

# What Drives Used Car Value?

## Predicting Used Car Prices in 2022

Robert Davies  
General Assembly DSI21



## GOALS

# Goals

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

## DESCRIBING THE MODEL

## RISKS AND LIMITATIONS

## IMPACT AND NEXT STEPS

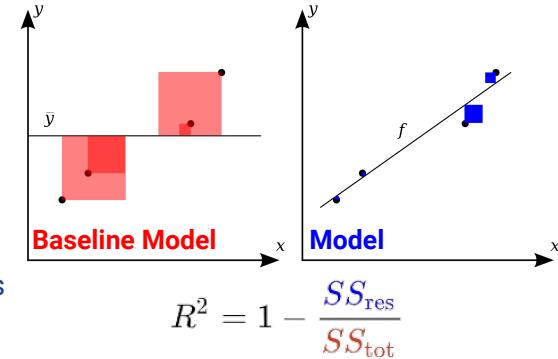
- The **primary** goal was to generate a predictive model for used car price in the UK in 2022, which is also interpretable.
- The **secondary** goal was to use the resulting model to test the hypothesis:

*When all other car attributes are equal, a Dacia branded car is cheaper than a Volvo branded car.*



# Success Criteria

- The **primary** metric used for model evaluation was **R<sup>2</sup> score**.
  - This can be interpreted as the fraction of target (car price) variance explained by the model which was not explained by the baseline (average) model.
  - An R<sup>2</sup> score of 1 is perfect prediction, whilst an R<sup>2</sup> score of 0 equates to a model equivalent to the average.
  - R<sup>2</sup> score is a unitless metric.
- The **secondary** metric of **RMSE** (Root Mean Square Error) was used to further describe model performance.
  - This is defined as the square root of mean squared errors.
  - As error values are squared, larger errors are penalised more than smaller errors.
  - The units of RMSE are in £.



## GOALS

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

## DESCRIBING THE MODEL

## RISKS AND LIMITATIONS

## IMPACT AND NEXT STEPS

# The Data - Where It Comes From

- 400,247 New and Used cars were scraped initially.
  - 369,730 Used cars at input to modelling after data cleaning
- Scraping problems encountered:
  - Limited search results
  - Cloudflare website security
- For each car, effort was made to scrape the data below. Ultimately, not all of it was used in modelling.

Engine Size	Drivetrain	MPG	Listing ID	
Fuel Type	ULEZ	Price	Brand	BHP
Dealer City	Mileage	Transmission	Seats	Year Made
Body		Dealer Rating	Doors	Listing URL
Owners				

Used in Modelling | Not Used in Modelling

## AutoTrader.co.uk

Volkswagen Tiguan  
2016 (66 reg)  
2.0 TDI BMT 150 4Motion SE Nav 5dr TOWBAR - LANE ASSIST - DRIVING MODES

£18,424  
£1,156 below market average

Car price	£18,424
Seller's admin fee	No fees
Total price	£18,424
Monthly price (CS)	£417 →

APPLE CARPLAY - WIFI - DAB  
Extra features →

carsa

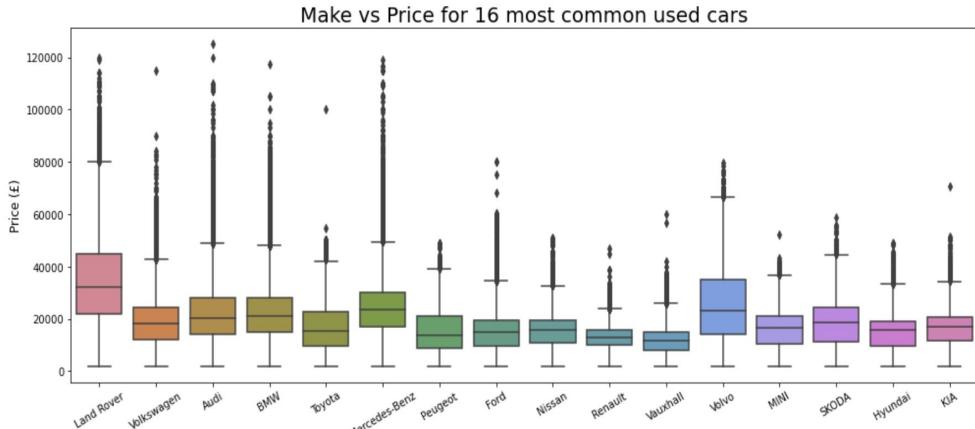
Carsa Mountsorrel  
[Visit website](#)  
[\(0116\) 484 7460](#)  
[Text](#) [Chat](#)  
 Visit this vehicle today at  
Carsa Mountsorrel  
[View map and address](#) →

High Beam Assistant, Metallic Paint, Cor-Net App-Connect, Lane Departure Warning System, Driving Modes, DAB Digital Radio/CD Player With AUX/USB/SD Input, Bluetooth Phone + Audio Streaming, Apple CarPlay, WiFi, Towbar, 3 Zone Climate Control, Four Wheel Drive, 18 inch alloys, Front and Rear Parking... [Read more](#)

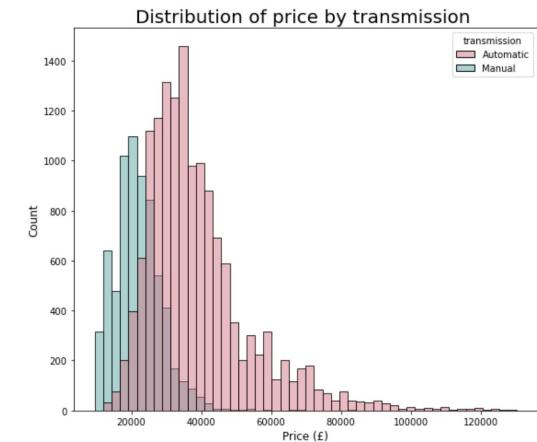
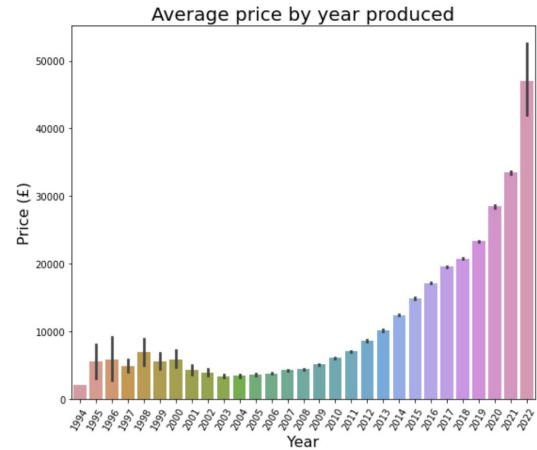
Specification

Available from multiple collection points.

# The Data - Brief Overview



- Land Rover, Mercedes-Benz and Volvo stand out to be the most expensive amongst the most common brands.
- It is evident which brands are targeting the high-end market and which are not.
- The relationship between year produced and price is not linear.
- Categorical predictors such as transmission may be useful in predicting car price



## GOALS

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

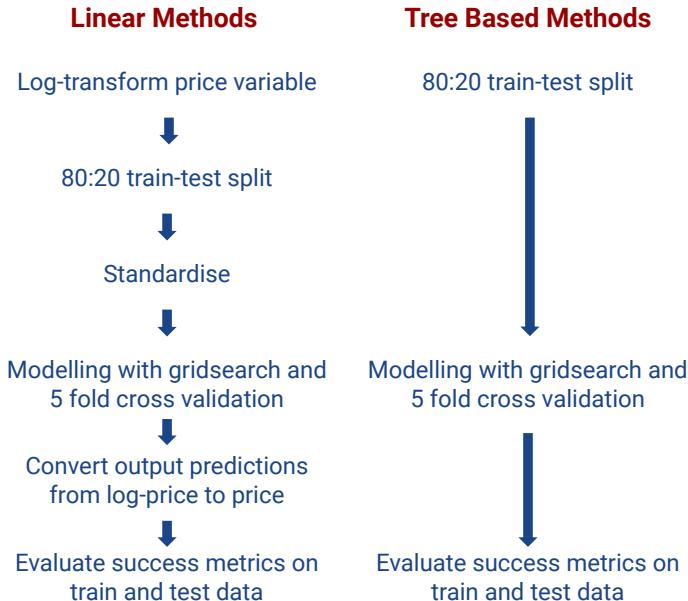
## DESCRIBING THE MODEL

## RISKS AND LIMITATIONS

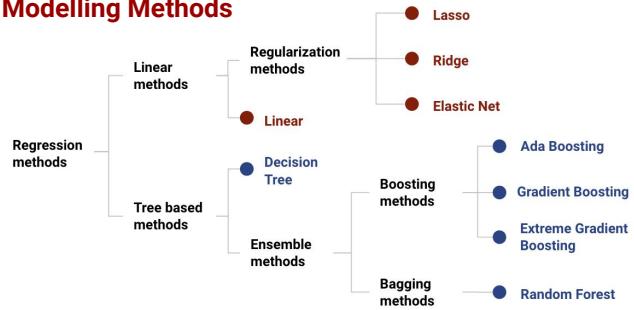
## IMPACT AND NEXT STEPS

# Selecting A Model

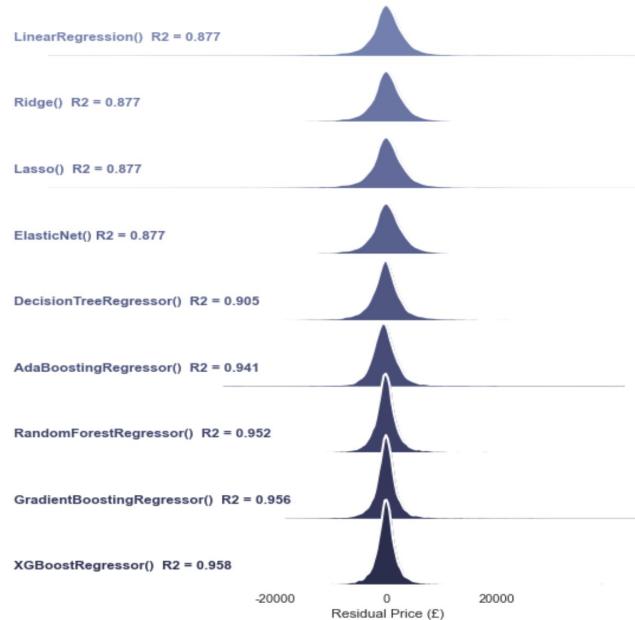
- Nine different models were tested on a random 50,000 car subset of the dataset.
- These models have been categorised into linear methods and tree based methods.



## Modelling Methods



## Price Residuals For Each Model Tested



## GOALS

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

## DESCRIBING THE MODEL

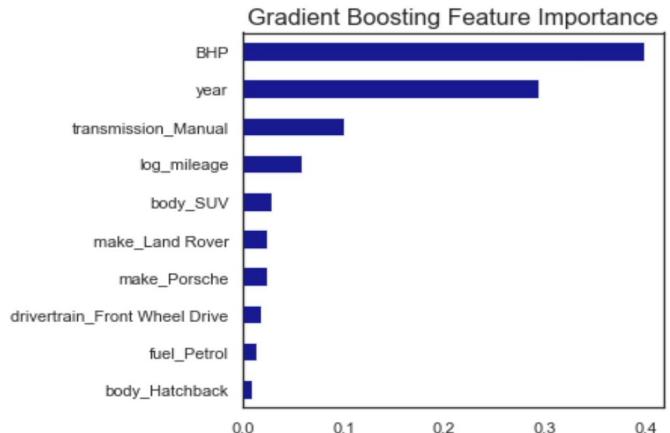
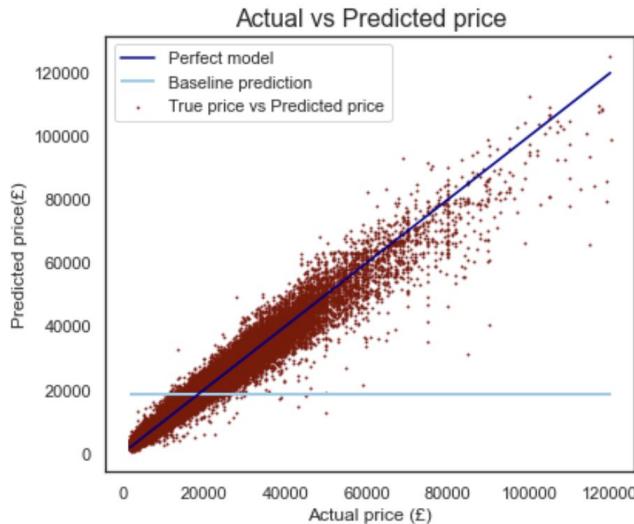
## RISKS AND LIMITATIONS

## IMPACT AND NEXT STEPS

# Gradient Boosting Model

- Gradient boosting regression is a which starts by simply predicting the average price of a car for every car and computing the price residuals.
- The model then uses a stepwise method to improve the result by iteratively modelling and reducing the residuals.
- After each step, new residuals are computed and the process repeats itself.
- The model is considered optimised when adding further steps fails to reduce the residuals.

Metric	R <sup>2</sup> Score	RMSE (£)
Test data	0.954	2496.72
Train data	0.955	2553.55



## GOALS

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

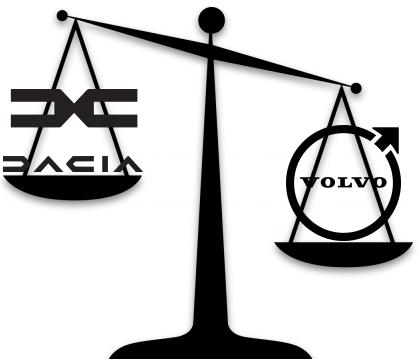
## DESCRIBING THE MODEL

## RISKS AND LIMITATIONS

## IMPACT AND NEXT STEPS

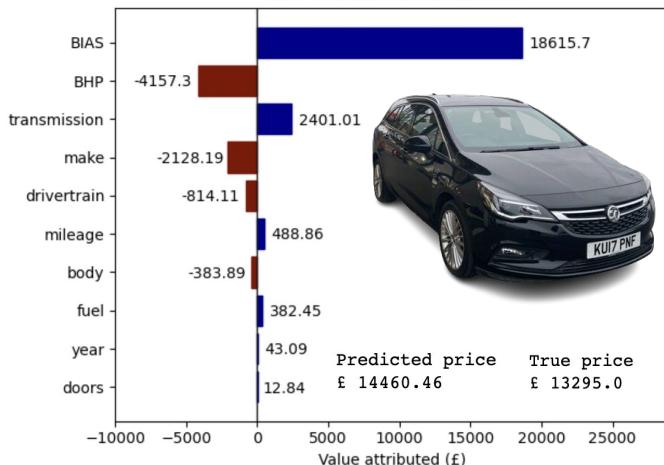
# Under The Hood

- One minor problem with Gradient Boosting Methods is that they aren't easy to interpret in detail.
- The Explain it Like I'm Five (ELI5) package from Sklearn has been used to obtain an itemised value breakdown for any car
- This has been made possible through the computation of permutation importance rather than feature importance, or Gini importance from the model.
- The original hypothesis has been answered - A Dacia is indeed cheaper than Volvo

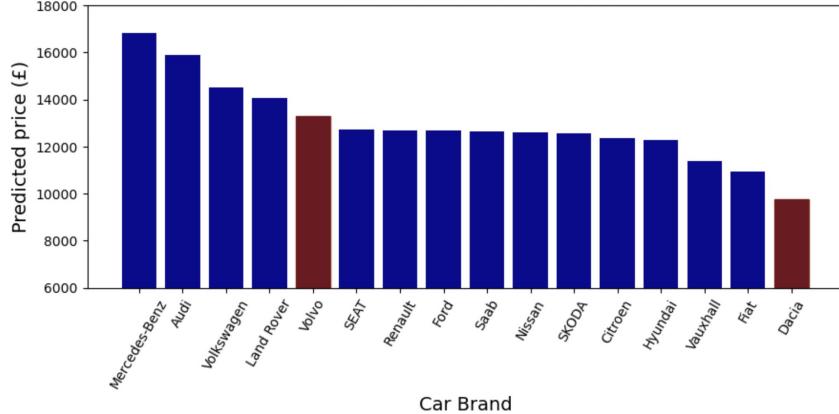


Car Summary: Vauxhall Astra 1.6 CDTi Elite Nav Sports Tourer Auto 5dr

Car Value Breakdown



Brand Influence on Car Price



## GOALS

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

## DESCRIBING THE MODEL

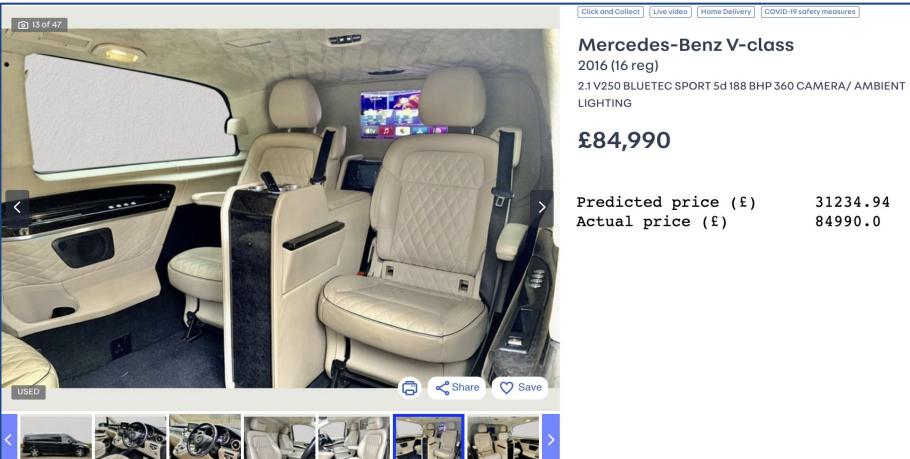
## RISKS AND LIMITATIONS

## IMPACT AND NEXT STEPS

# Risks And Limitations

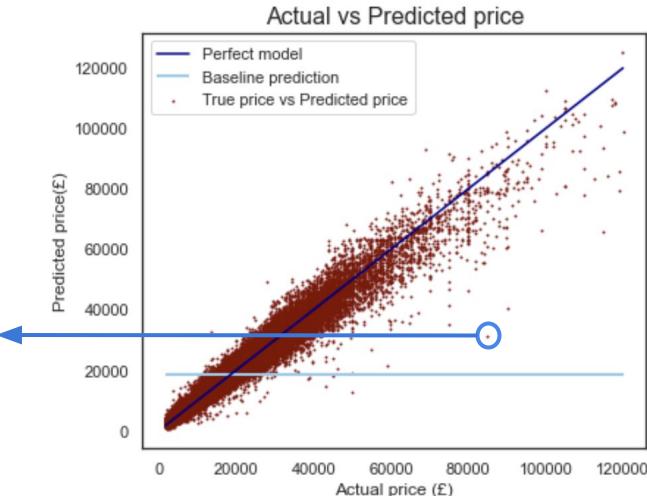
- The data may go stale quickly if the UK market changes.
- Outlier analysis has allowed me to sift out some erroneous rows caused by poor data entry. Non-remarkable data entry errors would harm the model but are difficult to identify.
- 'Car spec' is not captured in the dataset. Besides possibly the BHP variable, any non-standard upgrades are not captured.

## Largest Residual



## Suspicious SKODA

Erroneous SKODA (?) entry		Description
name	skoda 911	skoda 911
name_subtitle	911 996 GT3 3.6 Coupe Manual Petrol 5dr	911 996 GT3 3.6 Coupe Manual Petrol 5dr
make	SKODA	SKODA
price	119990.0	119990.0
year	2004	2004
body	Hatchback	Hatchback
mileage	6400.0	6400.0
engine	1.9	1.9
BHP	105.0	105.0



## GOALS

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

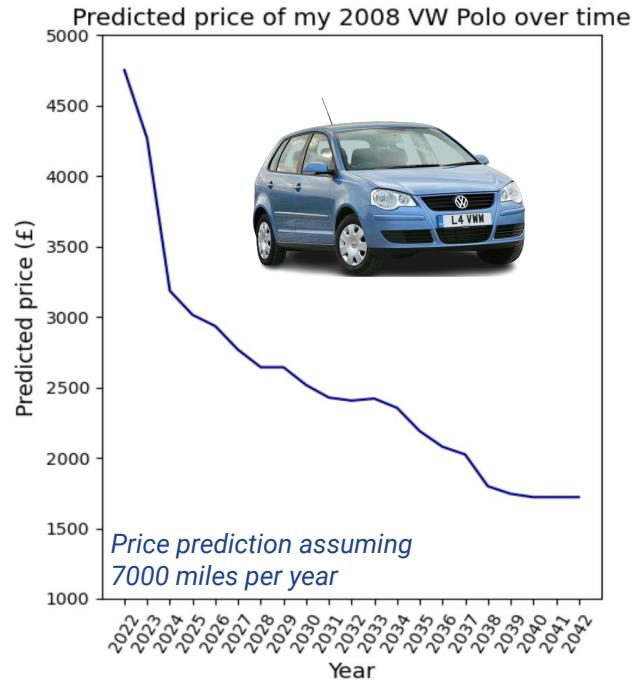
## DESCRIBING THE MODEL

## RISKS AND LIMITATIONS

## IMPACT AND NEXT STEPS

# Conclusions, Impact and Further Work

- All project goals have been met! With a 0.954 R<sup>2</sup> score it can be said that the used car model can be used to accurately predict used car price.
- Using permutation analysis with ELI5 has made it possible to explain why the Gradient Boosting Regressor does what it does.
- The model has been manipulated to rank common car brands in terms of influence on price and this could be extended to the other variables.
- With the model, we can predict the price of used cars in future.
- This model may be improved in future by implementing deep learning modelling methods.
- The model may be improved through enrichment of the predictor variables. It seems intuitive that car boot volume would impact the price of a car.
- It would be interesting to benchmark the model price predictions against AutoTrader's own value gauge. This isn't available on every listing and would require a new web scrape.



## GOALS

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

## DESCRIBING THE MODEL

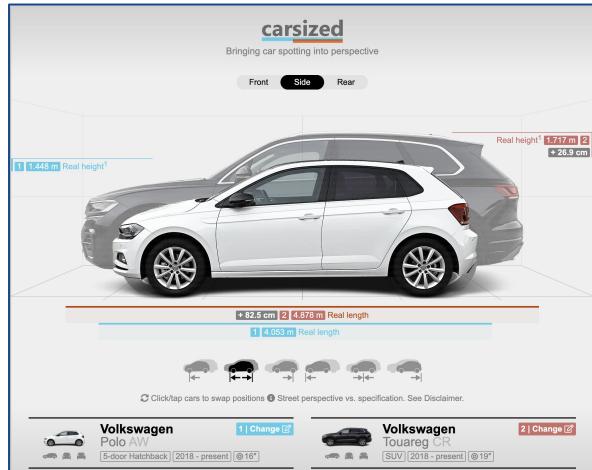
## RISKS AND LIMITATIONS

## IMPACT AND NEXT STEPS

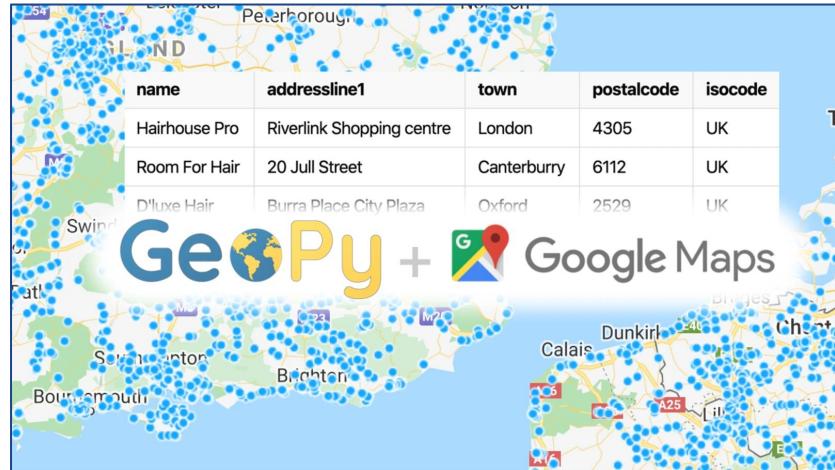
# Further Work - Enriching Predictors

- Some of it has been done already. Additional sources besides AutoTrader have been used to enrich the predictors. This sequence is awful but reflects how it was done.

## Carsized.com



## GeoPy with Google Maps API



Fuel Type Dealer County	Drivetrain Mileage	Price Transmission	Brand BHP	Year Made Cargo Volume	Body
----------------------------	-----------------------	-----------------------	--------------	---------------------------	------

Original | Additional

- 319,633** Complete Used Cars at input to further work
- 50,097 didn't have dealer county or cargo volume data

## GOALS

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

## DESCRIBING THE MODEL

## RISKS AND LIMITATIONS

## IMPACT AND NEXT STEPS

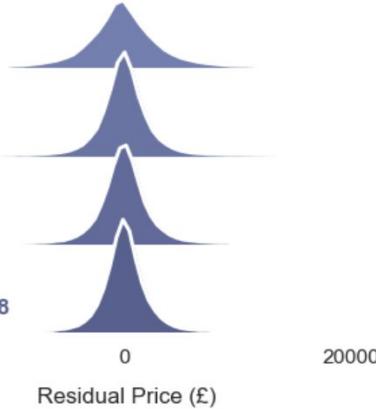
# Further Work

Fuel Type Dealer County	Drivetrain Mileage	Price Doors Transmission	Brand BHP	Year Made Cargo Volume	Body
----------------------------	-----------------------	--------------------------------	--------------	---------------------------	------

Original Continuous | Original Categorical | Additional

## Price Residuals

Continuous Predictors Only  $R^2 = 0.877$



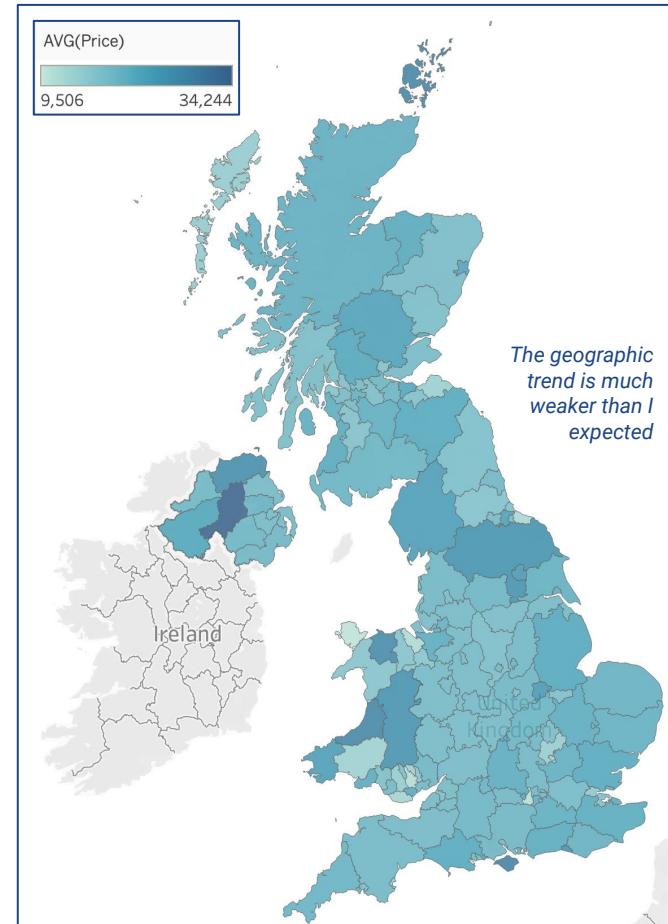
Original Predictors  $R^2 = 0.954$

Original Predictors + County  $R^2 = 0.956$

Original Predictors + Boot Volume  $R^2 = 0.968$

Metric	$R^2$ Score		RMSE (£)	
	Original	Upgraded	Original	Upgraded
Test data	0.954	0.968	2477.30	2077.61
Train data	0.958	0.968	2356.69	2069.74

## UK Counties Average Used Car Price Choropleth



# Questions?



# Car Value Breakdown with ELI5 - Example 1

y (score 35023.935) top features

Contribution?	Feature	Value
+18627.367	<BIAS>	1.000
+10981.325	BHP	301.775
+5272.278	log_mileage	9.507
+2718.439	year	2019.000
+2213.293	make_Mercedes-Benz	1.000
+2037.989	transmission_Manual	0.000
+1472.747	drivetrain_Front Wheel Drive	0.000
	... 30 more positive ...	
	... 16 more negative ...	
-661.444	make_Porsche	0.000
-837.760	fuel_Diesel	0.000
-1214.711	body_Hatchback	1.000
-1220.384	make_Land Rover	0.000
-2073.431	fuel_Petrol	1.000
-2345.195	body_SUV	0.000

## Alternative Slide

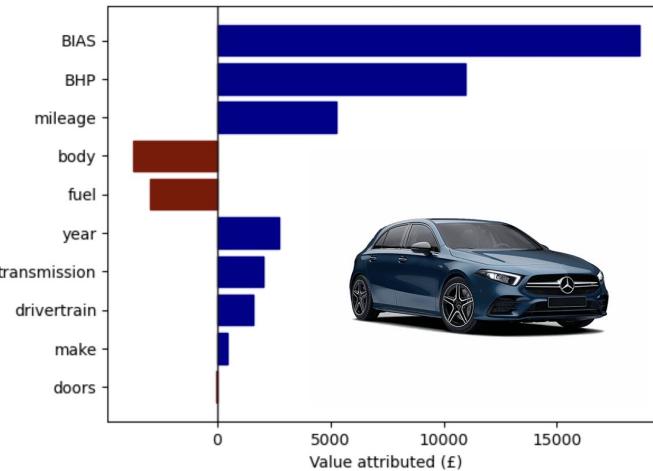
Car Summary: Mercedes-Benz A Class A35 4Matic Premium Plus 5dr Auto 2.0

Predicted price  
£ 35023.93

True price  
£ 38398.0

<https://www.autotrader.co.uk/car-details/202110028032322>

Car Value Breakdown



Predictor	Attribute	Price Contribution (£)	Cumulative Car Price (£)
BIAS		0	18627.367022
BHP		301.775148	10981.324640
mileage		13447.0	5272.278382
body	Hatchback	-3720.365950	31639.152757
fuel	Petrol	-2944.782282	28694.370475
year	2019	2718.438685	31412.809160
transmission	Automatic	2037.988511	33450.797671
drivetrain	Four Wheel Drive	1590.931155	35041.728826
make	Mercedes-Benz	478.548663	35359.518707
doors	5dr	-17.793867	35023.934959

# Car Value Breakdown with ELI5 - Example 1

## Alternative Slide

With permutation importance from ELI5 it is possible to quantify how each used car predictor impacted the final price prediction. The top row, labelled BIAS is the average car price.



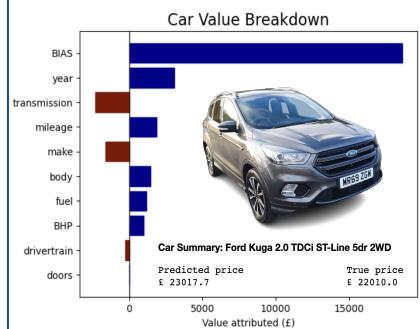
A function has been created to collate contributions from previously dummified categorical variables with accumulated price impacts displayed in a Pandas DataFrame.



Finally the car value breakdown is displayed in a horizontal bar chart with all ten predictor price contributions on the y-axis.

y (score 23017.701) top features		
Contribution?	Feature	Value
+18627.367	<BIAS>	1.000
+3100.195	year	2019.000
+1856.395	log_mileage	9.465
+1052.183	fuel_Petrol	0.000
+1037.665	body_Hatchback	0.000
+993.617	BHP	147.929
+403.330	body_SUV	1.000
	... 27 more positive ...	
	... 27 more negative ...	
-167.954	make_Mercedes-Benz	0.000
-198.167	make_Audi	0.000
-257.477	make_Land Rover	0.000
-285.174	drivetrain_Front Wheel Drive	1.000
-910.598	make_Ford	1.000
-2321.886	transmission_Manual	1.000

	Attribute	Price Contribution (£)	Cumulative Car Price (£)
Predictor			
BIAS		0	18627.367022
year		2019	3100.195227
transmission	Manual	-2321.885676	19405.676572
mileage		12898.0	1856.395476
make	Ford	-1595.570532	19665.501516
body	SUV	1469.615395	21135.116911
fuel	Diesel	1152.982247	22288.099158
BHP		147.928994	993.616996
drivetrain	Front Wheel Drive	-283.486230	23281.716054
doors	5dr	19.470730	23017.700554



# Car Value Breakdown with ELI5 - Example 2

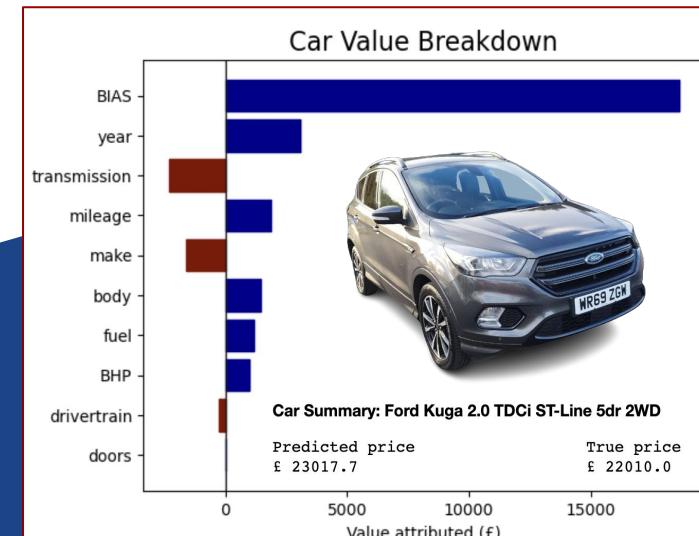
y (score 23017.701) top features

Contribution?	Feature	Value
+18627.367	<BIAS>	1.000
+3100.195	year	2019.000
+1856.395	log_mileage	9.465
+1052.183	fuel_Petrol	0.000
+1037.665	body_Hatchback	0.000
+993.617	BHP	147.929
+403.330	body_SUV	1.000
... 27 more positive ...		
... 17 more negative ...		
-167.954	make_Mercedes-Benz	0.000
-198.167	make_Audi	0.000
-257.477	make_Land Rover	0.000
-285.174	drivetrain_Front Wheel Drive	1.000
-910.598	make_Ford	1.000
-2321.886	transmission_Manual	1.000



Predictor	Attribute	Price Contribution (£)	Cumulative Car Price (£)
BIAS		0	18627.367022
year		2019	3100.195227
transmission	Manual	-2321.885676	19405.676572
mileage		12898.0	1856.395476
make	Ford	-1596.570532	17809.106040
body	SUV	1469.615395	21135.116911
fuel	Diesel	1152.982247	22288.099158
BHP		147.928994	993.616896
drivetrain	Front Wheel Drive	-283.486230	23281.716054
doors	5dr	19.470730	23017.700554

## Alternative Slide



## GOALS

## SUCCESS CRITERIA

## THE DATA

## OVERALL APPROACH

## DESCRIBING THE MODEL

## RISKS AND LIMITATIONS

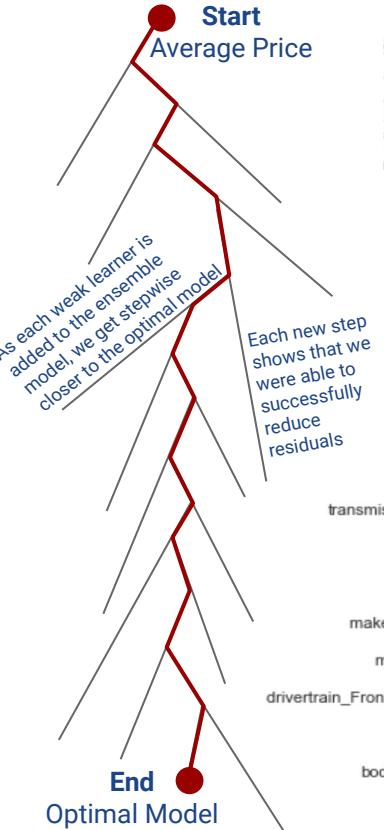
## IMPACT AND NEXT STEPS

# Gradient Boosting Model

- Gradient boosting regression is an ensemble method which starts by predicting the average price of a car for every car and computing the price residual.
- It then constructs a sequence of predictors which are designed to minimise these residuals.
- Each subsequent predictor is weighted by a small fraction to obtain a 'weak learner' before adding it to the ensemble. In this way, we take many small steps towards the optimal model.
- The model is considered optimised when adding further weak learners fail to reduce the residuals.

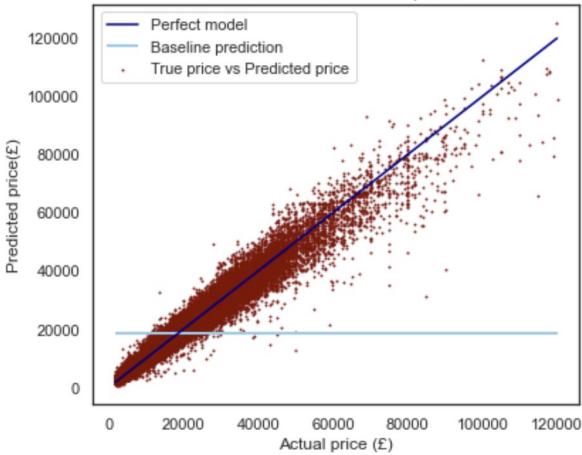
Metric	R <sup>2</sup> Score	RMSE (£)
Test data	0.954	2496.72
Train data	0.955	2553.55

Simplified model schematic

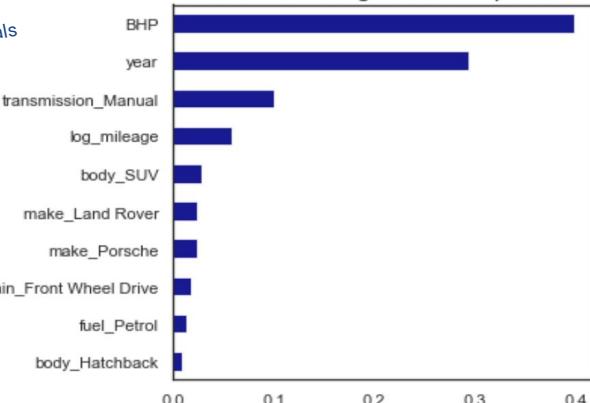


## Alternative Explanation

Actual vs Predicted price



Gradient Boosting Feature Importance



# AutoTrader Search Results Alternative Explanation

autotrader.co.uk/results-car-search

## Another outlier example

## Alternative Slide

