# UDACITY

‹ Back to Data Analyst Nanodegree

# Wrangle OpenStreetMap Data

| REVIEW |
|---|
| CODE REVIEW 4 |
| HISTORY |

▼ data.py  3

```python
1  import csv
2  import codecs
3  import pprint
4  import re
5  import xml.etree.cElementTree as ET
6  import codecs
7  import datetime
8  from collections import defaultdict
9
10
11 # Code to more efficiently parse the OSM xml and return a list of XML elements
12 def get_element(osm_file, tags=('node', 'way', 'relation')):
13     #Yield element if it is the right type of tag
14
15     context = ET.iterparse(osm_file, events=('start', 'end'))
16     _, root = next(context)
17     for event, elem in context:
18         if event == 'end' and elem.tag in tags:
19             yield elem
20             root.clear()
21
22 # Dictionary of terms to replace, i.e. "bbq" --> "barbecue"
23 cuisine_dict = {
24     "bbq": "barbecue",
```

AWESOME

Good work finding the inconsistent/incorrect issues with the cuisine names.

```
25      "bar&grill": "grill",
26      "donuts": "donut",
27      "bagels": "bagel",
28      "steak": "steak_house",
29      "burgers": "burger",
30      "salads": "salad",
31      "sandwiches": "sandwich",
32      "sandwhiches": "sandwich",
33      "sandwhich": "sandwich",
34      "hot_dogs": "hot_dog",
35      "chicago_dogs": "hot_dog",
36      "coffee": "coffee_shop",
37      "texmex": "tex_mex",
38      "roast_beel": "roast_beef",
39      "puertorican": "puerto_rican",
40      "jewish_deli": "jewish",
41      "fish_tacos": "tacos",
42      "juice_bar": "juice",
43      "mexcian_food": "mexican",
44      "mexican_food": "mexican",
45      "wrap": "wraps",
46      "french_fries": "fries",
47      "drive_thru_coffee": "coffee_shop",
48      "healthy_cuisine_-_greek": "greek",
49      "tailand": "thai",
50      "thailand": "thai"
51 }
52
53 # These are cuisine terms we will throw out and not allow, based on the large-ish samp
54 # NOTE the some terms - "vegan", "vegetarian", "drive_through" - should be included i
55 #   besides cuisine tags, and the code detects them and adds the appropriate sub-tags
56 disallowed = ["", "food", "vegan", "vegetarian", "vegetarian_and_vegan", "mon", "ret"
57               "lunch", "diner", "dinner", "bar", "pub", "drive_through"]
58
59 # strips the string of whitespace and leading underscore, and converts to lower-case
60 # i.e. "_pizza" --> "pizza";  "Fries " --> "fries"
61 def str_space_undersc_lower(string):
62     if string == "" or string == None:
63         return ""
64     string = string.strip()[1:] if string.strip()[0] == '_' else string.strip()
65     return string.lower()
66
67 # Replaces spaces between cuisine terms with an underscore
68 # i.e. "coffee shop" --> "coffee_shop"
69 def repl_space(string):
70     return "_".join(string.split())
71
72 # splits strings on comma and/or semicolon
73 def spl(string):
74     return re.split("[,;]", string)
75
76 # Replace string with its substitute in the cuisine dictionary, if it is in the dict
77 # If its not in dict and its not in the disallowed list and if it is 3+ characters lo
78 #   then use it as is.
79 def dict_lookup(string):
80     if string in cuisine_dict:
80         return cuisine_dict[string]
82     elif string in disallowed:
```

```python
83          return None
84      elif len(string) < 3:
85          return None
86      else:
87          return string
88
89  # check if have text indicating they are using text for a drive through for a cuisine
90  def has_drive_through_text(text):
91      lower_text = text.lower()
92      return ("drive_through" in lower_text) or ("drivethrough" in lower_text)
93
94  # check if have text for vegan in cuisine text
95  def has_vegan_text(text):
96      lower_text = text.lower()
97      return "vegan" in lower_text
98
99  # check if have text for vegetarian in cuisine text
100 def has_vegetarian_text(text):
101     lower_text = text.lower()
102     return "vegetarian" in lower_text
103
104 # call dict_lookup after cleaning with str_space_undersc_lower and then repl_space
105 def process(string):
106     return dict_lookup(repl_space(str_space_undersc_lower(string)))
107
108 # clean up each element after splitting them, and then rejoin
109 def shape_cuisine_element(cuisine_str):
110     return ";".join(filter(None, map(process, re.split('[,;]', cuisine_str))))
111
112
113 node_attribs = ['id', 'user', 'uid', 'version', 'lat', 'lon', 'timestamp', 'changeset
114 node_tag_keys = ['id', 'key', 'value', 'type']
115 way_attribs = ['changeset', 'id', 'timestamp', 'uid', 'user', 'version']
116
117 OSM_PATH = "phoenix.osm"
118 NODES_PATH = "nodes.csv"
119 NODE_TAGS_PATH = "nodes_tags.csv"
120 WAYS_PATH = "ways.csv"
121 WAY_NODES_PATH = "ways_nodes.csv"
122 WAY_TAGS_PATH = "ways_tags.csv"
123
124 LOWER_COLON = re.compile(r'^([a-z]|_)+:([a-z]|_)+')
125 PROBLEMCHARS = re.compile(r'[=\+/&<>;\'"\?%#$@\,\. \t\r\n]')
126
127 # Make sure the fields order in the csvs matches the column order in the sql table sc
128 NODE_FIELDS = ['id', 'lat', 'lon', 'user', 'uid', 'version', 'changeset', 'timestamp'
129 NODE_TAGS_FIELDS = ['id', 'key', 'value', 'type']
130 WAY_FIELDS = ['id', 'user', 'uid', 'version', 'changeset', 'timestamp']
131 WAY_TAGS_FIELDS = ['id', 'key', 'value', 'type']
132 WAY_NODES_FIELDS = ['id', 'node_id', 'position']
133
134 def get_valid_type_key(thestring):
135     if PROBLEMCHARS.search(thestring):
136         return False, None, None
137     elif ':' in thestring:
138         return True, thestring[: thestring.index(':')], thestring[thestring.index(':'
139     else:
140         return True, "regular", thestring
141
142 # Populate the node attribute dictionary, using the "lat", "lon", "uid", and "changes
143 # Note: we assume valid string attributes if not lat, lon, uid or changeset
```

```python
144 def pop_attrib(tag):
145     thedict = {}
146     for attrib in node_attribs:
147         if attrib == 'lat' or attrib == 'lon':
148             thedict[attrib] = float(tag.get(attrib))
149         elif attrib == 'id' or attrib == 'uid' or attrib == 'changeset':
150             thedict[attrib] = int(tag.get(attrib))
151         else:
152             thedict[attrib] = tag.get(attrib)
153     return thedict
154
155 # populate the way attributes dictionary
156 def get_way_attribs(tag):
157     attrib_dict = {}
158     for attrib in way_attribs:
159         if attrib == 'changeset' or attrib == 'id' or attrib == 'uid':
160             attrib_dict[attrib] = int(tag.get(attrib))
161         else:
162             attrib_dict[attrib] = tag.get(attrib)
163     return attrib_dict
164
165 # populate the way nodes dictionary
166 def pop_way_nodes(tag):
```

**SUGGESTION**

Great job commenting for code that is not intuitively readable. Although it is not required for this project, particularly for arguments. Here is the resource

https://www.python.org/dev/peps/pep-0257/

```python
167     nodes = []
168     tagid = int(tag.get('id'))
169     nd_tags = tag.findall('nd')
170     for index, nd_tag in enumerate(nd_tags):
171         newdict = {}
172         newdict['id'] = tagid
173         newdict['node_id'] = int(nd_tag.get('ref'))
174         newdict['position'] = index
175         nodes.append(newdict)
176     return nodes
177
178 # Populate the drive_through dictionary, using the passed in tag id and tag type
179 # Note this is used when we find that the cuisine tag is indicating this is a drive th
180 def pop_drive_through(tag_id, tag_type):
181     return {
182         "id": tag_id, "key": "drive_through", "value": "yes", "type": tag_type
183     }
184
185 # Populate the diet dictionary, using the passed in information for tag id, tag type,
186 #  that this is indicating a vegan or a veretagrian restaurant (or both)
187 # Note this is used when we find the cuisine tag is indicating this is a vegan and/or
188 def pop_diet(tag_id, tag_type, has_vegan, has_vegetarian):
189
190     if has_vegan and has_vegetarian:
191         text = "vegan;vegetarian"
192     elif has_vegan:
193         text = "vegan"
194     elif has_vegetarian:
```

```python
195                 text = "vegetarian"
196
197         if has_vegan or has_vegetarian:
198             return { "id": tag_id, "key": "diet", "value": text, "type": tag_type }
199         else:
200             return None
201
202 # Populate the tags dictionary based on the keys and values in the tag's subtags.
203 # If this is a cuisine tag, we do some extra checking and cleaning up.
204 # In particular, if the cuisine tag has indications that this is a vegan or vegetaria
205 # restaurant, then we add a dictionary for a "diet" subtag.  If it indicates a drive
206 # restaurant, then we add a dictionary for a "drive_through" subtag.
207 # We also do the cuisine tag cleanups based on formatting (stripping underscores, whi
208 #  and a dictionary lookup to replace terms.
209 def pop_tags(tag):
210     thelist = []
211     for subtag in tag.iter("tag"):
212         newdict = {}
213         valid, typ, key = get_valid_type_key(subtag.get('k'))
214         val = subtag.get('v')
215         if valid and val != None and val != "":
216
217             cleaned_val = get_subtag_value(subtag)
218             if cleaned_val != None and cleaned_val != '':
219                 newdict['id'] = int(tag.get('id'))
220                 newdict['key'] = key
221                 newdict['value'] = cleaned_val
222                 newdict['type'] = typ
223                 thelist.append(newdict)
224
225             if key == "cuisine":
226                 need_drive_through_dict = has_drive_through_text(val)
227                 need_dict_for_vegan = has_vegan_text(val)
228                 need_dict_for_veggie = has_vegetarian_text(val)
229
230                 if need_drive_through_dict and len(tag.findall("./tag/[@k='drive_thro
231                     thelist.append(pop_drive_through(int(tag.get('id')), typ))
232                 if (need_dict_for_vegan or need_dict_for_veggie) and len(tag.findall(
233                     thelist.append(pop_diet(int(tag.get('id')), typ, need_dict_for_ve
234     return thelist
235
236 def get_subtag_value(subtag):
237     val = subtag.get('v')
238     if subtag.get('k') == "cuisine":
239         return shape_cuisine_element(val)
240     else:
241         return val
242
243 # Make the dictionary for way tags
244 def make_way_dict(tag):
245     return { "way": get_way_attribs(tag), 'way_nodes': pop_way_nodes(tag), 'way_tags'
246
247 # Make the dictionary for node tags
248 def make_node_dict(tag):
249     return { "node": pop_attrib(tag), 'node_tags': pop_tags(tag) }
250
251 # Clean up the element, by creating a dictionary for a way or node with cleaned up va
252 def shape_element(element, node_attr_fields=NODE_FIELDS, way_attr_fields=WAY_FIELDS,
253                   problem_chars=PROBLEMCHARS, default_tag_type='regular'):
254
255     if element.tag == 'node':
```

```
256          return make_node_dict(element)
```

**AWESOME**

Great job calling all the functions within the `shape_element` function to update cleaned data back to the m

```
257      elif element.tag == 'way':
258          return make_way_dict(element)
259
260  # Parse the OSM xml and get a dictionary for every element, then write the dictionary
261  def process_map(file_in, validate):
262      #Iteratively process each XML element and write to csv(s)
263
264      with open(NODES_PATH, 'w', newline="", encoding='utf-8') as nodes_file, open(NODE
265          open(WAYS_PATH, 'w', newline="", encoding='utf-8') as ways_file, open(WAY_NOD
266          open(WAY_TAGS_PATH, 'w', newline="", encoding='utf-8') as way_tags_file:
267
268          nodes_writer = csv.DictWriter(nodes_file, fieldnames=NODE_FIELDS, delimiter='
269          node_tags_writer = csv.DictWriter(nodes_tags_file, fieldnames=NODE_TAGS_FIELD
270          ways_writer = csv.DictWriter(ways_file, fieldnames=WAY_FIELDS, delimiter='|')
271          way_nodes_writer = csv.DictWriter(way_nodes_file, fieldnames=WAY_NODES_FIELDS
272          way_tags_writer = csv.DictWriter(way_tags_file, fieldnames=WAY_TAGS_FIELDS, d
273
274          nodes_writer.writeheader()
275          node_tags_writer.writeheader()
276          ways_writer.writeheader()
277          way_nodes_writer.writeheader()
278          way_tags_writer.writeheader()
279
280          for element in get_element(file_in, tags=('node', 'way')):
281              el = shape_element(element)
282              if el:
283                  #if validate:
284                  #    validate_element(el, validator)
285
286                  if element.tag == 'node':
287                      nodes_writer.writerow(el['node'])
288                      node_tags_writer.writerows(el['node_tags'])
289                  elif element.tag == 'way':
290                      ways_writer.writerow(el['way'])
291                      way_nodes_writer.writerows(el['way_nodes'])
292                      way_tags_writer.writerows(el['way_tags'])
293
294
```

▶ README.txt    **1**

▶ references_urls.txt

▶ overpass_api_usage.txt

▶ audit.py

RETURN TO PATH

Student FAQ