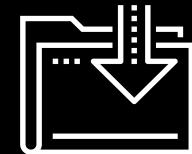




Introduction to Data

Data Boot Camp
Lesson 1.1





WELCOME

The Rise of Data



Why is data analytics such a
hot skill these days?

1

Explosive
growth in
digitised
data
(creation)



2

Explosive growth in analytic tools (synthesis)



Best Restaurants in San Francisco, CA

Showing 1-20 of 7848

3

Accelerating search for actionable insight (value)



1. Aracy Cafe

113 reviews

\$\$ · American (New), Venues & Event

This restaurant takes reservations

Get 7% Cash Back when you dine here

I seriously hope all the hipster jerks from the city don't find this place and ruin it. I hesitate to post this review so as not to bring them to the presence of this place. But since it's... [read more](#)

2. Derm Restaurant

34 reviews

\$\$ · Thai

This restaurant takes reservations

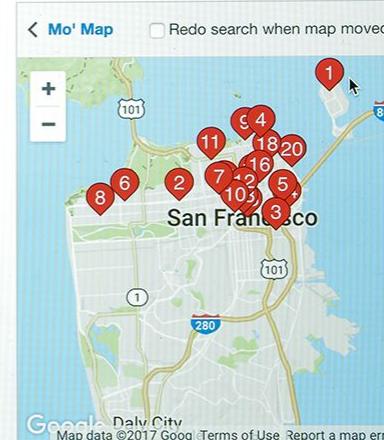
This restaurant accepts pickup orders

401 13th St
San Francisco, CA 94130
(415) 985-7117

[Find a Table](#)[Enroll in Cash Back](#)

Laurel Heights

3226 Geary Blvd
San Francisco, CA 94118
(415) 379-4549

[Find a Table](#)[Start Order](#)

Ads by Google

salutemarinabay.com ▾

Salute E Vita - Italian Food - Waterfront Dining

Enjoy authentic Italian dishes with breathtaking views of San Francisco Bay.
1900 Esplanade Dr. Richmond, CA



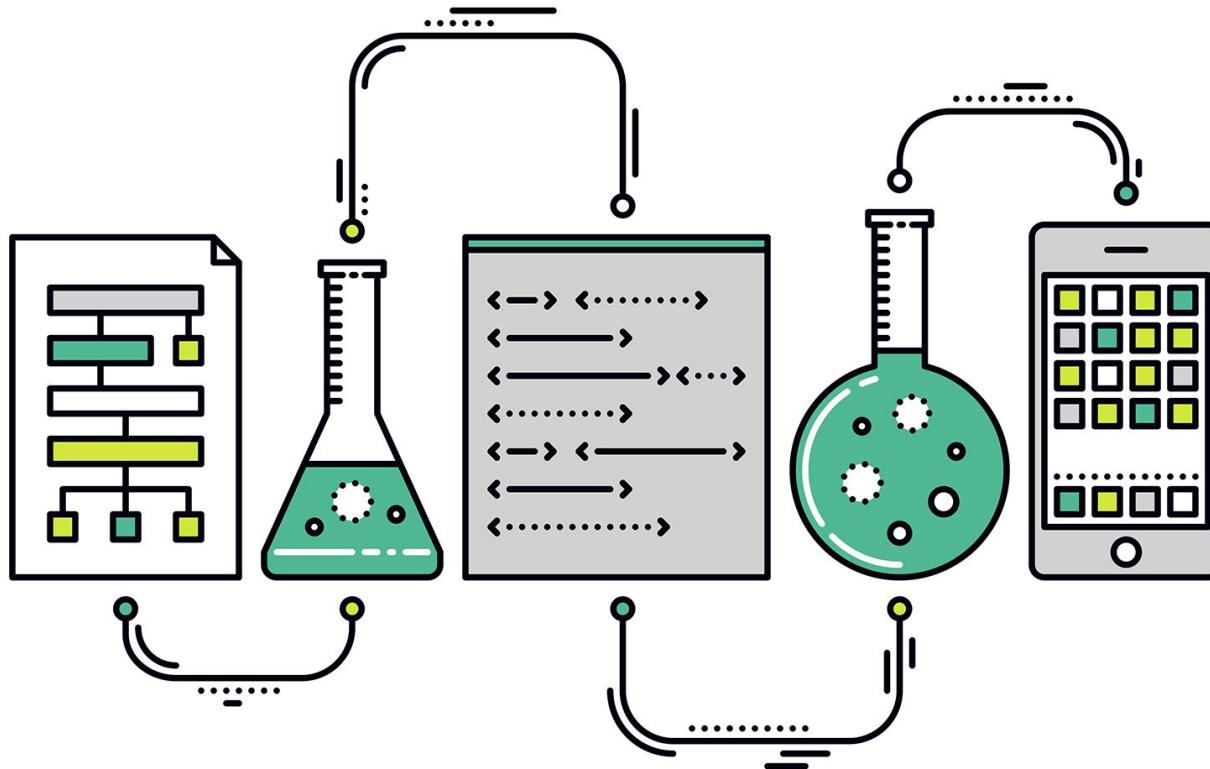


What does the term **data analytics** mean?

Perhaps you
are picturing
an Excel
spreadsheet.

5.3649	5.00E+05	4.93E-05
3.1631	5.26E-05	4.7659
8.1376	4.93E+05	3.1631
9.0365	5.20E+05	4.7659
5.4273	4.92E+05	8.1376
11.251	5.06E+05	9.0365
9.36018	5.06E+05	5.4273
2.5538	5.05E+05	11.251
1.6593	4.93E+05	9.36018
1.66169	5.12E+05	2.5538
1.66169	5.10E+05	1.6593
1.66169	5.12E+05	1.66169
1.66169	5.10E+05	1.66169

Data Analytics Involves Spreadsheets and Formulas





Fundamentally, data analytics is
about storytelling and truth-telling.

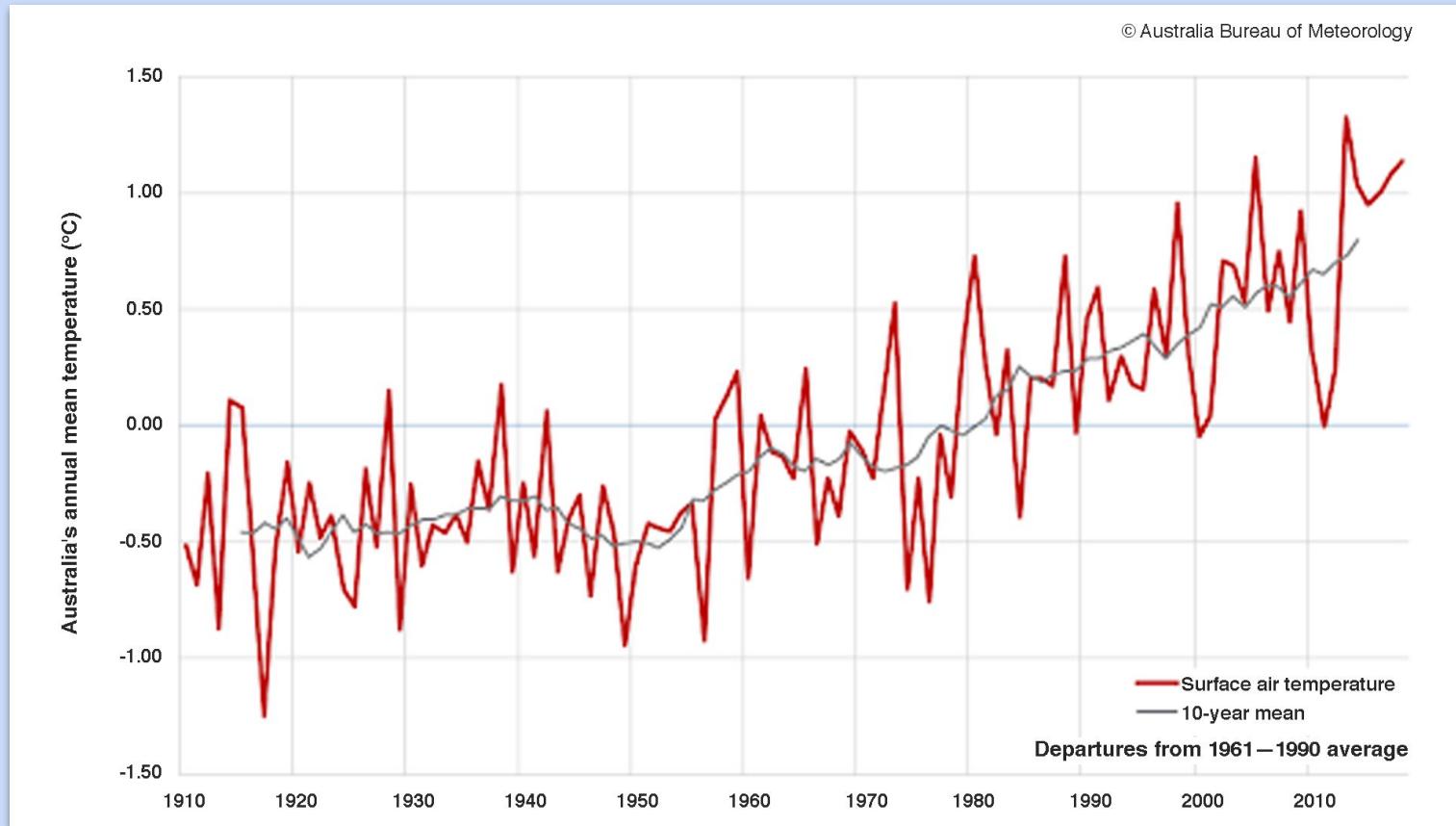
Data as Storytelling

Data as Storytelling

Australia GDP



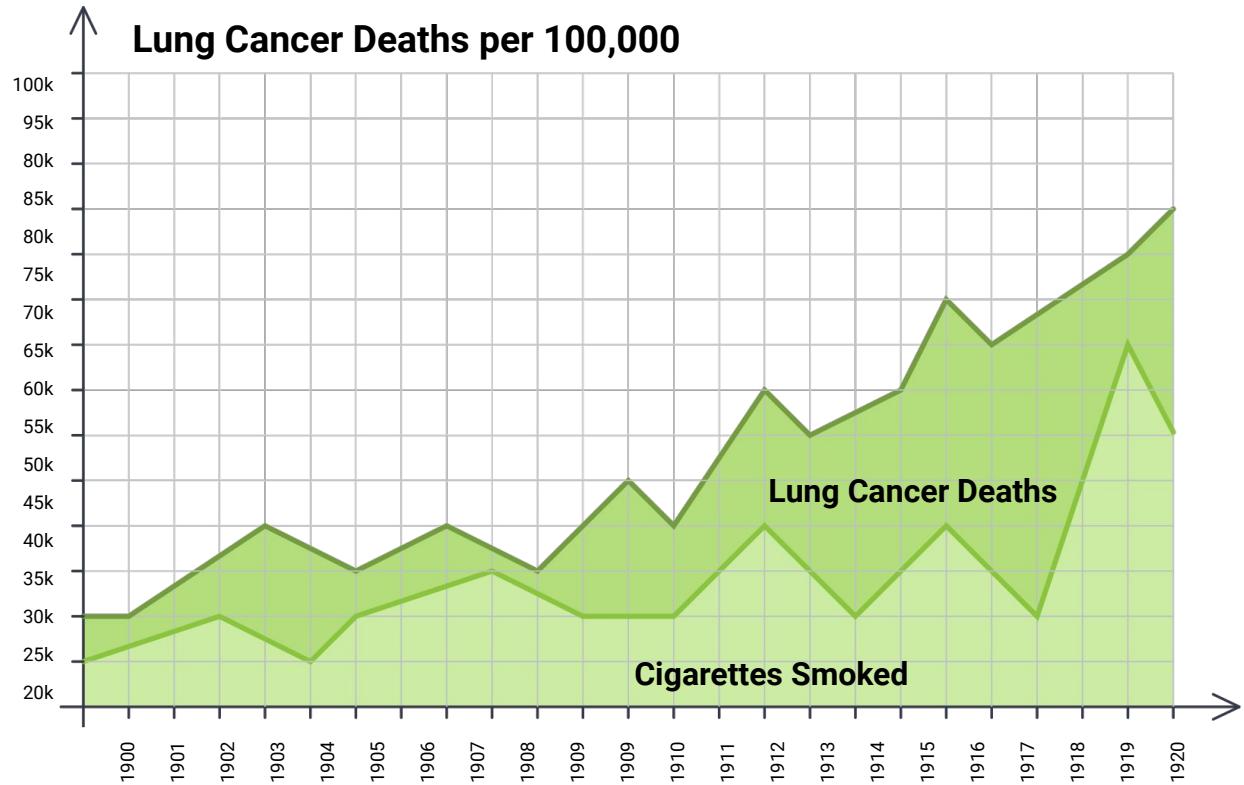
Data = Drama: Australia's Annual Mean Temperature (°C)



Data as Truth-Telling

Data as Truth-Telling

Unearthing Relationships

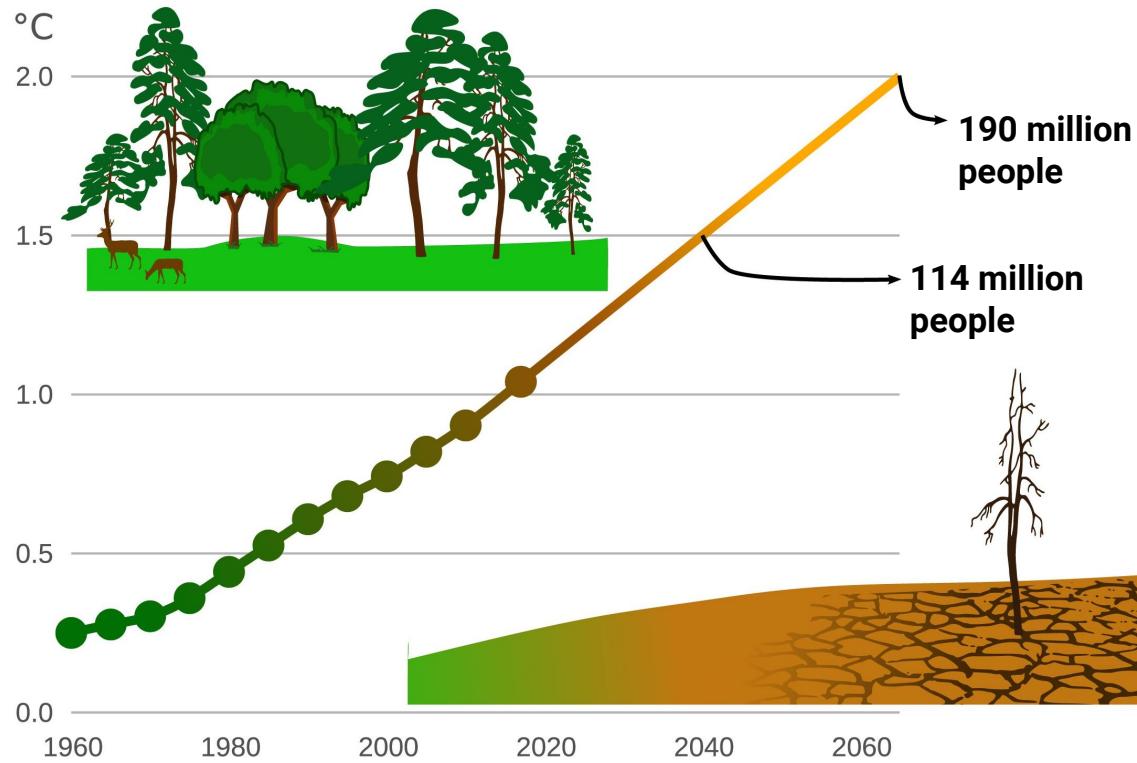




Data as
Truth-Telling

Making
Predictions

Exposure to Extreme Drought Is Increasing



Data as
Truth-Telling
—
Stating
Significance

Course Overview

Tools for Truths, Skills for Stories

Our Goals:



Truth-telling
Storytelling

Our Means:



Microsoft Excel

SQL

Python

MongoDB

pandas

HTML/CSS

Matplotlib/Seaborn

JavaScript

APIs

D3.js

Beautiful Soup

Leaflet.js/Google

Machine Learning

Maps

Tableau

Hadoop

Course Overview

Each class will include the following:



Overview of lesson topics



Instructor lecture



Instructor demonstration



Class discussions



In-class activities



Project work

Weekly Breakdown by Subject

Weeks 1–2

Week 1

Week 2

Week 3

Week 4

Week 5

Week 6

Week 7

Week 8

Week 9

Week 10

Week 11

Week 12

Week 13

Week 14

Week 15

Week 16

Week 17

Week 18

Week 19

Week 20

Week 21

Week 22

Week 23

Week 24

Intro to Data Analytics & Excel: Introduction to the high-level concepts of data analytics and real-world data crunching with Excel formulas, pivot tables, and conditional formatting.

Python Data Analytics and Visualisation: Thorough crash course on Python programming, followed by multiple weeks of data processing using advanced libraries like NumPy, pandas, Matplotlib, Seaborn, and BeautifulSoup.

Deep Dive into Databases: Immersion into introductory and advanced work with SQL (PostgreSQL) and noSQL databases (MongoDB).

Web-Based Data Visualisation: Introduction to the fundamental tools of web development (HTML, CSS, JavaScript) and advanced libraries that are useful for data visualisation (Plotly.js, Leaflet.js), and R.

Final Projects & Advanced Topics: Introduction to Tableau and advanced topics like Hadoop and machine learning. Develop a real-world data visualisation project.

Weeks 3–8

Week 1

Week 2

Week 3

Week 4

Week 5

Week 6

Week 7

Week 8

Week 9

Week 10

Week 11

Week 12

Week 13

Week 14

Week 15

Week 16

Week 17

Week 18

Week 19

Week 20

Week 21

Week 22

Week 23

Week 24

Intro to Data Analytics & Excel: Introduction to the high-level concepts of data analytics and real-world data crunching with Excel formulas, pivot tables, and conditional formatting.

Python Data Analytics and Visualisation: Thorough crash course on Python programming, followed by multiple weeks of data processing using advanced libraries like NumPy, pandas, Matplotlib, Seaborn, and BeautifulSoup.

Deep Dive into Databases: Immersion into introductory and advanced work with SQL (PostgreSQL) and noSQL databases (MongoDB).

Web-Based Data Visualisation: Introduction to the fundamental tools of web development (HTML, CSS, JavaScript) and advanced libraries that are useful for data visualisation (Plotly.js, Leaflet.js), and R.

Final Projects & Advanced Topics: Introduction to Tableau and advanced topics like Hadoop and machine learning. Develop a real-world data visualisation project.

Weeks 9–13

Week 1

Intro to Data Analytics & Excel: Introduction to the high-level concepts of data analytics and real-world data crunching with Excel formulas, pivot tables, and conditional formatting.

Week 2

Week 3

Week 4

Week 5

Week 6

Week 7

Week 8

Week 9

Week 10

Week 11

Week 12

Week 13

Week 14

Week 15

Week 16

Week 17

Week 18

Week 19

Week 20

Week 21

Week 22

Week 23

Week 24

Python Data Analytics and Visualisation: Thorough crash course on Python programming, followed by multiple weeks of data processing using advanced libraries like NumPy, pandas, Matplotlib, Seaborn, and BeautifulSoup.

Deep Dive into Databases: Immersion into introductory and advanced work with SQL (PostgreSQL) and noSQL databases (MongoDB).

Web-Based Data Visualisation: Introduction to the fundamental tools of web development (HTML, CSS, JavaScript) and advanced libraries that are useful for data visualisation (Plotly.js, Leaflet.js), and R.

Final Projects & Advanced Topics: Introduction to Tableau and advanced topics like Hadoop and machine learning. Develop a real-world data visualisation project.

Weeks 14–17

Week 1

Intro to Data Analytics & Excel: Introduction to the high-level concepts of data analytics and real-world data crunching with Excel formulas, pivot tables, and conditional formatting.

Week 2

Week 3

Week 4

Week 5

Week 6

Week 7

Week 8

Week 9

Week 10

Week 11

Week 12

Week 13

Week 14

Python Data Analytics and Visualisation: Thorough crash course on Python programming, followed by multiple weeks of data processing using advanced libraries like NumPy, pandas, Matplotlib, Seaborn, and BeautifulSoup.

Week 15

Week 16

Week 17

Week 18

Week 19

Week 20

Week 21

Week 22

Week 23

Week 24

Deep Dive into Databases: Immersion into introductory and advanced work with SQL (PostgreSQL) and noSQL databases (MongoDB).

Web-Based Data Visualisation: Introduction to the fundamental tools of web development (HTML, CSS, JavaScript) and advanced libraries that are useful for data visualisation (Plotly.js, Leaflet.js), and R.

Final Projects & Advanced Topics: Introduction to Tableau and advanced topics like Hadoop and machine learning. Develop a real-world data visualisation project.

Weeks 18–24

Week 1

Intro to Data Analytics & Excel: Introduction to the high-level concepts of data analytics and real-world data crunching with Excel formulas, pivot tables, and conditional formatting.

Week 2

Week 3

Week 4

Week 5

Week 6

Week 7

Week 8

Week 9

Week 10

Week 11

Week 12

Week 13

Week 14

Week 15

Week 16

Week 17

Week 18

Week 19

Week 20

Week 21

Week 22

Week 23

Week 24

Python Data Analytics and Visualisation: Thorough crash course on Python programming, followed by multiple weeks of data processing using advanced libraries like NumPy, pandas, Matplotlib, Seaborn, and BeautifulSoup.

Deep Dive into Databases: Immersion into introductory and advanced work with SQL (PostgreSQL) and noSQL databases (MongoDB).

Web-Based Data Visualisation: Introduction to the fundamental tools of web development (HTML, CSS, JavaScript) and advanced libraries that are useful for data visualisation (Plotly.js, Leaflet.js), and R.

Final Projects & Advanced Topics: Introduction to Tableau and advanced topics like Hadoop and machine learning. Develop a real-world data visualisation project.

Example Activity



Example Activity: Create a Heat Map

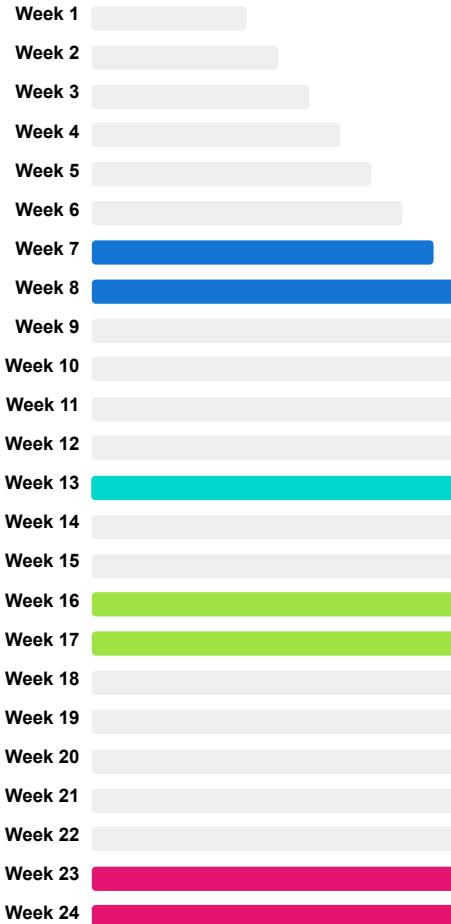
In this activity, you will use public data to map the density of water hydrants in Western Australia. You can compare this data with data from other sources like the Australian Census, the Australian Bureau of Statistics, Google Maps, and more to find insights on access to public utilities and services.

Suggested Time:

20 minutes

Projects

Projects



Project 1: Exploratory Data Analysis

Project Specialisation Tracks

Tackle each project with a specialised track!

Choose your focus:

- Finance
- Healthcare
- Custom

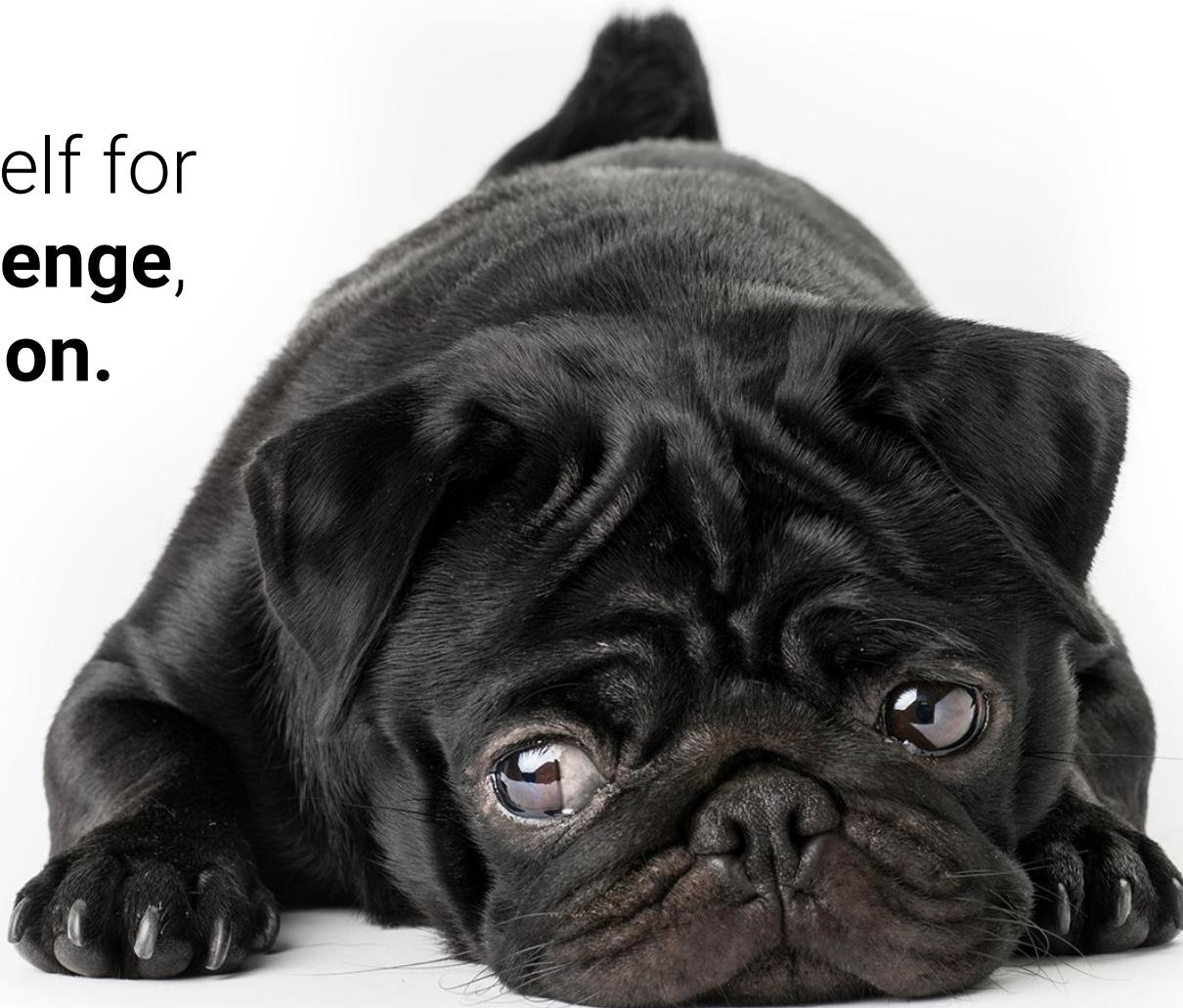
Project 2: Extract, Transform & Load

Project 3: Data Visualisations

Final Project: Machine Learning

Helpful Tips

Brace yourself for
doubt, challenge,
and **confusion.**



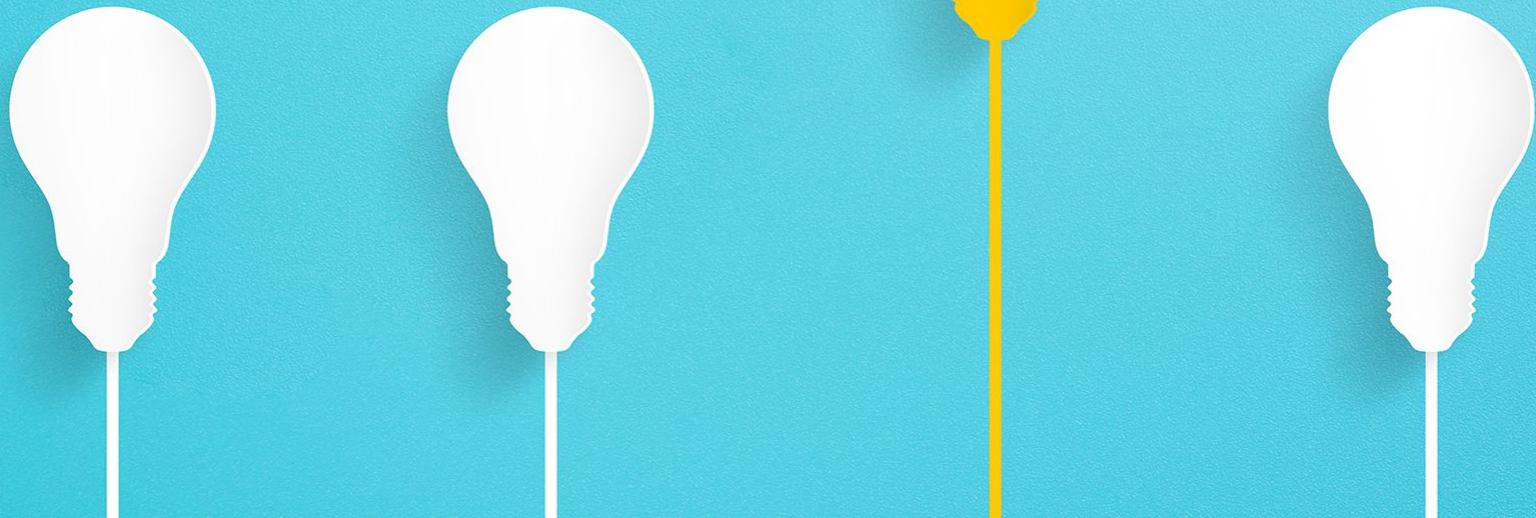
There is no shortcut.
You've got to **put in the hours!**



Form a community
with your classmates.



Enjoy the **novice experience**
and expect lots of
light bulb moments.



A close-up photograph of a fluffy orange and white cat lying on its back on a light-colored, textured surface. The cat's eyes are closed, and it has a content expression. Its front paws are raised towards its head, and its pink paw pads are visible. The background is slightly blurred.

Celebrate your successes!



Group Activity:

Form groups of 3 or 4 people. Introduce yourself to your group.

Suggested Time:

10 minutes

Break





Group Activity: The Great Debate

Find the group you formed before the break.
Together, ponder the following question.

Suggested Time:

15 minutes

Group Activity: The Great Debate

Which do Australians prefer:
Italian or Thai food?



Group Activity: The Great Debate

With your group, develop a strategy for answering this question with as much confidence possible. Specifically, answer questions like:



What data will you attempt to gather?



What relationships will you be searching for?



How will you ensure that your answer is most likely true?

Assumptions:



You are given 5 hours and a budget of \$10 to accomplish this.



Your answer will be tested by randomly selecting 9 Australians who will each be asked the question—with zero qualifiers.



You have only your group.

The Great Debate (Analysed)

Step 1: Decompose the “Ask”

Step 1: Decompose the 'Ask'

Which do **Australians** prefer:
Italian or Thai food?



Step 1: Decompose the “Ask”

Which do **Australians** prefer: Italian or Thai food?



Who exactly is an **Australian**?



Are **Australians** just homeowners?



Do **Australians** just live in big cities?



Are **Australians** just millennials?



How can we get a representative sample
of Australians?

Step 1: Decompose the 'Ask'

Which do Australians **prefer**:
Italian or Thai food?



Step 1: Decompose the “Ask”

Which do Australians **prefer**: Italian or Thai food?



How do we define “preference”?



Do people prefer the foods they eat the most frequently?



Do people prefer the foods they wish they could eat if cost was not an issue?



How uniform is the preference? Is it regionalised? Is it different by demographic?



Inherently, preference is subjective. We are going to need to make it objective.

Step 1: Decompose the 'Ask'

Which do Australians prefer:
Italian or Thai food?





Italian and Thai are broad categories to pursue. We will have to narrow the scope.

Step 1: Decompose the “Ask”

Which do Australians prefer: **Italian or Thai food?**

01

How do we categorise foods? Is pizza Italian? Is Bing Boy Thai?

02

How do we categorise food? Does making pasta at home constitute Italian? Or, are we just talking about restaurants?

03

Are we just talking about best experiences? Or, are we including poorer versions of these foods?

Step 2: Identify Data Sources

Step 2: Identify Data Sources

Why poll an audience when there are already enormous databases of information about Australians' food preferences—readily available online?



Step 2: Identify Data Sources

As everyday consumers, we are **regularly** getting a pulse of everyday Australian food preferences to inform our own decisions. Perhaps we can make use of the same approach.



Step 2a: Identify Data Sources

Web services like Yelp provide an almost-encyclopedic amount of information about the eating preferences of Australians.



Thai Food

Sydney, New South W.

For Businesses

Write a Review

Log In

Sign Up

Restaurants ▾

Home Services ▾

Auto Services ▾

More ▾

 157

\$\$ • Thai

"Solid 4 stars, good but not great **Thai** food. Maybe I was expecting more out of Yelps top choice for **Thai** in Sydney. Of the 7 dishes we ordered, everything was tasty but not overwhelming. Food comes out really quick but otherwise" [more](#)

 56

2. Chat Thai (02) 9221 0600
188 Pitt Street Sydney

\$\$ • Thai

"Oh wow! This is some phenomenal **Thai** food! I was in **Thai** town trying multiple spots and this was my favorite. They have excellent mains, I really enjoyed the pad see ew! But the real star is their

 Map data ©2020 Google Terms of Use Report a map error

Step 3: Define Strategy and Metrics

Step 3: Define Strategy and Metrics

Here, we created a blueprint for what we're targeting:

Australians	<ul style="list-style-type: none">Ideally, we need thousands of records from Australians in hundreds of localities (large samples).
Preference	<ul style="list-style-type: none">The number of Yelp reviews (more = preference).The average aggregated ratings (higher = preference).
Italian and Thai Food	<ul style="list-style-type: none">The top 20 Italian restaurants and the top 20 Thai restaurants in every city.



Step 3: Define Strategy and Metrics



Food Type

Yelp Best Thai Restaurant Sydney New South Wales, A Log In Sign Up

1. Home Thai Restaurant
★★★☆☆ 157 Review Count

\$ \$ • Thai

"Solid 4 stars, good but not great **Thai** food. Maybe I was expecting more out of Yelps top choice for **Thai** in Sydney. Of the 7 dishes we ordered, everything was tasty but not overwhelming. Food comes out really quick but otherwise" [more](#)

2. Chat Thai
★★★☆☆ Rating

\$ \$ • Thai

"Really fun trendy **Thai** restaurant in China town area of Sydney. Enjoyed the curries and noodles. Waited about 25 mins on Saturday at 8 pm." [more](#)

Lots of Data!

Step 3: Define Strategy and Metrics

Repeat this analysis for as many cities as possible.

Melbourne, VIC	
Italian	Thai
Restaurant	Restaurant
Restaurant	Restaurant
Restaurant	VS.
Restaurant	Restaurant
Restaurant	Restaurant

Sydney, NSW	
Italian	Thai
Restaurant	Restaurant
Restaurant	Restaurant
Restaurant	VS.
Restaurant	Restaurant
Restaurant	Restaurant

Brisbane, QLD	
Italian	Thai
Restaurant	Restaurant
Restaurant	Restaurant
Restaurant	VS.
Restaurant	Restaurant
Restaurant	Restaurant

Adelaide, SA	
Italian	Thai
Restaurant	Restaurant
Restaurant	Restaurant
Restaurant	VS.
Restaurant	Restaurant
Restaurant	Restaurant

Darwin, NT	
Italian	Thai
Restaurant	Restaurant
Restaurant	Restaurant
Restaurant	VS.
Restaurant	Restaurant
Restaurant	Restaurant

Perth, WA	
Italian	Thai
Restaurant	Restaurant
Restaurant	Restaurant
Restaurant	VS.
Restaurant	Restaurant
Restaurant	Restaurant

Step 4: Build a Data Retrieval Plan

Step 4: Build a Data Retrieval Plan

We could retrieve this data by brute force, but it would be:

- Extremely time consuming.
- Skewed by our city familiarity.
- Labour intensive.

The image displays three separate instances of a Yelp search interface, each with a red header and a white search bar. Each instance shows the Yelp logo on the left, followed by a search bar containing "Find Thai" and another search bar containing "Near [city], [state]". The first instance is for Melbourne, VIC; the second for Sydney, NSW; and the third for Brisbane, QLD. Each interface includes a red search button with a magnifying glass icon on the right.

Thank You, Yelp!

Thankfully, we can take advantage of the **Yelp Fusion API** to programmatically run our queries. (#ThankGoodnessForProgramming)

The screenshot shows the Yelp Fusion API documentation page. The left sidebar has sections for General (Create App, Email / Notifications, Display Requirements, Terms of Use, FAQ), Yelp Fusion (Introduction, Business Endpoints, Business Search, Phone Search), and a search bar. The main content area is titled '/businesses/search' and describes the endpoint for returning up to 1000 businesses based on search criteria. It includes a note about reviews and a 'Request' section with the GET URL `https://api.yelp.com/v3/businesses/search`. The 'Parameters' section lists 'term' (string) and 'location' (string) with their descriptions.

/businesses/search

This endpoint returns up to 1000 businesses based on the provided search criteria. It has some basic information about the business. To get detailed information and reviews, please use the Business ID returned here and refer to [/businesses/{id}](#) and [/businesses/{id}/reviews](#) endpoints.

Note: at this time, the API does not return businesses without any reviews.

Request

```
GET https://api.yelp.com/v3/businesses/search
```

Parameters

These parameters should be in the query string.

Name	Type	Description
term	string	Optional. Search term, for example "food" or "restaurants". The term may also be business names, such as "Starbucks". If term is not included the endpoint will default to searching across businesses from a small number of popular categories.
location	string	Required if either latitude or longitude is not provided. This string indicates the geographic area to be used when searching for businesses. Examples: "New York City", "NYC", "350 5th Ave, New York, NY 10118". Businesses returned in the response may not be strictly within the specified location.

Thank You, Yelp!

```
{"id": "WNpF7jhHH9U12t4dhuDlEA",
"alias": "pure-thai-berala-berala",
"name": "Pure Thai Berala",
'image_url': 'https://s3-media1.fl.yelpcdn.com/bphoto/mlkp03PLFa4gVIlWj82cXg/o.jpg',
'is_closed': False,
'url': 'https://www.yelp.com.au/biz/pure-thai-berala-berala?adjust_creative=1GwZyE0zIjSujpm-api_v3_business_search&utm_source=1GwZyE0zIjSujpHtlMnodQ',
'review_count': 3,
'categories': [{alias: 'thai', title: 'Thai'}],
'rating': 3.5,
'coordinates': {'latitude': -33.871492, 'longitude': 151.031695},
'transactions': [],
'location': {'address1': '160 Woodburn Rd',
'address2': '',
'address3': '',
'city': 'Berala',
'zip_code': '2141',
'country': 'AU',
'state': 'NSW',
'display_address': ['160 Woodburn Rd', 'Berala New South Wales 2141']},
'phone': '+61296492882',
'display_phone': '(02) 9649 2882',
'distance': 9435.618278117496}
```



Step 4: Build a Data Retrieval Plan

We will build a Python script to randomly select over 700 post codes from the Australian Bureau of Statistics and then acquire review data from the top 20 Thai and top 20 Italian restaurants for each post code by using the Yelp API.



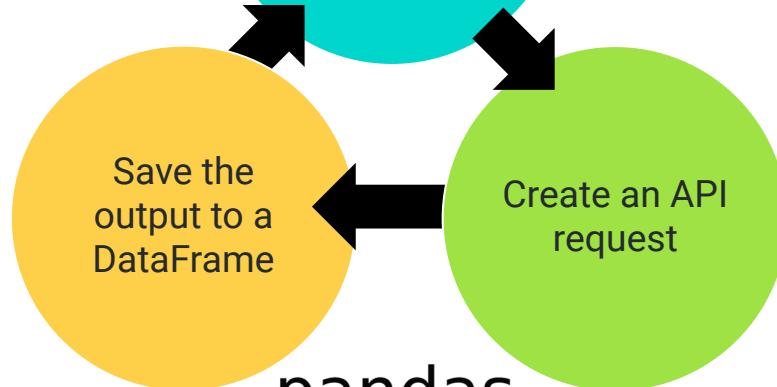
3001		2009		4000	
Italian	Thai	Italian	Thai	Italian	Thai
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant

5012		0810		6001	
Italian	Thai	Italian	Thai	Italian	Thai
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant
Restaurant	Restaurant	Restaurant	Restaurant	Restaurant	Restaurant



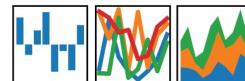
Step 5: Retrieve the Data

Pulling with Python



pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Pulling with Python

```
try:

    # Loop through all records to calculate the review count and weighted review value
    for business in yelp_reviews_italian["businesses"]:

        italian_review_count = italian_review_count + business["review_count"]
        italian_weighted_review = italian_weighted_review + (business["review_count"] * business["rating"])

    for business in yelp_reviews_thai["businesses"]:
        thai_review_count = thai_review_count + business["review_count"]
        thai_weighted_review = thai_weighted_review + (business["review_count"] * business["rating"])

    # Append the data to the appropriate column of the data frames
    italian_data = italian_data.append({
        'Postal Code': row["Postal Code"],
        'Italian Review Count':italian_review_count,
        'Italian Average Rating':(italian_weighted_review / italian_review_count),
        'Italian Weighted Rating':italian_weighted_review},ignore_index=True)

    thai_data = thai_data.append({
        'Postal Code': row["Postal Code"],
        'Thai Review Count':thai_review_count,
        'Thai Average Rating':(thai_weighted_review / thai_review_count),
        'Thai Weighted Rating':thai_weighted_review},ignore_index=True)

except:
    print("Uh oh")
```



This code...

Pulling with Python

1

<https://api.yelp.com/v3/businesses/search?term=Italian&location=2055>
<https://api.yelp.com/v3/businesses/search?term=Thai&location=2055>

2

<https://api.yelp.com/v3/businesses/search?term=Italian&location=2059>
<https://api.yelp.com/v3/businesses/search?term=Thai&location=2059>

3

<https://api.yelp.com/v3/businesses/search?term=Italian&location=2060>
<https://api.yelp.com/v3/businesses/search?term=Thai&location=2060>

4

<https://api.yelp.com/v3/businesses/search?term=Italian&location=2000>
<https://api.yelp.com/v3/businesses/search?term=Thai&location=2000>

5

<https://api.yelp.com/v3/businesses/search?term=Italian&location=2001>
<https://api.yelp.com/v3/businesses/search?term=Thai&location=2001>

6

<https://api.yelp.com/v3/businesses/search?term=Italian&location=2020>
<https://api.yelp.com/v3/businesses/search?term=Thai&location=2020>

7

<https://api.yelp.com/v3/businesses/search?term=Italian&location=2129>
<https://api.yelp.com/v3/businesses/search?term=Thai&location=2129>



**...will make all of
these URLs.**

Pulling with Python

GET https://api.yelp.com/v3/businesses/search?term=Italian&location=37764...

Headers (1)

Key	Value	Description	...	Bulk Edit	Presets
Authorization	Bearer gl6k6JmewUhjMVBv0I2x4Bz_NRIEggSjIjGbTaejmzbvBJXg 36F...				
New key	Value	Description			

Body

Pretty Raw Preview JSON

```
1 {  
2   "businesses": [  
3     {  
4       "id": "two-brothers-italian-pizza-kodak",  
5       "name": "Two Brothers Italian Pizza",  
6       "image_url": "https://s3-media3.fl.yelpcdn.com/bphoto/364BqQt0qtVHV1f0t_xznA/o.jpg",  
7       "is_closed": false,  
8       "url": "https://www.yelp.com/biz/two-brothers-italian-pizza-kodak?adjust_creative=1GwZyE0zIjSujpHtlMnodQ&utm_campaign=yelp_api_v3&utm_medium=  
9         _api_v3_business_search&utm_source=1GwZyE0zIjSujpHtlMnodQ",  
10      "review_count": 8,  
11      "categories": [  
12        {  
13          "alias": "pizza",  
14          "title": "Pizza"  
15        },  
16        {  
17          "alias": "italian",  
18          "title": "Italian"  
19        },  
20        {  
21          "alias": "pastashops",  
22          "title": "Pasta Shops"  
23        },  
24      ],  
25      "rating": 2,  
26      "coordinates":  
27        {  
28          "latitude": 35.9638662447754,  
29          "longitude": -83.5926620147413  
30        },  
31      "transactions": [],  
32      "location": {  
33        "address1": "1000 W Broad St",  
34        "address2": null,  
35        "city": "Columbus",  
36        "state": "OH",  
37        "zip_code": "43228",  
38        "country": "US",  
39        "display_address": ["1000 W Broad St", "Columbus, OH 43228"]  
40      }  
41    }  
42  ]  
43}  
44}
```



Each of these URLs holds a piece of our answer.

Step 6: Assemble and Clean the Data

Cleaning with Pandas

No data comes out exactly the way you want it to.
In our case, we needed multiple steps to aggregate the data along our channels of interest.

```
# Combine DataFrames into a single DataFrame  
combined_data = pd.merge(thai_data, italian_data, on="Postal Code")  
combined_data.head()
```

	Post Code	Thai Review Count	Thai Average Rating	Thai Weighted Rating	Italian Review Count	Italian Average Rating	Italian Weighted Rating
0	0801	97	4.1134	399	63	3.78571	238.5
1	4000	256	4.11133	1052.2	266	3.81955	1016
2	5012	378	3.64286	1377	66	3.2197	212.5
3	3001	222	4.16892	925.5	420	3.77857	1587
4	2009	2842	3.94053	11199	2829	3.92824	11113

Step 7: Analyse for Trends

Analyse for Trends (Table)

It's close:

```
# Model 1: Head-to-Head Review Counts
italian_summary = pd.DataFrame({"Review Counts": italian_data["Italian Review Count"].sum(),
                                 "Rating Average": italian_data["Italian Average Rating"].mean(),
                                 "Review Count Wins": combined_data["Review Count Wins"].value_counts()["Italian"],
                                 "Rating Wins": combined_data["Rating Wins"].value_counts()["Italian"]}, index=["Italian"])

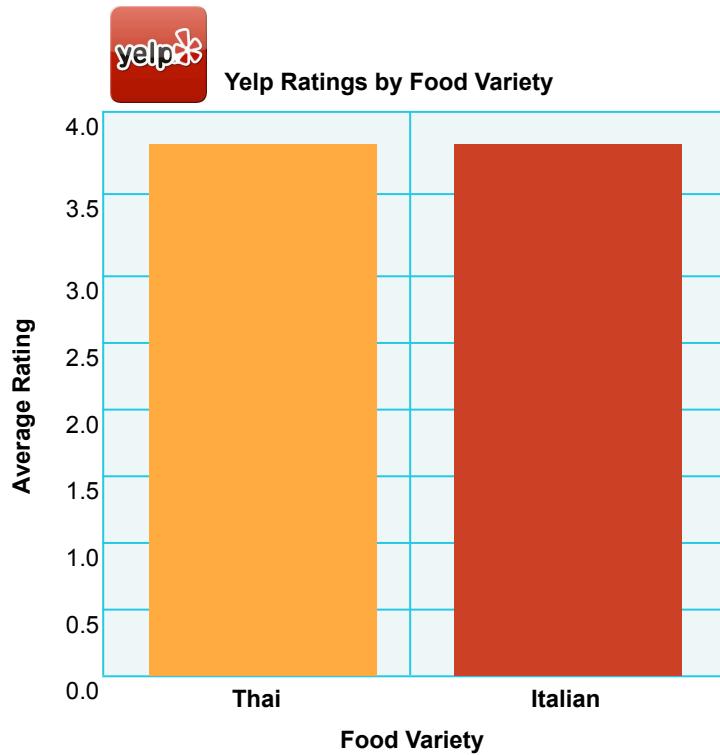
thai_summary = pd.DataFrame({"Review Counts": thai_data["Thai Review Count"].sum(),
                             "Rating Average": thai_data["Thai Average Rating"].mean(),
                             "Review Count Wins": combined_data["Review Count Wins"].value_counts()["Thai"],
                             "Rating Wins": combined_data["Rating Wins"].value_counts()["Thai"]}, index=["Thai"])

final_summary = pd.concat([thai_summary, italian_summary])
final_summary
```

	Review Counts	Rating Average	Review Count Wins	Rating Wins
Thai	2229761.0	3.987034	347	461
Italian	2670762.0	3.964877	600	486

Analyse for Trends (Ratings)

Yelpers rate Italian and Thai relatively **equally**.

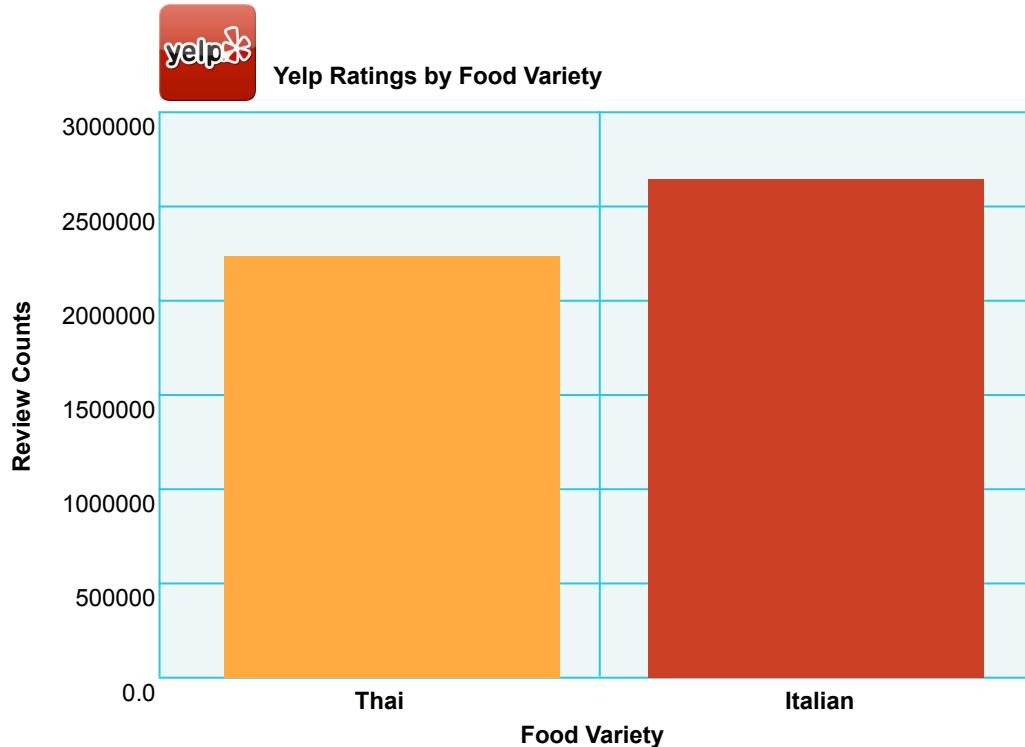


=



Analyse for Trends (Ratings)

Yelpers seem to **review significantly more Italian** restaurants.



Analyse for Trends (Statistical Analysis)

We have an intuitive sense that the numbers are close, but to quantify our intuition, we use a Student's t-test. After performing the t-test, we can quantifiably state that the differences are not statistically significant

Metric	Italian	Thai	p-value (t-test)
Average Rating	3.965	3.987	0.4696
Review Counts	2,671k	2,229k	0.0000235

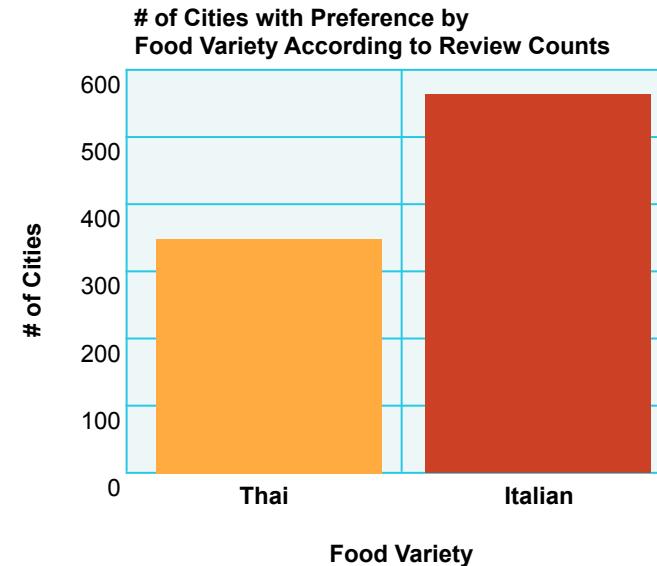
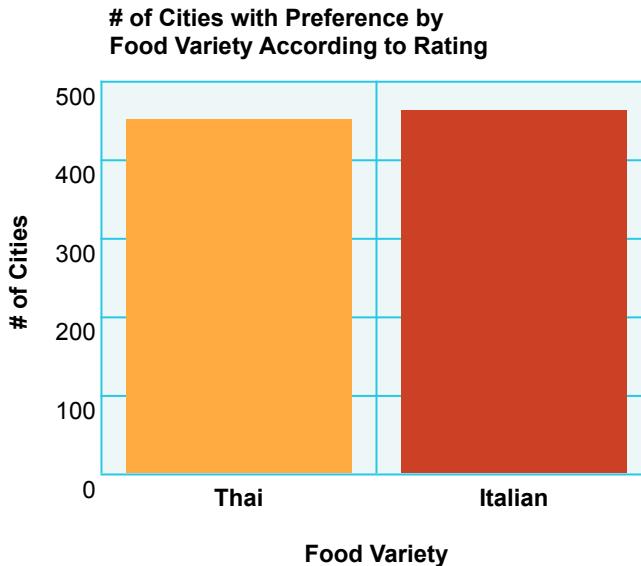


The difference in review counts is **not statistically significant**.

Analyse for Trends (Winner Take All)

Just for fun, let's throw in an analysis that aggregates the data from all the cities by using a winner-take-all approach.

It's sort of even when compared by rating, but **Italian is the clear winner** when comparing the number of reviews.



Step 8: Acknowledge Limitations

Limitations of Analysis

Yelp demographics might not match the Australian demographic.



Limitations of Analysis

Restaurant experiences do not equate to home-cooked meals.



Limitations of Analysis

The fine dining effect?



Step 9: Make the Call

Making the Call

The “proper” conclusion:

Based on our analysis, it's clear that Australians' preferences for Italian and Thai food are similar in nature. As a whole, Australians rate Thai and Italian restaurants at statistically similar scores (average score: 3.8, p-value: 0.4696). However, there is substantial evidence that Australians write more reviews of Italian restaurants than Thai restaurants (+441k, p-value: 0.0000235).



This might indicate there is an increased interest in visiting Italian restaurants for the experience. Or, it might merely suggest that Yelp users enjoy writing reviews of Italian restaurants more than Thai restaurants.

Making the Call

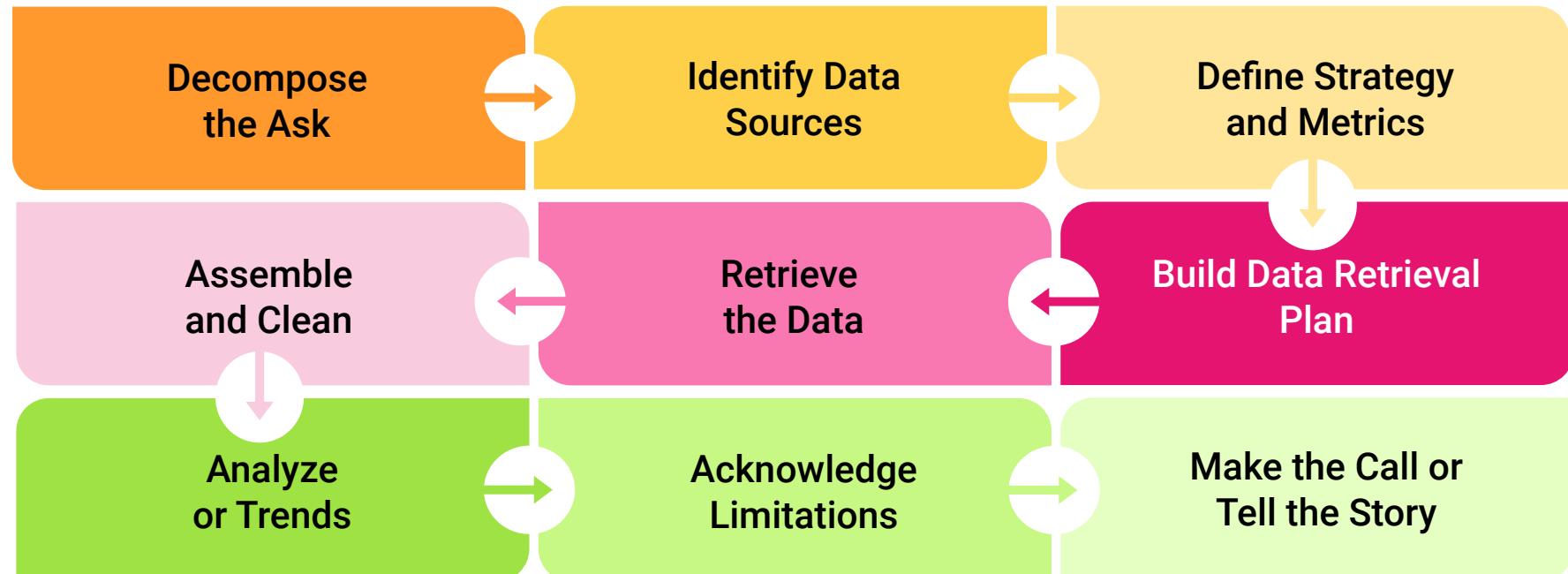
The “let’s be real” conclusion: Italian (but, it’s going to be close).



An Analytics Paradigm

Analytics Paradigm

Regardless of type or industry, this paradigm provides a repeatable pathway for effective data problem-solving.





Group Activity: Predicting Gentrification

Using the analytics paradigm as a framework, outline a strategy to identify which neighbourhoods in our city are showing signs of gentrification.

Suggested Time:

15 minutes

Group Activity: Predicting Gentrification

Specifically, how would you answer these questions:

-  What observable signs can we detect to suggest gentrification is happening?
-  What means can we use to determine how long the trend has been happening?
-  What proxies might we use to identify gentrification in non-obvious ways?
-  How might you create a visualisation of this data to best tell the story?

Pay special attention to details, like:

-  What data will you use to build your model?
-  How will you retrieve the data?
-  What does your final story look like?



Time's Up! Let's Review.

Prepare for Next Class

By next class:

01

Make certain that you have Excel installed.

02

Make certain that you have Slack installed and that you're actively checking it.

03

Figure out where the Git repository for our class is.

04

Figure out where the class videos will be posted.

Questions?

