

# Interpretability

ACTL3143 & ACTL5111 Deep Learning for Actuaries  
Patrick Laub



# Lecture Outline

- **Interpretability**
- Aspects of Interpretability
- Inherent Interpretability
- Post-hoc Interpretability
- Illustrative Example
- Illustrative Example (Fixed)



# Interpretability and Trust

Suppose a neural network informs us to increase the premium for Bob.

- Why are we getting such a conclusion from the neural network, and should we trust it?
- How can we explain our pricing scheme to Bob and the regulators?
- Should we be concerned with moral hazards, discrimination, unfairness, and ethical affairs?

We need to trust the model to employ it! With interpretability, we can trust it!



# Interpretability

## Interpretability Definition

*Interpretability refers to the ease with which one can understand and comprehend the model's algorithm and predictions.*

Interpretability of black-box models can be crucial to ascertaining trust.



# Lecture Outline

- Interpretability
- **Aspects of Interpretability**
- Inherent Interpretability
- Post-hoc Interpretability
- Illustrative Example
- Illustrative Example (Fixed)



# First Dimension of Interpretability

## Inherent Interpretability

*The model is interpretable by design.*

## Post-hoc Interpretability

*The model is not interpretable by design, but we can use other methods to explain the model.*



# Second Dimension of Interpretability

Global Interpretability:

- *The ability to understand how the model works.*
- Example: how each feature impacts the overall mean prediction.

Local Interpretability:

- *The ability to interpret/understand each prediction.*
- Example: how Bob's mean prediction has increased the most.



# Lecture Outline

- Interpretability
- Aspects of Interpretability
- **Inherent Interpretability**
- Post-hoc Interpretability
- Illustrative Example
- Illustrative Example (Fixed)





Perspective | [Published: 13 May 2019](#)

# Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead

[Cynthia Rudin](#) 

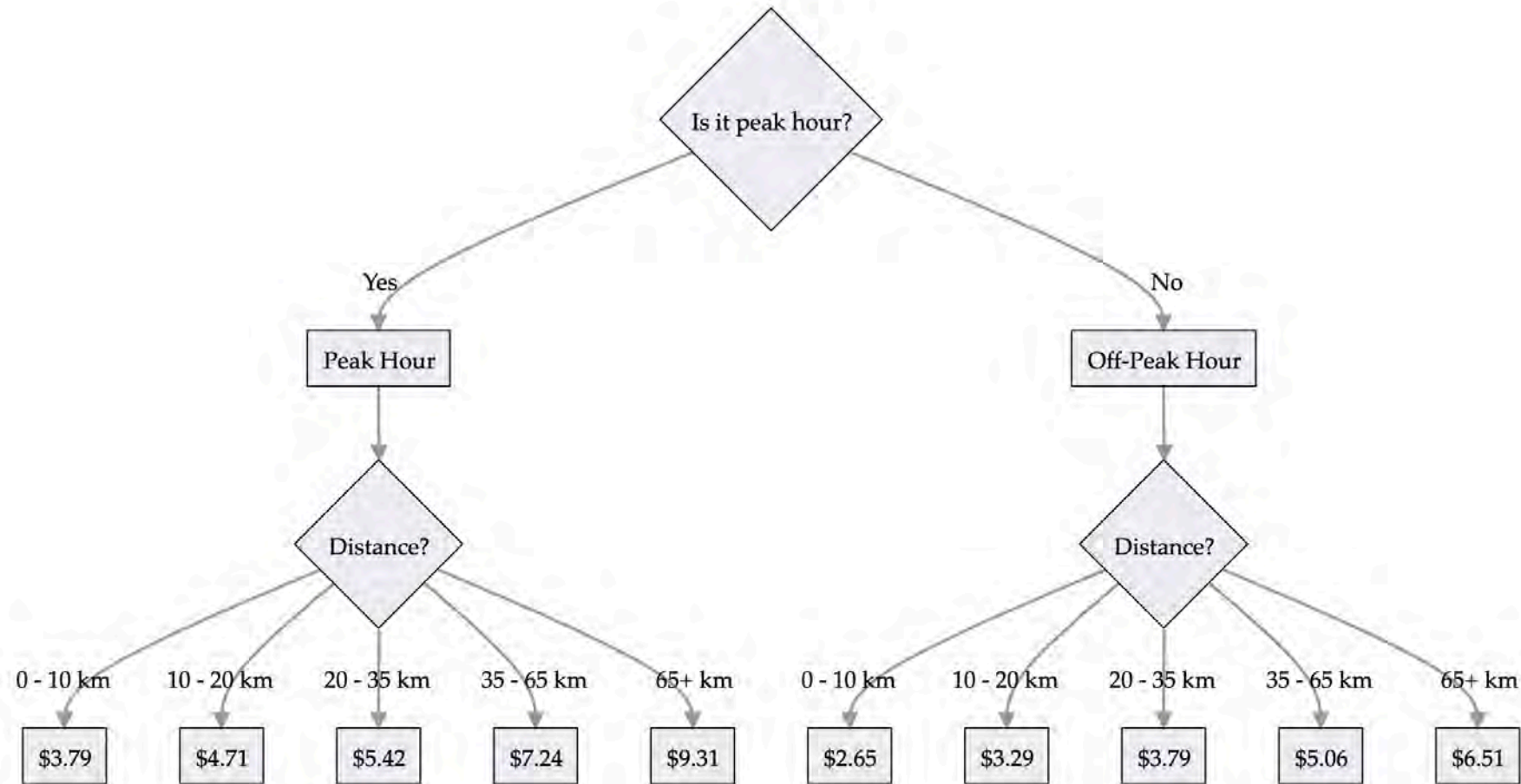
[Nature Machine Intelligence](#) 1, 206–215 (2019) | [Cite this article](#)

66k Accesses | 2230 Citations | 485 Altmetric | [Metrics](#)

## Abstract

Black box machine learning models are currently being used for high-stakes decision making throughout society, causing problems in healthcare, criminal justice and other domains. Some people hope that creating methods for explaining these black box models will alleviate some of the problems, but trying to explain black box models, rather than creating models that are interpretable in the first place, is likely to perpetuate bad practice and can potentially cause great harm to society. The way forward is to design models that are inherently interpretable. This Perspective clarifies the chasm between explaining black boxes and using inherently interpretable models, outlines several key reasons why explainable black boxes should be avoided in high-stakes decisions, identifies challenges to interpretable machine learning, and provides several example applications where interpretable models could potentially replace black box models in criminal justice, healthcare and computer vision.

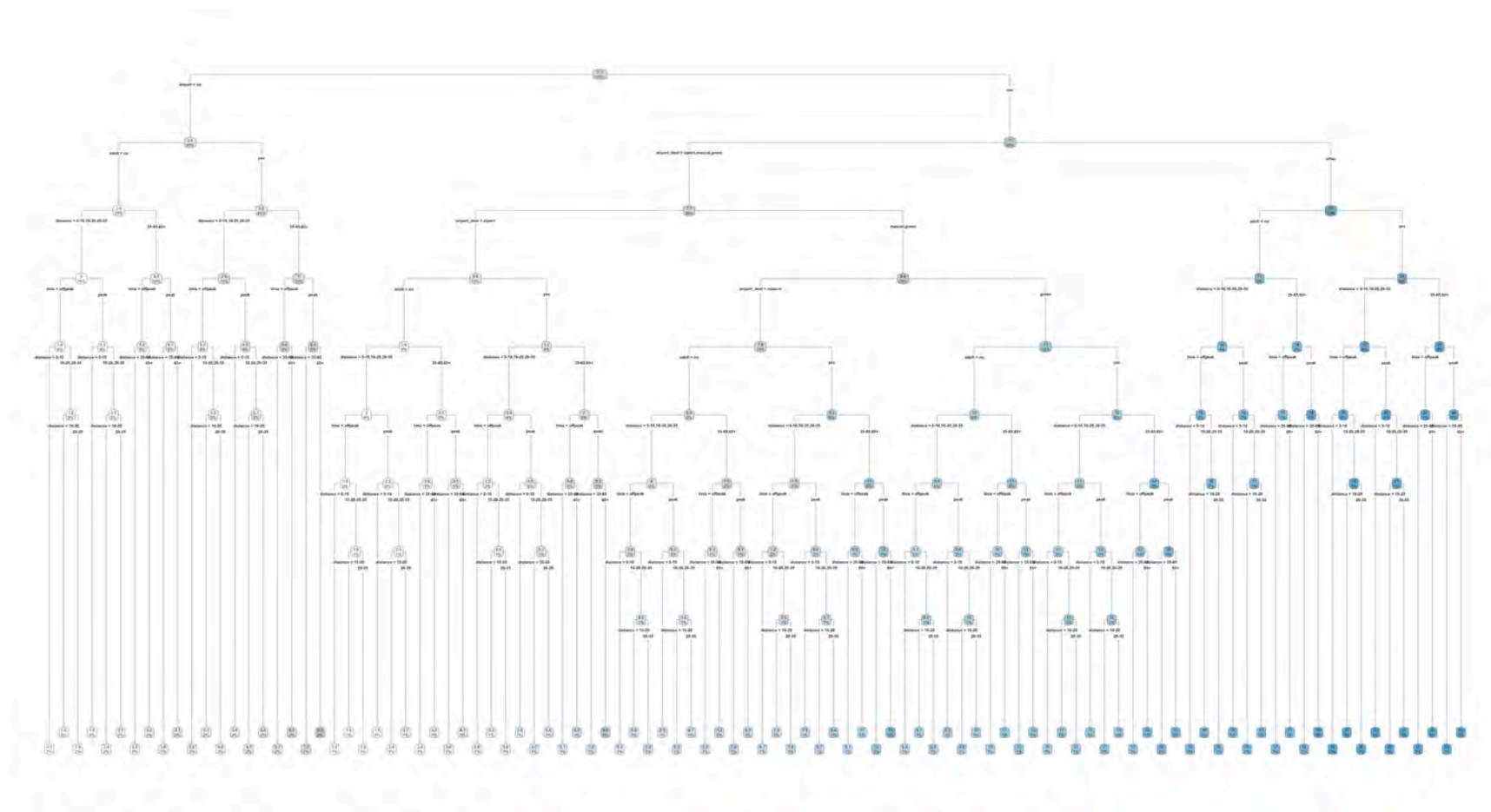
# Trees are interpretable!



Train prices



# Trees are interpretable?



# Linear models

A GLM has the form

$$\hat{y} = g^{-1}(\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p)$$

where  $\beta_0, \dots, \beta_p$  are the model parameters.

Global & local interpretations are easy to obtain.



# LocalGLMNet

Imagine:

$$\hat{y}_i = g^{-1}(\beta_0(\mathbf{x}_i) + \beta_1(\mathbf{x}_i)x_{i1} + \cdots + \beta_p(\mathbf{x}_i)x_{ip})$$

A GLM with local parameters  $\beta_0(\mathbf{x}_i), \dots, \beta_p(\mathbf{x}_i)$  for each observation  $\mathbf{x}_i$ .

The local parameters are the output of a neural network.

# Lecture Outline

- Interpretability
- Aspects of Interpretability
- Inherent Interpretability
- **Post-hoc Interpretability**
- Illustrative Example
- Illustrative Example (Fixed)



# Permutation importance

- Inputs: fitted model  $m$ , tabular dataset  $D$ .
- Compute the reference score  $s$  of the model  $m$  on data  $D$  (for instance the accuracy for a classifier or the  $R^2$  for a regressor).
- For each feature  $j$  (column of  $D$ ):
  - For each repetition  $k$  in  $1, \dots, K$ :
    - Randomly shuffle column  $j$  of dataset  $D$  to generate a corrupted version of the data named  $\tilde{D}_{k,j}$ .
    - Compute the score  $s_{k,j}$  of model  $m$  on corrupted data  $\tilde{D}_{k,j}$ .
  - Compute importance  $i_j$  for feature  $f_j$  defined as:

$$i_j = s - \frac{1}{K} \sum_{k=1}^K s_{k,j}$$

Source: scikit-learn documentation, [permutation\\_importance function](#).



# Permutation importance

```
1 def permutation_test(model, X, y, num_reps=1, seed=42):
2     """
3     Run the permutation test for variable importance.
4     Returns matrix of shape (X.shape[1], len(model.evaluate(X, y))).
5     """
6     rnd.seed(seed)
7     scores = []
8
9     for j in range(X.shape[1]):
10         original_column = np.copy(X[:, j])
11         col_scores = []
12
13         for r in range(num_reps):
14             rnd.shuffle(X[:, j])
15             col_scores.append(model.evaluate(X, y, verbose=0))
16
17         scores.append(np.mean(col_scores, axis=0))
18         X[:, j] = original_column
19
20     return np.array(scores)
```





# LIME

*Local Interpretable Model-agnostic Explanations* employs an interpretable surrogate model to explain locally how the black-box model makes predictions for individual instances.

E.g. a black-box model predicts Bob's premium as the highest among all policyholders. LIME uses an interpretable model (a linear regression) to explain how Bob's features influence the black-box model's prediction.



See “Why Should I Trust You?”: Explaining the Predictions of Any Classifier.

# Globally vs. Locally Faithful

## Globally Faithful

*The interpretable model's explanations accurately reflect the behaviour of the black-box model across the entire input space.*

## Locally Faithful

*The interpretable model's explanations accurately reflect the behaviour of the black-box model for a specific instance.*

LIME aims to construct an interpretable model that mimics the black-box model's behaviour in a *locally faithful* manner.



# LIME Algorithm

Suppose we want to explain the instance  $\mathbf{x}_{\text{Bob}} = (1, 2, 0.5)$ .

1. Generate perturbed examples of  $\mathbf{x}_{\text{Bob}}$  and use the trained gamma MDN  $f$  to make predictions:

$$\begin{aligned} \mathbf{x}_{\text{Bob}}^{(1)} &= (1.1, 1.9, 0.6), & f(\mathbf{x}_{\text{Bob}}^{(1)}) &= 34000 \\ \mathbf{x}_{\text{Bob}}^{(2)} &= (0.8, 2.1, 0.4), & f(\mathbf{x}_{\text{Bob}}^{(2)}) &= 31000 \\ \vdots & & \vdots & \end{aligned}$$

We can then construct a dataset of  $N_{\text{Examples}}$  perturbed examples:

$$\mathcal{D}_{\text{LIME}} = \left( \{ \mathbf{x}_{\text{Bob}}^{(i)}, f(\mathbf{x}_{\text{Bob}}^{(i)}) \} \right)_{i=0}^{N_{\text{Examples}}}.$$



# LIME Algorithm

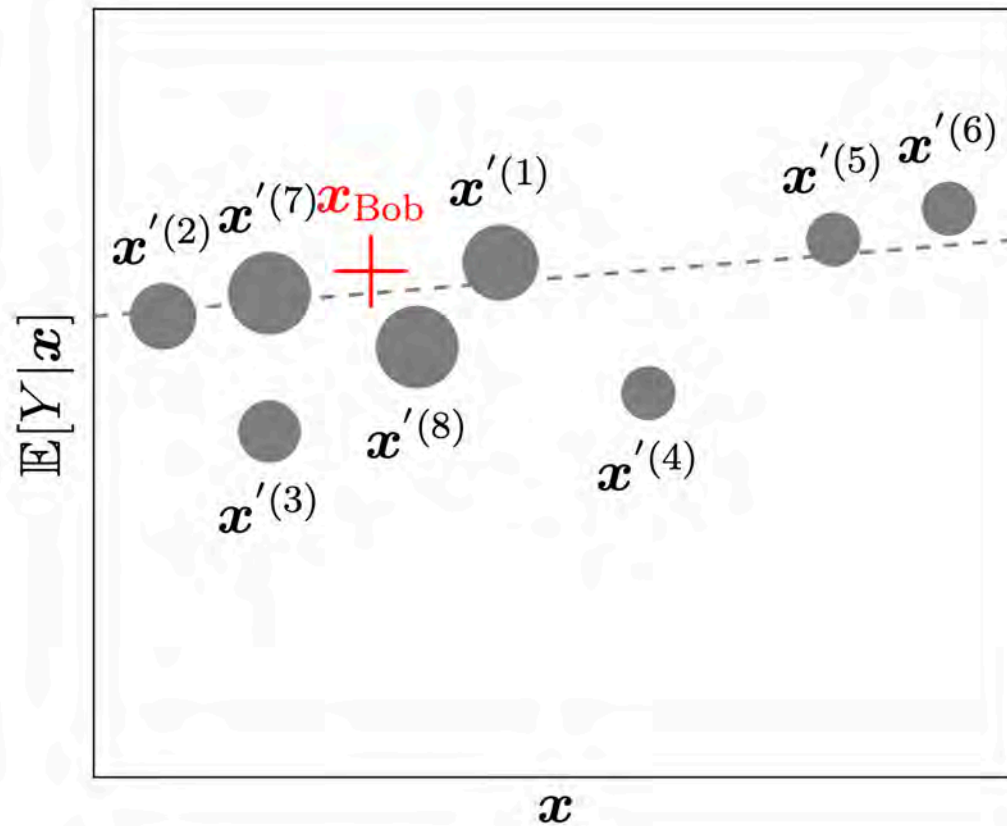
2. Fit an interpretable model  $g$ , i.e., a linear regression using  $\mathcal{D}_{\text{LIME}}$  and the following loss function:

$$\mathcal{L}_{\text{LIME}}(f, g, \pi_{\mathbf{x}_{\text{Bob}}}) = \sum_{i=1}^{N_{\text{Examples}}} \pi_{\mathbf{x}_{\text{Bob}}}(\mathbf{x}_{\text{Bob}}^{(i)}) \cdot \left( f(\mathbf{x}_{\text{Bob}}^{(i)}) - g(\mathbf{x}_{\text{Bob}}^{(i)}) \right)^2,$$

where  $\pi_{\mathbf{x}_{\text{Bob}}}(\mathbf{x}_{\text{Bob}}^{(i)})$  represents the distance from the perturbed example  $\mathbf{x}_{\text{Bob}}^{(i)}$  to the instance to be explained  $\mathbf{x}_{\text{Bob}}$ .



# “Explaining” to Bob



The bold red cross is the instance being explained. LIME samples instances (grey nodes), gets predictions using  $f$  (gamma MDN) and weighs them by the proximity to the instance being explained (represented here by size). The dashed line  $g$  is the learned local explanation.

# SHAP Values

The SHapley Additive exPlanations (SHAP) value helps to quantify the contribution of each feature to the prediction for a specific instance.

The SHAP value for the  $j$ th feature is defined as

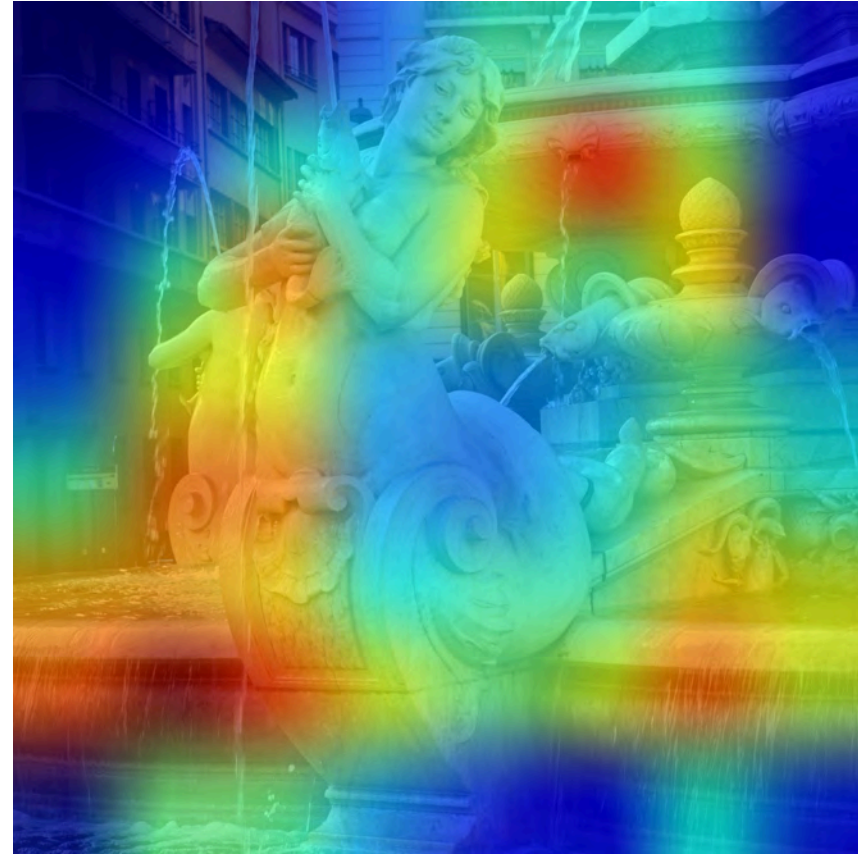
$$\text{SHAP}^{(j)}(\mathbf{x}) = \sum_{U \subset \{1, \dots, p\} \setminus \{j\}} \frac{1}{p} \binom{p-1}{|U|}^{-1} (\mathbb{E}[Y | \mathbf{x}^{(U \cup \{j\})}] - \mathbb{E}[Y | \mathbf{x}^{(U)}]),$$

where  $p$  is the number of features. A positive SHAP value indicates that the variable increases the prediction value.

# Grad-CAM



Original image



Grad-CAM

See Chollet (2021), Deep Learning with Python, Section 9.4.3.

# Lecture Outline

- Interpretability
- Aspects of Interpretability
- Inherent Interpretability
- Post-hoc Interpretability
- **Illustrative Example**
- Illustrative Example (Fixed)





# First attempt at NLP task

```
1 df_raw["SUMMARY_EN"]
```

```
0      V1, a 2000 Pontiac Montana minivan,
made a lef...
1      The crash occurred in the eastbound
lane of a ...
2      This crash occurred just after the
noon time h...
...
6946   The crash occurred in the eastbound
lanes of a ...
6947   This single-vehicle crash occurred
in a rural ...
6948   This two vehicle daytime collision
occurred mi...
Name: SUMMARY_EN, Length: 6949, dtype:
object
```

```
1 df_raw["NUM_VEHICLES"].value_counts()\
2     .sort_index()
```

```
NUM_VEHICLES
1      1822
2      4151
3+      976
Name: count, dtype: int64
```



# Bag of words for the top 1,000 words

```

1 vect = CountVectorizer(max_features=1_000, stop_words="english")
2 vect.fit(X_train["SUMMARY_EN"])
3
4 X_train_bow = vectorise_dataset(X_train, vect)
5 X_val_bow = vectorise_dataset(X_val, vect)
6 X_test_bow = vectorise_dataset(X_test, vect)
7
8 vectorise_dataset(X_train, vect, dataframe=True).head()

```

	10	105	113	12	15	150	16	17	18	180	...	yield	zoi
2532	0	0	0	0	0	0	0	0	0	0	...	0	0
6209	0	0	0	0	0	0	0	0	0	0	...	0	0
2561	1	0	1	0	0	0	0	0	0	0	...	0	0
6664	0	0	0	0	0	0	0	0	0	0	...	0	0
4214	0	0	0	0	0	0	0	0	0	0	...	0	0

5 rows × 1008 columns



# Trained a basic neural network on that

```
1 num_features = X_train_bow.shape[1]
2 num_cats = df_raw["NUM_VEHICLES"].nunique()
3 model = build_model(num_features, num_cats)
4 es = EarlyStopping(patience=1, restore_best_weights=True, monitor="val_accuracy")
```

```
1 model.fit(X_train_bow, y_train, epochs=10,
2         callbacks=[es], validation_data=(X_val_bow, y_val), verbose=0)
3 model.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 100)	100,900
dense_1 (Dense)	(None, 3)	303

Total params: 303,611 (1.16 MB)  
 Trainable params: 101,203 (395.32 KB)  
 Non-trainable params: 0 (0.00 B)  
 Optimizer params: 202,408 (790.66 KB)

```
1 model.evaluate(X_train_bow, y_train, verbose=0)
```

[0.004478050395846367, 1.0, 1.0]

```
1 model.evaluate(X_val_bow, y_val, verbose=0)
```

[8.163677215576172, 0.9784172773361206, 0.9985611438751221]



# Permutation importance algorithm

Taken directly from scikit-learn documentation:

- Inputs: fitted predictive model  $m$ , tabular dataset (training or validation)  $D$ .
- Compute the reference score  $s$  of the model  $m$  on data  $D$  (for instance the accuracy for a classifier or the  $R^2$  for a regressor).
- For each feature  $j$  (column of  $D$ ):
  - For each repetition  $k$  in  $1, \dots, K$ :
    - Randomly shuffle column  $j$  of dataset  $D$  to generate a corrupted version of the data named  $\tilde{D}_{k,j}$ .
    - Compute the score  $s_{k,j}$  of model  $m$  on corrupted data  $\tilde{D}_{k,j}$ .
  - Compute importance  $i_j$  for feature  $f_j$  defined as:

$$i_j = s - \frac{1}{K} \sum_{k=1}^K s_{k,j}$$

Source: scikit-learn documentation, [permutation\\_importance function](#).



# Find important inputs

```
1 def permutation_test(model, X, y, num_reps=1, seed=42):
2     """
3     Run the permutation test for variable importance.
4     Returns matrix of shape (X.shape[1], len(model.evaluate(X, y))).
5     """
6     rnd.seed(seed)
7     scores = []
8
9     for j in range(X.shape[1]):
10         original_column = np.copy(X[:, j])
11         col_scores = []
12
13         for r in range(num_reps):
14             rnd.shuffle(X[:, j])
15             col_scores.append(model.evaluate(X, y, verbose=0))
16
17         scores.append(np.mean(col_scores, axis=0))
18         X[:, j] = original_column
19
20     return np.array(scores)
```



# Run the permutation test

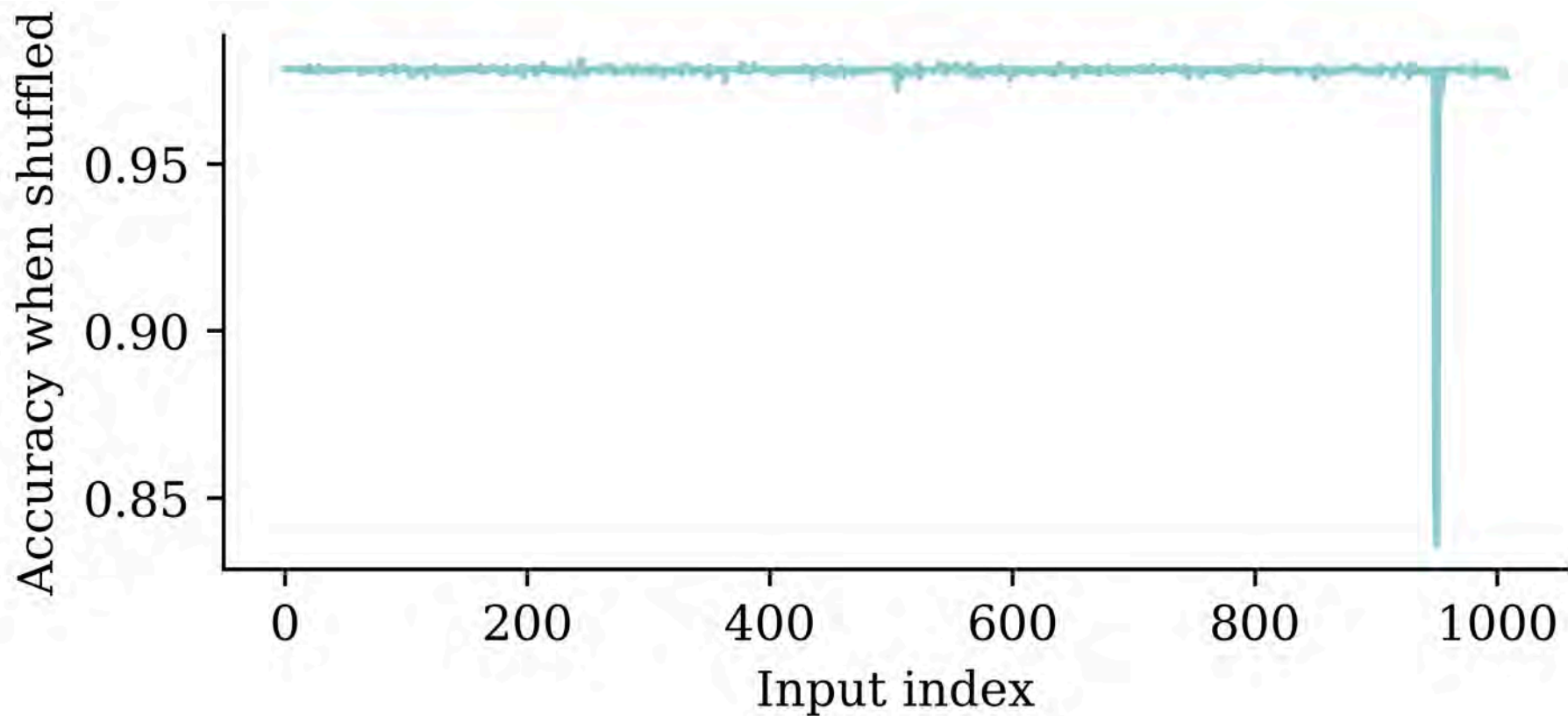
```
1 all_perm_scores = permutation_test(model, X_val_bow, y_val)
2 all_perm_scores
```

```
array([[ 8.16,  0.98,  1.  ],
       [ 8.16,  0.98,  1.  ],
       [ 8.16,  0.98,  1.  ],
       ...,
       [ 8.48,  0.98,  1.  ],
       [ 8.47,  0.98,  1.  ],
       [14.59,  0.98,  1.  ]])
```



# Plot the permuted accuracies

```
1 perm_scores = all_perm_scores[:,1]
2 plt.plot(perm_scores)
3 plt.xlabel("Input index")
4 plt.ylabel("Accuracy when shuffled");
```



# Find the most significant inputs

```

1 vocab = vect.get_feature_names_out()
2 input_cols = list(vocab) + weather_cols
3
4 best_input_inds = np.argsort(perm_scores)[:100]
5 best_inputs = [input_cols[idx] for idx in best_input_inds]
6
7 print(best_inputs)

```

```

['v3', 'v2', 'vehicle', 'involved', 'event', 'harmful', 'stated', 'motor', 'WEATHER8',
'left', 'v1', 'just', 'turn', 'traveling', 'approximately', 'medication', 'hurry',
'encroaching', 'chevrolet', 'parked', 'continued', 'saw', 'road', 'rest', 'distraction',
'ahead', 'wet', 'hit', 'WEATHER5', 'WEATHER6', 'WEATHER4', 'old', 'v4', '48', 'rested',
'direction', 'occurred', 'kph', 'clear', 'right', 'miles', 'uphill', 'WEATHER3', 'denied',
'80', 'attempt', 'thinking', 'assumption', 'damage', 'runs', 'time', 'consisted',
'following', 'towed', 'WEATHER1', 'alcohol', 'mph', 'pickup', 'lane', 'conversing', 'make',
'started', 'maneuver', 'stopped', 'store', 'car', 'local', 'dry', 'median', 'south',
'driver', 'higher', 'pre', 'northbound', 'impacting', 'congested', 'health', 'southwest',
'gmc', 'observed', 'parking', 'partially', 'heart', 'shoulders', 'shoulder', 'southeast',
'came', 'heard', 'gap', 'southbound', 'conditions', 'contacted', 'change', 'compact', '2002',
'cell', 'causing', 'oldsmobile', 'half', 'hard']

```





# How about a simple decision tree?

```
1 from sklearn import tree
2
3 clf = tree.DecisionTreeClassifier(random_state=0, max_leaf_nodes=3)
4 clf.fit(X_train_bow[:, best_input_inds], y_train);
```

```
1 print(clf.score(X_train_bow[:, best_input_inds], y_train))
2 print(clf.score(X_val_bow[:, best_input_inds], y_val))
```

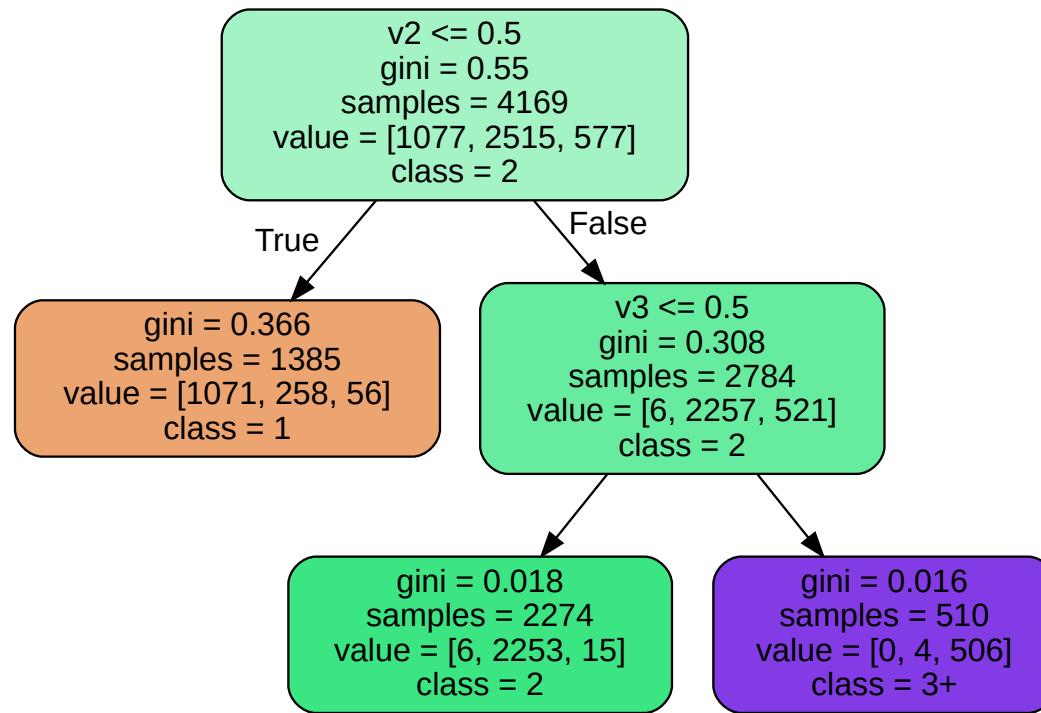
```
0.9186855360997841
0.9316546762589928
```

The decision tree ends up giving pretty good results.



# Decision tree

```
1 tree.plot_tree(clf, feature_names=best_inputs, filled=True);
```



```
1 print(np.where(clf.feature_importances_ > 0)[0])
2 [best_inputs[ind] for ind in np.where(clf.feature_importances_ > 0)[0]]
```

```
[0 1]
```

```
['v3', 'v2']
```



# Lecture Outline

- Interpretability
- Aspects of Interpretability
- Inherent Interpretability
- Post-hoc Interpretability
- Illustrative Example
- **Illustrative Example (Fixed)**



# This is why we replace “v1”, “v2”, “v3”

There was a slide in the NLP deck titled “Just ignore this for now...” That was going through each summary and replacing the words “V1”, “V2”, “V3” with random numbers. This was done to see if the model was overfitting to these words.

```
1 model.fit(X_train_bow, y_train, epochs=10,  
2         callbacks=[es], validation_data=(X_val_bow, y_val), verbose=0);
```

Retraining on the fixed dataset gives us a more realistic (lower) accuracy.

```
1 model.evaluate(X_train_bow, y_train, verbose=0)
```

```
[0.07491350173950195, 0.986567497253418, 0.9997601509094238]
```

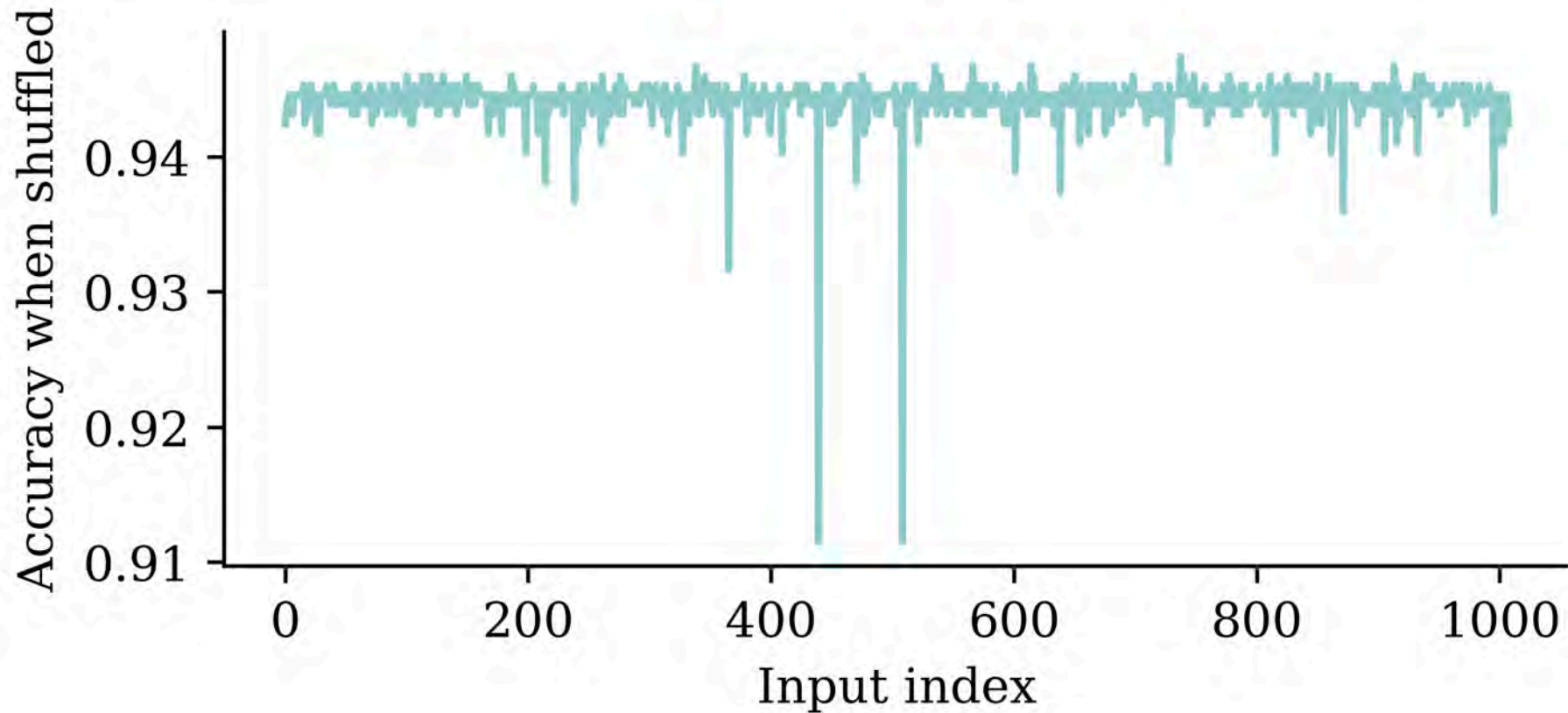
```
1 model.evaluate(X_val_bow, y_val, verbose=0)
```

```
[2.865464925765991, 0.9446043372154236, 0.9956834316253662]
```



# Permutation importance accuracy plot

```
1 perm_scores = permutation_test(model, X_val_bow, y_val)[: ,1]  
2 plt.plot(perm_scores)  
3 plt.xlabel("Input index"); plt.ylabel("Accuracy when shuffled");
```



# Find the most significant inputs

```

1 vocab = vect.get_feature_names_out()
2 input_cols = list(vocab) + weather_cols
3
4 best_input_inds = np.argsort(perm_scores)[:100]
5 best_inputs = [input_cols[idx] for idx in best_input_inds]
6
7 print(best_inputs)

```

```

['harmful', 'involved', 'event', 'year', 'struck', 'contacted', 'old', 'impact', 'coded',
'motor', 'rear', 'towed', 'stop', 'driven', 'turning', 'single', 'forward', 'chevrolet',
'lanes', 'crash', 'parked', 'continued', 'WEATHER4', 'WEATHER1', 'travel', 'divided',
'brakes', 'include', 'came', 'stopped', 'final', 'factor', 'clear', '2004', 'speed', '2002',
'reason', 'passing', 'pickup', 'crossing', 'driving', 'ahead', 'did', 'control', '10',
'highway', 'WEATHER5', 'weather', 'traveling', 'surveillance', 'stated', 'spin', 'heart',
'road', 'pole', 'occupants', 'afternoon', 'moderate', 'mirror', 'including', 'pushed',
'55mph', 'intersection', 'WEATHER8', '1993', 'cross', 'curve', 'slight', 'dark', 'day',
'sight', 'distance', 'seat', 'saw', 'said', 'drivers', 'earlier', 'early', 'dodge',
'corrected', 'state', 'contributed', '16', 'trying', 'trip', 'buick', 'trailer', 'traffic',
'tractor', 'cherokee', '42', 'taken', '41', 'cloudy', 'collision', 'congested', 'contact',
'eastbound', 'roadways', 'roadway']

```



# How about a simple decision tree?

```
1 clf = tree.DecisionTreeClassifier(random_state=0, max_leaf_nodes=3)
2 clf.fit(X_train_bow[:, best_input_inds], y_train);
```

```
1 print(clf.score(X_train_bow[:, best_input_inds], y_train))
2 print(clf.score(X_val_bow[:, best_input_inds], y_val))
```

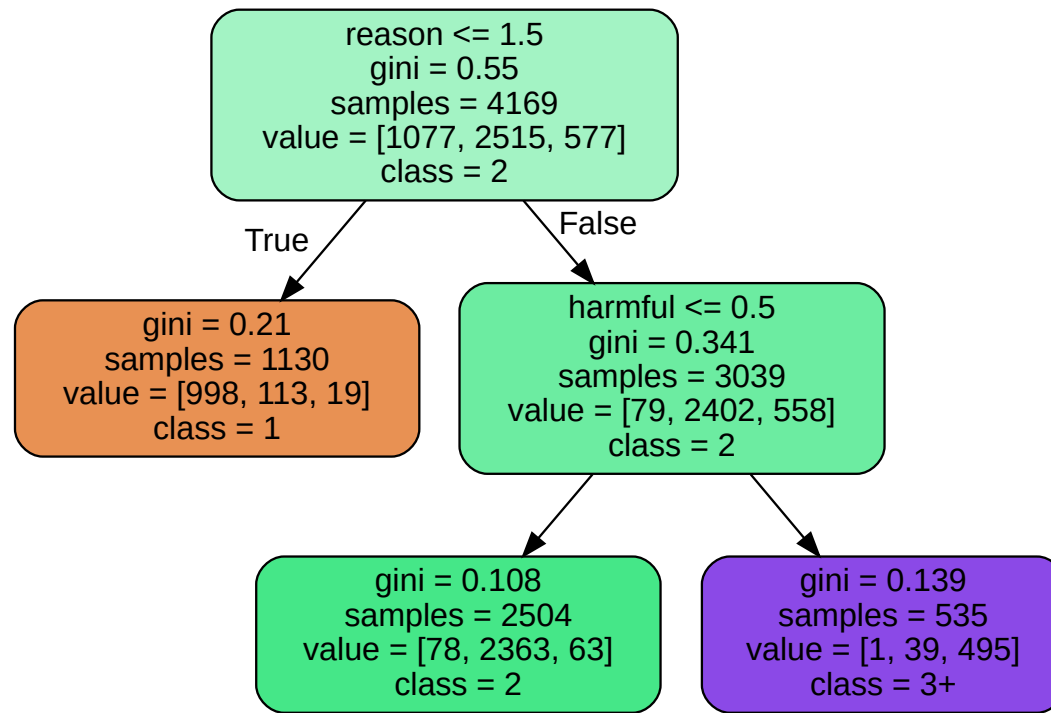
0.9249220436555529

0.9381294964028777



# Decision tree

```
1 tree.plot_tree(clf, feature_names=best_inputs, filled=True);
```



```
1 print(np.where(clf.feature_importances_ > 0)[0])
2 [best_inputs[ind] for ind in np.where(clf.feature_importances_ > 0)[0]]
```

```
[ 0 36]
```

```
['harmful', 'reason']
```





# Package Versions

```
1 from watermark import watermark
2 print(watermark(python=True, packages="keras,matplotlib,numpy,pandas,seaborn,scipy,torch"))
```

```
Python implementation: CPython
Python version       : 3.11.9
IPython version      : 8.24.0
```

```
keras      : 3.3.3
matplotlib: 3.9.0
numpy      : 1.26.4
pandas     : 2.2.2
seaborn    : 0.13.2
scipy      : 1.11.0
torch      : 2.3.1
tensorflow: 2.16.1
tf_keras   : 2.16.0
```



# Glossary

- global interpretability
- Grad-CAM
- inherent interpretability
- LIME
- local interpretability
- permutation importance
- post-hoc interpretability
- SHAP values

