

# Exploratory Data Analysis

MKT 566

Instructor: Davide Proserpio

# Before we start

- Groups
- Homework

# What we will learn

How to use visualization to explore your data in a systematic way  
(also called **Exploratory Data Analysis** or **EDA**)

- Generate questions about your data
- Search for answers by visualizing, transforming, and modelling your data
- Use what you learn to refine your questions and/or generate new questions.

(Partially based on [Chapter 7 of R for Data Science](#))

# EDA Goal

- There is no rule about which questions you should ask to guide your research.
- However, two types of questions will always be useful for making discoveries within your data. You can loosely word these questions as:
  - What type of **variation** occurs within my variables?
  - What type of **covariation** occurs between my variables?

# Covariation

- **Covariation** is the tendency for the values of two or more variables to vary together in a related way
- The best way to spot covariation is to **visualize the relationship between two or more variables**
- How you do that should again depend on the type of variables involved

# Visualizing covariation

Example with the [marketing](#) dataset from the library ‘datarium’

```
> head(marketing)
  youtube facebook newspaper sales
1  276.12    45.36    83.04  26.52
2   53.40    47.16    54.12  12.48
3   20.64    55.08    83.16  11.16
4  181.80    49.56    70.20  22.20
5  216.96    12.96    70.08  15.48
6   10.44    58.68    90.00   8.64
```

# A categorical and continuous variable

How can we visualize sales by ad spend?

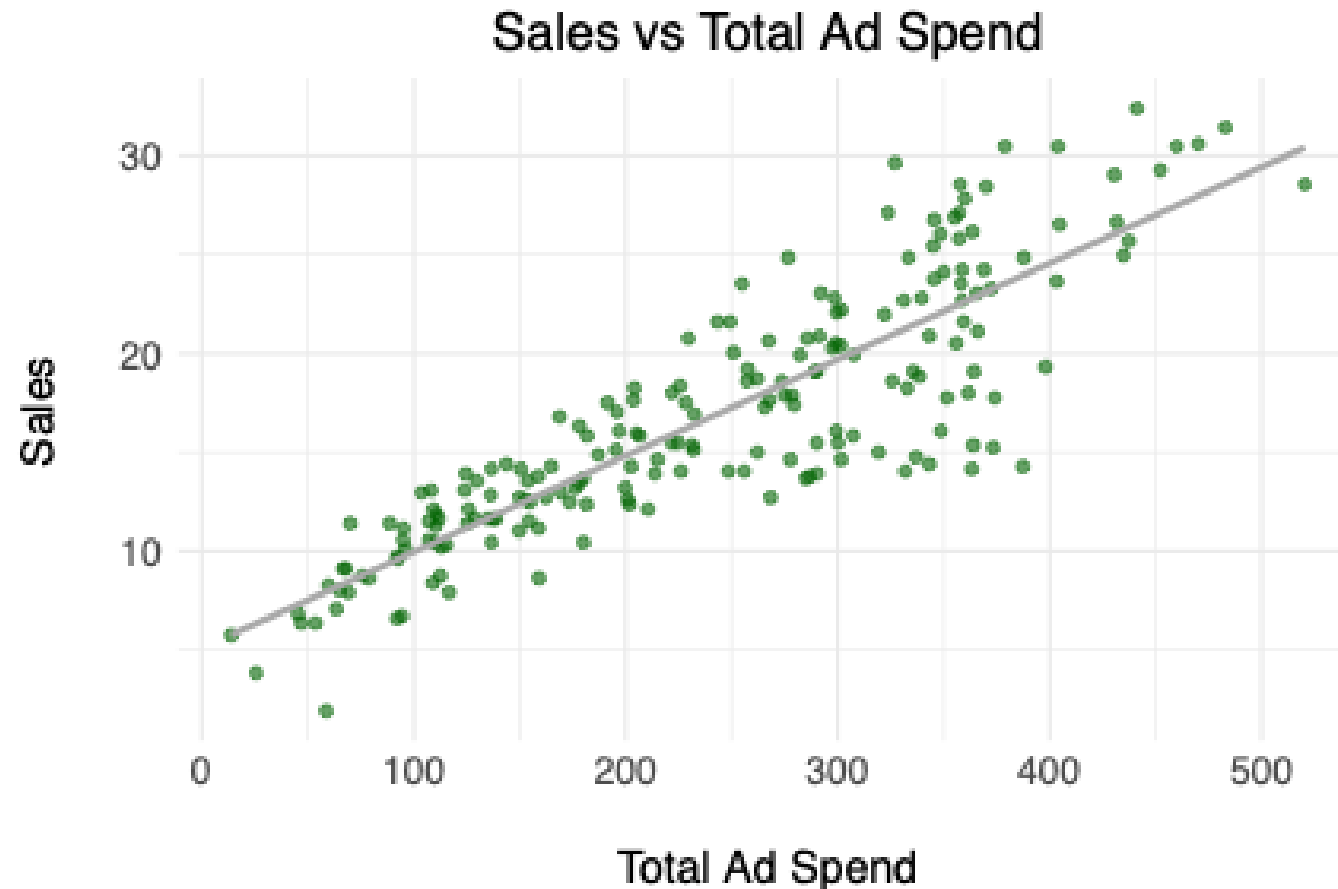
# Two continuous variables

How can we visualize sales by ad spend?



# Two continuous variables

How can we visualize sales by ad spend?

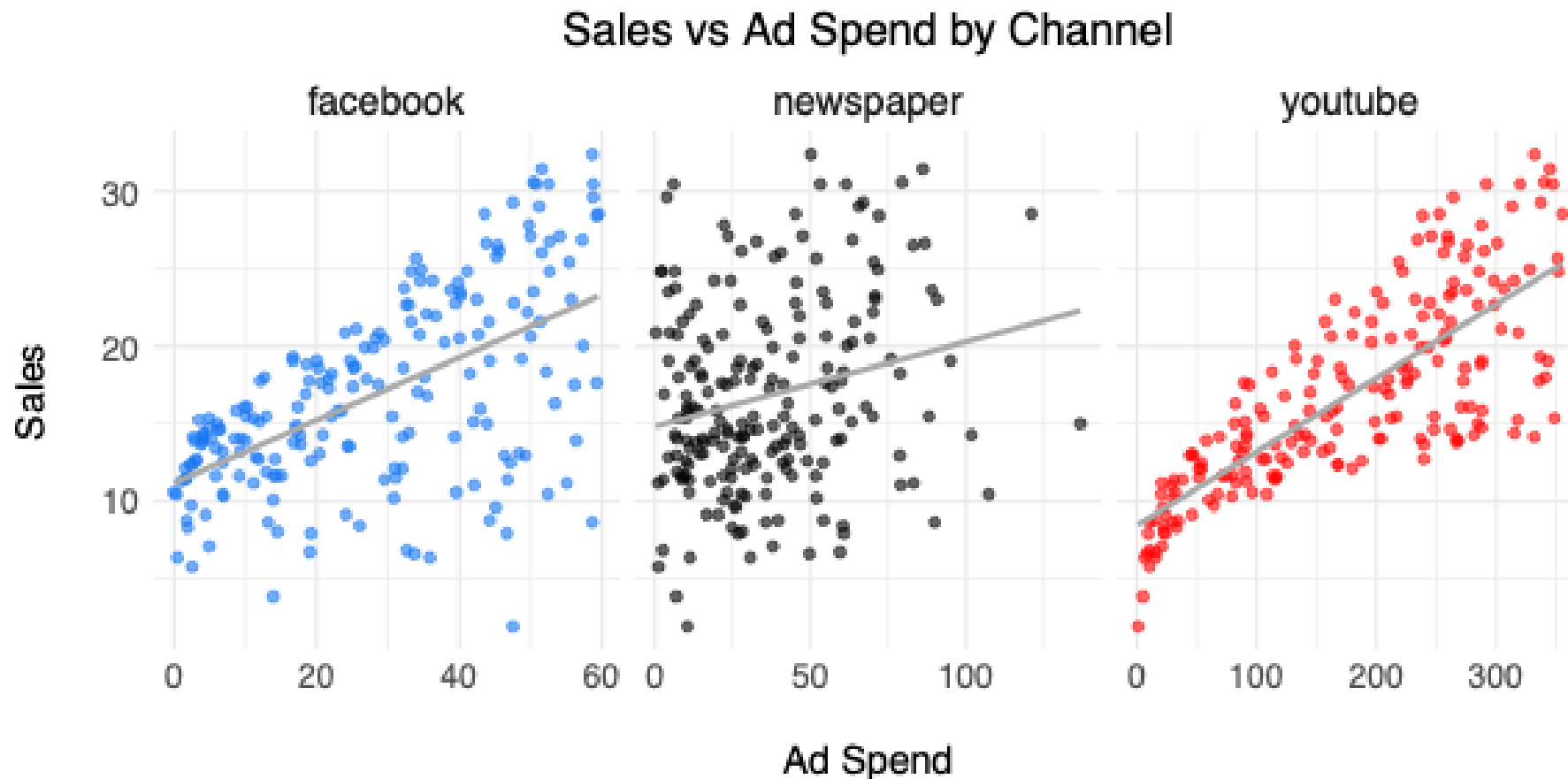


# Two continuous variables

Can we do a more informative viz?

# Two continuous variables

Can we do a more informative viz?



# One categorical and one continuous variable

Let's use the simulated marketing dataset we explored last week

```
> head(df)
```

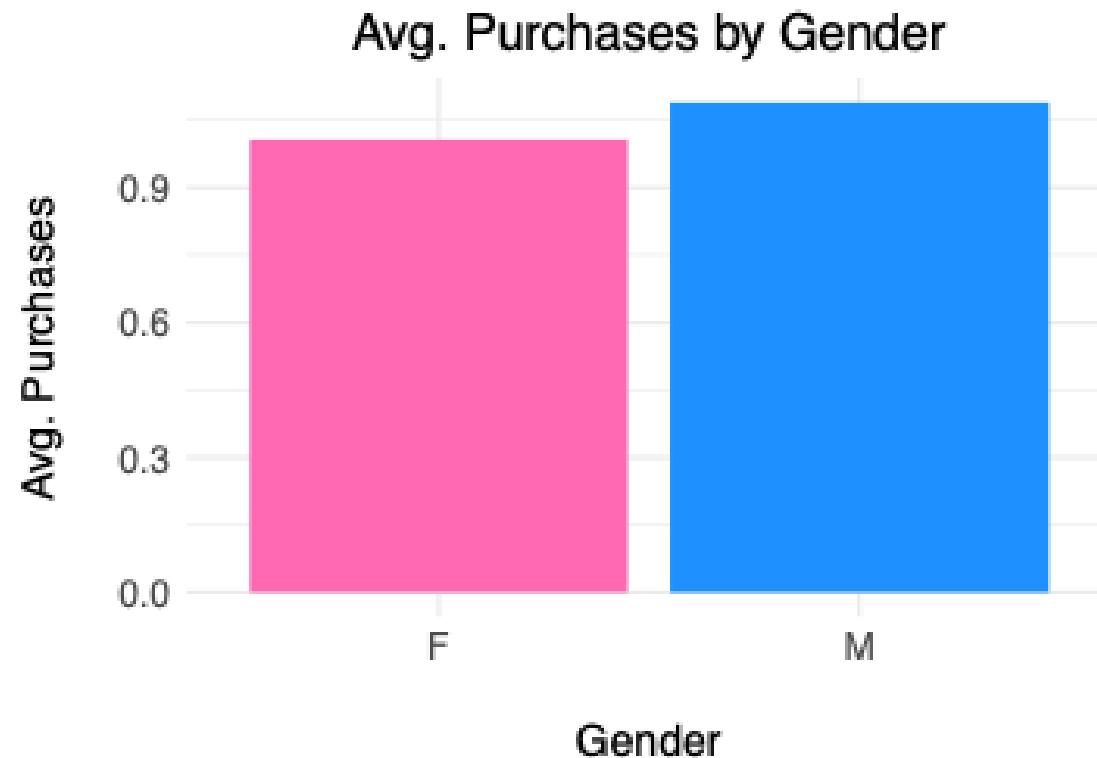
	CustomerID	Age	Gender	Device	Channel	Ad_Spend	Clicks	Purchases	Revenue
	<int>	<int>	<char>	<char>	<char>	<num>	<int>	<int>	<num>
1:	1	54	M	Mobile	Social	718.60	95	6	149.16
2:	2	18	F	Mobile	Search	233.00	34	1	22.22
3:	3	42	F	Mobile	Search	122.51	18	0	0.00
4:	4	27	F	Desktop	Social	198.78	19	1	13.22
5:	5	53	F	Mobile	Social	145.19	19	4	150.48
6:	6	35	M	Desktop	Video	125.74	9	0	0.00

# One categorical and one continuous variable

Which viz can we use to explore the relationship between purchases and gender?

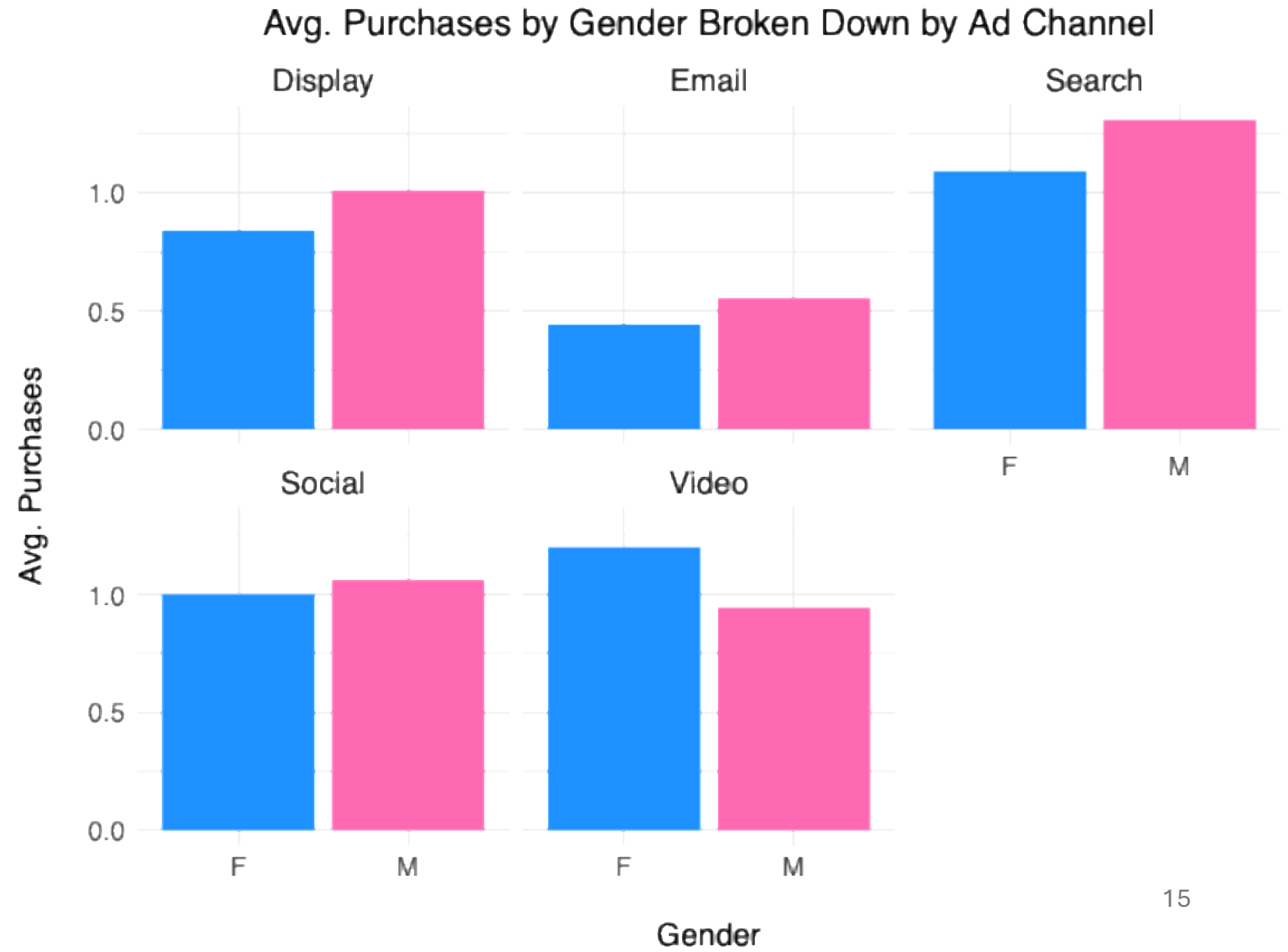
# One categorical and one continuous variable

Which viz can we use to explore the relationship between purchases and gender?



# One categorical and one continuous variable

Let's add an additional dimension: Ad Channel



# RateBeer data viz exercise

Code:

- RateBeer case: [html](#) and [R Markdown, dataset](#)

```
> head(beer)
```

	beer_name	beer_beerId	beer_brewerId	beer_ABV	beer_style	review_appearance	review_aroma	review_palate	review_taste	review_overall	review_time	review_profileName
	<char>	<char>	<int>	<char>	<char>	<char>	<char>	<char>	<char>	<char>	<int>	<char>
1:	John Harvards Fancy Lawnmower Beer	64125	8481	5.4	Klsch	2/5	4/10	2/5	4/10	8/20	1157587200	hopdog
2:	Barley Island Dirty &quot;Old&quot; Helen Sour Ale	114513	3228	-	Sour Ale/Wild Ale	4/5	8/10	4/5	8/10	17/20	1266019200	MI2CA
3:	Barley Island Sinister Minister Belgian Black Ale	77833	3228	6.7	Traditional Ale	3/5	6/10	3/5	6/10	16/20	1237420800	emacgee
4:	Barley Island Sinister Minister Belgian Black Ale	77833	3228	6.7	Traditional Ale	4/5	6/10	4/5	6/10	14/20	1229040000	after4ever
5:	Barley Island Sinister Minister Belgian Black Ale	77833	3228	6.7	Traditional Ale	3/5	6/10	3/5	6/10	12/20	1222041600	Sparky
6:	Barley Island Sinister Minister Belgian Black Ale	77833	3228	6.7	Traditional Ale	3/5	5/10	3/5	6/10	13/20	1221264000	jsquire