

Data Analyst Nanodegree: P2_Data Wrangle OpenStreetMaps Data

Zhihui Xie

xiewisdom@gmail.com

Project Summary

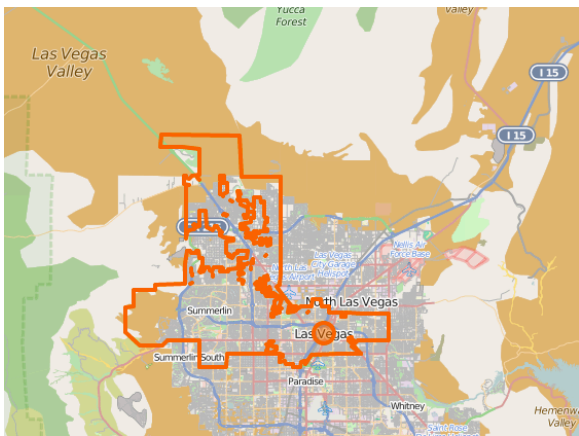
OpenStreetMap, an open project to create a free map around the world, is a powerful tool for viewing maps and humanitarian aid. The initial data of the map were collected using a handheld GPS and a notebook, digital camera, or a voice recorder. These data heavily rely on the human input, which may cause inconsistent input, misspellings or error (e.g. inconsistent input of Street, street, st., St.). This project is going to focus on map of Las Vegas, Nevada, USA and wrangling the data: 1) overview data (code: mapparser.py); 2) check the “k” value for each “<tag>” (code: tags.py); 3) find out number of unique users contributed to the map (code: users.py); 4) fix unexpected street types (e.g. street, st., St. to be Street) (code: audit.py); 5) transform the shape of data and insert data into mongodb (code: data.py and mongodb.py); 6) use MongoDB queries to find number and names of hotels, and shopping malls (code: query.by). The codes mentioned above are included in a separated file and tested by sample dataset (a part of the whole dataset) before run for the whole dataset.

Data source

I chose Las Vegas for data wrangling because it is one of the most popular tourist destination and my dream place for vacation. I hope I can have a vacation there this holiday and would like to know it well. I hope this map to be improved in quality!

I explored the data from the following link:

URL: <https://www.openstreetmap.org/relation/170117#map=11/36.2501/-115.2531>



Is there a list of Web sites, books, forums, blog posts, github repositories etc that you referred to or used in this submission (Add N/A if you did not use such resources)? The above images were cited from OpenStreetMap and wiki:

<https://www.openstreetmap.org/relation/170117#map=11/36.2551/-115.2387>
http://en.wikipedia.org/wiki/Las_Vegas#mediaviewer/File:LasVegasMontage6.jpg

Problems encountered in the map

In this dataset, there are 7 types of problem for street types (table 1). Some of them are inconsistent input, for example, “AVE”, “Ave” and “Ave.”. Some are abbreviation: “Ln.”, “PkwY”, “St” and “Rd”. Using attached auditing code (audit.py), these inconsistent input or abbreviation were fixed (table 1).

In addition to street type, the keys were also validated (using code: tags.py). In this dataset, there are 240883 keys including 'lower' regular expression ('^([a-z]|_)*\$'), 263369 keys including 'lower_colon' regular expression ('^([a-z]|_)*:([a-z]|_)*\$') and 2 including 'problemchars' expression ('[=\\+/&<>;\\\"\\'/?%#\$@\\., \\t\\r\\n]'). The keys with 'problemchars' expression will be ignored during data shape transformation (code: data.py).

Table 1 Fixing problem of street types

Types of problem	I			II		III		IV	V	VI	VII
Before auditing	AVE	Ave	Ave.	Blvd	Blvd.	Dr	Dr.	Ln.	PkwY	St	Rd
After auditing	Avenue			Boulevard		Drive		Lane	Parkway	Street	Road

Overview of the data

To overview the data, the following items were recorded:

- size of the file (las-vegas_nevada.osm): 166.9 MB;
- number of unique users (code: users.py): contributed by 618 unique users;
- number of nodes and ways (code: mappraser.py): nodes: 723370, ways: 78394;
- number of all type of amenities (attached code: query.py, function: get_query (data, db)):

```
{'grave_yard': 1, 'taxi': 3, 'marketplace': 1, 'casino': 33, 'public_building': 38, 'self_storage': 1, 'cinema': 6, 'fitness_center': 2, 'recycling': 6, 'car_wash': 20, 'telephone': 12, 'library': 23, 'conciierge': 1, 'clinic': 1, 'oil_tank': 2, 'college': 6, 'parking': 746, 'spa': 1, 'post_office': 62, 'cafe': 49, 'yes': 3, 'koolsville tattoo': 1, 'toilets': 53, 'ice_cream': 2, 'police': 18, 'townhall': 7, 'food_court': 1, 'hospital': 66, 'veterinary': 4, 'antiques': 1, 'pharmacy': 24, 'kindergarten': 2, 'shelter': 51, 'fountain': 269, 'bicycle_rental': 1, 'prison': 5, 'fuel': 139, 'smog': 1, 'bicycle_parking': 3, 'bbq': 6, 'car_rental': 8, 'fast_food': 158, 'fire_station': 70, 'hotel': 1, 'theatre': 26, 'post_box': 6, 'social_facility': 2, 'arts_centre': 2, 'bat': 1, 'pub': 16, 'waste_basket': 1, 'flower_shop': 1, 'dentist': 4, 'doctors': 7, 'adult day care': 1, 'bank': 48, 'courthouse': 3, 'place_of_worship': 367, 'school': 535, 'bar': 50, 'restaurant': 200, 'swingerclub': 1, 'parking_entrance': 5, 'university': 6, 'atm': 11, 'nightclub': 10,
```

```
'drinking_water': 12, 'swimming_pool': 21, 'mall': 2, 'lounge': 1, 'vending_machine': 4,
'bus_station': 6, 'bench': 16, 'childcare': 1, 'whirlpool': 1, 'finish_line': 1}
```

Hotels and malls in the datasets

I am interested in finding all the hotels and shopping malls in Las Vegas. To do that, I try to find hotels and malls listed as amenity using two pipelines (pipeline1 and pipeline3 in attached code: query.py, function: get_pipeline()) to query database inserted into MongoDB. However, only one hotel without name and two malls were found, which does not make sense. I went back to check the original dataset and noticed that “hotel” and “mall” were listed as “tourism” and “shop”, respectively. Thus, I edited the query pipelines (pipeline2 and pipeline4) and got the following results: Hotels as amenity: 1, Hotels as tourism: 103; Malls as amenity: 2, Malls as shop: 8.

Pipelines:

```
#hotels as amenity
pipeline1 = [{"$match": {"amenity": "hotel"}},
{"$project": {"_id": "$name", "cuisine": "$cuisine", "phone": "$phone"}}]
#hotels as tourism
pipeline2 = [{"$match": {"tourism": "hotel"}},
{"$project": {"_id": "$name", "address": "$address", "phone": "$phone"}}]
#mall as amenity
pipeline3 = [{"$match": {"amenity": "mall"}},
{"$project": {"_id": "$name", "cuisine": "$cuisine", "phone": "$phone"}}]
#mall as shop
pipeline4 = [{"$match": {"shop": "mall"}},
{"$project": {"_id": "$name", "address": "$address", "phone": "$phone"}}]
```

Outputs:

```
Hotels as amenity: 1 {u'ok': 1.0, u'result': [{u'_id': u'SLS Hotel and
Casino', u'address': {u'city': u'Las Vegas', u'housenumber': u'2535',
u'89109', u'state': u'NV', u'street': u'Las Vegas Boulevard S', u'phone': u'(702)
737-2111', {u'_id': u'Four Seasons', {u'_id': u'Tuscany', {u'_id': u'Embassy Suites',
{u'_id': u'Silver Sevens Hotel & Casino', u'address': {u'city': u'Las Vegas', u'housenumber':
u'4100', u'postcode': u'89169', u'street': u'Paradise Road', u'phone': u'(702)
733-7000', {u'_id': u'Suncoast Hotel and Casino', u'address': {u'state': u'NV'}}, {u'_id': u'Lake
Mead Lodge', u'address': {u'state': u'NV'}}, {u'_id': u'Encore Casino at Wynn', u'address': {u'city': u'Las
Vegas', u'housenumber': u'3131', u'postcode': u'89109', u'state': u'NV',
u'street': u'Las Vegas Boulevard S'}, {u'_id': u'Hacienda Hotel and Casino', u'phone': u'(702)
293-5000', {u'_id': u'Hilton Grand Vacations Suites on the Las Vegas Strip', u'address': {u'city': u'Las
Vegas', u'housenumber': u'2650', u'postcode': u'89109', u'state': u'NV',
u'street': u'Las Vegas Boulevard', u'phone': u'702-765-8300', {u'_id': u'Hilton Grand Vacations Suites -
Las Vegas (Convention Center)', u'address': {u'city': u'Las Vegas', u'housenumber': u'455',
u'postcode': u'89109', u'state': u'NV', u'street': u'Karen Avenue', u'phone':
u'702-946-9210', {u'_id': u'Embassy Suites Convention Center Las Vegas', u'address': {u'city': u'Las
Vegas', u'housenumber': u'3600', u'postcode': u'89169', u'state': u'NV',
u'street': u'Paradise Road', u'phone': u'702-893-8000', {u'_id': u'Renaissance Las Vegas Hotel',
u'address': {u'city': u'Las Vegas', u'housenumber': u'3400', u'postcode': u'89169',
u'state': u'NV', u'street': u'Paradise Road', u'phone': u'702-784-5700', {u'_id':
u'Extended Stay America', {u'_id': u'Boulder Station Hotel and Casino', {u'_id': u'Element Las Vegas
Summerline, 10555 Discovery Drive', {u'_id': u'Main Street Station Casino and Hotel', {}, {u'_id':
u'Americas Best Value Inn', {u'_id': u'Summer Bay Resort', {u'_id': u'Candlewood Suites Las Vegas',
{u'_id': u'Quality Inn', {u'_id': u'Boulder Dam Hotel', {u'_id': u'Boulder Inn and Suites', {u'_id':
u'Circus Circus, 2880 Las Vegas Boulevard, NV 89109 Las Vegas, Vereinigte Staaten von Amerika', {u'_id':
u'Fortune Hotel & Suites', u'address': {u'city': u'Las Vegas', u'housenumber': u'325',
u'postcode': u'89169', u'street': u'East Flamingo Road', u'phone': u'702-830-4010',
```

{u'_id': u'Super 8'},	{u'_id': u'Stratosphere Hotel and Casino'},	{u'_id': u'Residence Inn Las Vegas
Convention Center',	u'address': {u'city': u'Las Vegas',	u'housenumber': u'3225',
u'postcode': u'NV 89109',	u'street': u'Paradise Road'}},	{u'_id': u'Courtyard Las Vegas
Convention Center',	u'address': {u'city': u'Las Vegas',	u'housenumber': u'3275',
u'postcode': u'NV 89109',	u'street': u'Paradise Road'}},	{u'_id': u'Hilton Garden Inn Henderson'},
{u'_id': u'Palms Casino Resort'},	{u'_id': u'Gold Coast Hotel & Casino'},	{}, {u'_id': u'Tropicana
Hotel and Casino',	u'address': {u'city': u'Las Vegas',	u'county': u'Clark',
u'housenumber': u'3801',	u'postcode': u'89109',	u'state': u'NV',
u'Las Vegas Boulevard South'}},	{u'_id': u'Luxor Hotel and Casino',	u'address': {u'county': u'Clark',
u'housenumber': u'3900',	u'postcode': u'89119',	u'state': u'NV',
u'Las Vegas Boulevard South'}},	{u'_id': u'Mandalay Bay Casino & Resort'},	{u'_id': u'Polo Towers'},
{u'_id': u'Encore'},	{u'_id': u'Mirage Hotel and Casino',	u'address': {u'city': u'Las Vegas',
u'county': u'Clark',	u'housenumber': u'3400',	u'postcode': u'89109',
u'state': u'NV',	u'street': u'Las Vegas Boulevard South'}},	{u'_id': u'Treasure Island Hotel and
Casino',	u'address': {u'city': u'Las Vegas',	u'county': u'Clark',
u'housenumber': u'3300',	u'postcode': u'89109',	u'state': u'NV',
u'street': u'Las Vegas		
Boulevard South'}},	{u'_id': u'Best Western Mardi Gras Inn'},	{u'_id': u'Plaza Hotel',
{u'city': u'Las Vegas',	u'housenumber': u'1',	u'postcode': u'89101',
u'street': u'Main Street'},	u'phone': u'(702) 386-2110'},	{u'_id': u'Paris Hotel and Casino',
u'address': {u'city': u'Las Vegas',	u'county': u'Clark',	u'housenumber': u'3655',
u'postcode': u'89109',	u'state': u'NV',	u'street': u'Las Vegas Boulevard South'}},
{u'_id': u'Fiesta Henderson'},	{u'_id': u'Hilton Garden Inn'},	{u'_id': u'Crestwood Suites'},
{u'_id': u'Manor Suites',	u'address': {u'housenumber': u'7230'},	{u'_id': u'Tahiti Village'},
u'address': {u'housenumber': u'55'},	{u'_id': u'Mandarin Oriental'},	{u'_id': u'Micdrotel Inn',
u'address': {u'city': u'Las Vegas',	u'housenumber': u'3730',	u'postcode': u'89109',
u'state': u'NV',	u'street': u'Las Vegas Boulevard'}},	{u'_id': u'Vdara',
u'Las Vegas',	u'country': u'US',	u'housenumber': u'2600',
u'89109',	u'state': u'NV',	u'street': u'West Harmon Avenue'}},
Cromwell',	u'address': {u'city': u'Las Vegas',	u'county': u'Clark',
u'housenumber': u'3595',	u'postcode': u'89109',	u'state': u'NV',
u'Las Vegas Boulevard South'}},	{u'_id': u'Marriott's Grand Chateau"},	{u'_id': u'The Carriage House',
u'address': {u'housenumber': u'105',	u'street': u'East Harmon Avenue',	u'phone':
u'702.798.1020'},	{u'_id': u'The California Hotel and Casino'},	{u'_id': u'Trump Hotel'},
u'Bellagio Hotel and Casino',	u'address': {u'city': u'Las Vegas',	u'county': u'Clark',
u'housenumber': u'3600',	u'postcode': u'89109',	u'state': u'NV',
u'Las Vegas Boulevard South'}},	{u'_id': u'Excalibur Hotel and Casino',	u'address': {u'city': u'Las Vegas',
u'county': u'Clark',	u'housenumber': u'3850',	u'postcode': u'89109',
u'state': u'NV',	u'street': u'Las Vegas Boulevard South'}},	{u'_id': u'Monte Carlo Hotel and Casino',
u'address': {u'city': u'Las Vegas',	u'county': u'Clark',	u'housenumber': u'3770',
u'postcode': u'89109',	u'state': u'NV',	u'street': u'Las Vegas Boulevard South'}},
{u'_id': u'New York New York Hotel and Casino',	u'address': {u'city': u'Las Vegas',	u'county':
u'Clark',	u'housenumber': u'3790',	u'postcode': u'89109',
u'street': u'Las Vegas Boulevard South'}},	{u'_id': u'Caesars Hotel and Casino',	u'address': {u'city': u'Las
Vegas',	u'county': u'Clark',	u'housenumber': u'3570',
u'89109',	u'state': u'NV',	u'street': u'Las Vegas Boulevard'}},
u'Resort & Casino',	u'address': {u'city': u'Las Vegas',	u'county': u'Clark',
u'housenumber': u'3535',	u'postcode': u'89109',	u'state': u'NV',
u'Las Vegas Boulevard South',	u'phone': u'(800) 634-6441'},	{u'_id': u'Venetian Hotel and Casino',
u'address': {u'city': u'Las Vegas',	u'county': u'Clark',	u'housenumber': u'3355',
u'postcode': u'89109',	u'state': u'NV',	u'street': u'Las Vegas Boulevard South'}},
{u'_id': u'Mandalay Bay Hotel and Casino',	u'address': {u'state': u'NV'}},	{u'_id': u'Flamingo Hotel and
Casino',	u'address': {u'city': u'Las Vegas',	u'county': u'Clark',
u'3555',	u'postcode': u'89109',	u'state': u'NV',
u'Boulevard South'}},	{u'_id': u'MGM Grand Hotel and Casino',	u'address': {u'city': u'Las Vegas',
u'county': u'Clark',	u'housenumber': u'3799',	u'postcode': u'89109',
u'state': u'NV',	u'street': u'Las Vegas Boulevard South'}},	{u'_id': u'Bally's Hotel and Casino",
u'address': {u'city': u'Las Vegas',	u'county': u'Clark',	u'housenumber': u'3645',
u'postcode': u'89109',	u'state': u'NV',	u'street': u'Las Vegas Boulevard South'}},
{u'_id': u'Rio All Suite Hotel and Casino'},	{u'_id': u'The Hotel at Mandalay Bay'},	{u'_id': u'Hooters'},
{u'_id': u'Harrah's Hotel and Casino",	u'address': {u'city': u'Las Vegas',	u'county': u'Clark',
u'housenumber': u'3475',	u'postcode': u'89109',	u'state': u'NV',
u'Las Vegas Boulevard South'}},	{u'_id': u'the D Las Vegas Hotel & Casino',	u'address': {u'city': u'Las
Vegas',	u'housenumber': u'301',	u'postcode': u'89101',
		u'street':

```

u'Fremont Street'}, u'phone': u'(702) 388-2400'}, {u'_id': u'La Quinta Inn'}, {u'_id': u'Hard Rock
Hotel and Casino'}, {u'_id': u'Wyndham Grand Desert'}, {u'_id': u'LVH Las Vegas Hotel & Casino'},
u'address': {u'city': u'Las Vegas', u'housenumber': u'3000', u'postcode': u'89109',
u'state': u'NV', u'street': u'Paradise Road'}, u'phone': u'702-732-5111'}, {u'_id': u'Planet
Hollywood Resort and Casino', u'address': {u'city': u'Las Vegas', u'country': u'US',
u'county': u'Clark', u'housenumber': u'3667', u'postcode': u'89109',
u'state': u'NV', u'street': u'South Las Vegas Boulevard'}, u'phone': u'+1 866 919 7472'},
{u'_id': u'The Westin Casuarina'}, {u'_id': u'Circus Circus Las Vegas Hotel and Casino', u'address':
{u'city': u'Las Vegas', u'housenumber': u'2880', u'postcode': u'89109',
u'state': u'NV', u'street': u'Las Vegas Boulevard'}}, {u'_id': u'Stratosphere Hotel and Casino',
u'address': {u'city': u'Las Vegas', u'housenumber': u'2000', u'postcode': u'89104',
u'state': u'NV', u'street': u'Las Vegas Boulevard South'}}, {u'_id': u'J. W. Marriot Resort'},
{u'_id': u'Hampton Inn and Suites'}, {u'_id': u'The Palazzo', u'address': {u'city': u'Las Vegas',
u'county': u'Clark', u'housenumber': u'3325', u'postcode': u'89109',
u'state': u'NV', u'street': u'Las Vegas Boulevard South'}}, {}, {u'_id': u'Gold Strike'},
{u'_id': u'Spring Hill Suites'}, {u'_id': u'Town Place Suites'}, {u'_id': u'Wynn Hotel & Casino'},
{u'_id': u'The Grand'}, {u'_id': u'Hyatt Place'}, {u'_id': u'Super 8'}, {u'_id': u'Golden Gate Hotel and
Casino'}, {u'_id': u'Palace Station'}, {}, {u'_id': u'Blue Moon'}, {u'_id': u'Eastside Cannery
Casino Hotel', u'address': {u'city': u'Whitney', u'housenumber': u'5255',
u'postcode': u'89122', u'street': u'Boulder Highway'}}, {u'_id': u'Hampton Inn', u'address':
{u'city': u'Las Vegas', u'housenumber': u'7100', u'postcode': u'89128',
u'street': u'Cascade Valley Court'}}, {u'_id': u'La Quinta Inn & Suites', u'address': {u'city': u'Las Vegas',
u'housename': u'La Quinta Inn & Suites', u'housenumber': u'7101', u'postcode':
u'89128', u'street': u'Cascade Valley Court'}}, {u'_id': u'Best Western McCarran Inn'},
{u'_id': u'SpringHill Suites by Marriott', u'address': {u'country': u'US'}}}] Malls as amenity: 2 {u'ok': 1.0, u'result':
[{u'_id': u'Crystals'}, {u'_id': u'Miracle Mile Shops'}]} Malls as shop: 8 {u'ok': 1.0, u'result': [{u'_id': u'Fashion Show Mall'},
{}, {u'_id': u'Las Vegas Outlet Center Annex'}, {}, {}, {}, {u'_id': u'Fashion Show
Mall'}, {u'_id': u'Downtown Container Park', u'address': {u'city': u'Las Vegas',
u'housenumber': u'707', u'postcode': u'89101', u'street': u'Fremont Street'}}]}

```

Acknowledgement

In addition, to find users who contribute to this map, I setup pipeline5 to query MongoDB and showed top 5 contributors as follow:

Pipelines:

#top 5 contributors to the map

```

pipeline5 = [{"$match": {"created.user": {"$exists": 1}}},
{"$group": {"_id": "$created.user", "count": {"$sum": 1}}},
{"$sort": {"count": -1}},
{"$limit": 5}]

```

Outputs:

```

Top 5 contributors to the map: 5 {u'ok': 1.0, u'result': [{u'_id': u'alimamo', u'count': 254892},
{u'_id': u'woodpeck_fixbot', u'count': 79690}, {u'_id': u'nmixer', u'count': 70150}, {u'_id':
u'gMitchellD', u'count': 48266}, {u'_id': u'robgeb', u'count': 43484}]

```

Thank you for all contributors to this map and Udacity staffs whom help building and evaluating this project.

By Zhihui Xie

12/11/2014