



The usedcars.csv Dataset

The best way to learn the process of data exploration is with an example.

- We will explore the usedcars.csv dataset, which contains actual data about used cars recently advertised for sale on a popular U.S. website.
- Following along with the examples, we have to be sure that this file has been downloaded and saved to our R working directory.
- Since the dataset is stored in the CSV form, we can use the read.csv() function to load the data into an R data frame:

> usedcars <- read.csv("usedcars.csv", stringsAsFactors = FALSE)

- Given the usedcars data frame, we will now assume the role of a data scientist who has the task of understanding the used car data.
 - Although data exploration is a fluid process, the steps can be imagined as a sort of investigation in which questions about the data are answered.
 - The exact questions may vary across projects, but the types of questions are always similar. 19.03.2025 Eksploracja danych

Exploring the Structure of Data $(1/2)$)
---	---

- One of the first questions to ask in an investigation of a new dataset should be about how the dataset is organized.
 - If you are fortunate, your source will provide a data dictionary, which is a document that describes the dataset's features. Assuming that the used car data does not come with this documentation, we'll need to create one on our own.
 - The str() function provides a method to display the construction of R data structures such as data frames, vectors, or lists. It can be used to create the basic outline for our data dictionary:

> str(usedcars)			
'data.frame': 150 obs. of 6 variables:			
\$ year : int 2011 2011 2011 2012 2010 2011 2010 2011 2010			
\$ model : chr "SEL" "SEL" "SEL"			
\$ price : int 21992 20995 19995 17809 17500 17495 17000 16995 16995 16995			
\$ mileage : int 7413 10926 7351 11613 8367 25125 27393 21026 32655 36116			
\$ color : chr "Yellow" "Gray" "Silver" "Gray"			
\$ transmission: chr "AUTO" "AUTO" "AUTO" "AUTO"			

19.03.2025

Eksploracja danych



Exploring Numeric Variables (1/2)				
To investigate the numeric variables in the used car data, we will employ a common set of measurements to describe values known as summary statistics.				
The summary() function displays several common summary statistics. Let's take a look at a single feature, year:				
> summary(usedcars\$year) Min. 1st Qu. Median Mean 3rd Qu. Max. 2000 2008 2009 2010 2012				
 Even if you aren't already familiar with summary statistics, you may be able to guess some of them from the heading before the summary() output. 				
• Ignoring the meaning of the values for now, the fact that we see numbers such as 2000, 2008, and 2009 could lead us to believe that the year variable indicates the year of manufacture rather than the year the advertisement was posted, since we know the vehicles were recently listed for sale.				
19.03.2025 Eksploracja danych				













The Five-Number Summary

The five-number summary is a set of five statistics that roughly depict the spread of a feature's values.

All five of the statistics are included in the output of the summary() function. Written in order, they are:

- 1. Minimum (Min.)
- 2. First quartile, or Q1 (1st Qu.)
- 3. Median, or Q2 (Median)
- 4. Third quartile, or Q3 (3rd Qu.)
- 5. Maximum (Max.)
- As you would expect, minimum and maximum are the most extreme feature values, indicating the smallest and largest values, respectively. R provides the min() and max() functions to calculate these values on a vector of data.
- The span between the minimum and maximum value is known as the range. In R, the range() function returns both the minimum and maximum value.

19.03.2025

Eksploracja danych





The quantile() Function					
The quantile() function provides a robust tool to identify quantiles for a set of values.					
• By default, the quantile() function returns the five-number summary. Applying the function to the used car data results in the same statistics as done earlier:					
<pre>> quantile(usedcars\$price)</pre>	92.0				
 If we specify an additional probs parameter using a vector denoting cut points, we can obtain arbitrary quantiles, such as the 1st and 99th percentiles: 					
<pre>> quantile(usedcars\$price, probs = c(0 1% 99% 5428.69 20505.00</pre>	1.01, 0.99))				
 The seq() function is used to generate vectors of evenly-spaced values. This makes it easy to c other slices of data, such as the quintiles (five groups), as shown in the following command: 					
> quantile(usedcars\$price, seq(from = 0% 20% 40% 60% 80% 10 3800.0 10759.4 12993.8 13992.0 149	0, to = 1, by = 0.20)) 00% 99.0 21992.0				
19.03.2025	Eksploracja danych				





Visualizing Numeric Values – Boxplots (1/3)					
Tisualizing numeric variables can be helpful in diagnosing data problems.					
 A common visualization of the five-number summary is boxplot, also known as a box-and-whiskers plot. 					
 The boxplot displays the center and spread of a numeric variable in a format that allows you to quickly obtain a sense of the range and skew of a variable or compare it to other variables. 					
Ict's take a look at a boxplot for the used car price and mileage data.					
 To obtain a boxplot for a variable, we will use the boxplot() function. We will also specify a pair of extra parameters, main and ylab, to add a title to the figure and label the y axis (the vertical axis), respectively. 					
The commands to create the price and mileage boxplots are:					
> boxplot(usedcars\$price, main="Boxplot of Used Car Prices", ylab="Price (\$)")					
> boxplot(usedcars\$mileage, main="Boxplot of Used Car Mileage", ylab="Odometer (mi.)")					
19.03.2025 Eksploracja danych					













Measuring Spread – Variance and Standard Deviation (1/2)

- Distributions allow us to characterize a large number of values using a smaller number of parameters.
- The normal distribution, which describes many types of real-world data, can be defined with just two: center and spread.
 - The center of normal distribution is defined by its mean value, which we have used earlier.
 - The spread is measured by a statistic called the standard deviation.
- In order to calculate the standard deviation, we must first obtain the variance, which is defined as the average of the squared differences between each value and the mean value.
 - In mathematical notation, the variance of a set of *n* values of *x* is defined by the following formula:

$$Var(X) = \sigma^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu)^2$$

• The Greek letter μ denotes the mean of the values, and the variance itself is denoted by σ^2 (sigma squared).

19.03.2025

Eksploracja danych























Correlation Coefficient - Definition					
For pairs of variables measured on an interval or ratio scale, a correlation coefficient r can be calculated as follows:					
$r = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{n\sigma_x \sigma_y}$ x and y - variables, x_i - individual values of x, y_i - individual values of y, \bar{x} - the mean of the x variable, \bar{y} - the mean of the y variable, σ_x and σ_y - the standard deviations of the variables x and y, respectively, n - the number of observations.	 The values of <i>r</i> are within a range of - 1.0 ÷ 1.0 and quantify the linear relationship between the variables. Positive numbers for <i>r</i> indicate a positive correlation between the pair of variables, and negative numbers indicate a negative correlation. A value of <i>r</i> close to 0 indicates little or no relationship between the variables. 				
19.03.2025 Eksploracja	danych				





20







The CrossTable() Output			
The preceding command results in the following table:	Cell Contents N		
 The rows in the table indicate the three models of used cars: SE, SEL, and SES (plus an additional row for the total across all models). 	Chi-square contribution N / Row Total N / Col Total N / Table Total 		
• The columns indicate whether or not the car's color is conservative (plus a column totaling across both types of color).	Total Observations in Table: 150 usedcars\$conservative usedcars\$model FALSE TRUE Row Total		
What we are most interested in is the row proportion for conservative cars for each model.	0.346 0.654 0.520 0.529 0.515 0.180 0.340 		
 The row proportions tell us that 0.654 (65 percent) of SE cars are colored conservatively in comparison to 0.696 (70 percent) of SEL cars and 0.653 (65 percent) of SES. 	0.086 0.044 0.306 0.044 0.304 0.696 0.153 0.137 0.162 0.047 0.107 		
 These differences are relatively small, suggesting that there are no substantial differences in the types of colors chosen by the model of the car. 19.03.2025 Eksploracja danych 	0.007 0.004 0.347 0.653 0.327 0.333 0.323 0.113 0.213 0.113 0.213 Column Total 51 99 150 0.340 0.660 		