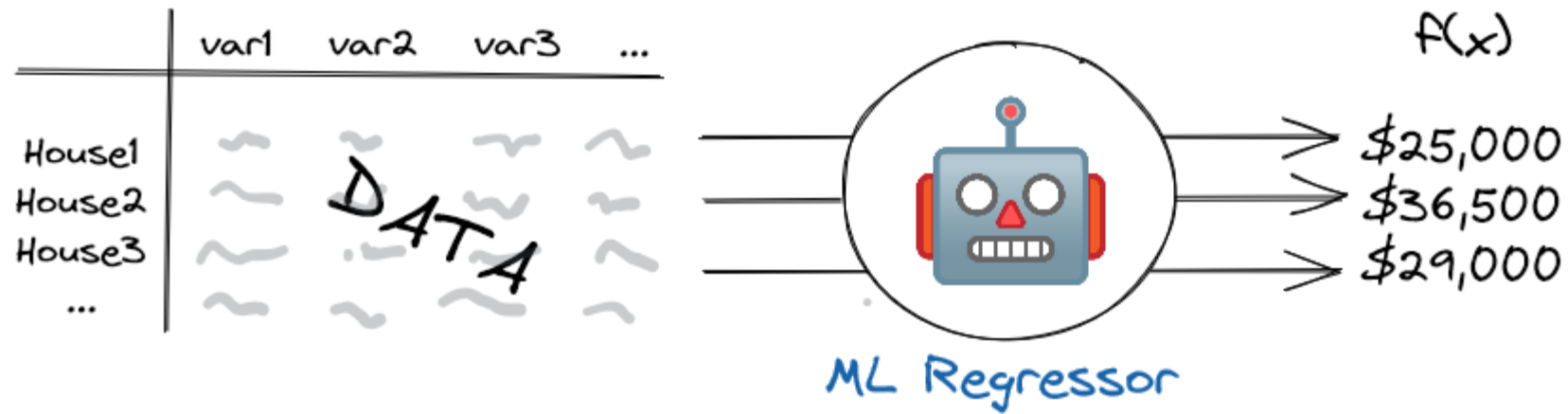
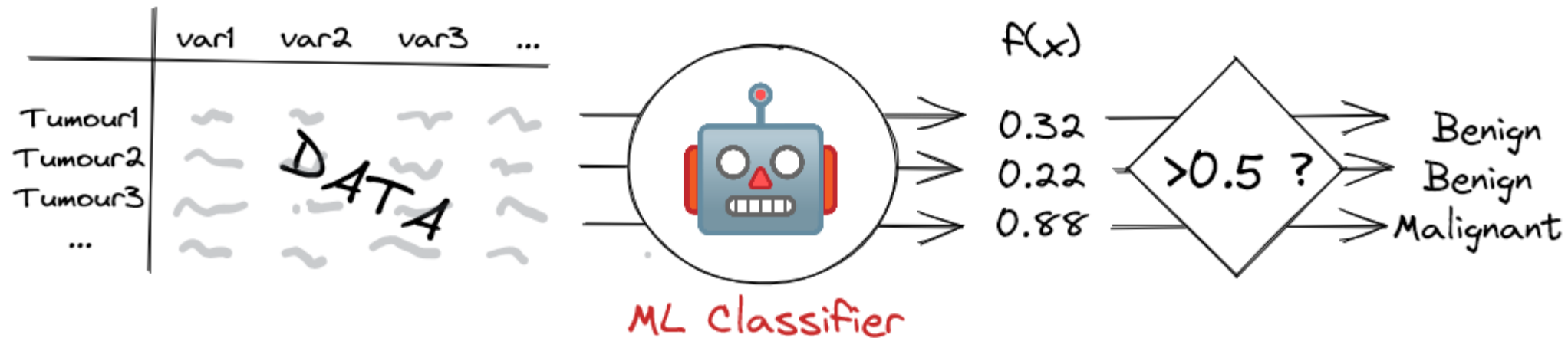


SHAP Analysis Extra Note

Regression:



Classification:



SHAP

- SHAP = *SHapley Additive exPlanations*
- Popularized use of Shapley values in ML
 - Also used in earlier work by Lipovetsky & Conklin (2001), Strumbelj et al. (2009), Datta et al. (2016)
- SHAP uses Shapley values to explain individual predictions

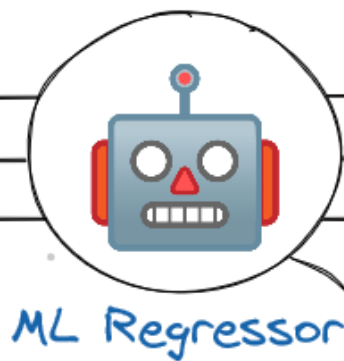
A machine learning model's prediction, $f(x)$, can be represented as the sum of its computed *SHAP values*, plus a fixed *base value*, such that:

$$f(x) = \text{base value} + \text{sum}(\text{SHAP values})$$

Regression:

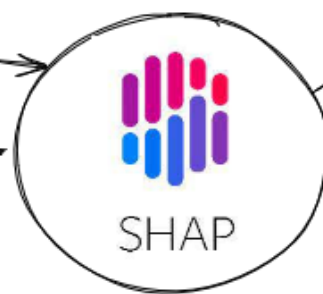
	var1	var2	var3	...
House1	~	~	~	~
House2	~	~	~	~
House3	~	~	~	~
...	~	~	~	~

DATA



$$f(x) = \text{base value} + \text{sum}(\text{SHAP values})$$

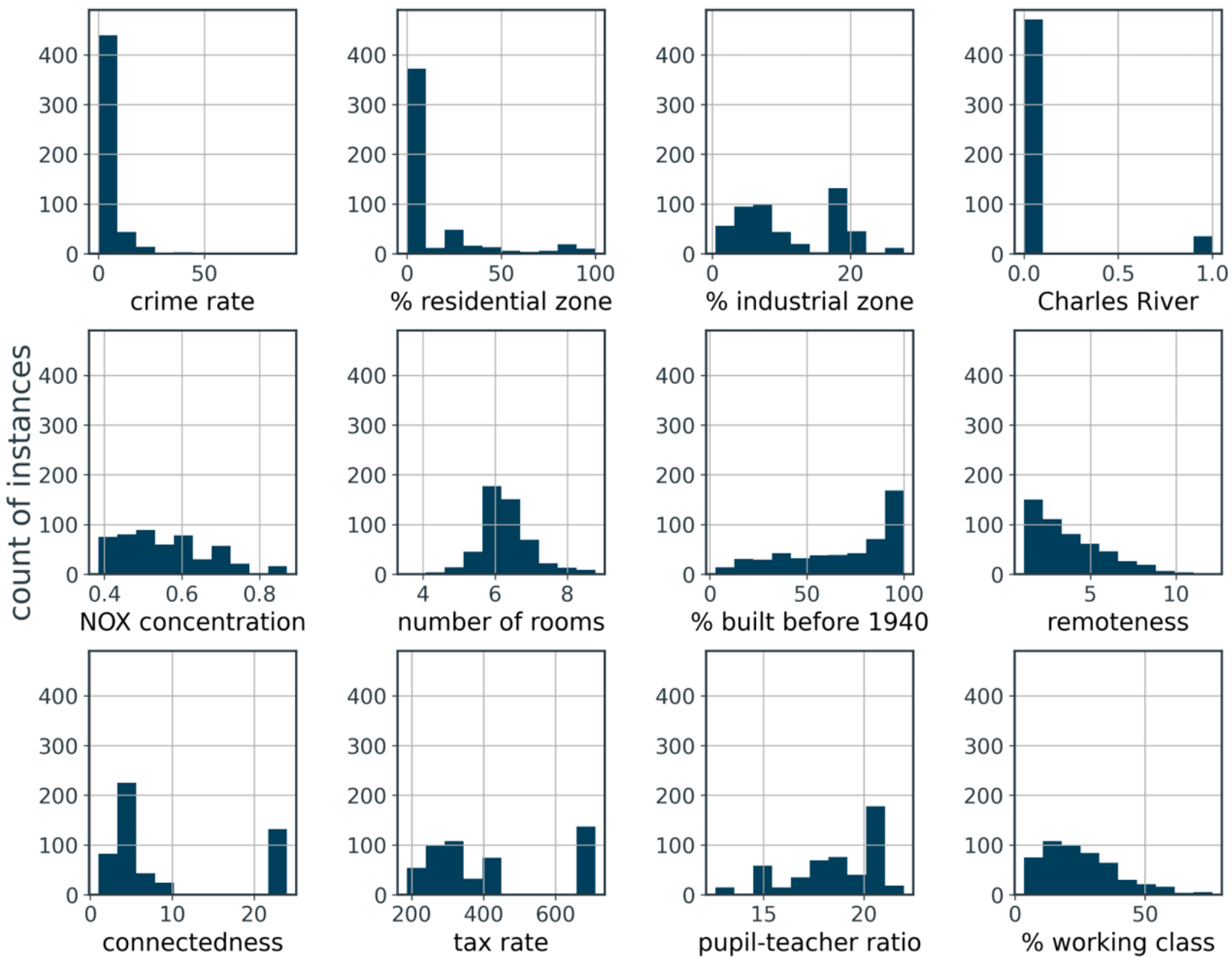
$$\begin{aligned} \$25,000 &= \$22,500 + \text{sum}(\text{House1 SHAP values}) \\ \$36,500 &= \$22,500 + \text{sum}(\text{House2 SHAP values}) \\ \$29,000 &= \$22,500 + \text{sum}(\text{House3 SHAP values}) \end{aligned}$$



	var1	var2	var3	...
House1	~	~	~	~
House2	~	~	~	~
House3	~	~	~	~
...	~	~	~	~

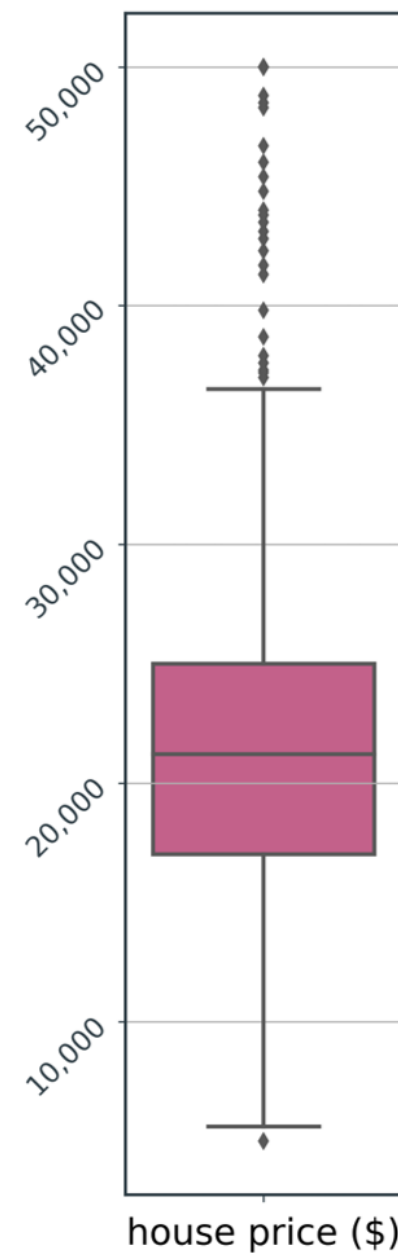
SHAP Values

Model Input Variables



predicts...

Target Variable



The model's predicted house price for this example.

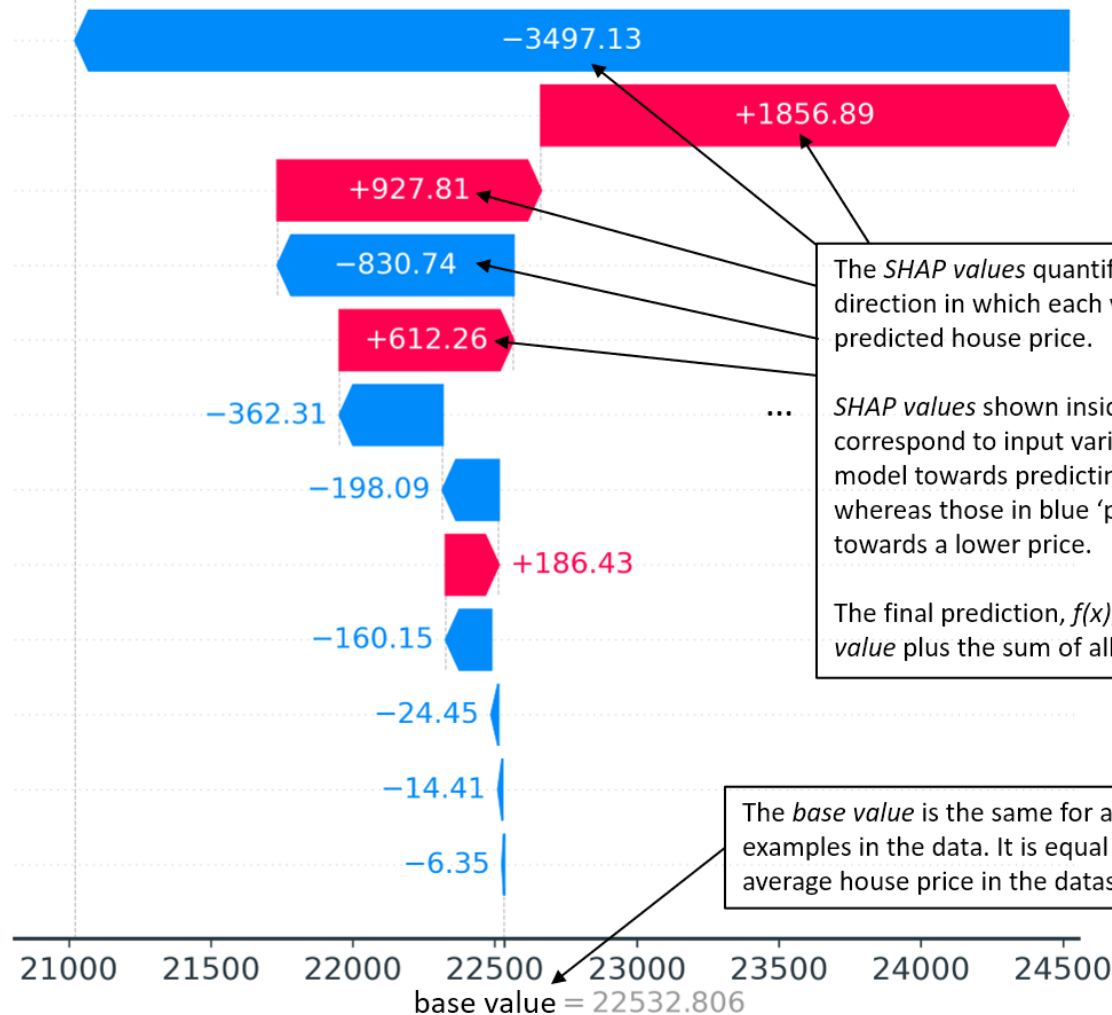
$$f(x) = 21022.57 = \text{base value} + \text{sum}(\text{SHAP values})$$

A **waterfall plot** provides a detailed breakdown of how each input variable contributes towards the predicted house price for a single instance of the data.

These are the input variables, ranked from top to bottom by how much impact they have on the model's prediction for this example from the data.

The grey numbers denote the values of the variables for this particular instance.

5.878 = number of rooms
16.2 = % working class
21.4 = % built before 1940
6.498 = remoteness
0.409 = NOX concentration
345 = tax rate
4 = connectedness
0.058 = crime rate
6.07 = % industrial zone
18.9 = pupil-teacher ratio
12.5 = % residential zone
0 = Charles River



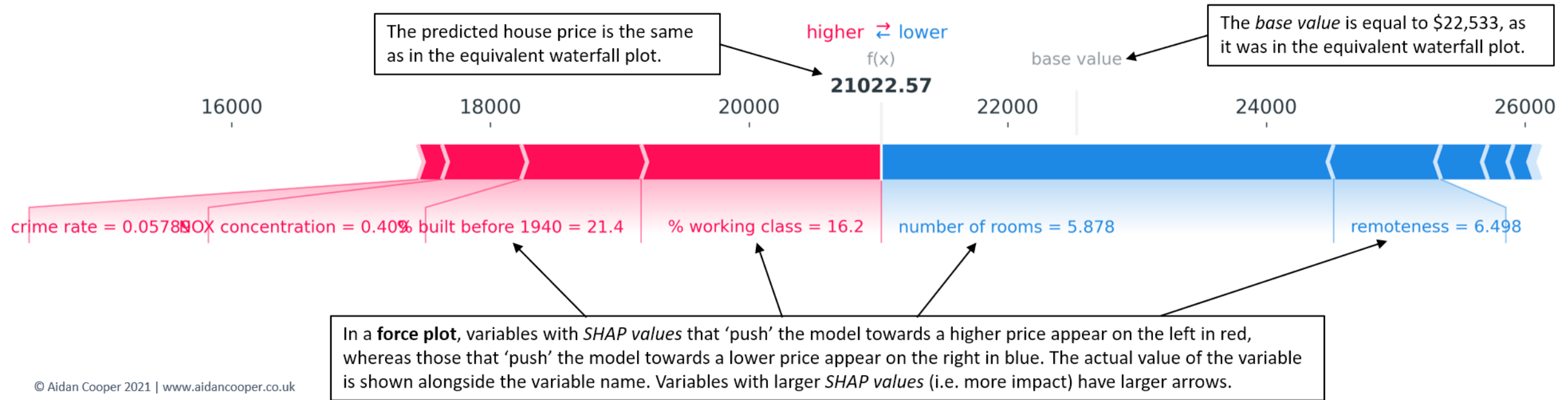
The *SHAP values* quantify the amount and direction in which each variable impacts the predicted house price.

SHAP values shown inside red arrows correspond to input variables that 'push' the model towards predicting a higher price, whereas those in blue 'push' the model towards a lower price.

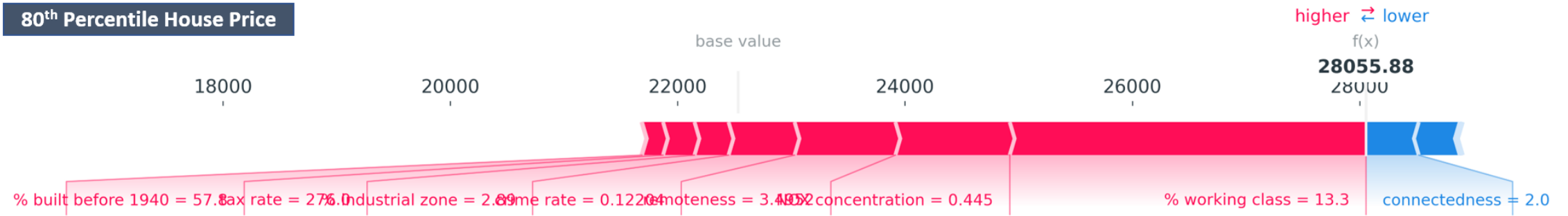
The final prediction, $f(x)$, is equal to the *base value* plus the sum of all the *SHAP values*.

The *base value* is the same for all examples in the data. It is equal to the average house price in the dataset.

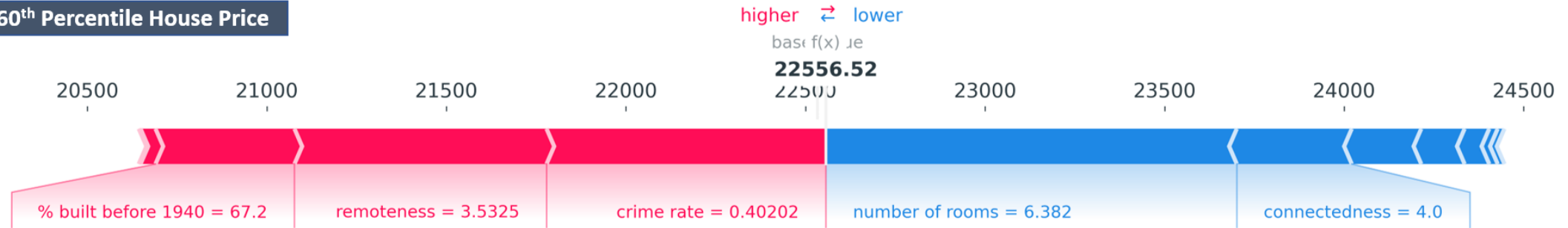
Force Plot



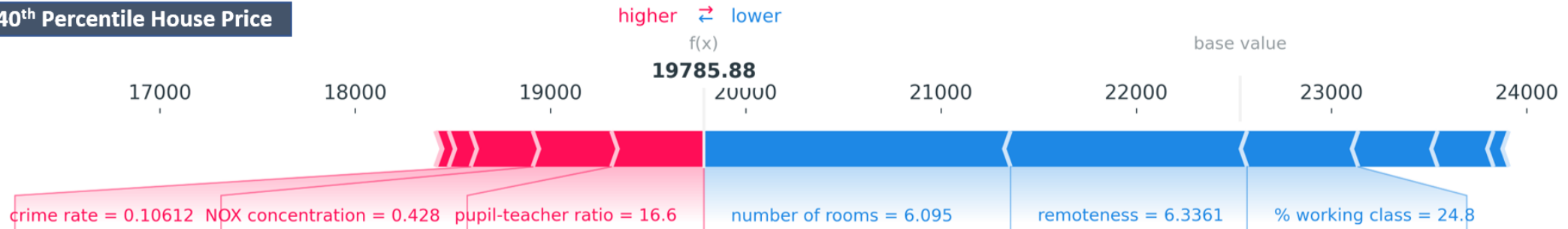
80th Percentile House Price



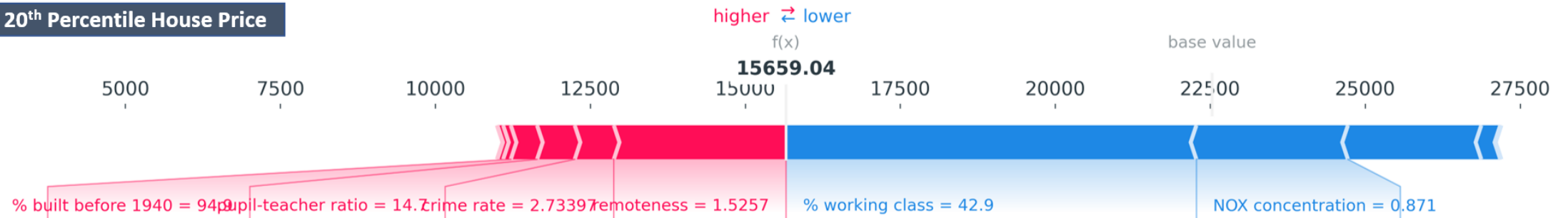
60th Percentile House Price



40th Percentile House Price

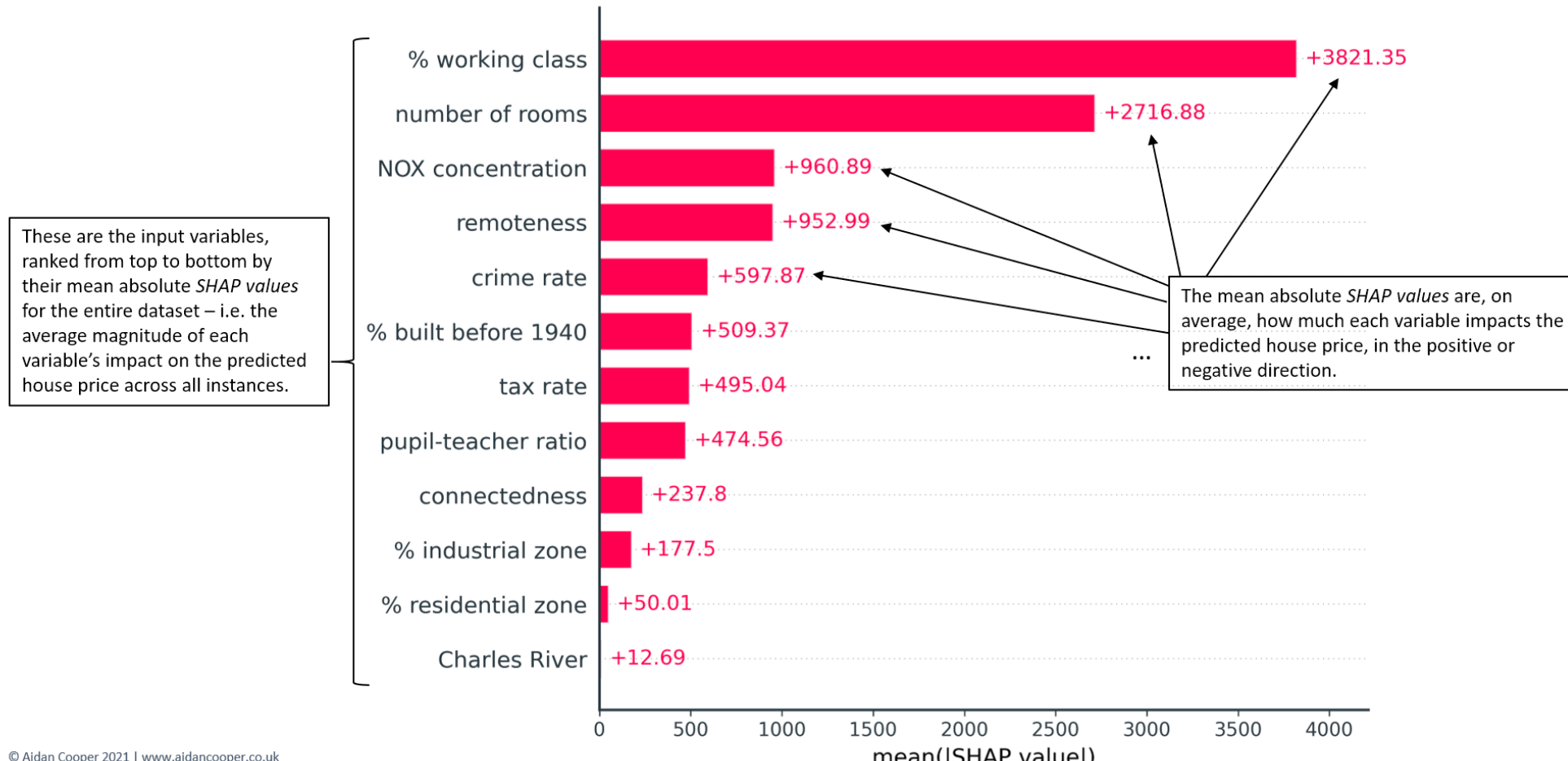


20th Percentile House Price

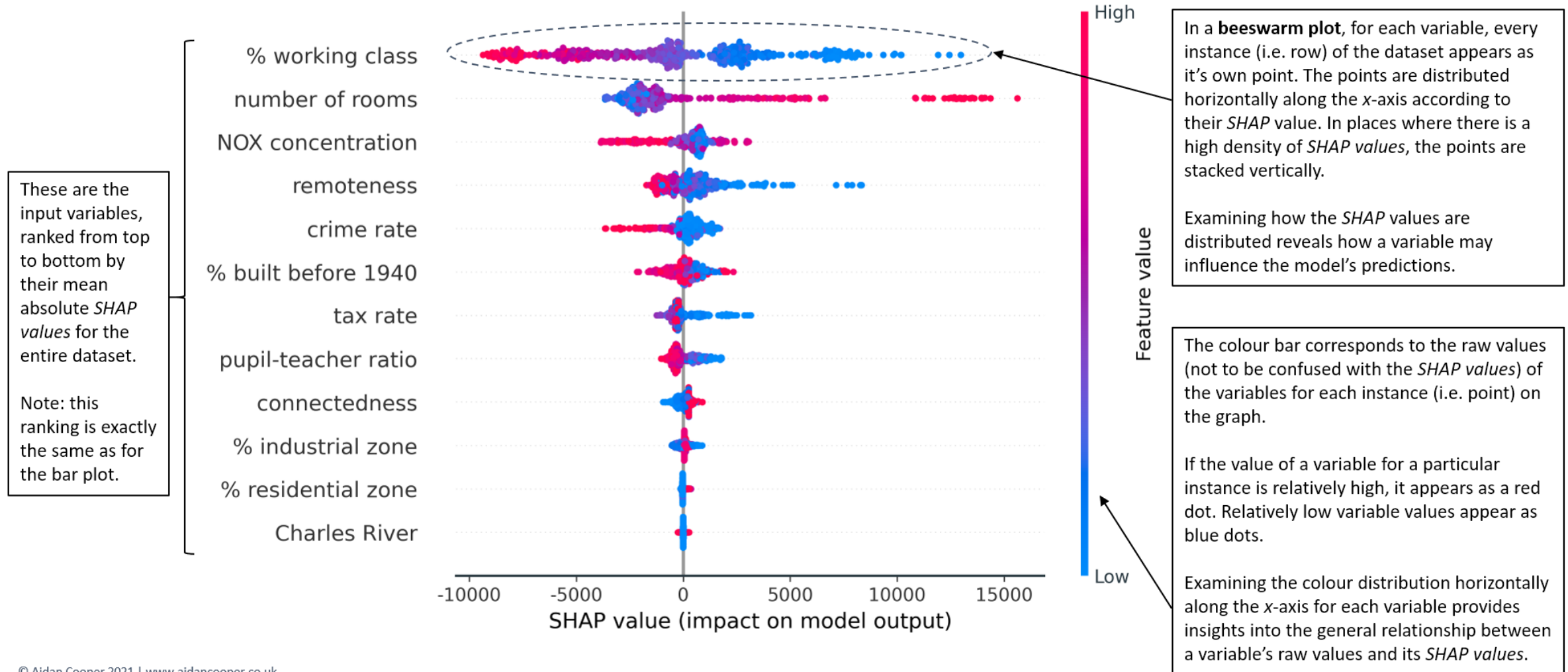


Global interpretability: understanding drivers of predictions across the population

- Bar Plot

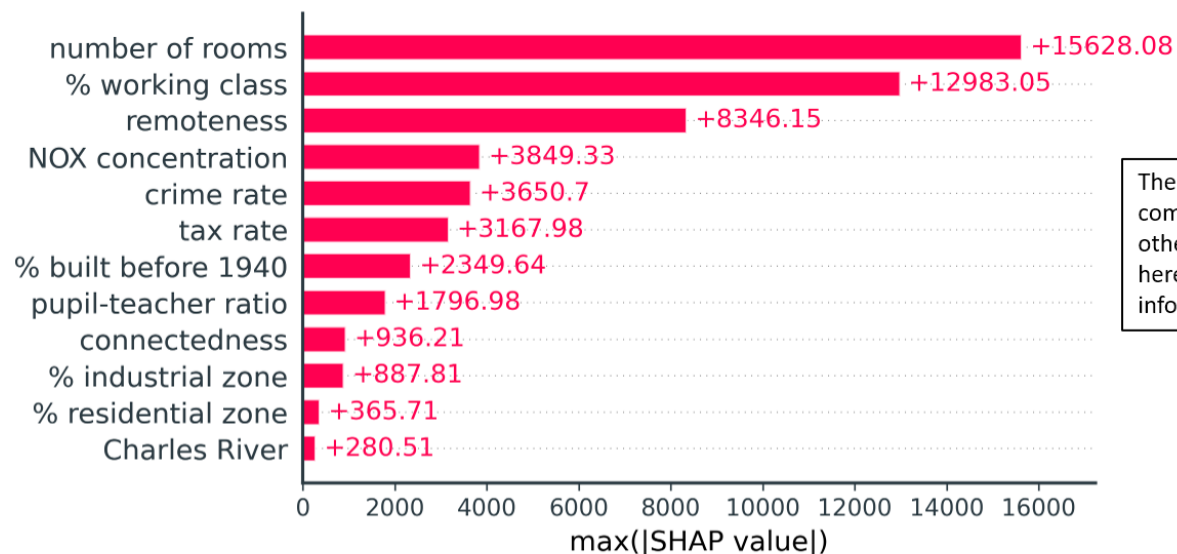


Beeswarm plots



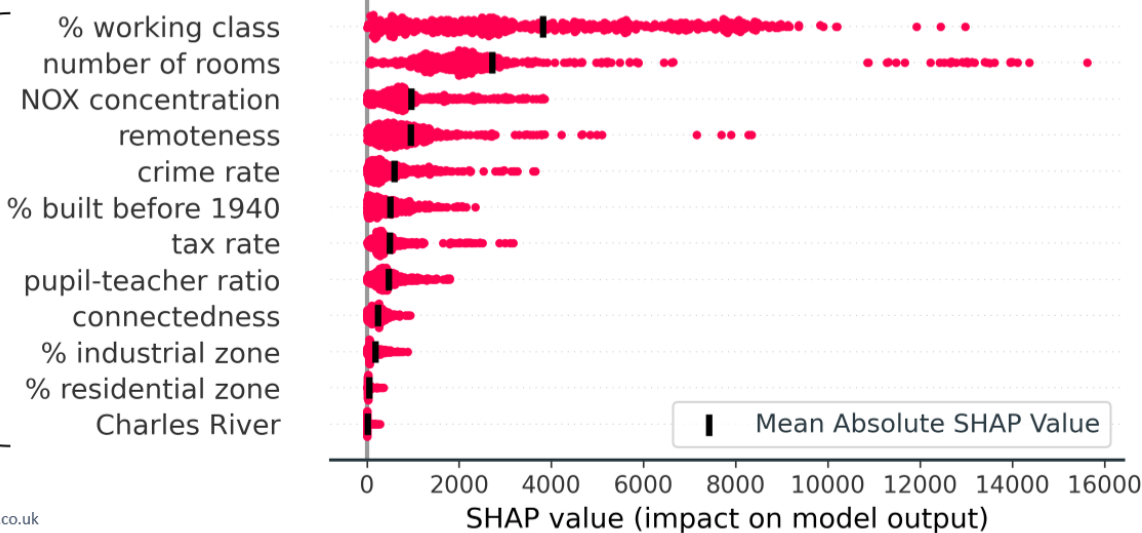
Whereas before, variables were ranked by their **mean** absolute SHAP value, here they are ranked by their **max** absolute SHAP value for the entire dataset.

Note that the ranking changes in places.



The mean absolute SHAP value is the most commonly used ranking for variables, but other statistics such as the max value (shown here) or median value may also be informative.

These are the input variables, ranked from top to bottom by their mean absolute SHAP values for the entire dataset.



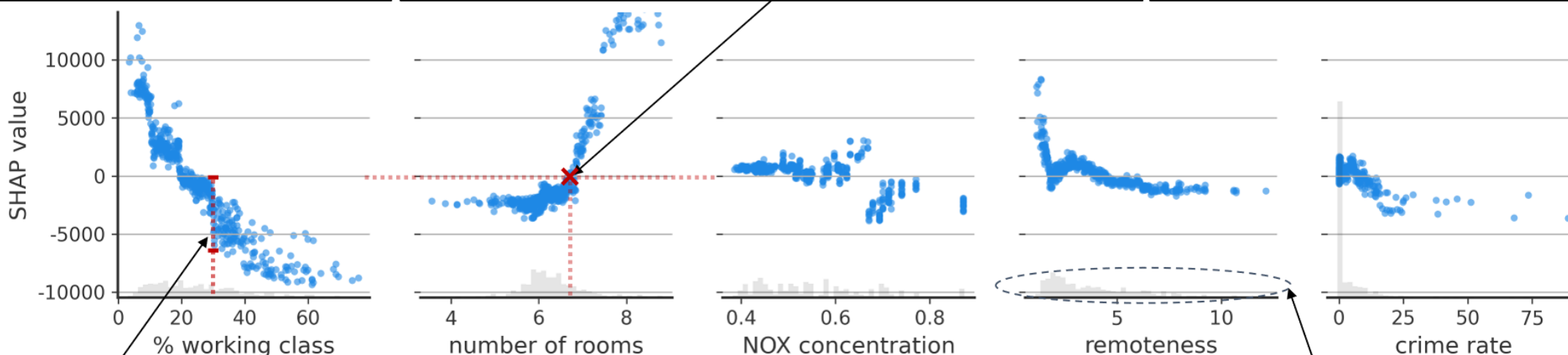
A balance can be struck between the simplicity of a bar plot and the information-rich complexity of a beeswarm plot, by creating a beeswarm plot for the **absolute** SHAP values.

This still shows us the ranking and relative influence of variables on the model's predictions, but also allows for further insights. E.g. the highest observed SHAP values actually occur for the 2nd ranked variable, *number of rooms*.

In a **dependence plot**, every instance (i.e. row) of the dataset appears as it's own point. The points are presented as a scatterplot of a variable's *SHAP values* versus the variables underlying raw values.

SHAP values above the $y=0$ line lead to predictions of higher house prices, whereas those below it are associated with lower house price predictions. The raw variable value at which the distribution of *SHAP values* cross the $y=0$ line tells you the threshold at which the model switches from predicting lower to higher house prices. For *number of rooms*, this is at approximately 6.8 rooms, as marked by the **X**.

With all five plots on the same y -scale, the extent of the vertical distribution of the *SHAP values* indicates how much relative influence each variable has on predictions. *% working class* has a much wider range of *SHAP values* than *crime rate*.



The vertical spread of *SHAP values* at a fixed raw variable value is due to *interaction effects* with other variables. For example, here we see that houses with a *% working class* of 30% can have *SHAP values* that range from \$0 to -\$6,500 depending on the other data for those particular instances.

The shapes of the distributions of points provide insights into the relationship between a variable's values and its *SHAP values*. For *% working class*, we see a negative, linear relationship across the full range of variable values. For *number of rooms*, we see that *SHAP values* are mostly flat between 4 and 6.5 rooms, but then increase sharply for higher room counts.

The inset histograms just above the x-axis display the distributions of raw variable values. We should be cautious not to overinterpret regions of the dependence plot where the underlying data is sparse (e.g. *crime rates* over 25%).

