

Methods for feature importance and feature interaction detection

Step 1: Feature importance detection

Here SubmodularPick in Lime model is used to save a representative sample of explanation objects for global interpretation.

Codes:

```
lime.lime_tabular.LimeTabularExplainer(training_data, mode='classification',  
training_labels=None, feature_names=None, categorical_features=None,  
categorical_names=None, kernel_width=None, kernel=None, verbose=False,  
class_names=None, feature_selection='auto', discretize_continuous=True,  
discretizer='quartile', sample_around_instance=False, random_state=None,  
training_data_stats=None)
```

Description:

Explains predictions on tabular data. For numerical features, perturb them by sampling from a Normal(0,1) and doing the inverse operation of mean-centering and scaling, according to the means and stds in the training data. For categorical features, perturb by sampling according to the training distribution, and making a binary feature that is 1 when the value is the same as the instance being explained.

Arguments:

<i>training_data</i>	numpy 2d array
<i>mode</i>	“classification” or “regression”
<i>training_labels</i>	labels for training data. Not required, but may be used by discretizer.
<i>feature_names</i>	list of names (strings) corresponding to the columns in the training data.
<i>categorical_features</i>	list of indices (ints) corresponding to the categorical columns. Everything else will be considered continuous. Values in these columns MUST be integers.
<i>categorical_names</i>	map from int to list of names, where categorical_names[x][y] represents the name of the yth value of column x.
<i>kernel_width</i>	kernel width for the exponential kernel. If None, defaults to sqrt (number of columns) * 0.75
<i>kernel</i>	similarity kernel that takes euclidean distances and kernel width as input and outputs weights in (0,1). If None, defaults to an exponential kernel.
<i>verbose</i>	if true, print local prediction values from linear model
<i>class_names</i>	list of class names, ordered according to whatever the classifier is using. If not present, class names will be ‘0’, ‘1’, ...
<i>feature_selection</i>	feature selection method. can be ‘forward_selection’, ‘highest_weights’, ‘lasso_path’, ‘none’ or ‘auto’. See function ‘explain_instance_with_data’ in lime_base.py for details on what each of the options does.
<i>discretize_continuous</i>	if True, all non-categorical features will be discretized into quartiles.

<i>Discretizer</i>	only matters if <code>discretize_continuous</code> is True. Options are 'quartile', 'decile' or 'entropy'
<i>sample_around_instance</i>	if True, will sample continuous features in perturbed samples from a normal centered at the instance being explained. Otherwise, the normal is centered on the mean of the feature data.
<i>random_state</i>	an integer or <code>numpy.RandomState</code> that will be used to generate random numbers. If None, the random state will be initialized using the internal numpy seed.
<i>training_data_stats</i>	a dict object having the details of training data statistics. If None, training data information will be used, only matters if <code>discretize_continuous</code> is True. Must have the following keys: 'means', 'mins', 'maxs', 'stds', 'feature_values', 'feature_frequencies'

Usage recommendations:

Most parameters can follow the defaults. 'training_data', 'feature_names', 'class_names' should be initialized. 'random_state' also had better to be fixed for stability.

e.g. `explainer = lime.lime_tabular.LimeTabularExplainer(test_data, feature_names=feature_names, class_names=class_names, mode = 'regression', random_state=20)`

`lime.submodular_pick.SubmodularPick(explainer, data, predict_fn, method='sample', sample_size=1000, num_exps_desired=5, num_features=10, **kwargs)`

Description:

A representative sample of explanation objects using SP-LIME is saved, as well as saving all generated explanations. Firstly, a collection of candidate explanations is generated. From these candidates, `num_exps_desired` are chosen using submodular pick.

Arguments:

<i>data</i>	a numpy array where each row is a single input into <code>predict_fn</code>
<i>predict_fn</i>	prediction function. For classifiers, this should be a function that takes a numpy array and outputs prediction probabilities. For regressors, this takes a numpy array and returns the predictions. For <code>ScikitClassifiers</code> , this is <code>classifier.predict_proba()</code> . For <code>ScikitRegressors</code> , this is <code>regressor.predict()</code> . The prediction function needs to work on multiple feature vectors (the vectors randomly perturbed from the <code>data_row</code>).
<i>method</i>	The method to use to generate candidate explanations <code>method == 'sample'</code> will sample the data uniformly at random. The sample size is given by <code>sample_size</code> . Otherwise if <code>method == 'full'</code> then explanations will be generated for the entire data.
<i>sample_size</i>	The number of instances to explain if <code>method == 'sample'</code>
<i>num_exps_desired</i>	The number of explanation objects returned
<i>num_features</i>	maximum number of features present in explanation

Returns:

<i>sp_explanations</i>	A list of explanation objects that has a high coverage explanation: All the candidate explanations saved for potential future use.
------------------------	---

Usage recommendations:

‘explainer’ is the output of `lime.lime_tabular.LimeTabularExplainer`. If training data is not too large, `method = ‘full’` is better. If data size is very large, part of the instances can be selected for training. For ‘num_exps_desired’, there is no special recommendation and no less than the default number is OK. For ‘num_features’, the best situation is that all the features or at least 95% of them are selected.

e.g.

```
sp = submodular_pick.SubmodularPick(explainer, data, predict_fn=model.predict, method='full', num_exps_desired=5, num_features=8)
```

#visualization

```
[exp.as_pyplot_figure(label=exp.available_labels()[0]) for exp in sp.sp_explanations]
```

feature_importance_table(*sbmodular*, *feature_names=None*, *threshold=0.1*)

Description:

To visualize the feature importance in a table form. ‘threshold’ is to determine whether the feature is important. The table has 3 columns: ‘features’, ‘weights’, ‘counts’. ‘weights’ shows the importance of each feature. ‘counts’ shows the number of features recognized as important ones among all picked data instances.

Arguments:

<i>sbmodular</i>	submodular pick in lime model
<i>feature_names</i>	the names of features
<i>threshold</i>	the threshold to pick up important features

Returns:

<i>sorted_count</i>	the table showing the feature importance in ascending order
---------------------	---

Usage recommendations:

‘threshold’ should be set carefully. The recommendation is a value that is a little larger than the least important feature according to `sp.sp_explanations`.

e.g.

```
table = feature_importance_table(sp.sp_explanations, feature_names=feature_names, threshold=1)
```

Step 2: Feature interactions detection

Here NID or GradientNID algorithm is used to detect feature interactions.

application_NID(*inputs*, *targets*, *std_scale=False*, *detector="GradientNID"*, *arch=[256, 128, 64]*, *batch_size=100*, *device=torch.device("cpu")*, *weight_samples=False*, *add_linear=False*, *l1_const=None*, *grad_gpu=-1*, *seed=42*, *number=20*, *feature_names=None*, *decimal=4*)

Description:

To train the dataset with ‘NID’ and ‘GradientNID’ to get the feature interaction information.

Arguments:

<i>inputs</i>	numpy 2d array. It can be the whole dataset or part of them
<i>targets</i>	the target data corresponding to inputs
<i>std_scale</i>	standardizing the processed data by sklearn.preprocessing.StandardScaler()
<i>detector</i>	"NID" or "GradientNID"
<i>arch</i>	architecture of the neural networks, default is [256, 128,64]
<i>batch_size</i>	the batch size of the neural network, default is 100
<i>device</i>	the object of the device to which torch.Tensor is allocated, default is torch.device("cpu")
<i>weight_samples</i>	weight sampling with gaussian kernel, default is False
<i>add_linear</i>	adding the full connection layer, default is False
<i>l1_const</i>	the const parameter for l1 regularization, default is None
<i>grad_gpu</i>	if it is -1, the device is cpu, else the device is gpu. Default is -1
<i>seed</i>	the seed for random state, default is 42
<i>number</i>	the number of output interactions, default is 20
<i>feature_names</i>	the names of the features
<i>decimal</i>	the number of decimal places of the result, default is 4

Returns:

<i>interactions</i>	the feature interaction information
<i>mlp_loss</i>	the loss of NN

Usage recommendations:

Most settings can follow the defaults. 'detector' is recommended as 'GradientNID' because the high-order interactions are always the noises. Here 'number' is important and the recommendation is to set it as more than the number of features.

e.g.

```
interactions, mlp_loss = application_NID(inputs=inputs, targets=targets,  
feature_names=feature_names, detector="GradientNID", number=10)
```

Step 3: Combination and retraining

Here all explanations above are combined to create a new dataset and then it will be retrained to improve the performance.

combination(interactions, training_data, test_data, feature_names, delete=None, num=10)

Description:

To create new features with feature interactions by multiplication and combine them with the original dataset. If there exist some unimportant features, they should be deleted.

Arguments:

<i>interactions</i>	the feature interactions interpretation
<i>training_data</i>	the original training data
<i>test_data</i>	the original test data
<i>feature_names</i>	the names of the features

<i>delete</i>	the unimportant features to be deleted, default is None. If it is not None, for example, the 1st and 2nd features are deleted, delete=[1,2].
<i>num</i>	the number of selected features

Returns:

<i>train_new</i>	the new training data
<i>test_new</i>	the new test data
<i>feature_names_new</i>	the new feature names

Usage recommendations:

‘delete’ needs to be treated with caution. Only when the feature has much lower weight in lime explanations, it can be deleted. According to ‘num’, it had better follow the number of original features.

e.g.

train_new, test_new, feature_names_new = combination(interactions, train_data, test_data, feature_names, delete=[1], num=11)

retraining(table, model, training_data, training_target, test_data, test_target, task='regression', method='rmse and r2', multi_class='macro', rmse=None, r2=None, accuracy=None, f1=None, n=5):

Description:

After the new dataset is created, lime model will be utilized again to detect the feature importance.

Arguments:

<i>table</i>	the feature importance table
<i>model</i>	the model used before
<i>training_data</i>	the new training data
<i>training_target</i>	the targets of training data
<i>test_data</i>	the new test data
<i>test_target</i>	the targets of the test data
<i>task</i>	'regression', 'binary_classification' or 'multi_classification'
<i>method</i>	'rmse and r2', 'rmse' or 'r2' in regression task; 'acc and f1', 'acc' or 'f1' in classification tasks
<i>multi_class</i>	'micro', 'macro', 'samples' or 'weighted', default is 'macro'
<i>rmse</i>	root mean squared error
<i>r2</i>	r2 score
<i>accuracy</i>	accuracy score
<i>f1</i>	f1 score
<i>n</i>	the number of candidate features to be deleted

Returns:

<i>rmse_best</i>	the optimized rmse score
<i>r2_best</i>	the optimized r2 score
<i>acc_best</i>	the optimized accuracy score
<i>f1_best</i>	the optimized f1 score
<i>feature_deleted</i>	the deleted unimportant features

Usage recommendations:

For 'method', select the suitable method according to the task. For 'n', the number of candidate features to be deleted should be less than half of the feature numbers.

e.g.

```
rmse_best, r2_best, feature_deleted = retraining(table, model, training_data, training_target,
test_data, test_target, task='regression', method='rmse and r2', rmse=rmse, r2=r2, n=6)
```

Example:

Dataset; IKV_dataset - Young's Modulus

Codes:

```
train, test = train_test_split(dataset, train_size=0.8, random_state=20)
# extracting the training data, test data and the feature names
train_data = train.iloc[:, 1:8].values
train_target = train.iloc[:, 8:].values
test_data = test.iloc[:, 1:8].values
test_target = test.iloc[:, 8:].values
feature_names=list(dataset.columns)[1:8]
target_names = list(dataset.columns)[8:]

# random forest model to predict the young's modulus
model_1 = RandomForestRegressor(n_estimators=300, random_state=20)
model_1.fit(train_data, train_target[:,0])
Young_pred = model_1.predict(test_data)
# RMSE and R2 score
rmse = np.sqrt(mean_squared_error(Young_pred, test_target[:,0]))
r2 = r2_score(Young_pred, test_target[:,0])

# local explanation in LIME algorithm
explainer = lime.lime_tabular.LimeTabularExplainer(test_data,
feature_names=feature_names, class_names=target_names[0], verbose = False, mode =
'regression', random_state=20)
# submodular pick based on the local explainer above to create the global explainer
sp = submodular_pick.SubmodularPick(explainer, test_data, model_1.predict, method = 'full')
# visualizing the feature importance for each picked data instance
```

```

[exp.as_pyplot_figure() for exp in sp.sp_explanations]

a = feature_importance_table(sp.sp_explanations, feature_names=feature_names,
threshold=3)

interactions, mlp_loss = application_NID(dataset.iloc[:, 1:8].values, dataset.iloc[:, 8].values,
feature_names=feature_names, detector="GradientNID", number=20)

# create new datasets with feature interactions and delete unimportant features

train_new, test_new, feature_names_new = combination(interactions, train_data, test_data,
feature_names, num=11)

# reuse the lime model for the new dataset

model_12 = RandomForestRegressor(n_estimators=300, random_state=20)

model_12.fit(train_new, train_target[:,0])

explainer_12 = lime.lime_tabular.LimeTabularExplainer(test_new,
feature_names=feature_names_new, class_names=target_names[0], verbose = False, mode =
'regression', random_state=20)

# submodular pick for the new dataset

sp_new = submodular_pick.SubmodularPick(explainer_12, test_new, model_12.predict,
num_features=len(feature_names_new), method = 'full', num_exps_desired=10)

# visualizing

[exp.as_pyplot_figure() for exp in sp_new.sp_explanations]

a12 = feature_importance_table(sp_new.sp_explanations,
feature_names=feature_names_new, threshold=3)

#retraining to optimize both rmse and r2

rmse_best, r2_best, feature_deleted = retraining(table=a12, model=model_12,
training_data=train_new, training_target=train_target[:,0], test_data=test_new,
test_target=test_target[:,0], task='regression', method='rmse and r2', rmse=rmse, r2=r2, n=8)

#retraining to optimize only rmse

rmse_best, feature_deleted = retraining(table=a12, model=model_12,
training_data=train_new, training_target=train_target[:,0], test_data=test_new,
test_target=test_target[:,0], task='regression', method='rmse', rmse=rmse, r2=r2, n=8)

#retraining to optimize only r2

r2_best, feature_deleted = retraining(table=a12, model=model_12, training_data=train_new,
training_target=train_target[:,0], test_data=test_new, test_target=test_target[:,0],
task='regression', method='r2', rmse=rmse, r2=r2, n=8)

```