



۱. (۵٪) [مطالعه و تحلیل] تفاوت‌های میان یادگیری ماشین، داده‌کاوی و بازشناسی الگو را ذکر کنید.

۲. (۵٪) [مطالعه و تحلیل] تفاوت روش‌های تمایزی (Discriminative) و تولیدی (Generative) در یادگیری ماشین را از نقطه نظر احتمالاتی شرح دهید.

۳. (۴۰٪) [پیاپی‌سازی: یک دسته‌بند ساده و معیارهای ارزیابی] در این مسئله شما یک دسته‌بند ساده را به منظور تشخیص سه زبان فارسی، عربی و کردی در یک متن طراحی می‌کنید. داده‌های مربوط به این تمرین در فایل ANN-HW1-Data.xlsx قرار گرفته است که حاوی ۵۰ جمله برای هر زبان است. از این داده، برای هر زبان، ۸۰٪ جملات را برای آموزش و مابقی ۲۰٪ را برای آزمون جدا کنید. الف) (۱۰٪) برای تشخیص این سه زبان، برای هر جمله تعداد پنج نویسه (کاراکتر) «پ، ژ، گ، چ، ع» را به تعداد کل نویسه‌های آن جمله تقسیم کنید و بر اساس مقدار آن در مورد نوع زبان تصمیم بگیرید. برای این منظور، میانگین این معیار را برای همه جملات آموزشی در هر زبان حساب کنید و از آن به عنوان آستانه کمک بگیرید. بعد از تعیین مقادیر آستانه، داده‌های مجموعه آزمون را برای ارزیابی روش خود به سیستم ارائه دهید و نتایج حاصل را با معیارهای صحت (Accuracy)، دقت (Precision)، یادآوری (Recall) و F-Measure گزارش کنید.

ب) (۲۰٪) می‌خواهیم کارایی سیستم را مقداری بهبود دهیم و به جای یک عدد به عنوان ویژگی، یک بردار ویژگی استفاده کنیم. به این منظور، از کل تعداد نویسه‌های متن به عنوان ویژگی استفاده می‌کنیم. به منظور طراحی این دسته‌بند ابتدا لازم است برداری با نام Character Frequency (CF) را معرفی کنیم. تعداد عناصر (مولفه‌های) این بردار برابر تعداد کل نویسه‌های موجود در تمامی سه زبان است. آنگاه برای هر جمله، عناصر این بردار برابر با فراوانی نرمال شده تعداد نویسه‌های آن متن استفاده می‌کنیم. به عنوان مثال فرض کنید نویسه‌های (A,B,C,D,E) تمامی نویسه‌های مجاز در سه زبان باشند. در این صورت، بردار CF متناظر با متن فرضی "ABBDCDEDEABB" به صورت (2, 4, 1, 3, 2) خواهد بود.



حال با استفاده از بردار CF، بردار دیگری با نام Normal CF (NCF) با استفاده از فرمول  $NCF(i) = CF(i)/\text{Sum}(CF)$  تعریف می‌شود. بعد از نرمال کردن، بردار این جمله به صورت (0.17, 0.33, 0.08, 0.25, 0.17) خواهد بود (تقسیم مولفه‌ها بر ۱۲).

با داشتن این بردار برای هر جمله، میانگین همه بردارهای داده آموزش را برای هر زبان محاسبه کنید و از آن به عنوان معیار مقایسه داده‌های آزمون استفاده کنید. بدین صورت که برای هر داده آزمون، هر زبانی که میانگین بردارهای NCF آموزش آن، کم‌ترین فاصله اقلیدسی را با NCF جمله آزمون داشته باشد، به عنوان زبان آن جمله تشخیص داده می‌شود.

معیارهای صحت (Accuracy)، دقت (Precision)، یادآوری (Recall) و F-Measure را در این حالت نیز محاسبه کنید و با نتایج بخش الف مقایسه کنید و مشاهده و تحلیل خود را گزارش کنید.

ج) (۱۰٪) نتایج بخش ب را با ارزیابی مبتنی بر روش 5-fold Cross-Validation گزارش کنید. برای این کار کل مجموعه داده را استفاده کنید (و نه فقط مجموعه آموزش).

۴. (۵۰٪) [پیاده‌سازی: پرسپترون] مسئلهٔ بازشناسی نویسه را برای الگوهای بیان شده در مثال ۲-۱۵ در فصل دوم کتاب، با استفاده از ساختار شبکه نشان داده شده در شکل ۲-۲۱ پیاده‌سازی کنید. بدین منظور الگوهای آموزش شکل ۲-۲۰ و الگوهای آزمون (تست) شکل ۲-۲۲ که به صورت فایل‌های متنی به همراه تمرین ارائه شده است، به کار بگیرید. برای موارد زیر نتیجه را گزارش کنید. گزارش‌ها درصد خطای بازشناسی را به صورت زیر محاسبه و گزارش کنید.

$$Error Rate = \frac{N_{err}}{N} \times 100 = \frac{\text{تعداد الگوهایی که اشتباهی بازشناسی شده اند}}{\text{تعداد کل الگوها}} \times 100$$

نکته: برای بهتر دیدن الگوهای ارائه شده در فایل‌های متنی، می‌توانید از فونت Courier New یا CourierPS در NotePad استفاده کنید.



۱-۴) [پرسپترون] آموزش شبکه را با استفاده از یادگیری پرسپترون انجام دهید و نتیجه بررسی موارد زیر را گزارش کنید.

۱-۴-۱) (۱۰٪) [مقادیر اولیه وزن‌ها] درصد خطای بازشناسی را برای داده‌های آزمون، به ازای اجراهای مختلفی از برنامه با ۳ مقدار مختلف اولیه وزن‌ها و بایاس‌های شبکه پیدا کنید (یک بار صفر و دو بار تصادفی). درصد خطای بازشناسی را برای داده‌های آموزش و آزمون با استفاده از بهترین مقادیر اولیه محاسبه کنید. با استفاده از نتایج، در مورد اثر مقادیر اولیه بر کارایی و سرعت شبکه بحث کنید.

۱-۴-۲) (۵٪) [مقدار آستانه] برنامه را برای ۳ مقدار مختلف آستانه  $\theta$  تکرار کنید و نتایج بازشناسی را برای داده‌های آزمون ارائه کنید. آیا مقدار بیشتر  $\theta$  تأثیری بر تعداد دفعاتی که شبکه اشتباه می‌کند، دارد؟

۱-۴-۳) (۵٪) [نرخ یادگیری] شبکه را برای ۳ مقدار مختلف از نرخ یادگیری، ۰.۱، ۰.۵ و ۰.۹ آموزش دهید و نتیجه بازشناسی بر روی داده‌های آزمون را به صورت نمودار رسم کنید. زمان همگرایی شبکه را در هر حالت گزارش کنید. با توجه به نتایج حاصل، در مورد تأثیر این پارامتر بر عملکرد شبکه بحث کنید.

۱-۴-۲) (۱۰٪) برای حالتی که برای الگوی تست، بیش از یکی از دسته‌ها انتخاب می‌شوند، چه راهکاری را پیشنهاد می‌کنید. روش خود را پیاده‌سازی کرده و نتیجه را گزارش کنید.

۱-۴-۳) (۱۰٪) برای افزایش توانایی شبکه به ویژه در برخورد با داده‌های نویزی مجموعه تست چه روشی را پیشنهاد می‌کنید.

۱-۴-۴) (۱۰٪) [استخراج ویژگی] در بخش‌های قبل از این تمرین، مقدار کل پیکسل‌های هر نویسه به عنوان ورودی شبکه استفاده شد. در این بخش، از روش تصویر کردن (projection) برای استخراج



ویژگی استفاده می‌شود و مقدار ویژگی‌ها (به جای مقادیر پیکسل‌ها) به عنوان ورودی به شبکه داده می‌شود. در این روش، به ازای هر ردیف (و هر ستون) از هر نویسه، مجموع پیکسل‌های روشن (با مقدار یک) آن ردیف (یا ستون) شمارش شده و مقدار حاصل جمع به عنوان ویژگی در نظر گرفته می‌شود. با توجه به ابعاد نویسه‌ها که  $9 \times 7$  هستند، تعداد ویژگی‌های هر نویسه  $9+7=16$  خواهد بود. مقادیر بدست آمده برای ویژگی‌ها را به گونه‌ای نرمال کنید که همه مقادیر بین صفر و یک قرار بگیرند (بر بیشترین مقدار تقسیم کنید).

شبکه را برای آموزش با ویژگی‌های حاصل تغییر دهید (با بهترین پارامترهای ممکن از تجربیات بخش‌های قبلی، شامل مقدار آستانه و نرخ یادگیری و نحوه مقداردهی اولیه وزن‌ها) و نتایج حاصل را (نرخ خطا روی مجموعه آزمون) گزارش کنید. آیا کارایی شبکه در این حالت، به نسبت حالتی که مقدار خود پیکسل‌ها استفاده شود، بهبود می‌یابد یا خیر؟ تحلیل خود را از این نتایج بیان کنید.