# Taxi demand prediction in New York City

In [1]:

```
!pip3 install graphviz
!pip3 install dask
!pip3 install toolz
!pip3 install cloudpickle
!pip3 install folium
!pip3 install gpxpy
```

```
Collecting graphviz
  Downloading
https://files.pythonhosted.org/packages/1f/e2/ef2581b5b86625657afd32030f90cf2717456c1d2b711ba074bf0
f1a/graphviz-0.10.1-py2.py3-none-any.whl
scikit-umfpack 0.3.2 has requirement numpy>=1.15.3, but you'll have numpy 1.15.1 which is
incompatible.
menpo 0.8.1 has requirement matplotlib<2.0,>=1.4, but you'll have matplotlib 2.2.3 which is
incompatible.
menpo 0.8.1 has requirement pillow<5.0,>=3.0, but you'll have pillow 5.2.0 which is incompatible.
menpo 0.8.1 has requirement scipy<1.0,>=0.16, but you'll have scipy 1.1.0 which is incompatible.
fastai 1.0.42 has requirement torch>=1.0.0, but you'll have torch 0.4.1 which is incompatible.
Installing collected packages: graphviz
Successfully installed graphviz-0.10.1
You are using pip version 10.0.1, however version 19.0.3 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.
Requirement already satisfied: dask in /usr/local/lib/python3.6/site-packages (1.1.1)
scikit-umfpack 0.3.2 has requirement numpy>=1.15.3, but you'll have numpy 1.15.1 which is
incompatible.
menpo 0.8.1 has requirement matplotlib<2.0,>=1.4, but you'll have matplotlib 2.2.3 which is
incompatible.
menpo 0.8.1 has requirement pillow<5.0,>=3.0, but you'll have pillow 5.2.0 which is incompatible.
menpo 0.8.1 has requirement scipy<1.0,>=0.16, but you'll have scipy 1.1.0 which is incompatible.
fastai 1.0.42 has requirement torch>=1.0.0, but you'll have torch 0.4.1 which is incompatible.
You are using pip version 10.0.1, however version 19.0.3 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.
Requirement already satisfied: toolz in /usr/local/lib/python3.6/site-packages (0.9.0)
scikit-umfpack 0.3.2 has requirement numpy>=1.15.3, but you'll have numpy 1.15.1 which is
incompatible.
menpo 0.8.1 has requirement matplotlib<2.0,>=1.4, but you'll have matplotlib 2.2.3 which is
incompatible.
menpo 0.8.1 has requirement pillow<5.0,>=3.0, but you'll have pillow 5.2.0 which is incompatible.
menpo 0.8.1 has requirement scipy<1.0,>=0.16, but you'll have scipy 1.1.0 which is incompatible.
fastai 1.0.42 has requirement torch>=1.0.0, but you'll have torch 0.4.1 which is incompatible.
You are using pip version 10.0.1, however version 19.0.3 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.
Requirement already satisfied: cloudpickle in /usr/local/lib/python3.6/site-packages (0.7.0)
scikit-umfpack 0.3.2 has requirement numpy>=1.15.3, but you'll have numpy 1.15.1 which is
incompatible.
menpo 0.8.1 has requirement matplotlib<2.0,>=1.4, but you'll have matplotlib 2.2.3 which is
incompatible.
menpo 0.8.1 has requirement pillow<5.0,>=3.0, but you'll have pillow 5.2.0 which is incompatible.
menpo 0.8.1 has requirement scipy<1.0,>=0.16, but you'll have scipy 1.1.0 which is incompatible.
fastai 1.0.42 has requirement torch>=1.0.0, but you'll have torch 0.4.1 which is incompatible.
You are using pip version 10.0.1, however version 19.0.3 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.
Collecting folium
  Downloading
https://files.pythonhosted.org/packages/43/77/0287320dc4fd86ae8847bab6c34b5ec370e836a79c7b0c16680a3
770/folium-0.8.3-py2.py3-none-any.whl (87kB)
    100% |████████████████████████████████| 92kB 5.0MB/s
Requirement already satisfied: numpy in /usr/local/lib/python3.6/site-packages (from folium)
(1.15.1)
Requirement already satisfied: requests in /usr/local/lib/python3.6/site-packages (from folium)
(2.21.0)
Collecting branca>=0.3.0 (from folium)
  Downloading
https://files.pythonhosted.org/packages/63/36/1c93318e9653f4e414a2e0c3b98fc898b4970e939afeedeee6075
703/branca-0.3.1-py3-none-any.whl
Requirement already satisfied: six in /usr/local/lib/python3.6/site-packages (from folium)
```

```
                                                                                    --                  -       -
(1.11.0)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.6/site-packages (from folium)
(2.10)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.6/site-packages (from
requests->folium) (2018.11.29)
Requirement already satisfied: chardet<3.1.0,>=3.0.2 in /usr/local/lib/python3.6/site-packages
(from requests->folium) (3.0.4)
Requirement already satisfied: urllib3<1.25,>=1.21.1 in /usr/local/lib/python3.6/site-packages
(from requests->folium) (1.22)
Requirement already satisfied: idna<2.9,>=2.5 in /usr/local/lib/python3.6/site-packages (from
requests->folium) (2.8)
Requirement already satisfied: MarkupSafe>=0.23 in /usr/local/lib/python3.6/site-packages (from
jinja2->folium) (1.0)
scikit-umfpack 0.3.2 has requirement numpy>=1.15.3, but you'll have numpy 1.15.1 which is
incompatible.
menpo 0.8.1 has requirement matplotlib<2.0,>=1.4, but you'll have matplotlib 2.2.3 which is
incompatible.
menpo 0.8.1 has requirement pillow<5.0,>=3.0, but you'll have pillow 5.2.0 which is incompatible.
menpo 0.8.1 has requirement scipy<1.0,>=0.16, but you'll have scipy 1.1.0 which is incompatible.
fastai 1.0.42 has requirement torch>=1.0.0, but you'll have torch 0.4.1 which is incompatible.
Installing collected packages: branca, folium
Successfully installed branca-0.3.1 folium-0.8.3
You are using pip version 10.0.1, however version 19.0.3 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.
Collecting gpxpy
  Downloading
https://files.pythonhosted.org/packages/6e/d3/ce52e67771929de455e76655365a4935a2f369f76dfb0d70c20a3
463/gpxpy-1.3.5.tar.gz (105kB)
    100% |████████████████████████████████| 112kB 3.9MB/s
Building wheels for collected packages: gpxpy
  Running setup.py bdist_wheel for gpxpy ... done
  Stored in directory:
/root/.cache/pip/wheels/d2/f0/5e/b8e85979e66efec3eaa0e47fbc5274db99fd1a07befd1b2aa4
Successfully built gpxpy
scikit-umfpack 0.3.2 has requirement numpy>=1.15.3, but you'll have numpy 1.15.1 which is
incompatible.
menpo 0.8.1 has requirement matplotlib<2.0,>=1.4, but you'll have matplotlib 2.2.3 which is
incompatible.
menpo 0.8.1 has requirement pillow<5.0,>=3.0, but you'll have pillow 5.2.0 which is incompatible.
menpo 0.8.1 has requirement scipy<1.0,>=0.16, but you'll have scipy 1.1.0 which is incompatible.
fastai 1.0.42 has requirement torch>=1.0.0, but you'll have torch 0.4.1 which is incompatible.
Installing collected packages: gpxpy
Successfully installed gpxpy-1.3.5
You are using pip version 10.0.1, however version 19.0.3 is available.
You should consider upgrading via the 'pip install --upgrade pip' command.
```

In [2]:

```python
#Importing Libraries
# pip3 install graphviz
#pip3 install dask
#pip3 install toolz
#pip3 install cloudpickle
# https://www.youtube.com/watch?v=ieW3G7ZzRZ0
# https://github.com/dask/dask-tutorial
# please do go through this python notebook: https://github.com/dask/dask-
tutorial/blob/master/07_dataframe.ipynb
import dask.dataframe as dd#similar to pandas
%matplotlib inline

import pandas as pd#pandas to create small dataframes

# pip3 install foliun
# if this doesnt work refere install_folium.JPG in drive
import folium #open street map

# unix time: https://www.unixtimestamp.com/
import datetime #Convert to unix time

import time #Convert to unix time

# if numpy is not installed already : pip3 install numpy
import numpy as np#Do aritmetic operations on arrays

# matplotlib: used to plot graphs
import matplotlib
```

```
# matplotlib.use('nbagg') : matplotlib uses this protocall which makes plots more user intractive
like zoom in and zoom out
matplotlib.use('nbagg')
import matplotlib.pylab as plt
import seaborn as sns#Plots
from matplotlib import rcParams#Size of plots

# this lib is used while we calculate the stight line distance between two (lat,lon) pairs in mile
s
import gpxpy.geo #Get the haversine distance

from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
import math
import pickle
import os

# download migwin: https://mingw-w64.org/doku.php/download/mingw-builds
# install it in your system and keep the path, migw_path ='installed path'
mingw_path = 'C:\\Program Files\\mingw-w64\\x86_64-5.3.0-posix-seh-rt_v4-rev0\\mingw64\\bin'
os.environ['PATH'] = mingw_path + ';' + os.environ['PATH']

# to install xgboost: pip3 install xgboost
# if it didnt happen check install_xgboost.JPG
import xgboost as xgb

# to install sklearn: pip install -U scikit-learn
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error
from sklearn.metrics import mean_absolute_error
import warnings
warnings.filterwarnings("ignore")
```

```
/usr/local/lib/python3.6/site-packages/ipykernel_launcher.py:29: UserWarning:
This call to matplotlib.use() has no effect because the backend has already
been chosen; matplotlib.use() must be called *before* pylab, matplotlib.pyplot,
or matplotlib.backends is imported for the first time.

The backend was *originally* set to 'module://ipykernel.pylab.backend_inline' by the following cod
e:
  File "/usr/local/lib/python3.6/runpy.py", line 193, in _run_module_as_main
    "__main__", mod_spec)
  File "/usr/local/lib/python3.6/runpy.py", line 85, in _run_code
    exec(code, run_globals)
  File "/usr/local/lib/python3.6/site-packages/ipykernel_launcher.py", line 16, in <module>
    app.launch_new_instance()
  File "/usr/local/lib/python3.6/site-packages/traitlets/config/application.py", line 658, in
launch_instance
    app.start()
  File "/usr/local/lib/python3.6/site-packages/ipykernel/kernelapp.py", line 486, in start
    self.io_loop.start()
  File "/usr/local/lib/python3.6/site-packages/tornado/platform/asyncio.py", line 132, in start
    self.asyncio_loop.run_forever()
  File "/usr/local/lib/python3.6/asyncio/base_events.py", line 422, in run_forever
    self._run_once()
  File "/usr/local/lib/python3.6/asyncio/base_events.py", line 1432, in _run_once
    handle._run()
  File "/usr/local/lib/python3.6/asyncio/events.py", line 145, in _run
    self._callback(*self._args)
  File "/usr/local/lib/python3.6/site-packages/tornado/platform/asyncio.py", line 122, in
_handle_events
    handler_func(fileobj, events)
  File "/usr/local/lib/python3.6/site-packages/tornado/stack_context.py", line 300, in
null_wrapper
    return fn(*args, **kwargs)
  File "/usr/local/lib/python3.6/site-packages/zmq/eventloop/zmqstream.py", line 450, in
_handle_events
    self._handle_recv()
  File "/usr/local/lib/python3.6/site-packages/zmq/eventloop/zmqstream.py", line 480, in
_handle_recv
    self._run_callback(callback, msg)
  File "/usr/local/lib/python3.6/site-packages/zmq/eventloop/zmqstream.py", line 432, in
_run_callback
    callback(*args, **kwargs)
  File "/usr/local/lib/python3.6/site-packages/tornado/stack_context.py", line 300, in
null_wrapper
    return fn(*args, **kwargs)
```

```
  File "/usr/local/lib/python3.6/site-packages/ipykernel/kernelbase.py", line 283, in dispatcher
    return self.dispatch_shell(stream, msg)
  File "/usr/local/lib/python3.6/site-packages/ipykernel/kernelbase.py", line 233, in
dispatch_shell
    handler(stream, idents, msg)
  File "/usr/local/lib/python3.6/site-packages/ipykernel/kernelbase.py", line 399, in
execute_request
    user_expressions, allow_stdin)
  File "/usr/local/lib/python3.6/site-packages/ipykernel/ipkernel.py", line 208, in do_execute
    res = shell.run_cell(code, store_history=store_history, silent=silent)
  File "/usr/local/lib/python3.6/site-packages/ipykernel/zmqshell.py", line 537, in run_cell
    return super(ZMQInteractiveShell, self).run_cell(*args, **kwargs)
  File "/usr/local/lib/python3.6/site-packages/IPython/core/interactiveshell.py", line 2662, in ru
n_cell
    raw_cell, store_history, silent, shell_futures)
  File "/usr/local/lib/python3.6/site-packages/IPython/core/interactiveshell.py", line 2785, in _r
un_cell
    interactivity=interactivity, compiler=compiler, result=result)
  File "/usr/local/lib/python3.6/site-packages/IPython/core/interactiveshell.py", line 2901, in ru
n_ast_nodes
    if self.run_code(code, result):
  File "/usr/local/lib/python3.6/site-packages/IPython/core/interactiveshell.py", line 2961, in ru
n_code
    exec(code_obj, self.user_global_ns, self.user_ns)
  File "<ipython-input-2-dcdf77e8f017>", line 10, in <module>
    get_ipython().run_line_magic('matplotlib', 'inline')
  File "/usr/local/lib/python3.6/site-packages/IPython/core/interactiveshell.py", line 2131, in ru
n_line_magic
    result = fn(*args,**kwargs)
  File "<decorator-gen-107>", line 2, in matplotlib
  File "/usr/local/lib/python3.6/site-packages/IPython/core/magic.py", line 187, in <lambda>
    call = lambda f, *a, **k: f(*a, **k)
  File "/usr/local/lib/python3.6/site-packages/IPython/core/magics/pylab.py", line 99, in
matplotlib
    gui, backend = self.shell.enable_matplotlib(args.gui)
  File "/usr/local/lib/python3.6/site-packages/IPython/core/interactiveshell.py", line 3049, in en
able_matplotlib
    pt.activate_matplotlib(backend)
  File "/usr/local/lib/python3.6/site-packages/IPython/core/pylabtools.py", line 311, in
activate_matplotlib
    matplotlib.pyplot.switch_backend(backend)
  File "/usr/local/lib/python3.6/site-packages/matplotlib/pyplot.py", line 231, in switch_backend
    matplotlib.use(newbackend, warn=False, force=True)
  File "/usr/local/lib/python3.6/site-packages/matplotlib/__init__.py", line 1422, in use
    reload(sys.modules['matplotlib.backends'])
  File "/usr/local/lib/python3.6/importlib/__init__.py", line 166, in reload
    _bootstrap._exec(spec, module)
  File "/usr/local/lib/python3.6/site-packages/matplotlib/backends/__init__.py", line 16, in
<module>
    line for line in traceback.format_stack()


/usr/local/lib/python3.6/site-packages/sklearn/ensemble/weight_boosting.py:29: DeprecationWarning:
numpy.core.umath_tests is an internal NumPy module and should not be imported. It will be removed
in a future NumPy release.
  from numpy.core.umath_tests import inner1d
```

# Data Information

Source of Data: Data can be downloaded from here: http://www.nyc.gov/html/tlc/html/about/trip_record_data.shtml. Here, we have used Jan- 2015 and Jan- 2016 data.


## Information on taxis:

*Yellow Taxi: Yellow Medallion Taxicabs*

These are the famous NYC yellow taxis that provide transportation exclusively through street-hails. The number of taxicabs is limited by a finite number of medallions issued by the TLC. You access this mode of transportation by standing in the street and hailing an available taxi with your hand. The pickups are not pre-arranged.

*For Hire Vehicles (FHVs)*

*For Hire Vehicles (FHVs)*

FHV transportation is accessed by a pre-arrangement with a dispatcher or limo company. These FHVs are not permitted to pick up passengers via street hails, as those rides are not considered pre-arranged.

*Green Taxi: Street Hail Livery (SHL)*

The SHL program will allow livery vehicle owners to license and outfit their vehicles with green borough taxi branding, meters, credit card machines, and ultimately the right to accept street hails in addition to pre-arranged rides.

Credits: Quora

*Footnote:*
In the given notebook we are considering only the yellow taxis for the time period between Jan - Mar 2015 & Jan - Mar 2016

# Data Collection

We Have collected all yellow taxi trips data from jan-2015 to dec-2016(Will be using only 2015 data)

| file name | file name size | number of records | number of features |
| --- | --- | --- | --- |
| yellow_tripdata_2016-01 | 1. 59G | 10906858 | 19 |
| yellow_tripdata_2016-02 | 1. 66G | 11382049 | 19 |
| yellow_tripdata_2016-03 | 1. 78G | 12210952 | 19 |
| yellow_tripdata_2016-04 | 1. 74G | 11934338 | 19 |
| yellow_tripdata_2016-05 | 1. 73G | 11836853 | 19 |
| yellow_tripdata_2016-06 | 1. 62G | 11135470 | 19 |
| yellow_tripdata_2016-07 | 884Mb | 10294080 | 17 |
| yellow_tripdata_2016-08 | 854Mb | 9942263 | 17 |
| yellow_tripdata_2016-09 | 870Mb | 10116018 | 17 |
| yellow_tripdata_2016-10 | 933Mb | 10854626 | 17 |
| yellow_tripdata_2016-11 | 868Mb | 10102128 | 17 |
| yellow_tripdata_2016-12 | 897Mb | 10449408 | 17 |
| yellow_tripdata_2015-01 | 1.84Gb | 12748986 | 19 |
| yellow_tripdata_2015-02 | 1.81Gb | 12450521 | 19 |
| yellow_tripdata_2015-03 | 1.94Gb | 13351609 | 19 |
| yellow_tripdata_2015-04 | 1.90Gb | 13071789 | 19 |
| yellow_tripdata_2015-05 | 1.91Gb | 13158262 | 19 |
| yellow_tripdata_2015-06 | 1.79Gb | 12324935 | 19 |
| yellow_tripdata_2015-07 | 1.68Gb | 11562783 | 19 |
| yellow_tripdata_2015-08 | 1.62Gb | 11130304 | 19 |
| yellow_tripdata_2015-09 | 1.63Gb | 11225063 | 19 |
| yellow_tripdata_2015-10 | 1.79Gb | 12315488 | 19 |
| yellow_tripdata_2015-11 | 1.65Gb | 11312676 | 19 |
| yellow_tripdata_2015-12 | 1.67Gb | 11460573 | 19 |

In [3]:

```
#Looking at the features
# dask dataframe  : # https://github.com/dask/dask-tutorial/blob/master/07_dataframe.ipynb
month = dd.read_csv('/floyd/input/nyc/yellow_tripdata_2015-01.csv')
print(month.columns)
```

```
Index(['VendorID', 'tpep_pickup_datetime', 'tpep_dropoff_datetime',
       'passenger_count', 'trip_distance', 'pickup_longitude',
       'pickup_latitude', 'RateCodeID', 'store_and_fwd_flag',
       'dropoff_longitude', 'dropoff_latitude', 'payment_type', 'fare_amount',
       'extra', 'mta_tax', 'tip_amount', 'tolls_amount',
       'improvement_surcharge', 'total_amount'],
      dtype='object')
```

In [4]:

```
# However unlike Pandas, operations on dask.dataframes don't trigger immediate computation,
# instead they add key-value pairs to an underlying Dask graph. Recall that in the diagram below,
# circles are operations and rectangles are results.

# to see the visulaization you need to install graphviz
# pip3 install graphviz if this doesnt work please check the install_graphviz.jpg in the drive
month.visualize()
```
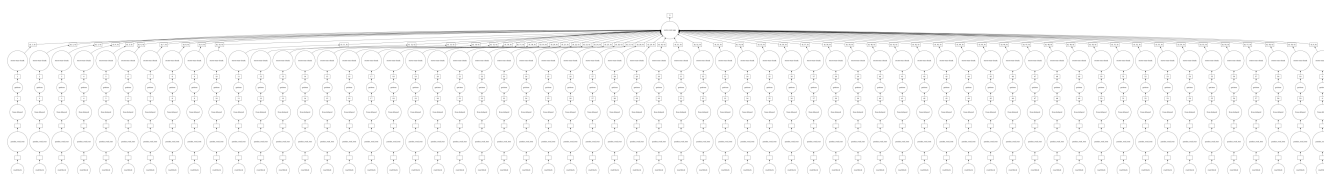
Out[4]:



In [5]:

```
month.fare_amount.sum().visualize()
```

Out[5]:



# Features in the dataset:

| Field Name | | Description |
|---|---|---|
| VendorID | 1. <br> 2. | A code indicating the TPEP provider that provided the record. <br> Creative Mobile Technologies <br> VeriFone Inc. |
| tpep_pickup_datetime | | The date and time when the meter was engaged. |
| tpep_dropoff_datetime | | The date and time when the meter was disengaged. |
| Passenger_count | | The number of passengers in the vehicle. This is a driver-entered value. |
| Trip_distance | | The elapsed trip distance in miles reported by the taximeter. |
| Pickup_longitude | | Longitude where the meter was engaged. |
| Pickup_latitude | | Latitude where the meter was engaged. |
| RateCodeID | 1. <br> 2. <br> 3. <br> 4. <br> 5. <br> 6. | The final rate code in effect at the end of the trip. <br> Standard rate <br> JFK <br> Newark <br> Nassau or Westchester <br> Negotiated fare <br> Group ride |
| Store_and_fwd_flag | | This flag indicates whether the trip record was held in vehicle memory before sending to the vendor, aka "store and forward," because the vehicle did not have a connection to the server. Y= store and forward trip N= not a store and forward trip |
| Dropoff_longitude | | Longitude where the meter was disengaged. |
| Dropoff_ latitude | | Latitude where the meter was disengaged. |
| Payment_type | 1. <br> 2. <br> 3. <br> 4. <br> 5. <br> 6. | A numeric code signifying how the passenger paid for the trip. <br> Credit card <br> Cash <br> No charge <br> Dispute <br> Unknown <br> Voided trip |
| Fare_amount | | The time-and-distance fare calculated by the meter. |
| Extra | | Miscellaneous extras and surcharges. Currently, this only includes. the $0.50$ and $1$ rush hour and overnight charges. |
| MTA_tax | | 0.50 MTA tax that is automatically triggered based on the metered rate in use. |
| Improvement_surcharge | | 0.30 improvement surcharge assessed trips at the flag drop. the improvement surcharge began being levied in 2015. |

| | |
|---|---|
| Tip_amount | Tip amount – This field is automatically populated for credit card tips.Cash tips are not included. |
| Tolls_amount | Total amount of all tolls paid in trip. |
| Total_amount | The total amount charged to passengers. Does not include cash tips. |

# ML Problem Formulation

**Time-series forecasting and Regression**

- *To find number of pickups, given location cordinates(latitude and longitude) and time, in the query reigion and surrounding regions.*

To solve the above we would be using data collected in Jan - Mar 2015 to predict the pickups in Jan - Mar 2016.

# Performance metrics

1. Mean Absolute percentage error.
2. Mean Squared error.

# Data Cleaning

In this section we will be doing univariate analysis and removing outlier/illegitimate values which may be caused due to some error

In [6]:

```
#table below shows few datapoints along with all our features
month.head(5)
```

Out[6]:

| | VendorID | tpep_pickup_datetime | tpep_dropoff_datetime | passenger_count | trip_distance | pickup_longitude | pickup_latitude | RateCode |
|---|---|---|---|---|---|---|---|---|
| **0** | 2 | 2015-01-15 19:05:39 | 2015-01-15 19:23:42 | 1 | 1.59 | -73.993896 | 40.750111 | |
| **1** | 1 | 2015-01-10 20:33:38 | 2015-01-10 20:53:28 | 1 | 3.30 | -74.001648 | 40.724243 | |
| **2** | 1 | 2015-01-10 20:33:38 | 2015-01-10 20:43:41 | 1 | 1.80 | -73.963341 | 40.802788 | |
| **3** | 1 | 2015-01-10 20:33:39 | 2015-01-10 20:35:31 | 1 | 0.50 | -74.009087 | 40.713818 | |
| **4** | 1 | 2015-01-10 20:33:39 | 2015-01-10 20:52:58 | 1 | 3.00 | -73.971176 | 40.762428 | |

## 1. Pickup Latitude and Pickup Longitude

It is inferred from the source https://www.flickr.com/places/info/2459115 that New York is bounded by the location cordinates(lat,long) - (40.5774, -74.15) & (40.9176,-73.7004) so hence any cordinates not within these cordinates are not considered by us as we are only concerned with pickups which originate within New York.

In [8]:

```
# Plotting pickup cordinates which are outside the bounding box of New-York
# we will collect all the points outside the bounding box of newyork city to outlier_locations
outlier_locations = month[((month.pickup_longitude <= -74.15) | (month.pickup_latitude <= 40.5774)| \
                    (month.pickup_longitude >= -73.7004) | (month.pickup_latitude >= 40.9176))]

# creating a map with the a base location
# read more about the folium here: http://folium.readthedocs.io/en/latest/quickstart.html

# note: you dont need to remember any of these, you dont need indeepth knowledge on these maps and
plots

map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')

# we will spot only first 100 outliers on the map, plotting all the outliers will take more time
sample_locations = outlier_locations.head(10000)
for i,j in sample_locations.iterrows():
    if int(j['pickup_latitude']) != 0:
        folium.Marker(list((j['pickup_latitude'],j['pickup_longitude']))).add_to(map_osm)
```

**Observation:-** As you can see above that there are some points just outside the boundary but there are a few that are in either South america, Mexico or Canada
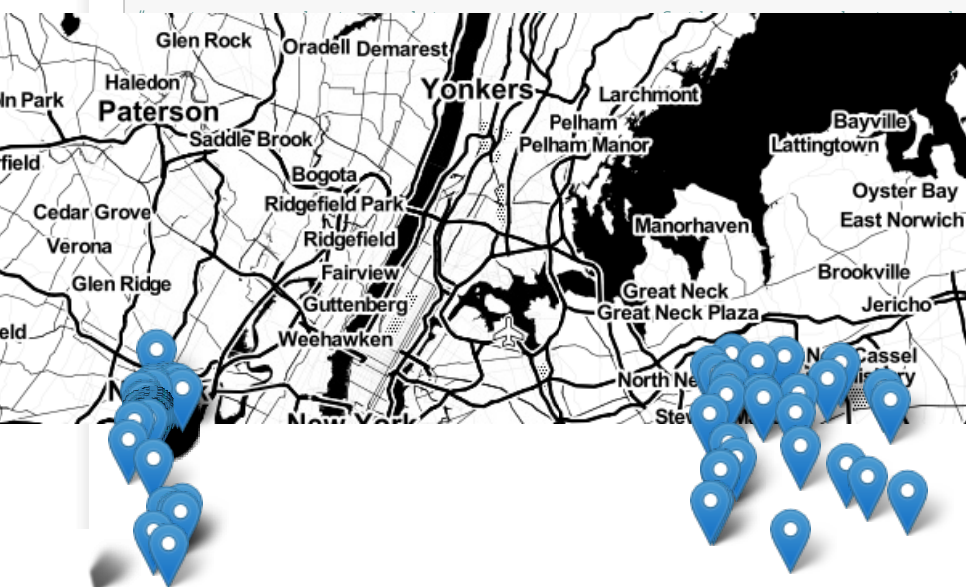
## 2. Dropoff Latitude & Dropoff Longitude

It is inferred from the source https://www.flickr.com/places/info/2459115 that New York is bounded by the location cordinates(lat,long) - (40.5774, -74.15) & (40.9176,-73.7004) so hence any cordinates not within these cordinates are not considered by us as we are only concerned with dropoffs which are within New York.

In [9]:

```
# Plotting dropoff cordinates which are outside the bounding box of New-York
# we will collect all the points outside the bounding box of newyork city to outlier_locations
outlier_locations = month[((month.dropoff_longitude <= -74.15) | (month.dropoff_latitude <= 40.5774
)| \
                (month.dropoff_longitude >= -73.7004) | (month.dropoff_latitude >= 40.9176))]

# creating a map with the a base location
# read more about the folium here: http://folium.readthedocs.io/en/latest/quickstart.html
```

indeepth knowledge on these maps and

tamen Toner')

l the outliers will take more time

ngitude']))).add_to(map_osm)

**Observation:-** The observations here are similar to those obtained while analysing pickup latitude and longitude

### 3. Trip Durations:

According to NYC Taxi & Limousion Commision Regulations **the maximum allowed trip duration in a 24 hour interval is 12 hours.**

In [148]:

```python
#The timestamps are converted to unix so as to get duration(trip-time) & speed also pickup-times in unix are used while binning

# in out data we have time in the formate "YYYY-MM-DD HH:MM:SS" we convert thiss sting to python time formate and then into unix time stamp
# https://stackoverflow.com/a/27914405
def convert_to_unix(s):
    return time.mktime(datetime.datetime.strptime(s, "%Y-%m-%d %H:%M:%S").timetuple())


# we return a data frame which contains the columns
# 1.'passenger_count' : self explanatory
# 2.'trip_distance' : self explanatory
# 3.'pickup_longitude' : self explanatory
# 4.'pickup_latitude' : self explanatory
# 5.'dropoff_longitude' : self explanatory
# 6.'dropoff_latitude' : self explanatory
# 7.'total_amount' : total fair that was paid
# 8.'trip_times' : duration of each trip
# 9.'pickup_times : pickup time converted into unix time
# 10.'Speed' : velocity of each trip
def return_with_trip_times(month):
    duration = month[['tpep_pickup_datetime','tpep_dropoff_datetime']].compute()
    #pickups and dropoffs to unix time
    duration_pickup = [convert_to_unix(x) for x in duration['tpep_pickup_datetime'].values]
    duration_drop = [convert_to_unix(x) for x in duration['tpep_dropoff_datetime'].values]
    #calculate duration of trips
    durations = (np.array(duration_drop) - np.array(duration_pickup))/float(60)

    #append durations of trips and speed in miles/hr to a new dataframe
    #new_frame =
month[['passenger_count','trip_distance','pickup_longitude','pickup_latitude','dropoff_longitude',
off_latitude','total_amount']].compute()
    new_frame =
month[['VendorID','payment_type','passenger_count','trip_distance','pickup_longitude','pickup_latit
ude','dropoff_longitude','dropoff_latitude','RateCodeID','store_and_fwd_flag','total_amount']].com
pute()

    new_frame['trip_times'] = durations
```

```
    new_frame['pickup_times'] = duration_pickup
    new_frame['Speed'] = 60*(new_frame['trip_distance']/new_frame['trip_times'])

    return new_frame

# print(frame_with_durations.head())
#  passenger_count trip_distance pickup_longitude pickup_latitude dropoff_longitude
dropoff_latitude total_amount trip_times pickup_times Speed
#   1                1.59      -73.993896        40.750111    -73.974785         40.750618
17.05    18.050000 1.421329e+09 5.285319
#   1                3.30      -74.001648        40.724243    -73.994415         40.759109
.80    19.833333 1.420902e+09 9.983193
#   1                1.80      -73.963341        40.802788    -73.951820         40.824413
10.80    10.050000 1.420902e+09 10.746269
#   1                0.50      -74.009087        40.713818    -74.004326         40.719986
4.80     1.866667 1.420902e+09 16.071429
#   1                3.00      -73.971176        40.762428    -74.004181         40.742653
6.30    19.316667 1.420902e+09 9.318378
frame_with_durations = return_with_trip_times(month)
```

In [12]:

```
# the skewed box plot shows us the presence of outliers
sns.boxplot(y="trip_times", data =frame_with_durations)
plt.show()
```



In [13]:

```
#calculating 0-100th percentile to find a the correct percentile value for removal of outliers
for i in range(0,100,10):
    var =frame_with_durations["trip_times"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print ("100 percentile value is ",var[-1])
```

```
0 percentile value is -1211.0166666666667
10 percentile value is 3.8333333333333335
20 percentile value is 5.383333333333334
30 percentile value is 6.816666666666666
40 percentile value is 8.3
50 percentile value is 9.95
60 percentile value is 11.866666666666667
70 percentile value is 14.283333333333333
80 percentile value is 17.633333333333333
90 percentile value is 23.45
100 percentile value is  548555.6333333333
```

In [14]:

```
#looking further from the 99th percecntile
for i in range(90,100):
    var =frame_with_durations["trip_times"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print ("100 percentile value is ",var[-1])
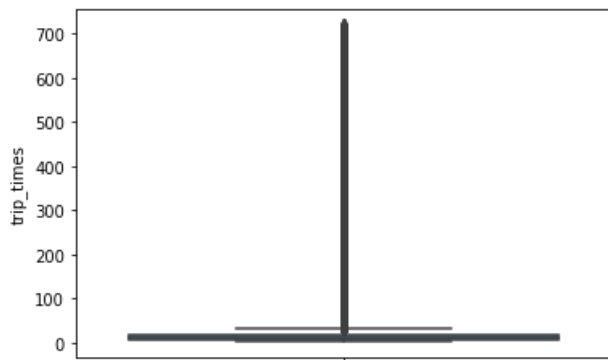```

```
90 percentile value is 23.45
```

```
91 percentile value is 24.35
92 percentile value is 25.383333333333333
93 percentile value is 26.55
94 percentile value is 27.933333333333334
95 percentile value is 29.583333333333332
96 percentile value is 31.683333333333334
97 percentile value is 34.46666666666667
98 percentile value is 38.71666666666667
99 percentile value is 46.75
100 percentile value is  548555.6333333333
```

In [149]:

```python
#removing data based on our analysis and TLC regulations
frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_times>1) &
(frame_with_durations.trip_times<720)]
```

In [16]:

```python
#box-plot after removal of outliers
sns.boxplot(y="trip_times", data =frame_with_durations_modified)
plt.show()
```



In [17]:

```python
#pdf of trip-times after removing the outliers
sns.FacetGrid(frame_with_durations_modified,size=6) \
      .map(sns.kdeplot,"trip_times") \
      .add_legend();
plt.show();
```

Observation: Above PDF plot shows that almost all of the trip durations are very less and approximately less than 100, extremely few trip durations are above 100.

In [18]:

```python
#converting the values to log-values to chec for log-normal
import math
frame_with_durations_modified['log_times']=[math.log(i) for i in frame_with_durations_modified['trip_times'].values]
```
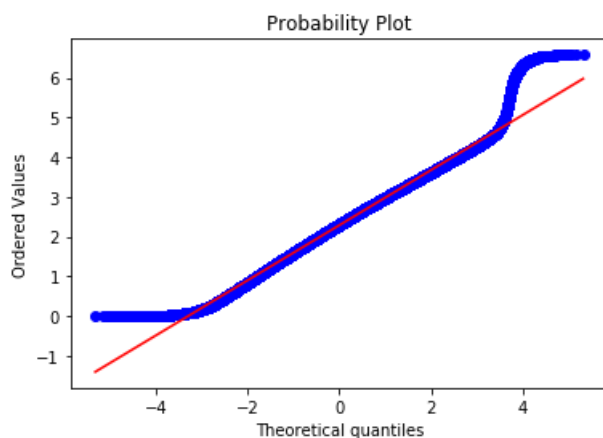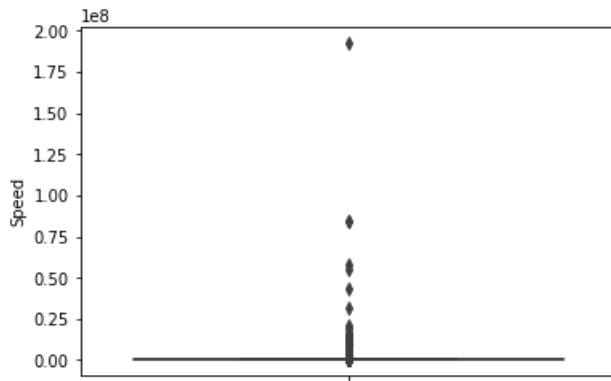
In [19]:

```python
#pdf of log-values
sns.FacetGrid(frame_with_durations_modified,size=6) \
       .map(sns.kdeplot,"log_times") \
       .add_legend();
plt.show();
```



In [20]:

```python
#Q-Q plot for checking if trip-times is log-normal
import scipy
scipy.stats.probplot(frame_with_durations_modified['log_times'].values, plot=plt)
plt.show()
```



## 4. Speed

In [21]:

```
# check for any outliers in the data after trip duration outliers removed
# box-plot for speeds with outliers
frame_with_durations_modified['Speed'] =
60*(frame_with_durations_modified['trip_distance']/frame_with_durations_modified['trip_times'])
sns.boxplot(y="Speed", data =frame_with_durations_modified)
plt.show()
```

```
#calculating speed values at each percntile 0,10,20,30,40,50,60,70,80,90,100
for i in range(0,100,10):
    var =frame_with_durations_modified["Speed"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print("100 percentile value is ",var[-1])
```

```
0 percentile value is 0.0
10 percentile value is 6.409495548961425
20 percentile value is 7.80952380952381
30 percentile value is 8.929133858267717
40 percentile value is 9.98019801980198
50 percentile value is 11.06865671641791
60 percentile value is 12.286689419795222
70 percentile value is 13.796407185628745
80 percentile value is 15.963224893917962
90 percentile value is 20.186915887850468
100 percentile value is  192857142.85714284
```

```
#calculating speed values at each percntile 90,91,92,93,94,95,96,97,98,99,100
for i in range(90,100):
    var =frame_with_durations_modified["Speed"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print("100 percentile value is ",var[-1])
```

```
90 percentile value is 20.186915887850468
91 percentile value is 20.91645569620253
92 percentile value is 21.752988047808763
93 percentile value is 22.721893491124263
94 percentile value is 23.844155844155843
95 percentile value is 25.182552504038775
96 percentile value is 26.80851063829787
97 percentile value is 28.84304932735426
98 percentile value is 31.591128254580514
99 percentile value is 35.7513566847558
100 percentile value is  192857142.85714284
```

```
#calculating speed values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100
for i in np.arange(0.0, 1.0, 0.1):
    var =frame_with_durations_modified["Speed"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
```

```
print("100 percentile value is ",var[-1])
```

```
99.0 percentile value is 35.7513566847558
99.1 percentile value is 36.31084727468969
99.2 percentile value is 36.91470054446461
99.3 percentile value is 37.588235294117645
99.4 percentile value is 38.33035714285714
99.5 percentile value is 39.17580340264651
99.6 percentile value is 40.15384615384615
99.7 percentile value is 41.338301043219076
99.8 percentile value is 42.86631016042781
99.9 percentile value is 45.3107822410148
100 percentile value is  192857142.85714284
```

Observations: Here, 100th percentile value of a speed is 192 Million miles/hr which is (BIZZARE). Furthermore, 99.9th percentile value of speed is 45.31miles/hr. So, we are removing all the data points where speed is greater than 45.31miles/hr.

In [150]:

```
#removing further outliers based on the 99.9th percentile value
frame_with_durations_modified=frame_with_durations[(frame_with_durations.Speed>0) &
(frame_with_durations.Speed<45.31)]
```

In [27]:

```
### Box plot of speed after removing outliers and erroneous points.
fig = plt.figure(figsize = (10,6))
ax = sns.boxplot("Speed", data = frame_with_durations_modified, orient = "v")

plt.tick_params(labelsize = 20)
plt.ylabel("Speed(Miles/hr)", fontsize = 20)
plt.show()
```
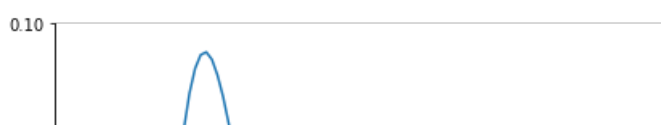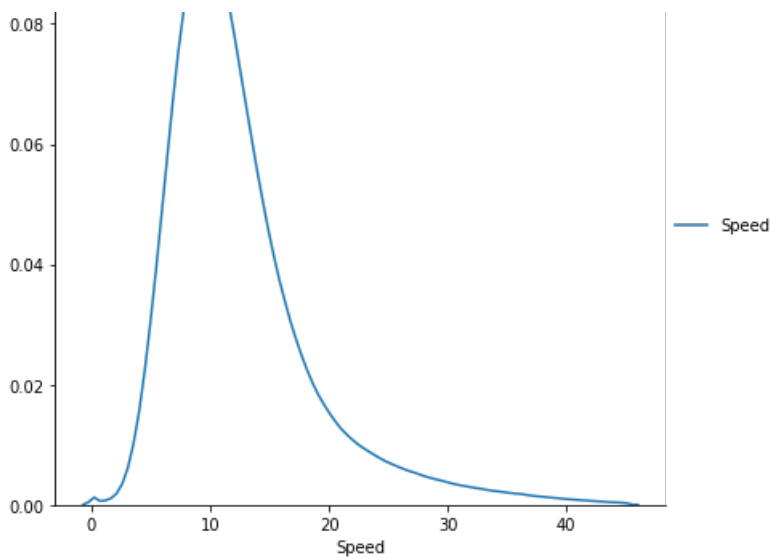


In [31]:

```
def dist_of_params(frame,variable,title):
    sns.FacetGrid(frame,size=6) \
      .map(sns.kdeplot, variable) \
      .add_legend()
    plt.show()
```
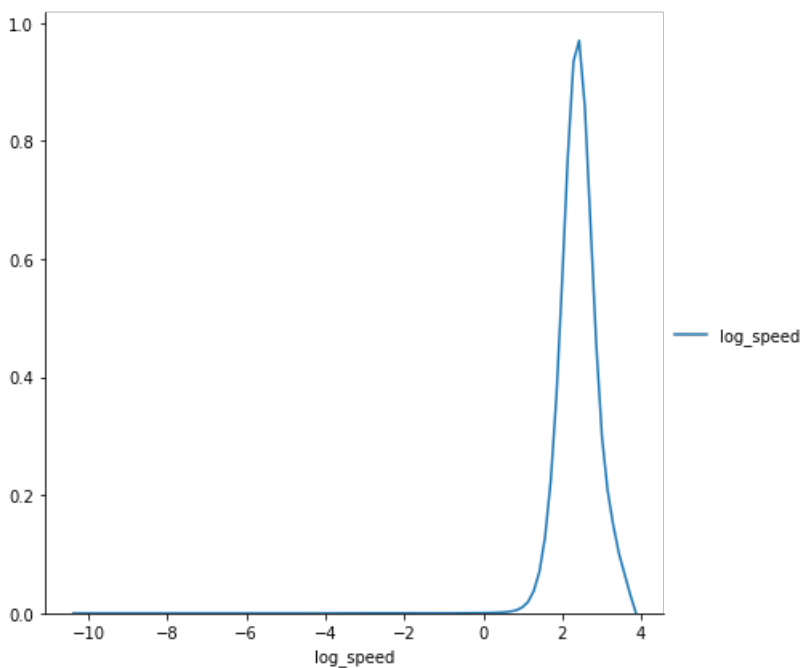
In [32]:

```
#trip speed
dist_of_params(frame_with_durations_modified,'Speed','average Speed of cab trips distribution')
```

```python
#log speed
log_trip_speed = frame_with_durations_modified.Speed.values
frame_with_durations_modified['log_speed'] = np.log(log_trip_speed)
dist_of_params(frame_with_durations_modified,'log_speed','log of speed for cab trips distribution'
)
```

```python
#avg.speed of cabs in New-York
Average_speed = sum(frame_with_durations_modified["Speed"]) /
float(len(frame_with_durations_modified["Speed"]))
Average_speed
```

Out[35]:

12.450173996027528

```python
print("Speed of Taxis around NYC per 10 minutes = "+str(Average_speed/6)+" per 10 minutes.")
```
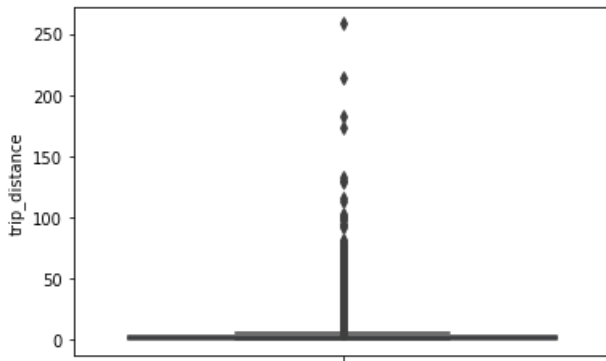
Speed of Taxis around NYC per 10 minutes = 2.0750289993379214 per 10 minutes.

**The avg speed in Newyork speed is 12.45miles/hr, so a cab driver can travel 2 miles per 10min on avg.**

## 4. Trip Distance

In [37]:

```python
# up to now we have removed the outliers based on trip durations and cab speeds
# lets try if there are any outliers in trip distances
# box-plot showing outliers in trip-distance values
sns.boxplot(y="trip_distance", data =frame_with_durations_modified)
plt.show()
```



In [38]:

```python
#calculating trip distance values at each percntile 0,10,20,30,40,50,60,70,80,90,100
for i in range(0,100,10):
    var =frame_with_durations_modified["trip_distance"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print("100 percentile value is ",var[-1])
```

```
0 percentile value is 0.01
10 percentile value is 0.66
20 percentile value is 0.9
30 percentile value is 1.1
40 percentile value is 1.39
50 percentile value is 1.69
60 percentile value is 2.07
70 percentile value is 2.6
80 percentile value is 3.6
90 percentile value is 5.97
100 percentile value is  258.9
```

In [39]:

```python
#calculating trip distance values at each percntile 90,91,92,93,94,95,96,97,98,99,100
for i in range(90,100):
    var =frame_with_durations_modified["trip_distance"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print("100 percentile value is ",var[-1])
```

```
90 percentile value is 5.97
91 percentile value is 6.45
92 percentile value is 7.07
93 percentile value is 7.85
94 percentile value is 8.72
95 percentile value is 9.6
96 percentile value is 10.6
97 percentile value is 12.1
98 percentile value is 16.03
99 percentile value is 18.17
100 percentile value is  258.9
```

```
#calculating trip distance values at each percntile
99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100
for i in np.arange(0.0, 1.0, 0.1):
    var =frame_with_durations_modified["trip_distance"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
print("100 percentile value is ",var[-1])
```

```
99.0 percentile value is 18.17
99.1 percentile value is 18.37
99.2 percentile value is 18.6
99.3 percentile value is 18.83
99.4 percentile value is 19.13
99.5 percentile value is 19.5
99.6 percentile value is 19.96
99.7 percentile value is 20.5
99.8 percentile value is 21.22
99.9 percentile value is 22.57
100 percentile value is  258.9
```
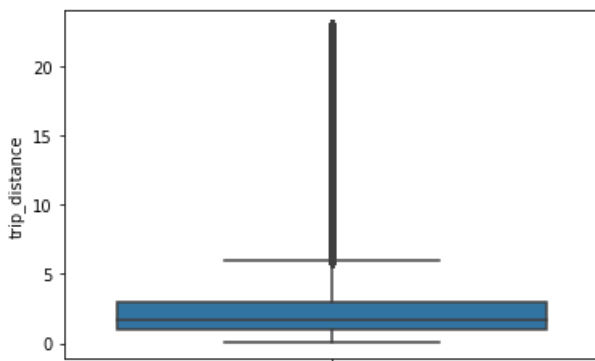
Observation: Here, 99.9th percentile of trip distance is 22.58miles, however, 100th percentile value is 258.9miles, which is very high. So, we are removing all the data points where trip distance is greater than 23miles.

```
#removing further outliers based on the 99.9th percentile value
frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_distance>0) &
(frame_with_durations.trip_distance<23)]
```
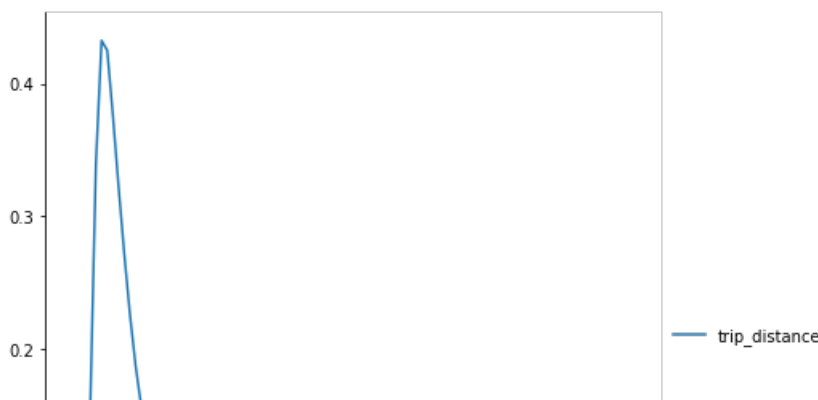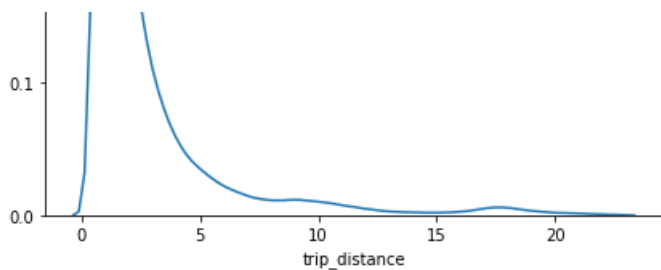
```
#box-plot after removal of outliers
sns.boxplot(y="trip_distance", data = frame_with_durations_modified)
plt.show()
```
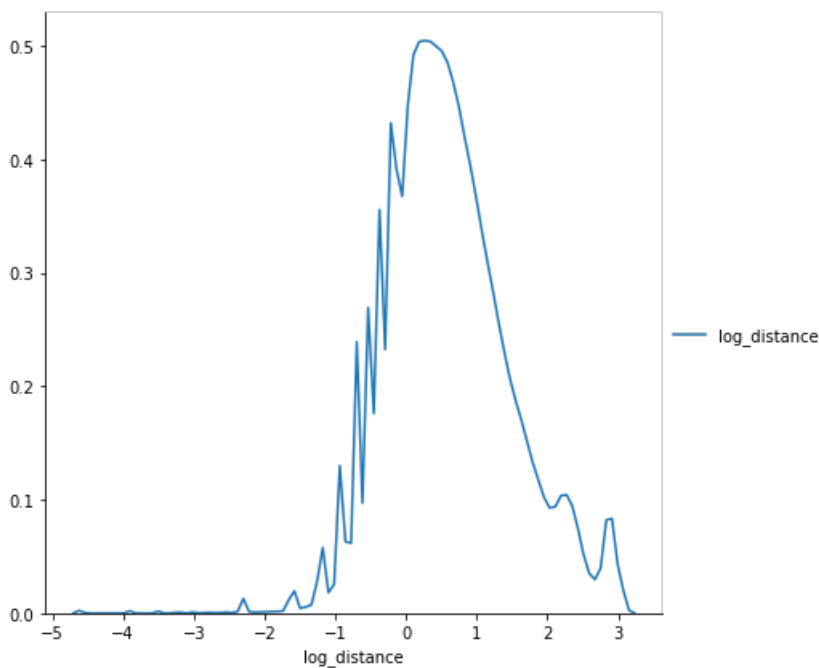
```
#trip distances
dist_of_params(frame_with_durations_modified,'trip_distance','distance for cab trips distribution'
)
```

```
#log trip distances
log_trip_distance = frame_with_durations_modified.trip_distance.values
frame_with_durations_modified['log_distance'] = np.log(log_trip_distance)
dist_of_params(frame_with_durations_modified,'log_distance','log of distance for cab trips
distribution')
```



**Remove all outliers/erronous points.**

```
#removing all outliers based on our univariate analysis above
def remove_outliers(new_frame):


    a = new_frame.shape[0]
    print ("Number of pickup records = ",a)
    temp_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_frame.dropoff_longitude
<= -73.7004) &\
                    (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitude <=
40.9176)) & \
                    ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_latitude >=
40.5774)& \
                    (new_frame.pickup_longitude <= -73.7004) & (new_frame.pickup_latitude <=
40.9176))]
    b = temp_frame.shape[0]
    print ("Number of outlier coordinates lying outside NY boundaries:",(a-b))


    temp_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times < 720)]
    c = temp_frame.shape[0]
    print ("Number of outliers from trip times analysis:",(a-c))


    temp_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_distance < 23)]
    d = temp_frame.shape[0]
```

```
      print ("Number of outliers from trip distance analysis:",(a-d))

      temp_frame = new_frame[(new_frame.Speed <= 65) & (new_frame.Speed >= 0)]
      e = temp_frame.shape[0]
      print ("Number of outliers from speed analysis:",(a-e))

      temp_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_amount >0)]
      f = temp_frame.shape[0]
      print ("Number of outliers from fare analysis:",(a-f))


      new_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_frame.dropoff_longitude <
= -73.7004) &\
                      (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitude <=
40.9176)) & \
                      ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_latitude >=
40.5774)& \
                      (new_frame.pickup_longitude <= -73.7004) & (new_frame.pickup_latitude <=
40.9176))]

      new_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times < 720)]
      new_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_distance < 23)]
      new_frame = new_frame[(new_frame.Speed < 45.31) & (new_frame.Speed > 0)]
      new_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_amount >0)]

      print ("Total outliers removed",a - new_frame.shape[0])
      print ("---")
      return new_frame
```

In [153]:

```
print ("Removing outliers in the month of Jan-2015")
print ("----")
frame_with_durations_outliers_removed = remove_outliers(frame_with_durations)
print("fraction of data points that remain after removing outliers",
float(len(frame_with_durations_outliers_removed))/len(frame_with_durations))
```

```
Removing outliers in the month of Jan-2015
----
Number of pickup records =  12748986
Number of outlier coordinates lying outside NY boundaries: 293919
Number of outliers from trip times analysis: 23889
Number of outliers from trip distance analysis: 92597
Number of outliers from speed analysis: 24473
Number of outliers from fare analysis: 5275
Total outliers removed 377910
---
fraction of data points that remain after removing outliers 0.9703576425607495
```

## 5. Vendor ID

In [154]:

```
#Get unique vendors
vendor_unique = frame_with_durations_outliers_removed .VendorID.unique()
Vendors = []
#Get the number of trips for each vendor
Vendors.append(frame_with_durations_outliers_removed[frame_with_durations_outliers_removed
['VendorID'] == 1].shape[0])
Vendors.append(frame_with_durations_outliers_removed[frame_with_durations_outliers_removed
['VendorID'] == 2].shape[0])
```

In [156]:

```
#Names of vendors
print("Unique vendors\n {}: Creative Mobile Technologies\n {}: VeriFone
Inc\n".format(vendor_unique[1],vendor_unique[0]))
print ("Vendor 1 and Vendor 2 count:\n",Vendors)
```

```
Unique vendors
 1: Creative Mobile Technologies
```
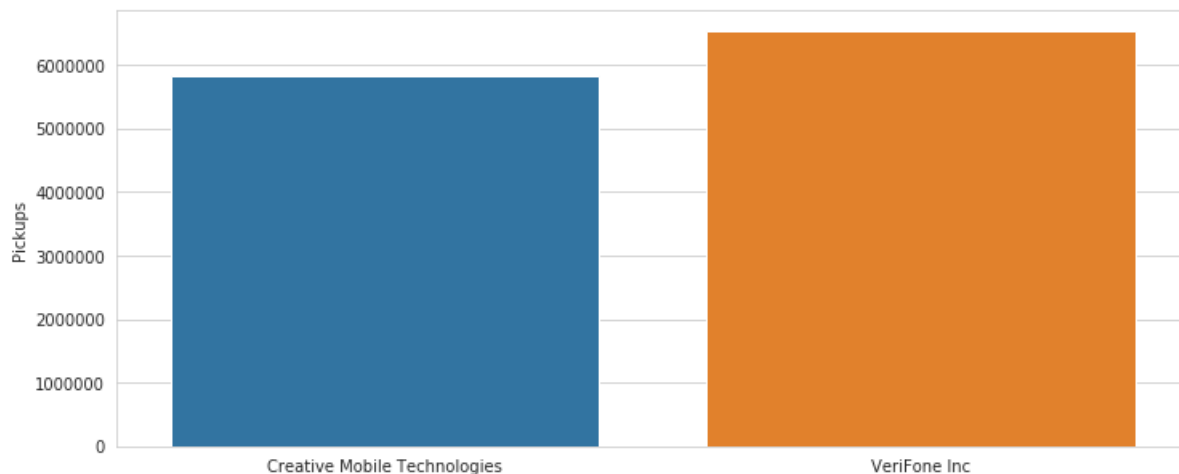
```
  2: VeriFone Inc

Vendor 1 and Vendor 2 count:
 [5837321, 6533755]
```

```python
sns.set_style("whitegrid")

fig, ax = plt.subplots()
fig.set_size_inches(12, 5)

ax = sns.barplot(x = ['Creative Mobile Technologies', 'VeriFone Inc'] , y = np.array(Vendors))
ax.set(ylabel='Pickups')
plt.show()
```

```python
print ("Creative Mobile Technologies share :{} % VeriFone Inc share :{} % ".format(100*float(Vendo
rs[1])/sum(Vendors),100*float(Vendors[0])/sum(Vendors)))
```

```
Creative Mobile Technologies share :52.814767284591895 % VeriFone Inc share :47.185232715408105 %
```

Observation: Both the vendors share almost same number of pickups in JAN month

## 6. Passenger counts

```python
#Get the number of pickups per each passenger count
passenger_count = []
passenger_uniq = set(frame_with_durations_outliers_removed['passenger_count'].values)
for i in passenger_uniq:

passenger_count.append(frame_with_durations_outliers_removed[frame_with_durations_outliers_removed
['passenger_count'] == i].shape[0])

#sort the passenger pickups based on the passenger count
pass_count_pickups = [x for _,x in sorted(zip(passenger_uniq,passenger_count))]
pass_count = [int(y) for y,_ in sorted(zip(passenger_uniq,passenger_count))]

print (pass_count)
print (pass_count_pickups)
```

```
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
[5836, 8703909, 1764123, 515199, 246597, 687392, 448008, 6, 2, 4]
```
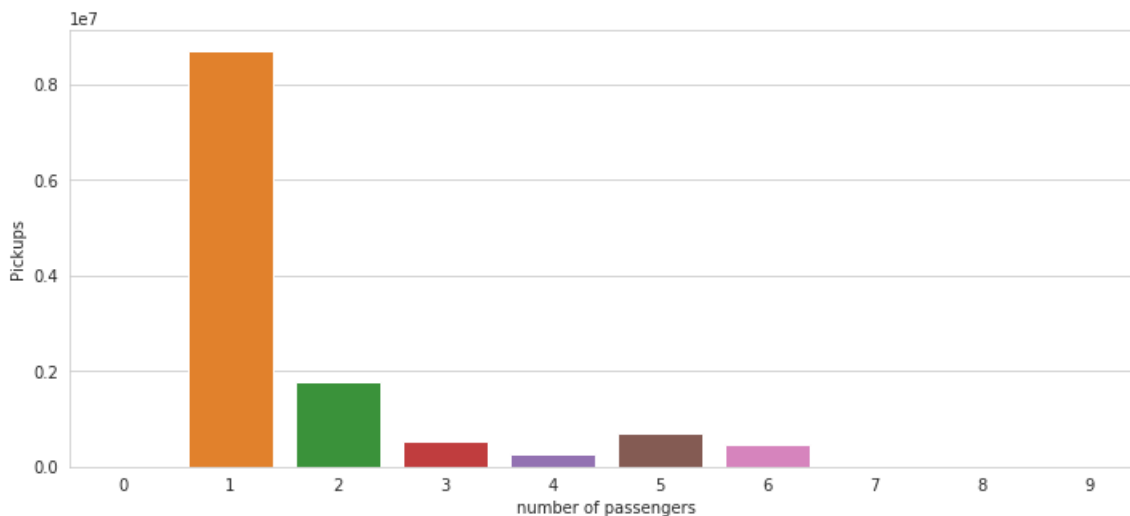
```python
sns.set_style("whitegrid")

fig, ax = plt.subplots()
```

```
fig, ax = plt.subplots()
fig.set_size_inches(12, 5)

ax = sns.barplot(x = np.array(pass_count) , y = np.array(pass_count_pickups))
ax.set(ylabel='Pickups',xlabel = 'number of passengers')
#sns.plt.title('Number of pickups for every type of passenger count')
plt.show()
```

```
for a,b in zip(pass_count,pass_count_pickups):
    print ("Pickups for passenger count {} : {}% of total pickups".format(int(a),round(float(b)*100
/sum(pass_count_pickups),4)))
```

```
Pickups for passenger count 0 : 0.0472% of total pickups
Pickups for passenger count 1 : 70.3569% of total pickups
Pickups for passenger count 2 : 14.2601% of total pickups
Pickups for passenger count 3 : 4.1645% of total pickups
Pickups for passenger count 4 : 1.9933% of total pickups
Pickups for passenger count 5 : 5.5564% of total pickups
Pickups for passenger count 6 : 3.6214% of total pickups
Pickups for passenger count 7 : 0.0% of total pickups
Pickups for passenger count 8 : 0.0% of total pickups
Pickups for passenger count 9 : 0.0% of total pickups
```
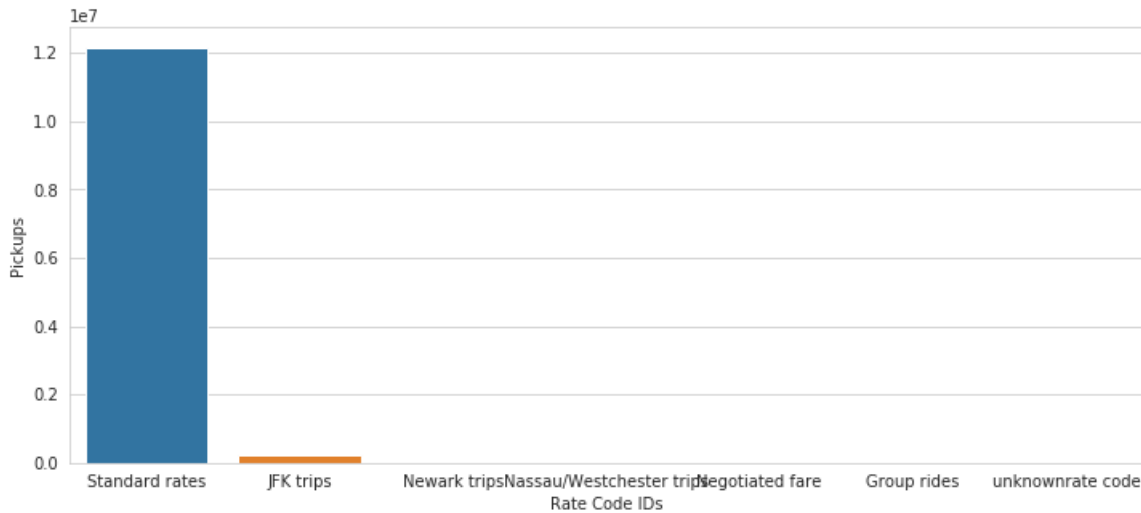
## 7. Rate Code ID

```
#Unique rate codes
rate_codes_uniq = set(frame_with_durations_outliers_removed['RateCodeID'].values)
rate_count = []
for i in rate_codes_uniq:
    rate_count.append(frame_with_durations_outliers_removed[frame_with_durations_outliers_removed['
RateCodeID'] == i].shape[0])
#sort the pickups based on the rate codes
rate_code_pickups = [x for _,x in sorted(zip(rate_codes_uniq,rate_count))]
rate_code = [int(y) for y,_ in sorted(zip(rate_codes_uniq,rate_count))]
Trips = ['Standard rates','JFK trips','Newark trips','Nassau/Westchester trips','Negotiated fare',
'Group rides','unknownrate code']
```

```
sns.set_style("whitegrid")

fig, ax = plt.subplots()
fig.set_size_inches(12, 5)

ax = sns.barplot(x = np.array(Trips) , y = np.array(rate_code_pickups))
ax.set(ylabel='Pickups',xlabel = 'Rate Code IDs')
plt.show()
```

```
for a,b in zip(Trips,rate_code_pickups):
    print ("Pickups for/with {} accounts {}% of total pickups".format(a,round(float(b)*100/sum(rate
_code_pickups),4)))
```

```
Pickups for/with Standard rates accounts 98.2535% of total pickups
Pickups for/with JFK trips accounts 1.637% of total pickups
Pickups for/with Newark trips accounts 0.0067% of total pickups
Pickups for/with Nassau/Westchester trips accounts 0.0131% of total pickups
Pickups for/with Negotiated fare accounts 0.0871% of total pickups
Pickups for/with Group rides accounts 0.0003% of total pickups
Pickups for/with unknownrate code accounts 0.0022% of total pickups
```

## 8. Store and Forward Flags

```
SWflag_uniq = set(frame_with_durations_outliers_removed['store_and_fwd_flag'].values)
SWflag_count = []
for i in SWflag_uniq :

SWflag_count.append(frame_with_durations_outliers_removed[frame_with_durations_outliers_removed['s
tore_and_fwd_flag'] == i].shape[0])
print (SWflag_count)
SWflag_uniq = list(SWflag_uniq)
print (SWflag_uniq)
```

```
[105726, 12265350]
['Y', 'N']
```

```
for a,b in zip(SWflag_uniq,SWflag_count):
    print ("Pickups for/with store and forward = {} accounts {}% of total pickups".format(a,round(f
loat(b)*100/sum(SWflag_count),4)))
```

```
Pickups for/with store and forward = Y accounts 0.8546% of total pickups
Pickups for/with store and forward = N accounts 99.1454% of total pickups
```
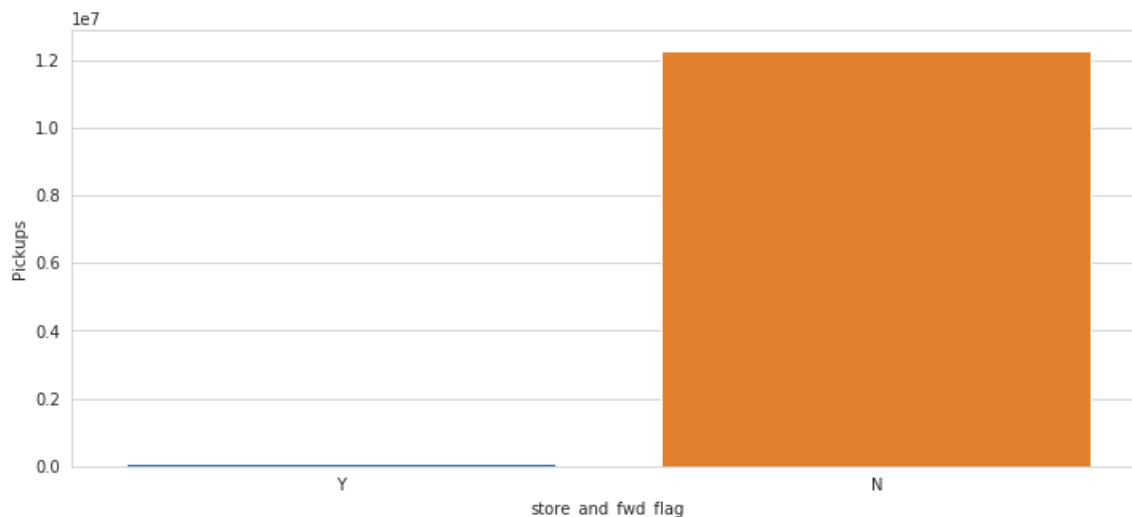
```
sns.set_style("whitegrid")

fig, ax = plt.subplots()
fig.set_size_inches(12, 5)

ax = sns.barplot(x = np.array(SWflag_uniq) , y = np.array(SWflag_count))
```

```
ax.set(ylabel='Pickups',xlabel = 'store_and_fwd_flag')
#sns.title('Number of pickups for store forward flags')
plt.show()
```



## 9. Payments Types

```
payment_uniq = set(frame_with_durations_outliers_removed['payment_type'].values)
payment_count = []
for i in payment_uniq :

payment_count.append(frame_with_durations_outliers_removed[frame_with_durations_outliers_removed['
payment_type'] == i].shape[0])
print (payment_count)
payment_uniq = list(payment_uniq)
print (payment_uniq)
```

```
[7674572, 4663218, 25562, 7723, 1]
[1, 2, 3, 4, 5]
```

```
pay_pickups = [x for _,x in sorted(zip(payment_uniq,payment_count))]
pay = [int(y) for y,_ in sorted(zip(payment_uniq,payment_count))]
pay_type = ['Credit card','Cash','No charge','Dispute','Unknow']
```

```
sns.set_style("whitegrid")

fig, ax = plt.subplots()
fig.set_size_inches(12, 5)

ax = sns.barplot(x = np.array(pay_type) , y = np.array(pay_pickups))
ax.set(ylabel='Pickups',xlabel = 'payment types')
#sns.plt.title('Number of pickups for each payment types')
plt.show()
```

```
for a,b in zip(pay_type,pay_pickups):
    print ("Pickups with payment types {} accounts to :{}% of total pickups".format(a,round(float(b
)*100/sum(pay_pickups),4)))
```
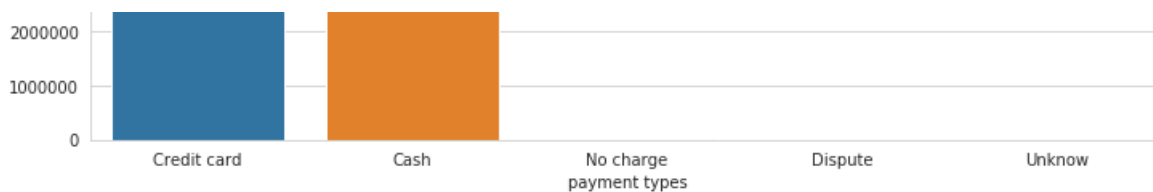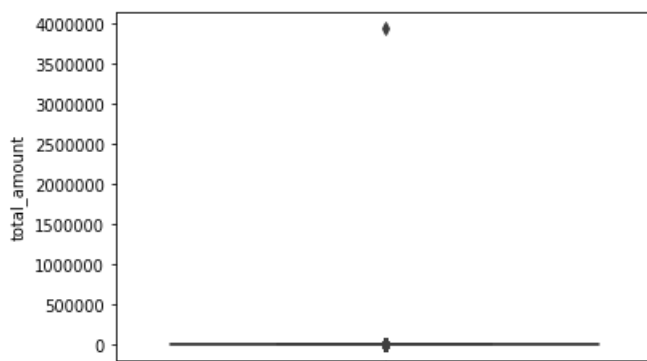
```
Pickups with payment types Credit card accounts to :62.0364% of total pickups
Pickups with payment types Cash accounts to :37.6945% of total pickups
Pickups with payment types No charge accounts to :0.2066% of total pickups
Pickups with payment types Dispute accounts to :0.0624% of total pickups
Pickups with payment types Unknow accounts to :0.0% of total pickups
```

## 10. Total Fare

In [56]:

```
# up to now we have removed the outliers based on trip durations, cab speeds, and trip distances
# lets try if there are any outliers in based on the total_amount
# box-plot showing outliers in fare
sns.boxplot(y="total_amount", data =frame_with_durations_modified)
plt.show()
```



In [57]:

```
#calculating total fare amount values at each percntile 0,10,20,30,40,50,60,70,80,90,100
for i in range(0,100,10):
    var = frame_with_durations_modified["total_amount"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print("100 percentile value is ",var[-1])
```

```
0 percentile value is -242.55
10 percentile value is 6.3
20 percentile value is 7.8
30 percentile value is 8.8
40 percentile value is 9.8
50 percentile value is 11.16
60 percentile value is 12.8
70 percentile value is 14.8
80 percentile value is 18.3
90 percentile value is 25.8
100 percentile value is  3950611.6
```

In [58]:

```
#calculating total fare amount values at each percntile 90,91,92,93,94,95,96,97,98,99,100
for i in range(90,100):
```

```
    var = frame_with_durations_modified["total_amount"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print("100 percentile value is ",var[-1])
```

```
90 percentile value is 25.8
91 percentile value is 27.3
92 percentile value is 29.3
93 percentile value is 31.8
94 percentile value is 34.8
95 percentile value is 38.53
96 percentile value is 42.6
97 percentile value is 48.13
98 percentile value is 58.13
99 percentile value is 66.13
100 percentile value is  3950611.6
```

In [59]:

```
#calculating total fare amount values at each percntile
99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100
for i in np.arange(0.0, 1.0, 0.1):
    var = frame_with_durations_modified["total_amount"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
print("100 percentile value is ",var[-1])
```

```
99.0 percentile value is 66.13
99.1 percentile value is 68.13
99.2 percentile value is 69.6
99.3 percentile value is 69.6
99.4 percentile value is 69.73
99.5 percentile value is 69.75
99.6 percentile value is 69.76
99.7 percentile value is 72.58
99.8 percentile value is 75.35
99.9 percentile value is 88.28
100 percentile value is  3950611.6
```

**Observation:-** As even the 99.9th percentile value doesnt look like an outlier,as there is not much difference between the 99.8th percentile and 99.9th percentile, we move on to do graphical analyis

In [60]:

```
#below plot shows us the fare values(sorted) to find a sharp increase to remove those values as ou
tliers
# plot the fare amount excluding last two values in sorted data
plt.plot(var[:-2])
plt.show()
```



In [61]:

```
# a very sharp increase in fare values can be seen
# plotting last three total fare values, and we can observe there is share increase in the values
plt plot(var[-3:])
```

```
plt.plot(var[ 3.])
plt.show()
```

```
#now looking at values not including the last two points we again find a drastic increase at aroun
d 1000 fare value
# we plot last 50 values excluding last two values
plt.plot(var[-50:-2])
plt.show()
```



## Plot pickup per day in jan 2015 month

In [63]:

```
def plot_2015(hour_pickup_month_2015,month):
    plt.figure(figsize=(8,4))
    hours_lis = [s for s in range(0,24)]
    plt.plot(hours_lis,hour_pickup_month_2015,'xkcd:magenta',label = 'average pickups per hour
'+month+' 2015')
    plt.plot(hours_lis,hour_pickup_month_2015, 'ro',  markersize=2)

    plt.xticks([s for s in range(0,24)])
    plt.xlabel('Hours of a day')
    plt.ylabel('Number of pickups')
    plt.title('Pickups for every hour for whole of month')
    plt.legend()
    plt.grid(True)
    plt.show()

# ??
def compute_pickups_per_day(month):

    time = month[['tpep_pickup_datetime','tpep_dropoff_datetime']].compute()
    time['tpep_pickup_datetime'] = pd.to_datetime(time['tpep_pickup_datetime'])
    time["pickup_day"] = time['tpep_pickup_datetime'].dt.strftime('%u').astype(int)
    time["pickup_hour"] = time['tpep_pickup_datetime'].dt.strftime('%H').astype(int)
    return time

def pickup_in_day(frame):
    hour_pickups = []
    temp = []
```
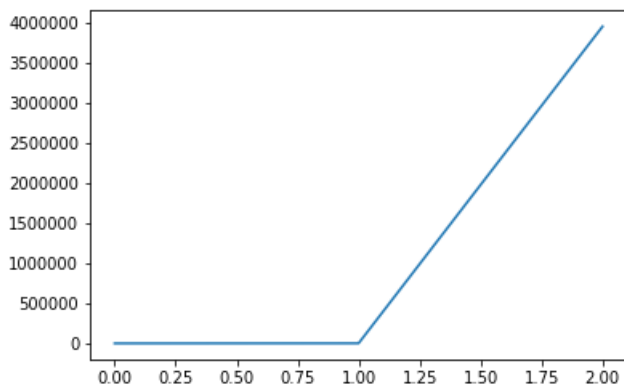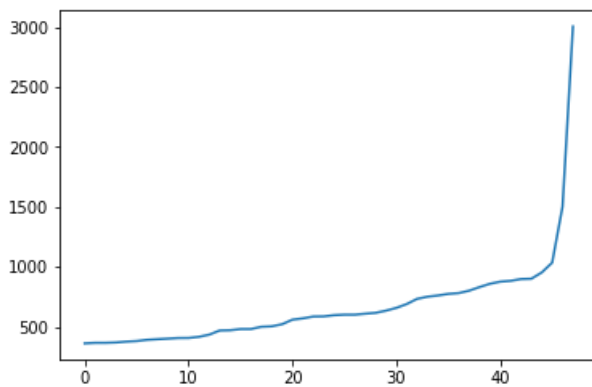
```
        temp = []
    for i in range(1,8):
        for j in range(0,24):
            temp.append(frame[(frame.pickup_day == i) & (frame.pickup_hour == j)].shape[0])
        hour_pickups.append(temp)
        temp = []

    days = ['Monday','Tuesday','Wednesday','Thursday','Friday','Saturday','Sunday']

    hour_pickup_month = []
    for j in range(0,24):
        hour_pickup_month.append(frame[frame.pickup_hour == j].shape[0])

    return hour_pickup_month
```

In [64]:

```
# jan month
times_of_day_dframe_jan_2015 = compute_pickups_per_day(month)

hour_pickup_month_jan_2015 = pickup_in_day(times_of_day_dframe_jan_2015)

plot_2015(hour_pickup_month_jan_2015, 'jan')
```



## Data-preperation

### Clustering/Segmentation

In [67]:

```
#trying different cluster sizes to choose the right K in K-means
coords = frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']].values
neighbours=[]

def find_min_distance(cluster_centers, cluster_len):
    nice_points = 0
    wrong_points = 0
    less2 = []
    more2 = []
    min_dist=1000
    for i in range(0, cluster_len):
        nice_points = 0
        wrong_points = 0
        for j in range(0, cluster_len):
            if j!=i:
                distance = gpxpy.geo.haversine_distance(cluster_centers[i][0], cluster_centers[i][1
,cluster_centers[j][0], cluster_centers[j][1])
                min_dist = min(min_dist,distance/(1.60934*1000))
                if (distance/(1.60934*1000)) <= 2:
                    nice_points +=1
                else:
                    wrong_points += 1
        less2.append(nice_points)
        more2.append(wrong_points)
    neighbours.append(less2)
```

```
    print ("On choosing a cluster size of ",cluster_len,"\nAvg. Number of Clusters within the vici
nity (i.e. intercluster-distance < 2):", np.ceil(sum(less2)/len(less2)), "\nAvg. Number of
Clusters outside the vicinity (i.e. intercluster-distance > 2):", np.ceil(sum(more2)/len(more2)),"
\nMin inter-cluster distance = ",min_dist,"\n---")

def find_clusters(increment):
    kmeans = MiniBatchKMeans(n_clusters=increment, batch_size=10000,random_state=42).fit(coords)
    frame_with_durations_outliers_removed['pickup_cluster'] =
kmeans.predict(frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']])
    cluster_centers = kmeans.cluster_centers_
    cluster_len = len(cluster_centers)
    return cluster_centers, cluster_len

# we need to choose number of clusters so that, there are more number of cluster regions
#that are close to any cluster center
# and make sure that the minimum inter cluster should not be very less
for increment in range(10, 100, 10):
    cluster_centers, cluster_len = find_clusters(increment)
    find_min_distance(cluster_centers, cluster_len)
```

```
On choosing a cluster size of  10
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 2.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 8.0
Min inter-cluster distance =  1.0945442325142543
---
On choosing a cluster size of  20
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 4.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 16.0
Min inter-cluster distance =  0.7131298007387813
---
On choosing a cluster size of  30
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 22.0
Min inter-cluster distance =  0.5185088176172206
---
On choosing a cluster size of  40
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 32.0
Min inter-cluster distance =  0.5069768450363973
---
On choosing a cluster size of  50
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 12.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 38.0
Min inter-cluster distance =  0.365363025983595
---
On choosing a cluster size of  60
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 14.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 46.0
Min inter-cluster distance =  0.34704283494187155
---
On choosing a cluster size of  70
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 16.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 54.0
Min inter-cluster distance =  0.30502203163244707
---
On choosing a cluster size of  80
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 18.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 62.0
Min inter-cluster distance =  0.29220324531738534
---
On choosing a cluster size of  90
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 21.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 69.0
Min inter-cluster distance =  0.18257992857034985
---
```

**Inference:**

- The main objective was to find a optimal min. distance(Which roughly estimates to the radius of a cluster) between the clusters
  which we got was 40

```
In [68]:
```

```
# if check for the 50 clusters you can observe that there are two clusters with only 0.3 miles apa
rt from each other
# so we choose 40 clusters for solve the further problem

# Getting 40 clusters using the kmeans
kmeans = MiniBatchKMeans(n_clusters=40, batch_size=10000,random_state=0).fit(coords)
frame_with_durations_outliers_removed['pickup_cluster'] =
kmeans.predict(frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']])
```

**Plotting the cluster centers:**

In [69]:

```
# Plotting the cluster centers on OSM
cluster_centers = kmeans.cluster_centers_
cluster_len = len(cluster_centers)
map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')
for i in range(cluster_len):
    folium.Marker(list((cluster_centers[i][0],cluster_centers[i][1])), popup=(str(cluster_centers[i
][0])+str(cluster_centers[i][1]))).add_to(map_osm)
map_osm
```

Out[69]:

**Plotting the clusters:**

In [70]:

```
#Visualising the clusters on a map
def plot_clusters(frame):
    city_long_border = (-74.03, -73.75)
    city_lat_border = (40.63, 40.85)
    fig, ax = plt.subplots(ncols=1, nrows=1)
    ax.scatter(frame.pickup_longitude.values[:100000], frame.pickup_latitude.values[:100000], s=10,
lw=0,
               c=frame.pickup_cluster.values[:100000], cmap='tab20', alpha=0.2)
    ax.set_xlim(city_long_border)
    ax.set_ylim(city_lat_border)
    ax.set_xlabel('Longitude')
    ax.set_ylabel('Latitude')
    plt.show()

plot_clusters(frame_with_durations_outliers_removed)
```

## Time-binning

```python
#Refer:https://www.unixtimestamp.com/
# 1420070400 : 2015-01-01 00:00:00
# 1422748800 : 2015-02-01 00:00:00
# 1425168000 : 2015-03-01 00:00:00
# 1427846400 : 2015-04-01 00:00:00
# 1430438400 : 2015-05-01 00:00:00
# 1433116800 : 2015-06-01 00:00:00

# 1451606400 : 2016-01-01 00:00:00
# 1454284800 : 2016-02-01 00:00:00
# 1456790400 : 2016-03-01 00:00:00
# 1459468800 : 2016-04-01 00:00:00
# 1462060800 : 2016-05-01 00:00:00
# 1464739200 : 2016-06-01 00:00:00

def add_pickup_bins(frame,month,year):
    unix_pickup_times=[i for i in frame['pickup_times'].values]
    unix_times = [[1420070400,1422748800,1425168000,1427846400,1430438400,1433116800],\
                  [1451606400,1454284800,1456790400,1459468800,1462060800,1464739200]]

    start_pickup_unix=unix_times[year-2015][month-1]
    # https://www.timeanddate.com/time/zones/est
    # (int((i-start_pickup_unix)/600)+33) : our unix time is in gmt to we are converting it to est
    tenminutewise_binned_unix_pickup_times=[(int((i-start_pickup_unix)/600)+33) for i in unix_picku
p_times]
    frame['pickup_bins'] = np.array(tenminutewise_binned_unix_pickup_times)
    return frame
```

```python
# clustering, making pickup bins and grouping by pickup cluster and pickup bins
frame_with_durations_outliers_removed['pickup_cluster'] =
kmeans.predict(frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']])
jan_2015_frame = add_pickup_bins(frame_with_durations_outliers_removed,1,2015)
jan_2015_groupby =
jan_2015_frame[['pickup_cluster','pickup_bins','trip_distance']].groupby(['pickup_cluster','pickup_
bins']).count()
```

```python
# we add two more columns 'pickup_cluster'(to which cluster it belogns to)
# and 'pickup_bins' (to which 10min intravel the trip belongs to)
jan_2015_frame.head()
```

| | passenger_count | trip_distance | pickup_longitude | pickup_latitude | dropoff_longitude | dropoff_latitude | total_amount | trip_times | picku |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1.59 | -73.993896 | 40.750111 | -73.974785 | 40.750618 | 17.05 | 18.050000 | 1.42 |

| 1 | passenger_count | trip_distance | pickup_longitude | pickup_latitude | dropoff_longitude | dropoff_latitude | total_amount | trip_times | pick |
|---|---|---|---|---|---|---|---|---|---|
| 2 | 1 | 1.80 | -73.963341 | 40.802788 | -73.951820 | 40.824413 | 10.80 | 10.050000 | 1.42( |
| 3 | 1 | 0.50 | -74.009087 | 40.713818 | -74.004326 | 40.719986 | 4.80 | 1.866667 | 1.42( |
| 4 | 1 | 3.00 | -73.971176 | 40.762428 | -74.004181 | 40.742653 | 16.30 | 19.316667 | 1.42( |

In [74]:

```python
# hear the trip_distance represents the number of pickups that are happend in that particular 10mi
n intravel
# this data frame has two indices
# primary index: pickup_cluster (cluster number)
# secondary index : pickup_bins (we devid whole months time into 10min intravels 24*31*60/10 =4464
bins)
jan_2015_groupby.head()
```

Out[74]:

| | | trip_distance |
|---|---|---|
| pickup_cluster | pickup_bins | |
| 0 | 33 | 104 |
| | 34 | 200 |
| | 35 | 208 |
| | 36 | 141 |
| | 37 | 155 |

In [75]:

```python
# upto now we cleaned data and prepared data for the month 2015,

# now do the same operations for months Jan, Feb, March of 2016
# 1. get the dataframe which inlcudes only required colums
# 2. adding trip times, speed, unix time stamp of pickup_time
# 4. remove the outliers based on trip_times, speed, trip_duration, total_amount
# 5. add pickup_cluster to each data point
# 6. add pickup_bin (index of 10min intravel to which that trip belongs to)
# 7. group by data, based on 'pickup_cluster' and 'pickuo_bin'

# Data Preparation for the months of Jan,Feb and March 2016
def datapreparation(month,kmeans,month_no,year_no):

    print ("Return with trip times..")

    frame_with_durations = return_with_trip_times(month)

    print ("Remove outliers..")
    frame_with_durations_outliers_removed = remove_outliers(frame_with_durations)

    print ("Estimating clusters..")
    frame_with_durations_outliers_removed['pickup_cluster'] =
kmeans.predict(frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']])
    #frame_with_durations_outliers_removed_2016['pickup_cluster'] =
kmeans.predict(frame_with_durations_outliers_removed_2016[['pickup_latitude',
'pickup_longitude']])

    print ("Final groupbying..")
    final_updated_frame = add_pickup_bins(frame_with_durations_outliers_removed,month_no,year_no)
    final_groupby_frame = final_updated_frame[['pickup_cluster','pickup_bins','trip_distance']].gro
upby(['pickup_cluster','pickup_bins']).count()

    return final_updated_frame,final_groupby_frame

month_jan_2016 = dd.read_csv('/floyd/input/nyc/yellow_tripdata_2016-01.csv')
month_feb_2016 = dd.read_csv('/floyd/input/nyc/yellow_tripdata_2016-02.csv')
month_mar_2016 = dd.read_csv('/floyd/input/nyc/yellow_tripdata_2016-03.csv')

jan_2016_frame,jan_2016_groupby = datapreparation(month_jan_2016,kmeans,1,2016)
feb_2016_frame,feb_2016_groupby = datapreparation(month_feb_2016,kmeans,2,2016)
mar_2016_frame,mar_2016_groupby = datapreparation(month_mar_2016,kmeans,3,2016)
```

```
Return with trip times..
Remove outliers..
Number of pickup records =  10906858
Number of outlier coordinates lying outside NY boundaries: 214677
Number of outliers from trip times analysis: 27190
Number of outliers from trip distance analysis: 79742
Number of outliers from speed analysis: 21047
Number of outliers from fare analysis: 4991
Total outliers removed 297784
---
Estimating clusters..
Final groupbying..
Return with trip times..
Remove outliers..
Number of pickup records =  11382049
Number of outlier coordinates lying outside NY boundaries: 223161
Number of outliers from trip times analysis: 27670
Number of outliers from trip distance analysis: 81902
Number of outliers from speed analysis: 22437
Number of outliers from fare analysis: 5476
Total outliers removed 308177
---
Estimating clusters..
Final groupbying..
Return with trip times..
Remove outliers..
Number of pickup records =  12210952
Number of outlier coordinates lying outside NY boundaries: 232444
Number of outliers from trip times analysis: 30868
Number of outliers from trip distance analysis: 87318
Number of outliers from speed analysis: 23889
Number of outliers from fare analysis: 5859
Total outliers removed 324635
---
Estimating clusters..
Final groupbying..
```

## Smoothing

In [79]:

```python
# Gets the unique bins where pickup values are present for each each reigion

# for each cluster region we will collect all the indices of 10min intravels in which the pickups
are happened
# we got an observation that there are some pickpbins that doesnt have any pickups
def return_unq_pickup_bins(frame):
    values = []
    for i in range(0,40):
        new = frame[frame['pickup_cluster'] == i]
        list_unq = list(set(new['pickup_bins']))
        list_unq.sort()
        values.append(list_unq)
    return values
```

In [80]:

```python
# for every month we get all indices of 10min intravels in which atleast one pickup got happened

#jan
jan_2015_unique = return_unq_pickup_bins(jan_2015_frame)
jan_2016_unique = return_unq_pickup_bins(jan_2016_frame)

#feb
feb_2016_unique = return_unq_pickup_bins(feb_2016_frame)

#march
mar_2016_unique = return_unq_pickup_bins(mar_2016_frame)
```

In [81]:

```python
# for each cluster number of 10min intravels with 0 pickups
for i in range(40):
    print("for the ",i,"th cluster number of 10min intavels with zero pickups: ",4464 -
len(set(jan_2015_unique[i])))
    print('-'*60)
```

```
for the  0 th cluster number of 10min intavels with zero pickups:  40
------------------------------------------------------------
for the  1 th cluster number of 10min intavels with zero pickups:  1985
------------------------------------------------------------
for the  2 th cluster number of 10min intavels with zero pickups:  29
------------------------------------------------------------
for the  3 th cluster number of 10min intavels with zero pickups:  354
------------------------------------------------------------
for the  4 th cluster number of 10min intavels with zero pickups:  37
------------------------------------------------------------
for the  5 th cluster number of 10min intavels with zero pickups:  153
------------------------------------------------------------
for the  6 th cluster number of 10min intavels with zero pickups:  34
------------------------------------------------------------
for the  7 th cluster number of 10min intavels with zero pickups:  34
------------------------------------------------------------
for the  8 th cluster number of 10min intavels with zero pickups:  117
------------------------------------------------------------
for the  9 th cluster number of 10min intavels with zero pickups:  40
------------------------------------------------------------
for the  10 th cluster number of 10min intavels with zero pickups:  25
------------------------------------------------------------
for the  11 th cluster number of 10min intavels with zero pickups:  44
------------------------------------------------------------
for the  12 th cluster number of 10min intavels with zero pickups:  42
------------------------------------------------------------
for the  13 th cluster number of 10min intavels with zero pickups:  28
------------------------------------------------------------
for the  14 th cluster number of 10min intavels with zero pickups:  26
------------------------------------------------------------
for the  15 th cluster number of 10min intavels with zero pickups:  31
------------------------------------------------------------
for the  16 th cluster number of 10min intavels with zero pickups:  40
------------------------------------------------------------
for the  17 th cluster number of 10min intavels with zero pickups:  58
------------------------------------------------------------
for the  18 th cluster number of 10min intavels with zero pickups:  1190
------------------------------------------------------------
for the  19 th cluster number of 10min intavels with zero pickups:  1357
------------------------------------------------------------
for the  20 th cluster number of 10min intavels with zero pickups:  53
------------------------------------------------------------
for the  21 th cluster number of 10min intavels with zero pickups:  29
------------------------------------------------------------
for the  22 th cluster number of 10min intavels with zero pickups:  29
------------------------------------------------------------
for the  23 th cluster number of 10min intavels with zero pickups:  163
------------------------------------------------------------
for the  24 th cluster number of 10min intavels with zero pickups:  35
------------------------------------------------------------
for the  25 th cluster number of 10min intavels with zero pickups:  41
------------------------------------------------------------
for the  26 th cluster number of 10min intavels with zero pickups:  31
------------------------------------------------------------
for the  27 th cluster number of 10min intavels with zero pickups:  214
------------------------------------------------------------
for the  28 th cluster number of 10min intavels with zero pickups:  36
------------------------------------------------------------
for the  29 th cluster number of 10min intavels with zero pickups:  41
------------------------------------------------------------
for the  30 th cluster number of 10min intavels with zero pickups:  1180
------------------------------------------------------------
for the  31 th cluster number of 10min intavels with zero pickups:  42
------------------------------------------------------------
for the  32 th cluster number of 10min intavels with zero pickups:  44
------------------------------------------------------------
for the  33 th cluster number of 10min intavels with zero pickups:  43
------------------------------------------------------------
for the  34 th cluster number of 10min intavels with zero pickups:  39
------------------------------------------------------------
for the  35 th cluster number of 10min intavels with zero pickups:  42
```

```
for the   35 th cluster number of 10min intavels with zero pickups:   42
-------------------------------------------------------------
for the   36 th cluster number of 10min intavels with zero pickups:   36
-------------------------------------------------------------
for the   37 th cluster number of 10min intavels with zero pickups:   321
-------------------------------------------------------------
for the   38 th cluster number of 10min intavels with zero pickups:   36
-------------------------------------------------------------
for the   39 th cluster number of 10min intavels with zero pickups:   43
-------------------------------------------------------------
```

there are two ways to fill up these values

- Fill the missing value with 0's
- Fill the missing values with the avg values
    - Case 1:(values missing at the start)
      Ex1: \_ \_ \_ x =>ceil(x/4), ceil(x/4), ceil(x/4), ceil(x/4)
      Ex2: \_ \_ x => ceil(x/3), ceil(x/3), ceil(x/3)
    - Case 2:(values missing in middle)
      Ex1: x \_ \_ y => ceil((x+y)/4), ceil((x+y)/4), ceil((x+y)/4), ceil((x+y)/4)
      Ex2: x \_ \_ \_ y => ceil((x+y)/5), ceil((x+y)/5), ceil((x+y)/5), ceil((x+y)/5), ceil((x+y)/5)
    - Case 3:(values missing at the end)
      Ex1: x \_ \_ \_ => ceil(x/4), ceil(x/4), ceil(x/4), ceil(x/4)
      Ex2: x \_ => ceil(x/2), ceil(x/2)

In [82]:

```python
# Fills a value of zero for every bin where no pickup data is present
# the count_values: number pickps that are happened in each region for each 10min intravel
# there wont be any value if there are no picksups.
# values: number of unique bins

# for every 10min intravel(pickup_bin) we will check it is there in our unique bin,
# if it is there we will add the count_values[index] to smoothed data
# if not we add 0 to the smoothed data
# we finally return smoothed data
def fill_missing(count_values,values):
    smoothed_regions=[]
    ind=0
    for r in range(0,40):
        smoothed_bins=[]
        for i in range(4464):
            if i in values[r]:
                smoothed_bins.append(count_values[ind])
                ind+=1
            else:
                smoothed_bins.append(0)
        smoothed_regions.extend(smoothed_bins)
    return smoothed_regions
```

In [83]:

```python
# Fills a value of zero for every bin where no pickup data is present
# the count_values: number pickps that are happened in each region for each 10min intravel
# there wont be any value if there are no picksups.
# values: number of unique bins

# for every 10min intravel(pickup_bin) we will check it is there in our unique bin,
# if it is there we will add the count_values[index] to smoothed data
# if not we add smoothed data (which is calculated based on the methods that are discussed in the
above markdown cell)
# we finally return smoothed data
def smoothing(count_values,values):
    smoothed_regions=[] # stores list of final smoothed values of each reigion
    ind=0
    repeat=0
    smoothed_value=0
    for r in range(0,40):
        smoothed_bins=[] #stores the final smoothed values
        repeat=0
        for i in range(4464):
            if repeat!=0: # prevents iteration for a value which is already visited/resolved
```

```
                    repeat-=1
                    continue
            if i in values[r]: #checks if the pickup-bin exists
                smoothed_bins.append(count_values[ind]) # appends the value of the pickup bin if it
exists
            else:
                if i!=0:
                    right_hand_limit=0
                    for j in range(i,4464):
                        if j not in values[r]: #searches for the left-limit or the pickup-bin
value which has a pickup value
                            continue
                        else:
                            right_hand_limit=j
                            break
                    if right_hand_limit==0:
                    #Case 1: When we have the last/last few values are found to be missing,hence we
have no right-limit here
                        smoothed_value=count_values[ind-1]*1.0/((4463-i)+2)*1.0

                        for j in range(i,4464):
                            smoothed_bins.append(math.ceil(smoothed_value))
                        smoothed_bins[i-1] = math.ceil(smoothed_value)
                        repeat=(4463-i)
                        ind-=1
                    else:
                    #Case 2: When we have the missing values between two known values
                        smoothed_value=(count_values[ind-1]+count_values[ind])*1.0/((right_hand_lim
t-i)+2)*1.0
                        for j in range(i,right_hand_limit+1):
                            smoothed_bins.append(math.ceil(smoothed_value))
                        smoothed_bins[i-1] = math.ceil(smoothed_value)
                        repeat=(right_hand_limit-i)
                else:
                    #Case 3: When we have the first/first few values are found to be missing,hence
we have no left-limit here
                    right_hand_limit=0
                    for j in range(i,4464):
                        if  j not in values[r]:
                            continue
                        else:
                            right_hand_limit=j
                            break
                    smoothed_value=count_values[ind]*1.0/((right_hand_limit-i)+1)*1.0
                    for j in range(i,right_hand_limit+1):
                            smoothed_bins.append(math.ceil(smoothed_value))
                    repeat=(right_hand_limit-i)
            ind+=1
        smoothed_regions.extend(smoothed_bins)
    return smoothed_regions
```

In [84]:

```
#Filling Missing values of Jan-2015 with 0
# here in jan_2015_groupby dataframe the trip_distance represents the number of pickups that are h
appened
jan_2015_fill = fill_missing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)

#Smoothing Missing values of Jan-2015
jan_2015_smooth = smoothing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)
```

In [85]:

```
def countZeros(num):
    count = 0
    for i in num:
        if i == 0:
            count += 1
    return count
```

In [86]:

```
print("Number of values filled with zero in zero fill data= "+str(countZeros(jan_2015_fill)))
```

Number of values filled with zero in zero fill data= 9451

In [87]:

```
print("Sanity check for number of zeros in smoothed data = "+str(countZeros(jan_2015_smooth)))
```

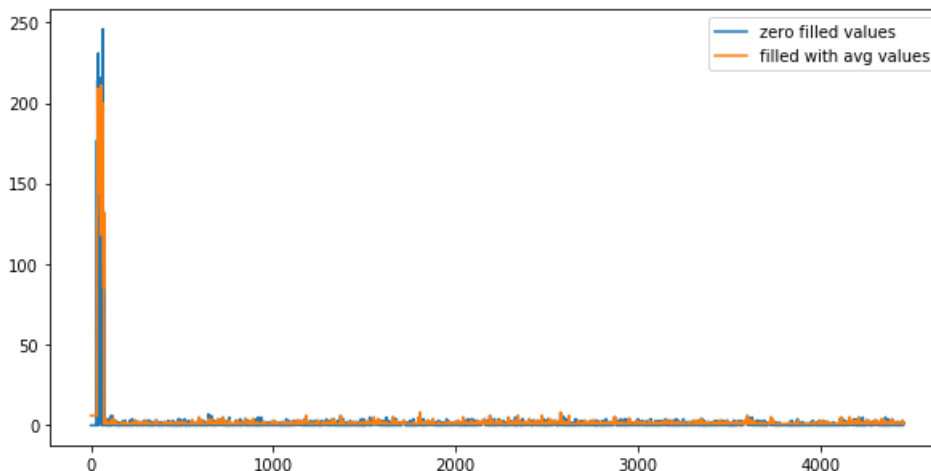Sanity check for number of zeros in smoothed data = 0

In [88]:

```
# number of 10min indices for jan 2015= 24*31*60/10 = 4464
# number of 10min indices for jan 2016 = 24*31*60/10 = 4464
# number of 10min indices for feb 2016 = 24*29*60/10 = 4176
# number of 10min indices for march 2016 = 24*30*60/10 = 4320
# for each cluster we will have 4464 values, therefore 40*4464 = 178560 (length of the
jan_2015_fill)
print("number of 10min intravels among all the clusters ",len(jan_2015_fill))
```

number of 10min intravels among all the clusters  178560

In [89]:

```
# Smoothing vs Filling
# sample plot that shows two variations of filling missing values
# we have taken the number of pickups for cluster region 2
plt.figure(figsize=(10,5))
plt.plot(jan_2015_fill[4464:8920], label="zero filled values")
plt.plot(jan_2015_smooth[4464:8920], label="filled with avg values")
plt.legend()
plt.show()
```



In [90]:

```
# why we choose, these methods and which method is used for which data?

# Ans: consider we have data of some month in 2015 jan 1st, 10 _ _ _ 20, i.e there are 10 pickups
that are happened in 1st
# 10st 10min intravel, 0 pickups happened in 2nd 10mins intravel, 0 pickups happened in 3rd 10min
intravel
# and 20 pickups happened in 4th 10min intravel.
# in fill_missing method we replace these values like 10, 0, 0, 20
# where as in smoothing method we replace these values as 6,6,6,6,6, if you can check the number o
f pickups
# that are happened in the first 40min are same in both cases, but if you can observe that we look
ing at the future values
# wheen you are using smoothing we are looking at the future number of pickups which might cause a
data leakage.

# so we use smoothing for jan 2015th data since it acts as our training data
# and we use simple fill_misssing method for 2016th data.
```

```
# Jan-2015 data is smoothed, Jan,Feb & March 2016 data missing values are filled with zero
jan_2015_smooth = smoothing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)
jan_2016_smooth = fill_missing(jan_2016_groupby['trip_distance'].values,jan_2016_unique)
feb_2016_smooth = fill_missing(feb_2016_groupby['trip_distance'].values,feb_2016_unique)
mar_2016_smooth = fill_missing(mar_2016_groupby['trip_distance'].values,mar_2016_unique)

# Making list of all the values of pickup data in every bin for a period of 3 months and storing t
hem region-wise
regions_cum = []

# a =[1,2,3]
# b = [2,3,4]
# a+b = [1, 2, 3, 2, 3, 4]

# number of 10min indices for jan 2015= 24*31*60/10 = 4464
# number of 10min indices for jan 2016 = 24*31*60/10 = 4464
# number of 10min indices for feb 2016 = 24*29*60/10 = 4176
# number of 10min indices for march 2016 = 24*31*60/10 = 4464
# regions_cum: it will contain 40 lists, each list will contain 4464+4176+4464 values which repres
ents the number of pickups
# that are happened for three months in 2016 data

for i in range(0,40):
    regions_cum.append(jan_2016_smooth[4464*i:4464*(i+1)]+feb_2016_smooth[4176*i:4176*(i+1)]+mar_20
16_smooth[4464*i:4464*(i+1)])

# print(len(regions_cum))
# 40
# print(len(regions_cum[0]))
# 13104
```
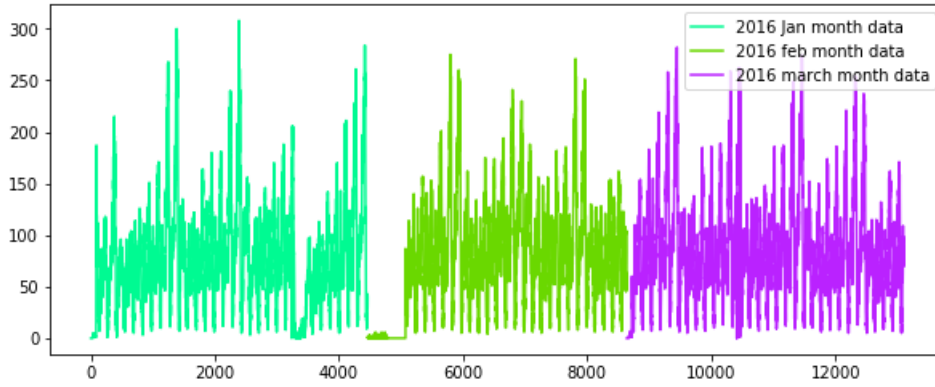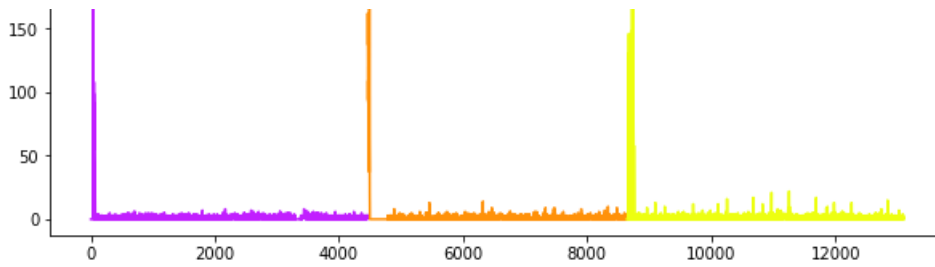
## Time series and Fourier Transforms

```
def uniqueish_color():
    """There're better ways to generate unique colors, but this isn't awful."""
    return plt.cm.gist_ncar(np.random.random())
first_x = list(range(0,4464))
second_x = list(range(4464,8640))
third_x = list(range(8640,13104))
for i in range(40):
    plt.figure(figsize=(10,4))
    plt.plot(first_x,regions_cum[i][:4464], color=uniqueish_color(), label='2016 Jan month data')
    plt.plot(second_x,regions_cum[i][4464:8640], color=uniqueish_color(), label='2016 feb month dat
a')
    plt.plot(third_x,regions_cum[i][8640:], color=uniqueish_color(), label='2016 march month data')
    plt.legend()
    plt.show()
```
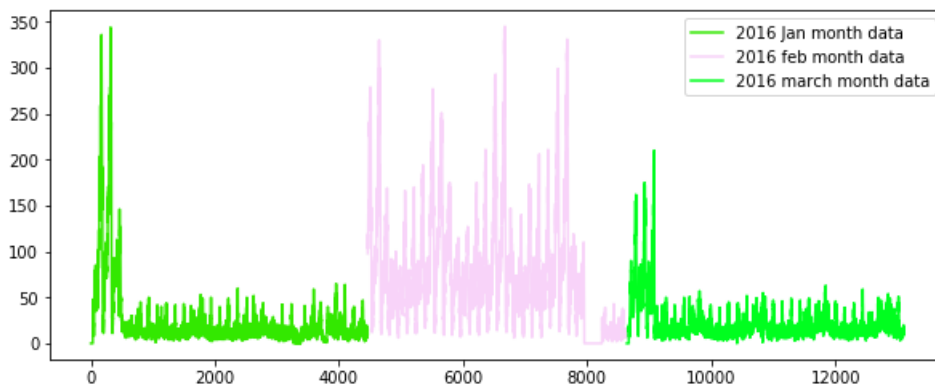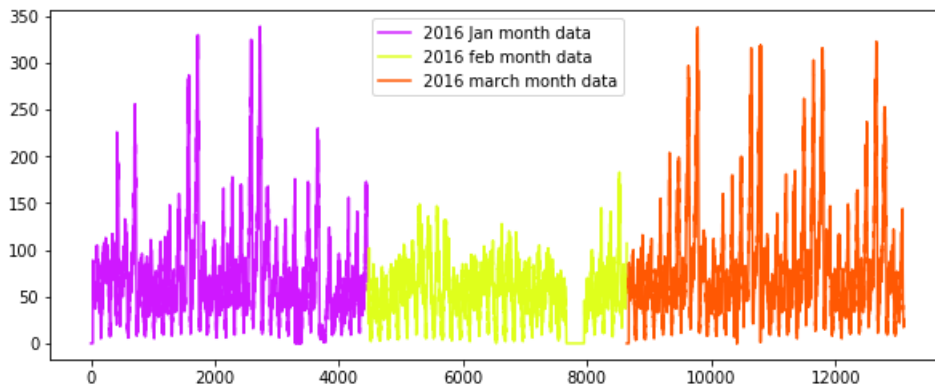
2016 Jan month data
2016 feb month data
2016 march month data



2016 Jan month data
2016 feb month data
2016 march month data



2016 Jan month data
2016 feb month data
2016 march month data



2016 Jan month data
2016 feb month data
2016 march month data

Legend: 2016 Jan month data, 2016 feb month data, 2016 march month data


Legend: 2016 Jan month data, 2016 feb month data, 2016 march month data


Legend: 2016 Jan month data, 2016 feb month data, 2016 march month data


Legend: 2016 Jan month data, 2016 feb month data, 2016 march month data


Legend: 2016 Jan month data, 2016 feb month data, 2016 march month data

Legend:
- 2016 Jan month data
- 2016 feb month data
- 2016 march month data



Legend:
- 2016 Jan month data
- 2016 feb month data
- 2016 march month data



Legend:
- 2016 Jan month data
- 2016 feb month data
- 2016 march month data



Legend:
- 2016 Jan month data
- 2016 feb month data
- 2016 march month data



Legend:
- 2016 Jan month data
- 2016 feb month data

Legend (second plot):
- 2016 Jan month data
- 2016 feb month data
- 2016 march month data



Legend (third plot):
- 2016 Jan month data
- 2016 feb month data
- 2016 march month data



Legend (fourth plot):
- 2016 Jan month data
- 2016 feb month data
- 2016 march month data



Legend (fifth plot):
- 2016 Jan month data
- 2016 feb month data
- 2016 march month data



Legend (sixth plot):
- 2016 Jan month data
- 2016 feb month data
- 2016 march month data

2016 Jan month data
2016 feb month data
2016 march month data



2016 Jan month data
2016 feb month data
2016 march month data



2016 Jan month data
2016 feb month data
2016 march month data



2016 Jan month data
2016 feb month data
2016 march month data

```
Y    = np.fft.fft(np.array(jan_2016_smooth)[0:4460])
freq = np.fft.fftfreq(4460, 1)
n = len(freq)
plt.figure()
plt.plot( freq[:int(n/2)], np.abs(Y)[:int(n/2)] )
plt.xlabel("Frequency")
plt.ylabel("Amplitude")
plt.show()
```

```
#Preparing the Dataframe only with x(i) values as jan-2015 data and y(i) values as jan-2016
ratios_jan = pd.DataFrame()
ratios_jan['Given']=jan_2015_smooth
ratios_jan['Prediction']=jan_2016_smooth
ratios_jan['Ratios']=ratios_jan['Prediction']*1.0/ratios_jan['Given']*1.0
```

# Modelling: Baseline Models

Now we get into modelling in order to forecast the pickup densities for the months of Jan, Feb and March of 2016 for which we are using multiple models with two variations

1. Using Ratios of the 2016 data to the 2015 data i.e $R_t = P_t^{2016}/P_t^{2015}$
2. Using Previous known values of the 2016 data itself to predict the future values

## Simple Moving Averages

The First Model used is the Moving Averages Model which uses the previous n values in order to predict the next value

Using Ratio Values - $R_t = (R_{t-1} + R_{t-2} + R_{t-3} \ldots R_{t-n})/n$

```python
def MA_R_Predictions(ratios,month):
    predicted_ratio=(ratios['Ratios'].values)[0]
    error=[]
    predicted_values=[]
    window_size=3
    predicted_ratio_values=[]
    for i in range(0,4464*40):
        if i%4464==0:
            predicted_ratio_values.append(0)
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_ratio_values.append(predicted_ratio)
        predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
        error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(ratios['Prediction'].values)[i],1))))
        if i+1>=window_size:
            predicted_ratio=sum((ratios['Ratios'].values)[(i+1)-window_size:(i+1)])/window_size
        else:
            predicted_ratio=sum((ratios['Ratios'].values)[0:(i+1)])/(i+1)


    ratios['MA_R_Predicted'] = predicted_values
    ratios['MA_R_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err
```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 3 is optimal for getting the best results using Moving Averages using previous Ratio values therefore we get $R_t = (R_{t-1} + R_{t-2} + R_{t-3})/3$

Next we use the Moving averages of the 2016 values itself to predict the future value using $P_t = (P_{t-1} + P_{t-2} + P_{t-3}\ldots P_{t-n})/n$

In [96]:

```python
def MA_P_Predictions(ratios,month):
    predicted_value=(ratios['Prediction'].values)[0]
    error=[]
    predicted_values=[]
    window_size=1
    predicted_ratio_values=[]
    for i in range(0,4464*40):
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
        if i+1>=window_size:
            predicted_value=int(sum((ratios['Prediction'].values)[(i+1)-window_size:(i+1)])/window_size)
        else:
            predicted_value=int(sum((ratios['Prediction'].values)[0:(i+1)])/(i+1))

    ratios['MA_P_Predicted'] = predicted_values
    ratios['MA_P_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err
```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 1 is optimal for getting the best results using Moving Averages using previous 2016 values therefore we get $P_t = P_{t-1}$

## Weighted Moving Averages

The Moving Avergaes Model used gave equal importance to all the values in the window used, but we know intuitively that the future is more likely to be similar to the latest values and less similar to the older values. Weighted Averages converts this analogy into a mathematical relationship giving the highest weight while computing the averages to the latest previous value and decreasing weights to the subsequent older ones

Weighted Moving Averages using Ratio Values - $R_t = (N * R_{t-1} + (N-1) * R_{t-2} + (N-2) * R_{t-3}\ldots 1 * R_{t-n})/(N * (N+1)/2)$

```python
def WA_R_Predictions(ratios,month):
    predicted_ratio=(ratios['Ratios'].values)[0]
    alpha=0.5
    error=[]
    predicted_values=[]
    window_size=5
    predicted_ratio_values=[]
    for i in range(0,4464*40):
        if i%4464==0:
            predicted_ratio_values.append(0)
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_ratio_values.append(predicted_ratio)
        predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
        error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(ratios['Pred
iction'].values)[i],1))))
        if i+1>=window_size:
            sum_values=0
            sum_of_coeff=0
            for j in range(window_size,0,-1):
                sum_values += j*(ratios['Ratios'].values)[i-window_size+j]
                sum_of_coeff+=j
            predicted_ratio=sum_values/sum_of_coeff
        else:
            sum_values=0
            sum_of_coeff=0
            for j in range(i+1,0,-1):
                sum_values += j*(ratios['Ratios'].values)[j-1]
                sum_of_coeff+=j
            predicted_ratio=sum_values/sum_of_coeff

    ratios['WA_R_Predicted'] = predicted_values
    ratios['WA_R_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].v
alues))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err
```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 5 is optimal for getting the best results using Weighted Moving Averages using previous Ratio values therefore we get
$R_t = (5 * R_{t-1} + 4 * R_{t-2} + 3 * R_{t-3} + 2 * R_{t-4} + R_{t-5})/15$

Weighted Moving Averages using Previous 2016 Values - $P_t = (N * P_{t-1} + (N-1) * P_{t-2} + (N-2) * P_{t-3}. \ldots 1 * P_{t-n})/(N * (N+1)/2)$

```python
def WA_P_Predictions(ratios,month):
    predicted_value=(ratios['Prediction'].values)[0]
    error=[]
    predicted_values=[]
    window_size=2
    for i in range(0,4464*40):
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
        if i+1>=window_size:
            sum_values=0
            sum_of_coeff=0
            for j in range(window_size,0,-1):
                sum_values += j*(ratios['Prediction'].values)[i-window_size+j]
                sum_of_coeff+=j
            predicted_value=int(sum_values/sum_of_coeff)

        else:
            sum_values=0
            sum_of_coeff=0
            for j in range(i+1,0,-1):
                sum_values += j*(ratios['Prediction'].values)[j-1]
                sum_of_coeff+=j
            predicted_value=int(sum_values/sum_of_coeff)

    ratios['WA_P_Predicted'] = predicted_values
```

```
    ratios['WA_P_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].v
alues))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err
```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 2 is optimal for getting the best results using Weighted Moving Averages using previous 2016 values therefore we get $P_t = (2 * P_{t-1} + P_{t-2})/3$

## Exponential Weighted Moving Averages

https://en.wikipedia.org/wiki/Moving_average#Exponential_moving_average Through weighted averaged we have satisfied the analogy of giving higher weights to the latest value and decreasing weights to the subsequent ones but we still do not know which is the correct weighting scheme as there are infinetly many possibilities in which we can assign weights in a non-increasing order and tune the the hyperparameter window-size. To simplify this process we use Exponential Moving Averages which is a more logical way towards assigning weights and at the same time also using an optimal window-size.

In exponential moving averages we use a single hyperparameter alpha $(\alpha)$ which is a value between 0 & 1 and based on the value of the hyperparameter alpha the weights and the window sizes are configured.
For eg. If $\alpha = 0.9$ then the number of days on which the value of the current iteration is based is~ $1/(1 - \alpha) = 10$ i.e. we consider values 10 days prior before we predict the value for the current iteration. Also the weights are assigned using $2/(N + 1) = 0.18$ ,where N = number of prior values being considered, hence from this it is implied that the first or latest value is assigned a weight of 0.18 which keeps exponentially decreasing for the subsequent values.

$$R_t' = \alpha * R_{t-1} + (1 - \alpha) * R_{t-1}'$$

In [99]:

```
def EA_R1_Predictions(ratios,month):
    predicted_ratio=(ratios['Ratios'].values)[0]
    alpha=0.6
    error=[]
    predicted_values=[]
    predicted_ratio_values=[]
    for i in range(0,4464*40):
        if i%4464==0:
            predicted_ratio_values.append(0)
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_ratio_values.append(predicted_ratio)
        predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
        error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(ratios['Pred
iction'].values)[i],1))))
        predicted_ratio = (alpha*predicted_ratio) + (1-alpha)*((ratios['Ratios'].values)[i])

    ratios['EA_R1_Predicted'] = predicted_values
    ratios['EA_R1_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].v
alues))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err
```

$$P_t' = \alpha * P_{t-1} + (1 - \alpha) * P_{t-1}'$$

In [100]:

```
def EA_P1_Predictions(ratios,month):
    predicted_value= (ratios['Prediction'].values)[0]
    alpha=0.3
    error=[]
    predicted_values=[]
    for i in range(0,4464*40):
        if i%4464==0:
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
```

```
            error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
        predicted_value =int((alpha*predicted_value) + (1-alpha)*((ratios['Prediction'].values)[i])
)

    ratios['EA_P1_Predicted'] = predicted_values
    ratios['EA_P1_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].v
alues))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err
```

In [101]:

```
mean_err=[0]*10
median_err=[0]*10
ratios_jan,mean_err[0],median_err[0]=MA_R_Predictions(ratios_jan,'jan')
ratios_jan,mean_err[1],median_err[1]=MA_P_Predictions(ratios_jan,'jan')
ratios_jan,mean_err[2],median_err[2]=WA_R_Predictions(ratios_jan,'jan')
ratios_jan,mean_err[3],median_err[3]=WA_P_Predictions(ratios_jan,'jan')
ratios_jan,mean_err[4],median_err[4]=EA_R1_Predictions(ratios_jan,'jan')
ratios_jan,mean_err[5],median_err[5]=EA_P1_Predictions(ratios_jan,'jan')
```

## Comparison between baseline models

We have chosen our error metric for comparison between models as **MAPE (Mean Absolute Percentage Error)** so that we can
know that on an average how good is our model with predictions and **MSE (Mean Squared Error)** is also used so that we have a
clearer understanding as to how well our forecasting model performs with outliers so that we make sure that there is not much of a
error margin between our prediction and the actual value

In [102]:

```
print ("Error Metric Matrix (Forecasting Methods) - MAPE & MSE")
print ("---------------------------------------------------------------------------------------
----------")
print ("Moving Averages (Ratios) -                        MAPE: ",mean_err[0]," MSE: ",me
ian_err[0])
print ("Moving Averages (2016 Values) -                   MAPE: ",mean_err[1]," MSE: ",m
dian_err[1])
print ("---------------------------------------------------------------------------------------
----------")
print ("Weighted Moving Averages (Ratios) -               MAPE: ",mean_err[2]," MSE: ",me
dian_err[2])
print ("Weighted Moving Averages (2016 Values) -          MAPE: ",mean_err[3]," MSE: ",me
dian_err[3])
print ("---------------------------------------------------------------------------------------
----------")
print ("Exponential Moving Averages (Ratios) -         MAPE: ",mean_err[4]," MSE: ",media
n_err[4])
print ("Exponential Moving Averages (2016 Values) -    MAPE: ",mean_err[5]," MSE: ",media
n_err[5])
```

```
Error Metric Matrix (Forecasting Methods) - MAPE & MSE
-----------------------------------------------------------------------------------------------
-
Moving Averages (Ratios) -                        MAPE:  0.22785156353133512      MSE:  1196.
953853046595
Moving Averages (2016 Values) -                   MAPE:  0.15583458712025738      MSE:  254.
6309363799283
-----------------------------------------------------------------------------------------------
-
Weighted Moving Averages (Ratios) -               MAPE:  0.22706529144871415      MSE:
1053.083529345878
Weighted Moving Averages (2016 Values) -          MAPE:  0.1479482182992932       MSE:
224.81054547491038
-----------------------------------------------------------------------------------------------
-
Exponential Moving Averages (Ratios) -         MAPE:  0.2275474636148534       MSE:
1019.3071012544802
Exponential Moving Averages (2016 Values) -    MAPE:  0.1475381297798153       MSE:
222.35159610215055
```

**Plese Note:-** The above comparisons are made using Jan 2015 and Jan 2016 only

From the above matrix it is inferred that the best forecasting model for our prediction would be:- $P'_t = \alpha * P_{t-1} + (1 - \alpha) * P'_{t-1}$ i.e Exponential Moving Averages using 2016 Values

# Regression Models

## Train-Test Split

Before we start predictions using the tree based regression models we take 3 months of 2016 pickup data and split it such that for every region we have 70% data in train and 30% in test, ordered date-wise for every region

In [103]:

```python
# Preparing data to be split into train and test, The below prepares data in cumulative form which
will be later split into test and train
# number of 10min indices for jan 2015= 24*31*60/10 = 4464
# number of 10min indices for jan 2016 = 24*31*60/10 = 4464
# number of 10min indices for feb 2016 = 24*29*60/10 = 4176
# number of 10min indices for march 2016 = 24*31*60/10 = 4464
# regions_cum: it will contain 40 lists, each list will contain 4464+4176+4464 values which repres
ents the number of pickups
# that are happened for three months in 2016 data

# print(len(regions_cum))
# 40
# print(len(regions_cum[0]))
# 12960

# we take number of pickups that are happened in last 5 10min intravels
number_of_time_stamps = 5

# output varaible
# it is list of lists
# it will contain number of pickups 13099 for each cluster
output = []


# tsne_lat will contain 13104-5=13099 times lattitude of cluster center for every cluster
# Ex: [[cent_lat 13099times],[cent_lat 13099times], [cent_lat 13099times].... 40 lists]
# it is list of lists
tsne_lat = []


# tsne_lon will contain 13104-5=13099 times logitude of cluster center for every cluster
# Ex: [[cent_long 13099times],[cent_long 13099times], [cent_long 13099times].... 40 lists]
# it is list of lists
tsne_lon = []

# we will code each day
# sunday = 0, monday=1, tue = 2, wed=3, thur=4, fri=5,sat=6
# for every cluster we will be adding 13099 values, each value represent to which day of the week
that pickup bin belongs to
# it is list of lists
tsne_weekday = []

# its an numbpy array, of shape (523960, 5)
# each row corresponds to an entry in out data
# for the first row we will have [f0,f1,f2,f3,f4] fi=number of pickups happened in i+1th 10min int
ravel(bin)
# the second row will have [f1,f2,f3,f4,f5]
# the third row will have [f2,f3,f4,f5,f6]
# and so on...
tsne_feature = []


tsne_feature = [0]*number_of_time_stamps
for i in range(0,40):
    tsne_lat.append([kmeans.cluster_centers_[i][0]]*13099)
    tsne_lon.append([kmeans.cluster_centers_[i][1]]*13099)
    # jan 1st 2016 is thursday, so we start our day from 4: "(int(k/144))%7+4"
    # our prediction start from 5th 10min intravel since we need to have number of pickups that ar
e happened in last 5 pickup bins
```

```
e happened in last 5 pickup bins
    tsne_weekday.append([int(((int(k/144))%7+4)%7) for k in range(5,4464+4176+4464)])
    # regions_cum is a list of lists [[x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3..x13104],
[x1,x2,x3..x13104], [x1,x2,x3..x13104], .. 40 lsits]
    tsne_feature = np.vstack((tsne_feature, [regions_cum[i][r:r+number_of_time_stamps] for r in ran
ge(0,len(regions_cum[i])-number_of_time_stamps)]))
    output.append(regions_cum[i][5:])
tsne_feature = tsne_feature[1:]
```

```
len(tsne_lat[0])*len(tsne_lat) == tsne_feature.shape[0] == len(tsne_weekday)*len(tsne_weekday[0]) =
= 40*13099 == len(output)*len(output[0])
```

True

```
alpha=0.3
predicted_values=[]
predict_list = []
tsne_flat_exp_avg = []
fr_am_final = pd.DataFrame(columns= ['f_1','a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5'])
for r in range(0,40):
    YJan = np.fft.fft(np.array(regions_cum[r][0:4464]))
    freqJan = np.fft.fftfreq((4464), 1)


    YFeb = np.fft.fft(np.array(regions_cum[r])[4464:(4176+4464)])
    freqFeb = np.fft.fftfreq((4176), 1)


    YMar = np.fft.fft(np.array(regions_cum[r])[(4176+4464):(4176+4464+4464)])
    freqMar = np.fft.fftfreq((4464), 1)

    fr_am_jan = pd.DataFrame()
    fr_am_feb = pd.DataFrame()
    fr_am_mar = pd.DataFrame()

    fr_am_jan['Frequency'] = freqJan
    fr_am_jan['Amplitude'] = YJan
    fr_am_feb['Frequency'] = freqFeb
    fr_am_feb['Amplitude'] = YFeb
    fr_am_mar['Frequency'] = freqMar
    fr_am_mar['Amplitude'] = YMar

    fr_am_list_jan = []
    fr_am_list_feb = []
    fr_am_list_mar = []

    fr_am_jan_sorted = fr_am_jan.sort_values(by=["Amplitude"], ascending=False)[:5].reset_index(dro
p=True).T
    fr_am_feb_sorted = fr_am_feb.sort_values(by=["Amplitude"], ascending=False)[:5].reset_index(dro
p=True).T
    fr_am_mar_sorted = fr_am_mar.sort_values(by=["Amplitude"], ascending=False)[:5].reset_index(dro
p=True).T

    for i in range(0,5):
        fr_am_list_jan.append(float(fr_am_jan_sorted[i]['Frequency']))
        fr_am_list_jan.append(float(fr_am_jan_sorted[i]['Amplitude']))

        fr_am_list_feb.append(float(fr_am_feb_sorted[i]['Frequency']))
        fr_am_list_feb.append(float(fr_am_feb_sorted[i]['Amplitude']))

        fr_am_list_mar.append(float(fr_am_mar_sorted[i]['Frequency']))
        fr_am_list_mar.append(float(fr_am_mar_sorted[i]['Amplitude']))

    fr_am_new_jan = pd.DataFrame([fr_am_list_jan]*4464)
    fr_am_new_feb = pd.DataFrame([fr_am_list_feb]*4176)
    fr_am_new_mar = pd.DataFrame([fr_am_list_mar]*4464)

    fr_am_new_jan.columns = ['f_1','a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5',]
    fr_am_new_feb.columns = ['f_1','a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5',]
```

```
        fr_am_new_mar.columns = ['f_1','a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5',]


    fr_am_final = fr_am_final.append(fr_am_new_jan, ignore_index=True)
    fr_am_final = fr_am_final.append(fr_am_new_feb, ignore_index=True)
    fr_am_final = fr_am_final.append(fr_am_new_mar, ignore_index=True)


    for i in range(0,13104):
        if i==0:
            predicted_value= regions_cum[r][0]
            predicted_values.append(0)
            continue
        predicted_values.append(predicted_value)
        predicted_value =int((alpha*predicted_value) + (1-alpha)*(regions_cum[r][i]))
    predict_list.append(predicted_values[5:])
    predicted_values=[]
fr_am_final.drop(['f_1'],axis=1,inplace=True)

fr_am_final = fr_am_final # (fr_am_final - fr_am_final.mean()) / (fr_am_final.max() -
fr_am_final.min())
fr_am_final = fr_am_final.fillna(0)
```

In [106]:

```
# train, test split : 70% 30% split
# Before we start predictions using the tree based regression models we take 3 months of 2016 pick
up data
# and split it such that for every region we have 70% data in train and 30% in test,
# ordered date-wise for every region
print("size of train data :", int(13099*0.7))
print("size of test data :", int(13099*0.3))
```

```
size of train data : 9169
size of test data : 3929
```

In [107]:

```
train_features =  [tsne_feature[i*13099:(13099*i+9169)] for i in range(0,40)]
test_features = [tsne_feature[(13099*(i))+9169:13099*(i+1)] for i in range(0,40)]
fr_am_final_train = pd.DataFrame(columns=['a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5'])
fr_am_final_test = pd.DataFrame(columns=['a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5'])
for i in range(0,40):
    fr_am_final_train = fr_am_final_train.append(fr_am_final[i*13099:(13099*i+9169)] )
fr_am_final_train.reset_index(inplace=True)
for i in range(0,40):
    fr_am_final_test = fr_am_final_test.append(fr_am_final[(13099*(i))+9169:13099*(i+1)])
fr_am_final_test.reset_index(inplace=True)
```

In [108]:

```
print("Number of data clusters",len(train_features), "Number of data points in trian data",
len(train_features[0]), "Each data point contains", len(train_features[0][0]),"features")
print("Number of data clusters",len(train_features), "Number of data points in test data",
len(test_features[0]), "Each data point contains", len(test_features[0][0]),"features")
```

```
Number of data clusters 40 Number of data points in trian data 9169 Each data point contains 5 fea
tures
Number of data clusters 40 Number of data points in test data 3930 Each data point contains 5 feat
ures
```

In [109]:

```
tsne_train_flat_lat = [i[:9169] for i in tsne_lat]
tsne_train_flat_lon = [i[:9169] for i in tsne_lon]
tsne_train_flat_weekday = [i[:9169] for i in tsne_weekday]
tsne_train_flat_output = [i[:9169] for i in output]
tsne_train_flat_exp_avg = [i[:9169] for i in predict_list]
```

In [110]:

```
tsne_test_flat_lat = [i[9169:] for i in tsne_lat]
tsne_test_flat_lon = [i[9169:] for i in tsne_lon]
tsne_test_flat_weekday = [i[9169:] for i in tsne_weekday]
tsne_test_flat_output = [i[9169:] for i in output]
tsne_test_flat_exp_avg = [i[9169:] for i in predict_list]
```

In [111]:

```
train_new_features = []
for i in range(0,40):
    train_new_features.extend(train_features[i])
test_new_features = []
for i in range(0,40):
    test_new_features.extend(test_features[i])
```

In [112]:

```
tsne_train_lat = sum(tsne_train_flat_lat, [])
tsne_train_lon = sum(tsne_train_flat_lon, [])
tsne_train_weekday = sum(tsne_train_flat_weekday, [])
tsne_train_output = sum(tsne_train_flat_output, [])
tsne_train_exp_avg = sum(tsne_train_flat_exp_avg,[])
```

In [113]:

```
tsne_test_lat = sum(tsne_test_flat_lat, [])
tsne_test_lon = sum(tsne_test_flat_lon, [])
tsne_test_weekday = sum(tsne_test_flat_weekday, [])
tsne_test_output = sum(tsne_test_flat_output, [])
tsne_test_exp_avg = sum(tsne_test_flat_exp_avg,[])
```

In [114]:

```
columns = ['ft_5','ft_4','ft_3','ft_2','ft_1']
df_train = pd.DataFrame(data=train_new_features, columns=columns)
df_train['lat'] = tsne_train_lat
df_train['lon'] = tsne_train_lon
df_train['weekday'] = tsne_train_weekday
df_train['exp_avg'] = tsne_train_exp_avg

print(df_train.shape)
```

(366760, 9)

In [115]:

```
df_test = pd.DataFrame(data=test_new_features, columns=columns)
df_test['lat'] = tsne_test_lat
df_test['lon'] = tsne_test_lon
df_test['weekday'] = tsne_test_weekday
df_test['exp_avg'] = tsne_test_exp_avg
print(df_test.shape)
```

(157200, 9)

In [116]:

```
df_test.head()
```

Out[116]:

|   | ft_5 | ft_4 | ft_3 | ft_2 | ft_1 | lat | lon | weekday | exp_avg |
|---|------|------|------|------|------|-----------|------------|---------|---------|
| 0 | 143 | 145 | 119 | 113 | 124 | 40.776228 | -73.982119 | 4 | 121 |
| 1 | 145 | 119 | 113 | 124 | 121 | 40.776228 | -73.982119 | 4 | 120 |
| 2 | 119 | 113 | 124 | 121 | 131 | 40.776228 | -73.982119 | 4 | 127 |
| 3 | 113 | 124 | 121 | 131 | 110 | 40.776228 | -73.982119 | 4 | 115 |

| | ft_5 | ft_4 | ft_3 | ft_2 | ft_1 | lat | lon | weekday | exp_avg |
|---|---|---|---|---|---|---|---|---|---|
| 4 | 126 | 124 | 133 | 112 | 116 | 40.776228 | -73.982119 | 4 | 115 |

In [117]:

```
df_test_lm = pd.concat([df_test, fr_am_final_test], axis=1)
df_train_lm = pd.concat([df_train, fr_am_final_train], axis=1)

df_test_lm.head()
print(df_test.shape)
print(fr_am_final_test.shape)
```

```
(157200, 9)
(157200, 10)
```

In [118]:

```
df_test_lm.head()
```

Out[118]:

| | ft_5 | ft_4 | ft_3 | ft_2 | ft_1 | lat | lon | weekday | exp_avg | index | a_1 | f_2 | a_2 | f_3 | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 143 | 145 | 119 | 113 | 124 | 40.776228 | -73.982119 | 4 | 121 | 9169 | 387761.0 | 0.006944 | 91160.781939 | -0.006944 | 91160.7819 |
| 1 | 145 | 119 | 113 | 124 | 121 | 40.776228 | -73.982119 | 4 | 120 | 9170 | 387761.0 | 0.006944 | 91160.781939 | -0.006944 | 91160.7819 |
| 2 | 119 | 113 | 124 | 121 | 131 | 40.776228 | -73.982119 | 4 | 127 | 9171 | 387761.0 | 0.006944 | 91160.781939 | -0.006944 | 91160.7819 |
| 3 | 113 | 124 | 121 | 131 | 110 | 40.776228 | -73.982119 | 4 | 115 | 9172 | 387761.0 | 0.006944 | 91160.781939 | -0.006944 | 91160.7819 |
| 4 | 124 | 121 | 131 | 110 | 116 | 40.776228 | -73.982119 | 4 | 115 | 9173 | 387761.0 | 0.006944 | 91160.781939 | -0.006944 | 91160.7819 |

In [119]:

```
# specify parameters and distributions to sample from
def report(results, n_top=3):
    for i in range(1, n_top + 1):
        candidates = np.flatnonzero(results['rank_test_score'] == i)
        for candidate in candidates:
            print("Model with rank: {0}".format(i))
            print("Mean validation score: {0:.3f} (std: {1:.3f})".format(
                    results['mean_test_score'][candidate],
                    results['std_test_score'][candidate]))
            print("Parameters: {0}".format(results['params'][candidate]))
            print("")
```

## Using Linear Regression

In [120]:

```
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import GridSearchCV,RandomizedSearchCV
```

In [122]:

```
from sklearn.preprocessing import StandardScaler
```

In [123]:

```
#standardizing the data
train_std = StandardScaler().fit_transform(df_train)
test_std = StandardScaler().fit_transform(df_test)
```

In [124]:

```
#hyper-paramater tuning
lr_reg=LinearRegression()
parameters = {'fit_intercept':[True,False], 'normalize':[True,False], 'copy_X':[True, False]}
grid = GridSearchCV(lr_reg,parameters, cv=None)
grid.fit(train_std, tsne_train_output)

print(grid.best_estimator_)
print(grid.best_params_)
```

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=True)
{'copy_X': True, 'fit_intercept': True, 'normalize': True}
```

In [125]:

```
#applying linear regression with best hyper-parameter
lr_reg=LinearRegression(copy_X=True, fit_intercept=True, normalize=True,n_jobs=-1).fit(train_std, t
sne_train_output)
```

In [127]:

```
y_pred = lr_reg.predict(df_test)
lr_test_predictions = [round(value) for value in y_pred]
y_pred = lr_reg.predict(df_train)
lr_train_predictions = [round(value) for value in y_pred]
```

## Using Random Forest Regressor

In [131]:

```
#hyper-paramater tuning
from scipy.stats import randint as sp_randint
values = [10,20,30]
reg1 = RandomForestRegressor(n_jobs=-1)
hyper_parameter = {"n_estimators": values,"max_depth": [3, None],"max_features": ['sqrt' , 'log2' ]
,"min_samples_split": sp_randint(2, 11),
             "min_samples_leaf": sp_randint(1, 11)}
best_parameter = RandomizedSearchCV(reg1,param_distributions=hyper_parameter,
                                 n_iter=20,n_jobs=-1)
best_parameter.fit(train_std, tsne_train_output)
report(best_parameter.cv_results_)
```

```
Model with rank: 1
Mean validation score: 0.946 (std: 0.013)
Parameters: {'max_depth': None, 'max_features': 'log2', 'min_samples_leaf': 10,
'min_samples_split': 6, 'n_estimators': 30}

Model with rank: 2
Mean validation score: 0.946 (std: 0.013)
Parameters: {'max_depth': None, 'max_features': 'log2', 'min_samples_leaf': 10,
'min_samples_split': 7, 'n_estimators': 30}

Model with rank: 3
Mean validation score: 0.946 (std: 0.013)
Parameters: {'max_depth': None, 'max_features': 'log2', 'min_samples_leaf': 7,
'min_samples_split': 8, 'n_estimators': 30}
```

In [132]:

```
# Training a hyper-parameter tuned random forest regressor on our train data
# find more about LinearRegression function here http://scikit-
learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html
# -----------------------
# default paramters
# sklearn.ensemble.RandomForestRegressor(n_estimators=10, criterion='mse', max_depth=None, min_sam
ples_split=2,
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_
impurity_decrease=0.0,
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None,
verbose=0, warm_start=False)
```

```python
# some of methods of RandomForestRegressor()
# apply(X) Apply trees in the forest to X, return leaf indices.
# decision_path(X) Return the decision path in the forest
# fit(X, y[, sample_weight]) Build a forest of trees from the training set (X, y).
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict regression target for X.
# score(X, y[, sample_weight]) Returns the coefficient of determination R^2 of the prediction.
# ----------------------
# video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-
using-decision-trees-2/
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-en
sembles/
# ----------------------

regr1 =
RandomForestRegressor(max_features='log2',min_samples_leaf=10,min_samples_split=6,n_estimators=30,
n_jobs=-1)
regr1.fit(df_train, tsne_train_output)
```

Out[132]:

```
RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
           max_features='log2', max_leaf_nodes=None,
           min_impurity_decrease=0.0, min_impurity_split=None,
           min_samples_leaf=10, min_samples_split=6,
           min_weight_fraction_leaf=0.0, n_estimators=30, n_jobs=-1,
           oob_score=False, random_state=None, verbose=0, warm_start=False)
```

In [133]:

```python
# Predicting on test data using our trained random forest model

# the models regr1 is already hyper parameter tuned
# the parameters that we got above are found using grid search

y_pred = regr1.predict(df_test)
rndf_test_predictions = [round(value) for value in y_pred]
y_pred = regr1.predict(df_train)
rndf_train_predictions = [round(value) for value in y_pred]
```

In [134]:

```python
#feature importances based on analysis using random forest
print (df_train.columns)
print (regr1.feature_importances_)
```

```
Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'lat', 'lon', 'weekday',
       'exp_avg'],
      dtype='object')
[0.00748843 0.05074405 0.10547138 0.14064703 0.29347633 0.00170107
 0.00174276 0.00093264 0.39779631]
```

## Using XgBoost Regressor

In [139]:

```python
## Hyperparameter Tuning
x_model = xgb.XGBRegressor(n_jobs=-1)
param_dist = {"max_depth": [3, 4,5],"min_child_weight": [3, 4,5,6],"gamma":
[0,0.1,0.2],"colsample_bytree":[0.7,0.8,0.9],
              "nthread":[3,4,5]}

# run randomized search
n_iter_search = 20
random_search = RandomizedSearchCV(x_model, param_distributions=param_dist,
                                   n_iter=n_iter_search,n_jobs=-1)

random_search.fit(df_train, tsne_train_output)

report(random_search.cv_results_)
```

```
Model with rank: 1
Mean validation score: 0.947 (std: 0.013)
Parameters: {'nthread': 4, 'min_child_weight': 6, 'max_depth': 4, 'gamma': 0, 'colsample_bytree':
0.9}

Model with rank: 2
Mean validation score: 0.946 (std: 0.013)
Parameters: {'nthread': 5, 'min_child_weight': 5, 'max_depth': 5, 'gamma': 0.1,
'colsample_bytree': 0.8}

Model with rank: 3
Mean validation score: 0.946 (std: 0.013)
Parameters: {'nthread': 3, 'min_child_weight': 5, 'max_depth': 5, 'gamma': 0, 'colsample_bytree':
0.7}
```

In [141]:

```python
# Training a hyper-parameter tuned Xg-Boost regressor on our train data

# find more about XGBRegressor function here
http://xgboost.readthedocs.io/en/latest/python/python_api.html?#module-xgboost.sklearn
# -----------------------
# default paramters
# xgboost.XGBRegressor(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True,
objective='reg:linear',
# booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1, max_delta_step=0, subsamp
le=1, colsample_bytree=1,
# colsample_bylevel=1, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, base_score=0.5,
random_state=0, seed=None,
# missing=None, **kwargs)

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbo
se=True, xgb_model=None)
# get_params([deep]) Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is no
t thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----------------------
# video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-
using-decision-trees-2/
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-en
sembles/
# -----------------------

x_model = xgb.XGBRegressor(
 learning_rate =0.1,
 max_depth=4,
 min_child_weight=6,
 gamma=0,
 subsample=0.8,
 reg_alpha=200, reg_lambda=200,
 colsample_bytree=0.9,nthread=4)
x_model.fit(df_train, tsne_train_output)
```

Out[141]:

```
XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
       colsample_bytree=0.9, gamma=0, importance_type='gain',
       learning_rate=0.1, max_delta_step=0, max_depth=4,
       min_child_weight=6, missing=None, n_estimators=100, n_jobs=1,
       nthread=4, objective='reg:linear', random_state=0, reg_alpha=200,
       reg_lambda=200, scale_pos_weight=1, seed=None, silent=True,
       subsample=0.8)
```

In [142]:

```python
#predicting with our trained Xg-Boost regressor
# the models x_model is already hyper parameter tuned
# the parameters that we got above are found using grid search

y_pred = x_model.predict(df_test)
xgb_test_predictions = [round(value) for value in y_pred]
```

```
xgb_test_predictions = [round(value) for value in y_pred]
y_pred = x_model.predict(df_train)
xgb_train_predictions = [round(value) for value in y_pred]
```

## Calculating the error metric values for various models

In [144]:

```
train_mape=[]
test_mape=[]

train_mape.append((mean_absolute_error(tsne_train_output,df_train['ft_1'].values))/(sum(tsne_train_
output)/len(tsne_train_output)))
train_mape.append((mean_absolute_error(tsne_train_output,df_train['exp_avg'].values))/(sum(tsne_tra
in_output)/len(tsne_train_output)))
train_mape.append((mean_absolute_error(tsne_train_output,rndf_train_predictions))/(sum(tsne_train_o
utput)/len(tsne_train_output)))
train_mape.append((mean_absolute_error(tsne_train_output,
xgb_train_predictions))/(sum(tsne_train_output)/len(tsne_train_output)))
train_mape.append((mean_absolute_error(tsne_train_output,
lr_train_predictions))/(sum(tsne_train_output)/len(tsne_train_output)))

test_mape.append((mean_absolute_error(tsne_test_output, df_test['ft_1'].values))/(sum(tsne_test_out
put)/len(tsne_test_output)))
test_mape.append((mean_absolute_error(tsne_test_output,
df_test['exp_avg'].values))/(sum(tsne_test_output)/len(tsne_test_output)))
test_mape.append((mean_absolute_error(tsne_test_output,
rndf_test_predictions))/(sum(tsne_test_output)/len(tsne_test_output)))
test_mape.append((mean_absolute_error(tsne_test_output,
xgb_test_predictions))/(sum(tsne_test_output)/len(tsne_test_output)))
test_mape.append((mean_absolute_error(tsne_test_output,
lr_test_predictions))/(sum(tsne_test_output)/len(tsne_test_output)))
```

In [145]:

```
print ("Error Metric Matrix (Tree Based Regression Methods) -  MAPE")
print ("------------------------------------------------------------------------------------
----------")
print ("Baseline Model -                              Train: ",train_mape[0]," Test: ",test_map
[0])
print ("Exponential Averages Forecasting -           Train: ",train_mape[1]," Test: ",test_map
e[1])
print ("Linear Regression -                          Train: ",train_mape[3]," Test: ",test_mape
3])
print ("Random Forest Regression -                   Train: ",train_mape[2]," Test: ",test_mape
[2])
```

```
Error Metric Matrix (Tree Based Regression Methods) -  MAPE
------------------------------------------------------------------------------------------------
-
Baseline Model -                              Train:  0.14870666996426116      Test:
0.14225522601041551
Exponential Averages Forecasting -           Train:  0.14121603560900353      Test:
0.13490049942819257
Linear Regression -                          Train:  0.14013022153536772      Test:
0.13348453650930908
Random Forest Regression -                   Train:  0.12043079748405067      Test:
0.1326694865641158
```

## Error Metric Matrix

In [146]:

```
print ("Error Metric Matrix (Tree Based Regression Methods) -  MAPE")
print ("------------------------------------------------------------------------------------
----------")
print ("Baseline Model -                              Train: ",train_mape[0]," Test: ",test_map
[0])
print ("Exponential Averages Forecasting -           Train: ",train_mape[1]," Test: ",test_map
e[1])
print ("Linear Regression -                          Train: ",train_mape[4]," Test: ",test_mape
```

```
4])
print ("Random Forest Regression -                    Train: ",train_mape[2],"     Test: ",test_mape
[2])
print ("XgBoost Regression -                          Train: ",train_mape[3],"      Test: ",test_map
[3])
print ("-------------------------------------------------------------------------------------
----------")
```

```
Error Metric Matrix (Tree Based Regression Methods) -  MAPE
--------------------------------------------------------------------------------------------
-
Baseline Model -                             Train:  0.14870666996426116        Test:
0.14225522601041551
Exponential Averages Forecasting -           Train:  0.14121603560900353         Test:
0.13490049942819257
Linear Regression -                          Train:  62.44342852233198       Test:
62.326781185156506
Random Forest Regression -                   Train:  0.12043079748405067        Test:
0.1326694865641158
XgBoost Regression -                         Train:  0.14013022153536772        Test:  0.1334845365(
30908
--------------------------------------------------------------------------------------------
-
```

## Observation:

```
1) Our goal is predict number of pickups in New York City.
2) We are working on Jan 2015, Jan 2016, Feb 2016 and ,March 2016 Yellow taxi Data.
 ### Exploratory Data Analysis
3) For EDA we used Jan 2015 file.
4) Using Dask package in Dataframe we read the csv file.It contains 1274898 samples and 19 features
.
5) Using GRaphviz we have visualize the data in which circles are operation and rectangles are resu
lts.
 ### Data Cleaning
6) Using folium we plot pickup lattitude and pickup longitude, in which some points just outside th
e boundary.
7) Similarly we plot Dropout lattitude and Dropout longitude.
 ### Trip Durations
8) The timestamps are converted to unix so as to get duration(trip-time) & speed also pickup-times
in unix are used while binning
9) Append durations of trips and speed in miles/hr to a new dataframe
10) The skewed box plot shows us the presence of outliers
11) Calculate 0-100th percentile to find a the correct percentile value for removal of outliers
12) Remove data based on our analysis and TLC regulations and box-plot after removal of outliers
13) PDF plot shows that almost all of the trip durations are very less and approximately less than
100, extremely few trip durations are above 100.
14) Convert the values to log-values to chec for log-normal and plot pdf and Q-Q plot of log values
.
 ### Speed
15) Check for any outliers in the data after trip duration outliers removed and plot box-plot for s
peeds with outliers
16) Calculate speed values at each percntile
17) Here, 100th percentile value of a speed is 192 Million miles/hr which is (BIZZARE).
    Furthermore, 99.9th percentile value of speed is 45.31miles/hr. So, we are removing all the dat
a points where speed is greater than 45.31miles/hr.
18) Remove further outliers based on the 99.9th percentile value and plot Box plot of speed after r
emoving outliers and erronous points.
19) Plot PDF of trip speed and pdf of log speed
20) Calculate average speed of cabs in New - York.
 ### Trip Distance
21) Plot box-plot showing outliers in trip-distance values
22) Calculate trip distance values at each percntile and remove the outlier.
23) PLot BOx plot and Pdf for trip durations and log values of trip durations.
 ### Remove all outliers/erronous points.
24) Remove all outliers based on our univariate analysis above
 ### Vendor ID
25) Get unique vendors and get the number of trips for each vendor
26) PLot box plot and both the vendors share almost same number of pickups in JAN month
 ### Passenger counts
27) Get the number of pickups per each passenger count and plot box plot
 ### Rate Code ID
```

28) Get Unique rate codes **and** sort the pickups based on the rate codes **and** plot box plot
 ### *Store and Forward Flags*
29) Get Unique Store **and** Forward Flags **and** sort the pickups based on the Store **and** Forward Flags **and** plot box plot
 ### *Payments Types*
30) Get Unique Payments Types **and** sort the pickups based on the Payments Types **and** plot box plot
 ### *Total fare*
31) PLot box-plot **for** Total amount calculate total fare amount values at each percntile
 ### *Plot pickup per day in jan 2015 month*
32) pickups **and** distances travelled at what time of the day

In [ ]:

### *Data preparation*
### *Clustering/Segmentation*
33) Remove pickup lattitude **and** pickup longitude **and** find minimum distance using gpxpy
34) We need to choose number of clusters so that, there are more number of cluster regions that are close to any cluster center
35) The main objective was to find a optimal min. distance(Which roughly estimates to the radius of a cluster) between the clusters which we got was 40
36) Get 40 clusters using the kmeans
37) Plot the cluster centre using foliuma **and** visualize the clusters on a map.
### *Time Binning*
38) clustering, making pickup bins **and** grouping by pickup cluster **and** pickup bins
39) we add two more columns 'pickup_cluster'(to which cluster it belogns to) **and** 'pickup_bins' (to which 10min intravel the trip belongs to)
40) Data Preparation **for** the months of Jan,Feb **and** March 2016
### *Smoothing*
41) Gets the unique bins where pickup values are present **for** each each reigion
42) **for** every month we get all indices of 10min intravels **in** which atleast one pickup got happened
43) **for** each cluster number of 10min intravels **with** 0 pickups
44) There are two ways to fill up these values 1)Fill the missing value **with** 0's 2) Fill the missing values with the avg values
45) Plot smoothing verses filling
46) Jan-2015 data **is** smoothed, Jan,Feb & March 2016 data missing values are filled **with** zero
 ### *Time series and Fourier Transforms*
47) PLot time series **and** Fourier Transform
 ### *Modelling: Baseline Models*
48) Now we get into modelling **in** order to forecast the pickup densities **for** the months of Jan, Feb **and** March of 2016 **for** which we are using multiple models **with** two variations

1. Using Ratios of the 2016 data to the 2015 data i.e  $R_t = P2016_t / P2015_t$
2. Using Previous known values of the 2016 data itself to predict the future values
49) Compare baseline Simple Moving average,Weighted Moving Averages ,Exponential Moving Averages **for** ratios **and** previous values.

In [ ]:

### *Regression Models*
### *Train Test split*
50)Take 3 months of 2016 pickup data **and** split it such that **for** every region we have 70% data **in** train **and** 30% **in** test, ordered date-wise **for** every region
51) Calculate number of bins **for** 10 min indices **for** jan15 **and** jan, feb,march of 16.
52)Add top 5 frequency **and** amplitude of jan , feb **and** march
53)Before we start predictions using the tree based regression models we take 3 months of 2016 pickup data
**and** split it such that **for** every region we have 70% data **in** train **and** 30% **in** test,
### *Using LInear REgression*
54)After hyperparameter tuning train the model **and** get prediction values
### *Using Random Forest*
55)After hyperparameter tuning train the model **and** get prediction values
### *Using XGBoost*
54)After hyperparameter tuning train the model **and** get prediction values
###*Calculating the error metric values for various models*
55) Get MAPE **for** train **and** test **for** all models

**Observation: By observing best Error Matrix model is Random Forest Regressor.**