

Data Visualization

Table of contents

1 Data Visualization	2
2 ggplot2: An Overview	2
3 Installing & Loading ggplot2	3
4 The Grammar of Graphics in ggplot2	3
5 The mpg Dataset	4
5.1 Dataset Overview	4
5.2 Dataset <code>head()</code>	5
5.3 Dataset <code>summary()</code>	5
5.4 Visualization Objectives	6
5.5 Univariate Analysis: Highway Mileage	6
5.6 Univariate Analysis: Displacement	7
5.7 Bivariate Analysis: Displacement vs. Highway Mileage	8
5.8 Bivariate Analysis: Manufacturer Comparison	9
5.9 Multivariate Analysis: Engine Size, Cylinders, and Fuel Economy	10
6 The diamonds Dataset	11
6.1 Dataset Overview	12
6.2 Dataset <code>head()</code>	13
6.3 Dataset <code>summary()</code>	13
6.4 Visualization Objectives	13
6.5 Univariate Analysis: Price	14
6.6 Univariate Analysis: Carat	14
6.7 Bivariate Analysis: Price vs. Carat	15
6.8 Bivariate Analysis: Cut vs. Price	16
6.9 Multivariate Analysis	17
6.10 Exercise	18

6.11 A bar chart of the <code>cut</code>	19
6.11.1 categorized by <code>clarity</code> , stacked	20
6.11.2 categorized by <code>clarity</code> , clustered	21
6.11.3 categorized by <code>clarity</code> , normalized	22
7 Useful ggplot2 References	23

1 Data Visualization

Data visualization is a key step in the data analysis process that helps to turn abstract data into actionable insights. Data visualization helps in making sense of data, revealing hidden patterns, and communicating results effectively.



2 ggplot2: An Overview

A visualization package in R that uses the grammar of graphics.



3 Installing & Loading ggplot2

- First, we need install the `ggplot2` library

```
1 if (!require(ggplot2)) install.packages("ggplot2")
```

(Typically, required only *once* for a workspace/environment)

- After installation, let us load the `ggplot2` library

```
1 library(ggplot2)
```

4 The Grammar of Graphics in ggplot2

- Explains how to compose graphs from a dataset, aesthetic mappings, and geometric objects.
- A flexible template that forms the foundation of any `ggplot2` graph.

```
1 ggplot(data = <DATA>) +  
2   <GEOM_FUNCTION>(mapping = aes(<MAPPINGS>))
```

5 The mpg Dataset

The `mpg` dataset contains fuel economy data for 234 cars ranging from 1999 to 2008.



Figure 1: US Environmental Protection Agency (EPA)

5.1 Dataset Overview

A look at the dataset's variables: manufacturer, model, displacement, year, number of cylinders, and fuel economy measurements.

Variable	Description
manufacturer	Name of the car manufacturer
model	Model of the car
displ	Engine displacement, in liters
year	Year of manufacture
cyl	Number of cylinders
trans	Type of transmission
drv	Type of drive train (f = front-wheel, r = rear-wheel, 4 = 4wd)
cty	City miles per gallon
hwy	Highway miles per gallon
fl	Fuel type
class	Vehicle class

5.2 Dataset head()

```
1 head(mpg)
```

manufacturer	model	displ	year	cyl	trans	drv	cty	hwy	fl	class
audi	a4	1.8	1999	4	auto(l5)	f	18	29	p	compact
audi	a4	1.8	1999	4	manual(m5)	f	21	29	p	compact
audi	a4	2.0	2008	4	manual(m6)	f	20	31	p	compact
audi	a4	2.0	2008	4	auto(av)	f	21	30	p	compact
audi	a4	2.8	1999	6	auto(l5)	f	16	26	p	compact
audi	a4	2.8	1999	6	manual(m5)	f	18	26	p	compact

5.3 Dataset summary()

```
1 summary(mpg)
```

manufacturer	model	displ	year	cyl	trans	drv	cty	hwy	fl	class
Length:234	Length:234	Min. :1.600	Min. :1999	Min. :4.000	Length:234	Length:234	Min. :12.00	Length:234	Length:234	
							9.00			
Class :character	Class :character	1st Qu.:2.400	1st Qu.:1999	1st Qu.:4.000	Class :character	Class :character	1st Qu.:14.00	1st Qu.:18.00	Class :character	Class :character

manufacturer	model	displ	year	cyl	trans	drv	cty	hwy	fl	class
Mode	Mode	Median	Median	Median	Mode	Mode	Median	Median	Mode	Mode
:character	:character	:3.300	:2004	:6.000	:character	:character	:17.00	:24.00	:character	:character
NA	NA	Mean	Mean	Mean	NA	NA	Mean	Mean	NA	NA
		:3.472	:2004	:5.889			:16.86	:23.44		
NA	NA	3rd	3rd	3rd	NA	NA	3rd	3rd	NA	NA
		Qu.:4.600	Qu.:2000	Qu.:8.000			Qu.:19.000	Qu.:27.00		
NA	NA	Max.	Max.	Max.	NA	NA	Max.	Max.	NA	NA
		:7.000	:2008	:8.000			:35.00	:44.00		

5.4 Visualization Objectives

- **Fuel Economy Trends:** Identify trends in fuel efficiency over time.
- **Displacement vs. Mileage:** Understand how engine size affects fuel economy.
- **Manufacturer Comparison:** Compare different manufacturers based on fuel economy.

5.5 Univariate Analysis: Highway Mileage

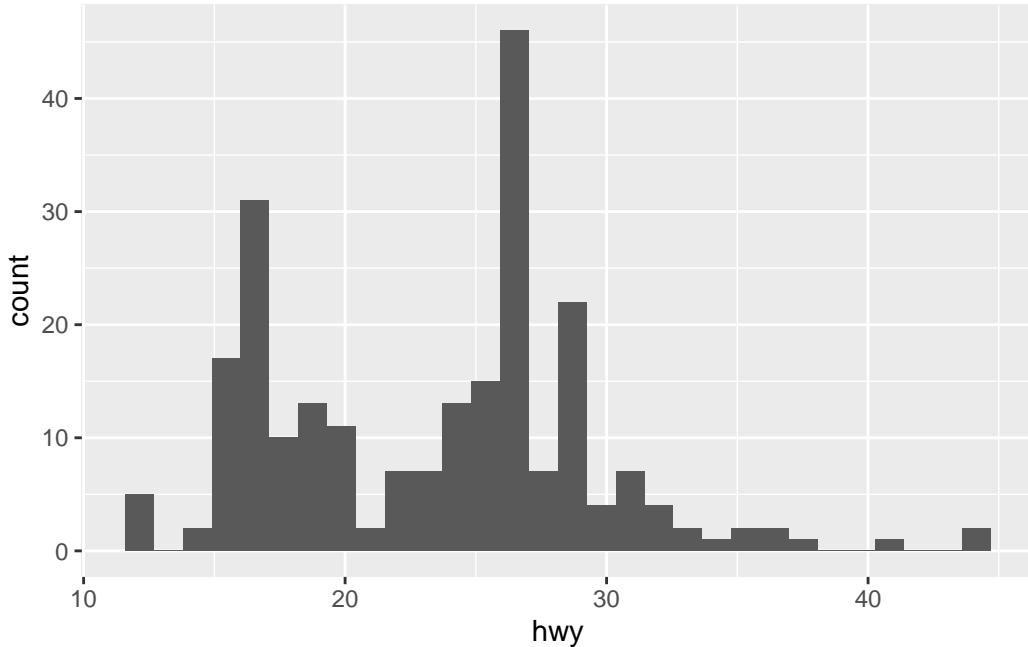
Analyzing the highway miles per gallon (hwy) distribution.

```

1 ggplot(mpg, aes(x = hwy)) +
2   geom_histogram()

```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

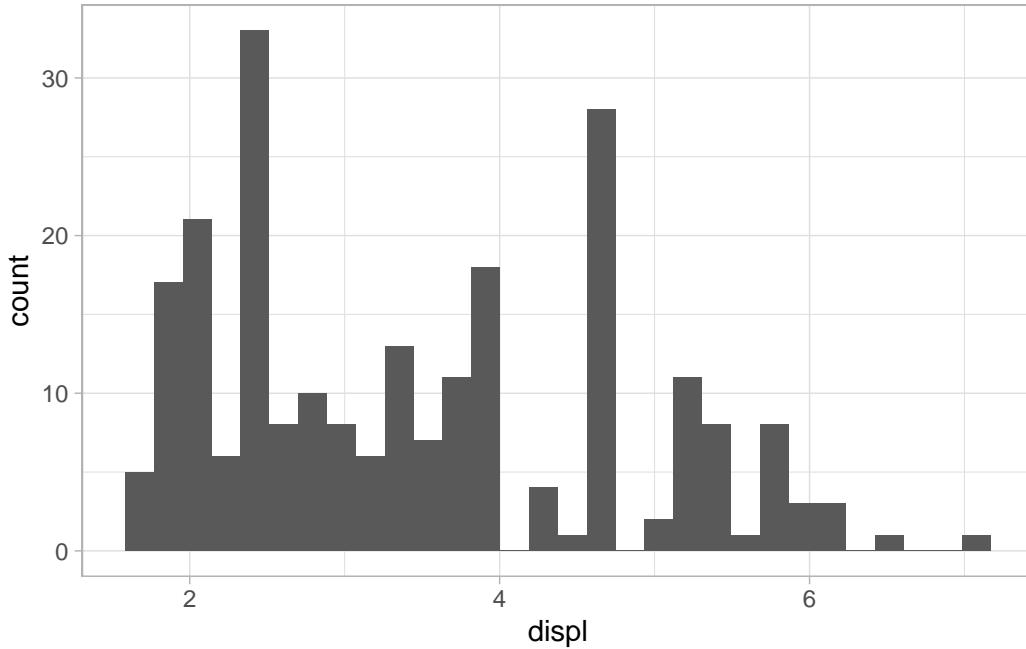


- **Insight:** The majority of cars have highway mileage ratings between 25 to 35 miles per gallon.

5.6 Univariate Analysis: Displacement

Exploring engine displacement (displ) across all cars.

```
1 ggplot(mpg, aes(x = displ)) +  
2   geom_histogram() +  
3   theme_light()  
  
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

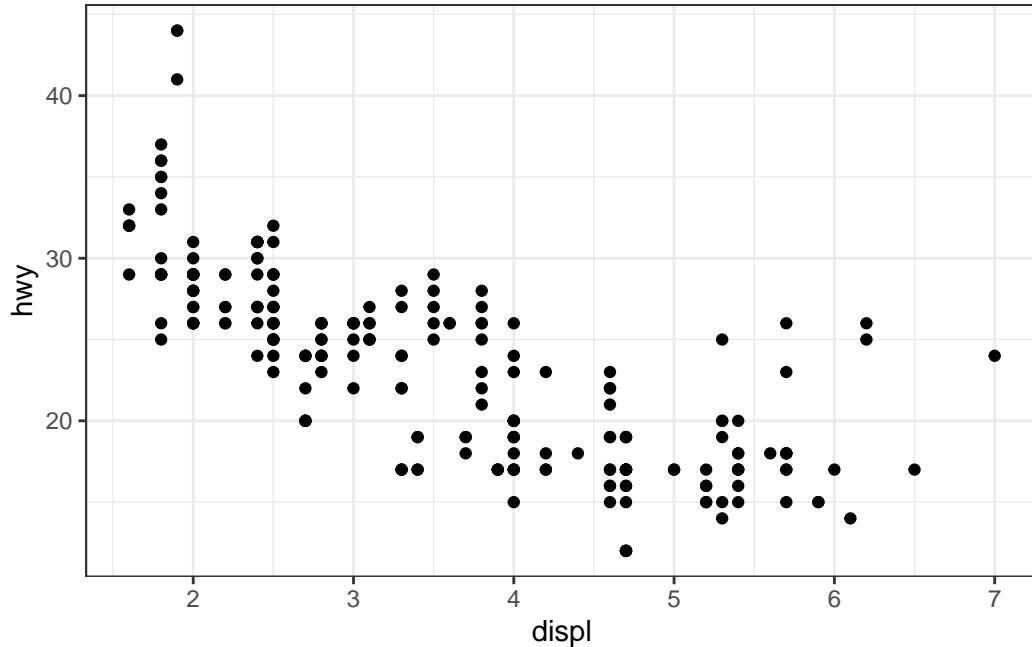


- **Insight:** Smaller engine displacements are more common, suggesting a concentration of more fuel-efficient vehicles.

5.7 Bivariate Analysis: Displacement vs. Highway Mileage

Investigating the relationship between engine displacement (`displ`) and highway mileage (`hwy`).

```
1 ggplot(mpg, aes(x = displ, y = hwy)) +  
2   geom_point() +  
3   theme_bw()
```

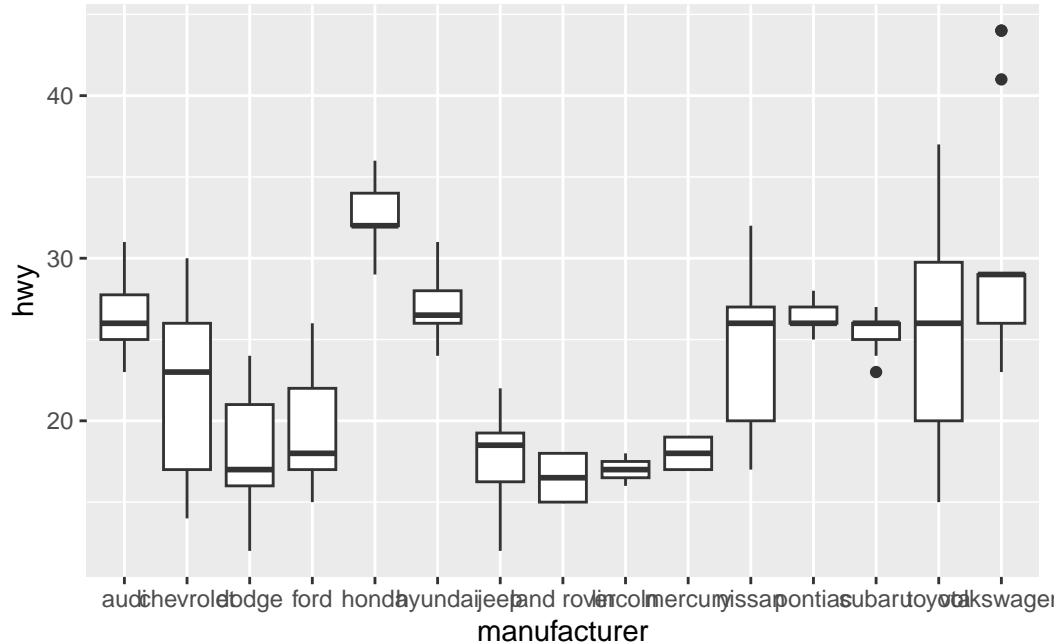


- **Insight:** There is a clear negative trend indicating that larger engine sizes are associated with lower highway mileage.

5.8 Bivariate Analysis: Manufacturer Comparison

Comparing fuel economy across manufacturers with boxplots.

```
1 ggplot(mpg, aes(x = manufacturer, y = hwy)) +  
2   geom_boxplot()
```



- **Insight:** There is significant variability in highway mileage among manufacturers, with some consistently outperforming others.

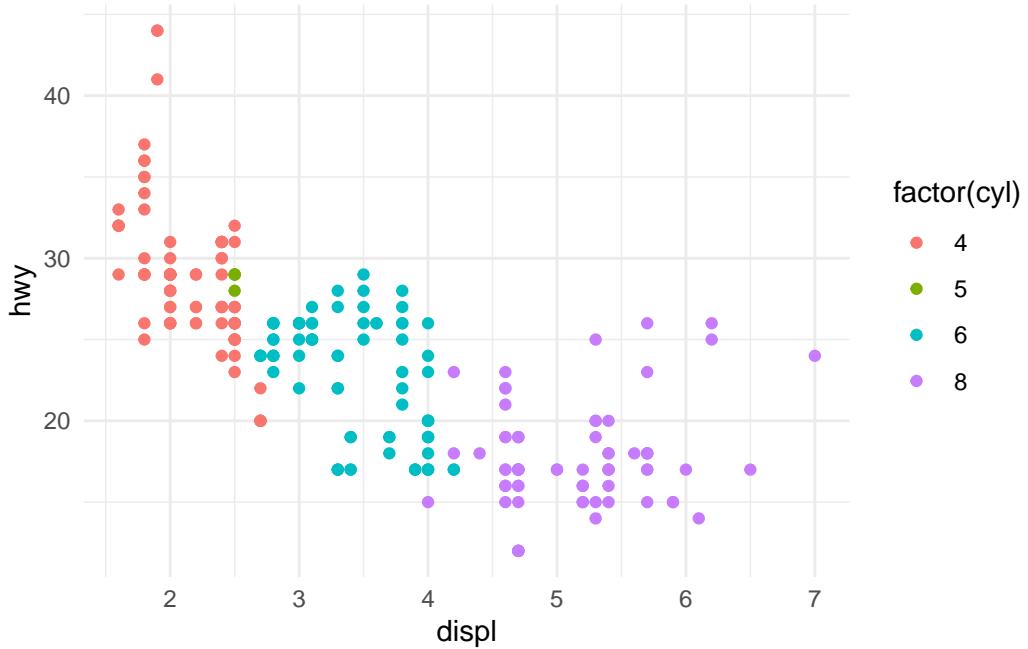
5.9 Multivariate Analysis: Engine Size, Cylinders, and Fuel Economy

Analyzing how engine size and the number of cylinders relate to fuel economy.

```

1 ggplot(mpg, aes(x = displ, y = hwy, color = factor(cyl))) +
2   geom_point() +
3   theme_minimal()

```



- **Insight:** Cars with fewer cylinders tend to have better highway mileage, despite the size of the engine.

6 The diamonds Dataset

The diamonds dataset contains prices and attributes of approximately 54,000 round-cut diamonds.



6.1 Dataset Overview

A quick overview of the dataset's structure and variables: `carat`, `cut`, `color`, `clarity`, `depth`, `table`, `price`, `x` (length), `y` (width), and `z` (depth).

Variable	Description
<code>carat</code>	Weight of the diamond
<code>cut</code>	Quality of the cut (Fair, Good, Very Good, Premium, Ideal)
<code>color</code>	Diamond color, from J (worst) to D (best)
<code>clarity</code>	How clear the diamond is (I1, SI2, SI1, VS2, VS1, VVS2, VVS1, IF)
<code>depth</code>	Depth percentage = $\left(\frac{z}{\frac{x+y}{2}} \right) \times 100$
<code>table</code>	Width of the top of the diamond relative to the widest point
<code>price</code>	Price in US dollars
<code>x</code>	Length in mm
<code>y</code>	Width in mm
<code>z</code>	Depth in mm

6.2 Dataset head()

```
1 head(diamonds)
```

carat	cut	color	clarity	depth	table	price	x	y	z
0.23	Ideal	E	SI2	61.5	55	326	3.95	3.98	2.43
0.21	Premium	E	SI1	59.8	61	326	3.89	3.84	2.31
0.23	Good	E	VS1	56.9	65	327	4.05	4.07	2.31
0.29	Premium	I	VS2	62.4	58	334	4.20	4.23	2.63
0.31	Good	J	SI2	63.3	58	335	4.34	4.35	2.75
0.24	Very Good	J	VVS2	62.8	57	336	3.94	3.96	2.48

6.3 Dataset summary()

```
1 summary(diamonds)
```

carat	cut	color	clarity	depth	table	price	x	y	z
Min.	Fair :	D:	SI1	Min.	Min.	Min. :	Min. :	Min. :	Min. :
:0.2000	1610	6775	:13065	:43.00	:43.00	326	0.000	0.000	0.000
1st	Good :	E:	VS2	1st	1st	1st	1st	1st	1st
Qu.:0.4000906	9797	:12258	Qu.:61.00	Qu.:56.00	Qu.:950	Qu.:950	Qu.:4.710	Qu.:4.720	Qu.:2.910
Median	Very	F:	SI2 :	Median	Median	Median	Median	Median	Median
:0.7000	Good:1209542	9194	:1209542	:61.80	:57.00	:2401	:5.700	:5.710	:3.530
Mean	Premium	G:11292	S1 :	Mean	Mean	Mean :	Mean :	Mean :	Mean :
:0.7979	:13791	8171	:13791	:61.75	:57.46	3933	5.731	5.735	3.539
3rd	Ideal	H:	VVS2	3rd	3rd	3rd	3rd	3rd	3rd
Qu.:1.040021551	8304	:5066	Qu.:62.50	Qu.:59.00	Qu.:5324	Qu.:5324	Qu.:6.540	Qu.:6.540	Qu.:4.040
Max.	NA	I:	VVS1	Max.	Max.	Max.	Max.	Max.	Max.
:5.0100		5422	:3655	:79.00	:95.00	:18823	:10.740	:58.900	:31.800
NA	NA	J:	(Other):	NA	NA	NA	NA	NA	NA
		2808	2531						

6.4 Visualization Objectives

- **Price Distribution:** Analyzing the distribution of diamond prices.

- **Carat Size Analysis:** Understanding how carat size affects price.
- **Cut Quality:** Examining the relationship between diamond cut and other attributes.

6.5 Univariate Analysis: Price

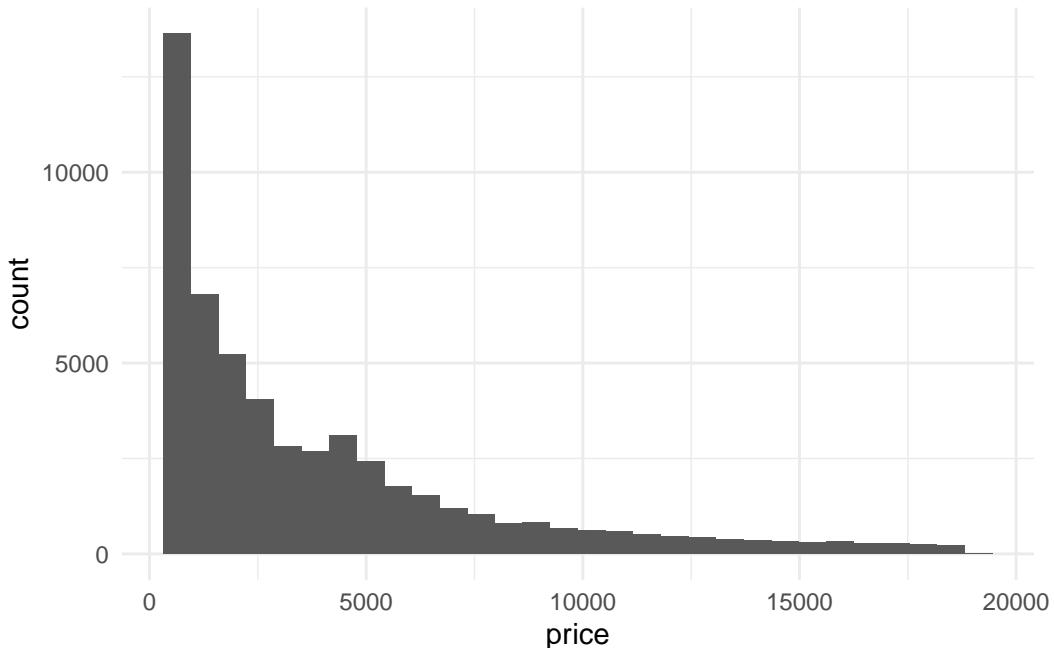
Visualizing the distribution of diamond prices using histograms and density plots.

```

1 ggplot(diamonds, aes(x = price)) +
2   geom_histogram() +
3   theme_minimal()

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

```



6.6 Univariate Analysis: Carat

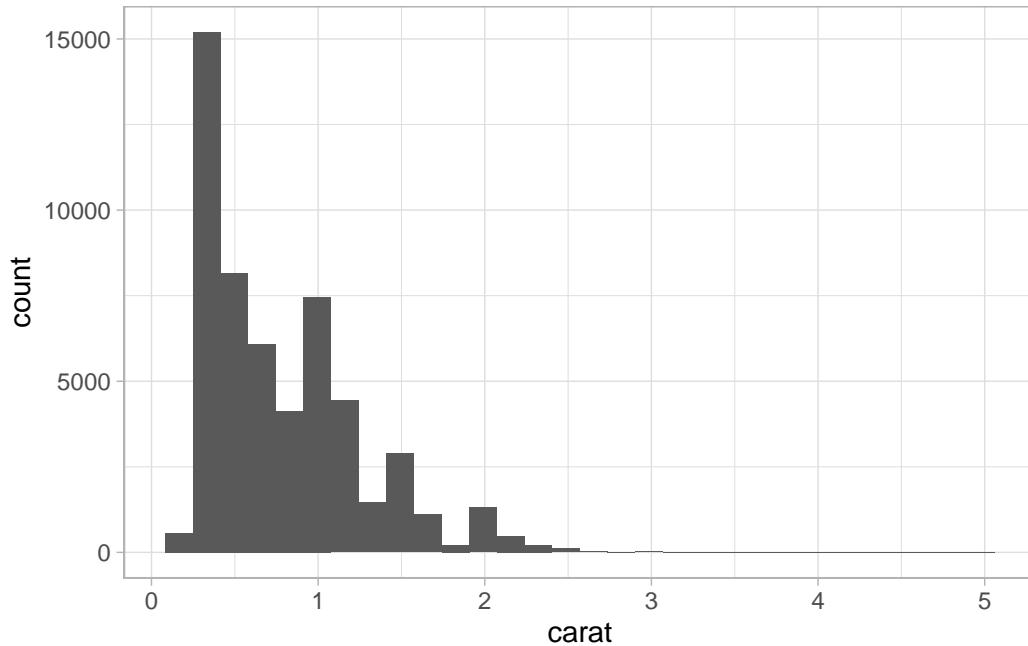
Understanding carat distribution with a histogram.

```

1 ggplot(diamonds, aes(x = carat)) +
2   geom_histogram() +
3   theme_light()

```

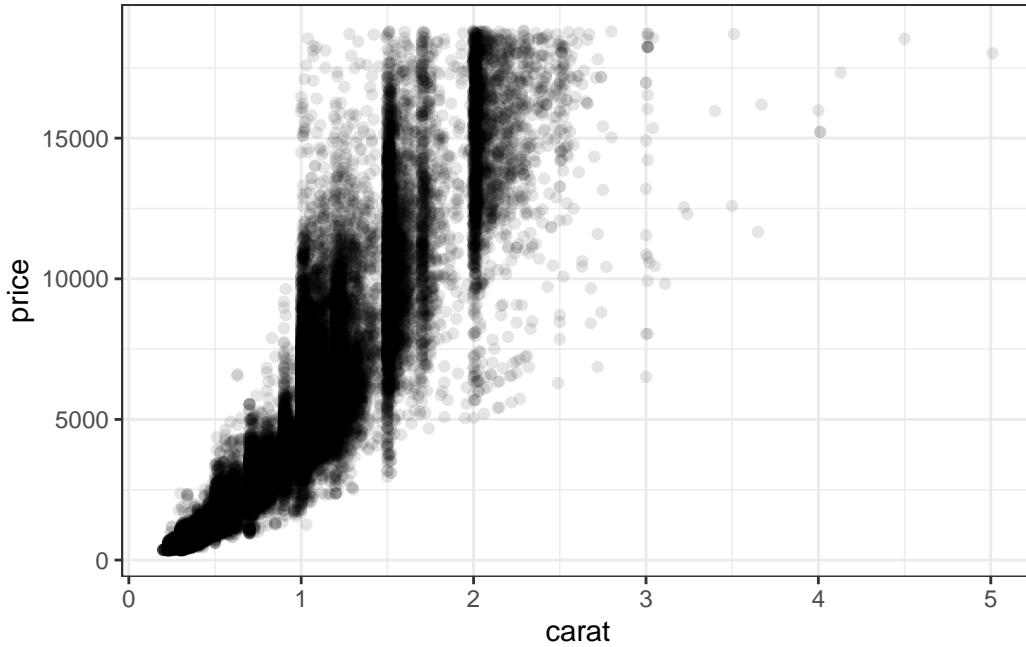
```
`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



6.7 Bivariate Analysis: Price vs. Carat

Exploring the relationship between price and carat size with scatter plots.

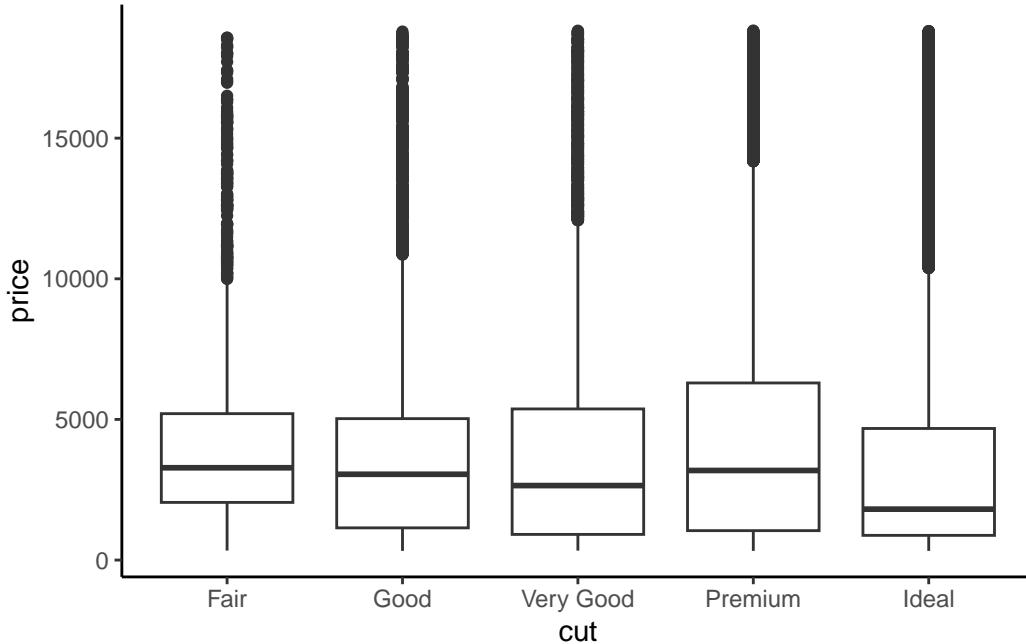
```
1 ggplot(diamonds, aes(x = carat, y = price)) +  
2   geom_point(alpha = 0.1) +  
3   theme_bw()
```



6.8 Bivariate Analysis: Cut vs. Price

Comparing average price across different cuts using boxplots.

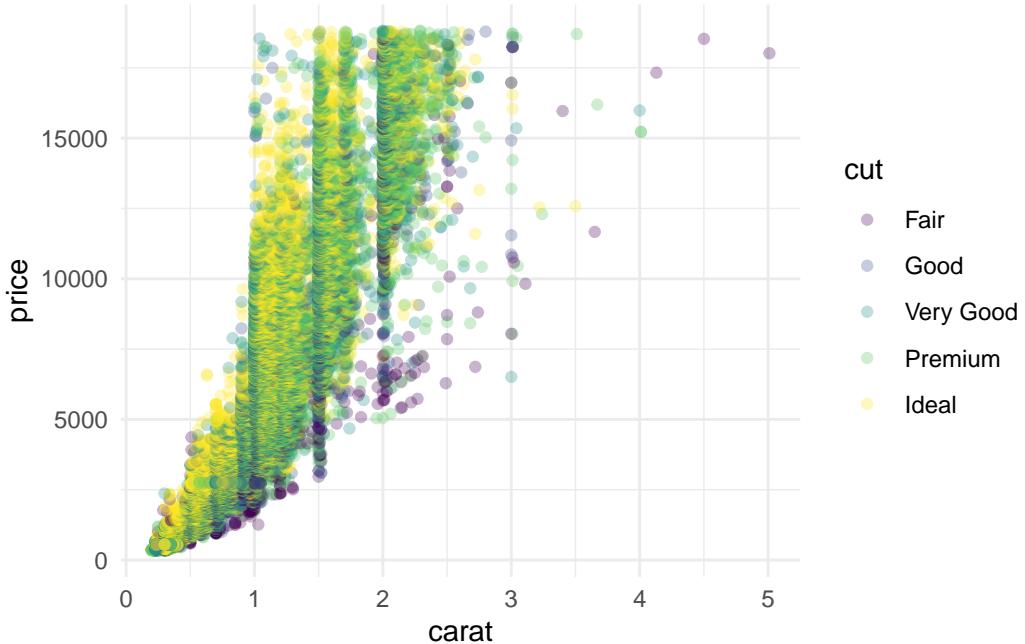
```
1 ggplot(diamonds, aes(x = cut, y = price)) +  
2   geom_boxplot() +  
3   theme_classic()
```



6.9 Multivariate Analysis

Incorporating multiple attributes into our visual analysis, like cut and color versus price.

```
1 ggplot(diamonds, aes(x = carat, y = price, color = cut)) +  
2   geom_point(alpha = 0.3) +  
3   theme_minimal()
```



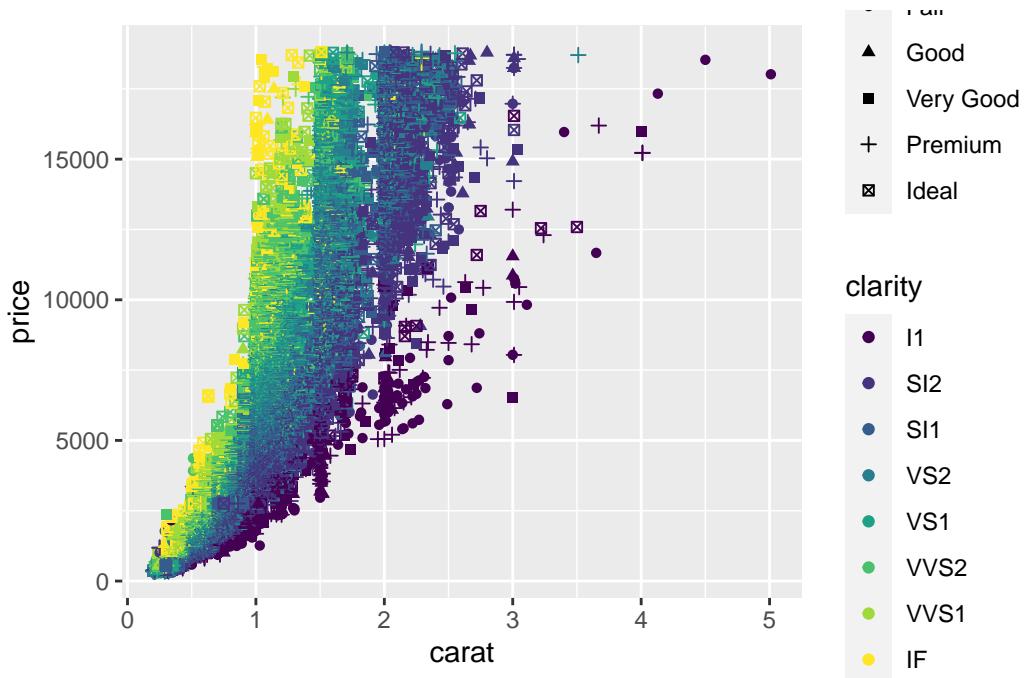
6.10 Exercise

Generate a plot to show the relationship between the following variables from the diamonds dataset: `carat`, `cut`, `clarity`, and `price`.

- Solution

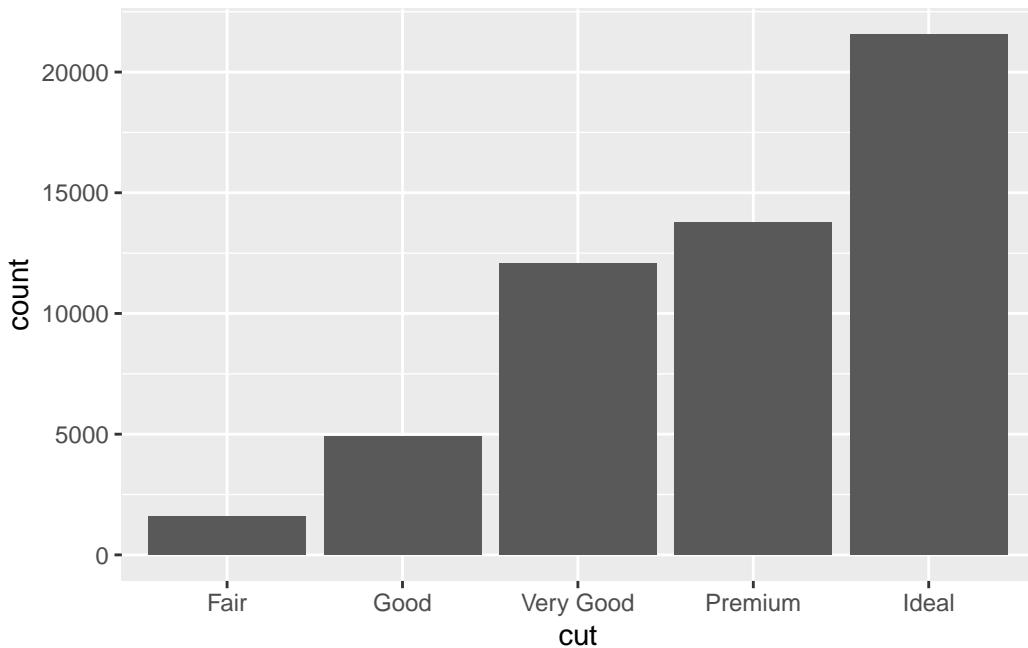
```
1 ggplot(diamonds) +
2   geom_point(aes(x = carat, y = price, color = clarity, shape = cut))
```

Warning: Using shapes for an ordinal variable is not advised



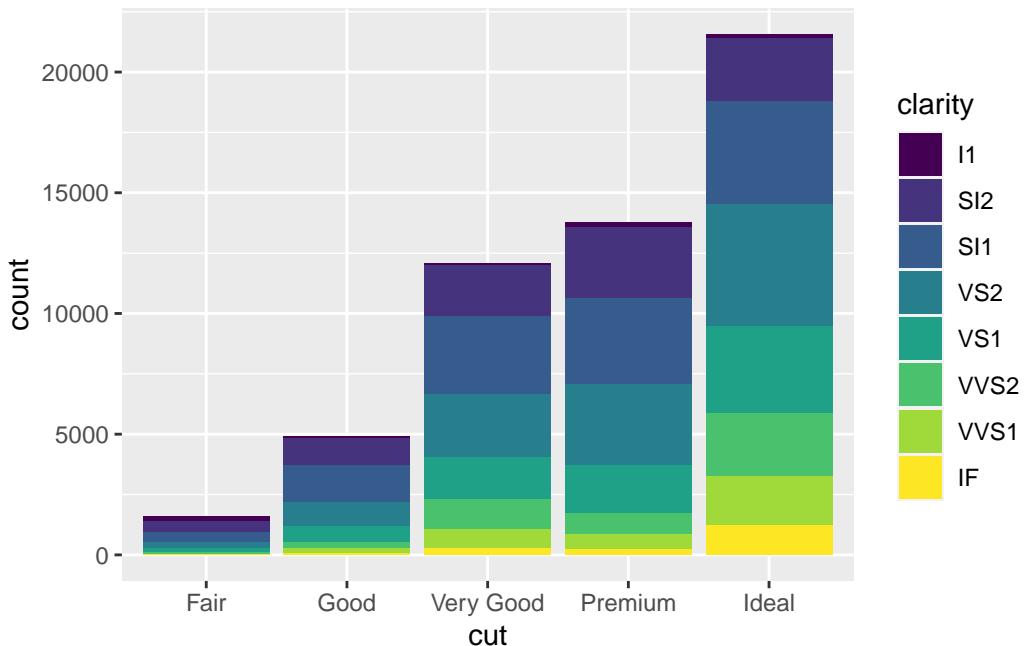
6.11 A bar chart of the cut

```
1 ggplot(data = diamonds) +  
2   geom_bar(mapping = aes(x = cut))
```



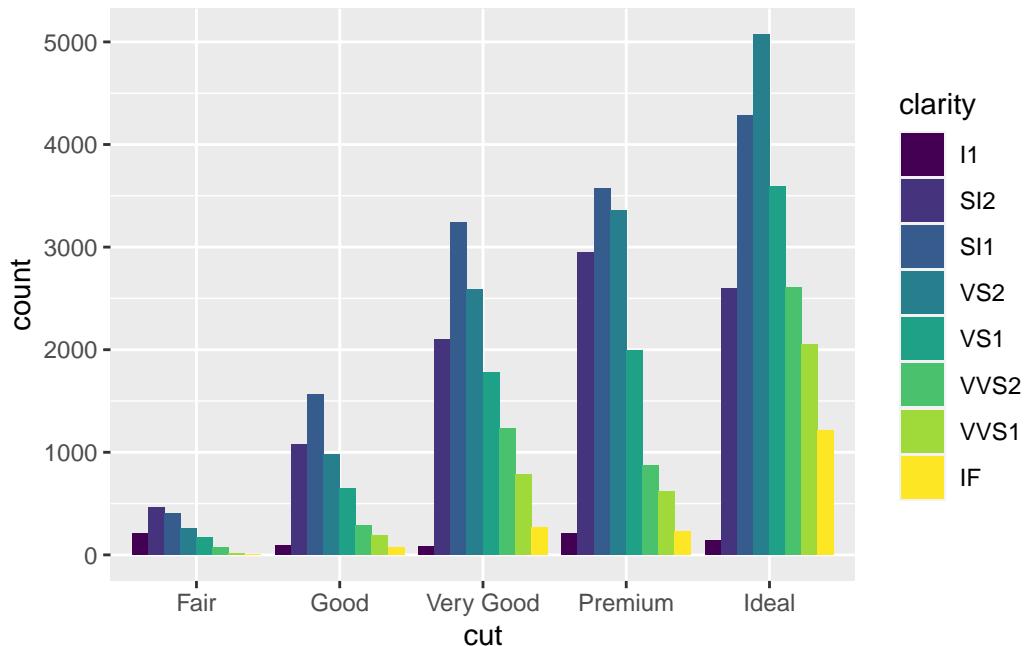
6.11.1 categorized by clarity, stakced

```
1 ggplot(data = diamonds) +  
2   geom_bar(mapping = aes(x = cut, fill = clarity))
```



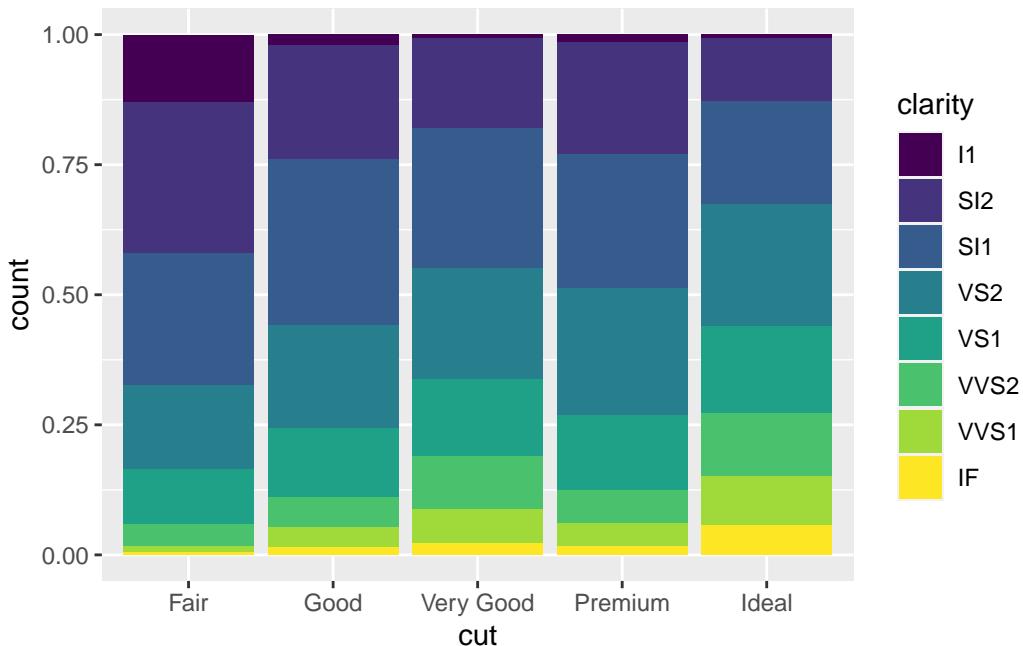
6.11.2 categorized by clarity, clustered

```
1 ggplot(data = diamonds) +  
2   geom_bar(mapping = aes(x = cut, fill = clarity), position = "dodge")
```



6.11.3 categorized by clarity, normalized

```
1 ggplot(data = diamonds) +  
2   geom_bar(mapping = aes(x = cut, fill = clarity), position = "fill")
```



7 Useful ggplot2 References

- [ggplot2 Documentation](#)
- [ggplot2 Cheat Sheet](#)
- [R Graphics Cookbook](#)
- [ggplot2 Book](#)
- [ggplot2 Extension Gallery](#)
- [ggplot2 Example Gallery](#)
- [ggplot2 Visualization Catalog](#)