

JiXi 绩溪

Package for Tuning (1st Generation)

11/04/2023

Background

The purpose of this package is to provide a sophisticated framework for Brute Force tuning. *The Criteria of First Generation Tuning is to use Brute Force to train every specified discrete combination.*

The package takes in X and y data for train, validate and test as DataFrame, as well as a dictionary of {hyperparameters name -> string: hyperparameter values as a list}, and autogenerates all combinations of these hyperparameters to be tuned.

JiXi allows tuning combinations to proceed 1) in the order of how they would be, if generated by nested loops, and also 2) randomly shuffled. A third and fourth sophisticated method of tuning order whereby 3) first tune the centre combination of the discrete parameter field, then tune the 'layer' immediately touching that, and then the 'layer' immediately touching the previous layer and 4) to first tune all combinations lying on the horizontal and also the diagonal relative to the centre combination, before going layer by layer.

The advantage of 3) and 4) is that, if used alongside YiLong, can guarantee that each hyperparameter's values will quickly see at least one combination, which allows the Data Scientist to discard certain parameter values clearly producing bad results and save time.

Class

<u>Class</u>	<u>Purpose</u>
JiXi	Object that performs brute force tuning with four different order choices

Methods:

<u>Methods</u>	<u>Purpose</u>
<i>JiXi()</i>	Initialisation
<code>read_in_data(train_x, train_y, val_x, val_y, test_x, test_y)</code>	<p>Read in Train Test Split data</p> <p>Parameters:</p> <p><code>train_x</code> – <code>pd.DataFrame</code> <code>train_y</code> - <code>pd.Series</code> <code>val_x</code> - <code>pd.DataFrame</code> <code>val_y</code> - <code>pd.Series</code> <code>test_x</code> - <code>pd.DataFrame</code> <code>test_y</code> – <code>pd.Series</code></p>
<code>read_in_model(model, type)</code>	<p>Read in the underlying model class that we want to tune to get optimal parameters for</p> <p>Parameters:</p> <p><code>model</code> – any model class that allows <code>.fit()</code> and <code>.predict()</code></p> <p><code>type</code> – str – either “Classification” or “Regression”</p>
<code>set_hyperparameters(parameter_choices)</code>	<p>Read in the different values of each hyperparameters we want to try. Function will automatically generate each combination</p> <p>Parameters:</p> <p><code>parameter_choices</code> – dict of str:list – str is hyperparameter name (strictly as defined in model class), and list is sorted values of hyperparameter which we want to try out.</p>

<code>set_non_tuneable_hyperparameters(non_tuneable_hyperparameter_choice)</code>	<p>Reads in values for non-tuneable hyperparameters (i.e. doesn't need to clog up the tuning output csv)</p> <p>Parameters: non_tuneable_hyperparameter_choices – dict of str:int</p>
<code>set_features(ningxiang_output)</code>	<p>Reads in feature combinations for tuning</p> <p>Parameters: ningxiang_output – dict of tuple:float</p>
<code>set_tuning_result_saving_address(address)</code>	<p>Set saving address for tuning output csv</p> <p>Parameters: address – str – does not need to include '.csv'</p>
<code>change_tuning_style(type, seed = None, outer_most_layer = 2, randomise = True)</code>	<p>Set which type of tuning order to use.</p> <p>'a': as if nested (according to order of dictionary input to set_hyperparameters())</p> <p>'b': (reset to 'a') before random shuffle using inputted seed, or default seed 19421221</p> <p>'c': (reset to a) before setting to layer by layer order</p> <p>'d': (reset to a) (reset to c) before setting to diag-hor -> layer by layer. Automatically randomised by default seed</p> <p>Parameters: type – str – 'a' or 'b' or 'c' or 'd'</p>

	<p>seed – int – for ‘b’ and ‘c’</p> <p>outer_most_layer – the outer most layer for ‘c’ and ‘d’ to actually order for, before remaining are all random</p> <p>randomise – bool – whether or not to randomise ‘c’</p>
tune(key_stats_only = <code>False</code>)	<p>Begin tuning process</p> <p>If key_stats_only = True then don’t calculate non important stats</p> <p>Parameters:</p> <p>key_stats_only – bool</p>
tune_parallel(part, splits, key_stats_only = <code>False</code>)	<p>Begin tuning process, splitting all combinations into <i>splits</i> parts and tune the <i>part</i>-th part.</p> <p>If key_stats_only = True then don’t calculate non important stats</p> <p>Parameters:</p> <p>key_stats_only – bool</p>
read_in_tuning_result_df(address)	<p>Read in existing DataFrame from .csv consisting of tuning result.</p> <p>Automatically populates result array and checked array if csv columns match parameter choices</p> <p>Parameters:</p> <p>address – str – include ‘.csv’</p>
set_tuning_best_model_saving_address(address)	<p>Set address for exporting best model as a pickle</p>

	Parameters: address – str - – does not need to include ‘.pickle’
view_best_combo_and_score()	View the current best combination and its validation score

Objects:

<u>Objects</u>	<u>Purpose</u>
train_x	DataFrame
train_y	Series
val_x	DataFrame
val_y	Series
test_x	DataFrame
test_y	Series
tuning_result	DataFrame
model	model class
parameter_choices	Dictionary -str:list – str is hyperparameter name (strictly as defined in model class), and list is sorted values of hyperparameter which we want to try out.
hyperparameters	list
feature_n_ningxiang_score_dict	Dictionary -str:float – str is hyperparameter name (strictly as defined in model class), and float is its NingXiang score
non_tuneable_parameter_choices	Dictionary -str:str/float/int - str is hyperparameter name (strictly as defined in model class), and values are valid hyperparameter values for model
checked	np.array
result	np.array
tuning_result_saving_address	str
best_model_saving_address	str

best_score = -np.inf	int
best_combo	list
best_clf	model object
clf_type	str – ‘Regression’ or ‘Classification’
combos	List of lists
n_items	list - denoting how many values in each hyperparameter dimensions
<pre> regression_extra_output_columns = ['Train r2', 'Val r2', 'Test r2', 'Train RMSE', 'Val RMSE', 'Test RMSE', 'Train MAPE', 'Val MAPE', 'Test MAPE', 'Time'] </pre>	list (pre-setted)
<pre> classification_extra_output_columns = ['Train accu', 'Val accu', 'Test accu', 'Train balanced_accu', 'Val balanced_accu', 'Test balanced_accu', 'Train f1', 'Val f1', 'Test f1', 'Train precision', 'Val precision', 'Test precision', 'Train recall', 'Val recall', 'Test recall', 'Time'] </pre>	list (pre-setted)

Dependencies

pandas

numpy

sklearn