

بخش دو تمرین دو

در این تمرین با کمک ۳ فیچری که در صورت سوال گفته شد، الگوریتم رگرسیون خطی را اجرا می‌کنیم. چون فیچر Heating type کتگوریکال است، از one hot encoding استفاده می‌کنیم که منجر می‌شود تعداد فیچرهای ورودی به الگوریتم برابر ۱۵ بشود. همچنین فیچر Living space نیز به عنوان لیبیل‌ها در نظر گرفته می‌شود.

داده‌ها را به نسبت ۸ به ۲ برای ترین و تست تقسیم می‌کنیم و قبل از این که داده‌ها را به الگوریتم بدهیم، داده‌ها را با کمک standard scaler نرمال‌سازی می‌کنیم.

مقدار mse بدون استفاده از پکیج (استفاده از الگوریتم پیاده‌سازی شده):

MSE on the test set: 1125.6715078175303

مقدار mse با استفاده از پکیج:

MSE on the test set: 1125.7051500327655

تفاوت کم بین مقدار mse الگوریتم پیاده‌سازی شده و پکیج استفاده شده نشان‌دهنده این است که الگوریتم به خوبی پیاده‌سازی شده است.

همان فیچرهایی که به مدل رگرسیون خطی به عنوان ورودی داده شده بود را به عنوان ورودی به Ridge regression و Lasso regression نیز می‌دهیم.

مقدار mse با استفاده از پکیج Ridge regression :

10 Fold:Mean MSE: 87382.327

مقدار mse با استفاده از پکیج Lasso regression :

10 Fold:Mean MSE: 86213.890

رگرسیون ridge و lasso به دلیل وجود ترم محدودکننده در شرط بهینه‌سازی اجازه نمی‌دهند که وزن‌ها از حدی بیشتر شوند. به این دلیل روی داده‌های غیرخطی خطا بالا می‌رود و مقدار وزن‌ها از حدی بیشتر نمی‌شود که این منجر به افزایش مقدار mse می‌شود و در این دیتاست مدل underfit شده است. اما مزیت مدل رگرسیون ridge و lasso به تفسیرپذیری بالای آن‌ها می‌باشد.

در واقع وجود ترم محدودکننده در شرط بهینه‌سازی موجب این می‌شود که خیلی از فیچر‌ها صفر می‌شوند و می‌تواند فیچرهای خوب مدل را به ما نشان دهد. به طور کلی یک trade off بین پیچیدگی مدل و تفسیرپذیری مدل وجود دارد و هرچقدر مدل پیچیده‌تر باشد معمولاً خطای کمتری دارد و مدل‌ها با تفسیرپذیری بالا خطای زیادی دارند.