

Tarunvir Singh  
Professor Bailey  
Statistics  
12/05/2025

## What factors best predict engine size?

### Introduction

In this project, I see how different car features relate to engine size using the Automobile dataset from the UCI Machine Learning Repository. The purpose of the analysis is to identify which variables are most effective at predicting engine size and to compare several linear regression models. After reviewing the dataset, I listed a variety of possible numerical predictors such as horsepower, wheelbase, curb weight, fuel mileage, and price. I then saw each predictor using scatter plots and selected the variables that showed the strongest relationships with engine size. The models will be assessed using  $R^2$  and AIC, where higher  $R^2$  and lower AIC indicate better predictive performance. I chose engine size as the response variable because it is a central element of a vehicle's performance and design, and many car characteristics such as horsepower, fuel efficiency, and price. They are strongly affected by the size of the engine. This makes engine size a meaningful outcome for analyzing how different vehicle features relate to one another.

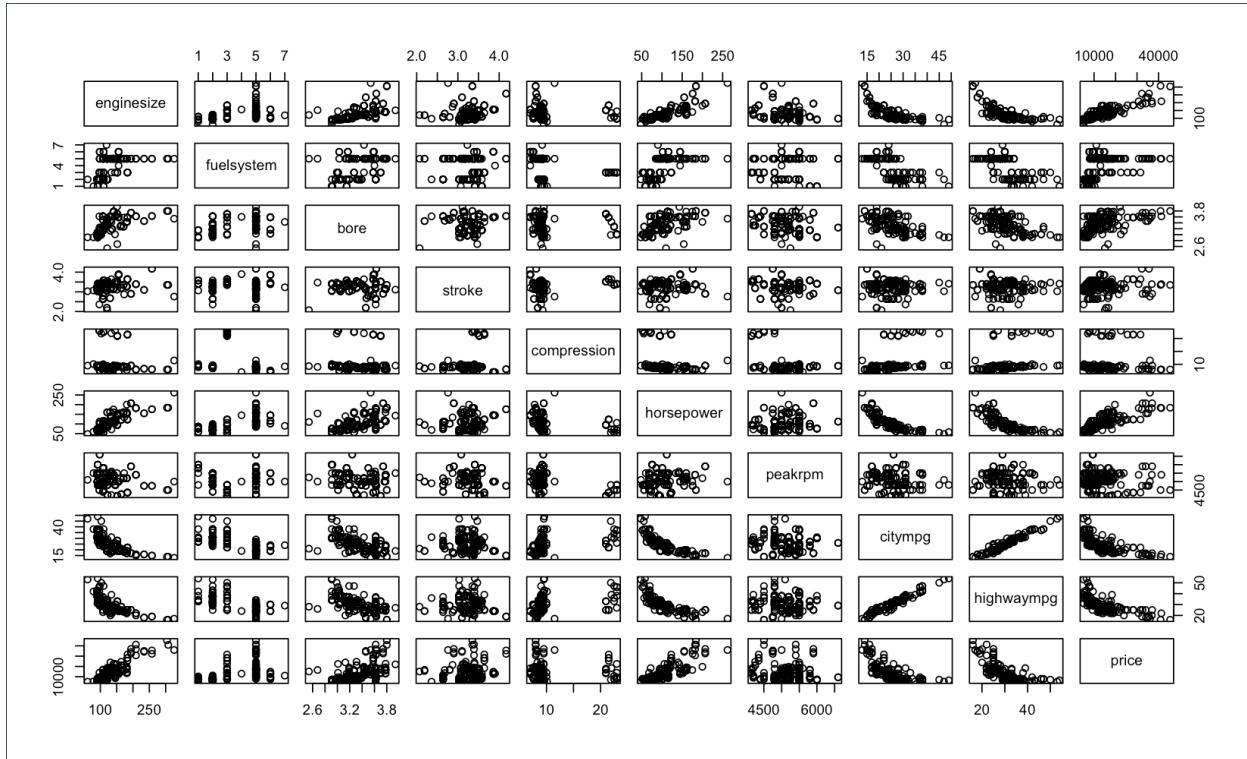
### Methods Section

Below is a snippet of the first six rows of the imported Automobile dataset (import85 file).

	symboling	make	fuel	aspiration	numdoors	bodystyle	drivewheels	enginelocation	wheelbase	length
1	3	alfa-romero	gas	std	two	convertible	rwd	front	88.6	168.8
2	3	alfa-romero	gas	std	two	convertible	rwd	front	88.6	168.8
3	1	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	171.2
4	2	audi	gas	std	four	sedan	fwd	front	99.8	176.6
5	2	audi	gas	std	four	sedan	4wd	front	99.4	176.6
6	2	audi	gas	std	two	sedan	fwd	front	99.8	177.3
	width	height	curbweight	enginetype	numcylinders	enginesize	fuelsystem	bore	stroke	compression
1	64.1	48.8	2548	dohc	four	130	mpfi	3.47	2.68	9.0
2	64.1	48.8	2548	dohc	four	130	mpfi	3.47	2.68	9.0
3	65.5	52.4	2823	ohcv	six	152	mpfi	2.68	3.47	9.0
4	66.2	54.3	2337	ohc	four	109	mpfi	3.19	3.40	10.0
5	66.4	54.3	2824	ohc	five	136	mpfi	3.19	3.40	8.0
6	66.3	53.1	2507	ohc	five	136	mpfi	3.19	3.40	8.5
	peakrpm	citympg	highwaympg	price						
1	5000	21	27	13495						
2	5000	21	27	16500						
3	5000	19	26	16500						
4	5500	24	30	13950						
5	5500	18	22	17450						
6	5500	19	25	15250						

I selected engine size as the response variable and examined a wide range of possible predictors, including:

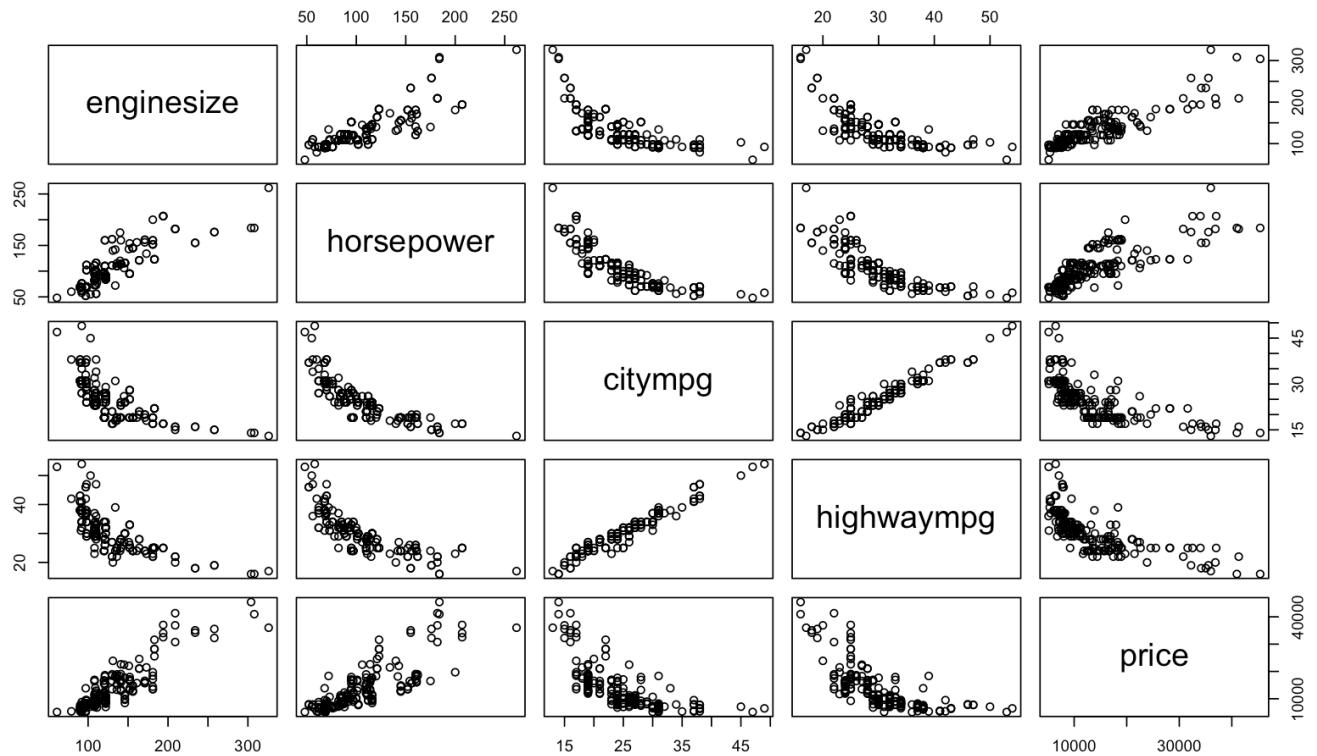
- horsepower
- fuelsystem
- bore
- stroke
- Compression
- peakrpm
- citympg
- highwaympg
- price



To narrow these down, I created scatter plots to visually check the strength, direction, and shape of each relationship with engine size. Scatter plots are useful because they show whether the relationship looks linear, curved, weak, or too scattered to model.

After reviewing every graph, I found that horsepower, city mpg, highway mpg, and price were the strongest candidates. These variables showed clear patterns, either positive linear trends (horsepower, price) or inverse trends (city and highway mpg). In contrast, variables like bore, stroke, compression, peakrpm, wheelbase, and curb weight showed relationships that were either too weak or too dispersed to be useful in a linear model.

Here are the models for them in Scatter Plots



### Model 1 - Single Predictor

I created one-predictor regression models for the four strongest candidates: horsepower, city mpg, highway mpg, and price. I compared their performance using Adjusted R<sup>2</sup> and AIC.

Model	Predictor	Adjusted R <sup>2</sup>	AIC
m1.hp	Horsepower	0.7131	1749.688
m1.ct	Citympg	0.5106	1852.726
m1.hw	Highwaympg	0.5416	1840.130
m1.pr	Price	0.7888	1690.528

Among the four, price stood out clearly. It had the highest Adjusted R<sup>2</sup> and the lowest AIC, making it the strongest single predictor of engine size.

For Model 1, I used price as the single predictor

**Coefficients:**

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	6.742e+01	2.650e+00	25.44	<2e-16	***
price	4.570e-03	1.705e-04	26.80	<2e-16	***

The R gave me

**Intercept** 67.42 Predicted engine size when price = 0 (not meaningful, but needed for the line)

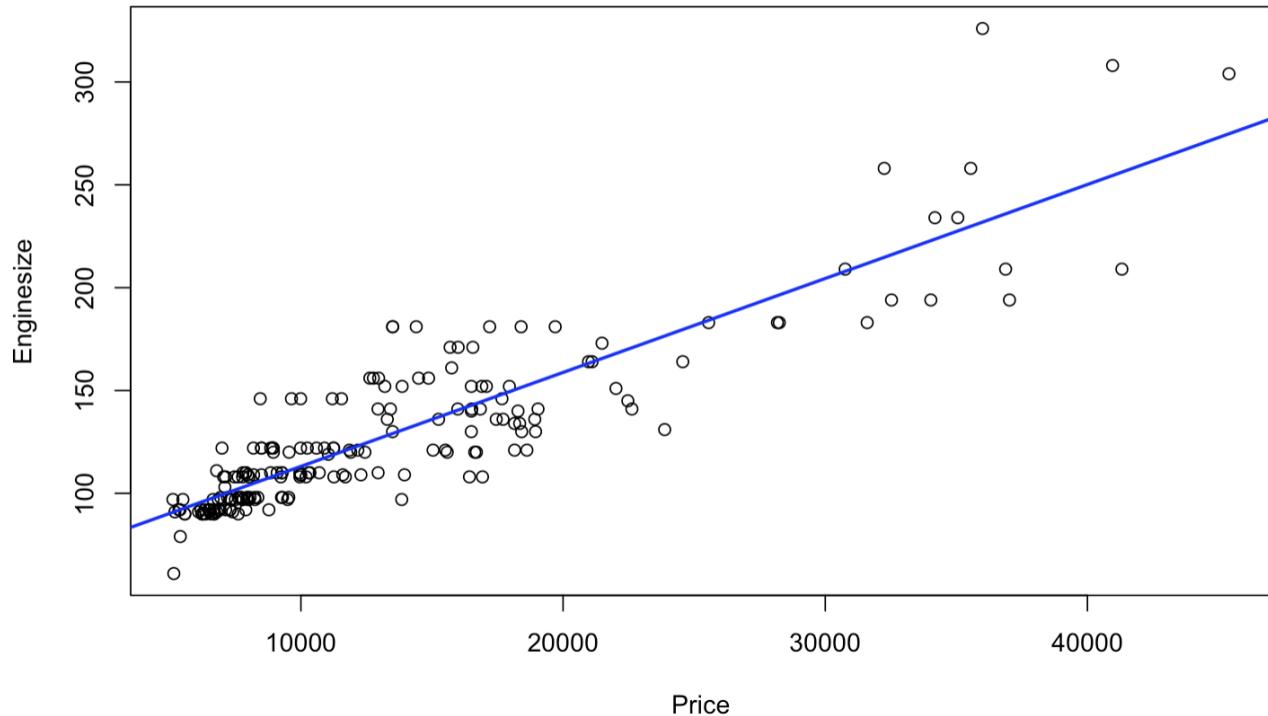
**price** 0.00457 Engine size increases by 0.00457 units for each \$1 increase in price

Equation: enginesize = 67.42 + 0.00457(price)

Price is statistically significant ( $p < 2 \times 10^{-16}$ ), meaning the positive relationship between price and engine size is highly reliable. Even as a single predictor, price explains a large amount of variation in engine size.

Scatterplot of engine size vs. price with fitted regression line for Model 1.

**Model 1 Enginesize vs Price**



### **Model 2 - Multiple Predictors**

Next, I built several multi-predictor models to see whether adding variables would improve predictive performance.

<b>Model</b>	<b>Predictor</b>	<b>Adjusted R<sup>2</sup></b>	<b>AIC</b>
m2.pr.hp	Price + horsepower	0.8329	1646.379
m2.pr.hw	Price + Highwaympg	0.808	1673.186
m2.pr.hp.hw.ct	Price + Horsepower + Highwaympg + Citympg	0.84	1639.909
m2.pr.hp.sq	price + price ^2 + horsepower + horsepower^2	0.8417	1637.826

The best-performing model was m2.pr.hp.sq, which included price, horsepower, and their squared terms. This model had the highest Adjusted R<sup>2</sup> (0.8417) and the lowest AIC (1637.826)

Full Equation from R

Coefficients:	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.763e+01	9.146e+00	6.301	2.04e-09 ***
price	7.245e-04	8.071e-04	0.898	0.37051
I(price^2)	5.147e-08	1.716e-08	2.998	0.00308 **
horsepower	4.754e-01	1.991e-01	2.388	0.01792 *
I(horsepower^2)	-6.302e-05	7.323e-04	-0.086	0.93152

$$\text{Enginesize} = 57.63 + 0.0007245 \text{ (price)} + 0.00000005147 \text{ (price}^2\text{)} + 0.4754 \text{ (horsepower)} - 0.00006302 \text{ (horsepower}^2\text{)}$$

Interpretation of coefficients:

- **Horsepower** has a statistically significant positive effect ( $p = 0.0179$ ), meaning vehicles with higher horsepower generally have larger engines.
- **Price<sup>2</sup>** is also significant ( $p = 0.00308$ ), suggesting a curved relationship where more expensive cars tend to have much larger engines.
- **Horsepower<sup>2</sup>** was not significant, but allowing curvature still helped model fit.
- **Price** was not significant alone because the curved effect explains more of the relationship.

Overall, this model shows that engine size increases more quickly at higher levels of horsepower and price, which means the relationships are nonlinear. To add on, when a predictor and its squared term are both included, they often overlap, which can make the linear term look less important. But the squared terms are very important and make the AIC much better, so this is still the best model. This confirms that engine size is better explained by curved (nonlinear) relationships rather than a straight line.

## **Comparing Model 1 and Model 2**

Model 1 (m1) uses only price to predict engine size. This model explains a good amount of the variability in engine size, with an  $R^2$  of 0.7888. Its AIC value is 1690.528, which shows how well the model fits the data based on its overall accuracy.

Model 2 (m2) includes multiple predictors. The best-performing version of Model 2 is the m2.pr.hp.sq model, which included price, price<sup>2</sup>, horsepower, and horsepower<sup>2</sup>. This model has an  $R^2$  of 0.8417, which is higher than Model 1. This means it explains more of the variability in engine size. It also has a much lower AIC of 1637.826, which shows that this model fits the data better than the one-predictor model.

Overall, Model 2 clearly performs better. It explains more of the variation in engine size and has the lowest AIC, meaning it should also do a better job predicting future data compared to Model 1. Even though Model 1 is simpler, Model 2 provides the strongest and most reliable fit.

## **Conclusion**

In this study, I analyzed how different car features relate to engine size and compared several regression models to find the strongest predictors. After analyzing scatter plots and testing multiple models, I found that price and horsepower were consistently important variables in explaining engine size. While the one-predictor model using price gave a clear linear relationship, the multi-predictor models provided a more complete and accurate picture overall.

The best-performing model combined price, horsepower, and their squared terms, showing that engine size is influenced by both linear and nonlinear effects (parabola or curve). This confirms that engine size does not increase at a constant rate; instead, both price and horsepower show accelerating effects at higher values, which is why the squared terms improved the model.

Overall, the analysis shows that cars with higher prices and greater horsepower tend to have larger engines, and that using several predictors leads to better predictions than relying on a single variable. This project demonstrates the value of building and comparing different regression models to understand which factors matter most.