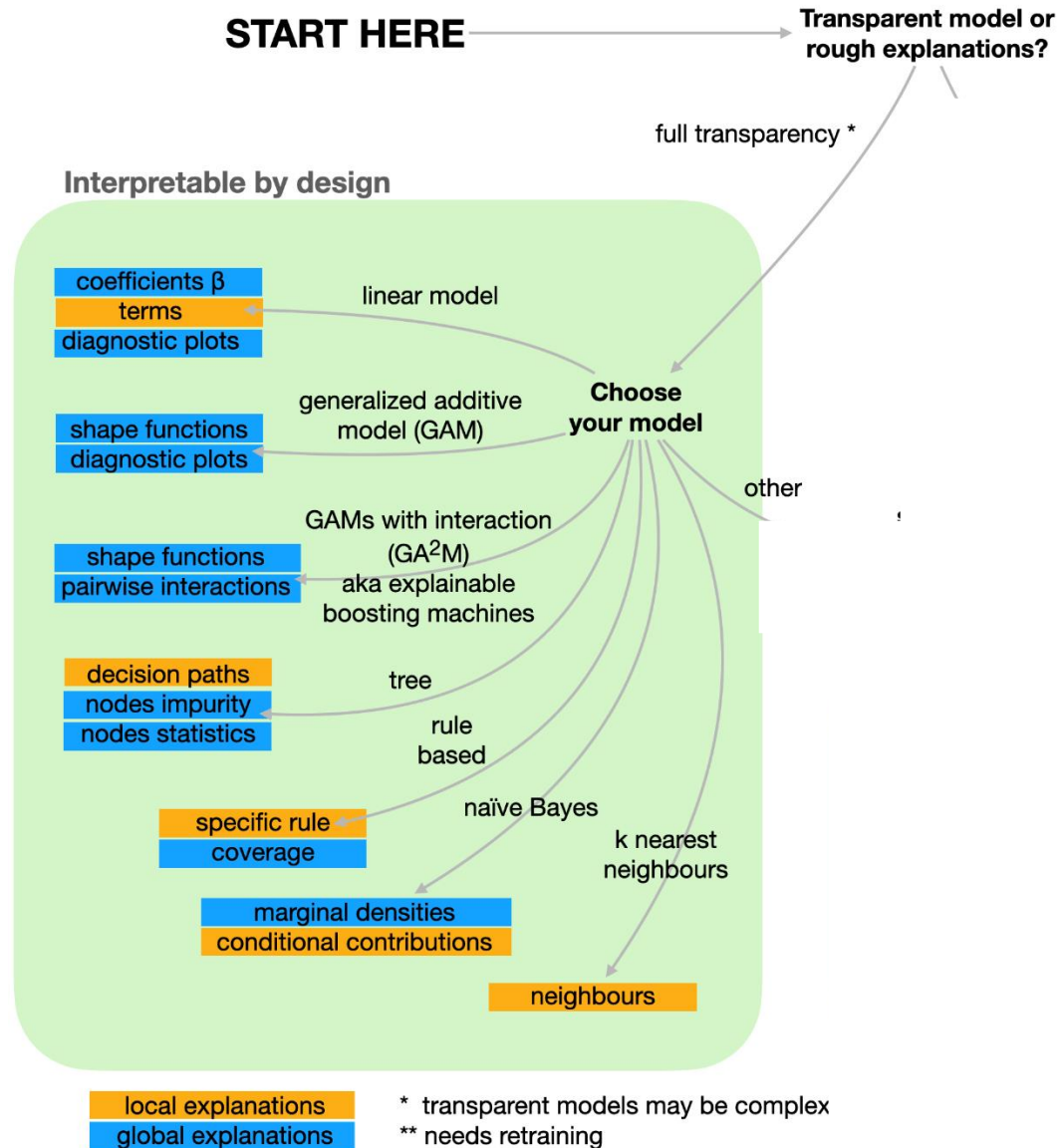


# Landscape of R packages for eXplainable Artificial Intelligence

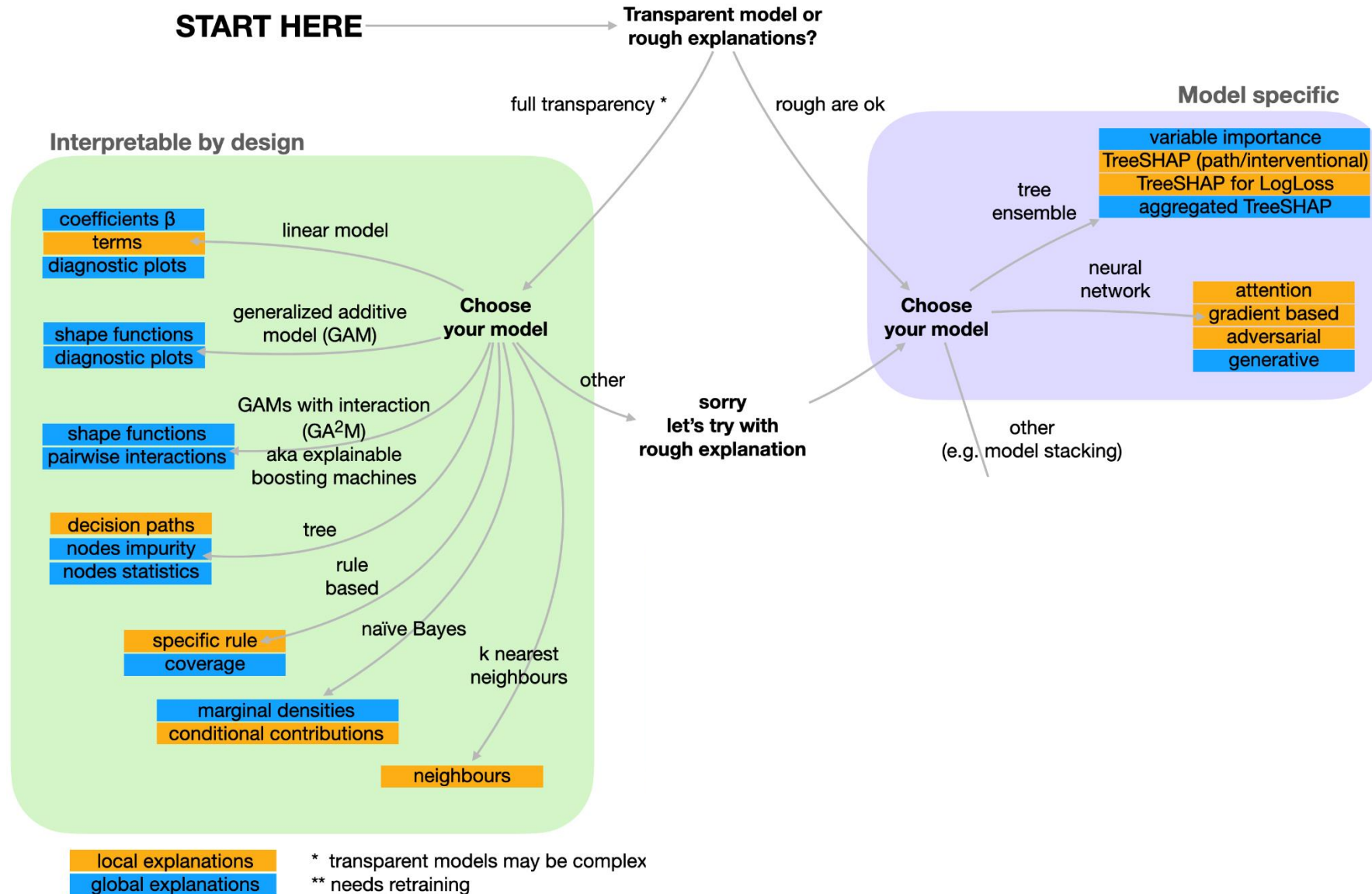
Szymon Maksymiuk and Alicja Gosiewska

12.10.2020

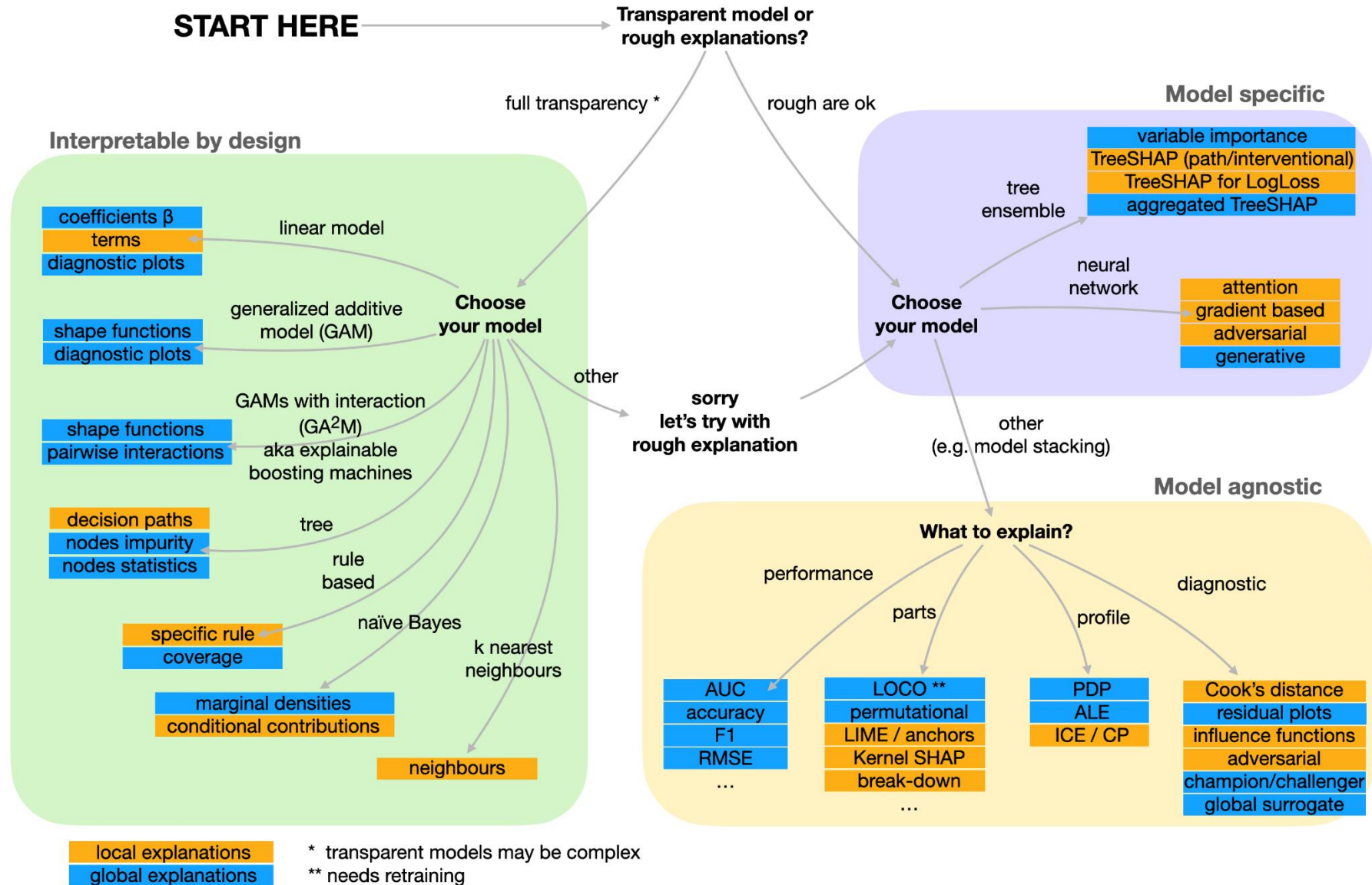
# Taxonomy in terms of models



# Taxonomy in terms of models



# Taxonomy in terms of models



# Taxonomy in terms of explanation object

**Global Explanations**

**Local Explanations**

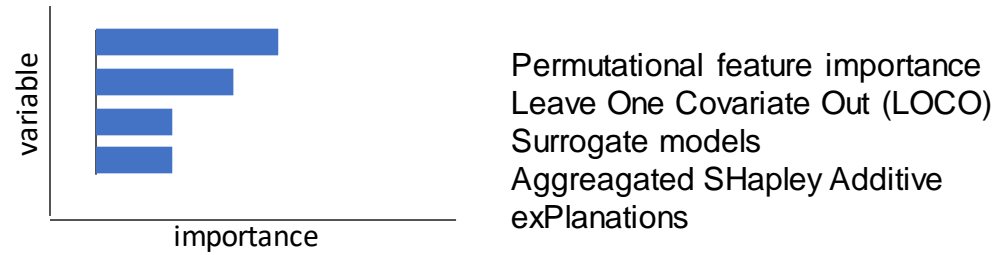


# Taxonomy in terms of explanation object

## Global Explanations

## Local Explanations

### Model Parts

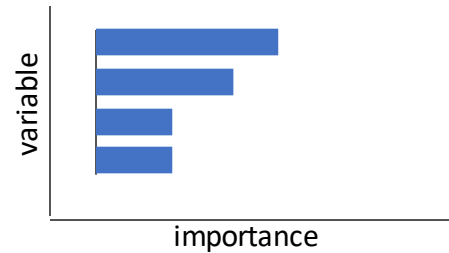


# Taxonomy in terms of explanation object

## Global Explanations

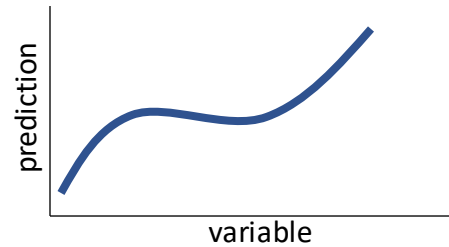
## Local Explanations

### Model Parts



Permutational feature importance  
Leave One Covariate Out (LOCO)  
Surrogate models  
Aggregated SHapley Additive  
exPlanations

### Model Profile



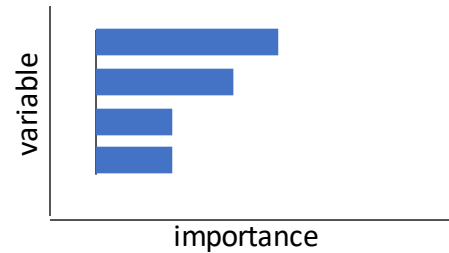
Partial Dependence Profiles (PDP)  
Accumulated Local Effects (ALE)

# Taxonomy in terms of explanation object

## Global Explanations

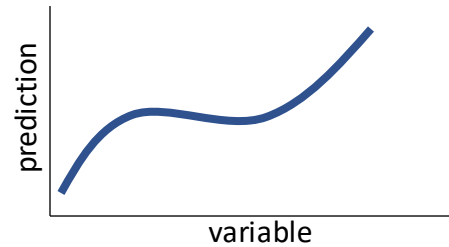
## Local Explanations

### Model Parts



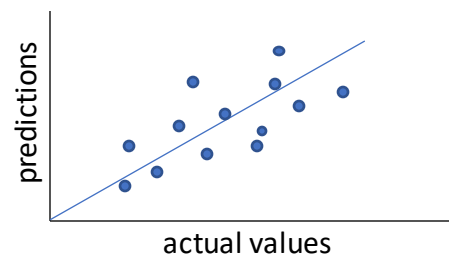
Permutational feature importance  
Leave One Covariate Out (LOCO)  
Surrogate models  
Aggreagated SHapley Additive  
exPlanations

### Model Profile



Partial Dependence Profiles (PDP)  
Accumulated Local Effects (ALE)

### Model Diagnostics



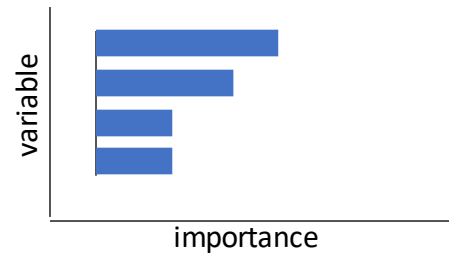
Residual plots  
Variable vs. prediction plots  
Demographic parity



# Taxonomy in terms of explanation object

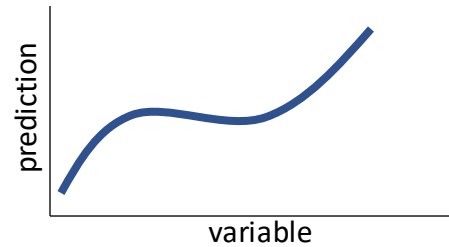
## Global Explanations

### Model Parts



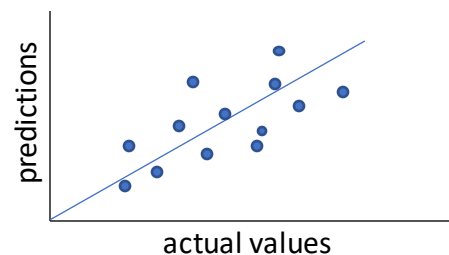
Permutational feature importance  
Leave One Covariate Out (LOCO)  
Surrogate models  
Aggregated SHapley Additive  
exPlanations

### Model Profile



Partial Dependence Profiles (PDP)  
Accumulated Local Effects (ALE)

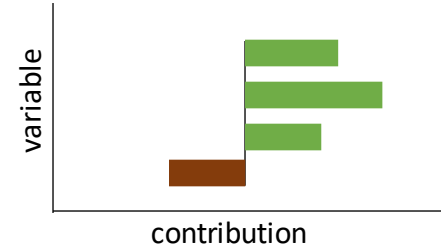
### Model Diagnostics



Residual plots  
Variable vs. prediction plots  
Demographic parity

## Local Explanations

### Predict Parts

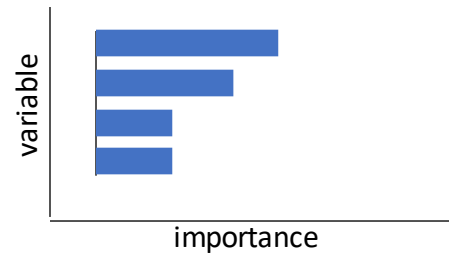


Break Down (BD)  
SHapley Additive exPlanations (SHAP)  
Local Interpretable Model agnostic  
Explanations (LIME)

# Taxonomy in terms of explanation object

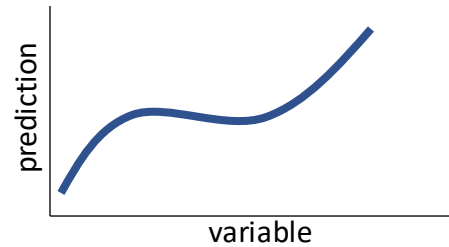
## Global Explanations

### Model Parts



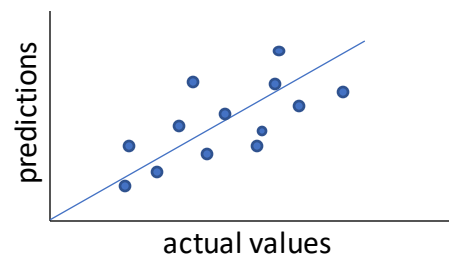
Permutational feature importance  
Leave One Covariate Out (LOCO)  
Surrogate models  
Aggregated SHapley Additive  
exPlanations

### Model Profile



Partial Dependence Profiles (PDP)  
Accumulated Local Effects (ALE)

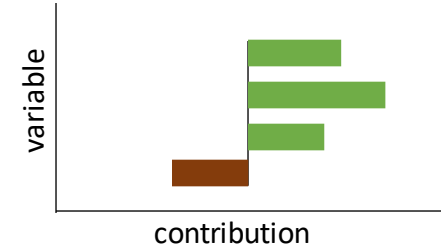
### Model Diagnostics



Residual plots  
Variable vs. prediction plots  
Demographic parity

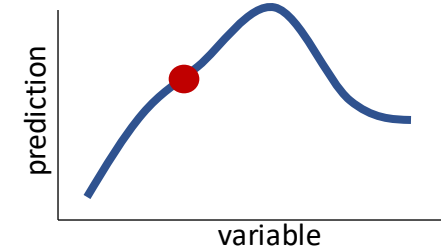
## Local Explanations

### Predict Parts



Break Down (BD)  
SHapley Additive exPlanations (SHAP)  
Local Interpretable Model agnostic  
Explanations (LIME)

### Predict Profile

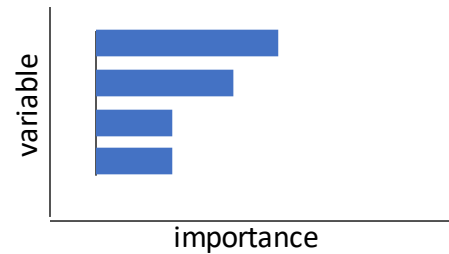


Ceteris Paribus (CP) / Individual  
Conditional Expectations (ICE)

# Taxonomy in terms of explanation object

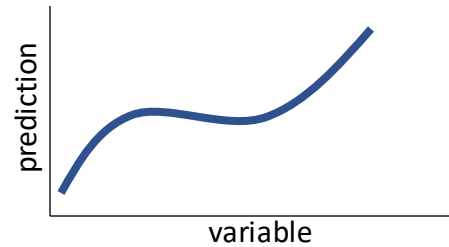
## Global Explanations

### Model Parts



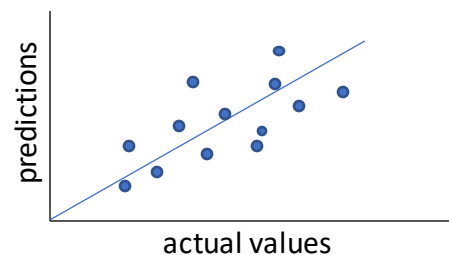
Permutational feature importance  
Leave One Covariate Out (LOCO)  
Surrogate models  
Aggregated SHapley Additive  
exPlanations

### Model Profile



Partial Dependence Profiles (PDP)  
Accumulated Local Effects (ALE)

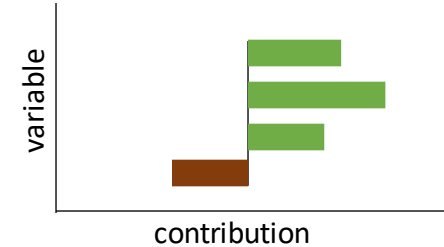
### Model Diagnostics



Residual plots  
Variable vs. prediction plots  
Demographic parity

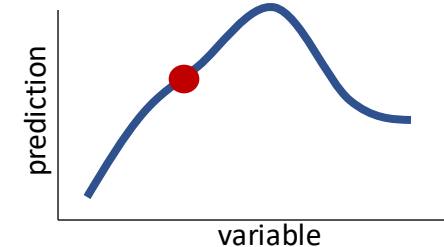
## Local Explanations

### Predict Parts



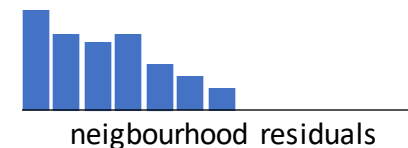
Break Down (BD)  
SHapley Additive exPlanations (SHAP)  
Local Interpretable Model agnostic  
Explanations (LIME)

### Predict Profile



Ceteris Paribus (CP) / Individual  
Conditional Expectations (ICE)

### Predict Diagnostics



Local residual density plot

# Six Python libraries for XAI

# eli5 (Python)

- One of the oldest XAI libraries.
- Supports the most common Python frameworks and packages: scikit-learn, Keras, xgboost, LightGBM, CatBoost, lightning, and sklearn-crfsuite.
- Good for text, image, and tabular data.

## Feature Importance

Weight	Feature
0.4278	Sex=female
0.1949	Pclass=3
0.0665	Embarked=S
0.0510	Pclass=2
0.0420	SibSp
0.0417	Cabin=
0.0385	Embarked=C
0.0358	Ticket=1601
0.0331	Age
0.0323	Fare
0.0220	Pclass=1
0.0143	Parch

# eli5 (Python)

## LIME

y=alt.atheism (probability 0.000, score -9.663) top features

Contribution <sup>2</sup>	Feature
-0.360	<BIAS>
-9.303	Highlighted in text (sum)

as i recall from my bout with kidney stones, there isn't any medication that can do anything about them except relieve the pain. either they pass, or they have to be broken up with sound, or they have to be extracted surgically. when i was in, the x-ray tech happened to mention that she'd had kidney stones and children, and the childbirth hurt less.

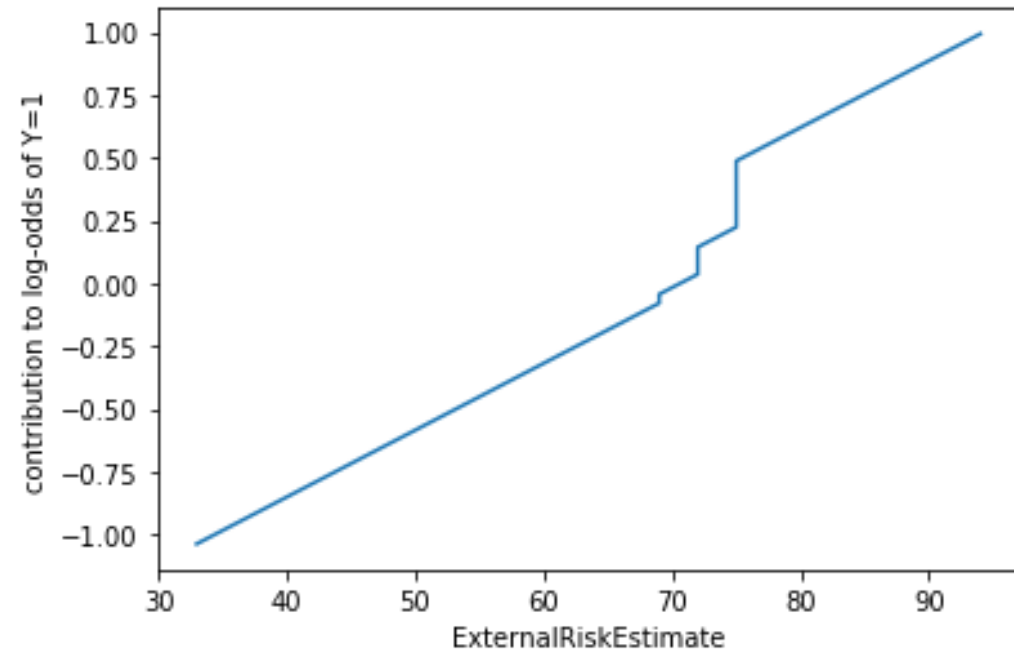
## Grad-CAM



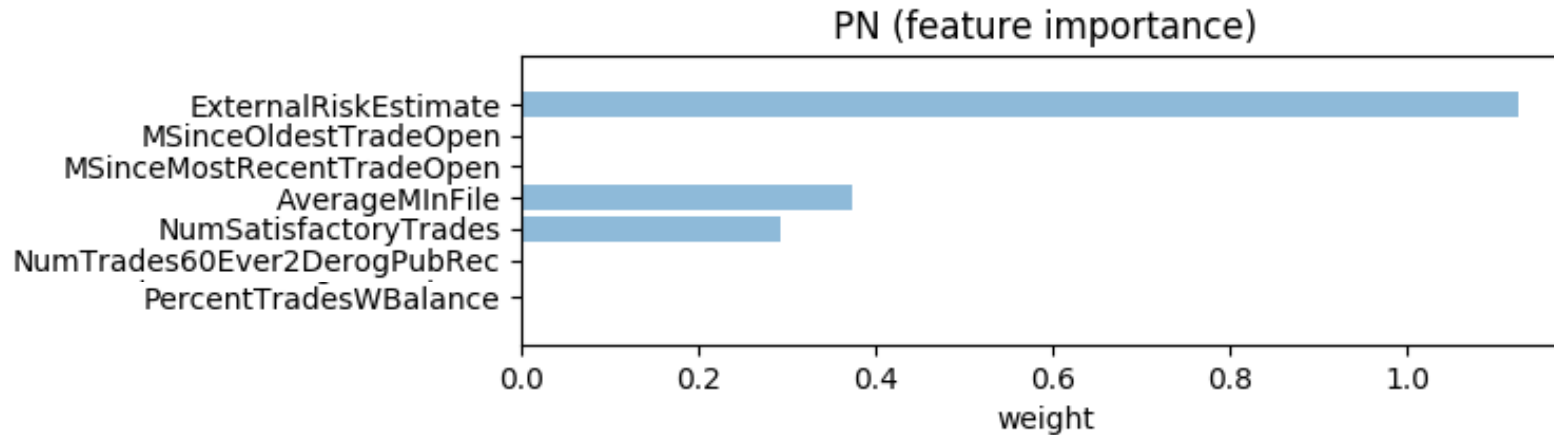
# aix360 (Python)

- Explainability of both models and data sets.
- A wide range of methods: local, global, model-specific, model-agnostic.
- Many tutorials and example notebooks

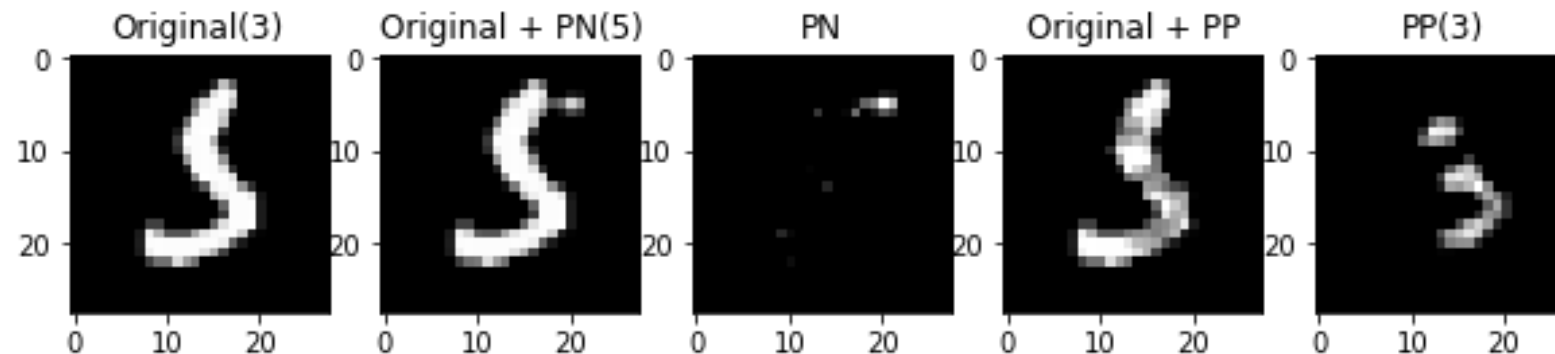
**Partial Dependence Profile**



# aix360 (Python)



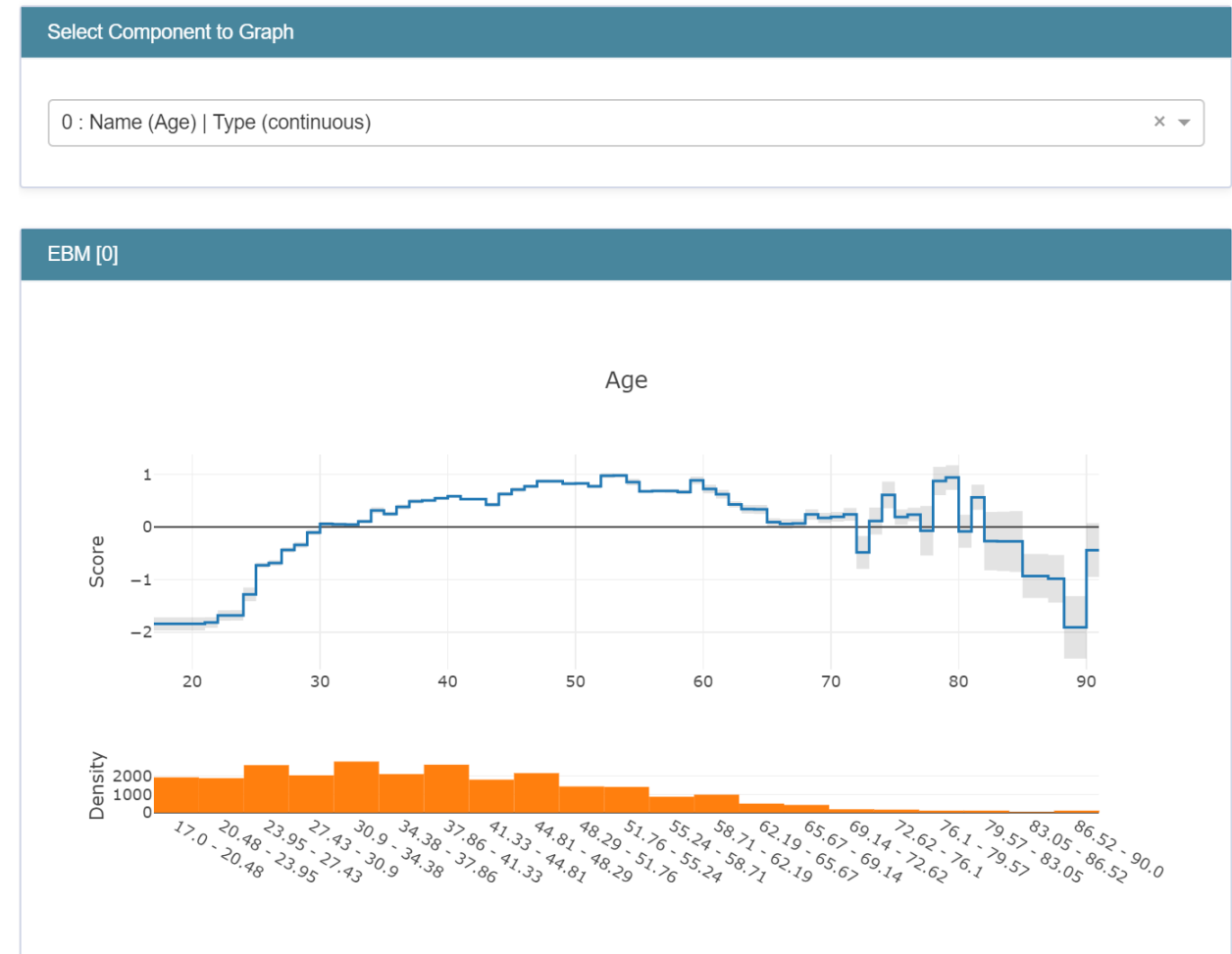
## Plot Pertinent Negative (PN) and Pertinent Positive (PP) explanations





# interpret (Python)

- Provides interpretable models, such as Explainable Boosting and Decision Rule List.
- Many model-agnostic methods: LIME, SHAP, PDP.
- Provides comparisons of multiple models on interactive dashboards.



Interpret ML Dashboard

Overview

	Data
1	0.78
2	0.69
3	0.85
4	0.72
5	0.81
6	0.76
7	0.83
8	0.74
9	0.80
10	0.77
11	0.82
12	0.75
13	0.84
14	0.73
15	0.86
16	0.71
17	0.87
18	0.70
19	0.88
20	0.68
21	0.89
22	0.67
23	0.90
24	0.66
25	0.91
26	0.65
27	0.92
28	0.64
29	0.93
30	0.63
31	0.94
32	0.62
33	0.95
34	0.61
35	0.96
36	0.60
37	0.97
38	0.59
39	0.98
40	0.58
41	0.99
42	0.57
43	1.00
44	0.56
45	1.01
46	0.55
47	1.02
48	0.54
49	1.03
50	0.53
51	1.04
52	0.52
53	1.05
54	0.51
55	1.06
56	0.50
57	1.07
58	0.49
59	1.08
60	0.48
61	1.09
62	0.47
63	1.10
64	0.46
65	1.11
66	0.45
67	1.12
68	0.44
69	1.13
70	0.43
71	1.14
72	0.42
73	1.15
74	0.41
75	1.16
76	0.40
77	1.17
78	0.39
79	1.18
80	0.38
81	1.19
82	0.37
83	1.20
84	0.36
85	1.21
86	0.35
87	1.22
88	0.34
89	1.23
90	0.33
91	1.24
92	0.32
93	1.25
94	0.31
95	1.26
96	0.30
97	1.27
98	0.29
99	1.28
100	0.28
101	1.29
102	0.27
103	1.30
104	0.26
105	1.31
106	0.25
107	1.32
108	0.24
109	1.33
110	0.23
111	1.34
112	0.22
113	1.35
114	0.21
115	1.36
116	0.20
117	1.37
118	0.19
119	1.38
120	0.18
121	1.39
122	0.17
123	1.40
124	0.16
125	1.41
126	0.15
127	1.42
128	0.14
129	1.43
130	0.13
131	1.44
132	0.12
133	1.45
134	0.11
135	1.46
136	0.10
137	1.47
138	0.09
139	1.48
140	0.08
141	1.49
142	0.07
143	1.50
144	0.06
145	1.51
146	0.05
147	1.52
148	0.04
149	1.53
150	0.03
151	1.54
152	0.02
153	1.55
154	0.01
155	1.56
156	0.00
157	1.57
158	-0.01
159	1.58
160	-0.02
161	1.59
162	-0.03
163	1.60
164	-0.04
165	1.61
166	-0.05
167	1.62
168	-0.06
169	1.63
170	-0.07
171	1.64
172	-0.08
173	1.65
174	-0.09
175	1.66
176	-0.10
177	1.67
178	-0.11
179	1.68
180	-0.12
181	1.69
182	

Performance
-------------

Global

Local
-------

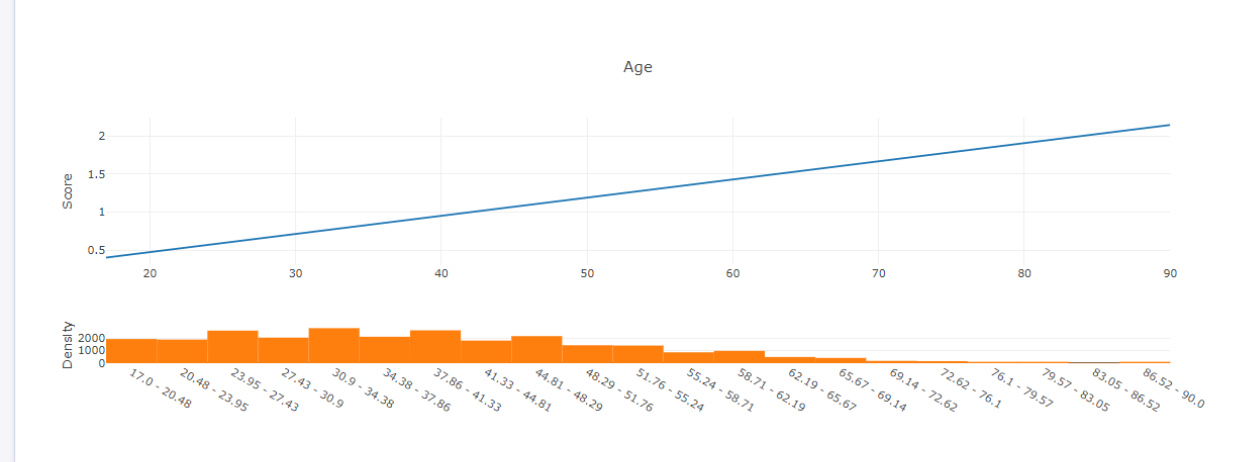
Select Explanation



Select Components to Graph

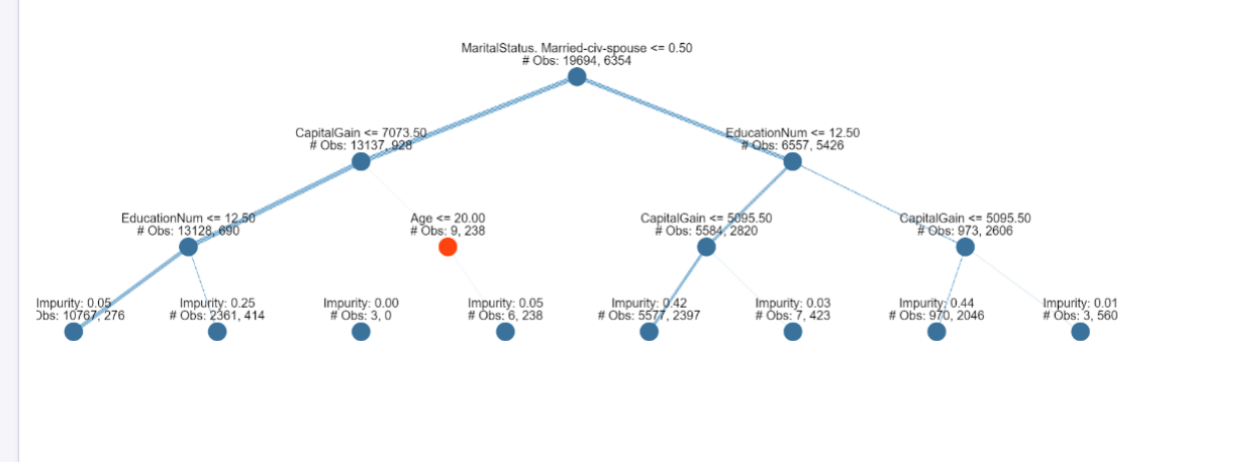
<input type="checkbox"/>	Name	Type	# Unique	% Non-zero	Importance	SelectID
<input checked="" type="checkbox"/>	Age	continuous	73	1		0
<input type="checkbox"/>	fnlwgt	continuous	18385	1		1
<input type="checkbox"/>	EducationNum	continuous	16	1		2
<input type="checkbox"/>	CapitalGain	continuous	116	0.084		3
<input type="checkbox"/>	CapitalLoss	continuous	88	0.047		4
<input type="checkbox"/>	HoursPerWeek	continuous	93	1		5
<input type="checkbox"/>	WorkClass. ?	categorical	2	0.056		6

LR [0]



LR (Overall)
0.99

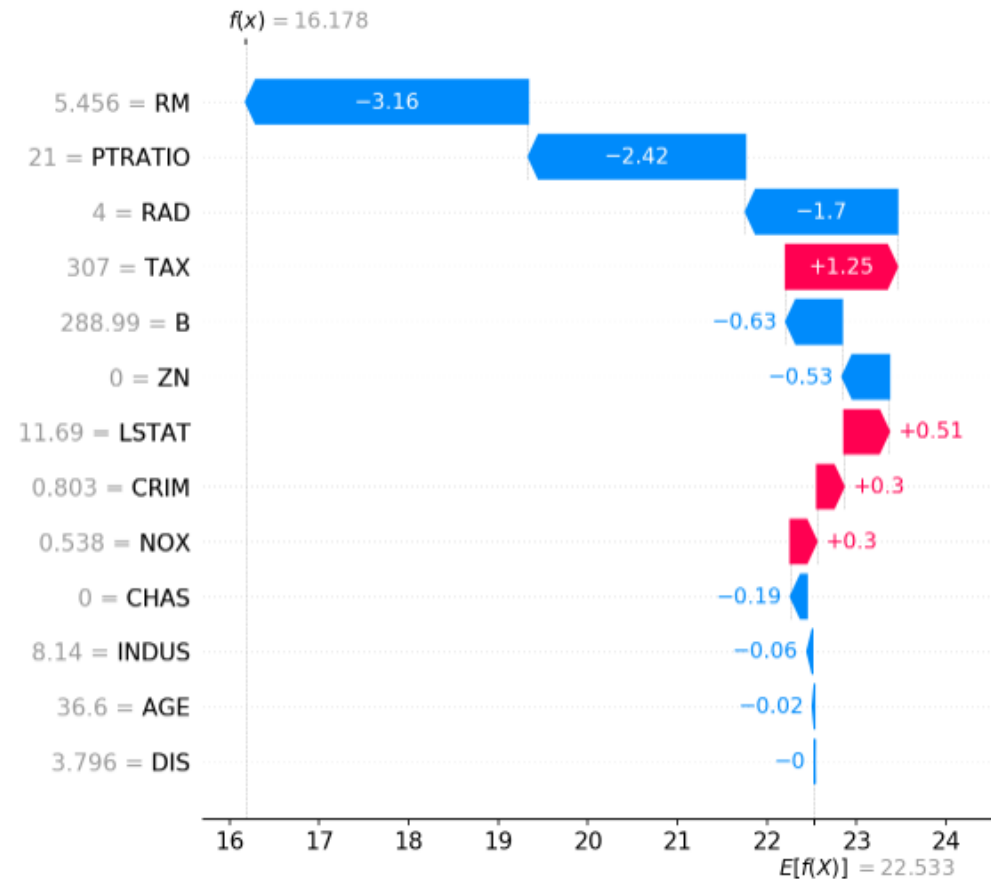
Tree [0]



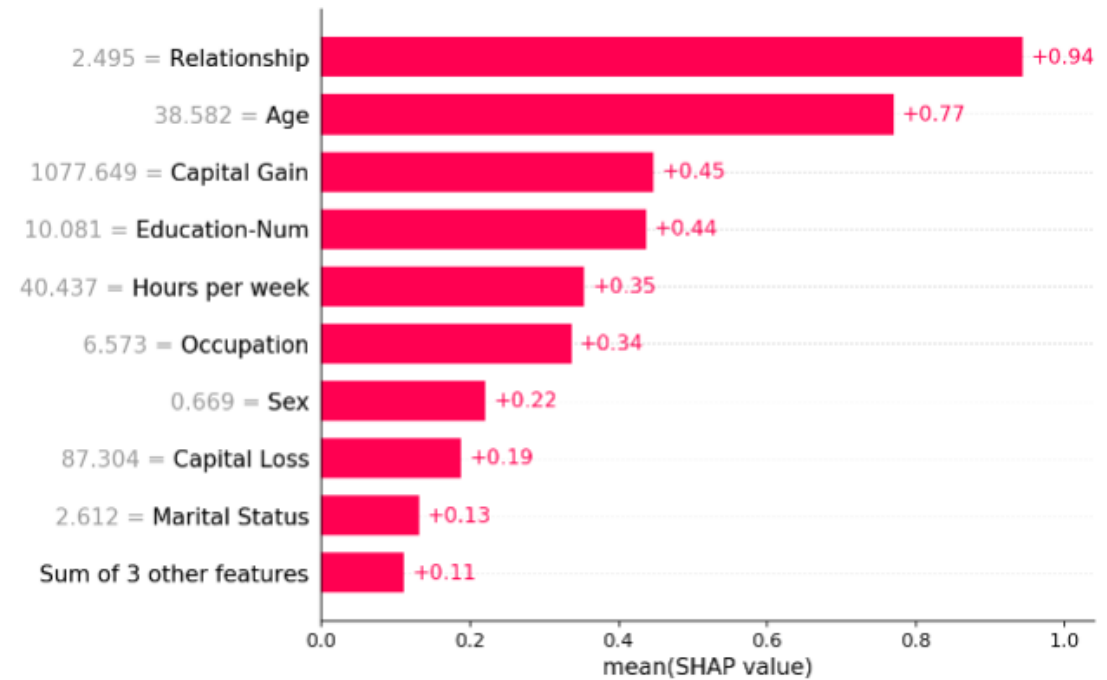
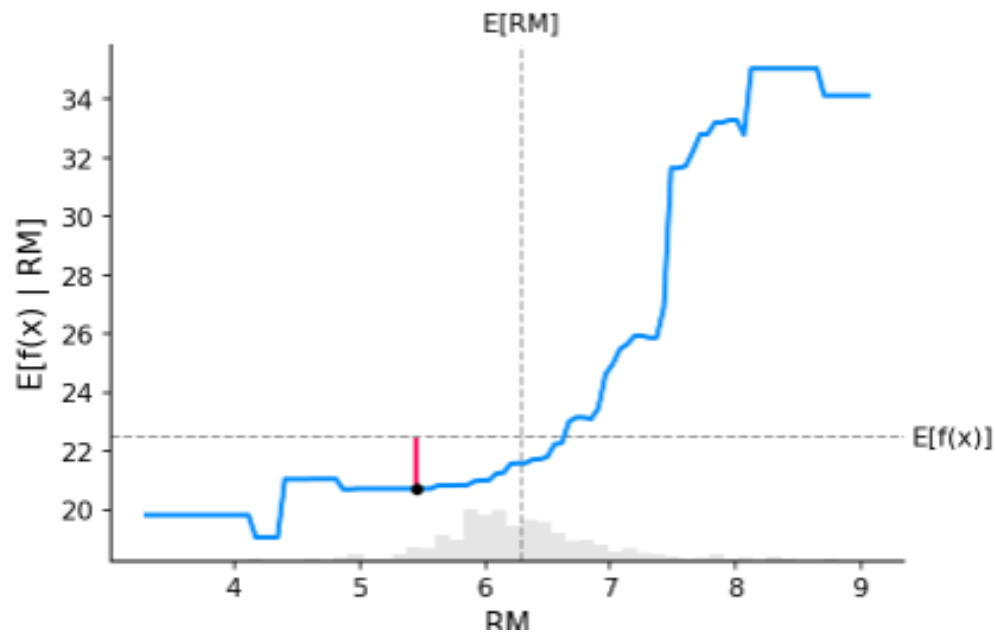
Tree (Overall)

# shap (Python)

- Popular and mature
- Strong theoretical background based on game theory
- Model agnostic
- Perfect for tabular image and text data
- Implements methods that comply with different task

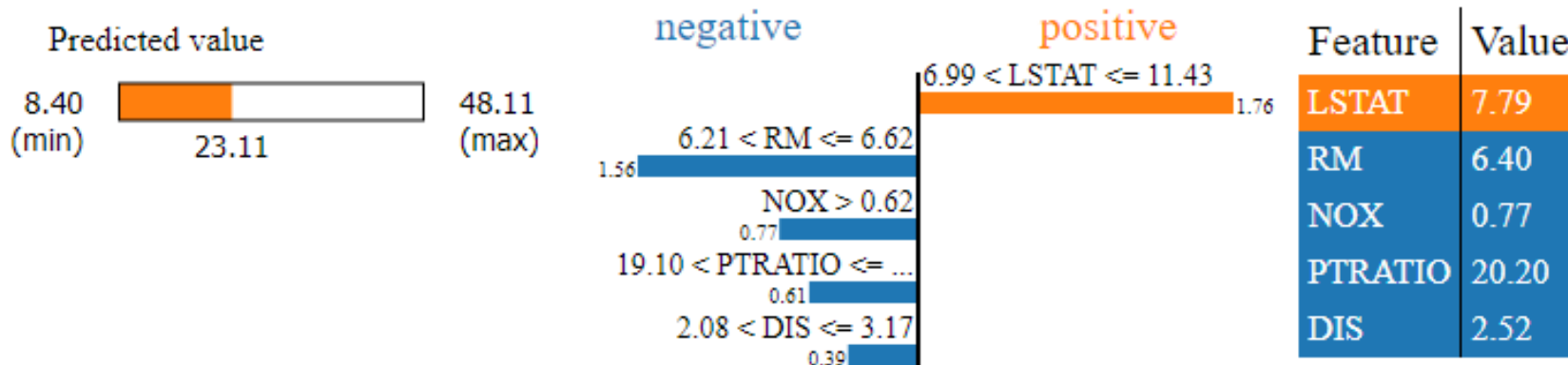


# shap (Python)



# lime (Python)

- Popular and mature
- Perfect for non-tabular data and models with hundreds of features
- Assumes local linearity of data

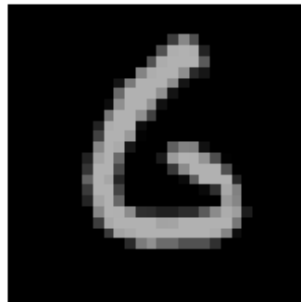


# lime (Python)

Positive for 0  
Actual 6



Positive for 1  
Actual 6



Positive for 2  
Actual 6



Positive for 3  
Actual 6



Positive for 4  
Actual 6



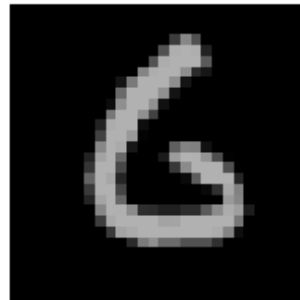
Positive for 5  
Actual 6



Positive for 6  
Actual 6



Positive for 7  
Actual 6



Positive for 8  
Actual 6

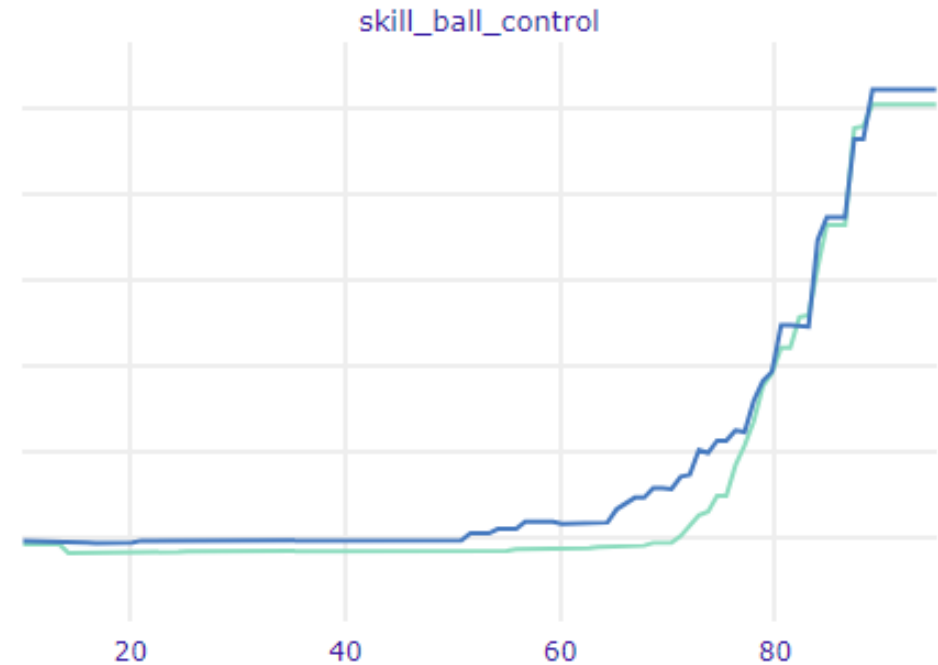


Positive for 9  
Actual 6



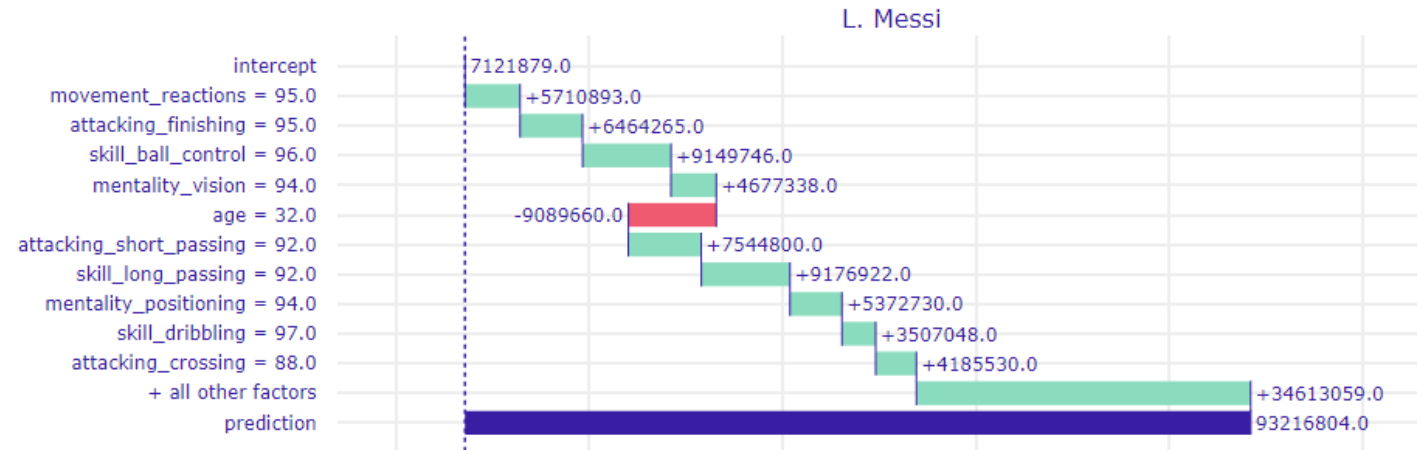
# dalex

- Combain library for plenty of explanation methods
- Based on R DALEX package
- Sufficient documentation with theoritical background for all implemented methods

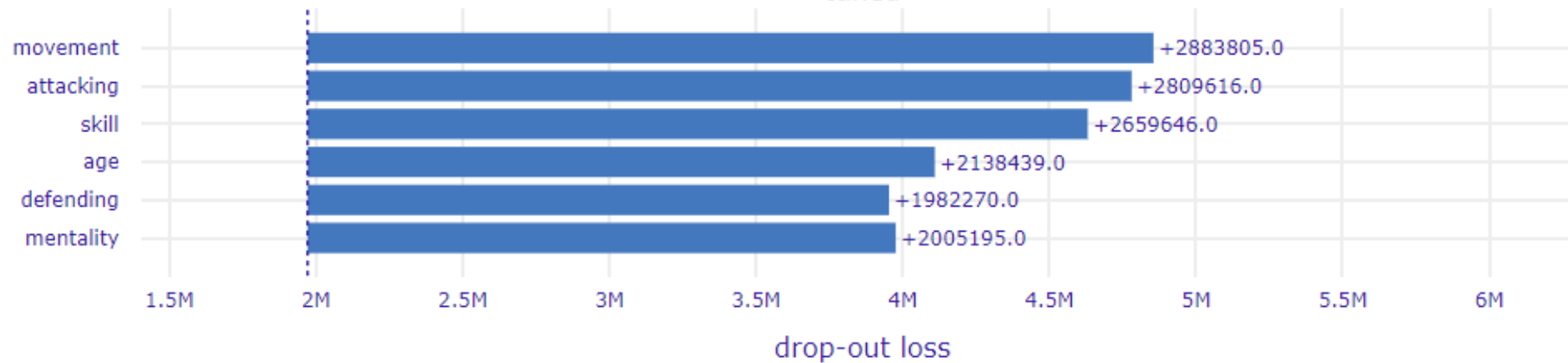


# dalex

Break Down



tuned

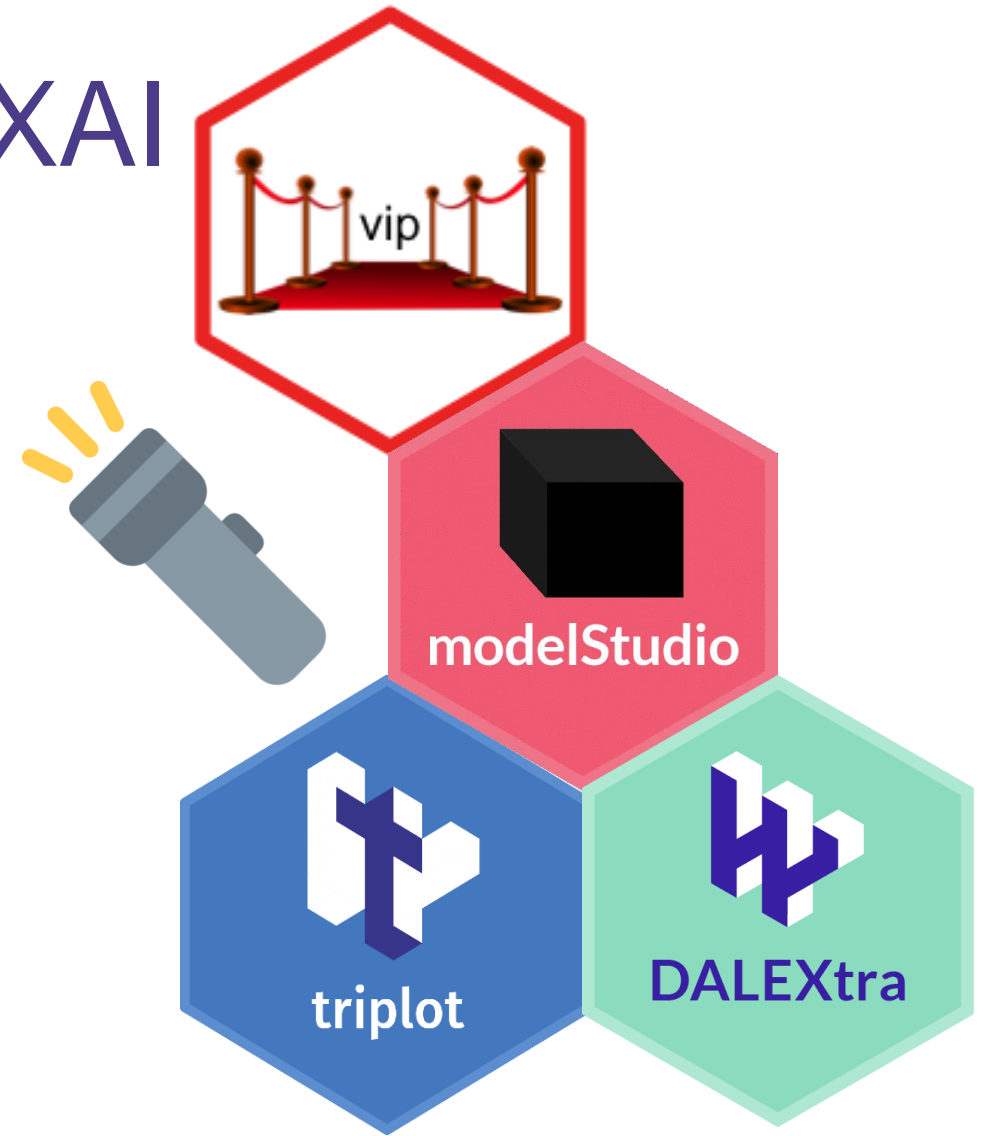




# Hot Five R packages for XAI

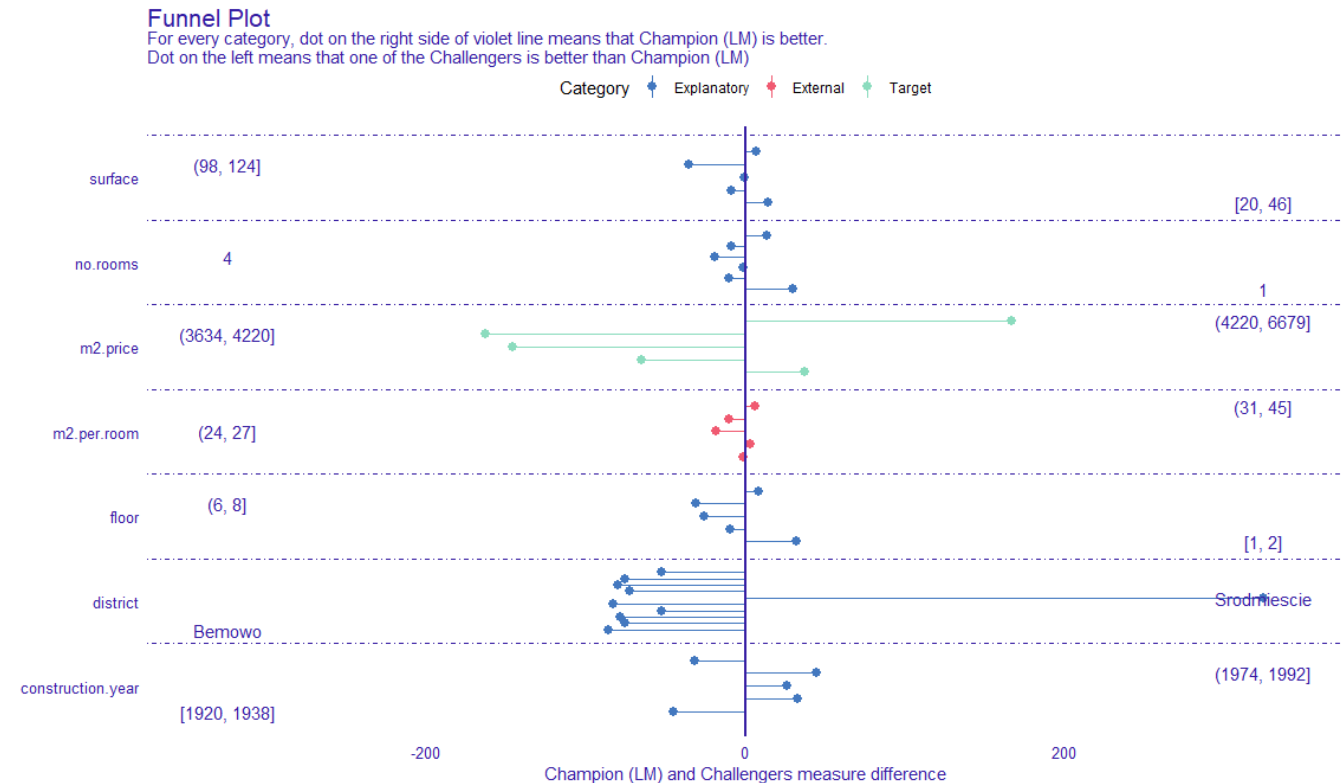
# Hot Five R Packages for XAI

- Outstand other packages in specific areas.
- The only one that implement interesting conceptions.
- Answer different XAI issues.



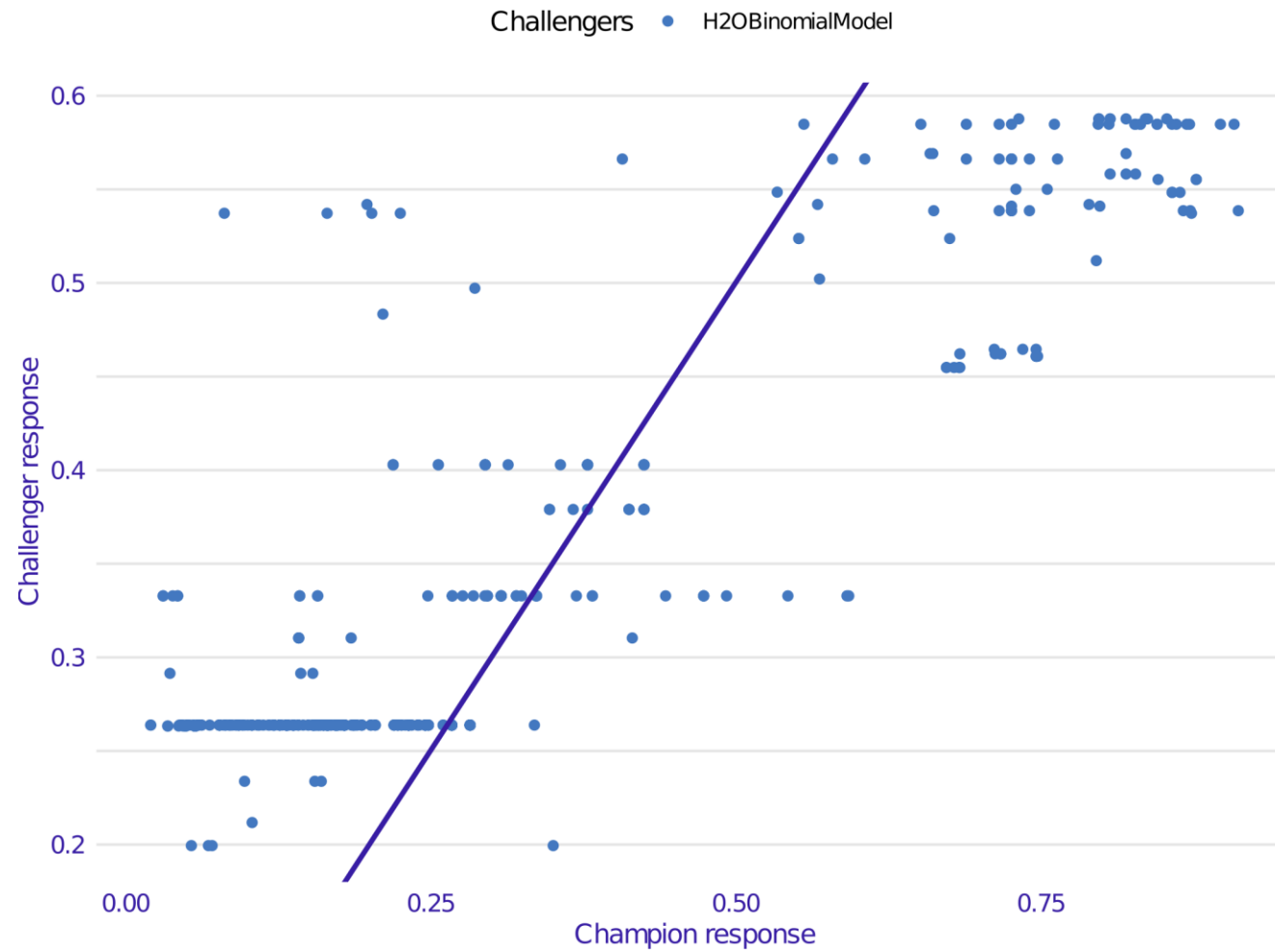
# DALEXtra

- Provides various methods for model diagnostics, including Funnel Plot.
- Allows users to compare more than one predictive model at once.
- Integration with various Machine Learning frameworks, also from different languages.  
For example, mlr, caret, scikit-learn, and more  
Other XAI tools do not!



Use-case on the Titanic dataset is under  
<https://mi2datalab.github.io/XAI-tools/DALEXtra.html>.

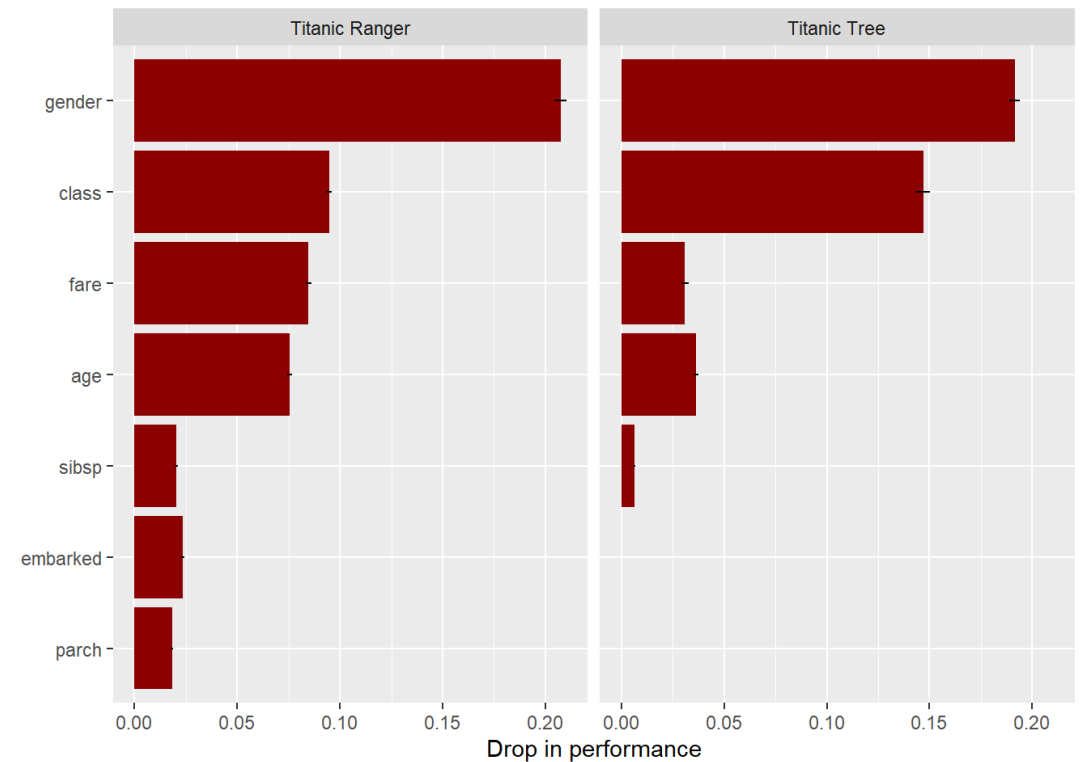
# DALEXtra



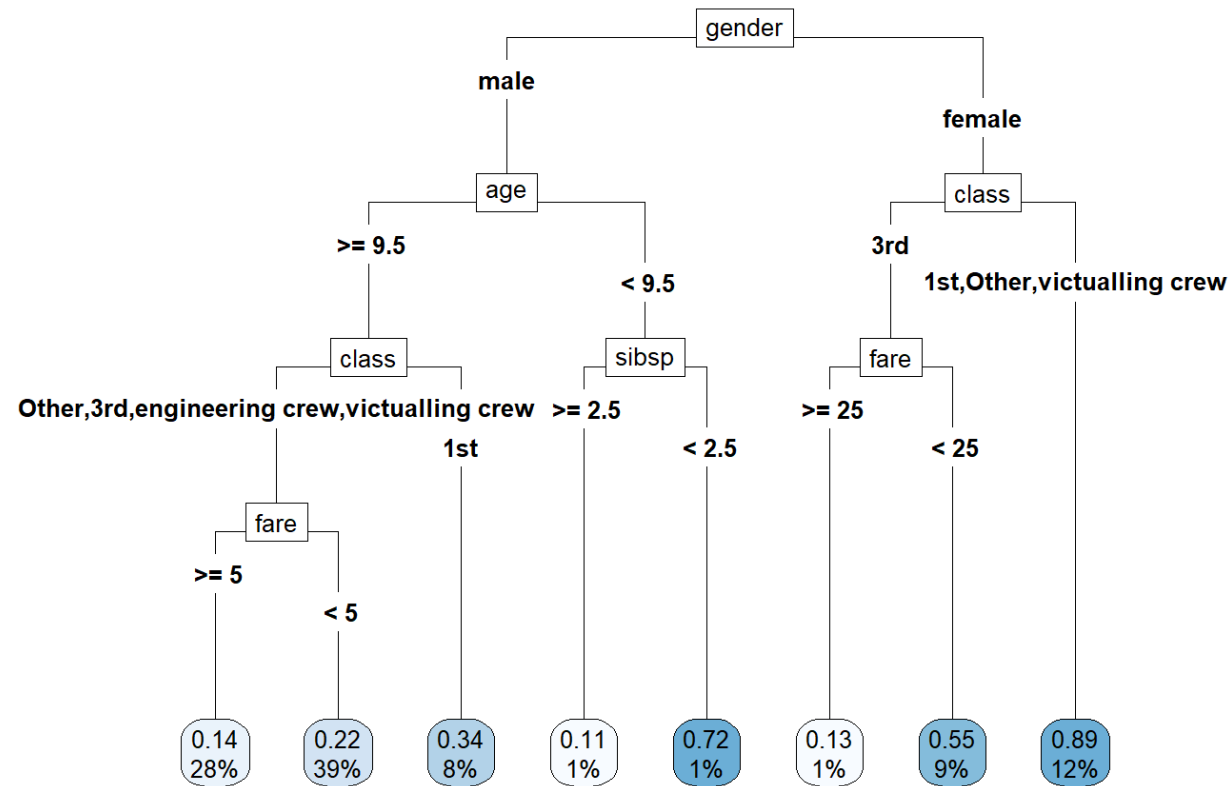
# flashlight

- Provides variable importance, PDP, ALE, residual, target and predicted value profiles. Local explanations: SHAP, BreakDown, and ICE.
- Observations can be weighted and taken into consideration while computing explanations  
No such feature in DALEX-verse or iml!

Use-case on the Titanic dataset is under  
<https://mi2datalab.github.io/XAI-tools/flashlight.html>.



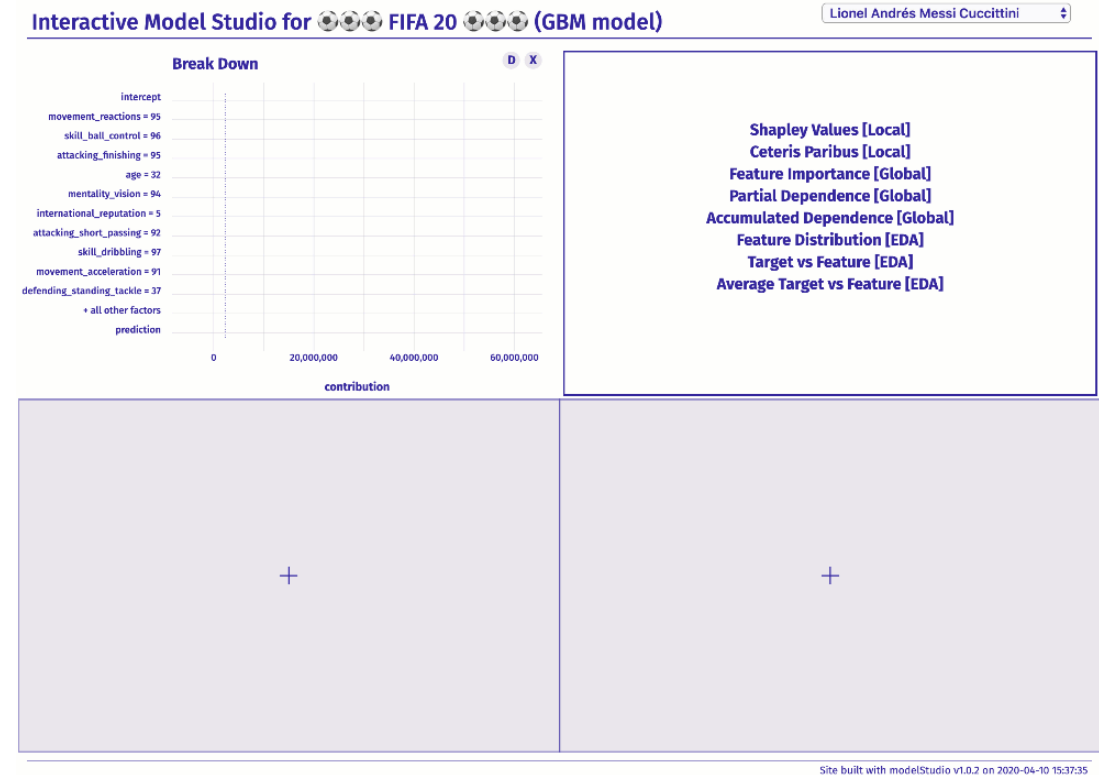
# flashlight



# modelStudio

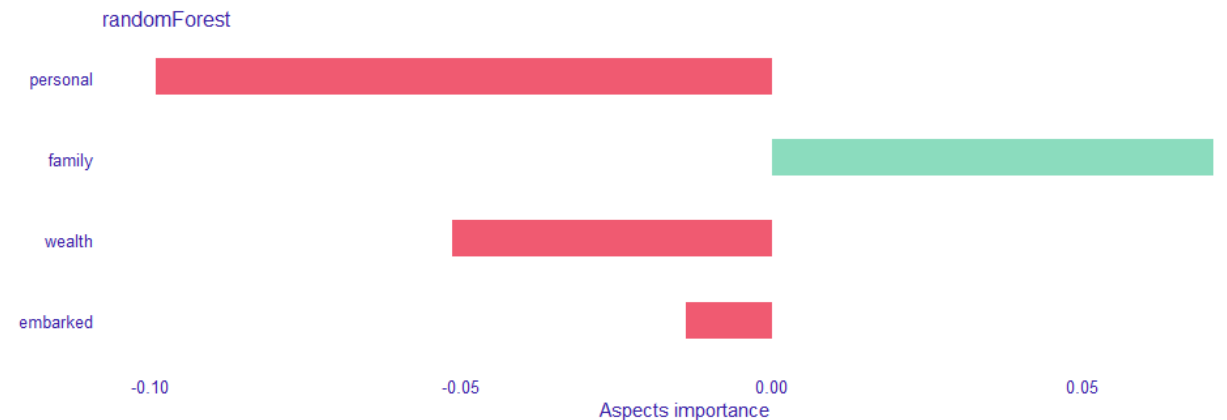
- Flexible standalone interface to model explanations.
- Ability to smoothly navigate between many different explanations at the same time
- Usage does not require fluency in programming.

Use-case on the Titanic dataset is under  
[https://mi2datalab.github.io/XAI-tools/modelStudio\\_titanic.html](https://mi2datalab.github.io/XAI-tools/modelStudio_titanic.html)



# tripplot

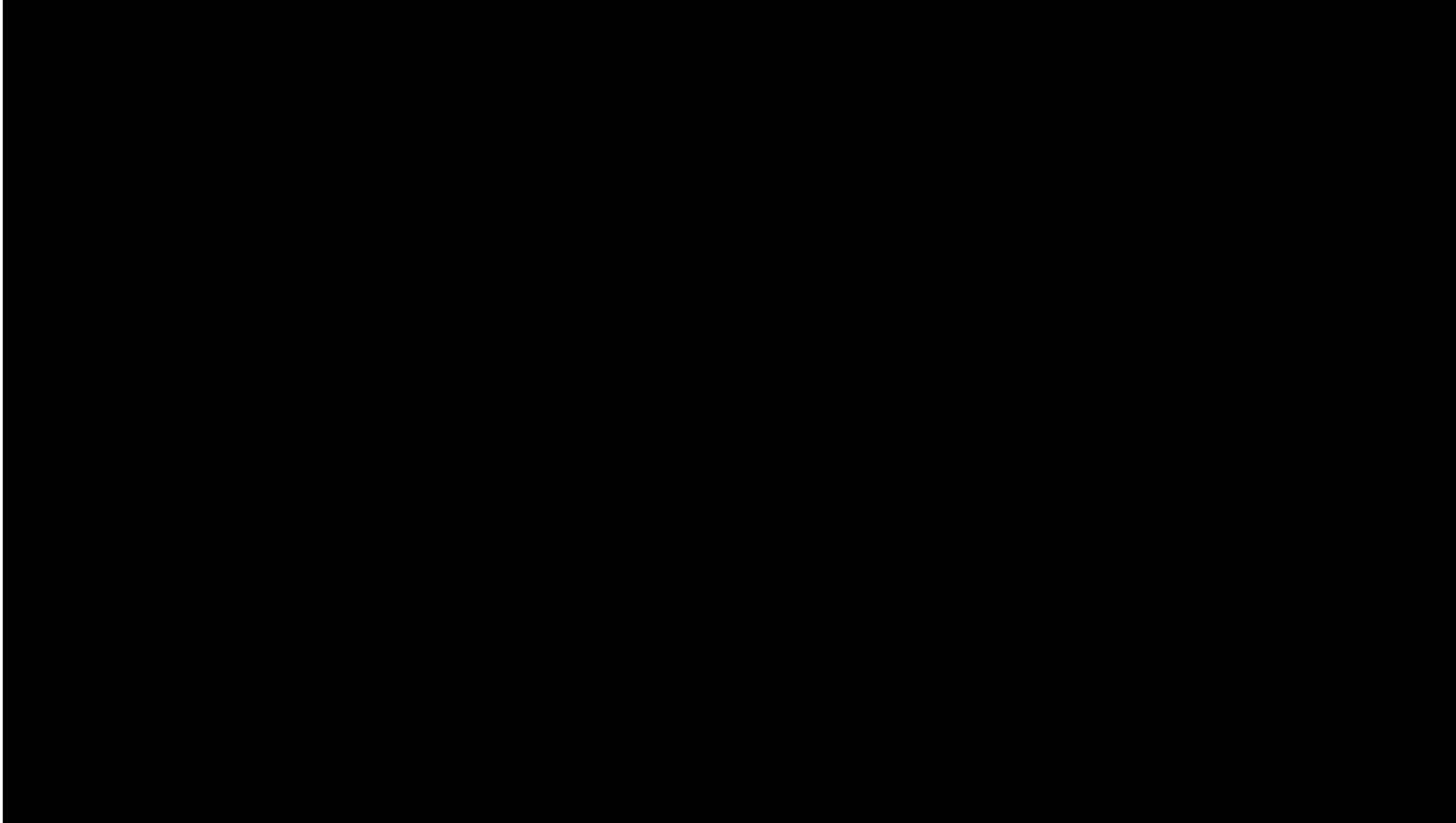
- A powerful tool that helps explaining models with correlated features
- Allows taking into consideration the whole groups of dependent variables;
- Model agnostic and it works for local and global explanations;



Use-case on the Titanic dataset is under  
<https://mi2datalab.github.io/XAI-tools/tripplot.html>



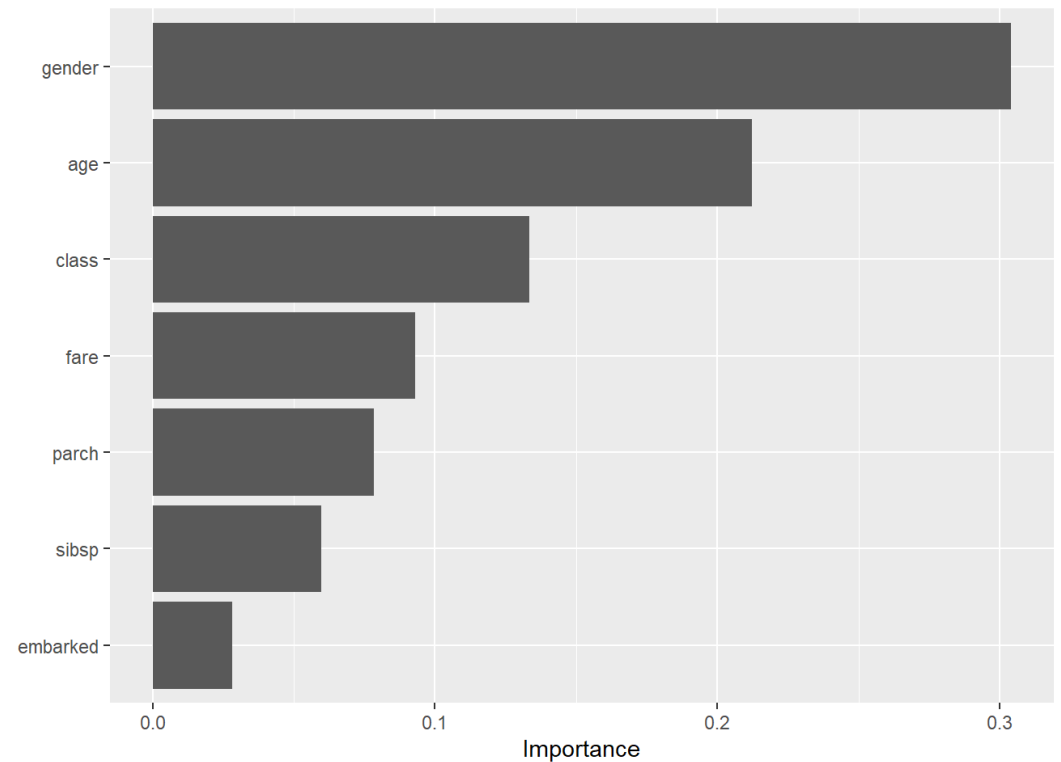
# tripplot



# vip

- Many different ways of computing variable importance, including, permutational feature importance, Shapley-based variable importance, and Variance-based variable importance.
- Offers both, model agnostics and model specific type of explanations, which is unusual in R packages for XAI.

Use-case on the Titanic dataset is under  
<https://mi2datalab.github.io/XAI-tools/vip.html>



# Takeouts

- There are a lot of R and Python packages for XAI with various functionalities.
- There is no single tool that will always be the best.
- We presented **just 11**, for more see our preprint:

[xai-toolkits.drwhy.ai](https://xai-toolkits.drwhy.ai)

You will see comparison of 27 packages and 6 python libraries for XAI. For example:

- What types of explanations the packages contain.
- What ML framework packages are compatible with.
- How to use each package.
- And more 😊



# Teaser

	Package	Global Explanations			Local Explanations		
		Model parts	Model profile	Model diagnostics	Predict parts	Predict profile	Predict diagnostics
R	ALEPlot	-	✓	-	-	-	-
	auditor	-	-	✓	-	-	-
	DALEX/DALEXtra	✓	✓	✓	✓	✓	✓
	EIX	✓	-	✓	✓	-	-
	ExplainPrediction	✓	-	-	✓	-	-
	fairness	-	-	✓	-	-	-
	fastshap	✓	✓	-	✓	-	-
	flashlight	✓	✓	-	✓	✓	-
	forestmodel	✓	-	-	-	-	-
	fscaret	✓	-	-	-	-	-
	ICEbox	-	✓	-	-	-	-
	iml	✓	✓	-	✓	✓	-
	lime	-	-	-	✓	-	-
	live	-	-	-	-	✓	-
	mcr	-	-	✓	-	-	-
	modelDown	✓	✓	✓	-	-	-
	modelStudio	✓	✓	-	✓	✓	-
	pdp	✓	✓	-	-	-	-
	randomForestExplainer	✓	-	-	-	-	-
	shapper	-	-	-	✓	-	-
	smbinning	✓	-	✓	-	-	-
	survxai	✓	-	✓	✓	✓	-
	vip	✓	-	-	-	-	-
	vivo	✓	-	-	✓	-	-
Python	aix360	✓	✓	✓	✓	-	-
	eli5	✓	-	-	✓	-	-
	interpret	✓	✓	-	✓	-	-
	lime	-	-	-	✓	-	-
	shap	✓	✓	-	✓	-	-
	skater	✓	✓	-	✓	-	-

		R				Python		Java
	Package	mlr	mlr3	parsnip	caret	keras	scikit-learn	h2o
R	ALEPlot	★	★	★	★	●	●	★
	auditor	★	★	✓	✓	●	●	★
	DALEX	★	★	✓	✓	●	●	★
	DALEXtra	✓	✓	✓	✓	✓	✓	✓
	EIX <sup>1</sup>	-	-	-	-	-	-	-
	ExplainPrediction	★	★	★	★	●	●	★
	fairness	★	★	★	★	●	●	★
	fastshap	★	★	★	★	●	●	★
	flashlight	★	★	★	★	●	●	★
	forestmodel <sup>2</sup>	-	-	-	-	-	-	-
	fscaret	-	-	-	✓	-	-	-
	iBreakDown	★	★	✓	✓	●	●	★
	ICEbox	★	★	★	★	●	●	★
	iml	✓	★	★	✓	●	●	★
	ingredients	★	★	✓	✓	●	●	★
	lime	✓	★	✓	✓	●	●	✓
	live	★	★	★	★	●	●	★
	mcr <sup>3</sup>	-	-	-	-	-	-	-
	modelDown	★	★	✓	✓	●	●	★
	modelStudio	✓	✓	✓	✓	✓	✓	✓
	pdp	★	★	★	★	●	●	★
	randomForestExplainer <sup>4</sup>	-	-	-	-	-	-	-
	shapper	★	★	★	★	●	●	★
	smbinning <sup>5</sup>	-	-	-	-	-	-	-
	survxai	★	★	★	-	-	-	-
	vip	★	★	✓	✓	●	●	★
	vivo	★	★	✓	✓	●	●	★
Python	aix360 <sup>6</sup>	-	-	-	-	-	✓	-
	eli5	-	-	-	-	✓	✓	-
	interpret <sup>7</sup>	-	-	-	-	-	✓	-
	lime	-	-	-	-	✓	✓	★
	shap	-	-	-	-	✓	✓	★
	skater	●	●	●	●	✓	✓	★

## **pdp: Partial Dependence Plots**

A general framework for constructing partial dependence (i.e., marginal effect) plots from various types machine learning models in R.

## **lime: Local Interpretable Model-Agnostic Explanations**

When building complex models, it is often difficult to explain why the model should be trusted. While global measures such as accuracy are useful, they cannot be used for explaining why a model made a specific prediction. 'lime' (a port of the 'lime' 'Python' package) is a method for explaining the outcome of black box models by fitting a local model around the point in question and perturbations of this point. The approach is described in more detail in the article by Ribeiro et al. (2016) <[arXiv:1602.04938](#)>.

## **DALEX: moDeL Agnostic Language for Exploration and eXplanation**

Unverified black box model is the path to the failure. Opacity leads to distrust. Distrust leads to ignorance. Ignorance leads to rejection. DALEX package x-rays any model and helps to explore and explain its behaviour. Machine Learning (ML) models are widely used and have various applications in classification or regression. Models created with boosting, bagging, stacking or similar techniques are often used due to their high performance. But such black-box models usually lack of direct interpretability. DALEX package contains various methods that help to understand the link between input variables and model output. Implemented methods help to explore model on the level of a single instance as well as a level of the whole dataset. All model explainers are model agnostic and can be compared across different models. DALEX package is the cornerstone for 'DrWhy.AI' universe of packages for visual model exploration. Find more details in (Biecek 2018) <[arXiv:1806.08915](#)>.

