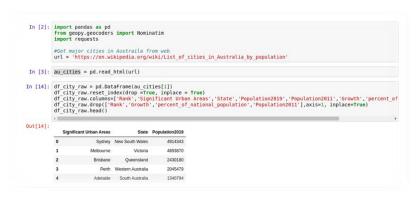
## **Data for Australian city grouping**

To help build a model to compare Australian cities with each other, we propose to follow the process below and gather the data required for the project.

1. Finalize the list of cities to be considered for analysis. This is needed because it is not practically possible to cover every inch of Australia and analyze it for similarity. Plus, there isn't data available for every small city/town or settlement. For the purpose of this project, we intend to get the list of cities available on Wikipedia <a href="https://en.wikipedia.org/wiki/List of cities in Australia by population">https://en.wikipedia.org/wiki/List of cities in Australia by population</a> This data also gives us the recent population of each city.

## Here is quick look at the data



2. Get the latitude and longitude for each city chosen for analysis in step 1 using **geopy** API Take a quick look at latitude and longitude for each City



3. Use Foursquare API to get the venues for each of the cities chosen in step1 using its latitude and longitude found in step 2

## (Point 3 continued...)

## A quick look at it

```
In [21]: radius = 500
limit=100
version = 20203112

for city_lat, city_lng in zip(df[df['City']="Kalgoorlie=Boulder']['Latitude'], df[df['City']="Kalgoorlie=Boulder']
url = "https://api.foursquare.com/v2/venues/search?client_id={}&client_secret={}&radius={}&lliet_{}, {}&v={}
url = "https://api.foursquare.com/v2/venues/search?client_id={}&client_secret={}&radius={}&lliet_{}, {}&v={}
url = "https://api.foursquare.com/v2/venues/search?client_id={}&client_secret={}&radius={}&lliet_{}, {}&v={}
url = "https://soilet_id={}&client_secret={}&radius={}&lliet_{}, {}&v={}
url = "https://soilet_id={}&client_secret={}&radius={}&lliet_{}, {}&v={}
url = "https://soilet_id={}&client_secret={}&radius={}&lliet_{}, {}&v={}
url = "https://soilet_id={}&client_secret={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={}&radius={
```

4. Foursquare API venues give us very good data about what kind of places a city has. It is called as venues.

Here is sneak peak in to one venue detail

```
{'categories': [{'icon': {'prefix': 'https://ss3.4sqi.net/img/categories_v2/nightlife/pub_',
    'suffix': '.png'},
   'id': '4bf58dd8d48988d11b941735',
   'name': 'Pub',
   'pluralName': 'Pubs',
   'primary': True,
   'shortName': 'Pub'}],
  'hasPerk': False,
  'id': '4da41978540ea1cd84db9dde',
  'location': {'cc': 'AU',
   'country': 'Australia',
   'distance': 78,
   'formattedAddress': ['Australia'],
   'labeledLatLngs': [{'label': 'display',
    'lat': -26.54196138365552,
    'lng': 151.8393920350731}],
```

(Point 4 Continued...)

'lat': -26.54196138365552,

'Ing': 151.8393920350731},

'name': 'Kingaroy RSL',

'referralId': 'v-1611223824'}

The venue categories and population are used as a measure of similarity and dissimilarity. The rationale behind this is a traveler to the city is most interested in the places to visit. These places are given by Foursquare API and are called **venues**. The population of the city plays a major role in deciding many important aspects of a city, including public transport, amenities, crime, poverty, etc. Thus, city population is one important feature to decide similarity and dissimilarity among cities.