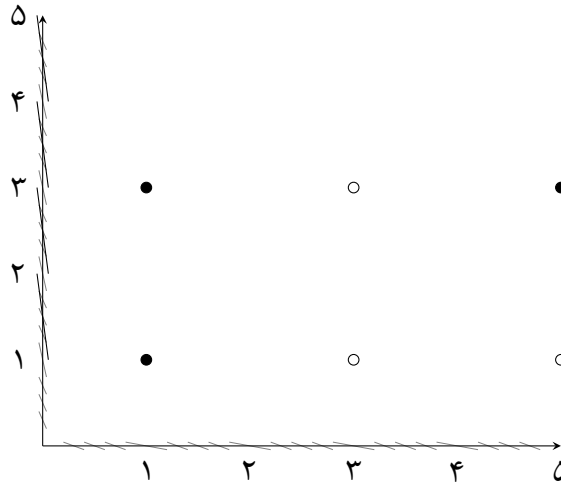


**سوال ۱: (30 Points) یادگیری جمعی**

۱. داده‌های شکل زیر را در نظر بگیرید و با فرض این که هر دسته بند ضعیف یک خط افقی یا عمودی باشد، آداپوست را دو مرحله بر روی این داده‌ها اجرا کرده و وزن نمونه‌ها، مرز دسته‌بندی و خطای آموزش را در هر مرحله به دست آورید.



۲. فرض کنید با دو دسته‌بند ضعیف مواجه هستیم. دسته‌بند اول دسته‌بندی است که بر روی داده‌های آموزش *over fit* می‌کند و قادر به ارائه عملکرد مثبت روی داده‌های دیده نشده نیست. همچنین دسته‌بند دوم هم با وجود قدرت تعمیم‌دهی بالا اما توانایی یادگیری مدل‌های پیچیده را ندارد. برای تقویت عملکرد هر کدام، از طریق روش‌های یادگیری جمعی، کدام از یک تکنیک‌های *bagging* و یا *boosting* را پیشنهاد می‌دهید؟ توضیح دهید.

۳. فرض کنید که  $h_1$  تا  $h_t$  و  $\alpha_1$  تا  $\alpha_t$  فرضیه‌ها و ضرایب فرضیه‌های به دست آمده از اجرای آداپوست باشند و داشته باشیم

$$H_t(x) = \frac{1}{\gamma} (\alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_t h_t(x))$$

و همچنین loss زیر را تعریف کنیم

$$Loss(y, H_t(x)) = \sum_{i=1}^N e^{-y^{(i)} H_t(x^{(i)})}$$

• اولاً نشان دهید عبارتی که در تکرار  $t$  از آداپوست بایستی کمینه شود عبارت زیر است

$$\sum_{i=1}^N w_t^{(i)} e^{-\frac{1}{\gamma} \alpha_t y^{(i)} h_t(x^{(i)})}$$

• ثانیاً، نشان دهید که کمینه کردن عبارت بالا معادل است با کمینه کردن عبارت زیر

$$\sum_{i=1}^N w_t^{(i)} * I(y^{(i)} \neq h_t(x^{(i)}))$$

• حال با مشتق‌گیری از  $loss$  بر حسب  $\alpha_t$  و صفر قرار دادن مشتق مقدار بهینه  $\alpha_t$  که برابر با زیر است را به دست آورید:

$$\epsilon_t = \frac{\sum_{i=1}^N w_t^{(i)} * I(y^{(i)} \neq h_t(x^{(i)}))}{\sum_{i=1}^N w_t^{(i)}}$$

$$\alpha_t = \ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)$$

## سوال ۲: (35 Points) شبکه پرسپترونی چندلایه

۱. یک شبکه عصبی را با  $k$  لایه در نظر بگیرید. اگر از تابع همانی به عنوان activation function در تمام لایه ها استفاده شود، در رابطه با خطی یا غیرخطی بودن مدل توضیح دهید.

۲. همانطور که میدانید تمام Boolean function ها را می توان با یک شبکه عصبی دو لایه مدل کرد:

- یک تابع بولین با  $N$  متغیر را در نظر بگیرید. بیشترین تعداد نورون ها در لایه مخفی چند است و در چه حالتی رخ می دهد.
- آیا می توان با افزایش لایه ها، تعداد نورون ها در لایه مخفی را نسبت به اندازه ورودی خطی کرد؟ با رسم شبکه عصبی نشان دهید.

۳. تعداد ۱۰ نمونه از دو کلاس  $C1$  و  $C2$  در زیر آمده است.

$$c1: [0/1, -0/2], [0/2, 0/1], [-0/15, 0/2], [1/1, 0/8], [1/2, 1/1]$$

$$c2: [1/1, -0/1], [1/25, 0/15], [0/9, 0/1], [0/1, 1/2], [0/2, 0/9]$$

- آیا این نقاط به صورت خطی تفکیک پذیر هستند؟ نشان دهید.
- اگر به صورت خطی تفکیک پذیر نیستند، یک شبکه عصبی با activation function پله برای آن طراحی کنید تا بتواند نقاط کلاسه های  $c1$ ,  $c2$  را به خوبی دسته بندی کند.

۴. نشان دهید اگر تابع هزینه در یک شبکه عصبی cross entropy باشد:

- خروجی شبکه عصبی که با استفاده از وزن های بهینه بدست می آید برابر تخمین  $p(w_i|x)$  است.
- اگر activation function برابر تابع سیگموئید باشد، آنگاه:

$$\delta_i^L(i) = \frac{\partial \varepsilon(i)}{\partial v_j^L(i)}$$

$$\delta_j^L(i) = a(1 - \hat{y}_j(i))y_j(i)$$

## سوال ۳: (35 + 20 Points) تخمین ساختار دوم پروتئین (Protein Secondary Structure)

یکی از مسائل مطرح در حوزه بیوانفورماتیک، تخمین ساختار دوم پروتئین بر اساس ساختار اول آن یعنی رشته آمینواسیدی پروتئین است. برای توصیف ساختار دوم پروتئین می توان رشته آمینواسید آن را به زیررشته هایی تقسیم کرد و برای هر زیر رشته به صورت محلی ساختار دوم را بررسی کرد. در این مسئله از داده هایی با فرمت dssp استفاده می شود. برای آشنایی با فایل dssp می توانید به این لینک مراجعه کنید. هر فایل dssp مربوط به یک رشته پروتئین است. در این فایل که ساختار ماتریسی دارد، هر سطر مربوط به یک آمینواسید از پروتئین است که به ترتیب محل قرار گیری در رشته آمده است. هر چند سطر با هم به عنوان یک زیر رشته در نظر گرفته شده است. ستون های موجود در این ماتریس در لینک بالا توضیح داده شده است. ستون های ابتدایی مربوط به شماره آمینواسید و نام آن است و ستون های بعدی هر کدام یکی از ویژگی های ساختار دوم را مشخص می کند. در این تمرین از دو ستون این ماتریس استفاده می شود که برای راحتی کار شما، در هر فایل تنها همین دو ستون قرار دارند و ستون های دیگر حذف شده اند. لازم به ذکر است که داده هایی که در اختیار شما قرار گرفته اند مربوط به ۵۴۲۶ پروتئین است که هر کدام در فایل جداگانه ای قرار دارد. ستون اول در فایل های داده شده مربوط به نام آمینواسید است که به عنوان ویژگی در فرآیند آموزش استفاده می شود. ستون دوم مربوط به ساختار آن آمینو اسید است که با یک حرف مشخص شده است و به عنوان برچسب در فرآیند آموزش استفاده می شود. همانطور که در بالا توضیح داده شد، هر چند سطر از این ماتریس که برچسب یکسانی دارند، مربوط به یک زیررشته هستند. در این مسئله می خواهیم با داشتن زیررشته، ساختار دوم آن را تخمین بزنیم. ۹/۰ داده ها را برای آموزش و ۱/۰ را برای آزمون کنار بگذارید. در مورد نحوه ای انجام هر کدام از موارد زیر و مشاهداتتان از نتایج توضیح دهید.

۱. ابتدا در پیش پردازش زیررشته های به طول ۷ را نگه دارید و سایر زیررشته ها در نظر نگیرید.
۲. همانطور که با ساختار ویژگی ها آشنا شده اید می دانید که از نوع Categorical هستند و باید به نوع Numerical تبدیل شوند.
۳. با استفاده از پکیج sklearn و روش Random Forest دسته بندی را انجام دهید. در این روش hyperparameterهای عمق درخت و تعداد درخت وجود دارند که باید به درستی انتخاب شوند. مقادیر مناسب برای این پارامترها را انتخاب کنید و علت انتخاب خود را در گزارش بیاورید.
۴. پس از انتخاب مقادیر پارامترها، accuracy این روش را بروی داده های آموزش و آزمون گزارش کنید. همچنین precision و recall را برای داده های آموزش و آزمون برای هر دسته به صورت جداگانه محاسبه کنید. confusion matrix را برای داده های آموزش و آزمون به طور جداگانه محاسبه کرده و در گزارش بیاورید.
۵. با استفاده از پکیج sklearn و روش Multi Layer Perceptron دسته بندی را انجام دهید. همانند قسمت قبل accuracy دسته بندی، precision و recall برای هر دسته و همچنین confusion matrix را در گزارش بیاورید.
۶. (امتیازی) در بخش پیش پردازش برای یکسان سازی تعداد ویژگی ها برای هر نمونه، تنها زیررشته های به طول ۷ را نگه داشتیم و در نتیجه بخشی از داده را دور ریختیم. روشی ارائه دهید تا بتوانیم از زیررشته های با طول کمتر و بیشتر از ۷ هم استفاده کنیم. پس از تشکیل دوباره نمونه ها، بخش ۴ و ۵ را تکرار کنید.