

`.explain_predict(X)`

The Objective

- A machine learning model:
 - *model = ...* → *fit(X,y)*
- A prediction:
 - *model.predict(X)* → „Iris-virginica“
- A prediction that explains itself:
 - *model.predict_explain(X)* → „The prediction is Iris-virginica, because ...“

Start with KNeighborsClassifier

- A personalized machine learning model:
 - *model = my_KneighborsClassifier.fit(X,y)*
- *Get prediction that explains itself:*
 - *model.predict_explain(X) → Prediction
Confidence
Explanation
Features_Distribution*

model.predict_explain(X)

- Prediction: *„Iris-versicolor“*
- Confidence: *False*

model.predict_explain(X)

- *Explanation:*
 - „The prediction 'Iris-virginica' is rather unsure:

model.predict_explain(X)

- *Explanation:*
 - *„The prediction 'Iris-virginica' is rather unsure:*
 - *On the one hand the 5 nearest neighbours have diverse target values (2x value 'Iris-versicolor', 3x value 'Iris-virginica').*

model.predict_explain(X)

- *Explanation:*
 - „The prediction 'Iris-virginica' is rather unsure:
 - On the one hand the 5 nearest neighbours have diverse target values (2x value 'Iris-versicolor', 3x value 'Iris-virginica').
 - But on the other hand the nearest neighbour has the same target value too."

model.predict_explain(X)

- Features_Distribution:
 - *The features given for predicting the target value are rather far from any other observations already known.*

model.predict_explain(X)

- Features_Distribution:
 - *The features given for predicting the target value are rather far from any other observations already known.*
 - *No feature has the exact same values in the range of the 5 nearest neighbours.*

model.predict_explain(X)

- Features_Distribution:
 - *The features given for predicting the target value are rather far from any other observations already known.*
 - *No feature has the exact same values in the range of the 5 nearest neighbours.*
 - *However, the feature 'sepal_length' differs remarkably ('5.6' vs. '6.0') throughout the inspected 5 nearest neighbours.*

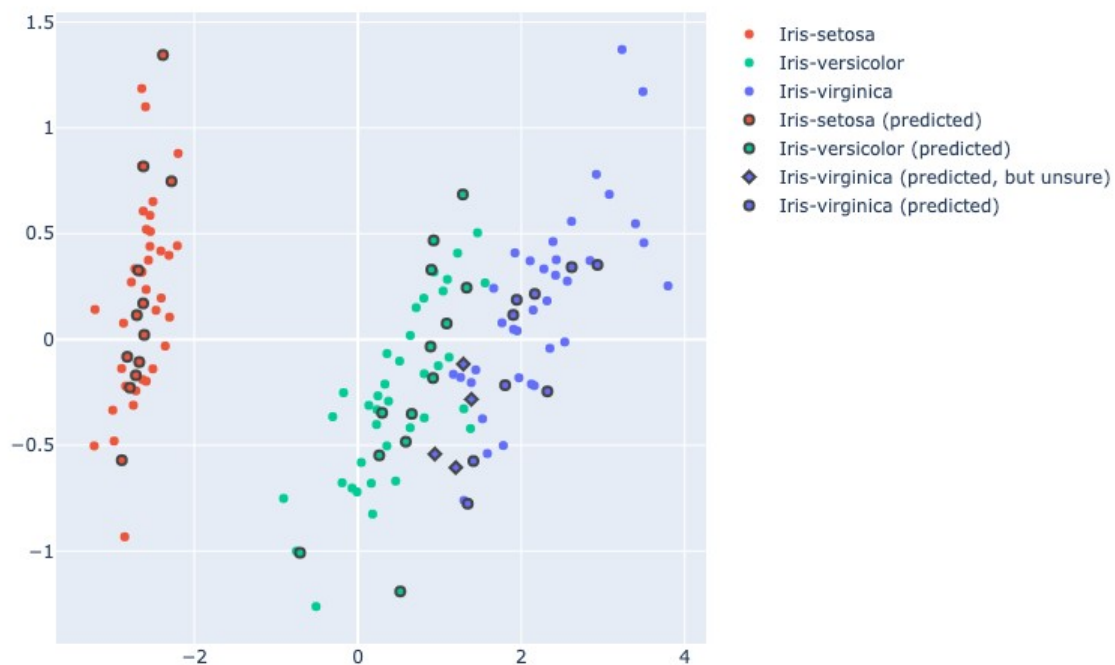
model.predict_explain(X)

- Features_Distribution:
 - *The features given for predicting the target value are rather far from any other observations already known.*
 - *No feature has the exact same values in the range of the 5 nearest neighbours.*
 - *However, the feature 'sepal_length' differs remarkably ('5.6' vs. '6.0') throughout the inspected 5 nearest neighbours.*
 - *There seems to be an intersection of the target values {'Iris-versicolor', 'Iris-virginica'}."*

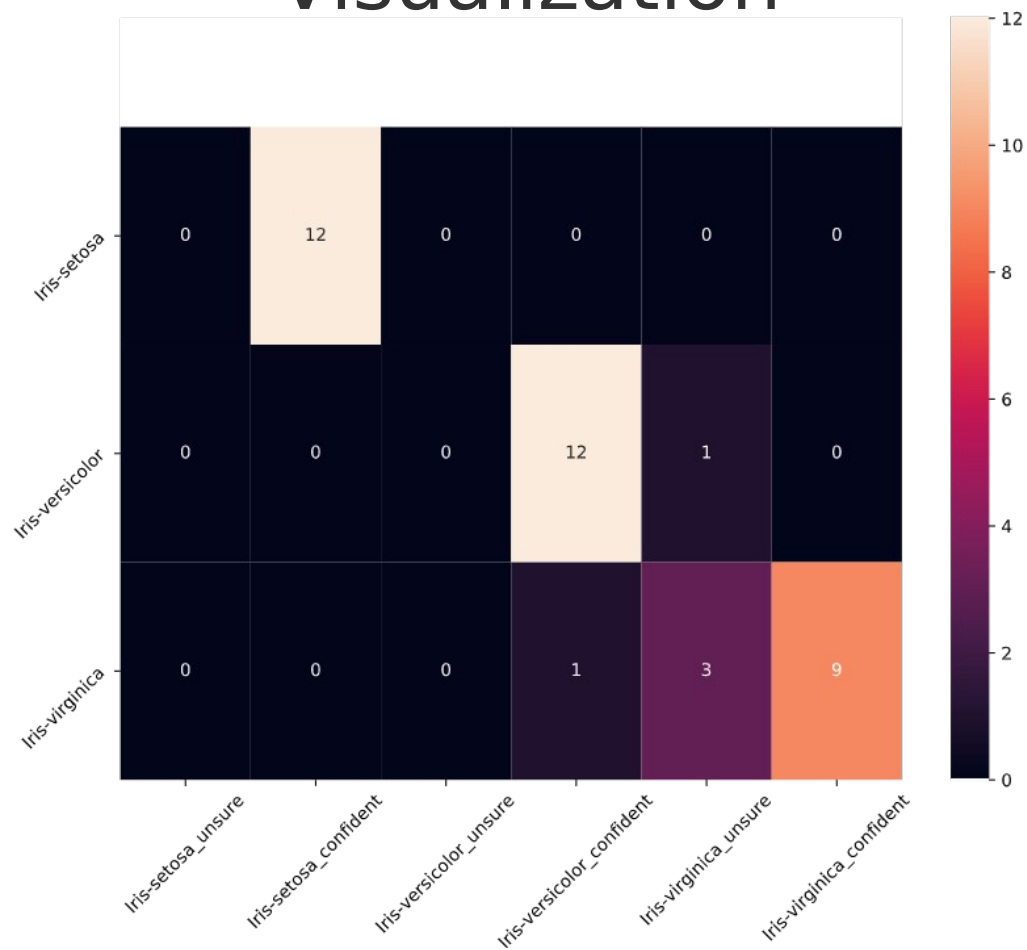
Visualization

Interactive

Dimensionality reduction for y_predict_explain: PCA visualization



Visualization



Libraries

```
from eli5.sklearn import PermutationImportance
```

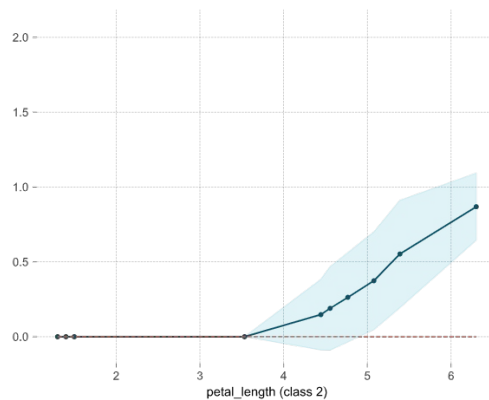
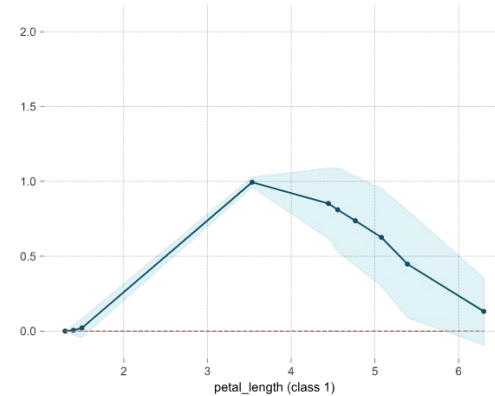
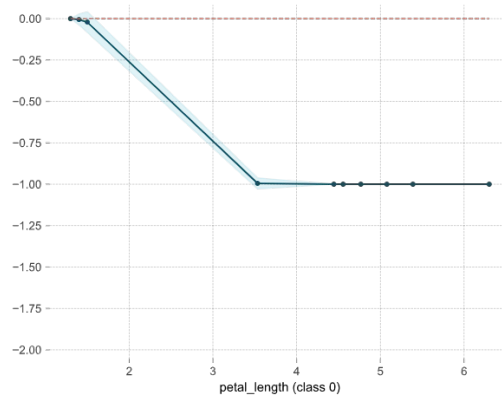
Weight	Feature
0.6316 ± 0.2183	petal_length
0.1000 ± 0.0394	petal_width
0.0158 ± 0.0537	sepal_width
-0.0000 ± 0.0744	sepal_length

Libraries

PDP for feature "petal_length"

Number of unique grid points: 10

```
from pdpbox import pdp
```



Libraries

```
import shap
```

Using 112 background data samples could cause slower run times. Consider using `shap.sample(data, K)` or `shap.kmeans(data, K)` to summarize the background as K samples.
Iris-virginica



Libraries

```
import lime, lime.lime_tabular
```

Iris-versicolor

Prediction probabilities

Iris-virginica	0.00
Iris-versicolor	1.00
Iris-setosa	0.00

NOT Iris-versicolor

Iris-versicolor

1.60 < petal_length <=...	0.52
sepal_length <= 5.10	0.08
0.30 < petal_width <=...	0.05
sepal_width <= 2.80	0.02

Näxt steps (?)

- Rollout `.predict_explain(X)` on `RandomForestClassifier`
- Dive into explainability models in more detail
- Try it out on a data set
 - Get one from the other groups
 - Get a new one

Merci!