# Data Wrangle OpenstreetMaps Data

**Lasso.py**                                                                                 ③

```python
1  '''
2      The lasso module contains all the functions used in data wrangling the Open Street Map Da
3  '''
4
5  import xml.etree.cElementTree as ET
6  import re
7  import codecs
8  import json
9  from collections import defaultdict
10
11
12 #  Regex
```

`AWESOME`

Outstanding commenting etiquette!

```python
13 RE_PROBLEM_CHARS = re.compile(r'[=\+/&<>;\'"\?%#$@\,\. \t\r\n]')
14 RE_POSTAL_CODE = re.compile(r"^([a-zA-Z]\d[a-zA-Z] ?\d[a-zA-Z]\d)$")
15 # http://stackoverflow.com/questions/16614648/canadian-postal-code-regex
16
17
18 # Default Dicts used in audit functions
```

`AWESOME`

You have avoided in-line comments which often extend the character-per-line preference of 80. Wel
done!

```python
19 def def_dict_2():
20     '''
21     defaultdict two dict deep with default of set
22     '''
23     return defaultdict(lambda: defaultdict(set))
24
25
26 def def_dict_3():
27     '''
28     defaultdict three dict deep with default of set
29     '''
30     return defaultdict(lambda: defaultdict(lambda: defaultdict(set)))
31
32
33 #######################
34 ### Auditing the data ###
35 #######################
```

```python
36
37
38  def dictify_element_and_children(element, atr_d=def_dict_2(), st_atr_d=def_dict_3(), s_st_d=d
```

```python
39          '''
40          From each element in the xml tree creats/adds summary dictionaries to better understand t
41
42          return:
43              1 - attrib_dict (atr_d): should return all potential attributes with a set of all ans
44              2 - sub_tag_attrib_dict (st_atr_d): return sub_tag: sub_tag.attrib.keys()
45              3 - sub_subtag_dict (s_st_d): return sub_tag: sub_tag.children
46              4 - tag_k_v_dict: for tag sub_tag:
47          '''
48          for key, val in element.attrib.items():
49              atr_d[element.tag][key].add(val)
50          for sub_tag in element.iter():
51              child_set = {el.tag for el in list(sub_tag)}
52              if child_set != set():
53                  s_st_d[element.tag][sub_tag.tag].update(child_set)
54              for key, val in sub_tag.attrib.items():
55                  st_atr_d[element.tag][sub_tag.tag][key].add(val)
56              if sub_tag.tag == 'tag':
57                  tag_k_v_dict[element.tag][sub_tag.attrib['k']].add(sub_tag.attrib['v'])
58
59          return atr_d, st_atr_d, s_st_d, tag_k_v_dict
60
61
62  def summarizes_data_2_tags_deep(filename):
63      '''
64      uses dictify_element_and_children to loop through entire xml tree creating summary data.
65      '''
66      atr_d = def_dict_2()
67      st_atr_d = def_dict_3()
68      s_st_d = def_dict_2()
69      tag_k_v_dict = def_dict_2()
70
71      for _, element in ET.iterparse(filename):
72          dictify_element_and_children(element, atr_d, st_atr_d, s_st_d, tag_k_v_dict)
73      return atr_d, st_atr_d, s_st_d, tag_k_v_dict
74
75
76  def check_keys_list(dict_key_list):
77      '''
78      checks a list of dictionary keys for problem characters
79      '''
80      problem_keys = []
81      for key in dict_key_list:
82          if RE_PROBLEM_CHARS.search(key):
83              problem_keys.append(key)
84      return problem_keys
85
86
87  def process_audit_address_type(tag_k_v_dict, directions=()):
88      '''
89      loops though all street addesses putting all the street types in a set
90      if the last word in the steet addess is a direction (E, W, N, S)
```

```python
 91            it uses the second last word in the street address
 92        if not
 93            it uses the last word in the street address
 94        '''
 95    street_types = set()
 96    street_list = wrap_up_tag_k_v_dict(tag_k_v_dict, 'addr:street')
 97
 98    for val in list(street_list):
 99        street_name = val
100        street_split = street_name.split()
101        if street_split[-1] in directions:
102            street_types.add(street_split[-2])
103        else:
104            street_types.add(street_split[-1])
105
106    return street_types
107
108
109 def wrap_up_tag_k_v_dict(tag_k_v_dict, key):
110    '''
111    used to look at the key value pairs of nodes, ways, and relations together
112    '''
113    return tag_k_v_dict['node'][key] | tag_k_v_dict['relation'][key] | tag_k_v_dict['way'][ke
114
115
116 ############################
117 ### Load Data into MongoDB ###
118 ############################
119
120
121 # Variable maps to swap out non normal values for normal values
122 # street direction map
123 ST_DIR_MAP = {'S': 'South',
124               's': 'South',
125               'South': 'South',
126               'E': 'East',
127               'e': 'East',
128               'East': 'East',
129               'W': 'West',
130               'w': 'West',
131               'West': 'West',
132               'N': 'North',
133               'n': 'North',
134               'North': 'North'}
135
136 # Street type map
137 ST_TYPE_MAP = {'AVenue': 'Avenue',
138               'Ave': 'Avenue',
139               'Crescent': 'Cresent',
140               'Dr': 'Drive',
141               'Dr.': 'Drive',
142               'Rd': 'Road',
143               'St': 'Street',
144               'St.': 'Street',
145               'Steet': 'Street'}
146
147 # Province Map
148 PROV_MAP = {'ON': 'ON',
149            'Ontario': 'ON',
150            'on': 'ON',
151            'ontario': 'ON'}
152
153 # City Map
```

```python
153  # City Map
154  CITY_MAP = {'City of Cambridge': 'Cambridge',
155              'City of Kitchener': 'Kitchener',
156              'kitchener': 'Kitchener',
157              'City of Waterloo': 'Waterloo',
158              'waterloo': 'Waterloo',
159              'St. Agatha': 'Saint Agatha'}
160
161
162  def map_subin(val, val_map):
163      '''
164      Subs in a value from a map given the map and a value
165      '''
166      # returns the maped value if it's in the map
167      if val in val_map.keys():
168          return val_map[val]
169      else:
170          # returns the original value if it's not in the map
171          return val
172
173
174  def update_street(street):
175      '''
176      uses the map_subin function to subin corrected street types and street directions
177      '''
178      # split the street address into a list of words
179      st_list = street.split()
180
181      if st_list[-1] in ST_DIR_MAP.keys():
182          # if the last word is a direction sub both direction(-1) & type(-2)
183          st_list[-1] = map_subin(st_list[-1], ST_DIR_MAP)
184          st_list[-2] = map_subin(st_list[-2], ST_TYPE_MAP)
185      else:
186          # otherwise sub in the street type
187          st_list[-1] = map_subin(st_list[-1], ST_TYPE_MAP)
188
189      return ' '.join(st_list)
190
191
192  def update_address(key, val, addr_dict):
193      if key == 'addr:street':
194          addr_dict[key[5:]] = update_street(val)
195      elif key == 'addr:state':
196          if not addr_dict.get('province'):
197              addr_dict['province'] = map_subin(val, PROV_MAP)
198      elif key == 'addr:province':
199          addr_dict[key[5:]] = map_subin(val, PROV_MAP)
200      elif key == 'addr:city':
201          addr_dict[key[5:]] = map_subin(val, CITY_MAP)
202      else:
203          addr_dict[key[5:]] = val
204      return addr_dict
205
206
207  def tag_subtag_process(sub_tag, address, tags):
208      key = sub_tag.attrib['k']
209      val = sub_tag.attrib['v']
210
211      # addr: tags are sent to update_address function
212      if key[0:5] == 'addr:':
213          address = update_address(key, val, address)
214
215      # merge 'fixme' and 'FIXME' into 'FIXME'
```

```python
216        elif key in ['fixme', 'FIXME']:
217            if tags.get('FIXME'):
218                tags['FIXME'] += '\nFIXME: ' + val
219            else:
220                tags['FIXME'] = val
221
222        # all other tag tags get added as k:v pairs
223        else:
224            tags[key] = val
225
226    return address, tags
227
228
229 def subtag_process(xml_tree):
230    '''
231    adds sub tags of an osm xml element to lists and dicts for easy joining to the JSON struc
232    '''
233    # dicts and lists for constucted values
234    node_refs = []
235    members = []
236    address = {}
237    tags = {}
238
239    # looping though each sub tag of xml_tree
240    for sub_tag in xml_tree.iter():
241
242        # sub tag of 'tag' sent to tag function
243        if sub_tag.tag == 'tag':
244            address, tags = tag_subtag_process(sub_tag, address, tags)
245
246        # sub tag of 'nd' appended in order to list
247        elif sub_tag.tag == 'nd':
248            node_refs.append(int(sub_tag.attrib['ref']))
249
250        # sub tag of 'member' appended in order as a list of dicts
251        elif sub_tag.tag == 'member':
252            mem = {}
253            for key, val in sub_tag.attrib.items():
254                if val:
255                    if key == 'ref':
256                        mem[key] = int(val)
257                    else:
258                        mem[key] = val
259            members.append(mem)
260
261    return node_refs, members, address, tags
262
263
264 def shape_xml_tree(xml_tree):
265    '''
266    takes an xml element (node, way, or relation) and converts it into a json element includi
267    '''
268    # returns an empty element if the xml_tree is not a node, way or relation
269    if xml_tree.tag not in ['node', 'way', 'relation']:
270        return {}
271
272    # This is the element we will return at the end
273    element = {}
274
275    # building out the xml_tree attributes
276    element['type'] = xml_tree.tag
277    element['id'] = int(xml_tree.attrib.get('id'))
```

```python
278
279         # location info from the start tag is converted into a list of two floats for easy coordi
280         if xml_tree.tag == 'node':
281             pos = [float(xml_tree.attrib.get('lat')), float(xml_tree.attrib.get('lon'))]
282             element['pos'] = pos
283
284         # creation info is saved in a dictionary under the creation key
285         element['created'] = {}
286         for key, val in xml_tree.attrib.items():
287             if key in ["uid", "version", "changeset"]:
288                 element['created'][key] = int(val)
289             if key in ["user", "timestamp"]:
290                 element['created'][key] = val
291
292         # sub tags are processed in the subtag_process function
293         node_refs, members, address, tags = subtag_process(xml_tree)
294
295         # append all the subtag values
296         if node_refs:
297             element['nd'] = node_refs
298         if members:
299             element['member'] = members
300         if address:
301             element['addr'] = address
302         if tags:
303             element['tag'] = tags
304
305         return element
306
307
308 def process_map(file_in, pretty=False):
309     '''
310     takes xml file with tag 'node', 'way', or 'relation'
311
312     Unpackes into json compatible dict and list structure.
313     saves the json to file for easy import into MongoDB
314
315     {'type':    xml_tree.tag,
316
317      'id':      int(xml_tree('id')),
318
319      'pos':     [float(xml_tree('lat')),
320                  float(xml_tree('lon'))],
321
322      'created': {'version':     int(xml_tree('uid')),
323                  'changeset':   int(xml_tree('changeset')),
324                  'timestamp':   xml_tree('timestamp'),
325                  'user':        xml_tree('user'),
326                  'uid':         int(xml_tree('uid'))},
327
328      'address': {'housenumber': tag_tag['addr:housenumber'],
329                  'postcode': tag_tag['addr:postcode'],
330                  'street': tag_tag['addr:street'], ...},
331
332      'member':  [{'type': member_tag('type'),
333                  'ref':  int(member_tag('ref')),
334                  'role': member_tag('role')},
335                  {.........................}],
336
337      'node_refs':[int(nd_tag['ref']),
338                  int(nd_tag['ref']), ... ],
339
340      'tag': {tag['k']:  tag_tag['v'],
```

```python
341                tag['k']:  tag_tag['v'],
342                ... }
343        }
344      '''
345
346      file_out = "{0}.json".format(file_in)
347      data = []
348      with codecs.open(file_out, "w") as file_out:
349          for _, xml_tree in ET.iterparse(file_in):
350              element = shape_xml_tree(xml_tree)
351              if element:
352                  data.append(element)
353                  if pretty:
354                      file_out.write(json.dumps(element, indent=4)+"\n")
355                  else:
356                      file_out.write(json.dumps(element) + "\n")
357      return data
358
359 if __name__ == '__main__':
360      pass
361
```