# NBS Data: Exploring 1

After the `load` and then `clean` of sample data, we now have a working ([Pandas](#) dataframe to explore. I'll poke around with the data here, and the once we know what type of analysis we are doing we can implement it in the `do` and `function` files.

The following fetches our data into a dataframe called `basicdf`:

```python
In [1]:  from pandas import *

         #Move to working directory
         import os
         os.chdir('C:\Users\Aman\Documents\Projects\NextBigSound\NBS-Data-Sample')

         #Set data directory
         datadir = 'data'
         path = os.path.join(os.getcwd(), datadir)

         #Retrieve dataframes from the hdf5 file
         store = HDFStore(os.path.join(path, 'NBSData.h5'))
         basicdf = store['basic']
         store.close()
```

basicdf has the following data:

```python
In [2]:  print basicdf.columns
         print 'df index is: ', basicdf.index.names

         array([Artist.ID, Day, Facebook.fans.d, Facebook.fans.t, Last.fm.plays.d,
                Last.fm.plays.t, MySpace.fans.d, MySpace.fans.t, MySpace.plays.d,
                MySpace.plays.t, Twitter.fans.d, Twitter.fans.t, Twitter.statuses.d,
                Twitter.statuses.t, YouTube.fans.d, YouTube.fans.t, YouTube.plays.d,
                YouTube.plays.t, Pandora.fans.d, Pandora.fans.t, Rdio.fans.d,
                Rdio.fans.t, Rdio.plays.d, Rdio.plays.t, SoundCloud.fans.d,
                SoundCloud.fans.t, SoundCloud.plays.d, SoundCloud.plays.t,
                iTunes.Album.Units.d, iTunes.Track.Units.d, Vevo.plays.d,
                Vevo.plays.t, SiteCatalyst.Visits.d, MediaGuide.Radio.Spins.d,
                Spotify.plays.d, Wikipedia.views.d], dtype=object)
         df index is:  ['artist', 'date']
```

We can slice the data conveniently.

The following shows us all the data for an artist with id=1035.

```python
In [3]:  basicdf.ix[1035,:]
```

```
Out[3]:  <class 'pandas.core.frame.DataFrame'>
         Index: 536 entries, 2010-08-31 17:00:00 to 2012-02-17 16:00:00
         Data columns:
         Artist.ID                  536   non-null values
         Day                        536   non-null values
         Facebook.fans.d            514   non-null values
         Facebook.fans.t            524   non-null values
         Last.fm.plays.d            187   non-null values
         Last.fm.plays.t            293   non-null values
         MySpace.fans.d             519   non-null values
         MySpace.fans.t             527   non-null values
         MySpace.plays.d            512   non-null values
         MySpace.plays.t            522   non-null values
         Twitter.fans.d             522   non-null values
         Twitter.fans.t             528   non-null values
```

```
        Twitter.statuses.d          171  non-null values
        Twitter.statuses.t          174  non-null values
        YouTube.fans.d              460  non-null values
        YouTube.fans.t              489  non-null values
        YouTube.plays.d             403  non-null values
        YouTube.plays.t             413  non-null values
        Pandora.fans.d              176  non-null values
        Pandora.fans.t              181  non-null values
        Rdio.fans.d                 114  non-null values
        Rdio.fans.t                 130  non-null values
        Rdio.plays.d                209  non-null values
        Rdio.plays.t                230  non-null values
        SoundCloud.fans.d           212  non-null values
        SoundCloud.fans.t           218  non-null values
        SoundCloud.plays.d          75   non-null values
        SoundCloud.plays.t          78   non-null values
        iTunes.Album.Units.d        534  non-null values
        iTunes.Track.Units.d        534  non-null values
        Vevo.plays.d                0    non-null values
        Vevo.plays.t                0    non-null values
        SiteCatalyst.Visits.d       536  non-null values
        MediaGuide.Radio.Spins.d    295  non-null values
        Spotify.plays.d             0    non-null values
        Wikipedia.views.d           532  non-null values
        dtypes: float64(36)
```

Apparently the sample dataset spans from (for this user, anyway) from the end of August last year to February 17th of this year 2012.

```
    2010-08-31 17:00:00 to 2012-02-17 16:00:00
```

We've converted the data into a timeseries. This gives us a lot of flexibility.

```
In [4]: i1035twitterstatus = basicdf.ix[1035,:]['Twitter.statuses.d']
        #stored into a new var for easy access
        i1035twitterstatus.plot()
        i1035twitterstatus.describe()
```

```
Out[4]: count    171.000000
        mean       4.929825
        std        8.980095
        min        0.000000
        25%        1.000000
        50%        2.000000
        75%        5.000000
        max       60.000000
```



```
In [5]: #zooming in on what happened just this year.
        yrstart, yrend = datetime(2012, 01,01), datetime(2012, 2, 17)
```
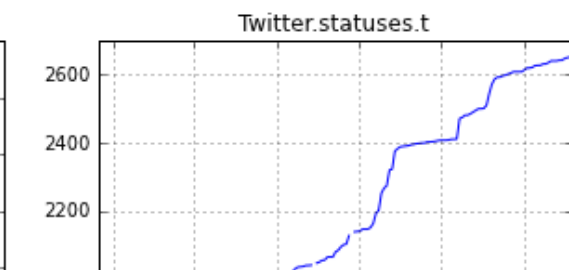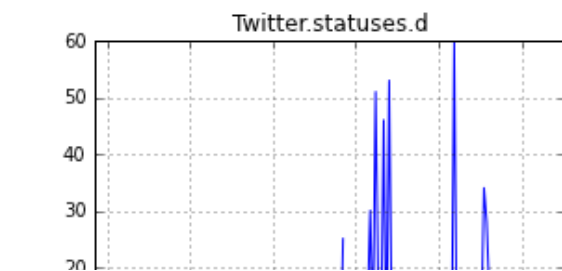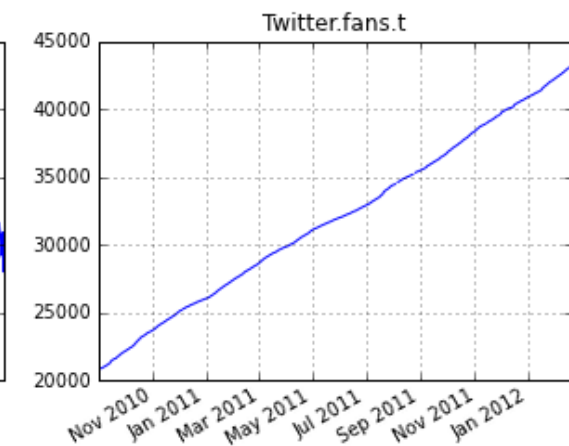
```
i1035twitterstatus[yrstart:yrend].plot()
```

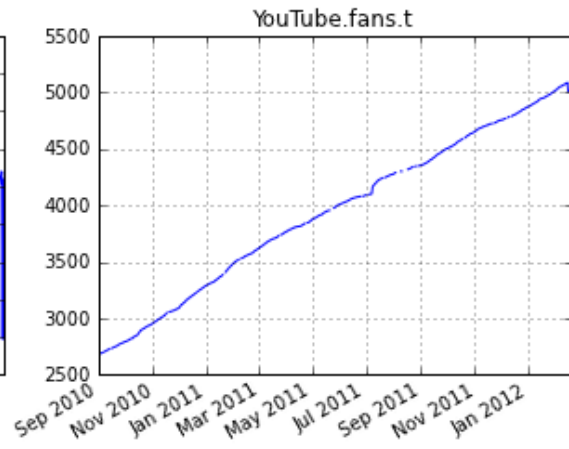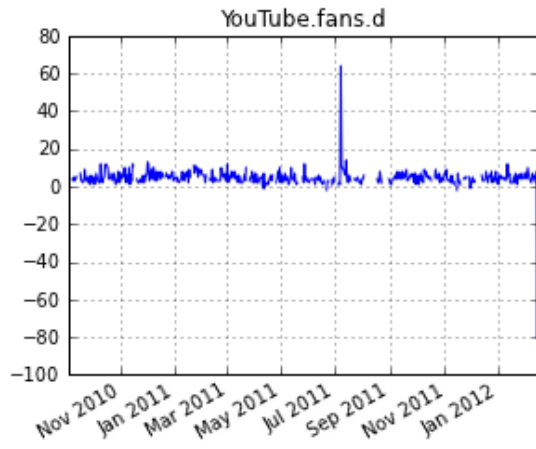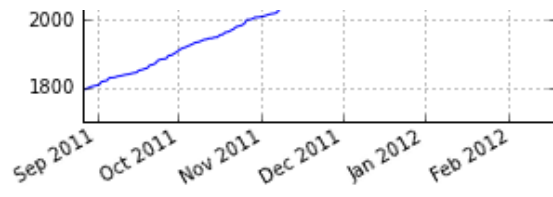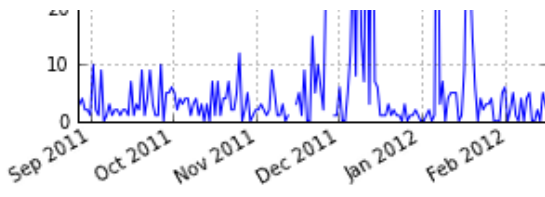Out[5]:   <matplotlib.axes.AxesSubplot at 0x8c851f0>



In [6]:
```
art1035 = basicdf.ix[1035,:]
del art1035['Artist.ID']
del art1035['Day']   #Don't need to plot this column

#Plotting parameters
ncols = len(art1035.columns)
plt.figure(figsize=(10,100))
subplots_adjust(hspace=0.4)

#Plot each column
for ii, col in enumerate(art1035):
    plt.subplot(20,2, ii+1)
    plt.title(str(col))
    if sum(art1035[col].notnull()) == 0:
        plt.text(0.5, 0.5, 'No Data points', fontsize=18, ha='center', va='top')
    else:
        art1035[col].plot()
```
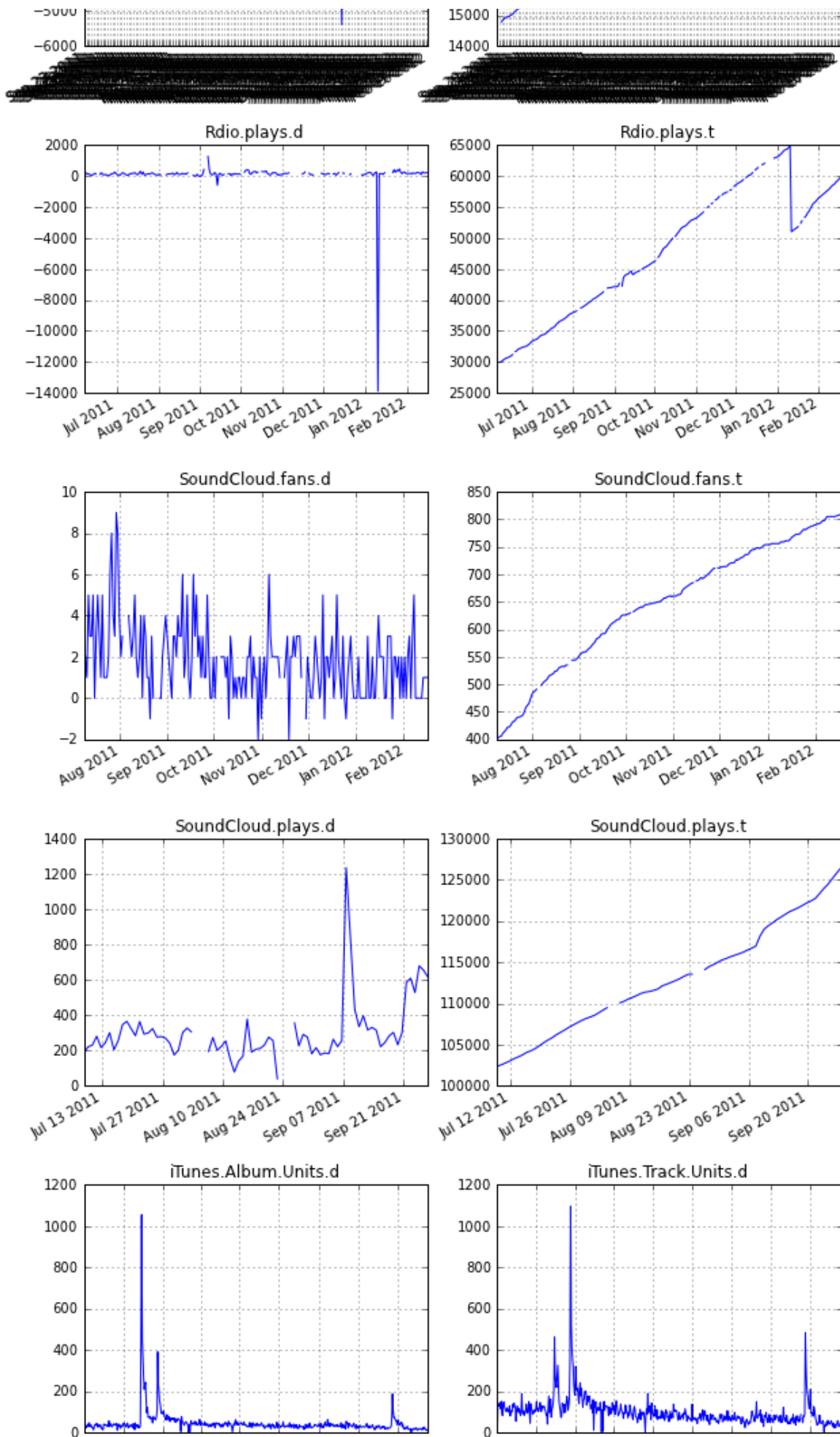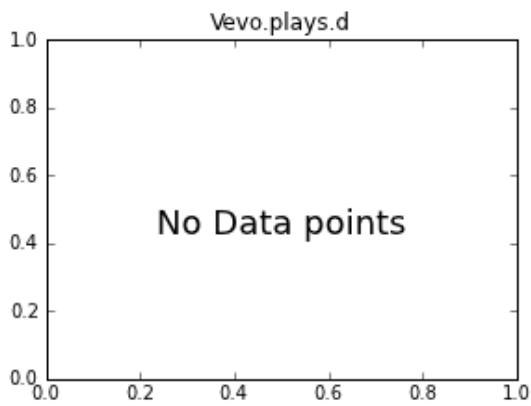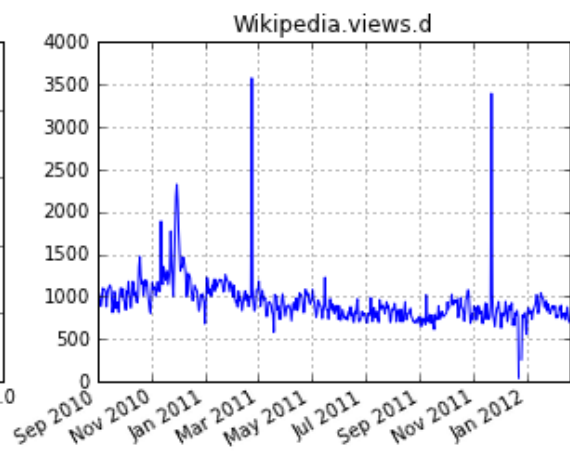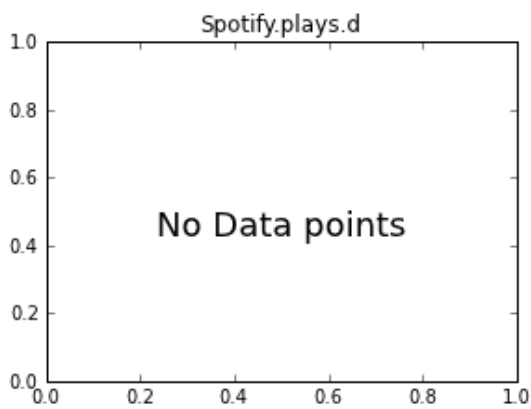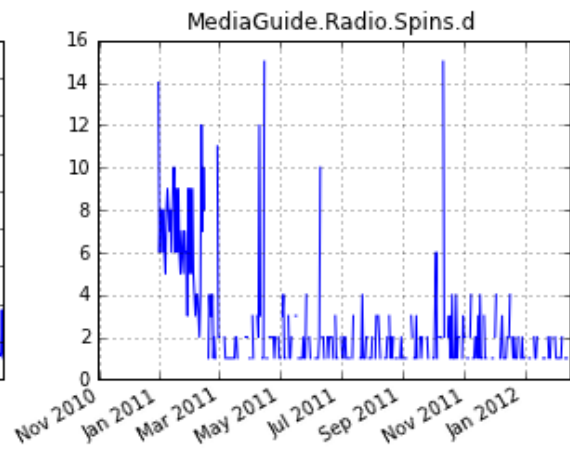
MySpace.fans.d

MySpace.fans.t

MySpace.plays.d

MySpace.plays.t

Twitter.fans.d

Twitter.fans.t

Twitter.statuses.d

Twitter.statuses.t

YouTube.fans.d

YouTube.fans.t

YouTube.plays.d

YouTube.plays.t

Pandora.fans.d

Pandora.fans.t

Rdio.fans.d

Rdio.fans.t

Rdio.plays.d



Rdio.plays.t



SoundCloud.fans.d



SoundCloud.fans.t



SoundCloud.plays.d



SoundCloud.plays.t



iTunes.Album.Units.d



iTunes.Track.Units.d