



## درس یادگیری ماشین

مدرس: دکتر سامان هراتی زاده

نیم سال اول ۱۴۰۱-۱۴۰۰

الگوریتم: oneR, Prism, ID3

تحويل آنالین: اعلام خواهد شد

تمرین شماره ۱

مهلت ارسال تمرین: ۱۴۰۰/۹/۳

### هدف: ارائه ی یک مدل پیش بینی برای تشخیص بیماری کرونا

داده ها : مجموعه داده ی پیوست شده شامل علائم و ویژگی هایی است که در احتمال ابتلا به بیماری کرونا موثر هستند. در این مجموعه داده ۵ علامت سرفه، تب، گلودرد، تنگی نفس و سردرد و ویژگی هایی نظیر جنسیت، سن بالای ۶۰ و دلیل تست ( تماس با فرد بیمار، سفر خارجی و غیره) برای ۴۰۰۰ نمونه گزارش شده است. هدف ما پیش بینی نتیجه ی تست کرونا (Corona Result) است. تمامی ویژگی های این مجموعه داده اسمی هستند.

۱. الگوریتم oneR را برای پیش بینی ابتلا به بیماری کرونا پیاده سازی کنید. قواعد ایجاد شده توسط این مدل را همراه با support و confidence آنها در داده آموزش نمایش دهید. نتایج ارزیابی مدل را برای بهترین قاعده ی استخراج شده گزارش کنید.

۲. با پیاده سازی الگوریتم Prism، مدلی برای تشخیص ابتلا به بیماری کرونا ارائه دهید. برای الگوریتم شرط توقف بر اساس تعداد ویژگی اضافه کنید و عدد بهینه ی تعداد ویژگی را برای حل مشکل بیش برآزش را پیدا کنید. نمودار دقت بر حسب تعداد ویژگی ها را رسم کنید. قواعد ایجاد شده توسط این مدل را نمایش دهید. نتایج ارزیابی مدل را گزارش کنید.

۳. الگوریتم ID3 را برای تشخیص بیماری کرونا پیاده سازی کنید. درخت را حداکثر تا ارتفاع ۴ ادامه دهید. قواعد استخراج توسط این الگوریتم را را به همراه نتایج ارزیابی آن را گزارش کنید. (نیازی به ترسیم درخت خروجی به صورت اتوماتیک و گرافیکی توسط سیستم نیست. اما باید امکان پیمایش درخت از ریشه تا برگ وجود داشته باشد.)

۴. (امتیازی) با توجه به اینکه با یک مسئله ی طبقه بندی Imbalanced رو به رو هستیم، به نظر شما Imbalanced بودن دیتا چه مشکلاتی را برای یک مدل پیش بینی ایجاد می کند. راه حل شما برای حل این مسئله چیست؟ راه حل خود را با روش های موجود برای حل این معضل مقایسه کنید.

۵. (امتیازی) الگوریتم درخت تصمیم را با کتابخانه ی scikit-learn پیاده سازی کنید و درخت حاصل از آن را ترسیم نمایید و نتایج آن را با مدل پیاده سازی شده توسط خودتان مقایسه کنید.

آ) تمامی مدل‌ها را با مترهای accuracy, precision, recall و f-measure ارزیابی و میانگین آن‌ها را برحسب 5-Fold-CrossValidation اعلام کنید.<sup>۱</sup>

- $\text{Precision} = \text{TruePositives} / (\text{TruePositives} + \text{FalsePositives})$
- $\text{Recall} = \text{TruePositives} / (\text{TruePositives} + \text{FalseNegatives})$
- $\text{F-measure} = (2 * \text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$

ب) نتایج بدست آمده از پیاده‌سازی این الگوریتم‌ها را در گزارش خود شرح دهید. گزارش شما باید شامل توضیح الگوریتم و شرح مختصری از قسمت‌های مختلف کد شما باشد. دقت داشته باشید که بخشی از نمره شما به گزارش شما تعلق می‌گیرد، بنابراین زمان کافی به آن اختصاص دهید.

- برای پیاده‌سازی می‌توانید از زبان پایتون و کتابخانه‌های pandas, NumPy, SciPy, math و یا کتابخانه‌های مشابه استفاده کنید. دقت داشته‌باشید که استفاده از کتابخانه‌ی آماده برای سؤالاتی که بایستی الگوریتم را پیاده‌سازی کنید مجاز نیست.

- الگوریتم‌ها را خودتان پیاده‌سازی کنید. به کدهای کپی شده از گیت‌هاب و دوستانان نمره‌ای تعلق نمی‌گیرد. در صورت مشاهده‌ی کدهای مشابه بین چند نفر به هیچ یک از افراد نمره‌ای تعلق نخواهد گرفت.

- کد خود را کامنت گذاری کنید. برای کامنت گذاری مناسب و تمیز بودن کد امتیاز مثبت در نظر گرفته خواهد شد.

- در صورت تاخیر تا ۲۴ ساعت ۲۰٪ تا ۴۸ ساعت ۴۵٪ و بیش از ۴۸ ساعت تمامی نمره تمرین را از دست خواهید داد.

- برای تحویل فایل‌های مربوطه را در یک فایل زیپ با نام ML\_StudentID\_FullName\_HW #1 قرار دهید.

- در صورت وجود هرگونه سوال یا ابهام می‌توانید با دستیاران آموزشی از طریق Piazza یا گروه WhatsApp در ارتباط باشید.

---

<sup>۱</sup> پیشنهاد می‌شود برای K-Fold-CrossValidation تابعی بنویسید که برای تمرین‌های بعدی هم قابل استفاده باشد.