

به نام خدا



دانشگاه صنعتی امیرکبیر  
( پلی تکنیک تهران )



DEPARTMENT OF COMPUTER  
ENGINEERING AND IT

دانشگاه صنعتی امیرکبیر

دانشکده‌ی مهندسی کامپیوتر

پروژه نهایی درس یادگیری ماشین

دکتر احسان ناظر فرد

طراح سوال:

سید اردلان قریشی

محمد رضا امامی ناصری

بهمن ۱۳۹۹

### توضیحات مهم:

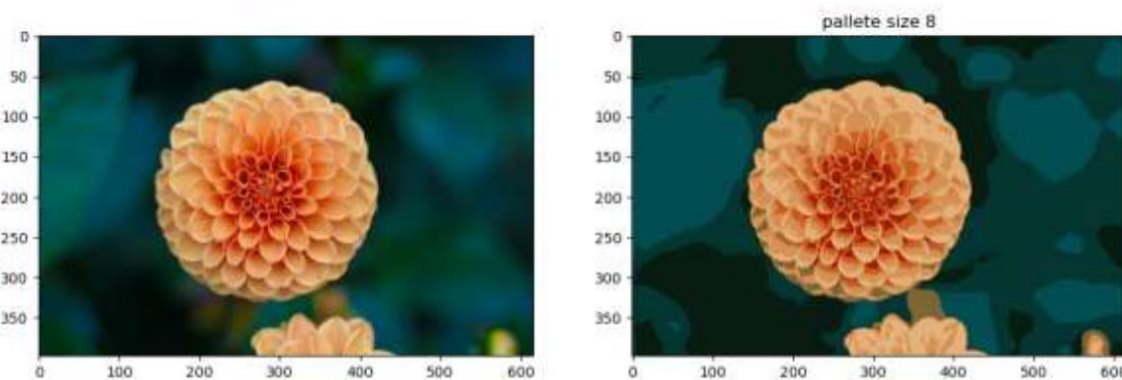
- تمامی مستندات خود شامل گزارش و کدهای خود را در یک فایل فشرده با فرمت zip ذخیره کرده و با عنوان #StudentId\_Project.zip بارگذاری نمایید (به عنوان مثال 99131000\_Project.zip).
- مهلت انجام تمرین تا ساعت ۸:۰۰ روز ۱۵ اسفند می باشد و به هیچ وجه تمدید نمی شود. پس از این ساعت امکان بارگذاری و ویرایش تکالیف در سامانه وجود نخواهد داشت.
- پروژه بدون گزارش فاقد ارزش می باشد و نمره ای به آن تعلق نمی یابد.
- تا حد ممکن سعی کنید اصول لازم برای گزارش مهندسی را رعایت نمایید (به بهترین گزارش نمره تشویقی تعلق می گیرد).
- مطابق قوانین دانشگاه هرگونه کپی برداری ممنوع می باشد و در صورت مشاهده نمره ی هر دو طرف صفر در نظر گرفته می شود.
- در صورت داشتن هرگونه ابهام می توانید از طریق ایمیل زیر سوال خود را مطرح نمایید:

[MLAUTFALL99@gmail.com](mailto:MLAUTFALL99@gmail.com)

۱- با استفاده از یک کتابخانه‌ی آماده که در آن الگوریتم خوشه‌بندی<sup>۱</sup> K-means وجود دارد، موارد زیر را پیاده‌سازی نمایید.

۱,۱) تصاویر bee.jpg و parrots.jpg را خوانده و نمایش دهید. هر تصویر از تعدادی پیکسل<sup>۲</sup> ساخته شده است و رنگ هر پیکسل با استفاده از ترکیب سه رنگ قرمز، سبز و آبی (RGB) ساخته می‌شود؛ به همین دلیل بعد از خواندن تصویر، مشاهده می‌کنید که تصویر خوانده شده یک ماتریس با مشخصات  $W \times H \times 3$  است به طوری که W و H اشاره به عرض<sup>۳</sup> و طول<sup>۴</sup> تصویر دارد و ۳ نشان دهنده‌ی هر کدام از سه رنگ RGB است. بنابراین، پیکسل‌های تصویر، داده‌های مورد نیاز مسئله می‌باشند که هر کدام دارای سه ویژگی هستند. پیکسل‌ها را با تعداد خوشه‌های ۲، ۳، ۴، ۵، ۶، ۱۰، ۱۵، ۲۰ خوشه‌بندی کنید.

- بعد از هر بار خوشه‌بندی تصاویر، رنگ پیکسل‌ها را با رنگ مرکز خوشه‌ای که در آن قرار می‌گیرند جایگزین کنید و تصویر حاصل را نمایش دهید.



شکل ۱ تصویر سمت چپ تصویر واقعی و تصویر سمت راست با استفاده از الگوریتم k-means با تعداد خوشه ۸ ایجاد شده است

۱,۲) در این بخش مجموعه داده‌ی Shill Bidding Dataset.csv را بارگذاری کنید.

الف) یکی از روش‌های تعیین تعداد خوشه‌های بهینه در الگوریتم k-means استفاده از روش elbow است؛ این روش را توضیح دهید.

ب) تعداد خوشه‌ها را از ۱ تا ۱۰ تغییر دهید و الگوریتم را اجرا نمایید. با توجه به روش elbow بهترین تعداد خوشه، برای خوشه‌بندی مجموعه داده را مشخص نمایید و دلیل انتخاب خود را

<sup>1</sup> clustering

<sup>2</sup> pixel

<sup>3</sup> width

<sup>4</sup> height

توضیح دهید. نمودار هزینه بر حسب تعداد خوشه<sup>۵</sup> را رسم کنید. (برای تابع هزینه می‌توانید از inertia یا distortion استفاده نمایید).

ج) معیار purity را به ازای تعداد خوشه برابر با ۲ ( $k=2$ ) محاسبه نمایید.

۲- با استفاده از الگوریتم خوشه‌بندی DBSCAN برای هر یک از مجموعه داده‌های موجود در پوشه‌ی مربوط به این سوال، نمونه‌ها را همراه با خوشه‌ی نسبت داده شده<sup>۶</sup> رسم کنید. به این نکته توجه کنید که داده‌ها می‌توانند متعلق به هیچ خوشه‌ای نباشند و می‌توانند هنگام نمایش به عنوان نویز<sup>۷</sup> تلقی شده و نمایش داده شوند. پس از اجرای الگوریتم خوشه‌بندی برای هر یک از مجموعه داده‌ها معیار purity را به دست آورده و به طور کیفی تاثیر نوع مجموعه داده بر کیفیت خوشه‌بندی را مقایسه و تحلیل کنید (در این سوال استفاده از کتابخانه آزاد است).

مجموعه داده‌ها: Compound – pathbased – rings – spiral – D31

۳- در این بخش می‌خواهیم دو الگوریتم value iteration و policy iteration را برای محیط<sup>۸</sup> Frozen Lake مانند شکل زیر پیاده‌سازی نماییم. (برای آشنایی بیشتر می‌توانید به مستندات ارائه شده در پانویس<sup>۹</sup> مراجعه کنید).

S	F	F	F	H
F	F	H	H	F
F	F	F	F	F
H	H	F	H	F
F	F	F	F	G

<sup>۵</sup> graph of cost sequence

<sup>۶</sup> assigned cluster

<sup>۷</sup> noise

<sup>۸</sup> environment

<sup>۹</sup> [Link1](#), [Link2](#), [Link3](#)

در این محیط عامل<sup>۱۰</sup> با شروع حرکت از خانه‌ی شروع<sup>۱۱</sup> (S) می‌خواهد به خانه‌ی هدف<sup>۱۲</sup> (G) برسد. در این بین خانه‌های یخ زده<sup>۱۳</sup> (F) هم وجود دارد. همچنین گودال<sup>۱۴</sup>‌هایی (H) نیز در نقشه (محیط) موجود است. عامل باید از طریق خانه‌های یخ‌زده حرکت کرده و به خانه‌ی هدف برسد. توجه کنید که عامل به هدف محیط شامل احتمال گذار وضعیت‌ها<sup>۱۵</sup> و میزان پاداش<sup>۱۶</sup> دسترسی دارد. در این بخش قصد داریم سیاست بهینه<sup>۱۷</sup> را برای محیط ۵ در ۵ شکل بالا به دست آوریم. هنگامی که عامل به خانه‌ی G یا H برسد یک اپیزود<sup>۱۸</sup> تمام می‌شود. در صورتی که عامل در خانه G قرار بگیرد، پاداش +۱ می‌گیرد. در مابقی موارد عامل پاداشی دریافت نمی‌کند. در پیاده‌سازی الگوریتم‌ها ۰.۸۵=γ در نظر بگیرید (انتخاب شرط خاتمه‌ی مناسب به عهده‌ی شما می‌باشد).

**الف)** الگوریتم value iteration را پیاده‌سازی کرده و مقادیر  $V^*$  را به دست آورید. زمان اجرا و تعداد تکرار<sup>۱۹</sup> مورد نیاز را نمایش دهید. سیاست بهینه را به دست آورید و آن را به صورت یک جدول متشکل از حروف U (بالا)، R (راست)، D (پایین) و L (چپ) نمایش دهید.

**ب)** الگوریتم policy iteration را نیز مانند حالت قبل پیاده‌سازی کرده و زمان اجرا و تعداد تکرار آن را با مورد قبل مقایسه کنید. سیاست بهینه را مانند قسمت قبل نمایش دهید.

**نکته:** برای پیاده‌سازی قسمت‌های **الف** و **ب** می‌توانید از ابزار gym استفاده نمایید. برای آشنایی بیشتر به مستندات<sup>۲۰</sup> آن رجوع کنید.

**۴-** مجموعه داده‌ی SeoulBikeData.csv در فایل مجموعه داده‌ها قرار داده شده است. با استفاده از آن موارد زیر را انجام دهید (استفاده از کتابخانه در تمامی بخش‌های سوال مجاز است).

**الف)** پیش‌پردازش<sup>۲۱</sup>‌های لازم را انجام دهید.

**ب)** همبستگی<sup>۲۲</sup> بین ویژگی‌ها را استخراج کرده و با توجه به آن بهترین ویژگی‌ها را انتخاب کنید. در این مرحله شما باید تعداد ویژگی‌های انتخاب شده را با توجه به یک مدل رگرسیون خطی پایه مورد بررسی قرار داده و بهترین K را پیدا کنید.

<sup>10</sup> agent

<sup>11</sup> start

<sup>12</sup> goal

<sup>13</sup> frozen

<sup>14</sup> hole

<sup>15</sup> state-transition probability

<sup>16</sup> reward

<sup>17</sup> optimal policy

<sup>18</sup> episode

<sup>19</sup> iteration

<sup>20</sup> gym doc ([Link](#))

<sup>21</sup> preprocess

<sup>22</sup> correlation

ج) با استفاده از داده پیش‌پردازش شده مدل رگرسیون لسو<sup>۲۳</sup> را آموزش دهید. نقش پارامتر  $\alpha$  در این مدل را بررسی کرده و با جستجو، بهترین مقدار آن را به دست آورید.

د) ویژگی‌های انتخاب شده در بخش ب و ج را با هم مقایسه کنید. چه نتیجه‌ای می‌گیرید؟

ه) برای بهبود عملکرد مدل چه پیشنهادی دارید؟

۵- مجموعه داده‌ی heart\_failure\_clinical\_records\_dataset.csv را بارگذاری کنید. شما باید با استفاده از ویژگی‌های موجود، هر فرد را بر اساس مقادیر ستون DEATH\_EVENT دسته‌بندی کنید (استفاده از کتابخانه در تمامی بخش‌های سوال مجاز است).

الف) پیش‌پردازش‌های لازم را انجام دهید. این مجموعه داده شامل مقادیر گم شده<sup>۲۴</sup> است. روش‌های مختلفی برای برطرف کردن این مشکل پیشنهاد شده است. درباره‌ی آن‌ها تحقیق کرده و با ذکر دلیل یکی از این روش‌ها را انتخاب کرده و مقادیر گم شده مجموعه داده را برطرف کنید.

ب) بهترین K ویژگی را با توجه به اهمیت آن‌ها انتخاب کنید. همانند بخش ب سوال ۴، باید تعداد ویژگی‌های انتخاب شده را با توجه به یک مدل دسته‌بند پایه<sup>۲۵</sup> مورد بررسی قرار دهید و بهترین K را پیدا کنید<sup>۲۶</sup>.

ج) ۳ مدل مختلف رای‌گیری<sup>۲۷</sup> که هر کدام شامل ۳ دسته‌بند است را برای بهترین K ویژگی آموزش دهید و بهترین مدل را انتخاب کنید.

د) دسته‌بندهای مورد استفاده در بهترین مدل را با استفاده از داده‌های به دست آمده در بخش ب (بهترین K ویژگی) به صورت مجزا آموزش دهید. از مقایسه عملکرد دسته‌بندها به صورت تکی و گروهی چه نتیجه‌ای می‌گیرید؟

با آرزوی موفقیت!

<sup>23</sup> Lasso regression

<sup>24</sup> missing value

<sup>25</sup> base classifier

<sup>26</sup> more info ([Link](#))

<sup>27</sup> voting