adaptive Random Forest classifier (جنگل تصادفی تطبیقی)	۴
سوال پنجم: نمودار خطا در حین آموزش	۵
فعالیت بیشتر: تحلیل عملکرد مدل بر اساس ماتریس در هم ریختگی	۵
بخش دوم: Singular Value Decomposition	۶
سوال اول: ارتباط ماتریس های تجزیه $M$ و $covariance(Y)$	۶
سوال دوم: ارتباط ماتریس های تجزیه $M$ و $Y$	۸
سوال سوم: یافتن ماتریس $M$	۸
سوال چهارم: تولید نمونه و انتقال تحت ماتریس محاسبه شده	۹
بخش سوم: Recommender Systems	۱۱
سوال اول: محاسبه مشتق نسبت به $r_{ui}$	۱۱
سوال دوم: پیاده سازی SGD	۱۲

### بخش اول: Data Stream

در این بخش پاسخ‌های مربوط به مبحث دسته‌بندی Data Stream آورده شده است که مشتمل بر پنج سوال می‌باشد.

#### سوال اول: عملکرد Adaptive Random Forest classifier (جنگل تصادفی تطبیقی)

دسته‌بند جنگل تصادفی یا Random Forest classifier یک دسته‌بند تجمعی (ensemble) بوده و با ترکیب چندین درخت تصمیم ضعیف (weak classifier) بصورت bagging یک دسته‌بند مطلوب را ایجاد می‌کند. bagging می‌تواند هم بصورت data bagging انجام پذیرد و هم بصورت feature bagging. بطور معمول درخت تصمیم ایجاد شده بصورت data bagging می‌باشد که هر درخت تصمیم صرفاً یک زیرمجموعه‌ای از مجموعه داده آموزشی را دیده و یاد می‌گیرد و در نهایت برای پیش‌بینی نمونه به همه‌ی درخت‌های تصمیم داده شده و نتایج با هم ادغام می‌گردد (رای گیری).

جنگل تصادفی تطبیقی یا Adaptive Random Forest classifier یک نسخه بهبود یافته جنگل تصادفی برای مقابله با حجم انبوه داده‌ها بصورت جریان داده می‌باشد. این الگوریتم بدین صورت کار می‌کند که در ابتدای کار هیچ درخت تصمیمی وجود دارد و اصطلاحاً جنگل خالی است. به محض ورود هر دسته داده (که می‌تواند یک داده نیز باشد)، یک درخت تصمیم جدید آموزش داده شده و به جنگل اضافه می‌شود. سپس تمام درخت‌های موجود با یک معیار ارزیابی نظیر Accuracy سنجیده می‌شوند؛ در صورتی که عملکرد هر درخت زیر حد آستانه تعیین شده باشد نشان می‌دهد که توزیع یاد گرفته‌ی آن outdate شده و عملکرد خوبی ندارد و پتانسیل آن را دارد که از جنگل حذف شود. البته در نسخه‌های مختلف این الگوریتم شروطی نیز برای تعیین تعداد درخت تصمیم موجود و جایگزینی درخت جدید (و نگهداری آن) صرفاً با درخت منسوخ ارائه شده است که چندین رویکرد برای بهبود عملکرد این مرحله وجود دارد که اینجا اشاره نمی‌شود. حال عملکرد کل جنگل موجود (دقت و خطا) ارزیابی می‌شود. اگر عملکرد کل زیر حد آستانه مشخص شده باشد یعنی عملکرد مدل کلی کاملاً نامطلوب بوده و می‌تواند حاکی از آن باشد که در توزیع جریان داده تغییرات اساسی به وجود

آمده است؛ در این حالت کل جنگل خالی شده و فرآیند آموزش درخت های کاملاً جدید با داده های کلا جدید ورودی از جریان داده مجدداً ادامه می یابد و با ورود دسته داده های جدید آموزش بصورت پیوسته ادامه یافته و این امکان به جنگل داده می شود که خود را با جریان داده و توزیع آن تطبیق دهد.

### سوال دوم: بارگذاری مجموعه داده mnist

برای بارگذاری مجموعه داده mnist از پکیج keras و زیر مجموعه ی dataset استفاده شده است که بصورت پیش فرض ۶۰ هزار داده آموزشی وجود دارد که بصورت تصاویر ۲۸\*۲۸ می باشد را در مقیاس ۰-۲۵۵ باز میگرداند. از جایی که scikit-multiflow بصورت بردار ویژگی پیاده سازی شده و امکان ورودی ماتریس وجود ندارد، من داده ها را بصورت بردارهای ۷۸۴ تایی تغییر شکل داده و نرمال نیز کرده ام. با توجه به همگرایی و همچنین محدودیت سخت افزاری بنده از ۲۰ هزار نمونه برای آموزش و حل این بخش از سوالات استفاده کرده ام.

### سوال سوم: ایجاد جریان داده با استفاده از DataStream

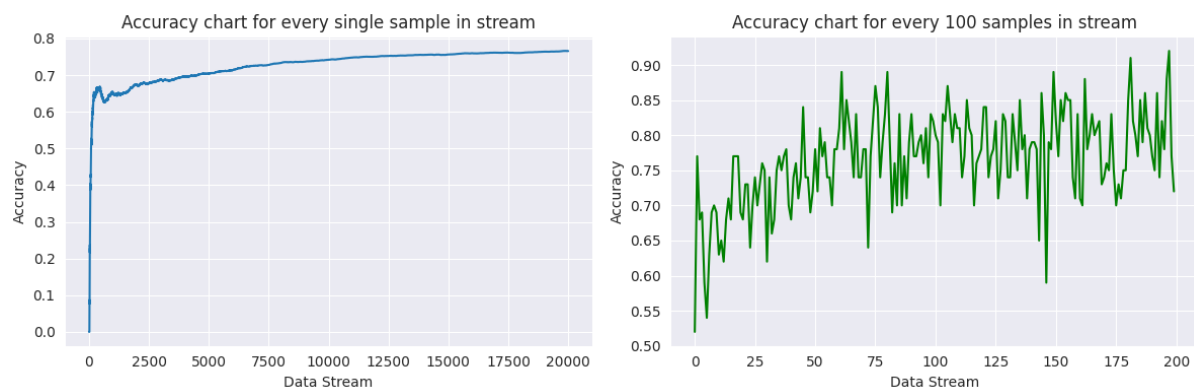
برای ایجاد جریان داده، مجموعه داده آموزشی مورد نظر را به همراه تعداد کلاس هدف که اینجا ۱۰ می باشد، به کلاس DataStream ورودی داده و شی مد نظر ساخته می شود جریان داده ای از مجموعه داده ی ما می باشد.

### سوال چهارم: آموزش مدل Adaptive Random Forest classifier (جنگل تصادفی تطبیقی)

برای آموزش مدل، از کلاس AdaptiveRandomForestClassifier از کتابخانه scikit-multiflow استفاده شده است. طبق سعی و خطا تعداد بهینه درخت های تصمیم موجود در جنگل تصادفی ۱۵ بدست آمده است؛ البته با توجه به ماهیت مجموعه داده که تصویری بوده و دسته بندی آن بصورت بردار ویژگی flat شده می باشد، تغییر پارامتر های مختلف در ایجاد جنگل تصادفی تطبیقی تاثیر چندانی در بهبود عملکرد مدل نداشت. آموزش این مدل بصورت incremental بوده و در هر مرحله یک نمونه وارد ورودی در نظر گرفته شده است. در پیاده سازی از متد های has\_more\_samples و next\_sample که مربوط به کلاس DataStream میباشد استفاده شده است.

## سوال پنجم: نمودار خطا در حین آموزش

دو نمودار خطای مدنظر تولید شده و بصورت زیر قابل مشاهده است. همانطور که دیده می‌شود پس از مشاهده تقریباً ۱۰ هزار نمونه آموزشی مدل به همگرایی نسبی رسیده و عملکرد آن در ادامه افزایش محسوس نداشته است. در خصوص عملکرد ۱۰۰ تایی نیز قابل نتیجه ای است که میانگین عملکرد موجود برابر با همان میانگین اصلی و به ازای نمونه های آموزشی تک و در حد ۸۰-۷۸٪ می باشد.



## فعالیت بیشتر: تحلیل عملکرد مدل بر اساس ماتریس در همریختگی

در این قسمت و به عنوان فعالیت بیشتر و اضافی ماتریس در همریختگی مربوط به مدل محاسبه شده و به صورت زیر می باشد؛ همانطور که دیده می‌شود عملکرد مدل بیشتر در مقابله با تصاویر اعداد ۳ و ۵ و ۸ با چالش مواجه بوده و خطا

Confusion Matrix for test data (10K samples) | Accuracy: 79.4%

0	0.884	0.000	0.031	0.027	0.000	0.015	0.019	0.002	0.022	0.000
1	0.000	0.880	0.103	0.004	0.000	0.001	0.004	0.000	0.007	0.001
2	0.018	0.010	0.902	0.015	0.006	0.005	0.014	0.016	0.012	0.003
3	0.008	0.008	0.072	0.821	0.002	0.045	0.002	0.022	0.010	0.011
4	0.007	0.005	0.032	0.007	0.742	0.037	0.009	0.006	0.008	0.147
5	0.037	0.008	0.047	0.179	0.010	0.646	0.011	0.012	0.036	0.013
6	0.029	0.004	0.093	0.007	0.027	0.014	0.816	0.003	0.006	0.000
7	0.001	0.015	0.044	0.007	0.016	0.005	0.001	0.844	0.007	0.061
8	0.023	0.038	0.076	0.149	0.015	0.072	0.004	0.013	0.569	0.041
9	0.011	0.009	0.009	0.024	0.042	0.030	0.002	0.068	0.007	0.799
	0	1	2	3	4	5	6	7	8	9

دارد. دلیل این امر آن است که ما تصاویر دو بعدی را در قالب بردار های ویژگی یک بعدی به مدل داده ایم (ماهیت درخت های تصمیم و جریان داده در کلاس گفته شده یک بعدی است) و از جایی که تصاویر این سه رقم در بردار های ویژگی و پیکسل های هم پوشان شان بیشتر است، درخت های موجود هر یک نظر و رای مختلف داشته و در نتیجه عملکرد مدل افزایش نمی یابد.

## بخش دوم

### بخش دوم: Singular Value Decomposition

در این بخش پاسخ‌های مربوط به مبحث دسته‌بندی Singular Value Decomposition آورده شده است که مشتمل بر چهار سوال و پرسش می‌باشد.

#### سوال اول: ارتباط ماتریس‌های تجزیه $M$ و $covariance(Y)$

$X_{2,n}$  sampled by  $N(0,1)_{20}$  : تجزیه ماتریس‌ها

$$Y = MX, \quad cov(Y) = C$$

$$X = U_X \Sigma_X V_X^T, \quad M = U_M \Sigma_M V_M^T$$

$$Y = U_Y \Sigma_Y V_Y^T, \quad C = U_C \Sigma_C V_C^T$$

توزیع نرمال  $X \rightarrow \Sigma_X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  (کواریانس)

$\Sigma_Y = M \Sigma_X M^T$  (کواریانس) تحت تبدیل (ضرب) ماتریسی

$$\Sigma_Y = (U_M \Sigma_M V_M^T) \Sigma_X (U_M \Sigma_M V_M^T)^T$$

$$= U_M \Sigma_M V_M^T \Sigma_X V_M \Sigma_M^T U_M^T$$

$\Sigma_M$  orthogonal

$$= U_M \Sigma_M \Sigma_M^T U_M^T = \Sigma_Y = U_C \Sigma_C V_C^T$$

(کواریانس)

$$\Sigma_C = \Sigma_M \Sigma_M^T = \Sigma_M^2$$

چون ماتریس‌های  $\Sigma$  ماتریس‌های قطری هستند و از طرفی چون تجزیه SVD برای هر ماتریس یکتا می‌باشد (نازغ از خطا محاسبات و تقریب) می‌توان طبق اسلایدهای درس نتیجه گرفت که باید ماتریس‌های  $U$  و  $V$  را به هم مرتبط کرد.

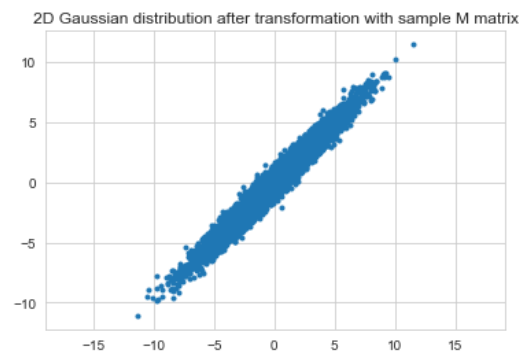
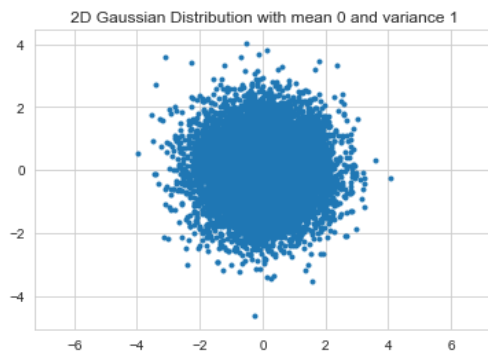
left singular در رابطه‌ای  $U_M \Sigma_M \Sigma_M^T U_M^T = U_C \Sigma_C V_C^T$  با هم

باید باشند لذا رابطه‌ای بین  $U_M$  و  $U_C$  بصورت تساوی حاصل می‌شود:

$$U_M = U_C$$

علاوه بر اثبات فوق و نشان دادن نحوه رسیدن به ارتباط فوق، بنده با شبیه سازی نیز صحت چنین رابطه را ارزیابی کردم که حاکی از درستی آن بود و نتیجه بصورت زیر است (یک ماتریس دلخواه  $M$  تعریف شده و بروی نمونه های گرفته شده از توزیع نرمال اعمال شد):

$$\text{Sample } M = \begin{bmatrix} 2.82 & -0.28 \\ 2.82 & 0.28 \end{bmatrix}$$



```
np.diag(S_C),np.diag(S_M**2)
✓ 0.0s
array([[16.16690675,  0.        ],
       [ 0.         ,  0.16081937]]),
array([[16.   ,  0.   ],
       [ 0.   ,  0.16]]))
```

```
U_C,U_M
✓ 0.0s
(array([[ -0.7071487 , -0.70706486],
        [-0.70706486,  0.7071487 ]]),
 array([[ -0.70710678, -0.70710678],
        [-0.70710678,  0.70710678]]))
```

## سوال دوم: ارتباط ماتریس‌های تجزیه M و Y

$$\begin{aligned}
 & \rightarrow Y = MX \\
 & U_Y \Sigma_Y V_Y^T = (U_M \Sigma_M V_M^T) (U_X \Sigma_X V_X^T) \\
 & \Sigma_X \rightarrow \text{چون حامل از توزیع نرمال} \\
 & \text{است، پس قطری بوده و پراکنش حول هر بُعد یکسان لذا} \\
 & \text{می‌توان } \Sigma_X \text{ را بصورت } c \times I \text{ نوشت که } I \text{ ماتریس} \\
 & \text{همانی و } c \text{ یک عدد است که می‌تواند جایزه شود:} \\
 & U_Y \Sigma_Y V_Y^T = (U_M \Sigma_M V_M^T) (U_X c \times I V_X^T) \\
 & = U_M c \times \Sigma_M \underbrace{V_M^T U_X V_X^T}_{\text{یکتا بودن تبزیه SVD}} \\
 & \text{و در نتیجه نام گذاری مجدد:} \\
 & U_Y \Sigma_Y V_Y^T = U_M c \times \Sigma_M V_{MX}^T \\
 & \boxed{\Sigma_Y = c \times \Sigma_M} \quad , \quad \boxed{\Sigma_X = c \times I}
 \end{aligned}$$

همانند قسمت قبل، برای این بخش نیز صحت رابطه‌ی بدست آمده را با ماتریس M دلخواه نشان می‌دهیم:

```

S_Y, S_M * np.mean(S_X)
✓ 0.0s
(array([402.06107626, 40.10540379]), array([401.5575755 , 40.15575755]))
    
```

## سوال سوم: یافتن ماتریس M

طبق محاسبات و روابط حاصل از دو قسمت قبلی و همچنین نتیجه‌گیری‌های صفحه ماتریس M بصورت زیر بدست می‌آید:

$$M = \begin{bmatrix} 6.0186 & 1.003 \\ 6.0210 & 1.998 \end{bmatrix}$$



$$* X = U_X \Sigma_X V_X^T$$

$$* M = U_M \Sigma_M V_M^T, \quad M^T = V_M \Sigma_M U_M^T$$

$$* Y_1 = U_{Y_1} \Sigma_{Y_1} V_{Y_1}^T, \quad Y_1 = M X, \quad \text{cov}(Y_1) = C_1$$

$$* Y_2 = U_{Y_2} \Sigma_{Y_2} V_{Y_2}^T, \quad Y_2 = M^T X, \quad \text{cov}(Y_2) = C_2$$

$$* C_1 = U_{C_1} \Sigma_{C_1} V_{C_1}^T, \quad C_2 = U_{C_2} \Sigma_{C_2} V_{C_2}^T$$

$$1. \Sigma_M = \sqrt{\Sigma_{C_1}} \quad \text{طبق الف}$$

$$2. U_M = U_{C_1} \quad \text{طبق الف}$$

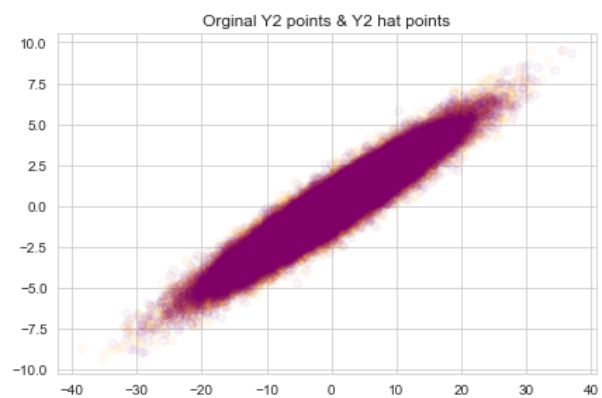
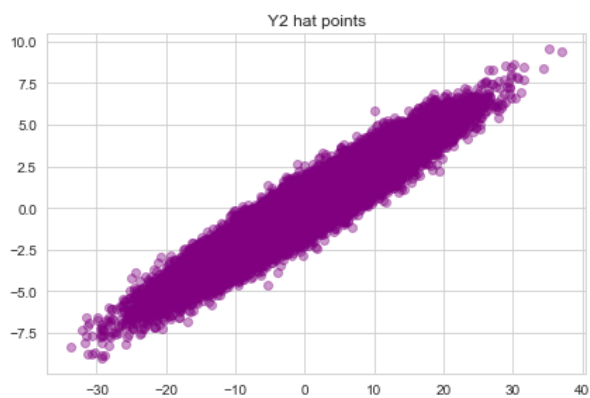
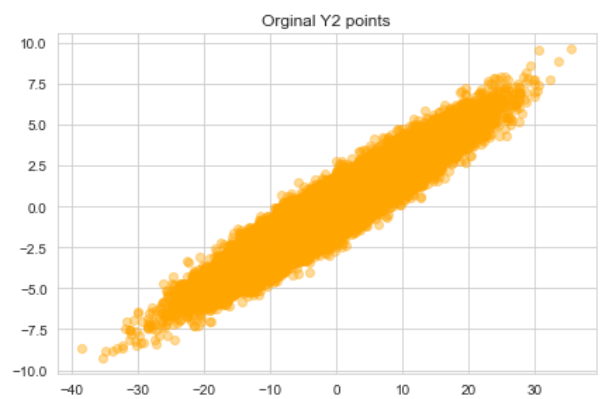
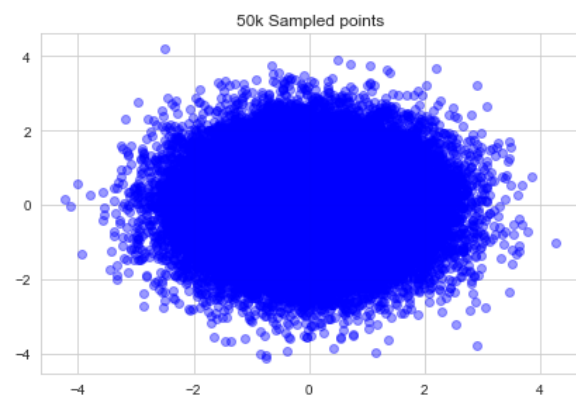
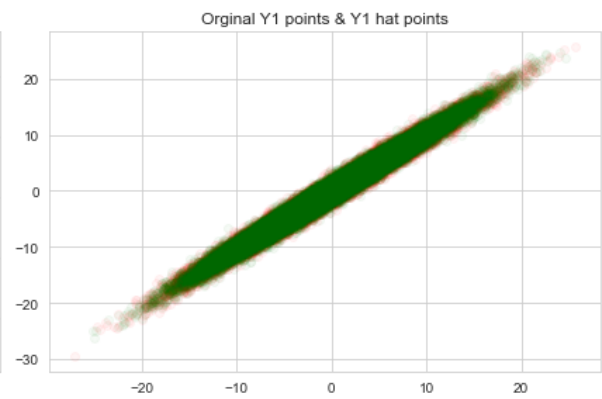
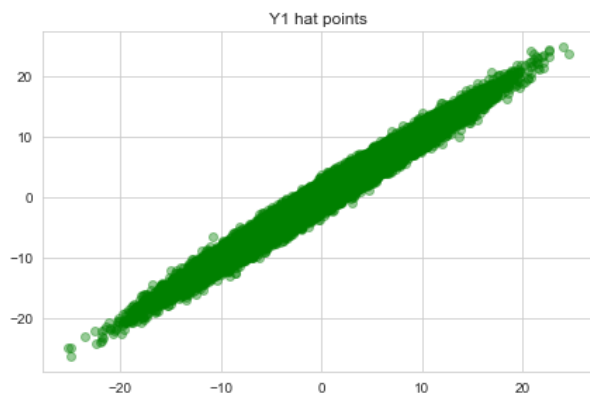
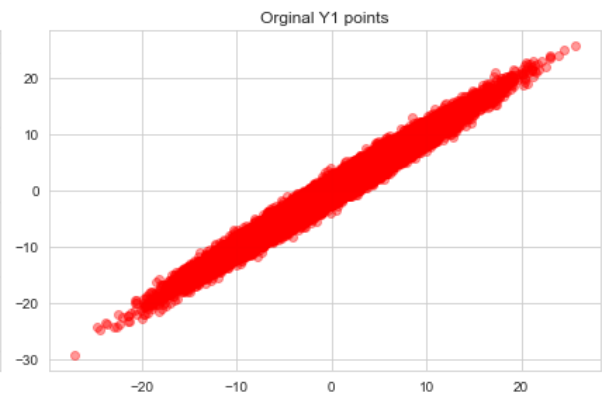
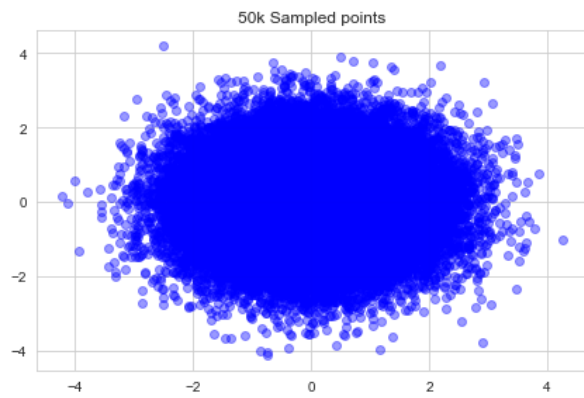
$$3. V_M = U_{C_2}$$

$$M = U_{C_1} \sqrt{\Sigma_{C_1}} U_{C_2}$$

### سوال چهارم: تولید نمونه و انتقال تحت ماتریس محاسبه شده

طبق صورت سوال، در این بخش ۵۰ هزار نمونه از توزیع گوسی با میانگین صفر و ماتریس کواریانس I تولید شده که در نمودار ها قابل مشاهده است. برای هر کدام از  $Y_1$  و  $Y_2$  چهار نمودار رسم شده است که اولی نمونه های تولید شده از توزیع گوسی ( $X$ )، دومی نمونه های اورجینال داده شده در پوشه ( $Y_1$  یا  $Y_2$ )، سومی نمونه های انتقال یافته از توزیع گوسی تحت ماتریس M محاسبه شده ( $\hat{Y}_1$  یا  $\hat{Y}_2$ ) و چهارمی نیز رسم توأمان هر دو می باشد.

حال اگر بخواهیم Y ها را با  $\hat{Y}$  ها مقایسه کنیم، ملاحظه میکنیم که شکل ها روی هم افتاده و بر هم منطبق شده اند که این دو نکته را نشان میدهد؛ اولین نکته اینکه ماتریس M بدست آورده شده صحیح می باشد و دومین نکته این است که یک ماتریس تبدیلی نظیر M تاثیر برابر و مشابهی بروی نمونه های حاصل از یک توزیع گوسی داشته است و روند نمونه گیری حاصل از توزیع گوسی تاثیری در رفتار ماتریس تبدیل M ندارد.



## بخش سوم: Recommender Systems

در این بخش پاسخ‌های مربوط به مبحث دسته‌بندی Recommender Systems آورده شده است که مشتمل بر دو سوال می‌باشد.

### سوال اول: محاسبه مشتق نسبت به $r_{ui}$

$$* \log(f(m))' = \frac{f'(m)}{f(m)} \Rightarrow \frac{\text{sigmoid}(f(m))'}{\text{sigmoid}(f(m))} \quad (I)$$

$$* \text{sigmoid}(f(m))' = \left( \frac{1}{1+e^{-f(m)}} \right)' = \frac{-(1+e^{-f(m)})'}{(1+e^{-f(m)})^2} \quad (II)$$

$$* -(1+e^{-f(m)})' = (-e^{-f(m)})' = -(-f(m))' e^{-f(m)} = (f(m))' e^{-f(m)} \quad (III)$$

$$* \log(f(m))' = \frac{(f(m))' e^{-f(m)}}{(1+e^{-f(m)})^2} = \frac{(1+e^{-f(m)}) \times (f(m))' e^{-f(m)}}{(1+e^{-f(m)})^2}$$

$$= \frac{(f(m))' e^{-f(m)}}{1+e^{-f(m)}} \quad f(m) = r_{ui} - q_i \cdot p_u^T$$

Extend =  $f(m) \left( \frac{1}{1+e^{-f(m)}} \right) (e^{-f(m)})$   
 $\text{sigmoid}(f(m))$

$$E' = \frac{(r_{ui} - q_i \cdot p_u^T) \cdot e^{-f(m)}}{1+e^{-f(m)}} + \lambda \left( \sum_i \|q_i\|_2^2 + \sum_u \|p_u\|_2^2 \right)$$

$$\frac{\partial E}{\partial r_{ui}} = \frac{1 \cdot e^{-f(m)} \cdot (-f(m))'}{1+e^{-f(m)}} + \lambda \left( \sum_i \|q_i\|_2^2 + \sum_u \|p_u\|_2^2 \right)$$

$r_{ui}$  مستقل از  $\lambda$

$$= \frac{1 \cdot e^{-f(m)} \cdot (-f(m))'}{1+e^{-f(m)}} = \frac{1 - \text{sigmoid}(r_{ui} - q_i \cdot p_u^T)}{1+e^{-f(m)}}$$

## سوال دوم: پیاده سازی SGD

\* برای پیاده سازی SGD و بروزرسانی وزن های ماتریس  $P$  و  $Q$ ، نیاز داریم لگاریتم نسبت به  $P$  و  $Q$  نیز محاسبه کنیم:

$$\frac{\partial E}{\partial P} = \frac{(y_{ui} - q_i p_u) \cdot e^{- (y_{ui} - q_i p_u^T)}}{1 + e^{- (y_{ui} - q_i p_u^T)}} + \lambda (\sum_u \|q_i\|^2 + \sum_u \|p_u\|^2)$$

$$= -q_i \cdot \exp(-(y_{ui} - q_i p_u^T)) \cdot \text{sigmoid}(y_{ui} - q_i p_u^T) + \lambda \sum_u 2 \|p_u\|$$

← بطور مشابه برای  $Q$

$$\frac{\partial E}{\partial Q} = -p_u \cdot \exp(-(y_{ui} - q_i p_u^T)) \cdot \text{sigmoid}(y_{ui} - q_i p_u^T) + \lambda \sum_i 2 \|q_i\|$$

برای پیاده سازی این قسمت از فرمت و قالب داده شده استفاده گردیده و کل قسمت های مشخص شده پر شده است. برای اهداف رسم نمودار های خطا یا پیش بینی و ارزیابی چند فیلد به کلاس اضافه شده است. همچنین برای راحتی و انسجام بیشتر چند تابع کمکی نیز تعریف و از آن استفاده شده است که عبارت اند از

۱. `read_data` برای خواندن دیتاست و ایجاد ماتریس  $R$
۲. `train_test_split` برای تفکیک و استخراج عناصر آموزش و آزمون از ماتریس  $R$  (۲۰ درصد از آیتام های هر کاربر برای آزمون در نظر گرفته شده است)
۳. `get_batches` برای ایجاد batch های تصادفی برای SGD به ازای کاربران
۴. `learning_rate_scheduler` برای تغییر نرخ آموزش با پیشروی تکرار های آموزشی
۵. `sum_squared_norms` برای محاسبه نرم ماتریس های  $Q$  و  $P$  برای بدست آوردن مقدار خطا در هر تکرار آموزشی

از مفاهیم یادگیری و نزول در راستای گرادیان می دانیم که نرخ یادگیری و گام حرکتی ثابت به هیچ عنوان تصمیم مطلوب و بهینه‌ای نمی باشد و طی اثبات های بهینه سازی میتوان نشان داد که نزول در راستای گرادیان و برای یک نقطه شروع تصادفی بهتر است با نرخ های بزرگ یادگیری و با سرعت بالا آغاز شده و به مرور گام های آموزشی نیز نرخ یادگیری کاهش یابد؛ با توجه به این نکته، بنده از نرخ یادگیری ۰.۲ شروع کرده و طی ۵۰ تکرار آموزشی آن را تقریباً به ۰.۰۱ بصورت خطی کاهش داده ام.

همچنین برای نرخ regularization نیز در زمان پیاده سازی و با سعی خطا به این نتیجه رسیدم که اگر برابر با ۰.۰۰۱ تنظیم شود همگرایی، دقت نهایی و روند آموزشی بهتر و مطلوب تری حاصل می شود و میتوان از فضای نهفته با ابعاد بالاتر نیز استفاده کرد و به همین جهت از این نرخ به جای ۰.۱ استفاده کرده ام. (همچنین در latent factor های بالا نرخ ۰.۱ باعث overflow شدن عناصر ماتریس P و Q در زمان بروزرسانی وزنهای می شد)

در یادگیری های انجام شده از batch size برابر با ۵۱۲ استفاده شده و در نتیجه هر تکرار آموزشی مشتمل بر ۵۹ بار بروزرسانی وزن شده است. (۲۹۸۵۸/۵۱۲) بنده در هر گام آموزشی طبق گرادیان محاسبه شده وزن ماتریس ها را بروز کرده و همچنین مقدار خطای حاصل را طبق تابع هزینه داده محاسبه کرده و در نهایت رسم کرده ام که در ادامه قابل مشاهده است. (میانگین خطای batch های هر تکرار)

بنده آموزش را به ازای ۳۲ و ۶۴ بعد برای factor انجام داده و برای هر کدام نمودار خطای آموزشی، معیار ۱۰ آیتم برتر (در عنوان نمودار درج شده است) و نمودار فراوانی تعداد آیتم های مشترک در مجموعه داده آزمون به همراه تغییرات نرخ یادگیری رسم کرده ام. به ازای تعداد فاکتور های ۳۲ معیار ۱۰ آیتم برتر برابر با ۰.۰۱۵۸ و به ازای تعداد فاکتور های ۶۴ معیار مذکور برابر با ۰.۰۲۲۷ حاصل شده است.

