

**by Rica Enriquez, July 1, 2015**

### Import the database for querying below

## Section 1. Problems Encountered in the Map

When importing the data, all labels were included. Many of these labels were only available for a small number of samples. To avoid underutilized upper labels and lower labels are removed. If a subset of the dataset is used in the future, the labels will be re-evaluated.

## Removing underutilized labels

```
In [2]: all_labels = ['real_ale', 'fhrrs', 'anglican', 'dance', 'sS052', 'demodifie
            'is_in', 'max_age', 'created_by', 'fax', 'cctv', 'mph',
            'school', 'level', 'notes', 'bus_stop', 'ncn_1', 'real_c
            'cost',
            'exit_to', 'leaf_type', 'access', 'fast_food', 'eligibil
            'used_to_be', 'permit_holders', 'microbrewery', 'survey'
            'lwn',
            'vehicle', 'have_riverbank', 'type', 'start_date', 'entr
            'campaigned_for_by',
            'give_way', 'visibility', 'site', 'phone', 'traffic_caln
            'history', 'estate', 'lock', 'currency', 'ncn_ref', 'por
            'gate', 'uk_postcode_centroid', 'FIXME', 'description',
            'mental'
```

```

'rental',
'natural', 'lcn_ref', 'wheelchair', 'outdoor_seating', '
'postal_code', 'motorcycle', 'int_ref', 'pitch', 'covere
'material', 'foot', 'tourism', 'smoothness', 'fixme', 'r
'embankment',
'crossing', 'kerb', 'name_1', 'frequency', 'naptan', 'ac
'highways_agency', 'ref', 'brewery', 'highway', 'barrier
'electrified',
'was_called', 'old_amenity', 'accommodation', 'tenant',
'box_type', 'turn', 'place', 'high_capacity', 'support',
'priority',
'motorcar', 'park_ride', 'enforcement', 'noname', 'est_w
'population', 'multi_storey', 'royal_cypher', 'aeroway',
'occupier', 'nqa', 'sidewalk', 'hgv', 'lit', 'takeaway',
'aerodrome', 'url', 'medical', 'tactile_paving', 'shop',
'last_survey', 'gauge', 'mapillary', 'wood', 'fuel', 'ia
'artist', 'motorboat', 'public_transport', 'power_source
'lamp_type',
'route_ref', 'parking', 'sport', 'power_supply', 'capaci
'boundary',
'email', 'screen', 'denomination', 'key', 'substation',
'genus', 'comment', 'maintainer', 'wall', 'loading_gauge
'ele', 'alt_description', 'boat', 'speech_output', 'mkgn
'direction', 'lanes', 'building_1', 'craft', 'official_r
'parking_space', 'replaced', 'overtaking', 'layer', 'ons
'guided_busway', 'beer_garden', 'waterway', 'cuisine', '
'collection_times',
'status', 'wires', 'cyclestreets_id', 'fence_type', 'fru
'oneway',
'landmark', 'left', 'taxi', 'livestock', 'proposed', 'hc
'guest_house',
'isced', 'toilets', 'generator', 'TODO', 'bench', 'sourc
'historic',
'lcn', 'psv', 'furniture', 'vending', 'tower', 'internet
'local_ref', 'man_made', 'religion', 'artwork_type', 'pc
'footway',
'industry', 'taxon', 'supervised', 'step_count', 'female
'opening_hours',
'museum', 'width', 'occupier3', 'occupier2', 'admin_leve
'construction',
'diaper', 'courts', 'old_name', 'real_fire', 'circuits',
'crossing_ref',
'cinema', 'carriageway_ref', 'maxheight', 'cafe', 'cable
'interior_decoration', 'cycleway', 'department', 'denota
'diet',
'min_age', 'maxstay', 'opened', 'building', 'yelp', 'wif

```

```

# Number of documents in the collection

```

```

N = db.cambridge.find().count()

```

```

removed = []

```

```

kept = []

```

```

# Remove labels used in less than 1000 documbets

```

```

for label in all_labels:
    pipeline = [{"$group": {"_id": "$" + label, "count": {"$sum": 1}}}, {"
    result = list(db.cambridge.aggregate(pipeline))
    if len(result) > 0:
        n = result[0]["count"]
        if n >= N - 1000:
            db.cambridge.update({}, {"$unset": {label: ""}}, multi=True)
            removed.append(label)
        else:
            kept.append(label)
print len(removed), "labels were removed and", len(all_labels) - len(remov
print "The labels kept are:", kept

```

336 labels were removed and 11 labels were kept.

The labels kept are: ['address', 'amenity', 'entrance', 'natural', 'foot', 'rator']

## Removing underutilized sublabels

Similarly, only lower labels that had at least 500 documents were kept. The "kept" list from the Mon dictionary of the list of kept upper labels and their sublabels. This dictionary is then looped through each sublabel is used in, using MongoDB. If there are less than 500 documents, that sublabel is removed. If a label contains any sublabels, it is also removed. The final structure of the collection is printed out.

```

In [3]: kept_sublabels = {'building': ['name', 'level', 'levels', 'min_level', 'ma
    'maxspeed': ['type', 'ype'],
    'name': ['cy', 'eo', 'ru', 'sr', 'uk', 'zh', 'en', 'zh_f
    'service': ['bicycle:pump', 'bicycle:chain_tool'], 'acce
    'source': ['crossing', 'addr', 'name', 'phone', 'populat
    'location', 'start_date', 'fhhs:id', 'opening
    'ele', 'cost', 'taxon', 'database', 'position
    'traffic_calming', 'designation', 'operator',
    'bicycle:backward', 'tourist_bus:backward', '
    'hgv', 'outline', 'maxspeed:date', 'addr:hous
    'bridge', 'addr:postcode', 'tracktype', 'heig
    'address': ['street', 'postcode', 'housenumber', 'houser
    'interpolation', 'flat', 'flats', 'place', '
    'ref': ['university_of_cambridge', 'observado']}

# Remove sublabels used in less than 500 documents
removed_sub = {}
kept_sub = {}
for label in kept_sublabels.keys():
    for sublabel in kept_sublabels[label]:
        pipeline = [{"$group": {"_id": ".".join(["$", label, ".", sublabel]
            {"$match": {"_id": None}}}]
        result = list(db.cambridge.aggregate(pipeline))
        if len(result) > 0:
            n = result[0]["count"]
            if n >= N - 500:
                db.cambridge.update({}, {"$unset": {".".join([label, '.', s

```

```

        try:
            removed_sub[label].append(sublabel)
        except:
            removed_sub[label] = [sublabel]
    else:
        try:
            kept_sub[label].append(sublabel)
        except:
            kept_sub[label] = [sublabel]

# Remove the upper labels that no longer have sublabels
for label in removed_sub.keys():
    if label not in kept_sub.keys():
        db.cambridge.update({}, {"$unset": {label: ""}}, multi=True)

# Print the final structure of the collection
final_labels = {}
for label in kept:
    if label in kept_sub:
        final_labels[label] = kept_sub[label]
    elif label not in removed_sub:
        final_labels[label] = None

print "The final structure of the collection is:"
pprint.pprint(final_labels)

```

The final structure of the collection is:

```

{'address': ['street',
             'postcode',
             'houzenumber',
             'houzename',
             'city',
             'country',
             'interpolation'],
 'amenity': None,
 'barrier': None,
 'entrance': None,
 'foot': None,
 'highway': None,
 'landuse': None,
 'natural': None,
 'operator': None,
 'source': ['name']}

```

## Remove postcodes that do not start with "CB"

All postcodes in Cambridge start with CB. Listing the postcodes in the collection, documents with a since it is for a place in Royston and the Stevenage postcode area. Since there are only two docum approach to extend the collection the surrounding area. Additionally the postcode "CB1" is incompl Documents with this are also removed.

```
In [4]: db.cambridge.remove({"address.postcode": "CB1"})
        db.cambridge.remove({"address.postcode": "SG8 5TF"})
```

```
Out[4]: {u'n': 0, u'ok': 1}
```

## Update Cities to Cambridge

Some cities were "cambridge" and not "Cambridge", overspecified to "Girton" or "South Cambridge", or "Cambridge" and "Girton". We used the following logic to show that each entry is in Cambridge. Therefore, the "city" is updated to "Cambridge".

```
In [5]: db.cambridge.update({"address.city": "cambridge"}, {"$set": {"address.city": "cambridge"},
                             multi=True)
db.cambridge.update({"address.city": "South Cambridgeshire"}, {"$set": {"address.city": "South Cambridgeshire"},
                     multi=True)
db.cambridge.update({"address.city": "Girton"}, {"$set": {"address.city": "Girton"}, multi=True)
db.cambridge.update({"address.city": "11"}, {"$set": {"address.city": "Cambridge"}, multi=True)
```

```
Out[5]: {u'n': 0, u'nModified': 0, u'ok': 1, 'updatedExisting': False}
```

## Pare down barrier, entrance, highway, landuse, and operator

Some entries for these sublabels were the same, but in a different format. They were updated to be

```
In [6]: # Pare down barrier
pipeline = [{"$group": {"_id": "$barrier", "count": {"$sum": 1}}},
            {"$sort": {"_id": -1}}]
result = list(db.cambridge.aggregate(pipeline))
pprint.pprint(result)

db.cambridge.update({"barrier": "fence;wall"}, {"$set": {"barrier": "fence;wall"}}, upsert=True, multi=True)
db.cambridge.update({"barrier": "fence;wall"}, {"$set": {"barrier": "fence;wall"}}, upsert=True, multi=True)
db.cambridge.update({"barrier": "fedr"}, {"$set": {"barrier": None}}, upsert=True, multi=True)
db.cambridge.update({"barrier": "bollards"}, {"$set": {"barrier": "bollards"}}, upsert=True, multi=True)

# Pare down entrance
pipeline = [{"$group": {"_id": "$entrance", "count": {"$sum": 1}}},
            {"$sort": {"_id": -1}}]
result = list(db.cambridge.aggregate(pipeline))
pprint.pprint(result)

db.cambridge.update({"entrance": "secondary_entrance"}, {"$set": {"entrance": "secondary_entrance"}}, upsert=True, multi=True)
db.cambridge.update({"entrance": "main_entrance; porters"}, {"$set": {"entrance": "main_entrance; porters"}}, upsert=True, multi=True)
db.cambridge.update({"entrance": "porters;main_entrance"}, {"$set": {"entrance": "porters;main_entrance"}}, upsert=True, multi=True)
db.cambridge.update({"entrance": "main"}, {"$set": {"entrance": "main"}}, upsert=True, multi=True)
db.cambridge.update({"entrance": "emergency"}, {"$set": {"entrance": "emergency"}}, upsert=True, multi=True)
db.cambridge.update({"entrance": "main_entrance;porters;"}, {"$set": {"entrance": "main_entrance;porters;"}}, upsert=True, multi=True)
```

```

        upsert=False, multi=True)

# Pare down highway
pipeline = [{"$group": {"_id": "$highway", "count": {"$sum": 1}}},
            {"$sort": {"_id": -1}}]
result = list(db.cambridge.aggregate(pipeline))
pprint.pprint(result)

db.cambridge.update({"highway": "bus_stand"}, {"$set": {"highway": "bus_st

# Pare down landuse
pipeline = [{"$group": {"_id": "$barrier", "count": {"$sum": 1}}},
            {"$sort": {"_id": -1}}]
result = list(db.cambridge.aggregate(pipeline))
pprint.pprint(result)
db.cambridge.update({"landuse": "institutional"}, {"$set": {"landuse": "

        multi=True)

# Pare down operator
pipeline = [{"$group": {"_id": "$operator", "count": {"$sum": 1}}},
            {"$sort": {"_id": -1}}]
result = list(db.cambridge.aggregate(pipeline))
pprint.pprint(result)
db.cambridge.update({"operator": "YourSpace"}, {"$set": {"operator": "Your
        multi=True)
db.cambridge.update({"operator": "Your Space"}, {"$set": {"operator": "You
        multi=True)
db.cambridge.update({"operator": "Trinity College"},
                    {"$set": {"operator": "Trinity College (University of
db.cambridge.update({"operator": "St John's College"},
                    {"$set": {"operator": "St John's College (University c
db.cambridge.update({"operator": "Lucy Cavendish College"},
                    {"$set": {"operator": "Lucy Cavendish College (Univers
        multi=True)
db.cambridge.update({"operator": "Lloyds"}, {"$set": {"operator": "Lloyds
db.cambridge.update({"operator": "King's College"},
                    {"$set": {"operator": "King's College (University of C
db.cambridge.update({"operator": "King's College (University Of Cambridge)
                    {"$set": {"operator": "King's College (University of C
db.cambridge.update({"operator": "Needham Institute"}, {"$set": {"operator
        upsert=False, multi=True)
db.cambridge.update({"operator": "Gonville and Caius College (University c
                    {"$set": {"operator": "Gonville & Caius College (Unive
        multi=True)
db.cambridge.update({"operator": "EDF"}, {"$set": {"operator": "EDF Energy
db.cambridge.update({"operator": "Clare College"},
                    {"$set": {"operator": "Clare College (University of Ca
db.cambridge.update({"operator": "Christ's College"},
                    {"$set": {"operator": "Christ's College (University of

{u'_id': u'kissing_gate', u'count': 12},
{u'_id': u'kerb', u'count': 34},
{u'_id': u'height_restrictor', u'count': 1},
{u'_id': u'hedge', u'count': 490},
{u'_id': u'gravestone', u'count': 1},
{u'_id': u'gate', u'count': 807}

```

```
{u'_id': u'gate', u'count': 80},
{u'_id': u'full-height_turnstile', u'count': 1},
{u'_id': u'fence', u'count': 5405},
{u'_id': u'entrance', u'count': 155},
{u'_id': u'ditch', u'count': 1},
{u'_id': u'cycle_barrier', u'count': 39},
{u'_id': u'chicane', u'count': 2},
{u'_id': u'cattle_grid', u'count': 36},
{u'_id': u'car_trap', u'count': 9},
{u'_id': u'boom', u'count': 1},
{u'_id': u'bollard', u'count': 186},
{u'_id': u'block', u'count': 1},
{u'_id': u'barrier', u'count': 3},
{u'_id': None, u'count': 298180}]
[{'_id': u'Your Space Apartments', u'count': 5},
```

## Section 2. Overview of the Data

A statistical overview of the dataset with the MongoDB queries used to obtain such statistics are bel

```
In [7]: print "The size of 'cambridge_england.osm' is", os.stat("cambridge_england.osm").st_size
print "The size of 'cambridge_england.osm.json' is", os.stat("cambridge_england.osm.json").st_size
N2 = db.cambridge.find().count()
print "There are", N, "documents in the original set and", N2, "documents in the aggregated set."
pipeline = [{"$group": {"_id": "$created.user", "count": {"$sum": 1}}}]
print "There are", len(list(db.cambridge.aggregate(pipeline))), "unique users."
print "There are", db.cambridge.find({"type": "node"}).count(), "nodes."
print "There are", db.cambridge.find({"type": "way"}).count(), "ways."
pipeline = [{"$group": {"_id": "$created.user", "count": {"$sum": 1}}},
             {"$sort": {"count": -1}},
             {"$limit": 1}]
print list(db.cambridge.aggregate(pipeline))[0]['_id'], "contributed the most to this collection with",
        list(db.cambridge.aggregate(pipeline))[0]['count'], "documents."
pipeline = [{"$group": {"_id": "$created.user", "count": {"$sum": 1}}},
             {"$group": {"_id": "$count", "num_users": {"$sum": 1}}},
             {"$sort": {"_id": 1}},
             {"$limit": 1}]
print list(db.cambridge.aggregate(pipeline))[0]['num_users'], "users contributed the most to this collection."
pipeline = [{"$group": {"_id": "$amenity", "count": {"$sum": 1}}},
             {"$match": {"_id": {"$ne": None}}},
             {"$sort": {"count": -1}},
             {"$limit": 5}]
print list(db.cambridge.aggregate(pipeline))[0]["_id"], list(db.cambridge.find({"amenity": list(db.cambridge.aggregate(pipeline))[0]["_id"]}).count()), "are the top two amenities."
```

The size of 'cambridge\_england.osm' is 61.828774 MB.

The size of 'cambridge\_england.osm.json' is 87.071371 MB.

There are 306428 documents in the original set and 306428 documents in the aggregated set.

There are 453 unique users.

There are 257058 nodes.

There are 49351 ways.

smb1001 contributed the most to this collection with 81443 documents.

120 users contributed once.

university bicycle\_parking are the top two amenities.

## Section 3. Additional Ideas

Some addresses have house names. It would be interesting to know if there's a certain postal code

Additionally, it would be interesting where the top amenities are located.



```

In [8]: # Find the top postcodes with housenames
pipeline = [{"$match": {"address.housename": {"$exists": True}}},
            {"$group": {"_id": "$address.postcode", "count": {"$sum":
            {"$sort": {"count": -1}}}]
pprint.pprint(list(db.cambridge.aggregate(pipeline)))

# Find the top operators with housenames
pipeline = [{"$match": {"address.housename": {"$exists": True}}},
            {"$group": {"_id": "$operator", "count": {"$sum": 1}}},
            {"$sort": {"count": -1}}]
pprint.pprint(list(db.cambridge.aggregate(pipeline)))

# Find the top amenities with postcodes and sort by postcode
pipeline = [{"$match": {"amenity": {"$ne": None}}},
            {"$match": {"amenity": {"$ne": "university"}}},
            {"$match": {"address.postcode": {"$ne": None}}},
            {"$group": {"_id": {"amenity": "$amenity",
                                "postcode": "$address.postcode"},
                        "count": {"$sum": 1}}},
            {"$sort": {"count": -1}},
            {"$group": {"_id": "$_id.amenity",
                        "info": {"$push": {
                                "postcode": "$_id.postcode",
                                "count": "$count"},},
                        "count": {"$sum": "$count"}}},
            {"$sort": {"count": -1}},
            {"$limit": 5}]
pprint.pprint(list(db.cambridge.aggregate(pipeline)))
[{"_id": "CB3 0JB", "count": 1},
 {"_id": "CB2 1UA", "count": 1},
 {"_id": "CB4 1HQ", "count": 1},
 {"_id": "CB2 3LL", "count": 1}]
[{"_id": None, "count": 514},
 {"_id": "University of Cambridge", "count": 53},
 {"_id": "Clare College (University of Cambridge)", "count": 6},
 {"_id": "Selwyn College (University of Cambridge)", "count": 4},
 {"_id": "Anglia Ruskin University", "count": 4},
 {"_id": "Riverside ECHG", "count": 3},
 {"_id": "CRM Students", "count": 2},
 {"_id": "Murray Edwards College (University of Cambridge)", "count": 2},
 {"_id": "Newnham College (University of Cambridge)", "count": 2},
 {"_id": "Girton College (University of Cambridge)", "count": 2},
 {"_id": "Christ's College (University of Cambridge)", "count": 2},
 {"_id": "British Antarctic Survey", "count": 1},
 {"_id": "Clare Hall (University of Cambridge)", "count": 1},
 {"_id": "St John's College (University of Cambridge)", "count": 1},
 {"_id": "Evangelical Lutheran Church of England", "count": 1},
 {"_id": "Murray Edwards College (University of Cambridge)", "count": 1}]

```

## Section 4. Conclusions

The Cambridge, England OpenStreetMap dataset is full of information. However, this can also be cleaned up. Some labels and sublabels were removed. It may be useful to move the information to an existing label rather than creating new ones. For example, "have\_riverbank" and "trees" can be moved to the kept "natural" label. Additionally, there are some "cuisine" tags that were prematurely removed. The top cuisine for this area couldn't be examined with Mongo

