

## تمرین کامپیوتری سوم

سیستم‌های عامل - بهار 1400

دانشکده مهندسی برق و کامپیوتر

مسئولان تمرین:

پیش‌زمینه

استاد:

مبینا شاه‌بنده، محمدصابر ابراهیم نژاد

دکتر مهدی کارگهی

### مقدمه



در این تمرین شما به تحلیل داده‌هایی که از مشخصات و قیمت فروش خانه‌ها جمع‌آوری

شده‌است می‌پردازید. در این تمرین به شبیه‌سازی یکی از روش‌های رایج در یادگیری ماشین<sup>1</sup>

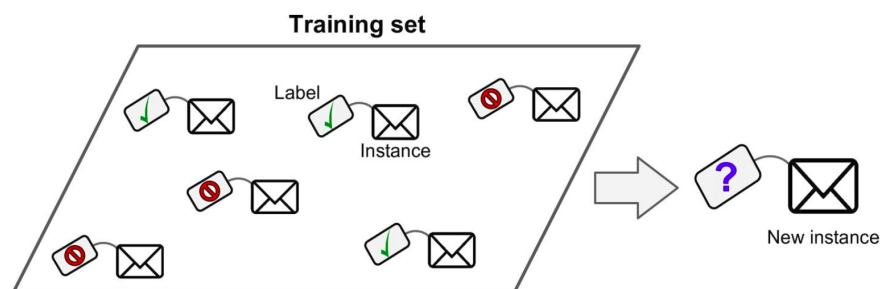
پرداخته می‌شود. به عنوان یکی از شاخه‌های وسیع و پرکاربرد هوش مصنوعی، یادگیری ماشین به تنظیم و اکتشاف شیوه‌ها و

الگوریتم‌هایی می‌پردازد که بر اساس آن‌ها رایانه‌ها و سامانه‌ها توانایی یادگیری و پیش‌بینی پیدا می‌کنند.

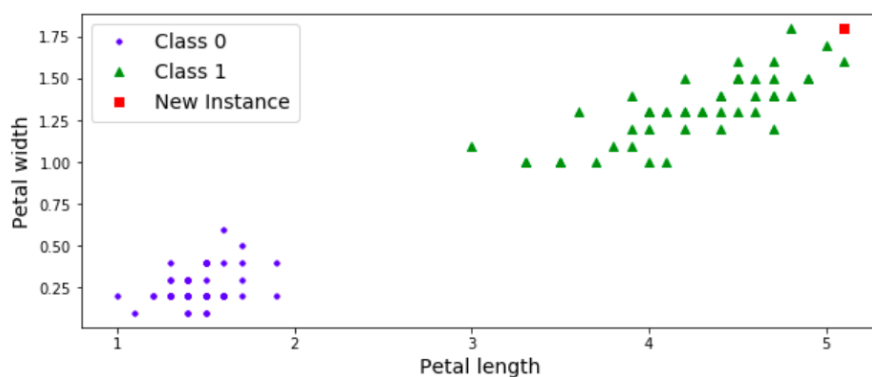
<sup>1</sup> Machine Learning

## طبقه‌بندی<sup>2</sup>

در حوزه یادگیری ماشین، طبقه‌بندی نوعی یادگیری محسوب می‌شود و طبقه‌بندی مسئله شناسایی تعلق مشاهده جدید، به یکی از دسته‌ها بر اساس مجموعه‌ای از مشاهدات می‌باشد که عضویت در دسته‌هایشان مشخص می‌باشد.



برای مثال تصور کنید که می‌خواهید نام یک گل را بر اساس طول و عرض گلبرگ‌های آن تشخیص دهید. بدین منظور لازم است که یک طبقه‌بند<sup>3</sup> برای این منظور آموزش ببیند (توانایی تشخیص نوع گل را پیدا کند) و پس از آن بر اساس ویژگی‌هایی که یک گل را توصیف می‌کند (طول و عرض در این مثال) به طبقه‌بند داده شود. این طبقه‌بند براساس مشاهداتی که در گذشته داشته است (در مرحله آموزش) تعلق این گل را به یکی از دسته‌ها تشخیص می‌دهد.



<sup>2</sup> Classification

<sup>3</sup> Classifier

## طبقه‌بندی بر اساس میانگین و انحراف معیار<sup>4</sup>

در حوزه یادگیری ماشین نمونه‌هایی که قصد پیش‌بینی نوع و یا یک ویژگی آن‌ها وجود دارد، با استفاده از تعدادی ویژگی عددی و قابل اندازه‌گیری در قالب بردار ویژگی<sup>5</sup> توصیف می‌شوند. یکی از روش‌های طبقه‌بندی داده‌ها، استفاده از میانگین و انحراف معیار است که در آن ابتدا میانگین و انحراف معیار یک ویژگی از داده‌ها بدست آمده و سپس بر اساس آن، بازه تخمینی ساخته می‌شود. این بازه تخمین به صورت زیر بدست می‌آید:

$$(mean - std, mean + std)$$

برای مثال تصور کنید که طبقه‌بند توانایی تشخیص دو نوع گل از یکدیگر را بر اساس طول گلبرگ آنها دارد. حال فرض کنید گل‌هایی که در کلاس  $\infty$  هستند، دارای میانگین طول گلبرگ 0.7 و انحراف معیار طول گلبرگ 0.1 هستند و این مقادیر برای گل‌های کلاس  $\beta$  به ترتیب برابر با 0.3 و 0.2 باشد. حال این طبقه‌بند با متغیرهای آماری مذکور، قصد تشخیص نمونه‌ای که دارای بردار ویژگی زیر است را دارد:

<i>Length</i>	<i>Width</i>
0.78	0.15

ستون‌های *Length* و *Width* همانطور که از نام آن‌ها برمی‌آید معرف طول و عرض گلبرگ مربوط به گل‌ها است. با توجه به اینکه طبقه‌بندی بر اساس طول گلبرگ است، مقدار این ویژگی را برای این نمونه بررسی می‌کنیم. بازه تخمینی برای کلاس  $\infty$ ، برابر است با:  $(0.1, 0.7 + 0.1, 0.7 - 0.1)$ . با توجه به اینکه طول گلبرگ این نمونه برابر با 0.78 است، پس در بازه تخمین کلاس  $\infty$  حضور دارد؛ بنابراین کلاس این نمونه کلاس  $\infty$  تعیین می‌شود. در صورتی که مقدار ویژگی مورد نظر نمونه‌ای

<sup>4</sup> Mean-Standard Deviation Classification

<sup>5</sup> Feature Vector

در هیچ یک از دو بازه تخمین دو کلاس (یا چند کلاس) حضور نداشته باشد، یکی از کلاس ها به صورت پیش فرض در نظر گرفته می شود.

## مجموعه داده<sup>6</sup>



مجموعه داده‌ای که در این تمرین به شما داده شده‌است در قالب CSV<sup>7</sup> است. CSV نام یک قالب برای پرونده‌های متنی است که در آن مقادیر با استفاده از نماد کاما (,) از یکدیگر جدا می‌شوند. این قالب یکی از روش‌های پرطرفدار برای تبادل اطلاعات است.



## اطلاعات خانه ها

اطلاعات مربوط به خانه ها در پرونده dataset.csv در اختیار شما قرار داده شده‌است. در ادامه درباره‌ی هر ویژگی و نوع داده‌ی<sup>8</sup> مربوط به آن، توضیح مختصری آمده‌است.

نام ویژگی	توضیح	نوع داده
LotArea	مساحت خانه در واحد فیت مربع	عدد صحیح

<sup>6</sup> Dataset

<sup>7</sup> Comma-Separated Values

<sup>8</sup> Data Type

OverallQual	کیفیت کلی متريال های به کار رفته در ساخت خانه	عدد صحيح (1 تا 10)
OverallCond	رتبه بندی وضعیت کلی خانه	عدد صحيح (1 تا 10)
YearBuilt	سال ساخت خانه	عدد صحيح
TotalBsmtSF	مساحت فضای زیرزمین در واحد فیت مربع	عدد صحيح
GrLivArea	مساحت فضای روزمین در واحد فیت مربع	عدد صحيح
GarageCars	ظرفیت پارکینگ خانه در واحد تعداد خودرو	عدد صحيح
GarageArea	مساحت پارکینگ در واحد فیت مربع	عدد صحيح
SalePrice	قیمت فروش خانه در واحد دلار	عدد صحيح

## منابع

<https://www.kaggle.com/lespin/house-prices-dataset>

<https://www.vectorstock.com/>