



## تمرین سری چهارم

مهلت تحویل: پایان روز ۱۶ آذر ۱۴۰۳  
طراح تمرین: تکین جزایری - مهدی شفیعی

### نکات ضروری

۱. لازم است کد تمرین به همراه گزارش تحویل داده شود که شامل پاسخ سوالات نظری، توضیحات و تحلیل کد و تصاویر و نمودارهای لازم می باشد.
۲. کد تمرین باید در زبان پایتون، به صورت notebook و شامل نتایج اجرا باشد.
۳. فایل گزارش باید به فرمت pdf باشد. حجم گزارش به هیچ عنوان معیار نمره دهی نیست، لذا لطفا در حد نیاز توضیح دهید.
۴. استفاده از کتابخانه های از پیش آماده مجاز است، مگر اینکه نقیض آن گفته شود. در گزارش کتابخانه ها ذکر شوند.
۵. نمره هر تمرین از ۱۰۰ است (مگر آنکه نمره امتیازی در تمرین مشخص شده باشد). تا سه روز بعد از مهلت اعلام شده برای تمرین امکان بارگذاری وجود دارد اما به ازای هر روز ۱۰ درصد از نمره تمرین کسر خواهد شد.
۶. مشاهده شباهت نامتعارف در گزارش و کد به منزله تقلب می باشد و طرفین مشمول کسر نمره خواهند شد.
۷. برای این تمرین، می توانید گزارش و پاسخ هایتان را به زبان های فارسی و یا انگلیسی تحویل دهید.

## مسئله ۱.

مجموعه داده زیر شامل ۱۲ نمونه است که هر کدام دارای دو ویژگی و یک برچسب هستند.

(آ) با استفاده از این مجموعه داده و دسته‌بندی ساده<sup>۱</sup>، برچسب داده  $(X_1, X_2) = (0, 7)$  را تخمین بزنید. مراحل و محاسبات انجام شده را تشریح کنید.

No.	$X_1$	$X_2$	Y	No.	$X_1$	$X_2$	Y
1	0.6	6.1	1	7	-1.3	7.6	2
2	-1.8	5.5	1	8	1.3	7.9	3
3	-0.8	6.1	1	9	1.1	4.9	3
4	0	6.3	1	10	1.1	6.3	3
5	-0.3	7.8	2	11	2.9	7.1	3
6	-0.8	7.1	2	12	1.1	6.3	3

[۱۰ نمره]

(ب) فرض کنید برخلاف قسمت (الف)، هزینه متفاوتی برای دسته‌بندی‌های اشتباه<sup>۲</sup> در نظر بگیریم. ماتریس هزینه<sup>۳</sup> را که در آن درایه  $z, i$  نشان‌دهنده هزینه تشخیص کلاس  $i$  ام برای نمونه‌های کلاس  $z$  ام است، به صورت در نظر بگیرید. این ماتریس بیانگر چه تغییر رویکردی در حل این مسئله دسته‌بندی است؟  
با اعمال این ماتریس هزینه، برچسب داده  $(X_1, X_2) = (0, 7)$  را مجدداً تخمین بزنید. مراحل و محاسبات انجام شده را تشریح کنید.

$$L = \begin{bmatrix} 0 & 1 & 3 \\ 1 & 0 & 3 \\ 1 & 1 & 0 \end{bmatrix}$$

[۱۰ نمره]

## مسئله ۲.

هدف این مسئله، پیاده سازی الگوریتم‌های خوشه بندی بدون استفاده از کتابخانه‌هایی است که آن‌ها را به صورت آماده در اختیار می‌گذارند. سه مجموعه داده xclara، engytime و target در فضای دوبعدی به پیوست این تمرین در اختیار شما قرار داده شده است.<sup>۴</sup> نمودار نقطه‌ای<sup>۵</sup> آن‌ها را بدون استفاده از برچسب‌های موجود رسم کنید. (در تمامی قسمت‌های این مسئله، نتایج نهایی، نمودارهای رسم شده و تحلیل‌ها را در گزارش بیاورید.)

<sup>۱</sup>Naive Bayes Classifier

<sup>۲</sup>Misclassifications

<sup>۳</sup>Loss Matrix

<sup>۴</sup>منبع مجموعه داده‌ها:

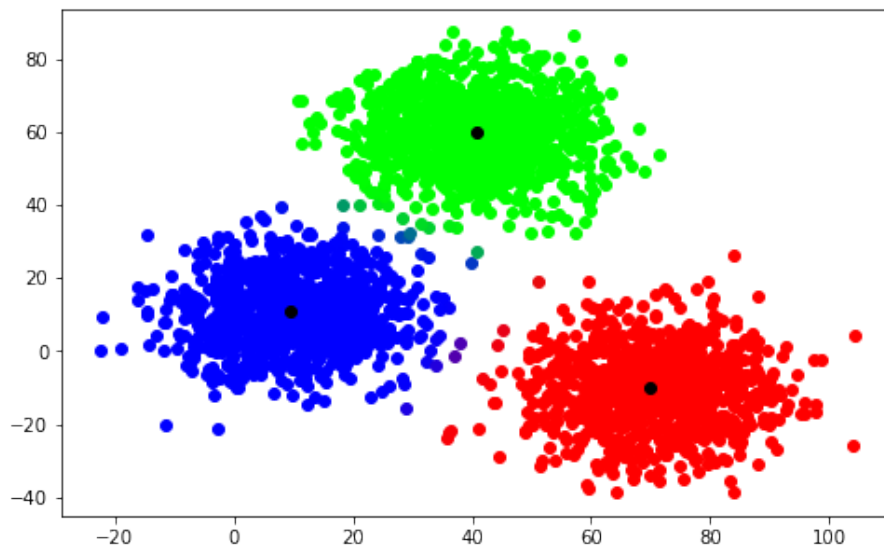
<https://github.com/milaan9/Clustering-Datasets>

<sup>۵</sup>Scatter Plot

(آ) الگوریتم  $k$ -mean را بدون استفاده مستقیم از کتابخانه‌ها پیاده‌سازی کنید. سپس برای هر مجموعه داده به ازای  $k$  از ۱ تا ۱۰، این الگوریتم را اجرا کنید و نتیجه را به صورت یک نمودار نقطه‌ای ترسیم کنید که در آن، داده‌های هر کلاس با رنگ‌های متفاوت از یکدیگر جدا شده باشند. در نهایت، برای هر مجموعه داده نمودار مقدار نهایی تابع هزینه ( $J$ ) نسبت به  $k$  ترسیم کنید و بر اساس آن مناسب‌ترین مقدار برای  $k$  را تعیین کنید.

[۱۵ نمره]

(ب) الگوریتم EM برای مدل ترکیبی گاوسی<sup>۶</sup> (GMM) را بدون استفاده مستقیم از کتابخانه‌ها پیاده‌سازی کنید و آن را روی مجموعه داده‌های xclara (با سه خوشه)، engytime (با دو خوشه) و target (با دو و سه خوشه) اجرا کنید. پس از هر یک از اجراها، نمودار نقطه‌ای داده‌ها را ترسیم کنید به طوری که هر خوشه دارای یک رنگ خاص باشد و هر داده بسته به احتمال تعلق به هر خوشه، ترکیبی از رنگ خوشه‌ها را داشته باشد. به طور مثال، نتیجه می‌تواند مشابه تصویر زیر باشد:



[۲۰ نمره]

(ج) (امتیازی) علاوه بر روش‌های خوشه‌بندی  $k$ -mean، GMM و سلسله‌مراتبی که با آن‌ها آشنا شده‌اید، الگوریتم مشهور دیگری به نام DBSCAN<sup>۷</sup> برای خوشه‌بندی در سال ۱۹۹۶ ارائه شده است که ایده آن متفاوت از روش‌های ذکر شده است. این روش یک الگوریتم مبتنی بر چگالی<sup>۸</sup> است؛ یعنی نقاط نزدیک به هم را در یک خوشه قرار می‌دهد و نقاطی را نیز که نزدیک به نقاط دیگر نیستند و در منطقه کم‌چگالی قرار گرفته‌اند، داده پرت در نظر می‌گیرد. این الگوریتم را بدون استفاده مستقیم از کتابخانه‌ها پیاده‌سازی کنید و با تغییر پارامترهای مدل، تلاش کنید به نتایج مناسبی روی هر یک از سه مجموعه داده برسید. برای هر اجرا، نمودار نقطه‌ای نمونه‌ها را ترسیم کنید.

[۲۰ نمره]

### مسئله ۳.

در این تمرین شما باید روی مجموعه داده‌ای کار کنید که از طریق فایل Loan\_Data.csv به شما داده شده است و شامل اطلاعاتی درباره متقاضیان بیمه است. هدف نهایی این است که پیش‌بینی کنید آیا بیمه به یک متقاضی تعلق می‌گیرد یا خیر. این تمرین شامل مراحل آشنایی داده و تحلیل کاوشگرانه داده<sup>۹</sup> (EDA)، پیش‌پردازش داده و پیاده‌سازی یک مدل رگرسیون لجیستیک<sup>۱۰</sup> است.

<sup>۶</sup>Gaussian Mixture Models

<sup>۷</sup>Density-Based Spatial Clustering of Applications with Noise

<sup>۸</sup>Density-Based

<sup>۹</sup>Exploratory Data Analysis

<sup>۱۰</sup>Logistic Regression

(آ) وجود مقادیر گم‌شده را بررسی کنید. سطرهای تکراری را در صورت وجود حذف کرده و ۱۰ سطر اول آن را چاپ کنید.

[۵ نمره]

(ب) تحلیل کاوشگرانه داده:

- (i) تعداد مقادیر یکتا در ستون `property_area` را محاسبه کرده و نمایش دهید.
- (ii) ماتریس همبستگی<sup>۱۱</sup> را رسم کنید. (می‌توانید از نقشه گرمایی<sup>۱۲</sup> استفاده کنید).
- (iii) برای بررسی رابطه بین ویژگی‌های `Gender` و `Loan_Status` نمودار میله‌ای رسم کنید.
- (iv) یک نمودار دایره‌ای برای ستون `Education` رسم کنید.
- (v) هیستوگرام ستون `ApplicantIncome` را با ده `bin` رسم کنید.

[۱۰ نمره]

(ج) پیش‌پردازش: تمام پیش‌پردازش‌های لازم (حذف مقادیر `null`، نرمالسازی، تبدیل داده‌های کیفی و ...) را انجام دهید و داده‌ها را به ۲۰ درصد برای آموزش و ۸۰ درصد برای تست تقسیم کنید.

[۵ نمره]

(د) با استفاده از کتابخانه‌های آماده، یک مدل رگرسیون لجیستیک آموزش دهید و ماتریس درهم‌ریختگی<sup>۱۳</sup> را محاسبه کنید.

[۱۰ نمره]

(ه) مدل رگرسیون لجیستیک را با استفاده از قالب زیر پیاده‌سازی کنید. در کد زیر باید توابع `predict`، `loss_derivate`، `loss` و حلقه `for` در تابع `fit` را کامل کنید. سپس `precision` و `recall` و `f1-score` را بر روی داده‌های تست محاسبه کنید.

```
import numpy as np
class GDLogisticRegression:
    def __init__(self, n_features, max_iter=50000, lr=0.0001, tol=1e-6, momentum=0.9):
        self.N = n_features
        self.beta = np.zeros((self.N+1,))
        self.max_iter = max_iter
        self.lr = lr
        self.tol = tol
        self.momentum = momentum

    def loss(self, X, y):
        pass

    def loss_derivative(self, X, y):
        pass
```

---

<sup>11</sup>Correlation Matrix

<sup>12</sup>Heatmap

<sup>13</sup>Confusion Matrix

```
def predict(self, X_test, threshold=0.5):  
    pass  
  
def fit(self, X_train, y_train):  
    X_train_new = np.concatenate((X_train, np.ones((X_train.shape  
        [0], 1))), axis=1)  
    last_loss = 0  
    momentum = 0  
  
    for _ in range(self.max_iter):  
        pass
```

[۱۵ نمره]

---

[بارم کل: ۲۰ + ۱۰۰ نمره]

با آرزوی موفقیت