

الگوریتم RESSEL

پروژه پایان ترم کلاس یادگیری ماشین: توصیف مدل و شیوه پیاده سازی

استاد: جناب آقای دکتر شاقوزی

محقق: سجاد صابری

معرفی روش RESSEL

روش RESSEL (Reliable Semi-Supervised Ensemble Learning) ، یک نوع روش یادگیری نیمه نظارتی است که ترکیبی از یادگیری نیمه نظارتی و یادگیری شورایی (Ensemble Learning) را به کار می بندد. در این روش، از داده های برچسب دار و بدون برچسب استفاده می شود تا با بهره گیری از داده های بدون برچسب، مدل یادگیری نظارتی بهبود یابد. هدف اصلی RESSEL این است که با استفاده از مجموعه ای از طبقه بندی ها، از نمونه های بدون برچسب به شکل موثرتری برای بهبود دقت پیش بینی ها استفاده شود.

هدف اصلی RESSEL

در این روش، مجموعه کوچکی از داده های برچسب دار و یک مجموعه بزرگ از داده های بدون برچسب انتخاب می شوند. این کار با خودآموزی (Self-Training) انجام می شود؛ به این معنا که مدل پس از یادگیری اولیه، نمونه هایی از داده های بدون برچسب که با اعتماد بیشتری طبقه بندی می شوند را به عنوان داده های برچسب دار به مجموعه آموزشی اضافه می کند. هدف این مسیر، بهبود پیش بینی و دسته بندی صحیح داده ها با بازده بالاتر در محیط واقعی است.

پارامترهای در نظر گرفته شده

در کد پیاده‌سازی RESSEL به زبان پایتون، چندین پارامتر کلیدی در نظر گرفته شده است که برای کنترل فرایند یادگیری مهم هستند:

1. **ensemble_size** تعداد مدل‌های جمعی : تعداد مدل‌های جمعی یا طبقه‌بندهایی که قرار است در الگوریتم

جمعی استفاده شوند. تعداد بیشتر باعث ایجاد تنوع بیشتر و بهبود دقت کلی می‌شود.

2. **Iterations** تعداد تکرارهای خودآموزی: تعداد تکرارهایی که در آن داده‌های بدون برچسب با اطمینان

بالا به مجموعه داده‌های برچسب‌دار اضافه می‌شوند.

3. **batch_size** اندازه دسته داده‌های بدون برچسب: تعداد نمونه‌هایی که در هر تکرار از داده‌های

بدون برچسب با اعتماد بالا انتخاب شده و به داده‌های برچسب‌دار اضافه می‌شوند.

4. **unlabeled_fraction** نسبت داده‌های بدون برچسب: نسبت داده‌های بدون برچسب که در هر تکرار

برای خودآموزی در نظر گرفته می‌شود.

5. **oob_threshold** آستانه خطای (OOB): آستانه‌ای که برای توقف زود هنگام (Early Stopping)

براساس خطای OOB (خطای خارج از کیسه) استفاده می‌شود. اگر خطای OOB به کمتر از این مقدار برسد،

مدل متوقف می‌شود تا از بیش‌برازش جلوگیری شود.

مراحل پیاده‌سازی در پایتون

1. **بارگذاری داده‌ها**: ابتدا داده‌های **Spambase** که شامل ویژگی‌های ایمیل‌ها برای تشخیص هرزنامه است،

بارگذاری می‌شود. داده‌ها به دو مجموعه **داده‌های برچسب‌دار** (برای آموزش اولیه) و **داده‌های بدون برچسب**

(برای خودآموزی) تقسیم می‌شوند.

2. **پیش‌پردازش داده‌ها**: داده‌ها با استفاده از ابزار **MinMaxScaler** مقیاس‌بندی شده و از

SelectFromModel برای انتخاب ویژگی‌های مهم استفاده می‌شود. مدل جنگل تصادفی

(RandomForest) برای انتخاب ویژگی‌ها به کار می‌رود.

3. **تعریف مدل پایه: DecisionTreeClassifier** به عنوان مدل پایه استفاده می شود. این طبقه بند به عنوان

واحد اصلی در روش جمعی به کار می رود.

4. **استفاده از BaggingClassifier :** برای ایجاد مجموعه ای از مدل ها، از **BaggingClassifier**

استفاده می شود که چندین طبقه بند تصمیم گیری را روی داده های آموزشی بوت استرپ شده آموزش می دهد و از OOB برای ارزیابی مدل استفاده می کند.

5. **خودآموزی:** در هر تکرار، مدل با استفاده از داده های برچسب دار اولیه آموزش می بیند و سپس از داده های

بدون برچسب استفاده می کند تا نمونه های با ضریب اعتماد بالا را به عنوان داده های برچسب دار به مجموعه آموزشی اضافه کند.

6. **محاسبه خطای OOB :** پس از هر تکرار، خطای OOB محاسبه می شود. این خطا برای ارزیابی عملکرد مدل

روی داده هایی که در نمونه برداری بوت استرپ حضور نداشته اند به کار می رود.

7. **توقف زودهنگام:** اگر خطای OOB به زیر مقدار تعیین شده برسد، فرایند خودآموزی متوقف می شود.

در نهایت، این کد به داده های Spambase اعمال می شود که شامل ویژگی های ایمیل هاست. این کد از روش خودآموزی

برای افزایش دقت طبقه بندهای جمعی استفاده می کند و از داده های بدون برچسب بهره می برد تا مدل بهتری ایجاد کند. در

قسمت انتهایی کد، فرایند ارزیابی مدل انجام می شود تا عملکرد نهایی مدل بر روی داده های تست شده بررسی گردد.