

EXPLAN: Explaining Black-box Classifiers using Adaptive Neighborhood Generation

Peyman Rasouli, Ingrid Chieh Yu

Introduction

EXPLAN is a local model-agnostic explanation method applicable to tabular data classification problems. It is a module-based algorithm consisted of dense data generation, representative data selection, data balancing, and rule-based interpretable model.

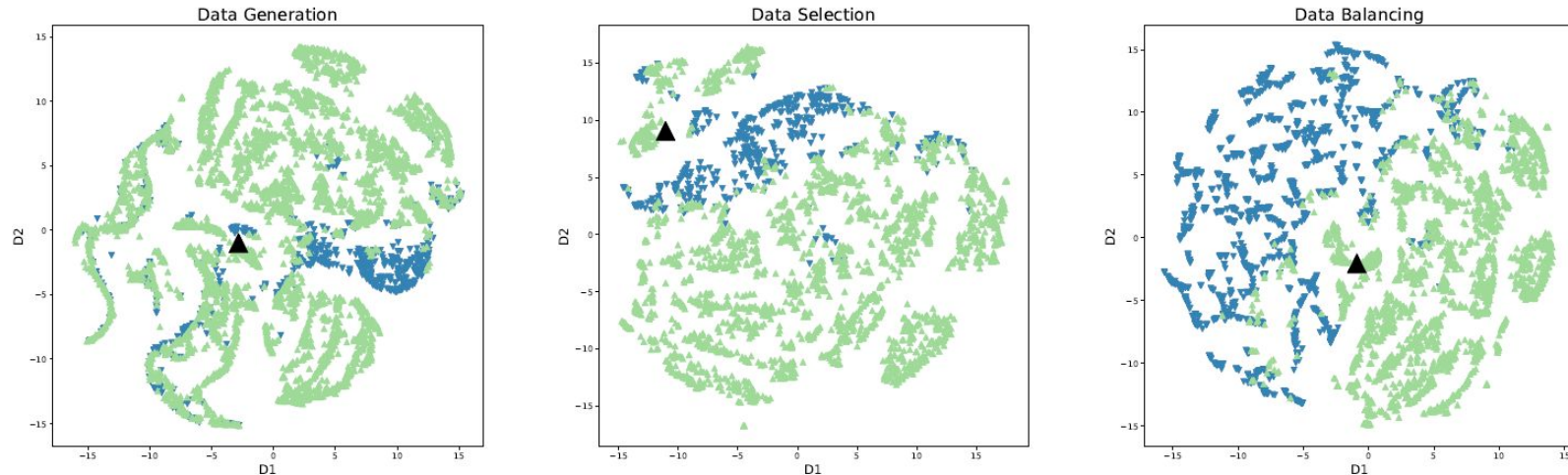


Fig. 2. Visualization of EXPLAN's neighborhood construction process.

Algorithm

Algorithm 1 EXPLAN Explanation Method

Input: $\{x, f, \mathcal{D}, \mathcal{N}, \tau\}$

/ x : instance to explain, f : black-box model, \mathcal{D} : distribution of training data, \mathcal{N} : # initial neighborhood samples, τ : # minimum samples per class */*

Output: $\{\mathcal{C}, e\}$

/ \mathcal{C} : interpretable model, e : explanation of x */*

- 1: **function** EXPLAN($x, f, \mathcal{D}, \mathcal{N}, \tau$)
 - 2: $\mathcal{Z} \leftarrow \text{DATAGENERATION}(x, f, \mathcal{D}, \mathcal{N})$
 - 3: $\mathcal{Z}' \leftarrow \text{DATA SELECTION}(x, f, \mathcal{Z}, \tau)$
 - 4: $\mathcal{X} \leftarrow \text{DATA BALANCING}(f, \mathcal{Z}')$
 - 5: $\mathcal{C}, e \leftarrow \text{INTERPRETABLE MODEL}(x, f, \mathcal{X})$
 - 6: **return** \mathcal{C}, e
-

Algorithm 2 Dense Data Generation

```
1: procedure DATAGENERATION( $x, f, \mathcal{D}, \mathcal{N}$ )
2:   procedure RANDOMDATAGENERATION( $\mathcal{D}, \mathcal{N}$ )
3:      $\mathcal{S} \leftarrow \text{DataSampling}(\mathcal{D}, \mathcal{N})$ 
4:     return  $\mathcal{S}$  /* random data points */
5:   procedure SURROGATEMODELCONSTRUCTION( $f, \mathcal{S}$ )
6:      $\mathcal{T} \leftarrow \text{RandomForestConstructor}(\mathcal{S}, f(\mathcal{S}))$ 
7:     return  $\mathcal{T}$  /* RF surrogate model */
8:   procedure CONTRIBUTIONEXTRACTION( $x, \mathcal{S}, \mathcal{T}$ )
9:      $\mathcal{V}(x) \leftarrow \text{TreeInterpreter}(\mathcal{T}, x)$ 
10:    for all  $s \in \mathcal{S}$  do
11:       $\mathcal{V}(s) \leftarrow \text{TreeInterpreter}(\mathcal{T}, s)$ 
12:    return  $\mathcal{V}$  /* feature importance */
13:   procedure SAMPLEMANIPULATION( $x, \mathcal{S}, \mathcal{T}, \mathcal{V}$ )
14:      $l_x \leftarrow \mathcal{T}(x)$ 
15:      $\mathcal{Z} \leftarrow \{\}$ 
16:     for all  $s \in \mathcal{S}$  do
17:        $l_s \leftarrow \mathcal{T}(s)$ 
18:       for  $j \leftarrow 1, \mathcal{F}$  do /*  $\mathcal{F}$ : feature dimension */
19:         if ( $s_j \neq x_j$ ) then
20:           if ( $\mathcal{V}_{s_j}^{l_x} = \mathcal{V}_{x_j}^{l_x}$ )  $\wedge$  ( $\mathcal{V}_{s_j}^{l_s} = \mathcal{V}_{x_j}^{l_s}$ ) then
21:              $s_j \leftarrow x_j$ 
22:        $\mathcal{Z} \leftarrow \mathcal{Z} \cup s$ 
23:     return  $\mathcal{Z}$  /* meaningful dense data w.r.t  $x$  */
24:   return  $\mathcal{Z}$ 
```

Draw perturbed samples from
training set distribution



Train surrogate random forest



Observation-level feature importance
from random forest explanations



Values of features in s with the same
contributions as in x set to values from x

Algorithm 3 Representative Data Selection

```
1: procedure DATASELECTION( $x, f, \mathcal{Z}, \tau$ )
2:    $n_c \leftarrow 2$  /*  $n_c$ : number of clusters */
3:    $\mathcal{Z}' \leftarrow \{\}$ 
4:   for all  $l \in \mathcal{L}$  do /*  $\mathcal{L}$ : set of labels */
5:      $\mathcal{G}_l \leftarrow \{z \in \mathcal{Z} \mid f(z) = l\}$ 
6:      $\mathcal{G}_l \leftarrow x \cup \mathcal{G}_l$ 
7:     while True do
8:        $c_x, c_{\neg x} \leftarrow \text{AgglomerativeClustering}(\mathcal{G}_l, n_c)$ 
9:       if  $|c_x| \geq \tau$  then
10:         $\mathcal{G}_l \leftarrow c_x$ 
11:       else
12:        break
13:      $\mathcal{Z}' \leftarrow \mathcal{Z}' \cup \mathcal{G}_l$ 
14: return  $\mathcal{Z}'$  /* representative data set */
```

Data Balancing

Generate samples in classes with small representations via equation:

$$x_{new} = x_i + \lambda \times (x_{zi} - x_i)$$

Rule-based Interpretable model

Build a classifier (decision tree) on obtained dataset and values of black-box model

Experiments

Experimental setup

- Matched against LIME, LORE, and Anchor explanation methods on:
 - fidelity comparison
 - neighborhood analysis
 - explanation comparison
- Black-boxes: Neural Network, Logistic Regression, Gradient Boosting
- Datasets:

TABLE I
DESCRIPTION OF THE DATA SETS.

Data set	# Instances	# Features	Class imbalance
<i>Adult</i>	49K	14	$\leq 50K$: 76% - $> 50K$: 24%
<i>German</i>	1K	20	Good: 70% - Bad: 30%
<i>COMPAS</i>	7K	52	Medium-Low: 72% - High: 28%

TABLE II
COMPARISON OF *fidelity_x* SCORES.

Data set	Black-box	EXPLAN	LIME	LORE
<i>Adult</i>	GB	0.994±.1	0.838±.4	0.980±.1
	LR	0.992±.1	0.940±.2	0.989±.1
	NN	0.992±.1	0.859±.3	0.977±.2
<i>German</i>	GB	1.000±.0	0.910±.3	0.950±.2
	LR	0.990±.1	0.940±.2	0.910±.3
	NN	1.000±.0	0.930±.3	0.990±.1
<i>COMPAS</i>	GB	1.000±.0	0.911±.3	0.999±.0
	LR	1.000±.0	0.925±.3	0.981±.1
	NN	0.999±.0	0.915±.3	0.986±.1

instance-level

TABLE III
COMPARISON OF *fidelity_x* SCORES.

Data set	Black-box	EXPLAN	LIME	LORE
<i>Adult</i>	GB	0.971±.0	0.738±.0	0.996±.0
	LR	0.990±.0	0.793±.1	0.995±.0
	NN	0.980±.0	0.804±.0	0.993±.0
<i>German</i>	GB	0.942±.0	0.223±.1	0.979±.0
	LR	0.972±.0	0.179±.1	0.944±.2
	NN	0.981±.0	0.037±.1	0.987±.0
<i>COMPAS</i>	GB	0.984±.0	0.897±.0	0.982±.1
	LR	0.988±.0	0.919±.0	0.975±.1
	NN	0.988±.0	0.896±.0	0.974±.1

neighborhood-level

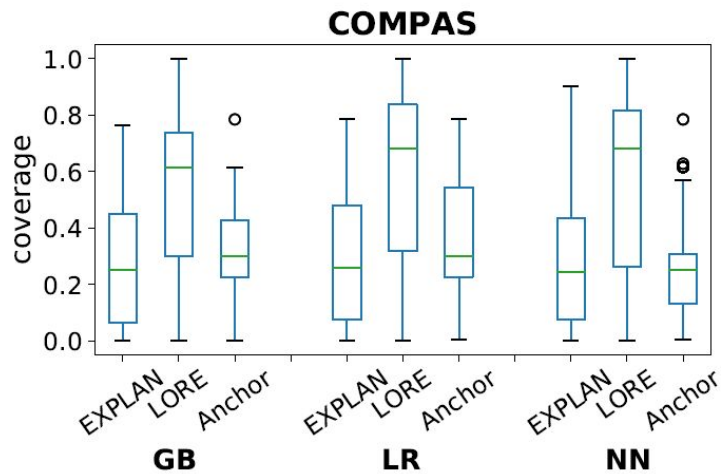
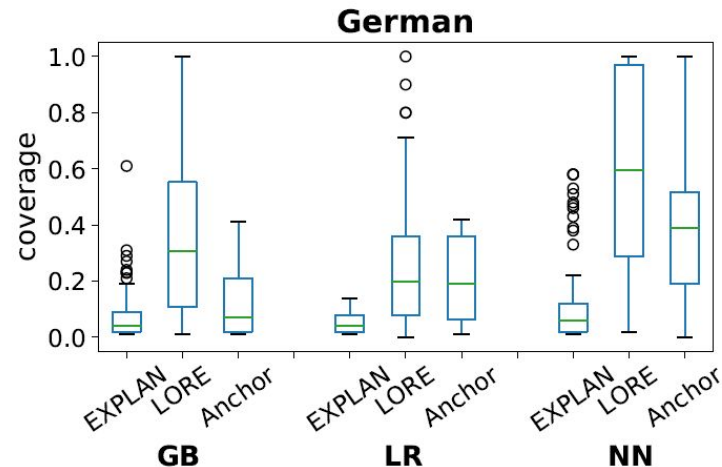
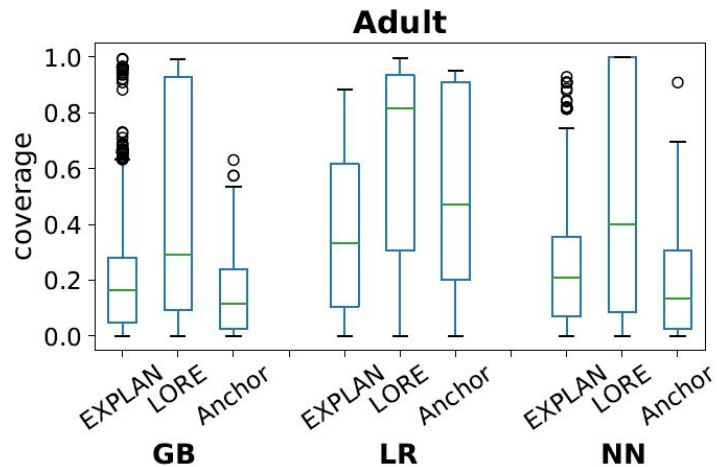


TABLE IV
COMPARISON OF $precision_e$ SCORES.

Data set	Black-box	EXPLAN	LORE	Anchor
<i>Adult</i>	GB	0.924 \pm .1	0.852 \pm .2	0.980 \pm .1
	LR	0.966 \pm .1	0.894 \pm .2	0.963 \pm .0
	NN	0.895 \pm .2	0.815 \pm .2	0.971 \pm .1
<i>German</i>	GB	0.897 \pm .2	0.816 \pm .2	0.984 \pm .0
	LR	0.937 \pm .1	0.835 \pm .3	0.994 \pm .0
	NN	0.950 \pm .1	0.879 \pm .2	0.976 \pm .1
<i>COMPAS</i>	GB	0.914 \pm .2	0.855 \pm .2	0.963 \pm .0
	LR	0.912 \pm .2	0.862 \pm .2	0.963 \pm .0
	NN	0.898 \pm .2	0.851 \pm .2	0.979 \pm .0

TABLE V
COMPARISON OF FEATURE FREQUENCY VARIANCE.

Data set	Black-box	EXPLAN	LORE	Anchor
<i>Adult</i>	GB	1.191 \pm .2	2.263 \pm .5	1.497 \pm .1
	LR	1.087 \pm .2	2.259 \pm .5	1.523 \pm .1
	NN	1.174 \pm .5	2.273 \pm .5	1.509 \pm .1
<i>German</i>	GB	0.469 \pm .0	2.296 \pm .6	0.472 \pm .0
	LR	0.513 \pm .0	2.293 \pm .7	0.469 \pm .0
	NN	0.467 \pm .0	2.334 \pm .6	0.468 \pm .0
<i>COMPAS</i>	GB	0.529 \pm .1	1.342 \pm .4	0.682 \pm .1
	LR	0.535 \pm .1	1.335 \pm .4	0.678 \pm .1
	NN	0.533 \pm .1	1.314 \pm .4	0.681 \pm .1

TABLE VI
COMPARISON OF JACCARD MEASURE OF STABILITY.

Data set	Black-box	EXPLAN	LORE	Anchor
<i>Adult</i>	GB	0.827 \pm .1	0.821 \pm .1	0.755 \pm .1
	LR	0.859 \pm .1	0.799 \pm .1	0.671 \pm .1
	NN	0.856 \pm .1	0.728 \pm .1	0.744 \pm .2
<i>German</i>	GB	0.702 \pm .1	0.694 \pm .2	0.754 \pm .1
	LR	0.729 \pm .1	0.698 \pm .2	0.819 \pm .2
	NN	0.846 \pm .1	0.779 \pm .1	0.884 \pm .1
<i>COMPAS</i>	GB	0.888 \pm .1	0.858 \pm .2	0.859 \pm .1
	LR	0.886 \pm .1	0.859 \pm .2	0.854 \pm .1
	NN	0.848 \pm .1	0.807 \pm .2	0.822 \pm .1

Explanation generated by EXPLAN for an input from *Adult* data set and GB black-box

$x = \{$ **age:** 30;
 workclass: Private;
 marital-status: Never-married;
 occupation: Prof-speciality;
 relationship: Unmarried;
 race: White;
 sex: Male;
 capital-gain: 0;
 capital-loss: 0;
 hours-per-week: 40;
 native-country: United-States
 class: \leq 50K

$e = \{$ **age:** \leq 30,
 capital-gain \leq 0
 hours-per-week \leq 44
 -> class: \leq 50K

Thank you for your attention

Time for Q&A