

Chen, Meixu; Fahrner, Dominik; Arribas-Bel, Daniel; Rowe, Francisco

## Article

# A reproducible notebook to acquire, process and analyse satellite imagery: Exploring long-term urban changes

REGION

## Provided in Cooperation with:

European Regional Science Association (ERSA)

*Suggested Citation:* Chen, Meixu; Fahrner, Dominik; Arribas-Bel, Daniel; Rowe, Francisco (2020) : A reproducible notebook to acquire, process and analyse satellite imagery: Exploring long-term urban changes, REGION, ISSN 2409-5370, European Regional Science Association (ERSA), Louvain-la-Neuve, Vol. 7, Iss. 2, pp. R15-R46, <https://doi.org/10.18335/region.v7i2.295>

This Version is available at:

<https://hdl.handle.net/10419/235822>

### Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

### Terms of use:

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*



<https://creativecommons.org/licenses/by-nc/4.0>

## A reproducible notebook to acquire, process and analyse satellite imagery: Exploring long-term urban changes

Meixu Chen, Dominik Fahrner, Daniel Arribas-Bel, Francisco Rowe<sup>1</sup>

<sup>1</sup> University of Liverpool, Liverpool, UK

Received: 22 December 2019/Accepted: 1 December 2020

**Abstract.** Satellite imagery is often used to study and monitor changes in natural environments and the Earth surface. The open availability and extensive temporal coverage of Landsat imagery has enabled to monitor changes in temperature, wind, vegetation and ice melting speed for a period of up to 46 years. Yet, the use of satellite imagery to study cities has remained underutilised in Regional Science, partly due to the lack of a practical methodological approach to capture data, extract relevant features and monitor changes in the urban environment. This notebook offers a framework to demonstrate how to batch-download high-resolution satellite imagery; and enable the extraction, analysis and visualisation of features of the built environment to capture long-term urban changes.

**Key words:** satellite imagery, image segmentation, urbanisation, cities, urban change, computational notebooks

### 1 Introduction

Sustainable urban habitats are a key component of many global challenges. Efficient management and planning of cities are pivotal to all 17 UN Sustainable Development Goals (SDGs). Over 90% of the projected urban population growth by 2050 will occur in less developed countries (United Nations 2019). Concentrated in cities, this growth offers an opportunity for social progress and economic development but it also imposes major challenges for urban planning. Prior work on urbanisation has identified the benefits of agglomeration and improvements in health and education, which tend to outweigh the costs of congestion, pollution and poverty (Glaeser, Henderson 2017). Yet research has remained largely focused on Western cities (e.g. Burchfield et al. 2006), developing a good understanding of urban areas in high-income, developed countries (Glaeser, Henderson 2017). Much less is known about the long-term evolution of urban habitats in less developed countries. Analysis of historical census data exist exploring changes at discrete points over time such as slum detection (e.g. Giada et al. 2003, Kit, Lüdeke 2013, Kohli et al. 2016). Less applications can be identified tracking changes in urban settings over a continuous temporal scale (Ibrahim et al. 2020). This gap is partly due to the lack of comprehensive and consistent data sources capturing the long-term dynamics of urban structures in less developed countries.

Cities in Asia provide a unique setting to explore the challenges triggered by rapid urbanisation. The share of urban population in Asia is currently at a turning point transitioning to exceed the share of rural population. Currently Asia is home to over 53% of the urban population globally and the share of urban population is projected to increase to 66% by 2050 (United Nations 2019). Developing tools to monitor and understand the past and current urbanisation process is key to guide appropriate urban planning and policy strategies.

Recent technological developments can help overcome the paucity in spatially-detailed urban data in less developed countries. The combination of geospatial technology, cheap computing and new machine learning algorithms has ushered in an age of new forms of data, producing brand new data sets and repurposing existing sources. Satellite imagery represents a key source of information. Photographs from the sky have existed for decades, but their use in the context of socioeconomic urban research has been limited. Image data has been hard to process and understand for social scientists. Yet recent developments in machine learning and artificial intelligence have made images computable and turned these data into brand new information to be explored by quantitative urban researchers. Further, satellite data has become more abundant and openly accessible in the past decade, and offers new possibilities for data exploration through increasing spatial and temporal resolution. This, together with more computational power being available, allows to process these data in an efficient and meaningful way.

This notebook illustrates an easy-to-use analytical framework based on Python tools which enables batch download, image feature extraction, analysis and visualisation of high-resolution satellite imagery to capture long-term urban changes. Our purpose is to fill in the absence of a systematic and reproducible framework to acquire, process and analyse satellite imagery in urban built environment related to the field of Regional Science. The source of satellite data and administrative boundaries data are from NASA's Landsat satellite programme and ArcGIS Online. The Python libraries used in this notebook are the following:

- [Landsat images in Google Cloud Storage](#): The Google Cloud Storage is accessed using an API to download Landsat imagery (version used: 0.4.9)
- [Matplotlib](#): A Python 2D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms.
- [Numpy](#): Adding support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions
- [Pandas](#): Provides high-performance, easy-to-use data structures and data analysis tools
- [GeoPandas](#): Python library that simplifies working with geospatial data (version used: 0.6.2)
- [Folium](#): Python library that enables plotting interactive maps using leaflet (version used: 0.10.0)
- [Glob](#): Unix style pathname pattern expansion
- [GDAL](#): Library for geospatial data processing (version used: 2.4.4)
- [Landsat578](#): Simple Landsat imagery download tool
- [L8qa](#): Landsat processing toolbox (version used: 0.1.1)
- [Rasterio](#): Library for raster data processing (version used: 1.1.3)
- [Scikit-image](#): Collection of algorithms for image processing
- [Wget](#): Pure python download utility (version used: 3.2)
- [OpenCV](#): Library for image processing
- [scikit-learn](#): Machine learning in Python. Simple and efficient tools for data mining and data analysis.

We can import them all as follows:

```
[1]: %matplotlib inline

#load external libraries
import matplotlib.pyplot as plt
from matplotlib import colors
import pandas as pd
import numpy as np
import geopandas as gpd
import folium
import os, shutil
import glob
import gdal
import wget
from landsat import google_download
from google_download import GoogleDownload
from l8qa.qa import write_cloud_mask
import rasterio
import rasterio as rio
from rasterio import merge
from rasterio.plot import show
from rasterio.mask import mask
from skimage import io, exposure, transform, data
from skimage.color import rgb2hsv, rgb2gray
from skimage.feature import local_binary_pattern
from sklearn.cluster import KMeans
import matplotlib.cm as cm
from sklearn import preprocessing
from rasterio.enums import Resampling
import seaborn as sns
import itertools

wdir= os.getcwd()
```

The remainder of this paper is structured as follows. The next section introduces the Landsat satellite imagery, study area Shanghai, and process on how to batch download and pre-process satellite data. Section 3 proposes our methods to extract different features including colour, texture, vegetation and built-up from imagery. Section 4 performs a clustering method on the extracted features, and section 5 interprets the results and gain insights from them. Finally, section 6 concludes by providing a summary of our work and avenues for further research using our proposed framework.

## 2 Data and Study Area

### 2.1 Landsat Imagery

We draw data from the NASA's Landsat satellite programme. It is the longest standing programme for Earth observation (EO) imagery (NASA 2019). Landsat satellites have been orbiting the Earth for 46 years providing increasingly higher resolution imagery. Landsat Missions 1-3 offer coarse imagery of 80m covering the period from 1972 to 1983. Landsat Missions 4-5 provides images of 30m resolution covering the period from 1983 to 2013 and Landsat Missions 7-8 are currently collecting enhanced images at 15m capturing Cirrus and Panchromatic bands, in addition to the traditional RGB, Near-, Shortwave-Infrared, and Thermal bands. The Landsat 6 mission was unsuccessful due to the transporting rocket not reaching orbit. Landsat imagery is openly available and offers extensive temporal coverage stretching for 46 years. Table 1 provides a summary overview of the operation, revisit time and image resolution for the Landsat programme, with other Earth observation satellite missions being shown in Table 2.

Additional Earth observation programmes exist. These programmes also offer freely accessible imagery at a higher resolution.

### 2.2 Study Area

In this analysis, we examine urban changes in Shanghai, China. Shanghai has experienced rapid population growth. Between 2000 and 2010, Shanghai's population rose by 7.4

Table 1: Overview of Landsat missions, their revisit time and spatial resolution

Mission	Operational time	Revisit time	Resolution
Landsat 1	1972-1978	18 d	80 m
Landsat 2	1975-1982	18 d	80 m
Landsat 3	1978-1983	18 d	80 m
Landsat 4	1983-1993	16 d	30 m
Landsat 5	1984-2013	16 d	30 m
Landsat 7	1999-present	16 d	15 m
Landsat 8	2013-present	16 d	15 m

Table 2: Overview of other Earth observation satellites, their revisit time and spatial resolution

Provider	Programme	Operational time	Revisit time	Resolution
European Space Agency	Sentinel	2015-present	5 d	10m
Planet Labs	Rapideye PlanetscopeSkysat	2009-present	4/5 d to daily	up to 0.8 m
NASA	Orbview 3	2003-2007	<3 d	1-4 m
NASA	EO-1	2003 -2017	-	10-30 m

million from 16.4 million to 23.8 million. It has an annual growth rate of 3.8 percent over 10 years. While the pace of population expansion has been less acute, Shanghai's population has continued to grow. In 2018, an estimated 24.24 million people were living in Shanghai experiencing a population expansion of approximately 8 million since 2010. The city is therefore a well suited example to explore long-term changes in urbanisation.

To extract satellite imagery, a first step is to identify the shape of the geographical area of interest. To this end, we use a polygon shapefile (<https://www.arcgis.com/home/item.html?id=105f92bd1fe54d428bea35eade65691b>). These polygons represent the Shanghai metropolitan area, so they include the city centre and surrounding areas. These polygons will be used as a bounding box to identify and extract relevant satellite images. We need to ensure the shapefile is in the same coordinate reference system (CRS) as the satellite imagery (WGS84 or EPSG:4326).

```
[2]: # Specify the path to your shapefile
directory = os.path.dirname(wdir)
shp = 'shang_dis_merged/shang_dis_merged.shp'
```

```
[3]: # Certify that the shapefile is in the right coordinate system, otherwise reproject
# it into the right CRS
def shapefile_crs_check(file):
    global bbox
    bbox = gpd.read_file(file)
    crs = bbox.crs
    data = crs.get("init", "")
    if 'epsg:4326' in data:
        print('Shapefile in right CRS')
    else:
        bbox = bbox.to_crs({'init': 'epsg:4326'})
    f,ax = plt.subplots(figsize=(5,5))
    plt.title('Fig.1: Shapefile of Shanghai urban area',y=-0.2)
    bbox.plot(ax=ax)
```

```
[4]: shapefile_crs_check(shp)
```

```
[4]: Shapefile in right CRS
image/png<Figure size 360x360 with 1 Axes>
```

The world reference system (WRS) from NASA is a system to identify individual satellite imagery scenes using path-row tuples instead of absolute latitude/longitude



Figure 1: Shapefile of Shanghai urban area

coordinates. The latitudinal centre of the image corresponds to the row, the longitudinal centre to the path. This system allows to uniformly catalogue satellite data across multiple missions and provides an easy to use reference system for the end user. It is necessary to note that the WRS was changed between Landsat missions, due to a difference in swath patterns of the more recent Landsat satellites (NASA 2019). The WRS1 is used for Landsat missions 1-3 and the WRS2 for Landsat missions 4,5,7,8. In order to obtain path-row tuples of relevant satellite images for an area of interest (AOI), it is necessary to intersect the WRS shapefile (either WRS1 or WRS2, depending on the Landsat satellite you would like to obtain data from) with the AOI shapefile. The resulting path-row tuples will later be used to locate and download the corresponding satellite images from the Google Cloud Storage. The output of the intersection between WRS and AOI files can be visualised using an interactive widget. The map below shows our area of interest in purple and the footprints of the relevant Landsat images on top of an OpenStreetMap basemap.

```
[5]: # Download the WRS 2 file to later intersect the shapefile with the WRS path/row
# tuples to identify relevant Landsat scenes
#
def sat_path():

    url = 'https://prd-wret.s3.us-west-2.amazonaws.com/assets/palladium/production/
    ...s3fs-public/atoms/files/WRS2_descending_0.zip'
    # Create folder for WRS2 file
    if os.path.exists(os.path.join('Landsat_images', 'wrs2')):
        print('folder exists')
    else:
        os.makedirs(os.path.join('Landsat_images', 'wrs2'))

    WRS_PATH = os.path.join('Landsat_images', 'WRS2_descending_0.zip')
    LANDSAT_PATH = os.path.dirname(WRS_PATH)

    # The WRS file is only needed once thus we add this loop
    if os.path.exists(WRS_PATH):
        print('File already exists')
    # Downloads the WRS file from the URL given and unzips it
    else:
        wget.download(url, out = LANDSAT_PATH)
        shutil.unpack_archive(WRS_PATH, os.path.join(LANDSAT_PATH, 'wrs2'))
```

```
[6]: %%time
# WARNING: this will take time the first time it's executed
# depending on your connection
sat_path()
```

```
[6]: folder exists
File already exists
Wall time: 1e+03 mu s
```

```
[7]: # Intersect the shapefile with the WRS2 shapefile to determine relevant path/row tuples
def get_pathrow():
    global paths,rows,path,row, wrs_intersection

    wrs=gpd.GeoDataFrame.from_file(os.path.join('Landsat_images','wrs2',
        'WRS2_descending.shp'))
    wrs_intersection=wrs[wrs.intersects(bbox.geometry[0])]
    paths,rows=wrs_intersection['PATH'].values, wrs_intersection['ROW'].values

    for i, (path,row) in enumerate(zip(paths,rows)):
        print('Image', i+1, ' -path:', path, 'row:', row)
```

```
[8]: get_pathrow()
```

```
[8]: Image 1 -path: 118 row: 38
Image 2 -path: 119 row: 38
```

```
[9]: # Visualise the output of the intersection with the shapefile using Folium

# Get the center of the map
xy = np.asarray(bbox.centroid[0].xy).squeeze()
center = list(xy[:-1])

# Select a zoom
zoom = 8

# Create the most basic OSM folium map
m = folium.Map(location = center, zoom_start = zoom, control_scale=True)

# Add the bounding box (bbox) GeoDataFrame in red using a lambda function
m.add_child(folium.GeoJson(bbox.__geo_interface__, name = 'Area of Interest',
    style_function = lambda x: {'color': 'purple', 'alpha': 0}))

loc = 'Fig 2.: Landsat satellite tiles that cover the Area of Interest'
title_html = '''
    <figcaption align="center" style="font-size:12px"><b>{}</b></figcaption>
    '''.format(loc)
m.get_root().html.add_child(folium.Element(title_html))

# Iterate through each polygon of paths and rows intersecting the area
for i, row in wrs_intersection.iterrows():
    # Create a string for the name containing the path and row of this Polygon
    name = 'path: %03d, row: %03d' % (row.PATH, row.ROW)
    # Create the folium geometry of this Polygon
    g = folium.GeoJson(row.geometry.__geo_interface__, name=name)
    # Add a folium Popup object with the name string
    g.add_child(folium.Popup(name))
    # Add the object to the map
    g.add_to(m)
m
```

```
[9]: text/html<folium.folium.Map at 0x1f0ea0d7dd8>
```

```
[10]: +fvtextcolorcomment_color# Display number of images and Path/Row of the image
for i, (path,row) in enumerate(zip(paths,rows)):
    print('Image', i+1, ' -path:', path, 'row:', row)
```

```
[10]: Image 1 -path: 118 row: 38
Image 2 -path: 119 row: 38
```

Note that here you have two options: 1) continuing and executing the code reported in the next two sections on data download and image cropping, or 2) skipping these sections and proceeding to the image mosaicing sections. We recommend 2) as the processing of unzipping every folder may take long causing the JupyterLab instance to crash.

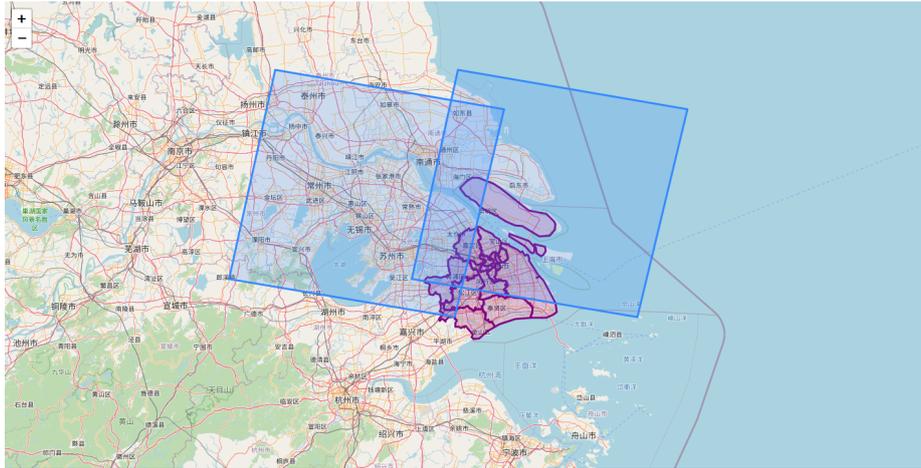


Figure 2: Landsat satellite tiles that cover the Area of Interest

### 2.3 Data download and pre-processing

We now have relevant path and row tuples for our area of analysis. So we can proceed to download satellite images, which are stored on the Google Cloud. To download images, we specify certain parameters: time frame, cloudcover in percentage (0-100 %) and satellite mission (1-5,7,8). The here used Landsat578 API automatically searches the Google Cloud for scenes with the specified parameters and downloads matching images. In order to search the Google Cloud for relevant images, a list of available needs to be downloaded when the code is run for the first time. The list provides basic information of the satellite images and since Landsat data acquisition is ongoing, is updated continuously. Thus, if data from the latest acquisition date is required, it is recommended to re-download the file list before running the code.

We use satellite imagery from a Landsat 5 scene taken in 1984 and a Landsat 8 taken in 2019 to determine neighbourhood changes over time. Landsat 5 scenes can be obtained from two different sensors, the Multispectral Scanner System and the Thematic Mapper, which provide 4 and 7 bands, respectively. The Multispectral Scanner System (MSS) is used in Landsat 1-3 and was superseded by the Thematic Mapper (TM). The MSS provides a green and red band (Band numbers: 1,2) and two infrared bands (Band numbers: 3,4), while the TM provides bands covering red, blue and green (Band numbers: 1,2,3), near-infrared (Band numbers: 4), short-wave infrared (Band numbers: 5,7) and thermal infrared (6). Each downloaded scene contains all bands with one image per band. The different bands can then be stacked in order to highlight various Earth surface processes. In this exercise, scenes from the MSS and TM are downloaded, but only data from the TM is used for analysis.

The Operational Land Imager (OLI) aboard Landsat 8 provides multispectral bands (bands 1-7 and 9) with a resolution of 30 metres and a panchromatic band (band 8) with a resolution of 15 metres (Barsi et al. 2014a). The Thermal Infrared Sensor (TIRS) provides thermal infrared images (bands 10 and 11) with a resolution of 100 meters (Barsi et al. 2014b). The Landsat 8 satellite has a swath width of 185 km for the OLI and TIRS instruments, so one scene usually captures the extent of a city. In other cases, the geographical area of interest may extend beyond one image so that multiple images may be needed (Barsi et al. 2014b, Knight, Kvaran 2014). Given the revisit time of 16 days, usually cloud free images can be retrieved for most cities on a bi-weekly or monthly basis (Roy et al. 2014). The folder and filename of each scene provides information about the satellite, instrument, path/row tuple and date.

Table 3 and Table 4 show which general information of the downloaded scenes can be inferred from the folder and file names of each individual scene:

Table 3: Overview of folder naming convention for Landsat images

Parameter	Meaning
L	Landsat
X	Sensor (“C”=OLI/TIRS combined, “O”=OLI-only, “T”=TIRS-only, “E”=ETM+, “T”=TM, “M”=MSS)
PPP	WRS path
RRR	WRS row
YYYY	Year
DDD	Julian day of year
GSI	Ground station identifier
VV	Archive version number

*Note:* Folder names are structured as LXPPPRRRYYYYDDDGSIIVV

Table 4: Overview of file naming convention for Landsat images

Parameter	Meaning
L	Landsat
X	Sensor (“C”=OLI/TIRS combined, “O”=OLI-only, “T”=TIRS-only, “E”=ETM+, “T”=TM, “M”=MSS)
SS	Satellite (“0”=Landsat 7, “08”=Landsat 8)
LLL	Processing correction level (L1TP/L1GT/L1GS)
PPP	WRS path
RRR	WRS row
YYYYMMDD	Acquisition year, month, day
yyyymmdd	Processing year, month, day
CC	Collection number (01, 02, ...)
TX	Collection category (“RT”=Real-Time, “T1”=Tier 1, “T2”=Tier 2)

*Note:* File names are structured as LXSS\_LLLL.PPPRRR.YYYYMMDD.yyyyymmdd.CC.TX

### 2.3.1 Landsat imagery download

We will now download two Landsat satellite images, one from 1984 and one from 2019. The starting year was chosen due to the increase in spatial resolution to 30 metres with Landsat 4, whereas the end year was chosen at random. The specific dates were selected as the cloud cover was below 5%, ensuring an unobstructed view of the urban area.

```
[11]: # Download Tile list from Google - only needs to be done when first running the code
# NOTE this cell is using the ! magic, which runs command line processes from a Jupyter
# notebook. Make sure the 'landsat' tool, from the 'landsat578' package is installed
# and available

# Path to index file
Index_PATH = os.path.join(directory, '/index.csv.gz')
if os.path.exists(Index_PATH):
    print('File already exists')
else:
    !landsat --update-scenes yes

[12]: # Define Download function to acquire scenes from the Google API
def landsat_download(start_date, end_date, sat,path,row,cloud,output):
    g=GoogleDownload(start=start_date, end=end_date, satellite=sat, path=path,
    ...row=row, max_cloud_percent=cloud, output_path=output)
    g.download()

[13]: # Specify start/end date (in YYYY-MM-DD format), the cloud coverage of the image (in %)
# and the satellite you would like to acquire images from (1-5,7,8). In this case we
# acquire a recent scene from Landsat 8 with a cloud coverage of 5 %.

start_date = '2019-01-01'
end_date = '2019-02-20'
```

```
cloud = 5
satellites = [8]
output = os.path.join(directory, '/Lansat_images/')
```

```
[14]: # Loop through the specified satellites for each path and row tuple
for sat in satellites:
    for i, (path,row) in enumerate(zip(paths,rows)):
        print('Image', i+1, '-path:', path, 'row:', row)
        landsat_download(start_date, end_date,sat,path,row,cloud,output)
```

```
[15]: # The above step is repeated to acquire a Landsat 5 scene from 1984 with 5 % cloud
# coverage.
start_date = '1984-04-22'
end_date = '1984-04-24'
cloud = 5
satellites = [5]
output = os.path.join(directory, '/Lansat_images/')
```

```
[16]: # Loop through the specified satellites for each path and row tuple
for sat in satellites:
    for i, (path,row) in enumerate(zip(paths,rows)):
        print('Image', i+1, '-path:', path, 'row:', row)
        landsat_download(start_date, end_date,sat,path,row,cloud,output)
```

```
[17]: # Delete Scenes that were acquired using the MSS:
outdir = os.listdir(output)
for i in outdir:
    if 'LM' in os.path.basename(i):
        try:
            shutil.rmtree(os.path.abspath(os.path.join(output,os.path.basename(i))))
        except OSError as e:
            print ("Error: %s - %s." % (e.filename, e.strerror))
```

### 2.3.2 Image Cropping

Satellite imagery is large. The size per image can easily equate to 1 GB. It often makes the data processing and analysis computationally expensive. Cropping the obtained scenes to the relevant region of the image enables faster processing and analysing by significantly reducing the size of the input.

```
[18]: # Define cropping function using command line gdalwarp.
## Note: The BQA band is the quality assessment band, which has a different no data
## value (1) than the other bands (0), which makes it necessary to us a different
## cropping function.
def crop(inraster,outraster,shape):
    !gdalwarp -cutline {shape} -srcnodata 0 -crop_to_cutline {inraster} {outraster}
def crop_bqa(inraster,outraster,shape):
    !gdalwarp -cutline {shape} -srcnodata 1 -crop_to_cutline {inraster} {outraster}
```

```
[19]: # Loop through every folder and a create an image cropped to the extent of the shapefile
# save it with the original name and the extension _Cropped
for t in range(0,12):
    for filename in glob.glob((output/'**/*_B{.tif}').format(t), recursive=True):
        inraster = filename
        outraster = filename[:-4] '_Cropped.tif'
        crop(inraster, outraster, shp)
for filename in glob.glob(output/'**/*.tif'):
    if 'BQA.TIF' in i:
        inraster = i
        outraster = i[:-4] '_Cropped.tif'
        crop_bqa(inraster,outraster,shp)
```

### 2.3.3 Image mosaic

As indicated above, a single Landsat scene may not cover the full extent of a city due to the satellite's flight path as can be observed from the interactive map. Creating a mosaic of two or more images is thus often needed to produce a single image that covers the entirety of the area under analysis.

```
[20]: # Read in the relevant Landsat 8 files
output = 'Landsat_images/'
images = sorted(os.listdir(output))
dirpath1 = os.path.join(output, images[0])
dirpath2 = os.path.join(output, images[1])
mosaic_n = os.path.join(output, 'Mosaic/')
search = 'L*_Cropped.tif'
query1 = os.path.join(dirpath1, search)
query2 = os.path.join(dirpath2, search)
files1 = glob.glob(query1)
files2 = glob.glob(query2)
files1.sort()
files2.sort()
if os.path.exists(mosaic_n):
    print('Output Folder exists')
else:
    os.makedirs(mosaic_n)
```

```
[21]: # Match bands together and create a mosaic. Since the BQA band and the cloudmask have
# different denominations than the other bands, these images have to be merged
# together separately.
def mosaic_new(scene1, scene2):
    src_mosaic = []
    string_list = []
    for i, j in zip(scene1, scene2):
        for k in range(1, 12):
            string_list.append('B{}_Cropped'.format(k))
    for l in range(0, 11):
        if string_list[l] in os.path.basename(i) and os.path.basename(j):
            src1 = rasterio.open(i)
            src2 = rasterio.open(j)
            src_mosaic = [src1, src2]
            mosaic, out_trans = rasterio.merge.merge(src_mosaic)
            out_meta = src1.meta.copy()
            out_meta.update({"driver": "GTiff", 'height': mosaic.shape[1],
                            'width': mosaic.shape[2], 'transform': out_trans})
            outdata = os.path.join(mosaic_n, 'B{}_mosaic.tif'.format(l))
            with rasterio.open(outdata, 'w', **out_meta) as dest:
                dest.write(mosaic)
    # Mosaic Quality Assessment Band
    if 'BQA_Cropped' in os.path.basename(i) and os.path.basename(j):
        bqa1 = rasterio.open(i)
        bqa2 = rasterio.open(j)
        bqa_mosaic = [bqa1, bqa2]
        mosaic_, out_trans = rasterio.merge.merge(bqa_mosaic, nodata=1)
        out_meta = bqa1.meta.copy()
        out_meta.update({"driver": "GTiff", 'height': mosaic_.shape[1],
                        'width': mosaic_.shape[2], 'transform': out_trans})
        outdata = os.path.join(mosaic_n, 'BQA_mosaic.tif')
        with rasterio.open(outdata, 'w', **out_meta) as dest:
            dest.write(mosaic_)
    # Mosaic of Cloudmask
    search = 'cloudmask.tif'
    query3 = os.path.join(dirpath1, search)
    query4 = os.path.join(dirpath2, search)
    files3 = glob.glob(query3)
    files4 = glob.glob(query4)
    for i, j in zip(files3, files4):
        if 'cloudmask' in os.path.basename(i) and os.path.basename(j):
            cloudmask1 = rasterio.open(i)
            cloudmask2 = rasterio.open(j)
            cloud_mosaic = [cloudmask1, cloudmask2]
            mosaic_c, out_trans = rasterio.merge.merge(cloud_mosaic, nodata=1)
            out_meta = cloudmask1.meta.copy()
            out_meta.update({"driver": "GTiff", 'height': mosaic_c.shape[1],
                            'width': mosaic_c.shape[2], 'transform': out_trans})
            outdata = os.path.join(mosaic_n, 'Cloudmask_mosaic.tif')
            with rasterio.open(outdata, 'w', **out_meta) as dest:
                dest.write(mosaic_c)
```

```
[22]: mosaic_new(files1, files2)
```

```
[23]: # Read in the relevant files for the Landsat 5 scenes
images = sorted(os.listdir(output))
dirpath_o1 = os.path.join(output, images[2])
dirpath_o2 = os.path.join(output, images[3])
mosaic_o = os.path.join(output, 'Mosaic_old/')
query_o1 = os.path.join(dirpath_o1, search)
query_o2 = os.path.join(dirpath_o2, search)
files_o1 = glob.glob(query_o1)
files_o2 = glob.glob(query_o2)
files_o1.sort()
files_o2.sort()
if os.path.exists(mosaic_o):
    print('Output Folder exists')
else:
    os.makedirs(mosaic_o)
```

```
[24]: # Match bands together and create a mosaic. Since the BQA band and the cloudmask have
# different denominations than the other bands, these images have to be merged together
# separately.
def mosaic_old(scene_o1, scene_o2):
    src_mosaic = []
    string_list = []
    for i, j in zip(scene_o1, scene_o2):

        for k in range(1, 8):
            string_list.append('B{}_Cropped'.format(k))
        for l in range(0, 7):
            if string_list[l] in os.path.basename(i) and os.path.basename(j):
                src1 = rasterio.open(i)
                src2 = rasterio.open(j)
                src_mosaic = [src1, src2]
                mosaic, out_trans = rasterio.merge.merge(src_mosaic)
                out_meta = src1.meta.copy()
                out_meta.update({"driver": "GTiff", 'height': mosaic.shape[1],
                                'width': mosaic.shape[2], 'transform': out_trans})
                outdata = os.path.join(mosaic_o, 'B{}_mosaic.tif'.format(l))
                with rasterio.open(outdata, 'w', **out_meta) as dest:
                    dest.write(mosaic)

        # Mosaic Quality Assessment Band
        if 'BQA_Cropped' in os.path.basename(i) and os.path.basename(j):
            bqa1 = rasterio.open(i)
            bqa2 = rasterio.open(j)
            bqa_mosaic = [bqa1, bqa2]
            mosaic_, out_trans = rasterio.merge.merge(bqa_mosaic, nodata=1)
            out_meta = bqa1.meta.copy()
            out_meta.update({"driver": "GTiff", 'height': mosaic_.shape[1],
                            'width': mosaic_.shape[2], 'transform': out_trans})
            outdata = os.path.join(mosaic_o, 'BQA_mosaic.tif')
            with rasterio.open(outdata, 'w', **out_meta) as dest:
                dest.write(mosaic_)

        # Mosaic of Cloudmask
        search = 'cloudmask.tif'
        query_o3 = os.path.join(dirpath_o1, search)
        query_o4 = os.path.join(dirpath_o2, search)
        files_o3 = glob.glob(query_o3)
        files_o4 = glob.glob(query_o4)
        for i, j in zip(files_o3, files_o4):
            if 'cloudmask' in os.path.basename(i) and os.path.basename(j):
                cloudmask1 = rasterio.open(i)
                cloudmask2 = rasterio.open(j)
                cloud_mosaic = [cloudmask1, cloudmask2]
                mosaic_c, out_trans = rasterio.merge.merge(cloud_mosaic, nodata=1)
                out_meta = cloudmask1.meta.copy()
                out_meta.update({"driver": "GTiff", 'height': mosaic_c.shape[1],
                                'width': mosaic_c.shape[2], 'transform': out_trans})
                outdata = os.path.join(mosaic_o, 'Cloudmask_mosaic.tif')
                with rasterio.open(outdata, 'w', **out_meta) as dest:
                    dest.write(mosaic_c)
```

```
[25]: mosaic_old(files_o1, files_o2)
```

### 2.3.4 Natural-colour (True-colour) composition

Our downloaded data from Landsat 8 and Landsat 5 have different band designations. Combining different satellite bands are useful to identify features of the urban environment: vegetation, built-up areas, ice and water. We create a standard natural-colour composition image using Red, Green and Blue satellite bands. This colour composition best reflects the natural environment. For instance, trees are green; snow and clouds are white; and water is blue. Landsat 8 has 11 bands with bands 4, 3 and 2 corresponding to Red, Green and Blue respectively. Landsat 5 has 7 bands with bands 3, 2 and 1, corresponding to Red, Green and Blue. We perform layer stacking to produce a true colour image composition to gain understanding of the local area before extracting and analysing features of the urban environment.

```
[26]: # Normalise the bands to so that they can be combined to a single image
def normalize(array):
    """Normalizes numpy arrays into scale 0.0 - 1.0"""
    array_min, array_max = array.min(), array.max()
    return ((array - array_min)/(array_max - array_min))

[27]: # Adjust the intensity of each band for visualisation.
# This is a way of rescaling each band by clipping the pixels that are outside the
# specified range to the range we defined. By adjusting the gamma, we change the
# brightness of the image with gamma >1 resulting in a brighter image. However
# there are more complex methods such as top of the atmosphere corrections, which
# subtracts any atmospheric interference from the image.
# For the purpose of this notebook, this way is sufficient.
def rescale_intensity(image):
    p2, p98 = np.percentile(image, (0.2, 98))
    img_exp = exposure.rescale_intensity(image, in_range=(p2, p98))
    img_gamma = exposure.adjust_gamma(img_exp, gamma=2.5, gain=1)
    return(img_gamma)

[28]: # Downsample image resolution with factor 0.5 for displaying purposes.
def downsample(file):
    downscale_factor=0.5
    data = file.read(1,
        out_shape=(
            file.count,
            int(file.height * downscale_factor),
            int(file.width * downscale_factor)
        ),
        resampling=Resampling.bilinear
    )
    # scale image transform
    transform = file.transform * file.transform.scale(
        (file.width / data.shape[-1]),
        (file.height / data.shape[-2])
    )
    return data

[29]: # Use rasterio to open the Red, Blue and Green bands of the mosaic image from 1984
# to create an RGB image
# **NOTE**: The Mosaic names do not correspond to the actual band designations as
# python starts counting at 0!
with rasterio.open('Landsat_images/Mosaic_old/B0_mosaic.tif') as band1_old:
    b1_old=downsample(band1_old)
with rasterio.open('Landsat_images/Mosaic_old/B1_mosaic.tif') as band2_old:
    b2_old=downsample(band2_old)
with rasterio.open('Landsat_images/Mosaic_old/B2_mosaic.tif') as band3_old:
    b3_old=downsample(band3_old)

[30]: # Normalise the bands so that they can be combined to a single image
red_old_n = normalize(b3_old)
green_old_n = normalize(b2_old)
blue_old_n = normalize(b1_old)

# Apply the function defined before to make more natural-looking image
red_adj = rescale_intensity(red_old_n)
green_adj = rescale_intensity(green_old_n)
blue_adj = rescale_intensity(blue_old_n)
```



Figure 3: True colour Landsat image of the Shanghai urban area from 1984

```
# Stack the three different bands together
rgb_2 = np.dstack((red_adj,green_adj,blue_adj))

# Visualise the true color image
fig,ax = plt.subplots(figsize=(10,10))
ax.imshow(rgb_2)
plt.title('Fig.3: True color Landsat image of the Shanghai urban area from 1984',
          y=-0.1, fontsize=12)
plt.show()
plt.close()
del rgb_2,b1_old,b2_old,b3_old,red_adj,green_adj,blue_adj
```

[30]: image/png<Figure size 720x720 with 1 Axes>

```
[31]: # Use rasterio to open the Red, Blue and Green bands of the mosaic image from 2019
# to create an RGB image
# **NOTE**: The Mosaic names do not correspond to the actual band designations as
# python starts counting at 0!!
with rasterio.open('Landsat_images/Mosaic/B1_mosaic.tif') as band2_new:
    b2_new = downsample(band2_new)
with rasterio.open('Landsat_images/Mosaic/B2_mosaic.tif') as band3_new:
    b3_new = downsample(band3_new)
with rasterio.open('Landsat_images/Mosaic/B3_mosaic.tif') as band4_new:
    b4_new = downsample(band4_new)
```

```
[32]: # Normalise the bands so that they can be combined to a single image
red_new_n = normalize(b4_new)
green_new_n = normalize(b3_new)
blue_new_n = normalize(b2_new)

# Apply the function defined before to make more natural-looking image
red_rescale = rescale_intensity(red_new_n)
green_rescale = rescale_intensity(green_new_n)
blue_rescale = rescale_intensity(blue_new_n)

# Stack the three different bands together
rgb = np.dstack((red_rescale, green_rescale, blue_rescale))

# Here we adjust the gamma (brightness) for the stacked image to achieve a more
# natural looking image.
rgb_adjust = exposure.adjust_gamma(rgb, gamma = 1.5, gain=1)

# Visualise the true color image
fig,ax = plt.subplots(figsize=(10,10))
```



Figure 4: True colour Landsat image of the Shanghai urban area from 2019

```
ax.imshow(rgb_adjust)
plt.title('Fig.4: True color Landsat image of the Shanghai urban area from 2019',
          y=-0.1, fontsize=12)
plt.show()
plt.close()
del rgb,red_new_n,green_new_n,blue_new_n,red_rescale,green_rescale,blue_rescale,
     rgb_adjust
```

[32]: Figure size 720x720 with 1 Axes

When comparing the true colour Landsat satellite images in Figures 3 and 4, the urbanisation of Shanghai between 1984 and 2019 is apparent. In the following steps, we will analyse and quantify these urban changes.

### 3 Feature extraction

Since the above two maps show that urban neighbourhoods of Shanghai have undergone dramatic changes over time in colour, texture, greenery, buildings, etc., the next stage is to gain valuable information out of satellite images and interpret these changes. Since the images we have downloaded are on a city-wide scale, which covers more than a thousand kilometre spatial resolution and less detailed. Therefore, feature extraction is performed to get a reduced representation of the initial image but informative and sufficiently accurate for subsequent analysis and interpretation.

We examine four sets of features based on the above two true colour maps and the scale, where the colour, texture, greenery, and buildings changed a lot during the past 25 years in Shanghai. Specifically, colour and texture features extracted from true colour imagery (i.e. RGB bands composition represented by bands 1-3 and bands 2-4 in 1984 and 2019), and vegetation features and built-up features extracted from Red, near infrared (NIR) and shortwave infrared (SWIR) bands, represented by bands 3-5 and bands 4-6 in 1984 and 2019. More detailed information about the meaning of each band can be found at [https://www.usgs.gov/faqs/what-are-best-landsat-spectral-bands-use-my-research?qt-news\\_science\\_products=0#qt-news\\_science\\_products](https://www.usgs.gov/faqs/what-are-best-landsat-spectral-bands-use-my-research?qt-news_science_products=0#qt-news_science_products). In this analysis, colour features measure the colour moments of true colour imagery to interpret colour distribution; texture features apply LBP (Local binary patterns) texture spectrum model to show spatial distribution of intensity values in an image; vegetation features calculate the NDVI (Normalised difference vegetation index) to capture the amount of vegetation, and built-up features calculate NDBI (Normalised difference built-up index) to highlight artificially constructed areas.



Figure 5: Spatial distribution of all administrative divisions of Shanghai

The administrative divisions of Shanghai have experienced tremendous changes in the last tens of years (Ministry of Civil Affairs of the People’s Republic of China 2018), thus, we will conduct feature extraction of imagery on the current administrative boundaries to explore if satellite imagery can be used to reflect and interpret urban changes. The figure below shows the spatial distribution of each administrative area with relative labels in Shanghai.

```
[33]: # read administrative boundary shapefile of Shanghai
poly = gpd.read_file(shp)

f, ax = plt.subplots(1, figsize = (9,9))
poly.plot(ax = ax)
# create a new column, in order to plot polygon labels (i.e. name) in the map
poly['coords']=poly['geometry'].apply(lambda x:x.representative_point().coords[:])
poly['coords']=[coords[0] for coords in poly['coords']]
for idx, row in poly.iterrows():
    ax.annotate(text=row['Name'],xy=row['coords'],va='center',ha='center',alpha = 0.8,
                fontsize = 8)
plt.axis('equal')
plt.axis('off')
f.suptitle('Fig.5: Spatial distribution of all administrative divisions of Shanghai',
           y=-0.1,fontsize = 12)
```

```
[33]: Text(0.5, -0.1, 'Fig.5: Spatial distribution of all administrative divisions of
Shanghai')
image/png<Figure size 648x648 with 1 Axes>
```

Figure 5 shows that administrative divisions of ‘Chongming’ in the north appear three geometries. Therefore, it is necessary to check if they belong to a single administrative unit.

```
[34]: poly.loc[poly['Name']== 'Chongming','Name']
```

```
[34]: 0    Chongming
3    Chongming
5    Chongming
Name: Name, dtype: object
```

Chongming administrative division consist of three separate geometries, which may confuse our further analysis. As a result, we dissolved these geometries into a single geometric feature and take a look at the new dataset. The below table shows that the Chongming administrative division now consists of multipolygons which includes all polygons as a whole.

```
[35]: # Dissolve geometries with the identical names together
poly = poly.dissolve(by = 'Name').reset_index()
# Have a look at the name of all administrative unit and we can see that chongming
# districts have been dissolved into a single administrative unit
poly['Name'].values
```

```
[35]: array(['Baoshan', 'Changning', 'Chongming', 'Fengxian', 'Hongkou',
        'Huangpu', 'Jiading', 'Jinshan', 'Minhang', 'Pudong New', 'Putuo',
        'Qingpu', 'Songjiang', 'Xuhui', 'Yangpu', 'Zhabei'], dtype=object)
```

### 3.1 Image processing

Further pre-processing of satellite imagery is needed before feature extraction. This pre-processing involves three steps:

1. Masking (cropping) of raster files (i.e., Blue, Green, Red, Nir and SWIR bands) into each administrative district polygon;
2. Image enhancement to improve the quality and content of the original image; and,
3. Band stacking based on each neighbourhood unit.

```
[36]: # open raster files
file_list_old = sorted(glob.glob('Landsat_images/Mosaic_old' + "/*.tif", recursive = True))
files_old = [rio.open(filename) for filename in file_list_old]
```

```
[37]: file_list = sorted(glob.glob('Landsat_images/Mosaic' + "/*.tif"))
files = [rio.open(filename) for filename in file_list]
```

Before cropping all raster files into each polygon in the vector file (i.e. Shanghai administrative area shapefile), we have to ensure they have the same coordinate reference system (CRS). Once matched, the cropping process is prepared to go.

```
[38]: poly.crs
```

```
[38]: {'init': 'epsg:4326'}
```

```
[39]: # check the crs of one band of satellite imagery
files[0].crs
```

```
[39]: CRS.from_epsg(32651)
```

```
[40]: # reproject the vector file to make it consistent with raster files
poly = poly.to_crs('EPSG:32651')
```

```
[41]: # get each neighbourhood geographic boundary based on administrative area data
geo = [poly.__geo_interface__['features'][i]['geometry']
        for i in range(len(poly))]
```

```
[42]: # clip R,G,B bands separately by each poly, so get pixel values in each poly and save
# them into a list
out_image = [[] for i in range(5)]
img_old = [[] for i in range(5)]

# x: Blue,Green,Red,NIR and SWIR bands, y: 16 polygons from vector file
for x,y in itertools.product(range(5),range(len(geo))):
    # out_image[0] means masked Blue band polygon
    out_image[x].append(mask(files_old[0:5][x], [geo[y]], crop=True))
    # image enhancement: normalisation and Histogram Equalization
    img_old[x].append(exposure.equalize_hist(normalize(out_image[x][y][0][0])))
del out_image,files_old
```

```
[43]: # clip R,G,B bands separately by each poly, so get pixel values in each poly and save
# them into a list
out_image = [[] for i in range(5)]
img_new = [[] for i in range(5)]

# x: Blue,Green,Red,NIR and SWIR bands, y: 16 polygons from vector file
```

```

for x,y in itertools.product(range(5),range(len(geo))):
    # out_image[0] means masked blue polygon
    out_image[x].append(mask(files[0:5][x], [geo[y]], crop=True))
    # image enhancement: normalisation and Histogram Equalization
    img_new[x].append(exposure.equalize_hist(normalize(out_image[x][y][0][0])))
del out_image,files

```

```
[44]: # have a look at the pixel values of one geographic area in blue band
img_new[0][0]
```

```
[44]: array([[0.48515378, 0.48515378, 0.48515378, ..., 0.48515378, 0.48515378,
         0.48515378],
        [0.48515378, 0.48515378, 0.48515378, ..., 0.48515378, 0.48515378,
         0.48515378],
        [0.48515378, 0.48515378, 0.48515378, ..., 0.48515378, 0.48515378,
         0.48515378],
        ...,
        [0.48515378, 0.48515378, 0.48515378, ..., 0.48515378, 0.48515378,
         0.48515378],
        [0.48515378, 0.48515378, 0.48515378, ..., 0.48515378, 0.48515378,
         0.48515378],
        [0.48515378, 0.48515378, 0.48515378, ..., 0.48515378, 0.48515378,
         0.48515378]])
```

```
[45]: # stack R,G,B bands together for later feature extraction
bb = [img_old[0][x].astype(np.float) for x in range(len(geo))]
bg = [img_old[1][x].astype(np.float) for x in range(len(geo))]
br = [img_old[2][x].astype(np.float) for x in range(len(geo))]
```

```
[46]: rgb_old = [np.dstack((br[x],bg[x],bb[x])) for x in range(len(geo))]
```

```
[47]: bb = [img_new[0][x].astype(np.float) for x in range(len(geo))]
bg = [img_new[1][x].astype(np.float) for x in range(len(geo))]
br = [img_new[2][x].astype(np.float) for x in range(len(geo))]
```

```
[48]: rgb_new = [np.dstack((br[x],bg[x],bb[x])) for x in range(len(geo))]
```

### 3.2 Colour features

Colour features are used to extract the characteristics of colours from satellite imagery. A commonly used method to extract colour features is to compute colour moments of an image. Colour moments provide a measurement of colour similarity between images (Keen 2005). Basically, colour probability distributions of an image are characterised by a range of unique moments. The mean, standard deviation and skewness these three central moments are generally used to identify colour distribution. Here we extract colour features on HSV (Hue, Saturation and Value) colour space because it corresponds to human vision and has been widely used in computer vision. HSV colour space can be converted from RGB colour channels, Hue represents the colour portion, saturation represents the amount of grey in a particular colour (0 is grey), and Value represents the brightness of the colour (0 is black). Therefore, the true-colour imagery is characterised by a total of nine moments - three moments for each HSV channel in the same units.

```
[49]: # interpret the color probability distribution by computing low order color
# moments(1,2,3)
def color_moments(img):
    if img is None:
        return
    # Convert RGB to HSV colour space
    img_hsv = rgb2hsv(img)
    # Split the channels - h,s,v
    h, s, v = [img_hsv[:, :, i] for i in [0,1,2]]
    # Initialize the colour feature
    color_feature = []
    # N = h.shape[0] * h.shape[1]
    # The first central moment - average
    h_mean = np.mean(h) # np.sum(h)/float(N)
    s_mean = np.mean(s) # np.sum(s)/float(N)
    v_mean = np.mean(v) # np.sum(v)/float(N)

```

Table 5: Partial colour features identified in 1984

Name	h_mean	s_mean	v_mean	h_std	s_std	v_std	h_skew	s_skew	v_skew
Baoshan	0.27216	0.05208	0.64415	0.32709	0.07246	0.18531	0.35671	0.09006	0.19945
Changning	0.22141	0.05156	0.65989	0.28837	0.07518	0.17446	0.33080	0.09250	0.18763
Chongming	0.15381	0.01739	0.74231	0.27216	0.03563	0.10235	0.33292	0.05118	0.12060
Fengxian	0.33961	0.11292	0.60576	0.32194	0.12281	0.24362	0.34723	0.14439	0.25767
Hongkou	0.24951	0.06370	0.65073	0.30983	0.08744	0.18781	0.34797	0.10619	0.20013

```

color_feature.extend([h_mean, s_mean, v_mean])
# The second central moment - standard deviation
h_std = np.std(h) # np.sqrt(np.mean(abs(h - h.mean())**2))
s_std = np.std(s) # np.sqrt(np.mean(abs(s - s.mean())**2))
v_std = np.std(v) # np.sqrt(np.mean(abs(v - v.mean())**2))
color_feature.extend([h_std, s_std, v_std])
# The third central moment - the third root of the skewness
h_skewness = np.mean(abs(h - h.mean())**3)
s_skewness = np.mean(abs(s - s.mean())**3)
v_skewness = np.mean(abs(v - v.mean())**3)
h_thirdMoment = h_skewness**(1./3)
s_thirdMoment = s_skewness**(1./3)
v_thirdMoment = v_skewness**(1./3)
color_feature.extend([h_thirdMoment, s_thirdMoment, v_thirdMoment])

return color_feature

```

```
[50]: # create and initialize a data table to store colour features
color_mom_old = pd.DataFrame(color_moments(rgb_old[0]))
# add the rest columns by assigning 9 color moments in each poly
for i in range(1, len(rgb_old)):
    color_mom_old[i] = color_moments(rgb_old[i])
    i = i+1

```

```
[51]: # create and initialize a data table
color_mom_new = pd.DataFrame(color_moments(rgb_new[0]))
# add the rest columns by assigning 9 color moments in each poly
for i in range(1, len(rgb_new)):
    color_mom_new[i] = color_moments(rgb_new[i])
    i = i+1

```

```
[52]: # Data manipulation
color_old_var = color_mom_old.T
# assign column names
color_old_var.columns =
    ['h_mean', 's_mean', 'v_mean', 'h_std', 's_std', 'v_std', 'h_skew', 's_skew', 'v_skew']
# set geographic name as index
color_old_var = color_old_var.set_index(poly.Name)

```

```
[53]: color_new_var = color_mom_new.T
color_new_var.columns =
    ['h_mean', 's_mean', 'v_mean', 'h_std', 's_std', 'v_std', 'h_skew', 's_skew', 'v_skew']
color_new_var = color_new_var.set_index(poly.Name)

```

As we have created two new tables for colour features in the year 1984 and 2019, it would be helpful to have a view of the tables and see how they look like. Table 5 and Table 6 show nine variables (column) representing colour features within five administrative division of Shanghai (row).

```
[54]: # check the information of colour feature
color_old_var.head().style.set_caption('Table 5: Partial colour features
... identified in 1984')
```

```
[54]: text/html<pandas.io.formats.style.Styler at 0x1f0ed9c4be0>
```

```
[55]: color_new_var.head().style.set_caption('Table 6: Partial colour features
... identified in 2019')
```

```
[55]: text/html<pandas.io.formats.style.Styler at 0x1f081f31518>
```

Table 6: Partial colour features identified in 2019

Name	h_mean	s_mean	v_mean	h_std	s_std	v_std	h_skew	s_skew	v_skew
Baoshan	0.23107	0.03587	0.63894	0.29785	0.05205	0.18019	0.33678	0.06787	0.19559
Changning	0.23185	0.03129	0.64924	0.30469	0.04847	0.16700	0.34439	0.06297	0.18248
Chongming	0.15731	0.01647	0.74240	0.28250	0.03184	0.10177	0.34529	0.04336	0.11980
Fengxian	0.29554	0.08620	0.60543	0.30273	0.09770	0.24360	0.32939	0.11570	0.25723
Hongkou	0.23994	0.03758	0.63861	0.30383	0.05505	0.18294	0.33962	0.07055	0.19762

### 3.3 Texture features

To extract texture features, we use a Local Binary Pattern (LBP) approach. LBP searches for pixels adjacent to a central point and tests whether these surrounding pixels are greater or less than the central pixel and generate a binary classification (Pedregosa et al. 2011) ([https://scikit-image.org/docs/dev/auto\\_examples/features\\_detection/plot\\_local\\_binary\\_pattern.html](https://scikit-image.org/docs/dev/auto_examples/features_detection/plot_local_binary_pattern.html)). In theory, eight adjacent neighbour pixels in greyscale are set to compare with one central pixel value by 3 \* 3 neighbourhood threshold, and consider the result as 1 or 0 (Ojala et al. 1996). Thus, these eight surrounding binary numbers correspond to LBP code for the central pixel value, determining the texture pattern of that threshold. Texture features are then the distribution of a collection of LBPs over an image.

```
[56]: # convert a RGB image into Grayscale, which takes less space for analysis
gray_images_old = [rgb2gray(rgb_old[i]) for i in range(len(rgb_old))]
gray_images_new = [rgb2gray(rgb_new[i]) for i in range(len(rgb_new))]
```

```
[57]: # settings for LBP
radius = 1 # radius = 1 refers to a 3*3 patch/window scale
n_points = 8 * radius # the number of circularly symmetric neighbour set points
method = 'uniform' # finer quantization of the angular space which is gray scale and
# rotation invariant

lbps_old = [local_binary_pattern(gray_images_old[i], n_points, radius, method)
... for i in range(len(rgb_old))]
lbps_new = [local_binary_pattern(gray_images_new[i], n_points, radius, method)
... for i in range(len(rgb_new))]
```

```
[58]: # n_bins are the same in each neighbourhood
n_bins = int(lbps_old[0].max()+1)
# define a function to count the number of points in a given bin of LBP distribution
# histogram
def count_hist(x):
    return np.histogram(lbps_old[x].ravel(), density=True, bins=n_bins, range=(0, n_bins))
# Assign counts to a new list, return the histogram vector features in this cell (polygon)
hist_features_old = [count_hist(i)[0] for i in range(len(rgb_old))]
```

```
[59]: # Extract texture features of another year based on same method
n_bins = int(lbps_new[0].max()+1)

def count_hist(x):
    return np.histogram(lbps_new[x].ravel(), density=True, bins=n_bins, range=(0, n_bins))

# Assign counts to a new list, return the histogram vector features in this cell (polygon)
hist_features_new = [count_hist(i)[0] for i in range(len(rgb_new))]
```

Same with operations on colour features, this time we build two new tables (Table 7 and 8) for texture features, with each row present administrative division and each column represent texture feature.

```
[60]: # The histogram features are the texture features
texture_old_var = pd.DataFrame([hist_features_old[a] for a in range(len(rgb_old))])
texture_old_var.columns = ['LBP'+ str(i) for i in range(n_bins)]
texture_old_var = texture_old_var.set_index(poly.Name)
# Have a look at the table with texture features of administrative division of
# Shanghai in 1984
texture_old_var.head().style.set_caption('Table 7: Partial texture features
... identified in 1984')
```

Table 7: Partial texture features identified in 1984

Name	LBP0	LBP1	LBP2	LBP3	LBP4	LBP5	LBP6	LBP7	LBP8	LBP9
Baoshan	0.03509	0.04196	0.04071	0.06839	0.07839	0.06748	0.04034	0.04110	0.52005	0.06648
Changning	0.03609	0.04608	0.04196	0.05979	0.06004	0.06442	0.03754	0.04339	0.53942	0.07129
Chongming	0.02582	0.02995	0.02176	0.03406	0.03958	0.03679	0.02404	0.02916	0.70944	0.04941
Fengxian	0.05551	0.06647	0.05123	0.07200	0.07377	0.07393	0.05291	0.06510	0.38211	0.10698
Hongkou	0.04202	0.05056	0.04354	0.05933	0.05676	0.07072	0.03969	0.04649	0.51043	0.08047

Table 8: Partial texture features identified in 2019

Name	LBP0	LBP1	LBP2	LBP3	LBP4	LBP5	LBP6	LBP7	LBP8	LBP9
Baoshan	0.04306	0.04774	0.04077	0.05862	0.06808	0.05681	0.03741	0.04557	0.52419	0.07777
Changning	0.04264	0.05012	0.03767	0.05137	0.05842	0.06153	0.03524	0.04708	0.53945	0.07648
Chongming	0.02547	0.02964	0.02333	0.03522	0.04762	0.03641	0.02359	0.02865	0.70442	0.04565
Fengxian	0.05121	0.06105	0.05288	0.08141	0.09716	0.07923	0.05152	0.06044	0.36993	0.09517
Hongkou	0.04703	0.05417	0.04294	0.05439	0.05501	0.06805	0.03879	0.04894	0.50732	0.08335

```
[60]: text/html<pandas.io.formats.style.Styler at 0x1f081ebe630>
```

```
[61]: # The histogram features are the texture features
texture_new_var = pd.DataFrame([hist_features_new[a] for a in range(len(rgb_new))])
texture_new_var.columns = ['LBP' + str(i) for i in range(n_bins)]
texture_new_var = texture_new_var.set_index(poly.Name)
# Have a look at the table with texture features of administrative division of
# Shanghai in 2019
texture_new_var.head().style.set_caption('Table 8: Partial texture features
... identified in 2019')
```

```
[61]: text/html<pandas.io.formats.style.Styler at 0x1f081ebeac8>
```

### 3.4 Vegetation and built-up features

Vegetation features and built-up features can be measured by calculating fundamental NDVI and NDBI indices in each administrative area respectively. The Normalized Difference Vegetation Index (NDVI) is a normalized index, using Red and NIR bands to display the amount of vegetation (NASA 2000). The use of NDVI maximizes the reflectance properties of vegetation by minimizing NIR and maximizing the reflectance in the red wavelength. The measure is used to distinguish vegetation in regions, as more vegetation will affect the ratio of visible light absorbed and near-infrared light reflected. The formula is as follows:

$$\text{NDVI} = (\text{NIR} - \text{Red}) / (\text{NIR} + \text{Red})$$

The output value of this index is between -1.0 and 1.0. Close to 0 represents no vegetation, close to 1 indicates the highest possible density of green leaves, and close to -1 indicates water bodies.

The Normalized Difference Built-up Index (NDBI) uses the NIR and SWIR bands to highlight artificially constructed areas (built-up areas) where there is a typically a higher reflectance in the shortwave infrared region than the near infrared region (Zha et al. 2003). The index is a ratio type that reduces the effects of differences in terrain illumination and atmospheric effects. The formula is as follows:

$$\text{NDBI} = (\text{SWIR} - \text{NIR}) / (\text{SWIR} + \text{NIR})$$

Also, the output value of this index is between -1 to 1. Higher values represent built-up areas whereas negative values represent water bodies.

After calculating these two indices, vegetation features and built-up features can be measured by calculating average values of index values within each administrative area.

Table 9: Partial vegetation and built-up features identified in 1984

Name	veg_mean	builtup_mean
Baoshan	-0.002218	0.000611
Changning	-0.002147	0.000582
Chongming	-0.000805	0.000190
Fengxian	-0.007201	0.001499
Hongkou	-0.004648	-0.000313

## 3.4.1 Vegetation features

```
[62]: # identify red and NIR band to each neighbourhood unit in 1984
red_old, nir_old = img_old[2],img_old[3]
# Calculate ndvi, assign 0 to nodata pixels
ndvi_old = [np.where((nir_old[i] red_old[i])==0, 0,
                    (nir_old[i]-red_old[i])/(nir_old[i] red_old[i]))
            for i in range(len(poly))]
```

```
[63]: # identify red and NIR band to each neighbourhood unit in 1984
red_new, nir_new = img_new[2],img_new[3]
# Calculate ndvi, assign 0 to nodata pixels
ndvi_new = list(map(lambda i: np.where((nir_new[i] red_new[i])==0, 0,
                    (nir_new[i]-red_new[i])/(nir_new[i] red_new[i])),
                    list(range(len(poly)))
                ))
```

```
[64]: veg_old_var = pd.DataFrame([np.mean(ndvi_old[i]) for i in range(len(poly))],
                                index = poly.Name, columns = ['veg_mean'])
```

```
[65]: veg_new_var = pd.DataFrame([np.mean(ndvi_new[i]) for i in range(len(poly))],
                                index = poly.Name, columns = ['veg_mean'])
```

## 3.4.2 Built-up features

```
[66]: # identify red and NIR band to each neighbourhood unit in 1984
nir_old, swir_old = img_old[3],img_old[4]
# Calculate ndbi, assign 0 to nodata pixels
ndbi_old = [np.where((nir_old[i] swir_old[i])==0., 0,
                    (swir_old[i] - nir_old[i])/(nir_old[i] swir_old[i]))
            for i in range(len(poly))]
```

```
[67]: # identify red and NIR band to each neighbourhood unit in 1984
nir_new, swir_new = img_new[3],img_new[4]
# Calculate ndbi, assign 0 to nodata pixels
ndbi_new = list(map(lambda i: np.where((nir_new[i] swir_new[i])==0., 0,
                    (swir_new[i] - nir_new[i])/(nir_new[i] swir_new[i])),
                    list(range(len(poly)))
                ))
```

```
[68]: builtup_old_var = pd.DataFrame([np.mean(ndbi_old[i]) for i in range(len(poly))],
                                    index = poly.Name, columns = ['builtup_mean'])
```

```
[69]: builtup_new_var = pd.DataFrame([np.mean(ndbi_new[i]) for i in range(len(poly))],
                                    index = poly.Name, columns = ['builtup_mean'])
```

Table 9 and Table 10 created as shown below contain both vegetation features (NDVI) and builtup features (NDBI), with the mean value of vegetation features and built-up features (two columns) calculated at each administrative division (row).

```
[70]: veg_built_old = pd.concat([veg_old_var,builtup_old_var], axis = 1)
veg_built_old.head().style.set_caption('Table 9: Partial vegetation and built-up
... features identified in 1984')
```

```
[70]: text/html<pandas.io.formats.style.Styler at 0x1f081e1b1d0>
```

Table 10: Partial vegetation and built-up features identified in 2019

Name	veg_mean	builtup_mean
Baoshan	-0.001801	0.001938
Changning	-0.001515	0.000774
Chongming	-0.000705	0.000318
Fengxian	-0.008185	-0.000408
Hongkou	-0.002057	-0.000277

```
[71]: veg_built_new = pd.concat([veg_new_var,builtup_new_var], axis = 1)
veg_built_new.head().style.set_caption('Table 10: Partial vegetation and built-up
... features identified in 2019')
```

```
[71]: text/html<pandas.io.formats.style.Styler at 0x1f0ed9c4828>
```

#### 4 Feature clustering

Now we have four types of features: colour, texture, vegetation and built-up area for Shanghai in 1984 and 2019. These features are the embodiment of urban changes and vary greatly due to rapid urbanisation and development. Therefore, the subsequent task is to identify systematic patterns from these integrated features for analysis of urban changes, such as whether several administrative areas share similar patterns. A clustering method is required within this context to group these geographical divisions that are similar within each other but different between them. Considering the ease of computation and fast implementation, we use generalised and the most popular k-means clustering to identify representative types of neighbourhoods based on multiple features. K-means clustering partitions the data by creating k groups of equal variance, minimising the within-cluster sum of squares (Pedregosa et al. 2011). We can perform K-means using the package `scikit-learn`, which is a powerful machine learning package for Python.

```
[72]: # merge all features together
features_old_var = pd.concat([color_old_var,texture_old_var,veg_old_var,
    builtup_old_var], axis = 1)
features_old_var.head().style.set_caption('Table 11: Four types of features
... (21 in total) identified in 1984')
```

```
[72]: text/html<pandas.io.formats.style.Styler at 0x1f0ed9c4908>
```

```
[73]: # merge all features together
features_new_var = pd.concat([color_new_var,texture_new_var,veg_new_var,
    builtup_new_var], axis=1)
features_new_var.head().style.set_caption('Table 12: Four types of features
... (21 in total) identified in 2019')
```

```
[73]: text/html<pandas.io.formats.style.Styler at 0x1f081f31438>
```

Table 11 and Table 12 reveal the integrated 21 features across our four sets of image features and their differences at geographical division in magnitude between 1984 and 2019. Since k-means clustering is one of the machine learning algorithms, which generally expect data transformation for preprocessing before fitting the algorithm. We therefore use one of the most popular rescale methods to standardise these features to lie between 0 and 1 based on `MinMaxScaler()` function in `scikit-learn` package. The motivation of this method relies on the robustness to very small standard deviation. This preprocess ensures individual features of dataset have the same scale that standard normally distributed.

```
[74]: # Last preprocessing step before machine learning: data rescaling
min_max_scaler = preprocessing.MinMaxScaler()
np_scaled = min_max_scaler.fit_transform(features_old_var)
oldvar_scale = pd.DataFrame(np_scaled)
oldvar_scale.columns = features_old_var.columns
```

Table 11: Four types of features (21 in total) identified in 1984

Name	h_mean	s_mean	v_mean	h_std	s_std	v_std	h_skew
Baoshan	0.272161	0.052081	0.644148	0.327094	0.072457	0.185309	0.356713
Changning	0.221412	0.051564	0.659894	0.288368	0.075177	0.174455	0.330803
Chongming	0.153807	0.017394	0.742309	0.272162	0.035627	0.102347	0.332916
Fengxian	0.339613	0.112915	0.605758	0.321941	0.122805	0.243621	0.347226
Hongkou	0.249526	0.063704	0.650725	0.309825	0.087439	0.187805	0.347968

Name	s_skew	v_skew	LBP0	LBP1	LBP2	LBP3	LBP4
Baoshan	0.090057	0.199446	0.035093	0.041960	0.040705	0.068394	0.078389
Changning	0.092504	0.187627	0.036086	0.046078	0.041956	0.059792	0.060040
Chongming	0.051184	0.120603	0.025822	0.029946	0.021757	0.034058	0.039580
Fengxian	0.144392	0.257670	0.055508	0.066468	0.051230	0.072002	0.073767
Hongkou	0.106194	0.200131	0.042018	0.050562	0.043542	0.059326	0.056759

Name	LBP5	LBP6	LBP7	LBP8	LBP9	veg_mean	builtup_mean
Baoshan	0.067483	0.040339	0.041101	0.520053	0.066483	-0.002218	0.000611
Changning	0.064422	0.037538	0.043385	0.539416	0.071285	-0.002147	0.000582
Chongming	0.036787	0.024035	0.029158	0.709444	0.049413	-0.000805	0.000190
Fengxian	0.073928	0.052907	0.065099	0.382110	0.106981	-0.007201	0.001499
Hongkou	0.070718	0.039691	0.046490	0.510429	0.080465	-0.004648	-0.000313

```
[75]: min_max_scaler = preprocessing.MinMaxScaler()
np_scaled = min_max_scaler.fit_transform(features_new_var)
newvar_scale = pd.DataFrame(np_scaled)
newvar_scale.columns = features_new_var.columns
```

Above two steps are the results of data transformation in 1984 and 2019. To identify robust and consistent clustering results, we merge them into a single one based on their common geographical units (see Table 13). The column names ended with ‘\_x’ and ‘\_y’ represent features extracted in 1984 and 2019, respectively. This table is the one prepared for the final k-mean clustering analysis. The dominant parameter in k-means clustering is the number of clusters (i.e., k), determining the optimal numbers of clusters is therefore a fundamental issue. We select a direct and popular elbow method as an example to assess the resulting partitions, testing nine different solutions varying k from 2 to 10. Basically, the idea of elbow method is to define clusters to minimise the total intra-cluster variation or total within-cluster sum of square (WSS). The optimal number can be determined by plotting the curve of WSS according to different k clusters and the location of a bend is considered as an indicator of the appropriate number for k.

```
[76]: merged_var = pd.merge(oldvar_scale, newvar_scale, left_index = True, right_index = True)
merged_var.head().style.set_caption('Table 13: Integrated preprocessed features
... identified in 1984 and 2019 seperately')
```

```
[76]: text/html<pandas.io.formats.style.Styler at 0x1f081e1b6d8>
```

```
[77]: # elbow analysis
cluster_range = range( 2, 11 )
cluster_errors = []

for num_clusters in cluster_range:
    clusters = KMeans( num_clusters )
    clusters.fit( merged_var )
    cluster_errors.append( clusters.inertia_ )
clusters_df = pd.DataFrame( { "num_clusters":cluster_range,
                             "cluster_errors": cluster_errors } )

plt.figure(figsize=(12,6))
plt.title('Fig.6: Elbow method to determine the optimal k for k-mean clustering',y=-0.2)
plt.plot( clusters_df.num_clusters, clusters_df.cluster_errors, marker = "o" )
```

```
[77]: [<matplotlib.lines.Line2D at 0x1f0817c2550>]image/png<Figure size 864x432 with 1 Axes>
```

Table 12: Four types of features (21 in total) identified in 2019

Name	h_mean	s_mean	v_mean	h_std	s_std	v_std	h_skew
Baoshan	0.231070	0.035865	0.638941	0.297847	0.052048	0.180189	0.336779
Changning	0.231849	0.031294	0.649237	0.304689	0.048471	0.167002	0.344394
Chongming	0.157306	0.016473	0.742402	0.282495	0.031843	0.101771	0.345289
Fengxian	0.295539	0.086197	0.605431	0.302731	0.097695	0.243596	0.329388
Hongkou	0.239944	0.037582	0.638613	0.303830	0.055047	0.182938	0.339615

Name	s_skew	v_skew	LBP0	LBP1	LBP2	LBP3	LBP4
Baoshan	0.067866	0.195591	0.043059	0.047740	0.040768	0.058617	0.068075
Changning	0.062974	0.182483	0.042641	0.050118	0.037668	0.051370	0.058422
Chongming	0.043360	0.119800	0.025468	0.029637	0.023332	0.035217	0.047621
Fengxian	0.115702	0.257234	0.051206	0.061050	0.052882	0.081410	0.097157
Hongkou	0.070552	0.197624	0.047032	0.054172	0.042940	0.054392	0.055014

Name	LBP5	LBP6	LBP7	LBP8	LBP9	veg_mean	builtup_mean
Baoshan	0.056809	0.037410	0.045565	0.524189	0.077767	-0.001801	0.001938
Changning	0.061528	0.035235	0.047082	0.539452	0.076482	-0.001515	0.000774
Chongming	0.036412	0.023590	0.028650	0.704424	0.045649	-0.000705	0.000318
Fengxian	0.079233	0.051521	0.060437	0.369931	0.095174	-0.008185	-0.000408
Hongkou	0.068051	0.038789	0.048937	0.507320	0.083353	-0.002057	-0.000277

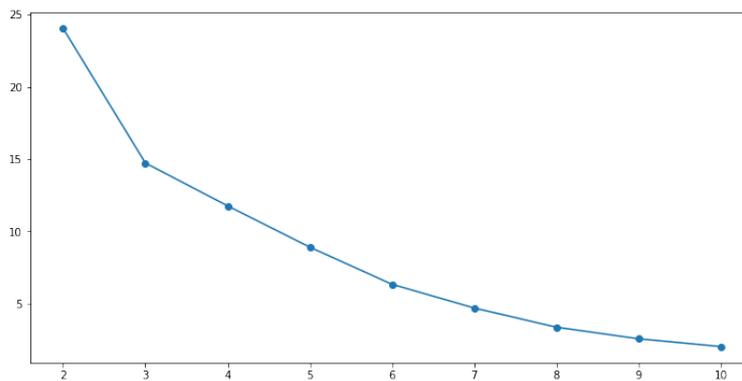


Figure 6: Elbow method to determine the optimal k for k-mean clustering

Figure 6 indicates that 2 and 6 (i.e. knee in the plot) can be the optimal numbers of k clusters for the features extracted from both years of satellite imagery. Considering the context of the paper, the number of 6 is finally assigned to k to fit the kmeans clustering model, varying labels are subsequently matched to features dataset.

```
[78]: np.random.seed(0)
      k = 6
      cls = pd.Series(KMeans(n_clusters=k, max_iter = 1000, n_init = 1000,
                           random_state = 24).fit_predict(merged_var))
```

After implementing k-means clustering on our constructed dataset, the label of each cluster is assigned to the last columns of data for further interpretation (as shown in Table 14).

```
[79]: # Assign the each cluster number to the merged data
      merged_var = merged_var.assign(lbls=cls)
      merged_var.index = features_old_var.index
      # last columns represent class labels
      merged_var.head().style.set_caption('Table 14: Assign cluster number to each
      ... administrative area')
```

```
[79]: text/html<pandas.io.formats.style.Styler at 0x1f081e58550>
```

Table 13: Integrated preprocessed features identified in 1984 and 2019 seperately

	h_mean_x	s_mean_x	v_mean_x	h_std_x	s_std_x	v_std_x	h_skew_x	s_skew_x	v_skew_x	LBP0_x	LBP1_x	LBP2_x	LBP3_x	LBP4_x
0	0.636975	0.359694	0.281144	1.000000	0.422465	0.580121	1.000000	0.417053	0.568198	0.286043	0.328943	0.588494	0.717523	0.837853
1	0.363843	0.354334	0.396450	0.295018	0.453664	0.504220	0.000000	0.443305	0.483028	0.316677	0.441716	0.627346	0.537772	0.441710
2	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.081547	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
3	1.000000	0.990530	0.000000	0.906183	1.000000	0.987872	0.633860	1.000000	0.987806	0.915923	1.000000	0.915397	0.792921	0.738052
4	0.515153	0.480226	0.329302	0.685631	0.594329	0.597572	0.662480	0.590183	0.573137	0.499704	0.564467	0.676601	0.528034	0.370871

	LBP5_x	LBP6_x	LBP7_x	LBP8_x	LBP9_x	veg_mean_x	builtup_mean_x	h_mean_x	s_mean_x	v_mean_x	h_std_y	s_std_y	v_std_y	h_skew_y
0	0.613055	0.564716	0.332298	0.714937	0.032829	0.647652	0.766995	0.442195	0.278123	0.301061	0.463123	0.308221	0.547531	0.504942
1	0.551933	0.467707	0.395850	0.744082	0.042063	0.654899	0.764480	0.446860	0.212560	0.370616	0.659603	0.254002	0.455458	0.773771
2	0.000000	0.000000	0.000000	1.000000	0.000000	0.790941	0.730085	0.000000	0.000000	1.000000	0.022307	0.002006	0.000000	0.805338
3	0.741784	1.000000	1.000000	0.507310	0.110712	0.142243	0.844884	0.828667	1.000000	0.074677	0.603370	1.000000	0.990252	0.244056
4	0.677669	0.542270	0.482237	0.700451	0.059718	0.401200	0.686044	0.495392	0.302755	0.298842	0.634913	0.353666	0.566728	0.605083

	s_skew_y	v_skew_y	LBP0_y	LBP1_y	LBP2_y	LBP3_y	LBP4_y	LBP5_y	LBP6_y	LBP7_y	LBP8_y	LBP9_y	veg_mean_y	builtup_mean_y
0	0.338747	0.541837	0.526221	0.472274	0.548142	0.506561	0.410501	0.406220	0.494801	0.499487	0.726866	0.060848	0.856451	1.000000
1	0.271126	0.448129	0.513894	0.534312	0.450686	0.349686	0.216767	0.500209	0.416926	0.544291	0.749995	0.058415	0.893850	0.699181
2	0.000000	0.000000	0.007423	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000	1.000000	0.581435
3	1.000000	0.982533	0.766492	0.819500	0.928936	1.000000	0.994141	0.852811	1.000000	0.938629	0.493097	0.093827	0.020274	0.393903
4	0.375880	0.556374	0.643383	0.640062	0.616412	0.415111	0.148374	0.630101	0.544153	0.599054	0.701302	0.071432	0.822844	0.427538

## 5 Interpretation

To understand the analysis result, the mean of each feature across each cluster can be calculated to uncover the feature differences among clusters. A categorical bar-plot shown below presents how the average of all features changed between 1984 and 2019. Besides, a choropleth map is created to visualise the spatial distribution of categories/clusters by varying colours.

```
[80]: # calculate the mean of features for each class
k6_mean = merged_var.groupby('lbls').mean()
k6_mean.style.set_caption('Table 15: Mean values of each feature at each cluster for
... different years')
```

```
[80]: text/html<pandas.io.formats.style.Styler at 0x1f081654b00>
```

Table 15 displays the mean values of all features in two years at varying groups. For more interpretability, a few data munging steps are required to generate visual representations.

```
[81]: # Rearrange our data in a way that every row is one feature in a class
k6_mean = k6_mean.stack()
k6_mean.head()
```

```
[81]: lbls
0  h_mean_x    0.803863
   s_mean_x    0.749195
   v_mean_x    0.146109
   h_std_x     0.895068
   s_std_x     0.749907
dtype: float64
```

```
[82]: # convert multi-indices into single index

k6_mean = k6_mean.reset_index()
# renmae the columns
k6_mean = k6_mean.rename(columns = {'lbls': 'Class', 'level_1': 'Features', 0: 'Values'})
# rename feature names in Feature column
old = k6_mean.loc[k6_mean['Features'].str.contains('x') == True, :]
new = k6_mean.loc[k6_mean['Features'].str.contains('y') == True, :]
# add a new column to represent time
old = old.assign(Time = 1984)
new = new.assign(Time = 2019)
# remove '_x' and '_y' in the table to make feature names for both years are the same
old['Features'] = old['Features'].str.replace('_x', '')
new['Features'] = new['Features'].str.replace('_y', '')
```

```
[83]: # create a new dataframe to store the mean of each feature each cluster with time

data = pd.concat([old,new])
data.head().style.set_caption('Table 16: Tidy table represents mean values of features
... for each cluster at different years')
```

```
[83]: text/html<pandas.io.formats.style.Styler at 0x1f08cee1d68>
```

Table 16 reveals different categorical information, with each row represents the number of class, the feature name, the mean value of the feature and the year when the feature is extracted. We can then visualise this table in the bar-plot in Figure 7 to understand the pattern from image features.

```
[84]: # visualise the distribution of mean values by features, class and time

g = sns.catplot( data = data, x = 'Features', y = 'Values', row = 'Class',
                hue = 'Time', kind = 'bar', aspect = 5, height = 3, palette = 'Accent')
g.fig.suptitle('Fig.7: Visual representation of patterns extracted from k-mean
... clustering', y = -0.1, fontsize = 18)
```

```
[84]: Text(0.5, -0.1, 'Fig.7: Visual representation of patterns extracted from k-mean
clustering')
image/png<Figure size 1141.5x1296 with 6 Axes>
```

Table 14: Assign cluster number to each administrative area

Name	h_mean_x	s_mean_x	v_mean_x	h_std_x	s_std_x	v_std_x	h_skew_x	s_skew_x	v_skew_x	LBP0_x	LBP1_x	LBP2_x	LBP3_x	LBP4_x
Baoshan	0.636975	0.359694	0.281144	1.000000	0.422465	0.580121	1.000000	0.417053	0.568198	0.286043	0.328943	0.588494	0.717523	0.837853
Changning	0.363843	0.354334	0.396450	0.295018	0.453664	0.504220	0.000000	0.443305	0.483028	0.316677	0.441716	0.627346	0.537772	0.441710
Chongming	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.081547	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
Fengxian	1.000000	0.990530	0.000000	0.906183	1.000000	0.987872	0.633860	1.000000	0.987806	0.915923	1.000000	0.915397	0.792921	0.738052
Hongkou	0.515153	0.480226	0.329302	0.685631	0.594329	0.597572	0.662480	0.590183	0.573137	0.499704	0.564467	0.676601	0.528034	0.370871

Name	LBP5_x	LBP6_x	LBP7_x	LBP8_x	LBP9_x	veg_mean_x	builtup_mean_x	h_mean_x	s_mean_x	v_mean_x	h_std_y	s_std_y	v_std_y	h_skew_y
Baoshan	0.613055	0.564716	0.332298	0.714937	0.032829	0.647652	0.766995	0.442195	0.278123	0.301061	0.463123	0.308221	0.547531	0.504942
Changning	0.551933	0.467707	0.395850	0.744082	0.042063	0.654899	0.764480	0.446860	0.212560	0.370616	0.659603	0.254002	0.455458	0.773771
Chongming	0.000000	0.000000	0.000000	1.000000	0.000000	0.790941	0.730085	0.000000	0.000000	1.000000	0.022307	0.002006	0.000000	0.805338
Fengxian	0.741784	1.000000	1.000000	0.507310	0.110712	0.142243	0.844884	0.828667	1.000000	0.074677	0.603370	1.000000	0.990252	0.244056
Hongkou	0.677669	0.542270	0.482237	0.700451	0.059718	0.401200	0.686044	0.495392	0.302755	0.298842	0.634913	0.353666	0.566728	0.605083

Name	s_skew_y	v_skew_y	LBP0_y	LBP1_y	LBP2_y	LBP3_y	LBP4_y	LBP5_y	LBP6_y	LBP7_y	LBP8_y	LBP9_y	veg_mean_y	builtup_mean_y	lbls
Baoshan	0.338747	0.541837	0.526221	0.472274	0.548142	0.506561	0.410501	0.406220	0.494801	0.499487	0.726866	0.060848	0.856451	1.000000	1
Changning	0.271126	0.448129	0.513894	0.534312	0.450686	0.349686	0.216767	0.500209	0.416926	0.544291	0.749995	0.058415	0.893850	0.699181	1
Chongming	0.000000	0.000000	0.007423	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000	1.000000	0.581435	2
Fengxian	1.000000	0.982533	0.766492	0.819500	0.928936	1.000000	0.994141	0.852811	1.000000	0.938629	0.493097	0.093827	0.020274	0.393903	3
Hongkou	0.375880	0.556374	0.643383	0.640062	0.616412	0.415111	0.148374	0.630101	0.544153	0.599054	0.701302	0.071432	0.822844	0.427538	1

Table 15: Mean values of each feature at each cluster for different years

0	0.803863	0.749195	0.146109	0.895068	h_std_x	s_std_x	v_std_x	h_skew_x	s_skew_x	v_skew_x	LBP0_x	LBP1_x	LBP2_x	LBP3_x	LBP4_x
1	0.426195	0.381664	0.384955	0.543494	0.481686	0.523486	0.450601	0.478586	0.504149	0.382849	0.444543	0.574545	0.510795	0.408613	0.408613
2	0.113097	0.088470	0.837305	0.210939	0.118478	0.127151	0.316506	0.107534	0.124469	0.070301	0.089850	0.162624	0.139945	0.089317	0.089317
3	0.949856	0.995265	0.004559	0.876415	0.991048	0.993936	0.676387	0.975688	0.993903	0.882717	0.957562	0.957699	0.878867	0.798240	0.798240
4	0.721478	0.979475	0.114017	0.567528	0.998762	0.933619	0.244557	0.971796	0.925791	1.000000	0.990273	0.845726	0.668671	0.407605	0.407605
5	0.471719	0.427622	0.373358	0.634773	0.554207	0.559062	0.554391	0.560426	0.426212	0.502654	0.647752	0.516126	0.345374	0.345374	0.345374
lbls	h_mean_x	s_mean_x	v_mean_x	h_std_x	s_std_x	v_std_x	h_skew_x	s_skew_x	v_skew_x	LBP0_x	LBP1_x	LBP2_x	LBP3_x	LBP4_x	
0	0.714188	0.767025	0.567415	0.603671	0.062757	0.336476	0.964021	0.782732	0.400459	0.582989	0.125378	0.621145	0.621750	0.781063	0.466975
1	0.491368	0.480169	0.415924	0.747533	0.045784	0.579325	0.782732	0.400459	0.233562	0.370746	0.449325	0.279210	0.479726	0.520610	0.520610
2	0.114726	0.140034	0.083018	0.941070	0.009137	0.739326	0.768474	0.081734	0.042801	0.835523	0.144022	0.062326	0.114256	0.757332	0.757332
3	0.870892	0.996516	0.943707	0.495601	0.106946	0.071121	0.808680	0.914333	0.849361	0.087764	0.539909	0.845232	0.995126	0.122028	0.122028
4	0.847049	0.819912	0.961713	0.538442	0.114071	1.000000	0.000000	0.883471	0.529215	0.000000	1.000000	0.546832	0.871779	0.698488	0.698488
5	0.299768	0.528165	0.568006	0.000000	1.000000	0.549489	0.733977	0.428816	0.207932	0.329358	0.435917	0.246778	0.505767	0.409755	0.409755
lbls	LBP5_x	LBP6_x	LBP7_x	LBP8_x	LBP9_x	veg_mean_x	builtup_mean_x	h_mean_x	h_mean_y	s_mean_x	s_mean_y	v_mean_x	v_mean_y	h_std_x	h_std_y
0	0.714188	0.767025	0.567415	0.603671	0.062757	0.336476	0.964021	0.621423	0.582989	0.125378	0.621145	0.621750	0.781063	0.466975	0.466975
1	0.491368	0.480169	0.415924	0.747533	0.045784	0.579325	0.782732	0.400459	0.233562	0.370746	0.449325	0.279210	0.479726	0.520610	0.520610
2	0.114726	0.140034	0.083018	0.941070	0.009137	0.739326	0.768474	0.081734	0.042801	0.835523	0.144022	0.062326	0.114256	0.757332	0.757332
3	0.870892	0.996516	0.943707	0.495601	0.106946	0.071121	0.808680	0.914333	0.849361	0.087764	0.539909	0.845232	0.995126	0.122028	0.122028
4	0.847049	0.819912	0.961713	0.538442	0.114071	1.000000	0.000000	0.883471	0.529215	0.000000	1.000000	0.546832	0.871779	0.698488	0.698488
5	0.299768	0.528165	0.568006	0.000000	1.000000	0.549489	0.733977	0.428816	0.207932	0.329358	0.435917	0.246778	0.505767	0.409755	0.409755
lbls	s_skew_y	v_skew_y	LBP0_y	LBP1_y	LBP2_y	LBP3_y	LBP4_y	LBP5_y	LBP6_y	LBP7_y	LBP8_y	LBP9_y	veg_mean_y	builtup_mean_y	
0	0.629142	0.768836	0.559714	0.585206	0.760955	0.841492	0.838157	0.722830	0.794093	0.652172	0.600350	0.070401	0.714016	0.835832	0.835832
1	0.300671	0.470766	0.465958	0.451277	0.510348	0.446913	0.309060	0.470514	0.466278	0.466282	0.748782	0.053443	0.852035	0.622854	0.622854
2	0.077196	0.114218	0.050600	0.054194	0.140384	0.172060	0.158486	0.147885	0.115339	0.054215	0.939642	0.006619	0.959421	0.572280	0.572280
3	0.870699	0.991266	0.883246	0.909750	0.964468	0.940470	0.808686	0.863404	0.987573	0.969315	0.488089	0.104292	0.010137	0.196952	0.196952
4	0.575903	0.866917	0.981285	0.905575	0.907557	0.660840	0.487015	1.000000	0.950033	0.913027	0.521429	0.107485	0.713309	0.986139	0.986139
5	0.280930	0.499205	0.621626	0.525659	0.570189	0.416572	0.116072	0.221324	0.445842	0.718360	0.000000	1.000000	0.895112	0.568988	0.568988



Figure 7: Visual representation of patterns extracted from k-mean clustering

```
[85]: # plot clustering results for two different years
f, ax = plt.subplots(1, figsize=(10, 12))
# plot cluster results
poly = poly.drop('coords', axis = 1)
poly.assign(lbls=cls)\
    .plot(column='lbls', categorical=True, linewidth=1, alpha=0.5, ax=ax,
          legend = True, cmap = 'Accent', edgecolor = 'black')
# add labels for geographical units
poly['coords']=poly['geometry'].apply(lambda x:x.representative_point().coords[:])
poly['coords']=[coords[0] for coords in poly['coords']]
for idx, row in poly.iterrows():
    ax.annotate(text=row['Name'],xy=row['coords'],va='center',ha='center',
               alpha = 0.8, fontsize = 10)
plt.title('Fig.8: Spatial distribution of classification results', y=-0.01)
# remove axes and set aspect ratio so that the data units are the same in every direction
ax.axis('off')
ax.axis('equal')
```

```
[85]: (290053.0696196473, 407301.6741094636, 3389866.639388826, 3533566.430983904)
image/png<Figure size 720x864 with 1 Axes>
```

Table 16: Tidy table represents mean values of features for each cluster at different years

	Class	Features	Values	Time
0	0	h_mean	0.803863	1984
1	0	s_mean	0.749195	1984
2	0	v_mean	0.146109	1984
3	0	h_std	0.895068	1984
4	0	s_std	0.749907	1984

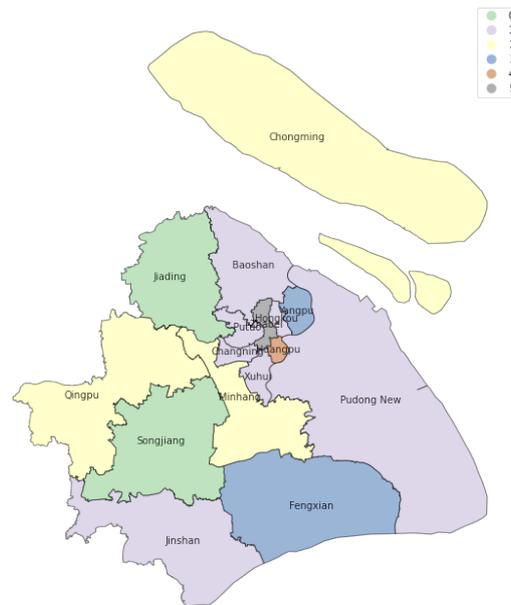


Figure 8: Spatial distribution of classification results

From Figures 7 and 8 we can see a few striking differences across clusters, or classes. For class 4, only one administrative area (i.e. Huangpu area) is grouped, displayed in the middle of north-east areas. The mean values for this class are mostly high in both years except a couple of features such as v\_mean, LBP4 and LBP9 features. The brightness (v\_mean) for this area is highly low and it became completely black over time. H\_mean value is high in both years, demonstrating that the dominating colour is blue, which represent water. This corresponds to the famous area of The Bund, with its river skyline, which is part of this polygon. The vegetation built-up features indicate that this area has experienced a remarkable change, from more vegetation and few buildings to less vegetation and completely constructed/urbanisation.

Class 0 and Class 1 are relatively consistent compared to other classes, implying that the urban areas in purple and green colours almost remained unchanged during the past 35 years. Besides, these two classes have similar transformation such as more vegetation coverage and less buildings for the current year of 2019. However, Class 0 has more brightness and more green colour based on v\_mean, h\_mean and veg\_mean features, and Class 1 has higher h\_mean, h\_std, h\_skew and built-up\_mean, implying these two areas have water covered and were highly constructed.

Class 2 distributed at north and middle-west areas in the map, which is extremely diverse and unique among all categories. It has the highest brightness features and LBP8 texture features, while the rest mean values of colour and texture features are highly low, especially for LBP9 where almost zero values in both years. The values for h\_mean, s\_mean and v\_mean display that the primary colour for these areas is red with little grey and much brightness, representing that these areas include more bare ground or soil and thus probably rural areas. Adversely, Class 5 has zero values for LBP8 but highest

values for LBP9 in both years. It contains only one administrative area (i.e. Zhabei area), surrounded by Class 4 and Class 0. Similarly, the area in Class 5 has more vegetation but slightly less built-up areas over the past years. Class 3 contains two areas distributed at the south and surrounded by Class 1 from the map. The feature values in Class 3 are mostly extremely high, while the `veg_mean` and `built-up_mean` for current year are the least, thus indicating that these areas have more water over the time.

## 6 Conclusion

Urbanisation has significantly changed the interaction between humans and the surrounding environment, which poses new challenges in a multitude of fields including construction and city planning, hazard mitigation or disease control. It is essential to quantify and assess urbanisation over time to enable policy makers and planners to make informed decisions about future urban changes. The sustainability of urban spaces will become particularly important in the light of future climate change. Satellite imagery could play a vital role in assessing cities for their livability by i.e. quantifying the greenspace to built environment ratio. This notebook shows the potential of open source satellite imagery to exploring urban changes and proposes a simple method framework for automatic data collection and features extraction to determine urbanisation over time using Python as a tool.

## References

- Barsi JA, Lee K, Kvaran G, Markham BL, Pedelty JA (2014a) The spectral response of the Landsat-8 operational land imager. *Remote Sensing* 6: 10232–10251. [CrossRef](#).
- Barsi JA, Schott JR, Hook SJ, Raqueno NG, Markham BL, Radocinski RG (2014b) Landsat-8 thermal infrared sensor (TIRS) vicarious radiometric calibration. *Remote Sensing* 6: 11607–11626. [CrossRef](#).
- Burchfield M, Overman HG, Puga D, Turner MA (2006) Causes of sprawl: A portrait from space. *The Quarterly Journal of Economics* 121: 587–633. [CrossRef](#).
- Giada S, De Groeve T, Ehrlich D, Soille P (2003) Information extraction from very high resolution satellite imagery over Lukole refugee camp, Tanzania. *International Journal of Remote Sensing* 24: 4251–4266. [CrossRef](#).
- Glaeser E, Henderson JV (2017) Urban economics for the developing world: An introduction. *Journal of Urban Economics* 98: 1–5. [CrossRef](#).
- Ibrahim MR, Haworth J, Cheng T (2020) Understanding cities with machine eyes: A review of deep computer vision in urban analytics. *Cities* 96. [CrossRef](#).
- Keen N (2005) Color moments. School of informatics, University of Edinburgh
- Kit O, Lüdeke M (2013) Automated detection of slum area change in Hyderabad, India using multitemporal satellite imagery. *Journal of Photogrammetry and Remote Sensing* 83: 130–137. [CrossRef](#).
- Knight EJ, Kvaran G (2014) Landsat-8 operational land imager design, characterization and performance. *Remote Sensing* 6: 10286–10305. [CrossRef](#).
- Kohli D, Sliuzas R, Stein A (2016) Urban slum detection using texture and spatial metrics derived from satellite imagery. *Journal of Spatial Science* 61: 405–426. [CrossRef](#).
- Ministry of Civil Affairs of the People’s Republic of China (2018) Change of administrative divisions at or above the county level. Available at: <http://202.108.98.30/description?dcpid=1> [Accessed 10 Oct. 2019]
- NASA (2000) Normalized difference vegetation index (NDVI). Available at: [https://earth-observatory.nasa.gov/features/MeasuringVegetation/measuring\\_vegetation\\_2.php](https://earth-observatory.nasa.gov/features/MeasuringVegetation/measuring_vegetation_2.php) [Accessed 30 Oct. 2019]

- NASA (2019) Landsat science. Available at: <https://landsat.gsfc.nasa.gov/> [Accessed 10 Sep. 2019]
- Ojala T, Pietikäinen M, Harwood D (1996) A comparative study of texture measures with classification based on feature distributions. *Pattern Recognition* 19: 51–59. [CrossRef](#).
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J (2011) Scikit-learn: Machine learning in Python. *Journal of machine learning research* 12: 2825–2830
- Roy DP, Wulder MA, Loveland TR, Woodcock CE, Allen RG, Anderson MC, Helder D, Irons JR, Johnson DM, Kennedy R, Scambos TA, Schaaf CB, Schott JR, Sheng Y, Vermote EF, Belward AS, Bindschadler R, Cohen WB, Gao F, Hipple JD, Hostert P, Huntington J, Justice CO, Kilic A, Kovalsky V, Lee ZP, Lymburner L, Masek JG, McCorkel J, Shuai Y, Trezza R, J. Vogelmann J, R.H. Wynne RH, Zhu Z (2014) Landsat-8: Science and product vision for terrestrial global change research. *Remote sensing of Environment* 145: 154–172. [CrossRef](#).
- United Nations (2019) World urbanization prospects 2018: Highlights. United Nations, Department of Economic and Social Affairs, Population Division (ST/ESA/SER.A/421)
- Zha Y, Gao J, Ni S (2003) Use of normalized difference built-up index in automatically mapping urban areas from TM imagery. *International Journal of Remote Sensing* 24: 583–594. [CrossRef](#).

